

# DS4400 Notes

DS4400 Notes 01/24

## 1. Convex functions:

A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex iff  $\forall \theta_1, \theta_2 \in \mathbb{R}^d$  and  $\forall \alpha \in [0, 1]$  we have  $f(\alpha\theta_1 + (1-\alpha)\theta_2) \leq \alpha f(\theta_1) + (1-\alpha)f(\theta_2)$

In the special case ( $d = 1$ )  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f$  is convex iff  $\forall \theta, f''(\theta) \geq 0$

When the function is convex, **local min**  $\equiv$  **global min**. When the system is not convex, we might find only a **local min** but not a **global min**

## 2. Dealing with non convex function:

In gradient descent:

- (a) use larger  $\rho$  in the beginning and gradually decrease  $\rho$  with iteration.
- (b) Run SGD/GD with multiple random initializations  $\theta_1^{(0)}, \theta_2^{(0)} \dots$  and keep the best solution.

## 3. $\min_{\theta} \sum_{i=1}^N (y_i - \theta^T x_i)^2 \triangleq J(\theta)$

In linear regression,  $J(\theta)$  is convex.

## 4. Robustness of Regression to outliers:

- (a) Run outlier detection algorithm, remove detected outliers, then run Linear Regression on remaining points.
- (b) Robust Regression cost function.  
 $\min_{\theta} \sum_{i=1}^N e_i^2$ ,  $e_i \triangleq y_i - \theta^T x_i$   
 $e^2$  is extremely unhappy with large errors.

we might use  $|e|$  to replace the function. This might be more tolerance. Then,  $\min_{\theta} \sum_{i=1}^N |y_i - \theta^T x_i|$

## 5. Exercise: $D = \{(x_1, y_1 = 100) \dots (x_1 = 0, y_1 = 0 = 100), (x_{11}, y_{11} = 0), (x_{12}, y_{12} = 0)\}$

$$e^2: 10(\theta - 100)^2 + 2\theta^2 \rightarrow$$

$$\frac{\partial}{\partial \theta} = 20(\theta - 100) + 4\theta = 0 \rightarrow \theta = 83.3$$

$$|e|: \min_{\theta} \sum_{i=1}^{12} |\theta - y_i| = 10|\theta - 100| + 2\theta$$

$$(\theta \leq 100) = \min_{\theta} 10(100 - \theta) + 2\theta$$

$$= 1000 - 8\theta \rightarrow \theta = 100$$

$$(\theta \geq 100) = \min_{\theta} 10(\theta - 100) + 2\theta$$

$$= 12\theta - 1000 \rightarrow \theta = 100$$

## 6. How to solve l1-norms cost functions?

- (a) No closed form
- (b) we need to be careful with gradient descent
- (c) We need to use convex programming toolboxes (convex optimizations)

## 7. Huber loss function

$$l_{\delta}(e) = \begin{cases} \frac{1}{2}e^2 & |e| \leq \delta \\ \delta|e| - \frac{\delta^2}{2} & |e| \geq \delta \end{cases}$$

$$\frac{\partial l_{\delta}(e)}{\partial e} = \begin{cases} e & -\delta \leq e \leq \delta \\ \delta & e > \delta \\ -\delta & e < -\delta \end{cases}$$

in huber loss function, we don't have closed form solution but we can run gradient descent now.

## 8. Definition: Overfitting:

Learning a system from training data that does very well on training data itself (e.g, very low regression error on training data), but performs poorly on test data.

9. **Definition:** Overfitting in Linear Regression

$$\Phi^T \Phi \theta = \Phi^T Y$$

$$\Rightarrow \theta^* = (\Phi^T \Phi)^{-1} \Phi^T Y$$

$$\text{rank}(\Phi^T \Phi) \leq \min\{\text{rk}(\Phi^T), \text{rk}(\Phi)\} = \text{rk}(\Phi) \leq \min\{N, d\}$$

$\Phi^T \Phi$  is  $d \times d$  matrix, then rank is  $\leq d$ .

Therefore, when  $N < d$  it is not invertible which means we have multiple solutions and results in overfitting.

DS4400 Notes 01/28

1. **Definition:** Overfitting

Refers to situation where the learned model does well on training data and poorly on testing data.

As  $d$  (dimension of system) increases, then training error goes down (can be exactly ZERO for sufficiently large  $d$ )

2. In Linear regression:

$$\min \sum_{i=1}^n (\theta^T \phi(x_i) - y_i)^2$$

set the derivative to 0 and we find

$$\Phi^T \Phi \theta = \Phi^T Y$$

Then  $\theta^* = (\Phi^T \Phi)^{-1} \Phi^T Y$

**When is it the case that  $\Phi^T \Phi$  is not invertible?**

Since  $\Phi^T \Phi \in \mathbb{R}^{N \times d}$

$$\text{rk}(\Phi^T \Phi) \leq \text{rk}(\Phi) \leq \min\{N, d\}$$

$\Phi^T \Phi \in \mathbb{R}^{d \times d}$  is invertible when  $\text{rk}(\Phi^T \Phi) = d$ . Therefore, when  $N < d$ ,  $\text{rk}(\Phi^T \Phi) = N$ ,  $\Phi^T \Phi$  is not invertible. There will be infinitely many solutions for  $\theta$ .

**Generally, need sufficient # samples**

3. Test overfitting.

If  $\Phi^T \Phi$  is not invertible,

$$\exists v \neq 0, \Phi^T \Phi v = 0$$

$$\begin{aligned} \Rightarrow \theta^* + \alpha v &\text{ is also a solution for any } \alpha \in \mathbb{R} \\ \Phi^T \Phi (\theta^* + \alpha v) &= \Phi^T \Phi \theta^* + \Phi^T \Phi (\alpha v) \\ &= \Phi^T \Phi \theta^* + \alpha \Phi^T \Phi v \\ &= \Phi^T \Phi \theta^* = \Phi^T Y \end{aligned}$$

We can find large  $\alpha$  so that  $\theta^*$  have extremely large entries.

**Generally, if the entries are very large (abs) we might have overfitting**

4. Treat overfitting

We want to change regression optimization to prevent  $\theta$  from very large terms.

then we change the cost function:

$$\min_{\theta} \sum_{i=1}^N (\theta^T \phi(x_i) - y_i)^2 + \lambda \sum_{j=1}^d \theta_j^2$$

$\lambda$ : regularization parameter ( $> 0$ )

$\sum_{j=1}^d \theta_j^2$ : regularizer.

$\lambda \rightarrow 0$ : back to overfitting

$\lambda \rightarrow \infty$ :  $\theta^* = 0$ , underfitting

- (a) closed-form

$$\frac{\partial J}{\partial \theta}$$

$$\begin{aligned} &= 2\Phi^T (\Phi \theta - Y) + \lambda \frac{\partial \sum_{j=1}^d \theta_j^2}{\partial \theta} \\ &= 2\Phi^T (\Phi \theta - Y) + 2\lambda \theta \end{aligned}$$

Let it be zero:

$$\Phi^T \Phi \theta + \lambda \theta = \Phi^T Y$$

$$(\Phi^T \Phi + \lambda I_d) \theta = \Phi^T Y$$

$$\text{Then } \theta^* = (\Phi^T \Phi + \lambda I_d)^{-1} \Phi^T Y$$

- (b) Gradient descent

Find initial  $\theta^{(0)}$

$$\theta^t = \theta^{(t-1)} - \rho \frac{\partial J}{\partial \theta} |_{\theta^{(t-1)}}$$

$$= \theta^{(t-1)} - 2\Phi^T (\Phi \theta^{(t-1)} - Y) + 2\lambda \theta^{(t-1)}$$

5. Hyperparameter Tuning

GD: set learning rate  $\rho$

Robust Reg: Huber loss  $\delta$

overfitting and regularization:  $\lambda$

$\rho, \delta, \lambda$  = hyperparameters

**How to pick hyperparameters?**

**BAD APPROACH 1:**

- (a) pick some set of possible  $\lambda_i \in \{\lambda_1, \lambda_2, \dots\}$   
 Run regression with  $\lambda_i$  and find  $\theta_i^*$   
 Measure regression error:

$$\epsilon_{tr}(\lambda) = \sum_{i=1}^N ((\theta^*(\lambda))^T x_i - y_i)^2$$

To sum: just find  $\lambda$  for which  $\epsilon_{tr}(\lambda)$  is minimum

**This approach is setting  $\lambda$  back to 0**

**Test data needed!!!**

- (a) We need to Train  $\lambda_i$  on **training set** to minimize the cost function

$$2\Phi^T(\Phi\theta - Y) + 2\lambda\theta$$

to find  $\theta_i^*$

- (b) Measure regression error on the **hold-out set**  $D^{ho}$

$$\epsilon_{tr} = \sum_{x_i, y_i \in D^{ho}} (y_i - (\theta^*(\lambda))^T x_i)^2$$

DS4400 Notes 01/31

### 1. Hyperparameter Tunning:

$$\min_{\theta} \sum_{i=1}^N (\theta^T \phi(x_i) - y_i)^2 + \lambda \sum_{j=1}^d \theta_j^2$$

- For  $\lambda \in \{\lambda_1, \lambda_2, \dots, \lambda_p\}$ 
  - Tran using  $D^{tr}$  with  $\lambda \rightarrow \theta^*(\lambda)$
  - Measure validation error

$$\epsilon^{tr}(\lambda) = \sum_{x_i, y_i \in D^{ho}} (y_i - (\theta^*(\lambda))^T x_i)^2$$

- select  $\lambda$  which minimizes

$$\epsilon^{ho}(\lambda) \rightarrow \lambda^* = \min_{\{\lambda_1, \lambda_2, \dots, \lambda_p\}} \epsilon^{ho}(\lambda)$$

### 2. Problems:

- Take much longer time since we are training the models multiple times
- Each training is using a subset of the data set, then each training is amplifying the problem of overfitting.

### 3. K-fold cross validation

divide Data set to k equally large sets  $\{D_1, D_2, \dots, D_k\} \in D$

- For  $\lambda \in \{\lambda_1, \lambda_2, \dots, \lambda_p\}$ 
  - For  $i = 1, 2, \dots, k$ 
    - \* train on  $\bigcup_{j \neq i} D^j$  and get  $\theta_i^*(\lambda)$
    - \* compute validate error on  $D^i \rightarrow \epsilon_i^{ho}(\lambda)$
    - compute average of  $\{\epsilon_i^{ho}(\lambda)\}$ :  
 $\epsilon^{ho} = \frac{1}{k} \sum_{i=1}^k \epsilon^{ho}(\lambda)$
  - select  $\lambda^* = \min_{\{\lambda_1, \lambda_2, \dots, \lambda_p\}} \epsilon^{ho}(\lambda)$

Once we find the best  $\lambda$ , train the model on the whole set.

### 4. PROBABILITY REVIEW

- Random Variable: a variable that takes values corresponding to outcome of a random phenomenon.
- Discrete r.v.: discrete values
- continuous r.v. continus range of values
- Condition:  $P(X|Y) = \frac{P(X \cap Y)}{P(Y)}$

$$P(X, Y) = P(X|Y)P(Y)$$

$$P(X, Y) = P(Y|X)P(X)$$

**Chain rule:**

$$P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots P(X_N|X_1, X_2, \dots, X_{N-1})$$

- Marginalization

$p(x, y)$  known

$$p(x) = \sum_y p(x, Y = y) = \int p(x, y) dy$$

- Bayes Rule:

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)} = \frac{P(X|Y)P(Y)}{P(X)}$$

- Independence:  
r.v. are independent ( $X \perp\!\!\!\perp Y$ ) iff  
 $P(X|Y) = p(X), P(Y|X) = p(Y)$   
or  $P(x, y) = P(x)p(y)$
- conditional independence example:  
 $X$  = height of person,  $Y$  = vocabulary,  
 $X$  is not independent of  $Y$  since babies may have less vocabulary and with lower heights.  
However,  $X$  = height,  $Y$  = vocab,  $Z$  = age. Then  $(X \perp\!\!\!\perp Y) | Z$

$$P(X, Y|Z) = P(X|Z)P(Y|Z)$$

$$\Rightarrow P(X|Y, Z) = P(X|Z)$$

- Expectation:  
 $E(X) = \sum xp(x)$  or  $\int xp(x)dx$   
 $E(f(X)) = \sum f(x)p(x)$  or  $\int f(x)p(x)dx$   
Given  $X \perp\!\!\!\perp Y$ ,  $E[XY] = E[X]E[Y]$   
hint:  $(E[XY] = E[f(x, y)])$
- IID r.v: independent and identically distributed  
 $p(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = p(X_1 = x_1)p(X_2 = x_2)\dots p(X_n = x_n)$   
and each experiment is identical.  
 $P(X_1 = \theta) = P(X_2 = \theta) = \dots = P(X_n = \theta)$