

# Numerical Optimization

Shan-Hung Wu  
*shwu@cs.nthu.edu.tw*

Department of Computer Science,  
National Tsing Hua University, Taiwan

Machine Learning

# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

# Outline

- ① **Numerical Computation**
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ **Duality\***

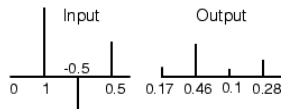
# Numerical Computation

- Machine learning algorithms usually require a high amount of numerical computation in involving real numbers
- However, real numbers cannot be represented precisely using a finite amount of memory
- Watch out the *numeric errors* when implementing machine learning algorithms

# Overflow and Underflow I

- Consider the **softmax function**  
 $\text{softmax} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ :

$$\text{softmax}(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{j=1}^d \exp(x_j)}$$



- Commonly used to transform a group of real values to “probabilities”
- Analytically, if  $x_i = c$  for all  $i$ , then  $\text{softmax}(\mathbf{x})_i = 1/d$
- Numerically, this may not occur when  $|c|$  is large
  - A positive  $c$  causes overflow
  - A negative  $c$  causes underflow and divide-by-zero error
- How to avoid these errors?

# Overflow and Underflow II

- Instead of evaluating  $\text{softmax}(\mathbf{x})$  directly, we can transform  $\mathbf{x}$  into

$$\mathbf{z} = \mathbf{x} - \max_i x_i \mathbf{1}$$

and then evaluate  $\text{softmax}(\mathbf{z})$

- $\text{softmax}(\mathbf{z})_i = \frac{\exp(x_i - m)}{\sum \exp(x_j - m)} = \frac{\exp(x_i) / \exp(m)}{\sum \exp(x_j) / \exp(m)} = \frac{\exp(x_i)}{\sum \exp(x_j)} = \text{softmax}(\mathbf{x})_i$
- No overflow, as  $\exp(\text{largest attribute of } \mathbf{x}) = 1$
- Denominator is at least 1, no divide-by-zero error
- What are the numerical issues of  $\log \text{softmax}(\mathbf{z})$ ? How to stabilize it?  
[Homework]

# Poor Conditioning I

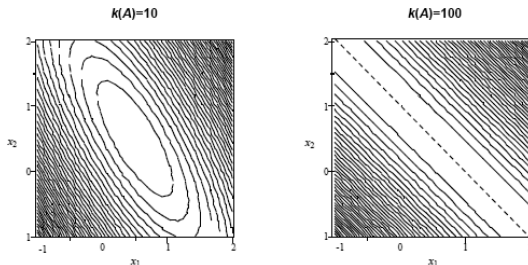
- The “conditioning” refer to how much the input of a function can change given a small change in the output
- Suppose we want to solve  $\mathbf{x}$  in  $f(\mathbf{x}) = \mathbf{A}\mathbf{x} = \mathbf{y}$ , where  $\mathbf{A}^{-1}$  exists
- The *condition number* of  $\mathbf{A}$  can be expressed by

$$\kappa(\mathbf{A}) = \max_{i,j} \left| \frac{\lambda_i}{\lambda_j} \right|$$

- We say the problem is *poorly* (or *ill-*) *conditioned* when  $\kappa(\mathbf{A})$  is large
- Hard to solve  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{y}$  precisely given a rounded  $\mathbf{y}$ 
  - $\mathbf{A}^{-1}$  amplifies pre-existing numeric errors

# Poor Conditioning II

- The contours of  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{A}\mathbf{x} + \mathbf{b}^\top \mathbf{x} + c$ , where  $\mathbf{A}$  is symmetric:



- When  $\kappa(\mathbf{A})$  is large,  $f$  stretches space differently along different attribute directions
  - Surface is flat in some directions but steep in others
- Hard to solve  $f'(\mathbf{x}) = \mathbf{0} \Rightarrow \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$



# Outline

- 1 Numerical Computation
- 2 Optimization Problems**
- 3 Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- 4 Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- 5 Constrained Optimization
- 6 Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- 7 Duality\*

# Optimization Problems

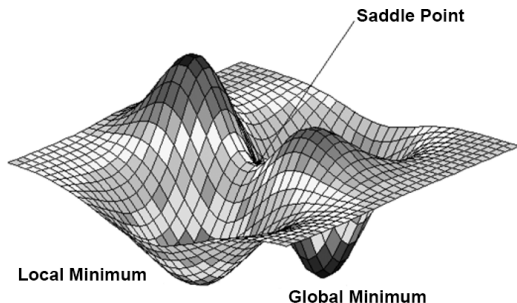
- An *optimization problem* is to minimize a *cost function*  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \mathbb{C} \end{aligned}$$

where  $\mathbb{C} \subseteq \mathbb{R}^d$  is called the *feasible set* containing *feasible points*

- Or, maximizing an *objective function*
- Maximizing  $f$  equals to minimizing  $-f$
- If  $\mathbb{C} = \mathbb{R}^d$ , we say the optimization problem is unconstrained
- $\mathbb{C}$  can be a set of function *constrains*, i.e.,  $\mathbb{C} = \{\mathbf{x} : g^{(i)}(\mathbf{x}) \leq 0\}_i$
- Sometimes, we single out equality constrains  
 $\mathbb{C} = \{\mathbf{x} : g^{(i)}(\mathbf{x}) \leq 0, h^{(j)}(\mathbf{x}) = 0\}_{i,j}$ 
  - Each equality constrain can be written as two inequality constrains

# Minimums and Optimal Points



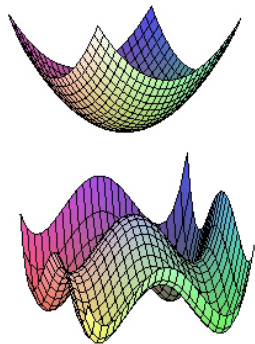
- **Critical points:**  $\{\mathbf{x} : f'(\mathbf{x}) = \mathbf{0}\}$ 
  - Minima:  $\{\mathbf{x} : f'(\mathbf{x}) = \mathbf{0} \text{ and } \mathbf{H}(f)(\mathbf{x}) \succ \mathbf{O}\}$ , where  $\mathbf{H}(f)(\mathbf{x})$  is the Hessian matrix (containing curvatures) of  $f$  at point  $\mathbf{x}$
  - Maxima:  $\{\mathbf{x} : f'(\mathbf{x}) = \mathbf{0} \text{ and } \mathbf{H}(f)(\mathbf{x}) \prec \mathbf{O}\}$
  - Plateau or saddle points:  $\{\mathbf{x} : f'(\mathbf{x}) = \mathbf{0} \text{ and } \mathbf{H}(f)(\mathbf{x}) = \mathbf{O} \text{ or indefinite}\}$
- $y^* = \min_{\mathbf{x} \in \mathbb{C}} f(\mathbf{x}) \in \mathbb{R}$  is called the **global minimum**
  - Global minima vs. local minima
- $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{C}} f(\mathbf{x})$  is called the **optimal point**

# Convex Optimization Problems

- An optimization problem is **convex** iff
  - ①  $f$  is convex by having a “convex hull” surface, i.e.,

$$H(f)(x) \succeq \mathbf{0}, \forall x$$

- ②  $g_i(x)$ 's are convex and  $h_j(x)$ 's are affine
- Convex problems are “easier” since
    - Local minima are necessarily global minima
    - No saddle point
    - We can get the global minimum by solving  $f'(x) = \mathbf{0}$



# Analytical Solutions vs. Numerical Solutions I

- Consider the problem:

$$\arg \min_x \frac{1}{2} (\|A\mathbf{x} - \mathbf{b}\|^2 + \lambda \|\mathbf{x}\|^2)$$

- Analytical solutions?
- The cost function  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top (A^\top A + \lambda I)\mathbf{x} - \mathbf{b}^\top A\mathbf{x} + \frac{1}{2}\|\mathbf{b}\|^2$  is convex
- Solving  $f'(\mathbf{x}) = \mathbf{x}^\top (A^\top A + \lambda I) - \mathbf{b}^\top A = 0$ , we have

$$\mathbf{x}^* = (A^\top A + \lambda I)^{-1} A^\top \mathbf{b}$$

# Analytical Solutions vs. Numerical Solutions II

- Problem ( $\mathbf{A} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{b} \in \mathbb{R}^n$ ,  $\lambda \in \mathbb{R}$ ):

$$\arg \min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} (\|\mathbf{Ax} - \mathbf{b}\|^2 + \lambda \|\mathbf{x}\|^2)$$

- Analytical solution:  $\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^\top \mathbf{b}$
- In practice, we may not be able to solve  $f'(\mathbf{x}) = \mathbf{0}$  analytically and get  $\mathbf{x}$  in a closed form
  - E.g., when  $\lambda = 0$  and  $n < d$
- Even if we can, the computation cost may be too high
  - E.g, inverting  $\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I} \in \mathbb{R}^{d \times d}$  takes  $O(d^3)$  time
- **Numerical methods**: since numerical errors are inevitable, why not just obtain an **approximation** of  $\mathbf{x}^*$ ?
- Start from  $\mathbf{x}^{(0)}$ , iteratively calculating  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$  such that  $f(\mathbf{x}^{(1)}) \geq f(\mathbf{x}^{(2)}) \geq \dots$ 
  - Usually require much less time to have a good enough  $\mathbf{x}^{(t)} \approx \mathbf{x}^*$

# Outline

- 1 Numerical Computation
- 2 Optimization Problems
- 3 Unconstrained Optimization**
  - Gradient Descent
  - Newton's Method
- 4 Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- 5 Constrained Optimization
- 6 Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- 7 Duality\*

# Unconstrained Optimization

- Problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}),$$

where  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is not necessarily convex



# General Descent Algorithm

---

Input:  $\mathbf{x}^{(0)} \in \mathbb{R}^d$ , an initial guess

repeat

    Determine a **descent direction**  $\mathbf{d}^{(t)} \in \mathbb{R}^d$  ;

**Line search**: choose a **step size** or **learning rate**  $\eta^{(t)} > 0$  such that  $f(\mathbf{x}^{(t)} + \eta^{(t)}\mathbf{d}^{(t)})$  is minimal along the ray  $\mathbf{x}^{(t)} + \eta^{(t)}\mathbf{d}^{(t)}$  ;

**Update rule**:  $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \eta^{(t)}\mathbf{d}^{(t)}$  ;

until *convergence criterion is satisfied*;

---

- Convergence criterion:  $\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\| \leq \varepsilon$ ,  $\|\nabla f(\mathbf{x}^{(t+1)})\| \leq \varepsilon$ , etc.
- Line search step could be skipped by letting  $\eta^{(t)}$  be a small constant

# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ **Unconstrained Optimization**
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

# Gradient Descent I

- By Taylor's theorem, we can approximate  $f$  locally at point  $\mathbf{x}^{(t)}$  using a linear function  $\tilde{f}$ , i.e.,

$$f(\mathbf{x}) \approx \tilde{f}(\mathbf{x}; \mathbf{x}^{(t)}) = f(\mathbf{x}^{(t)}) + \nabla f(\mathbf{x}^{(t)})^\top (\mathbf{x} - \mathbf{x}^{(t)})$$

for  $\mathbf{x}$  close enough to  $\mathbf{x}^{(t)}$

- This implies that if we pick a close  $\mathbf{x}^{(t+1)}$  that decreases  $\tilde{f}$ , we are likely to decrease  $f$  as well
- We can pick  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \nabla f(\mathbf{x}^{(t)})$  for some small  $\eta > 0$ , since

$$\tilde{f}(\mathbf{x}^{(t+1)}) = f(\mathbf{x}^{(t)}) - \eta \|\nabla f(\mathbf{x}^{(t)})\|^2 \leq \tilde{f}(\mathbf{x}^{(t)})$$

# Gradient Descent II

---

**Input:**  $\mathbf{x}^{(0)} \in \mathbb{R}^d$  an initial guess, a small  $\eta > 0$

**repeat**

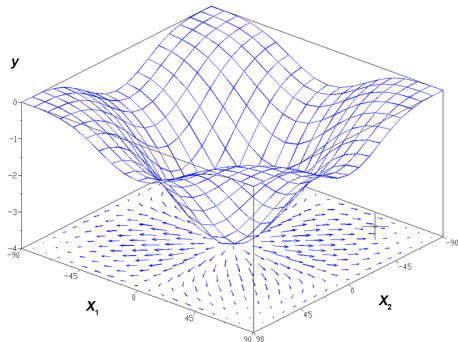
    |  $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \eta \nabla f(\mathbf{x}^{(t)})$  ;

**until** *convergence criterion is satisfied*;

---

# Is Negative Gradient a Good Direction? I

- Update rule:  
 $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \eta \nabla f(\mathbf{x}^{(t)})$
- Yes, as  $\nabla f(\mathbf{x}^{(t)}) \in \mathbb{R}^d$  denotes the steepest ascent direction of  $f$  at point  $\mathbf{x}^{(t)}$
- $-\nabla f(\mathbf{x}^{(t)}) \in \mathbb{R}^d$  the steepest descent direction
- But why?



# Is Negative Gradient a Good Direction? II

- Consider the slope of  $f$  in a given direction  $\mathbf{u}$  at point  $\mathbf{x}^{(t)}$
- This is the **directional derivative** of  $f$ , i.e., the derivative of function  $f(\mathbf{x}^{(t)} + \varepsilon \mathbf{u})$  with respect to  $\varepsilon$ , evaluated at  $\varepsilon = 0$
- By the chain rule, we have  $\frac{\partial}{\partial \varepsilon} f(\mathbf{x}^{(t)} + \varepsilon \mathbf{u}) = \nabla f(\mathbf{x}^{(t)} + \varepsilon \mathbf{u})^\top \mathbf{u}$ , which equals to  $\nabla f(\mathbf{x}^{(t)})^\top \mathbf{u}$  when  $\varepsilon = 0$

## Theorem (Chain Rule)

Let  $\mathbf{g} : \mathbb{R} \rightarrow \mathbb{R}^d$  and  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , then

$$(f \circ \mathbf{g})'(x) = f'(\mathbf{g}(x))\mathbf{g}'(x) = \nabla f(\mathbf{g}(x))^\top \begin{bmatrix} g'_1(x) \\ \vdots \\ g'_n(x) \end{bmatrix}.$$

# Is Negative Gradient a Good Direction? III

- To find the direction that decreases  $f$  fastest at  $\mathbf{x}^{(t)}$ , we solve the problem:

$$\arg \min_{\mathbf{u}, \|\mathbf{u}\|=1} \nabla f(\mathbf{x}^{(t)})^\top \mathbf{u} = \arg \min_{\mathbf{u}, \|\mathbf{u}\|=1} \|\nabla f(\mathbf{x}^{(t)})\| \|\mathbf{u}\| \cos \theta$$

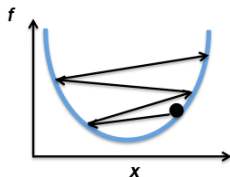
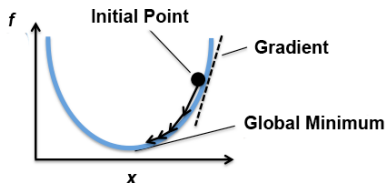
where  $\theta$  is the angle between  $\mathbf{u}$  and  $\nabla f(\mathbf{x}^{(t)})$

- This amounts to solve

$$\arg \min_{\mathbf{u}} \cos \theta$$

- So,  $\mathbf{u}^* = -\nabla f(\mathbf{x}^{(t)})$  is the steepest descent direction of  $f$  at point  $\mathbf{x}^{(t)}$

# How to Set Learning Rate $\eta$ ? I



- Too small an  $\eta$  results in slow descent speed and many iterations
- Too large an  $\eta$  may overshoot the optimal point along the gradient and goes uphill
- One way to set a better  $\eta$  is to leverage the *curvatures* of  $f$ 
  - The more curvy  $f$  at point  $x^{(t)}$ , the smaller the  $\eta$



# How to Set Learning Rate $\eta$ ? II

- By Taylor's theorem, we can approximate  $f$  locally at point  $\mathbf{x}^{(t)}$  using a quadratic function  $\tilde{f}$ :

$$f(\mathbf{x}) \approx \tilde{f}(\mathbf{x}; \mathbf{x}^{(t)}) = f(\mathbf{x}^{(t)}) + \nabla f(\mathbf{x}^{(t)})^\top (\mathbf{x} - \mathbf{x}^{(t)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(t)})^\top \mathbf{H}(f)(\mathbf{x}^{(t)}) (\mathbf{x} - \mathbf{x}^{(t)})$$

for  $\mathbf{x}$  close enough to  $\mathbf{x}^{(t)}$

- $\mathbf{H}(f)(\mathbf{x}^{(t)}) \in \mathbb{R}^{d \times d}$  is the (symmetric) Hessian matrix of  $f$  at  $\mathbf{x}^{(t)}$
- Line search at step  $t$ :

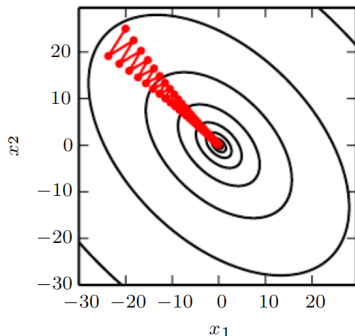
$$\begin{aligned} \arg \min_{\eta} \tilde{f}(\mathbf{x}^{(t)} - \eta \nabla f(\mathbf{x}^{(t)})) = \\ \arg \min_{\eta} f(\mathbf{x}^{(t)}) - \eta \nabla f(\mathbf{x}^{(t)})^\top \nabla f(\mathbf{x}^{(t)}) + \frac{\eta^2}{2} \nabla f(\mathbf{x}^{(t)})^\top \mathbf{H}(f)(\mathbf{x}^{(t)}) \nabla f(\mathbf{x}^{(t)}) \end{aligned}$$

- If  $\nabla f(\mathbf{x}^{(t)})^\top \mathbf{H}(f)(\mathbf{x}^{(t)}) \nabla f(\mathbf{x}^{(t)}) > 0$ , we can solve  $\frac{\partial}{\partial \eta} \tilde{f}(\mathbf{x}^{(t)} - \eta \nabla f(\mathbf{x}^{(t)})) = 0$  and get:

$$\eta^{(t)} = \frac{\nabla f(\mathbf{x}^{(t)})^\top \nabla f(\mathbf{x}^{(t)})}{\nabla f(\mathbf{x}^{(t)})^\top \mathbf{H}(f)(\mathbf{x}^{(t)}) \nabla f(\mathbf{x}^{(t)})}$$

# Problems of Gradient Descent

- Gradient descent is designed to find the steepest descent direction at step  $\mathbf{x}^{(t)}$ 
  - **Not aware of the conditioning** of the Hessian matrix  $\mathbf{H}(f)(\mathbf{x}^{(t)})$
- If  $\mathbf{H}(f)(\mathbf{x}^{(t)})$  has a large condition number, then  $f$  is curvy in some directions but flat in others at  $\mathbf{x}^{(t)}$
- E.g., suppose  $f$  is a quadratic function whose Hessian has a large condition number
- A step in gradient descent may overshoot the optimal points along flat attributes
  - “Zig-zags” around a narrow valley
- Why not take conditioning into account when picking descent directions?



# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ **Unconstrained Optimization**
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

# Newton's Method I

- By Taylor's theorem, we can approximate  $f$  locally at point  $\mathbf{x}^{(t)}$  using a quadratic function  $\tilde{f}$ , i.e.,

$$f(\mathbf{x}) \approx \tilde{f}(\mathbf{x}; \mathbf{x}^{(t)}) = f(\mathbf{x}^{(t)}) + \nabla f(\mathbf{x}^{(t)})^\top (\mathbf{x} - \mathbf{x}^{(t)}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(t)})^\top \mathbf{H}(f)(\mathbf{x}^{(t)}) (\mathbf{x} - \mathbf{x}^{(t)})$$

for  $\mathbf{x}$  close enough to  $\mathbf{x}^{(t)}$

- If  $f$  is strictly convex (i.e.,  $\mathbf{H}(f)(\mathbf{a}) \succ \mathbf{O}, \forall \mathbf{a}$ ), we can find  $\mathbf{x}^{(t+1)}$  that minimizes  $\tilde{f}$  in order to decrease  $f$
- Solving  $\nabla \tilde{f}(\mathbf{x}^{(t+1)}; \mathbf{x}^{(t)}) = \mathbf{0}$ , we have

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \mathbf{H}(f)(\mathbf{x}^{(t)})^{-1} \nabla f(\mathbf{x}^{(t)})$$

- $\mathbf{H}(f)(\mathbf{x}^{(t)})^{-1}$  as a “corrector” to the negative gradient

# Newton's Method II

---

**Input:**  $\mathbf{x}^{(0)} \in \mathbb{R}^d$  an initial guess,  $\eta > 0$

**repeat**

$\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \eta \mathbf{H}(f)(\mathbf{x}^{(t)})^{-1} \nabla f(\mathbf{x}^{(t)})$  ;

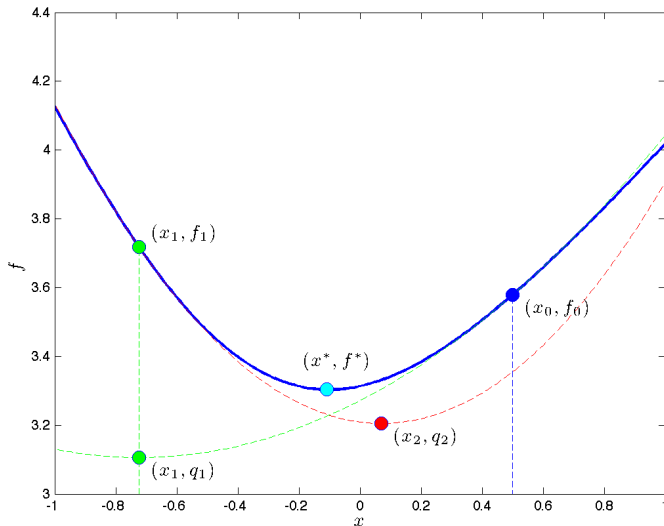
**until** *convergence criterion is satisfied*;

---

- In practice, we multiply the shift by a small  $\eta > 0$  to make sure that  $\mathbf{x}^{(t+1)}$  is close to  $\mathbf{x}^{(t)}$

# Newton's Method III

- If  $f$  is positive definite quadratic, then only one step is required



# General Functions

- Update rule:  $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \eta \mathbf{H}(f)(\mathbf{x}^{(t)})^{-1} \nabla f(\mathbf{x}^{(t)})$
- What if  $f$  is not strictly convex?
  - $\mathbf{H}(f)(\mathbf{x}^{(t)}) \preceq \mathbf{O}$  or indefinite
- The **Levenberg–Marquardt extension**:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \left( \mathbf{H}(f)(\mathbf{x}^{(t)}) + \alpha \mathbf{I} \right)^{-1} \nabla f(\mathbf{x}^{(t)}) \text{ for some } \alpha > 0$$

- With a large  $\alpha$ , degenerates into gradient descent of learning rate  $1/\alpha$

---

**Input:**  $\mathbf{x}^{(0)} \in \mathbb{R}^d$  an initial guess,  $\eta > 0$ ,  $\alpha > 0$

**repeat**

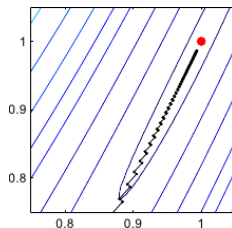
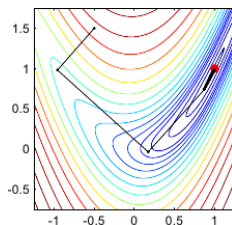
$\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \eta \left( \mathbf{H}(f)(\mathbf{x}^{(t)}) + \alpha \mathbf{I} \right)^{-1} \nabla f(\mathbf{x}^{(t)})$  ;

**until** *convergence criterion is satisfied*;

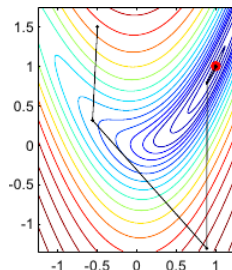
---

# Gradient Descent vs. Newton's Method

- Steps of Gradient descent when  $f$  is a Rosenbrock's banana:



- Steps of Newton's method:
  - Only 6 steps in total





# Problems of Newton's Method

- Computing  $\mathbf{H}(f)(\mathbf{x}^{(t)})^{-1}$  is slow
  - Takes  $O(d^3)$  time at each step, which is ***much slower*** than  $O(d)$  of gradient descent
- Imprecise  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \mathbf{H}(f)(\mathbf{x}^{(t)})^{-1} \nabla f(\mathbf{x}^{(t)})$  due to numerical errors
  - $\mathbf{H}(f)(\mathbf{x}^{(t)})$  may have a large condition number
- Attracted to ***saddle points*** (when  $f$  is not convex)
  - The  $\mathbf{x}^{(t+1)}$  solved from  $\nabla \tilde{f}(\mathbf{x}^{(t+1)}; \mathbf{x}^{(t)}) = \mathbf{0}$  is a critical point

# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent**
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

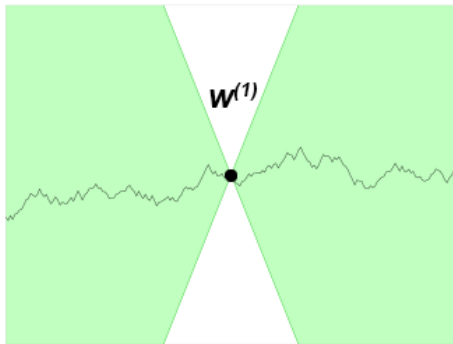
# Who is Afraid of Non-convexity?

- In ML, the function to solve is usually the cost function  $C(\mathbf{w})$  of a model  $\mathbb{F} = \{f : f \text{ parametrized by } \mathbf{w}\}$
- Many ML models have convex cost functions in order to take advantages of convex optimization
  - E.g., perceptron, linear regression, logistic regression, SVMs, etc.
- However, in deep learning, the cost function of a neural network is typically **not** convex
  - We will discuss techniques that tackle non-convexity later

# Assumption on Cost Functions

- In ML, we usually assume that the (real-valued) cost function is *Lipschitz continuous* and/or have Lipschitz continuous derivatives
- I.e., the rate of change of  $C$  is bounded by a *Lipschitz constant*  $K$ :

$$|C(\mathbf{w}^{(1)}) - C(\mathbf{w}^{(2)})| \leq K \|\mathbf{w}^{(1)} - \mathbf{w}^{(2)}\|, \forall \mathbf{w}^{(1)}, \mathbf{w}^{(2)}$$

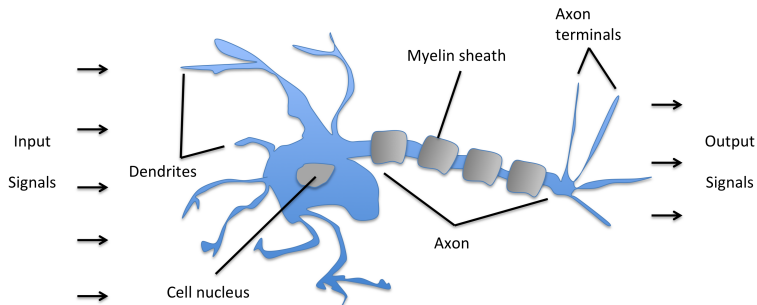


# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent**
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

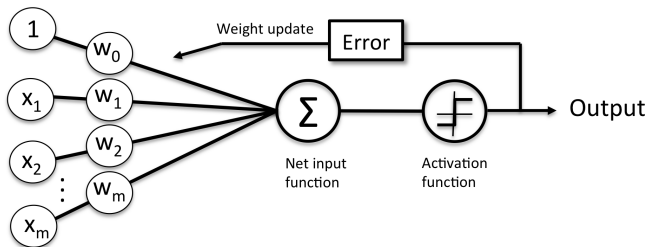
# Perceptron & Neurons

- **Perceptron**, proposed in 1950's by Rosenblatt, is one of the first ML algorithms for binary classification
- Inspired by McCulloch-Pitts (MCP) neuron, published in 1943
  - Our brains consist of interconnected **neurons**
  - Each neuron takes signals from other neurons as input
  - If the accumulated signal exceeds a certain threshold, an output signal is generated



# Model

- Binary classification problem:
  - Training dataset:  $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_i$ , where  $\mathbf{x}^{(i)} \in \mathbb{R}^D$  and  $y^{(i)} \in \{1, -1\}$
  - Output: a function  $f(\mathbf{x}) = \hat{y}$  such that  $\hat{y}$  is close to the true label  $y$
- Model:  $\{f : f(\mathbf{x}; \mathbf{w}, b) = \text{sign}(\mathbf{w}^\top \mathbf{x} - b)\}$ 
  - $\text{sign}(a) = 1$  if  $a \geq 0$ ; otherwise 0
  - For simplicity, we use shorthand  $f(\mathbf{x}; \mathbf{w}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$  where  $\mathbf{w} = [-b, w_1, \dots, w_D]^\top$  and  $\mathbf{x} = [1, x_1, \dots, x_D]^\top$



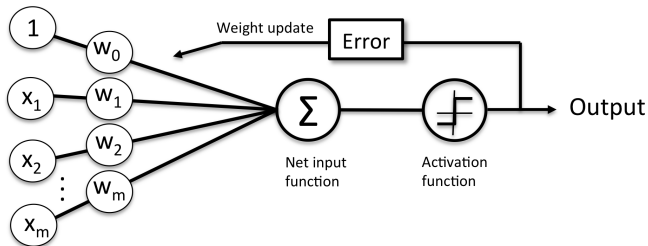
# Iterative Training Algorithm I

- 1 Initiate  $\mathbf{w}^{(0)}$  and learning rate  $\eta > 0$
- 2 Epoch: for each example  $(\mathbf{x}^{(t)}, y^{(t)})$ , update  $\mathbf{w}$  by

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta(y^{(t)} - \hat{y}^{(t)})\mathbf{x}^{(t)}$$

where  $\hat{y}^{(t)} = f(\mathbf{x}^{(t)}; \mathbf{w}^{(t)}) = \text{sign}(\mathbf{w}^{(t)\top} \mathbf{x}^{(t)})$

- 3 Repeat epoch several times (or until converge)



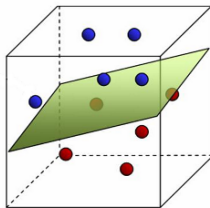


# Iterative Training Algorithm II

- Update rule:

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta(y^{(t)} - \hat{y}^{(t)})\mathbf{x}^{(t)}$$

- If  $\hat{y}^{(t)}$  is correct, we have  $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)}$
- If  $\hat{y}^{(t)}$  is incorrect, we have  $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + 2\eta y^{(t)}\mathbf{x}^{(t)}$ 
  - If  $y^{(t)} = 1$ , the updated prediction will more likely to be positive, as  $\text{sign}(\mathbf{w}^{(t+1)\top}\mathbf{x}^{(t)}) = \text{sign}(\mathbf{w}^{(t)\top}\mathbf{x}^{(t)} + c)$  for some  $c > 0$
  - If  $y^{(t)} = -1$ , the updated prediction will more likely to be negative
- Does **not** converge if the dataset cannot be separated by a hyperplane



# Outline

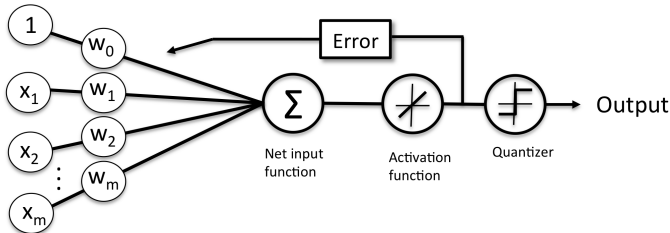
- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ **Optimization in ML: Stochastic Gradient Descent**
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

# ADaptive LInear NEuron (Adaline)

- Proposed in 1960's by Widrow et al.
- Defines and minimizes a **cost function** for training:

$$\arg \min_{\mathbf{w}} C(\mathbf{w}; \mathbb{X}) = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N \left( y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)} \right)^2$$

- Links numerical optimization to ML
- Sign function is only used for binary prediction **after** training



# Training Using Gradient Descent

- Update rule:

$$\begin{aligned}\mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \eta \nabla C(\mathbf{w}^{(t)}) \\ &= \mathbf{w}^{(t)} + \eta \sum_i (y^{(i)} - \mathbf{w}^{(t)\top} \mathbf{x}^{(i)}) \mathbf{x}^{(i)}\end{aligned}$$

- Since the cost function is convex, the training iterations will converge

# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent**
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

# Cost as an Expectation

- In ML, the cost function to minimize is usually a sum of *losses* over training examples
- E.g., in Adaline: sum of square losses (functions)

$$\arg \min_{\mathbf{w}} C(\mathbf{w}; \mathbb{X}) = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N \left( y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)} \right)^2$$

- Let examples be i.i.d. samples of random variables  $(\mathbf{x}, y)$
- We effectively minimize the estimate of  $\mathbf{E}[C(\mathbf{w})]$  *over the distribution  $P(\mathbf{x}, y)$* :

$$\arg \min_{\mathbf{w}} \mathbf{E}_{\mathbf{x}, y \sim P} [C(\mathbf{w})]$$

- $P(\mathbf{x}, y)$  may be unknown
- Since the problem is stochastic by nature, why not make the training algorithm stochastic too?

# Stochastic Gradient Descent

---

Input:  $\mathbf{w}^{(0)} \in \mathbb{R}^d$  an initial guess,  $\eta > 0$ ,  $M \geq 1$

repeat

    epoch:

        Randomly partition the training set  $\mathbb{X}$  into the *minibatches*

$\{\mathbb{X}^{(j)}\}_j$ ,  $|\mathbb{X}^{(j)}| = M$ ;

        foreach  $j$  do

$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla C(\mathbf{w}^{(t)}; \mathbb{X}^{(j)})$  ;

        end

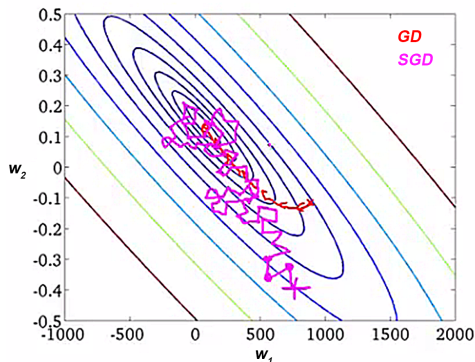
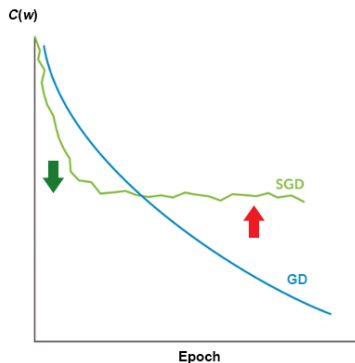
until *convergence criterion is satisfied*;

---

- $C(\mathbf{w}; \mathbb{X}^{(j)})$  is still an estimate of  $\mathbb{E}_{\mathbf{x}, y \sim P}[C(\mathbf{w})]$ 
  - $\mathbb{X}^{(j)}$  are samples of the same distribution  $P(\mathbf{x}, y)$
- It's common to set  $M = 1$  on a single machine
  - E.g., update rule for Adaline:  $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta (y^{(t)} - \mathbf{w}^{(t)\top} \mathbf{x}^{(t)}) \mathbf{x}^{(t)}$ , which is similar to that of Perceptron

# SGD vs. GD

- Each iteration can run *much faster* when  $M \ll N$
- Converges faster (in both #epochs and time) with large datasets
- Supports *online* learning
- But may wander around the optimal points
  - In practice, we set  $\eta = O(t^{-1})$





# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization**
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

# Constrained Optimization

- Problem:

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

subject to  $\mathbf{x} \in \mathbb{C}$

- $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is not necessarily convex
- $\mathbb{C} = \{\mathbf{x} : g^{(i)}(\mathbf{x}) \leq 0, h^{(j)}(\mathbf{x}) = 0\}_{i,j}$
- Iterative descent algorithm?

# Common Methods

- **Projective gradient descent**: if  $\mathbf{x}^{(t)}$  falls outside  $\mathbb{C}$  at step  $t$ , we “project” back the point to the tangent space (edge) of  $\mathbb{C}$
- **Penalty/barrier methods**: convert the constrained problem into one or more unconstrained ones
- And more...

# Karush-Kuhn-Tucker (KKT) Methods I

- Converts the problem

$$\begin{array}{ll} \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{subject to } \mathbf{x} \in \{\mathbf{x} : g^{(i)}(\mathbf{x}) \leq 0, h^{(j)}(\mathbf{x}) = 0\}_{i,j} \end{array}$$

into

$$\begin{array}{l} \min_{\mathbf{x}} \max_{\alpha, \beta, \alpha \geq 0} L(\mathbf{x}, \alpha, \beta) = \\ \min_{\mathbf{x}} \max_{\alpha, \beta, \alpha \geq 0} f(\mathbf{x}) + \sum_i \alpha_i g^{(i)}(\mathbf{x}) + \sum_j \beta_j h^{(j)}(\mathbf{x}) \end{array}$$

- $\min_{\mathbf{x}} \max_{\alpha, \beta} L$  means “minimize  $L$  with respect to  $\mathbf{x}$ , at which  $L$  is maximized with respect to  $\alpha$  and  $\beta$ ”
- The function  $L(\mathbf{x}, \alpha, \beta)$  is called the (generalized) **Lagrangian**
- $\alpha$  and  $\beta$  are called **KKT multipliers**

# Karush-Kuhn-Tucker (KKT) Methods II

- Converts the problem

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \{\mathbf{x} : g^{(i)}(\mathbf{x}) \leq 0, h^{(j)}(\mathbf{x}) = 0\}_{i,j} \end{aligned}$$

into

$$\begin{aligned} & \min_{\mathbf{x}} \max_{\alpha, \beta, \alpha \geq 0} L(\mathbf{x}, \alpha, \beta) = \\ & \min_{\mathbf{x}} \max_{\alpha, \beta, \alpha \geq 0} f(\mathbf{x}) + \sum_i \alpha_i g^{(i)}(\mathbf{x}) + \sum_j \beta_j h^{(j)}(\mathbf{x}) \end{aligned}$$

- Observe that for any feasible point  $\mathbf{x}$ , we have

$$\max_{\alpha, \beta, \alpha \geq 0} L(\mathbf{x}, \alpha, \beta) = f(\mathbf{x})$$

- The optimal feasible point is unchanged
- And for any infeasible point  $\mathbf{x}$ , we have

$$\max_{\alpha, \beta, \alpha \geq 0} L(\mathbf{x}, \alpha, \beta) = \infty$$

- Infeasible points will never be optimal (if there are feasible points)

# Alternate Iterative Algorithm

$$\min_{\mathbf{x}} \max_{\alpha, \beta, \alpha \geq \mathbf{0}} f(\mathbf{x}) + \sum_i \alpha_i g^{(i)}(\mathbf{x}) + \sum_j \beta_j h^{(j)}(\mathbf{x})$$

- “Large”  $\alpha$  and  $\beta$  create a “barrier” for feasible solutions

---

**Input:**  $\mathbf{x}^{(0)}$  an initial guess,  $\alpha^{(0)} = \mathbf{0}$ ,  $\beta^{(0)} = \mathbf{0}$

**repeat**

    Solve  $\mathbf{x}^{(t+1)} = \arg \min_{\mathbf{x}} L(\mathbf{x}; \alpha^{(t)}, \beta^{(t)})$  using some iterative algorithm starting at  $\mathbf{x}^{(t)}$ ;

**if**  $\mathbf{x}^{(t+1)} \notin \mathbb{C}$  **then**

        Increase  $\alpha^{(t)}$  to get  $\alpha^{(t+1)}$ ;

        Get  $\beta^{(t+1)}$  by increasing the magnitude of  $\beta^{(t)}$  and set  $\text{sign}(\beta_j^{(t+1)}) = \text{sign}(h^{(j)}(\mathbf{x}^{(t+1)}))$ ;

**end**

**until**  $\mathbf{x}^{(t+1)} \in \mathbb{C}$ ;

# KKT Conditions

## Theorem (KKT Conditions)

If  $\mathbf{x}^*$  is an optimal point, then there exists KKT multipliers  $\alpha^*$  and  $\beta^*$  such that the **Karush-Kuhn-Tucker (KKT) conditions** are satisfied:

Lagrangian stationarity:  $\nabla L(\mathbf{x}^*, \alpha^*, \beta^*) = 0$

Primal feasibility:  $g^{(i)}(\mathbf{x}^*) \leq 0$  and  $h^{(j)}(\mathbf{x}^*) = 0$  for all  $i$  and  $j$

Dual feasibility:  $\alpha^* \geq 0$

**Complementary slackness:**  $\alpha_i^* g^{(i)}(\mathbf{x}^*) = 0$  for all  $i$ .

- Only a necessary condition for  $\mathbf{x}^*$  being optimal
- Sufficient if the original problem is **convex**

# Complementary Slackness

- Why  $\alpha_i^* g^{(i)}(\mathbf{x}^*) = 0$ ?
- For  $\mathbf{x}^*$  being feasible, we have  $g^{(i)}(\mathbf{x}^*) \leq 0$
- If  $g^{(i)}$  is **active** (i.e.,  $g^{(i)}(\mathbf{x}^*) = 0$ ) then  $\alpha_i^* g^{(i)}(\mathbf{x}^*) = 0$
- If  $g^{(i)}$  is **inactive** (i.e.,  $g^{(i)}(\mathbf{x}^*) < 0$ ), then
  - To maximize the  $\alpha_i g^{(i)}(\mathbf{x}^*)$  term in the Lagrangian in terms of  $\alpha_i$  subject to  $\alpha_i \geq 0$ , we have  $\alpha_i^* = 0$
  - Again  $\alpha_i^* g^{(i)}(\mathbf{x}^*) = 0$
- So what?
  - $\alpha_i^* > 0$  implies  $g^{(i)}(\mathbf{x}^*) = 0$
  - Once  $\mathbf{x}^*$  is solved, we can quickly find out the active inequality constrains by checking  $\alpha_i^* > 0$



# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization**
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization**
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

# The Regression Problem

- Given a training dataset:  $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ 
  - $\mathbf{x}^{(i)} \in \mathbb{R}^D$ , called explanatory variables (attributes/features)
  - $y^{(i)} \in \mathbb{R}$ , called response/target variables (labels)
- Goal: to find a function  $f(\mathbf{x}) = \hat{y}$  such that  $\hat{y}$  is close to the true label  $y$
- Example: to predict the price of a stock tomorrow
- Could you define a model  $\mathbb{F} = \{f\}$  and cost function  $C[f]$ ?
- How about “relaxing” the Adaline by removing the sign function when making the final prediction?
  - Adaline:  $\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} - b)$
  - Regressor:  $\hat{y} = \mathbf{w}^\top \mathbf{x} - b$

# Linear Regression I

- Model:  $\mathbb{F} = \{f : f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^\top \mathbf{x} - b\}$ 
  - Shorthand:  $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x}$ , where  $\mathbf{w} = [-b, w_1, \dots, w_D]^\top$  and  $\mathbf{x} = [1, x_1, \dots, x_D]^\top$
- Cost function and optimization problem:

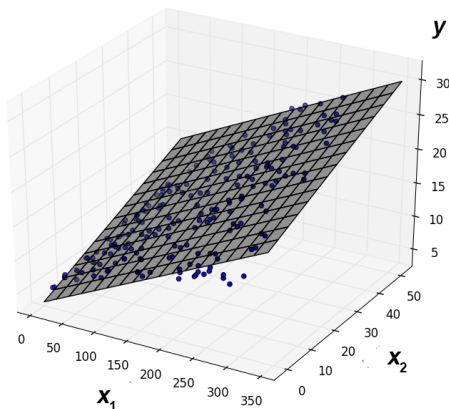
$$\arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N \|y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)}\|^2 = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

- $\mathbf{X} = \begin{bmatrix} 1 & \mathbf{x}^{(1)\top} \\ \vdots & \vdots \\ 1 & \mathbf{x}^{(N)\top} \end{bmatrix} \in \mathbb{R}^{N \times (D+1)}$  the design matrix
- $\mathbf{y} = [y^{(1)}, \dots, y^{(N)}]^\top$  the label vector

# Linear Regression II

$$\arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N \|\mathbf{y}^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)}\|^2 = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

- Basically, we fit a hyperplane to training data
  - Each  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} - b \in \mathbb{R}$  is a hyperplane in the graph



# Training Using Gradient Descent

$$\arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N \|\mathbf{y}^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)}\|^2 = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

- Batch:

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta \sum_{i=1}^N (\mathbf{y}^{(i)} - \mathbf{w}^{(t)\top} \mathbf{x}^{(i)}) \mathbf{x}^{(i)} = \mathbf{w}^{(t)} + \eta \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$$

- Stochastic (with minibatch size  $|M| = 1$ ):

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta (\mathbf{y}^{(t)} - \mathbf{w}^{(t)\top} \mathbf{x}^{(t)}) \mathbf{x}^{(t)}$$

# Evaluation Metrics of Regression Models

- Given a training/testing set  $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
- How to evaluate the predictions  $\hat{y}^{(i)}$  made by a function  $f$ ?
- Sum of Square Errors (SSE):  $\sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2$
- Mean Square Error (MSE):  $\frac{1}{N} \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2$
- Relative Square Error (RSE):

$$\frac{\sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2}{\sum_{i=1}^N (y^{(i)} - \bar{y})^2},$$

where  $\bar{y} = \frac{1}{N} \sum_i y^{(i)}$

- What does it mean? Compares  $f$  with a dummy prediction  $\bar{y}$
- Coefficient of Determination:  $R^2 = 1 - RSE \in [0, 1]$ 
  - Higher the better

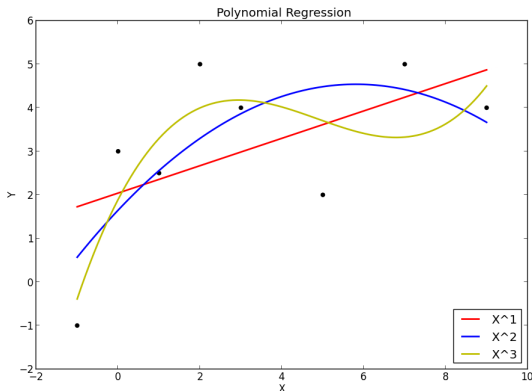
# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization**
  - Linear Regression
  - Polynomial Regression**
  - Generalizability & Regularization
- ⑦ Duality\*



# Polynomial Regression

- In practice, the relationship between explanatory variables and target variables may not be linear
- **Polynomial regression** fits a high-order polynomial to the training data
- How?



# Data Augmentation

- Suppose  $D = 2$ , i.e.,  $\mathbf{x} = [x_1, x_2]^\top$
- Linear model:

$$\mathbb{F} = \{f : f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x} + w_0 = w_0 + w_1 x_1 + w_2 x_2\}$$

- Quadratic model:

$$\mathbb{F} = \{f : f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_1 x_2 + w_5 x_2^2\}$$

- We can simply **augment** the data dimension to reduce a quadratic model to a linear one
  - A general technique in ML to “transform” a linear model into a nonlinear one
- How many variables to solve in  $\mathbf{w}$  for a polynomial regression problem of degree  $P$ ? [Homework]

# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization**
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\*

# Regularization

- There's another major difference between the ML algorithms and optimization techniques:
- We usually care about the *testing performance* rather than the training performance
  - E.g., in classification, we report the testing accuracy
- Goal: to learn a function that *generalizes* to unseen data well
- *Regularization*: techniques that improve the generalizability of the learned function
- How to regularize the linear regression?

$$\arg \min_w \frac{1}{2} \|y - Xw\|^2$$

# Regularized Linear Regression

- One way to improve the generalizability of  $f$  is to make it “flat:”

$$\begin{aligned} \arg \min_{\mathbf{w} \in \mathbb{R}^D, b} \frac{1}{2} \|\mathbf{y} - (X\mathbf{w} - b)\|^2 \\ \text{subject to } \|\mathbf{w}\|^2 \leq T \end{aligned} = \begin{aligned} \arg \min_{\mathbf{w} \in \mathbb{R}^{D+1}} \frac{1}{2} \|\mathbf{y} - (X\mathbf{w})\|^2 \\ \text{subject to } \mathbf{w}^\top \mathbf{S} \mathbf{w} \leq T \end{aligned}$$

- $\mathbf{S} = \text{diag}([0, 1, \dots, 1]^\top) \in \mathbb{R}^{(D+1) \times (D+1)}$  ( $b$  is not regularized)
- We will explain why this works later
- How to solve this problem?
- Using the KKT method, we have

$$\arg \min_{\mathbf{w}} \max_{\alpha, \alpha \geq 0} L(\mathbf{w}, \alpha) = \arg \min_{\mathbf{w}} \max_{\alpha, \alpha \geq 0} \frac{1}{2} \left( \|\mathbf{y} - X\mathbf{w}\|^2 + \alpha(\mathbf{w}^\top \mathbf{S} \mathbf{w} - T) \right)$$

# Alternate Iterative Algorithm

$$\arg \min_{\mathbf{w}} \max_{\alpha, \alpha \geq 0} L(\mathbf{w}, \alpha) = \arg \min_{\mathbf{w}} \max_{\alpha, \alpha \geq 0} \frac{1}{2} \left( \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \alpha(\mathbf{w}^\top \mathbf{S}\mathbf{w} - T) \right)$$

---

**Input:**  $\mathbf{w}^{(0)} \in \mathbb{R}^d$  an initial guess,  $\alpha^{(0)} = 0$ ,  $\delta > 0$

**repeat**

    Solve  $\mathbf{w}^{(t+1)} = \arg \min_{\mathbf{w}} L(\mathbf{w}; \alpha^{(t)})$  using some iterative algorithm starting at  $\mathbf{w}^{(t)}$ ;

**if**  $\mathbf{w}^{(t+1)\top} \mathbf{w}^{(t+1)} > T$  **then**

$\alpha^{(t+1)} = \alpha^{(t)} + \delta$  in order to increase  $L(\alpha; \mathbf{w}^{(t+1)})$ ;

**end**

**until**  $\mathbf{w}^{(t+1)\top} \mathbf{w}^{(t+1)} \leq T$ ;

---

- We could also solve  $\mathbf{w}^{(t+1)}$  analytically from  $\frac{\partial}{\partial \mathbf{x}} L(\mathbf{w}; \alpha^{(t)}) = \mathbf{0}$ :

$$\mathbf{w}^{(t+1)} = \left( \mathbf{X}^\top \mathbf{X} + \alpha^{(t)} \mathbf{S} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

# Outline

- ① Numerical Computation
- ② Optimization Problems
- ③ Unconstrained Optimization
  - Gradient Descent
  - Newton's Method
- ④ Optimization in ML: Stochastic Gradient Descent
  - Perceptron
  - Adaline
  - Stochastic Gradient Descent
- ⑤ Constrained Optimization
- ⑥ Optimization in ML: Regularization
  - Linear Regression
  - Polynomial Regression
  - Generalizability & Regularization
- ⑦ Duality\***

# Dual Problem

- Given a problem (called *primal problem*):

$$p^* = \min_{\mathbf{x}} \max_{\alpha, \beta, \alpha \geq 0} L(\mathbf{x}, \alpha, \beta)$$

- We define its *dual problem* as:

$$d^* = \max_{\alpha, \beta, \alpha \geq 0} \min_{\mathbf{x}} L(\mathbf{x}, \alpha, \beta)$$

- By the max-min inequality, we have  $d^* \leq p^*$  [Homework]
- $(p^* - d^*)$  is called the *duality gap*
  - $p^*$  and  $d^*$  are called the *primal* and *dual values*, respectively



# Strong Duality

- *Strong duality* holds if  $d^* = p^*$
- When will it happen?
- If the primal problem has solution and *convex*
- Why considering dual problem?

## Example

- Consider a primal problem:

$$\begin{array}{ll} \arg \min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{x}\|^2 \\ \text{subject to } \mathbf{Ax} \geq \mathbf{b}, \mathbf{A} \in \mathbb{R}^{n \times d} \end{array} = \arg \min_{\mathbf{x}} \max_{\alpha, \alpha \geq \mathbf{0}} \frac{1}{2} \|\mathbf{x}\|^2 - \alpha^\top (\mathbf{Ax} - \mathbf{b})$$

- Convex, so strong duality holds
- We can get the same solution via the dual problem:

$$\arg \max_{\alpha, \alpha \geq \mathbf{0}} \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x}\|^2 - \alpha^\top (\mathbf{Ax} - \mathbf{b})$$

- Solving  $\min_{\mathbf{x}} L(\mathbf{x}, \alpha)$  analytically, we have  $\mathbf{x}^* = \mathbf{A}^\top \alpha$
- Substituting this into the dual, we get

$$\arg \max_{\alpha, \alpha \geq \mathbf{0}} -\frac{1}{2} \|\mathbf{A}^\top \alpha\|^2 + \mathbf{b}^\top \alpha$$

- We now solve  $n$  variables instead of  $d$  (beneficial when  $n \ll d$ )