

S&DS 365 / 565  
**Data Mining and Machine Learning**

# **Cross Validation**

**Yale**

# Model Selection

For purposes of prediction, **minimizing test error** is priority.

Recall our two error metrics for evaluating predictions  $\hat{f}(x_i)$ :  
*(empirical risk)*

- Regression:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

*mean squared error*

- Classification: *misclassification*

$$Err = \frac{1}{n} \sum_{i=1}^n \mathbb{1} \left\{ \hat{f}(x_i) \neq y_i \right\}$$

# Bias-Variance Tradeoff

e.g. (Classification)

pop risk

$$\text{Define } R(f) = \mathbb{E} \mathbb{1}(f(X) \neq Y) = \mathbb{P}(f(X) \neq Y)$$

estimated  $\hat{f}$

Suppose that  $\hat{f}$  comes from some model class  $M$

true  $f^*$   
optimal error  $f^*_{\text{err}}$   
Let  $f^* \in \arg \min R(f)$  such that  $f \in M$

$$R(\hat{f}) - R(\bar{f}) = \underbrace{R(\hat{f}) - R(f^*)}_{\text{Excess Risk}} + \underbrace{R(f^*) - R(\bar{f})}_{\text{variance}} + \underbrace{R(\bar{f}) - R(f^*)}_{\text{bias}}$$

bias: ① bias on function class

I think my data fit into a linear model  
hope bias ↓ when degree of polynomial ↑

② KNN when  $K=1$ , no bias ↑  
when  $K \uparrow$  bias ↑ function more smooth

bias is NOT random

$\left\{ \begin{array}{ll} f^* & \text{true } f \\ \hat{f} & \text{estimated } f \\ \bar{f} & \text{optimal } f \end{array} \right\}$  we know

$$\text{optimal } f: \hat{f}^{\text{bias}} = \mathbb{1}[P(y=1|x) > P(y=0|x)]$$

proper learning: assume  $\hat{f} \in M$   
improper learning:  $\hat{f} \in$  bigger function class  
variance ↑ error ↑  
bias ↑ error ↑

variance: how much your function  $f$  fluctuate from one random data sample

to another

e.g. when  $K=1$ , change sample  $\rightarrow f$  change dramatically  
when  $K \uparrow$   $f$  change smaller

variance is random coz  $f$  is random  
we generate a  $f$  from a random sample

# Bias-Variance Tradeoff

often bias↑ variance↓



Eg (Regression case)

Model

Given  $Y = f(X) + \epsilon$ , (where  $E(\epsilon) = 0$  and  $\text{Var}(\epsilon) = \sigma^2$ , consider a predictor  $\hat{f}$ .

Assumption for calculate bias

Expected MSE for predicting a new  $Y$  at  $X = x$  can be decomposed into:

$$\begin{aligned} E(\text{MSE}) &= E[(Y - \hat{f}(x))^2] = \text{Var}(\hat{f}(x)) + [\text{Bias}(\hat{f}(x))]^2 + \sigma^2 \\ &= E\left[\left[f(x) + \epsilon - E\hat{f}(x) + E\hat{f}(x) - \hat{f}(x)\right]^2\right] \end{aligned}$$

excess risk

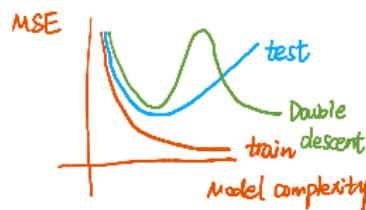
↓  
noise

2 random things, other are numbers

$\epsilon$  is independent of other thing, pull out  $\epsilon$  get  $\sigma^2$

$$\text{Var}(\hat{f}(x)) = E[(E\hat{f}(x) - \hat{f}(x))^2]$$

$$\text{Bias } \hat{f}(x) = [f(x) - E\hat{f}(x)]^2$$



# Bias-Variance Tradeoff

$$E[(Y - \hat{f}(x))^2] = \text{Var}(\hat{f}(x)) + [\text{Bias}(\hat{f}(x))]^2 + \sigma^2$$

- $\text{Var}(\hat{f})$  is the amount of variability in our predictor with different training set.
- $\text{Bias}(\hat{f})$  is the systematic error introduced by model approximation.  
*how much inflexibility of model*
- $\sigma^2$  is irreducible error, inherent in the error term  $\epsilon$ .  
*noise*

# Bias-Variance Tradeoff

$$E[(Y - \hat{f}(x))^2] = \text{Var}(\hat{f}(x)) + [\text{Bias}(\hat{f}(x))]^2 + \sigma^2$$

- $\text{Var}(\hat{f})$  is the amount of variability in our predictor with different training set. Increases with increasing model flexibility. (complexity)
- $\text{Bias}(\hat{f})$  is the systematic error introduced by model approximation. Decreases with increasing model flexibility.
- $\sigma^2$  is *irreducible error*, inherent in the error term  $\epsilon$ . Cannot get rid of this!

Need to balance bias and variance.

# Bias-Variance?

Let's take a step back. Compare this equation

$$E[(Y - \hat{f}(x))^2] = \text{Var}(\hat{f}(x)) + [\text{Bias}(\hat{f}(x))]^2 + \sigma^2$$

with what we are more much more familiar with: if  $Y \sim N(\theta, \sigma^2)$ , then with  $\hat{\theta} = Y$ , we have

$$\text{Var}(\hat{\theta}) = \sigma^2, \text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

$$kNN \quad k=1 \Rightarrow \text{Bias} = \theta - \theta = 0$$

# Bias-Variance?

Let's take a step back. Compare this equation

$$E[(Y - \hat{f}(x))^2] = \text{Var}(\hat{f}(x)) + [\text{Bias}(\hat{f}(x))]^2 + \sigma^2$$

with what we are more much more familiar with: if  $Y \sim N(\theta, \sigma^2)$ , then with  $\hat{\theta} = Y$ , we have

take  $\hat{\theta}$  is very small  $\hat{\theta}=0$

$$\text{Var}(\hat{\theta}) = 0$$

$$\text{Bias}(\hat{\theta}) = \theta - 0 = \theta$$

$$\text{Var}(\hat{\theta}) = \sigma^2, \text{Bias}(\hat{\theta}) = \theta$$

What does it mean to talk about bias/variance of  $\hat{f}$ ?

# Bias-Variance

How does this apply to the  $k$ -NN algorithm?

$k=1$  , no bias

$k \uparrow$    bias $\uparrow$    variance $\downarrow$

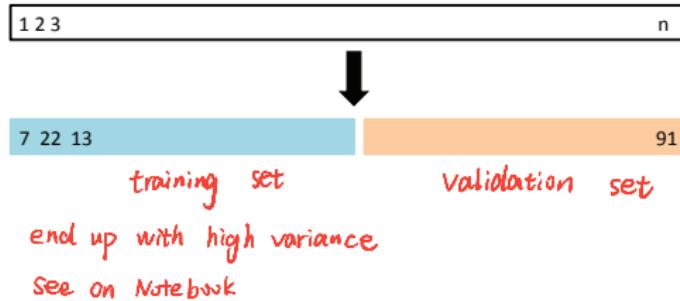
# Cross-Validation

**Cross-validation** is an intuitive, widely-applicable approach for:

- model assessment
- model selection

# Validation Sets

We've been doing this:



- ➊ Divide dataset randomly into a training set and a validation set.
- ➋ Fit the model on the training set.
- ➌ Use the validation set to obtain estimated test error.
- ➍ Repeat!

# Validation Sets

- highly variable validation error
- only uses a fraction of the training set

# Leave-One-Out Cross-Validation

How do we use more data to train with?

- Use a tiny validation set (e.g.  $(x_1, y_1)$ ) *one sample*
- Train with the rest (e.g.  $\{(x_2, y_2), \dots, (x_n, y_n)\}$ )

How do we feel about error rate evaluated on a single observation?

# Leave-One-Out Cross-Validation

How do we use more data to train with?

- Use a tiny validation set (e.g.  $(x_1, y_1)$ )
- Train with the rest (e.g.  $\{(x_2, y_2), \dots, (x_n, y_n)\}$ )

How do we feel about error rate evaluated on a single observation?

Not good, but we can iterate through the dataset, **each time using a different  $(x_i, y_i)$**  as the validation set and obtaining an error  $MSE_i$ .

train on rest samples ( $n-1$ )

# Leave-One-Out Cross-Validation

Obs	Iteration						
	1	2	3	4	...	n	
1	valid	train	train	train	...	train	
2	train	valid	train	train	...	train	
3	train	train	valid	train	...	train	
4	train	train	train	valid	...	train	
...	...	...	...	...	...	...	...
$n$	train	train	...	...	...	...	valid
MSE	$MSE_1$	$MSE_2$	$MSE_3$	$MSE_4$	...	$MSE_n$	

# Leave-One-Out Cross-Validation

n-fold cross validation

n is iteration times = sample size  
if sample size is small, do that

LOOCV estimate of test error is given by:

$$CV_{(n)} = \frac{1}{n} \sum_i MSE_i$$

# k-fold Cross-Validation

A potentially faster approach: *if sample size is large*

- Randomly divide the dataset into  $k$  folds.
- For  $b = 1, \dots, k$ :
  - ▶ Use  $b$ -th fold (“batch”) as validation set.
  - ▶ Use everything else as training set.
  - ▶ Compute validation error on  $b$ -th fold.
- Estimate test error using:

$$CV_{(k)} = \sum_b \frac{n_b}{n} MSE_b,$$

*usually  $n_b$  is same for  $k$  folds*

where  $n_b$  is the total # observations in the  $b$ -th fold, and  $n$  is the total # observations in the entire dataset.

# k-fold Cross-Validation

Obs	Iteration						
	1	2	3	4	...	$k$	
1	valid	train	train	train	...	train	fold 1
2	valid	train	train	train	...	train	
3	valid	train	train	train	...	train	
4	train	valid	train	train	...	train	
...	...	...	...	...	...	...	...
$n - 2$	train	train	...	...	...	valid	fold $\checkmark k$
$n - 1$	train	train	...	...	...	valid	
$n$	train	train	...	...	...	valid	
MSE	$MSE_1$	$MSE_2$	$MSE_3$	$MSE_4$	...	$MSE_k$	

# k-fold Cross-Validation

Obs	Iteration						
	1	2	3	4	...	$k$	
1	valid	train	train	train	...	train	fold 1
2	valid	train	train	train	...	train	
3	valid	train	train	train	...	train	
4	train	valid	train	train	...	train	
...	...	...	...	...	...	...	...
$n - 2$	train	train	...	...	...	valid	fold $v$
$n - 1$	train	train	...	...	...	valid	
$n$	train	train	...	...	...	valid	
MSE	$MSE_1$	$MSE_2$	$MSE_3$	$MSE_4$	...	$MSE_k$	

$n$ -fold CV is just LOOCV.

# A slick shortcut

Suppose that the fitted values can be written  $\hat{Y} = LY$  where  $L$  is an  $n \times n$  matrix.

$$\hat{Y} = \underbrace{X(X^T X)^{-1} X^T Y}_L$$

*only*  
This is the case for least squares regression

Then the leave-one-out-cross-validation error is

$$R_{LOO} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_{(-i)})^2$$

$$\frac{1}{n} \sum_{i=1}^n \left( \frac{Y_i - \hat{Y}_i}{1 - L_{ii}} \right)^2$$

where  $L_{ii}$  is the  $i$ th diagonal entry.

# A slick shortcut

Suppose that the fitted values can be written  $\hat{Y} = LY$  where  $L$  is an  $n \times n$  matrix.

This is the case for least squares regression

Then the leave-one-out-cross-validation error is

$$\text{Empirical Risk} = \text{Mean squared error}$$

$$R_{LOO} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_{(-i)})^2$$

↑ true  $Y$  of  $i$ th observ.      ↑ estimate of  $i$ th observ.  
    trained on excluding  $i$ th observ.

$$= \frac{1}{n} \sum_{i=1}^n \left( \frac{Y_i - \hat{Y}_i}{1 - L_{ii}} \right)^2$$

↑  $\hat{Y}$  trained on full data set

where  $L_{ii}$  is the  $i$ th diagonal entry.

So, no need to fit  $n$  regressions!

only need 1 regression

# Model Selection

So far, we've used it to estimate the test error. It's also useful for model selection.

Suppose we are interested in comparing the following 3 models:

Model 1:

Linear model  $\widehat{medv} = \widehat{\beta}_0 + \widehat{\beta}_1 lstat$

Model 2:

Log model  $\widehat{\log(medv)} = \widehat{\beta}_0 + \widehat{\beta}_1 lstat$

Model 3:

quadratic model  $\widehat{medv} = \widehat{\beta}_0 + \widehat{\beta}_1 lstat + \widehat{\beta}_2 lstat^2$

# Model Selection

So far, we've used it to estimate the test error. It's also useful for model selection.

Suppose we are interested in comparing the following 3 models:

Model 1:

$$\widehat{medv} = \widehat{\beta}_0 + \widehat{\beta}_1 lstat$$

Model 2:

$$\widehat{\log(medv)} = \widehat{\beta}_0 + \widehat{\beta}_1 lstat$$

Model 3:

$$\widehat{medv} = \widehat{\beta}_0 + \widehat{\beta}_1 lstat + \widehat{\beta}_2 lstat^2$$

We can use cross-validation to estimate the test error for each of these models, and select the model with the lowest test error.

# Model Selection: Workflow

In order to do model selection, you need a set of models  $\{M_1, \dots, M_m\}$ :

- ① for each model  $M_i$ :
  - ① calculate  $k$ -fold CV validation error *learn  $m \times k$  models*
- ② choose model with the lowest  $k$ -fold CV validation error  $M_{\min}$
- ③ apply  $M_{\min}$  to the full dataset

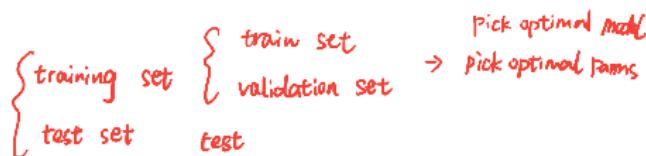
# Model Selection: Workflow

In order to do model selection, you need a set of models  $\{M_1, \dots, M_m\}$ :

- ① for each model  $M_i$ :
  - ① calculate  $k$ -fold CV validation error
- ② choose model with the lowest  $k$ -fold CV validation error  $M_{\min}$
- ③ apply  $M_{\min}$  to the full dataset *to learn params of  $M_{\min}$*   
*can't report this lowest cv error as model performance measure*
- ④ Don't necessarily pick the lowest always. If have a lot of models that are similar, might pick the **lowest complexity** model that's close enough to the smallest error. For KNN low-complexity means **higher  $k$** . For polynomials, **lower degree**  
*cos you use validation set to pick model and still use validation set for optimization*  
*may cause under report of error (over optimistic about models)*  
*pick model with 1 standard deviation of lowest CV error*

## A safer way

We have split into two datasets.



If you want an accurate assessment of the true test error need to split off a third set that is the **true** test set.

The result of model selection based on just validation will be **over optimistic**

To have an unbiased estimate of the final model error, need to use **third untouched** test dataset.

**untouched** means the dataset **was not ever used to pick a model**

With experience can develop a sense to connect validation set error to true error and can just ignore the test set