

UVA CS 6316/4501

– Fall 2016

Machine Learning

Lecture 15: K-nearest-neighbor Classifier / Bias-Variance Tradeoff

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

Rough Plan

- HW5 is due on Sat
- HW6 will have two coding Q – (image + audio)
- HW7 will be sample final questions
 - Final will be most contents after midterm
- Midterm grade will be released this evening
 - Mean 78 / Median 79 / Max 95
- Final will be in-class / close note / @ Dec5th

Where are we ? →

Five major sections of this course

- ❑ Regression (supervised)
- ❑ Classification (supervised)
- ❑ Unsupervised models
- ❑ Learning theory
- ❑ Graphical models

Where are we ? →

Three major sections for classification

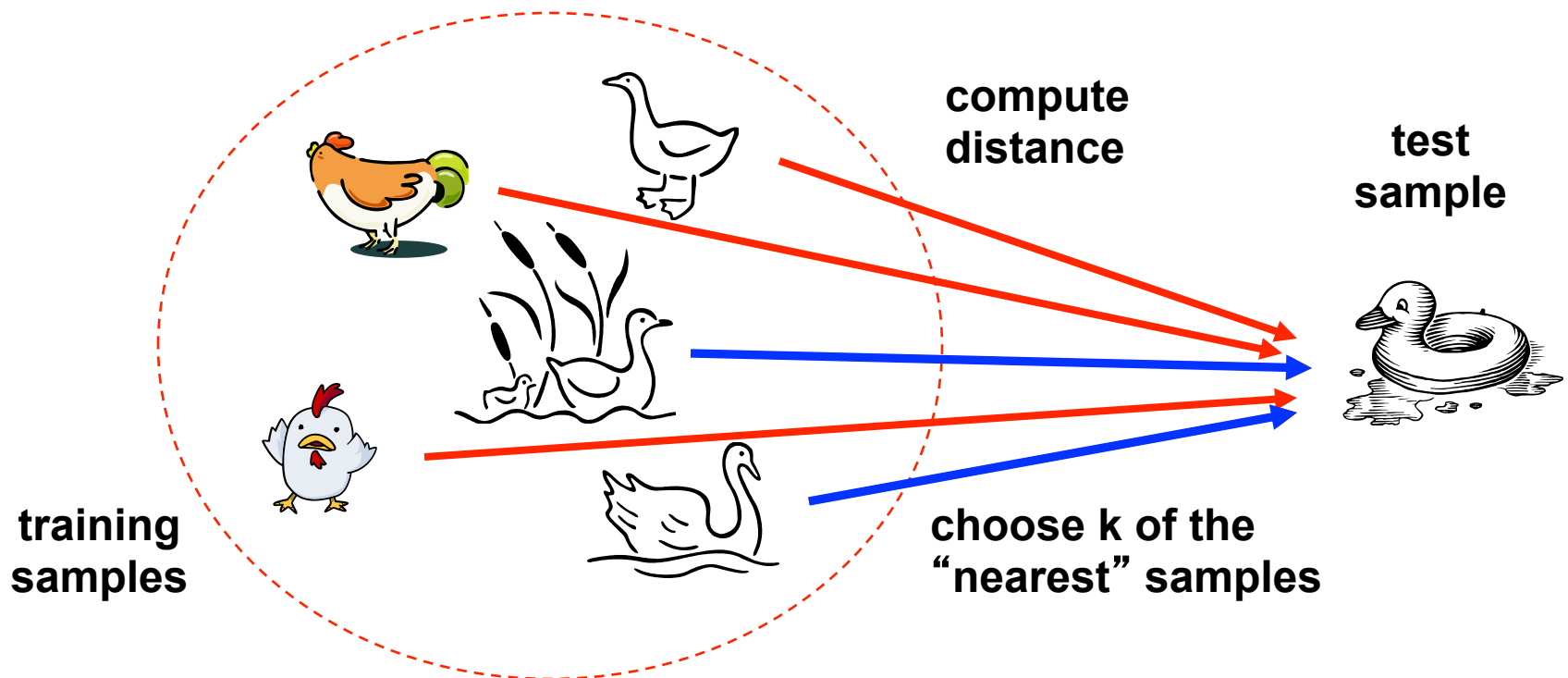
- We can divide the large variety of classification approaches into **roughly three major types**
 1. Discriminative
 - directly estimate a decision rule/boundary
 - e.g., **logistic regression**, support vector machine, decisionTree
 2. Generative:
 - build a generative statistical model
 - e.g., **naïve bayes classifier**, Bayesian networks
 3. Instance based classifiers
 - Use observation directly (no models)
 - e.g. **K nearest neighbors**

Today :

- 
- ✓ K-nearest neighbor
 - ✓ Model Selection / Bias Variance Tradeoff
- 

Nearest neighbor classifiers

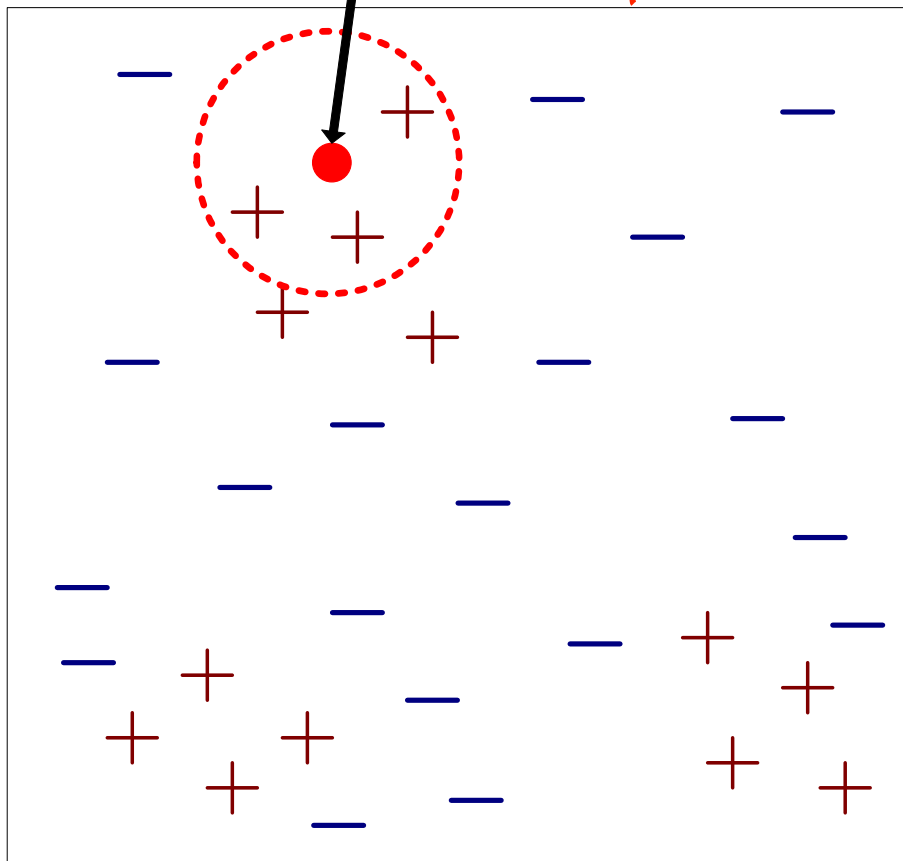
- Basic idea:
 - If it **walks** like a duck, **quacks** like a duck, then it's probably a duck



Nearest neighbor classifiers

Unknown record

$k=3$

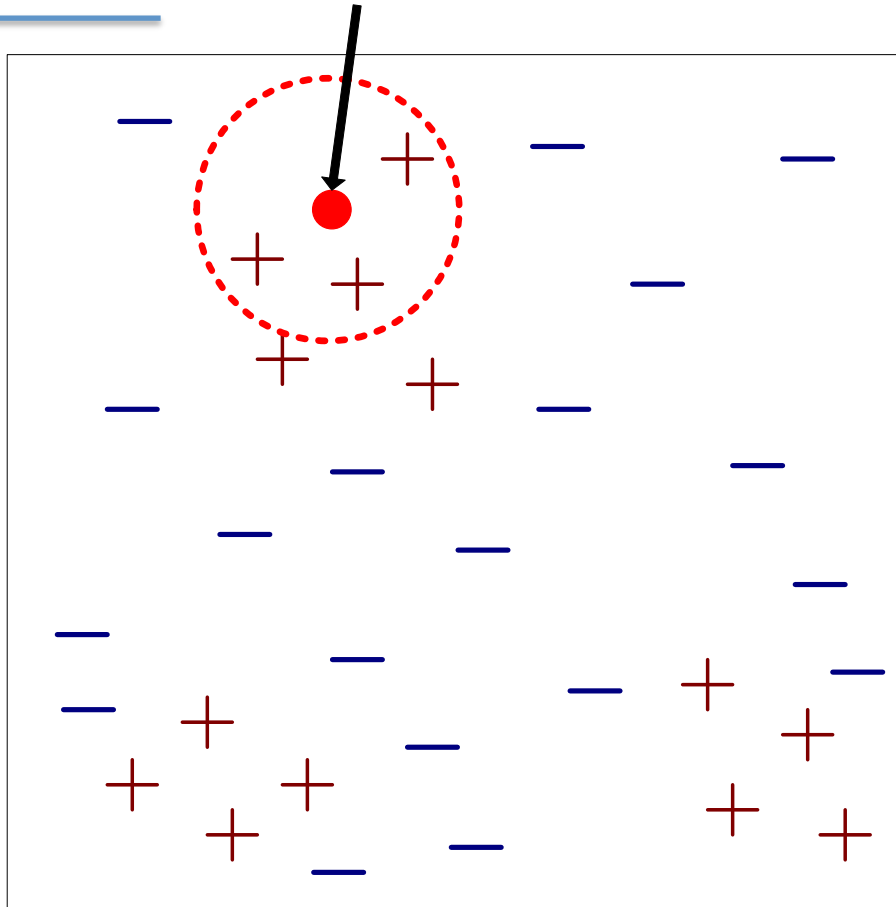


Requires **three** inputs:

1. The set of stored training samples
2. Distance metric to compute distance between samples
3. The value of k , i.e., the number of nearest neighbors to retrieve

Nearest neighbor classifiers

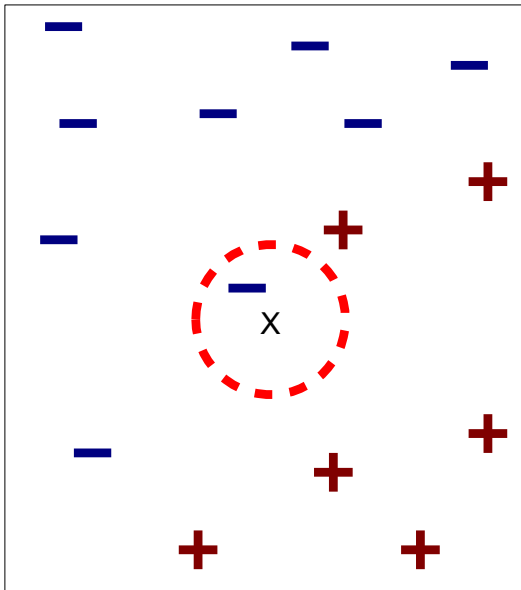
Unknown record



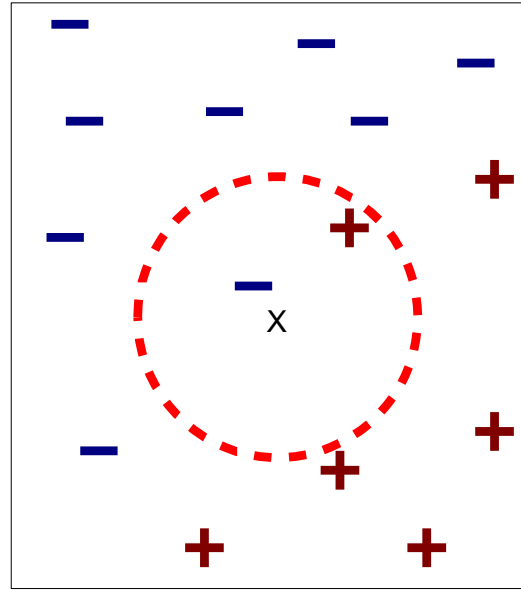
To classify **unknown** sample:

1. Compute distance to other training records
2. Identify k nearest neighbors
3. Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

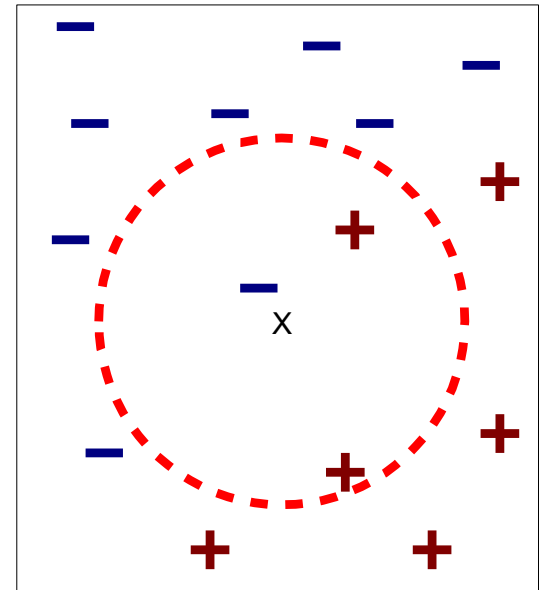
Definition of nearest neighbor



(a) 1-nearest neighbor



(b) 2-nearest neighbor



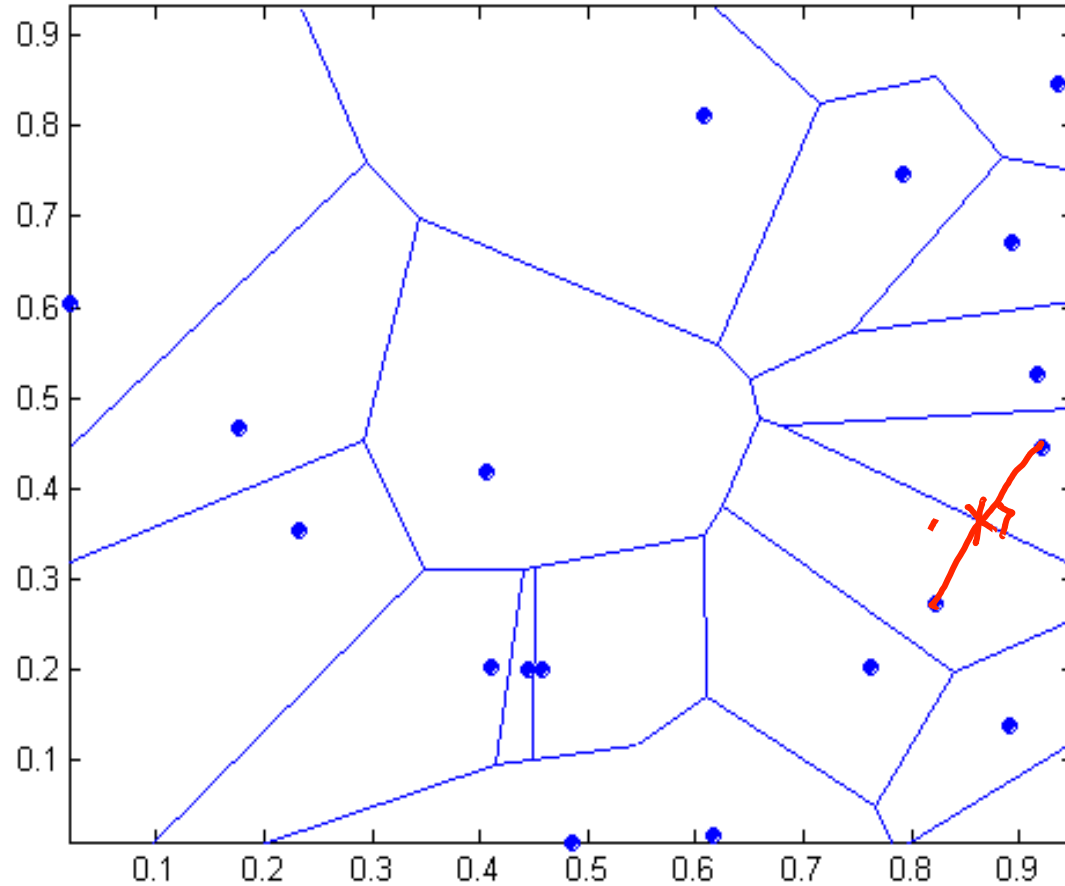
(c) 3-nearest neighbor

k -nearest neighbors of a sample x are datapoints that have the k smallest distances to x

K ?

1-nearest neighbor

Voronoi diagram:
partitioning of a
plane into
regions based
on distance to
points in a
specific subset
of the plane.



Nearest neighbor classification

- Compute **distance** between two points:
 - For instance, Euclidean distance

*e.g. cosine distance
for text*

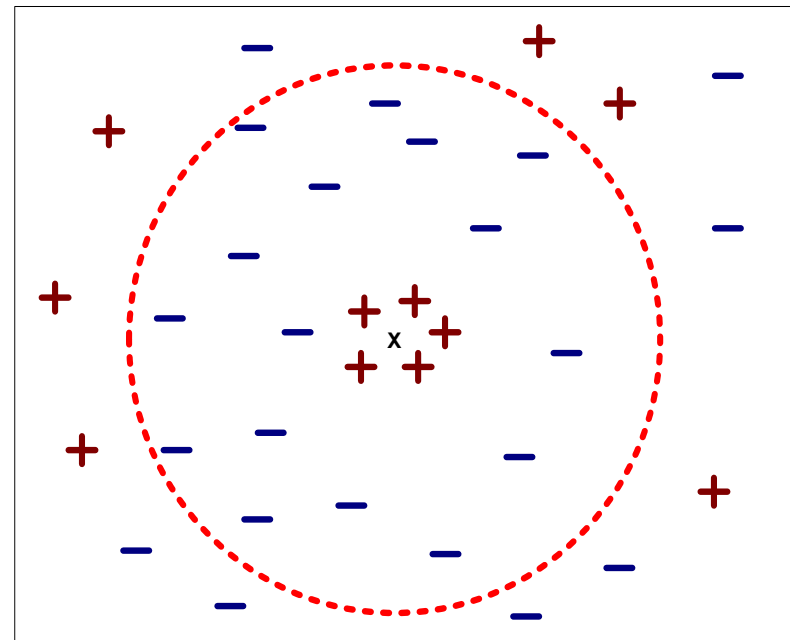
$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_i (x_i - y_i)^2}$$

- **Options** for determining the class from nearest neighbor list
 - Take **majority vote** of class labels among the k -nearest neighbors
 - **Weight the votes** according to distance
 - example: weight factor $w = 1 / d^2$

Nearest neighbor classification

- Choosing the value of k :
 - If k is too small, sensitive to noise points
 - If k is too large, neighborhood may include points from other classes

\downarrow P
 regression
 \downarrow P small
 \downarrow P large



k large

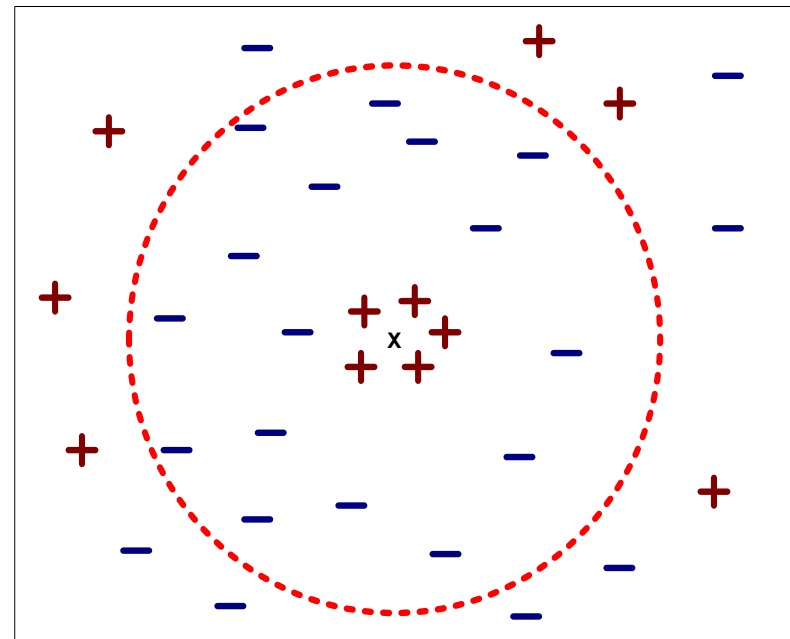
k small
flexible

Nearest neighbor classification

- Choosing the value of k :
 - If k is too small, sensitive to noise points
 - If k is too large, neighborhood may include points from other classes

$k \downarrow$ flexible / varies a lot

$k \uparrow$ smooth / varies little



Nearest neighbor classification

- Scaling issues
 - Attributes may have to be scaled to prevent distance measures from being dominated by one of the attributes
 - Example:
 - height of a person may vary from 1.5 m to 1.8 m
 - weight of a person may vary from 90 lb to 300 lb
 - income of a person may vary from \$10K to \$1M

Nearest neighbor classification

- k -Nearest neighbor classifier is a **lazy** learner
 - Does not build model explicitly.
 - Classifying unknown samples is relatively expensive. $X_{test} \rightarrow O(n) + O(\text{sort}-n-k)$
- k -Nearest neighbor classifier is a **local** model, vs. **global** model of linear classifiers.

Nearest neighbor classification

- k -Nearest neighbor classifier is a **lazy** learner

- Does **not** build **model** explicitly.

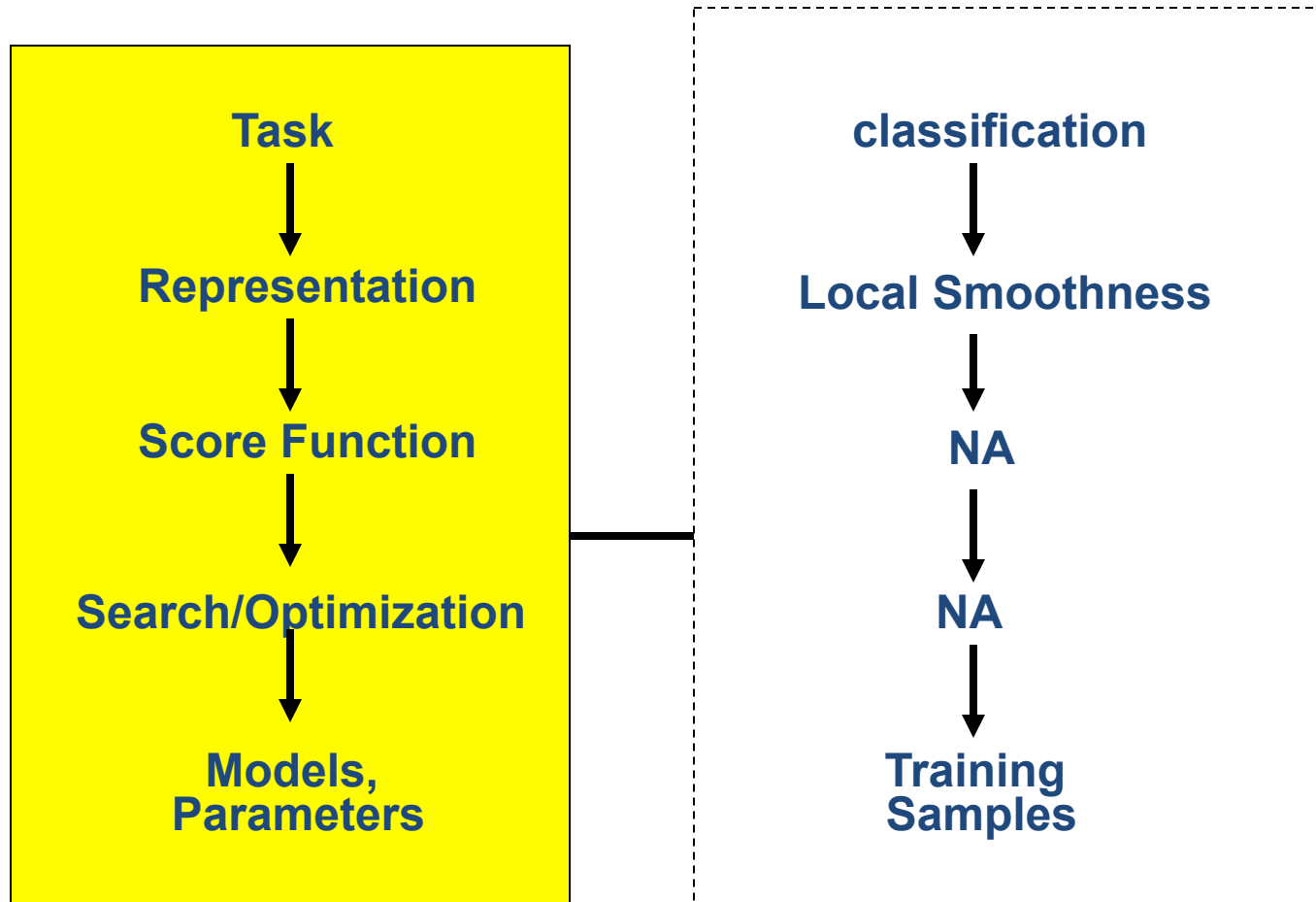
$k (X_{ts}, X_{tri})$

- Classifying unknown samples is relatively expensive.

testing { KNN: num_train / all train samples
SVM: num-support vectors points

- k -Nearest neighbor classifier is a **local** model, vs. **global** model of linear classifiers.

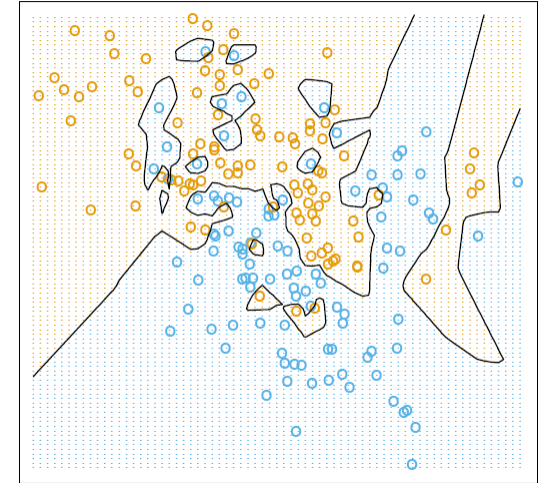
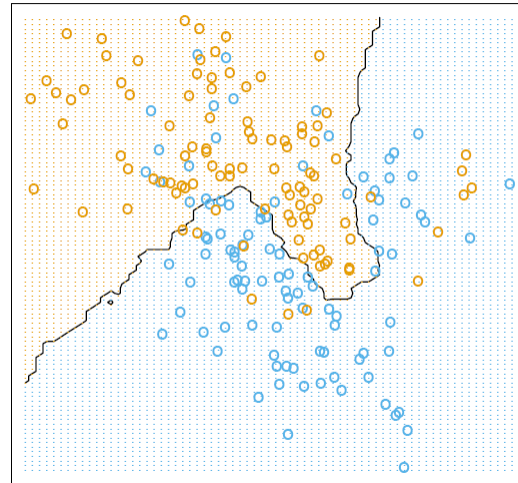
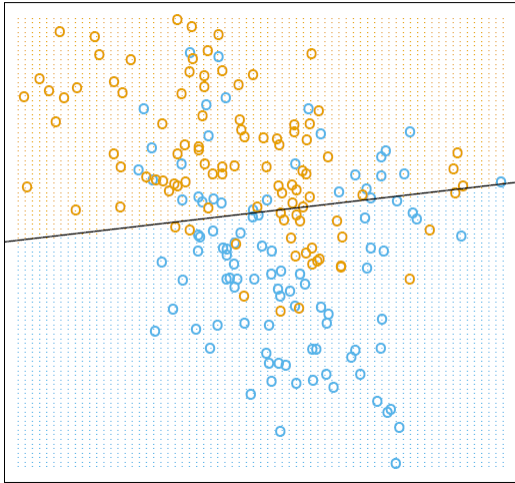
K-Nearest Neighbor



Decision boundaries in global vs. local models

K=15

K=1



linear regression

- global
- stable
- can be inaccurate

15-nearest neighbor

- K acts as a smoother

1-nearest neighbor

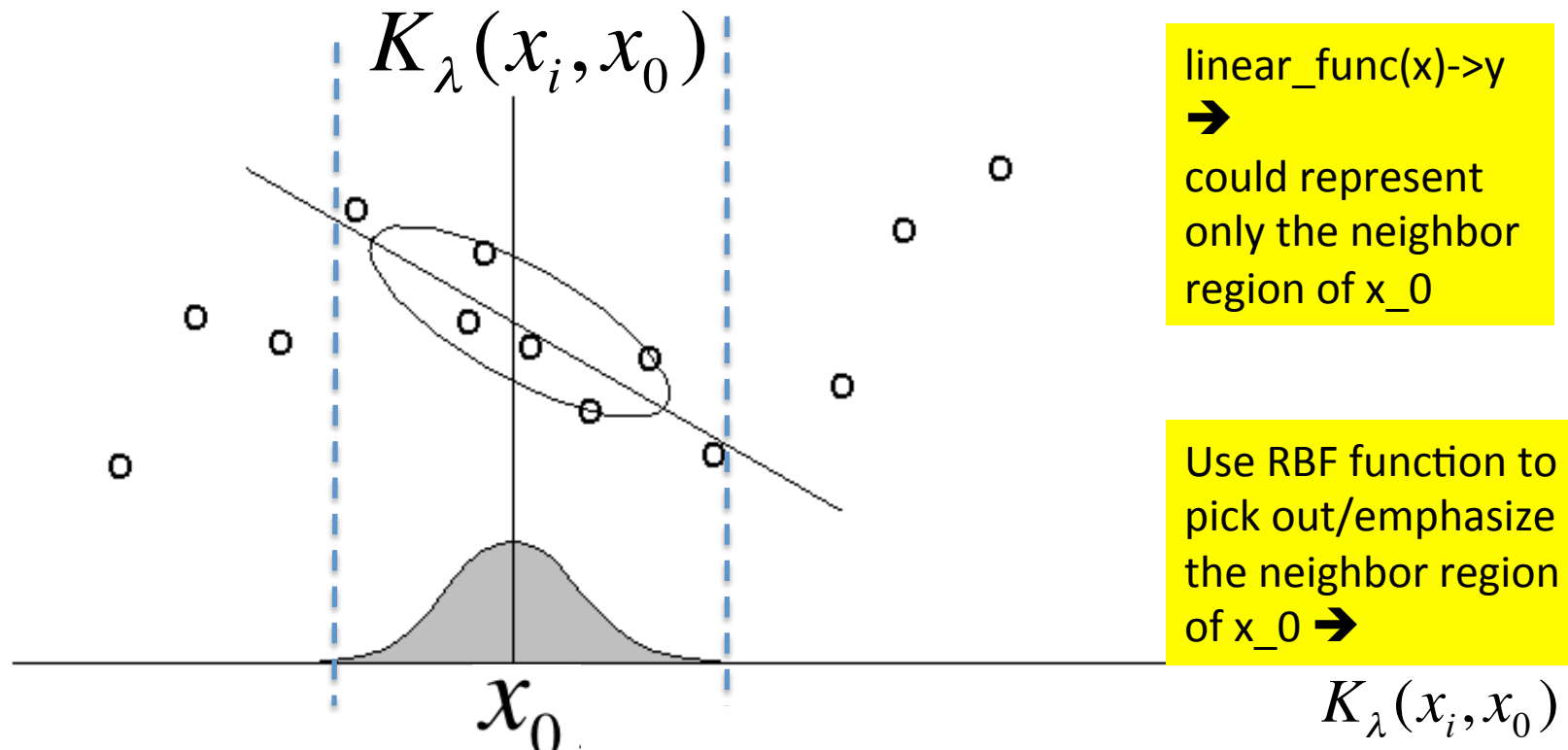
- local
- accurate
- unstable

What ultimately matters: **GENERALIZATION**

vs. KNN for regression (mean of KNN)

Vs. Locally weighted regression

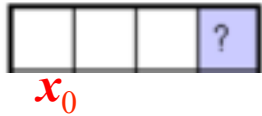
- *aka* locally weighted regression, locally linear regression, LOESS, ...



11/9/16 **Figure 2:** In locally weighted regression, points are weighted by proximity to the current x in question using a kernel. A regression is then computed using the weighted points.

Vs. Locally weighted regression

- Separate weighted least squares **at each target point \mathbf{x}_0** :




$$\min_{\alpha(x_0), \beta(x_0)} \sum_{i=1}^N K_{\lambda}(x_i, x_0) [y_i - \alpha(x_0) - \beta(x_0)x_i]^2$$

$$\hat{f}(x_0) = \hat{\alpha}(x_0) + \hat{\beta}(x_0)x_0$$

$$K_{\tau}(\mathbf{x}_i, \mathbf{x}_0) = \exp\left(-\frac{(\mathbf{x}_i - \mathbf{x}_0)^2}{2\tau^2}\right)$$

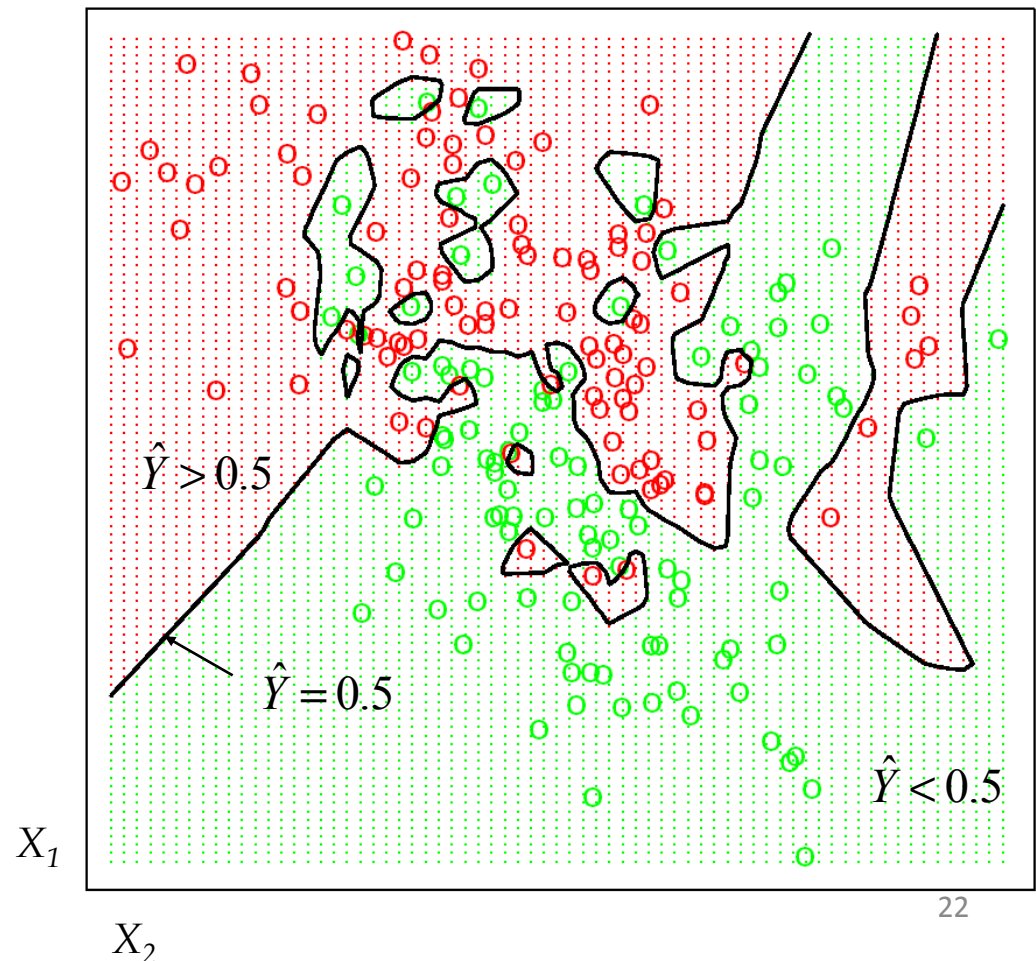
Today :

- ✓ K-nearest neighbor
- ✓ Model Selection / Bias Variance Tradeoff
-  ✓ EPE
- ✓ Decomposition of MSE
- ✓ Bias-Variance tradeoff
- ✓ High bias ? High variance ? How to respond ?

e.g. Training Error from KNN, Lesson Learned

- When $k = 1$,
- No misclassifications (on training): **Overtraining**
- Minimizing training error is not always good (e.g., 1-NN)

1-nearest neighbor averaging



Statistical Decision Theory

- Random input vector: X
- Random output variable: Y
- Joint distribution: $\Pr(X, Y)$
- Loss function $L(Y, f(X))$
- Expected prediction error (EPE):

$$\text{EPE}(f) = \mathbb{E}(L(Y, f(X))) = \int L(y, f(x)) \Pr(dx, dy)$$

$$\text{e.g.} = \int (y - f(x))^2 \Pr(dx, dy)$$

Consider population distribution

Expected prediction error (EPE)

Consider joint distribution

$$\text{EPE}(f) = \mathbb{E}(L(Y, f(X))) = \int L(y, f(x)) \Pr(dx, dy)$$

- For L2 loss: e.g. $= \int (y - f(x))^2 \Pr(dx, dy)$

under L2 loss, best estimator for EPE (Theoretically) is :

Conditional
mean

$$\hat{f}(x) = \mathbb{E}(Y | X = x)$$

e.g. KNN

NN methods are the direct implementation (approximation)

- For 0-1 loss: $L(k, \ell) = 1 - d_{kl}$

Bayes Classifier

$$\hat{f}(X) = C_k \text{ if}$$

$$\Pr(C_k | X = x) = \max_{g \in C} \Pr(g | X = x)$$

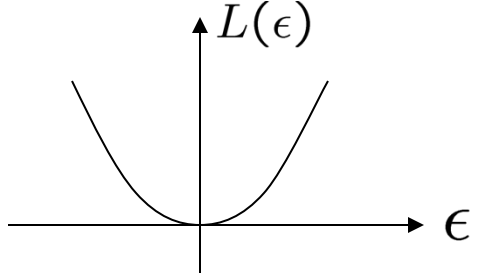
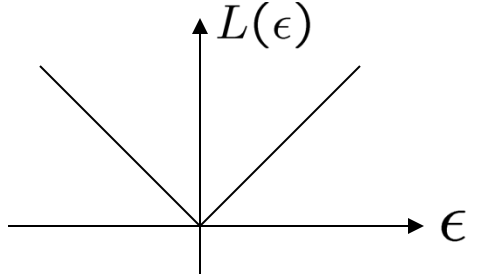
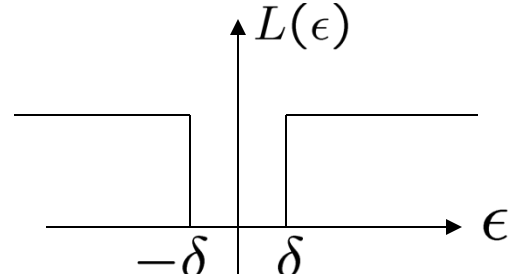
KNN FOR MINIMIZING EPE

- We know under L2 loss, best estimator for EPE (Theoretically) is :


Conditional
mean $f(x) = E(Y | X = x)$

- **Nearest neighbors** assumes that $f(x)$ is well approximated by a locally constant function.

Review : WHEN EPE USES DIFFERENT LOSS

Loss Function	Estimator $\hat{f}(x)$
L_2 	$\hat{f}(x) = E[Y X = x]$
L_1 	$\hat{f}(x) = \text{median}(Y X = x)$
$0-1$ 	$\hat{f}(x) = \arg \max_Y P(Y X = x)$ <p>(Bayes classifier / MAP)</p>

Today :

- ✓ K-nearest neighbor
- ✓ Model Selection / Bias Variance Tradeoff
 - ✓ EPE
-  ✓ Decomposition of MSE
 - ✓ Bias-Variance tradeoff
 - ✓ High bias ? High variance ? How to respond ?

Decomposition of EPE

– When additive error model:

– Notations $Y = f(X) + \epsilon, \epsilon \sim (0, \sigma^2)$

- Output random variable: Y
- Prediction function: f
- Prediction estimator: \hat{f}

$$\begin{aligned}
 EPE(\overset{f(x_0)}{\cancel{x_0}}) &= E[(Y - \hat{f})^2 | X = x_0] \\
 &= E[((Y - f) + (f - \hat{f}))^2 | X = x_0] \\
 &= E[\underbrace{(Y - f)^2}_{\epsilon} | X = x_0] + \underbrace{E[(f - \hat{f})^2 | X = x_0]}_{MSE} \\
 &= \sigma^2 + \text{Var}(\hat{f}) + \text{Bias}^2(\hat{f})
 \end{aligned}$$

MSE component of \hat{f} in estimating f

Bias-Variance Trade-off for EPE:

$$\text{EPE} (\overset{f(x_0)}{\cancel{x=0}}) = \text{noise}^2 + \text{bias}^2 + \text{variance}$$

Unavoidable error

Error due to incorrect assumptions

Error due to variance of training samples

The diagram illustrates the Bias-Variance Trade-off for Expected Prediction Error (EPE). The equation is $\text{EPE} (\overset{f(x_0)}{\cancel{x=0}}) = \text{noise}^2 + \text{bias}^2 + \text{variance}$. Three blue arrows point from the terms below to the corresponding terms in the equation: 'Unavoidable error' points to 'noise²', 'Error due to incorrect assumptions' points to 'bias²', and 'Error due to variance of training samples' points to 'variance'. The entire diagram is enclosed in a blue L-shaped frame on the left and right sides.

BIAS AND VARIANCE TRADE-OFF for MSE

(more general setting !!!):

more general setting of MSE

- θ : true value (normally unknown)
- $\hat{\theta}$: estimator
- $\bar{\theta} := E[\hat{\theta}]$ (mean, i.e. expectation of the estimator)

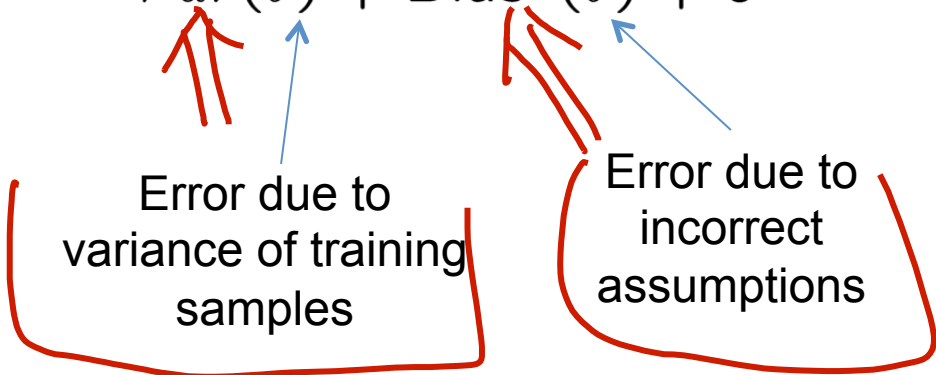
- Bias $E[(\bar{\theta} - \theta)^2]$
 - measures **accuracy** or **quality** of the estimator
 - low bias implies on average we will accurately estimate true **parameter** or **func** from training data
 - Variance $E[(\hat{\theta} - \bar{\theta})^2]$
 - Measures **precision** or **specificity** of the estimator
 - Low variance implies the estimator does not **change** much as **the training set varies**
- Handwritten note:* θ could be $\begin{cases} w & \text{for LR} \\ f & \text{for EPE} \end{cases}$

BIAS AND VARIANCE TRADE-OFF for MSE of parameter estimation:

e.g. In EPE case, $E[(f - \hat{f})^2 | X = X_0] \Rightarrow \text{MSE}(f(x_0))$

$$\begin{aligned}
 \text{MSE}(\hat{\theta}) &= E[(\hat{\theta} - \theta)^2] \\
 &= E[(\hat{\theta} - \bar{\theta}) + (\bar{\theta} - \theta)]^2 \\
 &= E[(\hat{\theta} - \bar{\theta})^2] + E[(\bar{\theta} - \theta)^2] + 2E[(\hat{\theta} - \bar{\theta})(\bar{\theta} - \theta)] \\
 &= \text{Var}(\hat{\theta}) + \text{Bias}^2(\hat{\theta}) + 0
 \end{aligned}$$

$\bar{\theta} = \text{mean}(\hat{\theta}) \Rightarrow$

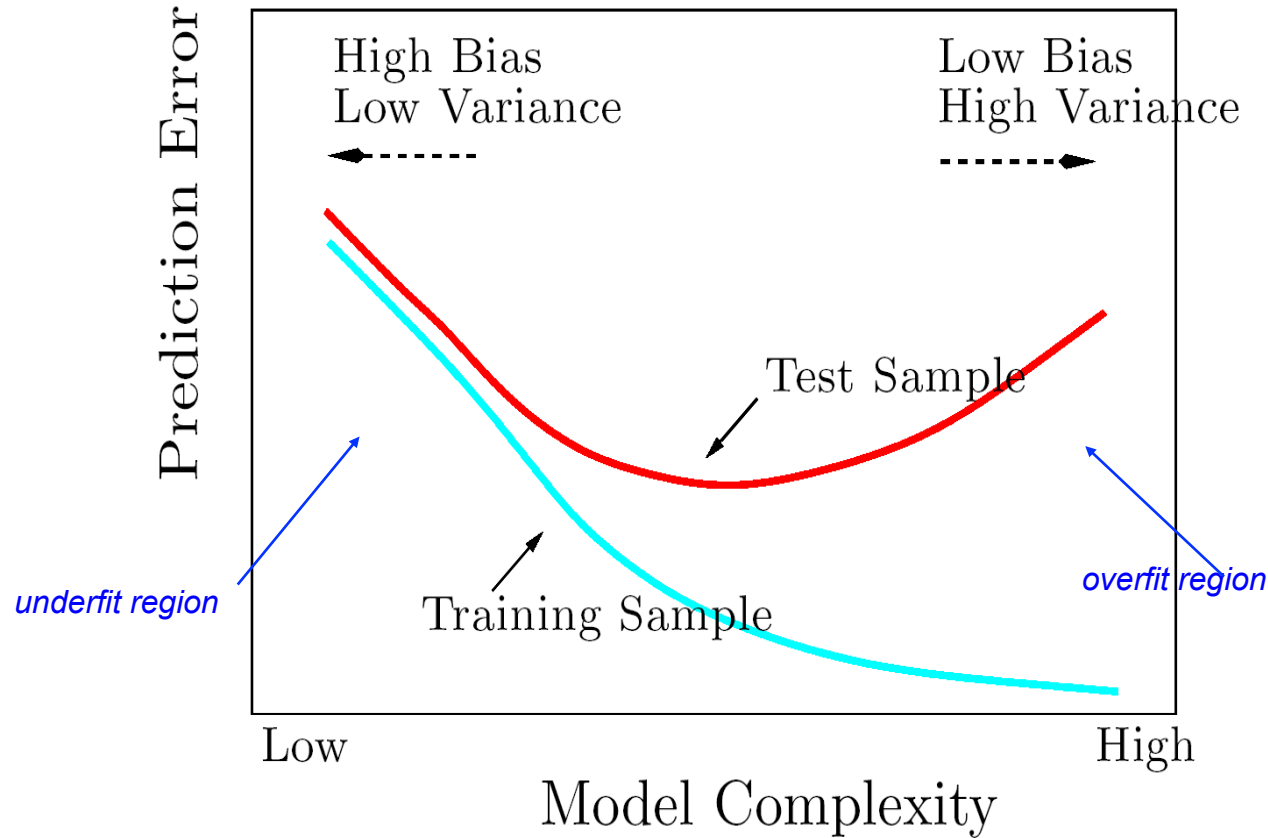


$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = \text{Bias}^2(\hat{\theta}) + \text{Var}(\hat{\theta})$$

Today :

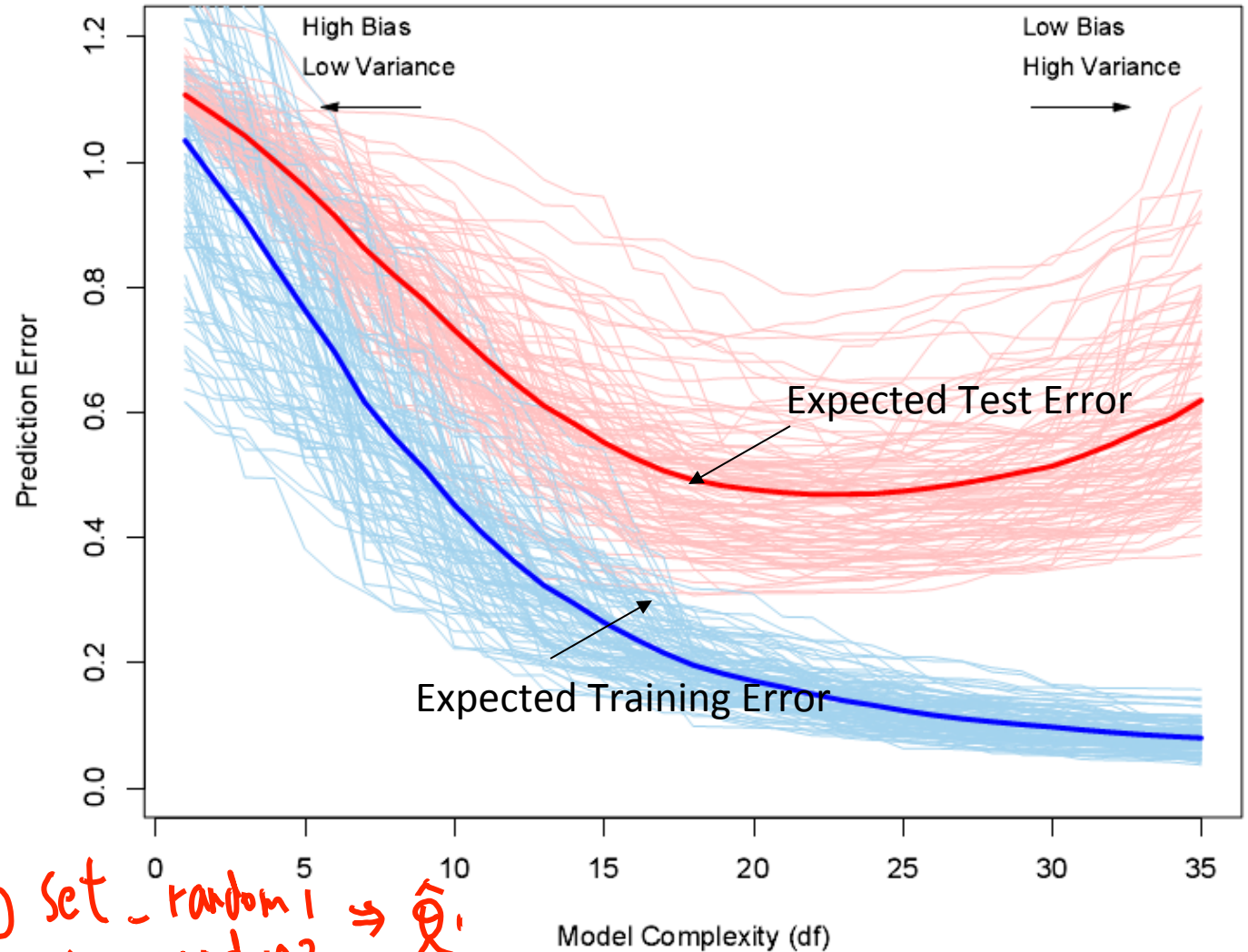
- ✓ K-nearest neighbor
- ✓ Model Selection / Bias Variance Tradeoff
 - ✓ EPE
 - ✓ Decomposition of MSE
 - ➔ ✓ Bias-Variance tradeoff
 - ✓ High bias ? High variance ? How to respond ?

Bias-Variance Tradeoff / Model Selection



(1) Training vs Test Error

- Training error can always be reduced when increasing model complexity,
- Expected Test error and CV error → good approximation of EPE



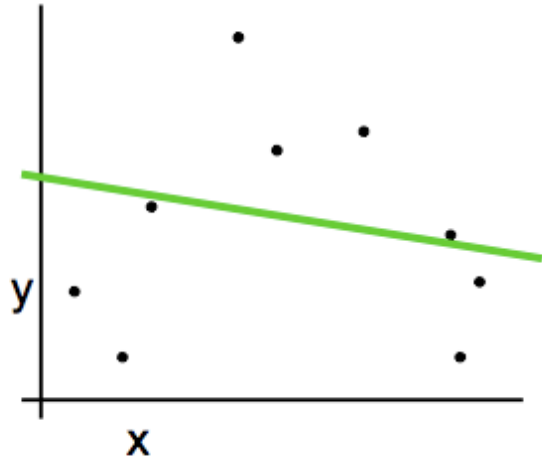
training set - random 1 $\Rightarrow \hat{\theta}_1$
 training set - random 2 $\Rightarrow \hat{\theta}_2$
 3 $\Rightarrow \hat{\theta}_3$

(2) Bias-Variance Trade-off

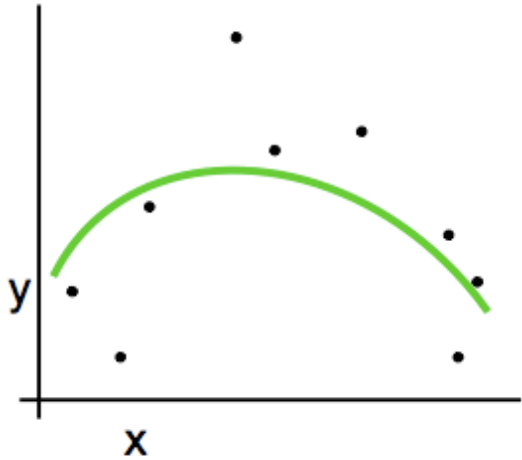
- Models with too few parameters are inaccurate because of a large bias (not enough flexibility).
 - poly regression: d small*
 - KNN: K large*
- Models with too many parameters are inaccurate because of a large variance (too much sensitivity to the sample randomness).

- poly regression: d large*
- KNN: K small*

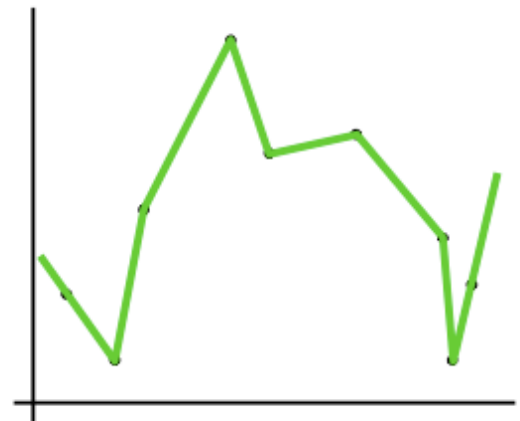
(2.1) Regression: Complexity versus Goodness of Fit



Highest Bias
Lowest variance
Model complexity = low

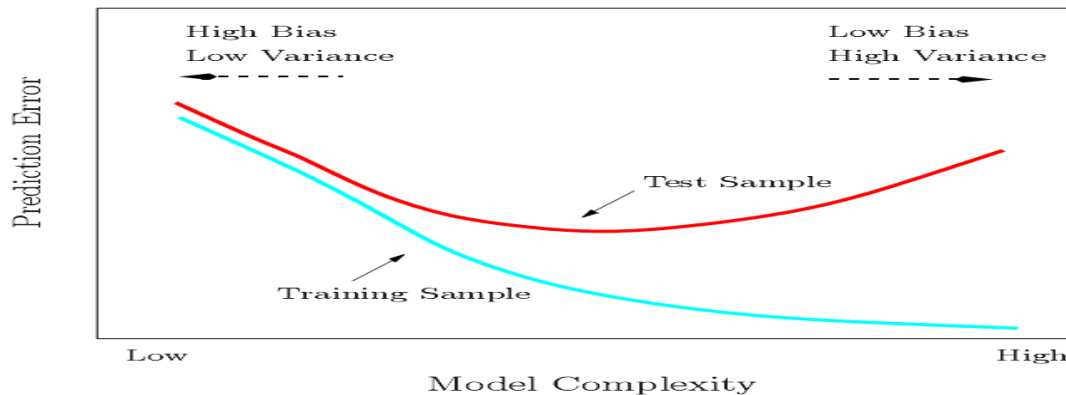


Medium Bias
Medium Variance
Model complexity = medium



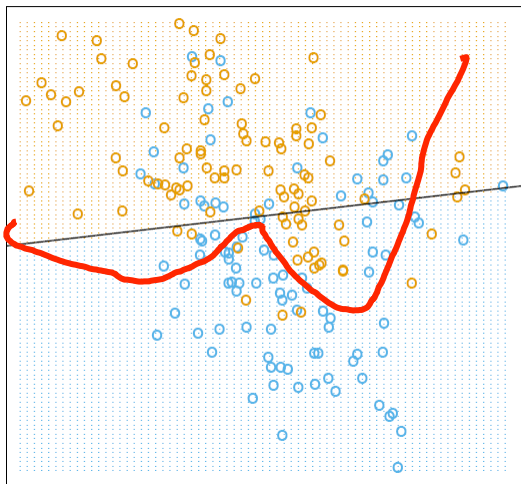
Smallest Bias
Highest variance
Model complexity = high

Low Variance / High Bias



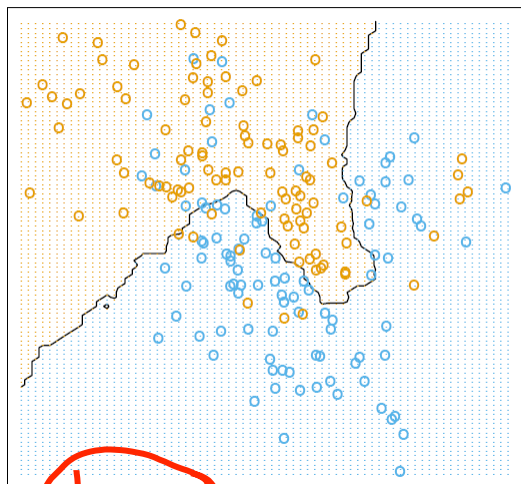
Low Bias / High Variance

(2.2) Classification, Decision boundaries in global vs. local models



Low Variance / High Bias

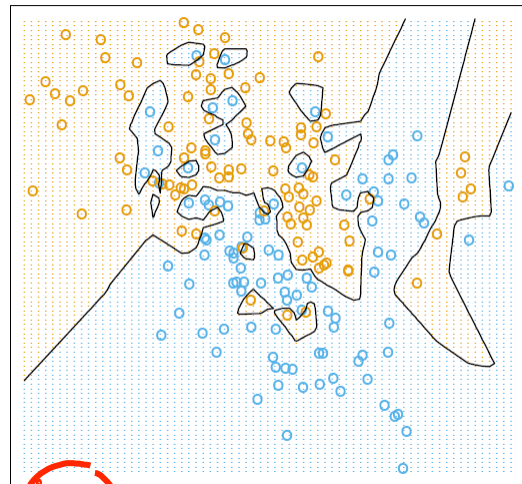
- linear regression
- global
- stable
- can be inaccurate



$k=15$

15-nearest neighbor

Low Variance / High Bias



$k=1$

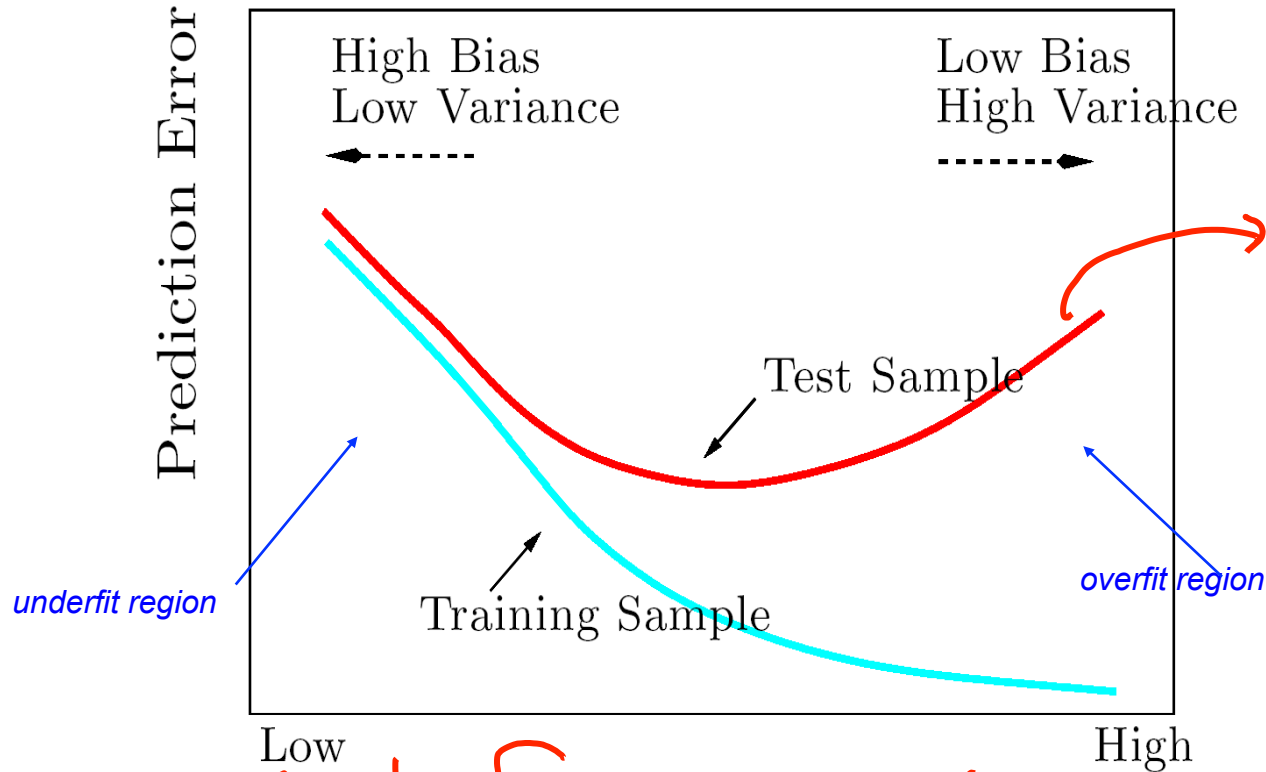
1-nearest neighbor

- KNN
- local
 - accurate
 - unstable

Low Bias / High Variance

What ultimately matters: **GENERALIZATION**

Bias-Variance Tradeoff / Model Selection



CV / test error
[good approximation
of EPE]

KNN: large k ← [Model Complexity] → small k
Regression: small d → large d

Model “bias” & Model “variance”

- Middle RED:
 - TRUE function θ [middle red]
- Error due to bias:
 - How far off in general from the middle red

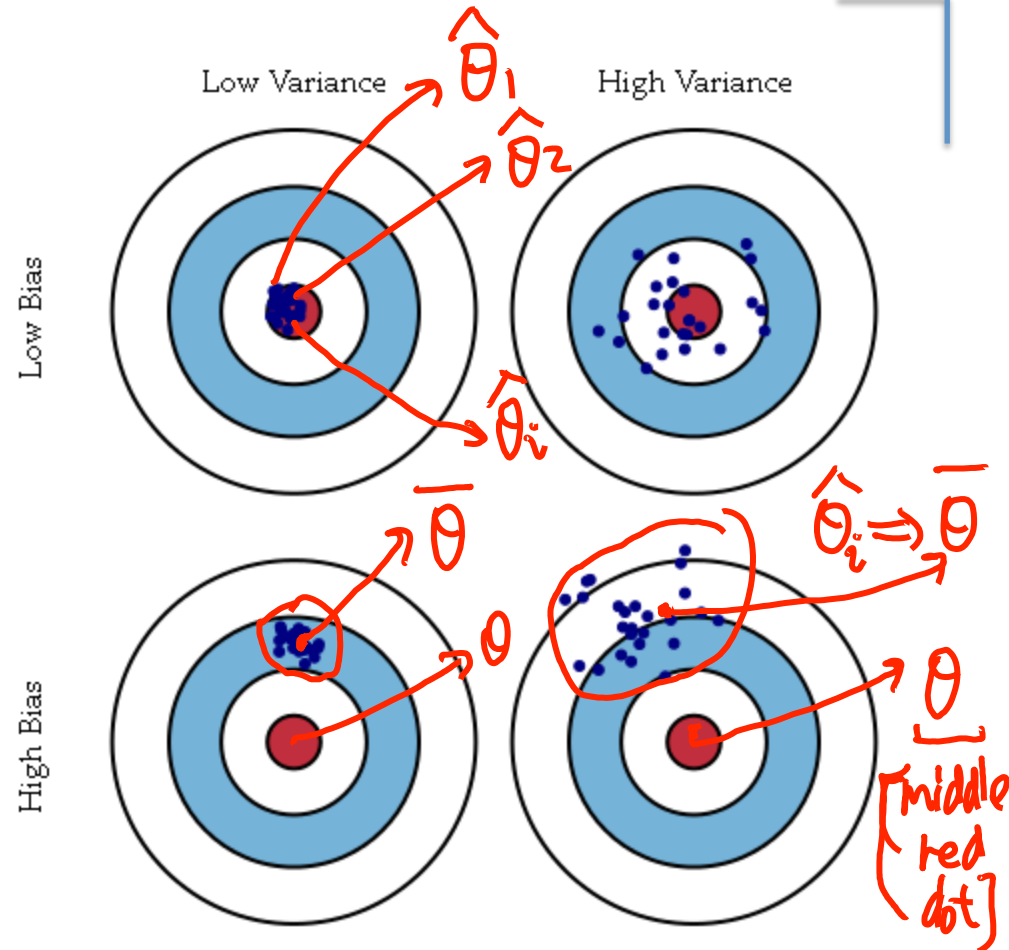
$$E(\theta - \bar{\theta})$$

mean of $\hat{\theta}$

- Error due to variance:
 - How wildly the blue points spread

$$E((\hat{\theta} - \bar{\theta})^2)$$


$\{\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3, \dots\}$ Blue dots



need to make assumptions that are able to generalize

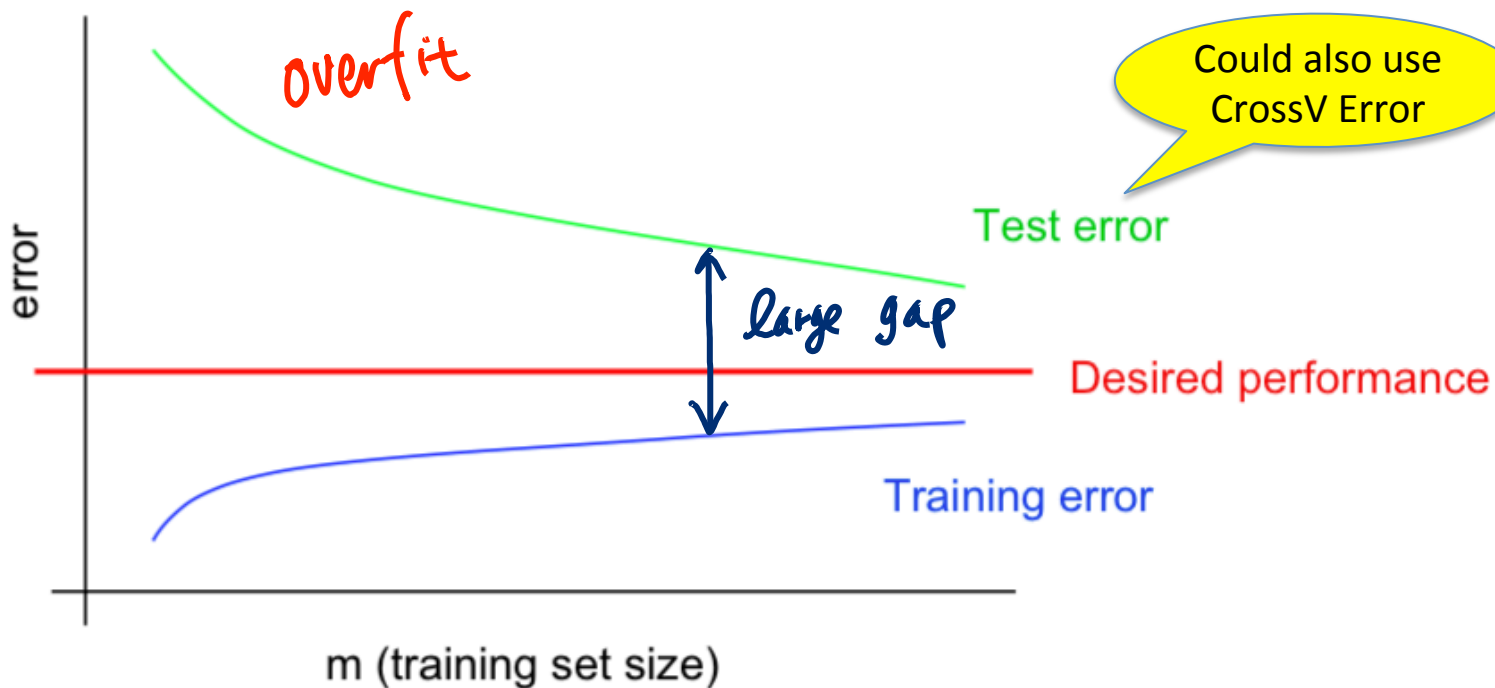
- Components of generalization error
 - **Bias:** how much the average model over all training sets differ from the true model?
 - Error due to inaccurate assumptions/simplifications made by the model
 - **Variance:** how much models estimated from different training sets differ from each other
- **Underfitting:** model is too “simple” to represent all the relevant class characteristics
 - High bias and low variance
 - High training error and high test error
- **Overfitting:** model is too “complex” and fits irrelevant characteristics (noise) in the data
 - Low bias and high variance
 - Low training error and high test error

Today :

- ✓ K-nearest neighbor
- ✓ Model Selection / Bias Variance Tradeoff
 - ✓ EPE
 - ✓ Decomposition of MSE
 - ✓ Bias-Variance tradeoff
-  ✓ High bias ? High variance ? How to respond ?

(1) High variance

Typical learning curve for high variance:



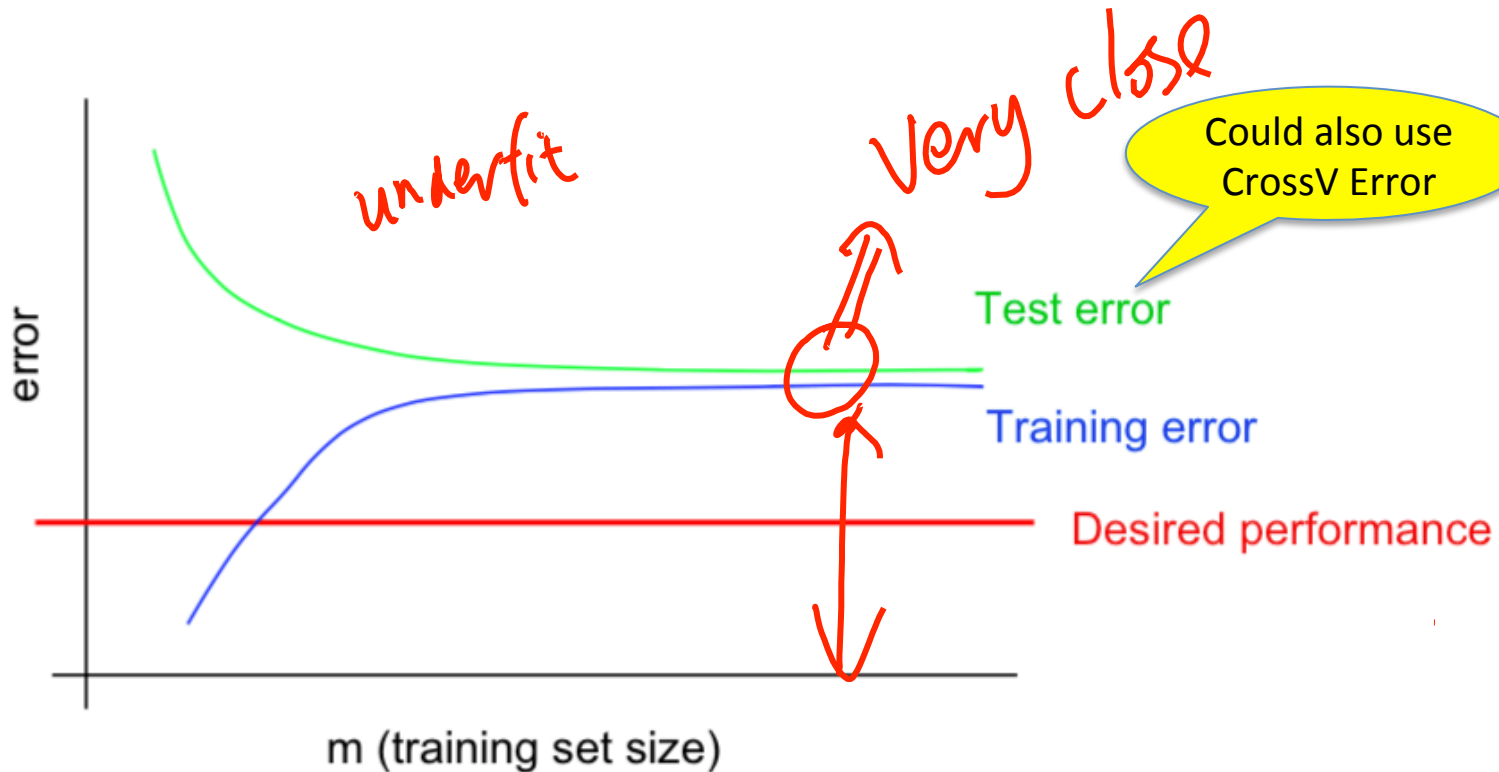
- Test error still decreasing as m increases. Suggests larger training set will help.
- Large gap between training and test error.
- **Low training error and high test error**

How to reduce variance?

- Choose a simpler classifier
- Regularize the parameters
- Get more training data
- Try smaller set of features

(2) High bias

Typical learning curve for high bias:



- Even training error is unacceptably high.
 - Small gap between training and test error.
- High training error and high test error**

How to reduce Bias ?

- E.g.
 - Get additional features
 - Try adding basis expansions, e.g. polynomial
 - Try more complex learner

(3) For instance, if trying to solve “spam detection” using (Extra)

L2 - logistic regression, implemented with gradient descent.

Fixes to try:

If performance is not as desired

- Try getting more training examples.
- Try a smaller set of features.
- Try a larger set of features.
- Try email header features.
- Run gradient descent for more iterations.
- Try Newton’s method.
- Use a different value for λ .
- Try using an SVM.

Why??

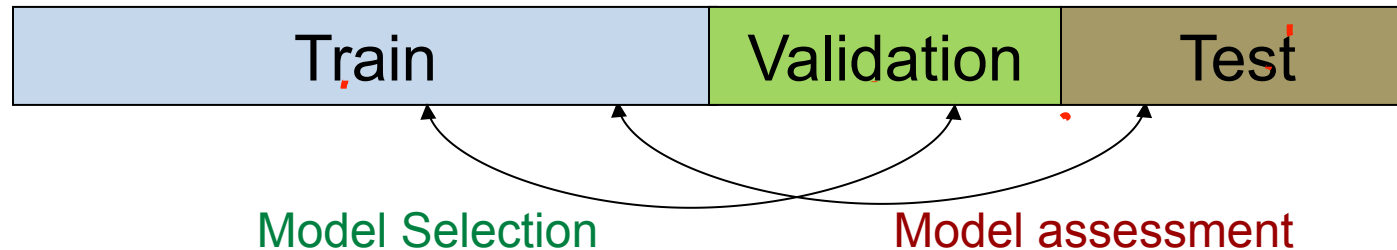
- Fixes high variance.
- Fixes high variance.
- Fixes high bias.
- Fixes high bias.
- Fixes optimization algorithm.
- Fixes optimization algorithm.
- Fixes optimization objective.
- Fixes optimization objective.

(4) Model Selection and Assessment

- Model Selection
 - Estimating performances of different models to choose the best one
- Model Assessment
 - Having chosen a model, estimating the prediction error on new data

Model Selection and Assessment (Extra)

- Data Rich Scenario: Split the dataset



- Insufficient data to split into 3 parts
 - Approximate validation step analytically
 - AIC, BIC, MDL, SRM
 - Efficient reuse of samples
 - Cross validation, bootstrap

Today Recap:

- ✓ K-nearest neighbor
- ✓ Model Selection / Bias Variance Tradeoff
 - ✓ EPE
 - ✓ Decomposition of MSE
 - ✓ Bias-Variance tradeoff
 - ✓ High bias ? High variance ? How to respond ?

References

- ❑ Prof. Tan, Steinbach, Kumar's "Introduction to Data Mining" slide
- ❑ Prof. Andrew Moore's slides
- ❑ Prof. Eric Xing's slides
- ❑ Hastie, Trevor, et al. *The elements of statistical learning*. Vol. 2. No. 1. New York: Springer, 2009.