Unconstrained Optimization Overview

Notation

In what follows, we will work with vectors in \mathbb{R}^n .

The vectors will be (usually) denoted by $x \in \mathbb{R}^n$.

We often consider sequences of vectors, $x_0, x_1, \ldots, x_k, \ldots$

The index k will usually indicate that x_k is the k-the vector in a sequence.

When we talk (relatively rarely) about components of vectors, we use *i* as an index, i.e., x_i will be the *i*-th component of $x \in \mathbb{R}^n$.

We denote by ||x|| the Euclidean norm of x.

We denote by $||x||_{\infty}$ the \mathcal{L}^{∞} norm giving the maximum of absolute values of components of x.

We ocasionally use the matrix morn ||A||, consistent with the Euclidean norm, defined by

$$||A|| = \sup_{||x||=1} ||Ax|| = \sqrt{\lambda_1}$$

Here λ_1 is the largest eigenvalue of $A^{\top}A$.

How to Recognize (Local) Minimum

How do we verify that $x^* \in \mathbb{R}^n$ is a minimizer of f?



How to Recognize (Local) Minimum

How do we verify that $x^* \in \mathbb{R}^n$ is a minimizer of f?



Technically, we should examine *all* points in the immediate vicinity if one has a smaller value (impractical).

Assuming the smoothness of f, we may benefit from the "stable" behavior of f around x^* .

Derivatives and Gradients

1

The gradient of $f : \mathbb{R}^n \to \mathbb{R}$, denoted by $\nabla f(x)$, is a column vector of first-order partial derivatives of the function concerning each variable:

$$abla f(x) = \left[rac{\partial f}{\partial x_1}, rac{\partial f}{\partial x_2}, \dots, rac{\partial f}{\partial x_n}
ight]^+,$$

Where each partial derivative is defined as the following limit:

$$\frac{\partial f}{\partial \mathbf{x}_{i}} = \lim_{\varepsilon \to 0} \frac{f(x_{1}, \dots, x_{i} + \varepsilon, \dots, x_{n}) - f(x_{1}, \dots, x_{i}, \dots, x_{n})}{\varepsilon}$$

Gradient



The gradient is a vector pointing in the direction of the most significant function increase from the current point.

Gradient

Consider the following function of two variables:

$$f(x_1, x_2) = x_1^3 + 2x_1x_2^2 - x_2^3 - 20x_1.$$

$$\nabla f(x_1, x_2) = \begin{bmatrix} 3x_1^2 + 2x_2^2 - 20 \\ 4x_1x_2 - 3x_2^2 \end{bmatrix}$$



Directional Derivatives vs Gradient

The rate of change in a direction p is quantified by a directional derivative, defined as

$$abla_{p}f(x) = \lim_{\varepsilon \to 0} \frac{f(x + \varepsilon p) - f(x)}{\varepsilon}$$

We can find this derivative by projecting the gradient onto the desired direction p using the dot product $\nabla_p f(x) = (\nabla f(x))^\top p$



(Here, we assume continuous partial derivatives.)

Geometry of Gradient

Consider the geometric interpretation of the dot product:

 $\nabla_{p}f(x) = (\nabla f(x))^{\top}p = ||\nabla f|| \, ||p|| \cos \theta$

Here θ is the angle between ∇f and p.

Geometry of Gradient

Consider the geometric interpretation of the dot product:

$$\nabla_{p}f(x) = (\nabla f(x))^{\top}p = ||\nabla f|| \, ||p|| \cos \theta$$

Here θ is the angle between ∇f and p.

The directional derivative is maximized by $\theta = 0$, i.e. when ∇f and p point in the same direction.



Hessian

Taking derivative twice, possibly w.r.t. different variables, gives the Hessian of \boldsymbol{f}

$$\nabla^2 f(x) = H(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

Note that the Hessian is a function which takes $x \in \mathbb{R}^n$ and gives a $n \times n$ -matrix of second derivatives of f.

٠

Hessian

Taking derivative twice, possibly w.r.t. different variables, gives the Hessian of \boldsymbol{f}

$$\nabla^2 f(x) = H(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

Note that the Hessian is a function which takes $x \in \mathbb{R}^n$ and gives a $n \times n$ -matrix of second derivatives of f.

We have

$$H_{ij}=\frac{\partial^2 f}{\partial x_i\partial x_j}.$$

If f has continuous second partial derivatives, then H is symmetric, i.e., $H_{ij} = H_{ji}$.

Let x be fixed and let g(t) = f(x + tp) and let $h_i(t) = \frac{\partial f}{\partial x_i}(x + tp)$ for $t \in \mathbb{R}$.

Let x be fixed and let g(t) = f(x + tp) and let $h_i(t) = \frac{\partial f}{\partial x_i}(x + tp)$ for $t \in \mathbb{R}$.

$$g'(t) = f(x+tp)' = [\nabla f(x+tp)]^\top p = \sum_{i=1}^n h_i(t)p_i$$

Let x be fixed and let g(t) = f(x + tp) and let $h_i(t) = \frac{\partial f}{\partial x_i}(x + tp)$ for $t \in \mathbb{R}$.

$$g'(t) = f(x+tp)' = [\nabla f(x+tp)]^{\top}p = \sum_{i=1}^{n} h_i(t)p_i$$

$$h'_{i}(t) = \left[\nabla \frac{\partial f}{\partial x_{i}}(x+tp)\right]^{\top} p = \sum_{j=1}^{n} \left(\frac{\partial f}{\partial x_{i}\partial x_{j}}(x+tp)\right) p_{j}$$
$$= [H(x+tp)p]_{i}$$

Let x be fixed and let g(t) = f(x + tp) and let $h_i(t) = \frac{\partial f}{\partial x_i}(x + tp)$ for $t \in \mathbb{R}$.

$$g'(t) = f(x+tp)' = [\nabla f(x+tp)]^{\top}p = \sum_{i=1}^{n} h_i(t)p_i$$

$$h'_{i}(t) = \left[\nabla \frac{\partial f}{\partial x_{i}}(x+tp)\right]^{\top} p = \sum_{j=1}^{n} \left(\frac{\partial f}{\partial x_{i}\partial x_{j}}(x+tp)\right) p_{j}$$
$$= [H(x+tp)p]_{i}$$

$$g''(t) = \sum_{i=1}^{n} h'_i(t) p_i = \sum_{i=1}^{n} [H(x+tp)p]_i p_i = p^{\top} H(x+tp)p$$

Let x be fixed and let g(t) = f(x + tp) and let $h_i(t) = \frac{\partial f}{\partial x_i}(x + tp)$ for $t \in \mathbb{R}$.

What exactly are g'(0) and g''(0)?

$$g'(t) = f(x+tp)' = [\nabla f(x+tp)]^{\top}p = \sum_{i=1}^{n} h_i(t)p_i$$

$$h'_{i}(t) = \left[\nabla \frac{\partial f}{\partial x_{i}}(x+tp)\right]^{\top} p = \sum_{j=1}^{n} \left(\frac{\partial f}{\partial x_{i}\partial x_{j}}(x+tp)\right) p_{j}$$
$$= [H(x+tp)p]_{i}$$

$$g''(t) = \sum_{i=1}^{n} h'_i(t) p_i = \sum_{i=1}^{n} [H(x+tp)p]_i p_i = p^{\top} H(x+tp)p$$

Thus,

$$g''(0) = p^\top H(x)p.$$

Principal Curvature Directions

Fix x and consider H = H(x). Consider unit eigenvectors \hat{v}_k of H:

 $H\hat{v}_k = \kappa_k\hat{v}_k$

For symmetric H, the unit eigenvectors form an orthonormal basis,

Principal Curvature Directions

Fix x and consider H = H(x). Consider unit eigenvectors \hat{v}_k of H:

$$H\hat{v}_k = \kappa_k \hat{v}_k$$

For symmetric H, the unit eigenvectors form an orthonormal basis, and there is a rotation matrix R such that

$$H = RDR^{-1} = RDR^{\top}$$

Here *D* is diagonal with $\kappa_1, \ldots, \kappa_n$ on the diagonal.

If $\kappa_1 \geq \cdots \geq \kappa_n$, the direction of \hat{v}_1 is the maximum curvature direction of f at x.



Consider $f(x) = x^{\top} H x$ where

$$H = \begin{pmatrix} 4/3 & 0 \\ 0 & 1 \end{pmatrix}$$

The eigenvalues are

 $\kappa_1 = 4/3$ $\kappa_2 = 1$

Their corresponding eigenvectors are $(1,0)^{\top}$ and $(0,1)^{\top}$.



Consider $f(x) = x^{\top} H x$ where

$$H = \begin{pmatrix} 4/3 & 0 \\ 0 & 1 \end{pmatrix}$$

The eigenvalues are

 $\kappa_1 = 4/3$ $\kappa_2 = 1$

Their corresponding eigenvectors are $(1,0)^{\top}$ and $(0,1)^{\top}$.

Note that

$$f(x) = \kappa_1 x_1^2 + \kappa_2 x_2^2$$

Considering a direction vector p we get

$$g(t) = f(0 + tp) = t^2 (\kappa_1 p_1^2 + \kappa_2 p_2^2)$$

which is a parabola with $g'' = 2 \left(\kappa_1 p_1^2 + \kappa_2 p_2^2 \right)$.



Consider $f(x) = x^{\top} Hx$ where

$$H = \begin{pmatrix} 4/3 & 1/3 \\ 1/3 & 3/3 \end{pmatrix}$$

Consider $f(x) = x^{\top} H x$ where

$$H = \begin{pmatrix} 4/3 & 1/3 \\ 1/3 & 3/3 \end{pmatrix}$$

The eigenvalues are

$$\kappa_1 = \frac{1}{6}(7 + \sqrt{5})$$
 $\kappa_2 = \frac{1}{6}(7 - \sqrt{5})$

Their corresponding eigenvectors are

$$\hat{\mathbf{v}}_{1} = \left(rac{1}{2}(1+\sqrt{5}),1
ight) \quad \hat{\mathbf{v}}_{2} = \left(rac{1}{2}(1-\sqrt{5}),1
ight)$$



Consider $f(x) = x^{\top} Hx$ where

$$H = \begin{pmatrix} 4/3 & 1/3 \\ 1/3 & 3/3 \end{pmatrix}$$

The eigenvalues are

$$\kappa_1 = \frac{1}{6}(7 + \sqrt{5})$$
 $\kappa_2 = \frac{1}{6}(7 - \sqrt{5})$

Their corresponding eigenvectors are

$$\hat{\mathbf{v}_1} = \left(rac{1}{2}(1+\sqrt{5}),1
ight) \quad \hat{\mathbf{v}_2} = \left(rac{1}{2}(1-\sqrt{5}),1
ight)$$

Note that

$$egin{aligned} \mathcal{H} = egin{pmatrix} \hat{v}_1 & \hat{v}_2 \end{pmatrix} egin{pmatrix} \kappa_1 & 0 \ 0 & \kappa_2 \end{pmatrix} egin{pmatrix} \hat{v}_1 & \hat{v}_2 \end{pmatrix}^ op \end{aligned}$$

Here $(\hat{v}_1 \ \hat{v}_2)$ is a 2 × 2 matrix whose columns are \hat{v}_1, \hat{v}_2 .



Hessian Visualization Example

Consider

$$f(x_1, x_2) = x_1^3 + 2x_1x_2^2 - x_2^3 - 20x_1.$$

And it's Hessian.

$$H(x_1, x_2) = \begin{bmatrix} 6x_1 & 4x_2 \\ 4x_2 & 4x_1 - 6x_2 \end{bmatrix}.$$



Theorem 1 (Taylor)

Suppose that $f : \mathbb{R}^n \to \mathbb{R}$ is twice continuously differentiable and that $p \in \mathbb{R}^n$. Then, we have

$$f(x+p) = f(x) + \nabla f(x)^{T} p + \frac{1}{2} p^{T} H(x) p + o(||p||^{2}).$$

Here $H = \nabla^2 f$ is the Hessian of f.

First-Order Necessary Conditions

Theorem 2

If x^* is a local minimizer and f is continuously differentiable in an open neighborhood of x^* , then $\nabla f(x^*) = 0$.





Note that $\nabla f(x^*) = 0$ does not tell us whether x^* is a minimizer, maximizer, or a saddle point.

Note that $\nabla f(x^*) = 0$ does not tell us whether x^* is a minimizer, maximizer, or a saddle point.

However, knowing the curvature in all directions from x^* might tell us what x^* is, right?

Note that $\nabla f(x^*) = 0$ does not tell us whether x^* is a minimizer, maximizer, or a saddle point.

However, knowing the curvature in all directions from x^* might tell us what x^* is, right?

Note that $\nabla f(x^*) = 0$ does not tell us whether x^* is a minimizer, maximizer, or a saddle point.

However, knowing the curvature in all directions from x^* might tell us what x^* is, right?

All comes down to the *definiteness* of $H := H(x^*)$.

- H is positive definite if p^THp > 0 for all p iff all eigenvalues of H are positive
- H is positive semi-definite if p[⊤]Hp ≥ 0 for all p iff all eigenvalues of H are nonnegative
- H is negative semi-definite if p[⊤]Hp ≤ 0 for all p iff all eigenvalues of H are nonpositive
- ► H is negative definite if p^T Hp < 0 for all p iff all eigenvalues of H are negative
- H is indefinite if it is not definite in the above sense iff H has at least one positive and one negative eigenvalue.

Definiteness



Second-Order Necessary Condition

Theorem 3 (Second-Order Necessary Conditions)

If x^* is a local minimizer of f and $\nabla^2 f$ is continuous in a neighborhood of x^* , then $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive semidefinite.

Theorem 4 (Second-Order Sufficient Conditions)

Suppose that $\nabla^2 f$ is continuous in a neighborhood of x^* and that $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive definite. Then x^* is a strict local minimizer of f.



Example

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Example

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Consider the gradient equal to zero:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1^3 + 6x_1^2 + 3x_1 - 2x_2 \\ 2x_2 - 2x_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Example

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Consider the gradient equal to zero:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1^3 + 6x_1^2 + 3x_1 - 2x_2 \\ 2x_2 - 2x_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

From the second equation, we have that $x_2 = x_1$. Substituting this into the first equation yields

$$x_1\left(2x_1^2+6x_1+1\right)=0.$$
Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Consider the gradient equal to zero:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1^3 + 6x_1^2 + 3x_1 - 2x_2 \\ 2x_2 - 2x_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

From the second equation, we have that $x_2 = x_1$. Substituting this into the first equation yields

$$x_1\left(2x_1^2+6x_1+1\right)=0.$$

The solution of this equation yields three points:

$$x_{\mathcal{A}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad x_{\mathcal{B}} = \begin{bmatrix} -\frac{3}{2} - \frac{\sqrt{7}}{2} \\ -\frac{3}{2} - \frac{\sqrt{7}}{2} \end{bmatrix}, \quad x_{\mathcal{C}} = \begin{bmatrix} \frac{\sqrt{7}}{2} - \frac{3}{2} \\ \frac{\sqrt{7}}{2} - \frac{3}{2} \end{bmatrix}$$

٠

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

To classify x_A, x_B, x_C , we need to compute the Hessian matrix:

$$H(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 6x_1^2 + 12x_1 + 3 & -2 \\ -2 & 2 \end{bmatrix}$$

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

To classify x_A, x_B, x_C , we need to compute the Hessian matrix:

$$H(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 6x_1^2 + 12x_1 + 3 & -2 \\ -2 & 2 \end{bmatrix}.$$

The Hessian, at the first point, is

$$H(x_{\mathcal{A}}) = \begin{bmatrix} 3 & -2 \\ -2 & 2 \end{bmatrix},$$

whose eigenvalues are $\kappa_1 \approx 0.438$ and $\kappa_2 \approx 4.561$. Because both eigenvalues are positive, this point is a local minimum.

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

To classify x_A, x_B, x_C , we need to compute the Hessian matrix:

$$H(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 6x_1^2 + 12x_1 + 3 & -2 \\ -2 & 2 \end{bmatrix}.$$

For the second point,

$$H(x_B) = \left[egin{array}{cc} 3(3+\sqrt{7}) & -2 \ -2 & 2 \end{array}
ight].$$

The eigenvalues are $\kappa_1 \approx 1.737$ and $\kappa_2 \approx 17.200$, so this point is another local minimum.

Consider the following function of two variables:

$$f(x_1, x_2) = 0.5x_1^4 + 2x_1^3 + 1.5x_1^2 + x_2^2 - 2x_1x_2.$$

To classify x_A, x_B, x_C , we need to compute the Hessian matrix:

$$H(x_1, x_2) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 6x_1^2 + 12x_1 + 3 & -2 \\ -2 & 2 \end{bmatrix}.$$

For the third point,

$$H(x_C) = \left[\begin{array}{rrr} 9 - 3\sqrt{7} & -2\\ -2 & 2 \end{array}\right]$$

The eigenvalues for this Hessian are $\kappa_1 \approx -0.523$ and $\kappa_2 \approx 3.586$, so this point is a saddle point.



Proofs of Some Theorems Optional

Taylor's Theorem

To prove the theorems characterizing minima/maxima, we need the following form of Taylor's theorem:

Theorem 5 (Taylor)

Suppose that $f : \mathbb{R}^n \to \mathbb{R}$ is continuously differentiable and that $p \in \mathbb{R}^n$. Then we have that.

$$f(x+p) = f(x) + \nabla f(x+tp)^T p,$$

for some $t \in (0, 1)$. Moreover, if f is twice continuously differentiable, we have that

$$f(x+p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x+tp)p,$$

for some $t \in (0, 1)$.

Proof of Theorem 2 (Optional)

We prove that if x^* is a local minimizer and f is continuously differentiable in an open neighborhood of x^* , then $\nabla f(x^*) = 0$.

Suppose for contradiction that $\nabla f(x^*) \neq 0$. Define the vector $p = -\nabla f(x^*)$ and note that $p^T \nabla f(x^*) = - \|\nabla f(x^*)\|^2 < 0$. Because ∇f is continuous near x^* , there is a scalar T > 0 such that

$$p^T
abla f(x^* + tp) < 0,$$
 for all $t \in [0, T]$

For any $ar{t} \in (0, T]$, we have by Taylor's theorem that

$$f\left(x^{*}+ar{t}p
ight)=f\left(x^{*}
ight)+ar{t}p^{T}
abla f\left(x^{*}+tp
ight), \hspace{0.5cm} ext{ for some }t\in(0,ar{t}).$$

Therefore, $f(x^* + \bar{t}p) < f(x^*)$ for all $\bar{t} \in (0, T]$. We have found a direction leading away from x^* along which f decreases, so x^* is not a local minimizer, and we have a contradiction.

Proof of Theorem 3 (Optional)

We prove that if x^* is a local minimizer of f and $\nabla^2 f$ is continuous in an open neighborhood of x^* , then $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive semidefinite.

We know that $\nabla f(x^*) = 0$. For contradiction, assume that $\nabla^2 f(x^*)$ is not positive semidefinite.

Then we can choose a vector p such that $p^T \nabla^2 f(x^*) p < 0$.

As $\nabla^2 f$ is continuous near x^* , $p^T \nabla^2 f(x^* + tp) p < 0$ for all $t \in [0, T]$ where T > 0.

By Taylor we have for all $\overline{t} \in (0, T]$ and some $t \in (0, \overline{t})$

$$f(x^* + \bar{t}p) = f(x^*) + \bar{t}p^T \nabla f(x^*) + \frac{1}{2} \bar{t}^2 p^T \nabla^2 f(x^* + tp) p < f(x^*).$$

Thus, x^* is not a local minimizer.

Proof of Theorem 4 (Optional)

We prove the following: Suppose that $\nabla^2 f$ is continuous in an open neighborhood of x^* and that $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive definite. Then x^* is a strict local minimizer of f.

Because the Hessian is continuous and positive definite at x^* , we can choose a radius r > 0 so that $\nabla^2 f(x)$ remains positive definite for all x in the open ball $\mathcal{D} = \{z \mid ||z - x^*|| < r\}$. Taking any nonzero vector p with ||p|| < r, we have $x^* + p \in \mathcal{D}$ and so

$$f(x^* + p) = f(x^*) + p^T \nabla f(x^*) + \frac{1}{2} p^T \nabla^2 f(z) p$$

= $f(x^*) + \frac{1}{2} p^T \nabla^2 f(z) p$,

where $z = x^* + tp$ for some $t \in (0, 1)$. Since $z \in D$, we have $p^T \nabla^2 f(z) p > 0$, and therefore $f(x^* + p) > f(x^*)$, giving the result.

Unconstrained Optimization Algorithms

Search Algorithms

We consider algorithms that

- Start with an initial guess x₀
- ▶ Generate a sequence of points *x*₀, *x*₁,...
- Stop when no progress can be made or when a minimizer seems approximated with sufficient accuracy.

To compute x_{k+1} the algorithms use the information about f at the previous iterates x_0, x_1, \ldots, x_k .

Search Algorithms

We consider algorithms that

- Start with an initial guess x₀
- Generate a sequence of points x_0, x_1, \ldots
- Stop when no progress can be made or when a minimizer seems approximated with sufficient accuracy.

To compute x_{k+1} the algorithms use the information about f at the previous iterates x_0, x_1, \ldots, x_k .

The *monotone* algorithms satisfy $f(x_{k+1}) < f(x_k)$.

Search Algorithms

We consider algorithms that

- Start with an initial guess x₀
- Generate a sequence of points x_0, x_1, \ldots
- Stop when no progress can be made or when a minimizer seems approximated with sufficient accuracy.

To compute x_{k+1} the algorithms use the information about f at the previous iterates x_0, x_1, \ldots, x_k .

The *monotone* algorithms satisfy $f(x_{k+1}) < f(x_k)$.

There are two overall strategies:

- Line search
- Trust region

Line Search Overview

To compute x_{k+1} , a line search algorithm chooses

- \blacktriangleright direction p_k
- step size α_k

and computes

$$x_{k+1} = x_k + \alpha_k p_k$$

Line Search Overview

To compute x_{k+1} , a line search algorithm chooses

direction p_k

• step size α_k

and computes

 $x_{k+1} = x_k + \alpha_k p_k$

The vector p_k should be a *descent* direction, i.e., a direction in which f decreases locally.

Line Search Overview

To compute x_{k+1} , a line search algorithm chooses

direction p_k

• step size α_k

and computes

 $x_{k+1} = x_k + \alpha_k p_k$

The vector p_k should be a *descent* direction, i.e., a direction in which f decreases locally.

 α_k is selected to approximately solve

$$\min_{\alpha>0}f(x_k+\alpha p_k)$$

However, typically, an exact solution is expensive and unnecessary. Instead, line search algorithms inspect a limited number of trial step lengths and find one that decreases f appropriately (see later).



A descent direction does not have to be followed to the minimum.



Trust Region

To compute x_{k+1} , a trust region algorithm chooses

• model function m_k whose behavior near x_k is similar to f

▶ a trust region $R \subseteq \mathbb{R}^n$ around x_k . Usually R is the ball defined by $||x - x_k|| \leq \Delta$ where $\Delta > 0$ is trust region radius.

and finds x_{k+1} solving

 $\min_{x\in R}m_k(x)$

Trust Region

To compute x_{k+1} , a trust region algorithm chooses

• model function m_k whose behavior near x_k is similar to f

▶ a trust region $R \subseteq \mathbb{R}^n$ around x_k . Usually R is the ball defined by $||x - x_k|| \leq \Delta$ where $\Delta > 0$ is trust region radius.

and finds x_{k+1} solving

 $\min_{x\in R}m_k(x)$

If the solution does not sufficiently decrease f, we shrink the trust region and re-solve.

Trust Region

To compute x_{k+1} , a trust region algorithm chooses

• model function m_k whose behavior near x_k is similar to f

▶ a *trust region* $R \subseteq \mathbb{R}^n$ around x_k . Usually R is the ball defined by $||x - x_k|| \leq \Delta$ where $\Delta > 0$ is *trust region radius*.

and finds x_{k+1} solving

 $\min_{x\in R}m_k(x)$

If the solution does not sufficiently decrease f, we shrink the trust region and re-solve.

The model m_k is usually derived from the Taylor's theorem.

$$m_k(x_k+p) = f_k + p^T \nabla f_k + \frac{1}{2} p^T B_k p$$

Where B_k approximates the Hessian of f at x_k .



Line Search Methods

Line Search

For setting the step size, we consider

- Armijo condition and backtracking algorithm
- strong Wolfe conditions and bracketing & zooming

Line Search

For setting the step size, we consider

- Armijo condition and backtracking algorithm
- strong Wolfe conditions and bracketing & zooming

For setting the direction, we consider

- Gradient descent
- Newton's method
- quasi-Newton methods (BFGS)
- (Conjugate gradients)

We start with the step size.

Step Size

Assume

 $x_{k+1} = x_k + \alpha_k p_k$

Where p_k is a descent direction

 $p_k^\top \nabla f_k < 0$



Step Size

Assume

 $x_{k+1} = x_k + \alpha_k p_k$

Where p_k is a descent direction $p_k^\top
abla f_k < 0$

Define

$$\phi(\alpha) = f(x_k + \alpha p_k)$$



Step Size

Assume

 $x_{k+1} = x_k + \alpha_k p_k$

Where p_k is a descent direction $p_k^\top \nabla f_k < 0$



Define

$$\phi(\alpha) = f(x_k + \alpha p_k)$$

We know that

 $\phi'(\alpha) = \nabla f(x_k + \alpha p_k)^\top p_k$ which means $\phi'(0) = \nabla f_k^\top p_k$

Note that $\phi'(0)$ must be negative as p_k is a descent direction.

Armijo Condition

The sufficient decrease condition (aka Armijo condition)

$$\phi(\alpha) \le \phi(\mathbf{0}) + \alpha \left(\mu_1 \phi'(\mathbf{0})\right)$$

where μ_1 is a constant such that $0 < \mu_1 \leq 1$



In practice, μ_1 is several orders smaller than 1, typically $\mu_1 = 10^{-4}$.

Backtracking Line Search Algorithm

Algorithm 1 Backtracking Line Search Input: $\alpha_{init} > 0, 0 < \mu_1 < 1, 0 < \rho < 1$ Output: α^* satisfying sufficient decrease condition 1: $\alpha \leftarrow \alpha_{init}$ 2: while $\phi(\alpha) > \phi(0) + \alpha \mu_1 \phi'(0)$ do 3: $\alpha \leftarrow \rho \alpha$

4: end while

The parameter ρ is typically set to 0.5. It can also be a variable set by a more sophisticated method (interpolation).

The α_{init} depends on the method for setting the descent direction p_k . For Newton and quasi-Newton, it is 1.0, but for other methods, it might be different.

There are two scenarios where the method does not perform well:

There are two scenarios where the method does not perform well:

The guess for the initial step is far too large, and the step sizes that satisfy sufficient decrease are smaller than the starting step by several orders of magnitude. Depending on the value of ρ, this scenario requires many backtracking evaluations.

There are two scenarios where the method does not perform well:

- The guess for the initial step is far too large, and the step sizes that satisfy sufficient decrease are smaller than the starting step by several orders of magnitude. Depending on the value of ρ, this scenario requires many backtracking evaluations.
- The guess for the initial step immediately satisfies sufficient decrease. However, the function's slope is still highly negative, and we could have decreased the function value by much more if we had taken a more significant step. In this case, our guess for the initial step is far too small.

There are two scenarios where the method does not perform well:

- The guess for the initial step is far too large, and the step sizes that satisfy sufficient decrease are smaller than the starting step by several orders of magnitude. Depending on the value of ρ, this scenario requires many backtracking evaluations.
- The guess for the initial step immediately satisfies sufficient decrease. However, the function's slope is still highly negative, and we could have decreased the function value by much more if we had taken a more significant step. In this case, our guess for the initial step is far too small.

Even if our original step size is not too far from an acceptable one, the basic backtracking algorithm ignores any information we have about the function values and gradients. It blindly takes a reduced step based on a preselected ratio ρ .
Backtracking Example

$$f(x_1, x_2) = 0.1x_1^6 - 1.5x_1^4 + 5x_1^2 + 0.1x_2^4 + 3x_2^2 - 9x_2 + 0.5x_1x_2$$

2.5 2 1.5 1 0.5 -3 -2 -1 0 1 2 3 x_1 x_1

 x_2

 $\mu_1 = 10^{-4}$ and $\rho = 0.7$.



We want to prevent too short of steps and to "motivate" the search to move closer to the minimum.

We want to prevent too short of steps and to "motivate" the search to move closer to the minimum.

We introduce the sufficient curvature condition

 $\left|\phi'(\alpha)\right| \le \mu_2 \left|\phi'(0)\right|$

where $\mu_1 < \mu_2 < 1$ is a constant.



We want to prevent too short of steps and to "motivate" the search to move closer to the minimum.

We introduce the sufficient curvature condition



Typical values of μ_2 range from 0.1 to 0.9, depending on the direction setting method.

We want to prevent too short of steps and to "motivate" the search to move closer to the minimum.

We introduce the sufficient curvature condition



Typical values of μ_2 range from 0.1 to 0.9, depending on the direction setting method.

As μ_2 tends to 0, the condition enforces $\phi'(\alpha) = 0$, which would yield an exact line search.

Strong Wolfe Conditions

Putting together Armijo and sufficient curvature conditions, we obtain *strong Wolfe conditions*

Strong Wolfe Conditions

Putting together Armijo and sufficient curvature conditions, we obtain *strong Wolfe conditions*

Sufficient decrease condition

 $\phi(\alpha) \leq \phi(0) + \mu_1 \alpha \phi'(0)$

Strong Wolfe Conditions

Putting together Armijo and sufficient curvature conditions, we obtain *strong Wolfe conditions*

Sufficient decrease condition

 $\phi(\alpha) \le \phi(0) + \mu_1 \alpha \phi'(0)$



Satisfiability of Strong Wolfe Conditions

Theorem 6

Suppose $f : \mathbb{R}^n \to \mathbb{R}$ is continuously differentiable. Let p_k be a descent direction at x_k , and assume that f is bounded below along the ray $\{x_k + \alpha p_k \mid \alpha > 0\}$. Then, if $0 < \mu_1 < \mu_2 < 1$, step length intervals exist that satisfy the strong Wolfe conditions.



Convergence of Line Search

Denote by θ_k the angle between p_k and $-\nabla f_k$, i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \, \|p_k\|}$$

Convergence of Line Search

Denote by θ_k the angle between p_k and $-\nabla f_k$, i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \, \|p_k\|}$$

Recall that f is L-smooth on a set \mathcal{N} for some L > 0 if

 $\|
abla f(x) -
abla f(ilde{x})\| \le L \|x - ilde{x}\|, \quad \text{ for all } x, ilde{x} \in \mathcal{N}$

Convergence of Line Search

Denote by θ_k the angle between p_k and $-\nabla f_k$, i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \, \|p_k\|}$$

Recall that f is L-smooth on a set \mathcal{N} for some L > 0 if

$$\|
abla f(x) -
abla f(ilde{x})\| \le L \|x - ilde{x}\|, \quad ext{ for all } x, ilde{x} \in \mathcal{N}$$

Theorem 7 (Zoutendijk)

Consider $x_{k+1} = x_k + \alpha_k p_k$, where p_k is a descent direction and α_k satisfies the strong Wolfe conditions. Suppose that f is bounded below in \mathbb{R}^n and that f is continuously differentiable in an open set \mathcal{N} containing the level set $\{x : f(x) \leq f(x_0)\}$. Assume also that f is L-smooth on \mathcal{N} . Then

$$\sum_{k\geq 0}\cos^2\theta_k \, \|\nabla f_k\|^2 < \infty.$$

Line Search Algorithm

How can we find a step size that satisfies strong Wolfe conditions?

How can we find a step size that satisfies *strong* Wolfe conditions? Use a bracketing and zoom algorithm, which proceeds in the following two phases:

Line Search Algorithm

How can we find a step size that satisfies strong Wolfe conditions?

Use a bracketing and zoom algorithm, which proceeds in the following two phases:

1. The bracketing phase finds an interval within which we are certain to find a point that satisfies the strong Wolfe conditions.

Line Search Algorithm

How can we find a step size that satisfies strong Wolfe conditions?

Use a bracketing and zoom algorithm, which proceeds in the following two phases:

- 1. The bracketing phase finds an interval within which we are certain to find a point that satisfies the strong Wolfe conditions.
- The zooming phase finds a point that satisfies the strong Wolfe conditions within the interval provided by the bracketing phase.

Algorithm 2 Bracketing

Input: $\alpha_1 > 0$ and α_{max} 1: Set $\alpha_0 \leftarrow 0$ 2: $i \leftarrow 1$ 3: repeat Evaluate $\phi(\alpha_i)$ 4: if $\phi(\alpha_i) > \phi(0) + \alpha_i \mu_1 \phi'(0)$ or $[\phi(\alpha_i) \ge \phi(\alpha_{i-1})$ and i > 15: then $\alpha^* \leftarrow \mathbf{zoom}(\alpha_{i-1}, \alpha_i)$ and stop 6: end if 7: Evaluate $\phi'(\alpha_i)$ 8: if $|\phi'(\alpha_i)| < \mu_2 |\phi'(0)|$ then 9: set $\alpha^* \leftarrow \alpha_i$ and stop 10: else if $\phi'(\alpha_i) > 0$ then 11: set $\alpha^* \leftarrow \mathbf{zoom}(\alpha_i, \alpha_{i-1})$ and stop 12: end if 13: Choose $\alpha_{i+1} \in (\alpha_i, \alpha_{\max})$ 14: $i \leftarrow i + 1$ 15: 16: **until** a condition is met

Explanation of Bracketing

Note that the sequence of trial steps α_i is monotonically increasing.

Note that the sequence of trial steps α_i is monotonically increasing.

Note that **zoom** is called when one of the following conditions is satisfied:

- α_i violates the sufficient decrease condition (lines 5 and 6)
- $\phi(\alpha_i) \ge \phi(\alpha_{i-1})$ (also lines 5 and 6)
- $\phi'(\alpha_i) \ge 0$ (lines 11 and 12)

The last step increases the α_i . May use, e.g., a constant multiple.

The following algorithm keeps two step lengths: $\alpha_{\rm lo}$ and $\alpha_{\rm hi}$

The following algorithm keeps two step lengths: α_{lo} and α_{hi}

The following invariants are being preserved:

The interval bounded by α_{lo} and α_{hi} always contains one or more intervals satisfying the strong Wolfe conditions. Note that we *do not* assume α_{lo} ≤ α_{hi}

The following algorithm keeps two step lengths: $\alpha_{\rm lo}$ and $\alpha_{\rm hi}$

The following invariants are being preserved:

- The interval bounded by α_{lo} and α_{hi} always contains one or more intervals satisfying the strong Wolfe conditions. Note that we *do not* assume α_{lo} < α_{hi}
- α_{lo} is, among all step lengths generated so far and satisfying the sufficient decrease condition, the one giving the smallest value of φ,

The following algorithm keeps two step lengths: $\alpha_{\rm lo}$ and $\alpha_{\rm hi}$

The following invariants are being preserved:

- The interval bounded by α_{lo} and α_{hi} always contains one or more intervals satisfying the strong Wolfe conditions. Note that we *do not* assume α_{lo} ≤ α_{hi}
- α_{lo} is, among all step lengths generated so far and satisfying the sufficient decrease condition, the one giving the smallest value of φ,

1: function $ZOOM(\alpha_{lo}, \alpha_{hi})$

2: repeat

- 3: Set α between α_{lo} and α_{hi} using interpolation (bisection, quadratic, etc.)
- 4: Evaluate $\phi(\alpha)$
- 5: if $\phi(\alpha) > \phi(0) + \alpha \mu_1 \phi'(0)$ or $\phi(\alpha) \ge \phi(\alpha_{\mathsf{lo}})$ then
- 6: $\alpha_{hi} \leftarrow \alpha$

else

7:

- 8: Evaluate $\phi'(\alpha)$
- 9: **if** $|\phi'(\alpha)| \le \mu_2 |\phi'(0)|$ then 10: Set $\alpha^* \leftarrow \alpha$ and stop
- 10: Set $\alpha^* \leftarrow$ 11: **end if**
- 12: **if** $\phi'(\alpha)(\alpha_{hi} \alpha_{lo}) \ge 0$ then

13: $\alpha_{hi} \leftarrow \alpha_{lo}$

- 14: **end if**
- 15: $\alpha_{\mathsf{lo}} \leftarrow \alpha$
- 16: end if
- 17: **until** a condition is met
- 18: end function

Bracketing & Zooming Example

We use quadratic interpolation; the bracketing chooses $\alpha_{i+1} = 2\alpha_i$, and the sufficient curvature factor is $\mu_2 = 0.9$.



Bracketing & Zooming Example

We use quadratic interpolation; the bracketing chooses $\alpha_{i+1} = 2\alpha_i$, and the sufficient curvature factor is $\mu_2 = 0.9$.



Bracketing is achieved in the first iteration by using a significant initial step of $\alpha_{init} = 1.2$ (left). Then, zooming finds an improved point through interpolation.

Bracketing & Zooming Example

We use quadratic interpolation; the bracketing chooses $\alpha_{i+1} = 2\alpha_i$, and the sufficient curvature factor is $\mu_2 = 0.9$.



Bracketing is achieved in the first iteration by using a significant initial step of $\alpha_{init} = 1.2$ (left). Then, zooming finds an improved point through interpolation.

The small initial step of $\alpha_{init} = 0.05$ (right) does not satisfy the strong Wolfe conditions, and the bracketing phase moves forward toward a flatter part of the function.

The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.

- The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.
- Practical line search algorithms also use the interpolating polynomials' properties to make educated guesses of where the next step length should lie.

- The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.
- Practical line search algorithms also use the interpolating polynomials' properties to make educated guesses of where the next step length should lie.
- A problem that can arise in the implementation is that as the optimization algorithm approaches the solution, two consecutive function values f (x_k) and f (x_{k-1}) may be indistinguishable in finite-precision arithmetic.

- The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.
- Practical line search algorithms also use the interpolating polynomials' properties to make educated guesses of where the next step length should lie.
- A problem that can arise in the implementation is that as the optimization algorithm approaches the solution, two consecutive function values f (x_k) and f (x_{k-1}) may be indistinguishable in finite-precision arithmetic.
- Some procedures also stop if the relative change in x is close to machine accuracy or some user-specified threshold.

- The interpolation of the zoom phase that determines α should be safeguarded to ensure that the new step length is not too close to the endpoints of the interval.
- Practical line search algorithms also use the interpolating polynomials' properties to make educated guesses of where the next step length should lie.
- A problem that can arise in the implementation is that as the optimization algorithm approaches the solution, two consecutive function values f (x_k) and f (x_{k-1}) may be indistinguishable in finite-precision arithmetic.
- Some procedures also stop if the relative change in x is close to machine accuracy or some user-specified threshold.
- The presented algorithm is implemented in https://docs.scipy.org/doc/scipy/reference/ generated/scipy.optimize.line_search.html

Unconstrained Optimization Algorithms

Descent Direction

First-Order Methods

Gradient Descent

Consider the *gradient descent* (aka *gradient descent*) method where

$$x_{k+1} = x_k + \alpha_k p_k$$
 $p_k = -\nabla f(x_k)$



Gradient Descent

Consider the *gradient descent* (aka *gradient descent*) method where

$$x_{k+1} = x_k + \alpha_k p_k$$
 $p_k = -\nabla f(x_k)$



Unfortunately, the gradient does not possess much information about the step size.

So usually, a normalized gradient is used to obtain the direction, and then a line search is performed:

$$x_{k+1} = x_k + \alpha_k p_k$$
 $p_k = -\frac{\nabla f(x_k)}{||\nabla f(x_k)||}$

The line search is *exact* if α_k minimizes $f(x_k + \alpha_k p_k)$. Not practical, we usually find α_k satisfying the strong Wolfe conditions.

Gradient Descent Algorithm with Line Search

Algorithm 3 Gradient Descent with Line Search

Input: x_0 starting point, $\varepsilon > 0$

Output: x^* approximation to a stationary point

1:
$$k \leftarrow 0$$

2: while
$$\|\nabla f\|_{\infty} > \varepsilon$$
 do

3:
$$p_k \leftarrow -\frac{\nabla f(x_k)}{\|\nabla f(x_k)\|}$$

4: Set
$$\alpha_{init}$$
 for line search

5:
$$\alpha_k \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$$

$$6: \qquad x_{k+1} \leftarrow x_k + \alpha_k p_k$$

7:
$$k \leftarrow k+1$$

8: end while
Gradient Descent Algorithm with Line Search

Algorithm 4 Gradient Descent with Line Search **Input:** x_0 starting point, $\varepsilon > 0$ **Output:** x^* approximation to a stationary point 1: $k \leftarrow 0$ 2: while $\|\nabla f\|_{\infty} > \varepsilon$ do $p_k \leftarrow -\frac{\nabla f(x_k)}{\|\nabla f(x_k)\|}$ 3: 4: Set α_{init} for line search 5: $\alpha_k \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$ 6: $x_{k+1} \leftarrow x_k + \alpha_k p_k$ $k \leftarrow k+1$ 7. 8: end while

Here α_{init} can be estimated from the previous step size α_{k-1} by demanding similar decrease in the objective:

$$\alpha_{init} \boldsymbol{p}_{k}^{\top} \nabla f_{k}^{\top} \approx \alpha_{k-1} \boldsymbol{p}_{k-1}^{\top} \nabla f_{k-1}^{\top} \quad \Rightarrow \quad \alpha_{init} = \alpha_{k-1} \frac{\alpha_{k-1} \boldsymbol{p}_{k-1}^{\top} \nabla f_{k-1}^{\top}}{\nabla \boldsymbol{p}_{k}^{\top} f_{k}^{\top}}$$

Gradient Descent Algorithm with Line Search

Algorithm 5 Gradient Descent with Line Search

Input: x_0 starting point, $\varepsilon > 0$ **Output:** x^* approximation to a stationary point

1:
$$k \leftarrow 0$$

2: while
$$\|\nabla f\|_{\infty} > \varepsilon$$
 do

3:
$$p_k \leftarrow -\frac{\nabla f(x_k)}{\|\nabla f(x_k)\|}$$

4: Set
$$\alpha_{init}$$
 for line search

5:
$$\alpha_k \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$$

$$6: \qquad x_{k+1} \leftarrow x_k + \alpha_k p_k$$

7:
$$k \leftarrow k+1$$

8: end while

Gradient Descent Algorithm with Line Search

Algorithm 6 Gradient Descent with Line Search

Input: x_0 starting point, $\varepsilon > 0$ **Output:** x^* approximation to a stationary point 1: $k \leftarrow 0$ 2: while $\|\nabla f\|_{\infty} > \varepsilon$ do $p_k \leftarrow -\frac{\nabla f(x_k)}{\|\nabla f(x_k)\|}$ 3: 4: Set α_{init} for line search 5: $\alpha_k \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$ 6: $x_{k+1} \leftarrow x_k + \alpha_k p_k$ $k \leftarrow k+1$ 7. 8: end while

Here α_{init} can be estimated from the previous step size α_{k-1} by demanding similar decrease in the objective:

$$\alpha_{init} \boldsymbol{p}_{k}^{\top} \nabla f_{k}^{\top} \approx \alpha_{k-1} \boldsymbol{p}_{k-1}^{\top} \nabla f_{k-1}^{\top} \quad \Rightarrow \quad \alpha_{init} = \alpha_{k-1} \frac{\alpha_{k-1} \boldsymbol{p}_{k-1}^{\top} \nabla f_{k-1}^{\top}}{\nabla \boldsymbol{p}_{k}^{\top} f_{k}^{\top}}$$



Note that p_{k+1} and p_k are always orthogonal.

$$f(x_1, x_2) = (1 - x_1)^2 + (1 - x_2)^2 + \frac{1}{2}(2x_2 - x_1^2)^2$$

Stopping: $||\nabla f||_{\infty} \leq 10^{-6}$.



The gradient descent can be prolonged.

Global Convergence with Line Search

Recall the Zoutendijk's theorem.

Denote by θ_k the angle between p_k and $-\nabla f_k$, i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \, \|p_k\|}$$

Recall that f is L-smooth on a set \mathcal{N} for some L > 0 if

$$\|
abla f(x) -
abla f(ilde{x})\| \leq L \|x - ilde{x}\|, \quad ext{ for all } x, ilde{x} \in \mathcal{N}$$

Theorem 8 (Zoutendijk)

Consider $x_{k+1} = x_k + \alpha_k p_k$, where p_k is a descent direction and α_k satisfies the strong Wolfe conditions. Suppose that f is bounded below in \mathbb{R}^n and that f is continuously differentiable in an open set \mathcal{N} containing the level set $\{x : f(x) \leq f(x_0)\}$. Assume also that f is L-smooth on \mathcal{N} . Then

$$\sum_{k\geq 0}\cos^2\theta_k \, \|\nabla f_k\|^2 < \infty.$$

Global Convergence of Gradient Descent

Assume that each α_k satisfies strong Wolfe conditions.

Global Convergence of Gradient Descent

Assume that each α_k satisfies strong Wolfe conditions.

Note that the angle θ_k between $p_k = -\nabla f_k$ and the negative gradient $-\nabla f_k$ equals 0. Hence, $\cos \theta_k = 1$.

Global Convergence of Gradient Descent

Assume that each α_k satisfies strong Wolfe conditions.

Note that the angle θ_k between $p_k = -\nabla f_k$ and the negative gradient $-\nabla f_k$ equals 0. Hence, $\cos \theta_k = 1$.

Thus, under the assumptions of Zoutendijk's theorem, we obtain

$$\sum_{k\geq 0}\cos^2\theta_k \|\nabla f_k\|^2 = \sum_{k\geq 0} \|\nabla f_k\|^2 < \infty$$

which implies that $\lim_{k\to\infty} ||\nabla f_k|| = 0$.

Local Linear Convergence of Gradient Descent

Theorem 9

Suppose that $f : \mathbb{R}^n \to \mathbb{R}$ is twice continuously differentiable, that the line search is exact, and that the descent converges to x^* where $\nabla f(x^*) = 0$ and the Hessian matrix $\nabla^2 f(x^*)$ is positive definite. Then

$$f(x_{k+1}) - f(x^*) \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}\right)^2 \left[f(x_k) - f(x^*)\right],$$

where $\lambda_1 \leq \cdots \leq \lambda_n$ are the eigenvalues of $\nabla^2 f(x^*)$.





Here $\ell_1 = 12, \ell_2 = 8, k_1 = 1, k_2 = 10, mg = 7$

Two Spring Problem - Gradient Descent



Gradient descent, line search, stop. cond. $||\nabla f||_{\infty} \leq 10^{-6}$.

Rosenbrock Function - Gradient Descent Rosenbrock: $f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$



Gradient descent, line search, stop. cond. $||\nabla f||_{\infty} \leq 10^{-6}$.

The method needs evaluation of ∇f at each x_k. If f is not differentiable at x_k, subgradients can be considered (out of the scope of this course).

- The method needs evaluation of ∇f at each x_k. If f is not differentiable at x_k, subgradients can be considered (out of the scope of this course).
- Slow, zig-zagging, provides insufficient information for line search initialization.

- The method needs evaluation of ∇f at each x_k. If f is not differentiable at x_k, subgradients can be considered (out of the scope of this course).
- Slow, zig-zagging, provides insufficient information for line search initialization.
- Susceptible to scaling of variables (see the paraboloid example).

- The method needs evaluation of ∇f at each x_k. If f is not differentiable at x_k, subgradients can be considered (out of the scope of this course).
- Slow, zig-zagging, provides insufficient information for line search initialization.
- Susceptible to scaling of variables (see the paraboloid example).
- THE basis for algorithms training neural networks a huge amount of specific adjustments are developed for working with huge numbers of variables in neural networks (trillions of weights).

Unconstrained Optimization Algorithms

Descent Direction

Second-Order Methods

Consider an objective $f : \mathbb{R}^n \to \mathbb{R}$.

Assume that f is twice differentiable.

Consider an objective $f : \mathbb{R}^n \to \mathbb{R}$.

Assume that f is twice differentiable.

Then, by the Taylor's theorem,

$$f(x_k+s) \approx f_k + \nabla f_k^{\top} s + \frac{1}{2} s^{\top} H_k s$$

where we denote the Hessian $\nabla^2 f(x_k)$ by H_k .

Consider an objective $f : \mathbb{R}^n \to \mathbb{R}$.

Assume that f is twice differentiable.

Then, by the Taylor's theorem,

$$f(\mathbf{x}_k + \mathbf{s}) \approx f_k + \nabla f_k^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top H_k \mathbf{s}$$

where we denote the Hessian $\nabla^2 f(x_k)$ by H_k . Define

$$q(s) = f_k +
abla f_k^\top s + rac{1}{2} s^\top H_k s$$

and minimize q w.r.t. s by setting $\nabla q(s) = 0$.

Consider an objective $f : \mathbb{R}^n \to \mathbb{R}$.

Assume that f is twice differentiable.

Then, by the Taylor's theorem,

$$f(\mathbf{x}_k + \mathbf{s}) \approx f_k + \nabla f_k^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top H_k \mathbf{s}$$

where we denote the Hessian $\nabla^2 f(x_k)$ by H_k .

Define

$$q(s) = f_k +
abla f_k^ op s + rac{1}{2}s^ op H_k s$$

and minimize q w.r.t. s by setting $\nabla q(s) = 0$. We obtain:

$$H_k s = -\nabla f_k$$

Denote by s_k the solution, and set $x_{k+1} = x_k + s_k$.

Algorithm 7 Newton's Method

Input: x_0 starting point, $\varepsilon > 0$

Output: x^* approximation to a stationary point

1:
$$k \leftarrow 0$$

2: while
$$\|\nabla f_k\|_{\infty} > \varepsilon$$
 do

3:
$$p_k \leftarrow -H_k^{-1} \nabla f(x_k)$$

4:
$$x_{k+1} \leftarrow x_k + p_k$$

5:
$$k \leftarrow k+1$$

6: end while

Newton's Method - Example

Newton's method finds the minimum of a quadratic function in a single step.



Note that the Newton's method is scale-invariant!

$$f(x_1, x_2) = (1 - x_1)^2 + (1 - x_2)^2 + \frac{1}{2} (2x_2 - x_1^2)^2$$

Stopping: $||\nabla f||_{\infty} \leq 10^{-6}$.



k = 0

k = 2

$$f(x_1, x_2) = (1 - x_1)^2 + (1 - x_2)^2 + \frac{1}{2}(2x_2 - x_1^2)^2$$

Stopping: $||\nabla f||_{\infty} \leq 10^{-6}$.



k = 2

Convergence Issues



Also, the computation of the Hessian is costly.

Theorem 10

Assume f is defined and twice differentiable on a convex set \mathcal{N} . Assume that ∇f is L-smooth on \mathcal{N} .

Let x_* be a minimizer of f(x) in \mathcal{N} and assume that $\nabla^2 f(x_*)$ is positive definite.

If $||x_0 - x_*||$ is sufficiently small, then $\{x_k\}$ converges quadratically to x_* .

Theorem 10 Assume f is defined and twice differentiable on a convex set \mathcal{N} . Assume that ∇f is L-smooth on \mathcal{N} . Let x_* be a minimizer of f(x) in \mathcal{N} and assume that $\nabla^2 f(x_*)$ is positive definite. If $||x_0 - x_*||$ is sufficiently small, then $\{x_k\}$ converges quadratically to x_* .

Note that the theorem implicitly assumes that $\nabla^2 f(x_k)$ is nonsingular for every k.

Theorem 10 Assume f is defined and twice differentiable on a convex set \mathcal{N} . Assume that ∇f is L-smooth on \mathcal{N} . Let x_* be a minimizer of f(x) in \mathcal{N} and assume that $\nabla^2 f(x_*)$ is positive definite. If $||x_0 - x_*||$ is sufficiently small, then $\{x_k\}$ converges quadratically to x_* .

Note that the theorem implicitly assumes that $\nabla^2 f(x_k)$ is nonsingular for every k.

As the theorem is concerned only with x_k approaching x^* , the continuity of $\nabla^2 f(x_k)$ and positive definiteness of $\nabla^2 f(x^*)$ imply that $\nabla^2 f(x_k)$ is positive definite for all sufficiently large k.

Theorem 10 Assume f is defined and twice differentiable on a convex set \mathcal{N} . Assume that ∇f is L-smooth on \mathcal{N} . Let x_* be a minimizer of f(x) in \mathcal{N} and assume that $\nabla^2 f(x_*)$ is positive definite. If $||x_0 - x_*||$ is sufficiently small, then $\{x_k\}$ converges quadratically to x_* .

Note that the theorem implicitly assumes that $\nabla^2 f(x_k)$ is nonsingular for every k.

As the theorem is concerned only with x_k approaching x^* , the continuity of $\nabla^2 f(x_k)$ and positive definiteness of $\nabla^2 f(x^*)$ imply that $\nabla^2 f(x_k)$ is positive definite for all sufficiently large k.

However, what happens if we start far away from a minimizer?

Newton's Method with Line Search

Algorithm 8 Newton's Method with Line Search

Input: x_0 starting point, $\varepsilon > 0$

Output: x^* approximation to a stationary point

1:
$$k \leftarrow 0$$

2: while
$$\|\nabla f_k\|_{\infty} > \varepsilon$$
 do

3:
$$p_k \leftarrow -H_k^{-1} \nabla f(x_k)$$

4: Set α_{init} for line search

5:
$$\alpha_k \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$$

$$6: \qquad x_{k+1} \leftarrow x_k + p_k$$

7:
$$k \leftarrow k+1$$

8: end while



Here $\ell_1 = 12, \ell_2 = 8, k_1 = 1, k_2 = 10, mg = 7$

Two Spring Problem - Newton's Method



Gradient descent, line search, stop. cond. $||\nabla f||_{\infty} \leq 10^{-6}$. Compare this with 32 iterations of gradient descent. Rosenbrock Function - Newton's Method Rosenbrock: $f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$



Gradient descent, line search, stop. cond. $||\nabla f||_{\infty} \leq 10^{-6}$. Compare this with 10,662 iterations of gradient descent.
Global Convergence with Line Search

Recall the Zoutendijk's theorem.

Denote by θ_k the angle between p_k and $-\nabla f_k$, i.e., satisfying

$$\cos \theta_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \, \|p_k\|}$$

Recall that f is L-smooth on a set \mathcal{N} for some L > 0 if

$$\|
abla f(x) -
abla f(ilde{x})\| \leq L \|x - ilde{x}\|, \quad ext{ for all } x, ilde{x} \in \mathcal{N}$$

Theorem 11 (Zoutendijk)

Consider $x_{k+1} = x_k + \alpha_k p_k$, where p_k is a descent direction and α_k satisfies the strong Wolfe conditions. Suppose that f is bounded below in \mathbb{R}^n and that f is continuously differentiable in an open set \mathcal{N} containing the level set $\{x : f(x) \leq f(x_0)\}$. Assume also that f is L-smooth on \mathcal{N} . Then

$$\sum_{k\geq 0}\cos^2\theta_k \, \|\nabla f_k\|^2 < \infty.$$

Assume that all α_k satisfy strong Wolfe conditions.

Assume that all α_k satisfy strong Wolfe conditions.

Assume that the Hessians H_k are positive definite with a uniformly bounded condition number:

 $||H_k|| ||H_k^{-1}|| \le M$ for all k

Assume that all α_k satisfy strong Wolfe conditions.

Assume that the Hessians H_k are positive definite with a uniformly bounded condition number:

 $||H_k|| ||H_k^{-1}|| \le M$ for all k

Then θ_k between $p_k = -H_k^{-1} \nabla f_k$ and $-\nabla f_k$ and satisfies

 $\cos \theta_k \geq 1/M$

Assume that all α_k satisfy strong Wolfe conditions.

Assume that the Hessians H_k are positive definite with a uniformly bounded condition number:

 $||H_k|| ||H_k^{-1}|| \le M$ for all k

Then θ_k between $p_k = -H_k^{-1} \nabla f_k$ and $-\nabla f_k$ and satisfies

 $\cos \theta_k \geq 1/M$

Thus, under the assumptions of Zoutendijk's theorem, we obtain

$$\frac{1}{M^2} \sum_{k \ge 0} \|\nabla f_k\|^2 \le \sum_{k \ge 0} \cos^2 \theta_k \, \|\nabla f_k\|^2 < \infty$$

which implies that $\lim_{k\to\infty} ||\nabla f_k|| = 0$.

Assume that all α_k satisfy strong Wolfe conditions.

Assume that the Hessians H_k are positive definite with a uniformly bounded condition number:

 $||H_k|| ||H_k^{-1}|| \le M$ for all k

Then θ_k between $p_k = -H_k^{-1} \nabla f_k$ and $-\nabla f_k$ and satisfies

 $\cos \theta_k \geq 1/M$

Thus, under the assumptions of Zoutendijk's theorem, we obtain

$$\frac{1}{M^2}\sum_{k\geq 0} \|\nabla f_k\|^2 \leq \sum_{k\geq 0} \cos^2 \theta_k \, \|\nabla f_k\|^2 < \infty$$

which implies that $\lim_{k\to\infty} ||\nabla f_k|| = 0$.

What if H_k is not positive definite or (nearly) singular?

Eigenvalue Modification

Consider $H_k = \nabla^2 f(x_k)$ and consider its diagonal form:

$$H_k = QDQ^T$$

Where D contains the eigenvalues of H_k on the diagonal, i.e., $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ and Q is an orthogonal matrix.

Eigenvalue Modification

Consider $H_k = \nabla^2 f(x_k)$ and consider its diagonal form:

$$H_k = Q D Q^T$$

Where D contains the eigenvalues of H_k on the diagonal, i.e., $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ and Q is an orthogonal matrix.

Observe that

- H_k is not positive definite iff $\lambda_i \leq 0$ for some *i*
- ▶ $||H_k||$ grows with max{ $\lambda_1, \ldots, \lambda_n$ } going to infinity.
- ► ||H_k⁻¹|| grows with min{λ₁,...,λ_n} going to 0 (i.e., the matrix becomes close to a singular matrix)

We want to prevent all three cases, i.e., make sure that for some reasonably large $\delta > 0$ we have $\lambda_i \ge \delta$ but not too large.

Eigenvalue Modification

Consider $H_k = \nabla^2 f(x_k)$ and consider its diagonal form:

$$H_k = QDQ^T$$

Where D contains the eigenvalues of H_k on the diagonal, i.e., $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ and Q is an orthogonal matrix.

Observe that

- H_k is not positive definite iff $\lambda_i \leq 0$ for some *i*
- ▶ $||H_k||$ grows with max $\{\lambda_1, \ldots, \lambda_n\}$ going to infinity.
- ||H_k⁻¹|| grows with min{λ₁,...,λ_n} going to 0 (i.e., the matrix becomes close to a singular matrix)

We want to prevent all three cases, i.e., make sure that for some reasonably large $\delta > 0$ we have $\lambda_i \ge \delta$ but not too large.

Two questions are in order:

- What is a reasonably large δ ?
- How to modify H_k so the minimum is large enough?

Consider an example:

$$abla f(x_k) = (1, -3, 2)$$
 and

$$\nabla^2 f(x_k) = \mathsf{diag}(10, 3, -1)$$

Consider an example:

$$\nabla f(x_k) = (1, -3, 2)$$
 and $\nabla^2 f(x_k) = diag(10, 3, -1)$

Now, the diagonalization is trivial:

$$abla^2 f(x_k) = Q$$
 diag $(10,3,-1) \; Q^ op \; Q = I$ is the identity matrix

Consider an example:

 $\nabla f(x_k) = (1, -3, 2)$ and $\nabla^2 f(x_k) = diag(10, 3, -1)$

Now, the diagonalization is trivial:

 $abla^2 f(x_k) = Q \operatorname{diag}(10,3,-1) \ Q^ op \quad Q = I$ is the identity matrix

What if we consider a minimum modification replacing the negative eigenvalue with a small number, say $\delta = 10^{-8}$?

Consider an example:

 $\nabla f(x_k) = (1, -3, 2)$ and $\nabla^2 f(x_k) = diag(10, 3, -1)$

Now, the diagonalization is trivial:

 $abla^2 f(x_k) = Q \operatorname{diag}(10,3,-1) \ Q^ op \quad Q = I$ is the identity matrix

What if we consider a minimum modification replacing the negative eigenvalue with a small number, say $\delta = 10^{-8}$? Obtain

$$B_k = Q \operatorname{diag}(10, 3, 10^{-8}) \ Q^{ op} = \operatorname{diag}(10, 3, 10^{-8})$$

Consider an example:

 $\nabla f(x_k) = (1, -3, 2)$ and $\nabla^2 f(x_k) = diag(10, 3, -1)$

Now, the diagonalization is trivial:

 $abla^2 f(x_k) = Q \operatorname{diag}(10,3,-1) \ Q^ op \quad Q = I$ is the identity matrix

What if we consider a minimum modification replacing the negative eigenvalue with a small number, say $\delta = 10^{-8}$? Obtain

$$B_k = Q ext{ diag}(10, 3, 10^{-8}) ext{ } Q^ op = ext{diag}(10, 3, 10^{-8})$$

If used in Newton's method, we obtain the following direction:

$$p_k = -B_k^{-1} \nabla f(x_k) = (1/10, -1, -(2 \cdot 10^8))$$

Thus, a very long vector almost parallel to the third dimension.

Consider an example:

 $\nabla f(x_k) = (1, -3, 2)$ and $\nabla^2 f(x_k) = diag(10, 3, -1)$

Now, the diagonalization is trivial:

 $abla^2 f(x_k) = Q \operatorname{diag}(10,3,-1) \ Q^ op \quad Q = I ext{ is the identity matrix}$

What if we consider a minimum modification replacing the negative eigenvalue with a small number, say $\delta = 10^{-8}$? Obtain

$$B_k = Q ext{ diag}(10, 3, 10^{-8}) ext{ } Q^ op = ext{diag}(10, 3, 10^{-8})$$

If used in Newton's method, we obtain the following direction:

$$p_k = -B_k^{-1} \nabla f(x_k) = (1/10, -1, -(2 \cdot 10^8))$$

Thus, a very long vector almost parallel to the third dimension.

Even though f decreases along p_k , it is far from the minimum of the quadratic approximation of f.

Note that the original Newton's direction is

 $-\text{diag}(1/10, 1/3, -1)(1, -3, 2)^{\top} = (-1/10, 1, 2)$ which is completely different.

Other strategies for eigenvalue modification can be devised.

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix B_k should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix B_k should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Strategies for eigenvalue modification include flipping negative eigenvalues to positive values, substituting negative eigenvalues with small positive ones, etc.

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix B_k should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Strategies for eigenvalue modification include flipping negative eigenvalues to positive values, substituting negative eigenvalues with small positive ones, etc.

There is no consensus on the best method for the modification.

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix B_k should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Strategies for eigenvalue modification include flipping negative eigenvalues to positive values, substituting negative eigenvalues with small positive ones, etc.

There is no consensus on the best method for the modification.

The implementation is based on computing $B_k = H_k + \Delta H_k$ for an appropriate modification matrix ΔH_k . What is ΔH_k in our example?

Other strategies for eigenvalue modification can be devised.

The criteria are rather loose. The resulting matrix B_k should be

- positive definite,
- ▶ of bounded norm (for all *k*),
- not too close to being singular.

(i.e., the eigenvalues should be sufficiently large)

Strategies for eigenvalue modification include flipping negative eigenvalues to positive values, substituting negative eigenvalues with small positive ones, etc.

There is no consensus on the best method for the modification.

The implementation is based on computing $B_k = H_k + \Delta H_k$ for an appropriate modification matrix ΔH_k .

What is ΔH_k in our example?

Various methods for computing ΔH_k have been devised in literature. Typically, it is based on some computationally cheaper decomposition than spectral decomposition (e.g., Cholesky).

Modified Newton's Method

Algorithm 9 Newton's Method with Line Search **Input:** x_0 starting point, $\varepsilon > 0$ **Output:** x^* approximation to a stationary point 1: $k \leftarrow 0$ 2: while $\|\nabla f_k\|_{\infty} > \varepsilon$ do $H_{\iota} \leftarrow \nabla^2 f(x_k)$ 3: if H_k is not sufficiently positive definite then 4: $H_k \leftarrow H_k + \Delta H_k$ so that H_k is sufficiently pos. definite 5: end if 6: Solve $H_k p_k = -\nabla f(x_k)$ for p_k 7: Set $x_{k+1} = x_k + \alpha_k p_k$, here α_k sat. the Wolfe cond. 8: $k \leftarrow k + 1$ 9. 10: end while

Newton's method is scale invariant.

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.
- Without modification, it may converge to an arbitrary stationary point (maximum, saddle point).

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.
- Without modification, it may converge to an arbitrary stationary point (maximum, saddle point).
- Computationally expensive:
 - \$\mathcal{O}(n^2)\$ second derivatives in the Hessian, each may be hard to compute.

Automated derivation methods help but still need store $O(n^2)$ results.

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.
- Without modification, it may converge to an arbitrary stationary point (maximum, saddle point).
- Computationally expensive:
 - \$\mathcal{O}(n^2)\$ second derivatives in the Hessian, each may be hard to compute.

Automated derivation methods help but still need store $O(n^2)$ results.

\$\mathcal{O}(n^3)\$ arithmetic operations to solve the linear system for the direction \$p_k\$.

May be mitigated by more efficient methods in case of sparse Hessians.

- Newton's method is scale invariant.
- Quadratic convergence in a close vicinity of a strict minimizer.
- Without modification, it may converge to an arbitrary stationary point (maximum, saddle point).
- Computationally expensive:
 - \$\mathcal{O}(n^2)\$ second derivatives in the Hessian, each may be hard to compute.

Automated derivation methods help but still need store $\mathcal{O}(n^2)$ results.

\$\mathcal{O}(n^3)\$ arithmetic operations to solve the linear system for the direction \$p_k\$.

May be mitigated by more efficient methods in case of sparse Hessians.

In a sense, Newton's method is an impractical "ideal" with which other methods are compared.

The efficiency issues (and the necessity of second-order derivatives) will be mitigated by using quasi-Newton methods.

Quasi-Newton Methods

Quasi-Newton Methods

Recall that Newton's method step p_k in $x_{k+1} = x_k + p_k$ comes from minimization of

$$q(p) = f_k + \nabla f_k^\top p + \frac{1}{2} p^\top H_k p$$

w.r.t. p by setting q'(p) = 0 and solving

$$-H_k p = -\nabla f_k$$

So Newton's method needs the second derivative (Hessian) that is computationally hard to obtain.

Gradient descent needs only the first derivatives but converges slowly.

Can we find a compromise?

Quasi-Newton methods use first derivatives to approximate the Hessian H_k in Newton's method with a matrix \tilde{H}_k .

BFGS

Denote by \tilde{H}_k the approximate of the Hessian $H_k = \nabla^2 f(x_k)$.

Suppose we just obtained the new point x_{k+1} after a line search starting from x_k in the direction p_k .

We can write the new quadratic approximation of f at x_{k+1} based on an updated Hessian approximation as follows:

$$q\left(p
ight)=f_{k+1}+
abla f_{k+1}^{ op}p+rac{1}{2}p^{ op} ilde{H}_{k+1}p.$$

Assume that f_{k+1} and ∇f_{k+1} are given, but we do not have the new approximate Hessian yet. Taking the gradient of this quadratic concerning p, we obtain

$$abla q(p) =
abla f_{k+1} + \tilde{H}_{k+1}p$$

Now we demand that the gradient ∇q of q w.r.t. p matches the gradient of f at x_{k+1} and at x_k .

$$q(p) = f_{k+1} + \nabla f_{k+1}^{\top} p + \frac{1}{2} p^{\top} \tilde{H}_{k+1} p$$
$$\nabla q(p) = \nabla f_{k+1} + \tilde{H}_{k+1} p$$

The gradient of the quadratic matching ∇f at x_k and x_{k+1} :



Note that $\nabla q(0) = \nabla f_{k+1}$ (just set p = 0 above).

$$q(p) = f_{k+1} + \nabla f_{k+1}^{\top} p + \frac{1}{2} p^{\top} \tilde{H}_{k+1} p$$
$$\nabla q(p) = \nabla f_{k+1} + \tilde{H}_{k+1} p$$

The gradient of the quadratic matching ∇f at x_k and x_{k+1} :



Note that $\nabla q(0) = \nabla f_{k+1}$ (just set p = 0 above). Just impose $\nabla q(-\alpha_k p_k) = \nabla f_{k+1} - \alpha_k \tilde{H}_{k+1} p_k = \nabla f_k$

$$q(p) = f_{k+1} + \nabla f_{k+1}^{\top} p + \frac{1}{2} p^{\top} \tilde{H}_{k+1} p$$
$$\nabla q(p) = \nabla f_{k+1} + \tilde{H}_{k+1} p$$

Just impose $\nabla q(-\alpha_k p_k) = \nabla f_{k+1} - \alpha_k \tilde{H}_{k+1} p_k = \nabla f_k$ Now, apparently, we have

$$\nabla f_{k+1} - \alpha_k \tilde{H}_{k+1} p_k = \nabla f_k \Rightarrow$$
$$\alpha_k \tilde{H}_{k+1} p_k = \nabla f_{k+1} - \nabla f_k$$

To simplify the notation, we define the step as

$$s_k = x_{k+1} - x_k = \alpha_k p_k,$$

and the difference in the gradient as

$$y_k = \nabla f_{k+1} - \nabla f_k.$$

Using this notation, we get the secant condition

$$\tilde{H}_{k+1}s_k = y_k$$

Now, we can obtain an approximate Hessian \tilde{H}_{k+1} by solving the secant condition $\tilde{H}_{k+1}s_k = y_k$.

Ideally, we want to

have H
 ^k_{k+1} symmetric positive definite
 To have a nice model for minimization around x_{k+1}.
 obtain H
 ^k_{k+1} from H
 ^k_k by

 $ilde{H}_{k+1} = ilde{H}_k + ext{something}$

To have a nice iterative algorithm.

Even if we demand symmetric positive definite solutions to the secant condition, there are infinitely many.

Note that the information about the solution is somehow present in s_k and y_k , so it is natural to compose the solution using these vectors.

We strive to choose \tilde{H}_{k+1} "close" to \tilde{H}_k .

Symmetric Rank One Update Consider $u_k = (y_k - \tilde{H}_k s_k)$

$$\tilde{H}_{k+1} = \tilde{H}_k + \frac{uu^\top}{u^\top s_k}$$

Now, the secant condition is satisfied:

$$ilde{H}_{k+1}s_k = ilde{H}_ks_k + rac{uu^{ op}s_k}{u^{ op}s_k} = ilde{H}_ks_k + u_k = y_k$$

Note that the updated matrix $\frac{uu^{\top}}{u^{\top}s_k}$ is of rank one and is a unique symmetric rank one matrix which makes \tilde{H}_{k+1} satisfy the secant condition.

To obtain a quasi-Newton method, it suffices to initialize \tilde{H}_0 , typically to the identity *I*, and use \tilde{H}_k instead of the Hessian $H_k = \nabla^2 f_k$ in Newton's method.

Even though \tilde{H}_k is a symmetric positive definite, the updated matrix \tilde{H}_{k+1} does not have to be a symmetric positive definite.
Rank One Update

Algorithm 10 Rank 1 update v1

$$k \leftarrow 0$$

$$\alpha_{init} \leftarrow 1$$

$$\tilde{V}_0 \leftarrow I \quad (or \quad \tilde{V}_0 \leftarrow 1/||\nabla f|| \cdot I)$$

while $||\nabla f_k||_{\infty} > \varepsilon$ do
 $s \leftarrow x_k - x_{k-1}$
 $y \leftarrow \nabla f_k - \nabla f_{k-1}$
 $\tilde{H}_k = \tilde{H}_{k-1} + \frac{uu^{\top}}{u^{\top}s_k}$
Solve for p_k in $\tilde{H}_k^{-1}p_k = -\nabla f_k$
 $\alpha \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$
 $x_{k+1} \leftarrow x_k + \alpha p_k$
 $k \leftarrow k+1$
end while

Symmetric Rank Two Update

Consider

$$ilde{H}_{k+1} = ilde{H}_k - rac{\left(ilde{H}_k s_k
ight) \left(ilde{H}_k s_k
ight)^ op}{s_k^ op ilde{H}_k s_k} + rac{y_k y_k^ op}{y_k^ op s_k}$$

Once again, verifying $\tilde{H}_{k+1}s_k = y_k$ is not difficult.

Lemma 1 If \tilde{H}_k is symmetric positive definite, then \tilde{H}_{k+1} is positive definite iff $y_k^\top s_k > 0$. $y_k^\top s_k > 0$ is called curvature condition

Now, it is not difficult to prove that if proper line search is performed, satisfying the strong Wolfe conditions, the curvature condition $y_k^{\top} s_k > 0$ will always be satisfied.

Thus, starting with a symmetric positive definite \tilde{H}_0 (e.g., a scalar multiple of I), every \tilde{H}_k is symmetric positive definite and satisfies the secant condition.

BFGS

Algorithm 11 BFGS v1

$$k \leftarrow 0$$

$$\alpha_{init} \leftarrow 1$$

$$\tilde{V}_{0} \leftarrow I \quad (or \quad \tilde{V}_{0} \leftarrow 1/\|\nabla f\| \cdot I)$$

while $\|\nabla f_{k}\|_{\infty} > \tau$ do
 $s \leftarrow x_{k} - x_{k-1}$
 $y \leftarrow \nabla f_{k} - \nabla f_{k-1}$
 $\tilde{H}_{k} \leftarrow \tilde{H}_{k-1} - \frac{(\tilde{H}_{k-1}s_{k})(\tilde{H}_{k-1}s_{k})^{\top}}{s_{k}^{\top}\tilde{H}_{k-1}s_{k}} + \frac{y_{k}y_{k}^{\top}}{y_{k}^{\top}s_{k}}$
Solve for p_{k} in $\tilde{H}_{k}^{-1}p_{k} = -\nabla f_{k}$
 $\alpha \leftarrow \text{linesearch}(p_{k}, \alpha_{\text{init}})$
 $x_{k+1} \leftarrow x_{k} + \alpha p_{k}$
 $k \leftarrow k + 1$
end while

Note that we still have to solve a linear system for p_k .

Sherman-Morrison-Woodbury Formula

Ideally, we would like to compute \tilde{H}_k^{-1} iteratively along the optimization, i.e.,

 $ilde{H}_{k+1}^{-1} = ilde{H}_k^{-1} + ext{something}$

To get such a "something" we use the following Sherman–Morrison–Woodbury (SMW) formula:

$$(A + UV^{T})^{-1} = A^{-1} - A^{-1}U(I + V^{T}U)^{-1}V^{T}A^{-1}$$

where

$$U = [u_1, u_2, \dots, u_k]$$
 $V = [v_1, v_2, \dots, v_k]$

SMW can be written as

$$\left(A + \sum_{i=1}^{k} u_i v_i^{T}\right)^{-1} = A^{-1} - A^{-1} U C^{-1} V^{T} A^{-1}$$

where

$$C_{ij} = \delta_{ij} + \mathbf{v}_i^T \mathbf{u}_j \quad i, j = 1, 2, \dots, k$$

Rank 1 – Iterative Inverse Hessian Approximation

Applying SMW to the rank one update

$$ilde{H}_{k+1} = ilde{H}_k + rac{\left(y_k - ilde{H}_k s_k
ight) \left(y_k - ilde{H}_k s_k
ight)^ op}{\left(y_k - ilde{H}_k s_k
ight)^ op s_k}$$

yields

$$\tilde{H}_{k+1}^{-1} = \tilde{H}_{k}^{-1} + \frac{\left(s_{k} - \tilde{H}_{k}^{-1}y_{k}\right)\left(s_{k} - \tilde{H}_{k}^{-1}y_{k}\right)^{\top}}{\left(s_{k} - \tilde{H}_{k}^{-1}y_{k}\right)^{\top}y_{k}}$$

Yes, only y and s swapped places.

This allows us to avoid solving for p_k in every iteration.

Rank One Update V2

Algorithm 12 Rank 1 update v1 $k \leftarrow 0$ $\alpha_{\text{init}} \leftarrow 1$ $\tilde{V}_0 \leftarrow I$ (or $\tilde{V}_0 \leftarrow 1/\|\nabla f\| \cdot I$) while $\|\nabla f_k\|_{\infty} > \tau$ do $s \leftarrow x_k - x_{k-1}$ $y \leftarrow \nabla f_k - \nabla f_{k-1}$ $\tilde{H}_k^{-1} \leftarrow \tilde{H}_{k-1}^{-1} + \frac{\left(s_k - \tilde{H}_{k-1}^{-1} y_k\right)\left(s_k - \tilde{H}_{k-1}^{-1} y_k\right)^{\top}}{\left(s_k - \tilde{H}_{k-1}^{-1} y_k\right)^{\top} y_k}$ $p_k \leftarrow -\tilde{H}_k^{-1} \nabla f_k$ $\alpha \leftarrow \text{linesearch}(p_k, \alpha_{\text{init}})$ $x_{k+1} \leftarrow x_k + \alpha p_k$ $k \leftarrow k + 1$ end while

BFGS

Applying SMW to the BFGS Hessian update

$$\tilde{H}_{k+1} = \tilde{H}_k - \frac{\left(\tilde{H}_k s_k\right) \left(\tilde{H}_k s_k\right)^\top}{s_k^\top \tilde{H}_k s_k} + \frac{y_k y_k^\top}{y_k^\top s_k}$$

yields

$$H_{k+1}^{-1} = \left(I - \frac{s_k y_k^\top}{s_k^\top y_k}\right) H_k^{-1} \left(I - \frac{y_k s_k^\top}{s_k^\top y_k}\right) + \frac{s_k s_k^\top}{s_k^\top y_k}$$

We avoid solving the linear system for p_k .

BFGS V2

Algorithm 13 BFGS v2

$$k \leftarrow 0$$

$$\alpha_{init} \leftarrow 1$$

$$\tilde{V}_{0} \leftarrow I \quad (or \quad \tilde{V}_{0} \leftarrow 1/\|\nabla f\| \cdot I)$$

while $\|\nabla f_{k}\|_{\infty} > \tau$ do
 $s \leftarrow x_{k} - x_{k-1}$
 $y \leftarrow \nabla f_{k} - \nabla f_{k-1}$
 $H_{k}^{-1} \leftarrow \left(I - \frac{s_{k}y_{k}^{\top}}{s_{k}^{\top}y_{k}}\right) H_{k-1}^{-1} \left(I - \frac{y_{k}s_{k}^{\top}}{s_{k}^{\top}y_{k}}\right) + \frac{s_{k}s_{k}^{\top}}{s_{k}^{\top}y_{k}}$
 $p_{k} \leftarrow -\tilde{H}_{k}^{-1}\nabla f_{k}$
 $\alpha \leftarrow \text{linesearch}(p_{k}, \alpha_{\text{init}})$
 $x_{k+1} \leftarrow x_{k} + \alpha p_{k}$
 $k \leftarrow k+1$
end while

Another View on BFGS (Optional)

We search for \tilde{H}_{k+1}^{-1} where \tilde{H}_{k+1} satisfies $\tilde{H}_{k+1}s_k = y_k$. Simply, search for a solution \tilde{V}_{k+1} for $\tilde{V}_{k+1}y_k = s_k$.

The idea is to use \tilde{V}_{k+1} close to \tilde{V}_k (in some sense):

$$\min_{ ilde{V}} \left\| ilde{V} - ilde{V}_k
ight\|$$

subject to $ilde{V} = ilde{V}^ op, \quad ilde{V} y_k = s_k$

Here the norm is weighted Frobenius norm:

$$\|A\| \equiv \left\| W^{1/2} A W^{1/2} \right\|_F,$$

where $\|\cdot\|_F$ is defined by $\|C\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n c_{ij}^2$. The weight W can be chosen as any matrix satisfying the relation $Wy_k = s_k$.

BFGS is obtained with $W = \overline{G}_k^{-1}$ where \overline{G}_k is the average Hessian defined by $\overline{G}_k = \left[\int_0^1 \nabla^2 f(x_k + \tau \alpha_k p_k) d\tau\right]$

Solving this gives precisely the BFGS formula for $\tilde{H}_{k+1}^{-1}.$

Quasi-Newton Methods Convergence

Quasi-Newton Methods Rate of Convergence

Quasi-Newton Methods - Practical Issues

Quasi-Newton Methods - Comments

When the number of design variables is extensive, working with the whole Hessian inverse approximation matrix might not be practical.

This motivates limited-memory quasi-Newton methods,

In addition, these methods also improve the computational efficiency of medium-sized problems (hundreds or thousands of design variables) with minimal sacrifice in accuracy.

L-BFGS

Recall that we compute iteratively the approximation to the inverse Hessian by

$$H_{k+1}^{-1} = \left(I - \frac{s_k y_k^\top}{s_k^\top y_k}\right) H_k^{-1} \left(I - \frac{y_k s_k^\top}{s_k^\top y_k}\right) + \frac{s_k s_k^\top}{s_k^\top y_k}$$

However, eventually, we are interested in

$$p_k = H_k^{-1} \nabla f$$

Note that given the sequences s_1, \ldots, s_k and y_1, \ldots, y_k and H_0^{-1} we can recursively compute H_{k+1}^{-1} for every k.

What if we limit the sequences in memory to just m last elements:

$$s_{k-m+1}, s_{k-m+2}, \dots, s_k$$
 $y_{k-m+1}, y_{k-m+2}, \dots, y_k$

In practice, m between 5 and 20 is usually sufficient. We also initialize the recurrence with the last iterate:

L-BFGS

Let us rewrite the BFGS update formula as follows:

$$\tilde{H}_{k+1}^{-1} = V_k^T \tilde{H}_k^{-1} V_k + \rho_k s_k s_k^\top$$

where

$$\rho_k = s_k^\top y_k \quad \text{and} \quad V_k = I - \rho_k s_k y_k^\top$$
$$s_k = x_{k+1} - x_k \quad \text{and} \quad y_k = \nabla f_{k+1} - \nabla f_k$$

By substitution we obtain

$$\begin{split} \tilde{H}_{k}^{-1} &= \left(V_{k-1}^{T} \cdots V_{k-m}^{T} \right) \tilde{H}_{k}^{0} \left(V_{k-m} \cdots V_{k-1} \right) \\ &+ \rho_{k-m} \left(V_{k-1}^{T} \cdots V_{k-m+1}^{T} \right) s_{k-m} s_{k-m}^{T} \left(V_{k-m+1} \cdots V_{k-1} \right) \\ &+ \rho_{k-m+1} \left(V_{k-1}^{T} \cdots V_{k-m+2}^{T} \right) s_{k-m+1} s_{k-m+1}^{T} \left(V_{k-m+2} \cdots V_{k} \right) \\ &+ \cdots \\ &+ \cdots \\ &+ \rho_{k-1} s_{k-1} s_{k-1}^{T} \end{split}$$

L-BFGS Algorithm

Algorithm 14 L-BFGS two-loop recursion **Input:** : $s_{k-1}, ..., s_{k-m}$ and $y_{k-1}, ..., y_{k-m}$ **Output:** : p_k the search direction $-\tilde{H}_k^{-1}\nabla f_k$ 1: $a \leftarrow \nabla f_{k}$ 2: for $i = k - 1, k - 2, \dots, k - m$ do 3: $\alpha_i \leftarrow \rho_i s_i^T q$ 4: $q \leftarrow q - \alpha_i y_i$ 5: end for 6: $r \leftarrow H^0_{\mu} q$ 7: for $i = k - m, k - m + 1, \dots, k - 1$ do 8: $\beta \leftarrow \rho_i v_i^T r$ 9: $r \leftarrow r + s_i(\alpha_i - \beta)$ 10: end for 11: stop with result $\tilde{H}_{\iota}^{-1} \nabla f_k = r$

L-BFGS Algorithm

Algorithm 15 L-BFGS

- 1: Choose starting point x_0 , integer m > 0
- 2: $k \leftarrow 0$
- 3: repeat

4: Choose
$$H_k^0$$
 e.g. $\frac{s_{k-1}^{-1}y_{k-1}}{y_{k-1}^{-1}y_{k-1}}$

- 5: Compute $p_k \leftarrow -H_k \nabla f_k$ using the previous algorithm
- 6: Compute $x_{k+1} \leftarrow x_k + \alpha_k p_k$, where α_k is chosen to satisfy the strong Wolfe conditions
- 7: **if** k > m **then**
- 8: Discard the vector pair $\{s_{k-m}, y_{k-m}\}$ from storage
- 9: end if
- 10: Compute and save $s_k \leftarrow x_{k+1} x_k$, $y_k \leftarrow \nabla f_{k+1} \nabla f_k$
- 11: $k \leftarrow k+1$

12: **until** convergence

$$f(x_1, x_2) = (1 - x_1)^2 + (1 - x_2)^2 + \frac{1}{2} (2x_2 - x_1^2)^2$$

Stopping: $||\nabla f||_{\infty} \leq 10^{-6}$.



In L-BFGS the memory length m was 5. The results are similar.



Here $\ell_1 = 12, \ell_2 = 8, k_1 = 1, k_2 = 10, mg = 7$



Rosenbrock: $f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$



Rosenbrock:

$$f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$$



Computational Complexity

Algorithm	Computational Complexity
Steepest Descent	$O(n^2)$ per iteration
Conjugate Gradients	O(n) per iteration
Newton's Method	$O(n^3)$ to compute Hessian and solve system
BFGS	$O(n^2)$ to update Hessian approximation

Table: Summary of the computational complexity for each optimization algorithm.

- Steepest Descent: Simple but often slow, requiring many iterations.
- Conjugate Gradients: Efficient for large sparse systems, fewer iterations.
- Newton's Method: Fast convergence but expensive per iteration.
- BFGS: Quasi-Newton, no Hessian needed, good speed and iteration count balance.

Constrained Optimization

Constrained Optimization Problem

Recall that the constrained optimization problem is

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{by varying} & x \\ \text{subject to} & g_j(x) \leq 0 \quad j = 1, \dots, n_g \\ & h_l(x) = 0 \quad l = 1, \dots, n_h \end{array}$$

 x^* is now a *constrained minimizer* if

$$f(x^*) \leq f(x)$$
 for all $x \in \mathcal{F}$

where \mathcal{F} is the feasibility region

$$\mathcal{F} = \{x \mid g_i(x) \le 0, h_j(x) = 0, j = 1, \dots, n_x, l = 1, \dots, n_h\}$$

Thus, to find a constrained minimizer, we have to inspect unconstrained minima of f inside of \mathcal{F} and points along the boundary of \mathcal{F} .

COP - Example



Equality Constraints

Let us restrict our problem only to the equality constraints:

minimize f(x)by varying xsubject to $h_j(x) = 0$ $j = 1, ..., n_h$

Assume that f and h_i have continuous second derivatives.

Now, we try to imitate the theory from the unconstrained case and characterize minima using gradients.

This time, we have to consider the gradient of f and h_i .

Half-Space of Decrease

Consider the first-order Taylor approximation of f at x

 $f(x+p) \approx f(x) + \nabla f(x)^{\top} p$

Note that if x^* is an unconstrained minimum of f, then

 $f(x^* + p) \ge f(x^*)$

for all p small enough.

Together with the Taylor approximation, we obtain

$$f(x^*) + \nabla f(x^*)^\top p \ge f(x^*)$$

and hence

$$\nabla f(x^*) \geq 0$$



The hyperplane defined by $\nabla f^{\top} p = 0$ contains directions p of zero variation in f.

In the unconstrained case, x^* is minimum only if $\nabla f(x^*) = 0$ because otherwise there would be a direction p satisfying $\nabla f(x^*)p < 0$, a *decrease direction*.

Decrease Direction in COP

In COP, p is a decrease direction in $x \in \mathcal{F}$ not only if $\nabla f(x)p < 0$, it also needs to be a *feasible direction*!

I.e., point into the feasible region.

How do we characterize feasible directions?

Consider Taylor approximation of h_j for all j:

$$h_j(x+p) \approx h_j(x) + \nabla h_j(x)^\top p$$

Assuming $x \in \mathcal{F}$, we have $h_j(x) = 0$ for all j and thus

$$h_j(x+p) \approx \nabla h_j(x)^\top p$$

As p is a feasible direction iff $h_j(x + p) = 0$, we obtain that p is a feasible direction iff

$$\nabla h_j(x)^\top p = 0$$
 for all j

Feasible Points and Directions



Here, the only feasible direction at x is p = 0.

Feasible Points and Directions



Here the feasible directions at x^* point along the red line, i.e.,

$$abla h_1(x^*)p = 0$$
 $abla h_2(x^*)p = 0$

Constrained Minima

Consider a direction p. Observe that

If h_j(x)[⊤]p ≠ 0, then moving a short step in the direction p violates the constraint h_j(x) = 0.

• If
$$h_j(x)^\top p = 0$$
 for all j and

- ∇f(x)p > 0, then moving a short step in the direction p increases f and stays in F.
- ∇f(x)p < 0, then moving a short step in the direction p decreases f and stays in F.</p>
- ∇f(x)p = 0, then moving a short step in the direction p does not change f and stays F.

To be a minimizer, x must be feasible and every direction satisfying $h_j(x)^\top p = 0$ for all j must also satisfy $\nabla f(x)p \ge 0$. Note that if p is a feasible direction, then -p is also. So finally, If x^* is a constrained minimizer, then

 $abla f(x^*)p = 0$ for all p such that $abla h_j(x^*)^\top p = 0$ for all j

Lagrange Multipliers



Left: f increases along p. Right: f does not change along p.

Observe that at an optimum, ∇f lies in the space spanned by the gradients of constraint functions.

There are Lagrange multipliers λ_1, λ_2 satisfying

$$\nabla f(x^*) = \lambda_1 \nabla h_1 + \lambda_2 \nabla h_2$$

Lagrange Multipliers

We know that if x^* is a constrained minimizer, then.

 $\nabla f(x)p = 0$ for all p such that $\nabla h_j(x)^\top p = 0$ for all j

But then, from the geometry of the problem, we obtain

Theorem 12

Consider the COP with only equality constraints and f and all h_j twice continuously differentiable.

Assume that x^* is a constrained minimizer and that x^* is regular, which means that $\nabla h_j(x^*)$ are linearly independent. Then there are $\lambda_1, \ldots, \lambda_{n_h} \in \mathbb{R}$ satisfying

$$abla f(x^*) = \sum_{j=1}^{n_h} \lambda_j
abla h_j(x^*)$$

The coefficients $\lambda_1, \ldots, \lambda_{n_h}$ are called *Lagrange multipliers*.
Lagrangian Function

Try to transform the constrained problem into an unconstrained one by moving the constraints $h_i(x) = 0$ into the objective.

Consider Lagrangian function $\mathcal{L}: \mathbb{R}^n \times \mathbb{R}^{n_h} \to \mathbb{R}$ defined by

 $\mathcal{L}(x,\lambda) = f(x) + h(x)^{\top} \lambda$ here $h(x) = (h_1(x), \dots, h_{n_h}(x))^{\top}$

Note that

$$abla_{\mathbf{x}}\mathcal{L} =
abla f(\mathbf{x}) + \sum_{j=1}^{n_h}
abla h_j(\mathbf{x})^\top \lambda_j$$
 $abla_{\lambda}\mathcal{L} = h(\mathbf{x})$

Now putting $\nabla \mathcal{L}(x) = 0$, we obtain precisely the above properties of the constrained minimizer:

$$h(x) = 0 \qquad ext{and} \qquad
abla f(x) = \sum_{j=1}^{n_h} -\lambda_j
abla h_j(x)^ op$$

However, we cannot use the unconstrained optimization methods here because searching for a minimizer in x asks for a maximizer in λ .

$$\begin{array}{ll} \underset{x_{1},x_{2}}{\text{minimize}} & f\left(x_{1},x_{2}\right) = x_{1} + 2x_{2} \\ \text{subject to} & h\left(x_{1},x_{2}\right) = \frac{1}{4}x_{1}^{2} + x_{2}^{2} - 1 = 0 \end{array}$$

The Lagrangian function

$$\mathcal{L}(x_1, x_2, \lambda) = x_1 + 2x_2 + \lambda \left(\frac{1}{4}x_1^2 + x_2^2 - 1\right)$$

Differentiating this to get the first-order optimality conditions,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x_1} &= 1 + \frac{1}{2}\lambda x_1 = 0 \qquad \frac{\partial \mathcal{L}}{\partial x_2} = 2 + 2\lambda x_2 = 0\\ \frac{\partial \mathcal{L}}{\partial \lambda} &= \frac{1}{4}x_1^2 + x_2^2 - 1 = 0. \end{aligned}$$

Solving these three equations for the three unknowns (x_1, x_2, λ) , we obtain two possible solutions:

$$x_A = (x_1, x_2) = (-\sqrt{2}, -\sqrt{2}/2), \quad \lambda_A = \sqrt{2}$$

 $x_B = (x_1, x_2) = (\sqrt{2}, \sqrt{2}/2), \quad \lambda_A = -\sqrt{2}$



Second-Order Sufficient Conditions

As in the unconstrained case, the first-order conditions characterize any "stable" point (minimum, maximum, saddle).

Consider Lagrangian Hessian:

$$H_{\mathcal{L}}(x,\lambda) = H_f(x) + \sum_{j=1}^{n_h} \lambda_j H_{h_j}(x)$$

Here H_f is the Hessian of f, and each H_{h_i} is the Hessian of h_j .

The second-order sufficient conditions are as follows: Assume x^* is regular and feasible. Also, assume that there is λ s.t.

$$abla f(x^*) = \sum_{j=1}^{n_h} -\lambda_j
abla h_j(x^*)^ op$$

and that

 $p^{\top}H_{\mathcal{L}}(x^*,\lambda)p > 0$ for all p satisfying $\nabla h_j(x^*)^{\top}p = 0$ for all j. Then, x^* is a constrained minimizer of f.

Inequality Constraints

Recall that the constrained optimization problem is

 $\begin{array}{ll} \text{minimize} & f(x) \\ \text{by varying} & x \\ \text{subject to} & g_i(x) \leq 0 \quad i = 1, \dots, n_g \\ & h_j(x) = 0 \quad j = 1, \dots, n_h \end{array}$

We say that a constraint $g_i(x) \le 0$ is *active* for x^* if $g_i(x^*) = 0$, otherwise it is *inactive* for x^* .

As before, if x^* is optimum, any small step in a feasible direction p must not decrease f, i.e.,

 $\nabla f(x^*)^\top p \geq 0$

How do we identify feasible directions for inequality constraints?

Feasible Directions

For inactive constraints, arbitrary direction p is feasible.

For active constraints $g_i(x) = 0$ we have

$$g_i(x+p) \approx g_i(x) + \nabla g_i(x)^\top p \leq 0, \quad i=1,\ldots,n_g$$

and p is feasible iff $\nabla g_i(x)^\top p \leq 0$ for all active constr. $g_i(x) = 0$.



Lagrange Multipliers

When can f (not) be decreased in a feasible direction?



Left: f decreases in the blue cone. Right: f does not decrease in any feasible direction.

At an optimum there are Lagrange multipliers $\sigma_1, \sigma_2 \ge 0$:

$$-\nabla f = \sigma_1 \nabla g_1 + \sigma_2 \nabla g_2$$

Lagrange Multipliers

We know that if x^* is a constrained minimizer, then.

 $\nabla f(x)p = 0$ for all p feasible

Using Farkas' lemma, one can prove the following

Theorem 13

Consider the COP with f and all g_i , h_j twice continuously differentiable.

Assume that x^* is a constrained minimizer and that x^* is regular which means that $\nabla g_i(x^*), \nabla h_j(x^*)$ are linearly independent. Then there are Lagrange multipliers $\lambda_1, \ldots, \lambda_{n_h} \in \mathbb{R}$ and $\sigma_1, \ldots, \sigma_{n_g} \in \mathbb{R}$ satisfying

$$abla f(x^*) = \sum_{j=1}^{n_h} \lambda_j
abla h_j(x^*) + \sum_{i=1}^{n_h} \sigma_i
abla g_i(x^*) \qquad \textit{where } \sigma_i \geq 0$$

Lagrangian Function

Note that inequality $g_i(x) \le 0$ can be equivalently expressed using a *slack variable* s_i by

$$g(x)+s_i^2=0$$

The Lagrangian function then generalizes from equality to inequality COP as follows.

$$\mathcal{L}(x,\lambda,\sigma,s) = f(x) + h(x)^{\top}\lambda + (g(x) + s \odot s)^{\top}\sigma$$

Here, $h(x) = (h_1(x), \ldots, h_{n_h}(x))^\top$, $g(x) = (g_1(x), \ldots, g_{n_g}(x))^\top$, $s = (s_1, \ldots, s_{n_g})$, and \odot is the component-wise multiplication.

Now compute the stable point of $\ensuremath{\mathcal{L}}$ by considering

$$abla_{x}\mathcal{L} = 0$$
 $abla_{\lambda}\mathcal{L} = 0$
 $abla_{\sigma}\mathcal{L} = 0$
 $abla_{s}\mathcal{L} = 0$

(see the whiteboard)

KKT

If x^* is a constrained minimizer and that x^* is regular. Then there are λ,σ,s satisfying

$$\frac{\partial f}{\partial x_{\ell}}(x) + \sum_{j=1}^{n_h} \lambda_j \frac{\partial h_j}{\partial x_{\ell}} + \sum_{j=1}^{n_g} \sigma_j \frac{\partial g_j}{\partial x_{\ell}} = 0 \qquad \ell = 1, \dots, n$$

$$h_j = 0 \qquad j = 1, \dots, n_h$$

$$g_i + s_i^2 = 0 \qquad i = 1, \dots, n_g$$

$$2\sigma_i s_i = 0 \qquad i = 1, \dots, n_g$$

$$\sigma_i \ge 0$$

So, solving the above system allows us to identify potential constrained minimizers.

To decide whether x^* solving KKT is a minimizer, check whether

$$p^{\top}H_{\mathcal{L}}(x,\lambda)p>0$$

For all feasible directions p (similarly to the equality case).

Example

$$\begin{array}{ll} \underset{x_{1},x_{2}}{\text{minimize}} & f\left(x_{1},x_{2}\right) = x_{1} + 2x_{2} \\ \text{subject to} & g\left(x_{1},x_{2}\right) = \frac{1}{4}x_{1}^{2} + x_{2}^{2} - 1 \leq 0. \end{array}$$

The Lagrangian function for this problem is

$$\mathcal{L}(x_1, x_2, \sigma, s) = x_1 + 2x_2 + \sigma \left(\frac{1}{4}x_1^2 + x_2^2 - 1 + s^2\right)$$



Example

$$\frac{\partial \mathcal{L}}{\partial x_1} = 1 + \frac{1}{2}\sigma x_1 = 0$$
$$\frac{\partial \mathcal{L}}{\partial x_2} = 2 + 2\sigma x_2 = 0$$
$$\frac{\partial \mathcal{L}}{\partial \sigma} = \frac{1}{4}x_1^2 + x_2^2 - 1 = 0$$
$$\frac{\partial \mathcal{L}}{\partial s} = 2\sigma s = 0.$$

Setting $\sigma = 0$ does not yield any solution. Setting s = 0 and $\sigma \neq 0$ we obtain

$$x_{A} = \begin{bmatrix} x_{1} \\ x_{2} \\ \sigma \end{bmatrix} = \begin{bmatrix} -\sqrt{2} \\ -\sqrt{2}/2 \\ \sqrt{2} \end{bmatrix}, \quad x_{B} = \begin{bmatrix} x_{1} \\ x_{2} \\ \sigma \end{bmatrix} = \begin{bmatrix} \sqrt{2} \\ \sqrt{2}/2 \\ -\sqrt{2} \end{bmatrix}$$

Now, σ must be non-negative, so only x_A is the solution. There is no feasible descent direction at x_A . We already know that the Hessian Lagrangian is positive definite, so this is a minimizer.

Penalty methods

The idea: Transform a constrained problem into an unconstrained one by adding a penalty to the objective function when constraints are violated or close to being violated.

Assuming an objective function f, the penalized objective is of the form

 $\hat{f}(x) = f(x) + \mu \pi(x)$

Here, μ is a fixed constant determining how strong the penalty should be, and π is the penalty function.

Now we may apply the unconstrained optimization methods (e.g., L-BFGS) to \hat{f} and obtain an approximation of a minimizer of f.

There are two types

- exterior penalizing infeasible x
- interior penalizing x close to being infeasible

Exterior Penalty Methods

Consider equality-constrained problems:

minimize
$$f(x)$$

by varying x
subject to $h_j(x) = 0$ $j = 1, \dots, n_h$

Consider quadratic penalty:

$$\hat{f}(x;\mu) = f(x) + \frac{\mu}{2} \sum_{j=1}^{n_h} h_j(x)^2$$

If f is continuously differentiable, \hat{f} is as well (w.r.t. x).

Quadratic Penalty



The true solution would be recovered for $\mu = \infty$.

However, large μ means large condition number of the Hessian of \hat{f} Intuitively, curvature of \hat{f} changes rapidly with the direction.

Need to choose μ carefully, iteratively.

Quadratic Penalty

The problems

- Small µ may result in so weak penalty that f unbounded below results in f unbounded as well
- As $\mu = \infty$ is impossible, the solution is always slightly infeasible
- \blacktriangleright Growing "curvature" of \hat{f} as μ grows making the Hessian of \hat{f} almost singular



Quadratic Penalty for Inequality Constraints



Minimizer approached from the infeasible side.

Example

$$\hat{f}(x;\mu) = x_1 + 2x_2 + \frac{\mu}{2} \max\left(0, \frac{1}{4}x_1^2 + x_2^2 - 1\right)^2$$



Augmented Lagrangian

Instead of minimizing f, we search for an optimal point of the Lagrangian.

Similarly, instead of minimizing \hat{f} we may augment the Lagrangian L with penalty and optimize the augmented Lagrangian

$$\hat{L}(x;\lambda,\mu) = f(x) + \sum_{i=1}^{n_h} \lambda_i h_i(x) + \frac{\mu}{2} \sum_{i=1}^{n_h} h_i(x)^2$$

Note the relationship between optimality conditions for L and \hat{L}

$$\nabla_{x}\hat{L}(x;\lambda,\mu) = \nabla f(x) + \sum_{i=1}^{n_{h}} (\lambda_{i} + \mu h_{i}(x)) \nabla h_{i} = 0$$
$$\nabla_{x}\mathcal{L}(x^{*},\lambda^{*}) = \nabla f(x^{*}) + \sum_{i=1}^{n_{h}} \lambda_{i}^{*} \nabla h_{i}(x^{*}) = 0.$$

Comparing these two conditions suggests an approximation:

$$\lambda_j^* \approx \lambda_j + \mu h_j.$$

Augmented Lagrangian Penalty Method

Inputs:

- ► x₀: Starting point
- $\lambda_0 = 0$: Initial Lagrange multiplier
- $\mu_0 > 0$: Initial penalty parameter
- $\rho > 1$: Penalty increase factor

Outputs:

- x*: Optimal point
- $f(x^*)$: Corresponding function value

Algorithm:

k = 0 not converged $x_{k+1} \leftarrow x$ minimizing $f(x; \lambda_k, \mu_k)$ $\lambda_{k+1} = \lambda_k + \mu_k h(x_k) \ \mu_{k+1} = \rho \mu_k \ k = k+1$

Comparison of Quadratic and Lagrangian Penalty

Compare

$$h_j pprox rac{1}{\mu} \left(\lambda_j^* - \lambda_j
ight).$$

with the corresponding approximation of h_j in the quadratic penalty method is

$$h_j \approx \frac{\lambda_j^+}{\mu}$$

Thus, the quadratic penalty relies solely on increasing μ .

However, the augmented Lagrangian also controls the numerator via estimating λ_j .

If λ_j is close to $\lambda_j^*,$ we may obtain a close solution for modest values of $\mu.$

Several variants of the Lagrangian penalty exist for inequality constraints; see Nocedal & Wright.

Interior Penalty Methods

Always seek to maintain feasibility as opposed to the exterior methods.

Instead of adding a penalty only when constraints are violated; add a penalty as the constraint is approached from the feasible region.

Desirable if the objective function is ill-defined outside the feasible region.

The interior methods are also referred to as *barrier methods* because the penalty function acts as a barrier preventing iterates from leaving the feasible region.

Barrier Methods

Minimize the augmented objective function.

$$\hat{f}(x;\mu) = f(x) + \mu \pi(x)$$

Here π is a penalty function.



Algorithms based on these penalties must be prevented from evaluating infeasible points.

Barrier methods



Solve a sequence of unconstrained problems for \hat{f} with $\mu \rightarrow 0$.

Every unconstrained optimization must start at an initial point feasible for the constrained problem.

The line search must check for feasibility and backtrack from steps to infeasible points.

Example

$$\hat{f}(x;\mu) = x_1 + 2x_2 - \mu \ln \left(-\frac{1}{4}x_1^2 - x_2^2 + 1\right)$$



As for exterior methods, the Hessian becomes increasingly ill-conditioned as $\mu \rightarrow 0.$

Various modifications exist that alleviate the above problem. These methods lead to a class of modern *interior point methods*. Summary of Penalty Methods

... not too efficient but simple

Quadratic Programming

The quadratic optimization problem with equality constraints is to

minimize $\frac{1}{2}x^{\top}Qx + q^{t}x$ by varying xsubject to Ax + b = 0

Here

- Q is a n × n symmetric matrix. For simplicity assume positive definite.
- A is a $m \times n$ matrix. Assume full rank.



Quadratic Programming

How to solve the quadratic program?

Consider the Lagrangian function

$$L(x,\lambda) = \frac{1}{2}x^{\top}Qx + q^{\top}x + \lambda^{\top}(Ax + b)$$

and its partial derivatives:

$$abla_x L(x) = Qx + q + A^\top \lambda = 0$$

 $abla_\lambda L(x) = Ax + b = 0$

As Q is positive definite, we know that a solution to the above system is a minimizer.

So in order to solve the quadratic program, it suffices to solve the system of linear equations.

The situation is much more complicated as some constraints can be active (i.e., equality) and some inactive (i.e., locally irrelevant).

There are methods based on the concept of *active set* which keeps track of active constraints that iteratively search for solutions.

We shall see an analogy in linear programming, now won't go any further.

Sequential Quadratic Programming