

# CS260: Machine Learning Algorithms

## Lecture 3: Optimization

Cho-Jui Hsieh  
UCLA

Jan 14, 2019

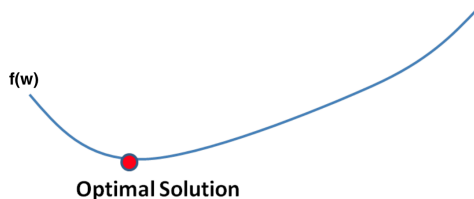
# Optimization

- Goal: find the minimizer of a function

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

For now we assume  $f$  is twice differentiable

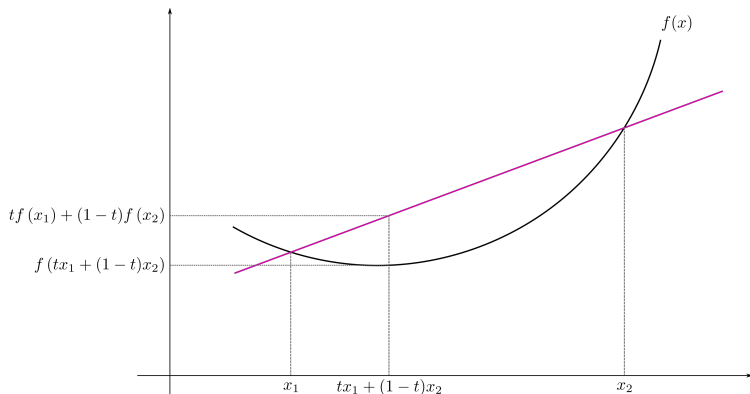
- Machine learning algorithm: find the hypothesis that **minimizes training error**



# Convex function

- A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex function  
 $\Leftrightarrow$  the function  $f$  is **below** any line segment between two points on  $f$ :

$$\forall \mathbf{x}_1, \mathbf{x}_2, \forall t \in [0, 1], \quad f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$

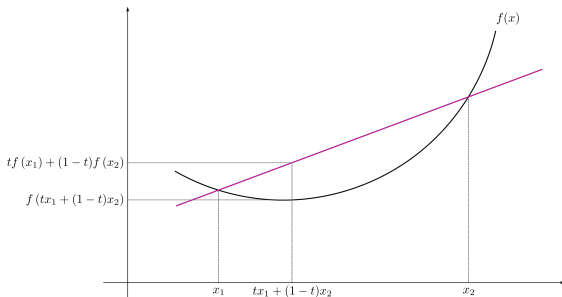


# Convex function

- A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex function

$\Leftrightarrow$  the function  $f$  is below any line segment between two points on  $f$ :

$$\forall \mathbf{x}_1, \mathbf{x}_2, \forall t \in [0, 1], \quad f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$

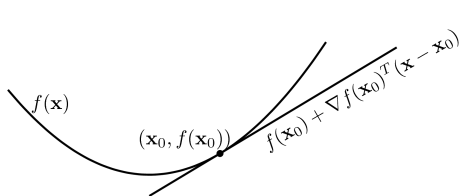


**Strict convex:**  $f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) < tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$

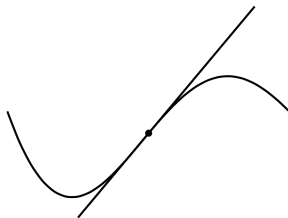
# Convex function

- Another equivalent definition for differentiable function:

$f$  is convex if and only if  $f(\mathbf{x}) \geq f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0), \quad \forall \mathbf{x}, \mathbf{x}_0$



convex function



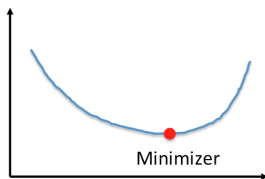
nonconvex function



# Convex function

- Strict convex function:
  - $\nabla f(\mathbf{w}^*) = 0 \Leftrightarrow \mathbf{w}^*$  is the unique global minimum  
most algorithms only converge to gradient = 0
  - Example: Linear regression when  $X^T X$  is invertible

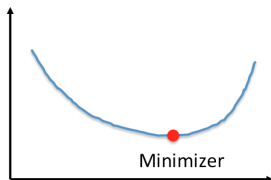
**Convex**



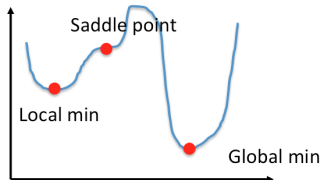
# Convex vs Nonconvex

- Convex function:
  - $\nabla f(\mathbf{w}^*) = 0 \Leftrightarrow \mathbf{w}^*$  is a global minimum
  - Example: linear regression, logistic regression, ...
- Non-convex function:
  - $\nabla f(\mathbf{w}^*) = 0 \Leftrightarrow \mathbf{w}^*$  is Global min, local min, or saddle point  
(also called **stationary points**)  
most algorithms only converge to **stationary points**
  - Example: neural network, ...

**Convex**



**Non-Convex**





# Gradient descent

# Gradient Descent

- Gradient descent: repeatedly do

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \alpha \nabla f(\mathbf{w}^t)$$

$\alpha > 0$  is the **step size**

# Gradient Descent

- Gradient descent: repeatedly do

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \alpha \nabla f(\mathbf{w}^t)$$

$\alpha > 0$  is the **step size**

- Generate the sequence  $\mathbf{w}^1, \mathbf{w}^2, \dots$

converge to stationary points (  $\lim_{t \rightarrow \infty} \|\nabla f(\mathbf{w}^t)\| = 0$  )

# Gradient Descent

- Gradient descent: repeatedly do

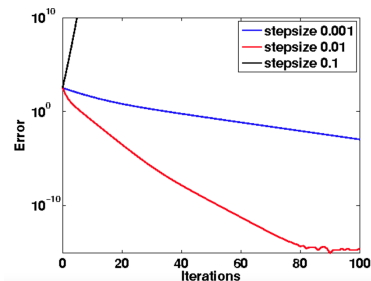
$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \alpha \nabla f(\mathbf{w}^t)$$

$\alpha > 0$  is the **step size**

- Generate the sequence  $\mathbf{w}^1, \mathbf{w}^2, \dots$

converge to stationary points (  $\lim_{t \rightarrow \infty} \|\nabla f(\mathbf{w}^t)\| = 0$  )

- Step size **too large**  $\Rightarrow$  **diverge**; **too small**  $\Rightarrow$  **slow convergence**



# Why gradient descent?

- Successive approximation view

At each iteration, form an approximation function of  $f(\cdot)$ :

$$f(\mathbf{w}^t + \mathbf{d}) \approx g(\mathbf{d}) := f(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)^T \mathbf{d} + \frac{1}{2\alpha} \|\mathbf{d}\|^2$$

Update solution by  $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \mathbf{d}^*$

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} g(\mathbf{d})$$

$$\nabla g(\mathbf{d}^*) = 0 \Rightarrow \nabla f(\mathbf{w}^t) + \frac{1}{\alpha} \mathbf{d}^* = 0 \Rightarrow \mathbf{d}^* = -\alpha \nabla f(\mathbf{w}^t)$$

# Why gradient descent?

- Successive approximation view

At each iteration, form an approximation function of  $f(\cdot)$ :

$$f(\mathbf{w}^t + \mathbf{d}) \approx g(\mathbf{d}) := f(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)^T \mathbf{d} + \frac{1}{2\alpha} \|\mathbf{d}\|^2$$

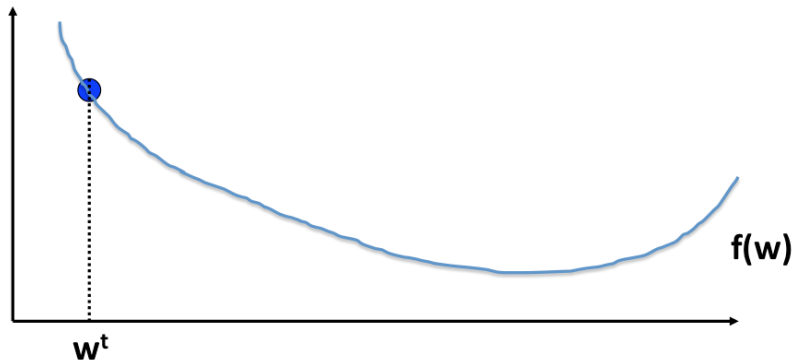
Update solution by  $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \mathbf{d}^*$

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} g(\mathbf{d})$$

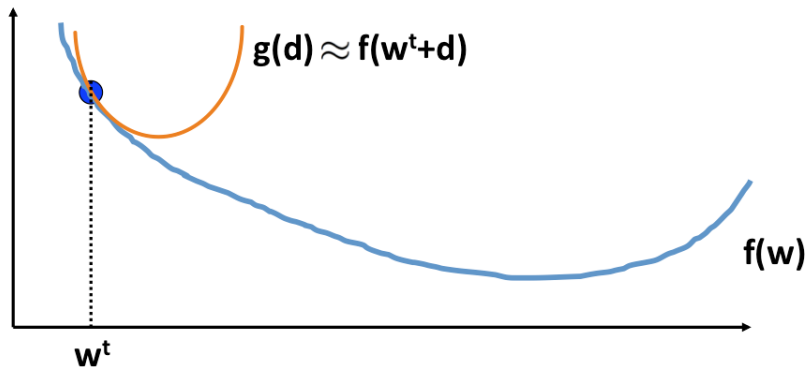
$$\nabla g(\mathbf{d}^*) = 0 \Rightarrow \nabla f(\mathbf{w}^t) + \frac{1}{\alpha} \mathbf{d}^* = 0 \Rightarrow \mathbf{d}^* = -\alpha \nabla f(\mathbf{w}^t)$$

- $\mathbf{d}^*$  will decrease  $f(\cdot)$  if  $\alpha$  (step size) is sufficiently small

# Illustration of gradient descent



# Illustration of gradient descent

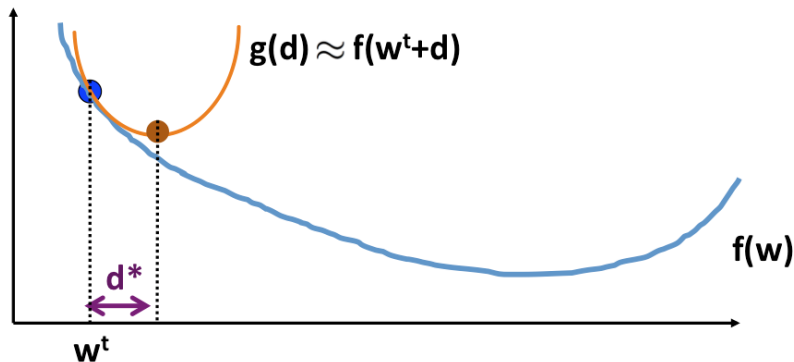


Form a quadratic approximation

$$f(\mathbf{w}^t + \mathbf{d}) \approx g(\mathbf{d}) = f(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)^T \mathbf{d} + \frac{1}{2\alpha} \|\mathbf{d}\|^2$$



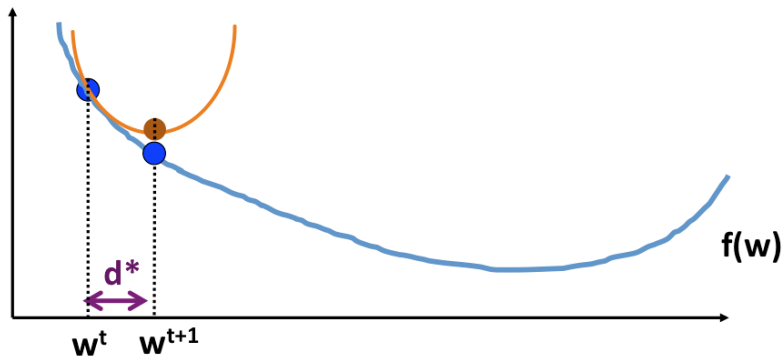
# Illustration of gradient descent



Minimize  $g(\mathbf{d})$ :

$$\nabla g(\mathbf{d}^*) = 0 \Rightarrow \nabla f(\mathbf{w}^t) + \frac{1}{\alpha} \mathbf{d}^* = 0 \Rightarrow \mathbf{d}^* = -\alpha \nabla f(\mathbf{w}^t)$$

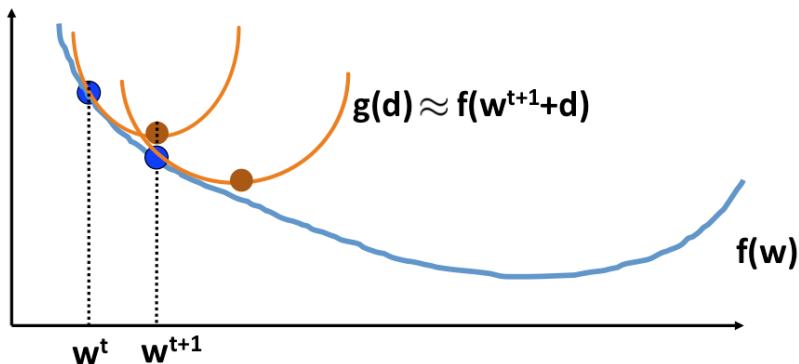
# Illustration of gradient descent



Update

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{d}^* = \mathbf{w}^t - \alpha \nabla f(\mathbf{w}^t)$$

# Illustration of gradient descent

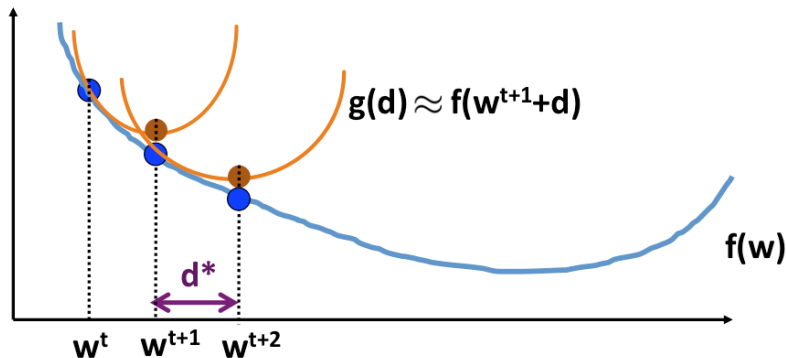


Form another quadratic approximation

$$f(\mathbf{w}^{t+1} + \mathbf{d}) \approx g(\mathbf{d}) = f(\mathbf{w}^{t+1}) + \nabla f(\mathbf{w}^{t+1})^T \mathbf{d} + \frac{1}{2\alpha} \|\mathbf{d}\|^2$$

$$\mathbf{d}^* = -\alpha \nabla f(\mathbf{w}^{t+1})$$

# Illustration of gradient descent

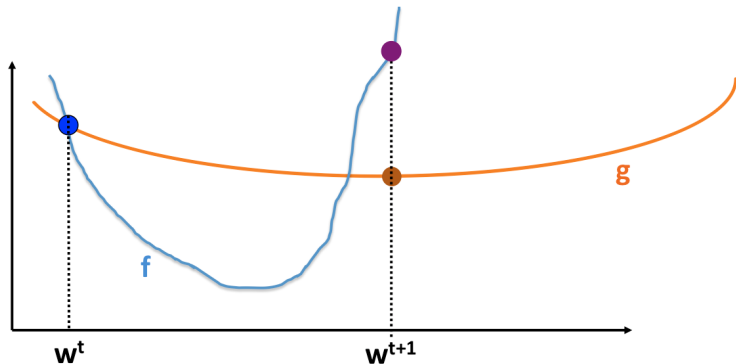


Update

$$w^{t+2} = w^{t+1} + \mathbf{d}^* = w^{t+1} - \alpha \nabla f(w^{t+1})$$

# When will it diverge?

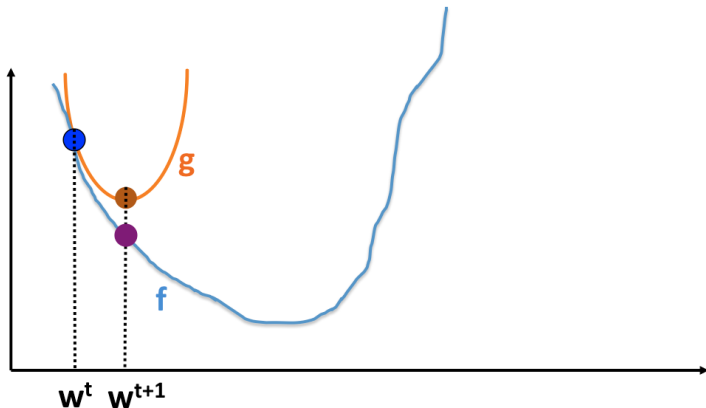
Can diverge ( $f(\mathbf{w}^t) < f(\mathbf{w}^{t+1})$ ) if  $g$  is **not** an upperbound of  $f$



$f(\mathbf{w}^t) < f(\mathbf{w}^{t+1})$ , diverge because  $g$ 's curvature is too small

# When will it converge?

Always converge ( $f(\mathbf{w}^t) > f(\mathbf{w}^{t+1})$ ) when  $g$  is an upper bound of  $f$



$f(\mathbf{w}^t) > f(\mathbf{w}^{t+1})$ , converge when  $g$ 's curvature is large enough

# Convergence

- Let  $L$  be the **Lipchitz constant**  
 $(\nabla^2 f(\mathbf{x}) \preceq L I \text{ for all } \mathbf{x})$
- **Theorem:** gradient descent converges if  $\alpha < \frac{1}{L}$

# Convergence

- Let  $L$  be the **Lipchitz constant**

$$(\nabla^2 f(\mathbf{x}) \preceq L I \text{ for all } \mathbf{x})$$

- Theorem:** gradient descent converges if  $\alpha < \frac{1}{L}$

- Why?

- When  $\alpha < 1/L$ , for any  $\mathbf{d}$ ,

$$\begin{aligned} g(\mathbf{d}) &= f(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)^T \mathbf{d} + \frac{1}{2\alpha} \|\mathbf{d}\|^2 \\ &> f(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)^T \mathbf{d} + \frac{L}{2} \|\mathbf{d}\|^2 \\ &\geq f(\mathbf{w}^t + \mathbf{d}) \end{aligned}$$



# Convergence

- Let  $L$  be the **Lipchitz constant**

$$(\nabla^2 f(\mathbf{x}) \preceq L I \text{ for all } \mathbf{x})$$

- Theorem:** gradient descent converges if  $\alpha < \frac{1}{L}$

- Why?

- When  $\alpha < 1/L$ , for any  $\mathbf{d}$ ,

$$\begin{aligned} g(\mathbf{d}) &= f(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)^T \mathbf{d} + \frac{1}{2\alpha} \|\mathbf{d}\|^2 \\ &> f(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)^T \mathbf{d} + \frac{L}{2} \|\mathbf{d}\|^2 \\ &\geq f(\mathbf{w}^t + \mathbf{d}) \end{aligned}$$

- So,  $f(\mathbf{w}^t + \mathbf{d}^*) < g(\mathbf{d}^*) \leq g(0) = f(\mathbf{w}^t)$

# Convergence

- Let  $L$  be the **Lipchitz constant**

$$(\nabla^2 f(\mathbf{x}) \preceq L I \text{ for all } \mathbf{x})$$

- Theorem:** gradient descent converges if  $\alpha < \frac{1}{L}$

- Why?

- When  $\alpha < 1/L$ , for any  $\mathbf{d}$ ,

$$\begin{aligned} g(\mathbf{d}) &= f(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)^T \mathbf{d} + \frac{1}{2\alpha} \|\mathbf{d}\|^2 \\ &> f(\mathbf{w}^t) + \nabla f(\mathbf{w}^t)^T \mathbf{d} + \frac{L}{2} \|\mathbf{d}\|^2 \\ &\geq f(\mathbf{w}^t + \mathbf{d}) \end{aligned}$$

- So,  $f(\mathbf{w}^t + \mathbf{d}^*) < g(\mathbf{d}^*) \leq g(0) = f(\mathbf{w}^t)$
- In formal proof, need to show  $f(\mathbf{w}^t + \mathbf{d}^*)$  is **sufficiently** smaller than  $f(\mathbf{w}^t)$

# Applying to Logistic regression

## gradient descent for logistic regression

- Initialize the weights  $\mathbf{w}_0$
- For  $t = 1, 2, \dots$ 
  - Compute the gradient

$$\nabla f(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}}$$

- Update the weights:  $\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla f(\mathbf{w})$
- Return the final weights  $\mathbf{w}$

# Applying to Logistic regression

## gradient descent for logistic regression

- Initialize the weights  $\mathbf{w}_0$
- For  $t = 1, 2, \dots$ 
  - Compute the gradient

$$\nabla f(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}}$$

- Update the weights:  $\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla f(\mathbf{w})$
- Return the final weights  $\mathbf{w}$

When to stop?

- Fixed number of iterations, or
- Stop when  $\|\nabla f(\mathbf{w})\| < \epsilon$

# Line Search

- In practice, we do not know  $L \dots$   
need to tune step size when running gradient descent

# Line Search

- In practice, we do not know  $L \dots$   
need to tune step size when running gradient descent
- Line Search: Select step size automatically (for gradient descent)

# Line Search

- The **back-tracking** line search:

- Start from some **large**  $\alpha_0$

- Try  $\alpha = \alpha_0, \frac{\alpha_0}{2}, \frac{\alpha_0}{4}, \dots$

Stop when  $\alpha$  satisfies some **sufficient decrease condition**

# Line Search

- The **back-tracking** line search:
  - Start from some **large**  $\alpha_0$
  - Try  $\alpha = \alpha_0, \frac{\alpha_0}{2}, \frac{\alpha_0}{4}, \dots$   
Stop when  $\alpha$  satisfies some **sufficient decrease condition**
  - A simple condition:  $f(\mathbf{w} + \alpha \mathbf{d}) < f(\mathbf{w})$



# Line Search

- The **back-tracking** line search:
  - Start from some **large**  $\alpha_0$
  - Try  $\alpha = \alpha_0, \frac{\alpha_0}{2}, \frac{\alpha_0}{4}, \dots$   
Stop when  $\alpha$  satisfies some **sufficient decrease condition**
  - A simple condition:  $f(\mathbf{w} + \alpha \mathbf{d}) < f(\mathbf{w})$   
often works in practice but doesn't work in theory

# Line Search

- The **back-tracking** line search:

- Start from some **large**  $\alpha_0$

- Try  $\alpha = \alpha_0, \frac{\alpha_0}{2}, \frac{\alpha_0}{4}, \dots$

Stop when  $\alpha$  satisfies some **sufficient decrease condition**

- A simple condition:  $f(\mathbf{w} + \alpha \mathbf{d}) < f(\mathbf{w})$

often works in practice but doesn't work in theory

- A (provable) sufficient decrease condition:

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

for a constant  $\sigma \in (0, 1)$

# Line Search

## gradient descent with backtracking line search

- Initialize the weights  $\mathbf{w}_0$
- For  $t = 1, 2, \dots$ 
  - Compute the gradient

$$\mathbf{d} = -\nabla f(\mathbf{w})$$

- For  $\alpha = \alpha_0, \alpha_0/2, \alpha_0/4, \dots$   
Break if  $f(\mathbf{w} + \alpha\mathbf{d}) \leq f(\mathbf{w}) + \sigma\alpha\nabla f(\mathbf{w})^T \mathbf{d}$
  - Update  $\mathbf{w} \leftarrow \mathbf{w} + \alpha\mathbf{d}$
- Return the final solution  $\mathbf{w}$

# Stochastic Gradient descent

# Large-scale Problems

- Machine learning: usually minimizing the training loss

$$\min_{\mathbf{w}} \left\{ \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{w}^T \mathbf{x}_n, y_n) \right\} := f(\mathbf{w}) \text{ (linear model)}$$

$$\min_{\mathbf{w}} \left\{ \frac{1}{N} \sum_{n=1}^N \ell(h_{\mathbf{w}}(\mathbf{x}_n), y_n) \right\} := f(\mathbf{w}) \text{ (general hypothesis)}$$

$\ell$ : loss function (e.g.,  $\ell(a, b) = (a - b)^2$ )

- Gradient descent:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \underbrace{\nabla f(\mathbf{w})}_{\text{Main computation}}$$

# Large-scale Problems

- Machine learning: usually minimizing the training loss

$$\min_{\mathbf{w}} \left\{ \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{w}^T \mathbf{x}_n, y_n) \right\} := f(\mathbf{w}) \text{ (linear model)}$$

$$\min_{\mathbf{w}} \left\{ \frac{1}{N} \sum_{n=1}^N \ell(h_{\mathbf{w}}(\mathbf{x}_n), y_n) \right\} := f(\mathbf{w}) \text{ (general hypothesis)}$$

$\ell$ : loss function (e.g.,  $\ell(a, b) = (a - b)^2$ )

- Gradient descent:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \underbrace{\nabla f(\mathbf{w})}_{\text{Main computation}}$$

- In general,  $f(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N f_n(\mathbf{w})$ ,  
each  $f_n(\mathbf{w})$  only depends on  $(\mathbf{x}_n, y_n)$

# Stochastic gradient

- Gradient:

$$\nabla f(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w})$$

- Each gradient computation needs to go through **all training samples**  
slow when millions of samples
- Faster way to compute “approximate gradient”?

# Stochastic gradient

- Gradient:

$$\nabla f(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w})$$

- Each gradient computation needs to go through **all training samples** slow when millions of samples
- Faster way to compute “approximate gradient”?
- Use **stochastic sampling**:
  - Sample a small subset  $B \subseteq \{1, \dots, N\}$
  - Estimated gradient

$$\nabla f(\mathbf{w}) \approx \frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w})$$

$|B|$ : batch size



# Stochastic gradient descent

## Stochastic Gradient Descent (SGD)

- Input: training data  $\{\mathbf{x}_n, y_n\}_{n=1}^N$
- Initialize  $\mathbf{w}$  (zero or random)
- For  $t = 1, 2, \dots$ 
  - Sample a **small batch**  $B \subseteq \{1, \dots, N\}$
  - Update parameter

$$\mathbf{w} \leftarrow \mathbf{w} - \eta^t \frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w})$$

# Stochastic gradient descent

## Stochastic Gradient Descent (SGD)

- Input: training data  $\{\mathbf{x}_n, y_n\}_{n=1}^N$
- Initialize  $\mathbf{w}$  (zero or random)
- For  $t = 1, 2, \dots$ 
  - Sample a **small batch**  $B \subseteq \{1, \dots, N\}$
  - Update parameter

$$\mathbf{w} \leftarrow \mathbf{w} - \eta^t \frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w})$$

**Extreme case:**  $|B| = 1 \Rightarrow$  **Sample one training data at a time**

# Logistic Regression by SGD

- Logistic regression:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N \underbrace{\log(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n})}_{f_n(\mathbf{w})}$$

## SGD for Logistic Regression

- Input: training data  $\{\mathbf{x}_n, y_n\}_{n=1}^N$
- Initialize  $\mathbf{w}$  (zero or random)
- For  $t = 1, 2, \dots$ 
  - Sample a batch  $B \subseteq \{1, \dots, N\}$
  - Update parameter

$$\mathbf{w} \leftarrow \mathbf{w} - \eta^t \frac{1}{|B|} \sum_{i \in B} \underbrace{\frac{-y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}}}_{\nabla f_n(\mathbf{w})}$$

# Why SGD works?

- Stochastic gradient is an **unbiased estimator** of full gradient:

$$\begin{aligned} E\left[\frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w})\right] &= \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w}) \\ &= \nabla f(\mathbf{w}) \end{aligned}$$

# Why SGD works?

- Stochastic gradient is an **unbiased estimator** of full gradient:

$$\begin{aligned} E\left[\frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w})\right] &= \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w}) \\ &= \nabla f(\mathbf{w}) \end{aligned}$$

- Each iteration updated by

gradient + **zero-mean noise**

# Stochastic gradient descent

- In gradient descent,  $\eta$  (step size) is a fixed constant
- Can we use fixed step size for SGD?

# Stochastic gradient descent

- In gradient descent,  $\eta$  (step size) is a fixed constant
- Can we use fixed step size for SGD?
- SGD with fixed step size **cannot converge to global/local minimizers**

# Stochastic gradient descent

- In gradient descent,  $\eta$  (step size) is a fixed constant
- Can we use fixed step size for SGD?
- SGD with fixed step size **cannot converge to global/local minimizers**
- If  $\mathbf{w}^*$  is the minimizer,  $\nabla f(\mathbf{w}^*) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w}^*) = 0$ ,



# Stochastic gradient descent

- In gradient descent,  $\eta$  (step size) is a fixed constant
- Can we use fixed step size for SGD?
- SGD with fixed step size **cannot converge to global/local minimizers**
- If  $\mathbf{w}^*$  is the minimizer,  $\nabla f(\mathbf{w}^*) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w}^*) = 0$ ,

but  $\frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w}^*) \neq 0$  if  $B$  is a subset

# Stochastic gradient descent

- In gradient descent,  $\eta$  (step size) is a fixed constant
- Can we use fixed step size for SGD?
- SGD with fixed step size **cannot converge to global/local minimizers**
- If  $\mathbf{w}^*$  is the minimizer,  $\nabla f(\mathbf{w}^*) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w}^*) = 0$ ,

$$\text{but } \frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w}^*) \neq 0 \quad \text{if } B \text{ is a subset}$$

(Even if we got minimizer, SGD will **move away** from it)

# Stochastic gradient descent, step size

- To make SGD converge:

Step size should decrease to 0

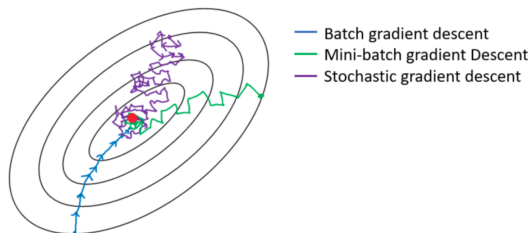
$$\eta^t \rightarrow 0$$

Usually with polynomial rate:  $\eta^t \approx t^{-a}$  with constant  $a$

# Stochastic gradient descent vs Gradient descent

## Stochastic gradient descent:

- pros:
  - cheaper computation per iteration
  - faster convergence in the beginning
- cons:
  - less stable, slower final convergence
  - hard to tune step size



(Figure from <https://medium.com/@ImadPhd/>

gradient-descent-algorithm-and-its-variants-10f652806a3)

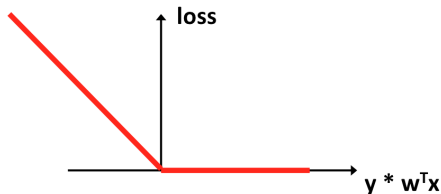
# Revisit perceptron Learning Algorithm

- Given a classification data  $\{\mathbf{x}_n, y_n\}_{n=1}^N$
- Learning a linear model:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{w}^T \mathbf{x}_n, y_n)$$

- Consider the loss:

$$\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = \max(0, -y_n \mathbf{w}^T \mathbf{x}_n)$$



What's the gradient?

# Revisit perceptron Learning Algorithm

$$\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = \max(0, -y_n \mathbf{w}^T \mathbf{x}_n)$$

Consider two cases:

- Case I:  $y_n \mathbf{w}^T \mathbf{x}_n > 0$  (prediction **correct**)
  - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
  - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$

# Revisit perceptron Learning Algorithm

$$\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = \max(0, -y_n \mathbf{w}^T \mathbf{x}_n)$$

Consider two cases:

- Case I:  $y_n \mathbf{w}^T \mathbf{x}_n > 0$  (prediction **correct**)
  - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
  - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
- Case II:  $y_n \mathbf{w}^T \mathbf{x}_n < 0$  (prediction **wrong**)
  - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = -y_n \mathbf{w}^T \mathbf{x}_n$
  - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = -y_n \mathbf{x}_n$

# Revisit perceptron Learning Algorithm

$$\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = \max(0, -y_n \mathbf{w}^T \mathbf{x}_n)$$

Consider two cases:

- Case I:  $y_n \mathbf{w}^T \mathbf{x}_n > 0$  (prediction **correct**)
  - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
  - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
- Case II:  $y_n \mathbf{w}^T \mathbf{x}_n < 0$  (prediction **wrong**)
  - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = -y_n \mathbf{w}^T \mathbf{x}_n$
  - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = -y_n \mathbf{x}_n$

SGD update rule: Sample an index  $n$

$$\mathbf{w}^{t+1} \leftarrow \begin{cases} \mathbf{w}^t & \text{if } y_n \mathbf{w}^T \mathbf{x}_n \geq 0 \text{ (predict correct)} \\ \mathbf{w}^t + \eta^t y_n \mathbf{x}_n & \text{if } y_n \mathbf{w}^T \mathbf{x}_n < 0 \text{ (predict wrong)} \end{cases}$$

Equivalent to Perceptron Learning Algorithm when  $\eta^t = 1$



# Conclusions

- Gradient descent
- Stochastic gradient descent

Questions?