# Natural Language Processing

**Lecture 07**

Dirk Hovy

dirk.hovy@unibocconi.it

@dirk_hovy

Bocconi

# Today's Goals

- Understand the difference between **sparse** and **dense** representations

- Learn about **word2vec** and **doc2vec**

- **Understand** the underlying algorithms

Bocconi

# Dense Distributed Representations

Bocconi

# Distributional Hypothesis

*"You shall know the meaning of a word by the company it keeps"*

Firth (1957)

Similar words have similar **contexts**

Represent **words** as **vectors**/points in space

Similar words have similar vectors

Bocconi

# An Example

# Latent Semantic Analysis

# Part 1
# Representing Words
# as Vectors

# Semantic Similarity

# Similarity Measures

**Cosine similarity**

$$\frac{A \cdot B}{\|A\|\|B\|}$$

**x**

*flats ≈ apartments*

0.84

0.13

*platypus*

−1.0

**y**

# Dot Product

- "combine" vectors to a scalar

*SUM*

$$x \cdot y = \sum_{i=1}^{D} x_i y_i$$

*MULTIPLY*

$$\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 6 \end{bmatrix} = \begin{matrix} 1 \\ 4 \\ 3 \end{matrix}$$

Bocconi

# Vector Norm

- add up square of each element, take $\sqrt{\ }$

$$\begin{bmatrix} 2 \\ 6 \end{bmatrix} \qquad = \sqrt{2^2 + 6^2} \quad = 6.324$$

Bocconi

# Nearest neighbors

# Word2Vec – Intuitively

```
place all words randomly on fridge

for each pair of words:
    if in same sentence:
        move closer together
    else:
        move further apart
```

Bocconi

house

wash

dog

one

very

always

never

monday

weekend

door

buy

two

billions

tuesday

Bocconi

two

one

billions

never

very

always

wash buy

house

weekend

tuesday

door

dog

monday

X-POS    Y-POS

two
one
billions
never
very
always
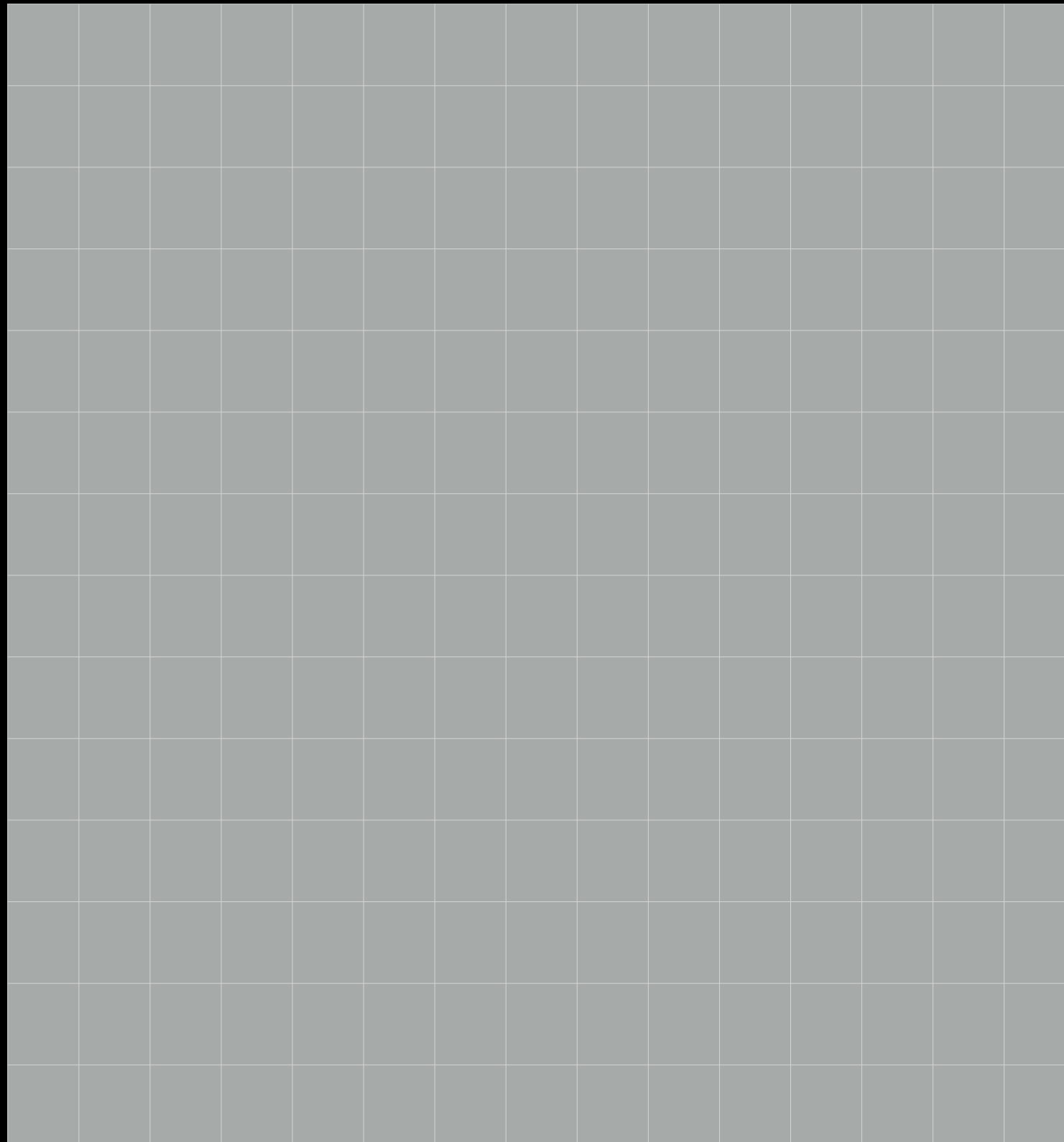wash
buy
house
door
dog
weekend
monday
tuesday

VECTORS

16

Bocconi

# Word2Vec – CBOW Model

**[16.3, 7.8]**
*flat*

*MATRIX OF*

*OUTPUT*  `garden`

*TARGET WORDS*

*ERROR*
*BACKPROPAGATION*  *SUM*

*INPUT*

*MATRIX OF*

rent    Renting our large apartment in great location

*CONTEXT WORDS*

Bocconi

OUTPUT

INPUT

rent      Renting our large apartment in great location

Bocconi

# Nuts and Bolts
# and
# Engineering Tricks

Bocconi

# Problem?

- We are trying to learn a conditional probability distribution over the vocabulary *for each word in the vocabulary*:

$$P(w_{out} | w_{in})$$

- With a large vocabulary comes large trouble…

**aardvark**

**...**

**Zzzyx**

$w_{out}$

**Bocconi**

# Trick 1: Negative Sampling

Sample small set of words, labeled as 0 (not a context word) or 1 (is a context word)



CHECK AGAINST

TRUE ANSWER

SIGMOID — 0.8 ← 0

0.2 — DOT PRODUCT

UPDATE

CONTEXT WORD   TARGET WORD

Bocconi

# Trick 2: Sub-Sampling

Sample a word:

```
      the
      the
        a
      the
      the
       in
      the
        a
      the
        a
 platypus
```

*SOLUTION:*

*REMOVE WORDS IN THE INPUT SENTENCE PROPORTIONAL TO THEIR FREQUENCY*

50000
40000
30000
20000
10000
0

Bocconi

# Trick 3: Hierarchical Softmax

Update to regular softmax: $O(|V|)$

Hierarchical softmax: $O(log|V|)$



*RANDOM WALK ALONG A TREE*

N0

0.65     0.35

N1

0.54     0.46     0.29     0.71

N2

0.71   0.29   0.83   0.17   1   0   0.8   0.2

horse   fridge   time   it   and   potato   zebra   white

$$P(\texttt{time}\,|\,C) = P_{N0}(right\,|\,C) \cdot P_{N1}(left\,|\,C) \cdot P_{N2}(right\,|\,C) = 0.25$$

source: http://building-babylon.net/2017/08/01/hierarchical-softmax/

Bocconi

# Vector Space Semantics

***king – man + woman ≈ queen***

# Caveat: Antonyms

*His kitchen was always very* _____



x

**clean**

**dirty**

*SAME CONTEXT, OPPOSITE MEANING!*

y

Bocconi

# Debiasing Vectors

# Part 2
# Representing Documents as Vectors

**Bocconi**

# Example 1: Songs

**Billboard HOT 100**

song 1

$\vdots$

song $n$

$C$

**Hiphop songs**

**Country songs**

# Example 2: Cities

jodel

city 1

$C$

city $n$

# Doc2Vec – Intuitively

```
place words & cities randomly on fridge

for each pair of (word, city):

    if word seen in city:

        move closer together

    else:

        move further apart
```

# Adding Labels

I.E., CITIES, REGIONS, PEOPLE, ...

NUMBERS two

one

billions

never

very

always

ADVERBS

VERBS

wash buy

TIMES

house weekend

tuesday

door NOUNS

dog

monday DAYS

Bocconi

# Words and Documents

# Preview:
# Better, Contextualized
# Document Embeddings

Bocconi

# Contextual Representations

I had to stay home **sick**

Totally **sick** move, bro

*pill*

*nurse*

***sick***

*great*

*awesome*

Bocconi

# Encoding Words

| 0.001 | 0.001 | ... | 0.13 | ... | 0.001 |
|-------|-------|-----|------|-----|-------|

*VOCABULARY*  aardvark  Aarhus  ...  like  ...  Zzzyx

**FeedForward w/ Softmax**

**BERT**

Encoder #12 or 24

...

Encoder #2

Encoder #1

| 1 | 2 | 3 | 4 | 5 | 512 |
|---|---|---|---|---|-----|
| [CLS] | Mice | [MASK] | cheese | –PAD– | –PAD– |
| *START* | | *MASK* | | *PADDING* | |

37

# Wrapping up…

**Bocconi**

# Representation Comparison

| | Discrete | Distributed |
|---|---|---|
| **#Dimensions** | Data-dependent | Pre-defined |
| **Content** | Count-based | Coefficients |
| **Density** | Sparse | Dense |
| **Strength** | Interpretability | Similarity |
| **Application** | Understanding | Performance |
| **School of thought** | Rationalism | Empiricism |

Bocconi

# Take home points

- Text can be represented as dense, continuous embedding vectors

- Embedding models learn similarity via co-occurrence

- Word and document embeddings reflect semantic similarity in high-dimensional space

- Good for similarity, visualization, and classification, bad for analysis

Bocconi