

CS 412

JAN 20TH – NEAREST NEIGHBORS

$$\delta = 0.05$$

Error Bounds

PAC

To be within ϵ with 95% probability requires N flips:

| | | |
|------------------|---------------|-----------|
| $\epsilon = .2$ | \rightarrow | $N = 10$ |
| $\epsilon = .1$ | \rightarrow | $N = 38$ |
| $\epsilon = .05$ | \rightarrow | $N = 149$ |
| $\epsilon = .03$ | \rightarrow | $N = 414$ |
| $\epsilon = .01$ | \rightarrow | $N = 931$ |

$$P(|\hat{\theta} - \theta^*| < \epsilon) \geq \delta$$

Summary

Probability important for:

- Estimating prediction uncertainty
- Analyzing random samples

Many important concepts:

- Random variables
- Marginalization
- Approximation bounds

Data Separation

Before your model construction begins,
you need to separate your data into at least two sets

- Training
- Testing
- *Validation (which we'll discuss later)*

You should not conduct any analysis of the data before separating it

- This is called *data snooping* and can result in an inaccurate estimate of our final reported error

The only way we can estimate the final reported error of our model
is by judging on the test set

Data Separation

The training set is used to train the data and is usually a smaller section of the data as a whole (much less than half)

By not analyzing or training on a large set of the data, what is the advantage?

Data Separation

The training set is used to train the data and is usually a smaller section of the data as a whole (much less than half)

By not analyzing or training on a large set of the data, what is the advantage?

When we've created our "choice" model, we can run it on the test data *only once* and get an estimate of the error without having to induce a penalty

Approach

A bank creates an investment ML algorithm that forecasts whether a certain bond value will go up or down. As input, it receives a group of financial features that have all been standardized for mean and variance. The bank separated the data into two groups, training and testing. After building the model on the training data, the model trained on the testing data can correctly determine whether the bond will go up or down with 52% accuracy.

When this model was applied to the actual market, it consistently was wrong less than 50% of the time. What went wrong?

Approach

A bank creates an investment ML algorithm that forecasts whether a certain bond value will go up or down. As input, it receives a group of financial features that have all been standardized for mean and variance. The bank separated the data into two groups, training and testing. After building the model on the training data, the model trained on the testing data can correctly determine whether the bond will go up or down with 52% accuracy.

When this model was applied to the actual market, it consistently was wrong less than 50% of the time. What went wrong?

The bank normalized both the training and testing sets, meaning the training set was actually impacted by the values in the test set.

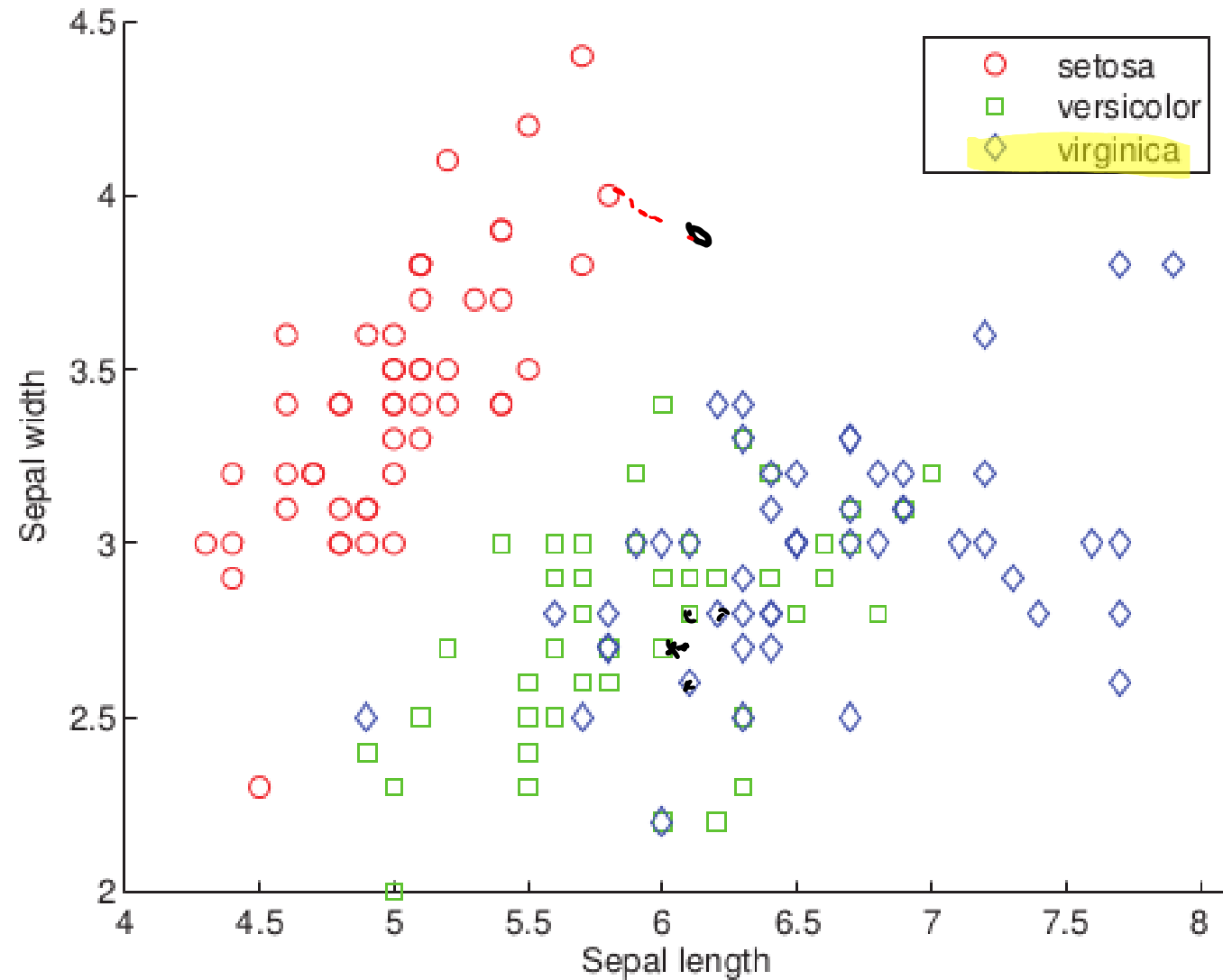
yes

no

Three interlocking gray gears are shown in a 3D perspective. They are arranged in a triangular pattern, with one gear at the top and two below it. The gears have a metallic texture and are set against a white background.

The image shows two boxes containing various geometric shapes and colors. The left box contains: a blue rectangle, a red circle, a grey oval, a green oval, a blue circle with a white center, a blue crescent moon, a yellow circle, a blue arrow, and a grey four-pointed star. The right box contains: a green diamond, a red five-pointed star, a yellow circle, a yellow circle with a white center, an orange parallelogram, a green triangle, a green parallelogram, a yellow five-pointed star, and a blue diamond.

left alone
until finished



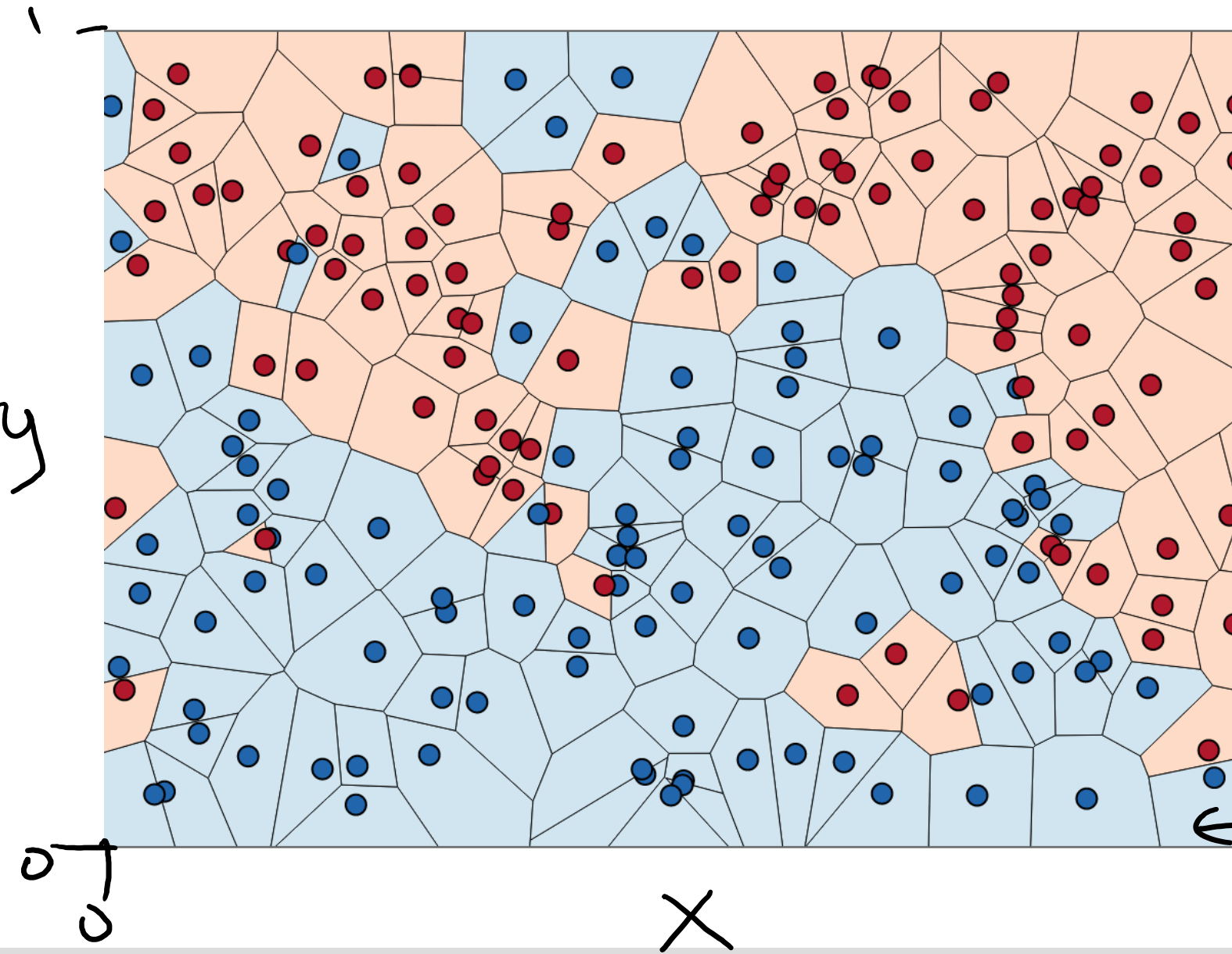
Approach

Given the following data, what are some approaches that we could use to make a simple model?

$if k = n$

Nearest Neighbor

The nearest neighbor model takes the closest point and assigns that value to any region in the space



← Decision Region

Nearest Neighbor

What are some observations we could make about this modeling technique?

Number of neighbors - k (usually odd)
for training we have 100% accuracy - overfitting

Nearest Neighbor

What are some observations we could make about this modeling technique?

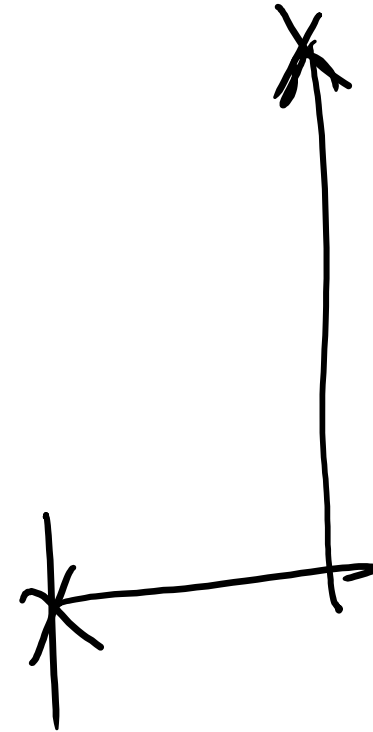
- Needs to store and run through all of the data
- Not clearly defined by what we mean by distance
- Does not deal well with regions of high overlap
- Categorical features?

Nearest Neighbor

What are some observations we could make about this modeling technique?

- Needs to store and run through all of the data
- Not clearly defined by what we mean by distance
 - Euclidian – straight line distance
 - Manhattan – sum of feature differences (going by street blocks)
 - Chebyshev – closest by maximum dimension
- Does not deal well with regions of high overlap
- Categorical features?
 - Need to create (n-1) dummy features

$$\text{Euclid} = \sqrt[n]{\sum_i^2}$$
$$Man = \sum |i|$$



k-Nearest Neighbor

For binary classification problems

- Take the majority classification of the k-nearest neighbors (where k is odd)

For numerical output problems

- Take the average of the k-nearest neighbors OR
- Take the weighted average of the k-nearest neighbors

n

k-Nearest Neighbor

For binary classification problems

- Take the majority classification of the k-nearest neighbors (where k is odd)

For numerical output problems

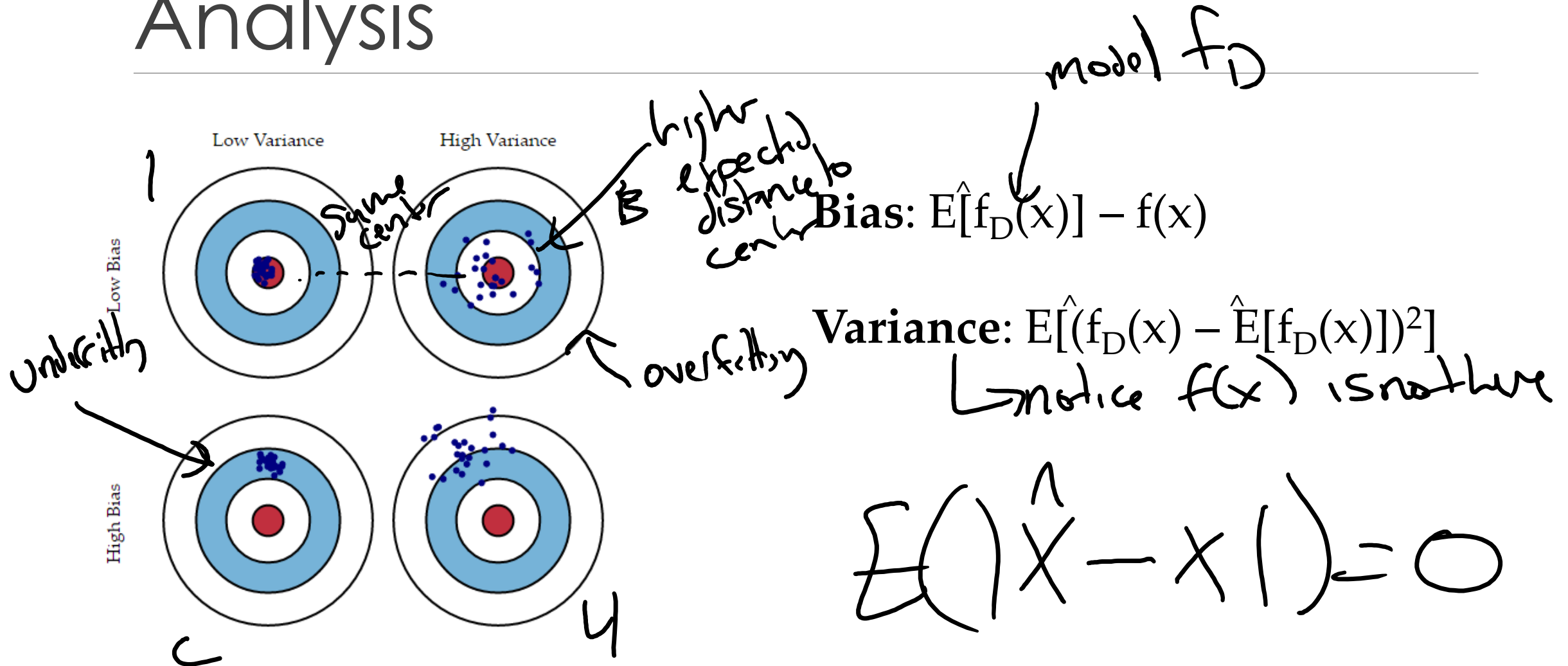
- Take the average of the k-nearest neighbors OR
- Take the weighted average of the k-nearest neighbors (Gaussian distribution)

Normal

$K=n$ $acc = \text{most common point}$

$K=1$ $acc = 100\%$ on the training set

Nearest Neighbor Regression Analysis



How do we evaluate?

train validation

After test data is split off

K-Fold Cross Validation:

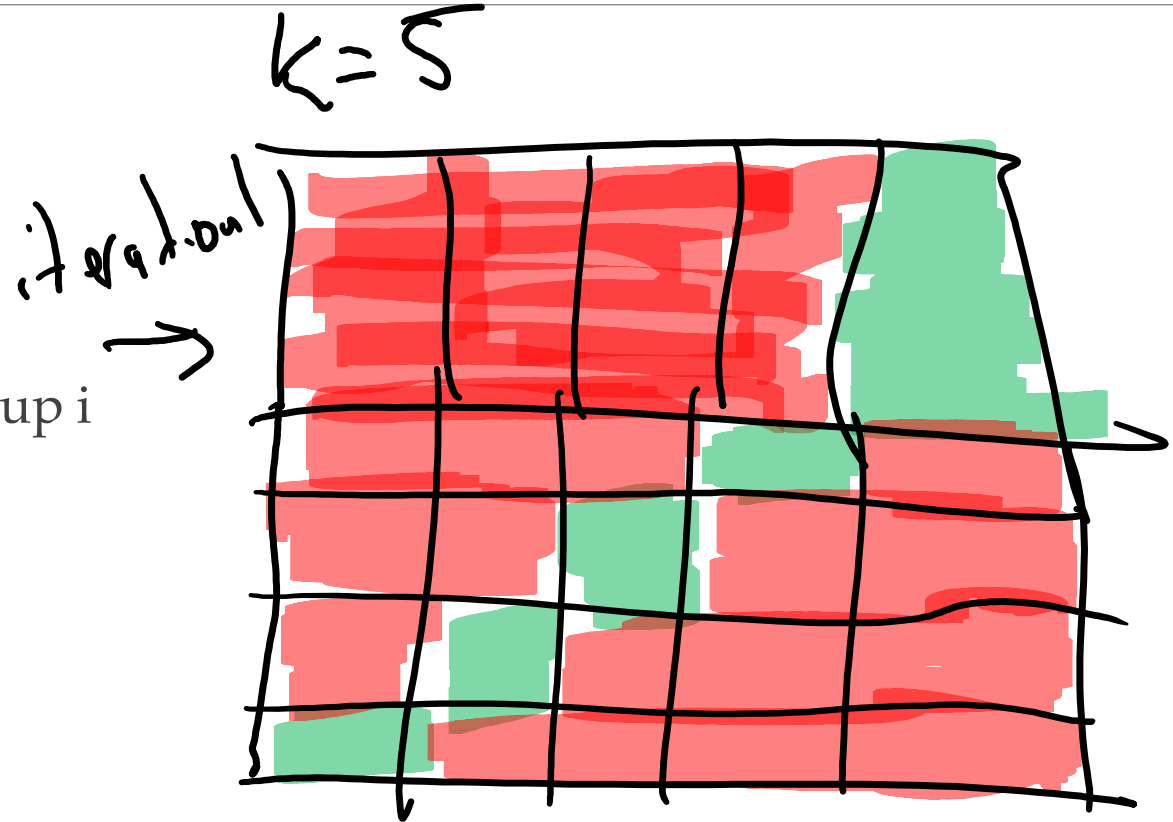
Split data into K groups.

For $i=1:K$

Train classifier on all except group i

Evaluate on group i

Return average evaluation



How do we evaluate?

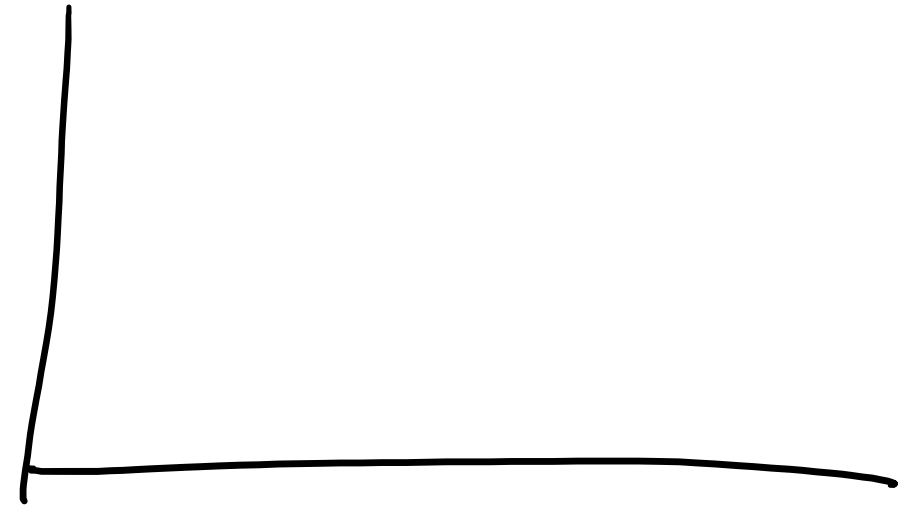
Leave One Out Cross Validation (LOOCV): = n -fold cross validation

For all data points (k) in the validation set:

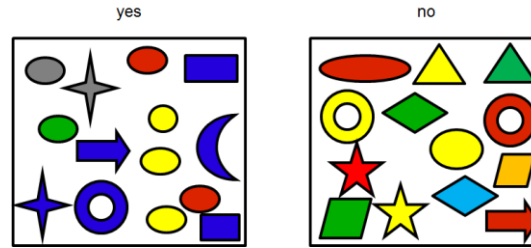
- Train on all other data points

- Test on the single point k

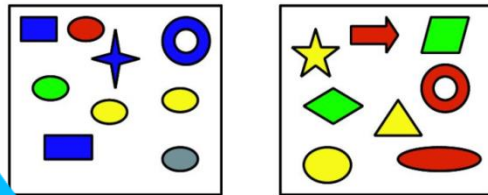
Return the average evaluation



Full Data



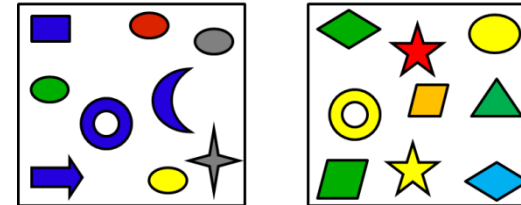
Training Data



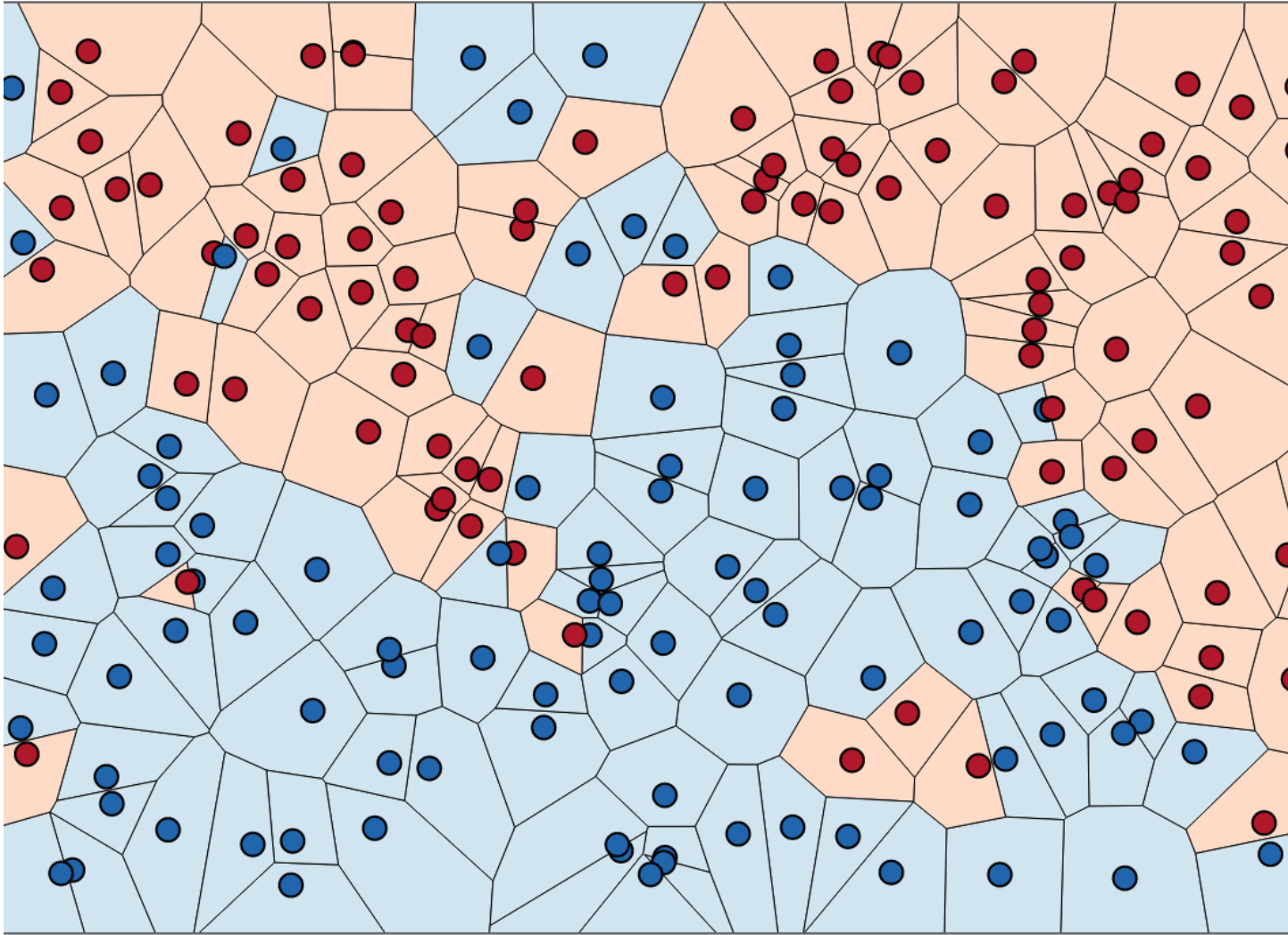
Classifier



Testing Data



separate into validation



Nearest Neighbor

The nearest neighbor model takes the closest point and assigns that value to any region in the space

Nearest Neighbor

What are some observations we could make about this modeling technique?

- Needs to store and run through all of the data
- Not clearly defined by what we mean by distance
 - Euclidian – straight line distance
 - Manhattan – sum of feature differences (going by street blocks)
 - Chebyshev – closest by maximum dimension
- Does not deal well with regions of high overlap
- Categorical features?
 - Need to create (n-1) dummy features

Eye color (B, ~~Green~~, Gray, Brown)

→ Blue (1, 0)
→ Gray (1, 0)
→ Brown (1, 0)

$\therefore (0, 0) \rightarrow \text{Green}$

k-Nearest Neighbor

For binary classification problems

- Take the majority classification of the k-nearest neighbors (where k is odd)

For numerical output problems

- Take the average of the k-nearest neighbors OR
- Take the weighted average of the k-nearest neighbors

k-Nearest Neighbor

For binary classification problems

- Take the majority classification of the k-nearest neighbors (where k is odd)

For numerical output problems

- Take the average of the k-nearest neighbors OR
- Take the weighted average of the k-nearest neighbors (Gaussian distribution)

k-Nearest Neighbor

For binary classification problems

- Take the majority classification of the k-nearest neighbors (where k is odd)

For numerical output problems

- Take the average of the k-nearest neighbors OR
- Take the weighted average of the k-nearest neighbors (Gaussian distribution)

We may also want to take representative samples

This is k-clustering, and we'll discuss this later

How do we evaluate a model?

EPE (*estimated prediction error*) for a given model f

$$\text{EPE}(f) = E(Y - f(X))^2 \quad \leftarrow \text{Squared error term}$$

Where X is the vector of input attributes,
 $f(X)$ is the predicted output and
 Y is the actual output

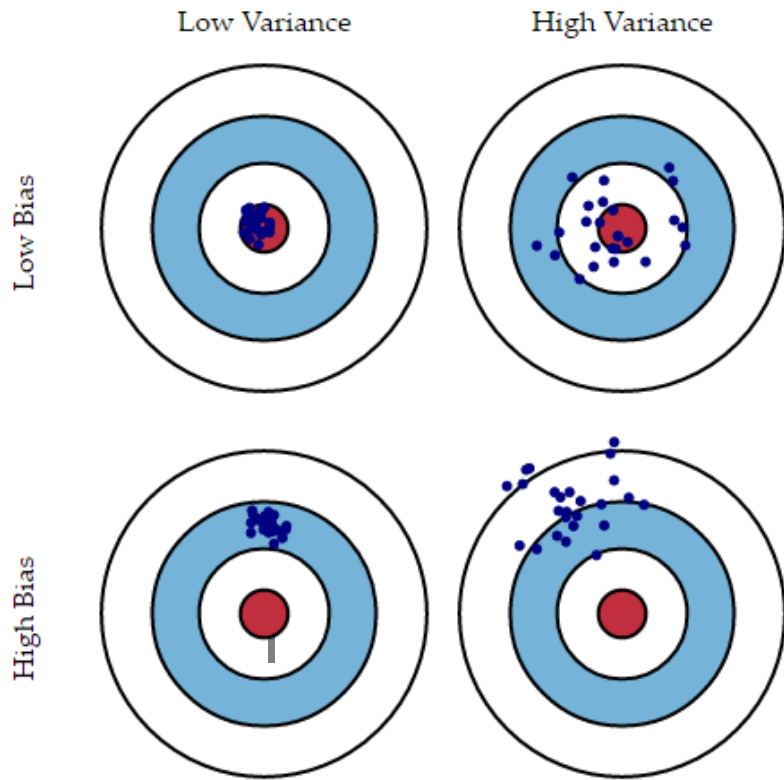
We will usually try to find the model which minimizes the
estimated error, but in the classification model, this is just accuracy

This may not show the whole picture!

$$Y = \{0, 1\}$$
$$X = \{0, 1\}$$

$$(Y - X)^2 = \{0, 1\}$$

Nearest Neighbor Regression Analysis



Bias: $E[\hat{f}_D(x)] - f(x)$

Variance: $E[(\hat{f}_D(x) - E[\hat{f}_D(x)])^2]$

PAC:

$$P(|\hat{\theta} - \theta^*| > \epsilon) < \delta$$

How do we evaluate?

K-Fold Cross Validation:

Split data into K groups.

For $i=1:K$

 Train classifier on all except group i

 Evaluate on group i

Return average evaluation

How do we evaluate?

Leave One Out Cross Validation (LOOCV):

For all data points (k) in the validation set:

- Train on all other data points

- Test on the single point k

Return the average evaluation

Homework 1

For HW 1, you'll be solving a binary classification problem with k-nearest neighbors for different levels of k

(1, 5)

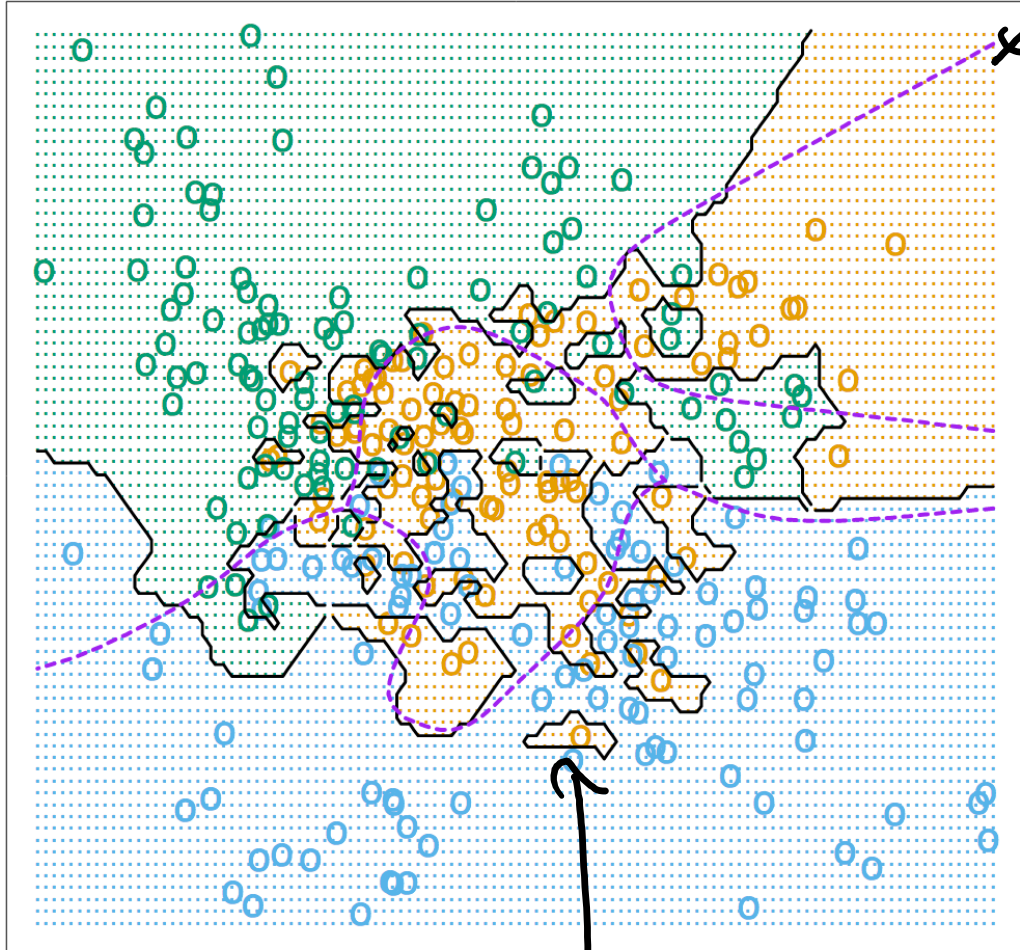
You will also be running your models with differing sizes of cross-validation and comparing runtimes

Make sure that you separate your test data before begin!

intersh
variance

mean
intersh

→ 1-Nearest Neighbor



Bayes Boundary

Results

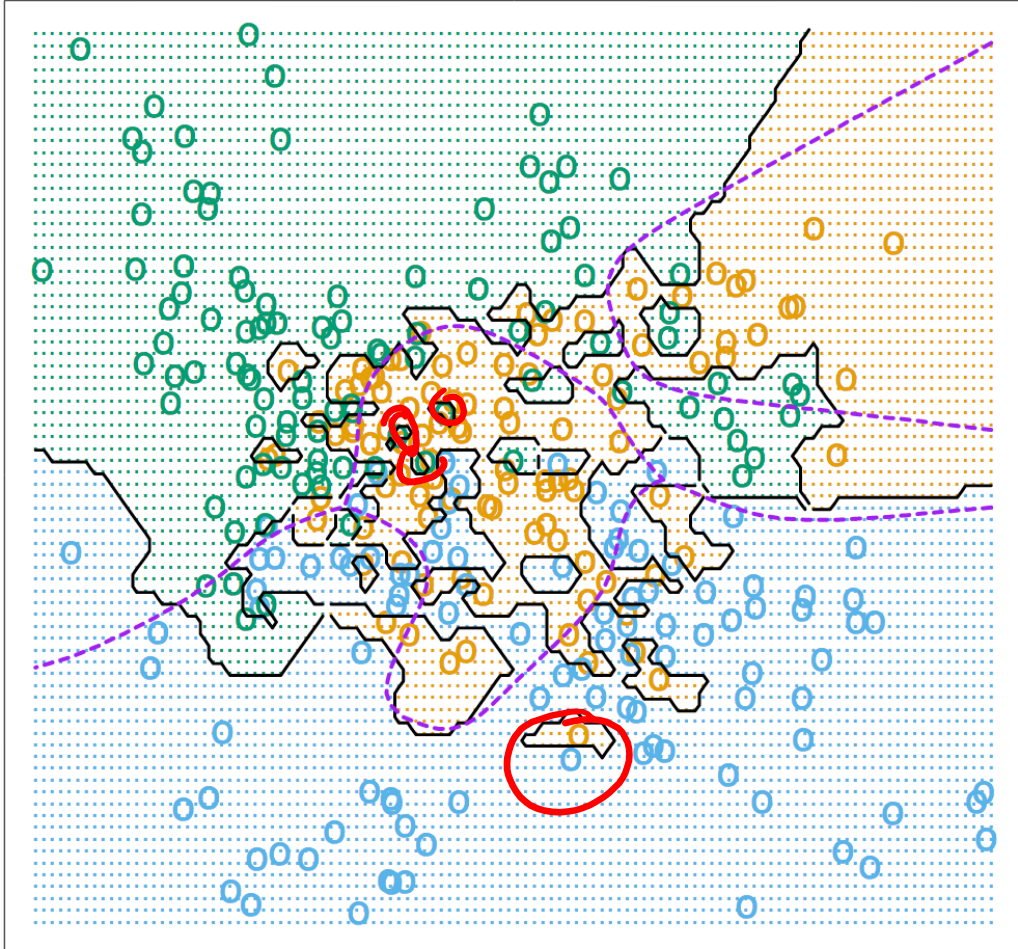
What are some observations we can make about this model?

$$k=1$$

3-classification

Over or under fit

1-Nearest Neighbor

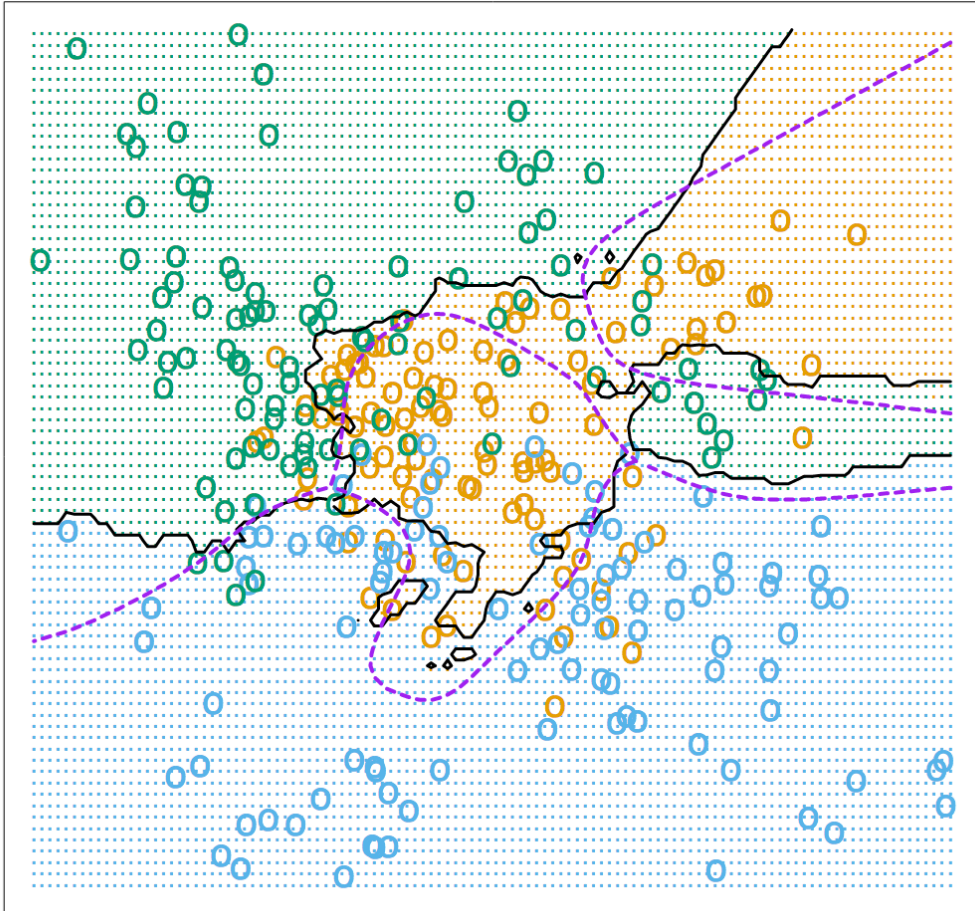


Results

What are some observations we can make about this model?

- 3 classifications
 - blue
 - orange
 - purple
- Purple line is called the Bayes decision boundary
- Does this suffer from overfitting?

15-Nearest Neighbors



Results

What are some observations we can make about this model?

- 3 classifications
 - blue
 - orange
 - purple
- Purple line is called the Bayes decision boundary
- Does this suffer from overfitting?

*bias has increased
variance*