

## Lec3 Supervised Learning (continued)

Outline:

Linear Regression (Recap)

Locally weighted regression

Probabilistic interpretation

Logistic Regression

Newton's method

Recap:

$(x^{(i)}, y^{(i)})$  -  $i^{\text{th}}$  example

$x^{(i)} \in \mathbb{R}^{n+1}$   $y^{(i)} \in \mathbb{R}$

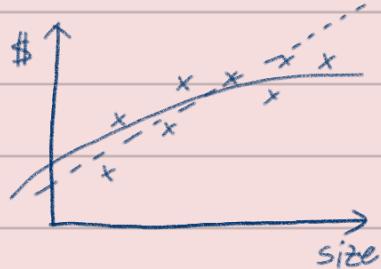
$m = \# \text{ of examples}$   $n = \# \text{ of features}$

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j = \theta^T x$$

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \text{인수인 } \theta \rightarrow \theta^T = [\theta_0, \theta_1, \dots, \theta_n]$$

모든 계수를 포함한 벡터.



$\theta_0 + \theta_1 x_1$  do we want linear fn?  
 $\theta_0 + \theta_1 x_1 + \theta_2 x_2^2$  or quadratic?  
 or sth else?

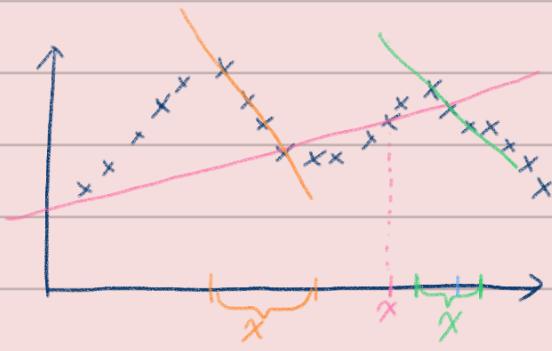
## Locally weighted regression

"parametric" learning algorithm:

Fit fixed set of parameters ( $\theta_i$ ) to data

"Non-parametric" learning algorithm:  $\rightarrow$  not good when data set is huge

Amount of data/parameters you need to keep grows (linearly)  
 with the size of data



To evaluate  $h$  at certain  $X$ :

LR: Fit  $\theta$  to minimize  $J(\theta)$ ,  
 return  $\theta^T x$

Locally weight regression:

Fit  $\theta$  to minimize  $\sum_{i=1}^m w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$

where  $w^{(i)}$  is a weight function

$$w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\sigma^2}\right) \quad \begin{array}{l} x^{(i)}: i^{\text{th}} \text{ training example} \\ x: where I want to make the prediction \end{array}$$

If  $|x^{(i)} - x|$  is small,  $w^{(i)} \approx 1$   
 If  $|x^{(i)} - x|$  is large,  $w^{(i)} \approx 0$

$\Rightarrow$  net effect: training example of  $\sigma$  distance  
 input with  $w^{(i)}$ ,  $w^{(i)} \approx 0$  if  $x^{(i)}$  far away  
 essentially sums over for the squared error  
 for the examples that are close to  $x$

$\sigma$ : bandwidth

LWR에서  $\hat{y}$ (예측)가真实값과 underfitting인 경우

작은수로 오류이 overfitting인 경우



underfitting

training data fit  
여러개의 test data  
모든이 예측 성능 나쁨

overfitting

training data에는 잘 맞음  
test data는 잘 안 맞음  
( $\because$  noise 많아 학습해)!

## Probabilistic interpretation

Why Least Squares?

Assume ①  $y^{(i)} = \theta^T X^{(i)} + \varepsilon^{(i)}$

$J(\theta)$  error term (unmodeled effect, random noise)  
e.g. the mood of the seller..!

probability function

$$= \text{"density"} \quad P(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)$$

: central limit theorem

상당히 정규분포  
very far from Gaussian  
but default noise distribution is Gaussian!

Assumption ② on  $\varepsilon$ : IID (independent & identically distributed)

this implies:

$$\text{"parameterized by } \theta\text{"} \quad P(y^{(i)} | X^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T X^{(i)})^2}{2\sigma^2}\right) \quad = \varepsilon^{(i)}$$

$$\text{i.e. } y^{(i)} | X^{(i)}; \theta \sim N(\theta^T X^{(i)}, \sigma^2)$$

= Given  $X^{(i)}$  &  $\theta$ , the house price  $y^{(i)}$ 's probability (density) is  
 $\theta^T X^{(i)}$ 를 mean으로,  $\sigma^2$ 을 variance로 가지는 정규분포(Gaussian)을 갖는다.  
이는 확률밀도함수  $P(y^{(i)} | X^{(i)}; \theta)$ 는, 모델이  $\theta$ 를 선택했을 때 가질 것이다.  
 $y^{(i)}$ 가 갖는 확률을 나타낸다

Likelihood of  $\theta$   $L(\theta) = P(\vec{y} | \mathbf{x}; \theta)$

'likelihood of parameter' ( $\theta$ )  
likelihood of data ( $\mathbf{x}$ )

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

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

Log Likelihood of  $\theta$   $l(\theta) = \log L(\theta)$

$$= \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp(...)$$

$$= \sum_{i=1}^m \left[ \log \frac{1}{\sqrt{2\pi}\sigma} + \log \exp(...) \right]$$

$$= m \log \frac{1}{\sqrt{2\pi}\sigma} + \boxed{\sum_{i=1}^m \frac{(y^{(i)} - \theta^T X^{(i)})^2}{2\sigma^2}}$$

MLE: Maximum Likelihood estimation : choose  $\theta$  to maximize  $L(\theta)$

$\Leftrightarrow$  choose  $\theta$  to maximize  $l(\theta)$

i.e. choose  $\theta$  to minimize  $\sum_{i=1}^m \frac{(y^{(i)} - \theta^T X^{(i)})^2}{2\sigma^2}$

( $\because$  첫번재 term은  $\theta$ 가 영향을 미친다)

which is  $J(\theta)$ .

$\therefore$  Least Squares?

## Classification

$y \in \{0, 1\}$  binary classification



for classification, linear regression is not a good algorithm.

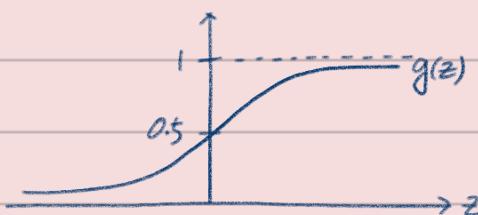
## Logistic regression

이진 분류에 대한 LR. LR는 비선형  $\ell(\theta)$ 가 global max를 가지는 경우에만 가능하다 (No local optima)

want  $h_{\theta}(x) \in [0, 1]$

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \quad h_{\theta}(x) = \theta^T x \text{ is sigmoid function } \approx \frac{1}{1 + e^{-z}}$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{"Sigmoid" or "logistic" function}$$



if tumor is malignant

$$\left. \begin{aligned} \text{e.g. } P(y=1 | X; \theta) &= h_{\theta}(x) \\ &\downarrow \text{size of the tumor} \end{aligned} \right\} \begin{array}{l} \text{Linear} \\ \text{Algebra} \\ \text{Compress} \end{array}$$

$$p(y|x; \theta) = h(x)^y (1 - h(x))^{1-y}$$

① Define Likelihood

$$L(\theta) = P(\vec{y}|X; \theta) = \prod_{i=1}^m P(y^{(i)}|x^{(i)}; \theta) = \prod_{i=1}^m h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

② Define log Likelihood

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

③ Choose  $\theta$  to maximize  $l(\theta)$

④ plug  $\theta$  to the model and use the new feature  $X \rightarrow$  predict  $y$ .

$$\text{① ④ } \rightarrow \theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} l(\theta) \quad <> \text{ square error } J(\theta) \text{ を } \text{최소화하는 (least squares)}$$

경사 하강법에 적용되는 원칙

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \text{ 를 }$$

모든  $\theta$ 를 최소화하는 원칙

$$\Rightarrow \theta_j := \theta_j + \alpha \cdot \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) \cdot x_j^{(i)}$$

## Newton's method

Linear Rmk Normal Eq을 사용하여 더 빠른 progress

Have  $f$ , want to find  $\theta$  s.t.  $f(\theta) = 0$

[Want to maximize  $l(\theta)$  i.e. want  $l'(\theta) = 0$ ]

$$\theta^{(1)} := \theta^{(0)} - \Delta$$

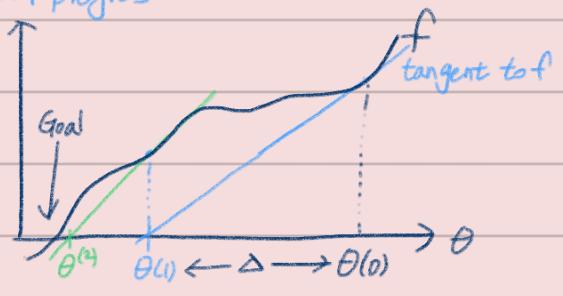
$$f'(\theta^{(0)}) = \frac{f(\theta^{(0)})}{\Delta}$$

$$\Delta = \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

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

$$\theta^{(t+1)} := \theta^{(t)} - \frac{l(\theta^{(t)})}{l'(\theta^{(t)})}$$



"Quadratic convergence" Newton's method

0.01 error  $\rightarrow$  0.0001 error  $\rightarrow$  0.00000001 error

becomes much inaccurate by 1 single iteration

# of parameters  $\leq 50$  ish otherwise go for Gradient Descent instead

when  $\theta$  is a vector: ( $\theta \in \mathbb{R}^{n+1}$ )  $\xrightarrow{\text{size } \theta}$

$$\theta^{(t+1)} := \theta^{(t)} + H^{-1} \nabla l \quad \text{vector of derivatives } \mathbb{R}^{n+1}$$

when  $H$  is the Hessian matrix

$$H_{ij} = \frac{\partial^2 l}{\partial \theta_i \partial \theta_j}$$

$\rightarrow$   $\theta$  is a high dimensional vector  $\Rightarrow$  Newton's method each step  $\Theta(\# \text{ iterations})$

$\therefore$  having to invert a pretty big matrix ( $H^{-1}$ )