# CSE/STAT 416

**Classification**

**Pemi Nguyen**
**Paul G. Allen School of Computer Science & Engineering**
**University of Washington**

**April 11, 2022**

**Slides by Hunter Schafer**

# Pre-Lecture Video 1

*Classification*

# Logistics

Homework 2 out, due this Friday

A bit more challenging than the first one so you should try to start early

COVID concerns: I and some of my staff members are still dealing with COVID situations, so we'd appreciate your patience as things can be slow.

Sorry for not releasing the re-recordings over the weekends. Will finish them today. Hopefully will help you for hw2. In the mean time you can watch Hunter's.

# Roadmap So Far

1. Housing Prices - Regression $(x, y)$ continuous value
   - Regression Model
   - Assessing Performance
   - Ridge Regression
   - LASSO

   $\}$ Regularization – Overfitting

2. Sentiment Analysis – Classification
   - Classification Overview
   - Logistic Regression

# Spam Filtering



2 output classes

Spam

Not Spam (ham)

**Input: x**
*Text of email*
*Sender*
*Subject*
*…*

**Output: y**
Spam
Ham

# Object Detection



Top Predictions

- Labrador retriever
- golden retriever
- redbone
- bloodhound
- Rhodesian ridgeback

**Input: x**
*Pixels*

**Output: y**
*Class*
*(+ Probability)*

Training Data → **x** → Feature extraction → **h(x)** → ML model → **ŷ**

**y** — *texts*

*transform*

**f̂**

Optimization Algorithm — *Gradient Descent*

Quality metric — ~~MSE~~

# Sentiment Classifier

In our example, we want to classify a restaurant review as positive or negative.



**Input: x**

**Output: y**
Predicted class

**Idea**: Use a list of ~~good~~ words and ~~bad~~ words, classifier by most frequent type of word

Positive Words: great, awesome, good, amazing, ...

Negative Words: bad, terrible, disgusting, sucks, ...

**Simple Threshold Classifier**

Input $x$: Sentence from review

Count the number of positive and negative words, in $x$

If num_positive > num_negative:

- $\hat{y} = +1$

Else:

- $\hat{y} = -1$

$=$ : arbitrary

2 positive

1 negative

Example: "Sushi was great, the food was awesome, but the service was terrible"

→ positive

9

# Limitations of Implementation 1

How do we get list of positive/negative words?

Words have different degrees of sentiment.
- Great > Good
  (2)    (1)
- How can we weigh them differently?

Single words are not enough sometimes...
- "Good" → Positive
- "Not Good" → Negative

Words depend on contexts (NLP?) not in this class

unigram: 1 word at a time

bigram: 2 words at a time | I have a cat

# Implementation 2: Linear Classifier

**Idea**: Use labelled training data to learn a weight for each word. Use weights to score a sentence.

| Word | Weight |
|------|--------|
| good | 1.0 |
| great | 1.5 |
| awesome | 2.7 |
| bad | 1.0 |
| terrible | -2.1 |
| awful | -3.3 |
| restaurant, the, we, where, … | 0.0 |
| … | … |

*positive* (good, great, awesome)

*negative* (bad, terrible, awful)

*neutral* (restaurant, the, we, where, …)

# Score a Sentence

| Word | Weight |
|------|--------|
| good | 1.0 |
| great | 1.5 |
| awesome | 2.7 |
| bad | -1.0 |
| terrible | -2.1 |
| awful | -3.3 |
| restaurant, the, we, where, … | 0.0 |
| … | … |

Input $x_i$:

$freq = 1$
$\parallel$

"Sushi was **awesome**, the food was **great**, but the service was **terrible**"

$\parallel$
$1$

$\parallel$
$1$

$Score(x) =$

$2.7 \times 1$  (awesome)

$+ 1.5 \times 1$  (great)

$+ (-2.1) \times 1$

$= 2.1 \Rightarrow positive$

Linear classifier, because output is linear weighted sum of inputs.

Will learn how to learn weights soon!

Think 👤

1 min

pollev.com/cs416

**What is the score of this sentence?**

| Word | Weight |
|------|--------|
| good | 1.0 |
| great | 1.5 |
| awesome | 2.7 |
| bad | -1.0 |
| terrible | -2.1 |
| awful | -3.3 |
| restaurant, the, we, where, … | 0.0 |
| … | … |

Input $x_i$:

"Sushi was **awful**, but the food was **great**, and the service was **great** as well".

$$Score(x)$$
$$= (-3.3) \times 1$$
$$+ 1.5 \times 2$$
$$= -0.3 \implies predicted: negative$$

1:00

13

**Idea**: Use labelled training data to learn a weight for each word. Use weights to score a sentence.

See last slide for example weights and scoring.

**Linear Classifier**

Input $x$: Sentence from review

Compute $Score(x)$

If $Score(x) > 0$:                    threshold

- $\hat{y} = +1$

Else:

- $\hat{y} = -1$

$Score = 0 \implies$ choose arbitrary

14

# Linear Classifier Notation

**Model**: $\hat{y}^{(i)} = sign\left(Score(x^{(i)})\right)$

$$Score(x_i) = w_0 h_0\left(x^{(i)}\right) + w_1 h_1\left(x^{(i)}\right) + \ldots + w_D h_D\left(x^{(i)}\right)$$

$$= \sum_{j=0}^{D} w_j h_j(x^{(i)}) \qquad \text{features} \quad \text{weights}$$

$$= w^T h(x^{(i)})$$

We will also use the notation

$$\hat{s}^{(i)} = Score\left(x^{(i)}\right) = w^T h\left(x^{(i)}\right)$$
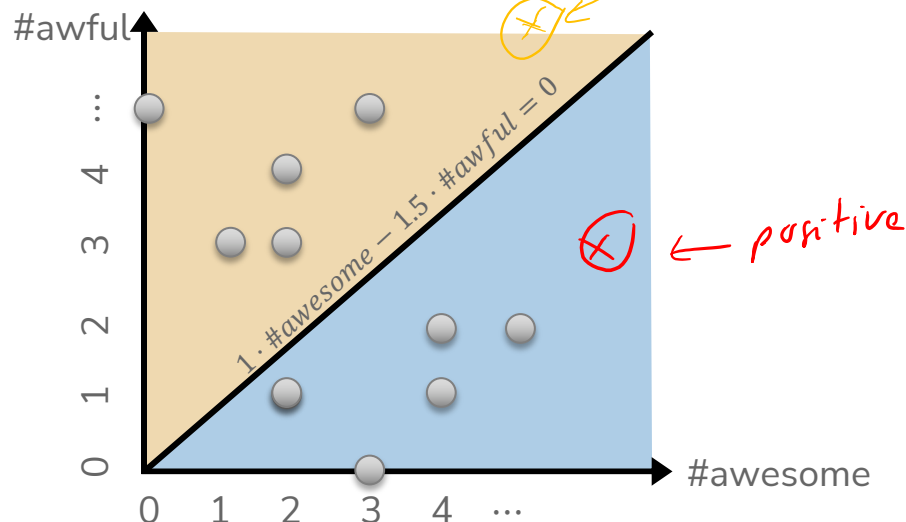
$$\hat{y}^{(i)} = sign(\hat{s}^{(i)})$$

# Decision Boundary

Consider if only two words had non-zero coefficients

| Word | Coefficient | Weight |
|------|-------------|--------|
|      | $w_0$       | 0.0    |
| awesome | $w_1$    | 1.0    |
| awful   | $w_2$    | -1.5   |

$\hat{s} = 1 \cdot \#awesome - 1.5 \cdot \#awful$



#awful

negative

$1 \cdot \#awesome - 1.5 \cdot \#awful = 0$

positive

#awesome

# Limitations of Implementation 2

Words are not single on their own, but depend on surrounding contexts:
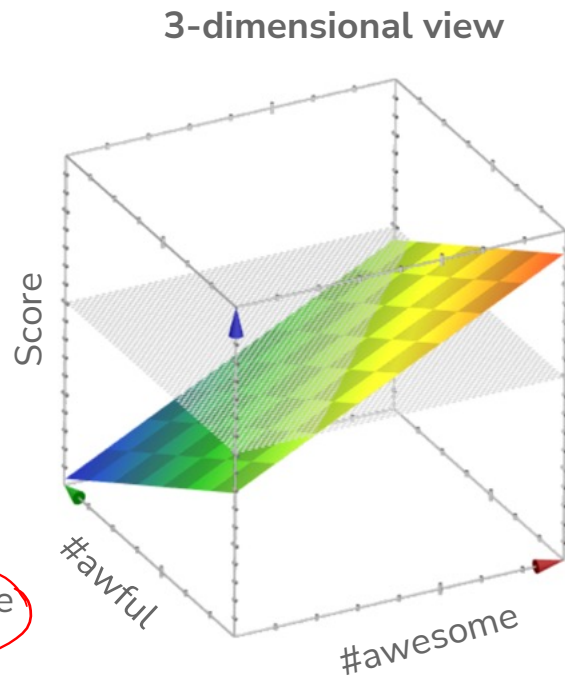
"Not not good" -> positive
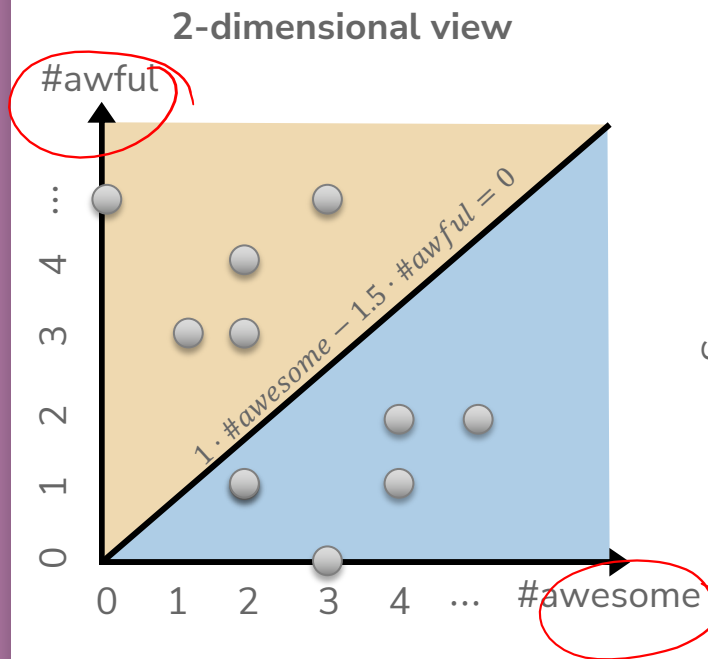
A bit sad -> negative (but not low negative score)

Linear classifier can't learn complex models well

# Decision Boundary

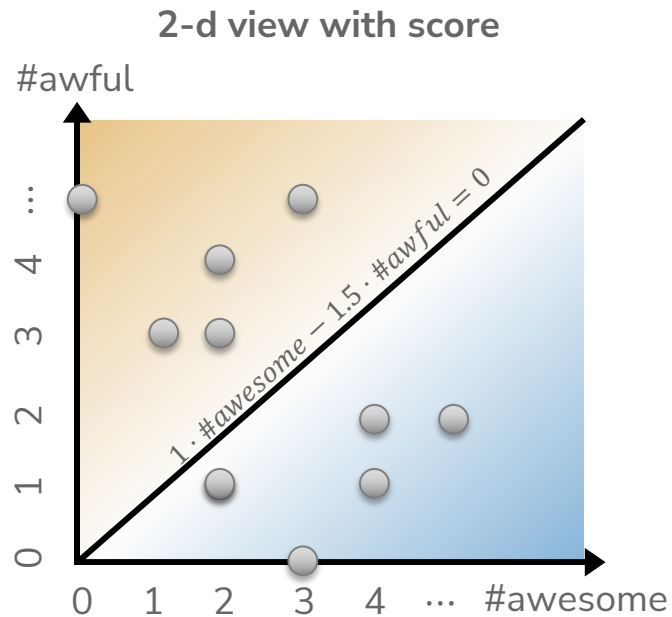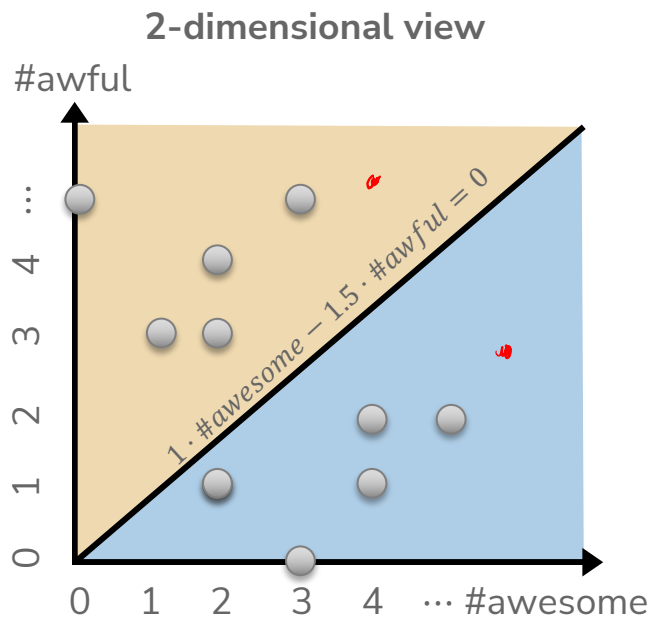$$Score(x) = 1 \cdot \#awesome - 1.5 \cdot \#awful$$

**2-dimensional view**



**3-dimensional view**



Generally, with classification we don't us a plot like the 3d view since it's hard to visualize, instead use 2d plot with decision boundary

18

# Decision Boundary

$$Score(x) = 1 \cdot \#awesome - 1.5 \cdot \#awful$$

**2-dimensional view**

#awful

... 4 3 2 1 0

$1 \cdot \#awesome - 1.5 \cdot \#awful = 0$

0 1 2 3 4 ... #awesome

**2-d view with score**

#awful

... 4 3 2 1 0

$1 \cdot \#awesome - 1.5 \cdot \#awful = 0$
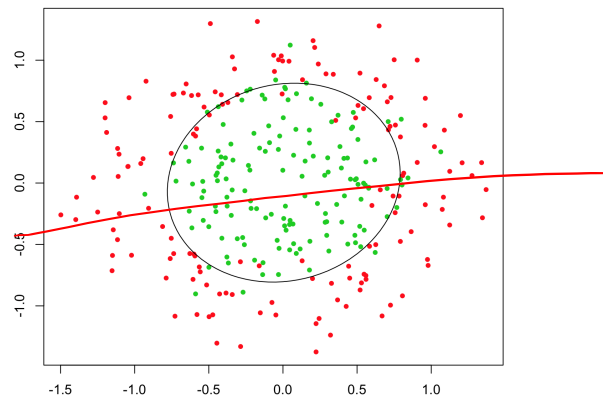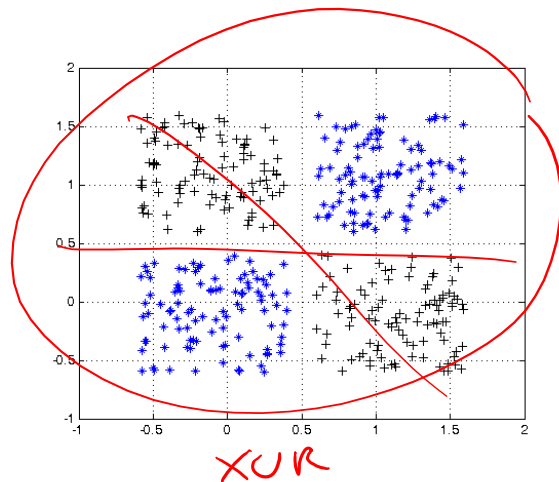
0 1 2 3 4 ... #awesome

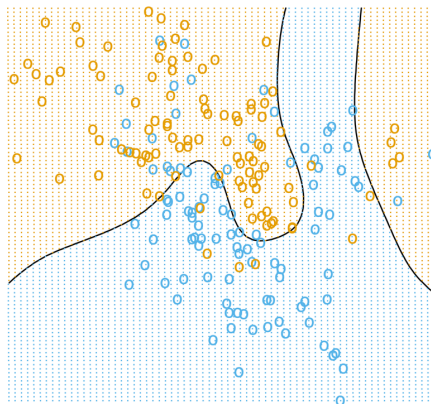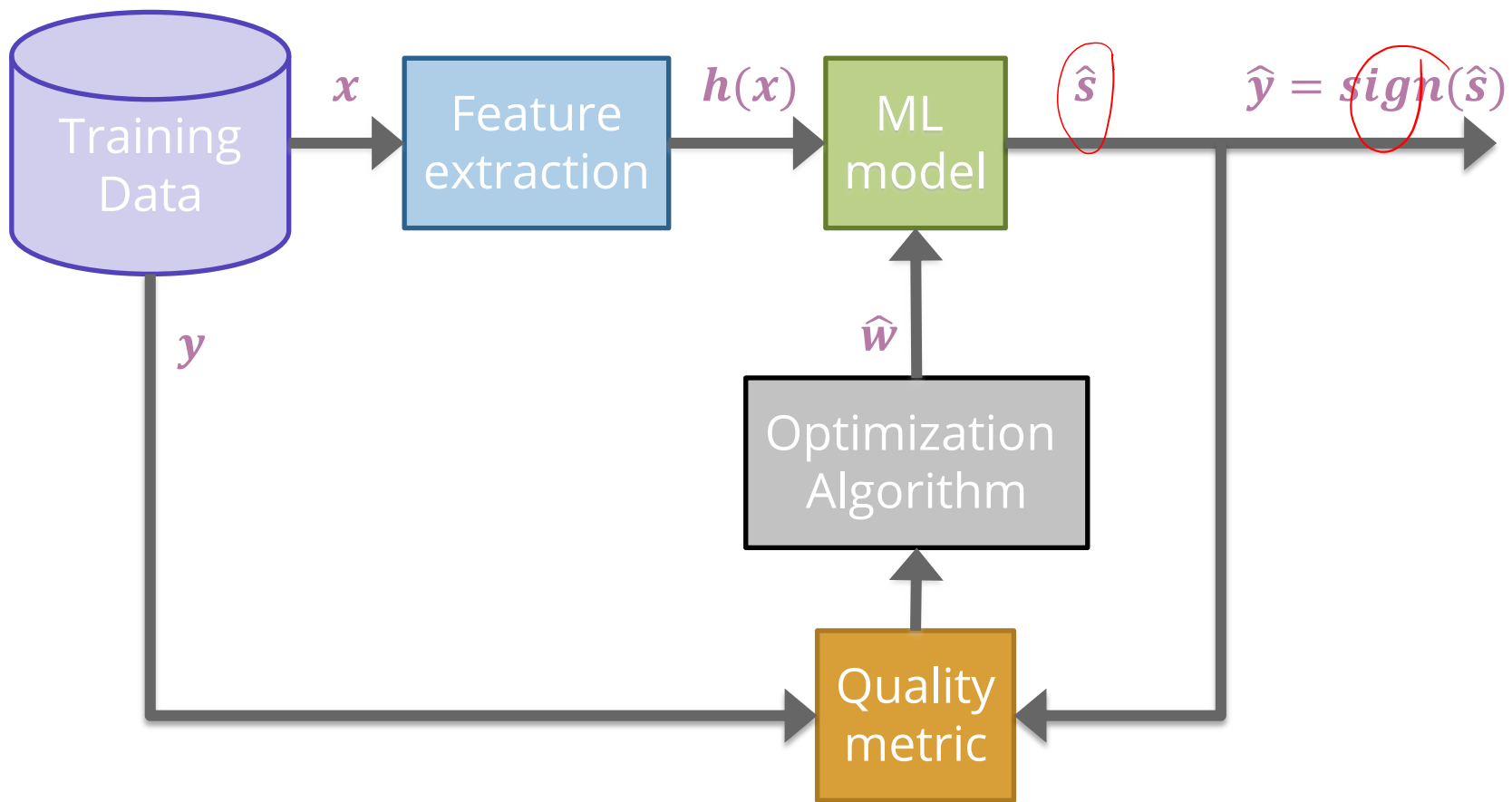# Class Session

# Complex Decision Boundaries?

What if we want to use a more complex decision boundary?

– Need more complex model/features!

$h(x)$



$X \cup R$

# Evaluating Classifiers

# Classification Error

Ratio of examples where there was a mistaken prediction

What's a mistake?

If the true label was positive ($y = +1$),
but we predicted negative ($\hat{y} = -1$)

If the true label was negative ($y = -1$),
but we predicted positive ($\hat{y} = +1$)

**Classification Error**

$$\frac{\#\ mistakes}{\#\ examples} = \frac{1}{n} \sum_i \mathbb{1}\left\{ y^{(i)} \neq \hat{y}^{(i)} \right\}$$

**Classification Accuracy**

$$1 - error = \frac{\#\ correct}{\#\ examples}$$

$$3 \qquad 3.5$$

$$0.01 \qquad 0.0001$$

# What's a good accuracy?

For binary classification:

Should at least beat random guessing...

Accuracy should be at least 0.5

*example*

For multi-class classification ($k$ classes):

~~test error~~

Should still beat random guessing

Accuracy should be at least $1/k$

- 3-class: 0.33
- 4-class: 0.25
- ...

**Besides that, higher accuracy means better, right?**

# Detecting Spam

Imagine I made a "Dummy Classifier" for detecting spam

> The classifier ignores the input, and always predicts spam.

> This actually results in 90% accuracy! Why?
> - Most emails are spam…

This is called the **majority class classifier.**

A classifier as simple as the majority class classifier can have a high accuracy if there is a **class imbalance**.

> A class imbalance is when one class appears much more frequently than another in the dataset

This might suggest that accuracy isn't enough to tell us if a model is a good model.

# Assessing Accuracy

Always digging in and ask critical questions of your accuracy.

Is there a **class imbalance**?

How does it compare to a baseline approach?
- Random guessing
- Majority class
- ...

Most important: **What does my application need?**
- What's good enough for user experience?
- What is the impact of a mistake we make?

Brain Break

# Confusion Matrix

For binary classification, there are only two types of mistakes

$$\hat{y} = +1, \quad y = -1$$

$$\hat{y} = -1, \quad y = +1$$

Generally we make a **confusion matrix** to understand mistakes.

**Predicted Label**

| | ➕ | ➖ |
|---|---|---|
| ➕ | True Positive (TP) | False Negative (FN) |
| ➖ | False Positive (FP) | True Negative (TN) |

**True Label**

# Confusion Matrix Example

**Predicted Label**

$n = 100$

| | | Predicted + | Predicted − |
|---|---|---|---|
| **True Label** | + | True Positive (TP) $\textcircled{50}$ | False Negative (FN) 10 |
| | − | False Positive (FP) 15 | True Negative (TN) $\textcircled{35}$ |

$$\text{Accuracy} = \frac{\#\ correct}{\#\ examples} = \frac{50 + 35}{100} = 85\%$$

# Which is Worse?

**What's worse, a false negative or a false positive?**

It entirely depends on your application!

**Detecting Spam**  ← *positive*

False Negative: Annoying

False Positive: Email lost  *worse*

**Medical Diagnosis**

False Negative: Disease not treated  *worse*

False Positive: Wasteful treatment

In almost every case, how treat errors depends on your context.

**Think**  ⌾

2 mins

In your group, discuss an example of the social implications using machine learning classifiers in making decisions.

# Errors and Fairness

We mentioned on the first day how ML is being used in many contexts that impact crucial aspects of our lives.

Models making errors is a given, what we do about that is a choice:

Are the errors consequential enough that we shouldn't use a model in the first place?

Do different demographic groups experience errors at different rates?
- If so, we would hopefully want to avoid that model!

Will talk more about how to define whether or a not a model is fair / discriminatory in a later lecture! Will use these notions of error as a starting point!

# Binary Classification Measures

Notation

$$C_{TP} = \#TP, \quad C_{FP} = \#FP, \quad C_{TN} = \#TN, \quad C_{FN} = \#FN$$

$$N = C_{TP} + C_{FP} + C_{TN} + C_{FN}$$

$$N_P = C_{TP} + C_{FN}, \quad N_N = C_{FP} + C_{TN}$$

**Error Rate**

$$\frac{C_{FP} + C_{FN}}{N}$$

**Accuracy Rate**

$$\frac{C_{TP} + C_{TN}}{N}$$

**False Positive rate (FPR)**

$$\frac{C_{FP}}{N_N}$$

**False Negative Rate (FNR)**

$$\frac{C_{FN}}{N_P}$$

**True Positive Rate or Recall**

$$\frac{C_{TP}}{N_P}$$

**Precision**

$$\frac{C_{TP}}{C_{TP} + C_{FP}}$$

**F1-Score**

$$2 \frac{Precision \cdot Recall}{Precison + Recall}$$

See more!

# Multiclass Confusion Matrix

Consider predicting (*Healthy, Cold, Flu*)

|  | Healthy | Cold | Flu |
|---|---|---|---|
| **Healthy** | 60 | 8 | 2 |
| **Cold** | 4 | 12 | 4 |
| **Flu** | 0 | 2 | 8 |

Predicted Label

True Label

Think

2 min

Suppose we trained a classifier and computed its confusion matrix on the training dataset. **Is there a class imbalance in the dataset and if so, which class has the highest representation?**

**Predicted Label**

| | Pupper | Doggo | Boofer |
|---|---|---|---|
| Pupper | 2 | 27 | 4 |
| Doggo | 4 | 25 | 4 |
| Boofer | 1 | 30 | 2 |

True Label

2:00

# Learning Theory

# How much data?

The more the merrier

    But data quality is also an extremely important factor

Theoretical techniques can bound how much data is needed

    Typically too loose for practical applications

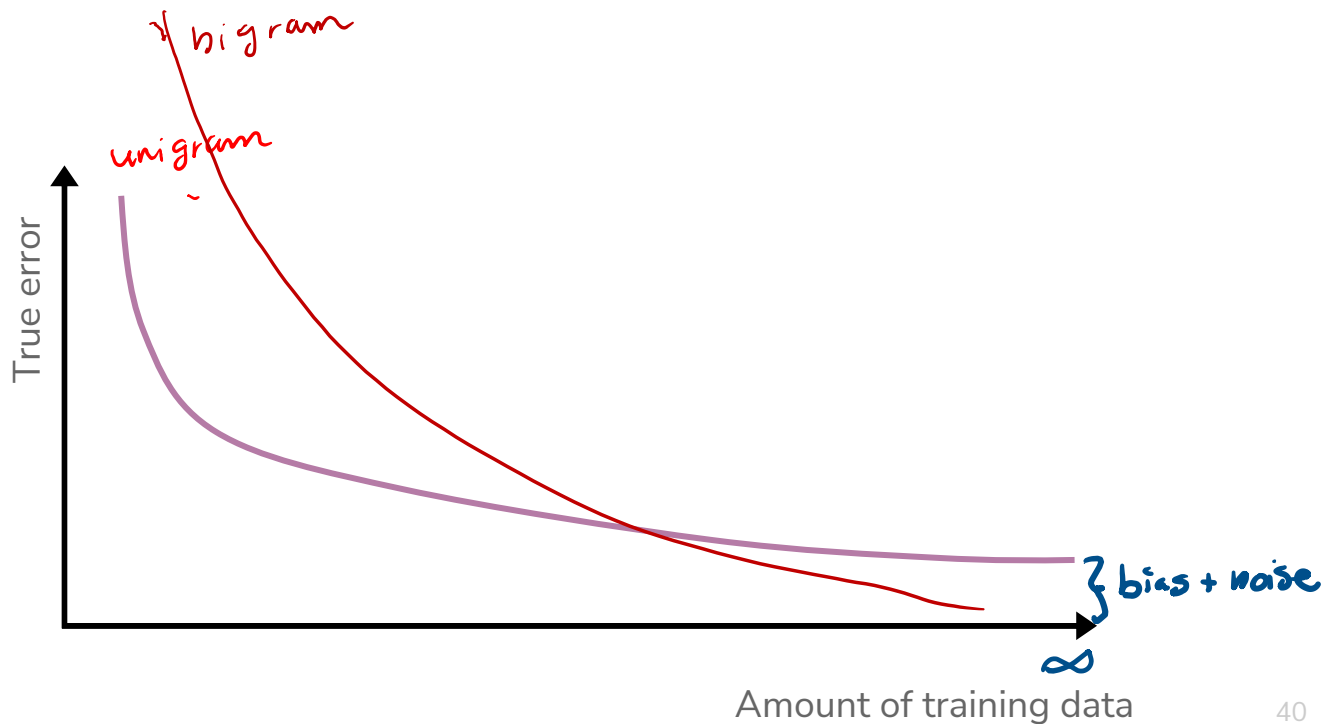    But does provide some theoretical guarantee

In practice

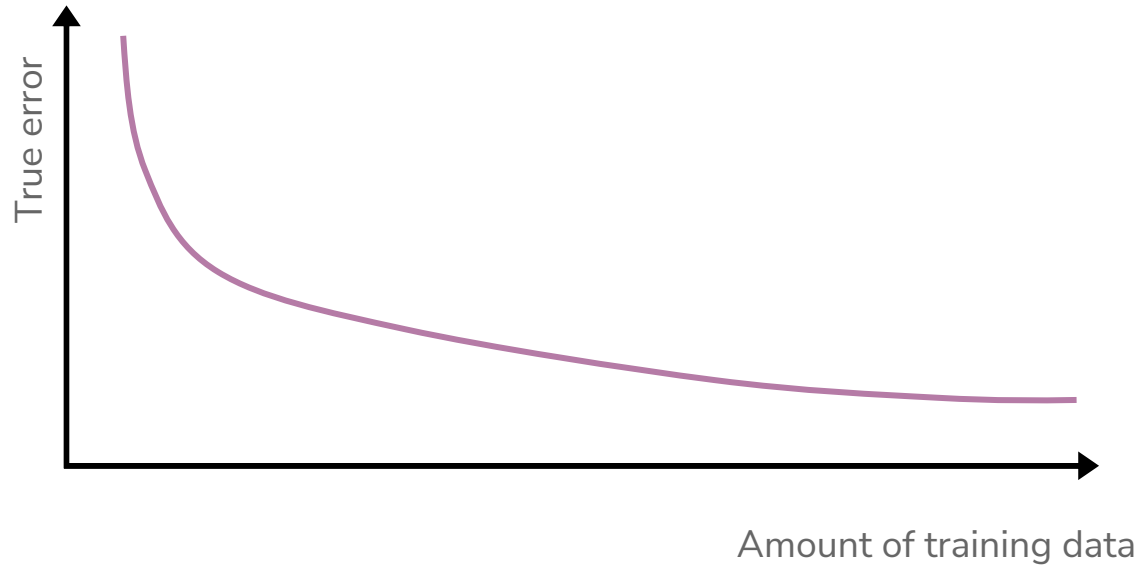    More complex models need more data

# Learning Curve

How does the true error of a model relate to the amount of training data we give it?

Hint: We've seen this picture before



*Amount of training data*

*True error*

bigram

unigram

bias + noise

∞

# Learning Curve

What if we use a more complex model?



True error

Amount of training data

# Threshold Model

# Change Threshold

What if I never want to make a false positive prediction?

*Always predict neg* $(\alpha = \infty)$

What if I never want to make a false negative prediction?
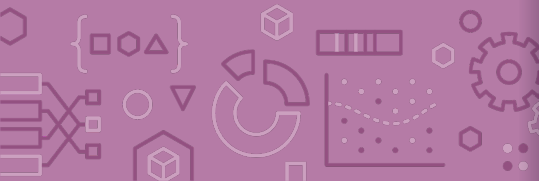
*Always predict pos* $(\alpha = -\infty)$

One way to control for our application is to change the scoring threshold. (Could also change intercept!)

If $Score(x) > \alpha$:   *0*
- Predict $\hat{y} = +1$

Else:
- Predict $\hat{y} = -1$

# Next Time

We will talk about learning classifiers that model the probability of seeing a particular class at a given input.

$$P(y|x)$$

Normally assume some structure on the probability (e.g. linear)
$$P(y|x,w) \approx w^T x$$

Use machine learning algorithm to learn approximate $\hat{w}$ such that
$$\hat{P}(y|x) = P(y|x,\hat{w})$$

And $P(y|x)$ and $\hat{P}(y|x)$ are close.

# Recap

**Theme**: Describe high level idea and metrics for classification

**Ideas:**

Applications of classification

Linear classifier

Decision boundaries

Classification error / Classification accuracy

Class imbalance

Confusion matrix

Learning theory