

DISTRIBUTIONAL SEMANTICS P1

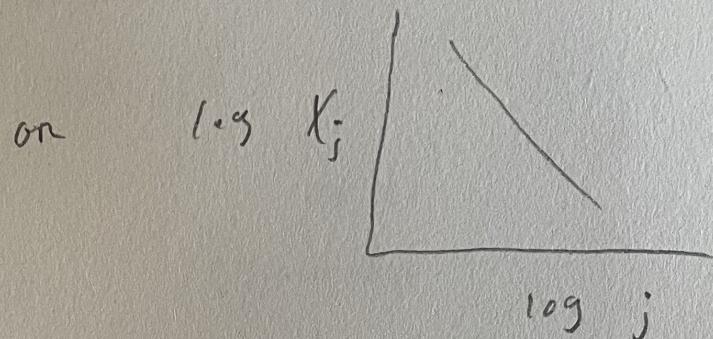
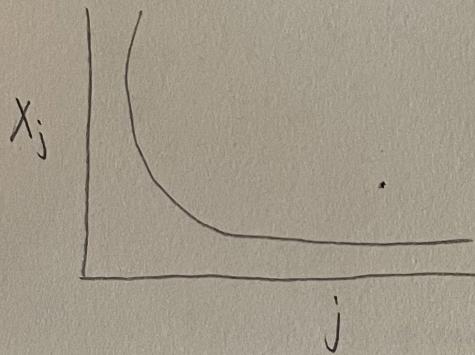
L4-1

- Many problems in NLP require us to assign a measure of similarity, or relevance, between pieces of text:
 - WORDS
 - DOCUMENTS
 - CORPORA
- KEY pieces of technology require this:
 - INFORMATION RETRIEVAL / SEARCH SYSTEMS
 - RECOMMENDER SYSTEMS
 - VIRTUAL ASSISTANTS / CHATBOTS
- = Thus far in this course we have USED WORD frequency (A.K.A. BOW A.K.A. WORD histograms) to represent pieces of text. THIS IS PROBLEMATIC:
 - OUR CORPUS $X \in \mathbb{Z}^{M \times N}$ IS ~~usually~~ SPARSE
 - WE DON'T HAVE A RELIABLE DISTANCE MEASURE IN OUR \mathbb{Z}^N INPUT VECTOR SPACE

- ANOTHER PROBLEM IS THAT $P(X)$ IS
HIGHLY NON-UNIFORM

L4-2

→ Recall ZIPF'S LAW: $P(X_j) = \frac{1/x_j^\alpha}{\sum_{j=1}^n 1/j^\alpha}$ $\alpha > 1$



→ IN MODULE #1 we Learned that

"Self information" is proportional to the inverse log of $P(X)$: $I(X) = -\log P(X)$

- High frequency words carry low information
- Low frequency words carry high information

→ MANY EARLY IR SYSTEMS were DESIGNED TO BE SENSITIVE TO LOW FREQUENCY WORDS, AND INSENSITIVE TO HIGH FREQUENCY WORDS.

→ THIS CAN BE ACCOMPLISHED BY WEIGHING THE WORD COUNTS IN X ACCORDING TO FREQUENCY

TF-IDF ALGORITHM

- TF-IDF MAPS FREQUENCY COUNTS IN X onto AN ASSOCIATED SET OF "WEIGHTS" ACCORDING TO:

$$W_{ij} = TF(X)_{ij} \cdot IDF(X)_j$$

- $TF :=$ "TERM-frequency"

$$TF(X)_{ij} = X_j^{(i)}$$

- $IDF :=$ "INVERSE DOCUMENT frequency"

$$IDF(X)_j = \left[\sum_{i=1}^M 1 \{ X_j^{(i)} > 0 \} \right]^{-1}$$

- MANY VARIANTS OF THE ABOVE. THE MOST WIDELY USED IS:

$$W_{ij} = (1 + \log X_j^{(i)}) \cdot \log \left(\frac{M}{\sum_{i=1}^M 1 \{ X_j^{(i)} > 0 \}} \right)$$

- DOCUMENT SEARCH USING TF-IDF

- GIVEN A TF-IDF WEIGHTED CORPUS MATRIX,
 $W \in \mathbb{R}^{M \times N}$, AND A QUERY, $q \in \mathbb{Z}^N$, OF WORD FREQUENCY COUNTS, FIND THE MOST RELEVANT DOCUMENT IN X . (i.e., THE MOST RELEVANT ROW OF X)

- LET THE SIMILARITY BETWEEN DOCUMENT i AND OUR QUERY BE $S_i(q)$

$$S_i(q) = \sum_{j=1}^N I\{q_j > 0\} \cdot W_{ij}$$

- THE MOST RELEVANT DOCUMENT IS THEN SELECTED ACCORDING TO:

$$\hat{i} = \underset{i \in \{1, \dots, M\}}{\operatorname{ARGMAX}} S_i(q)$$

POINTWISE MUTUAL INFORMATION (PMI)

L4-6

- CAN BE USED IN THE SAME WAY THAT
WE USE TF-IDF FOR IR

$$\begin{aligned}
 & \text{doc} \quad \text{word} \\
 - \quad \text{PMI}(i, j) &= \log \frac{P(i, j)}{P(i) P(j)} \\
 &= \log \frac{P(j|i) P(i)}{P(i) P(j)} \\
 &= \log \frac{P(j|i)}{P(j)} = \log P(j|i) - \log P(j) \\
 &= \log \frac{x_j^{(i)}}{\sum_{j'=1}^N x_{j'}^{(i)}} - \log \frac{\sum_{i'=1}^M x_i^{(i')}}{\sum_{i'=1}^M \sum_{j'=1}^N x_{j'}^{(i')}} \\
 &= \log \left(\frac{\text{Count}(i, j)}{\text{Count}(i, \text{All } j)} \right) - \log \left(\frac{\text{Count}(\text{All } i, j)}{\text{Count}(\text{All } i, \text{All } j)} \right)
 \end{aligned}$$

- POSITIVE PMI (PPMI)

$$\text{PPMI}(i, j) = \begin{cases} \text{PMI}(i, j), & P(j|i) > P(j) \\ 0, & \text{ELSE} \end{cases}$$

- TF-IDF / PMI ARE USED TO ADDRESS PROBLEMS RELATED TO ZIPF'S LAW (i.e., THE NON-UNIFORMITY IN $P(x)$)
- THESE METHODS DO NOT CHANGE THE SPARSITY OF OUR REPRESENTATION, THOUGH!

- OUR GOAL IS TO MAP W ONTO A "DENSE" REPRESENTATION: $W \in \mathbb{R}^N \rightarrow U \in \mathbb{R}^K$
 $K \ll N$

~~$X \in \mathbb{R}^N \rightarrow \mathbb{R}^K$~~

- THERE ARE TWO MATRIX FACTORIZATION METHODS THAT DO A GOOD JOB AT THIS IN PRACTICE

- TRUNCATED SVD

- NON-NEGATIVE MATRIX FACTORIZATION

LATENT SEMANTIC ANALYSIS (LSA)

L4-8

- A.K.A. LATENT SEMANTIC INDEXING (LSI)
- USES TRUNCATED SVD
- WE START WITH A MATRIX of WEIGHTED Frequency Counts, $W \in \mathbb{R}^{M \times N}$
- RECALL THE TRUNCATED SVD:

$$W = U \Sigma V^T$$

$M \times N$ $M \times K$ $K \times K$ $K \times N$

↳ K IS OFTEN CHOSEN USING

- ELBOW METHOD Performance on Σ from Full SVD
- JUST $K \ll N$
- NUMERICAL PACKAGES IN Python (Scilab, Scipy) CAN DO THIS FOR YOU!
- OUR COMPRESSED CORPUS IS GIVEN BY :

$$U = W V \Sigma^{-1} \rightarrow \text{RECALL THAT}$$

V IS ORTHOGONAL

Σ IS A DIAGONAL MATRIX

DOCUMENT SEARCH USING LSA

L4-9

- Using the Problem STATEMENT From

L4-5 :

- FIRST Compute the WEIGHTED Query VECTORS : TF-IDF / PMF for EXAMPLE

$$\rightarrow q \Rightarrow q^{(w)}$$

- Project $q^{(w)} \in \mathbb{R}^N$ to $q'^{(w)} \in \mathbb{R}^K$

$$q'^{(w)} = q^{(w)} V \Sigma^{-1}$$

- USE COSINE-SIMILARITY TO MEASURE THE RELEVANCE of document $x^{(i)}$ to q

$$S_i(q) = \cos(w^{(i)}, q'^{(w)})$$

$$= \frac{\sum_{j=1}^K q_j'^{(w)} w_j^{(i)}}{\sqrt{\sum_{j=1}^K q_j'^{(w)2}} \sqrt{\sum_{j=1}^K w_j^{(i)2}}}$$

- Document SELECTION:

$$\hat{i} = \underset{i \in \{1, \dots, M\}}{\operatorname{Argmax}} S_i(q'^{(w)})$$

NON-NEGATIVE MATRIX FACTORIZATION (NMF) L4-10

- NOTE THAT IN LSA, THE TRANSFORMED VECTOR SPACE IN \mathbb{R}^K IS DEFINED BY THE K-COLUMNS OF V . IN OTHER WORDS, EACH COLUMN OF V REPRESENTS A BASIS VECTOR IN OUR NEW \mathbb{R}^K SPACE. IMPORTANTLY, EACH OF THESE BASIS VECTORS IS A LINEAR COMBINATION OF THE N-BASIS VECTORS IN THE ORIGINAL \mathbb{R}^N SPACE. WITH LSA, THE COEFFICIENTS THAT GET APPLIED TO EACH ORIGINAL DIMENSION (RE: WORD) CAN BE NEGATIVE, AND THIS IS UNDESIRABLE.

EXAMPLE: $V_j^T = (0.5)\text{swim} + (0.7)\text{sun} + (-1.0)\text{tarn} + (-0.2)\text{dark}$

COULD EQUIVALENTLY BE REPRESENTED AS:

$$(0.5)\text{swim} + (1.4)\text{sun} + (0.01)\text{tarn} + (0.01)\text{dark}$$

- NMF FACTORS X USING A POSITIVITY CONSTRAINT ON THESE COEFFICIENTS THAT GET APPLIED TO THE WORDS.

NMF Cont...

L4-11

- NMF FACTORIZATION typically is performed ON A WEIGHTED version of X , similar to LSA

$$W = UV \quad \Rightarrow \quad \underbrace{V}_{M \times K} \text{ is ORTHOGONAL} \quad \underbrace{U}_{K \times N}$$

- OBJECTIVE Function:

$$\hat{U}, \hat{V} = \min_{U, V \geq 0} L(W; U, V) + \alpha p(\|U\|_F + \|V\|_F) + \alpha(1-p)(\|U\|_F^2 + \|V\|_F^2)$$

α = REGULARIZATION CONSTANT

p = SPARSITY CONSTANT

L = SOME OBJECTIVE FUNCTION

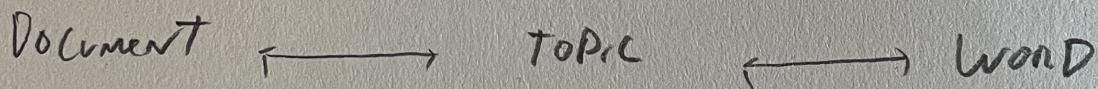
FROBENIUS LOSS: $L_F = \|W - UV\|_F^2$

KL-DIVERGENCE LOSS: $L_{KL} = \sum_{i=1}^M \sum_{j=1}^N W_{ij}^{(i)} \left[\log \frac{W_{ij}^{(i)}}{(UV)_{ij}} - 1 \right] + (UV)_{ij}$

TOPIC MODELING

L4-12

- Thus far we have primarily been concerned with modeling the relationship between documents and their constituent words.
- Topic modeling is a set of methods that involve an "intermediate" representation that separates documents from words. The intermediary is referred to as "Topics"



0

0

0

0

0

0

0

0

0

0

0

0

- WE HAVE ALREADY LEARNED two such topic models!

- LSA : $X = U \Sigma V^T \rightarrow$ "Topics" = columns of V

- NMF : $X = UV \rightarrow$ "Topics" = rows of V