# IB031 Úvod do strojového učení Tomáš Brázdil

#### Course Info

Resources:

- Lectures & tutorials (the main source)
- Many books, few perfect for introductory level One relatively good, especially the first part:
  A. Géron. Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems. O'Reilly Media; 3rd edition, 2022
- (Almost) infinitely many online courses, tutorials, materials, etc.

#### Evaluation

The evaluation is composed of three parts:

- Mid-term exam: Written exam from the material of the first half of the semester.
- End-term exam: The "big" one containing everything from the semester (with possibly more stress in the second half).
- Projects: During tutorials, you will work on larger projects (in pairs).
- Each part contributes the following number of points:
  - Mid-term exam: 25
  - End-term exam: 50
  - Project: 25

To pass, you need to obtain at least 60 points.

#### Distinguishing Properties of the Course

- Introductory, prerequisites are held to a minimum
- Formal and precise: Be prepared for a complete and "mathematical" description of presented methods.

#### Distinguishing Properties of the Course

- Introductory, prerequisites are held to a minimum
- Formal and precise: Be prepared for a complete and "mathematical" description of presented methods.
- I assume that you have basic knowledge of
  - Elementary understanding of mathematical notation (operations on sets, logic, etc.)
  - ► Linear algebra: Vectors in ℝ<sup>n</sup>, operations on vectors (including the dot product). Geometric interpretation!
  - Calculus: Functions of multiple real variables, partial derivatives, basic differential calculus.
  - Probability: Notion of probability distribution, random variables/vectors, expectation.

#### What Is Machine Learning?

Machine learning is the science (and art) of programming computers so they can learn from data.

#### What Is Machine Learning?

Machine learning is the science (and art) of programming computers so they can learn from data.

Here is a slightly more general definition:

#### Arthur Samuel, 1959

Machine learning is the field of study that allows computers to learn without being explicitly programmed.

#### What Is Machine Learning?

Machine learning is the science (and art) of programming computers so they can learn from data.

Here is a slightly more general definition:

#### Arthur Samuel, 1959

Machine learning is the field of study that allows computers to learn without being explicitly programmed.

And a more engineering-oriented one:

#### Tom Mitchell, 1997

A computer program is said to learn from experience E concerning some task T and some performance measure P if its performance on T, as measured by P, improves with experience E.

#### Example

In the context of spam filtering:

- The task T is to flag spam in new emails.
- The experience E is represented by a set of emails labeled either spam or ham by hand (the training data).
- The performance measure P could be the accuracy, which is the ratio of the number of correctly classified emails and all emails.

There are many more performance measures; we will study the basic ones later.

#### Example

In the context of spam filtering:

- The task T is to flag spam in new emails.
- The experience E is represented by a set of emails labeled either spam or ham by hand (the training data).
- The performance measure P could be the accuracy, which is the ratio of the number of correctly classified emails and all emails.

There are many more performance measures; we will study the basic ones later.

In the context of housing price prediction:

- The task T is to predict prices of new houses based on their basic parameters (size, number of bathrooms, etc.)
- The experience E is represented by information about existing houses.
- The performance measure P could be, e.g., an absolute difference between the predicted and real price.

## Examples (cont.)

In the context of game playing:

- ▶ The task *T* is to play chess.
- The experience E is represented by a series of self-plays where the computer plays against itself.
- The performance measure P is winning/losing the game. Here, the trick is to spread the delayed and limited feedback about the result of the game throughout the individual decisions in the game.

## Examples (cont.)

In the context of game playing:

- ▶ The task *T* is to play chess.
- The experience E is represented by a series of self-plays where the computer plays against itself.
- The performance measure P is winning/losing the game. Here, the trick is to spread the delayed and limited feedback about the result of the game throughout the individual decisions in the game.

In the context of customer behavior:

- The task T is to group customers with similar shopping habits in an e-shop.
- The experience E consists of lists of items individual customers bought in the shop.
- The performance measure P? Measure how "nicely" the customers are grouped. (whether people with similar habits, as seen by humans, fall into the same group).

#### Comparison of Programming and Learning

How to code the spam filter?

- Examine what spam mails typically contain: Specific words ("Viagra"), sender's address, etc.
- ▶ Write down a rule-based system that detects specific features.
- Test the program on new emails and (most probably) go back to look for more spam features.



#### Comparison of Programming and Learning

The machine learning way:

- Study the problem and collect lots of emails, labeling them spam or ham.
- Train a machine learning model that reads an email and decides whether it's spam or ham.
- Test the model and (most probably) go back to collect more data and adjust the model.



#### ML Solutions are Adaptive

Spam filter: Authors of spam might and will adapt to your spam filter (possibly change the wording to pass through).

ML systems can be adjusted to new situations by retraining on new data (unless the data becomes ugly).



#### ML for Human Understanding

Spam filter: A trained system can be inspected for notorious spam features.

Some models allow direct inspection, such as decision trees or linear/logistic regression models.



#### Usage of Machine Learning

Machine learning suits various applications, especially where traditional methods fall short. Here are some areas where it excels:

- Solving complex problems where fine-tuning and rule-based solutions are inadequate.
- Tackling complex issues that resist traditional problem-solving approaches.
- Adapting to fluctuating environments through retraining on new data.
- Gaining insights from large and complex datasets.

In summary, machine learning offers innovative solutions and adaptability for today's complex and ever-changing problems, (sometimes) providing insights beyond the reach of traditional approaches.

## Types of Learning

There are main categories based on information available during the training:

- Supervised learning
- Unsupervised learning
- Semi-supervised learning
- Self-supervised learning
- Reinforcement learning

#### Supervised Learning



Labels are available for all input data.

#### Supervised Learning



Labels are available for all input data.

Typical supervised learning tasks are

 Classification where the aim is to classify inputs into (typically few) classes

(e.g., the spam filter where the classes are spam/ham)

 Regression where a numerical value is output for a given input (e.g., housing prices)

#### Unsupervised Learning





No labels are available for input data.

#### Unsupervised Learning





No labels are available for input data.

Typical unsupervised learning tasks are

Clustering where inputs are grouped according to their features

(e.g., clients of a bank grouped according to their age, wealth, etc.)

 Association where interesting relations and rules are discovered among the features of inputs

(e.g., market basket mining where associations between various types of goods are being learned from the behavior of customers)

 Dimensionality reduction reduce high-dimensional data to few dimensions (e.g., images to few image features)

#### Semi-Supervised Learning



Labels for some data.

#### Semi-Supervised Learning



Labels for some data.

For example, Medical data, where elaborate diagnosis is available only for some patients.

Combines supervised and unsupervised learning: e.g., clusters all data and labels the unlabeled inputs with the most common labels in their clusters.

#### Self-Supervised Learning



Generate labels from (unlabeled) inputs.

The goal is to learn typical features of the data.

It can be later modified to generate images, classify, etc.

## Reinforcement Learning



Learn from performing actions and getting feedback from environment.

- ChatGPT (and similar generative models)
  - The basis forms a generative language model, i.e., a text-generating model trained on texts in a self-supervised way
  - Currently extended to multimodal versions (text, image, sound)

- ChatGPT (and similar generative models)
  - The basis forms a generative language model, i.e., a text-generating model trained on texts in a self-supervised way
  - Currently extended to multimodal versions (text, image, sound)
- Machine translation, image captioning
  - Google translate, etc.
  - Typically (semi)-supervised learning,

- ChatGPT (and similar generative models)
  - The basis forms a generative language model, i.e., a text-generating model trained on texts in a self-supervised way
  - Currently extended to multimodal versions (text, image, sound)
- Machine translation, image captioning
  - Google translate, etc.
  - Typically (semi)-supervised learning,
- Various image recognition and processing tasks
  - In medicine where it is slowly making its way into hospitals as assistance tools
  - Automotive, advertising, quality control etc., etc., etc.

- ChatGPT (and similar generative models)
  - The basis forms a generative language model, i.e., a text-generating model trained on texts in a self-supervised way
  - Currently extended to multimodal versions (text, image, sound)
- Machine translation, image captioning
  - Google translate, etc.
  - Typically (semi)-supervised learning,
- Various image recognition and processing tasks
  - In medicine where it is slowly making its way into hospitals as assistance tools
  - Automotive, advertising, quality control etc., etc., etc.
- Science
  - Chemistry & biology: E.g., prediction of features of chemical compounds (Alpha-fold)

- ChatGPT (and similar generative models)
  - The basis forms a generative language model, i.e., a text-generating model trained on texts in a self-supervised way
  - Currently extended to multimodal versions (text, image, sound)
- Machine translation, image captioning
  - Google translate, etc.
  - Typically (semi)-supervised learning,
- Various image recognition and processing tasks
  - In medicine where it is slowly making its way into hospitals as assistance tools
  - Automotive, advertising, quality control etc., etc., etc.

#### Science

- Chemistry & biology: E.g., prediction of features of chemical compounds (Alpha-fold)
- ▶ Various "table" data processing in finance, management, etc.
  - Often straightforward methods (linear/logistic regression)
  - Essential but not fancy

- ChatGPT (and similar generative models)
  - The basis forms a generative language model, i.e., a text-generating model trained on texts in a self-supervised way
  - Currently extended to multimodal versions (text, image, sound)
- Machine translation, image captioning
  - Google translate, etc.
  - Typically (semi)-supervised learning,
- Various image recognition and processing tasks
  - In medicine where it is slowly making its way into hospitals as assistance tools
  - Automotive, advertising, quality control etc., etc., etc.
- Science
  - Chemistry & biology: E.g., prediction of features of chemical compounds (Alpha-fold)
- ▶ Various "table" data processing in finance, management, etc.
  - Often straightforward methods (linear/logistic regression)
  - Essential but not fancy
- Game playing: More fancy than useful, learning models beating humans in several difficult games.

#### ML in Context



## Supervised Learning

#### Example - Fruit Recognition

**The goal:** Create an automatic system for fruit recognition, concretely apple, lemon, and mandarin.

**Inputs:** Measures of *height* and *width* of each fruit.

Suppose we have a dataset of dimensions of several fruits labeled with the correct class.





Data



Use similarity to solve the problem.

	height	width	fruit
0	3.91	5.76	Mandarin
1	7.09	7.69	Apple
2	10.48	7.32	Lemon
3	9.21	7.20	Lemon
4	7.95	5.90	Lemon
5	7.62	7.51	Apple
6	7.95	5.32	Mandarin
7	4.69	6.19	Mandarin
8	7.50	5.99	Lemon
9	7.11	7.02	Apple
10	4.15	5.60	Mandarin
11	7.29	8.38	Apple
12	8.49	6.52	Lemon
13	7.44	7.89	Apple
14	7.86	7.60	Apple
15	3.93	6.12	Apple
16	4.40	5.90	Mandarin
17	5.50	4.50	Lemon
18	8.10	6.15	Lemon
19	8.69	5.82	Lemon
## **KNN** Classification

Given a new fruit. What is it?

Find five closest examples



Where is the machine learning?

## **KNN** Classification

Given a new fruit. What is it?

Find five closest examples

Among the five closest:

- M = 4 mandarins
- A = 1 apples
- $\blacktriangleright$  L = 0 lemons



Where is the machine learning?

## **KNN** Classification

Given a new fruit. What is it?

Find five closest examples

Among the five closest:

- M = 4 mandarins
- A = 1 apples
- L = 0 lemons

It is a mandarin!



Where is the machine learning?

Learning in Fruit Classification with KNN





#### Fruit Classification Algorithm

## **Input:** A fruit *F* with dimensions *height*, *width* **Output:** *mandarin*, *lemon*, *apple*

- 1: Find K examples  $\{E_1, \ldots, E_K\}$  in the dataset whose dimensions are closest to the dimensions of the fruit F
- 2: Count the number of examples of each class in  $\{E_1, \ldots, E_K\}$

$$M$$
 mandarins in  $\{E_1, \ldots, E_K\}$ 

L lemons in 
$$\{E_1,\ldots,E_K\}$$

A apples in  $\{E_1, \ldots, E_K\}$ 

- 3: if  $M \ge L$  and  $M \ge A$  then return mandarin
- 4: else if  $L \ge A$  then return lemon
- 5: else return apple
- 6: end if

Does it work?

#### Testing the Model for Fruit Classification

Consider a test set of new instances (K = 5, d is Euclidean):



Perfect classification of new data! Just deploy and sell!!

# K Nearest Neighbors

#### Learning and Inference

Two crucial components of machine learning are the following:



#### Training Data

<i>x</i> <sub>11</sub> Xo1	<i>x</i> <sub>12</sub>		x <sub>1n</sub>	$c_1$
×21 :	×22 :	•	^2n :	:
x <sub>p1</sub>	<i>x</i> <sub>p2</sub>	• • •	x <sub>pn</sub>	Ср

#### Formally, we define training dataset

$$\mathcal{T} = \{ (\vec{x}_k, c_k) \mid k = 1, \dots, p \}$$

Here each  $\vec{x}_k \in \mathbb{R}^n$  is an input vector and  $c_k \in C$  is the correct class.

height	width	fruit
4.0	6.5	Mandarin
4.47	7.13	Mandarin
6.49	7.0	Apple
7.51	5.01	Lemon
8.34	4.23	Lemon

$$\mathcal{T} = \{(4.0, 6.5), M), \\ (4.47, 7.13), M), \\ (6.49, 7.0), A), \\ \dots \}$$

#### KNN: Learning

Consider the training set:

$$\mathcal{T} = \{(\vec{x}_k, c_k) \mid k = 1, \dots, p\}$$

and memorize it exactly as it is.

Store in a table.

Possibly use a clever representation allowing fast computation of nearest neighbors such as KDTrees (out of the scope of this lecture).

Also,

- determine the number of neighbors  $K \in \mathbb{N}$ ,
- ▶ and the distance measure *d*.

#### Inference in KNN

Assume a KNN "trained" by memorizing  $\mathcal{T} = \{(\vec{x}_k, c_k) \in \mathbb{R}^n \times C \mid k = 1, ..., p\}$ , a constant  $K \in \mathbb{N}$  and a distance measure d.

For d, consider Euclidean distance, but different norms may also be used to define different distance measures.

#### Inference in KNN

Assume a KNN "trained" by memorizing  $\mathcal{T} = \{(\vec{x}_k, c_k) \in \mathbb{R}^n \times C \mid k = 1, ..., p\}$ , a constant  $K \in \mathbb{N}$  and a distance measure d.

For d, consider Euclidean distance, but different norms may also be used to define different distance measures.

**Input:** A vector  $\vec{z} = (z_1, ..., z_n) \in \mathbb{R}^n$ **Output:** A class from *C* 

1: Find K indices of examples  $X = \{i_1, \ldots, i_K\} \subseteq \{1, \ldots, p\}$  with minimum distance to  $\vec{z}$ , i.e., satisfying

 $\max \big\{ d(\vec{z}, \vec{x_\ell}) \mid \ell \in X \big\} \leq \min \big\{ d(\vec{z}, \vec{x_\ell}) \mid \ell \in \{1, \dots, p\} \smallsetminus X \big\}$ 

- 2: For every  $c \in C$  count the number #c of elements  $\ell$  in X such that  $c_\ell = c$
- 3: Return some

```
c_{max} \in \underset{c \in C}{\operatorname{arg\,max}} \# c
```

A class  $c_{max} \in C$  which maximizes #c.

#### The resulting model

What exactly constitutes the model? The model consists of

- The trained parameters: In this case the memorized training data.
- ► The *hyperparameters* set "from the outside": In this case, the number of neighbors *K* and the distance measure *d*.

#### The resulting model

What exactly constitutes the model? The model consists of

- The trained parameters: In this case the memorized training data.
- The hyperparameters set "from the outside": In this case, the number of neighbors K and the distance measure d.
  Note that different settings of K lead to different classifiers (for the same d):



... to get an efficient solution:

... to get an efficient solution:

- Deal with issues in the data
  - Data almost always comes in weird formats, with inconsistencies, missing values, wrong values, etc.
  - Data rarely have the ideal form for a given learning model.

We need to ingest, validate, and preprocess the data.

... to get an efficient solution:

- Deal with issues in the data
  - Data almost always comes in weird formats, with inconsistencies, missing values, wrong values, etc.
  - Data rarely have the ideal form for a given learning model.

We need to ingest, validate, and preprocess the data.

- Deal with issues in the model
  - In KNN, the training memorizes the example, but at least the K can be tuned.

We need to tune the model.

... to get an efficient solution:

- Deal with issues in the data
  - Data almost always comes in weird formats, with inconsistencies, missing values, wrong values, etc.
  - Data rarely have the ideal form for a given learning model.

We need to ingest, validate, and preprocess the data.

- Deal with issues in the model
  - In KNN, the training memorizes the example, but at least the K can be tuned.

We need to tune the model.

Deal with the wrong model by testing and validation in as realistic conditions as possible.

... to get an efficient solution:

Deal with issues in the data

- Data almost always comes in weird formats, with inconsistencies, missing values, wrong values, etc.
- Data rarely have the ideal form for a given learning model.

We need to ingest, validate, and preprocess the data.

- Deal with issues in the model
  - In KNN, the training memorizes the example, but at least the K can be tuned.

We need to tune the model.

- Deal with the wrong model by testing and validation in as realistic conditions as possible.
- Deal with deployment real-world application issues involving, e.g., implementation in embedded devices with limited resources.

#### Models Considered in This Course

Throughout this course, we will meet the following models:

- KNN (already did)
- Decision trees
- (Naive) Bayes classifier
- Clustering: K-means and hierarchical
- Linear and logistic regression
- Support Vector Machines (SVM)
- Kernel linear models
- Neural networks (light intro to feed-forward networks)
- Ensemble methods + random forests
- (maybe some reinforcement learning)

#### Models Considered in This Course

Throughout this course, we will meet the following models:

- KNN (already did)
- Decision trees
- (Naive) Bayes classifier
- Clustering: K-means and hierarchical
- Linear and logistic regression
- Support Vector Machines (SVM)
- Kernel linear models
- Neural networks (light intro to feed-forward networks)
- Ensemble methods + random forests
- (maybe some reinforcement learning)

... but first, let us see the whole machine learning pipeline.

### Machine Learning Pipeline



Always start with

Always start with

The problem formulation & understanding. For example, diagnosis of diabetes from medical records. What info is (possibly) sufficient for such a diagnosis?

Always start with

The problem formulation & understanding. For example, diagnosis of diabetes from medical records. What info is (possibly) sufficient for such a diagnosis?

Find data sources.

In our example, the sources are hospitals. It would be best to persuade them to give you the data and sign a contract.

Always start with

- The problem formulation & understanding. For example, diagnosis of diabetes from medical records. What info is (possibly) sufficient for such a diagnosis?
- Find data sources.

In our example, the sources are hospitals. It would be best to persuade them to give you the data and sign a contract.

#### Collect the data.

In our example, the data is possibly small (just tables with results of tests). But for other diagnoses, you may include huge amounts of data from MRI, CT, etc. Then, the collection itself might be a serious technical problem.

Always start with

- The problem formulation & understanding. For example, diagnosis of diabetes from medical records. What info is (possibly) sufficient for such a diagnosis?
- Find data sources.

In our example, the sources are hospitals. It would be best to persuade them to give you the data and sign a contract.

Collect the data.

In our example, the data is possibly small (just tables with results of tests). But for other diagnoses, you may include huge amounts of data from MRI, CT, etc. Then, the collection itself might be a serious technical problem.

Integrate data from various sources.

A serious diagnostic system must be trained/tested on data from many hospitals. You must blend the data from various sources (different formats, etc.).

For simple "toy" machine learning projects, you may fetch prepared datasets from various databases on the internet.

For simple "toy" machine learning projects, you may fetch prepared datasets from various databases on the internet.

The data should be stored in an identified location and versioned. You will probably keep adding data and training models on the ever-changing datasets. You have to be able to keep track of the changes and map training data to particular models.

Tools such as ML Flow or Weights & Biases might be helpful.

For simple "toy" machine learning projects, you may fetch prepared datasets from various databases on the internet.

The data should be stored in an identified location and versioned. You will probably keep adding data and training models on the ever-changing datasets. You have to be able to keep track of the changes and map training data to particular models.

Tools such as ML Flow or Weights & Biases might be helpful.

#### **Data Separation**

At this point, you should randomize the ordering of the data and select a test set to be used in model evaluation!

The test data are supposed to simulate the actual conditions, i.e., they should be "unseen".

For simple "toy" machine learning projects, you may fetch prepared datasets from various databases on the internet.

The data should be stored in an identified location and versioned. You will probably keep adding data and training models on the ever-changing datasets. You have to be able to keep track of the changes and map training data to particular models.

Tools such as ML Flow or Weights & Biases might be helpful.

#### **Data Separation**

At this point, you should randomize the ordering of the data and select a test set to be used in model evaluation!

The test data are supposed to simulate the actual conditions, i.e., they should be "unseen".

#### Data Exploration

Compute basic statistics to identify missing values, outliers, etc.

## Clean Data

The cleaning usually comprises the following steps:

- Fix or remove incorrect or corrupted values.
- Identify outliers and decide what to do with them.
  Outliers may harm some training methods and are not "representative".
  However, sometimes, they naturally belong to the dataset, and expert insight is needed.
- Fix formatting.

For example, the Date may be expressed in many ways, and a simple  $\ensuremath{\mathsf{Yes}}\xspace/No$  answer.

 Resolve missing values (by either removing the whole examples or imputing)

Many methods have been developed for missing values imputation. It is a susceptible issue because new values may strongly bias the model.

Remove duplicates.

The above steps often affect the training and need expertise in the application domain.

Later in this course, we will discuss techniques for data cleaning.

ID	Age	Income	Gender	Customer_Satisfaction
1	38	46641.356413713	nan	Unsatisfied
2	42	49129.0615585107	female	Neutral
3	18	119965.049731014	Male	nan
4	18	66828.0762224329	nan	very unsatisfied
5	58	57422.2721106762	female	very unsatisfied
6	28	59502.8174855665	Other	Satisfied
7	18	42659.6675768587	Other	Neutral
8	18	54019.1173206374	Other	Satisfied
9	40	25429.1604541137	female	Unsatisfied
10	21	15595.5862129548	Other	Satisfied
11	18	58094.2328460069	Other	very unsatisfied
12	18	39097.3278583155	female	Very Satisfied
13	30		Other	Satisfied
14	50	30617.3914472273	Female	Very Satisfied
15	18		nan	Neutral
16	34	39902.4430953214	male	nan
17	49	68381.6997683133	Female	Very Satisfied
18	33	44796.0962271524	Other	Very Satisfied
19	47	39218.9560738814	Female	very unsatisfied
20		14544.9226784447	Other	Satisfied

#### Prepare Data

Unlike cleaning, which is application-dependent, data preparation/transformation is model-dependent. This usually subsumes:

**Scaling**: Settings values of inputs to a similar range.

Some models, especially those utilizing distance, are sensitive to large differences between input sizes.

Encoding: Encode non-numeric data using real-valued vectors. Many models, especially those based on geometry, work only with numeric data. Non-numeric data such as Yes/No, Short/Medium/Long must be encoded appropriately.

 Binning or Discretization Convert continuous features into discrete bins to capture patterns in ranges.

**Comment:** Sometimes **Normalization**, that is changing the distribution of inputs to resemble the normal distribution, is mentioned. However, this step is typically not essential for machine learning itself. However, it is important to use statistical inference to test the significance of learned parameters.

#### Prepare Data

 Feature selection Throw out input features that are too "similar" to other features.

For example, if the temperature is measured both in Celsius and in Kelvin, keep one of them. The relationship can, of course, be a more complex (non-linear) correlation.

#### Prepare Data

 Feature selection Throw out input features that are too "similar" to other features.

For example, if the temperature is measured both in Celsius and in Kelvin, keep one of them. The relationship can, of course, be a more complex (non-linear) correlation.

▶ Dimensionality reduction Transforming data from  $\mathbb{R}^n$  to  $\mathbb{R}^m$  where  $m \ll n$ .

Growing dimension means growing difficulty of training for all models. Some models cease to work for high-dimensional data. The reduction typically searches for a few important characteristic features of inputs.
#### Prepare Data

 Feature selection Throw out input features that are too "similar" to other features.

For example, if the temperature is measured both in Celsius and in Kelvin, keep one of them. The relationship can, of course, be a more complex (non-linear) correlation.

▶ Dimensionality reduction Transforming data from ℝ<sup>n</sup> to ℝ<sup>m</sup> where m << n.</p>

Growing dimension means growing difficulty of training for all models. Some models cease to work for high-dimensional data. The reduction typically searches for a few important characteristic features of inputs.

 Feature aggregation Introducing new features using operations on the original ones.

We will see kernel transformations later in this course, allowing simple models to solve complex problems.

# Train Model

Now the dataset has been cleaned; we may train a model.

# Train Model

Now the dataset has been cleaned; we may train a model.

Before training, we should split the dataset into

- training dataset on which the model will learn
- validation dataset on which we fine-tune hyperparameters



The resulting model is obtained after several iterations of the above process.

## **Evaluate Model**

Here, we use the test set that we separated during data fetching. In some cases, a brand new test set can be generated. patients are examined regularly, creating new records continuously. In some cases, it is tough to obtain new data. For example, new expensive and difficult measurements are needed to obtain new data.

## **Evaluate Model**

Here, we use the test set that we separated during data fetching. In some cases, a brand new test set can be generated. patients are examined regularly, creating new records continuously. In some cases, it is tough to obtain new data. For example, new expensive and difficult measurements are needed to obtain new data.

Critical issue: Make sure that you are truly testing

exactly the whole inference process.

Often, just a model is tested, and the testing and production inference engines are separated. This leads to truly nasty errors in the production!

We will discuss various generic metrics helpful in measuring the quality of the resulting model.

Deployment of machine learning models is a complex question, application dependent.

The recently emerging area of MLOps is concerned with the engineering side of the model deployment.

Deployment of machine learning models is a complex question, application dependent.

The recently emerging area of MLOps is concerned with the engineering side of the model deployment.

From the technical point of view, the typical issues solved by  $\mathsf{ML}$  Ops teams are

- how to extract/process data in real-time
- how much storage is required
- how to store/collect model (and data) artifacts/predictions
- how to set up APIs, tools, and software environments
- What the period of predictions (instantaneous or batch predictions) should be
- how to set up hardware requirements (or cloud requirements for on-cloud environments) by the computational resources required
- how to set up a pipeline for continuous training and parameter tuning

From the user's point of view:

- How to get a sensible and valuable user output?
  - Al researchers will be satisfied with tons of running text in terminals.
  - "Normal" people need a graphical interface with understandable output.
  - Experts working in other domains typically demand speed and clarity at the extreme.

From the user's point of view:

- How to get a sensible and valuable user output?
  - Al researchers will be satisfied with tons of running text in terminals.
  - "Normal" people need a graphical interface with understandable output.
  - Experts working in other domains typically demand speed and clarity at the extreme.
- How do you persuade users that the AI is working for them?
  - Especially if safety is at stake, you need to have outstanding arguments and explanations ready for end-users
  - In many areas, the devices need to be certified (medicine, automotive) for ML-based systems.

This complex subject will be only touched on in this course.

#### Monitor, collect Data

Deployed machine learning models must be constantly monitored. Because of the influx of new data, ML models work in highly dynamic environments.

For example, an image-processing medical diagnostic model suddenly misdiagnosed a patient because a nurse marked the sample with a marker pen.

Every customer has a different infrastructure and may produce data slightly differently.

Data for retraining and improvement should be stored.

Also, many areas allow the *active learning* where users provide feedback for (continuous) retraining of the models.

# Data

You receive data from a medical researcher concerning a project that you are eager to work on.

You receive data from a medical researcher concerning a project that you are eager to work on.

The data consists of a 1000 lines table with five columns:

012	232	33.5	0	10.7
020	121	16.9	2	210.1
027	165	24.0	0	427.6

. . .

The aim is to predict the last field given the others.

You receive data from a medical researcher concerning a project that you are eager to work on.

The data consists of a 1000 lines table with five columns:

232	33.5	0	10.7
121	16.9	2	210.1
165	24.0	0	427.6
	232 121 165	23233.512116.916524.0	23233.5012116.9216524.00

. . .

The aim is to predict the last field given the others.

The medical researcher does not elaborate further on the data, but they seem to be pretty easy to work with, right?

You receive data from a medical researcher concerning a project that you are eager to work on.

The data consists of a 1000 lines table with five columns:

012	232	33.5	0	10.7
020	121	16.9	2	210.1
027	165	24.0	0	427.6

. . .

The aim is to predict the last field given the others.

The medical researcher does not elaborate further on the data, but they seem to be pretty easy to work with, right?

After a few days, you have trained a model that predicts numbers resembling the ones in the table.

You contact the medical researcher and discuss the results.

Researcher: So, you got the data for all the patients?

**Researcher:** So, you got the data for all the patients? **Data Miner:** Yes. I haven't had much time for analysis, but I do have a few interesting results.

**Researcher:** So, you got the data for all the patients? **Data Miner:** Yes. I haven't had much time for analysis, but I do have a few interesting results.

**Researcher:** Amazing. There were so many data issues with this set of patients that I couldn't do much.

**Researcher:** So, you got the data for all the patients? **Data Miner:** Yes. I haven't had much time for analysis, but I do have a few interesting results.

**Researcher:** Amazing. There were so many data issues with this set of patients that I couldn't do much.

Data Miner: Oh? I didn't hear about any possible problems.

**Researcher:** So, you got the data for all the patients? **Data Miner:** Yes. I haven't had much time for analysis, but I do have a few interesting results.

**Researcher:** Amazing. There were so many data issues with this set of patients that I couldn't do much.

**Data Miner:** Oh? I didn't hear about any possible problems. **Researcher:** Well, first, there is field 5, the variable we want to predict. It's common knowledge among people who analyze this type of data that results are better if you work with the log of the values, but I didn't discover this until later. Was it mentioned to you?

**Researcher:** So, you got the data for all the patients? **Data Miner:** Yes. I haven't had much time for analysis, but I do have a few interesting results.

**Researcher:** Amazing. There were so many data issues with this set of patients that I couldn't do much.

**Data Miner:** Oh? I didn't hear about any possible problems. **Researcher:** Well, first, there is field 5, the variable we want to predict. It's common knowledge among people who analyze this type of data that results are better if you work with the log of the values, but I didn't discover this until later. Was it mentioned to you?

Data Miner: No.

**Researcher:** But surely you heard about what happened to field 4? It's supposed to be measured on a scale from 1 to 10, with 0 indicating a missing value, but because of a data entry error, all 10's were changed into 0's. Unfortunately, since some of the patients have missing values for this field, it's impossible to say whether a 0 in this field is a real 0 or a 10. Quite a few of the records have that problem.

**Researcher:** But surely you heard about what happened to field 4? It's supposed to be measured on a scale from 1 to 10, with 0 indicating a missing value, but because of a data entry error, all 10's were changed into 0's. Unfortunately, since some of the patients have missing values for this field, it's impossible to say whether a 0 in this field is a real 0 or a 10. Quite a few of the records have that problem.

Data Miner: Interesting. Were there any other problems?

**Researcher:** But surely you heard about what happened to field 4? It's supposed to be measured on a scale from 1 to 10, with 0 indicating a missing value, but because of a data entry error, all 10's were changed into 0's. Unfortunately, since some of the patients have missing values for this field, it's impossible to say whether a 0 in this field is a real 0 or a 10. Quite a few of the records have that problem.

**Data Miner:** Interesting. Were there any other problems? **Researcher:** Yes, fields 2 and 3 are basically the same, but I assume that you probably noticed that.

**Researcher:** But surely you heard about what happened to field 4? It's supposed to be measured on a scale from 1 to 10, with 0 indicating a missing value, but because of a data entry error, all 10's were changed into 0's. Unfortunately, since some of the patients have missing values for this field, it's impossible to say whether a 0 in this field is a real 0 or a 10. Quite a few of the records have that problem.

**Data Miner:** Interesting. Were there any other problems? **Researcher:** Yes, fields 2 and 3 are basically the same, but I assume that you probably noticed that.

**Data Miner:** Yes, but these fields were only weak predictors of field 5.

**Researcher:** Anyway, given all those problems, I'm surprised you were able to accomplish anything.

**Data Miner:** True, but my results are really quite good. Field 1 is a very strong predictor of field 5. I'm surprised that this wasn't noticed before.

**Researcher:** What? Field 1 is just an identification number.

Data Miner: Nonetheless, my results speak for themselves.

**Researcher:** Oh, no! I just remembered. We assigned ID numbers after we sorted the records based on field 5. There is a strong connection, but it isn't very sensible. Sorry.

**Researcher:** Anyway, given all those problems, I'm surprised you were able to accomplish anything.

**Data Miner:** True, but my results are really quite good. Field 1 is a very strong predictor of field 5. I'm surprised that this wasn't noticed before.

**Researcher:** What? Field 1 is just an identification number. **Data Miner:** Nonetheless, my results speak for themselves. **Researcher:** Oh, no! I just remembered. We assigned ID numbers after we sorted the records based on field 5. There is a strong connection, but it isn't very sensible. Sorry.

OK, what's the point?

You have to

Understand the task you want to solve and the data!

## Data Objects

*Data objects* represent entities we work with (e.g., classify them). For example, in cancer prediction, the data objects are patients. In fruit classification, the data objects are individual fruits. *Data objects* represent entities we work with (e.g., classify them).

For example, in cancer prediction, the data objects are patients. In fruit classification, the data objects are individual fruits.

Data objects are described by *attributes* (or *features* or *variables*). For example, the age, weight, genetic profile, and other patient

characteristics. Or the width and height of a fruit.

## Attributes vs Features vs Variables

The name differs from field to field.

#### Attributes vs Features vs Variables

The name differs from field to field.

So, the following names are usually used as synonyms:

- Attributes used mostly by database and data mining experts.
- Features used mostly by machine learning experts.
- Variables used mostly by statisticians.

## Attributes vs Features vs Variables

The name differs from field to field.

So, the following names are usually used as synonyms:

- Attributes used mostly by database and data mining experts.
- Features used mostly by machine learning experts.
- Variables used mostly by statisticians.

One may make some distinctions

- Attributes represent information about the object without any additional assumptions.
- Features assume that their values are somewhat characteristic of the object.
- Variables assume that there is some process behind them (typically a random process in the case of statistics).

# Data Types - Categorical Attributes

*Categorical attributes* (nominal attributes) are symbols or names of things.

- Each value represents some kind of category, code, or state.
- Values are not ordered and should not be used quantitatively (in computer science, the values are known as enumerations).

# Data Types - Categorical Attributes

*Categorical attributes* (nominal attributes) are symbols or names of things.

- Each value represents some kind of category, code, or state.
- Values are not ordered and should not be used quantitatively (in computer science, the values are known as enumerations).
- Examples:

 $\mathsf{hair\_color} \in \{\mathsf{black}, \mathsf{brown}, \mathsf{blond}, \mathsf{red}, \mathsf{auburn}, \mathsf{gray}, \mathsf{white}\}$ 

```
marital\_status \in {single, married, divorced, widowed}
```

 $\mathsf{customer\_ID} \in \{0, 1, 2, \ldots\}$ 

Even though the last one is usually expressed using numbers, it should not be used quantitatively.

# Data Types - Categorical Attributes

*Categorical attributes* (nominal attributes) are symbols or names of things.

- Each value represents some kind of category, code, or state.
- Values are not ordered and should not be used quantitatively (in computer science, the values are known as enumerations).
- Examples:

 $\mathsf{hair\_color} \in \{\mathsf{black}, \mathsf{brown}, \mathsf{blond}, \mathsf{red}, \mathsf{auburn}, \mathsf{gray}, \mathsf{white}\}$ 

```
marital_status \in {single, married, divorced, widowed}
```

 $\mathsf{customer\_ID} \in \{0, 1, 2, \ldots\}$ 

Even though the last one is usually expressed using numbers, it should not be used quantitatively.

*Binary attributes* are categorical attributes with only two values.

## DataTypes - Ordinal Attributes

*Ordinal attribute* is an attribute with values that have a meaningful order or ranking among them.
DataTypes - Ordinal Attributes

*Ordinal attribute* is an attribute with values that have a meaningful order or ranking among them.

#### **Examples:**

```
\mathsf{drink\_size} \in \{\mathsf{small}, \mathsf{medium}, \mathsf{large}\}
```

```
\mathsf{grades} \in \{\mathsf{A},\mathsf{B},\mathsf{C},\mathsf{D},\mathsf{E},\mathsf{F}\}
```

It can also be obtained by discretizing numeric quantities into series of intervals.

Ordinal attributes do not allow arithmetic operations.

DataTypes - Ordinal Attributes

*Ordinal attribute* is an attribute with values that have a meaningful order or ranking among them.

#### **Examples:**

```
\mathsf{drink\_size} \in \{\mathsf{small}, \mathsf{medium}, \mathsf{large}\}
```

```
\mathsf{grades} \in \{\mathsf{A},\mathsf{B},\mathsf{C},\mathsf{D},\mathsf{E},\mathsf{F}\}
```

It can also be obtained by discretizing numeric quantities into series of intervals.

Ordinal attributes do not allow arithmetic operations.

Categorical and ordinal attributes are called *qualitative* attributes. Next, we look at numeric, i.e., *quantitative* attributes.

# Data Types - Numeric Attributes

Numeric attributes are quantities represented by numbers.

# Data Types - Numeric Attributes

Numeric attributes are quantities represented by numbers.

Distinguish two types: Interval-scale and ratio-scale.

	INTERVAL SCALE	RATIO SCALE	
Measurement	Equal intervals between	Equal intervals with	
interval	consecutive points.	the presence of a true zero.	
Absolute	Lacks a true zero point	Possesses a true	
zero	Lacks a true zero point.	zero point.	
Statistical	Limited to addition	Allows for meaningful	
analysis	and subtraction	multiplication and division.	
Meaningful	Ratios are not meaningful	Ratios are meaningful	
ratios	due to the lack of zero.	due to the presence of zero.	
Examples	IQ scores,	Height, weight, income, etc.	
	Celsius temperature,		
	NPS data, etc.		

## Discrete vs Continuous Attributes

Often, two kinds of numeric attributes are distinguished:

## Discrete vs Continuous Attributes

Often, two kinds of numeric attributes are distinguished:

► Discrete

A finite or countably infinite range of values, i.e., integers may represent the values.

Some (but not all) authors count the qualitative (categorical, ordinal) attributes among the discrete attributes.

## Discrete vs Continuous Attributes

Often, two kinds of numeric attributes are distinguished:

► Discrete

A finite or countably infinite range of values, i.e., integers may represent the values.

Some (but not all) authors count the qualitative (categorical, ordinal) attributes among the discrete attributes.

#### Continuous

An uncountably infinite range of values, typically an interval. There are several more or less formal definitions of continuous attributes in the literature. For example:

- All non-discrete variables.
- Have an infinite number of values between any two values.
- ► Their values are measured (??).

Deeper characteristics of data (statistical properties, etc.) will be examined at tutorials.

- One of the widely used methods for machine learning.
- Intuitively simple, directly explainable.
- Basis for random forests (a powerful model).

- One of the widely used methods for machine learning.
- Intuitively simple, directly explainable.
- Basis for random forests (a powerful model).
- We will consider the ID3 algorithm. Quinlan, 1979
- Various adjustments that appear in C4.5, CART, etc.

Consider the weather forecast for tennis playing. How would you decide whether to play today?



Consider the weather forecast for tennis playing. How would you decide whether to play today?



How do we obtain such a tree based on experience/data?

Consider data represented as follows:

• A finite set of *attributes*  $\mathcal{A} = \{A_1, \ldots, A_n\}$ .

• Each attribute  $A \in A$  has its set of values V(A).

We start with trees on discrete datasets, that is, assume V(A) finite for all  $A \in A$ .

Consider data represented as follows:

- A finite set of *attributes*  $\mathcal{A} = \{A_1, \ldots, A_n\}$ .
- Each attribute  $A \in A$  has its set of values V(A).

We start with trees on discrete datasets, that is, assume V(A) finite for all  $A \in A$ .

Objects to be classified are described by vectors of values of all attributes:

$$\vec{x} = (x_1, \ldots, x_n) \in V(A_1) \times \cdots \times V(A_n)$$

Given  $\vec{x}$  and an attribute  $A_k$  we denote by  $A_k(\vec{x})$  the value  $x_k$  of the attribute  $A_k$  in  $\vec{x}$ .

Consider data represented as follows:

- A finite set of *attributes*  $\mathcal{A} = \{A_1, \ldots, A_n\}$ .
- Each attribute  $A \in A$  has its set of values V(A).

We start with trees on discrete datasets, that is, assume V(A) finite for all  $A \in A$ .

Objects to be classified are described by vectors of values of all attributes:

$$\vec{x} = (x_1, \ldots, x_n) \in V(A_1) \times \cdots \times V(A_n)$$

Given  $\vec{x}$  and an attribute  $A_k$  we denote by  $A_k(\vec{x})$  the value  $x_k$  of the attribute  $A_k$  in  $\vec{x}$ .

Consider a set *C* of *classes*.

We consider a multiclass classification in general, i.e., C is an arbitrary finite set.

The tennis problem:

The attributes are:

 $A_1 = Outlook, A_2 = Temperature, A_3 = Humidity, A_4 = Wind$ 

The tennis problem:

The attributes are:

 $A_1 = Outlook, A_2 = Temperature, A_3 = Humidity, A_4 = Wind$ 

The sets of values of the attributes:

- $\blacktriangleright$  V(A<sub>1</sub>) = {Sunny, Overcast, Rain}
- $\blacktriangleright V(A_2) = \{Hot, Mild, Cool\}$
- $\blacktriangleright V(A_3) = \{High, Normal\}$

$$\blacktriangleright V(A_4) = \{Strong, Weak\}$$

The tennis problem:

The attributes are:

 $A_1 = Outlook, A_2 = Temperature, A_3 = Humidity, A_4 = Wind$ 

The sets of values of the attributes:

- $\blacktriangleright$  V(A<sub>1</sub>) = {Sunny, Overcast, Rain}
- $\blacktriangleright V(A_2) = \{Hot, Mild, Cool\}$
- $\blacktriangleright V(A_3) = \{ High, Normal \}$
- $\blacktriangleright V(A_4) = \{Strong, Weak\}$

Consider

 $ec{x} = (Overcast, Hot, Normal, Weak)$  $\in V(A_1) imes V(A_2) imes V(A_3) imes V(A_4)$ 

The tennis problem:

The attributes are:

 $A_1 = Outlook, A_2 = Temperature, A_3 = Humidity, A_4 = Wind$ 

The sets of values of the attributes:

- $\blacktriangleright V(A_1) = \{Sunny, Overcast, Rain\}$
- $\blacktriangleright V(A_2) = \{Hot, Mild, Cool\}$

$$\blacktriangleright V(A_3) = \{ High, Normal \}$$

 $\blacktriangleright V(A_4) = \{ Strong, Weak \}$ 

Consider

 $ec{x} = (Overcast, Hot, Normal, Weak)$  $\in V(A_1) imes V(A_2) imes V(A_3) imes V(A_4)$ 

Then

 $A_3(\vec{x}) = Humidity(\vec{x}) = Normal$  $A_4(\vec{x}) = Wind(\vec{x}) = Weak$ 

The tennis problem:

The attributes are:

 $A_1 = Outlook, A_2 = Temperature, A_3 = Humidity, A_4 = Wind$ 

The sets of values of the attributes:

- $\blacktriangleright V(A_1) = \{Sunny, Overcast, Rain\}$
- $\blacktriangleright V(A_2) = \{Hot, Mild, Cool\}$
- $\blacktriangleright V(A_3) = \{High, Normal\}$
- $\blacktriangleright V(A_4) = \{Strong, Weak\}$

Consider

 $ec{x} = (Overcast, Hot, Normal, Weak)$  $\in V(A_1) imes V(A_2) imes V(A_3) imes V(A_4)$ 

Then

$$A_{3}(\vec{x}) = Humidity(\vec{x}) = Normal$$
$$A_{4}(\vec{x}) = Wind(\vec{x}) = Weak$$
$$\bullet C = \{Yes, No\}$$

Consider (directed, rooted) trees T = (T, E) where T is a set of nodes and  $E \subseteq T \times T$  is a set of directed edges.

Consider (directed, rooted) trees T = (T, E) where T is a set of nodes and  $E \subseteq T \times T$  is a set of directed edges.

Denote by  $T_{leaf} \subseteq T$  the set of all *leaves* of the tree and by  $T_{int}$  the set  $T \smallsetminus T_{leaf}$  of *internal nodes*.

Consider (directed, rooted) trees  $\mathcal{T} = (T, E)$  where T is a set of nodes and  $E \subseteq T \times T$  is a set of directed edges.

Denote by  $T_{leaf} \subseteq T$  the set of all *leaves* of the tree and by  $T_{int}$  the set  $T \smallsetminus T_{leaf}$  of *internal nodes*.

A decision tree is

• a tree 
$$\mathcal{T} = (T, E)$$
 where

▶ each leaf  $\tau \in T_{leaf}$  is assigned a class  $C(\tau) \in C$ ,

▶ each internal node  $\tau \in T_{int}$  is assigned an attribute  $A(\tau) \in A$ ,

and there is a bijection between edges from τ and values of the attribute A(τ). Given an edge (τ, τ') ∈ E we write V(τ, τ') to denote the value of the attribute A(τ) assigned to the edge.

Consider (directed, rooted) trees T = (T, E) where T is a set of nodes and  $E \subseteq T \times T$  is a set of directed edges.

Denote by  $T_{leaf} \subseteq T$  the set of all *leaves* of the tree and by  $T_{int}$  the set  $T \setminus T_{leaf}$  of *internal nodes*.

A decision tree is

▶ a tree T = (T, E) where

• each leaf  $au \in T_{leaf}$  is assigned a class  $C( au) \in C$ ,

▶ each internal node  $\tau \in T_{int}$  is assigned an attribute  $A(\tau) \in A$ ,

and there is a bijection between edges from τ and values of the attribute A(τ). Given an edge (τ, τ') ∈ E we write V(τ, τ') to denote the value of the attribute A(τ) assigned to the edge.

**Inference:** Given an input  $\vec{x}$ , we traverse the tree from the root to a leaf, always choosing edges labeled with values of attributes from  $\vec{x}$ . The output is the class labeling the leaf.

$$T = \{O, H, W, z_1, z_2, z_3, z_4, z_5\}$$
  

$$T_{leaf} = \{z_1, z_2, z_3, z_4, z_5\}, T_{int} = \{O, H, W\}$$
  

$$E = \{(O, H), (O, W), (H, z_1), (H, z_2), (O, z_3), (W, z_4), (W, z_5)\}$$
  

$$C(z_1) = C(z_3) = No, C(z_2) = C(z_4) = Yes$$
  

$$A(O) = Outlook, A(H) = Humidity, A(W) = Wind$$
  

$$V(O, H) = Sunny, V(O, z_3) = Overcast, V(O, W) = Rain$$
  

$$V(H, z_1) = High, V(H, z_2) = Normal$$
  

$$V(W, z_4) = Strong, V(W, z_5) = Weak$$

**Inference:** For (*Rain*, *Hot*, *High*, *Strong*) we reach *z*<sub>4</sub>, yielding *No*.

65

Outlook

Rain

Yes

Strong

No

Wind

Weak

Yes

Consider a training dataset

$$\mathcal{D} = \{ (\vec{x}_k, c_k) \mid k = 1, \dots, p \}$$

Here  $\vec{x}_k \in V(A_1) \times \cdots \times V(A_k)$  and  $c_k \in C$  for every k.

Technically  $\ensuremath{\mathcal{D}}$  can be a multiset containing several occurrences of the same vector.

Index	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

 $\mathcal{D} = \{((Sunny, Hot, High, Weak), No), \\ ((Sunny, Hot, High, Strong), No)$ 

. . .

((Rain, Mild, High, Strong), No)

The learning algorithm ID3 works as follows:

Start with the whole training dataset  $\mathcal{D}$ .

The learning algorithm ID3 works as follows:

- Start with the whole training dataset  $\mathcal{D}$ .
- If there is just a single class in D, create a single node decision tree that returns the class.

The learning algorithm ID3 works as follows:

- Start with the whole training dataset  $\mathcal{D}$ .
- If there is just a single class in D, create a single node decision tree that returns the class.
- Otherwise, identify an attribute A ∈ A which best classifies the examples in D. For every v ∈ V(A) we obtain

$$\mathcal{D}_{\mathbf{v}} = \{ \vec{x} \mid \vec{x} \in \mathcal{D}, A(\vec{x}) = \mathbf{v} \}$$

We aim to have each  $\mathcal{D}_{v}$  as pure as possible, that is, ideally, to contain examples of just a single class.

The learning algorithm ID3 works as follows:

- Start with the whole training dataset  $\mathcal{D}$ .
- If there is just a single class in D, create a single node decision tree that returns the class.
- Otherwise, identify an attribute A ∈ A which best classifies the examples in D. For every v ∈ V(A) we obtain

$$\mathcal{D}_{\mathbf{v}} = \{ \vec{x} \mid \vec{x} \in \mathcal{D}, A(\vec{x}) = \mathbf{v} \}$$

We aim to have each  $\mathcal{D}_{\nu}$  as pure as possible, that is, ideally, to contain examples of just a single class.

Finally,

• create a root node  $\tau$  of a decision tree,

The learning algorithm ID3 works as follows:

- Start with the whole training dataset  $\mathcal{D}$ .
- If there is just a single class in D, create a single node decision tree that returns the class.
- Otherwise, identify an attribute A ∈ A which best classifies the examples in D. For every v ∈ V(A) we obtain

$$\mathcal{D}_{\mathbf{v}} = \{ \vec{x} \mid \vec{x} \in \mathcal{D}, A(\vec{x}) = \mathbf{v} \}$$

We aim to have each  $\mathcal{D}_{v}$  as pure as possible, that is, ideally, to contain examples of just a single class.

Finally,

create a root node τ of a decision tree,

assign the attribute A to τ,

The learning algorithm ID3 works as follows:

- Start with the whole training dataset  $\mathcal{D}$ .
- If there is just a single class in D, create a single node decision tree that returns the class.
- Otherwise, identify an attribute A ∈ A which best classifies the examples in D. For every v ∈ V(A) we obtain

$$\mathcal{D}_{\mathbf{v}} = \{\vec{x} \mid \vec{x} \in \mathcal{D}, A(\vec{x}) = \mathbf{v}\}$$

We aim to have each  $\mathcal{D}_{v}$  as pure as possible, that is, ideally, to contain examples of just a single class.

Finally,

- create a root node τ of a decision tree,
- assign the attribute A to  $\tau$ ,
- For every v ∈ V(A), recursively construct a decision tree with a root τ<sub>v</sub> using D<sub>v</sub>,

The learning algorithm ID3 works as follows:

- Start with the whole training dataset  $\mathcal{D}$ .
- If there is just a single class in D, create a single node decision tree that returns the class.
- Otherwise, identify an attribute A ∈ A which best classifies the examples in D. For every v ∈ V(A) we obtain

$$\mathcal{D}_{\mathbf{v}} = \{\vec{x} \mid \vec{x} \in \mathcal{D}, A(\vec{x}) = \mathbf{v}\}$$

We aim to have each  $\mathcal{D}_{v}$  as pure as possible, that is, ideally, to contain examples of just a single class.

Finally,

- create a root node τ of a decision tree,
- assign the attribute A to  $\tau$ ,
- For every v ∈ V(A), recursively construct a decision tree with a root τ<sub>v</sub> using D<sub>v</sub>,
- for every  $v \in V(A)$  introduce an edge  $(\tau, \tau_v)$  assigned v.

- 1: function ID3(dataset  $\mathcal{D}$ , attribute set  $\mathcal{A}$ )
- 2: Create a root node  $\tau$  for the tree
- 3: if  $\mathcal{D} = \emptyset$  then
- 4: Return the single node  $\tau$  assigned with a default class.
- 5: else if all examples in  $\mathcal{D}$  are of the same class c then
- 6: Return the single-node tree, where au is assigned c
- 7: else if set of attributes  $\mathcal{A}$  is empty then
- 8: Return the single-node tree where  $\tau$  is assigned the most common class in  $\mathcal{D}$
- 9: else
- 10: Choose attribute  $A \in A$  best classifying examples in  $\mathcal{D}$ .
- 11: Set the decision attribute for  $\tau$  to A
- 12: for each value  $v \in D(A)$  do
- 13: Compute a decision tree ID3( $\mathcal{D}_{v}, \mathcal{A} \setminus \{A\}$ ) with root  $\tau_{v}$ ,
- 14: add a new edge  $(\tau, \tau_v)$  assigned v.
- 15: end for
- 16: end if
- 17: return au
- 18: end function

# Best Classifying Attribute

We aim to choose an attribute that best informs us about the class. As a result, we would possibly use as few attributes as possible and obtain a small tree containing only class-relevant decisions.

How to choose an attribute that best classifies examples in  $\mathcal{D}$ ?

There are several measures used in practice.

The most common are

- information gain
- Gini impurity decrease
The information gain is based on the notion of entropy.

The information gain is based on the notion of entropy.

We need some notation:

• Given a training dataset  $\mathcal{D}$  and a class  $c \in C$  we denote by  $p_c$  the proportion of examples with class c in  $\mathcal{D}$ .

The information gain is based on the notion of entropy.

We need some notation:

- ► Given a training dataset D and a class c ∈ C we denote by p<sub>c</sub> the proportion of examples with class c in D.
- We define the *entropy* of  $\mathcal{D}$  by

$$Entropy(\mathcal{D}) = \sum_{c \in C} -p_c \log_2 p_c$$

The information gain is based on the notion of entropy.

We need some notation:

- ► Given a training dataset D and a class c ∈ C we denote by p<sub>c</sub> the proportion of examples with class c in D.
- We define the *entropy* of  $\mathcal{D}$  by

$$Entropy(\mathcal{D}) = \sum_{c \in C} -p_c \log_2 p_c$$

The information gain of an attribute A is then defined by

$$extsf{Gain}(\mathcal{D}, \mathcal{A}) = extsf{Entropy}(\mathcal{D}) - \sum_{v \in V(\mathcal{A})} rac{|\mathcal{D}_v|}{|\mathcal{D}|} extsf{Entropy}(\mathcal{D}_v)$$

In every step of the ID3 algorithm, we choose an attribute *maximizing* the information gain for the current dataset  $\mathcal{D}$ .

The intuition behind information gain:

Consider C = {0,1} and p the proportion of examples of class 1. p measures the "uncertainty" of the class:



The intuition behind information gain:

Consider C = {0,1} and p the proportion of examples of class 1. p measures the "uncertainty" of the class:



∑<sub>v∈V(A)</sub> |D|/|D| Entropy(D<sub>v</sub>) is weighted uncertainty of classes in each D<sub>v</sub> (weighted by the relative size of D<sub>v</sub>).

The intuition behind information gain:

Consider C = {0,1} and p the proportion of examples of class 1. p measures the "uncertainty" of the class:



∑<sub>v∈V(A)</sub> |D|/|D| Entropy(D<sub>v</sub>) is weighted uncertainty of classes in each D<sub>v</sub> (weighted by the relative size of D<sub>v</sub>).
 Gain(D, A) measures reduction in uncertainty of classes by splitting D according to A.

• We define *Gini impurity* of  $\mathcal{D}$  by

$$Gini(\mathcal{D}) = 1 - \sum_{c \in C} p_c^2$$

• We define *Gini impurity* of  $\mathcal{D}$  by

$$Gini(\mathcal{D}) = 1 - \sum_{c \in C} p_c^2$$

The *impurity decrease* of an attribute A is then defined similarly to the gain in the entropy case

$$\mathit{ImpDec}(\mathcal{D}, A) = \mathit{Gini}(\mathcal{D}) - \sum_{v \in V(A)} \frac{|\mathcal{D}_v|}{|\mathcal{D}|} \mathit{Gini}(\mathcal{D}_v)$$

In every step of the ID3 algorithm, we choose an attribute maximizing the impurity decrease for the current dataset  $\mathcal{D}$ .

What is the intuition behind  $Gini(\mathcal{D})$  ?

What is the intuition behind  $Gini(\mathcal{D})$  ?

Assume we randomly independently choose objects from  $\mathcal{D}$ .

 $1 - \sum_{c \in C} p_c^2$  is the probability of choosing two objects of different classes in two consecutive independent trials. Indeed,  $p_c$  is the probability of choosing an object of class c,  $p_c^2$  the probability of choosing objects of the class c twice, and  $\sum_{c \in C} p_c^2$  the probability of choosing two objects of the same class.

In what follows (and at the exam), we will work only with the Gini impurity as it is easier to compute.

Consider our tennis example (see the table).

Consider the whole dataset 
$$D$$
.  
•  $p_{Yes} = 9/14$   
•  $p_{No} = 5/14$   
•  $Gini(D) = 1 - (9/14)^2 - (5/14)^2 = 0.45918$ 

Consider our tennis example (see the table).

Consider the whole dataset 
$$\mathcal{D}$$
.
 $p_{Yes} = 9/14$ 
 $p_{No} = 5/14$ 
 $Gini(\mathcal{D}) = 1 - (9/14)^2 - (5/14)^2 = 0.45918$ 
For  $A = Outlook$  we get
 $Gini(\mathcal{D}_{Sunny}) = 1 - (2/5)^2 - (3/5)^2 = 0.48$ 
 $Gini(\mathcal{D}_{Overcast}) = 1 - 1^2 - 0^2 = 0$ 
 $Gini(\mathcal{D}_{Rain}) = 1 - (3/5)^2 - (2/5)^2 = 0.48$ 

Consider our tennis example (see the table).

Consider the whole dataset 
$$\mathcal{D}$$
.
 $p_{Yes} = 9/14$ 
 $p_{No} = 5/14$ 
 $Gini(\mathcal{D}) = 1 - (9/14)^2 - (5/14)^2 = 0.45918$ 
For  $A = Outlook$  we get
 $Gini(\mathcal{D}_{Sunny}) = 1 - (2/5)^2 - (3/5)^2 = 0.48$ 
 $Gini(\mathcal{D}_{Overcast}) = 1 - 1^2 - 0^2 = 0$ 
 $Gini(\mathcal{D}_{Rain}) = 1 - (3/5)^2 - (2/5)^2 = 0.48$ 
Thus

 $ImpDec(\mathcal{D}, Outlook) =$  $0.459 - (5/14) \cdot 0.48 - (4/14) \cdot 0 - (5/14) \cdot 0.48$ = 0.117

• 
$$ImpDec(\mathcal{D}, Temperature) = 0.018$$

•  $ImpDec(\mathcal{D}, Humidity) = 0.091$ 

So the largest information gain is given by the Outlook.

Going further on, consider  $\mathcal{D} = \mathcal{D}_{Sunny}$ . We get

- $ImpDec(\mathcal{D}, Temperature) = 0.279$
- $ImpDec(\mathcal{D}, Humidity) = 0.48$
- $ImpDec(\mathcal{D}, Wind) = 0.013$

The best choice attribude after Sunny in Outlook is Humidity.

Going further on, consider  $\mathcal{D}=\mathcal{D}_{\textit{Sunny}}.$  We get

- ► ImpDec(D, Temperature) = 0.279
- $ImpDec(\mathcal{D}, Humidity) = 0.48$
- $ImpDec(\mathcal{D}, Wind) = 0.013$

The best choice attribude after Sunny in Outlook is Humidity.

Now consider  $\mathcal{D} = \mathcal{D}_{Rain}$ .

- $ImpDec(\mathcal{D}, Temperature) = 0.013$
- $ImpDec(\mathcal{D}, Humidity) = 0.013$
- $ImpDec(\mathcal{D}, Wind) = 0.48$

The best choice attribude after Rain in Outlook is Wind.

#### Continuous-Valued Attributes

What if values of an attribute *A* come from a continuous variable? *A* is a numerical attribute that can take any value in an interval, such as temperature, size, time, etc.

#### Continuous-Valued Attributes

What if values of an attribute A come from a continuous variable? A is a numerical attribute that can take any value in an interval, such as temperature, size, time, etc.

Consider an internal node  $\tau \in T_{int}$  assigned such a continuous attribute *A*. Then

- $\tau$  is assigned a threshold value called a *cut point*  $H \in \mathbb{R}$ ,
- there are two edges  $e_{true}$ ,  $e_{false}$  from  $\tau$ ,
- e<sub>true</sub> labeled with True and e<sub>false</sub> labeled with False.

During inference, when considering an example  $\vec{x}$  in the node  $\tau$ ,

- evaluate  $A(\vec{x}) \leq H$ ,
- ▶ if  $A(\vec{x}) \leq H$ , then follow  $e_{\text{true}}$ ,

else follow e<sub>false</sub>.



In training, the cut point is chosen from the attribute values in the training set using information gain/impurity decrease similar to discrete attributes.

# Iris Example



#### Attributes

Sepal.Length, Sepal.Width, Petal.Length, Petal.Width

**Classes** (Variety) Setosa, Versicolor, Virginica

#### Iris Example

\_

The dataset (150 examples):

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Variety
5.5	3.5	1.3	0.2	Setosa
6.8	2.8	4.8	1.4	Versicolor
6.7	3.1	4.7	1.5	Versicolor
6.9	3.1	5.1	2.3	Virginica
7.3	2.9	6.3	1.8	Virginica
5.4	3.7	1.5	0.2	Setosa
4.6	3.4	1.4	0.3	Setosa
6.2	2.8	4.8	1.8	Virginica
5.4	3.0	4.5	1.5	Versicolor
4.7	3.2	1.6	0.2	Setosa
6.7	3.3	5.7	2.1	Virginica
5.0	3.4	1.5	0.2	Setosa
5.0	3.0	1.6	0.2	Setosa
4.4	2.9	1.4	0.2	Setosa
6.0	3.4	4.5	1.6	Versicolor
5.1	3.5	1.4	0.2	Setosa
6.6	3.0	4.4	1.4	Versicolor
5.9	3.2	4.8	1.8	Versicolor
5.6	2.8	4.9	2.0	Virginica

#### Iris Example



#### Iris Example - Decision Tree



#### Iris Example - Decision Tree Boudaries



If the leaves are split further, the Depth = 2 boundary would be added.

How important are attributes for the trained tree  $\mathcal{T}$ ? Depends on

- how close they are to the root of  $\mathcal{T}$ ,
- how large information gain/decrease in impurity they give.

There are several formulae for computing the importance.

How important are attributes for the trained tree  $\mathcal{T}$ ? Depends on

• how close they are to the root of  $\mathcal{T}$ ,

how large information gain/decrease in impurity they give.

There are several formulae for computing the importance.

One example is *mean decrease impurity* defined as follows:

Consider a decision tree T trained on a dataset  $\mathcal{D}$  using the ID3.

How important are attributes for the trained tree  $\mathcal{T}$ ? Depends on

• how close they are to the root of  $\mathcal{T}$ ,

► how large information gain/decrease in impurity they give. There are several formulae for computing the importance.

One example is *mean decrease impurity* defined as follows:

Consider a decision tree T trained on a dataset  $\mathcal{D}$  using the ID3.

For every node  $\tau$  of  $\mathcal{T}$ , denote by  $\mathcal{D}[\tau]$  the subset of  $\mathcal{D}$  which was used in the ID3 procedure when the node  $\tau$  was created (line 2).

How important are attributes for the trained tree  $\mathcal{T}$ ? Depends on

• how close they are to the root of  $\mathcal{T}$ ,

► how large information gain/decrease in impurity they give. There are several formulae for computing the importance.

One example is *mean decrease impurity* defined as follows:

Consider a decision tree T trained on a dataset  $\mathcal{D}$  using the ID3.

For every node  $\tau$  of  $\mathcal{T}$ , denote by  $\mathcal{D}[\tau]$  the subset of  $\mathcal{D}$  which was used in the ID3 procedure when the node  $\tau$  was created (line 2).

Consider an attribute A and denote by  $T[A] \subseteq T_{int}$  the set of all nodes of  $\mathcal{T}$  assigned the attribute A by ID3 (line 11).

How important are attributes for the trained tree  $\mathcal{T}$ ? Depends on

• how close they are to the root of  $\mathcal{T}$ ,

► how large information gain/decrease in impurity they give. There are several formulae for computing the importance.

One example is *mean decrease impurity* defined as follows:

Consider a decision tree T trained on a dataset D using the ID3.

For every node  $\tau$  of  $\mathcal{T}$ , denote by  $\mathcal{D}[\tau]$  the subset of  $\mathcal{D}$  which was used in the ID3 procedure when the node  $\tau$  was created (line 2).

Consider an attribute A and denote by  $T[A] \subseteq T_{int}$  the set of all nodes of T assigned the attribute A by ID3 (line 11).

Then define the importance as the average decrease in Gini impurity (i.e., average ImpDec) in the nodes of T[A]:

$$\textit{GiniImportance}(\textit{A}) = \sum_{ au \in \mathcal{T}[\textit{A}]} rac{|\mathcal{D}[ au]|}{|\mathcal{D}|}\textit{ImpDec}(\mathcal{D}[ au],\textit{A})$$

# **Decision Trees**

Practical Issues

#### **Practical Issues**

- Data preprocessing
- Model tunning (overfitting and underfitting)
- Sensitivity to changes in data/hyperparameters
- Learning representation problems (the XOR)

Little preprocessing is needed for decision trees.

Little preprocessing is needed for decision trees.

Of course, ideally, clean up the data, but

 Missing values are not such a big issue (considering a concrete example, exclude the attributes with missing values)

Little preprocessing is needed for decision trees.

Of course, ideally, clean up the data, but

- Missing values are not such a big issue (considering a concrete example, exclude the attributes with missing values)
- Outliers are not such a big issue either; the division of nodes is done based on relative counts, not concrete values.

Little preprocessing is needed for decision trees.

Of course, ideally, clean up the data, but

- Missing values are not such a big issue (considering a concrete example, exclude the attributes with missing values)
- Outliers are not such a big issue either; the division of nodes is done based on relative counts, not concrete values.
- Decision trees can cope with continuous and categorical values directly (i.e., no encoding necessary)

Little preprocessing is needed for decision trees.

Of course, ideally, clean up the data, but

- Missing values are not such a big issue (considering a concrete example, exclude the attributes with missing values)
- Outliers are not such a big issue either; the division of nodes is done based on relative counts, not concrete values.
- Decision trees can cope with continuous and categorical values directly (i.e., no encoding necessary)

Imbalanced classes might cause problems because of small information gain/impurity decrease in splitting.

#### Imbalanced Classes

Consider a dataset  $\ensuremath{\mathcal{D}}$  where

- there are two classes,  $C = \{0, 1\}$ ,
- $\blacktriangleright$  10<sup>6</sup> examples have the class 1,
- ▶ 100 examples have the class 0.
Consider a dataset  $\ensuremath{\mathcal{D}}$  where

- there are two classes,  $C = \{0, 1\}$ ,
- ▶ 10<sup>6</sup> examples have the class 1,
- ▶ 100 examples have the class 0.

Then

- ▶  $p_1 = 10^6/(10^6 + 100) \approx 1$  and  $p_0 = 100/10^6 \approx 0$ ,
- thus the Gini impurity  $1 p_1^2 p_0^2 \approx 0$ .

Consider a dataset  $\ensuremath{\mathcal{D}}$  where

- there are two classes,  $C = \{0, 1\}$ ,
- ▶ 10<sup>6</sup> examples have the class 1,
- ▶ 100 examples have the class 0.

Then

▶ 
$$p_1 = 10^6/(10^6 + 100) pprox 1$$
 and  $p_0 = 100/10^6 pprox 0$ ,

▶ thus the Gini impurity  $1 - p_1^2 - p_0^2 \approx 0$ . Consider an attribute *A* with  $V(A) = \{a, b\}$ .

Splitting  $\mathcal{D}$  according to A gives to sets  $\mathcal{D}_a$  and  $\mathcal{D}_b$ .

Consider a dataset  $\ensuremath{\mathcal{D}}$  where

- there are two classes,  $C = \{0, 1\}$ ,
- ▶ 10<sup>6</sup> examples have the class 1,
- 100 examples have the class 0.

Then

▶ 
$$p_1 = 10^6/(10^6 + 100) \approx 1$$
 and  $p_0 = 100/10^6 \approx 0$ ,

▶ thus the Gini impurity  $1 - p_1^2 - p_0^2 \approx 0$ . Consider an attribute *A* with  $V(A) = \{a, b\}$ .

Splitting  $\mathcal{D}$  according to A gives to sets  $\mathcal{D}_a$  and  $\mathcal{D}_b$ .

What is the impurity decrease caused by the attribute?

$$ImpDec(\mathcal{D}, A) = Gini(\mathcal{D}) - \frac{|\mathcal{D}_{a}|}{|\mathcal{D}|}Gini(\mathcal{D}_{a}) - \frac{|\mathcal{D}_{b}|}{|\mathcal{D}|}Gini(\mathcal{D}_{b})$$

Consider a dataset  $\ensuremath{\mathcal{D}}$  where

- there are two classes,  $C = \{0, 1\}$ ,
- 10<sup>6</sup> examples have the class 1,
- 100 examples have the class 0.

Then

▶ 
$$p_1 = 10^6/(10^6 + 100) \approx 1$$
 and  $p_0 = 100/10^6 \approx 0$ ,  
▶ thus the Gini impurity  $1 - p_1^2 - p_0^2 \approx 0$ .

Consider an attribute A with  $V(A) = \{a, b\}$ .

Splitting  $\mathcal{D}$  according to A gives to sets  $\mathcal{D}_a$  and  $\mathcal{D}_b$ .

What is the impurity decrease caused by the attribute?

$$ImpDec(\mathcal{D}, A) = Gini(\mathcal{D}) - \frac{|\mathcal{D}_{a}|}{|\mathcal{D}|}Gini(\mathcal{D}_{a}) - \frac{|\mathcal{D}_{b}|}{|\mathcal{D}|}Gini(\mathcal{D}_{b})$$

For small  $|\mathcal{D}_{a}|$  (say  $\leq$  1000) we have small  $|\mathcal{D}_{a}|/|\mathcal{D}|$ 

For not so small  $\mathcal{D}_a$  we have  $Gini(\mathcal{D}_a) \approx 0$ .

In both cases, the impurity decrease is very small.

The behavior of the model on the training set:



**Right Fit** 

Underfitting







The behavior of the model on the training set:



The left one is strongly overfitting. It would possibly not work well on new data.

The behavior of the model on the training set:



- The left one is strongly overfitting. It would possibly not work well on new data.
- The right one is strongly underfitting. It would probably give poor classification results.

The behavior of the model on the training set:



- The left one is strongly overfitting. It would possibly not work well on new data.
- The right one is strongly underfitting. It would probably give poor classification results.
- The middle one seems good (but still needs to be tested on fresh data).

# Model Tuning - Overfitting in Decision Trees



See the overfitting on the left and the "nice" model on the right. Both overfitting and underfitting are best avoided. But how do we find out?

# Model Tuning (In General)

Recall from the first lecture:



The validation should be done on a **validation set** separated from the training set.

We will discuss more sophisticated techniques later.

# Model Tuning (In General)

Recall from the first lecture:



The validation should be done on a **validation set** separated from the training set.

We will discuss more sophisticated techniques later.

What hyperparameters to set? (see the next slide)

# Model Tuning (In General)

Recall from the first lecture:



The validation should be done on a **validation set** separated from the training set.

We will discuss more sophisticated techniques later.

What hyperparameters to set? (see the next slide)

What to observe? In the case of decision trees, one should observe the difference between performance measures (e.g., classification accuracy) on the training and validation sets.

The too-large difference implies an improperly fitting model.

There are several approaches available for decision trees.

Generally, the overfitting can be either prevented or resolved.

There are several approaches available for decision trees.

Generally, the overfitting can be either prevented or resolved.

Pre-pruning: Build the tree so it does not overfit by restricting its size.

There are several approaches available for decision trees.

Generally, the overfitting can be either prevented or resolved.

- Pre-pruning: Build the tree so it does not overfit by restricting its size.
- Post-pruning: Overfit with a large tree and remove subtrees to make it smaller.

There are several approaches available for decision trees.

Generally, the overfitting can be either prevented or resolved.

- Pre-pruning: Build the tree so it does not overfit by restricting its size.
- Post-pruning: Overfit with a large tree and remove subtrees to make it smaller.
- Ensemble methods: Fit several different trees and let them classify together (e.g., using majority voting).

There are several approaches available for decision trees.

Generally, the overfitting can be either prevented or resolved.

- Pre-pruning: Build the tree so it does not overfit by restricting its size.
- Post-pruning: Overfit with a large tree and remove subtrees to make it smaller.
- Ensemble methods: Fit several different trees and let them classify together (e.g., using majority voting).

The post-pruning approach has been more successful in practice than the pre-pruning because it is usually hard to say when to stop growing the tree.

We shall meet this controversy also in deep learning, where recent history shows a similar phenomenon.

The ensemble methods will be covered later when we discuss random forests.

Hyperparameters controlling the size of the tree:

► Maximum depth - do not grow the tree beyond the max depth The deeper the tree, the more complex models you can create ⇒ overfitting. Low depth may restrict expressivity.

Hyperparameters controlling the size of the tree:

- ► Maximum depth do not grow the tree beyond the max depth The deeper the tree, the more complex models you can create ⇒ overfitting. Low depth may restrict expressivity.
- Minimum number of examples to split a node if D[\[\tau]\] is small, \[\tau\$ becomes a leaf (labeled with the majority class)
   Higher values prevent a model from learning relations specific only for a few examples. Too high values may result in underfitting.

Hyperparameters controlling the size of the tree:

- ► Maximum depth do not grow the tree beyond the max depth The deeper the tree, the more complex models you can create ⇒ overfitting. Low depth may restrict expressivity.
- Minimum number of examples to split a node if D[τ] is small, τ becomes a leaf (labeled with the majority class)
   Higher values prevent a model from learning relations specific only for a few examples. Too high values may result in underfitting.
- Minimum number of examples required to be in a leaf Similar to the previous one. A higher number means we cannot have very specific branches concerned with particular combinations of values.

Hyperparameters controlling the size of the tree:

- ► Maximum depth do not grow the tree beyond the max depth The deeper the tree, the more complex models you can create ⇒ overfitting. Low depth may restrict expressivity.
- Minimum number of examples to split a node if D[τ] is small, τ becomes a leaf (labeled with the majority class)
   Higher values prevent a model from learning relations specific only for a few examples. Too high values may result in underfitting.
- Minimum number of examples required to be in a leaf Similar to the previous one. A higher number means we cannot have very specific branches concerned with particular combinations of values.
- Minimum information gain/impurity decrease A small impurity decrease means that the split does not contribute too much to the classification (their proportions after a split are similar to proportions before a split). However, keep in mind that it is *weighted average impurity* after the split.

## Post-Pruning - Reduced Error Pruning

Train a large tree and then remove nodes that make classification worse on the validation set.

## Post-Pruning - Reduced Error Pruning

Train a large tree and then remove nodes that make classification worse on the validation set.

Given a decision tree  $\mathcal{T}$  and its internal node  $\tau \in T_{int}$ , we denote by  $\mathcal{T}_{-\tau}$  the tree obtained from  $\mathcal{T}$  by removing the subtree rooted in  $\tau$ , i.e.,  $\tau$  is a leaf of  $\mathcal{T}_{-\tau}$ .

## Post-Pruning - Reduced Error Pruning

Train a large tree and then remove nodes that make classification worse on the validation set.

Given a decision tree  $\mathcal{T}$  and its internal node  $\tau \in T_{int}$ , we denote by  $\mathcal{T}_{-\tau}$  the tree obtained from  $\mathcal{T}$  by removing the subtree rooted in  $\tau$ , i.e.,  $\tau$  is a leaf of  $\mathcal{T}_{-\tau}$ .

- 1: Train  $\mathcal T$  to maximum fit on the *training dataset*.
- 2: while true do
- 3:  $Err[\mathcal{T}] \leftarrow$  the error of  $\mathcal{T}$  on the validation set.

4: for 
$$\tau \in T_{int}$$
 do

- 5:  $Err[\mathcal{T}_{-\tau}] \leftarrow$  the error of  $\mathcal{T}_{-\tau}$  on the validation set.
- 6: end for

7: **if** 
$$Err[\mathcal{T}] \leq \min\{Err[\mathcal{T}_{-\tau}] \mid \tau \in T_{int}\}\}$$
 then return  $\mathcal{T}$ 

8: **else** 

9: 
$$\mathcal{T} \leftarrow \operatorname{argmin}\{\operatorname{Err}[\mathcal{T}_{-\tau}] \mid \tau \in T_{\operatorname{int}}\}$$

10: end if

#### 11: end while

The error  $Err[\mathcal{T}]$  can be any measure of the "badness" of the decision tree  $\mathcal{T}$ . For example, 1 - Accuracy.

#### Other Pruning Methods

There are more pruning methods.

- Rule Post-Pruning:
  - Transform the tree into a set of rules. Rules correspond to paths in the tree; they have a form of implication: Specific values of attributes imply a class.
  - Remove the attribute conditions from the premises of the implications.

This gives a more refined pruning: Instead of removing the whole subtree, each path of the subtree can be pruned individually.

## Other Pruning Methods

There are more pruning methods.

- Rule Post-Pruning:
  - Transform the tree into a set of rules. Rules correspond to paths in the tree; they have a form of implication: Specific values of attributes imply a class.
  - Remove the attribute conditions from the premises of the implications.

This gives a more refined pruning: Instead of removing the whole subtree, each path of the subtree can be pruned individually.

Using cost complexity measure: Evaluate trees not only based on the classification error but also based on their size.

#### Other Pruning Methods

There are more pruning methods.

- Rule Post-Pruning:
  - Transform the tree into a set of rules. Rules correspond to paths in the tree; they have a form of implication: Specific values of attributes imply a class.
  - Remove the attribute conditions from the premises of the implications.

This gives a more refined pruning: Instead of removing the whole subtree, each path of the subtree can be pruned individually.

Using cost complexity measure: Evaluate trees not only based on the classification error but also based on their size.

Typically introduce regularization into the error functions: Given a decision tree  $\ensuremath{\mathcal{T}}$ 

 $Err_{\alpha}(\mathcal{T}) = Err(\mathcal{T}) + \alpha |\mathcal{T}|$ 

The original paper by Breiman et al. (1984) defined  $Err(\mathcal{T})$  to be the misclassification rate on the training dataset, and  $|\mathcal{T}|$  is the number of nodes of the tree  $\mathcal{T}$ .

Sensitivity to Small Changes and Randomness

 Decision trees are sensitive to small changes in data and hyperparameters.
 Several attributes may provide (almost) identical information gain but

divide the training dataset very differently.

Some implementations choose attributes partially in random (sci-kit-learn). You may get completely different trees. Sensitivity to Small Changes and Randomness

 Decision trees are sensitive to small changes in data and hyperparameters.

Several attributes may provide (almost) identical information gain but divide the training dataset very differently.

Some implementations choose attributes partially in random (sci-kit-learn). You may get completely different trees.

A solution is to train an ensemble of many decision trees and then use majority voting for classification.

This is the fundamental idea behind random forests (see later lectures).

## Iris - Illustration



Decision trees trained on the Iris dataset.

Iris Setosa is perfectly separated by many choices for the first split.

## Axis Sensitivity



The decision makes divisions along particular axes:

That is, rotated data may result in a completely different model.

That is why decision trees are often preceded by the *principal component analysis (PCA)* transformation, which aligns data along the axes of maximum data variance.

#### **XOR Training Problem**

Consider the following training dataset:



#### **XOR Training Problem**

Consider the following training dataset:



An ideal decision tree would look like this:



#### Attempts at Training on XOR

Max depth = 2:



#### Attempts at Training on XOR

Max depth = 2:



The problem: Both information gain and decrease in impurity consider only the relationship of a *single* attribute and the class.

However, there is no relationship between a single attribute and the class; both attributes need to be considered together!

## More Attempts at Training on XOR

Max depth = 3:



It's better but still fails occasionally.
- Simple to understand and interpret; trees can be visualized.
- Uses a white box model, where conditions are easily explained by boolean logic.

- Simple to understand and interpret; trees can be visualized.
- Uses a white box model, where conditions are easily explained by boolean logic.
- Can approximate an arbitrary (reasonable) boundary and capture complex non-linear relationships between attributes.
- Capable of handling multi-class problems.

- Simple to understand and interpret; trees can be visualized.
- Uses a white box model, where conditions are easily explained by boolean logic.
- Can approximate an arbitrary (reasonable) boundary and capture complex non-linear relationships between attributes.
- Capable of handling multi-class problems.
- Little data preparation, unlike other techniques requiring normalization, dummy variables, or missing value removal.
- Handles numerical and categorical data.
- Not sensitive to outliers since the splitting is based on the proportion of examples within the split ranges and not on absolute values.

- Simple to understand and interpret; trees can be visualized.
- Uses a white box model, where conditions are easily explained by boolean logic.
- Can approximate an arbitrary (reasonable) boundary and capture complex non-linear relationships between attributes.
- Capable of handling multi-class problems.
- Little data preparation, unlike other techniques requiring normalization, dummy variables, or missing value removal.
- Handles numerical and categorical data.
- Not sensitive to outliers since the splitting is based on the proportion of examples within the split ranges and not on absolute values.
- The cost of using a well-balanced tree is logarithmic in the number of data points used to train it.

- Overfitting: Trees can be over-complex and not generalize well, needing pruning or limits on tree depth.
- Instability: Small data variations can result in very different trees. This is mitigated in ensemble methods.
- Non-smooth predictions: Decision trees make piecewise constant approximations, which are not suitable for extrapolation.

- Overfitting: Trees can be over-complex and not generalize well, needing pruning or limits on tree depth.
- Instability: Small data variations can result in very different trees. This is mitigated in ensemble methods.
- Non-smooth predictions: Decision trees make piecewise constant approximations, which are not suitable for extrapolation.
- Difficulty expressing certain concepts, such as XOR, parity, or multiplexer problems (see the next slide).
- Bias in trees: Decision trees can create biased trees if some classes dominate. Balancing the dataset is recommended.

- Overfitting: Trees can be over-complex and not generalize well, needing pruning or limits on tree depth.
- Instability: Small data variations can result in very different trees. This is mitigated in ensemble methods.
- Non-smooth predictions: Decision trees make piecewise constant approximations, which are not suitable for extrapolation.
- Difficulty expressing certain concepts, such as XOR, parity, or multiplexer problems (see the next slide).
- Bias in trees: Decision trees can create biased trees if some classes dominate. Balancing the dataset is recommended.
- Learning optimal trees is NP-complete: Heuristic algorithms like greedy algorithms are used, which do not guarantee globally optimal trees. Ensemble methods can help.

# History of Decision Trees

- Hunt and colleagues use exhaustive search decision-tree methods (CLS) to model human concept learning in the 1960's.
- In the late 70's, Quinlan developed ID3 with the information gain heuristic to learn expert systems from examples.
- Simultaneously, Breiman, Friedman, and colleagues develop CART (Classification and Regression Trees), similar to ID3.
- In the 1980s, various improvements were introduced to handle noise, continuous features, missing features, and improved splitting criteria. Various expert-system development tools results.
- Quinlan's updated decision-tree package (C4.5) released in 1993.

#### Comment on Regression Trees

Decision trees can also be used to approximate functions. Assign a function value to the leaves instead of classes.



Here, "mse" is the mean-squared-error. We get to this notion later in connection with linear models and neural networks.

# Comment on Regression Trees



Intuitively, for every subinterval of  $x_1$ , the value (the red line) is at the average y over the subinterval.

How are the subintervals being set? We will answer this question later once we master the mean-squared error.

# Probabilistic Classification

Imagine that

- I look out of a window and see a bird,
- ▶ it is black, approx. 25 cm long and has a rather yellow beak.

Imagine that

I look out of a window and see a bird,

it is black, approx. 25 cm long and has a rather yellow beak.
 My daughter asks: What kind of bird is this?

Imagine that

- I look out of a window and see a bird,
- it is black, approx. 25 cm long and has a rather yellow beak.My daughter asks: What kind of bird is this?

My usual answer: This is *probably* a kind of blackbird (kos černý in Czech).

Imagine that

I look out of a window and see a bird,

it is black, approx. 25 cm long and has a rather yellow beak.My daughter asks: What kind of bird is this?

My usual answer: This is *probably* a kind of blackbird (kos černý in Czech).

Here *probably* means that out of my extensive catalog of four kinds of birds that I can recognize, "blackbird" gets the highest degree of belief based on *features* of this particular bird.

Frequentists might say that the largest proportion of birds with similar features I have ever seen were blackbirds.

Imagine that

I look out of a window and see a bird,

it is black, approx. 25 cm long and has a rather yellow beak.My daughter asks: What kind of bird is this?

My usual answer: This is *probably* a kind of blackbird (kos černý in Czech).

Here *probably* means that out of my extensive catalog of four kinds of birds that I can recognize, "blackbird" gets the highest degree of belief based on *features* of this particular bird.

Frequentists might say that the largest proportion of birds with similar features I have ever seen were blackbirds.

The degree of belief (Bayesian), or the relative frequency (frequentists), is the *probability*.

A finite or countably infinite set Ω of *possible outcomes*, Ω is called *sample space*.

Experiment: Roll one dice once. Sample space:  $\Omega = \{1, \dots, 6\}$ 

A finite or countably infinite set Ω of *possible outcomes*, Ω is called *sample space*.

Experiment: Roll one dice once. Sample space:  $\Omega = \{1, \ldots, 6\}$ 

Each element ω of Ω is assigned a "probability" value f(ω), here f must satisfy

• 
$$f(\omega) \in [0,1]$$
 for all  $\omega \in \Omega$ ,

$$\blacktriangleright \sum_{\omega \in \Omega} f(\omega) = 1.$$

If the dice is fair, then  $f(\omega) = \frac{1}{6}$  for all  $\omega \in \{1, \dots, 6\}$ .

 $\blacktriangleright$  A finite or countably infinite set  $\Omega$  of *possible outcomes*,  $\Omega$  is called *sample space*.

Experiment: Roll one dice once. Sample space:  $\Omega = \{1, \ldots, 6\}$ 

Each element  $\omega$  of  $\Omega$  is assigned a "probability" value  $f(\omega)$ , here f must satisfy

• 
$$f(\omega) \in [0, 1]$$
 for all  $\omega \in \Omega$ ,  
•  $\sum_{\alpha \in I} f(\omega) = 1$ 

• 
$$\sum_{\omega \in \Omega} f(\omega) = 1.$$

If the dice is fair, then  $f(\omega) = \frac{1}{6}$  for all  $\omega \in \{1, \dots, 6\}$ .

An event is any subset E of Ω.

• The *probability* of a given event  $E \subseteq \Omega$  is defined as

$$P(E) = \sum_{\omega \in E} f(\omega)$$

Let *E* be the event that an odd number is rolled, i.e.,  $E = \{1, 3, 5\}$ . Then  $P(E) = \frac{1}{2}$ .

 $\blacktriangleright$  A finite or countably infinite set  $\Omega$  of *possible outcomes*,  $\Omega$  is called *sample space*.

Experiment: Roll one dice once. Sample space:  $\Omega = \{1, \dots, 6\}$ 

Each element  $\omega$  of  $\Omega$  is assigned a "probability" value  $f(\omega)$ , here f must satisfy

• 
$$f(\omega) \in [0, 1]$$
 for all  $\omega \in \Omega$ ,  
•  $\sum_{\alpha \in I} f(\omega) = 1$ 

• 
$$\sum_{\omega \in \Omega} f(\omega) = 1.$$

If the dice is fair, then  $f(\omega) = \frac{1}{6}$  for all  $\omega \in \{1, \dots, 6\}$ .

An event is any subset E of Ω.

The *probability* of a given event  $E \subseteq \Omega$  is defined as

$$P(E) = \sum_{\omega \in E} f(\omega)$$

Let *E* be the event that an odd number is rolled, i.e.,  $E = \{1, 3, 5\}$ . Then  $P(E) = \frac{1}{2}$ .

**•** Basic laws:  $P(\Omega) = 1$ ,  $P(\emptyset) = 0$ , given disjoint sets A, B we have  $P(A \cup B) = P(A) + P(B)$ ,  $P(\Omega \setminus A) = 1 - P(A)$ .

Conditional Probability and Independence

► P(A | B) is the probability of A given B (assume P(B) > 0) defined by

 $P(A \mid B) = P(A \cap B)/P(B)$ 

(We assume that B is all and only information known.)

A fair dice: what is the probability that 3 is rolled assuming that an odd number is rolled? ... and assuming that an even number is rolled?

Conditional Probability and Independence

► P(A | B) is the probability of A given B (assume P(B) > 0) defined by

 $P(A \mid B) = P(A \cap B)/P(B)$ 

(We assume that B is all and only information known.)

A fair dice: what is the probability that 3 is rolled assuming that an odd number is rolled? ... and assuming that an even number is rolled?

• A and B are independent if  $P(A \cap B) = P(A) \cdot P(B)$ .

It is easy to show that if P(B) > 0, then A, B are independent iff P(A | B) = P(A). Random Variables and Random Vectors

- ► A random variable X is a function  $X : \Omega \to \mathbb{R}$ . A dice:  $X : \{1, ..., 6\} \to \{0, 1\}$  such that  $X(n) = n \mod 2$ .
- A random vector is a function  $X : \Omega \to \mathbb{R}^d$ .

#### Random Variables and Random Vectors

• A random variable X is a function  $X : \Omega \to \mathbb{R}$ . A dice:  $X : \{1, \dots, 6\} \to \{0, 1\}$  such that  $X(n) = n \mod 2$ .

• A random vector is a function  $X : \Omega \to \mathbb{R}^d$ .

We use  $X = (X_1, ..., X_d)$  where  $X_i$  is a random variable returning the *i*-th component of X.

Consider random variables X<sub>1</sub>, X<sub>2</sub> and Y. The variables X<sub>1</sub>, X<sub>2</sub> are *conditionally independent given* Y if for all x<sub>1</sub>, x<sub>2</sub> and y we have that

$$P(X_1 = x_1, X_2 = x_2 | Y = y) = P(X_1 = x_1 | Y = y) \cdot P(X_2 = x_2 | Y = y)$$

Let  $\Omega$  be a space of colored geometric shapes that are divided into two categories (1 and 0).

Assume a random vector  $X = (X_{color}, X_{shape}, X_{cat})$  where

• 
$$X_{color} : \Omega \to \{red, blue\},\$$

• 
$$X_{shape} : \Omega \rightarrow \{circle, square\},\$$

$$\blacktriangleright X_{cat}: \Omega \to \{\mathbf{1}, \mathbf{0}\}.$$

The following tables give probability distribution of values:

category 1:

category	<b>0</b> :	
----------	------------	--

	circle	square	
red	0.2	0.02	
blue	0.02	0.01	

	circle	square
red	0.05	0.3
blue	0.2	0.2

Example:  $P(red, circle, \mathbf{1}) = P(X_{color} = red, X_{shape} = circle, X_{cat} = \mathbf{1}) = 0.2$ 

Example:  $P(red, circle, \mathbf{1}) = P(X_{color} = red, X_{shape} = circle, X_{cat} = \mathbf{1}) = 0.2$ 

"Summing over" all possible values of some variable(s) gives the distribution of the rest:

$$P(red, circle) = P(X_{color} = red, X_{shape} = circle)$$
  
= P(red, circle, 1) + P(red, circle, 0)  
= 0.2 + 0.05 = 0.25

Example:  $P(red, circle, \mathbf{1}) = P(X_{color} = red, X_{shape} = circle, X_{cat} = \mathbf{1}) = 0.2$ 

"Summing over" all possible values of some variable(s) gives the distribution of the rest:

$$P(red, circle) = P(X_{color} = red, X_{shape} = circle)$$
  
= P(red, circle, 1) + P(red, circle, 0)  
= 0.2 + 0.05 = 0.25

P(red) = 0.2 + 0.02 + 0.05 + 0.3 = 0.57

Example:  $P(red, circle, 1) = P(X_{color} = red, X_{shape} = circle, X_{cat} = 1) = 0.2$ 

"Summing over" all possible values of some variable(s) gives the distribution of the rest:

$$P(red, circle) = P(X_{color} = red, X_{shape} = circle)$$
  
= P(red, circle, 1) + P(red, circle, 0)  
= 0.2 + 0.05 = 0.25

$$P(red) = 0.2 + 0.02 + 0.05 + 0.3 = 0.57$$

Thus also, all conditional probabilities can be computed:

$$P(\mathbf{1} \mid red, circle) = \frac{P(red, circle, \mathbf{1})}{P(red, circle)} = \frac{0.2}{0.25} = 0.8$$

Let  $\Omega$  be a sample space (a universum) of all objects that can be classified. We assume a probability P on  $\Omega$ .

Let  $\Omega$  be a sample space (a universum) of all objects that can be classified. We assume a probability P on  $\Omega$ .

We consider the problem of **binary classification**:

Let Y be the random variable for the category which takes values in {0,1}.

Let  $\Omega$  be a sample space (a universum) of all objects that can be classified. We assume a probability P on  $\Omega$ .

We consider the problem of **binary classification**:

- Let Y be the random variable for the category which takes values in {0,1}.
- Let X be the random vector describing n features of a given instance, i.e., X = (X<sub>1</sub>,...,X<sub>n</sub>)
  - Denote by  $\vec{x} \in \mathbb{R}^n$  values of X,
  - and by  $x_i \in \mathbb{R}$  values of  $X_i$ .

Let  $\Omega$  be a sample space (a universum) of all objects that can be classified. We assume a probability P on  $\Omega$ .

We consider the problem of **binary classification**:

- Let Y be the random variable for the category which takes values in {0,1}.
- Let X be the random vector describing n features of a given instance, i.e., X = (X<sub>1</sub>,...,X<sub>n</sub>)
  - Denote by  $\vec{x} \in \mathbb{R}^n$  values of X,
  - ▶ and by  $x_i \in \mathbb{R}$  values of  $X_i$ .

**Bayes classifier:** Given a vector of feature values  $\vec{x}$ ,

$$C^{Bayes}(\vec{x}) := \begin{cases} \mathbf{1} & \text{if } P(Y = \mathbf{1} \mid X = \vec{x}) \ge P(Y = \mathbf{0} \mid X = \vec{x}) \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

Intuitively,  $C^{Bayes}$  assigns to  $\vec{x}$  the most probable category it might be in.

## Bayesian Classification – Example

Imagine a conveyor belt with apples and apricots.

A machine is supposed to correctly distinguish apples from apricots based on their weight and diameter.

## Bayesian Classification – Example

Imagine a conveyor belt with apples and apricots.

A machine is supposed to correctly distinguish apples from apricots based on their weight and diameter.

That is,

 Y ∈ {1,0} (here our interpretation is 1 = apple, 0 = appricot)
 X = (X<sub>weight</sub>, X<sub>diam</sub>)

## Bayesian Classification – Example

Imagine a conveyor belt with apples and apricots.

A machine is supposed to correctly distinguish apples from apricots based on their weight and diameter.

That is,

 Y ∈ {1,0} (here our interpretation is 1 = apple, 0 = appricot)
 X = (X<sub>weight</sub>, X<sub>diam</sub>)

We are given a fruit of a diameter of 5cm that weighs 40g.
#### Bayesian Classification – Example

Imagine a conveyor belt with apples and apricots.

A machine is supposed to correctly distinguish apples from apricots based on their weight and diameter.

That is,

We are given a fruit of a diameter of 5cm that weighs 40g.

The Bayes classifier compares  $P(Y = \mathbf{1} | X = (40g, 5cm))$  with  $P(Y = \mathbf{0} | X = (40g, 5cm))$  and selects the more probable category given the features.

#### Crucial question: Is such a classifier good?

There are other classifiers, e.g., one which compares the weight divided by 10 with the diameter and decides based on the answer, or maybe a classifier that sums the weight and the diameter and compares the result with a constant, etc.

# **Bayes** Classifier

Let C be an arbitrary *classifier*, that is a function that to every feature vector  $\vec{x} \in \mathbb{R}^n$  assigns a class from  $\{0, 1\}$ .

# **Bayes Classifier**

Let C be an arbitrary *classifier*, that is a function that to every feature vector  $\vec{x} \in \mathbb{R}^n$  assigns a class from  $\{0, 1\}$ .

Define the error of the classifier C by

$$E_C = P(Y \neq C)$$

(Here we slightly abuse notation and apply *C* to samples, technically we apply the composition  $C \circ X$  of *C* and *X* which first determines the features using *X* and then classifies according to *C*).

# **Bayes** Classifier

Let C be an arbitrary *classifier*, that is a function that to every feature vector  $\vec{x} \in \mathbb{R}^n$  assigns a class from  $\{0, 1\}$ .

Define the error of the classifier C by

 $E_C = P(Y \neq C)$ 

(Here we slightly abuse notation and apply *C* to samples, technically we apply the composition  $C \circ X$  of *C* and *X* which first determines the features using *X* and then classifies according to *C*).

#### Theorem

The Bayes classifier  $C^{Bayes}$  minimizes  $E_C$ , that is

 $E_{C^{Bayes}} = \min_{C \text{ is a classifier}} E_C$ 

## Practical Use of Bayes Classifier

**The crucial problem:** The probability *P* is not known! In particular, where to get  $P(Y = \mathbf{1} | X = \vec{x})$ ? Note that  $P(Y = \mathbf{0} | X = \vec{x}) = 1 - P(Y = \mathbf{1} | X = \vec{x})$ 

# Practical Use of Bayes Classifier

**The crucial problem:** The probability *P* is not known! In particular, where to get  $P(Y = \mathbf{1} | X = \vec{x})$ ? Note that  $P(Y = \mathbf{0} | X = \vec{x}) = 1 - P(Y = \mathbf{1} | X = \vec{x})$ 

Given no other assumptions, this requires a table showing the probability of the category  ${\bf 1}$  for each possible feature vector  $\vec{x}.$ 

Where do you get these probabilities?

# Practical Use of Bayes Classifier

**The crucial problem:** The probability *P* is not known! In particular, where to get  $P(Y = \mathbf{1} | X = \vec{x})$ ? Note that  $P(Y = \mathbf{0} | X = \vec{x}) = 1 - P(Y = \mathbf{1} | X = \vec{x})$ 

Given no other assumptions, this requires a table showing the probability of the category **1** for each possible feature vector  $\vec{x}$ . Where do you get these probabilities?

In some cases, the probabilities might come from the knowledge of the solved problem (e.g., applications in physics might be supported by a theory giving the probabilities).

In most cases, however, P is estimated from sampled data by

$$\bar{P}(Y = \mathbf{1} \mid X = \vec{x}) = \frac{\text{number of samples with } Y = \mathbf{1} \text{ and } X = \vec{x}}{\text{number of samples with } X = \vec{x}}$$

(We use  $\overline{P}$  to denote an estimate of P from data.)

#### Estimating P

Consider a problem with  $X = (X_1, X_2, X_3)$  where each  $X_i$  returns either 0 or 1. What might the data look like?

#### Estimating P

Consider a problem with  $X = (X_1, X_2, X_3)$  where each  $X_i$  returns either 0 or 1. What might the data look like?

Part of the data table:

All data with  $X_1 = 1$ ,  $X_2 = 0$ ,  $X_3 = 1$ :

Y	$X_1$	$X_2$	$X_3$	
1	1	0	1	
1	0	1 1		
0	1	1 0		
0	0	0	1	
1	0	0	0	
0	<b>0</b> 1		1	

Y	$X_1$	$X_2$	$X_3$	
1	1	0	1	
1	1	0	1	
0	1	0	1	
0	1	0	1	
1	1	0	1	
1	1	0	1	

Estimate:  $\bar{P}(\mathbf{1} \mid 1, 0, 1) = 2/3$ 

### Estimating P

Consider a problem with  $X = (X_1, X_2, X_3)$  where each  $X_i$  returns either 0 or 1. What might the data look like?

Part of the data table:

All data with  $X_1 = 1$ ,  $X_2 = 0$ ,  $X_3 = 1$ :



Y	$X_1$	$X_2$	<i>X</i> <sub>3</sub>	
1	1	0	1	
1	1	0	1	
0	1	0	1	
0	1	0	1	
1	1	0	1	
1	1	0	1	

Estimate:  $\bar{P}(\mathbf{1} \mid 1, 0, 1) = 2/3$ 

The probability table and the necessary data are typically too large!

Concretely, if all  $X_1, \ldots, X_n$  are binary, there are  $2^n$  probabilities  $P(Y = \mathbf{1} \mid X = \vec{x})$ , one for each possible  $\vec{x} \in \{0, 1\}^n$ .

Let's Look at It the Other Way Round

Theorem (Bayes, 1764)

$$P(A \mid B) = \frac{P(B \mid A) \cdot P(A)}{P(B)}$$

Let's Look at It the Other Way Round

Theorem (Bayes, 1764)

$$P(A \mid B) = \frac{P(B \mid A) \cdot P(A)}{P(B)}$$

Proof.

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)} = \frac{\frac{P(A \cap B)}{P(A)} \cdot P(A)}{P(B)} = \frac{P(B \mid A) \cdot P(A)}{P(B)}$$

#### **Bayesian Classification**

Determine the category for  $\vec{x}$  by computing

$$P(Y = y \mid X = \vec{x}) = \frac{P(Y = y) \cdot P(X = \vec{x} \mid Y = y)}{P(X = \vec{x})}$$

for both  $y \in \{0, 1\}$  and deciding whether or not the following holds:

$$P(Y = \mathbf{1} \mid X = \vec{x}) \ge P(Y = \mathbf{0} \mid X = \vec{x})$$

### **Bayesian Classification**

Determine the category for  $\vec{x}$  by computing

$$P(Y = y \mid X = \vec{x}) = \frac{P(Y = y) \cdot P(X = \vec{x} \mid Y = y)}{P(X = \vec{x})}$$

for both  $y \in \{0, 1\}$  and deciding whether or not the following holds:

$$P(Y = \mathbf{1} \mid X = \vec{x}) \ge P(Y = \mathbf{0} \mid X = \vec{x})$$

So, to make the classifier, we need to compute the following:

- The prior P(Y = 1) (then P(Y = 0) = 1 P(Y = 1))
- ► The conditionals  $P(X = \vec{x} | Y = y)$  for  $y \in \{0, 1\}$  and for every  $\vec{x}$

# Estimating the Prior and Conditionals

• P(Y = 1) can be easily estimated from data by

$$\bar{P}(Y = \mathbf{1}) = \frac{\text{number of samples with } Y = \mathbf{1}}{\text{number of all samples}}$$

# Estimating the Prior and Conditionals

• P(Y = 1) can be easily estimated from data by

$$\bar{P}(Y = \mathbf{1}) = \frac{\text{number of samples with } Y = \mathbf{1}}{\text{number of all samples}}$$

If the dimension of features is small, P(X = x | Y = y) can be estimated from data similarly as P(Y = 1 | X = x) by

$$\bar{P}(X = \vec{x} \mid Y = y) = \frac{\text{number of samples with } Y = y \text{ and } X = \vec{x}}{\text{number of samples with } Y = y}$$

Unfortunately, for higher dimensional data too many samples are needed to estimate all  $P(X = \vec{x} | Y = y)$  (there are too many  $\vec{x}$ 's).

So where is the advantage of using the Bayes thm.??

# Estimating the Prior and Conditionals

• P(Y = 1) can be easily estimated from data by

$$\bar{P}(Y = \mathbf{1}) = \frac{\text{number of samples with } Y = \mathbf{1}}{\text{number of all samples}}$$

► If the dimension of features is small, P(X = x | Y = y) can be estimated from data similarly as P(Y = 1 | X = x) by

$$\bar{P}(X = \vec{x} \mid Y = y) = \frac{\text{number of samples with } Y = y \text{ and } X = \vec{x}}{\text{number of samples with } Y = y}$$

Unfortunately, for higher dimensional data too many samples are needed to estimate all  $P(X = \vec{x} | Y = y)$  (there are too many  $\vec{x}$ 's).

So where is the advantage of using the Bayes thm.??

We introduce *independence assumptions* about the features!

### Naive Bayes

We assume that features are (conditionally) independent given the category. That is for all x = (x<sub>1</sub>,...,x<sub>n</sub>) and y ∈ {0,1} we assume:

$$P(X = x | Y = y) = P(X_1 = x_1, \cdots, X_n = x_n | Y)$$
$$= \prod_{i=1}^n P(X_i = x_i | Y = y)$$

# Naive Bayes

We assume that features are (conditionally) independent given the category. That is for all x = (x<sub>1</sub>,...,x<sub>n</sub>) and y ∈ {0,1} we assume:

$$P(X = x | Y = y) = P(X_1 = x_1, \cdots, X_n = x_n | Y)$$
$$= \prod_{i=1}^n P(X_i = x_i | Y = y)$$

► Therefore, we only need to specify P(X<sub>i</sub> = x<sub>i</sub> | Y = y) for each possible pair of a feature-value x<sub>i</sub> and y ∈ {0,1}.

### Naive Bayes

We assume that features are (conditionally) independent given the category. That is for all x = (x<sub>1</sub>,...,x<sub>n</sub>) and y ∈ {0,1} we assume:

$$P(X = x | Y = y) = P(X_1 = x_1, \cdots, X_n = x_n | Y)$$
$$= \prod_{i=1}^n P(X_i = x_i | Y = y)$$

► Therefore, we only need to specify P(X<sub>i</sub> = x<sub>i</sub> | Y = y) for each possible pair of a feature-value x<sub>i</sub> and y ∈ {0,1}.

Note that if all  $X_i$  are binary (values in  $\{0, 1\}$ ), this requires specifying only 2n parameters:

$$P(X_i = 1 \mid Y = \mathbf{1})$$
 and  $P(X_i = 1 \mid Y = \mathbf{0})$  for each  $X_i$ 

as 
$$P(X_i = 0 | Y = y) = 1 - P(X_i = 1 | Y = y)$$
 for  $y \in \{0, 1\}$ .

Compared to specifying  $2^n$  parameters without any independence assumption.

#### Estimating the marginal probabilities

Estimate the probabilities  $P(X_i = x_i | Y = y)$  by

$$\bar{P}(X_i = x_i \mid Y = y) = \frac{\text{number of samples with } X_i = x_i \text{ and } Y = y}{\text{number of samples with } Y = y}$$

#### Estimating the marginal probabilities

Estimate the probabilities  $P(X_i = x_i | Y = y)$  by

 $\bar{P}(X_i = x_i \mid Y = y) = \frac{\text{number of samples with } X_i = x_i \text{ and } Y = y}{\text{number of samples with } Y = y}$ 

**Example:** Consider a problem with  $X = (X_1, X_2, X_3)$  where each  $X_i$  returns either 0 or 1. The data is

Y	$X_1$	$X_2$	$X_3$	
1	1 1		1	
1	0	1	1	
0	1	0	1	
0	0	0	1	
1	0	0	0	
0	1	1	1	

$$\bar{P}(X_1 = 1 | Y = \mathbf{1}) = 1/3 \qquad \bar{P}(X_1 = 1 | Y = \mathbf{0}) = 2/3 
\bar{P}(X_2 = 1 | Y = \mathbf{1}) = 1/3 \qquad \bar{P}(X_2 = 1 | Y = \mathbf{0}) = 1/3 
\bar{P}(X_3 = 1 | Y = \mathbf{1}) = 2/3 \qquad \bar{P}(X_3 = 1 | Y = \mathbf{0}) = 1$$

# Naive Bayes - Example

Consider classification of geometric shapes:

- $X_1 \in \{ small, medium, large \}$
- $X_2 \in \{red, blue, green\}$
- $X_3 \in \{ square, triangle, circle \}$

## Naive Bayes – Example

Consider classification of geometric shapes:

- $X_1 \in \{ small, medium, large \}$
- $X_2 \in \{red, blue, green\}$
- $X_3 \in \{ square, triangle, circle \}$

Assume that we have already estimated the following probabilities:

	Y = <b>1</b>	Y = <b>0</b>
$\bar{P}(Y)$	0.5	0.5
$\bar{P}(small \mid Y)$	0.4	0.4
$\overline{P}(medium \mid Y)$	0.1	0.2
$\bar{P}(large \mid Y)$	0.5	0.4
$\bar{P}(red \mid Y)$	0.9	0.3
$\bar{P}(blue \mid Y)$	0.05	0.3
$\bar{P}(green \mid Y)$	0.05	0.4
$\bar{P}(square \mid Y)$	0.05	0.4
$\bar{P}(triangle \mid Y)$	0.05	0.3
$\bar{P}(circle \mid Y)$	0.9	0.3

Does (*medium*, *red*, *circle*) belong to the category **1** ?

	Y = <b>1</b>	Y = <b>0</b>
$\bar{P}(Y)$	0.5	0.5
$\bar{P}(medium \mid Y)$	0.1	0.2
$\bar{P}(red \mid Y)$	0.9	0.3
$\bar{P}(circle \mid Y)$	0.9	0.3

	Y = <b>1</b>	Y = <b>0</b>
$\bar{P}(Y)$	0.5	0.5
$\bar{P}(medium \mid Y)$	0.1	0.2
$\bar{P}(red \mid Y)$	0.9	0.3
$\bar{P}(circle \mid Y)$	0.9	0.3

$$P(Y = \mathbf{1} | X = \vec{x}) =$$
  
= P(1) · P(medium | 1) · P(red | 1) · P(circle | 1) / P(X = \vec{x})  
= 0.5 · 0.1 · 0.9 · 0.9 / P(X = \vec{x}) = 0.0405/P(X = \vec{x})

	Y = <b>1</b>	Y = <b>0</b>
$\bar{P}(Y)$	0.5	0.5
$\bar{P}(medium \mid Y)$	0.1	0.2
$\bar{P}(red \mid Y)$	0.9	0.3
$\overline{P}(circle \mid Y)$	0.9	0.3

$$P(Y = \mathbf{1} | X = \vec{x}) =$$
  
= P(\mathbf{1}) \cdot P(medium | \mathbf{1}) \cdot P(red | \mathbf{1}) \cdot P(circle | \mathbf{1}) / P(X = \vec{x})  
\delta 0.5 \cdot 0.1 \cdot 0.9 \cdot 0.9 / P(X = \vec{x}) = 0.0405 / P(X = \vec{x})

$$P(Y = \mathbf{0} | X = \vec{x}) =$$
  
= P(\mathbf{0}) \cdot P(medium | \mathbf{0}) \cdot P(red | \mathbf{0}) \cdot P(circle | \mathbf{0}) / P(X = \vec{x})  
\delta 0.5 \cdot 0.2 \cdot 0.3 \cdot 0.3 / P(X = \vec{x}) = 0.009 / P(X = \vec{x})

(Note that we used the estimates  $\overline{P}$  of P to finish the computation above.)

	Y = <b>1</b>	Y = <b>0</b>
$\bar{P}(Y)$	0.5	0.5
$\bar{P}(medium \mid Y)$	0.1	0.2
$\bar{P}(red \mid Y)$	0.9	0.3
$\overline{P}(circle \mid Y)$	0.9	0.3

$$P(Y = \mathbf{1} | X = \vec{x}) =$$
  
= P(\mathbf{1}) \cdot P(medium | \mathbf{1}) \cdot P(red | \mathbf{1}) \cdot P(circle | \mathbf{1}) / P(X = \vec{x})  
\delta 0.5 \cdot 0.1 \cdot 0.9 \cdot 0.9 / P(X = \vec{x}) = 0.0405 / P(X = \vec{x})

$$P(Y = \mathbf{0} | X = \vec{x}) =$$
  
= P(\mathbf{0}) \cdot P(medium | \mathbf{0}) \cdot P(red | \mathbf{0}) \cdot P(circle | \mathbf{0}) / P(X = \vec{x})  
\delta 0.5 \cdot 0.2 \cdot 0.3 \cdot 0.3 / P(X = \vec{x}) = 0.009 / P(X = \vec{x})

(Note that we used the estimates  $\bar{P}$  of P to finish the computation above.) Apparently,

$$P(Y = \mathbf{1} \mid X = \vec{x}) \doteq 0.0405 / P(X = \vec{x}) > 0.009 / P(X = \vec{x}) \doteq P(\mathbf{0} \mid X = \vec{x})$$

So we classify  $\vec{x}$  to the category **1**.

#### Estimating Probabilities in Practice

We already know that  $P(X_i = x_i | Y = y)$  can be estimated by

$$\bar{P}(X_i = x_i \mid Y = y) = \ell_{y, x_i} / \ell_y$$

where

#### Estimating Probabilities in Practice

We already know that  $P(X_i = x_i | Y = y)$  can be estimated by

$$\bar{P}(X_i = x_i \mid Y = y) = \ell_{y,x_i} / \ell_y$$

where

**Problem:** If, by chance, a rare value  $x_i$  of a feature  $X_i$  never occurs in the training data, we get

$$\bar{P}(X_i = x_i \mid Y = y) = 0$$
 for both  $y \in \{\mathbf{0}, \mathbf{1}\}$ 

But then  $\overline{P}(X = x) = 0$  for x containing the value  $x_i$  for  $X_i$ , and thus  $\overline{P}(Y = y \mid X = x)$  is not well defined. Moreover,  $\overline{P}(Y = y) \cdot \overline{P}(X = x \mid Y = y) = 0$  (for  $y \in \{0, 1\}$ ) so even this cannot be used for classification.

# Probability Estimation Example

Estimated probabilities:

					Y = <b>1</b>	Y = <b>0</b>
				$\overline{P}(Y)$	0.5	0.5
Training	data:			$\bar{P}(small \mid Y)$	0.5	0.5
Size	Color	Shape	Class	$\bar{P}(medium \mid Y)$	0	0
small	red	circle	1	$\bar{P}(large \mid Y)$	0.5	0.5
large	red	circle	1	$\overline{P}(red \mid Y)$	1	0.5
small	red	triangle	0	$\overline{P}(blue \mid Y)$	0	0.5
large	blue	circle	0	$\bar{P}(green \mid Y)$	0	0
				$\bar{P}(square \mid Y)$	0	0
				$\overline{P}(triangle \mid Y)$	0	0.5
				$\overline{P}(circle \mid Y)$	1	0.5

Note that  $\overline{P}(medium \mid \mathbf{1}) = P(medium \mid \mathbf{0}) = 0$  and thus also  $\overline{P}(medium, red, circle) = 0$ .

So what is  $\overline{P}(1 \mid medium, red, circle)$ ?

# Smoothing

To account for estimation from small samples, probability estimates are adjusted or *smoothed*.

### Smoothing

- To account for estimation from small samples, probability estimates are adjusted or *smoothed*.
- Laplace smoothing adds one to every count of feature values

$$\tilde{P}(X_i = x_i \mid Y = y) = \frac{\ell_{y,x_i} + 1}{\ell_y + v_i}$$

where

To understand note that

$$\ell_y = \sum_{x_i \text{ is a value of } X_i} \ell_{y,x_i}$$

and thus

$$\begin{split} \bar{P}(X_i = x_i \mid Y = y) &= \ell_{y, x_i} / \sum_{x_i \text{ is a value of } X_i} \ell_{y, x_i} \\ \tilde{P}(X_i = x_i \mid Y = y) &= (\ell_{y, x_i} + 1) / \sum_{x_i \text{ is a value of } X_i} (\ell_{y, x_i} + 1) \end{split}$$

# Laplace Smoothing Example

Assume training set contains 10 samples of category 1:

- 4 small
- 0 medium
- ▶ 6 large

# Laplace Smoothing Example

Assume training set contains 10 samples of category 1:

- 4 small
- 0 medium
- 6 large

Estimate parameters as follows

- $\tilde{P}(small \mid \mathbf{1}) = (4+1)/(10+3) = 0.384$
- $\tilde{P}(medium \mid \mathbf{1}) = (0+1)/(10+3) = 0.0769$
- $\tilde{P}(large \mid \mathbf{1}) = (6+1)/(10+3) = 0.538$

## **Continuous Features**

 $\Omega$  may be (potentially) continuous,  $X_i$  may assign a continuum of values in  $\mathbb{R}$ .
#### **Continuous Features**

 $\Omega$  may be (potentially) continuous,  $X_i$  may assign a continuum of values in  $\mathbb{R}$ .

The probabilities are computed using *probability density p*: ℝ → ℝ<sup>+</sup>.
 A random variable X : Ω → ℝ<sup>+</sup> has a density *p* : ℝ → ℝ<sup>+</sup> if for every interval [*a*, *b*] we have

$$P(a \le X \le b) = \int_a^b p(x) dx$$

Usually,  $P(X_i | Y = y)$  is used to denote the *density* of  $X_i$  conditioned on Y = y.

#### **Continuous Features**

 $\Omega$  may be (potentially) continuous,  $X_i$  may assign a continuum of values in  $\mathbb{R}$ .

The probabilities are computed using *probability density p*: ℝ → ℝ<sup>+</sup>.
 A random variable X : Ω → ℝ<sup>+</sup> has a density *p* : ℝ → ℝ<sup>+</sup> if for every interval [*a*, *b*] we have

$$P(a \le X \le b) = \int_a^b p(x) dx$$

Usually,  $P(X_i | Y = y)$  is used to denote the *density* of  $X_i$  conditioned on Y = y.

- The densities  $P(X_i | Y = y)$  are usually estimated using Gaussian densities as follows:
  - Estimate the mean μ<sub>iy</sub> and the standard deviation σ<sub>iy</sub> based on training data.

Then put

$$ar{P}(X_i \mid Y = y) = rac{1}{\sigma_{iy}\sqrt{2\pi}} \exp\left(rac{-(X_i - \mu_{iy})^2}{2\sigma_{iy}^2}
ight)$$

Tends to work well despite rather a strong assumption of conditional independence of features.

- Tends to work well despite rather a strong assumption of conditional independence of features.
- Experiments show that it is quite competitive with other classification methods.
   Even if the probabilities are not accurately estimated, it often picks the correct maximum probability category.

- Tends to work well despite rather a strong assumption of conditional independence of features.
- Experiments show that it is quite competitive with other classification methods.
   Even if the probabilities are not accurately estimated, it often picks the
  - correct maximum probability category.
- Directly constructs a model from parameter estimates that are calculated from the training data.

- Tends to work well despite rather a strong assumption of conditional independence of features.
- Experiments show that it is quite competitive with other classification methods.
   Even if the probabilities are not accurately estimated, it often picks the

correct maximum probability category.

- Directly constructs a model from parameter estimates that are calculated from the training data.
- Typically handles outliers and noise well.

- Tends to work well despite rather a strong assumption of conditional independence of features.
- Experiments show that it is quite competitive with other classification methods.
   Even if the probabilities are not accurately estimated, it often picks the

correct maximum probability category.

- Directly constructs a model from parameter estimates that are calculated from the training data.
- Typically handles outliers and noise well.
- Missing values are easy to deal with (simply average overall missing values in feature vectors).

In the Naive Bayes, we have assumed that *all* features  $X_1, \ldots, X_n$  are independent.

In the Naive Bayes, we have assumed that all features  $X_1, \ldots, X_n$  are independent.

This is usually not realistic.

E.g. Variables "rain" and "grass wet" are (usually) strongly dependent.

In the Naive Bayes, we have assumed that *all* features  $X_1, \ldots, X_n$  are independent.

This is usually not realistic.

E.g. Variables "rain" and "grass wet" are (usually) strongly dependent.

What if we return some dependencies? (But now in a well-defined sense.)

In the Naive Bayes, we have assumed that *all* features  $X_1, \ldots, X_n$  are independent.

This is usually not realistic.

E.g. Variables "rain" and "grass wet" are (usually) strongly dependent.

What if we return some dependencies?

(But now in a well-defined sense.)

Bayesian networks are a graphical model that uses a directed acyclic graph to specify dependencies among variables.

#### Bayesian Networks - Frample P(S = T) | P(S = F)P(C = T) P(C = F)0.8 0.2 0.02 0.98 Chair Sport P(B = T|C,S) P(B = F|C,S)S С С P(W = T|C) P(W = F|C)т т 0.9 0.1 т 0.9 0.1 т F 0.2 0.8 Worker Back F 0.01 0.99 F Т 0.9 0.1 F F 0.01 0.99 P(A = T|B) P(A = F|B)в т 0.7 0.3 Ache F 0.1 0.9

Now, e.g.,

 $P(C, S, W, B, A) = P(C) \cdot P(S) \cdot P(W \mid C) \cdot P(B \mid C, S) \cdot P(A \mid B)$ 

#### Bavesian Networks - Frample



Now, e.g.,

 $P(C, S, W, B, A) = P(C) \cdot P(S) \cdot P(W \mid C) \cdot P(B \mid C, S) \cdot P(A \mid B)$ 

Now, we may, e.g., infer the probability P(C = T | A = T) that we sit in the wrong chair, assuming that our back aches.

#### Bavesian Networks - Frample



Now, e.g.,

 $P(C, S, W, B, A) = P(C) \cdot P(S) \cdot P(W \mid C) \cdot P(B \mid C, S) \cdot P(A \mid B)$ 

Now, we may, e.g., infer the probability P(C = T | A = T) that we sit in the wrong chair, assuming that our back aches. We have to store only 10 numbers as opposed to  $2^5 - 1$  possible probabilities for all vectors of values of C, S, W, B, A.

#### Bayesian Networks – Learning & Naive Bayes

Many algorithms have been developed for learning:

- the structure of the graph of the network,
- the conditional probability tables.

The methods are based on maximum-likelihood estimation, gradient descent, etc.

Automatic procedures are usually combined with expert knowledge.

#### Bayesian Networks – Learning & Naive Bayes

Many algorithms have been developed for learning:

- the structure of the graph of the network,
- the conditional probability tables.

The methods are based on maximum-likelihood estimation, gradient descent, etc.

Automatic procedures are usually combined with expert knowledge.

Can you express the naive Bayes for  $Y, X_1, \ldots, X_n$  using a Bayesian network?

# **Classifier Evaluation**

Assume binary classification into two classes  $\{0, 1\}$ .

Assume binary classification into two classes  $\{0,1\}$ .

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in \{0, 1\}$  for all k.

Assume binary classification into two classes  $\{0,1\}$ .

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in \{0, 1\}$  for all k.

Consider a sequence of predictions generated by a classifier:

$$h_1,\ldots,h_p\in\{0,1\}$$

Here each  $h_k$  has been predicted for the k-the example  $(\vec{x}_k, c_k)$ .

Assume binary classification into two classes  $\{0,1\}$ .

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in \{0, 1\}$  for all k.

Consider a sequence of predictions generated by a classifier:

$$h_1,\ldots,h_p\in\{0,1\}$$

Here each  $h_k$  has been predicted for the k-the example  $(\vec{x}_k, c_k)$ .

How good are the predictions  $h_1, \ldots, h_p$  w.r.t.  $c_1, \ldots, c_p$ ?

There are many possible metrics ...

Assume binary classification into two classes  $\{0,1\}$ .

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in \{0, 1\}$  for all k.

Consider a sequence of predictions generated by a classifier:

 $h_1,\ldots,h_p\in\{0,1\}$ 

Here each  $h_k$  has been predicted for the k-the example  $(\vec{x}_k, c_k)$ .

How good are the predictions  $h_1, \ldots, h_p$  w.r.t.  $c_1, \ldots, c_p$ ?

There are many possible metrics ...

I will call the class 1 *positive* and the class 0 *negative*. Note that the class 0 is not negative in the numerical sense but in the absence of something (e.g., predicted illness).

		Pred	icted
		1	0
Actual	1	TP	FN
Actual	0	FP	ΤN

		Pred	icted
		1	0
Actual	1	ΤP	FN
Actual	0	FP	ΤN

 $\blacktriangleright$  TP = number of correctly classified examples with actual class 1

 $\mathsf{TP} = |\{k \mid h_k = 1 \land c_k = 1\}|$ 

		Pred	icted
		1	0
Actual	1	ΤP	FN
Actual	0	FP	ΤN

- ► TP = number of correctly classified examples with actual class 1 TP =  $|\{k \mid h_k = 1 \land c_k = 1\}|$
- ► TN = number of correctly classified examples with actual class 0 TN = |{k | h<sub>k</sub> = 0 ∧ c<sub>k</sub> = 0}|

		Predicted			
		1	0		
Actual	1	ΤP	FN		
Actual	0	FP	ΤN		

- ► TP = number of correctly classified examples with actual class 1 TP =  $|\{k \mid h_k = 1 \land c_k = 1\}|$
- ► TN = number of correctly classified examples with actual class 0 TN = |{k | h<sub>k</sub> = 0 ∧ c<sub>k</sub> = 0}|
- ► FP = number of incorrectly classified examples with actual class 0 FP = |{k | h<sub>k</sub> = 1 ∧ c<sub>k</sub> = 0}|

		Predicted			
		1	0		
Actual	1	ΤP	FN		
Actual	0	FP	ΤN		

- ► TP = number of correctly classified examples with actual class 1 TP =  $|\{k \mid h_k = 1 \land c_k = 1\}|$
- ► TN = number of correctly classified examples with actual class 0 TN = |{k | h<sub>k</sub> = 0 ∧ c<sub>k</sub> = 0}|
- ► FP = number of incorrectly classified examples with actual class 0 FP = |{k | h<sub>k</sub> = 1 ∧ c<sub>k</sub> = 0}|

► FN = number of correctly classified examples with actual class 1 FN =  $|\{k \mid h_k = 0 \land c_k = 1\}|$ 

#### Example

Given a sample of 12 individuals, eight have cancer, and four are cancer-free.

#### Example

Given a sample of 12 individuals, eight have cancer, and four are cancer-free.

Assume that we have trained a classifier with the following results:

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	1	1	1	0	0	0	0
Predicted	0	0	1	1	1	1	1	1	1	0	0	0
Result	FN	FN	TP	TP	TP	TP	ΤP	TP	FP	ΤN	ΤN	ΤN

#### Example

Given a sample of 12 individuals, eight have cancer, and four are cancer-free.

Assume that we have trained a classifier with the following results:

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	1	1	1	0	0	0	0
Predicted	0	0	1	1	1	1	1	1	1	0	0	0
Result	FN	FN	TP	TP	TP	TP	TP	TP	FP	TN	ΤN	ΤN

Actual condition	Predicted condition					
	Cancer	Non-cancer				
Cancer	TP = 6	FN = 2				
Non-cancer	FP = 1	TN = 3				
Total	8 + 4 = 12					

#### Terminology

- TP aka hit
- TN aka correct rejection
- ▶ FP aka type I error, false alarm, overestimation
- FN aka type II error, miss, underestimation

Usually, TP, TN, FP, and FN are used to denote the individual examples of a particular kind and the number of these examples.

#### Terminology

- TP aka hit
- TN aka correct rejection
- ▶ FP aka type I error, false alarm, overestimation
- FN aka type II error, miss, underestimation

Usually, TP, TN, FP, and FN are used to denote the individual examples of a particular kind and the number of these examples.

In what follows, we also use

- P = TP + FN of all cases with the *actual* class 1
- ▶ N = TN + FP of all cases with the *actual* class 0
- ▶ PP = TP + FP of all cases with the *predicted* class 1

▶ PN = TN + FN of all cases with the *predicted* class 0 Note that P + N = PP + PN is the number of all cases.

#### Terminology

- TP aka hit
- TN aka correct rejection
- ▶ FP aka type I error, false alarm, overestimation
- FN aka type II error, miss, underestimation

Usually, TP, TN, FP, and FN are used to denote the individual examples of a particular kind and the number of these examples.

In what follows, we also use

- P = TP + FN of all cases with the *actual* class 1
- N = TN + FP of all cases with the actual class 0
- ▶ PP = TP + FP of all cases with the *predicted* class 1

▶ PN = TN + FN of all cases with the *predicted* class 0 Note that P + N = PP + PN is the number of all cases.

There is a large number of derived metrics. We consider some of the most used in practice.

#### Accuracy

$$\mathsf{Accuracy} = \frac{\mathsf{TP} + \mathsf{TN}}{\mathsf{P} + \mathsf{N}}$$

Intuitively, Accuracy is the proportion of correctly classified cases w.r.t. all cases.

#### Accuracy

$$\mathsf{Accuracy} = \frac{\mathsf{TP} + \mathsf{TN}}{\mathsf{P} + \mathsf{N}}$$

Intuitively, Accuracy is the proportion of correctly classified cases w.r.t. all cases.

Example: Consider our cancer predictor with the confusion matrix

Actual condition	Predicted condition				
	Cancer	Non-cancer			
Cancer	TP = 6	FN = 2			
Non-cancer	FP = 1	TN = 3			
Total	8 + 4 = 12				

#### Accuracy

$$\mathsf{Accuracy} = \frac{\mathsf{TP} + \mathsf{TN}}{\mathsf{P} + \mathsf{N}}$$

Intuitively, Accuracy is the proportion of correctly classified cases w.r.t. all cases.

Example: Consider our cancer predictor with the confusion matrix

Actual condition	Predicted condition					
	Cancer	Non-cancer				
Cancer	TP = 6	FN = 2				
Non-cancer	FP = 1	TN = 3				
Total	8 + 4 = 12					

The Accuracy is

$$\mathsf{ACC} = \frac{\mathsf{TP} + \mathsf{TN}}{\mathsf{P} + \mathsf{N}} = \frac{6+3}{12} = \frac{3}{4}$$
Accuracy can be misleading when the classes are imbalanced:

- Consider 100 cases, 90 in the class 0 and 10 in the class 1,
- consider a classifier that returns 1 for a single sample of class 1 and 0 for all other samples.

Accuracy can be misleading when the classes are imbalanced:

- Consider 100 cases, 90 in the class 0 and 10 in the class 1,
- consider a classifier that returns 1 for a single sample of class 1 and 0 for all other samples.

Actual	Predicted	
	Pos	Neg
Pos	1	9
Neg	0	90
Total	90 +	10 = 100

Accuracy can be misleading when the classes are imbalanced:

- Consider 100 cases, 90 in the class 0 and 10 in the class 1,
- consider a classifier that returns 1 for a single sample of class 1 and 0 for all other samples.

Actual	Predicted	
	Pos Neg	
Pos	1	9
Neg	0	90
Total	90 +	10 = 100

The Accuracy is 91/100 > 0.9. Pretty good, right?

Accuracy can be misleading when the classes are imbalanced:

- Consider 100 cases, 90 in the class 0 and 10 in the class 1,
- consider a classifier that returns 1 for a single sample of class 1 and 0 for all other samples.

Actual	Predicted	
	Pos Neg	
Pos	1	9
Neg	0	90
Total	90 +	10 = 100

The Accuracy is 91/100 > 0.9. Pretty good, right?

However, the classifier is pretty bad in the positive cases. In the case of cancer prediction, such a classifier would be a disaster.

#### Precision & Recall

To mitigate the defect of the Accuracy, we may compute the following metrics:

$$Precision = \frac{TP}{PP} \quad (= how often is predicted positive actually positive)$$

Precision is also known as positive predictive value (PPV)

#### Precision & Recall

To mitigate the defect of the Accuracy, we may compute the following metrics:

$$Precision = \frac{TP}{PP} \quad (= how often is predicted positive actually positive)$$

Precision is also known as positive predictive value (PPV)

$$Recall = \frac{TP}{P} \quad (= how often is actually positive predicted positive)$$

Recall is also known as true positive rate, sensitivity, hit rate, and power.

Precision & Recall - Example

**Example:** In our cancer example:

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 +	4 = 12

Precision & Recall - Example

**Example:** In our cancer example:

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 + 4 = 12	

 Precision measures how often is the patient predicted to be ill truly ill (in our case, 6/7) Precision & Recall - Example

**Example:** In our cancer example:

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 +	4 = 12

- Precision measures how often is the patient predicted to be ill truly ill (in our case, 6/7)
- Recall measures how often is an ill patient found to be ill (in our case, 6/8)

Precision & Recall - Imbalanced Classes

Consider 100 cases, 90 in the class 0 and 10 in the class 1,

Precision & Recall - Imbalanced Classes

- Consider 100 cases, 90 in the class 0 and 10 in the class 1,
- consider a classifier that returns 1 for a single sample of class 1 and 0 for all other samples.

Actual	Predicted	
	Pos	Neg
Pos	1	9
Neg	0	90
Total	90 +	10 = 100

Precision & Recall - Imbalanced Classes

- Consider 100 cases, 90 in the class 0 and 10 in the class 1,
- consider a classifier that returns 1 for a single sample of class 1 and 0 for all other samples.

Actual	Predicted	
	Pos	Neg
Pos	1	9
Neg	0	90
Total	90 +	10 = 100

$$\begin{aligned} &\mathsf{Precision} = 1\\ &\mathsf{Recall} = \frac{1}{10} \end{aligned}$$

You can see that the predictor is very precise (on the class 1) but useless due to the weak Recall.

Let us get back to our cancer example:

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 +	4 = 12

Consider Precision and Recall.

By now, you should remember what they measure.

Let us get back to our cancer example:

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 +	4 = 12

Consider Precision and Recall.

By now, you should remember what they measure.

Which of the two is more important in medicine?

Let us get back to our cancer example:

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 +	4 = 12

Consider Precision and Recall.

By now, you should remember what they measure.

Which of the two is more important in medicine?

Which of the two is more important for plagiarism detectors?

Let us get back to our cancer example:

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 +	4 = 12

Consider Precision and Recall.

By now, you should remember what they measure.

Which of the two is more important in medicine?

Which of the two is more important for plagiarism detectors?

Can we get a single number summarizing both Precision and Recall?

For example, to compare two classifiers.

## $F_1$ Score

 $F_1$  score is the harmonic mean of Recall and Precision:

$$\mathsf{F}_{1} = \frac{2}{\mathsf{Recall}^{-1} + \mathsf{Precision}^{-1}} = \frac{2\mathsf{TP}}{2\mathsf{TP} + \mathsf{FP} + \mathsf{FN}}$$

## F<sub>1</sub> Score

 $F_1$  score is the harmonic mean of Recall and Precision:

$$F_1 = \frac{2}{\text{Recall}^{-1} + \text{Precision}^{-1}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$

#### Compare the arithmetic (left) and harmonic (right) mean:



The harmonic mean prefers the two values closer to each other. For example, the harmonic mean of 2/3 and 1/3 is (approx) 0.44444.

## F<sub>1</sub> Score - Examples

Consider the cancer example:

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 + 4 = 12	

Here  $F_1 = \frac{2\text{TP}}{2\text{TP}+\text{FP}+\text{FN}} = (2 \cdot 6)/((2 \cdot 6) + 1 + 2) = 0.8$ .

## F<sub>1</sub> Score - Examples

Consider the cancer example:

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 + 4 = 12	

Here  $F_1 = \frac{2\text{TP}}{2\text{TP}+\text{FP}+\text{FN}} = (2 \cdot 6)/((2 \cdot 6) + 1 + 2) = 0.8.$ 

Our imbalanced example:

Actual	Predicted	
	Pos	Neg
Pos	1	9
Neg	0	90
Total	90 +	10 = 100

Here  $F_1 = \frac{2\text{TP}}{2\text{TP}+\text{FP}+\text{FN}} = (2 \cdot 1)/((2 \cdot 1) + 0 + 9) = 0.18$ . Note that the average of Precision and Recall is 0.55, which would give us a much less severe warning that the classifier is bad.

Note that the standard definitions of Precision and Recall for binary classifiers reveal only part of the truth.

Note that the standard definitions of Precision and Recall for binary classifiers reveal only part of the truth.

In particular, *true negatives are not used* in the definition of  $F_1$ .

Note that the standard definitions of Precision and Recall for binary classifiers reveal only part of the truth.

In particular, *true negatives are not used* in the definition of  $F_1$ .

Consider

Actual	Predicted	
	Pos	Neg
Pos	90	0
Neg	9	1
Total	90 +	10 = 100

Note that the standard definitions of Precision and Recall for binary classifiers reveal only part of the truth.

In particular, *true negatives are not used* in the definition of  $F_1$ .

Consider

Actual	Predicted	
	Pos	Neg
Pos	90	0
Neg	9	1
Total	90 +	10 = 100

Precision = 90/99 Recall = 90/90  $F_1 = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}} = (2 \cdot 90)/(2 \cdot 90 + 9 + 0) = 0.95$ 

Note that the standard definitions of Precision and Recall for binary classifiers reveal only part of the truth.

In particular, *true negatives are not used* in the definition of  $F_1$ .

Consider

Actual	Predicted	
	Pos	Neg
Pos	90	0
Neg	9	1
Total	90 +	10 = 100

Precision = 90/99 Recall = 90/90  

$$F_1 = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}} = (2 \cdot 90)/(2 \cdot 90 + 9 + 0) = 0.95$$

All great, except that the classifier sucks on the negative cases. If you are concerned with the negative cases, swap the classes and compute another set of metrics.

## $F_1$ Score

*F*<sub>1</sub> is often used as a summary score for binary classifiers instead of Accuracy.

Works better with imbalanced classes.

## F<sub>1</sub> Score

*F*<sub>1</sub> is often used as a summary score for binary classifiers instead of Accuracy.
 Works better with imbalanced classes.

- Criticised for giving Precision and Recall the same importance.
- Is not symmetric, ignores true negatives, i.e., is misleading for some cases of imbalanced classes.

## F<sub>1</sub> Score

*F*<sub>1</sub> is often used as a summary score for binary classifiers instead of Accuracy.
 Works better with imbalanced classes.

Criticised for giving Precision and Recall the same importance.

Is not symmetric, ignores true negatives, i.e., is misleading for some cases of imbalanced classes.

 Fowlkes-Mallows index is a geometric mean of Precision and Recall (used in clustering).

The geometric mean is between the arithmetic and harmonic mean. For example, the geometric mean of 2/3 and 1/3 is (approx) 0.4714.

### More Derived Metrics

Positive predictive value (PPV),	False omission
precision	rate (FOR)
$=\frac{\mathrm{TP}}{\mathrm{PP}}=1-\mathrm{FDR}$	$=\frac{FN}{PN}=1-NPV$
False discovery rate (FDR) = $\frac{FP}{PP} = 1 - PPV$	Negative predictive value (NPV) = $\frac{TN}{PN} = 1 - FOR$

You can see that the negative predictive value becomes the Precision when we swap the classes (and vice versa).

## More Derived Metrics

True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{TP}{P} = 1 - FNR$	False negative rate (FNR), miss rate = $\frac{FN}{P} = 1 - TPR$
False positive rate (FPR),	True negative rate (TNR),
probability of false alarm, fall-out	specificity (SPC), selectivity
$=\frac{FP}{N}=1-TNR$	$=\frac{\mathrm{TN}}{\mathrm{N}}=1-\mathrm{FPR}$

Note that *specificity* becomes Recall when we swap the classes (and vice versa).

For example, medical doctors communicate in terms of *sensitivity* and *specificity*.

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 + 4 = 12	

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 + 4 = 12	

 $\mathsf{TPR} = \mathsf{Sensitivity} = \mathsf{Recall} = \mathsf{TP}/\mathsf{P} = 6/8$ 

How often is positive predicted positive?

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 + 4 = 12	

TPR = Sensitivity = Recall = TP/P = 6/8

How often is positive predicted positive?

$$TNR = Specificity = TN/N = 3/4$$

How often is negative predicted negative?

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 + 4 = 12	

TPR = Sensitivity = Recall = TP/P = 6/8

How often is positive predicted positive?

$$TNR = Specificity = TN/N = 3/4$$

How often is negative predicted negative?

FPR = Prob. of false alarm = FP/N = 1/4

How often is negative predicted positive?

Actual condition	Predicted condition	
	Cancer	Non-cancer
Cancer	TP = 6	FN = 2
Non-cancer	FP = 1	TN = 3
Total	8 + 4 = 12	

TPR = Sensitivity = Recall = TP/P = 6/8

How often is positive predicted positive?

$$TNR = Specificity = TN/N = 3/4$$

How often is negative predicted negative?

FPR = Prob. of false alarm = FP/N = 1/4

How often is negative predicted positive?

$$FNR = Miss rate = FN/P = 2/8$$

How often is positive predicted negative?

# **Evaluating Multi-class Classifiers**
Assume classification into classes from a finite set C.

Assume classification into classes from a finite set C.

Consider a classification dataset:

 $\{(\vec{x}_k,c_k)\mid k=1,\ldots,p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in C$  for all k.

Assume classification into classes from a finite set C.

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in C$  for all k.

Consider a sequence of predictions generated by a classifier:

$$h_1,\ldots,h_p\in C$$

Here each  $h_k$  has been predicted for the k-the example  $(\vec{x}_k, c_k)$ .

Assume classification into classes from a finite set C.

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in C$  for all k.

Consider a sequence of predictions generated by a classifier:

$$h_1,\ldots,h_p\in C$$

Here each  $h_k$  has been predicted for the k-the example  $(\vec{x}_k, c_k)$ .

How good are the predictions  $h_1, \ldots, h_p$  w.r.t.  $c_1, \ldots, c_p$ ? There are many possible metrics ...

Consider an arbitrary (finite) number of classes in C.

### Confusion Matrix

Assume that  $C = \{1, \ldots, m\}$ .

# Confusion Matrix

Assume that  $C = \{1, \ldots, m\}$ .

Now, given two classes  $i, j \in C$  we denote by  $M_{ij}$  the number of samples of class *i* classified into the class *j*.

# Confusion Matrix

Assume that  $C = \{1, \ldots, m\}$ .

Now, given two classes  $i, j \in C$  we denote by  $M_{ij}$  the number of samples of class *i* classified into the class *j*.

Formally,

$$M_{ij} = |\{k \mid c_k = i \land h_k = j\}|$$

Actual		P	redicte	ed	
	1	•••	j	•••	т
1	$M_{11}$	•••	$M_{1j}$	•••	$M_{1m}$
÷	÷		÷		÷
i	$M_{i1}$		M <sub>ij</sub>	• • •	M <sub>im</sub>
:	÷		÷		÷
т	$M_{m1}$	•••	M <sub>mj</sub>	•••	$M_{mm}$

# Example

Actual	Predicted
big	big
big	big
small	big
medium	medium
big	small
big	big
small	small
small	small
medium	medium
medium	small
small	small
big	big
medium	small
small	medium
big	big

# Example

Actual	Predicted
big	big
big	big
small	big
medium	medium
big	small
big	big
small	small
small	small
medium	medium
medium	small
small	small
big	big
medium	small
small	medium
big	big

Actual	Predicted											
	big	big medium sma										
big	5	0	1									
medium	0	2	2									
small	1	1	3									

Note that the diagonal counts the correctly classified samples.

The off-diagonal elements correspond to misclassified samples.

We can easily generalize Accuracy, Precision, Recall, and  $F_1$ -score from the binary classification to multiple classes.

We can easily generalize Accuracy, Precision, Recall, and  $F_1$ -score from the binary classification to multiple classes.

Notation

$$\blacktriangleright M_{i\bullet} = \sum_{j=1}^m M_{ij}$$

We can easily generalize Accuracy, Precision, Recall, and  $F_1$ -score from the binary classification to multiple classes.

Notation

$$M_{i\bullet} = \sum_{j=1}^{m} M_{ij}$$
$$M_{\bullet j} = \sum_{i=1}^{m} M_{ij}$$

We can easily generalize Accuracy, Precision, Recall, and  $F_1$ -score from the binary classification to multiple classes.

Notation

$$M_{i\bullet} = \sum_{j=1}^{m} M_{ij}$$

$$M_{\bullet j} = \sum_{i=1}^{m} M_{ij}$$

$$M_{\bullet \bullet} = \sum_{i=1}^{m} \sum_{j=1}^{m} M_{ij}$$

We can easily generalize Accuracy, Precision, Recall, and  $F_1$ -score from the binary classification to multiple classes.

Notation

$$M_{i\bullet} = \sum_{j=1}^{m} M_{ij}$$

$$M_{\bullet j} = \sum_{i=1}^{m} M_{ij}$$

$$M_{\bullet \bullet} = \sum_{i=1}^{m} \sum_{j=1}^{m} M_{ij}$$

Now, the metrics:

Accuracy = 
$$\frac{\sum_{k=1}^{m} M_{kk}}{M_{\bullet \bullet}}$$

We can easily generalize Accuracy, Precision, Recall, and  $F_1$ -score from the binary classification to multiple classes.

Notation

$$M_{i\bullet} = \sum_{j=1}^{m} M_{ij}$$

$$M_{\bullet j} = \sum_{i=1}^{m} M_{ij}$$

$$M_{\bullet \bullet} = \sum_{i=1}^{m} \sum_{j=1}^{m} M_{ij}$$

Now, the metrics:

Accuracy = 
$$\frac{\sum_{k=1}^{m} M_{kk}}{M_{\bullet \bullet}}$$

For a given class  $i \in C$ :

$$Precision[i] = \frac{M_{ii}}{M_{\bullet i}} \qquad \text{Recall}[i] = \frac{M_{ii}}{M_{i\bullet}}$$

We can easily generalize Accuracy, Precision, Recall, and  $F_1$ -score from the binary classification to multiple classes.

Notation

$$M_{i\bullet} = \sum_{j=1}^{m} M_{ij}$$

$$M_{\bullet j} = \sum_{i=1}^{m} M_{ij}$$

$$M_{\bullet \bullet} = \sum_{i=1}^{m} \sum_{j=1}^{m} M_{ij}$$

Now, the metrics:

Accuracy = 
$$\frac{\sum_{k=1}^{m} M_{kk}}{M_{\bullet \bullet}}$$

For a given class  $i \in C$ :

$$\begin{aligned} &\mathsf{Precision}[i] = \frac{M_{ii}}{M_{\bullet i}} \quad \mathsf{Recall}[i] = \frac{M_{ii}}{M_{i\bullet}} \\ &\mathsf{F}_1[i] = \frac{2 * \mathsf{Precision}[i] * \mathsf{Recall}[i]}{\mathsf{Precision}[i] + \mathsf{Recall}[i]} \end{aligned}$$

Note that Precision, Recall, and  $F_1$  can be defined only for a given class!

# Example

Actual	Predicted										
	big	big medium small									
big	5	0	1								
medium	0	2	2								
small	1	1	3								

Compute the metrics.

# Example

Accuracy = $(5+2+3)/15 = 0.66$	5
Precision[big] = 5/6	Actual
Precision[medium] = 2/3	
Precision[small] = 3/6	medium
Recall[big] = 5/6	small
Recall[medium] $= 2/4$	
Recall[small] = 3/5	
$F_1[\text{big}] = \frac{2 * (5/6) * (5/6)}{(5/6) + (5/6)} = 5$	/6 = 0.83
$F_1$ [medium] = 0.57	
$F_1$ [medium] = 0.54	

How do you get a single number out of these? Average Precision, Recall, and  $F_1$  are usually computed, but one needs to be careful about the variance.

Actual	Predicted										
	big	big medium small									
big	5	0	1								
medium	0	2	2								
small	1	1	3								





Machine learning/data mining is needed to understand the matrix.

# Probabilistic Classifier Evaluation

Assume binary classification into two classes  $\{0, 1\}$ .

Assume binary classification into two classes  $\{0, 1\}$ .

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in C$  for all k.

Assume binary classification into two classes  $\{0, 1\}$ .

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in C$  for all k.

Consider a sequence of predictions generated by a classifier. Now the classifier returns *probability of class* 1 for a given input:

$$h_1,\ldots,h_p\in[0,1]$$

Here each  $h_k$  has been predicted for the k-the example  $(\vec{x}_k, c_k)$ .

Assume binary classification into two classes  $\{0, 1\}$ .

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in C$  for all k.

Consider a sequence of predictions generated by a classifier. Now the classifier returns *probability of class* 1 for a given input:

 $h_1,\ldots,h_p\in[0,1]$ 

Here each  $h_k$  has been predicted for the k-the example  $(\vec{x}_k, c_k)$ .

How to interpret the predictions  $h_1, \ldots, h_p$ ?

Assume binary classification into two classes  $\{0, 1\}$ .

Consider a classification dataset:

 $\{(\vec{x}_k, c_k) \mid k = 1, \ldots, p\}$ 

Here  $\vec{x}_k$  is a vector of attributes/features and  $c_k \in C$  for all k.

Consider a sequence of predictions generated by a classifier. Now the classifier returns *probability of class* 1 for a given input:

 $h_1,\ldots,h_p\in[0,1]$ 

Here each  $h_k$  has been predicted for the k-the example  $(\vec{x}_k, c_k)$ .

How to interpret the predictions  $h_1, \ldots, h_p$ ? How good are the predictions  $h_1, \ldots, h_p$  w.r.t.  $c_1, \ldots, c_p$ ?

Let us fix predictions  $h_1, \ldots, h_p$ .

Let us fix predictions  $h_1, \ldots, h_p$ .

Given a threshold  $\mathcal{T} \in [0,1]$  we define

$$h_k^T = \begin{cases} 1 & \text{if } h_k \ge T \\ 0 & \text{if } h_k < T \end{cases}$$

For every T we can compute all the metrics (Precision, Recall, etc.)

Let us fix predictions  $h_1, \ldots, h_p$ .

Given a threshold  $\,\mathcal{T}\in[0,1]$  we define

$$h_k^T = \begin{cases} 1 & \text{if } h_k \ge T \\ 0 & \text{if } h_k < T \end{cases}$$

For every T we can compute all the metrics (Precision, Recall, etc.)

Given a metric MET and a threshold T, we denote by MET[T] the metric MET evaluated on  $h_1^T, \ldots, h_p^T$ .

Let us fix predictions  $h_1, \ldots, h_p$ .

Given a threshold  $\,\mathcal{T}\in[0,1]$  we define

$$h_k^T = \begin{cases} 1 & \text{if } h_k \ge T \\ 0 & \text{if } h_k < T \end{cases}$$

For every T we can compute all the metrics (Precision, Recall, etc.)

Given a metric MET and a threshold T, we denote by MET[T] the metric MET evaluated on  $h_1^T, \ldots, h_p^T$ .

We obtain

$$\mathsf{TP}[T] = |\{k \mid h_k^T = 1 \land c_k = 1\}|$$

Let us fix predictions  $h_1, \ldots, h_p$ .

Given a threshold  $\,\mathcal{T}\in[0,1]$  we define

$$h_k^T = \begin{cases} 1 & \text{if } h_k \ge T \\ 0 & \text{if } h_k < T \end{cases}$$

For every  $\mathcal{T}$  we can compute all the metrics (Precision, Recall, etc.)

Given a metric MET and a threshold T, we denote by MET[T] the metric MET evaluated on  $h_1^T, \ldots, h_p^T$ .

We obtain

$$\mathsf{TP}[T] = |\{k \mid h_k^T = 1 \land c_k = 1\}|$$

and

 $\mathsf{TN}[T], \mathsf{FP}[T], \mathsf{FN}[T], \mathsf{Accuracy}[T], \mathsf{Precision}[T], \mathsf{Recall}[T], F_1[T], \dots$ 

Let us fix predictions  $h_1, \ldots, h_p$ .

Given a threshold  $\,\mathcal{T}\in[0,1]$  we define

$$h_k^T = \begin{cases} 1 & \text{if } h_k \ge T \\ 0 & \text{if } h_k < T \end{cases}$$

For every  $\mathcal{T}$  we can compute all the metrics (Precision, Recall, etc.)

Given a metric MET and a threshold T, we denote by MET[T] the metric MET evaluated on  $h_1^T, \ldots, h_p^T$ .

We obtain

$$\mathsf{TP}[T] = |\{k \mid h_k^T = 1 \land c_k = 1\}|$$

and

 $TN[T], FP[T], FN[T], Accuracy[T], Precision[T], Recall[T], F_1[T], \dots$ 

However, all metrics are now functions of the threshold T.

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05
T=0.5	TP	TP	ΤP	TP	ΤP	ΤN	TN	FN	FN	TN	TN	ΤN
T=0.42	TP	TP	ΤP	TP	ΤP	FP	FP	ΤP	FN	TN	TN	ΤN
T=0.1	TP	TP	ΤP	TP	ΤP	FP	FP	TP	ΤP	FP	FP	ΤN

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05
T=0.5	ΤP	TP	ΤP	TP	ΤP	ΤN	TN	FN	FN	TN	ΤN	ΤN
T=0.42	TP	TP	ΤP	TP	ΤP	FP	FP	ΤP	FN	TN	ΤN	ΤN
T=0.1	TP	TP	TP	TP	TP	FP	FP	TP	ΤP	FP	FP	ΤN

For example, consider T = 0.42, then

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05
T=0.5	ΤP	TP	ΤP	TP	ΤP	ΤN	TN	FN	FN	TN	ΤN	ΤN
T=0.42	ΤP	TP	ΤP	TP	TP	FP	FP	TP	FN	ΤN	ΤN	ΤN
T=0.1	TP	TP	ΤP	TP	ΤP	FP	FP	TP	ΤP	FP	FP	ΤN

For example, consider T = 0.42, then

 $\mathsf{TP}[T] = 6$   $\mathsf{FP}[T] = 2$   $\mathsf{FN}[T] = 1$   $\mathsf{TN}[T] = 3$ 

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05
T=0.5	ΤP	TP	TP	TP	TP	ΤN	TN	FN	FN	TN	ΤN	ΤN
T=0.42	ΤP	TP	ΤP	TP	TP	FP	FP	ΤP	FN	TN	ΤN	ΤN
T=0.1	TP	TP	TP	TP	TP	FP	FP	TP	ΤP	FP	FP	ΤN

For example, consider T = 0.42, then

 $\mathsf{TP}[T] = 6$   $\mathsf{FP}[T] = 2$   $\mathsf{FN}[T] = 1$   $\mathsf{TN}[T] = 3$ 

Accuracy
$$[T] = rac{3+6}{12}$$
 Precision $[T] = rac{6}{6+2}$  Recall $[T] = rac{6}{6+1}$
#### Thresholded Classifier Metrics

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05
T=0.5	ΤP	TP	ΤP	TP	ΤP	ΤN	ΤN	FN	FN	TN	ΤN	ΤN
T=0.42	ΤP	TP	ΤP	TP	TP	FP	FP	TP	FN	ΤN	ΤN	ΤN
T=0.1	TP	TP	ΤP	TP	ΤP	FP	FP	TP	ΤP	FP	FP	ΤN

For example, consider T = 0.42, then

 $\mathsf{TP}[T] = 6$   $\mathsf{FP}[T] = 2$   $\mathsf{FN}[T] = 1$   $\mathsf{TN}[T] = 3$ 

Accuracy[T] = 
$$\frac{3+6}{12}$$
 Precision[T] =  $\frac{6}{6+2}$  Recall[T] =  $\frac{6}{6+1}$   
 $F_1[T] = \frac{2 \cdot 6/8 \cdot 6/7}{6/8 + 6/7} = 0.8$ 

Consider two metrics for a given T:

$$\mathsf{TPR}[T] = \frac{\mathsf{TP}[T]}{\mathsf{P}[T]} \qquad (\mathsf{True Positive Rate})$$

Consider two metrics for a given T:

$$\begin{aligned} \mathsf{TPR}[T] &= \frac{\mathsf{TP}[\mathsf{T}]}{\mathsf{P}[T]} & (\mathsf{True Positive Rate}) \\ \mathsf{FPR}[T] &= \frac{\mathsf{FP}[T]}{\mathsf{N}[T]} & (\mathsf{False Positive Rate}) \end{aligned}$$

Consider two metrics for a given T:

$$TPR[T] = \frac{TP[T]}{P[T]}$$
(True Positive Rate)  
$$FPR[T] = \frac{FP[T]}{N[T]}$$
(False Positive Rate)

ROC curve is then a function  $\mathsf{ROC}:[0,1] \to [0,1]^2$  defined by

ROC(T) = (TPR[T], FPR[T])

Consider two metrics for a given T:

$$TPR[T] = \frac{TP[T]}{P[T]}$$
(True Positive Rate)  
$$FPR[T] = \frac{FP[T]}{N[T]}$$
(False Positive Rate)

ROC curve is then a function  $\mathsf{ROC}:[0,1]\to[0,1]^2$  defined by

ROC(T) = (TPR[T], FPR[T])

Observe that

ROC(0) = (1, 1)

Because the classifier with T = 0 simply classifies everything as positive, i.e., into the class 1.

Both TPR[T] and FPR[T] are non-increasing in T.

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

▶  $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

▶  $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5

▶  $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

▶  $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5

▶  $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5

•  $0.15 < T \le 0.36$ : TPR = 1 and FPR = 2/5

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

•  $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5

•  $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5

- ▶  $0.15 < T \le 0.36$ : TPR = 1 and FPR = 2/5
- $0.36 < T \le 0.42$ : TPR = 6/7 and FPR = 2/5

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

- $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5
- $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5
- ▶  $0.15 < T \le 0.36$ : TPR = 1 and FPR = 2/5
- $0.36 < T \le 0.42$ : TPR = 6/7 and FPR = 2/5

▶ 0.42 <  $T \le$  0.43: TPR = 5/7 and FPR = 2/5

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

- $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5
- $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5
- 0.15 <  $T \le$  0.36: TPR = 1 and FPR = 2/5
- $0.36 < T \le 0.42$ : TPR = 6/7 and FPR = 2/5
- $0.42 < T \le 0.43$ : TPR = 5/7 and FPR = 2/5
- 0.43 <  $T \le$  0.48: TPR = 5/7 and FPR = 1/5

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

- $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5
- $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5
- $0.15 < T \le 0.36$ : TPR = 1 and FPR = 2/5
- $0.36 < T \le 0.42$ : TPR = 6/7 and FPR = 2/5
- $0.42 < T \le 0.43$ : TPR = 5/7 and FPR = 2/5
- ▶ 0.43 <  $T \le 0.48$ : TPR = 5/7 and FPR = 1/5
- $0.48 < T \le 0.66$ : TPR = 5/7 and FPR = 0

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

- $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5
- $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5
- $0.15 < T \le 0.36$ : TPR = 1 and FPR = 2/5
- $0.36 < T \le 0.42$ : TPR = 6/7 and FPR = 2/5
- $0.42 < T \le 0.43$ : TPR = 5/7 and FPR = 2/5
- 0.43 <  $T \le$  0.48: TPR = 5/7 and FPR = 1/5
- $0.48 < T \le 0.66$ : TPR = 5/7 and FPR = 0
- $0.66 < T \le 0.86$ : TPR = 4/7 and FPR = 0

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

- $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5
- $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5
- $0.15 < T \le 0.36$ : TPR = 1 and FPR = 2/5
- $0.36 < T \le 0.42$ : TPR = 6/7 and FPR = 2/5
- $0.42 < T \le 0.43$ : TPR = 5/7 and FPR = 2/5
- 0.43 <  $T \le$  0.48: TPR = 5/7 and FPR = 1/5
- 0.48 <  $T \le 0.66$ : TPR = 5/7 and FPR = 0
- $0.66 < T \le 0.86$ : TPR = 4/7 and FPR = 0
- $0.86 < T \le 0.90$ : TPR = 3/7 and FPR = 0

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

- $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5
- $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5
- $0.15 < T \le 0.36$ : TPR = 1 and FPR = 2/5
- $0.36 < T \le 0.42$ : TPR = 6/7 and FPR = 2/5
- $0.42 < T \le 0.43$ : TPR = 5/7 and FPR = 2/5
- 0.43 <  $T \le$  0.48: TPR = 5/7 and FPR = 1/5
- 0.48 <  $T \le 0.66$ : TPR = 5/7 and FPR = 0
- $0.66 < T \le 0.86$ : TPR = 4/7 and FPR = 0
- $0.86 < T \le 0.90$ : TPR = 3/7 and FPR = 0
- $0.90 < T \le 0.95$ : TPR = 2/7 and FPR = 0

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

- $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5
- $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5
- 0.15 <  $T \le$  0.36: TPR = 1 and FPR = 2/5
- $0.36 < T \le 0.42$ : TPR = 6/7 and FPR = 2/5
- ▶ 0.42 <  $T \le 0.43$ : TPR = 5/7 and FPR = 2/5
- 0.43 <  $T \le$  0.48: TPR = 5/7 and FPR = 1/5
- 0.48 <  $T \le 0.66$ : TPR = 5/7 and FPR = 0
- $0.66 < T \le 0.86$ : TPR = 4/7 and FPR = 0
- $0.86 < T \le 0.90$ : TPR = 3/7 and FPR = 0
- $0.90 < T \le 0.95$ : TPR = 2/7 and FPR = 0

•  $0.95 < T \le 0.98$ : TPR = 1/7 and FPR = 0

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

- $0.05 < T \le 0.10$ : TPR = 1 and FPR = 4/5
- $0.10 < T \le 0.15$ : TPR = 1 and FPR = 3/5
- ▶  $0.15 < T \le 0.36$ : TPR = 1 and FPR = 2/5
- $0.36 < T \le 0.42$ : TPR = 6/7 and FPR = 2/5
- $0.42 < T \le 0.43$ : TPR = 5/7 and FPR = 2/5
- 0.43 <  $T \le$  0.48: TPR = 5/7 and FPR = 1/5
- ▶ 0.48 <  $T \le 0.66$ : TPR = 5/7 and FPR = 0
- $0.66 < T \le 0.86$ : TPR = 4/7 and FPR = 0
- $0.86 < T \le 0.90$ : TPR = 3/7 and FPR = 0
- $0.90 < T \le 0.95$ : TPR = 2/7 and FPR = 0
- $0.95 < T \le 0.98$ : TPR = 1/7 and FPR = 0
- $0.98 < T \le 1.00$ : TPR = 0 and FPR = 0

ROC

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05



#### Iris Dataset - A Classifier



 $\mathsf{Example}$  from the scikit-learn manual -  $\mathsf{SVM}$  classifier trained in Iris

### Using ROC and Threshold



Search for the best threshold at the elbow of the ROC curve.

# **ROC** - Explanation



The larger the *area under the ROC curve (ROC-AUC)*, the better. ROC-AUC ranges from 0 to 1. ROC-AUC  $\approx$  0.5 indicates random guessing.

#### **ROC-AUC**

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05



 $\mathsf{ROC}\text{-}\mathsf{AUC} = 0.8857$ 

### Iris - ROC-AUC



ROC-AUC = 0.79

How is the ROC-AUC connected with the samples?

How is the ROC-AUC connected with the samples?

Consider our cancer detection example:

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

How is the ROC-AUC connected with the samples?

Consider our cancer detection example:

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

AUC has a probabilistic explanation:

Consider the following experiment:

How is the ROC-AUC connected with the samples?

Consider our cancer detection example:

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

AUC has a probabilistic explanation:

Consider the following experiment:

Choose randomly a patient *i* from positive patients
Each positive patient has the same probability of being chosen.

How is the ROC-AUC connected with the samples?

Consider our cancer detection example:

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

AUC has a probabilistic explanation:

Consider the following experiment:

- Choose randomly a patient *i* from positive patients Each positive patient has the same probability of being chosen.
- Choose randomly a patient j from negative patients Each negative patient has the same probability of being chosen.

How is the ROC-AUC connected with the samples?

Consider our cancer detection example:

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

AUC has a probabilistic explanation:

Consider the following experiment:

- Choose randomly a patient *i* from positive patients Each positive patient has the same probability of being chosen.
- Choose randomly a patient *j* from negative patients Each negative patient has the same probability of being chosen.
- Check if  $h_i > h_j$ .

How is the ROC-AUC connected with the samples?

Consider our cancer detection example:

Index	1	2	3	4	5	6	7	8	9	10	11	12
Actual	1	1	1	1	1	0	0	1	1	0	0	0
Predicted	.98	.95	.9	.86	.66	.48	.43	.42	.36	.15	.1	.05

AUC has a probabilistic explanation:

Consider the following experiment:

- Choose randomly a patient *i* from positive patients Each positive patient has the same probability of being chosen.
- Choose randomly a patient j from negative patients Each negative patient has the same probability of being chosen.
- Check if  $h_i > h_j$ .

The ROC-AUC is the probability of succeeding in the  $h_i > h_j$  test.

### Summary

We have discussed various metrics that can be used to evaluate the quality of a classifier.

The metrics summarize the results of evaluation on a given dataset.

### Summary

We have discussed various metrics that can be used to evaluate the quality of a classifier.

The metrics summarize the results of evaluation on a given dataset.

We have discussed metrics for evaluating

 binary classifiers, Accuracy, Precision, Recall, F<sub>1</sub>, and few more

multi-class classifiers,

Accuracy, Precision, Recall, F1

 probabilistic classifiers, parametrized metrics, ROC-AUC

### Summary

We have discussed various metrics that can be used to evaluate the quality of a classifier.

The metrics summarize the results of evaluation on a given dataset.

We have discussed metrics for evaluating

binary classifiers,

Accuracy, Precision, Recall, F1, and few more

multi-class classifiers,

Accuracy, Precision, Recall, F1

 probabilistic classifiers, parametrized metrics, ROC-AUC

There are still several questions unanswered:

- When to use the metrics.
- ▶ How to estimate the influence of sampling the dataset.

#### Use of Evaluation Metrics

In our case, the following scenarios are typical:

Final test: Evaluate the model on the test set (separated at the beginning of training) and then compute the metrics. May inform the user about the quality of the model.
#### Use of Evaluation Metrics

In our case, the following scenarios are typical:

- Final test: Evaluate the model on the test set (separated at the beginning of training) and then compute the metrics. May inform the user about the quality of the model.
- Validation: Evaluate models on a separate validation set and use the metrics to compare models. There are (at least) two scenarios in which this happens:
  - Hyperparameter fine-tuning.
  - Comparison of different models (e.g., KNN and decision trees).

#### Use of Evaluation Metrics

In our case, the following scenarios are typical:

Final test: Evaluate the model on the test set (separated at the beginning of training) and then compute the metrics. May inform the user about the quality of the model.

 Validation: Evaluate models on a separate validation set and use the metrics to compare models. There are (at least) two scenarios in which this happens:

- Hyperparameter fine-tuning.
- Comparison of different models (e.g., KNN and decision trees).

Keep in mind that the metrics are artificial, and the results of the model are roughly summarized.

It would be best if you always strived to test the proper functionality of your model in as natural conditions as possible.

For example, a model for medical diagnosis should be evaluated by medical doctors who may observe many features of its behavior that are difficult to express quantitatively.

Machine learning models are typically trained on (pseudo) random samples of data objects.

For example, a set of patients treated by the concrete hospital.

Machine learning models are typically trained on (pseudo) random samples of data objects.

For example, a set of patients treated by the concrete hospital.

However, the purpose of testing/evaluation is to get information about the whole population (i.e., all possible patients).

Machine learning models are typically trained on (pseudo) random samples of data objects.

For example, a set of patients treated by the concrete hospital.

However, the purpose of testing/evaluation is to get information about the whole population (i.e., all possible patients).

How do we estimate how much specific properties of the given sample influence our model?

This is a challenging question; methods of inferential statistics are needed to get the answer.

Machine learning models are typically trained on (pseudo) random samples of data objects.

For example, a set of patients treated by the concrete hospital.

However, the purpose of testing/evaluation is to get information about the whole population (i.e., all possible patients).

How do we estimate how much specific properties of the given sample influence our model?

This is a challenging question; methods of inferential statistics are needed to get the answer.

We will consider these issues in some later lecture. Concretely,

- ► *Bias-variance* tradeoff
- Statistical tests for testing
  - significance of the metrics values,
  - paired t-tests for comparing models.

Let us consider two classifiers. How do you compare them?

Let us consider two classifiers. How do you compare them?

Accuracies and  $F_1$  scores can be compared easily (they are just numbers).

Let us consider two classifiers. How do you compare them?

Accuracies and  $F_1$  scores can be compared easily (they are just numbers).

How to compare (Precision<sub>1</sub>, Recall<sub>1</sub>) of the fist classifier with (Precision<sub>2</sub>, Recall<sub>2</sub>) of the second classifier?

Let us consider two classifiers. How do you compare them?

Accuracies and  $F_1$  scores can be compared easily (they are just numbers).

How to compare (Precision<sub>1</sub>, Recall<sub>1</sub>) of the fist classifier with (Precision<sub>2</sub>, Recall<sub>2</sub>) of the second classifier?

Thresholding

- Introduce a threshold  $0 \le t \le 1$
- Demand, one of the two metrics (typically the Recall), to be at least t. That is

 $\operatorname{Recall}_1 \geq t \qquad \operatorname{Recall}_2 \geq t$ 

Compare the values of the other metric numerically. In our case, decide whether

```
Precision_1 \ge Precision_2
```

(Still need to be concerned about the statistical significance.)

# Example

Actual condition	Predicted condition		Actual condition	Predicted condition	
	Canc.	Non-canc.		Canc.	Non-canc.
Cancer	6	2	Cancer	5	3
Non-canc.	1	3	Non-canc.	0	4
Total	8 + 4 = 12		Total	8 + 4 = 12	

$$\begin{aligned} &\mathsf{Precision}_1 = \frac{6}{7} \qquad \mathsf{Recall}_1 = \frac{6}{8} \\ &\mathsf{Precision}_2 = \frac{5}{5} = 1 \qquad \mathsf{Recall}_2 = \frac{5}{8} \end{aligned}$$

Consider a threshold t on the Recall.

The second classifier is better if the threshold t is 5/8, then the second classifier is better.

If the threshold t is 6/8, then the second classifier is unacceptable.

#### Numerical features

► Throughout this lecture we assume that all features are numerical, i.e., feature vectors belong to ℝ<sup>n</sup>.

# Numerical features

- ► Throughout this lecture we assume that all features are numerical, i.e., feature vectors belong to ℝ<sup>n</sup>.
- Most non-numerical features can be conveniently transformed to numerical ones.

#### For example:

Colors {blue, red, yellow} can be represented by

 $\{(1,0,0),(0,1,0),(0,0,1)\}$ 

(one-hot encoding)

- Words can be embedded into vector spaces by various means (word2vec etc.)
- A black-and-white picture of x × y pixels can be encoded as a vector of xy numbers that capture the shades of gray of the pixels.

(Even though this is not the best way of representing images.)

# **Basic Problems**

We consider two basic problems:

► (Binary) classification

**Our goal:** Classify inputs into two categories.



# **Basic Problems**

We consider two basic problems:

(Binary) classification

**Our goal:** Classify inputs into two categories.





**Our goal:** Find a (hypothesized) functional dependency in data.



# Linear Models Binary Classification

# Binary classification in $\mathbb{R}^n$

Our goal:

• Given a set D of training examples of the form  $(\vec{x}, c)$  where  $\vec{x} \in \mathbb{R}^n$  and  $c \in \{0, 1\}$ ,

## Binary classification in $\mathbb{R}^n$

#### Our goal:

- Given a set *D* of training examples of the form  $(\vec{x}, c)$  where  $\vec{x} \in \mathbb{R}^n$  and  $c \in \{0, 1\}$ ,
- construct a model  $h : \mathbb{R}^n \to \{0, 1\}$  that is consistent with D, i.e.,

 $h(\vec{x}) = c$  for all training examples  $(\vec{x}, c) \in D$ 

# Binary classification in $\mathbb{R}^n$

#### Our goal:

- Given a set *D* of training examples of the form  $(\vec{x}, c)$  where  $\vec{x} \in \mathbb{R}^n$  and  $c \in \{0, 1\}$ ,
- construct a model  $h : \mathbb{R}^n \to \{0, 1\}$  that is consistent with D, i.e.,

 $h(\vec{x}) = c$  for all training examples  $(\vec{x}, c) \in D$ 

Comments:

In practice, we often do not strictly demand h(x) = c for all training examples (x, c) ∈ D (often it is impossible)

We are more interested in good generalization, that is how well h classifies new instances that do not belong to D.
(Recall that we usually evaluate accuracy of the resulting hypothesized function h on a test set.)

# Models

We consider two kinds of hypothesis spaces:

Linear (affine) classifiers (this lecture)



# Models

We consider two kinds of hypothesis spaces:

Linear (affine) classifiers (this lecture)



 Non-linear classifiers (kernel SVM, neural networks) (later lectures)



#### Linear Classifier – Example



Length and Scalar Product of Vectors

• We consider vectors  $\vec{x} = (x_1, \ldots, x_n) \in \mathbb{R}^m$ .

#### Length and Scalar Product of Vectors

- We consider vectors  $\vec{x} = (x_1, \ldots, x_n) \in \mathbb{R}^m$ .
- Euclidean metric on vectors:  $||\vec{x}|| = \sqrt{\sum_{i=1}^{n} x_i^2}$ The distance between two vectors (points)  $\vec{x}, \vec{y}$  is  $||\vec{x} - \vec{y}||$ .

#### Length and Scalar Product of Vectors

• We consider vectors 
$$\vec{x} = (x_1, \ldots, x_n) \in \mathbb{R}^m$$
.

- Euclidean metric on vectors:  $||\vec{x}|| = \sqrt{\sum_{i=1}^{n} x_i^2}$ The distance between two vectors (points)  $\vec{x}, \vec{y}$  is  $||\vec{x} - \vec{y}||$ .
- Scalar product  $\vec{x} \cdot \vec{y}$  of vectors  $\vec{x} = (x_1, \dots, x_n)$  and  $\vec{y} = (y_1, \dots, y_n)$  defined by

$$\vec{x} \cdot \vec{y} = \sum_{i=1}^{n} x_i y_i$$

Recall that x̄ · ȳ = ||x̄|| ||ȳ|| cos θ where θ is the angle between x̄ and ȳ. That is x̄ · ȳ is the length of the projection of ȳ on x̄ multiplied by ||x̄||.

• Note that 
$$\vec{x} \cdot \vec{x} = ||\vec{x}||^2$$

## Linear Classifier

A *linear classifier*  $h[\vec{w}]$  is determined by a vector of *weights*  $\vec{w} = (w_0, w_1, \dots, w_n) \in \mathbb{R}^{n+1}$  as follows:

#### Linear Classifier

A *linear classifier*  $h[\vec{w}]$  is determined by a vector of *weights*  $\vec{w} = (w_0, w_1, \dots, w_n) \in \mathbb{R}^{n+1}$  as follows:

Given  $\vec{x} = (x_1, \ldots, x_n) \in \mathbb{R}^n$ ,

$$h[\vec{w}](\vec{x}) := \begin{cases} 1 & w_0 + \sum_{i=1}^n w_i \cdot x_i \ge 0\\ 0 & w_0 + \sum_{i=1}^n w_i \cdot x_i < 0 \end{cases}$$

#### Linear Classifier

A *linear classifier*  $h[\vec{w}]$  is determined by a vector of *weights*  $\vec{w} = (w_0, w_1, \dots, w_n) \in \mathbb{R}^{n+1}$  as follows:

Given  $\vec{x} = (x_1, \ldots, x_n) \in \mathbb{R}^n$ ,

$$h[\vec{w}](\vec{x}) := \begin{cases} 1 & w_0 + \sum_{i=1}^n w_i \cdot x_i \ge 0\\ 0 & w_0 + \sum_{i=1}^n w_i \cdot x_i < 0 \end{cases}$$

More succinctly:

$$h(\vec{x}) = sgn\left(w_0 + \sum_{i=1}^n w_i \cdot x_i\right) \quad \text{where} \quad sgn(y) = \begin{cases} 1 & y \ge 0\\ 0 & y < 0 \end{cases}$$

We define separating hyperplane determined by  $\vec{w}$  as the set of all  $\vec{x} \in \mathbb{R}^n$  satisfying  $w_0 + \sum_{i=1}^n w_i \cdot x_i = 0$ .









#### Linear Classifier – Geometry



#### Linear Classifier – Notation

Given 
$$\vec{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$$
 we define an *augmented feature vector*

 $\mathbf{\tilde{x}} = (x_0, x_1, \dots, x_n)$  where  $x_0 = 1$ 

#### Linear Classifier – Notation

Given 
$$\vec{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$$
 we define an *augmented feature vector*  
 $\tilde{\mathbf{x}} = (x_0, x_1, \dots, x_n)$  where  $x_0 = 1$ 

This makes the notation for the linear classifier more succinct:

$$h[\vec{w}](\vec{x}) = sgn(\vec{w} \cdot \tilde{\mathbf{x}})$$

# Linear Classifier – Learning



- classification in the plane using a linear classifier
- if a point is incorrectly classified, the learning algorithm turns the line (hyperplane) to improve the classification
Given a training set

$$D = \{ (\vec{x}_1, c_1), (\vec{x}_2, c_2) \}, \dots, (\vec{x}_p, c_p) \}$$

Here  $\vec{x}_k = (x_{k1} \dots, x_{kn}) \in \mathbb{R}^n$  and  $c_k \in \{0, 1\}$ .

Given a training set

$$D = \{ (\vec{x}_1, c_1), (\vec{x}_2, c_2) \}, \dots, (\vec{x}_p, c_p) \}$$

Here  $\vec{x}_k = (x_{k1} \dots, x_{kn}) \in \mathbb{R}^n$  and  $c_k \in \{0, 1\}$ .

Recall that  $\mathbf{\tilde{x}}_k = (x_{k0}, x_{k1} \dots, x_{kn})$  where  $x_{k0} = 1$ .

Given a training set

$$D = \{ (\vec{x}_1, c_1), (\vec{x}_2, c_2) \}, \dots, (\vec{x}_p, c_p) \}$$

Here  $\vec{x}_k = (x_{k1} \dots, x_{kn}) \in \mathbb{R}^n$  and  $c_k \in \{0, 1\}$ . Recall that  $\tilde{\mathbf{x}}_k = (x_{k0}, x_{k1} \dots, x_{kn})$  where  $x_{k0} = 1$ .

• A weight vector  $\vec{w} \in \mathbb{R}^{n+1}$  is **consistent with** D if  $h[\vec{w}](\vec{x}_k) = sgn(\vec{w} \cdot \tilde{\mathbf{x}}_k) = c_k$  for all k = 1, ..., p

Given a training set

$$D = \{ (\vec{x}_1, c_1), (\vec{x}_2, c_2) \}, \dots, (\vec{x}_p, c_p) \}$$

Here  $\vec{x}_k = (x_{k1} \dots, x_{kn}) \in \mathbb{R}^n$  and  $c_k \in \{0, 1\}$ . Recall that  $\tilde{\mathbf{x}}_k = (x_{k0}, x_{k1} \dots, x_{kn})$  where  $x_{k0} = 1$ .

• A weight vector  $\vec{w} \in \mathbb{R}^{n+1}$  is **consistent with** D if

$$h[\vec{w}](\vec{x}_k) = sgn(\vec{w} \cdot \tilde{\mathbf{x}}_k) = c_k$$
 for all  $k = 1, \dots, p$ 

*D* is **linearly separable** if there is a vector  $\vec{w} \in \mathbb{R}^{n+1}$  which is consistent with *D*.

Given a training set

$$D = \{ (\vec{x}_1, c_1), (\vec{x}_2, c_2) \}, \dots, (\vec{x}_p, c_p) \}$$

Here  $\vec{x}_k = (x_{k1} \dots, x_{kn}) \in \mathbb{R}^n$  and  $c_k \in \{0, 1\}$ . Recall that  $\tilde{\mathbf{x}}_k = (x_{k0}, x_{k1} \dots, x_{kn})$  where  $x_{k0} = 1$ .

• A weight vector  $\vec{w} \in \mathbb{R}^{n+1}$  is **consistent with** *D* if

$$h[\vec{w}](\vec{x}_k) = sgn(\vec{w} \cdot \tilde{\mathbf{x}}_k) = c_k$$
 for all  $k = 1, \dots, p$ 

*D* is **linearly separable** if there is a vector  $\vec{w} \in \mathbb{R}^{n+1}$  which is consistent with *D*.

• Our goal is to find a consistent  $\vec{w}$  assuming that D is linearly separable.

### Online learning algorithm:

Idea: Cyclically go through the training examples in D and adapt weights. Whenever an example is incorrectly classified, turn the hyperplane so that the example becomes closer to its correct half-space.

### Online learning algorithm:

Idea: Cyclically go through the training examples in D and adapt weights. Whenever an example is incorrectly classified, turn the hyperplane so that the example becomes closer to its correct half-space.

Compute a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$ 

### Online learning algorithm:

Idea: Cyclically go through the training examples in D and adapt weights. Whenever an example is incorrectly classified, turn the hyperplane so that the example becomes closer to its correct half-space.

Compute a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$ 

•  $\vec{w}^{(0)}$  is randomly initialized close to  $\vec{0} = (0, ..., 0)$ 

### Online learning algorithm:

Idea: Cyclically go through the training examples in D and adapt weights. Whenever an example is incorrectly classified, turn the hyperplane so that the example becomes closer to its correct half-space.

Compute a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$ 

- $\vec{w}^{(0)}$  is randomly initialized close to  $\vec{0} = (0, ..., 0)$
- ▶ In (t+1)-th step,  $\vec{w}^{(t+1)}$  is computed as follows:

$$ec{w}^{(t+1)} = ec{w}^{(t)} - \varepsilon \cdot \left(h[ec{w}^{(t)}](ec{x}_k) - c_k
ight) \cdot \widetilde{\mathbf{x}}_k$$

### Online learning algorithm:

Idea: Cyclically go through the training examples in D and adapt weights. Whenever an example is incorrectly classified, turn the hyperplane so that the example becomes closer to its correct half-space.

Compute a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$ 

- $\vec{w}^{(0)}$  is randomly initialized close to  $\vec{0} = (0, \dots, 0)$
- ▶ In (t+1)-th step,  $\vec{w}^{(t+1)}$  is computed as follows:

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon \cdot \left(h[\vec{w}^{(t)}](\vec{x}_k) - c_k\right) \cdot \tilde{\mathbf{x}}_k$$
$$= \vec{w}^{(t)} - \varepsilon \cdot \left(sgn\left(\vec{w}^{(t)} \cdot \tilde{\mathbf{x}}_k\right) - c_k\right) \cdot \tilde{\mathbf{x}}_k$$

Here  $k = (t \mod p) + 1$ , i.e., the examples are considered cyclically, and  $0 < \varepsilon \le 1$  is a **learning rate**.

### Online learning algorithm:

Idea: Cyclically go through the training examples in D and adapt weights. Whenever an example is incorrectly classified, turn the hyperplane so that the example becomes closer to its correct half-space.

Compute a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$ 

- $\vec{w}^{(0)}$  is randomly initialized close to  $\vec{0} = (0, ..., 0)$
- ▶ In (t+1)-th step,  $\vec{w}^{(t+1)}$  is computed as follows:

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon \cdot \left(h[\vec{w}^{(t)}](\vec{x}_k) - c_k\right) \cdot \tilde{\mathbf{x}}_k$$
$$= \vec{w}^{(t)} - \varepsilon \cdot \left(sgn\left(\vec{w}^{(t)} \cdot \tilde{\mathbf{x}}_k\right) - c_k\right) \cdot \tilde{\mathbf{x}}_k$$

Here  $k = (t \mod p) + 1$ , i.e., the examples are considered cyclically, and  $0 < \varepsilon \le 1$  is a **learning rate**.

### Theorem (Rosenblatt)

If D is linearly separable, then there is  $t^*$  such that  $\vec{w}^{(t^*)}$  is consistent with D.

Example

Training set:

$$D = \{((2, -1), 1), ((2, 1), 1), ((1, 3), 0)\}$$

That is

$$\vec{x}_1 = (2, -1) \qquad \qquad \tilde{x}_1 = (1, 2, -1) \\ \vec{x}_2 = (2, 1) \qquad \qquad \tilde{x}_2 = (1, 2, 1) \\ \vec{x}_3 = (1, 3) \qquad \qquad \tilde{x}_3 = (1, 1, 3)$$

$$egin{array}{rcl} c_1 &=& 1 \ c_2 &=& 1 \ c_3 &=& 0 \end{array}$$

Assume that the initial vector  $\vec{w}^{(0)}$  is  $\vec{w}^{(0)} = (0, -1, 1)$ . Consider  $\varepsilon = 1$ .

## Example: Separating by $\vec{w}^{(0)}$



Denoting  $\vec{w}^{(0)} =$  $(w_0, w_1, w_2) = (0, -1, 1)$ the blue separating line is given by  $w_0 + w_1 x_1 + w_2 x_2 = 0$ .

The red vector normal to the blue line is  $(w_1, w_2)$ .

The points on the side of  $(w_1, w_2)$  are assigned 1 by the classifier, the others zero. (In this case  $\vec{x_3}$  is assigned one and  $\vec{x_1}, \vec{x_2}$  are assigned zero, all of this is inconsistent with  $c_1 = 1, c_2 = 1, c_3 = 0.$ )

## Example: Computing $\vec{w}^{(1)}$

We have

$$\vec{w}^{(0)} \cdot \tilde{\mathbf{x}}_1 = (0, -1, 1) \cdot (1, 2, -1) = 0 - 2 - 1 = -3$$

thus

$$sgn\left(ec{w}^{(0)}\cdot\widetilde{\mathbf{x}}_{1}
ight)=0$$

and thus

$$sgn\left(ec{w}^{(0)}\cdot\widetilde{\mathbf{x}}_{1}
ight)-c_{1}=0-1=-1$$

(I.e.,  $\vec{x_1}$  is not correctly classified, and  $\vec{w}^{(0)}$  is not consistent with D.) Hence,

$$\vec{w}^{(1)} = \vec{w}^{(0)} - \left( sgn\left( \vec{w}^{(0)} \cdot \tilde{\mathbf{x}}_1 \right) - c_1 \right) \cdot \tilde{\mathbf{x}}_1$$

$$= \vec{w}^{(0)} + \tilde{\mathbf{x}}_1$$

$$= (0, -1, 1) + (1, 2, -1)$$

$$= (1, 1, 0)$$

# Example: Separating by $\vec{w}^{(1)}$



## Example: Computing $\vec{w}^{(2)}$

We have

$$\vec{w}^{(1)} \cdot \tilde{\mathbf{x}}_2 = (1, 1, 0) \cdot (1, 2, 1) = 1 + 2 = 3$$

thus

$$sgn\left(ec{w}^{(1)}\cdot\widetilde{\mathbf{x}}_{2}
ight)=1$$

and thus

$$sgn\left(ec{w}^{(1)}\cdot\widetilde{\mathbf{x}}_{2}
ight)-c_{2}=1-1=0$$

(I.e.,  $\vec{x_2}$  is currently correctly classified by  $\vec{w}^{(1)}$ . However, as we will see,  $\vec{x_3}$  is not well classified.) Hence,

$$\vec{w}^{(2)} = \vec{w}^{(1)} = (1, 1, 0)$$

## Example: Computing $\vec{w}^{(3)}$

We have

$$\vec{w}^{(2)} \cdot \tilde{\mathbf{x}}_3 = (1, 1, 0) \cdot (1, 1, 3) = 1 + 1 = 2$$

thus

$$sgn\left(ec{w}^{(2)}\cdot\widetilde{\mathbf{x}}_{3}
ight)=1$$

and thus

$$sgn\left(ec{w}^{(2)}\cdot\widetilde{\mathbf{x}}_{3}
ight)-c_{3}=1-0=1$$

(This means that  $\vec{x}_3$  is not well classified, and  $\vec{w}^{(2)}$  is not consistent with D.) Hence,

$$\vec{w}^{(3)} = \vec{w}^{(2)} - \left(sgn\left(\vec{w}^{(2)} \cdot \tilde{\mathbf{x}}_{3}\right) - c_{3}\right) \cdot \tilde{\mathbf{x}}_{3}$$
  
=  $\vec{w}^{(2)} - \tilde{\mathbf{x}}_{3}$   
=  $(1, 1, 0) - (1, 1, 3)$   
=  $(0, 0, -3)$ 

## Example: Separating by $\vec{w}^{(3)}$



## Example: Computing $\vec{w}^{(4)}$

We have

$$\vec{w}^{(3)} \cdot \tilde{\mathbf{x}}_1 = (0, 0, -3) \cdot (1, 2, -1) = 3$$

thus

$$sgn\left(ec{w}^{(3)}\cdot\mathbf{\widetilde{x}}_{1}
ight)=1$$

and thus

$$sgn\left(ec{w}^{(3)}\cdot\mathbf{\tilde{x}}_{1}
ight)-c_{1}=1-1=0$$

(I.e.,  $\vec{x_1}$  is currently correctly classified by  $\vec{w}^{(3)}$ . However, we shall see that  $\vec{x_2}$  is not.) Hence,

$$\vec{w}^{(4)} = \vec{w}^{(3)} = (0, 0, -3)$$

## Example: Computing $\vec{w}^{(5)}$

We have

$$\vec{w}^{(4)} \cdot \tilde{\mathbf{x}}_2 = (0, 0, -3) \cdot (1, 2, 1) = -3$$

thus

$$sgn\left(ec{w}^{(4)}\cdot\widetilde{\mathbf{x}}_{2}
ight)=0$$

and thus

$$sgn\left(ec{w}^{(4)}\cdot\widetilde{\mathbf{x}}_{2}
ight)-c_{2}=0-1=-1$$

(I.e.,  $\vec{x_2}$  is not correctly classified, and  $\vec{w}^{(4)}$  is not consistent with D.) Hence,

$$\vec{w}^{(5)} = \vec{w}^{(4)} - \left(sgn\left(\vec{w}^{(4)} \cdot \tilde{\mathbf{x}}_2\right) - c_2\right) \cdot \tilde{\mathbf{x}}_2$$
  
=  $\vec{w}^{(4)} + \tilde{\mathbf{x}}_2$   
=  $(0, 0, -3) + (1, 2, 1)$   
=  $(1, 2, -2)$ 

# Example: Separating by $\vec{w}^{(5)}$



## Example: The result

The vector  $\vec{w}^{(5)}$  is consistent with *D*:

$$sgn\left(\vec{w}^{(5)} \cdot \tilde{\mathbf{x}}_{1}\right) = sgn\left((1, 2, -2) \cdot (1, 2, -1)\right) = sgn(7) = 1 = c_{1}$$
$$sgn\left(\vec{w}^{(5)} \cdot \tilde{\mathbf{x}}_{2}\right) = sgn\left((1, 2, -2) \cdot (1, 2, 1)\right) = sgn(3) = 1 = c_{2}$$
$$sgn\left(\vec{w}^{(5)} \cdot \tilde{\mathbf{x}}_{3}\right) = sgn\left((1, 2, -2) \cdot (1, 1, 3)\right) = sgn(-3) = 0 = c_{3}$$

### Batch learning algorithm:

Compute a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$ 

### Batch learning algorithm:

Compute a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$ 

•  $\vec{w}^{(0)}$  is randomly initialized close to  $\vec{0} = (0, \dots, 0)$ 

### Batch learning algorithm:

Compute a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$   $\vec{w}^{(0)}$  is randomly initialized close to  $\vec{0} = (0, \dots, 0)$   $\ln (t+1)$ -th step,  $\vec{w}^{(t+1)}$  is computed as follows:  $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^{p} \left(h[\vec{w}^{(t)}](\vec{x}_{k}) - c_{k}\right) \cdot \tilde{\mathbf{x}}_{k}$  $= \vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^{p} \left(sgn\left(\vec{w}^{(t)} \cdot \tilde{\mathbf{x}}_{k}\right) - c_{k}\right) \cdot \tilde{\mathbf{x}}_{k}$ 

Here  $0 < \varepsilon \leq 1$  is a learning rate.

### Batch learning algorithm:

Compute a sequence of weight vectors  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$   $\vec{w}^{(0)}$  is randomly initialized close to  $\vec{0} = (0, \dots, 0)$   $\ln (t+1)$ -th step,  $\vec{w}^{(t+1)}$  is computed as follows:  $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^{p} \left(h[\vec{w}^{(t)}](\vec{x}_{k}) - c_{k}\right) \cdot \tilde{\mathbf{x}}_{k}$  $= \vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^{p} \left(sgn\left(\vec{w}^{(t)} \cdot \tilde{\mathbf{x}}_{k}\right) - c_{k}\right) \cdot \tilde{\mathbf{x}}_{k}$ 

Here  $0 < \varepsilon \leq 1$  is a learning rate.

## Linear Regression – Oaks in Wisconsin

This example is from How to Lie with Statistics by Darrell Huff (1954)

Age	DBH
(years)	(inch)
97	12.5
93	12.5
88	8.0
81	9.5
75	16.5
57	11.0
52	10.5
45	9.0
28	6.0
15	1.5
12	1.0
11	1.0



## Linear Regression – Oaks in Wisconsin

This example is from How to Lie with Statistics by Darrell Huff (1954)





## Linear Regression - Oaks in Wisconsin

This example is from How to Lie with Statistics by Darrell Huff (1954)





### Our goal:

• Given a set *D* of training examples of the form  $(\vec{x}, f)$  where  $\vec{x} \in \mathbb{R}^n$  and  $f \in \mathbb{R}$ ,

### Our goal:

• Given a set *D* of training examples of the form  $(\vec{x}, f)$  where  $\vec{x} \in \mathbb{R}^n$  and  $f \in \mathbb{R}$ ,

▶ construct a model function  $h : \mathbb{R}^n \to \mathbb{R}$  such that

 $h(\vec{x}) \approx f$  for all training examples  $(\vec{x}, f) \in D$ 

Here  $\approx$  means that the values are somewhat close to each other w.r.t. an appropriate *error function E*.

### Our goal:

• Given a set *D* of training examples of the form  $(\vec{x}, f)$  where  $\vec{x} \in \mathbb{R}^n$  and  $f \in \mathbb{R}$ ,

• construct a model function  $h : \mathbb{R}^n \to \mathbb{R}$  such that

 $h(\vec{x}) \approx f$  for all training examples  $(\vec{x}, f) \in D$ 

Here  $\approx$  means that the values are somewhat close to each other w.r.t. an appropriate *error function E*.

In what follows we use the squared error defined by

$$E = \frac{1}{2} \sum_{(\vec{x}, f) \in D} (h(\vec{x}) - f)^2$$

Our goal is to minimize E.

The main reason is that this function has nice mathematical properties (as opposed, e.g., to  $\sum_{(\vec{x},f)\in D} |h(\vec{x}) - f|$ ).

## Linear Function Approximation

▶ Given a set *D* of training examples:

$$D = \{ (\vec{x}_1, f_1), (\vec{x}_2, f_2), \dots, (\vec{x}_p, f_p) \}$$

Here  $\vec{x}_k = (x_{k1} \dots, x_{kn}) \in \mathbb{R}^n$  and  $f_k \in \mathbb{R}$ .

### Linear Function Approximation

Given a set D of training examples:

 $D = \{ (\vec{x}_1, f_1), (\vec{x}_2, f_2), \dots, (\vec{x}_p, f_p) \}$ 

Here  $\vec{x}_k = (x_{k1} \dots, x_{kn}) \in \mathbb{R}^n$  and  $f_k \in \mathbb{R}$ .

► Our goal: Find w so that h[w](x<sub>k</sub>) = w · x<sub>k</sub> is close to f<sub>k</sub> for every k = 1,..., p. Recall that x<sub>k</sub> = (x<sub>k0</sub>, x<sub>k1</sub>..., x<sub>kn</sub>) where x<sub>k0</sub> = 1.

### Linear Function Approximation

Given a set D of training examples:

 $D = \{ (\vec{x}_1, f_1), (\vec{x}_2, f_2), \dots, (\vec{x}_p, f_p) \}$ 

Here  $\vec{x}_k = (x_{k1} \dots, x_{kn}) \in \mathbb{R}^n$  and  $f_k \in \mathbb{R}$ .

► Our goal: Find w so that h[w](x<sub>k</sub>) = w · x<sub>k</sub> is close to f<sub>k</sub> for every k = 1,..., p. Recall that x<sub>k</sub> = (x<sub>k0</sub>, x<sub>k1</sub>..., x<sub>kn</sub>) where x<sub>k0</sub> = 1.

Squared Error Function:

$$E(\vec{w}) = \frac{1}{2} \sum_{k=1}^{p} (\vec{w} \cdot \tilde{\mathbf{x}}_{k} - f_{k})^{2} = \frac{1}{2} \sum_{k=1}^{p} \left( \sum_{i=0}^{n} w_{i} x_{ki} - f_{k} \right)^{2}$$
# Error function



Consider the gradient of the error function:

$$\nabla E(\vec{w}) = \left(\frac{\partial E}{\partial w_0}(\vec{w}), \dots, \frac{\partial E}{\partial w_n}(\vec{w})\right) = \sum_{k=1}^p (\vec{w} \cdot \tilde{\mathbf{x}}_k - f_k) \cdot \tilde{\mathbf{x}}_k$$

What is the gradient  $\nabla E(\vec{w})$ ? It is a vector in  $\mathbb{R}^{n+1}$  which points in the direction of the steepest *ascent* of *E* (its length corresponds to the steepness). Note that here the vectors  $\tilde{\mathbf{x}}_k$  are *fixed* parameters of *E*!

Consider the gradient of the error function:

$$\nabla E(\vec{w}) = \left(\frac{\partial E}{\partial w_0}(\vec{w}), \dots, \frac{\partial E}{\partial w_n}(\vec{w})\right) = \sum_{k=1}^p \left(\vec{w} \cdot \tilde{\mathbf{x}}_k - f_k\right) \cdot \tilde{\mathbf{x}}_k$$

What is the gradient  $\nabla E(\vec{w})$ ? It is a vector in  $\mathbb{R}^{n+1}$  which points in the direction of the steepest *ascent* of *E* (its length corresponds to the steepness). Note that here the vectors  $\tilde{\mathbf{x}}_k$  are *fixed* parameters of *E*!

#### Fact: If $\nabla E(\vec{w}) = \vec{0} = (0, ..., 0)$ , then $\vec{w}$ is a global minimum of E.

This follows from the fact that E is a convex paraboloid that has a unique extreme, which is a minimum.



Consider n = 1, which means that  $\vec{w} = (w_0, w_1)$  and we write x instead of  $\vec{x}$  since  $\vec{x} \in \mathbb{R}^n = \mathbb{R}^1 = \mathbb{R}$ .

Then the model is  $h[\vec{w}](x) = w_0 + w_1 \cdot x$ .

Consider n = 1, which means that  $\vec{w} = (w_0, w_1)$  and we write x instead of  $\vec{x}$  since  $\vec{x} \in \mathbb{R}^n = \mathbb{R}^1 = \mathbb{R}$ .

Then the model is  $h[\vec{w}](x) = w_0 + w_1 \cdot x$ .

Consider a concrete training set:

$$\mathcal{T} = \{(2,1), (3,2), (4,5)\} \\ = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\}$$

Consider n = 1, which means that  $\vec{w} = (w_0, w_1)$  and we write x instead of  $\vec{x}$  since  $\vec{x} \in \mathbb{R}^n = \mathbb{R}^1 = \mathbb{R}$ .

Then the model is  $h[\vec{w}](x) = w_0 + w_1 \cdot x$ .

Consider a concrete training set:

$$\mathcal{T} = \{(2,1), (3,2), (4,5)\} \\ = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\}$$

$$E(w_0, w_1) = \frac{1}{2} [(w_0 + w_1 \cdot 2 - 1)^2 + (w_0 + w_1 \cdot 3 - 2)^2 + (w_0 + w_1 \cdot 4 - 5)^2]$$

Consider n = 1, which means that  $\vec{w} = (w_0, w_1)$  and we write x instead of  $\vec{x}$  since  $\vec{x} \in \mathbb{R}^n = \mathbb{R}^1 = \mathbb{R}$ .

Then the model is  $h[\vec{w}](x) = w_0 + w_1 \cdot x$ .

Consider a concrete training set:

$$\mathcal{T} = \{(2,1), (3,2), (4,5)\} \\ = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\}$$

The augmented feature vectors are: (1, 2), (1, 3), (1, 4).

$$E(w_0, w_1) = \frac{1}{2} [(w_0 + w_1 \cdot 2 - 1)^2 + (w_0 + w_1 \cdot 3 - 2)^2 + (w_0 + w_1 \cdot 4 - 5)^2]$$

 $\frac{\partial E}{\partial w_0}$ 

Consider n = 1, which means that  $\vec{w} = (w_0, w_1)$  and we write x instead of  $\vec{x}$  since  $\vec{x} \in \mathbb{R}^n = \mathbb{R}^1 = \mathbb{R}$ .

Then the model is  $h[\vec{w}](x) = w_0 + w_1 \cdot x$ .

Consider a concrete training set:

$$\mathcal{T} = \{(2,1), (3,2), (4,5)\} \\ = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\}$$

$$E(w_0, w_1) = \frac{1}{2} [(w_0 + w_1 \cdot 2 - 1)^2 + (w_0 + w_1 \cdot 3 - 2)^2 + (w_0 + w_1 \cdot 4 - 5)^2]$$

$$\frac{\partial E}{\partial w_0} = (w_0 + w_1 \cdot 2 - 1) \cdot 1 + (w_0 + w_1 \cdot 3 - 2) \cdot 1 + (w_0 + w_1 \cdot 4 - 5) \cdot 1$$

Consider n = 1, which means that  $\vec{w} = (w_0, w_1)$  and we write x instead of  $\vec{x}$  since  $\vec{x} \in \mathbb{R}^n = \mathbb{R}^1 = \mathbb{R}$ .

Then the model is  $h[\vec{w}](x) = w_0 + w_1 \cdot x$ .

Consider a concrete training set:

$$\mathcal{T} = \{(2,1), (3,2), (4,5)\} \\ = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\}$$

$$E(w_0, w_1) = \frac{1}{2} [(w_0 + w_1 \cdot 2 - 1)^2 + (w_0 + w_1 \cdot 3 - 2)^2 + (w_0 + w_1 \cdot 4 - 5)^2]$$

$$\frac{\partial E}{\partial w_0} = (w_0 + w_1 \cdot 2 - 1) \cdot 1 + (w_0 + w_1 \cdot 3 - 2) \cdot 1 + (w_0 + w_1 \cdot 4 - 5) \cdot 1$$
$$\frac{\partial E}{\partial w_1}$$

Consider n = 1, which means that  $\vec{w} = (w_0, w_1)$  and we write x instead of  $\vec{x}$  since  $\vec{x} \in \mathbb{R}^n = \mathbb{R}^1 = \mathbb{R}$ .

Then the model is  $h[\vec{w}](x) = w_0 + w_1 \cdot x$ .

Consider a concrete training set:

$$\mathcal{T} = \{(2,1), (3,2), (4,5)\} \\ = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\}$$

$$E(w_0, w_1) = \frac{1}{2} [(w_0 + w_1 \cdot 2 - 1)^2 + (w_0 + w_1 \cdot 3 - 2)^2 + (w_0 + w_1 \cdot 4 - 5)^2]$$

$$\frac{\partial E}{\partial w_0} = (w_0 + w_1 \cdot 2 - 1) \cdot 1 + (w_0 + w_1 \cdot 3 - 2) \cdot 1 + (w_0 + w_1 \cdot 4 - 5) \cdot 1$$
$$\frac{\partial E}{\partial w_1} = (w_0 + w_1 \cdot 2 - 1) \cdot 2 + (w_0 + w_1 \cdot 3 - 2) \cdot 3 + (w_0 + w_1 \cdot 4 - 5) \cdot 4$$

Consider n = 1, which means that  $\vec{w} = (w_0, w_1)$  and we write x instead of  $\vec{x}$  since  $\vec{x} \in \mathbb{R}^n = \mathbb{R}^1 = \mathbb{R}$ .

Then the model is  $h[\vec{w}](x) = w_0 + w_1 \cdot x$ .

Consider a concrete training set:

$$\mathcal{T} = \{(2,1), (3,2), (4,5)\} \\ = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\}$$

$$E(w_0, w_1) = \frac{1}{2} [(w_0 + w_1 \cdot 2 - 1)^2 + (w_0 + w_1 \cdot 3 - 2)^2 + (w_0 + w_1 \cdot 4 - 5)^2]$$

$$\frac{\partial E}{\partial w_0} = (w_0 + w_1 \cdot 2 - 1) \cdot 1 + (w_0 + w_1 \cdot 3 - 2) \cdot 1 + (w_0 + w_1 \cdot 4 - 5) \cdot 1$$
$$\frac{\partial E}{\partial w_1} = (w_0 + w_1 \cdot 2 - 1) \cdot 2 + (w_0 + w_1 \cdot 3 - 2) \cdot 3 + (w_0 + w_1 \cdot 4 - 5) \cdot 4$$

$$\nabla E(\vec{w}) = \left(\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}\right) = (w_0 + w_1 \cdot 2 - 1) \cdot (1, 2) + (w_0 + w_1 \cdot 3 - 2) \cdot (1, 3) + (w_0 + w_1 \cdot 4 - 5) \cdot (1, 4)$$

#### **Gradient Descent:**

• Weights  $\vec{w}^{(0)}$  are initialized randomly close to  $\vec{0}$ .

#### **Gradient Descent:**

- Weights  $\vec{w}^{(0)}$  are initialized randomly close to  $\vec{0}$ .
- ► In (t + 1)-th step,  $\vec{w}^{(t+1)}$  is computed as follows:  $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon \cdot \nabla E(\vec{w}^{(t)})$

#### **Gradient Descent:**

- Weights  $\vec{w}^{(0)}$  are initialized randomly close to  $\vec{0}$ .
- ▶ In (t + 1)-th step,  $\vec{w}^{(t+1)}$  is computed as follows:

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon \cdot \nabla E(\vec{w}^{(t)})$$
  
=  $\vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^{p} \left( \vec{w}^{(t)} \cdot \tilde{\mathbf{x}}_{k} - f_{k} \right) \cdot \tilde{\mathbf{x}}_{k}$   
=  $\vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^{p} \left( h[\vec{w}^{(t)}](\vec{x}_{k}) - f_{k} \right) \cdot \tilde{\mathbf{x}}_{k}$ 

Here  $0 < \varepsilon \leq 1$  is a learning rate.

Note that the algorithm is almost similar to the batch perceptron algorithm!

#### **Gradient Descent:**

- Weights  $\vec{w}^{(0)}$  are initialized randomly close to  $\vec{0}$ .
- In (t+1)-th step,  $\vec{w}^{(t+1)}$  is computed as follows:

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon \cdot \nabla E(\vec{w}^{(t)})$$
  
=  $\vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^{p} \left( \vec{w}^{(t)} \cdot \tilde{\mathbf{x}}_{k} - f_{k} \right) \cdot \tilde{\mathbf{x}}_{k}$   
=  $\vec{w}^{(t)} - \varepsilon \cdot \sum_{k=1}^{p} \left( h[\vec{w}^{(t)}](\vec{x}_{k}) - f_{k} \right) \cdot \tilde{\mathbf{x}}_{k}$ 

Here  $0 < \varepsilon \leq 1$  is a learning rate.

Note that the algorithm is almost similar to the batch perceptron algorithm!

### Proposition

For sufficiently small  $\varepsilon > 0$  the sequence  $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \ldots$  converges (component-wisely) to the global minimum of E.

Training set:

$$D = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\} = \{(0, 0), (2, 1), (2, 2)\}$$

Note that input vectors are one dimensional, so we write them as numbers. That is

$$x_3 = 2$$
  $\tilde{\mathbf{x}}_3 = (1,2)$ 

$$f_1 = 0$$
  
 $f_2 = 1$   
 $f_3 = 2$ 

Assume that the initial vector  $\vec{w}^{(0)}$  is  $\vec{w}^{(0)} = (w_0^{(0)}, w_1^{(0)}) = (0, 2)$ . Consider  $\varepsilon = \frac{1}{10}$ .



Training set:  $D = \{(x_1, f_1), (x_2, f_2), (x_3, f_3)\} = \{(0, 0), (2, 1), (2, 2)\}$  Augmented input vectors:  $\mathbf{\tilde{x}}_1 = (1, 0), \ \mathbf{\tilde{x}}_2 = (1, 2), \ \mathbf{\tilde{x}}_1 = (1, 2)$ 

$$\nabla E(\vec{w}) = \left(\frac{\partial E}{\partial w_0}(\vec{w}), \frac{\partial E}{\partial w_1}(\vec{w})\right) = (w_0 + w_1 \cdot x_1 - f_1) \cdot \tilde{\mathbf{x}}_1 + (w_0 + w_1 \cdot x_2 - f_2) \cdot \tilde{\mathbf{x}}_2 + (w_0 + w_1 \cdot x_3 - f_3) \cdot \tilde{\mathbf{x}}_3$$

For  $\vec{w}^{(0)} = (0, 2)$  we have

$$\nabla E(\vec{w}^{(0)}) = (0 + 2 \cdot 0 - 0) \cdot (1, 0) + (0 + 2 \cdot 2 - 1) \cdot (1, 2) + (0 + 2 \cdot 2 - 2) \cdot (1, 2) = (3, 6) + (2, 4) = (5, 10)$$

Finally,  $\vec{w}^{(1)}$  is computed by

$$ec{w}^{(1)} = ec{w}^{(0)} - arepsilon \cdot 
abla E(ec{w}^{(0)}) = (0,2) - rac{1}{10} \cdot (5,10) = (-1/2,1)$$


















































#### Finding the Minimum in Dimension One

Assume n = 1. Then, the error function E is

$$E(w_0, w_1) = \frac{1}{2} \sum_{k=1}^{p} (w_0 + w_1 x_k - f_k)^2$$

#### Finding the Minimum in Dimension One

Assume n = 1. Then, the error function E is

$$E(w_0, w_1) = \frac{1}{2} \sum_{k=1}^{p} (w_0 + w_1 x_k - f_k)^2$$

Minimize E w.r.t.  $w_0$  a  $w_1$ :

$$\frac{\partial E}{\partial w_0} = 0 \quad \Leftrightarrow \quad w_0 = \bar{f} - w_1 \bar{x} \quad \Leftrightarrow \quad \bar{f} = w_0 + w_1 \bar{x}$$

where  $\bar{x} = \frac{1}{p} \sum_{k=1}^{p} x_k$  a  $\bar{f} = \frac{1}{p} \sum_{k=1}^{p} f_k$ 

#### Finding the Minimum in Dimension One

Assume n = 1. Then, the error function E is

$$E(w_0, w_1) = \frac{1}{2} \sum_{k=1}^{p} (w_0 + w_1 x_k - f_k)^2$$

Minimize E w.r.t.  $w_0$  a  $w_1$ :

$$\frac{\partial E}{\partial w_0} = 0 \quad \Leftrightarrow \quad w_0 = \bar{f} - w_1 \bar{x} \quad \Leftrightarrow \quad \bar{f} = w_0 + w_1 \bar{x}$$

where  $\bar{x} = \frac{1}{p} \sum_{k=1}^{p} x_k$  a  $\bar{f} = \frac{1}{p} \sum_{k=1}^{p} f_k$ 

$$\frac{\partial E}{\partial w_1} = 0 \quad \Leftrightarrow \quad w_1 = \frac{\frac{1}{p} \sum_{k=1}^{p} (f_k - \bar{f})(x_k - \bar{x})}{\frac{1}{p} \sum_{k=1}^{p} (x_k - \bar{x})^2}$$

i.e.  $w_1 = cov(f, x) / var(x)$ 











#### Maximum Likelihood vs Least Squares (Dim 1) Fix a training set $D = \{(x_1, f_1), (x_2, f_2), \dots, (x_p, f_p)\}$

Assume that each  $f_k$  has been generated randomly by

 $f_k = (\mathbf{w}_0 + \mathbf{w}_1 \cdot \mathbf{x}_k) + \epsilon_k$ 

where  $w_0, w_1$  are **unknown weights**, and  $\epsilon_k$  are independent, normally distributed noise values with mean 0 and some variance  $\sigma^2$ 



How "probable" is it to generate the correct  $f_1, \ldots, f_p$  ?

#### Maximum Likelihood vs Least Squares (Dim 1)



How "probable" is it to generate the correct  $f_1, \ldots, f_p$  ?

The following conditions are equivalent:

- $\blacktriangleright$  w<sub>0</sub>, w<sub>1</sub> minimize the squared error E
- ▶  $w_0, w_1$  maximize the likelihood (i.e., the "probability") of generating the correct values  $f_1, \ldots, f_p$  using  $f_k = (w_0 + w_1 \cdot x_k) + \epsilon_k$

Linear models are parametric, i.e., they have a fixed form with a small number of parameters that need to be learned from data (as opposed, e.g. to decision trees where the structure is not fixed in advance).

- Linear models are parametric, i.e., they have a fixed form with a small number of parameters that need to be learned from data (as opposed, e.g. to decision trees where the structure is not fixed in advance).
- Linear models are stable, i.e., small variations in the training data have only limited impact on the learned model. (tree models typically vary more with the training data).

- Linear models are parametric, i.e., they have a fixed form with a small number of parameters that need to be learned from data (as opposed, e.g. to decision trees where the structure is not fixed in advance).
- Linear models are stable, i.e., small variations in the training data have only limited impact on the learned model. (tree models typically vary more with the training data).
- Linear models are less likely to overfit (low variance) the training data but sometimes tend to underfit (high bias).

- Linear models are parametric, i.e., they have a fixed form with a small number of parameters that need to be learned from data (as opposed, e.g. to decision trees where the structure is not fixed in advance).
- Linear models are stable, i.e., small variations in the training data have only limited impact on the learned model. (tree models typically vary more with the training data).
- Linear models are less likely to overfit (low variance) the training data but sometimes tend to underfit (high bias).
- Linear models are prone to outliers.

# Unsupervised Learning

#### Clustering

# Often data form clusters based on some notion of similarity.



#### Clustering

Often data form clusters based on some notion of similarity.



This means that the data distribution is *multimodal*, i.e., contains several regions of higher probability mass.

#### Clustering

Often data form clusters based on some notion of similarity.

basic on  $\sum_{\substack{0 \\ -25 \\ -50 \\ -75 \\ -25$ 

10.0 7.5 5.0

2.5

This means that the data distribution is *multimodal*, i.e., contains several regions of higher probability mass.

We aim to group data into clusters of "similar" examples without using any additional information. (no supervision).



#### Motivation

...

Clustering is useful, e.g., in

- Customer segmentation based on their purchases.
- Data exploration identify patterns in data
- Semi-supervised learning cluster labeled examples with the unlabeled ones
- Search engines searching for images similar to a given image
- Image segmentation

## Segmentation



Consider a dataset

 $D = \{\vec{x}_1, \ldots, \vec{x}_p\}$ 

Note that no target class/value is provided.

*Clustering* is a partition  $\mathcal{U} = \{U_1, \ldots, U_K\}$  of D into K clusters.

Consider a dataset

 $D = \{\vec{x}_1, \ldots, \vec{x}_p\}$ 

Note that no target class/value is provided.

*Clustering* is a partition  $\mathcal{U} = \{U_1, \ldots, U_K\}$  of D into K clusters.

How do we identify the clusters?

Consider a dataset

 $D = \{\vec{x}_1, \ldots, \vec{x}_p\}$ 

Note that no target class/value is provided.

*Clustering* is a partition  $\mathcal{U} = \{U_1, \ldots, U_K\}$  of D into K clusters.

How do we identify the clusters?

Assume that we have a distance measure d measuring how far apart the objects being clustered. We want close objects to be clustered together.

Consider a dataset

 $D = \{\vec{x}_1, \ldots, \vec{x}_p\}$ 

Note that no target class/value is provided.

*Clustering* is a partition  $\mathcal{U} = \{U_1, \ldots, U_K\}$  of D into K clusters.

How do we identify the clusters?

Assume that we have a distance measure d measuring how far apart the objects being clustered. We want close objects to be clustered together.

#### For concreteness:

- ▶ We stick with numerical features, which means that the dataset  $D = {\vec{x_1}, ..., \vec{x_p}}$  contains vectors  $\vec{x_i} \in \mathbb{R}^n$ .
- Assume the Euclidean distance *d*.

Note that clustering may be based on completely different similarity/dissimilarity measures and non-numerical data.

## K-Means Clustering

The K-means clustering model consists of

• The number of clusters K

The K-means clustering model consists of

- ► The number of clusters K
- *K* cluster *prototypes*  $\vec{m}_1, \ldots, \vec{m}_K \in \mathbb{R}^n$

The K-means clustering model consists of

- The number of clusters K
- *K* cluster *prototypes*  $\vec{m}_1, \ldots, \vec{m}_K \in \mathbb{R}^n$
- ▶ An assignment  $q_{ij} \in \{0, 1\}$  for i = 1, ..., p and j = 1, ..., K of inputs  $\vec{x_i}$  to clusters  $U_j$  so that

$$\sum_j q_{ij} = 1$$
 for  $i = 1, \dots, p$ 

The K-means clustering model consists of

- The number of clusters K
- *K* cluster *prototypes*  $\vec{m}_1, \ldots, \vec{m}_K \in \mathbb{R}^n$
- ▶ An assignment  $q_{ij} \in \{0, 1\}$  for i = 1, ..., p and j = 1, ..., K of inputs  $\vec{x_i}$  to clusters  $U_j$  so that

$$\sum_j q_{ij} = 1$$
 for  $i = 1, \dots, p$ 

A given assignment  $\{q_{ij}\}$  induces a clustering

$$\mathcal{U} = \{U_1, \ldots, U_K\}$$
 where  $ec{x_i} \in U_j$  iff  $q_{ij} = 1$
### K-means clustering

The K-means clustering model consists of

- The number of clusters K
- *K* cluster *prototypes*  $\vec{m}_1, \ldots, \vec{m}_K \in \mathbb{R}^n$
- ▶ An assignment  $q_{ij} \in \{0, 1\}$  for i = 1, ..., p and j = 1, ..., K of inputs  $\vec{x_i}$  to clusters  $U_j$  so that

$$\sum_j q_{ij} = 1$$
 for  $i = 1, \dots, p$ 

A given assignment  $\{q_{ij}\}$  induces a clustering

$$\mathcal{U} = \{ U_1, \ldots, U_{\mathcal{K}} \}$$
 where  $ec{x_i} \in U_j$  iff  $q_{ij} = 1$ 

How good is a given model?

#### Error Function

Measure the distance of inputs  $\vec{x}_i$  to their cluster prototypes  $\vec{m}_i$ 



#### **Error Function**

Measure the distance of inputs  $\vec{x}_i$  to their cluster prototypes  $\vec{m}_i$ 



$$E(\{q_{ij}\},\{\vec{m}_j\}) = \sum_{i=1}^{r} \sum_{j=1}^{N} q_{ij} d(\vec{x}_i,\vec{m}_j)^2$$

#### **Error Function**

Measure the distance of inputs  $\vec{x}_i$  to their cluster prototypes  $\vec{m}_j$ 



We aim to *minimize* this error, i.e., to find proper positions of cluster prototypes and their assignment to minimize the total squared distance of examples to their prototypes.

The Problem: Minimize

$$E(\{q_{ij}\},\{\vec{m}_j\}) = \sum_{i=1}^{p} \sum_{j=1}^{K} q_{ij} d(\vec{x}_i,\vec{m}_j)^2 \text{ w.r.t. } \{q_{ij}\},\{\vec{m}_j\}$$

The Problem: Minimize

$$E(\{q_{ij}\},\{\vec{m}_j\}) = \sum_{i=1}^{p} \sum_{j=1}^{K} q_{ij} d(\vec{x}_i,\vec{m}_j)^2 \text{ w.r.t. } \{q_{ij}\},\{\vec{m}_j\}$$

Note that

If we fix {m<sub>j</sub>}, we can minimize E({q<sub>ij</sub>}, {m<sub>j</sub>}) by setting q<sub>ij</sub> = 1 iff m<sub>j</sub> is the *closest prototype* to x<sub>i</sub>.

The Problem: Minimize

$$E(\{q_{ij}\},\{\vec{m}_j\}) = \sum_{i=1}^{p} \sum_{j=1}^{K} q_{ij} d(\vec{x}_i,\vec{m}_j)^2 \text{ w.r.t. } \{q_{ij}\},\{\vec{m}_j\}$$

Note that

- ▶ If we fix  $\{\vec{m}_j\}$ , we can minimize  $E(\{q_{ij}\}, \{\vec{m}_j\})$  by setting  $q_{ij} = 1$  iff  $\vec{m}_j$  is the *closest prototype* to  $\vec{x}_i$ .
- If we fix {q<sub>ij</sub>}, we can minimize E({q<sub>ij</sub>}, {m<sub>j</sub>}) by letting each m<sub>j</sub> to minimize the total squared distance to its prototypes:

$$\sum_i q_{ij} d(\vec{x}_i, \vec{m}_j)^2$$

The Problem: Minimize

$$E(\{q_{ij}\},\{\vec{m}_j\}) = \sum_{i=1}^{p} \sum_{j=1}^{K} q_{ij} d(\vec{x}_i,\vec{m}_j)^2 \text{ w.r.t. } \{q_{ij}\},\{\vec{m}_j\}$$

Note that

- If we fix {m<sub>j</sub>}, we can minimize E({q<sub>ij</sub>}, {m<sub>j</sub>}) by setting q<sub>ij</sub> = 1 iff m<sub>j</sub> is the *closest prototype* to x<sub>j</sub>.
- If we fix {q<sub>ij</sub>}, we can minimize E({q<sub>ij</sub>}, {m<sub>j</sub>}) by letting each m<sub>j</sub> to minimize the total squared distance to its prototypes:

$$\sum_i q_{ij} d(\vec{x}_i, \vec{m}_j)^2$$

This is achieved by putting each prototype  $\vec{m}_j$  into the centroid of all inputs it represents:

$$\vec{m}_j = \frac{1}{\sum_{i=1}^p q_{ij}} \sum_{i=1}^p q_{ij} \vec{x}_i$$

Note that  $\sum_{i=1}^{p} q_{ij}$  is the size of the cluster represented by  $\vec{m_j}$ .

#### Algorithm 1 K-means clustering

- 1: Initialize K cluster centers  $\vec{m}_1, \vec{m}_2, \ldots, \vec{m}_K$  randomly
- 2: repeat
- 3: **for** each data point  $\vec{x_i}$  **do**

4: Assign  $\vec{x_i}$  to the nearest centroid, i.e., set  $q_{ij} = 1$  for

 $j = \arg\min_j d(\vec{x_i}, \vec{m_j})^2$ 

- 5: end for
- 6: **for** each cluster prototype  $\vec{m}_j$  **do**
- 7: Update  $\vec{m}_j$  to be the centroid of all points assigned to it

$$\vec{m}_j = \frac{1}{\sum_{i=1}^p q_{ij}} \sum_{i=1}^p q_{ij} \vec{x}_i$$

#### 8: end for

9: until convergence





Lines 3-5: Assign examples to the prototypes.



Lines 6-8: Move the prototypes to the centroids of their examples.



Lines 3-5: Assign examples to the prototypes.



Lines 6-8: Move the prototypes to the centroids of their examples.







### Convergence of K-means Clustering

Every step of K-means reduces the error  $E(\{q_{ij}\}, \{\vec{m}_j\})$ :

- ▶ We always assign an input vector to the closest prototype.
- We always move the prototype to be "closest" to the input vectors it represents.

### Convergence of K-means Clustering

Every step of K-means reduces the error  $E(\{q_{ij}\}, \{\vec{m}_j\})$ :

- ▶ We always assign an input vector to the closest prototype.
- We always move the prototype to be "closest" to the input vectors it represents.

Convergence can be tested by computing the error and checking whether it has not changed in the last step.

This will always happen after finitely many steps. There are only finitely many possible assignments to  $q_{ij}$ , and we always minimize the distance of inputs to their assigned centers.

#### Convergence of K-means Clustering

Every step of K-means reduces the error  $E(\{q_{ij}\}, \{\vec{m}_j\})$ :

- We always assign an input vector to the closest prototype.
- We always move the prototype to be "closest" to the input vectors it represents.

Convergence can be tested by computing the error and checking whether it has not changed in the last step.

This will always happen after finitely many steps. There are only finitely many possible assignments to  $q_{ij}$ , and we always minimize the distance of inputs to their assigned centers.

Example error development during training. Blue circles mean reassignment, and red circles mean moving prototypes.



#### Setting K - the Elbow Method

K-means clustering minimizes the *inertia* measure:

$$E(\{q_{ij}\},\{\vec{m}_j\}) = \sum_{i=1}^{p} \sum_{j=1}^{k} q_{ij} d(\vec{x}_i,\vec{m}_j)^2$$

That is the sum of squared distances of all examples of D to the cluster prototypes.

### Setting K - the Elbow Method

K-means clustering minimizes the *inertia* measure:

$$E(\{q_{ij}\},\{\vec{m}_j\}) = \sum_{i=1}^{p} \sum_{j=1}^{k} q_{ij} d(\vec{x}_i,\vec{m}_j)^2$$

That is the sum of squared distances of all examples of D to the cluster prototypes.

Note that the error does not consider the distance between the centers of the clusters.

Still, it is a valid measure that can be used to select the number of clusters.

#### Elbow Method

The following method for setting up the hyperparameters can be used in general. Let us illustrate the elbow method on K-means clustering with the inertia measure.

Consider the following data:



## Elbow Method



We could choose four clusters because adding more leads only to small decrements in the inertia.

#### **Bad Behavior**

Minimizing  $E(\{q_{ij}\}, \{\vec{m_j}\})$  starting from random positions of prototypes does not always produce "nice" results.

## **Bad Behavior**

Minimizing  $E(\{q_{ij}\}, \{\vec{m_j}\})$  starting from random positions of prototypes does not always produce "nice" results.

Some runs correspond to apparently bad solutions to the clustering problem even though a better solution exists.



Possible solution: Start the algorithm several times with random initialization of the prototypes.

- Prototype initialization is a big issue in K-means. There are various strategies. For example:
  - Start with all centers in a single corner.
  - Include randomness in the setting of centers throughout the algorithm.
  - Initialize sequentially, always fit prototypes, and then choose a new one as far away from the others as possible.
  - Use hierarchical clustering (next slides) to find clusters and initialize K-means with their centroids.

- Prototype initialization is a big issue in K-means. There are various strategies. For example:
  - Start with all centers in a single corner.
  - Include randomness in the setting of centers throughout the algorithm.
  - Initialize sequentially, always fit prototypes, and then choose a new one as far away from the others as possible.
  - Use hierarchical clustering (next slides) to find clusters and initialize K-means with their centroids.
- Empty clusters may occur need to resolve, e.g., by assigning the farthest point from any current prototype.

- Prototype initialization is a big issue in K-means. There are various strategies. For example:
  - Start with all centers in a single corner.
  - Include randomness in the setting of centers throughout the algorithm.
  - Initialize sequentially, always fit prototypes, and then choose a new one as far away from the others as possible.
  - Use hierarchical clustering (next slides) to find clusters and initialize K-means with their centroids.
- Empty clusters may occur need to resolve, e.g., by assigning the farthest point from any current prototype.
- As the squared error is behind the basic method, outliers may strongly affect its behavior (as in the linear regression case).

- Prototype initialization is a big issue in K-means. There are various strategies. For example:
  - Start with all centers in a single corner.
  - Include randomness in the setting of centers throughout the algorithm.
  - Initialize sequentially, always fit prototypes, and then choose a new one as far away from the others as possible.
  - Use hierarchical clustering (next slides) to find clusters and initialize K-means with their centroids.
- Empty clusters may occur need to resolve, e.g., by assigning the farthest point from any current prototype.
- As the squared error is behind the basic method, outliers may strongly affect its behavior (as in the linear regression case).
- Other problematic properties of data include
  - non-convex clusters
  - clusters of different sizes
  - non-linearly separable clusters
  - overlapping clusters









Consider a dataset

$$D = \{\vec{x}_1, \ldots, \vec{x}_p\}$$

Here  $\vec{x_i} \in \mathbb{R}^n$  for all i = 1, ..., p. Assume a distance d (e.g., Euclidean).

Consider a dataset

$$D = \{\vec{x}_1, \ldots, \vec{x}_p\}$$

Here  $\vec{x_i} \in \mathbb{R}^n$  for all i = 1, ..., p. Assume a distance d (e.g., Euclidean).

Idea:

- Start by merging the closest examples (w.r.t. d)
- Incrementally build larger clusters by merging smaller clusters.

Consider a dataset

$$D = \{\vec{x}_1, \ldots, \vec{x}_p\}$$

Here  $\vec{x_i} \in \mathbb{R}^n$  for all i = 1, ..., p. Assume a distance d (e.g., Euclidean).

Idea:

- Start by merging the closest examples (w.r.t. d)
- ▶ Incrementally build *larger clusters* by merging smaller clusters.

#### More concretely:

Maintain a set of clusters

Consider a dataset

$$D = \{\vec{x}_1, \ldots, \vec{x}_p\}$$

Here  $\vec{x_i} \in \mathbb{R}^n$  for all i = 1, ..., p. Assume a distance d (e.g., Euclidean).

Idea:

- Start by merging the closest examples (w.r.t. d)
- Incrementally build larger clusters by merging smaller clusters.

#### More concretely:

- Maintain a set of clusters
- lnitially, each  $\vec{x_i}$  in its own cluster

Consider a dataset

$$D = \{\vec{x}_1, \ldots, \vec{x}_p\}$$

Here  $\vec{x}_i \in \mathbb{R}^n$  for all i = 1, ..., p. Assume a distance d (e.g., Euclidean).

Idea:

- Start by merging the closest examples (w.r.t. d)
- Incrementally build larger clusters by merging smaller clusters.

#### More concretely:

- Maintain a set of clusters
- lnitially, each  $\vec{x_i}$  in its own cluster
- Repeat until only one cluster is left:
  - Pick two closest clusters
  - Merge them into a new cluster

Consider a dataset

$$D = \{\vec{x}_1, \ldots, \vec{x}_p\}$$

Here  $\vec{x_i} \in \mathbb{R}^n$  for all i = 1, ..., p. Assume a distance d (e.g., Euclidean).

Idea:

- Start by merging the closest examples (w.r.t. d)
- Incrementally build larger clusters by merging smaller clusters.

#### More concretely:

- Maintain a set of clusters
- lnitially, each  $\vec{x_i}$  in its own cluster
- Repeat until only one cluster is left:
  - Pick two closest clusters
  - Merge them into a new cluster

How do we determine the closest clusters?
# **Closest Clusters**



# **Closest Clusters**



# **Closest Clusters**











# Which One is Closer?



Consider two clusters  $U_j, U_k \subseteq D$ .

single\_linkage $(U_j, U_k)$ = min $\{d(\vec{x}, \vec{z}) \mid \vec{x} \in U_j, \vec{z} \in U_k\}$ 

 $\begin{aligned} & \text{complete\_linkage}(U_j, U_k) \\ &= \max\{d(\vec{x}, \vec{z}) \mid \vec{x} \in U_j, \vec{z} \in U_k\} \end{aligned}$ 

$$egin{array}{l} \mathsf{average\_linkage}(U_j, U_k) \ &= rac{1}{|U_j||U_k|} \sum_{ec{x} \in U_l} \sum_{ec{x} \in U_k} d(ec{x}, ec{z}) \end{array}$$

Each linkage can result in a different clustering.



Agglomerative Hierarchical Clustering Algorithm

Maintain a set of clusters Initially, each  $\vec{x_i}$  in its own cluster **repeat** 

Pick two closest clusters

Using the distance measure d and single, average, or complete linkage. Merge them into a new cluster until only one cluster is left

# Example



Point	x Coordinate	y Coordinate
p1	0.40	0.53
p2	0.22	0.38
p3	0.35	0.32
p4	0.26	0.19
p5	0.08	0.41
p6	0.45	0.30

	p1	p2	p3	p4	$p_5$	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

•1



d(3,6) = 0.11

which is the minimum distance between points.

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

•1



d(2,5) = 0.14

which is the second smallest distance.

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

•1



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



d(2,3) = 0.15 = $\min\{d(2,3), d(2,6), d(5,3), d(5,6)\}$ 

which is smaller than

$$d(1,2) = 0.24,$$
  
 $d(1,3) = 0.22,$   
 $d(4,2) = 0.2,$   
 $d(4,3) = 0.15,$   
 $d(4,1) = 0.37$   
the min. distances of points  
in all other pairs of clusters.

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



d(4,3) = 0.15= min{d(4,3), d(4,5), d(4,2), d(4,6)}

which is smaller than d(1,3) = 0.22, the distance of 1 to the cluster 3.

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00





	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

•1



d(2,5) = 0.14

which is second smallest distance.

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

The average distance between 4 and both points of  $\{3,6\}$  is

$$\frac{1}{2}(d(4,3)+d(4,6))=0.185$$

which is smaller than the average distance between all points of clusters 1, 2:

$$\frac{d(5,2) + d(5,3) + d(2,3) + d(2,6)}{4}$$

(equal to 0.205), and the average distance of 1 to any cluster.



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

•1



The average distance between clusters 2, 3 is 0.26 which is smaller than the average distance of 1 to any of the two clusters 1, 2(the average distances are 0.273 and 0.29).

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00





	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

•1





which is the minimum distance between points.

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

•1



d(2,5) = 0.14

which is the second smallest distance.

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00


	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



$$d(4,6) = 0.15 = \max\{d(4,3), d(4,6)\}$$

which is smaller than d(4,5) = 0.29, d(1,5) = 0.34, d(1,6) = 0.23, d(5,6) = 0.39, d(4,1) = 0.37the max distances of points in all other pairs of clusters.

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

•1



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



$$d(1,5) = 0.34$$

which is smaller than d(1, 4) = 0.37, d(5, 6) = 0.39, which are the maximum distances of points in all other pairs of clusters.

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00





# Properties of Agglomerative Hierarchical Clustering

- Provides hierarchy of clusters different cut levels provide different levels of coarseness of clusters
- Compared with k-means, it does not depend on the initialization and may provide better clusters than k-means.

# Properties of Agglomerative Hierarchical Clustering

- Provides hierarchy of clusters different cut levels provide different levels of coarseness of clusters
- Compared with k-means, it does not depend on the initialization and may provide better clusters than k-means.
- Lack of global objective function
  - The agglomerative hierarchical clustering uses local criteria to decide which clusters to merge.
- Agglomerative clustering has a "rich get richer" behavior that leads to uneven cluster sizes
- Merging decision cannot be undone bad for noisy data
- Computationally expensive.



#	A tibble:	50 × 20								
	state	homeo1	multi…²	income	med_i…³	poverty	fed_s…4	smoke	murder	robbery
	<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
1	Alabama	71.1	15.5	22984	42081	17.1	11.7	24.8	8.2	141.
2	Alaska	64.7	24.6	30726	66521	9.5	16.8	25	4.8	80.9
3	Arizona	67.4	20.7	25680	50448	15.3	9.85	20.4	7.5	144.
4	Arkansas	67.7	15.2	21274	39267	18	9.61	23.5	6.7	91.1
5	Californ…	. 57.4	30.7	29188	60883	13.7	8.89	15.2	6.9	176.
6	Colorado	67.6	25.6	30151	56456	12.2	9.15	19.9	3.7	84.6
7	Connecti	. 69.2	34.6	36775	67740	9.2	14.8	16.5	2.9	113
8	Delaware	73.6	17.7	29007	57599	11	8.89	20.7	4.4	155.
9	Florida	69.7	30	26551	47661	13.8	9.62	21.6	5	169.
10	Georgia	67.2	20.5	25134	49347	15.7	8.88	22.2	6.2	155.
#	# with 40 more rows, 10 more variables: agg_assault <dbl>, larceny <dbl>,</dbl></dbl>									
#	<pre># motor_theft <dbl>, soc_sec <dbl>, nuclear <dbl>, coal <dbl>,</dbl></dbl></dbl></dbl></pre>									
#	tr_death	s <dbl>,</dbl>	tr_death	hs_no_a	lc <dbl></dbl>	, unempl	<dbl>,  </dbl>	popdens	s2010 <d< td=""><td>ibl&gt;,</td></d<>	ibl>,

# and abbreviated variable names 'homeownership, 'multiunit, 'med\_income,

# 4fed\_spend









# **Cluster Validation**

# **Cluster Validity**

For supervised classification (= we have class labels) we have a variety of measures to evaluate how good our model is: Accuracy, Precision, Recall,  $F_1$ , etc.

# **Cluster Validity**

For supervised classification (= we have class labels) we have a variety of measures to evaluate how good our model is: Accuracy, Precision, Recall,  $F_1$ , etc.

For cluster analysis (=unsupervised learning), the analogous question is:

#### How to evaluate the "goodness" of the resulting clusters?

Keep in mind that the dataset can be large and high-dimensional. Visualization might be difficult. Random points:



Random points:







#### Hierarchical



1. Determining the **clustering tendency** of a set of data, i.e., distinguishing whether non-random structure exists in the data (e.g., to avoid overfitting).

- 1. Determining the **clustering tendency** of a set of data, i.e., distinguishing whether non-random structure exists in the data (e.g., to avoid overfitting).
- 2. Internal Validation: Evaluating how well the cluster analysis results fit the data without reference to external information.

- 1. Determining the **clustering tendency** of a set of data, i.e., distinguishing whether non-random structure exists in the data (e.g., to avoid overfitting).
- 2. Internal Validation: Evaluating how well the cluster analysis results fit the data without reference to external information.
- 3. External Validation: Compare the cluster analysis results to externally known class labels (class labels).

- 1. Determining the **clustering tendency** of a set of data, i.e., distinguishing whether non-random structure exists in the data (e.g., to avoid overfitting).
- 2. Internal Validation: Evaluating how well the cluster analysis results fit the data without reference to external information.
- 3. External Validation: Compare the cluster analysis results to externally known class labels (class labels).
- 4. Compare clusterings to determine which is better.

- 1. Determining the **clustering tendency** of a set of data, i.e., distinguishing whether non-random structure exists in the data (e.g., to avoid overfitting).
- 2. Internal Validation: Evaluating how well the cluster analysis results fit the data without reference to external information.
- 3. External Validation: Compare the cluster analysis results to externally known class labels (class labels).
- 4. Compare clusterings to determine which is better.
- 5. Determining the 'correct' number of clusters.

For 2, 3, and 4, we can further distinguish whether we want to evaluate the entire clustering or just individual clusters.

Numerical measures applied to judge various aspects of cluster validity are classified into the following three types.

- Internal Index: Used to measure the goodness of a clustering structure without respect to external information.
- External Index: Used to measure the extent to which cluster labels match externally supplied class labels.
- Relative Index: Used to compare two different clusterings or clusters.

#### Internal Index

Consider a dataset

 $D = \{\vec{x}_1, \ldots, \vec{x}_p\}$ 

Assume that a clustering algorithm produced a partition  $\mathcal{U} = \{U_1, \dots, U_K\}$  of D into K clusters.

No other information has been provided.

We aim to measure the clustering's "niceness" (??)

#### Internal Index

Consider a dataset

 $D = \{\vec{x}_1, \ldots, \vec{x}_p\}$ 

Assume that a clustering algorithm produced a partition  $\mathcal{U} = \{U_1, \dots, U_K\}$  of D into K clusters.

No other information has been provided.

We aim to measure the clustering's "niceness" (??)

Assume that we have a distance measure d measuring how far apart the objects being clustered.

#### For concreteness:

- ▶ We stick with numerical features, which means that the dataset  $D = {\vec{x_1}, ..., \vec{x_p}}$  contains vectors  $\vec{x_i} \in \mathbb{R}^n$ .
- Assume the Euclidean distance *d*.

Note that the validity measures may be based on completely different similarity/dissimilarity measures and non-numerical data.

Consider a dataset  $D = \{\vec{x}_1, \dots, \vec{x}_p\}$  and its clustering  $\mathcal{U} = \{U_1, \dots, U_K\}.$ 

Let us utilize the concept of distance to cluster prototypes and consider the distance between prototypes.

Consider a dataset  $D = \{\vec{x_1}, \dots, \vec{x_p}\}$  and its clustering  $\mathcal{U} = \{U_1, \dots, U_K\}.$ 

Let us utilize the concept of distance to cluster prototypes and consider the distance between prototypes.

To measure the dissimilarity of examples, we use the notion of *proximity* defined on pairs of vectors  $\vec{x}, \vec{z}$ :

 $proximity(\vec{x}, \vec{z})$ 

Consider a dataset  $D = \{\vec{x}_1, \dots, \vec{x}_p\}$  and its clustering  $\mathcal{U} = \{U_1, \dots, U_K\}.$ 

Let us utilize the concept of distance to cluster prototypes and consider the distance between prototypes.

To measure the dissimilarity of examples, we use the notion of *proximity* defined on pairs of vectors  $\vec{x}, \vec{z}$ :

 $proximity(\vec{x}, \vec{z})$ 

The proximity might be, e.g.,

- the distance  $d(\vec{x}, \vec{z})$ ,
- the square of the distance, that is  $d(\vec{x}, \vec{z})^2$ ,
- > any other notion of dissimilarity based on the application.

Consider a dataset  $D = \{\vec{x}_1, \dots, \vec{x}_p\}$  and its clustering  $\mathcal{U} = \{U_1, \dots, U_K\}.$ 

Let us utilize the concept of distance to cluster prototypes and consider the distance between prototypes.

To measure the dissimilarity of examples, we use the notion of *proximity* defined on pairs of vectors  $\vec{x}, \vec{z}$ :

 $proximity(\vec{x}, \vec{z})$ 

The proximity might be, e.g.,

- the distance  $d(\vec{x}, \vec{z})$ ,
- the square of the distance, that is  $d(\vec{x}, \vec{z})^2$ ,

any other notion of dissimilarity based on the application.
We consider the notions of *cohesion* (proximity of examples within clusters) and *separation* (proximity of clusters).



Prototype-based cohesion = the similarity of examples within a given cluster to a prototype of the cluster (e.g., centroid). Given a cluster  $U_i \in U$  and its prototype  $\vec{m_i} \in \mathbb{R}^n$ ,

$$cohesion(U_j) = \sum_{\vec{x} \in U_j} proximity(\vec{x}, \vec{m}_j)$$

Note that the prototype **does not** have to be an element of  $U_j$ . Intuitively, cohesion is the proximity of cluster's examples and a point somewhere "between" all examples of the cluster.



**Prototype-based separation** = dissimilarity of prototypes of different clusters.

Given a cluster  $U_j \in U$ , its prototype  $\vec{m}_j \in \mathbb{R}^n$ , and a prototype of all examples  $\vec{m} \in \mathbb{R}^n$  (e.g. the centroid of all examples)

separation
$$(U_j) = proximity(\vec{m}_j, \vec{m})$$

Intuitively, separation is the proximity of the cluster's examples to the dataset's center.

Summarize the prototype-based cohesion and separation as follows:

$$egin{aligned} \mathsf{cohesion}(\mathcal{U}) &= \sum_{j=1}^{K} \mathit{cohesion}(U_j) \ &= \sum_{j=1}^{K} \sum_{ec{x} \in U_j} \mathit{proximity}(ec{x},ec{m}_j) \end{aligned}$$

Summarize the prototype-based cohesion and separation as follows:

$$egin{aligned} \mathsf{cohesion}(\mathcal{U}) &= \sum_{j=1}^{K} \mathsf{cohesion}(U_j) \ &= \sum_{j=1}^{K} \sum_{ec{x} \in U_j} \mathsf{proximity}(ec{x},ec{m}_j) \end{aligned}$$

$$ext{separation}(\mathcal{U}) = \sum_{j=1}^{K} |U_j| ext{separation}(U_j) \ = \sum_{j=1}^{K} |U_j| ext{proximity}(ec{m}_j, ec{m})$$

If proximity  $(\vec{x}, \vec{z})$  is defined as  $d(\vec{x}, \vec{z})^2$  then the cohesion is the inertia.

Summarize the prototype-based cohesion and separation as follows:

$$egin{aligned} \mathsf{cohesion}(\mathcal{U}) &= \sum_{j=1}^{K} \mathsf{cohesion}(U_j) \ &= \sum_{j=1}^{K} \sum_{ec{x} \in U_j} \mathsf{proximity}(ec{x},ec{m}_j) \end{aligned}$$

$$egin{aligned} \mathsf{separation}(\mathcal{U}) &= \sum_{j=1}^{K} |U_j| \mathsf{separation}(U_j) \ &= \sum_{j=1}^{K} |U_j| \mathsf{proximity}(ec{m}_j, ec{m}) \end{aligned}$$

If  $proximity(\vec{x}, \vec{z})$  is defined as  $d(\vec{x}, \vec{z})^2$  then the cohesion is the inertia.

There is an interesting relationship between the above measures and the squared distances to the prototype of the whole dataset  $\vec{m}$ .
Consider a dataset  $D = \{\vec{x}_1, \dots, \vec{x}_p\}$  and its clustering  $\mathcal{U} = \{U_1, \dots, U_K\}$  of D.

Consider *proximity*  $(\vec{x}, \vec{z}) = d(\vec{x}, \vec{z})^2$  and all prototypes to be centroids.

Consider a dataset  $D = \{\vec{x}_1, \dots, \vec{x}_p\}$  and its clustering  $\mathcal{U} = \{U_1, \dots, U_K\}$  of D.

Consider *proximity*  $(\vec{x}, \vec{z}) = d(\vec{x}, \vec{z})^2$  and all prototypes to be centroids.

Let  $\vec{m} \in \mathbb{R}^n$  be the centroid of all examples:

$$\vec{m} = \frac{1}{|D|} \sum_{i=1}^{p} \vec{x}_i$$

Consider a dataset  $D = \{\vec{x}_1, \dots, \vec{x}_p\}$  and its clustering  $\mathcal{U} = \{U_1, \dots, U_K\}$  of D.

Consider *proximity*  $(\vec{x}, \vec{z}) = d(\vec{x}, \vec{z})^2$  and all prototypes to be centroids.

Let  $\vec{m} \in \mathbb{R}^n$  be the centroid of all examples:

$$\vec{m} = \frac{1}{|D|} \sum_{i=1}^{p} \vec{x_i}$$

Define

$$\mathsf{TSS} = \sum_{i=1}^{p} proximity(\vec{x_i}, \vec{m}) = \sum_{i=1}^{p} d(\vec{x_i}, \vec{m})^2$$

Consider a dataset  $D = \{\vec{x}_1, \dots, \vec{x}_p\}$  and its clustering  $\mathcal{U} = \{U_1, \dots, U_K\}$  of D.

Consider *proximity*  $(\vec{x}, \vec{z}) = d(\vec{x}, \vec{z})^2$  and all prototypes to be centroids.

Let  $\vec{m} \in \mathbb{R}^n$  be the centroid of all examples:

$$\vec{m} = \frac{1}{|D|} \sum_{i=1}^{p} \vec{x_i}$$

Define

$$\mathsf{TSS} = \sum_{i=1}^{p} proximity(\vec{x_i}, \vec{m}) = \sum_{i=1}^{p} d(\vec{x_i}, \vec{m})^2$$

The following holds:

$$\mathsf{TSS} = \mathsf{cohesion}(\mathcal{U}) + \mathsf{separation}(\mathcal{U})$$

Note that TSS is determined by D.

*Silhouette score* can be used to measure both qualities of clustering from the point of view of individual examples and from the point of view of the overall clustering.



silhouette
$$(\vec{x}) = \frac{b-a}{\max\{a, b\}}$$

Consider a clustering  $\mathcal{U} = \{U_1, \ldots, U_k\}$  and  $\vec{x} \in U_j$ .

Consider a clustering  $\mathcal{U}=\{U_1,\ldots,U_k\}$  and  $ec{x}\in U_j.$  If  $|U_j|>1$  we define

$$a(ec{x}) = rac{1}{|U_j|-1}\sum_{ec{z}\in U_j\smallsetminus\{ec{x}\}}d(ec{x},ec{z})$$

Consider a clustering  $\mathcal{U}=\{U_1,\ldots,U_k\}$  and  $ec{x}\in U_j.$  If  $|U_j|>1$  we define

$$a(ec{x}) = rac{1}{|U_j|-1}\sum_{ec{z}\in U_j\smallsetminus\{ec{x}\}}d(ec{x},ec{z})$$

$$b(\vec{x}) = \min_{k \neq j} \frac{1}{|U_k|} \sum_{\vec{z} \in U_k} d(\vec{x}, \vec{z})$$

Consider a clustering  $\mathcal{U}=\{U_1,\ldots,U_k\}$  and  $ec{x}\in U_j.$  If  $|U_j|>1$  we define

$$a(\vec{x}) = \frac{1}{|U_j| - 1} \sum_{\vec{z} \in U_j \smallsetminus \{\vec{x}\}} d(\vec{x}, \vec{z})$$

$$b(\vec{x}) = \min_{k \neq j} \frac{1}{|U_k|} \sum_{\vec{z} \in U_k} d(\vec{x}, \vec{z})$$

If  $|U_j| > 1$  we define

$$silhouette(ec{x}) = rac{b(ec{x}) - a(ec{x})}{\max\{a(ec{x}), b(ec{x})\}}$$

Else, we define *silhouette*( $\vec{x}$ ) = 0.





## Silhouette for Clusters and Clusterings

We have defined the silhouette for a single  $\vec{x} \in D$ .

To obtain the silhouette score for a whole cluster  $U_j$  or for D we summarize using simple averaging:

$$silhouette(U_j) = \frac{1}{|U_j|} \sum_{\vec{x} \in U_j} silhouette(\vec{x})$$

$$\textit{silhouette}(D) = rac{1}{|D|} \sum_{ec{x} \in D} \textit{silhouette}(ec{x})$$



#### Silhouette analysis for KMeans clustering on sample data with n\_clusters = 2

The colored graphs on the left are silhouette scores of the individual elements of clusters.



#### Silhouette analysis for KMeans clustering on sample data with n\_clusters = 3

The colored graphs on the left are silhouette scores of the individual elements of clusters.



#### Silhouette analysis for KMeans clustering on sample data with n\_clusters = 4

The colored graphs on the left are silhouette scores of the individual elements of clusters.



#### Silhouette analysis for KMeans clustering on sample data with n\_clusters = 5

The colored graphs on the left are silhouette scores of the individual elements of clusters.



#### Silhouette analysis for KMeans clustering on sample data with n\_clusters = 6

The colored graphs on the left are silhouette scores of the individual elements of clusters.

## External Index

Consider a *supervised learning* dataset

$$D = \{ (\vec{x_1}, c_1), \dots, (\vec{x_p}, c_p) \}$$

Here  $c_i \in C$  is a class of  $\vec{x_i}$ .

Assume that a clustering algorithm produced a partition  $U = \{U_1, \ldots, U_K\}$  of D into K clusters.

We measure how the clustering conforms with the given classes.

# Purity

Consider the clustering to be a classification model.

Define a classifier  $h: D \to C$  such that given  $\vec{x_i} \in U_i \in U$ 

 $h(\vec{x}_i)$  = the most frequent class in  $U_i$ 

Now we can measure the Accuracy of h.

# Purity

Consider the clustering to be a classification model.

Define a classifier  $h: D \to C$  such that given  $\vec{x_i} \in U_i \in U$ 

 $h(\vec{x}_i)$  = the most frequent class in  $U_i$ 

Now we can measure the Accuracy of h.

Accuracy of *h* is called *purity*.

Intuitively, it is the proportion of non-majority class elements in clusters.

# Purity

Consider the clustering to be a classification model.

Define a classifier  $h: D \to C$  such that given  $\vec{x_i} \in U_i \in U$ 

 $h(\vec{x}_i)$  = the most frequent class in  $U_i$ 

Now we can measure the Accuracy of h.

Accuracy of *h* is called *purity*.

Intuitively, it is the proportion of non-majority class elements in clusters.

Is it a good measure?

Probably not; many clusters lead to high purity (each element in its own cluster means purity = 1).

Given  $\vec{x_i}$ , denote by  $\mathcal{U}(\vec{x_i})$  the cluster  $U_j \in \mathcal{U}$  containing  $\vec{x_i}$ .

Given  $\vec{x_i}$ , denote by  $\mathcal{U}(\vec{x_i})$  the cluster  $U_j \in \mathcal{U}$  containing  $\vec{x_i}$ .

Distinguish the following types of *pairs* of examples:

 $\blacktriangleright$  TP = number of examples of the same class and the same cluster

$$\mathsf{TP} = |\{(i,j) \mid \mathcal{U}(\vec{x}_i) = \mathcal{U}(\vec{x}_j) \land c_i = c_j)\}$$

Given  $\vec{x_i}$ , denote by  $\mathcal{U}(\vec{x_i})$  the cluster  $U_j \in \mathcal{U}$  containing  $\vec{x_i}$ .

Distinguish the following types of *pairs* of examples:

► TP = number of examples of the same class and the same cluster TP =  $|\{(i,j) | U(\vec{x}_i) = U(\vec{x}_i) \land c_i = c_i)\}$ 

► TN = number of examples of different classes and different clusters  $TN = |\{(i,j) \mid U(\vec{x}_i) \neq U(\vec{x}_j) \land c_i \neq c_j)\}$ 

Given  $\vec{x_i}$ , denote by  $\mathcal{U}(\vec{x_i})$  the cluster  $U_j \in \mathcal{U}$  containing  $\vec{x_i}$ .

Distinguish the following types of *pairs* of examples:

► TP = number of examples of the same class and the same cluster  
TP = 
$$|\{(i,j) | U(\vec{x}_i) = U(\vec{x}_j) \land c_i = c_j)\}$$

► TN = number of examples of different classes and different clusters TN =  $|\{(i,j) | U(\vec{x}_i) \neq U(\vec{x}_j) \land c_i \neq c_j)\}$ 

FP = number of examples of the same class and different clusters FP = |{(i,j) | U(x<sub>i</sub>) = U(x<sub>j</sub>) ∧ c<sub>i</sub> ≠ c<sub>j</sub>)}

Given  $\vec{x_i}$ , denote by  $\mathcal{U}(\vec{x_i})$  the cluster  $U_j \in \mathcal{U}$  containing  $\vec{x_i}$ .

Distinguish the following types of *pairs* of examples:

► TP = number of examples of the same class and the same cluster  
TP = 
$$|\{(i,j) | U(\vec{x}_i) = U(\vec{x}_j) \land c_i = c_j)\}$$

► TN = number of examples of different classes and different clusters TN =  $|\{(i,j) | U(\vec{x}_i) \neq U(\vec{x}_j) \land c_i \neq c_j)\}$ 

- ► FP = number of examples of the same class and different clusters  $FP = |\{(i,j) \mid U(\vec{x_i}) = U(\vec{x_j}) \land c_i \neq c_j)\}$
- ► FN = number of examples of different classes and the same cluster  $FN = |\{(i,j) \mid U(\vec{x}_i) \neq U(\vec{x}_j) \land c_i = c_j)\}$

Now, we may apply all the measures from the supervised model.



$$TP = \binom{4}{2} + \binom{5}{2} + \binom{2}{2}$$
  
= 6 + 10 + 1 = 17  
$$TN = 4 * 5 + 1 * 2 = 22$$
  
$$FP = 1 * 4 + 5 * 2 = 14$$
  
$$FN = 1 * 5 + 4 * 2 = 13$$

		Cluster	
		same	diff
Class	same	TP=17	FN=13
	diff	FP=14	TN=22

# Rand Index

Accuracy (in this area known as *Rand index*) is

$$\mathsf{RandInd} = \mathsf{Accuracy} = \frac{\mathsf{TP} + \mathsf{TN}}{\mathsf{TP} + \mathsf{TN} + \mathsf{FP} + \mathsf{FN}}$$

# Rand Index

Accuracy (in this area known as *Rand index*) is

$$\mathsf{RandInd} = \mathsf{Accuracy} = \frac{\mathsf{TP} + \mathsf{TN}}{\mathsf{TP} + \mathsf{TN} + \mathsf{FP} + \mathsf{FN}}$$

In our example,

$$\mathsf{RandInd} = (17 + 22)/(17 + 22 + 14 + 13) = 0.59$$

Here, note that the Rand index considers the purity and the number of clusters.

# Rand Index

Accuracy (in this area known as *Rand index*) is

$$\mathsf{RandInd} = \mathsf{Accuracy} = \frac{\mathsf{TP} + \mathsf{TN}}{\mathsf{TP} + \mathsf{TN} + \mathsf{FP} + \mathsf{FN}}$$

In our example,

$$\mathsf{RandInd} = (17 + 22)/(17 + 22 + 14 + 13) = 0.59$$

Here, note that the Rand index considers the purity and the number of clusters.

Note that the Rand index can be used to compare two clusterings: Simply consider class labels to be indicators of clusters.

Similarly, we may compute the other measures such as Precision, Recall, and  $F_1$  with all their benefits and limitations.