

Day 5; 9/21/2018

Bias - Variance trade-off

Prediction for  $x = x_0$  squared loss, averaged over all training and future data

$$\text{Error}(x_0) = \underset{x \text{ is fixed, } \epsilon \text{ is random}}{E} [Y - \hat{h}_\theta(x_0)]^2 = E [(Y - h(x_0)) + (h(x_0) - E[\hat{h}_\theta(x_0)]) + (E[\hat{h}_\theta(x_0)] - \hat{h}_\theta(x_0))]^2$$

$$= E [Y - h(x_0)]^2 + \underbrace{E[h(x_0) - E[\hat{h}_\theta(x_0)]]^2}_{\text{irreducible error}} + \underbrace{E[E[\hat{h}_\theta(x_0)] - \hat{h}_\theta(x_0)]^2}_{\text{variance}^2}$$

cross-products cancel out      model mis-specification      variance of the sampling distribution of parameter estimates

[Ave model bias]

[Ave estimation uncertainty]

Specific to linear regression:  $Y = X\beta + \epsilon$ ,  $\text{Var}(\epsilon) = \sigma^2 I$

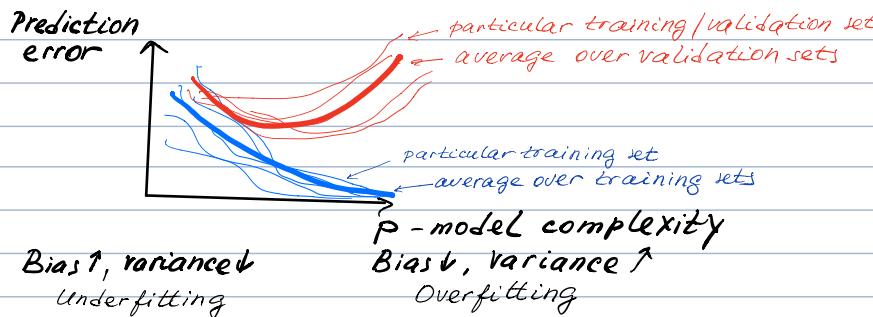
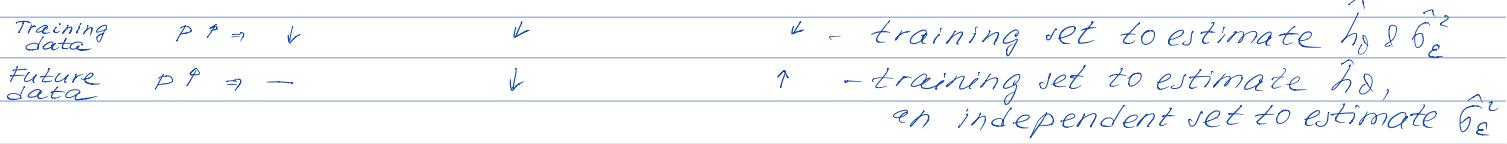
$$\text{Error}(x_0) = \sigma^2 + E[h(x_0) - E[\hat{h}_\theta(x_0)]]^2 + \|X(X^T X)^{-1} x_0\|_2^2 \sigma^2$$

Average error, over all observed values  $x_0$

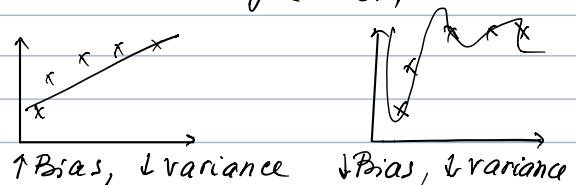
$$\frac{1}{N} \sum_{i=1}^N \text{Error}(x_i) = \sigma^2 + \frac{1}{N} \sum_{i=1}^N [h(x_i) - E[\hat{h}_\theta(x_i)]]^2 + \frac{P}{N} \sigma^2$$

How to estimate this error in practice?

$$\hat{\sigma}_e^2 = \frac{1}{N} \sum_{i=1}^N e_i^2$$



Example: polynomial regression



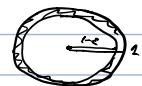
In high dimensions, we are likely to overfit: curse of dimensionality  
In  $p$  dimensions, each observation is unique



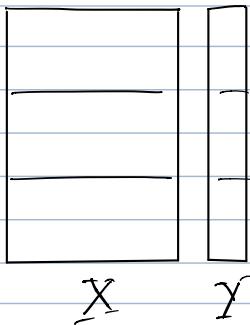
A sphere with radius 1 in  $p$  dimensions: what fraction of the volume lies between  $r=1$  and  $r=1-\epsilon$ ?  $V_p(r) = \text{constant}_p \cdot r^p$

$$\frac{V_p(1) - V_p(1-\epsilon)}{V_p(1)} = 1 - (1-\epsilon)^p \xrightarrow{p \gg 1} 1 \quad (\text{even for small } \epsilon)$$

$\Rightarrow$  most of volume is in the most distant points



Goal: minimize prediction error on future data



- model fitting / training
- model development: evaluate competing models
- model evaluation: evaluate final model

subset selection  
→ regularization

### Subset selection

- (1) Explore the space of candidate models - add feature engineering
- (2) Evaluate the candidate models on the dev set, keep the best performing
- (3) Evaluate the selected model on the evaluation set

Exploring the model space in (1):

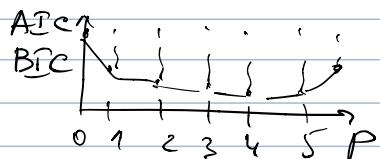
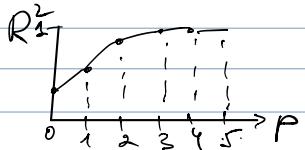
- Exhaustive search of  $\binom{p}{k}$  predictors
- Forward stepwise selection
- Backward stepwise selection

many heuristics exist

Example: forward stepwise selection

• Denote  $M_0$  the null model with no predictors

- For  $k = 0, \dots, p-1$ :
  - consider all  $p-k$  models that augment  $M_k$  with one predictor
  - $M_k \leftarrow$  the best among  $p-k$  models
- select the model among  $M_0, \dots, M_p$  that performs best on dev set



$$AIC = \frac{SSE}{N} + 2p$$

$$BIC = \frac{SSE}{N} + \log N p$$

$$AIC, BIC = -2 \log \text{likelihood} + kp = \frac{\sum (y_i - \hat{y}_i)^2}{\hat{\sigma}^2} + kp$$

$\hat{\sigma}^2 \xrightarrow{\text{estimated using the saturated model}}$

If  $N$  is relatively small - cross-validation

Iterate partitioning into train and dev sets



Problem: each iteration does not always have the same best model.

"Consensus" model may not perform as well

10-fold cross-validation: repeat 10 folds

Frequently used as an intermediate parameter tuning step

Recall: Normal equations  $\hat{\delta} = \arg \min_{\delta} (Y - X\delta)^T(Y - X\delta)$

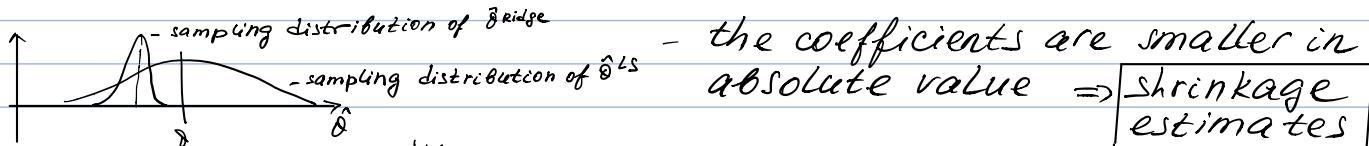
Analytical solution:  $\frac{d J(\delta)}{d \delta} = \frac{d}{d \delta} (Y - X\delta)^T(Y - X\delta) \stackrel{\text{set to 0}}{=} 0$

Properties of  $\hat{\delta}$ :  $E\{\hat{\delta}\} = E\left\{\left(X^T X\right)^{-1} X^T (X\delta + \epsilon)\right\} = \delta$  singular if  $N < P$   
or collinear columns

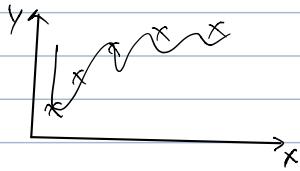
$$\text{Var}\{\hat{\delta}\} = \text{Var}\left\{\left(X^T X\right)^{-1} X^T (X\delta + \epsilon)\right\} = \sigma^2 (X^T X)^{-1} X^T I X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

I.e.:  $\hat{\delta}$  is a linear combination of  $y$ , is unbiased;  $\text{Var} = \sigma^2 (X^T X)^{-1}$   
 - can show that has min variance among all unbiased estimators  
 (Gauss-Markov theorem)

Can improve the estimation when  $X^T X$  is (nearly) singular:  
 define  $\hat{\delta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$       - biased, but smaller variance } invert  
 } pre-defined constant  $\lambda > 0$



Intuition:  
 polynomial regression



$$h_{\delta}(x) = \delta_0 + \delta_1 x + \delta_2 x^2 + \delta_3 x^3 + \delta_4 x^4 \quad (\times)$$

If we shrink  $\delta_3$  and  $\delta_4$  closer to 0  
 $\Rightarrow h_{\delta}(x)$  is closer to quadratic

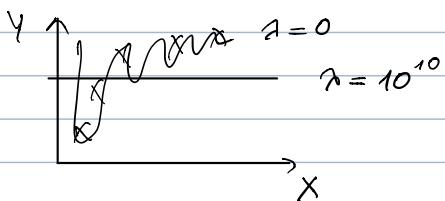
The same solution can be obtained using  
 (1) penalized estimation  
 (2) constrained optimization

(1)  $\hat{\delta} = \arg \min_{\delta} \underbrace{(Y - X\delta)^T(Y - X\delta)}_{\text{New loss function } J_{\lambda}(\delta)} + \lambda \|\delta\|_2^2 = \sum_{j=1}^p \delta_j^2$

Solution as function of  $\lambda$ :

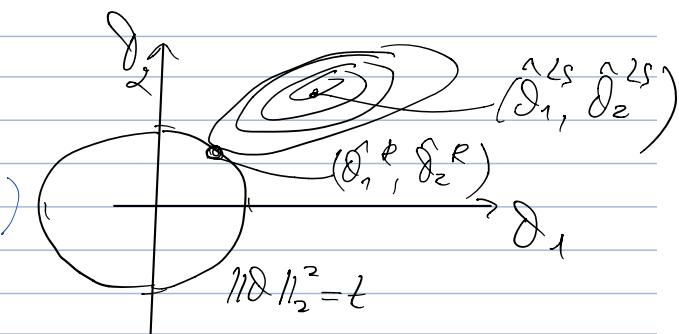
$$\frac{d}{d \delta} J_{\lambda}(\delta) = 0 \quad X^T X \delta - X^T y + \lambda \delta = 0 \quad \Rightarrow (X^T X + \lambda I) \delta - X^T y = 0$$

$\lambda$  is a hyperparameter



(2)  $\hat{\delta} = \arg \min_{\delta} (Y - X\delta)^T(Y - X\delta), \text{ s.t. } \|\delta\|_2^2 \leq t$

$$(t = \|\hat{\delta}(2)\|^2)$$



Even though we have an analytical solution, can also solve numerically by gradient descent

# Numeric parameter estimation in Ridge regression: Batch gradient descent

- Initialize  $\delta$
- Repeat until convergence }

- For  $j = 1, \dots, P$  }

$$\delta_j \leftarrow \delta_j + \lambda \frac{d}{d\delta_j} \left[ \sum_{i=1}^N (y_i - x^T \delta)^2 + \lambda \sum_{k=1}^P \delta_k^2 \right]$$

}

Details of the update:

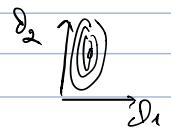
$$\delta_j \leftarrow \delta_j - \lambda \left[ 2 \sum_{i=1}^N (y_i - x^T \delta) x_{ij} + 2\lambda \delta_j \right]$$

$$\delta_j \leftarrow \delta_j (1 - 2\lambda) - \lambda \sum_{i=1}^N (y_i - x^T \delta) x_{ij}$$

usually  $< 1$       same as in least squares  
 $\Rightarrow$  shrinkage

Practical considerations for faster convergence:

- scaling features (parameter values are comparable)
- centering features (no intercept)
- trace plots of  $J(\delta)$



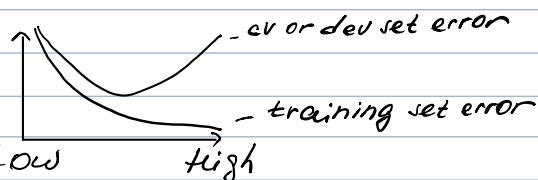
$J(\delta)$

# iterations

Q: how to select  $\lambda$ ? Try a range of  $\lambda$ ; use dev set or cross-valid. to find the optimal point

model complexity (# params)

$\lambda$	Low	High
Bias	High	Low
Variance	Low	High



$N \uparrow$  → fix high variance

$P \downarrow$  → -||-

$P \uparrow$  → fix high bias

$\lambda \downarrow$  → fix high bias

$\lambda \uparrow$  → fix high variance