

2/23/15

1

# (FINISH) REGRESSION

## ① Bayesian Linear Regression (aka Ridge Regression)

Likelihood  $[ Y | X=x \sim N(\vec{w}^T x, \sigma^2) ]$  same as last time

Prior  $[ \vec{w} \sim N(\vec{0}, \tau^2 I) \Leftrightarrow w_i \sim N(0, \tau^2) ]$  IID components of  $\vec{w}$

$$\Rightarrow P(\vec{w}) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi\tau^2}} e^{-\frac{(w_i-0)^2}{2\tau^2}}$$

$$\text{Now } \hat{\vec{w}}_{\text{MAP}} = \underset{\vec{w}}{\text{argmax}} \log P(\vec{w} | D)$$

$$= \underset{\vec{w}}{\text{argmax}} \log \left[ \frac{P(D | \vec{w}) P(\vec{w})}{P(D)} \right]$$

$$= \underset{\vec{w}}{\text{argmax}} [ \log P(D | \vec{w}) + \log P(\vec{w}) ]$$

$$= \underset{\vec{w}}{\text{argmax}} \left[ -\sum_{i=1}^n \frac{(y_i - \vec{w}^T \vec{x}_i)^2}{2\sigma^2} - \frac{n}{2} \log 2\pi\sigma^2 - \sum_{j=0}^d \frac{w_j^2}{2\tau^2} - \frac{d}{2} \log 2\pi\tau^2 \right]$$

$$= \underset{\vec{w}}{\text{argmin}} \left[ \frac{1}{n} \sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i)^2 + \frac{1}{n} \left( \frac{\sigma^2}{\tau^2} \right) \|\vec{w}\|_2^2 \right]$$

$$= \underset{\vec{w}}{\text{argmin}} \left[ \frac{1}{n} \sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i)^2 + \frac{\lambda}{n} \|\vec{w}\|_2^2 \right]$$

Avg. LOSS

trade-off

Regularization!

Bayesian Loss  $\equiv$  Avg-Fitting-Error + Norm-Square-of  $\vec{w}$

Smell right? Let's check.

$$\lambda \rightarrow 0 \quad \left( \frac{\sigma^2}{n} \rightarrow 0 \text{ or } \text{relatively uninformative prior} \right)$$
$$\Rightarrow \hat{\vec{w}}_{\text{MAP}} = \hat{\vec{w}}_{\text{MLE}}$$

$$\lambda \rightarrow \infty$$

$$\hat{\vec{w}}_{\text{MAP}} \rightarrow \vec{0}$$

$$\Rightarrow \hat{y}_i = \vec{w}^T \vec{x}_i \rightarrow 0 \quad \forall i$$

Hmm. that's a bit odd. Why always zero?

replace  $\sum_{j=0}^d w_j^2$  with  $\sum_{j=1}^d w_j^2$   
don't penalize norm of  $w_0$

Now as  $\lambda \rightarrow \infty$

$$w_j \rightarrow 0 \quad \forall j \neq 0$$

$$w_0 \rightarrow \frac{1}{n} \sum_{i=1}^n y_i \quad \left[ \text{Avg of dataset. Why?} \right]$$

HW 2!

So  $\hat{y}_i \rightarrow y_{\text{mean}}$  (Makes sense)

Okay so  $\hat{\vec{w}}_{\text{MAP}} = \underset{\vec{w}}{\text{argmin}} \frac{1}{n} \|Y - Xw\|_2^2 + \frac{\lambda}{n} \|w\|_2^2$

$$\text{So } \frac{\partial L(\vec{w})}{\partial \vec{w}} = \frac{1}{n} \left[ -2Y^T X + 2w^T X^T X \right] + \frac{2\lambda}{n} w^T = 0$$

2

$$\Rightarrow \frac{1}{n} X^T Y = \frac{1}{n} X^T X \vec{w} + \frac{\lambda}{n} \vec{w}$$

$$\Rightarrow (X^T X + \lambda I) \cdot \vec{w} = X^T Y$$

$$\Rightarrow \vec{w}_{MAP} = (X^T X + \lambda I)^{-1} X^T Y$$

Compare with  $\vec{w}_{MLE} = X^T Y = (X^T X)^{-1} X^T Y$

$$(X^T X)^{-1} \rightarrow (X^T X + \lambda I)^{-1}$$

might not be invertible

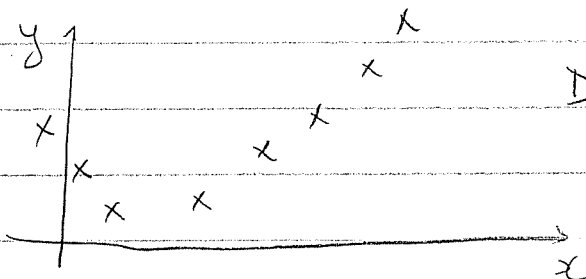
Shifts eigenvalues of  $X^T X$  up  $\uparrow$  by  $\lambda$

$$\left[ \begin{array}{l} \circ \circ (A + \lambda I) \cdot \vec{x} = A\vec{x} + \lambda\vec{x} \\ = (\text{eigenvalue} + \lambda) \cdot \vec{x} \end{array} \right]$$

This is a Bayesian Interpretation of a Linear Algebra Procedure!

### ② Polynomial Regression

Consider 1D data  $x \in \mathbb{R}^1$  for now (for simplicity)



Doesn't look like a linear relationship.

Can we do better?

Linear Regression:  $\hat{y} = w_0 + w_1 x$

Polynomial Regression:  $\hat{y} = w_0 + w_1 x + w_2 x^2 + \dots + w_d x^d$   
(of degree- $d$ )

$$= [w_0 \ w_1 \ \dots \ w_d] \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^d \end{bmatrix}$$

But we can easily rewrite this as.

$$= \vec{w}^T \underbrace{\Phi(x)}$$

higher-dim feature space  
"embedding"

$$\Phi: \mathbb{R}^1 \rightarrow \mathbb{R}^{(d+1)}$$

But  $\hat{y} = \vec{w}^T \Phi(x)$  is still linear in  $\vec{w}$ !

That's all that matters!

Linearity in params ( $\vec{w}$ ) not input ( $\vec{x}$ )

so same steps apply!

In general  $\vec{x} \in \mathbb{R}^d \rightarrow \Phi(\vec{x}) \in \mathbb{R}^D \quad D \neq d$

say  $\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \Phi(\vec{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_1^2 \\ x_2^2 \\ x_3^2 \\ x_1 x_2 \\ x_2 x_3 \\ x_3 x_1 \end{bmatrix}$

Quadratic Feature Augmentation /  
Quadratic "Kernel Map"  
[We will see kernels later]

$\hat{y} = \vec{w}^T \Phi(\vec{x})$  still linear in  $\vec{w}$ !

$\hat{y} = \vec{x}^T A \vec{x} + \vec{c}^T \vec{x}$  Quadratic in  $\vec{x}$

③

# NEW TOPIC: MODEL SELECTION

③ How do we pick  $d$  in 1-dim polynomial regression?  
( $x \in \mathbb{R}^d$ )

Alternatively,

Model 0:  $y^{(0)} = w_0$

Model 1:  $y^{(1)} = w_0 + w_1 x$

⋮

⋮

Model  $D$ :  $y^{(D)} = w_0 + w_1 x + w_2 x^2 + \dots + w_D x^D$

Model Selection: Which model class is better?

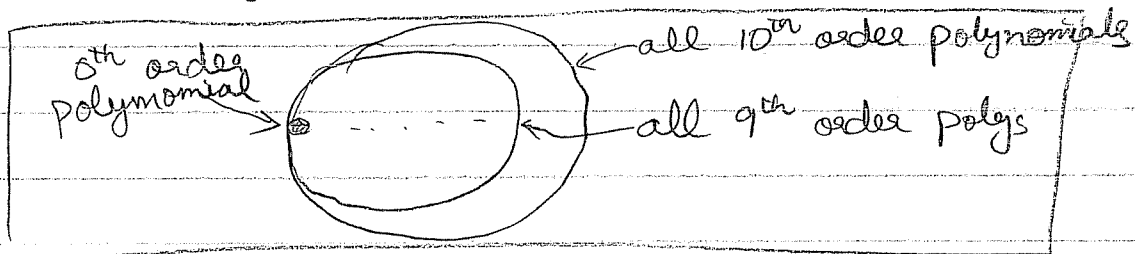
Notice that  $\vec{w}_{OLS}^{(d)} = \underset{\vec{w} \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{n} \sum_{\text{training data}} (y_i - \vec{w}^{(d)T} \vec{x}_i)^2$

So how about?

$d = \underset{\{1, \dots, 100\}}{\operatorname{argmin}} \frac{1}{n} \sum_{\text{training data}} (y_i - \vec{w}_{OLS}^{(d)T} \vec{x}_i)^2$

Won't work. Why?

Problem: Model classes are nested,  $d=100$  will always give lowest training error (or  $E_{\text{train}}$ )



Large Model Classes ( $\equiv$  Strong ML models) will always do better in terms of  $E_{\text{train}}$ .

But why is low train error not a good thing?  
Because that's not what we care about

④ ¶ What do we care about? Expected Error/Loss

$$X, Y \sim P(X, Y) \quad [\text{Unknown}]$$

$$\text{What we want} = \min_w E_{P(X, Y)} [L(y, \hat{y}(x; w))] \quad \begin{array}{l} \uparrow \\ \text{reality} \end{array} \quad \begin{array}{l} \uparrow \quad \uparrow \\ \text{Model Class} \quad \text{Parameters} \end{array}$$

$$= \min_w \int \int L(y, \hat{y}(x; w)) P(x, y) dx dy$$

$\begin{array}{cc} \uparrow & \uparrow \\ \text{Space of } X & \text{Space of } Y \\ \text{Say } \mathbb{R}^d & \text{Say } \mathbb{R}^1 \end{array}$

Problems

- $P(x, y)$  unknown
- Integral hard to solve
- Min of Integral

"Empirical Risk Minimization (ERM)"

$\equiv$  Lets approximate integral with samples  $(x_i, y_i) \sim P(X, Y)$

$$E_{\text{approx}} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{y}(x_i; \tilde{w}))$$



