

ECE 544: Pattern Recognition

A. G. Schwing

University of Illinois at Urbana-Champaign, 2021

Scribe & Exercises

Optimization Primal

Goals of this lecture

Goals of this lecture

- Understanding the basics of optimization

Goals of this lecture

- Understanding the basics of optimization

Reading Material

Goals of this lecture

- Understanding the basics of optimization

Reading Material

- S. Boyd and L. Vandenberghe; Convex Optimization; Chapters 2-4

Optimization problems that we have seen so far:

Optimization problems that we have seen so far:

- Linear Regression

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left(y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

Optimization problems that we have seen so far:

- Linear Regression

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left(y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Logistic Regression

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log \left(1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})) \right)$$

Optimization problems that we have seen so far:

- Linear Regression

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left(y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Logistic Regression

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log \left(1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})) \right)$$

Finding optimum:

Optimization problems that we have seen so far:

- Linear Regression

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left(y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Logistic Regression

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log \left(1 + \exp(-y^{(i)} \mathbf{w}^\top \phi(x^{(i)})) \right)$$

Finding optimum:

Analytically computable optimum vs. gradient descent

The Problem more generally:

The Problem more generally:

$$\begin{array}{ll} \min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\} \end{array}$$

The Problem more generally:

$$\begin{array}{ll} \min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\} \end{array}$$

Solution:

The Problem more generally:

$$\begin{array}{ll} \min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\} \end{array}$$

Solution:

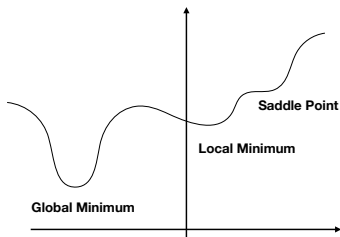
Solution \mathbf{w}^* has smallest value $f_0(\mathbf{w}^*)$ among all values that satisfy constraints

The Problem more generally:

$$\begin{array}{ll}\min_{\mathbf{w}} & f_0(\mathbf{w}) \\ \text{s.t.} & f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\}\end{array}$$

Solution:

Solution \mathbf{w}^* has smallest value $f_0(\mathbf{w}^*)$ among all values that satisfy constraints



Questions:

Questions:

- When can we find the optimum?

Questions:

- When can we find the optimum?
- Algorithms to search for the optimum?

Questions:

- When can we find the optimum?
- Algorithms to search for the optimum?
- How long does it take to find the optimum?

The pace of the class is?

The pace of the class is?

<https://www.strawpoll.me/19324883>



Questions:

Questions:

- When can we find the optimum?

Questions:

- When can we find the optimum?
- Algorithms to search for the optimum?

Questions:

- When can we find the optimum?
- Algorithms to search for the optimum?
- How long does it take to find the optimum?

When can we find the optimum?

When can we find the optimum?

When can we find the optimum?

- Least squares, linear and convex programs can be solved efficiently and reliably

When can we find the optimum?

- Least squares, linear and convex programs can be solved efficiently and reliably
- General optimization problems are very difficult to solve

When can we find the optimum?

- Least squares, linear and convex programs can be solved efficiently and reliably
- General optimization problems are very difficult to solve
- Often compromise between accuracy and computation time

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} \left(y^{(i)} - \phi(\mathbf{x}^{(i)})^\top \mathbf{w} \right)^2$$

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left(y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Linear program

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}} \left(y^{(i)} - \phi(\mathbf{x}^{(i)})^\top \mathbf{w} \right)^2$$

- Linear program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left(y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Linear program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

- Convex program

What's the form of 'least squares,' 'linear,' and 'convex' programs?

- Least squares program

$$\min_{\mathbf{w}} \frac{1}{2} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \left(y^{(i)} - \phi(x^{(i)})^\top \mathbf{w} \right)^2$$

- Linear program

$$\min_{\mathbf{w}} \mathbf{c}^\top \mathbf{w} \quad \text{s.t.} \quad \mathbf{A}\mathbf{w} \leq \mathbf{b}$$

- Convex program when all f_i **convex** (generalizes the above)

$$\min_{\mathbf{w}} f_0(\mathbf{w}) \quad \text{s.t.} \quad f_i(\mathbf{w}) \leq 0 \quad \forall i \in \{1, \dots, C\}$$

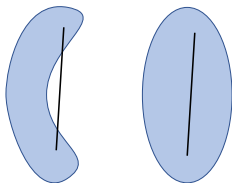
Convex set:

Convex set:

A set is convex if for any two points $\mathbf{w}_1, \mathbf{w}_2$ in the set, the line segment $\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2$ for $\lambda \in [0, 1]$ also lies in the set.

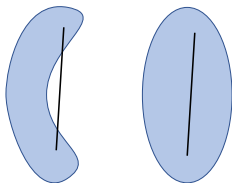
Convex set:

A set is convex if for any two points $\mathbf{w}_1, \mathbf{w}_2$ in the set, the line segment $\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2$ for $\lambda \in [0, 1]$ also lies in the set.



Convex set:

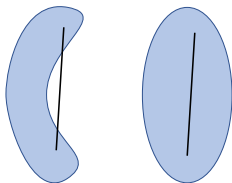
A set is convex if for any two points $\mathbf{w}_1, \mathbf{w}_2$ in the set, the line segment $\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2$ for $\lambda \in [0, 1]$ also lies in the set.



Example:

Convex set:

A set is convex if for any two points $\mathbf{w}_1, \mathbf{w}_2$ in the set, the line segment $\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2$ for $\lambda \in [0, 1]$ also lies in the set.



Example: Polyhedron

$$\{\mathbf{w} | \mathbf{Aw} \leq \mathbf{b}, \mathbf{Cw} = \mathbf{d}\}$$

Convex function

Convex function

A function f is convex if its domain is a convex set and for any points $\mathbf{w}_1, \mathbf{w}_2$ in the domain and any $\lambda \in [0, 1]$

Convex function

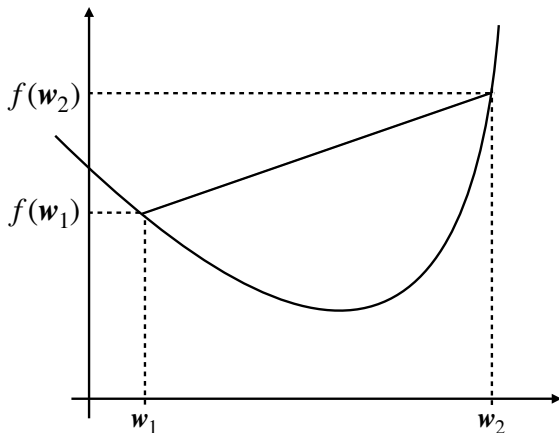
A function f is convex if its domain is a convex set and for any points $\mathbf{w}_1, \mathbf{w}_2$ in the domain and any $\lambda \in [0, 1]$

$$f((1 - \lambda)\mathbf{w}_1 + \lambda\mathbf{w}_2) \leq (1 - \lambda)f(\mathbf{w}_1) + \lambda f(\mathbf{w}_2)$$

Convex function

A function f is convex if its domain is a convex set and for any points $\mathbf{w}_1, \mathbf{w}_2$ in the domain and any $\lambda \in [0, 1]$

$$f((1 - \lambda)\mathbf{w}_1 + \lambda\mathbf{w}_2) \leq (1 - \lambda)f(\mathbf{w}_1) + \lambda f(\mathbf{w}_2)$$



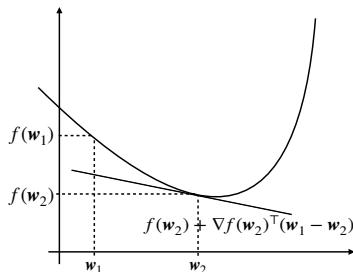
Recognizing convex functions

Recognizing convex functions

- If f is differentiable, then f is convex if and only if its domain is convex and $f(\mathbf{w}_1) \geq f(\mathbf{w}_2) + \nabla f(\mathbf{w}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2) \forall \mathbf{w}_1, \mathbf{w}_2$ in the domain

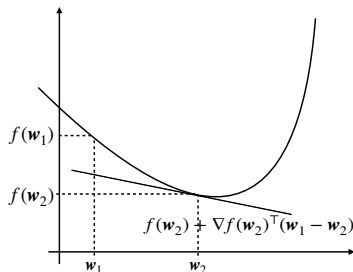
Recognizing convex functions

- If f is differentiable, then f is convex if and only if its domain is convex and $f(\mathbf{w}_1) \geq f(\mathbf{w}_2) + \nabla f(\mathbf{w}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2) \forall \mathbf{w}_1, \mathbf{w}_2$ in the domain



Recognizing convex functions

- If f is differentiable, then f is convex if and only if its domain is convex and $f(\mathbf{w}_1) \geq f(\mathbf{w}_2) + \nabla f(\mathbf{w}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2) \forall \mathbf{w}_1, \mathbf{w}_2$ in the domain

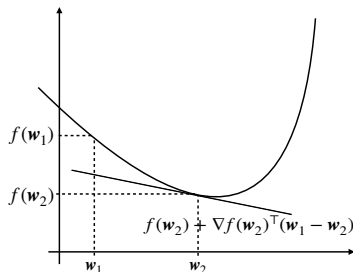


- If f is differentiable, then f is convex if and only if its domain is convex and $\forall \mathbf{w}_1, \mathbf{w}_2$ in the domain

$$(\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^\top (\mathbf{w}_1 - \mathbf{w}_2) \geq 0 \quad \text{monotone mapping}$$

Recognizing convex functions

- If f is differentiable, then f is convex if and only if its domain is convex and $f(\mathbf{w}_1) \geq f(\mathbf{w}_2) + \nabla f(\mathbf{w}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2) \quad \forall \mathbf{w}_1, \mathbf{w}_2$ in the domain



- If f is differentiable, then f is convex if and only if its domain is convex and $\forall \mathbf{w}_1, \mathbf{w}_2$ in the domain

$$(\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^\top (\mathbf{w}_1 - \mathbf{w}_2) \geq 0 \quad \text{monotone mapping}$$

- If f is twice differentiable, then f is convex if and only if its domain is convex and $\nabla^2 f(\mathbf{w}) \succeq 0 \quad \forall \mathbf{w}$ in the domain

Examples of convex functions

Examples of convex functions

- Exponential: $\exp(ax)$ convex on $x \in \mathbb{R} \forall a \in \mathbb{R}$

Examples of convex functions

- Exponential: $\exp(ax)$ convex on $x \in \mathbb{R} \forall a \in \mathbb{R}$
- Negative Logarithm: $-\log(x)$ is convex on $x \in \mathbb{R}_{++}$

Examples of convex functions

- Exponential: $\exp(ax)$ convex on $x \in \mathbb{R} \forall a \in \mathbb{R}$
- Negative Logarithm: $-\log(x)$ is convex on $x \in \mathbb{R}_{++}$
- Negative Entropy: $-H(x) = x \log(x)$ is convex on $x \in \mathbb{R}_{++}$

Examples of convex functions

- Exponential: $\exp(ax)$ convex on $x \in \mathbb{R} \forall a \in \mathbb{R}$
- Negative Logarithm: $-\log(x)$ is convex on $x \in \mathbb{R}_{++}$
- Negative Entropy: $-H(x) = x \log(x)$ is convex on $x \in \mathbb{R}_{++}$
- Norms: $\|\mathbf{w}\|_p$ for $p \geq 1$

Examples of convex functions

- Exponential: $\exp(ax)$ convex on $x \in \mathbb{R} \forall a \in \mathbb{R}$
- Negative Logarithm: $-\log(x)$ is convex on $x \in \mathbb{R}_{++}$
- Negative Entropy: $-H(x) = x \log(x)$ is convex on $x \in \mathbb{R}_{++}$
- Norms: $\|\mathbf{w}\|_p$ for $p \geq 1$
- Log-Sum-Exp: $\log(\exp(w_1) + \dots + \exp(w_d))$

Operations which preserve convexity

Operations which preserve convexity

- Non-negative weighted sums: $\alpha_i \geq 0$; if f_i convex $\forall i$, so is

$$g = \alpha_1 f_1 + \alpha_2 f_2 + \dots$$

Operations which preserve convexity

- Non-negative weighted sums: $\alpha_i \geq 0$; if f_i convex $\forall i$, so is

$$g = \alpha_1 f_1 + \alpha_2 f_2 + \dots$$

- Composition with an affine mapping: if f is convex, so is

$$g(\mathbf{w}) = f(\mathbf{A}\mathbf{w} + \mathbf{b})$$

Operations which preserve convexity

- Non-negative weighted sums: $\alpha_i \geq 0$; if f_i convex $\forall i$, so is

$$g = \alpha_1 f_1 + \alpha_2 f_2 + \dots$$

- Composition with an affine mapping: if f is convex, so is

$$g(\mathbf{w}) = f(\mathbf{Aw} + \mathbf{b})$$

- Pointwise maximum: if f_1, f_2 are convex, so is

$$g(\mathbf{w}) = \max\{f_1(\mathbf{w}), f_2(\mathbf{w})\}$$

Operations which preserve convexity

- Non-negative weighted sums: $\alpha_i \geq 0$; if f_i convex $\forall i$, so is

$$g = \alpha_1 f_1 + \alpha_2 f_2 + \dots$$

- Composition with an affine mapping: if f is convex, so is

$$g(\mathbf{w}) = f(\mathbf{A}\mathbf{w} + \mathbf{b})$$

- Pointwise maximum: if f_1, f_2 are convex, so is

$$g(\mathbf{w}) = \max\{f_1(\mathbf{w}), f_2(\mathbf{w})\}$$

Show that $\log(1 + \exp(x))$ is convex for $x \in \mathbb{R}$

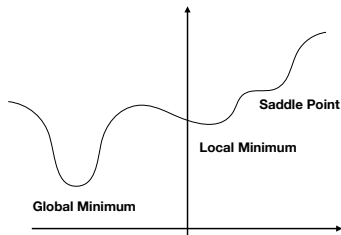
Optimality of convex optimization

Optimality of convex optimization

- A point \mathbf{w}^* is locally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$ in a neighborhood of \mathbf{w}^* ; globally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$

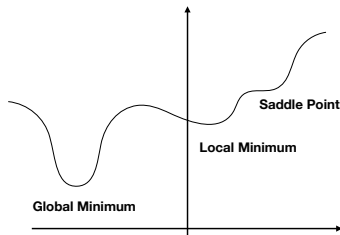
Optimality of convex optimization

- A point \mathbf{w}^* is locally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$ in a neighborhood of \mathbf{w}^* ; globally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$



Optimality of convex optimization

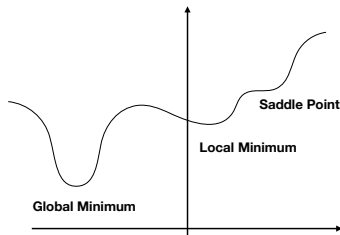
- A point \mathbf{w}^* is locally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$ in a neighborhood of \mathbf{w}^* ; globally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$



For convex problems global optimality follows directly from local optimality.

Optimality of convex optimization

- A point \mathbf{w}^* is locally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$ in a neighborhood of \mathbf{w}^* ; globally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$

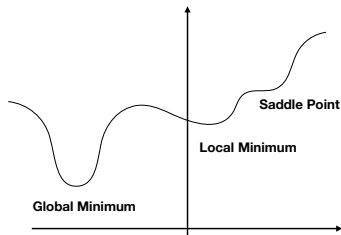


For convex problems global optimality follows directly from local optimality.

- For a local minimum of f , $\nabla f(\mathbf{w}^*) = 0$

Optimality of convex optimization

- A point \mathbf{w}^* is locally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$ in a neighborhood of \mathbf{w}^* ; globally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$

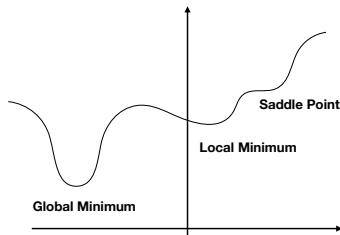


For convex problems global optimality follows directly from local optimality.

- For a local minimum of f , $\nabla f(\mathbf{w}^*) = 0$
- If f convex, then $\nabla f(\mathbf{w}^*) = 0$ sufficient for global optimality

Optimality of convex optimization

- A point \mathbf{w}^* is locally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$ in a neighborhood of \mathbf{w}^* ; globally optimal if $f(\mathbf{w}^*) \leq f(\mathbf{w}) \forall \mathbf{w}$



For convex problems global optimality follows directly from local optimality.

- For a local minimum of f , $\nabla f(\mathbf{w}^*) = 0$
- If f convex, then $\nabla f(\mathbf{w}^*) = 0$ sufficient for global optimality

This makes convex optimization special!

Algorithms to search for the optimum?

Descent methods

$$\min_{\mathbf{w}} f(\mathbf{w})$$

Intuition

Descent methods

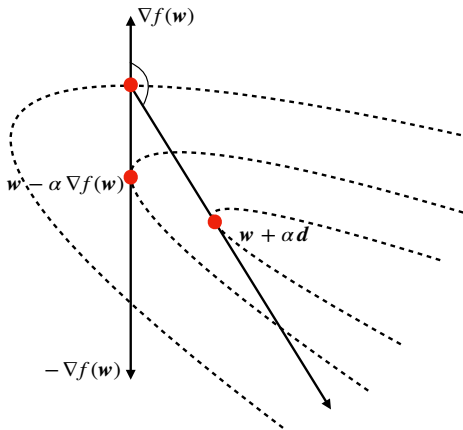
$$\min_{\mathbf{w}} f(\mathbf{w})$$

Intuition (find a stationary point with $\nabla f(\mathbf{w}) = 0$)

Descent methods

$$\min_{\mathbf{w}} f(\mathbf{w})$$

Intuition (find a stationary point with $\nabla f(\mathbf{w}) = 0$)



Iterative algorithm

Iterative algorithm

- Start with some guess \mathbf{w}

Iterative algorithm

- Start with some guess \mathbf{w}
- Iterate $k = 1, 2, 3, \dots$

Iterative algorithm

- Start with some guess \mathbf{w}
- Iterate $k = 1, 2, 3, \dots$
 - ▶ Select direction \mathbf{d}_k and stepsize α_k

Iterative algorithm

- Start with some guess \mathbf{w}
- Iterate $k = 1, 2, 3, \dots$
 - ▶ Select direction \mathbf{d}_k and stepsize α_k
 - ▶ $\mathbf{w} \leftarrow \mathbf{w} + \alpha_k \mathbf{d}_k$

Iterative algorithm

- Start with some guess \mathbf{w}
- Iterate $k = 1, 2, 3, \dots$
 - ▶ Select direction \mathbf{d}_k and stepsize α_k
 - ▶ $\mathbf{w} \leftarrow \mathbf{w} + \alpha_k \mathbf{d}_k$
 - ▶ Check whether we should stop (e.g., if $\nabla f(\mathbf{w}) \approx 0$)

Iterative algorithm

- Start with some guess \mathbf{w}
- Iterate $k = 1, 2, 3, \dots$
 - ▶ Select direction \mathbf{d}_k and stepsize α_k
 - ▶ $\mathbf{w} \leftarrow \mathbf{w} + \alpha_k \mathbf{d}_k$
 - ▶ Check whether we should stop (e.g., if $\nabla f(\mathbf{w}) \approx 0$)

Descent direction \mathbf{d}_k satisfies

Iterative algorithm

- Start with some guess \mathbf{w}
- Iterate $k = 1, 2, 3, \dots$
 - ▶ Select direction \mathbf{d}_k and stepsize α_k
 - ▶ $\mathbf{w} \leftarrow \mathbf{w} + \alpha_k \mathbf{d}_k$
 - ▶ Check whether we should stop (e.g., if $\nabla f(\mathbf{w}) \approx 0$)

Descent direction \mathbf{d}_k satisfies $\nabla f(\mathbf{w})^\top \mathbf{d}_k < 0$

How to select direction:

How to select direction:

- Steepest descent: $\mathbf{d}_k = -\nabla f(\mathbf{w}_k)$

How to select direction:

- Steepest descent: $\mathbf{d}_k = -\nabla f(\mathbf{w}_k)$
- Scaled gradient: $\mathbf{d}_k = -\mathbf{D}_k \nabla f(\mathbf{w}_k)$ for $\mathbf{D}_k \succ 0$

How to select direction:

- Steepest descent: $\mathbf{d}_k = -\nabla f(\mathbf{w}_k)$
- Scaled gradient: $\mathbf{d}_k = -\mathbf{D}_k \nabla f(\mathbf{w}_k)$ for $\mathbf{D}_k \succ 0$
 - ▶ E.g., Newton's method: $\mathbf{D}_k = [\nabla^2 f(\mathbf{w}_k)]^{-1}$

How to select direction:

- Steepest descent: $\mathbf{d}_k = -\nabla f(\mathbf{w}_k)$
- Scaled gradient: $\mathbf{d}_k = -\mathbf{D}_k \nabla f(\mathbf{w}_k)$ for $\mathbf{D}_k \succ 0$
 - ▶ E.g., Newton's method: $\mathbf{D}_k = [\nabla^2 f(\mathbf{w}_k)]^{-1}$
- ...

How to select stepsize:

How to select stepsize:

- Exact: $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$

How to select stepsize:

- Exact: $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$
- Constant: $\alpha_k = 1/L$ (for suitable L)

How to select stepsize:

- Exact: $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$
- Constant: $\alpha_k = 1/L$ (for suitable L)
- Diminishing: $\alpha_k \rightarrow 0$ but $\sum_k \alpha_k = \infty$ (e.g., $\alpha_k = 1/k$)

How to select stepsize:

- Exact: $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$
- Constant: $\alpha_k = 1/L$ (for suitable L)
- Diminishing: $\alpha_k \rightarrow 0$ but $\sum_k \alpha_k = \infty$ (e.g., $\alpha_k = 1/k$)
- Armijo Rule:

How to select stepsize:

- Exact: $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$
- Constant: $\alpha_k = 1/L$ (for suitable L)
- Diminishing: $\alpha_k \rightarrow 0$ but $\sum_k \alpha_k = \infty$ (e.g., $\alpha_k = 1/k$)
- Armijo Rule:

Start with $\alpha = s$ and continue with $\alpha = \beta s, \alpha = \beta^2 s, \dots$, until $\alpha = \beta^m s$ falls within the set of α with

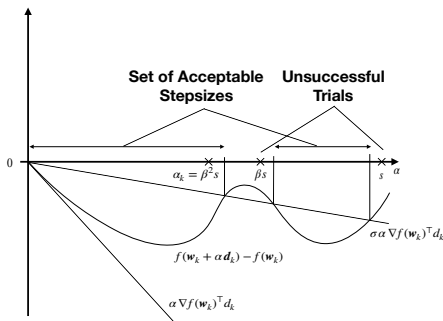
$$f(\mathbf{w}_k + \alpha \mathbf{d}_k) - f(\mathbf{w}_k) \leq \sigma \alpha \nabla f(\mathbf{w}_k)^\top \mathbf{d}_k$$

How to select stepsize:

- Exact: $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{w}_k + \alpha \mathbf{d}_k)$
- Constant: $\alpha_k = 1/L$ (for suitable L)
- Diminishing: $\alpha_k \rightarrow 0$ but $\sum_k \alpha_k = \infty$ (e.g., $\alpha_k = 1/k$)
- Armijo Rule:

Start with $\alpha = s$ and continue with $\alpha = \beta s$, $\alpha = \beta^2 s$, \dots , until $\alpha = \beta^m s$ falls within the set of α with

$$f(\mathbf{w}_k + \alpha \mathbf{d}_k) - f(\mathbf{w}_k) \leq \sigma \alpha \nabla f(\mathbf{w}_k)^\top \mathbf{d}_k$$



How long does it take to find the optimum?

Goal:

How many iterations k for

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon$$

Two important properties:

Two important properties:

- Lipschitz continuous gradient

Two important properties:

- Lipschitz continuous gradient
- Strong convexity

Properties: Lipschitz continuous gradient

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \end{aligned}$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \end{aligned}$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \\ & \leq \end{aligned}$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \\ & \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \end{aligned}$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \\ & \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \end{aligned}$$

$$\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2) =$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \\ & \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \end{aligned}$$

$$\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2) = L(\mathbf{w}_1 - \mathbf{w}_2) - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \\ & \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \end{aligned}$$

$$\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2) = L(\mathbf{w}_1 - \mathbf{w}_2) - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))$$

$$\begin{aligned} & (\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & = \end{aligned}$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \\ & \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \end{aligned}$$

$$\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2) = L(\mathbf{w}_1 - \mathbf{w}_2) - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))$$

$$\begin{aligned} & (\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & = L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \end{aligned}$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \\ & \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \end{aligned}$$

$$\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2) = L(\mathbf{w}_1 - \mathbf{w}_2) - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))$$

$$\begin{aligned} & (\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & = L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \geq \end{aligned}$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \\ & \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \end{aligned}$$

$$\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2) = L(\mathbf{w}_1 - \mathbf{w}_2) - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))$$

$$\begin{aligned} & (\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & = L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \geq 0 \end{aligned}$$

Properties: Lipschitz continuous gradient

$$\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Intuition:

Lipschitz continuous gradient, then $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex

Proof:

$$\begin{aligned} & (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \leq \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \text{Cauchy-Schwartz} \\ & \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \end{aligned}$$

$$\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2) = L(\mathbf{w}_1 - \mathbf{w}_2) - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))$$

$$\begin{aligned} & (\nabla g(\mathbf{w}_1) - \nabla g(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & = L\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 - (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^{\top}(\mathbf{w}_1 - \mathbf{w}_2) \\ & \geq 0 \quad \text{monotone mapping} \end{aligned}$$

If $g(\mathbf{w}) = \frac{L}{2} \|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex, then

If $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex, then

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{L}{2}\|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

If $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex, then

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{L}{2}\|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Proof:

If $g(\mathbf{w}) = \frac{L}{2} \|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex, then

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{L}{2} \|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Proof: plug definition of g into $g(\mathbf{w}_2) \geq g(\mathbf{w}_1) + \nabla g(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1)$ and re-arrange

If $g(\mathbf{w}) = \frac{L}{2} \|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex, then

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{L}{2} \|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Proof: plug definition of g into $g(\mathbf{w}_2) \geq g(\mathbf{w}_1) + \nabla g(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1)$
and re-arrange

Properties: Strong convexity

If $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex, then

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{L}{2}\|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Proof: plug definition of g into $g(\mathbf{w}_2) \geq g(\mathbf{w}_1) + \nabla g(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1)$ and re-arrange

Properties: Strong convexity

$$f(\mathbf{w}_2) \geq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{\sigma}{2}\|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

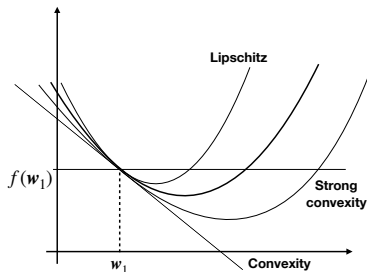
If $g(\mathbf{w}) = \frac{L}{2} \|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex, then

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{L}{2} \|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Proof: plug definition of g into $g(\mathbf{w}_2) \geq g(\mathbf{w}_1) + \nabla g(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1)$ and re-arrange

Properties: Strong convexity

$$f(\mathbf{w}_2) \geq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{\sigma}{2} \|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$



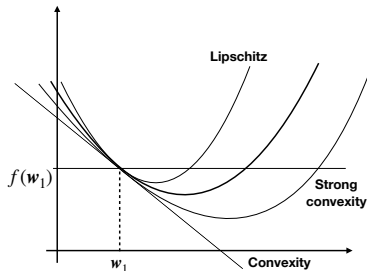
If $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex, then

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{L}{2}\|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Proof: plug definition of g into $g(\mathbf{w}_2) \geq g(\mathbf{w}_1) + \nabla g(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1)$ and re-arrange

Properties: Strong convexity

$$f(\mathbf{w}_2) \geq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{\sigma}{2}\|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$



if f twice differentiable

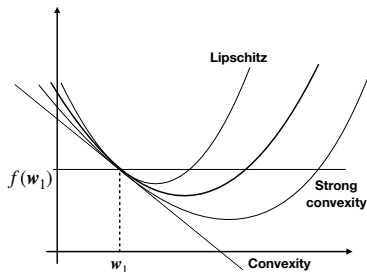
If $g(\mathbf{w}) = \frac{L}{2}\|\mathbf{w}\|_2^2 - f(\mathbf{w})$ convex, then

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{L}{2}\|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$

Proof: plug definition of g into $g(\mathbf{w}_2) \geq g(\mathbf{w}_1) + \nabla g(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1)$ and re-arrange

Properties: Strong convexity

$$f(\mathbf{w}_2) \geq f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}_2 - \mathbf{w}_1) + \frac{\sigma}{2}\|\mathbf{w}_2 - \mathbf{w}_1\|_2^2 \quad \forall \mathbf{w}_1, \mathbf{w}_2$$



if f twice differentiable

$$\sigma I \prec \nabla^2 f(\mathbf{w}) \prec LI \quad \forall \mathbf{w}$$

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$?

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$$

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$$

Hence

$$\alpha \mathbf{d}_k =$$

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$$

Hence

$$\alpha \mathbf{d}_k = -\frac{1}{L} \nabla f(\mathbf{w}_k)$$

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$$

Hence

$$\begin{aligned} \alpha \mathbf{d}_k &= -\frac{1}{L} \nabla f(\mathbf{w}_k) \\ f(\mathbf{w}_{k+1}) &\leq \end{aligned}$$

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$$

Hence

$$\begin{aligned} \alpha \mathbf{d}_k &= -\frac{1}{L} \nabla f(\mathbf{w}_k) \\ f(\mathbf{w}_{k+1}) &\leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2 \end{aligned}$$

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$$

Hence

$$\begin{aligned} \alpha \mathbf{d}_k &= -\frac{1}{L} \nabla f(\mathbf{w}_k) \\ f(\mathbf{w}_{k+1}) &\leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2 \quad \text{Bound on guaranteed progress} \end{aligned}$$

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$$

Hence

$$\begin{aligned} \alpha \mathbf{d}_k &= -\frac{1}{L} \nabla f(\mathbf{w}_k) \\ f(\mathbf{w}_{k+1}) &\leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2 \quad \text{Bound on guaranteed progress} \end{aligned}$$

Bound on sub-optimality from strong convexity:

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$$

Hence

$$\begin{aligned} \alpha \mathbf{d}_k &= -\frac{1}{L} \nabla f(\mathbf{w}_k) \\ f(\mathbf{w}_{k+1}) &\leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2 \quad \text{Bound on guaranteed progress} \end{aligned}$$

Bound on sub-optimality from strong convexity:

$$f(\mathbf{w}^*) \geq$$

How many iterations k such that

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \epsilon \quad \text{for} \quad \mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \mathbf{d}_k$$

How to pick $\alpha \mathbf{d}_k$? Minimize w.r.t. \mathbf{w}_{k+1} right-hand-side of upper bound

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$$

Hence

$$\begin{aligned} \alpha \mathbf{d}_k &= -\frac{1}{L} \nabla f(\mathbf{w}_k) \\ f(\mathbf{w}_{k+1}) &\leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2 \quad \text{Bound on guaranteed progress} \end{aligned}$$

Bound on sub-optimality from strong convexity:

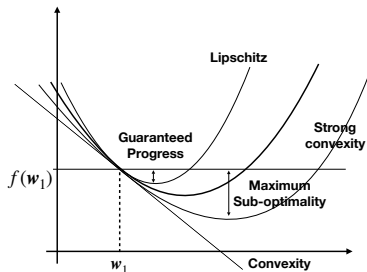
$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$



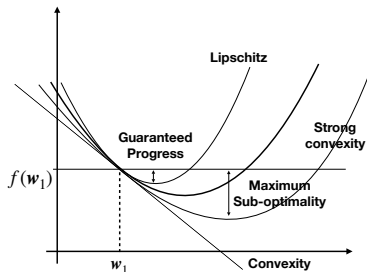
Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Distance to go:

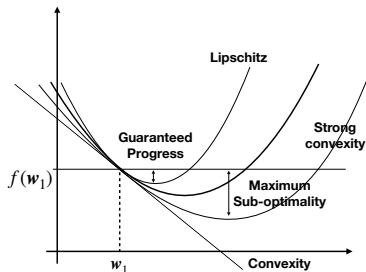


Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$



Distance to go:

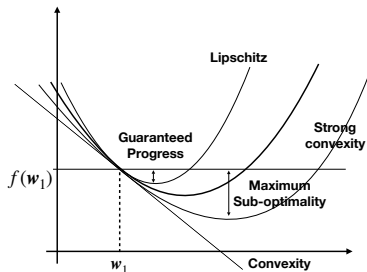
$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$



Distance to go:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

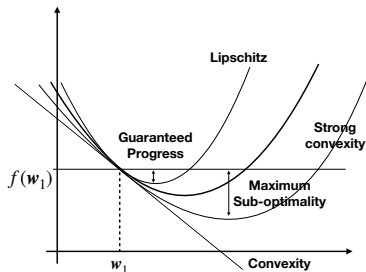
in 'guaranteed progress':

Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$



Distance to go:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

in 'guaranteed progress':

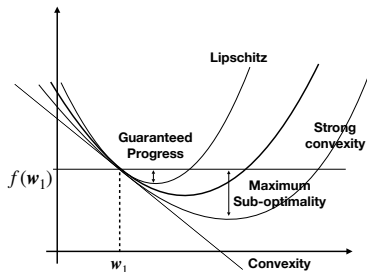
$$\begin{aligned} f(\mathbf{w}_k) - f(\mathbf{w}^*) &\leq f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*) - \frac{\sigma}{L} (f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*)) \\ &\leq \end{aligned}$$

Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$



Distance to go:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

in 'guaranteed progress':

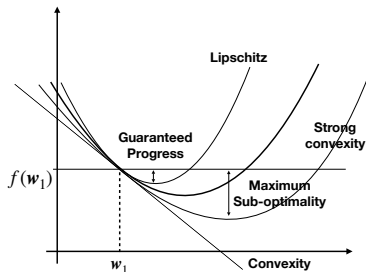
$$\begin{aligned} f(\mathbf{w}_k) - f(\mathbf{w}^*) &\leq f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*) - \frac{\sigma}{L} (f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*)) \\ &\leq \left(1 - \frac{\sigma}{L}\right)^k (f(\mathbf{w}_0) - f(\mathbf{w}^*)) \quad (\sigma < L) \end{aligned}$$

Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$



Distance to go:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

in 'guaranteed progress':

$$\begin{aligned} f(\mathbf{w}_k) - f(\mathbf{w}^*) &\leq f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*) - \frac{\sigma}{L} (f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*)) \\ &\leq \left(1 - \frac{\sigma}{L}\right)^k (f(\mathbf{w}_0) - f(\mathbf{w}^*)) \quad (\sigma < L) \end{aligned}$$

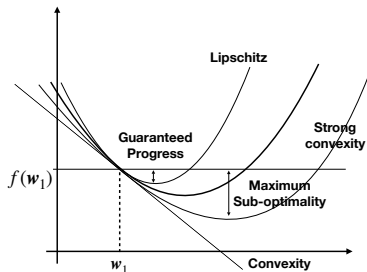
Rate:

Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$



Distance to go:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

in 'guaranteed progress':

$$\begin{aligned} f(\mathbf{w}_k) - f(\mathbf{w}^*) &\leq f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*) - \frac{\sigma}{L} (f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*)) \\ &\leq \left(1 - \frac{\sigma}{L}\right)^k (f(\mathbf{w}_0) - f(\mathbf{w}^*)) \quad (\sigma < L) \end{aligned}$$

Rate:

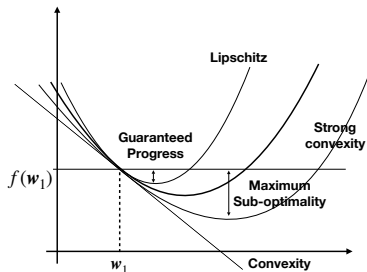
$$c \left(1 - \frac{\sigma}{L}\right)^k \leq \epsilon \quad \implies$$

Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$



Distance to go:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

in 'guaranteed progress':

$$\begin{aligned} f(\mathbf{w}_k) - f(\mathbf{w}^*) &\leq f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*) - \frac{\sigma}{L} (f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*)) \\ &\leq \left(1 - \frac{\sigma}{L}\right)^k (f(\mathbf{w}_0) - f(\mathbf{w}^*)) \quad (\sigma < L) \end{aligned}$$

Rate:

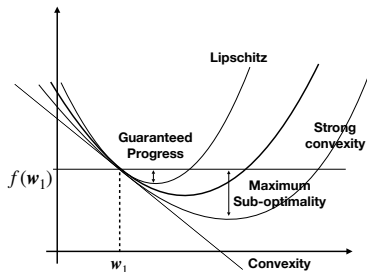
$$c \left(1 - \frac{\sigma}{L}\right)^k \leq \epsilon \implies k \geq O(\log(1/\epsilon))$$

Guaranteed progress:

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|_2^2$$

Maximum sub-optimality:

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) - \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$



Distance to go:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{2\sigma} \|\nabla f(\mathbf{w}_k)\|_2^2$$

in 'guaranteed progress':

$$\begin{aligned} f(\mathbf{w}_k) - f(\mathbf{w}^*) &\leq f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*) - \frac{\sigma}{L} (f(\mathbf{w}_{k-1}) - f(\mathbf{w}^*)) \\ &\leq \left(1 - \frac{\sigma}{L}\right)^k (f(\mathbf{w}_0) - f(\mathbf{w}^*)) \quad (\sigma < L) \end{aligned}$$

Rate:

$$c \left(1 - \frac{\sigma}{L}\right)^k \leq \epsilon \implies k \geq O(\log(1/\epsilon)) \quad (\text{sometimes } O(e^k))$$

No strong convexity assumption:

No strong convexity assumption:

Lipschitz bound and $\mathbf{w}_2 = \mathbf{w}_1 - \alpha \nabla f(\mathbf{w}_1)$ yields

No strong convexity assumption:

Lipschitz bound and $\mathbf{w}_2 = \mathbf{w}_1 - \alpha \nabla f(\mathbf{w}_1)$ yields

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) - (1 - \frac{L\alpha}{2})\alpha \|\nabla f(\mathbf{w}_1)\|_2^2$$

No strong convexity assumption:

Lipschitz bound and $\mathbf{w}_2 = \mathbf{w}_1 - \alpha \nabla f(\mathbf{w}_1)$ yields

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) - (1 - \frac{L\alpha}{2})\alpha \|\nabla f(\mathbf{w}_1)\|_2^2$$

Combined with convexity: $f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}^* - \mathbf{w}_1) \leq f(\mathbf{w}^*)$

No strong convexity assumption:

Lipschitz bound and $\mathbf{w}_2 = \mathbf{w}_1 - \alpha \nabla f(\mathbf{w}_1)$ yields

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) - (1 - \frac{L\alpha}{2})\alpha \|\nabla f(\mathbf{w}_1)\|_2^2$$

Combined with convexity: $f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}^* - \mathbf{w}_1) \leq f(\mathbf{w}^*)$

$$f(\mathbf{w}_2) \leq f(\mathbf{w}^*) + \nabla f(\mathbf{w}_1)(\mathbf{w}_1 - \mathbf{w}^*) - \frac{\alpha}{2} \|\nabla f(\mathbf{w}_1)\|_2^2$$

No strong convexity assumption:

Lipschitz bound and $\mathbf{w}_2 = \mathbf{w}_1 - \alpha \nabla f(\mathbf{w}_1)$ yields

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) - (1 - \frac{L\alpha}{2})\alpha \|\nabla f(\mathbf{w}_1)\|_2^2$$

Combined with convexity: $f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}^* - \mathbf{w}_1) \leq f(\mathbf{w}^*)$

$$f(\mathbf{w}_2) \leq f(\mathbf{w}^*) + \nabla f(\mathbf{w}_1)(\mathbf{w}_1 - \mathbf{w}^*) - \frac{\alpha}{2} \|\nabla f(\mathbf{w}_1)\|_2^2$$

Using $\mathbf{w}_2 - \mathbf{w}_1 = -\alpha \nabla f(\mathbf{w}_1)$ and rearranging terms gives

No strong convexity assumption:

Lipschitz bound and $\mathbf{w}_2 = \mathbf{w}_1 - \alpha \nabla f(\mathbf{w}_1)$ yields

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) - \left(1 - \frac{L\alpha}{2}\right)\alpha \|\nabla f(\mathbf{w}_1)\|_2^2$$

Combined with convexity: $f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}^* - \mathbf{w}_1) \leq f(\mathbf{w}^*)$

$$f(\mathbf{w}_2) \leq f(\mathbf{w}^*) + \nabla f(\mathbf{w}_1)(\mathbf{w}_1 - \mathbf{w}^*) - \frac{\alpha}{2} \|\nabla f(\mathbf{w}_1)\|_2^2$$

Using $\mathbf{w}_2 - \mathbf{w}_1 = -\alpha \nabla f(\mathbf{w}_1)$ and rearranging terms gives

$$f(\mathbf{w}_2) \leq f(\mathbf{w}^*) + \frac{1}{2\alpha} \left(\|\mathbf{w}_1 - \mathbf{w}^*\|_2^2 - \|\mathbf{w}_2 - \mathbf{w}^*\|_2^2 \right)$$

No strong convexity assumption:

Lipschitz bound and $\mathbf{w}_2 = \mathbf{w}_1 - \alpha \nabla f(\mathbf{w}_1)$ yields

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) - (1 - \frac{L\alpha}{2})\alpha \|\nabla f(\mathbf{w}_1)\|_2^2$$

Combined with convexity: $f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}^* - \mathbf{w}_1) \leq f(\mathbf{w}^*)$

$$f(\mathbf{w}_2) \leq f(\mathbf{w}^*) + \nabla f(\mathbf{w}_1)(\mathbf{w}_1 - \mathbf{w}^*) - \frac{\alpha}{2} \|\nabla f(\mathbf{w}_1)\|_2^2$$

Using $\mathbf{w}_2 - \mathbf{w}_1 = -\alpha \nabla f(\mathbf{w}_1)$ and rearranging terms gives

$$f(\mathbf{w}_2) \leq f(\mathbf{w}^*) + \frac{1}{2\alpha} \left(\|\mathbf{w}_1 - \mathbf{w}^*\|_2^2 - \|\mathbf{w}_2 - \mathbf{w}^*\|_2^2 \right)$$

Summing over all iterations

$$\sum_{i=1}^k (f(\mathbf{w}_i) - f(\mathbf{w}^*)) \leq$$

No strong convexity assumption:

Lipschitz bound and $\mathbf{w}_2 = \mathbf{w}_1 - \alpha \nabla f(\mathbf{w}_1)$ yields

$$f(\mathbf{w}_2) \leq f(\mathbf{w}_1) - \left(1 - \frac{L\alpha}{2}\right)\alpha \|\nabla f(\mathbf{w}_1)\|_2^2$$

Combined with convexity: $f(\mathbf{w}_1) + \nabla f(\mathbf{w}_1)^\top (\mathbf{w}^* - \mathbf{w}_1) \leq f(\mathbf{w}^*)$

$$f(\mathbf{w}_2) \leq f(\mathbf{w}^*) + \nabla f(\mathbf{w}_1)(\mathbf{w}_1 - \mathbf{w}^*) - \frac{\alpha}{2} \|\nabla f(\mathbf{w}_1)\|_2^2$$

Using $\mathbf{w}_2 - \mathbf{w}_1 = -\alpha \nabla f(\mathbf{w}_1)$ and rearranging terms gives

$$f(\mathbf{w}_2) \leq f(\mathbf{w}^*) + \frac{1}{2\alpha} \left(\|\mathbf{w}_1 - \mathbf{w}^*\|_2^2 - \|\mathbf{w}_2 - \mathbf{w}^*\|_2^2 \right)$$

Summing over all iterations

$$\sum_{i=1}^k (f(\mathbf{w}_i) - f(\mathbf{w}^*)) \leq \frac{1}{2\alpha} \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2$$

$f(\mathbf{w}_i)$ non-increasing:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{k} \sum_{i=1}^k (f(\mathbf{w}_i) - f(\mathbf{w}^*)) \leq \frac{1}{2k\alpha} \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2 \leq \epsilon$$

Consequently:

$f(\mathbf{w}_i)$ non-increasing:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{k} \sum_{i=1}^k (f(\mathbf{w}_i) - f(\mathbf{w}^*)) \leq \frac{1}{2k\alpha} \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2 \leq \epsilon$$

Consequently:

$$k \geq O(1/\epsilon)$$

$f(\mathbf{w}_i)$ non-increasing:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{1}{k} \sum_{i=1}^k (f(\mathbf{w}_i) - f(\mathbf{w}^*)) \leq \frac{1}{2k\alpha} \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2 \leq \epsilon$$

Consequently:

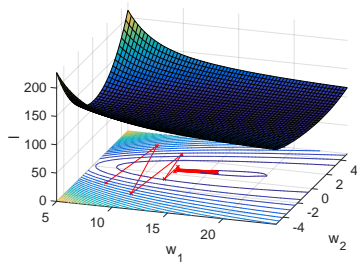
$$k \geq O(1/\epsilon)$$

Are these rates optimal?

Gradient with momentum

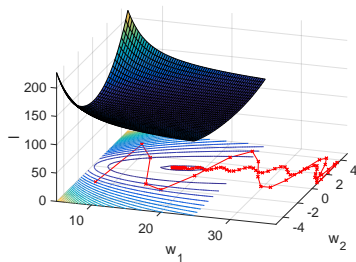
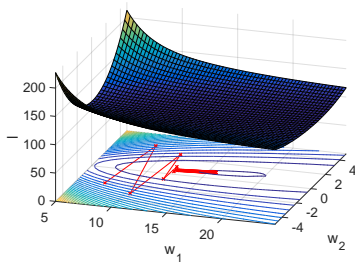
Gradient with momentum

Intuition:



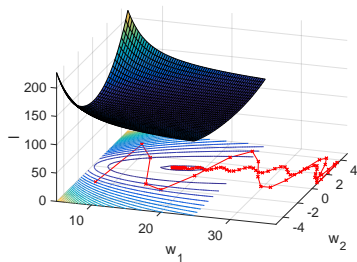
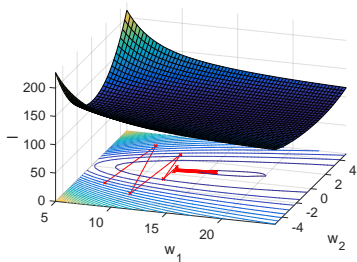
Gradient with momentum

Intuition:



Gradient with momentum

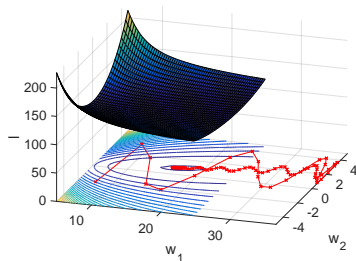
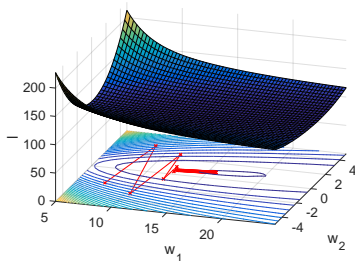
Intuition:



Video

Gradient with momentum

Intuition:

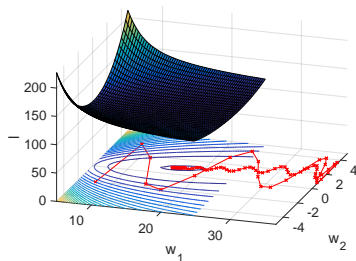
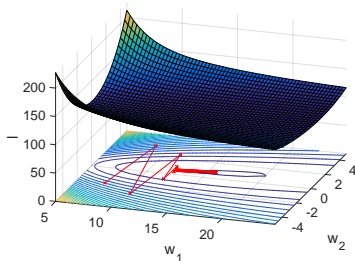


Video

- Polyak's method (aka heavy-ball)

Gradient with momentum

Intuition:



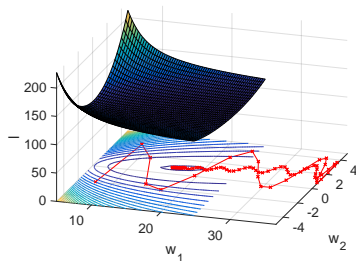
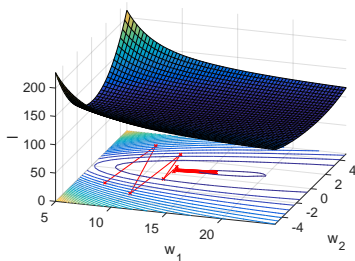
Video

- Polyak's method (aka heavy-ball)

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \beta_k (\mathbf{w}_k - \mathbf{w}_{k-1})$$

Gradient with momentum

Intuition:



Video

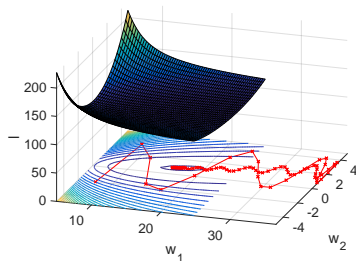
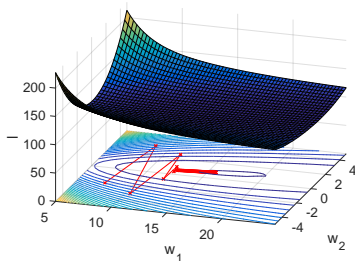
- Polyak's method (aka heavy-ball)

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \beta_k (\mathbf{w}_k - \mathbf{w}_{k-1})$$

- Momentum method in deep learning

Gradient with momentum

Intuition:



Video

- Polyak's method (aka heavy-ball)

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \beta_k (\mathbf{w}_k - \mathbf{w}_{k-1})$$

- Momentum method in deep learning

$$\mathbf{v}_{k+1} = \beta \mathbf{v}_k + \nabla f(\mathbf{w}_k)$$

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{v}_{k+1}$$

Recall the structure of our optimization problems:

Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

- So far we didn't consider the time for computing the gradient

Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

- So far we didn't consider the time for computing the gradient
- Iteration complexity is linear in the number of samples $|\mathcal{D}|$

Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

- So far we didn't consider the time for computing the gradient
- Iteration complexity is linear in the number of samples $|\mathcal{D}|$
- A large dataset makes gradient computation slow

Recall the structure of our optimization problems:

$$\min_{\mathbf{w}} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \ell(y_i, F(x^{(i)}, \mathbf{w}))$$

- So far we didn't consider the time for computing the gradient
- Iteration complexity is linear in the number of samples $|\mathcal{D}|$
- A large dataset makes gradient computation slow

How to deal with this?

Stochastic gradient descent

Stochastic gradient descent

Consider a subset of samples and approximate the gradient based on this batch of data.

Stochastic gradient descent

Consider a subset of samples and approximate the gradient based on this batch of data.

- Select a subset of samples \mathcal{B}_k

Stochastic gradient descent

Consider a subset of samples and approximate the gradient based on this batch of data.

- Select a subset of samples \mathcal{B}_k
- Gradient update using approximation

$$\nabla f(\mathbf{w}) \approx \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{B}_k} \nabla \ell(y^{(i)}, F(x^{(i)}, \mathbf{w}))$$

Stochastic gradient descent

Consider a subset of samples and approximate the gradient based on this batch of data.

- Select a subset of samples \mathcal{B}_k
- Gradient update using approximation

$$\nabla f(\mathbf{w}) \approx \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{B}_k} \nabla \ell(y^{(i)}, F(x^{(i)}, \mathbf{w}))$$

Convergence rates for stochastic gradient descent:

Stochastic gradient descent

Consider a subset of samples and approximate the gradient based on this batch of data.

- Select a subset of samples \mathcal{B}_k
- Gradient update using approximation

$$\nabla f(\mathbf{w}) \approx \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{B}_k} \nabla \ell(y^{(i)}, F(x^{(i)}, \mathbf{w}))$$

Convergence rates for stochastic gradient descent:

- Lipschitz continuous gradient and strongly convex: $k \geq O(1/\epsilon)$

Stochastic gradient descent

Consider a subset of samples and approximate the gradient based on this batch of data.

- Select a subset of samples \mathcal{B}_k
- Gradient update using approximation

$$\nabla f(\mathbf{w}) \approx \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{B}_k} \nabla \ell(y^{(i)}, F(x^{(i)}, \mathbf{w}))$$

Convergence rates for stochastic gradient descent:

- Lipschitz continuous gradient and strongly convex: $k \geq O(1/\epsilon)$
- Lipschitz continuous gradient: $k \geq O(1/\epsilon^2)$

Stochastic vs. deterministic (strongly convex)

Batch gradient descent:

- Convergence rate:
- Iteration complexity:

Stochastic gradient descent:

- Convergence rate:
- Iteration complexity:

Stochastic vs. deterministic (strongly convex)

Batch gradient descent:

- Convergence rate: $O(\log 1/\epsilon)$
- Iteration complexity:

Stochastic gradient descent:

- Convergence rate:
- Iteration complexity:

Stochastic vs. deterministic (strongly convex)

Batch gradient descent:

- Convergence rate: $O(\log 1/\epsilon)$
- Iteration complexity: linear in $|\mathcal{D}|$

Stochastic gradient descent:

- Convergence rate:
- Iteration complexity:

Stochastic vs. deterministic (strongly convex)

Batch gradient descent:

- Convergence rate: $O(\log 1/\epsilon)$
- Iteration complexity: linear in $|\mathcal{D}|$

Stochastic gradient descent:

- Convergence rate: $O(1/\epsilon)$
- Iteration complexity:

Stochastic vs. deterministic (strongly convex)

Batch gradient descent:

- Convergence rate: $O(\log 1/\epsilon)$
- Iteration complexity: linear in $|\mathcal{D}|$

Stochastic gradient descent:

- Convergence rate: $O(1/\epsilon)$
- Iteration complexity: independent of $|\mathcal{D}|$

Stochastic vs. deterministic (strongly convex)

Batch gradient descent:

- Convergence rate: $O(\log 1/\epsilon)$
- Iteration complexity: linear in $|\mathcal{D}|$

Stochastic gradient descent:

- Convergence rate: $O(1/\epsilon)$
- Iteration complexity: independent of $|\mathcal{D}|$

Can we get the best of both worlds?

Many related algorithms:

- SAG (Le Roux, Schmidt, Bach 2012)
- SDCA (Shalev-Shwartz and Zhang 2013)
- SVRG (Johnson and Zhang 2013)
- MISO (Mairal 2015)
- Finito (Defazio 2014)
- SAGA (Defazio, Bach, Lacoste-Julien 2014)
- ...

Many related algorithms:

- SAG (Le Roux, Schmidt, Bach 2012)
- SDCA (Shalev-Shwartz and Zhang 2013)
- SVRG (Johnson and Zhang 2013)
- MISO (Mairal 2015)
- Finito (Defazio 2014)
- SAGA (Defazio, Bach, Lacoste-Julien 2014)
- ...

Idea: variance reduction

Example: SVRG

- Initialize $\hat{\mathbf{w}}$
- For epoch 1, 2, 3, ...
 - ▶ Compute $\nabla f(\hat{\mathbf{w}}) = \sum_{i \in \mathcal{D}} \nabla \ell_i(\hat{\mathbf{w}})$
 - ▶ Initialize $\mathbf{w}_0 = \hat{\mathbf{w}}$
 - ▶ For t in length of epochs

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha [\nabla f(\hat{\mathbf{w}}) + \nabla \ell_{i(t)}(\mathbf{w}_{t-1}) - \nabla \ell_{i(t)}(\hat{\mathbf{w}})]$$

- ▶ Update $\hat{\mathbf{w}} = \mathbf{w}_t$
- Output $\hat{\mathbf{w}}$

Quiz:

Quiz:

- Stepsize/Learning rate rules?

Quiz:

- Stepsize/Learning rate rules?
- Descent directions?

Quiz:

- Stepsize/Learning rate rules?
- Descent directions?
- Properties of convex functions?

Quiz:

- Stepsize/Learning rate rules?
- Descent directions?
- Properties of convex functions?
- Convergence rates?

Quiz:

- Stepsize/Learning rate rules?
- Descent directions?
- Properties of convex functions?
- Convergence rates?
- Improvements?

Important topics of this lecture

Important topics of this lecture

- Convex optimization basics

Important topics of this lecture

- Convex optimization basics
- Algorithm choices

Important topics of this lecture

- Convex optimization basics
- Algorithm choices
- Rates

Important topics of this lecture

- Convex optimization basics
- Algorithm choices
- Rates

Up next:

- How to deal with constraints