## **Lecture 6: Intro to Classifiers**

**INFO 1998: Introduction to Machine Learning** 



#### Announcements

- No HW Due tonight!
- Mid-Semester Check-Ins today!
  - Last day next Monday
- HW 6 Released!



## Agenda

- 1. What is a Classifier?
- 2. K-Nearest Neighbors Classifier
- 3. Review of Underfitting v. Overfitting
- 4. Confusion Matrices



## What are Classifiers?



### What are Classifiers?

#### **Classifiers are able to help answer questions like...**

- "What species is this?"
- "What major is a student in based on their classes?"
- "Which Hogwarts House do I belong to?"
- "Am I going to pass this class?"



#### What are Classifiers?

- Classifiers predict the class/category of a set of data points. This class/category is based off of the target variable we are looking at.
- Difference between linear regression and classifiers
  - Linear regression is used to predict the value of a **continuous variable**
  - Classifiers are used to predict categorical or binary variables



# K-Nearest Neighbors Classifier



## What is the KNN Classifier?

- Lazy learner classifier
- Easy to interpret
- Fast to calculate
- Good for coarse analysis





#### **How Does It Work?**

Uses the k (a user specified value) nearest data points to predict the unknown one

- A simple assumption: the values nearest to a data point are similar to it
- k is a hyperparameter of the KNN model
  - a parameter which affects the training process





**Define** a *k* value (in this case k = 3)





Define a k value (in this case k = 3)
Pick a point to predict (blue diamond)



Define a k value (in this case k = 3)
Pick a point to predict (blue diamond)
Count the k closest points





- **Define** a k value (in this case k = 3)
- Pick a point to predict
- (blue diamond)
- **Count** the number of closest points
- **Increase** the radius until the number of points in circle adds up to 3



**Define** a k value (in this case k = 3) **Pick** a point to predict (blue diamond) **Count** the number of closest points **Increase** the radius until the number of points within the radius adds up to 3 **Predict** the blue diamond to be a blue circle!

3/3







## **Underfitting v. Overfitting**



## Underfitting

Underfitting means we have <u>high bias</u> and <u>low variance</u>.

- Lack of relevant variables/factor
- Imposing limiting assumptions
   Linearity
  - Assumptions on distribution
  - Wrong values for parameters





## Overfitting

Overfitting means we have <u>low bias</u> and <u>high variance</u>.

- Model fits too well to specific cases
- Model is over-sensitive to sample-specific noise
- Model introduces too many variables/complexities than needed





## Relationship Between k and Fit

The **k** value you use has a relationship to the fit of the model

A higher k gives a smoother line, but too large of a k and it is the average of all the data (or the label that is most common/likely)



k=3

k=7



## **Confusion Matrix**



## What is a Confusion Matrix?

Table used to describe the performance of a classifier on a set of binary test data for which the true values are known

	P <sup>'</sup> (Predicted)	n' (Predicted)
P (Actual)	True Positive	False Negative
n (Actual)	False Positive	True Negative



## Sensitivity

#### Called the true positive rate

Tells us how many positives are correctly identified as positives **Optimize for:** Initial diagnosis of fatal disease

	P <sup>'</sup> (Predicted)	n' (Predicted)
P (Actual)	True Positive	False Negative
n (Actual)	False Positive	True Negative

**Sensitivity** = True Positive/ (True Positive + False Negative)



## Specificity

#### Called the true negative rate

Tells us how many negatives are correctly identified as negatives **Optimize for:** testing for a disease with a risky treatment

	P <sup>'</sup> (Predicted)	n' (Predicted)
P (Actual)	True Positive	False Negative
n (Actual)	False Positive	True Negative

**Specificity** = True Negative/ (True Negative + False Positive)



## Question

Which is an example of when you would want higher specificity?

- A. DNA tests for a death penalty case
- B. Deciding which iPhone to buy
- C. Airport Extra Screening





Attendance!

### **Overall Accuracy**

Proportion of correct predictions

	P <sup>'</sup> (Predicted)	n' (Predicted)
р (Actual)	True Positive	False Negative
n (Actual)	False Positive	True Negative

## **Accuracy** = (True Positive + True Negative) / Total



#### **Overall Error Rate**

Proportion of incorrect predictions

	P <sup>'</sup> (Predicted)	n' (Predicted)
р (Actual)	True Positive	False Negative
n (Actual)	False Positive	True Negative

## **Error** = (False Positive + False Negative) / Total



### **Precision**

Proportion of correct positive predictions among all positive predictions

	P <sup>'</sup> (Predicted)	n' (Predicted)
р (Actual)	True Positive	False Negative
n (Actual)	False Positive	True Negative

Precision = True Positive /
(True Positive + False Positive)



## **Coming Up**

- Assignment 6: Due Monday 03/24 at 11:59pm!
- Mid-Semester Check-In: Now!
- Next Lecture: Supervised Learning Pt. 1