

CSE/STAT 416

Precision/Recall k-Nearest Neighbors

Vinitra Swamy
University of Washington
July 20, 2020



Reflection

Good sense of ensemble models!

Pacing of course seems to be good for the majority of the class

- Information density
- Lecture: brain breaks, poll-everywhere questions

Group project? -> happening soon, Homework 5

Thank you!

Ensemble Method

Instead of switching to a brand new type of model that is more powerful than trees, what if we instead tried to make the tree into a more powerful model.

What if we could combine many weaker models in such a way to make a more powerful model?

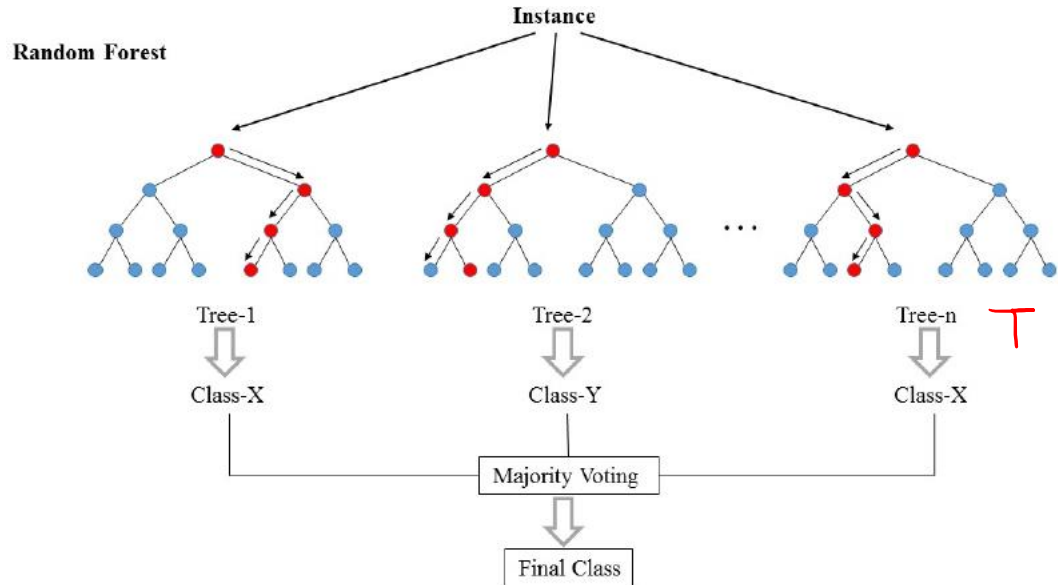
A **model ensemble** is a collection of (generally weak) models that are combined in such a way to create a more powerful model.

There are two common ways this is done with trees

- Random Forest (Bagging)
- AdaBoost (Boosting)

Random Forest

A **Random Forest** is a collection of T Decision Trees. Each decision tree casts a “vote” for a prediction and the ensemble predicts the majority vote of all of its trees.



AdaBoost

AdaBoost is a model similar to Random Forest (an ensemble of decision trees) with two notable differences that impact how we train it quite severely.

- Instead of using high depth trees that will overfit, we limit ourselves to **decision stumps**.
- Each model in the ensemble gets a weight associated to it, and we take a weighted majority vote

$$\hat{y} = \hat{F}(x) = \text{sign} \left(\sum_{t=1}^T \hat{w}_t \hat{f}_t(x) \right)$$

AdaBoost

Ada Glance

Train

$$\alpha_i = 1/N$$

for t in $[1, 2, \dots, T]$:

- Learn $\hat{f}_t(x)$ based on weights α_i
- Compute model weight \hat{w}_t
- Recompute weights α_i
- Normalize α_i

$$\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - \text{WeightedError}(\hat{f}_t)}{\text{WeightedError}(\hat{f}_t)} \right)$$

$$\alpha_i \leftarrow \begin{cases} \alpha_i e^{-\hat{w}_t}, & \text{if } \hat{f}_t(x_i) = y_i \\ \alpha_i e^{\hat{w}_t}, & \text{if } \hat{f}_t(x_i) \neq y_i \end{cases}$$

$$\alpha_i \leftarrow \frac{\alpha_i}{\sum_{j=1}^n \alpha_j}$$

Predict

$$\hat{y} = \hat{F}(x) = \text{sign} \left(\sum_{t=1}^T \hat{w}_t \hat{f}_t(x) \right)$$

Roadmap

1. Housing Prices - Regression
 - Regression Model
 - Assessing Performance
 - Ridge Regression
 - LASSO

2. Sentiment Analysis – Classification
 - Classification Overview
 - Logistic Regression
 - Decision Trees
 - Ensemble Methods

3. Document Retrieval – Clustering and Similarity
 - Precision / Recall
 - k-Nearest Neighbor
 - Kernel Methods
 - Locality Sensitive Hashing
 - Clustering
 - Hierarchical Clustering

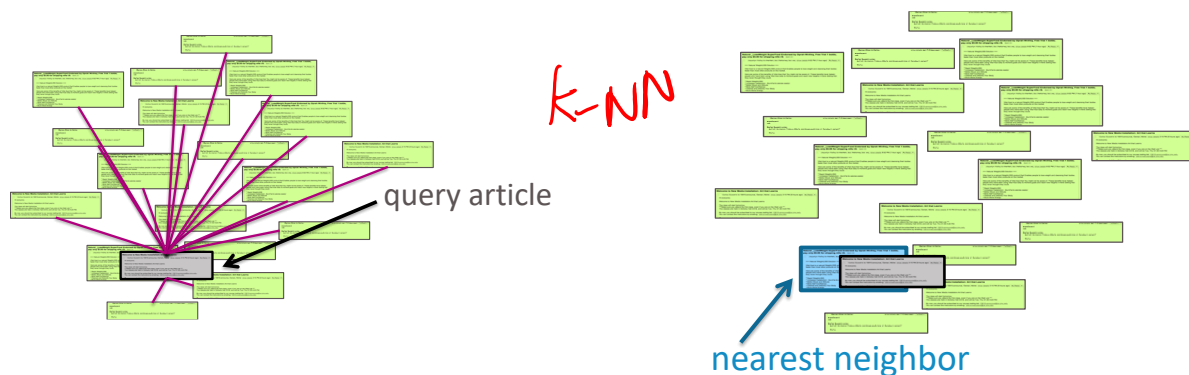
Supervised

Unsupervised

Document Retrieval

- Consider you had some time to read a book and wanted to find other books similar to that one.
- If we wanted to write a system to recommend books
 - How do we measure similarity?
 - How do we search over books?
 - How do we measure accuracy?

Big Idea: Define an **embedding** and a **similarity metric** for the books, and find the **“nearest neighbor”** to some query book.



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?

Confusion Matrix

For binary classification, there are only two types of mistakes

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

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

		Predicted Label	
		+	-
True Label	+	True Positive (TP)	False Negative (FN)
	-	False Positive (FP)	True Negative (TN)

Binary Classification Measures

Notation

- $C_{TP} = \#TP$, $C_{FP} = \#FP$, $C_{TN} = \#TN$, $C_{FN} = \#FN$
- $N = C_{TP} + C_{FP} + C_{TN} + C_{FN}$
- $N_P = C_{TP} + C_{FP}$, $N_N = C_{FP} + C_{TN}$

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


$$\frac{T_P}{N_P}$$

Precision


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

F1-Score

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

[See more!](#)

Assessing Accuracy

Often with binary classification, we treat the positive label as being the more important of the two. We then often then focus on these metrics:

Precision: Of the ones I predicted positive, how many of them were actually positive?

Recall: Of all the things that are truly positive, how many of them did I correctly predict as positive?

Precision

What fraction of the examples I predicted positive were correct?

Sentences predicted to be positive:

$\hat{y}_i = +1$

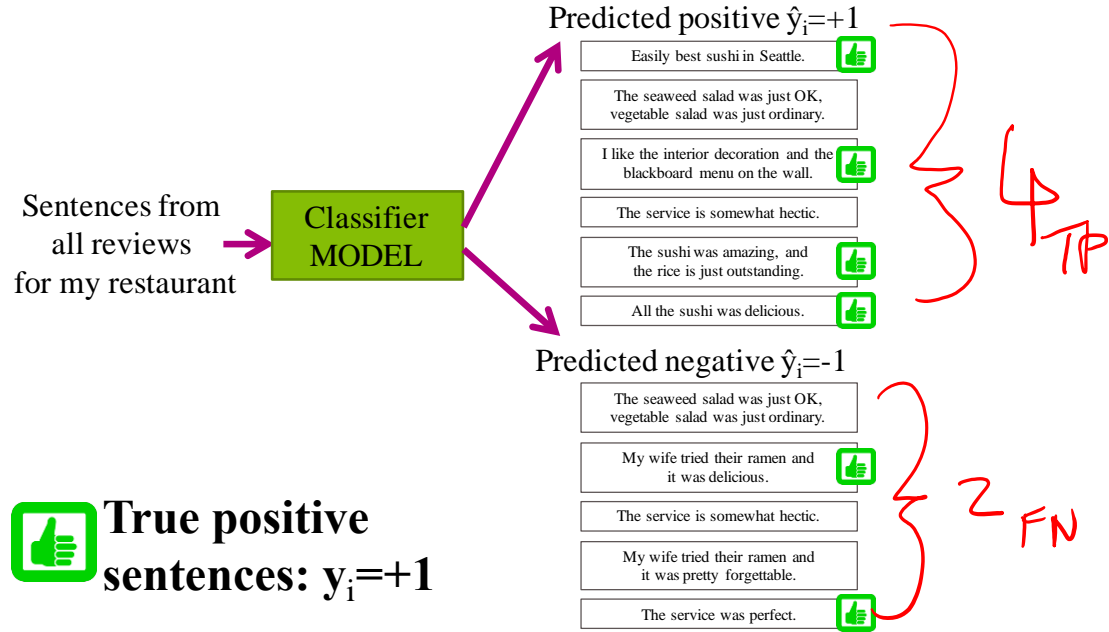
Easily best sushi in Seattle.	✓
The seaweed salad was just OK, vegetable salad was just ordinary.	✗
I like the interior decoration and the blackboard menu on the wall.	✓
The service is somewhat hectic.	✗
The sushi was amazing, and the rice is just outstanding.	✓
All the sushi was delicious.	✓

Only 4 out of 6 sentences predicted to be **positive** are actually **positive**

$$\text{precision} = \frac{C_{TP}}{C_{TP} + C_{FP}} = \frac{4}{4 + 2} = \frac{2}{3}$$

Recall

Of the truly positive examples, how many were predicted positive?



$$recall = \frac{C_{TP}}{N_P} = \frac{C_{TP}}{C_{TP} + C_{FN}} = \frac{4}{4+2} = \frac{2}{3}$$

Precision & Recall

An optimistic model will predict almost everything as positive

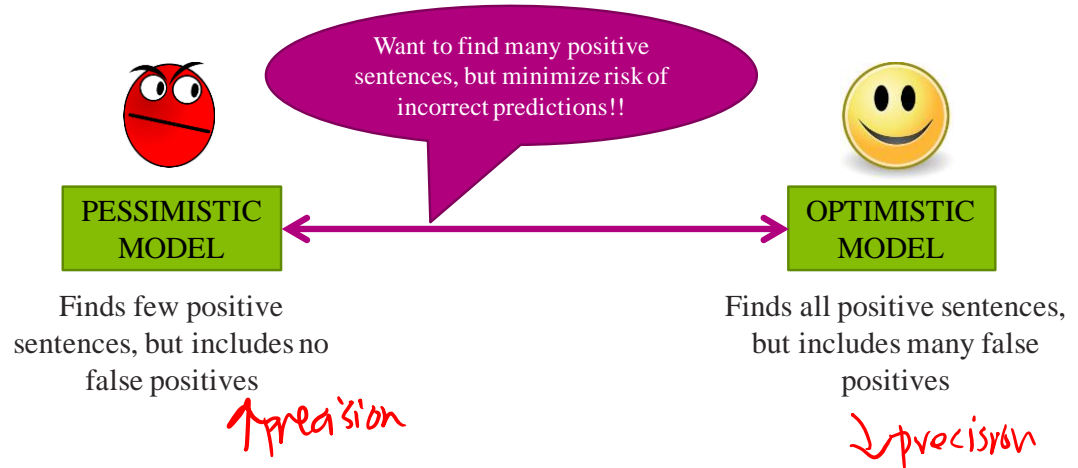


- High recall, low precision

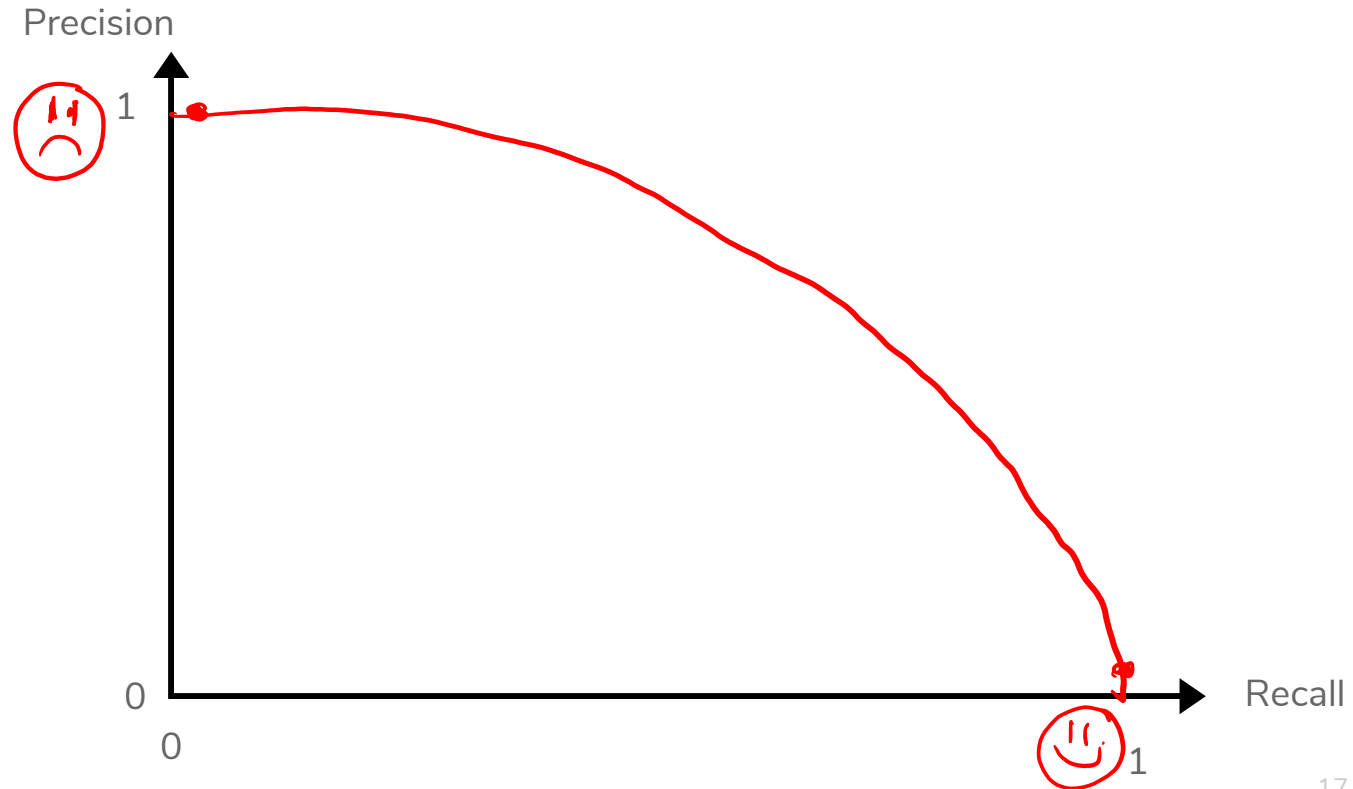
A pessimistic model will predict almost everything as negative



- High precision, low recall



Precision- Recall Curve



Controlling Precision/Recall

Depending on your application, precision or recall might be more important

- Ideally you will have high values for both, but generally increasing recall will decrease precision and vice versa.

For logistic regression, we can control for how optimistic the model is by changing the threshold for positive classification

Before

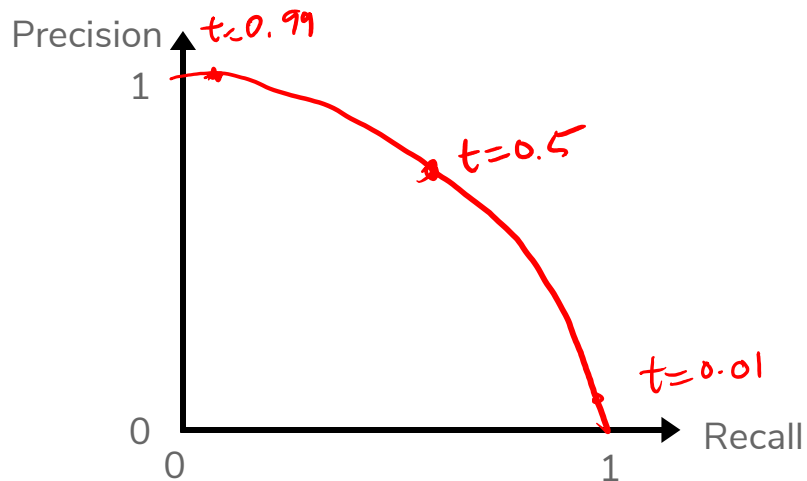
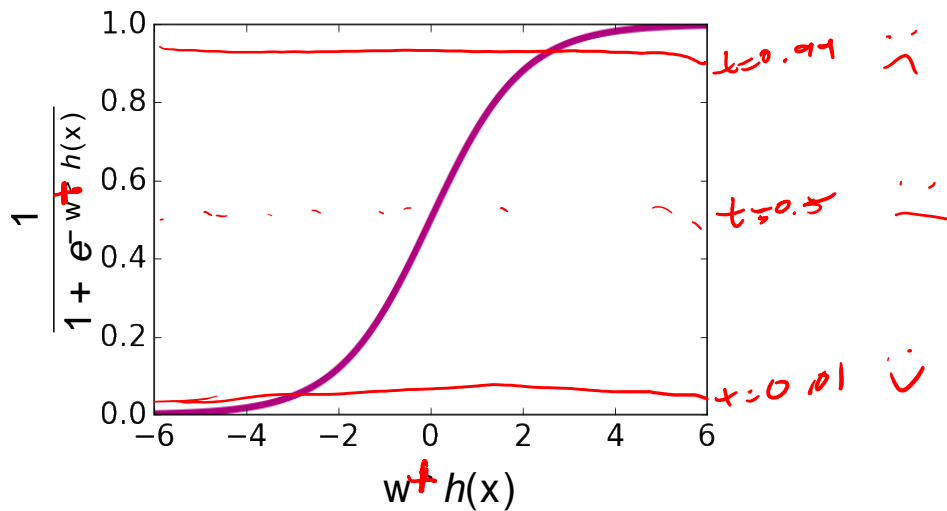
$$\hat{y}_i = +1 \text{ if } \hat{P}(y = +1|x_i) > 0.5 \text{ else } \hat{y}_i = -1$$

Now

$$\hat{y}_i = +1 \text{ if } \hat{P}(y = +1|x_i) > t \text{ else } \hat{y}_i = -1$$

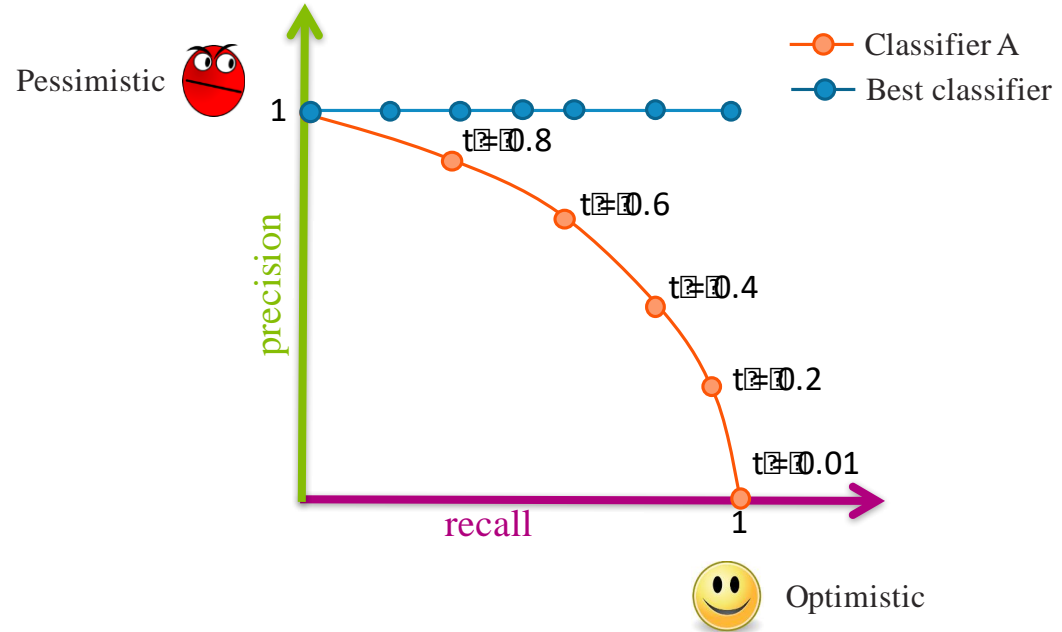
0.99 → 😊
0.01 → ☹️

Precision-Recall Tradeoff



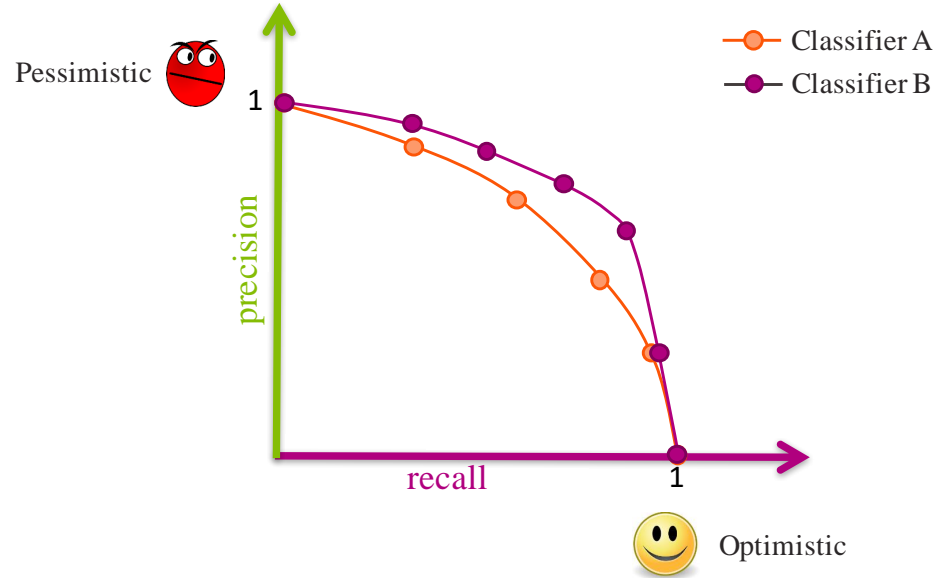
Precision-Recall Curve

Can try every threshold to get a curve like below



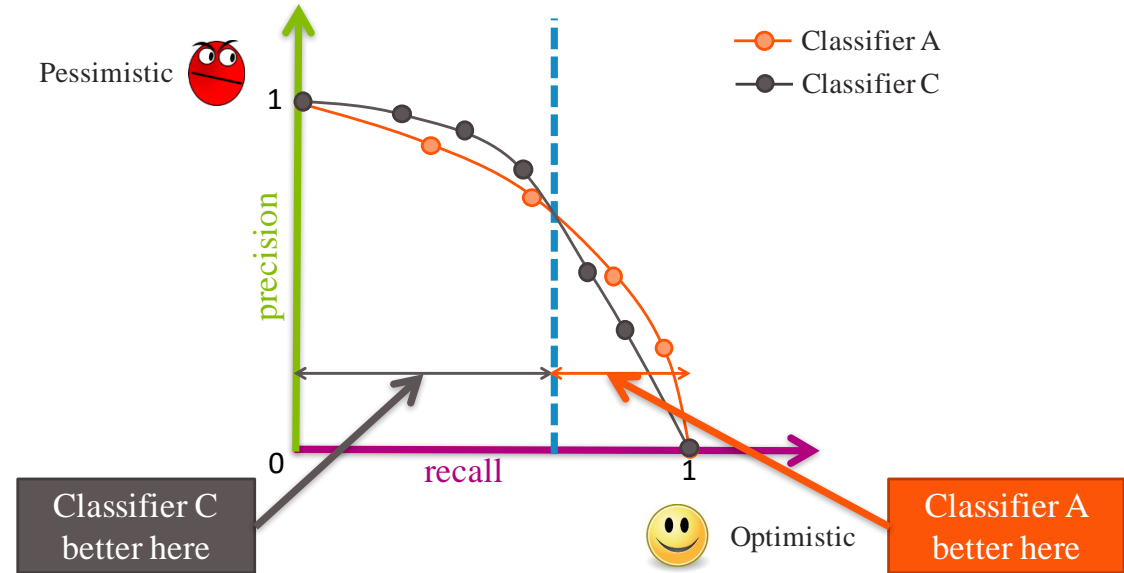
Precision-Recall Curve

Classifier B is strictly better than Classifier A



Precision-Recall Curve

Most times, the classifiers are incomparable



Compare Classifiers

Often come up with a single number to describe it

- F1-score, AUC, etc.
- Remember, what your application needs is most important





Also common to use precision at k

- If you show the top **k** most likely positive examples, how many of them are true positives

Showing
k=5 sentences
on website



Sentences model
most sure are positive

- Easily best sushi in Seattle. 
- My wife tried their ramen and it was pretty forgettable.
- The sushi was amazing, and the rice is just outstanding. 
- All the sushi was delicious. 
- The service was perfect. 



precision at $k = 0.8$

9:35



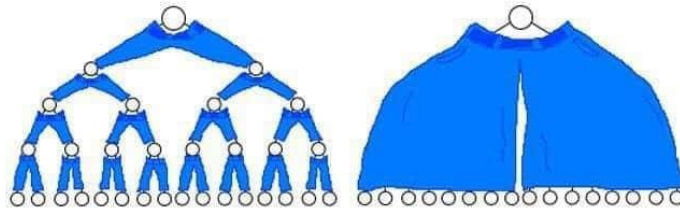
Brain Break

If a binary tree wore pants would he wear them

like this

or

like this?



Hugo Larochelle liked



Michael Hoffman
@michaelhoffman

I've tried this "artificial intelligence" plant identification app multiple times on the same tree and get a different answer each time. They must be using random forest

6:51 pm · 15 Jul 20 · [Twitter Web App](#)

205 Retweets and comments 1,952 Likes

Poll Everywhere

Think 

1 min

A model with high bias will have high precision and low recall.

- True
- False

pollev.com/cse416



1:00

Poll Everywhere

Pair 

2 min

A model with high bias will have high precision and low recall.

- True
- False

pollev.com/cse416

2:00

Nearest Neighbors

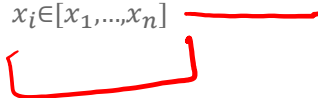
1-Nearest Neighbor

Input

- x_q : Query example (e.g. my book)
- x_1, \dots, x_n : Corpus of documents (e.g. Amazon books)

Output

- The document in corpus that is most similar to x_q

$$x^{NN} = \arg \min_{x_i \in [x_1, \dots, x_n]} \text{distance}(x_q, x_i)$$


It's very critical to properly define how we represent each document x_i and the similarity metric *distance*! Different definitions will lead to very different results.

1-Nearest Neighbor

How long does it take to find the 1-NN? About n operations

Input: x_q

$x^{NN} = \emptyset$

$nn_dist = \infty$

for $x_i \in [x_1, \dots, x_n]$:

$dist = distance(x_q, x_i)$

if $dist < nn_dist$:

$x^{NN} = x_i$

$nn_dist = dist$

Output: x^{NN}

$O(n)$

k-Nearest Neighbors

Input

- x_q : Query example (e.g. my book)
- x_1, \dots, x_n : Corpus of documents (e.g. Amazon books)

Output

- List of k documents most similar to x_q

Formally

$$X^{k\text{-NN}} = \{x^{NN_1}, x^{NN_2}, \dots, x^{NN_k}\}$$

for all x_i not in $X^{k\text{-NN}}$

$$\text{dist}(x_q, x_i) > \max_{j=1 \dots k} \text{dist}(x_q, x^{NN_j})$$

k-Nearest Neighbors

Same idea as 1-NN algorithm, but maintain list of k-NN

Input: x_q

$X^{k-NN} = [x_1, \dots, x_k]$

$nn_dists = [dist(x_1, x_q), dist(x_2, x_q), \dots, dist(x_k, x_q)]$

for $x_i \in [x_{k+1}, \dots, x_n]$:

$dist = distance(x_q, x_i)$

if $dist < \max(nn_dists)$:

remove largest dist from X^{k-NN} and nn_dists

add x_i to X^{k-NN} and $distance(x_q, x_i)$ to nn_dists

Output: X^{k-NN}

x_1

3NN

$[x_1, x_2, x_3]$

dist
from
 x_q

1

3

7

x_4

5

k-Nearest Neighbors

Can be used in many circumstances!

Retrieval

Return X^{k-NN}

Regression

$$\hat{y}_i = \frac{1}{k} \sum_{j=1}^k x^{NN_j}$$

Classification

$$\hat{y}_i = \text{majority_class}(X^{k-NN})$$

Important Points

While the formalization of these algorithms are fairly tedious, the intuition is fairly simple. Find the 1 or k nearest neighbors to a given document and return those as the answer.

This intuition relies on answering two important questions

- How do we represent the documents x_i ? *embedding / representation*
- How do we measure the distance $distance(x_q, x_i)$?

*↗
distance metric*

Document Representation

$D = \# \text{unique words in our corpus}$

$[\#I, \#like, \#dogs, \#cats]$

Like our previous ML algorithms, we will want to make a vector out of the document to represent it as a point in space.

Simplest representation is the bag-of-words representation.

- Each document will become a D dimension vector where D is the number of words in the entire corpus of documents
- The value of $x_i[j]$ will be the number of times word j appears in document i .
- This ignores order of words in the document, just the counts.

"I like dogs" $\rightarrow [1, 1, 1, 0]$

"I like cats" $\rightarrow [1, 1, 0, 1]$

"I like dogs dogs" $\rightarrow [1, 1, 2, 0]$

Bag of Words

Pros

- Very simple to describe
- Very simple to compute

Cons

- Common words like “the” and “a” dominate counts of uncommon words
- Often it’s the uncommon words that uniquely define a doc.

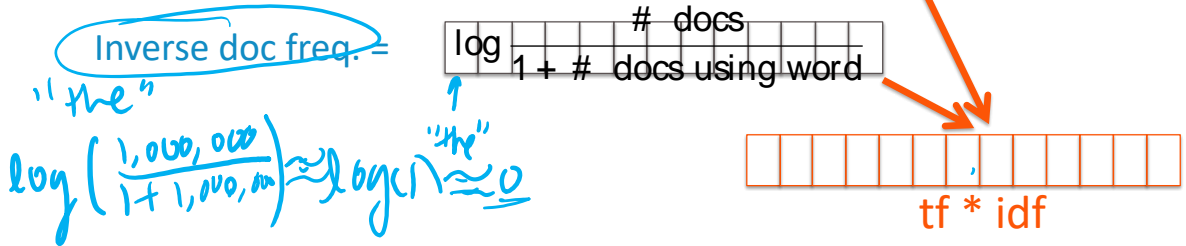
TF-IDF

Goal: Emphasize important words

- Appear frequently in the document (common locally)



- Appears rarely in the corpus (rare globally)



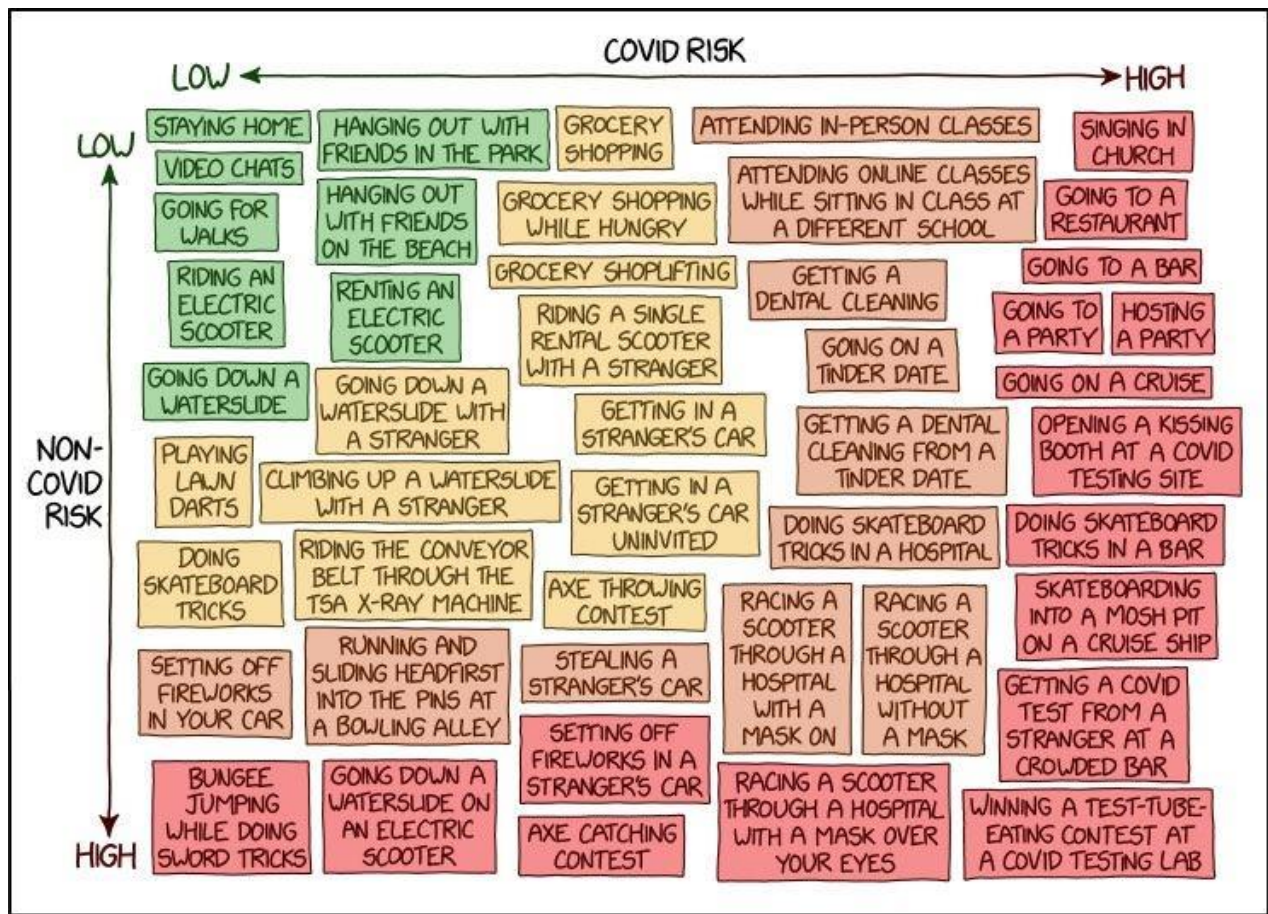
Do a pair-wise multiplication to compute the TF-IDF for each word

- Words that appear in every document will have a small IDF making the TF-IDF small!

10:12



Brain Break



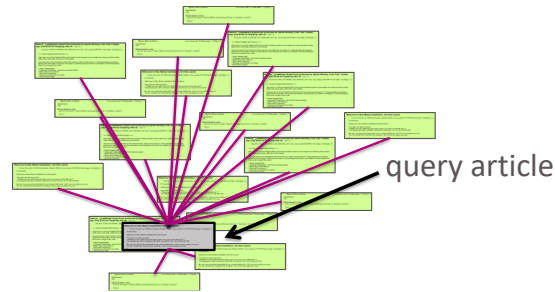
Document Retrieval

- Consider you had some time to read a book and wanted to find other books similar to that one.
- If we wanted to write a system to recommend books
 - How do we measure similarity?
 - How do we search over books?
 - How do we measure accuracy?

Bow TF-IDF

recomming!

Big Idea: Define an embedding and a similarity metric for the books, and find the “nearest neighbor” to some query book.

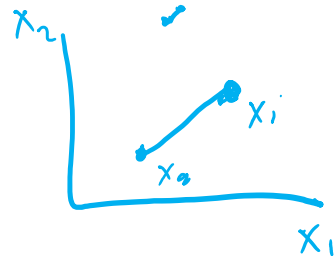


Distance

Now we will define what similarity/distance means

Want to define how “close” two vectors are. A smaller value for distance means they are closer, a large value for distance means they are farther away.

The simplest way to define distance between vectors is the **Euclidean distance**

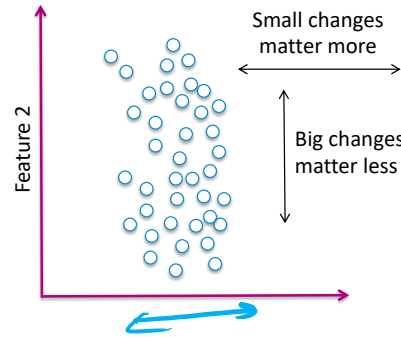


$$\text{distance}(x_i, x_q) = \|x_i - x_q\|_2$$

$$= \sqrt{\sum_{j=1}^D (x_i[j] - x_q[j])^2}$$

Weighted Distances

Some features vary more than others or are measured in different units. We can weight different dimensions differently to make the distance metric more reasonable.



Specify weights as
a function of
feature spread

For feature j :
$$\frac{1}{\max_i(x_i[j]) - \min_i(x_i[j])} = a_j$$

Weighted Euclidean distance

$$distance(x_i, x_q) = \sqrt{\sum_{j=1}^D a_j^2 (x_i[j] - x_q[j])^2}$$

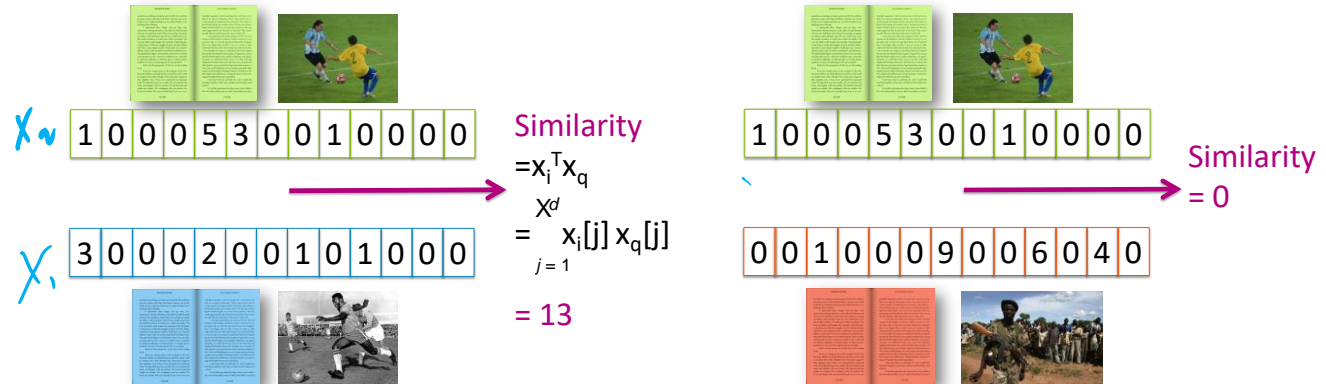
Similarity

Another natural similarity measure would use

$$x_i^T x_q = \sum_{j=1}^D \underbrace{x_i[j]} \underbrace{x_q[j]}$$

Notice this is a measure of similarity, not distance

- This means a bigger number is better



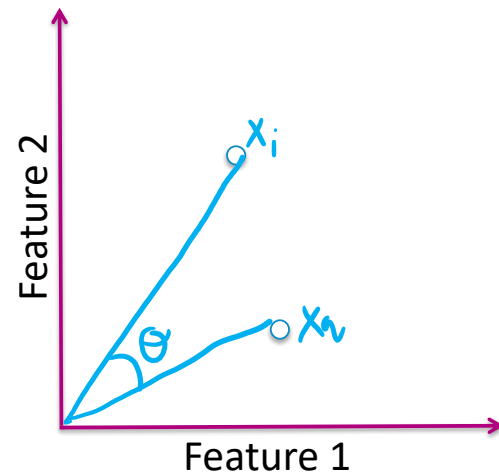
Cosine Similarity

Should we normalize the vectors before finding the similarity?

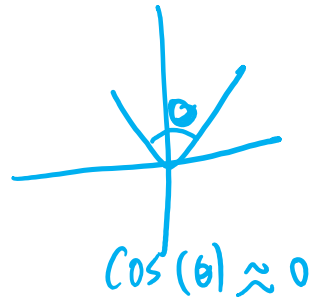
$$\text{similarity} = \frac{x_i^T x_q}{\|x_i\|_2 \|x_q\|_2} = \cos(\theta)$$

Note:

- Not a true distance metric
- Efficient for sparse vectors!



Cosine Similarity



In general

$$-1 \leq \text{cosine similarity} \leq 1$$

For positive features (like TF-IDF)

$$0 \leq \text{cosine similarity} \leq 1$$

Define

$$\text{distance} = 1 - \text{similarity}$$

To Normalize or Not To Normalize?

Not normalized



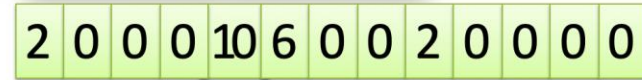
x_1



Similarity = 13



x_2

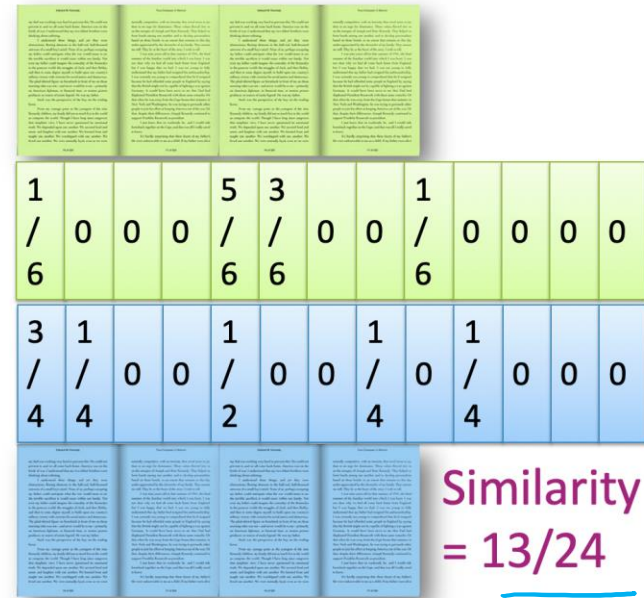
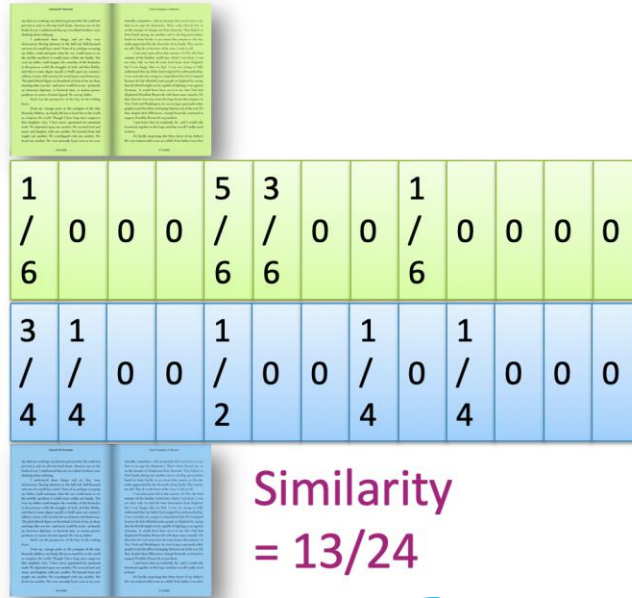


Similarity = 52



To Normalize or Not To Normalize?

Normalized



To Normalize or Not To Normalize?

Normalization is not desired when comparing documents of different sizes since it ignores length.



long document



short tweet

Normalizing can
make dissimilar
objects appear more
similar



long document



long document

Common
compromise:
Just cap maximum
word counts

In practice, can use multiple distance metrics and combine them using some defined weights

Think 

2 min

pollev.com/cse416

For the given documents, what are their Euclidean Distance and Cosine Similarity?

Assume we are using a bag of words representation

Document 1: "I really like dogs"

Document 2: "dogs are really really awesome"

Steps:

- Write out bag of words vectors
- Compute Euclidean distance
- Compute Cosine similarity

Poll Everywhere

Pair 

3 min

For the given documents, what are their Euclidean Distance and Cosine Similarity?

Assume we are using a bag of words representation

Document 1: “I really like dogs”

Document 2: “dogs are really really awesome”

Steps:

- Write out bag of words vectors
- Compute Euclidean distance
- Compute Cosine similarity

pollev.com/cse416

Think 

Document 1: "I really like dogs"

Document 2: "dogs are really really awesome"

Bag of words: (# I, # really, # like, # dogs, # are, # awesome)

$$x_1 = [1, 1, 1, 1, 0, 0] \quad x_2 = [0, 2, 0, 1, 1, 1]$$

euclidean distance: $(\|x_1 - x_2\|_2)$

$$\begin{aligned} \text{dist}(x_1, x_2) &= \sqrt{(1-0)^2 + (1-2)^2 + (1-0)^2 + (1-1)^2 + (0-1)^2 + (0-1)^2} \\ &= \sqrt{5} \end{aligned}$$

$$\text{cosine distance} = \left(1 - \frac{x_1^T x_2}{\|x_1\|_2 \|x_2\|_2}\right) = 1 - \frac{3}{\sqrt{4}\sqrt{7}} \approx 0.433$$

$$\text{dist}(x_1, x_2) = 1 - \frac{1 \cdot 0 + 1 \cdot 2 + 1 \cdot 0 + 1 \cdot 1 + 0 \cdot 1 + 0 \cdot 1}{\sqrt{1^2 + 1^2 + 1^2 + 1^2 + 0^2 + 0^2} \sqrt{0^2 + 2^2 + 0^2 + 1^2 + 1^2 + 1^2}}$$

Recap

Theme: Use nearest neighbors to recommend documents.

Ideas:

- Precision and Recall Curves
- Implement a nearest neighbor algorithm
- Compare and contrast different document representations
 - Emphasize important words with TF-IDF
- Compare and contrast different measurements of similarity
 - Euclidean and weighted Euclidean
 - Cosine similarity and inner-product similarity