

# Notes 1.

1. LMS algorithm

$$\text{Cost function: } J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\Rightarrow \text{gradient descent: } \theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta)$$

- (i) simultaneously performed for all values of  $j$
- (ii)  $\alpha$ : learning rate
- (iii) " - " : we are minimizing

Further, solve the partial derivative part:

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \underbrace{\sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}}_{\downarrow} \quad (\text{for every } j).$$

intuition: larger error  $\rightarrow$  larger  
look at every example in the entire  
training set on every step  
 $\Rightarrow$  is called "batch gradient descent"

An alternative:

loop {

for  $i = 1$  to  $m$ , {

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)} \quad (\text{for every } j)$$

}

}

"Stochastic gradient descent"

do not have to scan through the entire training set before taking a single step

+ The normal equations

2. Matrix derivatives.

D

$$f: \mathbb{R}^{m \times n} \mapsto \mathbb{R} \quad \nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \cdots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}$$

D tr: trace

$$(i) \text{ if } AB \text{ is square } \text{tr } AB = \text{tr } BA$$

$$(ii) \text{ tr } ABC = \text{tr } CAB = \text{tr } BCA$$

$$(iii) \text{ tr } ABCD = \dots$$

$$(iv) A \& B \text{ are square Matrices}$$

$$\text{tr } A = \text{tr } A^\top$$

$$\text{tr}(A+B) = \text{tr} A + \text{tr} B$$

$$\text{tr} aA = a \text{tr} A$$

⇒ Some facts of matrix derivatives.

$$\nabla_A \text{tr} AB = B^T$$

$$\nabla_A \text{tr} f(A) = (\nabla_A f(A))^T$$

$$\nabla_A \text{tr} ABC = CAB + C^T B A^T$$

$$\nabla_A |A| = |A|(A^{-1})^T$$

2.2 Least squares revisited.

design Matrix :  $X = \begin{bmatrix} (X^{(1)})^T \\ \vdots \\ (X^{(m)})^T \end{bmatrix}$  m samples & n features

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$X\theta - \vec{y} = \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix} \Rightarrow \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y})$$
$$= \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$
$$= J(\theta)$$

$$\Rightarrow \nabla_{\theta} J(\theta) = \dots = X^T X \theta - X^T \vec{y} \stackrel{\text{to minimize}}{=} 0$$

$$\Rightarrow \theta = (X^T X)^{-1} X^T \vec{y}$$

3 Probabilistic interpretation.  $\rightarrow$  solve: why is  $J$  (loss function) a reasonable choice

$$\text{Assumption } \epsilon^{(i)} = \theta^T x^{(i)} + \xi^{(i)}$$

( $\xi^{(i)}$  is an error term that captures either unmodeled effects or random noise.)  
 $\xi^{(i)} \sim N(0, \sigma^2)$

$$\Rightarrow P(y^{(i)} | x^{(i)}, \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right).$$

Namely,  $y^{(i)} | x^{(i)}; \theta \sim N(\theta^T x^{(i)}, \sigma^2)$

Likelihood function:

$$L(\theta) = L(\theta; X, \bar{y}) = P(\bar{y} | X; \theta)$$

$$\begin{aligned} &\text{independence} \\ &\text{assumption} \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta) \end{aligned}$$

$$= \frac{1}{\prod_{i=1}^m \sqrt{2\pi}\sigma} \exp\left(-\frac{\sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right).$$

Principle of maximum likelihood:  $\hat{\theta}(\theta) = \log L(\theta) = m \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$

$$\Leftrightarrow \text{minimize } \underbrace{\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2}_{\text{---}}$$

Interesting fact: our final choice of  $\theta$  did not depend on  $\sigma^2$

4. Locally weighted linear regression. — non-parametric algorithm.  
(LWR)

Fit  $\theta$  to minimize  $\sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$

\* fairly standard choice:

$$w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2T^2}\right).$$

( $x$ : the particular point at which we're trying to evaluate  $x$ )

interpretation: ① give a much higher weight to the training examples close to the query point  $x$

②  $T$ : bandwidth parameter

controls how quickly the weight of a training example falls off with distance of its  $x^{(i)}$  from the query point  $x$ .

## Part II. Classification and logistic regression.

### 5. Logistic Regression.

- Use linear regression:

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

where

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

logistic function / sigmoid function.

A useful property:  $g'(z) = g(z)(1 - g(z))$

$$\text{Assumptions: } P(y=1|x; \theta) = h_{\theta}(x) \quad \left. \begin{array}{l} \\ P(y=0|x; \theta) = 1 - h_{\theta}(x) \end{array} \right\} \Rightarrow P(y|x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

$$\Rightarrow L(\theta) = P(\bar{y} | X; \theta)$$

$$= \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta)$$

$$= \prod_{i=1}^m (h_{\theta}(x^{(i)}))^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

$$\Rightarrow L(\theta) = \log L(\theta) = \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))$$

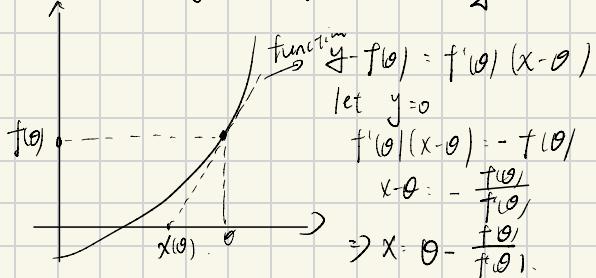
$$\Rightarrow \text{gradient descent: } \theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

(Stochastic)

### 6. Dimension. The perceptron learning algorithm.

$$g(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$

### 7. Another algorithm for maximizing $L(\theta)$ .



$\Rightarrow$  It performs the following update:

$$\theta := \theta - \frac{\nabla l(\theta)}{l'(\theta)}$$

Interpretation:

- $\textcircled{1}$  approximate the function  $f$  via a linear function that is tangent to  $f$  at the current guess  $\theta$
- $\textcircled{2}$  let the next guess for  $\theta$  be where that linear function is zero

$\Rightarrow$  Application: maximize  $l$

find point where  $l'(\theta) = 0$

let  $f(\theta) = l'(\theta)$

update rule:

$$\theta := \theta - \frac{l''(\theta)}{l'(\theta)}$$

$\Rightarrow$  generalize for vector-valued  $\theta$

$$\vec{\theta} := \vec{\theta} - H^{-1} \nabla l(\vec{\theta}),$$

$$(H: \text{(Hessian)} \quad H_{ij} = \frac{\partial^2 l(\theta)}{\partial \theta_i \partial \theta_j})$$

Result: faster convergence than batch gradient descent

Called: Fisher scoring ✓

### Part III Generalized Linear Models.

In regression:  $y|x; \mu \sim N(\mu, \sigma^2)$

In classification:  $y|x; \theta \sim \text{Bernoulli}(\theta)$

8. The exponential family (canonical parameter)

$$P(y; \theta) = b(y) \exp \left( \underbrace{\eta^\top T(y)}_{\text{sufficient statistic}} - \underbrace{a(\eta)}_{\log \text{partition function}} \right).$$

make sure integral equals to 1

$$\begin{aligned}
 \text{Bernoulli: } P(y|\phi) &= \phi^y (1-\phi)^{1-y} \\
 &= \exp(y \log \phi + (1-y) \log(1-\phi)) \\
 &= \exp\left(\log\left(\frac{\phi}{1-\phi}\right) y + \log(1-\phi)\right) \\
 \therefore \hat{\phi} &= \log \frac{y}{1-y}, \\
 \Rightarrow T(y) &= y \\
 a(y) &= \log(1+e^y), \\
 b(y) &=
 \end{aligned}$$

$$\begin{aligned}
 \text{Gaussian: } P(y|\mu) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y-\mu)^2\right) \\
 &= \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}y^2) \cdot \exp(\mu y - \frac{1}{2}\mu^2) \\
 \Rightarrow \hat{\mu} &= \mu \\
 T(y) &= y \\
 a(y) &= \mu^2/2 = y^2/2 \\
 b(y) &= (1/\sqrt{2\pi}) \exp(-y^2/2).
 \end{aligned}$$

## 9. Constructing GLMs.

To derive a GLM, have to make assumptions:

- ①  $y|x, \theta \sim \text{Exponential Family}(\eta)$
- ② the prediction  $h(x)$  outputs to satisfy  $h(x) = E[y|x]$
- ③ natural parameter  $\eta$  and inputs  $x$  are related linearly:  $\eta = \theta^T x$ .

## 9.1 Ordinary Least Squares

$$h_0(x) = E[y|x, \theta]$$

$$\text{Assumption 1: } \mu$$

$$\text{Assumption 2: } \eta$$

$$= \theta^T X$$

$$\uparrow$$

Assumption 3.

## 9.2 Logistic Regression.

$$h_0(x) = E[y|x, \theta]$$

$$= \phi$$

$$= 1/(1+e^{-\theta})$$

$$= \underbrace{1/(1+e^{-\theta^T x})}$$

Once we assume that  $y$  conditioned on  $x$  is Bernoulli,

it arises as a consequence of the definition of GLMs and exponential family distributions.

### 9.3 Softmax Regression.

$y \in \{1, 2, \dots, k\}$ . ( $y$  is response variable)

Multinomial distribution.

⇒ Parameterize  $\phi_1, \dots, \phi_{k-1}$  (since  $\sum_{i=1}^k \phi_i = 1$ ).

define  $T(y) \in \mathbb{R}^{k^1}$  as follows:

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, T(2) = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, T(3) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, T(k-1) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}, T(k) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Notation:  $(T(y))_i$  denotes the  $i$ -th element of the vector  $T(y)$ .

$$\Rightarrow (T(y))_i = \{y = i\}$$

$$E[(T(y))_i] = P(y = i) = \phi_i$$

$$\Rightarrow P(y, \phi) = \phi_1^{1 \{y=1\}} \phi_2^{1 \{y=2\}} \cdots \phi_k^{1 \{y=k\}}$$

$$= \phi_1^{1 \{y=1\}} \phi_2^{1 \{y=2\}} \cdots \phi_k^{1 \{y=k\}} e^{-\sum_{i=1}^k (T(y))_i}$$

$$= \phi_1^{(T(y))_1} \phi_2^{(T(y))_2} \cdots \phi_k^{(T(y))_k} e^{-\sum_{i=1}^k (T(y))_i}$$

$$= \exp \left\{ (T(y))_1 \log(\phi_1) + (T(y))_2 \log(\phi_2) + \cdots + (T(y))_k \log(\phi_k) \right\}$$

$$= \exp \left\{ (T(y))_1 \log(\phi_1 / \phi_k) + (T(y))_2 \log(\phi_2 / \phi_k) + \cdots + (T(y))_k \log(\phi_{k-1} / \phi_k) + \log(\phi_k) \right\}$$

$$= b(y) \exp(\beta^\top T(y) - a(\beta)).$$

where

$$\beta = \begin{bmatrix} \log(\phi_1 / \phi_k) \\ \log(\phi_2 / \phi_k) \\ \vdots \\ \log(\phi_{k-1} / \phi_k) \end{bmatrix}$$

$$a(\beta) = -\log(\phi_k)$$

$$b(y) = 1.$$

\* link function:  $\eta_i = \log \frac{\phi_i}{\phi_k}$

define:  $\partial_k \cdot \log(\phi_k / \phi_k) = 0$

$$e^{\eta_i} = \phi_i / \phi_k \Rightarrow \phi_k \cdot e^{\eta_i} = \phi_i$$

$$\Rightarrow \phi_k \sum_{j=1}^k e^{\eta_j} = \sum_{j=1}^k \phi_j = 1$$

$$\Rightarrow \phi_k = \frac{1}{\sum_{j=1}^k e^{\eta_j}}$$

$$\Rightarrow \phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}} \quad \text{softmax function}$$

Assumption 3.  $\eta_i = \theta_i^\top x$

$$\Rightarrow P(y=i | x; \theta) = \phi_i \\ = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}$$

$$= \frac{e^{\theta_i^\top x}}{\sum_{j=1}^k e^{\theta_j^\top x}} \quad \text{softmax regression}$$

Output:  $h(x) = E[T(y) | x; \theta]$

$$= E \left[ \begin{array}{c|c} \{y=1\} & \\ \{y=2\} & \\ \vdots & \\ \{y=k\} & \end{array} \middle| x; \theta \right]$$

$$= \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix} = \begin{bmatrix} \frac{\exp(\theta_1^\top x)}{\sum_{j=1}^k \exp(\theta_j^\top x)} \\ \vdots \\ \frac{\exp(\theta_{k-1}^\top x)}{\sum_{j=1}^k \exp(\theta_j^\top x)} \end{bmatrix}$$

(choose  $\theta$ )

parameter fitting:

$$l(\theta) = \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}; \theta)$$

$$= \sum_{i=1}^m \log \frac{1}{k} \left( \frac{e^{\theta_i^\top x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^\top x^{(i)}}} \right)^{1\{y^{(i)}=l\}}$$