

## Statistical Models & Machine Learning

- So far mostly visualization and point statistics.
- Modeling is another important part of EDA.
- Other classes cover theory of models, so here the focus is implementation.
- We saw performance was an issue for HWS with text; it can also be an issue with fitting models.
- The basic idea of modeling is that we have some inputs  $X_1, X_2, \dots$  and we want to predict an output  $Y$ .

or understand

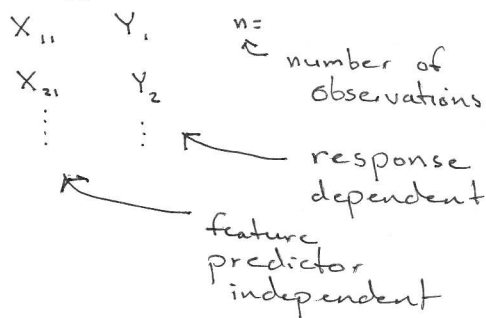
Tradeoff!

$$Y = f(X_1, X_2, \dots) + \varepsilon$$

↖ error term

$$\hat{Y} = \hat{f}(X_1, X_2, \dots)$$

### Training Set



### Test Set

### Parametric vs non-parametric

- Parametric models assume a <sup>shape</sup> form for  $f$ , then estimate parameters

Non-parametric models do not assume a shape for  $f$ , but do assume other properties, like smoothness. More obs needed

### Supervised vs unsupervised

(response) (no response)

### Regression vs classification

- We use mean squared error to measure error in regression models:

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

Other measurements are possible and can change how the model fits.

- We care about test error, not training error! This is where overfitting can be a problem.

} Revisit in CV example

- Bias-variance tradeoff:

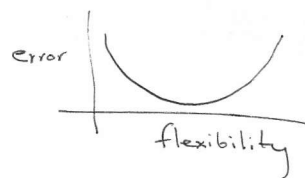
$$\mathbb{E} (y_0 - \hat{f}(x_0))^2 = \underbrace{\text{Var}(\hat{f}(x_0))}_{\text{model variance}} + \underbrace{[\text{Bias}(\hat{f}(x_0))]^2}_{\text{model bias}} + \underbrace{\text{Var}(\epsilon)}_{\text{error variance}}$$

Variance — how  $\hat{f}$  changes with different train set

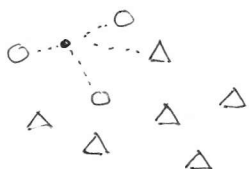
Bias — error due to approximation

bias inference

variance flexibility



- k-nearest neighbors



Use k-nearest training points to predict class

$$P(Y_0 = \Delta | X_{1,0} = x_{1,0}, X_{2,0} = x_{2,0}) = \frac{\# \Delta}{k}$$

Ties!

Bayes classifier approximation

The choice of k makes a difference!

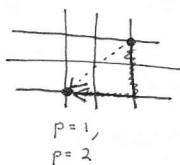
Similarly, how we measure distance makes a difference!

- Minkowski distances:

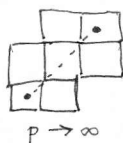
- Euclidean  $p=2$

- Manhattan  $p=1$

Manhattan, Euclidean



Chessboard



$$\left( \sum_{i=1}^n |a_i - b_i|^p \right)^{1/p}$$

$p \geq 1$

- How to actually implement kNN?
- What are the inputs? The outputs?

Input:  $x_1, x_2, \dots$  predictors  
 $k$  number of neighbors  
 $\{(x, y)_i\}$  training set

Output:  $\hat{y}_0$  predicted label

- So what are the steps?

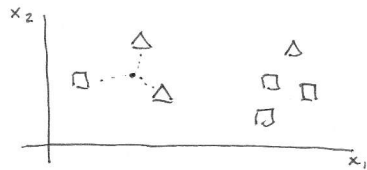
- Find  $k$  nearest neighbors to  $x_0$ 
  - How do we tell which points are "near"?
  - Compute distances, find smallest distances  
 $\text{dist}()$   $\text{sort}()$   $\text{order}()$
- Use neighbor labels to "vote"
  - Break ties  
 $\text{sample}()$   $\text{table}()$
- Return "winning" label

- Is kNN "fast"?

- Is kNN flexible? Where is bias-variance tradeoff?

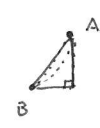
References: ISLR 2, 5.1 / TARP 8, 14

5-min reminder of kNN



Minkowski distance

$$\left( \sum_{i=1}^d |a_i - b_i|^p \right)^{1/p} \quad p \geq 1$$



Effect of k on model flexibility

How to implement kNN?

`predict_knn()`

Inputs: new point, train points, k, metric

Outputs: class label

Steps:

- |      |   |  |
|------|---|--|
| Find | { | 1. Compute distances from new point(s) to all train points |
|      |   | 2. Sort/order distances                                    |
| Vote | { | 3. Take top k train points                                 |
|      |   | 4. Count classes for top k points                          |
|      |   | 5. Return winning class, possibly breaking ties.           |

`dist()`  
`as.matrix()`

`order()`

`[]`

`table()`

`which.max() / max()`

Break the voting step (including ordering) into a separate function.

`apply(new-points, 1, vote) ...`

Error rates: how can we tell how well the model works?

$$\sum \frac{\text{actual} == \text{pred}}{\text{total}}$$

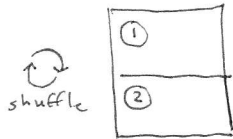
Fit on train → predict on train → underestimate

Fit on train → predict on test → Why not use test to train?

## □ Cross-validation (m-fold)

\*

Idea: fit multiple times and average

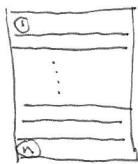


Iteration	1	2
Train	①	②
Test	②	①
Error	$e_1$	$e_2$

$$\hat{e} = \frac{e_1 + \dots + e_m}{m}$$

Two folds doesn't give us much to average.

Error rate will vary depending on how we choose folds, and model doesn't get to see all the data } High bias



n error estimates

Leave-one-out CV

Train ~~①~~ all others

Test ① ② ... ①

Now we have to fit the model n times — not feasible for large n or slow model.

Estimates are highly dependent! → high variance

Usually we compromise and use  $m=5$  or  $m=10$

## □ Error rates

Bias - variance tradeoff

## □ How to implement CV?

`cv_error_knn()`

Inputs: m, train set, k, metric

Outputs: CV error estimate

Steps:

1. Shuffle the rows or sample row ids.
2. Split rows into folds.
3. Loop over test folds; remaining are train folds
  - a. Fit model on train folds
  - b. Predict for test fold
  - c. Compute error rate
4. Return average of error rates.

`sample()`

`rep()`  
`split()`

`mean()`

- Cross-validation (Lecture 19)
- Implementing CV (Lecture 19)
- Performance considerations for CV
  - Compute distances within training set once, before looping over the folds
  - Split row indexes instead of actual observations, and look up appropriate rows for each fold (distances)
- Recap designing functions / ideas to code:
  - List inputs, outputs, and steps
    - └ This requires careful thought
  - Draw pictures or make simple examples of the problem. Work these out by hand
  - Divide-and-conquer: complicated steps can be written as separate functions
  - Test on simple/small examples first
- Function examples?
  - dist()    order()    table()    split()
  - sample()    rep()    ~~+~~    indexing
- Where to go from here?
  - STA 141B focuses on collecting data, interactive visualizations, and uses Python
  - STA 141C focuses on programming for statistics and uses R (possibly with some C)
  - STA 135/108
  - ECS 171 machine learning