

# THEORY OF LINEAR REGRESSION

This page intentionally left blank

Regression

- So far, we've been interested in learning  $P(Y|X)$  where  $Y$  has discrete values. This is a.k.a. 'classification'.
- What if  $Y$  is continuous? (a.k.a. 'regression')  
predict real-valued attributes (e.g. weight, height).
- Predict stock market tomorrow.

Setting

learn  $f: X \rightarrow Y$ , where  $Y$  is real-valued, given  $\{ \langle x^1, y^1 \rangle, \dots, \langle x^n, y^n \rangle \}$

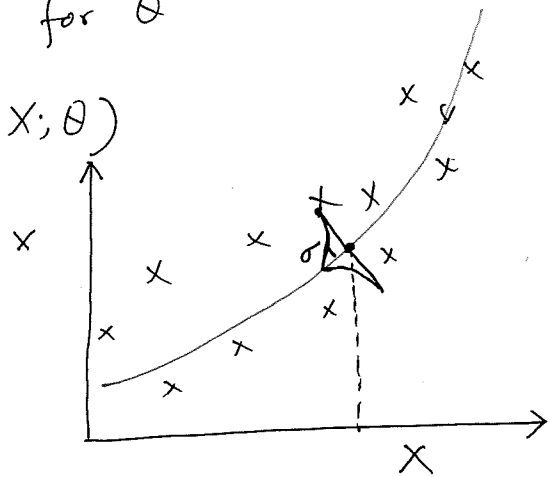
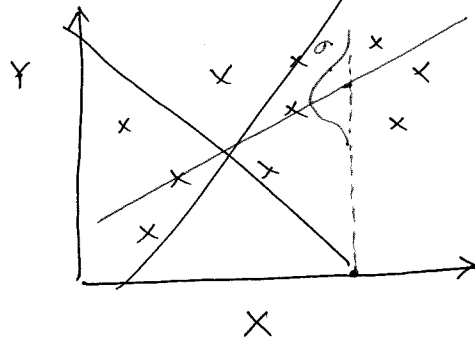
Approach

I. choose some parameterized form for  $P(Y|X; \theta)$

$\theta$  is a vector of parameters

II. Use MLE or MAP estimate for  $\theta$

I. Choose parameterized form for  $P(Y|X; \theta)$



Assume  $Y$  is some deterministic  $f(x)$ , plus some random noise

$$y = f(x) + \epsilon \quad \text{where } \epsilon \sim N(0, \sigma)$$

Therefore,  $Y$  is a random variable that follows the distribution

$$p(y|x) = N(f(x), \sigma)$$

And expected value of  $y$  for any given  $x$  is  $f(x)$

## Linear Regression

$$p(y|x) = N(f(x), \sigma)$$

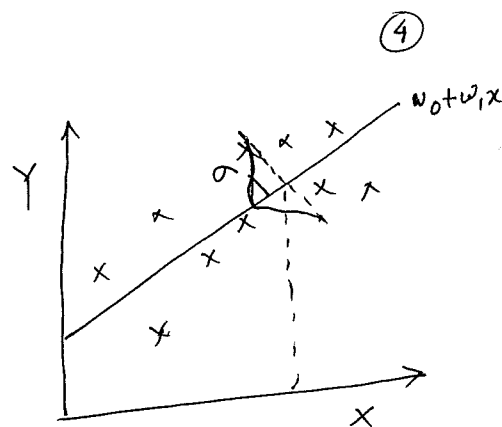
e.g. assume  $f(x)$  is a linear function of  $x$

$$p(y|x) = N(w_0 + w_1 x, \sigma)$$

$$E[y|x] = w_0 + w_1 x$$

$W = \langle w_0, w_1 \rangle$  denotes the parameters (weights)

$$p(y|x; W) = N(w_0 + w_1 x, \sigma)$$



## Training Linear Regression

$$p(y|x; W) = N(w_0 + w_1 x, \sigma)$$

Q. How can we learn  $W$  from the training data?

A. Learn Maximum Conditional Likelihood Estimate (MCLE)

$$W_{MCLE} = \arg \max_W \prod_l p(y^l | x^l, W)$$

$$\ln W_{MCLE} = \arg \max_W \sum_l \ln p(y^l | x^l, W)$$

l<sup>th</sup> training examples

$$= -\frac{1}{2} \left( \frac{y - f(x; W)}{\sigma} \right)^2$$

where  $p(y|x; W) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left( \frac{y - (w_0 + w_1 x)}{\sigma} \right)^2}$

$$p(y^l | x^l, W) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left( \frac{y^l - (w_0 + w_1 x^l)}{\sigma} \right)^2}$$

## Logistic Regression More Generally

$$p(y|x) = N(f(x), \sigma I)$$

$$f(x) = w_0 + \sum_{i=1}^n w_i x_i$$

$$p(y|x) = N\left(w_0 + \sum_{i=1}^n w_i x_i, \sigma I\right)$$

$$E[y|x] = w_0 + \sum_{i=1}^n w_i x_i$$

$$p(y|x; W) = N\left(w_0 + \sum_{i=1}^n w_i x_i, \sigma I\right)$$

$$\vec{x} = \langle x_1, x_2, \dots, x_n \rangle$$

$$W = \langle w_0, w_1, \dots, w_n \rangle$$

MCLE derivation

$$W_{\text{MCLE}} = \arg \max_W \sum_l \ln P(y^l | x^l, W)$$

$$\text{where } p(y^l | x^l, W) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{1}{2} \left( \frac{y^l - (w_0 + w_1 x^l)}{\sigma} \right)^2}$$

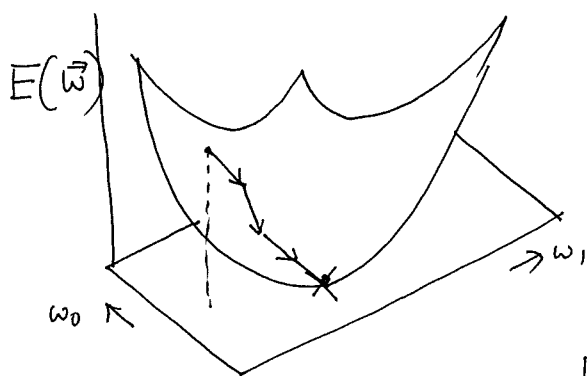
$$W_{\text{MCLE}} = \arg \max_W \sum_l \ln \frac{1}{\sqrt{2\pi}\sigma^2} - \frac{1}{2} \left( \frac{y^l - (w_0 + w_1 x^l)}{\sigma} \right)^2$$

$$W_{\text{MCLE}} = \arg \max_W \sum_l -\frac{1}{2\sigma^2} (y^l - (w_0 + w_1 x^l))^2$$

$$W_{\text{MCLE}} = \arg \min_W \sum_l \boxed{\frac{1}{2\sigma^2}} (y^l - (w_0 + w_1 x^l))^2$$

$$= \arg \min_W \sum_l \underbrace{(y^l - (w_0 + w_1 x^l))^2}_{\text{"squared error"}}$$

[often used in "curve fitting"]

Gradient Descent :

$$\text{Gradient } \nabla E(\vec{w}) = \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

$$\text{Training Rule: } \vec{w}^{(i+1)} \leftarrow \vec{w}^{(i)} - \eta \nabla E(\vec{w})$$

opposite  
of gradient

learning rate

$$\text{i.e. } \Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

# Minimizing Squared Error: Gradient Descent

⑥

$$\frac{\partial E}{\partial w_1} = \sum_l 2 (y^l - (w_0 + w_1 x^l)) (-x^l)$$

$$= -2 \sum_l (y^l - (w_0 + w_1 x^l)) x^l$$

$$\frac{\partial E}{\partial w_0} = -2 \sum_l (y^l - (w_0 + w_1 x^l))$$

So, our update rule is:

$$w_0 \leftarrow w_0 - \eta \left[ -2 \sum_l (y^l - (w_0 + w_1 x^l)) \right]$$

$$w_0 \leftarrow w_0 + 2\eta \left[ \sum_l (y^l - (w_0 + w_1 x^l)) \right]$$

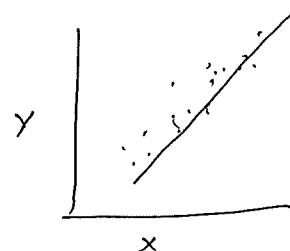
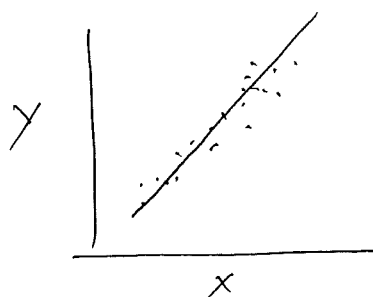
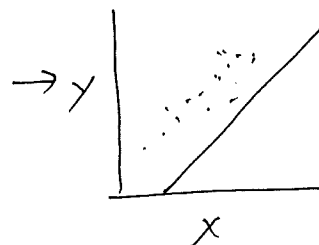
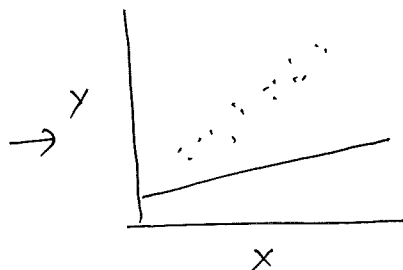
~~w\_0 \leftarrow w\_0~~

$$w_1 \leftarrow w_1 - \eta \left[ -2 \sum_l x^l (y^l - (w_0 + w_1 x^l)) \right]$$

$$w_1 \leftarrow w_1 + 2\eta \left[ \sum_l x^l (y^l - (w_0 + w_1 x^l)) \right]$$

In general,

$$w_i \leftarrow w_i + 2\eta \sum_l x_i^l \underbrace{\left( y^l - \sum_{j=1}^n w_j x_j^l \right)}_{\text{"misassignment"}}$$



"converged"

(7)

## Gradient Descent More Generally

Iterate until change  $< \epsilon$

$\forall i$  repeat

$$w_i \leftarrow w_i + \eta 2 \sum_l x_i^l \left( y^l - \left( w_0 + \sum_{j=1}^n w_j x_j^l \right) \right)$$

assume  $x_\phi = 1$  for  $w_0$

How about MAP estimate insted of MCLE?

$$W_{\text{MAP}} = \arg \max_W \ln N(W | \phi, I) + \sum_l \ln (P(Y^l | X^l, W))$$

$$= \arg \max_W \left( \underbrace{-c \sum_l w_i^2}_{\text{regulization term}} \right) + \sum_l \ln (P(Y^l | X^l, W))$$

regulization term  
pushes  $W$  to zero

## Bias and Variance

Given algorithm that outputs estimate  $\hat{\theta}$  for  $\theta$ ,

the bias of the estimator :  $E[\hat{\theta}] - \theta$

the variance of the estimator :  $E[(\hat{\theta} - E(\hat{\theta}))^2]$

e.g.  $\hat{\theta}^{\text{MLE}}$  estimator for probability  $\theta$  of heads, based on  $n$  independent coin flips

$$\hat{\theta}^{\text{MLE}} = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

What is the bias?  $\theta = E[\hat{\theta}^{\text{MLE}}]$

What is the variance?  $\frac{\theta(1-\theta)}{n} = \text{Var}[\hat{\theta}^{\text{MLE}}] = E[\underbrace{(\hat{\theta}^{\text{MLE}} - \theta)^2}_{\text{squared error}}]$

Let's use  $\hat{\theta}^{\text{MAP}}$  estimator

$$\hat{\theta}^{\text{MAP}} = \frac{\alpha_1 + \beta_1 - 1}{(\alpha_1 + \beta_1 - 1) + (\alpha_0 + \beta_0 - 1)}$$

Which estimator has higher bias? higher variance? (for finite  $n$ )

Bias  $\hat{\theta}^{\text{MAP}}$  has higher bias than  $\hat{\theta}^{\text{MLE}}$

Variance

$\hat{\theta}^{\text{MLE}}$  has higher variance.

since the fraction will vary more than the fraction in  $\hat{\theta}^{\text{MAP}}$

As  $n \rightarrow \infty$ ,  $\hat{\theta}^{\text{MAP}} \rightarrow \hat{\theta}^{\text{MLE}}$

So, bias approaches zero

So, variance, approaches zero

Train Set error :	1%	15%	15%	0.5%
CV Set error :	11%	16%	30%	1%
	high variance (overfitting)	high bias (underfitting)	(high bias high variance)	low bias low variance