

# Convex Optimization

## Lecture 6 - Applications in Machine Learning

Instructor: Yuanzhang Xiao

University of Hawaii at Manoa

Fall 2017

# Today's Lecture

- ① Regression / Prediction
- ② Classification
- ③ Unsupervised Learning
- ④ Reinforcement Learning

# Outline

① Regression / Prediction

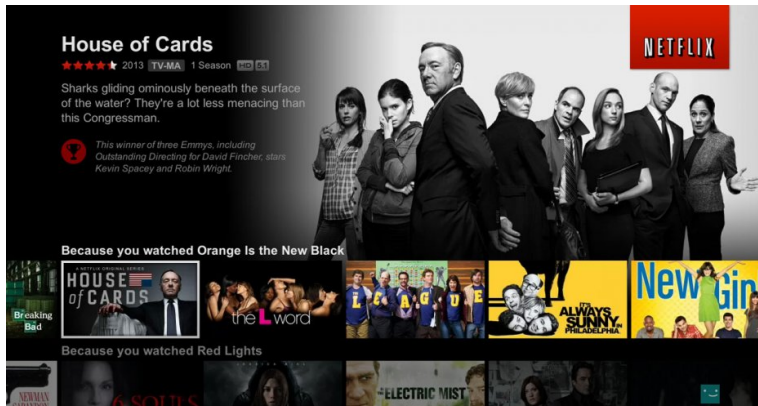
② Classification

③ Unsupervised Learning

④ Reinforcement Learning


# Supervised Learning - Examples

Netflix recommendation systems:



**House of Cards**  
★★★★★ 2013 TV-MA 1 Season HD 5.1

Sharks gliding ominously beneath the surface of the water? They're a lot less menacing than this Congressman.

 This winner of three Emmys, including Outstanding Directing for David Fincher, stars Kevin Spacey and Robin Wright.

**Because you watched Orange Is the New Black**

- Breaking Bad
- HOUSE of CARDS
- the L word
- LEAGUE
- ALWAYS SUNNY IN PHILADELPHIA
- New Girl

**Because you watched Red Lights**

- NEWMAN
- 6 SOULS
- ELECTRIC MIST
- Netflix logo

# Supervised Learning

basic elements:

- $a^{(i)} \in \mathbb{R}^n$ : **features**
  - gender, occupation, income, zip code
- $b^{(i)} \in \mathbb{R}$ : **target**
  - ratings of the movie, watched a movie or not
- $\{(a^{(i)}, b^{(i)}) \mid i = 1, \dots, m\}$ : **training set**

find **hypothesis**  $h : a \mapsto b$

- $b$  continuous (e.g., rating): **regression**
- $b$  discrete (e.g., watched or not): **classification**

# Linear Regression

linear hypothesis:

$$h_x(a) = x_1 a_1 + \cdots + x_n a_n$$

where  $x$  are **weights**

choose  $x$  to minimize **cost function**:

$$f_0(x) : \mathbb{R}^n \rightarrow \mathbb{R}$$

many choices of different cost functions  $\rightarrow$  different methods

# Ordinary Least Squares Regression

ordinary least squares regression:

$$\text{minimize } f_0(x) = \|Ax - b\|_2$$

where

$$A = \begin{bmatrix} a^{(1)T} \\ \vdots \\ a^{(m)T} \end{bmatrix}, \quad b = \begin{bmatrix} b^{(1)} \\ \vdots \\ b^{(m)} \end{bmatrix}$$

- convex optimization problem
- normal equation

$$A^T A x = A^T b$$

- if  $\text{rank}A = n$ , we have

$$x^* = (A^T A)^{-1} A^T b$$

# Chebyshev Regression

Chebyshev regression:

$$\text{minimize } f_0(x) = \|Ax - b\|_\infty$$

where  $\|y\|_\infty = \max\{|y_1|, \dots, |y_m|\}$

- convex optimization problem (may be hard to solve)
- equivalent formulation: a linear program

$$\begin{aligned} &\text{minimize } t \\ &\text{subject to } -t\mathbf{1} \leq Ax - b \leq t\mathbf{1} \end{aligned}$$

with optimization variables  $x$  and  $t \in \mathbb{R}$



# Sum of Absolute Residuals Regression

sum of absolute residuals regression:

$$\text{minimize } f_0(x) = \|Ax - b\|_1$$

where  $\|y\|_1 = \sum_{i=1}^m |y_i|$

- convex optimization problem (may be hard to solve)
- equivalent formulation: a linear program

$$\text{minimize } \mathbf{1}^T t$$

$$\text{subject to } -t \leq Ax - b \leq t$$

with optimization variables  $x$  and  $t \in \mathbb{R}^m$

# Regularized Regression

regularized regression:

$$\text{minimize } f_0(x) = \|Ax - b\| + \gamma\|x\|$$

- forces  $x$  to be small
- less sensitive to errors in features  $A$
- select few important features
  
- convex optimization problem (may be hard to solve)
- $\|Ax - b\|_1 + \gamma\|x\|_1$  can be reformulated as LP
- $\|Ax - b\|_2 + \gamma\|x\|_1$  can be reformulated as QP

# Regression Regularized by Cardinality

regression regularized by cardinality:

$$\text{minimize } f_0(x) = \|Ax - b\| + \text{card}(x)$$

where  $\text{card}(x)$  is the number of nonzero elements of  $x$

properties of  $\text{card}(\cdot)$ :

- quasiconcave on  $\mathbb{R}_+^n$ :

$$\text{card}(x + y) \geq \min \{ \text{card}(x), \text{card}(y) \}$$

- non-convex

# General Convex-Cardinality Problems

**convex-cardinality problem:** one that would be convex except for the appearance of  $\text{card}(\cdot)$  in objective or constraints

examples:

- regression regularized by cardinality
- minimum cardinality:

$$\begin{aligned} & \text{minimize} && \text{card}(x) \\ & \text{subject to} && \|Ax - b\|_2 \leq \epsilon \end{aligned}$$

- sparse modeling / regressor selection:

$$\begin{aligned} & \text{minimize} && \|Ax - b\|_2 \\ & \text{subject to} && \text{card}(x) \leq k \end{aligned}$$

**NP-hard** problems

# Solutions to General Convex-Cardinality Problems

exact solutions:

- fix sparsity pattern (which elements are nonzero), then solve the resulting convex problem
- $2^n$  convex problems to solve

convex heuristics:

- replace  $\text{card}(x)$  with  $\|x\|_1$

example: minimize  $\|Ax - b\|_2$  subject to  $\text{card}(x) \leq k$

- solve minimize  $\|Ax - b\|_2 + \gamma \text{card}(x)$
- adjust  $\gamma$  so that  $\text{card}(x) \leq k$
- fix sparsity pattern, and solve the resulting convex problem

theoretical guarantee that heuristic is exact for some problems

# Outline

① Regression / Prediction

② Classification

③ Unsupervised Learning

④ Reinforcement Learning

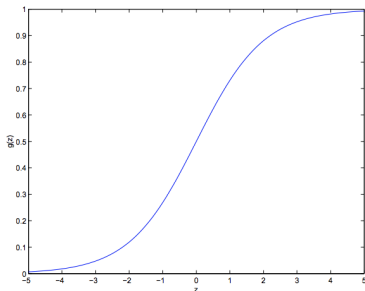
# Logistic Regression

classification:  $b \in \{0, 1\}$

logistic / sigmoid function:

$$h_x(a) = \frac{1}{1 + e^{-x^T a}}$$

illustration of  $g(z) = \frac{1}{1+e^{-z}}$ :



normalize  $x^T a$  to  $[0, 1]$

# Logistic Regression as Convex Optimization Problem

assume that

$$\text{prob}(b = 1|a; x) = h_x(a)$$

$$\text{prob}(b = 0|a; x) = 1 - h_x(a)$$

which can be written compactly as

$$\text{prob}(b|a; x) = [h_x(a)]^b [1 - h_x(a)]^{1-b}$$

assume training examples are independent, then likelihood of  $x$  is

$$L(x) = \text{prob}(b|A; x) = \prod_{i=1}^m [h_x(a^{(i)})]^{b^{(i)}} [1 - h_x(a^{(i)})]^{1-b^{(i)}}$$

easier to maximize the log-likelihood:

$$\ell(x) = \sum_{i=1}^m b^{(i)} \log [h_x(a^{(i)})] + (1 - b^{(i)}) \log [1 - h_x(a^{(i)})]$$



# Logistic Regression as Convex Optimization Problem

rearrange the expressions:

$$\begin{aligned}
 b^{(i)} \log [h_x(a^{(i)})] &= -b^{(i)} \log (1 + e^{-a^{(i)T}x}) \\
 (1 - b^{(i)}) \log [1 - h_x(a^{(i)})] &= (1 - b^{(i)}) \cdot \\
 &\quad [\log (e^{-a^{(i)T}x}) - \log (1 + e^{-a^{(i)T}x})]
 \end{aligned}$$

log-likelihood:

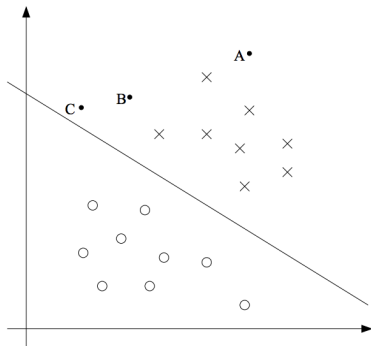
$$\ell(x) = \sum_{i=1}^m -\log (1 + e^{-a^{(i)T}x}) - (1 - b^{(i)}) (a^{(i)T}x)$$

the resulting convex optimization problem:

$$\text{maximize } \ell(x)$$

# Support Vector Machine (SVM)

a geometric view of classification problem:



- separating hyperplane:  $a^T x + y = 0$
- point A: high confidence about the output
- point C: low confidence about the output

# Support Vector Machine (SVM)

assume  $b \in \{-1, 1\}$  (note the difference from logistic regression)

hypothesis function:

$$h_{x,y}(a) = g(a^T x + y)$$

$$\text{where } g(z) = \begin{cases} 1 & z \geq 0 \\ -1 & z < 0 \end{cases}$$

distance between separating hyperplane  $a^T x + y = 0$  and the  $i$ th sample point  $a^{(i)}$ :

$$\gamma_i = b^{(i)} (a^{(i)T} x + y) / \|x\|_2$$

maximize the distance from the closest sample point: (**nonconvex**)

$$\begin{aligned} & \text{maximize } \gamma \\ & \text{subject to } \frac{b^{(i)} (a^{(i)T} x + y)}{\|x\|_2} \geq \gamma, \quad i = 1, \dots, m \end{aligned}$$

# SVM as Convex Optimization

equivalent formulation: (still **nonconvex**)

$$\begin{aligned} & \text{maximize} && \gamma \\ & \text{subject to} && b^{(i)} \left( a^{(i)T} x + y \right) \geq \gamma, \quad i = 1, \dots, m \\ & && \|x\|_2 = 1 \end{aligned}$$

**convex** reformulation:

$$\begin{aligned} & \text{maximize} && \gamma \\ & \text{subject to} && b^{(i)} \left( a^{(i)T} x + y \right) \geq \gamma, \quad i = 1, \dots, m \\ & && \|x\|_2 \leq 1 \end{aligned}$$

more commonly-used **convex** reformulation:

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|x\|_2^2 \\ & \text{subject to} && b^{(i)} \left( a^{(i)T} x + y \right) \geq 1, \quad i = 1, \dots, m \end{aligned}$$

# Dual Problem of SVM

Lagrangian:

$$L(x, y, \lambda) = \frac{1}{2} \|x\|_2^2 - \sum_{i=1}^m \lambda_i \left[ b^{(i)} \left( a^{(i)T} x + y \right) - 1 \right]$$

dual function:

$$g(\lambda) = \inf_{x, y} \frac{1}{2} x^T x - \left( \sum_{i=1}^m \lambda_i b^{(i)} a^{(i)} \right)^T x - \left( \sum_{i=1}^m \lambda_i b^{(i)} \right) y + \sum_{i=1}^m \lambda_i$$

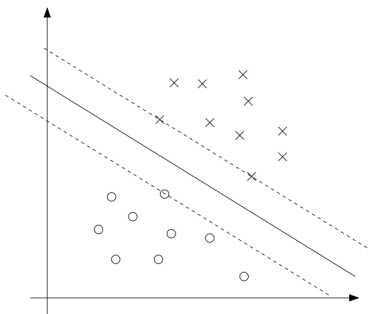
$$\Rightarrow x^*(\lambda) = \sum_{i=1}^m \lambda_i b^{(i)} a^{(i)} = A^T \text{diag}(b) \lambda, \quad \sum_{i=1}^m \lambda_i b^{(i)} = 0$$

dual problem:

$$\begin{aligned} & \text{maximize} && -\frac{1}{2} \lambda^T \text{diag}(b)^T A A^T \text{diag}(b) \lambda + 1^T \lambda \\ & \text{subject to} && \lambda \geq 0 \\ & && \sum_{i=1}^m \lambda_i b^{(i)} = 0 \end{aligned}$$

# Support Vectors

**support vectors:** (the data where the inequality is binding)



**a lot fewer support vectors than training data points**

# Solving SVM Efficiently

new training data  $(a^{(m+1)}, b^{(m+1)})$  coming in:

- if  $b^{(m+1)} (a^{(m+1)T} x^* + y^*) > 1$ , no need to update  $x^*$  and  $y^*$

if data is separable  $\rightarrow$  strong duality  $\rightarrow$  complementary slackness:

$$\lambda_i^* \left[ b^{(i)} \left( a^{(i)T} x^* + y^* \right) - 1 \right] = 0$$

$\Rightarrow \lambda_i^* > 0$  only if  $a^{(i)}$  is a support vector

- sparsity in the solution to the dual problem
- $x^* = \sum_i: a^{(i)} \text{ is a support vector } \lambda_i^* b^{(i)} a^{(i)}$
- classification:

$$a^T x^* + y^* = \sum_{i: a^{(i)} \text{ is a support vector}} \lambda_i^* b^{(i)} \left( a^T a^{(i)} \right) + y^*$$

# Outline

① Regression / Prediction

② Classification

③ Unsupervised Learning

④ Reinforcement Learning



# Unsupervised Learning

basic elements:

- $a^{(i)} \in \mathbb{R}^n$ : features
- no label
- $\{(a^{(i)}) \mid i = 1, \dots, m\}$ : training set

divide the data into  $k$  clusters

common approach:

- initially, “guess”  $k$  labels (i.e., centers of clusters)
- iterate between the following two steps:
  - ① assign data to different clusters
  - ② update the centers of clusters

# k-Means Clustering

k-means clustering:

- randomly select  $k$  cluster centroids  $\mu_1, \dots, \mu_k \in \mathbb{R}^n$
- iterate between the following two steps:
  - ① assign data to different clusters: for each  $j = 1, \dots, k$ ,

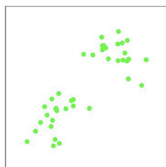
$$\mathcal{C}_j = \left\{ i : \|a^{(i)} - \mu_j\| \leq \|a^{(i)} - \mu_{j'}\|, \forall j' \right\}$$

- ② update the centers of clusters: for each  $j = 1, \dots, k$ ,

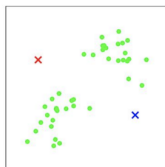
$$\text{minimize } \sum_{i \in \mathcal{C}_j} \|a^{(i)} - \mu_j\|$$

with optimization variable  $\mu_j$

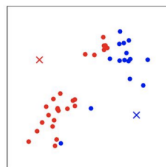
# Illustration of $k$ -Means Clustering



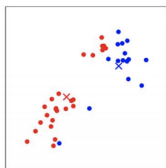
(a)



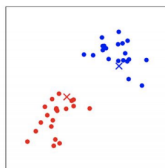
(b)



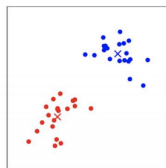
(c)



(d)



(e)



(f)

- (a), (b): data and initial guess
- (c), (d): first iteration
- (e), (f): second iteration

# Outline

- ① Regression / Prediction
- ② Classification
- ③ Unsupervised Learning
- ④ Reinforcement Learning

# Reinforcement Learning



autonomous driving



gaming (e.g.,  
AlphaGo)



robotics

# Markov Decision Process

basic elements:

- $s \in S$ : **states**
  - traffic, GPS position, lanes, distances from other cars, etc.
- $a \in A$ : **actions**
  - route, speed, lane switching, etc.
- $P(s'; s, a)$ : **state transition probabilities**
- $r(s, a)$ : **rewards**
  - time spent, "safety", etc.
- $\delta \in (0, 1)$ : **discount factor**

find a policy  $\pi : S \rightarrow A$  to maximize the total reward:

$$\mathbb{E}_{\pi} \{ r(s_0, a_0) + \delta r(s_1, a_1) + \delta^2 r(s_2, a_2) + \dots \}$$

sometimes state transition probabilities and rewards are **unknown**

# Bellman Equation

value function:

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \{ r(s_0 = s, \pi(s)) + \delta r(s_1, \pi(s_1)) + \delta^2 r(s_2, \pi(s_2)) + \dots \} \\ &= \underbrace{r(s, \pi(s))}_{\text{current reward}} + \delta \underbrace{\sum_{s' \in \mathcal{S}} P(s'; s, \pi(s)) V^\pi(s')}_{\text{expected future reward}} \end{aligned}$$

optimal value function:

$$V^*(s) = \max_{\pi} V^\pi(s)$$

Bell equation: optimal value functions satisfy

$$V^*(s) = \max_{a \in \mathcal{A}} r(s, a) + \delta \sum_{s' \in \mathcal{S}} P(s'; s, a) V^*(s')$$

given optimal value functions, easy to find optimal policy

# Solving Bellman Equation as Linear Program

from Bellman equation, the optimal values satisfy

$$V^*(s) \geq r(s, a) + \delta \sum_{s' \in S} P(s'; s, a) V^*(s'), \quad \forall a \in A$$

LP formulation of Bellman equation:

$$\text{minimize} \quad \sum_{s \in S} V(s)$$

$$\text{subject to} \quad V(s) \geq r(s, a) + \delta \sum_{s' \in S} P(s'; s, a) V(s'), \quad \forall a \in A, \quad \forall s \in S$$

with optimization variables  $V(s)$ ,  $\forall s \in S$