# *Distributional Models of Semantics*

Pawan Goyal

CSE, IIT Kharagpur

Week 7, Lecture 2

Words are treated as atomic symbols

# Vector Space Model without distributional similarity

Words are treated as atomic symbols

## One-hot representation

```
motel [0 0 0 0 0 0 0 0 0 0 1 0 0 0 0]  AND
hotel [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0] = 0
```

# Distributional Similarity Based Representations

*You know a word by the company it keeps*

# Distributional Similarity Based Representations

*You know a word by the company it keeps*

government debt problems turning into banking crises as has happened in

saying that Europe needs unified banking regulation to replace the hodgepodge

# Distributional Similarity Based Representations

You know a word by the company it keeps

government debt problems turning into banking crises as has happened in

saying that Europe needs unified banking regulation to replace the hodgepodge

These words will represent banking

## *Building a DSM step-by-step*

### *The "linguistic" steps*

Pre-process a corpus (to define targets and contexts)

⇓

Select the targets and the contexts

# Building a DSM step-by-step

## The "linguistic" steps

Pre-process a corpus (to define targets and contexts)

⇓

Select the targets and the contexts

## The "mathematical" steps

Count the target-context co-occurrences

⇓

Weight the contexts (optional)

⇓

Build the distributional matrix

⇓

Reduce the matrix dimensions (optional)

⇓

Compute the vector distances on the (reduced) matrix

# Many design choices

| Matrix type | Weighting | Dimensionality reduction | Vector comparison |
|---|---|---|---|
| word × document | probabilities | LSA | Euclidean |
| word × word | length normalization | PLSA | Cosine |
| word × search proximity | TF-IDF | LDA | Dice |
| adj. × modified noun | PMI | PCA | Jaccard |
| word × dependency rel. | Positive PMI | IS | KL |
| verb × arguments | PPMI with discounting | DCA | KL with skew |
| ⋮ | ⋮ | ⋮ | ⋮ |

# *Many design choices*

| Matrix type | | Weighting | | Dimensionality reduction | | Vector comparison |
| --- | --- | --- | --- | --- | --- | --- |
| word × document | | probabilities | | LSA | | Euclidean |
| word × word | | length normalization | | PLSA | | Cosine |
| word × search proximity | × | TF-IDF | × | LDA | × | Dice |
| adj. × modified noun | | PMI | | PCA | | Jaccard |
| word × dependency rel. | | Positive PMI | | IS | | KL |
| verb × arguments | | PPMI with discounting | | DCA | | KL with skew |
| ⋮ | | ⋮ | | ⋮ | | ⋮ |

### *General Questions*

- How do the rows (words, ...) relate to each other?

- How do the columns (contexts, documents, ...) relate to each other?

*A number of parameters to be fixed*

- Which type of context?
- Which weighting scheme?
- Which similarity measure?
- ...

A specific parameter setting determines a particular type of DSM (e.g. LSA, HAL, etc.)

# Documents as context: Word × document

|         | d1 | d2 | d3 | d4 | d5 | d6 | d7 | d8 | d9 | d10 |
|---------|----|----|----|----|----|----|----|----|----|-----|
| against | 0  | 0  | 0  | 1  | 0  | 0  | 3  | 2  | 3  | 0   |
| age     | 0  | 0  | 0  | 1  | 0  | 3  | 1  | 0  | 4  | 0   |
| agent   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   |
| ages    | 0  | 0  | 0  | 0  | 0  | 2  | 0  | 0  | 0  | 0   |
| ago     | 0  | 0  | 0  | 2  | 0  | 0  | 0  | 0  | 3  | 0   |
| agree   | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   |
| ahead   | 0  | 0  | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0   |
| ain't   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   |
| air     | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   |
| aka     | 0  | 0  | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0   |

# Words as context: Word × Word

|         | against | age  | agent | ages | ago  | agree | ahead | ain.t | air | aka | al  |
|---------|---------|------|-------|------|------|-------|-------|-------|-----|-----|-----|
| against | 2003    | 90   | 39    | 20   | 88   | 57    | 33    | 15    | 58  | 22  | 24  |
| age     | 90      | 1492 | 14    | 39   | 71   | 38    | 12    | 4     | 18  | 4   | 39  |
| agent   | 39      | 14   | 507   | 2    | 21   | 5     | 10    | 3     | 9   | 8   | 25  |
| ages    | 20      | 39   | 2     | 290  | 32   | 5     | 4     | 3     | 6   | 1   | 6   |
| ago     | 88      | 71   | 21    | 32   | 1164 | 37    | 25    | 11    | 34  | 11  | 38  |
| agree   | 57      | 38   | 5     | 5    | 37   | 627   | 12    | 2     | 16  | 19  | 14  |
| ahead   | 33      | 12   | 10    | 4    | 25   | 12    | 429   | 4     | 12  | 10  | 7   |
| ain't   | 15      | 4    | 3     | 3    | 11   | 2     | 4     | 166   | 0   | 3   | 3   |
| air     | 58      | 18   | 9     | 6    | 34   | 16    | 12    | 0     | 746 | 5   | 11  |
| aka     | 22      | 4    | 8     | 1    | 11   | 19    | 10    | 3     | 5   | 261 | 9   |
| al      | 24      | 39   | 25    | 6    | 38   | 14    | 7     | 3     | 11  | 9   | 861 |

# Words as contexts

### Parameters

- Window size
- Window shape - rectangular/triangular/other

# Words as contexts

### Parameters

- Window size
- Window shape - rectangular/triangular/other

### Consider the following passage

*Suspected communist rebels on 4 July 1989 killed Col. Herminio Taylo, police chief of Makati, the Philippines major financial center, in an escalation of street violence sweeping the Capitol area. The gunmen shouted references to the rebel New People's Army. They fled in a commandeered passenger jeep. The military says communist rebels have killed up to 65 soldiers and police in the Capitol region since January.*

# Words as contexts

## Parameters

- Window size
- Window shape - rectangular/triangular/other

## 5 words window (unfiltered): 2 words either side of the target word

*Suspected communist rebels on 4 July 1989 killed Col. Herminio Taylo, police chief of Makati, the Philippines major financial center, in an escalation of street violence sweeping the Capitol area. The gunmen shouted references to the rebel New People's Army. They fled in a commandeered passenger jeep. The military says communist rebels have killed up to 65 soldiers and police in the Capitol region since January.*

# Words as contexts

## Parameters

- Window size
- Window shape - rectangular/triangular/other

## 5 words window (filtered): 2 words either side of the target word

Suspected communist rebels on 4 July 1989 killed Col. Herminio Taylo, police chief of Makati, the Philippines major financial center, in an escalation of street violence sweeping the Capitol area. The gunmen shouted references to the rebel New People's Army. They fled in a commandeered passenger jeep. The military says communist rebels have killed up to 65 soldiers and police in the Capitol region since January.

# Context weighting: documents as context

## Indexing function F: Essential factors

- **Word frequency ($f_{ij}$):** How many times a word appears in the document? $F \propto f_{ij}$
- **Document length ($|D_i|$):** How many words appear in the document? $F \propto \frac{1}{|D_i|}$
- **Document frequency ($N_j$):** Number of documents in which a word appears. $F \propto \frac{1}{N_j}$

# Context weighting: documents as context

## Indexing function $F$: Essential factors

- **Word frequency ($f_{ij}$):** How many times a word appears in the document? $F \propto f_{ij}$
- **Document length ($|D_i|$):** How many words appear in the document? $F \propto \frac{1}{|D_i|}$
- **Document frequency ($N_j$):** Number of documents in which a word appears. $F \propto \frac{1}{N_j}$

## Indexing Weight: tf-Idf

- $f_{ij} * log(\frac{N}{N_j})$ for each term, normalize the weight in a document with respect to $L_2$-norm.

# Context weighting: words as context

## basic intuition

| word1 | word2 | freq(1,2) | freq(1) | freq(2) |
|-------|-------|-----------|---------|---------|
| dog | small | 855 | 33,338 | 490,580 |
| dog | domesticated | 29 | 33,338 | 918 |

# *Context weighting: words as context*

### *basic intuition*

| word1 | word2 | freq(1,2) | freq(1) | freq(2) |
|-------|-------|-----------|---------|---------|
| dog | small | 855 | 33,338 | 490,580 |
| dog | domesticated | 29 | 33,338 | 918 |

**Association measures** are used to give more weight to contexts that are more significantly associed with a targer word.

# *Context weighting: words as context*

## *basic intuition*

| word1 | word2 | freq(1,2) | freq(1) | freq(2) |
|-------|-------|-----------|---------|---------|
| dog | small | 855 | 33,338 | 490,580 |
| dog | domesticated | 29 | 33,338 | 918 |

**Association measures** are used to give more weight to contexts that are more significantly associed with a targer word.

- The less frequent the target and context element are, the higher the weight given to their co-occurrence count should be.

# Context weighting: words as context

### basic intuition

| word1 | word2 | freq(1,2) | freq(1) | freq(2) |
|-------|-------|-----------|---------|---------|
| dog | small | 855 | 33,338 | 490,580 |
| dog | domesticated | 29 | 33,338 | 918 |

**Association measures** are used to give more weight to contexts that are more significantly associated with a targer word.

- The less frequent the target and context element are, the higher the weight given to their co-occurrence count should be.
  $\Rightarrow$ Co-occurrence with frequent context element *small* is less informative than co-occurrence with rarer *domesticated*.

## *Context weighting: words as context*

**Association measures** are used to give more weight to contexts that are more significantly associoted with a targer word.

- The less frequent the target and context element are, the higher the weight given to their co-occurrence count should be.
  $\Rightarrow$ Co-occurrence with frequent context element *small* is less informative than co-occurrence with rarer *domesticated*.
- different measures - e.g., Mutual information, Log-likelihood ratio

$$PMI(w_1, w_2) = log_2 \frac{P_{corpus}(w_1, w_2)}{P_{ind}(w_1, w_2)}$$

$$PMI(w_1, w_2) = log_2 \frac{P_{corpus}(w_1, w_2)}{P_{ind}(w_1, w_2)}$$

$$PMI(w_1, w_2) = log_2 \frac{P_{corpus}(w_1, w_2)}{P_{corpus}(w_1)P_{corpus}(w_2)}$$

# Pointwise Mutual Information (PMI)

$$PMI(w_1, w_2) = log_2 \frac{P_{corpus}(w_1, w_2)}{P_{ind}(w_1, w_2)}$$

$$PMI(w_1, w_2) = log_2 \frac{P_{corpus}(w_1, w_2)}{P_{corpus}(w_1) P_{corpus}(w_2)}$$

$$P_{corpus}(w_1, w_2) = \frac{freq(w_1, w_2)}{N}$$

$$P_{corpus}(w) = \frac{freq(w)}{N}$$

## *Positive PMI*

All PMI values less than zero are replaced with zero.

# PMI: Issues and Variations

### Positive PMI
All PMI values less than zero are replaced with zero.

### Bias towards infrequent events
Consider $w_j$ having the maximum association with $w_i$,
$$P_{corpus}(w_i) \approx P_{corpus}(w_j) \approx P_{corpus}(w_i, w_j)$$

# PMI: Issues and Variations

## Positive PMI

All PMI values less than zero are replaced with zero.

## Bias towards infrequent events

Consider $w_j$ having the maximum association with $w_i$,
$P_{corpus}(w_i) \approx P_{corpus}(w_j) \approx P_{corpus}(w_i, w_j)$
PMI increases as the probability of $w_i$ decreases.

# PMI: Issues and Variations

### Positive PMI
All PMI values less than zero are replaced with zero.

### Bias towards infrequent events
Consider $w_j$ having the maximum association with $w_i$,
$P_{corpus}(w_i) \approx P_{corpus}(w_j) \approx P_{corpus}(w_i, w_j)$
PMI increases as the probability of $w_i$ decreases.
Also, consider a word $w_j$ that occurs once in the corpus, also in the context of $w_i$.

# PMI: Issues and Variations

## Positive PMI
All PMI values less than zero are replaced with zero.

## Bias towards infrequent events
Consider $w_j$ having the maximum association with $w_i$,
$P_{corpus}(w_i) \approx P_{corpus}(w_j) \approx P_{corpus}(w_i, w_j)$
PMI increases as the probability of $w_i$ decreases.
Also, consider a word $w_j$ that occurs once in the corpus, also in the context of $w_i$. A discounting factor proposed by Pantel and Lin:

$$\delta_{ij} = \frac{f_{ij}}{f_{ij}+1} \frac{min(f_i, f_j)}{min(f_i, f_j)+1}$$

$PMI_{new}(w_i, w_j) = \delta_{ij} PMI(w_i, w_j)$

# Distributional Vectors: Example

## Normalized Distributional Vectors using Pointwise Mutual Information

| | |
|---|---|
| **petroleum** | oil:0.032 gas:0.029 crude:0.029 barrels:0.028 exploration:0.027 barrel:0.026 opec:0.026 refining:0.026 gasoline:0.026 fuel:0.025 natural:0.025 exporting:0.025 |
| **drug** | trafficking:0.029 cocaine:0.028 narcotics:0.027 fda:0.026 police:0.026 abuse:0.026 marijuana:0.025 crime:0.025 colombian:0.025 arrested:0.025 addicts:0.024 |
| **insurance** | insurers:0.028 premiums:0.028 lloyds:0.026 reinsurance:0.026 underwriting:0.025 pension:0.025 mortgage:0.025 credit:0.025 investors:0.024 claims:0.024 benefits:0.024 |
| **forest** | timber:0.028 trees:0.027 land:0.027 forestry:0.026 environmental:0.026 species:0.026 wildlife:0.026 habitat:0.025 tree:0.025 mountain:0.025 river:0.025 lake:0.025 |
| **robotics** | robots:0.032 automation:0.029 technology:0.028 engineering:0.026 systems:0.026 sensors:0.025 welding:0.025 computer:0.025 manufacturing:0.025 automated:0.025 |