

WEEK 3

VECTOR SPACE MODELS.

- Why learn vector space models? → to identify similarity:
Where are you heading?
Where are you from?
different meaning
- What is your age?
How old are you?
same meaning

Applications:

- > capture the dependencies between words
- > information extraction to answer the questions
- > machine translation
- > chatbots

→ Represent words and documents as vectors

→ Representation that captures relative meaning

- Word by word & word by doc.

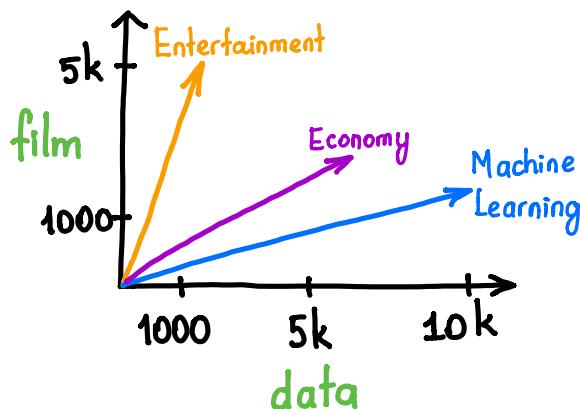
Vector representation designs (co-occurrence):

1 → word by word design: # of times they occur together within a distance k

	simple	raw	like	I
"I like simple data"				
"I prefer simple raw data" data	2	1	1	0

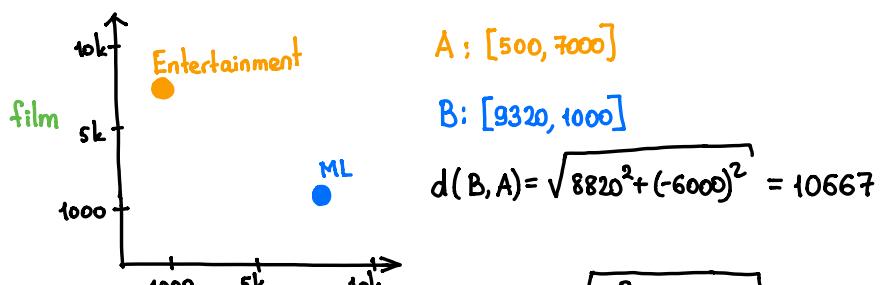
$\vec{w} \rightarrow$ word by document design: # of times a word occurs within a certain category.

	Entertainment	Economy	Machine Learning
data	500	6620	9320
film	7000	4000	1000



Measures of "similarity":
angle distance

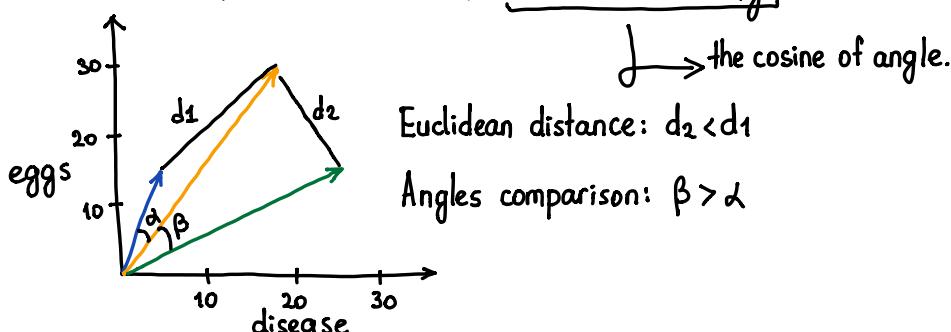
• Euclidean distance



$$d(\vec{v}, \vec{w}) = \sqrt{\sum_{i=1}^n (v_i - w_i)^2} \rightarrow \text{norm of } (\vec{v} - \vec{w})$$

• Cosine Similarity

Euclidean distance and cosine similarity,



Use cosine similarity (over euclidean distance) when corpora are in different sizes.

$$\rightarrow \text{Vector norm} \quad \|\vec{v}\| = \sqrt{\sum_{i=1}^n v_i^2}$$

$$\rightarrow \text{Dot product} \quad \vec{v} \cdot \vec{w} = \sum_{i=1}^n v_i \cdot w_i$$

$$\vec{v} \cdot \vec{w} = \|\vec{v}\| \|\vec{w}\| \cos(\beta)$$

$$\cos(\beta) = \frac{\vec{v} \cdot \vec{w}}{\|\vec{v}\| \|\vec{w}\|}$$

- When $\beta = 90^\circ \rightarrow \cos(\beta) = 0 \rightarrow \text{dissimilar/orthogonal}$
- When $\beta = 0^\circ \rightarrow \cos(\beta) = 1 \rightarrow \text{similar/identical}$

• Manipulating words in vector spaces.

USA → Washington DC

Russia → ? ($\text{Russia} + (\text{Washington DC} - \text{USA})$)

→ use known relationships to make predictions.

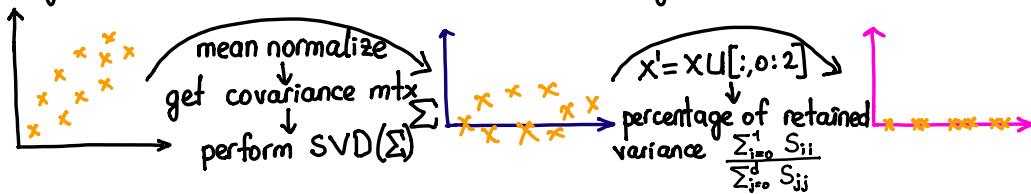
• PCA: Principal Component Analysis

- PCA is an algorithm used for dimensionality reduction
- Original space → Uncorrelated features → Dimensionality reduction
- Visualization to see words relationships in the vector space

• Eigenvectors & eigenvalues

- Eigenvector: uncorrelated features for the data (U)

- Eigenvalue: the amount of information retained by each feature (S)



- Eigenvectors give the direction of uncorrelated features.
- Eigenvalues are the variance of the new features.
- Dot product gives the projection on uncorrelated features.