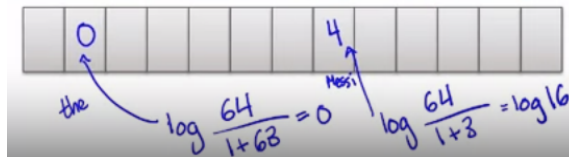


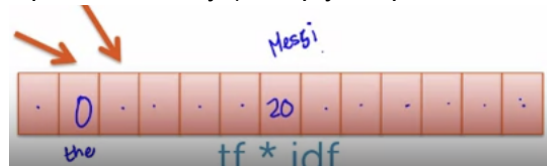
Document Retrieval	-
Overall process	<p>In order to retrieve articles that Carlos is interested in reading..</p> <ol style="list-style-type: none"> 1. How to measure similarity between articles 2. How to search over the articles that exist out there and retrieve the next article to recommend
Measuring Similarity	<p>Comparing the similarity between two documents where each word is an index and the vector contains the count of each word.</p> <div style="display: flex; align-items: center; justify-content: center; margin: 10px 0;"> <div style="display: flex; flex-direction: column; align-items: center;"> <div style="background-color: #d9ead3; border: 1px solid black; padding: 2px 5px; margin-bottom: 5px;">1 0 0 0 5 3 0 0 1 0 0 0 0</div> <div style="background-color: #d9ead3; border: 1px solid black; padding: 2px 5px; margin-bottom: 5px;">3 0 0 0 2 0 0 1 0 1 0 0 0</div> </div> <div style="margin: 0 10px; text-align: center;"> \longrightarrow </div> <div style="display: flex; flex-direction: column; align-items: center;"> <div>1×3</div> <div>$+$</div> <div>5×2</div> <div>$= 13$</div> </div> </div> <p>Shortfall with this method is, if the document doubles in length, then the number of similarities increases to 52. Therefore, we take the norm of a vector for even distribution.</p> <div style="display: flex; align-items: center; justify-content: center; margin: 10px 0;"> <div style="background-color: #d9ead3; border: 1px solid black; padding: 2px 5px; margin-bottom: 5px;">1 0 0 0 5 3 0 0 1 0 0 0 0</div> <div style="margin: 0 10px;"> $\sqrt{(1^2 + 5^2 + 3^2 + 1^2)}$ </div> </div> <div style="display: flex; align-items: center; justify-content: center; margin: 10px 0;"> <div style="background-color: #d9ead3; border: 1px solid black; padding: 2px 5px; margin-bottom: 5px;"> 1 5 3 1 / 0 0 0 / / 0 0 / 0 0 0 0 6 6 6 6 </div> </div>
Prioritizing important words with tf-idf	<p>Rare Words - words that appear infrequently but pertinent to the topic/article (futbol, messi)</p> <p>Important Words - common locally in the document but rare globally (soccer, field, goal). There can be a trade-off between local frequency and global rarity</p> <p>TF-IDF (Term Frequency - Inverse Document Frequency) - a function to weigh the rare words more than the common words</p> <ul style="list-style-type: none"> - Step 1: List out the Term Frequency



- Step 2: Apply the log function to create an inverse document frequency array



- Step 3: tf*idf array (multiply Step 1 and 2 together)



TF-IDF Logic

- Assumes log base 2
- Common words will be closer to 0 since # of docs the word is in will be high, causing log 1 which equals 0

- Inverse document frequency



$\log \frac{\# \text{ docs}}{1 + \# \text{ docs using word}}$

$\log \frac{\text{large \#}}{1 + \text{large \#}} \approx \log 1 = 0$
word in many docs
 $\log \frac{\text{large \#}}{1 + \text{small \#}} \rightarrow \text{large \#}$
rare word

Retrieving Similar Documents



Present someone with another relevant article recommendation

1-Nearest Neighbor Approach

- Input is the query article, output is the most similar article
- Algorithm: We search over each article in the corpus, calculate the similarity using TF-IDF and keep tabs on the best similarity score.

Search over each article  in corpus

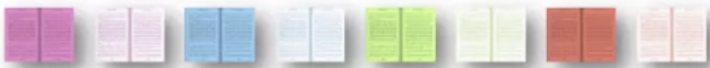
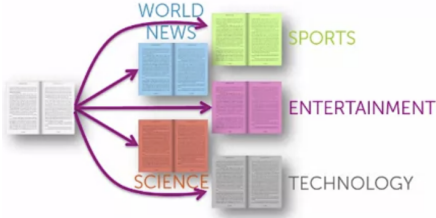
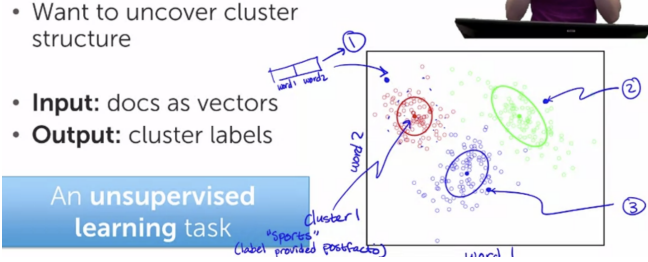
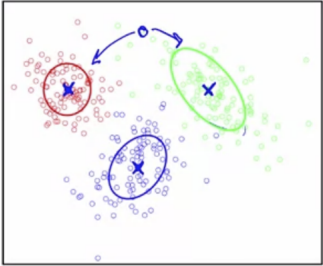
- Compute $s = \text{similarity}(\text{query article}, \text{article})$

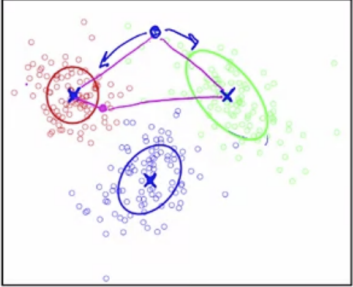
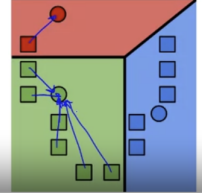
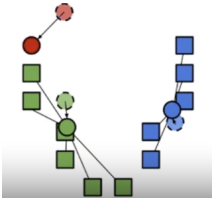
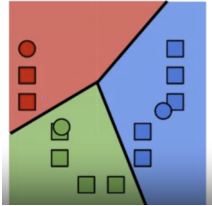
- If $s > \text{Best}_s$, record  =  and set $\text{Best}_s = s$

Return 

K-Nearest Neighbor Approach

- Input is the query article, output is list of k similar article
- Algorithm is nearly the same. We search over each article in the corpus, calculate the similarity using TF-IDF, and return the top n articles with the highest similarity score

	
<p>Clustering Documents Task Overview</p>	<p>How to quickly discover groups (clusters) of related articles instead of sorting through them one by one.</p> <p>Supervised Learning: Multiclass classification problem</p> <ul style="list-style-type: none"> - Query article is not labeled, but labels are provided (corpus is already grouped)  <p>Unsupervised Learning:</p> <ul style="list-style-type: none"> - No labels are provided and want to uncover the cluster structure <ul style="list-style-type: none"> • Want to uncover cluster structure • Input: docs as vectors • Output: cluster labels <p>An unsupervised learning task</p>  <ul style="list-style-type: none"> - The outputs shows the shapes of the clusters and a label is provided post-facto - Each cluster has a center. To measure new observations, we can do two things <ol style="list-style-type: none"> 1. Based on shape: observation is scored against the cluster center and the shape of the cluster is considered. In this case, it would be assigned to the oblong green shape  <ol style="list-style-type: none"> 2. Based on distance: observation is grouped to the cluster whose center is the shortest distance from the observation.

	
<p>Clustering Algorithm: K-Means</p>	<p>Clustering that uses distance from cluster center</p> <p>Step 1: Initialize cluster centers</p> <ul style="list-style-type: none"> - specify the number of clusters ahead of time. So we are looking at k number of clusters (3) <p>Step 2: Assign observation to the closest cluster center</p> <ul style="list-style-type: none"> - Use Voronoi Tessellation to find clusters where the items are connected to their closest centers  <p>Step 3: Revise cluster centers</p> <ul style="list-style-type: none"> - Revise as new observations come in, as I update my definition of the cluster center - Redraw the Voronoi Tessellation and reassign my observations to the nearest cluster center.  <p>Step 4: Iterate 1+2 until convergence.</p> 
<p>Other examples of clustering application</p>	<ol style="list-style-type: none"> 1. Image search 2. Grouping patients by medical conditions 3. Discover product categories from purchase histories 4. Structuring web search results 5. Discovering Similar Neighborhoods based on price of living

ML Blocking Diagram	<ol style="list-style-type: none">1. Training Data (doc id, document text table) then extract the features with tf-idf.2. The x (inputs) get put into a ML model that outputs (\hat{y}) cluster label3. To assess the accuracy of the cluster though we don't have a true cluster label to compare against, the y doesn't exist. Therefore we look at how coherent the clustering is by seeing if the distance from inputs (x) to cluster centers (\hat{w}) are minimized.4. ML algorithm iteratively updates the cluster center to minimize that distance
---------------------	--