# Natural Language Processing

**Text Classification**

Dirk Hovy

dirk.hovy@unibocconi.it

@dirk_hovy

# Text is an exploding data source
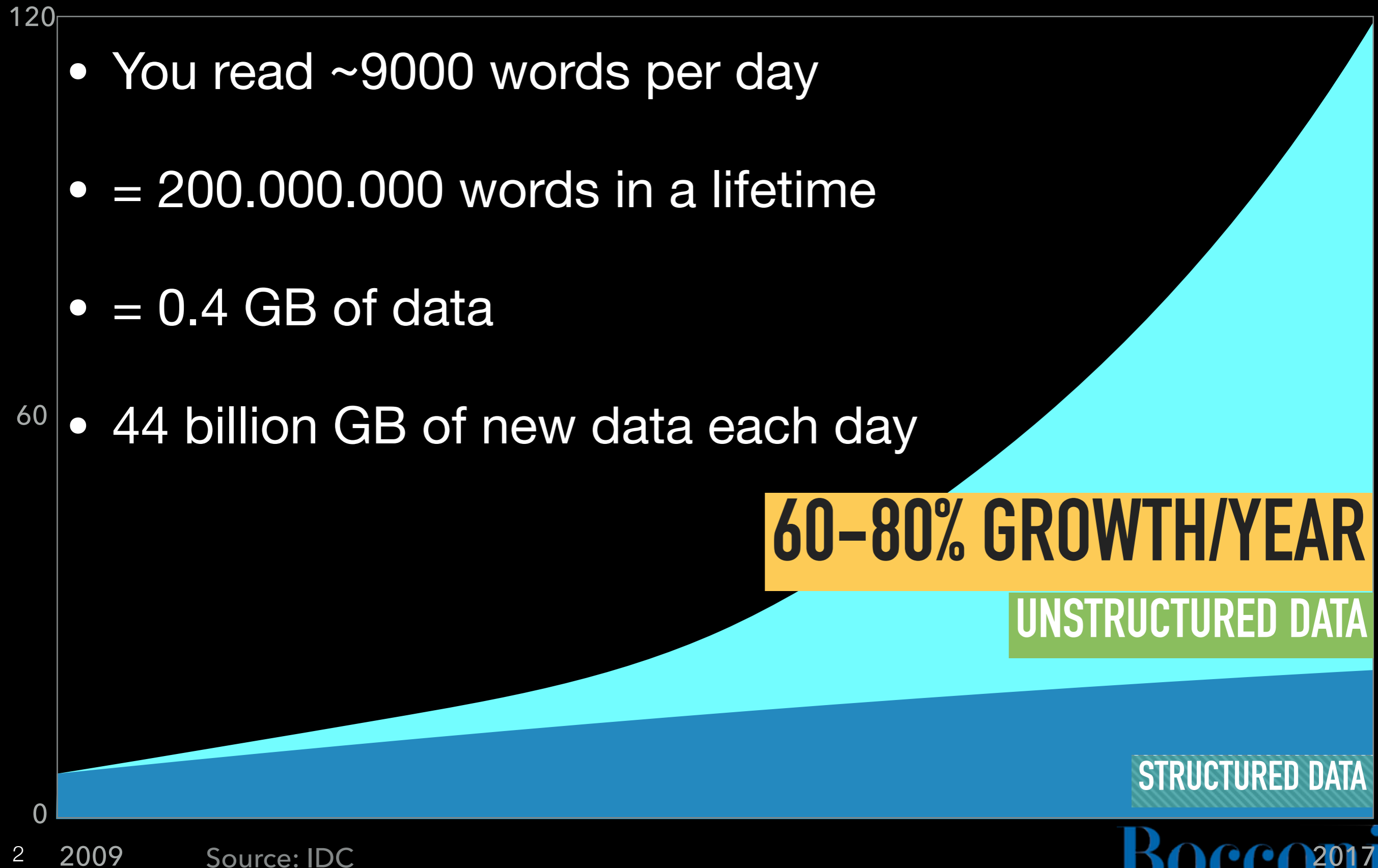
- You read ~9000 words per day

- = 200.000.000 words in a lifetime

- = 0.4 GB of data

- 44 billion GB of new data each day

**60–80% GROWTH/YEAR**

**UNSTRUCTURED DATA**

**STRUCTURED DATA**

120

60

0

2009

Source: IDC

2017

2

**Bocconi**

# NLP is booming



$5.400.000.000

$136.000.000

2016  2017  2018  2019  2020  2021  2022  2023  2024  2025
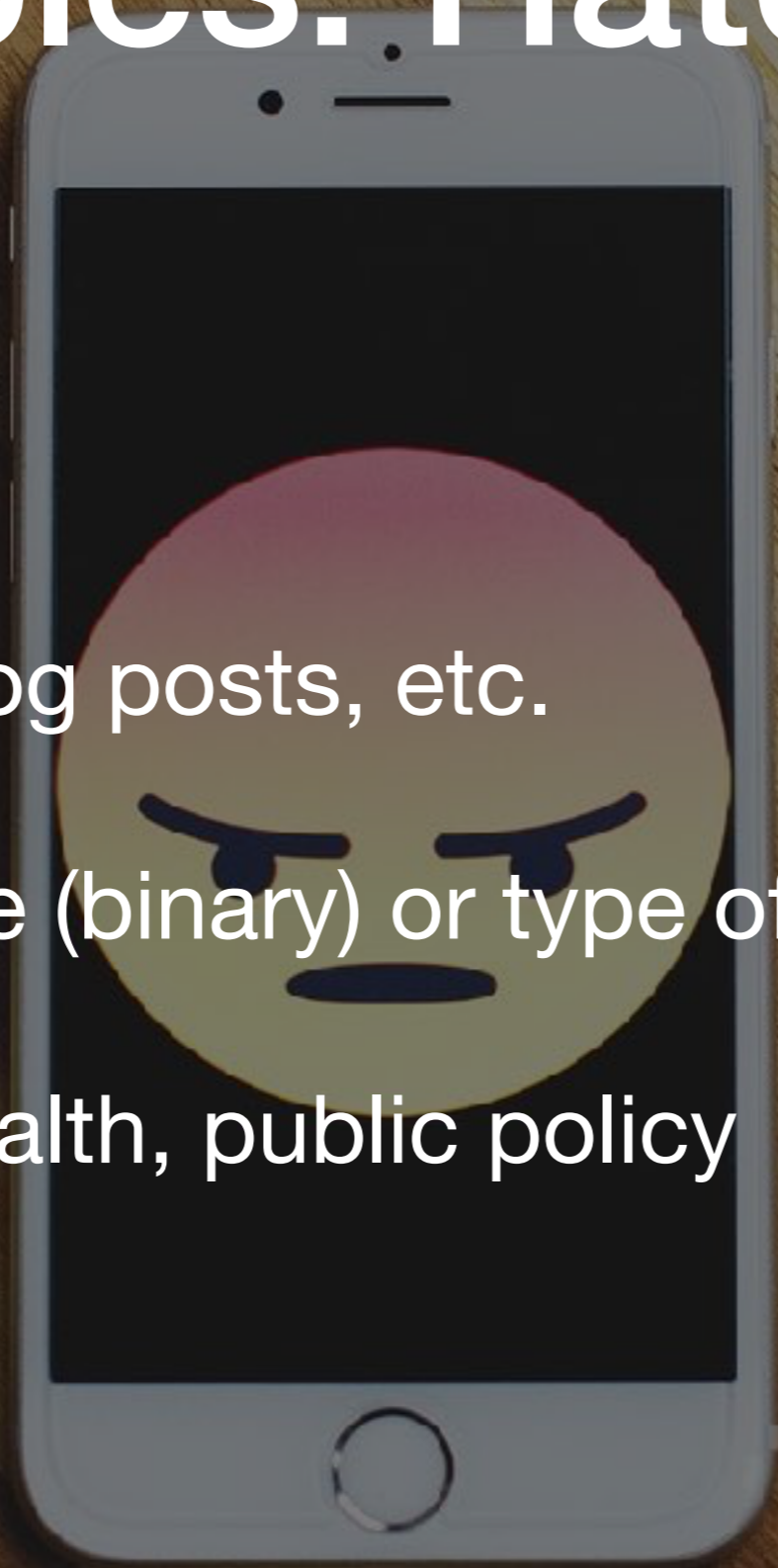
Bocconi

# Examples: Sentiment

- Input: reviews

- Output: positive, negative, neutral

- Use: business intelligence, market analysis

# Examples: Hate Speech

- Input: tweets, blog posts, etc.

- Output: presence (binary) or type of hate speech

- Use: platform health, public policy

Bocconi

# Examples: Mental Health

- Input: social media

- Output: presence of risk for mental health condition

- Use: psychologist support, risk screening

6

Bocconi

# Examples: Geolocation

## Author Attribute Prediction

- Input: tweet history

- Output: coordinates or predefined region

- Use: social media analysis, targeting

Bocconi

# Sentiment Analysis

**Bocconi**

# Classification Steps

- **preprocess** the data

- choose **text representation** (discrete or continuous)

- **select** a **model** (CV, metrics, regularization)

- **fit** the final **model**

Bocconi

# Let's start!

# Today's Goals

- Understand where NLP comes from

- Learn about the different steps of **preprocessing**

- Learn about **bag of words** (BOW) representations

- Learn about forms of **TF-IDF** and its possibilities

- Understand the difference between sparse and dense representations

- Learn about word2vec and doc2vec

Bocconi

# Pre-processing

# Pre-processing steps

```
<div id="text">I've been in New York
in 2011, but didn't like it. I
preferred Los Angeles.</div>
```

## GOAL: MINIMIZE VARIATION

Bocconi

# Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

  - numbers

  - lemmas vs. stems

- Remove unwanted words

  - stopwords

  - content words (use POS tagging!)

- join collocations

I've been in New York in 2011, but didn't like it. I preferred Los Angeles.

14

Bocconi

# Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

  - numbers

  - lemmas vs. stems

- Remove unwanted words

  - stopwords

  - content words (use POS tagging!)

- join collocations

I've been in New York in 2011, but didn't like it.

I preferred Los Angeles.

**Bocconi**

# Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

  - numbers

  - lemmas vs. stems

- Remove unwanted words

  - stopwords

  - content words (use POS tagging!)

- join collocations

```
I 've been in New York
in 2011 , but did n't
like it .

I preferred Los
Angeles .
```

Bocconi

# Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

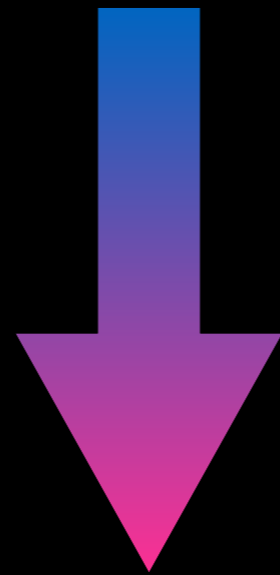  - numbers

  - lemmas vs. stems

- Remove unwanted words

  - stopwords

  - content words (use POS tagging!)

- join collocations

```
i 've been in new york
in 0000 , but did n't
like it .

i preferred los
angeles .
```

**Bocconi**

# Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

  - numbers

  - lemmas vs. stems

- Remove unwanted words

  - stopwords

  - content words (use POS tagging!)

- join collocations

```
i have be in new york in
0000 , but do not like
it .

i prefer los angeles .
```

**Bocconi**

# Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

  - numbers

  - lemmas vs. stems

- Remove unwanted words

  - stopwords

  - content words (use POS tagging!)

- join collocations

```
i new york 0000 , like .

i prefer los angeles .
```

Bocconi

# Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

  - numbers

  - lemmas vs. stems

- Remove unwanted words

  - stopwords

  - content words (use POS tagging!)

- join collocations

```
new york 0000 like

prefer los angeles
```

*CONTENT = (NOUN, VERB, NUM)*

Bocconi

# Pre-processing steps

- Remove formatting (e.g. HTML)

- Segment sentences

- Tokenize words

- Normalize words

  - numbers

  - lemmas vs. stems

- Remove unwanted words

  - stopwords

  - content words (use POS tagging!)

- join collocations

```
new_york 0000 like

prefer los_angeles
```

21

**Bocconi**

# Pre-processing steps

`<div id="text">I've been in New York in 2011, but didn't like it. I preferred Los Angeles.</div>`

*MINIMIMAL VARIATION*

*"BAG OF WORDS"*

new_york 0000 like

prefer los_angeles

Bocconi

# Telling Neighbors: Pointwise Mutual Information

# Some are not like the Others

*Frequent Bigrams*

*Collocations*

frying pan

social media

and thought

New York

of the

Bocconi

# Mutual Informativity

How well can we guess the Blank?

social _____          and _____

_____ media          _____ the

Bocconi

# Pointwise Mutual Information

*CHANCE OF SEEING THEM TOGETHER*

$$PMI(x, y) = log \frac{P(x, y)}{P(x)P(y)}$$

*...SEEING EITHER*

| x | y | c(x) | c(y) | c(xy) | P(x) | P(y) | P(x, y) | PMI(x; y) |
|---|---|---|---|---|---|---|---|---|
| moby | dick | 83 | 83 | 82 | 0.0003 | 0.0003 | 0.0003 | **3.48** |
| captain | ahab | 327 | 511 | 61 | 0.0013 | 0.0020 | 0.0002 | **1.97** |
| white | whale | 280 | 1150 | 106 | 0.0011 | 0.0045 | 0.0004 | **1.93** |
| under | the | 119 | 14175 | 45 | 0.0005 | 0.0553 | 0.0002 | **0.83** |
| is | a | 1690 | 4636 | 110 | 0.0066 | 0.0181 | 0.0004 | **0.56** |

c(X) = 256,149
c(XY) = 256,148

Bocconi

# Representing Text

# Ham or Spam?

From: offr4u@rsph.com
Subject: Unique wealth offerings
To: dirk.hovy@unibocconi.it

---

Greetings dear friend

We have an amazing offer 4U: Click here to get access to a free consultation for serious wealth benefits! Urgent: offer expires soon.
Works guaranteed! Triple your income.

Spam terms:

- 4U

- click

- amazing

- free

- guarantee

- offer

- urgent

- dear friend

- income

- serious

**Bocconi**

# Terminology



Data Matrix

Row or Feature Vector

Instances

Element (Scalar number)

Features or Dimensions

# Discrete Representations

# Bags of words (BOW)



*Count Words*

```
{
    'shakespeare': 6,
    'in': 20,
    'love': 6,
    'is': ...
}
```

*Vectorize Features*

| shakespeare | | in | | ... | | love | beer | | |
|---|---|---|---|---|---|---|---|---|---|
| 6 | ... | 20 | 0 | ... | 0 | 6 | 0 | ... | 0 |

$X$

Bocconi

# Quiz!

What happens if we allow _every_ *possible word* to constitute a feature?

Expensive computation, and vectors have too many zeros.
Limit to most frequent/informative words!

Bocconi

# Counting Trouble

*...And a Man named Zipf*

50000

40000

30000

20000

*50%*

10000

0

*THE OTHER 50%...*

33

Bocconi

# *N*-grams

**"As Gregor Samsa awoke one morning from uneasy dreams, he found himself transformed in his bed into a gigantic insect-like creature."**

Unigrams `As, Gregor, Samsa, awoke, one, morning, from, uneasy, dreams, ...`

Bigrams `As_Gregor, Gregor_Samsa, Samsa_awoke, awoke_one, one_morning, ...`

Trigrams `As_Gregor_Samsa, Gregor_Samsa_awoke, Samsa_awoke_one, awoke_one_morning, ...`

4-grams `As_Gregor_Samsa_awoke, Gregor_Samsa_awoke_one, Samsa_awoke_one_morning, ...`

Bocconi

# Finding Important Words:
# TF-IDF

Bocconi

# Some Words are Just More Interesting…

# Karen Spärck Jones

1935–2007

- Became a teacher before starting CS career at Cambridge

- Laid the foundation for modern NLP, Google Search, text classification

- Campaigned for more women in CS

- Namesake of prestigious CS prize

# Problems with Term Frequency

# Document and Term Frequency



$$IDF = log \ \frac{N}{df(w)}$$

FEATURE

DOCUMENT
FREQUENCY
(COUNT): 4

TERM FREQUENCY
(SUM): 9 TF

Bocconi

# Putting it Together

How often we saw the word

$$TFIDF(w) = TF(w) \cdot log \frac{N}{df(w)}$$

Adjusted by how many documents

BOCCONI

# Document and Term Frequency

Words in "Moby Dick"



| word | tf | idf | tfidf |
|------|------|----------|------------|
| ye | 467 | 4.257380 | 148.497079 |
| chapter | 171 | 5.039475 | 147.504638 |
| whale | 1150 | 3.262357 | 139.755743 |
| man | 525 | 3.982412 | 106.932953 |
| ahab | 511 | 4.019453 | 103.357774 |

# Dense Distributed Representations

# 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

# An Example

# Semantic Similarity



flats ≈ apartments

platypus

# Similarity Measures

**x**

**Cosine similarity**

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

*flats ≈ apartments*

0.84

0.13

*platypus*

-1.0

**y**

# Dot Product

- "combine" vectors to a scalar

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

*SUM*

*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 $\sqrt{\phantom{x}}$

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

# Nearest neighbors



k=5

k=3

# 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

# Word2Vec – CBOW Model

OUTPUT

garden

MATRIX OF

TARGET WORDS

ERROR

BACKPROPAGATION SUM

INPUT

MATRIX OF

rent    Renting a large apartment in great location

CONTEXT WORDS

52

Bocconi

OUTPUT

INPUT

rent    Renting a large apartment in great location

Bocconi

# Caveat: Antonyms

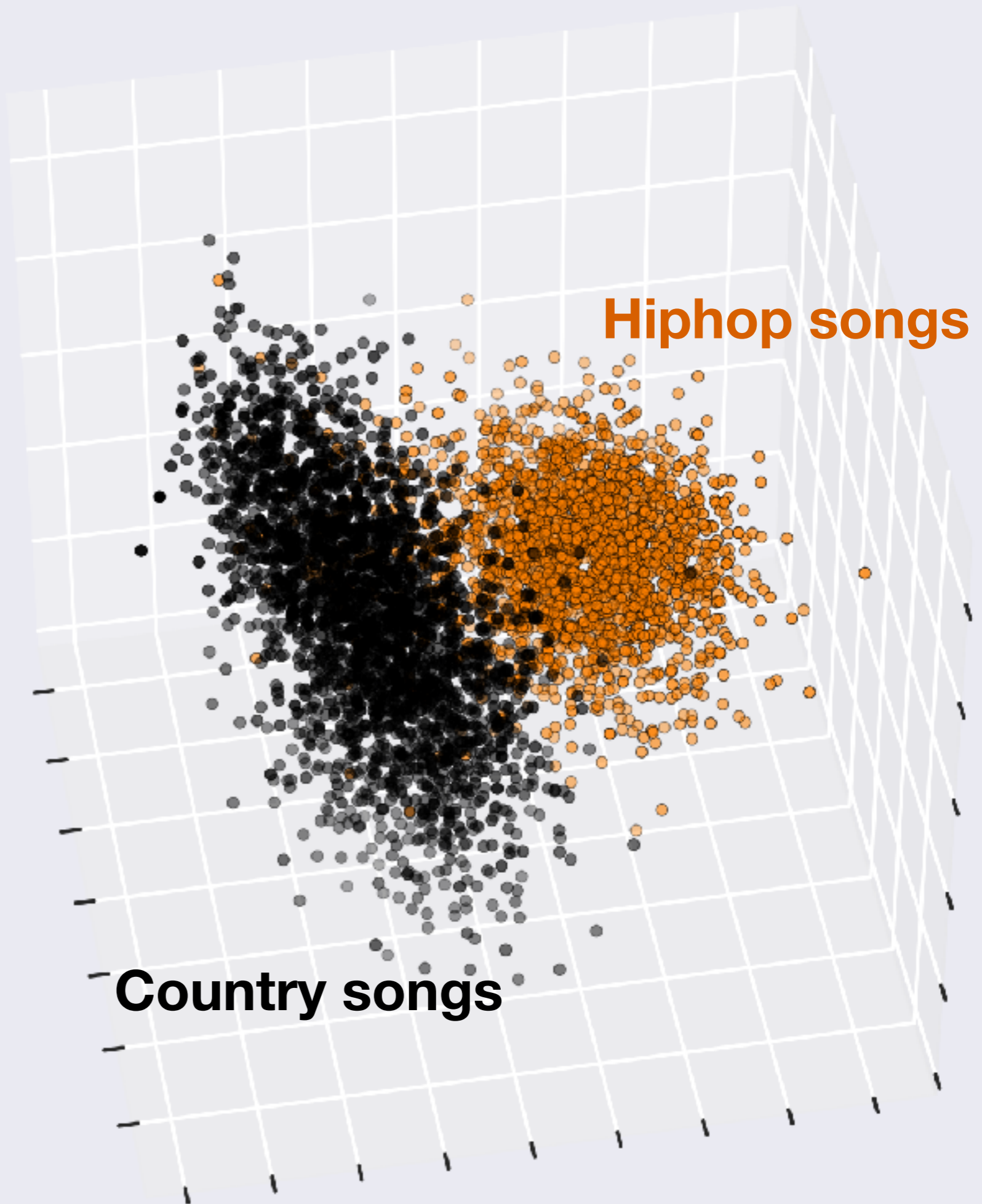*His kitchen was always very* _____



x

clean

dirty

*Same context, opposite meaning!*

y

Bocconi

# Part 2
# Representing Documents as Vectors

Billboard HOT 100

$C$

song 1

⋮

song *n*

**Hiphop songs**

**Country songs**
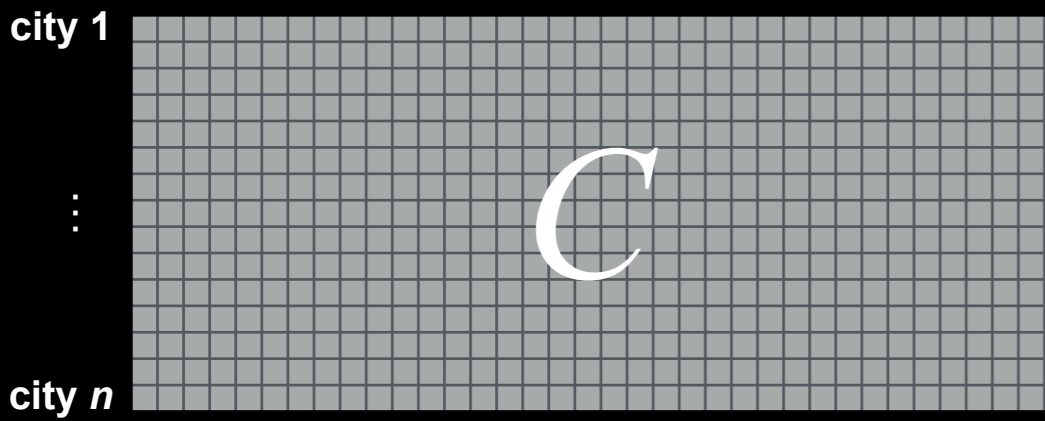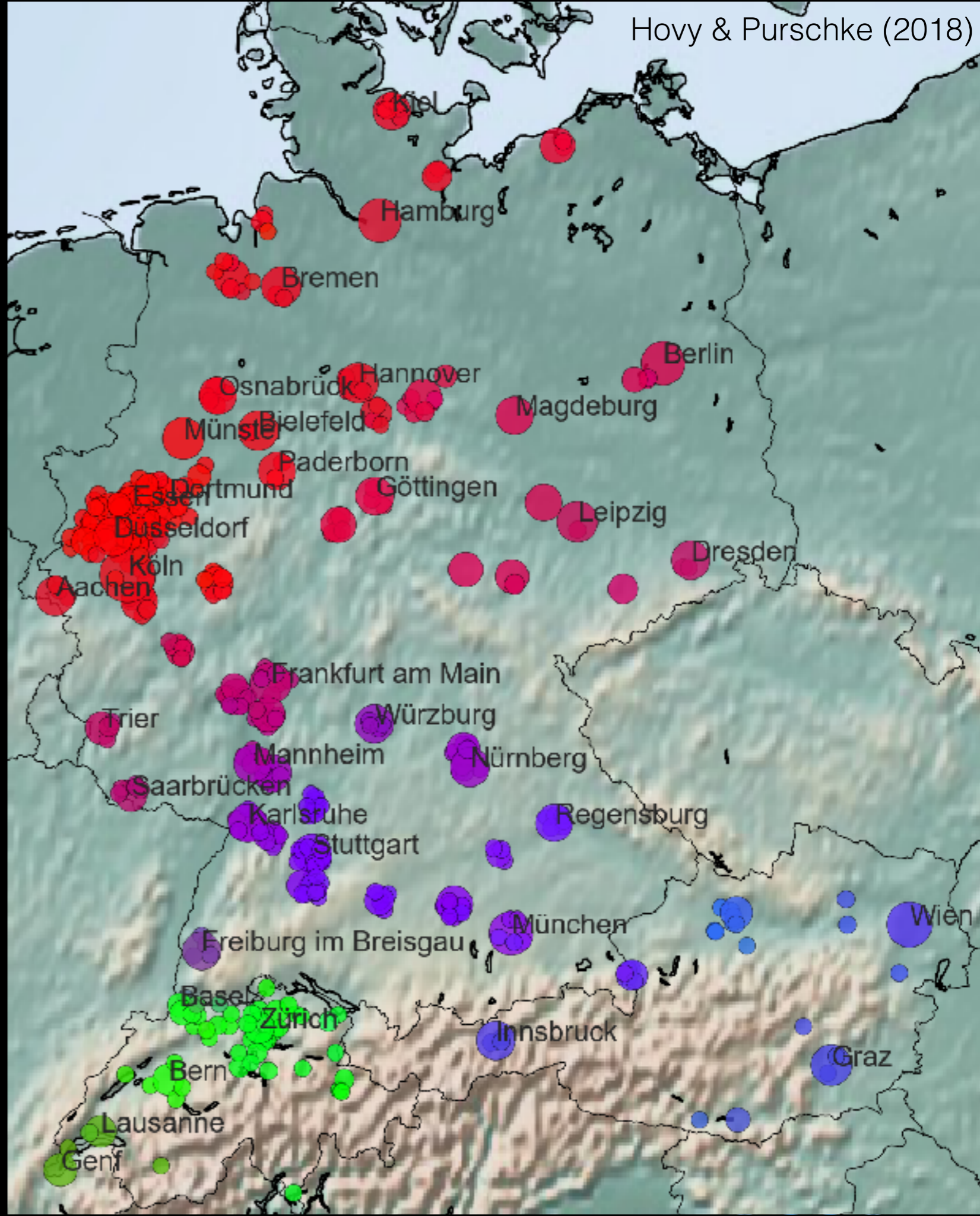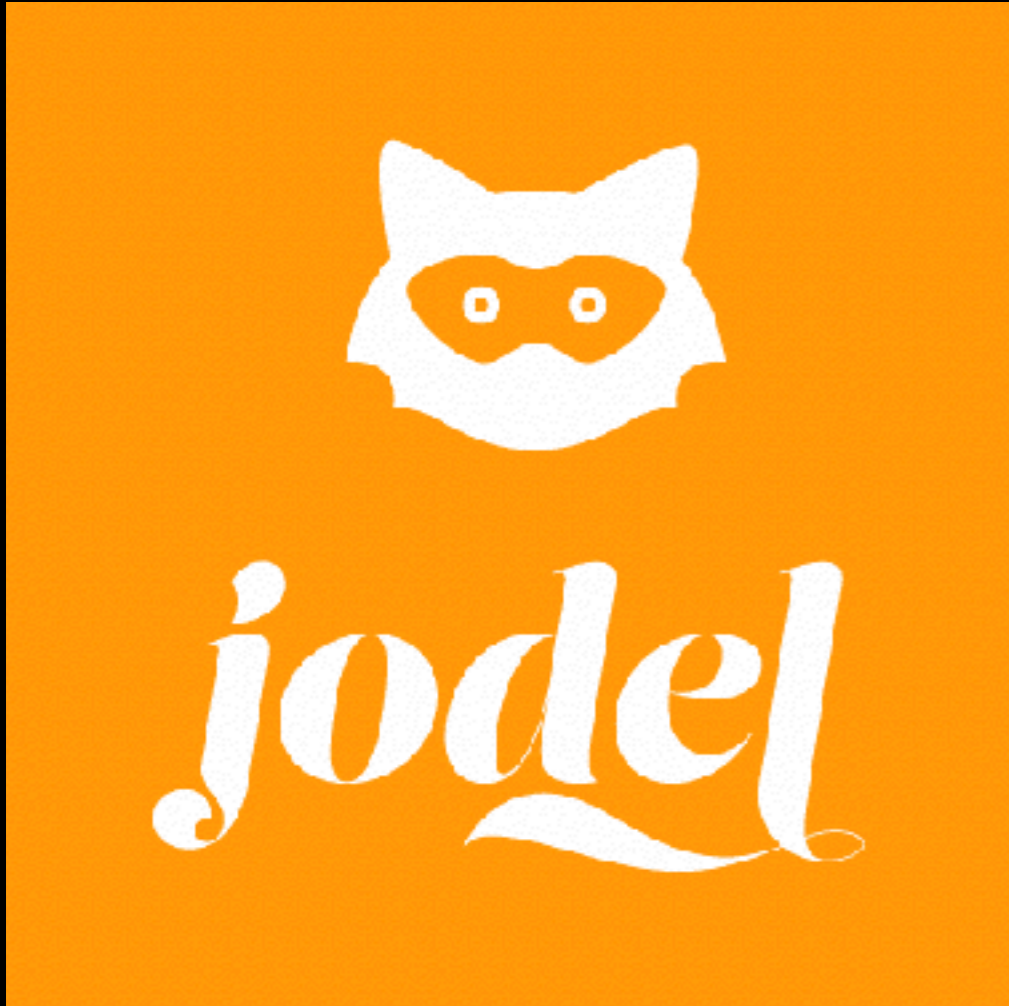
Example 2: Cities

Hovy & Purschke (2018)

# 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
```
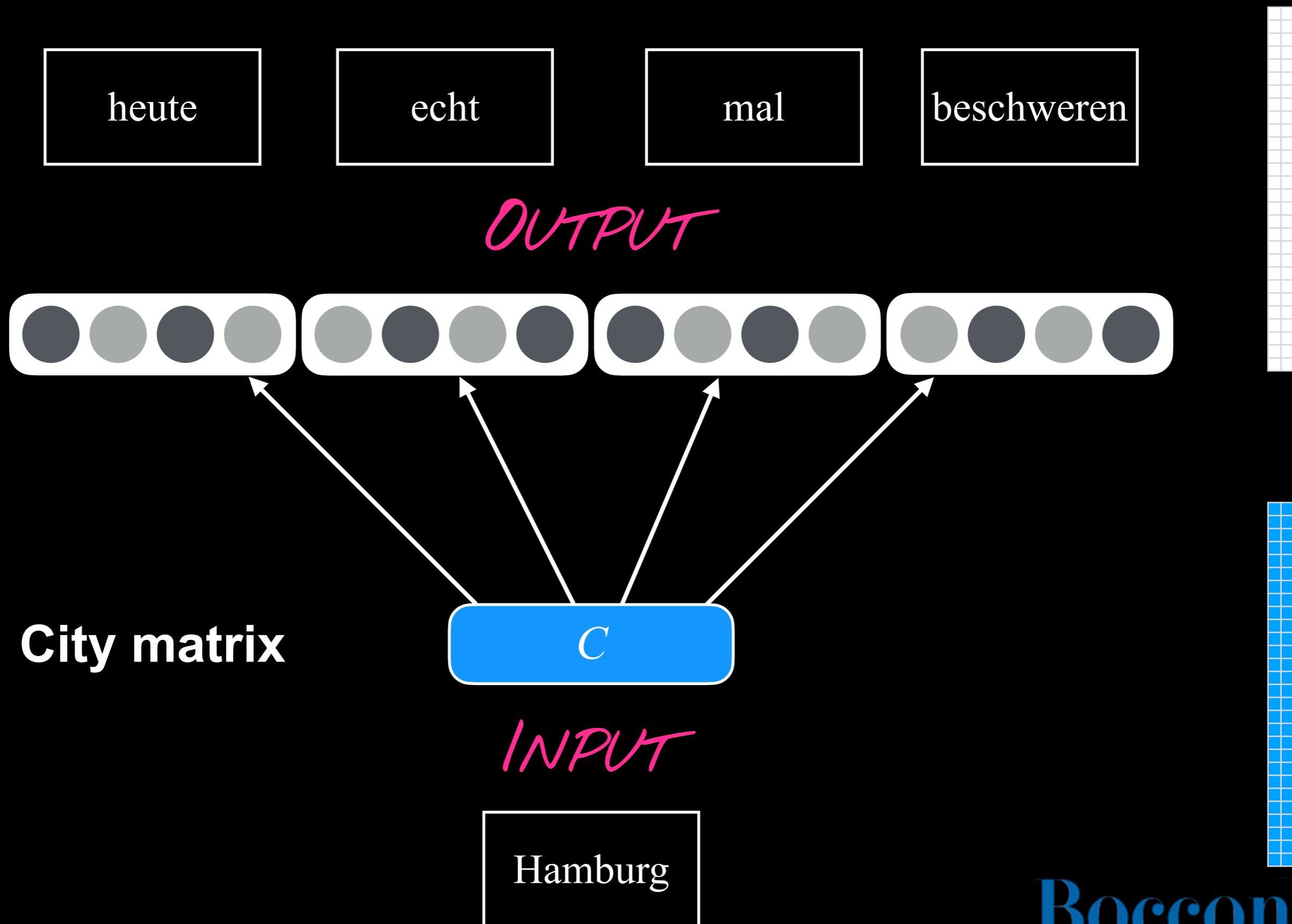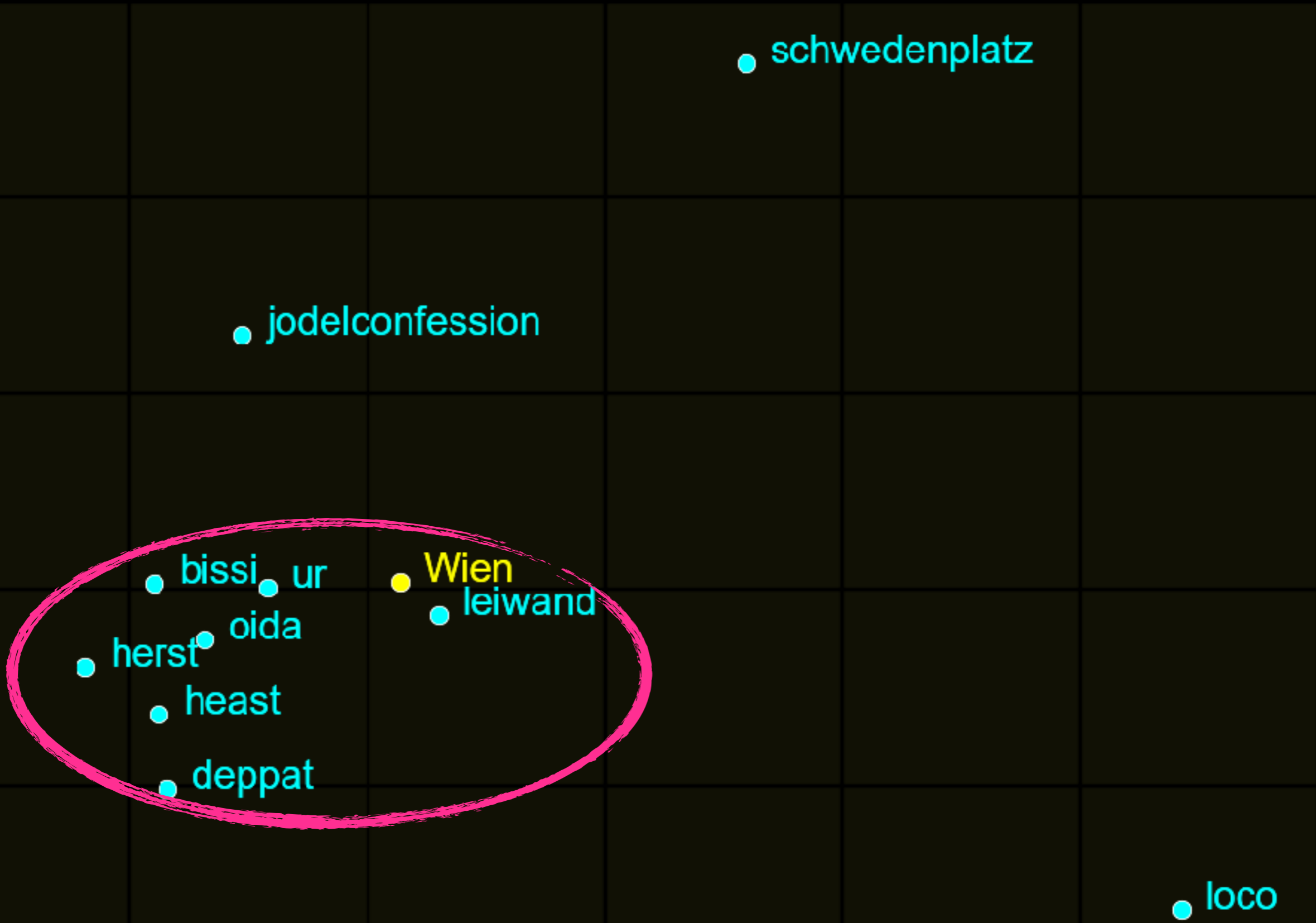
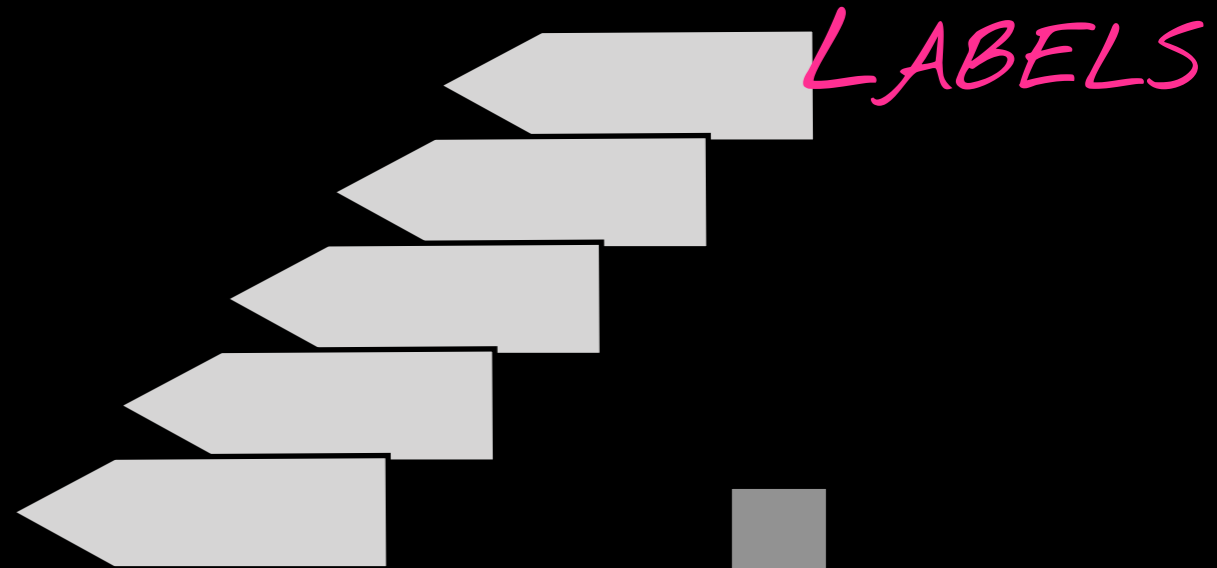# Doc2Vec – Model

heute     echt     mal     beschweren

*Output*

**City matrix**

$C$

*Input*

Hamburg

Bocconi

# Words and Documents



schwedenplatz

jodelconfession

bissi  ur  Wien  leiwand

oida

herst

heast

deppat

loco

# Wrapping up

# Comparison

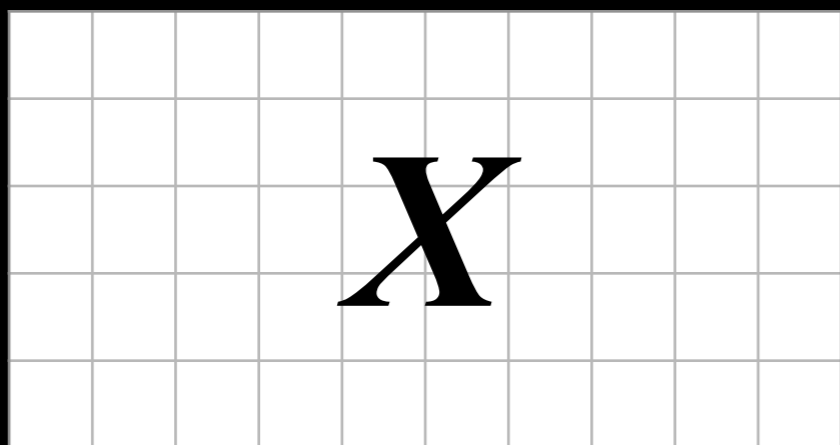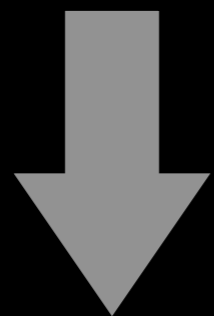| | Discrete | Distributed |
|---|---|---|
| **#Dimensions** | Data-dependent | Pre-defined |
| **Content** | Count-based | Coefficients |
| **Density** | Sparse | Dense |
| **Strength** | Interpretability | Similarity |
| **Application** | Understanding | Performance |
| **School of thought** | Rationalism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

Bocconi

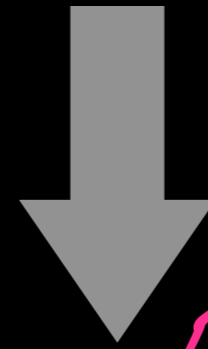# Text Classification

N Texts

Labels

N-by-D Matrix

$X$

N-by-1 Vector

$y$

Bocconi

# Fitting

$$f(\boldsymbol{X}) = y$$



D-BY-1 VECTOR

$w^T$

$X$

$y$

Bocconi

# Predicting

$$f(\mathbf{Z}) = \mathbf{Z}\, w^T = \hat{y}$$

K-BY-D
MATRIX

$\mathbf{Z}$

$w$

K-BY-K
VECTOR

$\hat{y}$

Bocconi

# Evaluating Performance

# Performance Problems

| x | y | ŷ |
|---|---|---|
| frog | 1 | 1 |
| deer | 1 | 1 |
| wolf | 1 | 1 |
| dog | 1 | 1 |
| bear | 1 | 1 |
| fish | 1 | 1 |
| bird | 1 | 0 |
| cat | 1 | 0 |
| stone | 0 | 1 |
| tree | 0 | 0 |

I HAVE A CLASSIFIER THAT'S 70% ACCURATE!

A 70% ACCURATE CLASSIFIER

Bocconi

# True and False

| predicted | | |
|---|---|---|
| **gold** | 1 | 0 |
| 1 | TP | FN |
| 0 | FP | TN |

*TARGET = ANIMAL*

| **x** | **y** | **ŷ** | |
|---|---|---|---|
| frog | 1 | 1 | |
| deer | 1 | 1 | |
| wolf | 1 | 1 | true positive |
| dog | 1 | 1 | |
| bear | 1 | 1 | |
| fish | 1 | 1 | |
| bird | 1 | 0 | |
| cat | 1 | 0 | false negative |
| stone | 0 | 1 | false positive |
| tree | 0 | 0 | true negative |

**accuracy** = (TP+TN) / (P + N)

**precision** = TP / (TP + FP)
**recall** = TP / (TP + FN)
**F1** = 2 (prec x rec) / (prec + rec)

*ACCURACY = 7/10 = 0.7*
*PRECISION = 6/7 = 0.86*
*RECALL = 6/8 = 0.75*
*F1 = 0.81*

**Bocconi**

| predicted | | |
|---|---|---|
| g o l d | 1 | 0 |
| 1 | TP | FN |
| 0 | FP | TN |

*TARGET = THING*

| x | y | ŷ | |
|---|---|---|---|
| frog | 0 | 0 | |
| deer | 0 | 0 | |
| wolf | 0 | 0 | true negative |
| dog | 0 | 0 | |
| bear | 0 | 0 | |
| fish | 0 | 0 | |
| bird | 0 | 1 | |
| cat | 0 | 1 | false positive |
| stone | 1 | 0 | false negative |
| tree | 1 | 1 | true positive |

**accuracy** = (TP+TN) / (P + N)

**precision** = TP / (TP + FP)
**recall** = TP / (TP + FN)
**F1** = 2 (prec x rec) / (prec + rec)

*ACCURACY = 7/10 = 0.7*
*PRECISION = 1/3 = 0.33*
*RECALL = 1/2 = 0.5*
*F1 = 0.4*

**Bocconi**

| predicted | | |
|---|---|---|
| g o l d | | 1 | 0 |
| | 1 | TP | FN |
| | 0 | FP | TN |

# *MICRO*Averaging

*WEIGH BY CLASS SIZE*

**accuracy** = (TP+TN) / (P + N)

**precision** = TP / (TP + FP)
**recall** = TP / (TP + FN)
**F1** = 2 (prec x rec) / (prec + rec)

*ANIMAL*                 *THING*

| **x** | **y** | **ŷ** | **x** | **y** | **ŷ** |
|---|---|---|---|---|---|
| frog | 1 | 1 | frog | 0 | 0 |
| deer | 1 | 1 | deer | 0 | 0 |
| wolf | 1 | 1 | wolf | 0 | 0 |
| dog | 1 | 1 | dog | 0 | 0 |
| bear | 1 | 1 | bear | 0 | 0 |
| fish | 1 | 1 | fish | 0 | 0 |
| bird | 1 | 1 | bird | 0 | 0 |
| cat | 1 | 0 | cat | 0 | 1 |
| stone | 0 | 1 | stone | 1 | 0 |
| tree | 0 | 0 | tree | 1 | 1 |

*ACC = (7+7)/(10+10) = 14/20 =0.7*
*PREC = (6+1)/(7+3) = 7/10 = 0.7*
*REC = (6+1)/(8+2) = 7/10 = 0.7*
*F1 = 0.7*

**Bocconi**

# *MACRO* **Averaging**

*WEIGH ALL CLASSES EQUALLY*

| | predicted | |
|---|---|---|
| **gold** | 1 | 0 |
| 1 | TP | FN |
| 0 | FP | TN |

**accuracy** = (TP+TN) / (P + N)

**precision** = TP / (TP + FP)

**recall** = TP / (TP + FN)

**F1** = 2 (prec x rec) / (prec + rec)

*ANIMAL*

| x | y | ŷ |
|---|---|---|
| frog | 1 | 1 |
| deer | 1 | 1 |
| wolf | 1 | 1 |
| dog | 1 | 1 |
| bear | 1 | 1 |
| fish | 1 | 1 |
| bird | 1 | 1 |
| cat | 1 | 0 |
| stone | 0 | 1 |
| tree | 0 | 0 |

*THING*

| x | y | ŷ |
|---|---|---|
| frog | 0 | 0 |
| deer | 0 | 0 |
| wolf | 0 | 0 |
| dog | 0 | 0 |
| bear | 0 | 0 |
| fish | 0 | 0 |
| bird | 0 | 0 |
| cat | 0 | 1 |
| stone | 1 | 0 |
| tree | 1 | 1 |

*ACC = (0.7 + 0.7) / 2 = 0.7*

*PREC = (0.86 + 0.33) / 2 = 0.6*

*REC = (0.5 + 0.75) / 2 = 0.63*

*F1 = 0.61*

**Bocconi**

| gold | | predicted | |
|---|---|---|---|
| | | 1 | 0 |
| | 1 | TP | FN |
| | 0 | FP | TN |

*PREDICT MAJORITY CLASS FOR ALL*

*TARGET = ANIMAL*

| x | y | ŷ |
|---|---|---|
| frog | 1 | 1 |
| deer | 1 | 1 |
| wolf | 1 | 1 |
| dog | 1 | 1 |
| bear | 1 | 1 |
| fish | 1 | 1 |
| bird | 1 | 1 |
| cat | 1 | 1 |
| stone | 0 | 1 |
| tree | 0 | 1 |

true positive

false positive

**accuracy** = (TP+TN) / (P + N)

**precision** = TP / (TP + FP)
**recall** = TP / (TP + FN)
**F1** = 2 (prec x rec) / (prec + rec)

*ACCURACY = 8/10 = 0.8*
*PRECISION = 8/10 = 0.8*
*RECALL = 8/8 = 1.0*
*F1 = 0.9*

*Bocconi*

# Metrics Overview

- **accuracy** can be too general

- **precision** and **recall** are per-class measures

- **precision** = how many of instances labeled as target class are actually *in* target class?

- **recall** = how many of *all* target class instances in data identified correctly?

- **F1** = symmetric mean of precision and recall

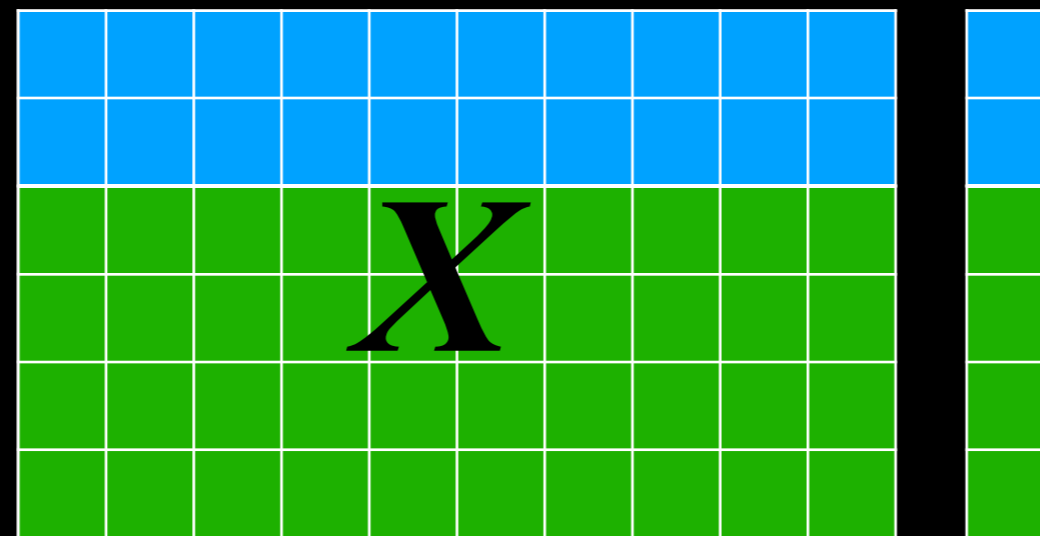**Bocconi**

# Beware: Overgeneralization

## *FALSE POSITIVES*

> June 6 2019
>
> Dear Ms Hovy,
>
> Congratulations on reaching retirement age!
>
> Also, you're on a no-fly list because of your political views and religious beliefs.

Bocconi

# Cross Validation

# Baselines

# Baseline: Total Recall

| g o l d | predicted | | |
|---|---|---|---|
| | | 1 | 0 |
| | 1 | TP | FN |
| | 0 | FP | TN |

*PREDICT MAJORITY CLASS FOR ALL*

**accuracy** = (TP+TN) / (P + N)

**precision** = TP / (TP + FP)

**recall** = TP / (TP + FN)

**F1** = 2 (prec x rec) / (prec + rec)

*TARGET = ANIMAL*

| x | y | ŷ |
|---|---|---|
| frog | 1 | 1 |
| deer | 1 | 1 |
| wolf | 1 | 1 |
| dog | 1 | 1 |
| bear | 1 | 1 |
| fish | 1 | 1 |
| bird | 1 | 1 |
| cat | 1 | 1 |
| stone | 0 | 1 |
| tree | 0 | 1 |

true positive

false positive

*ACCURACY = 8/10 = 0.8*

*PRECISION = 8/10 = 0.8*

*RECALL = 8/8 = 1.0*

*F1 = 0.9*

Bocconi

# Baseline: The Hulk

**(dumb but powerful)**

- Character 2–6 grams

- TFIDF weights

- L2-regularized Logistic Regression with balanced classes

- Can be further improved with dimensionality reduction

*ALWAYS CHECK AGAINST THIS BASELINE!*

Bocconi

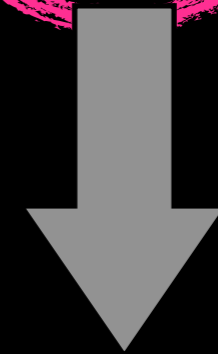# Regularization

# Regularization

$$y = X w^T + e$$

*D-BY-1 VECTOR*

$w^T$

$||w||$

# Regularization Norms

*L1 NORM*

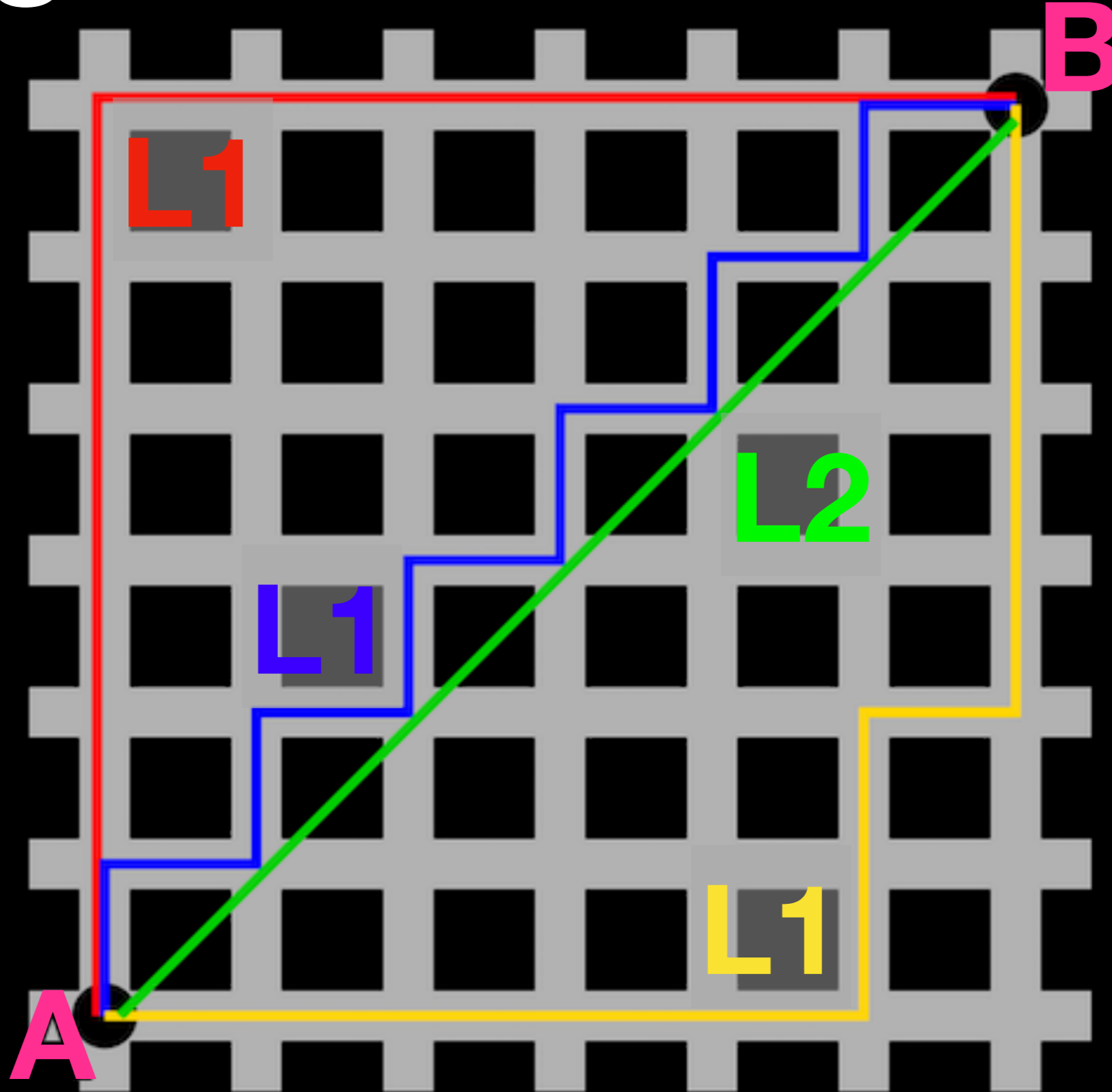$$||W||_1 = \sum_{i=1}^{N} |w_i|$$

*SPARSE*

*L2 NORM*

$$||W||_2 = \sqrt{\sum_{i=1}^{N} w_i^2}$$

*EVENLY DISTRIBUTED*

Bocconi
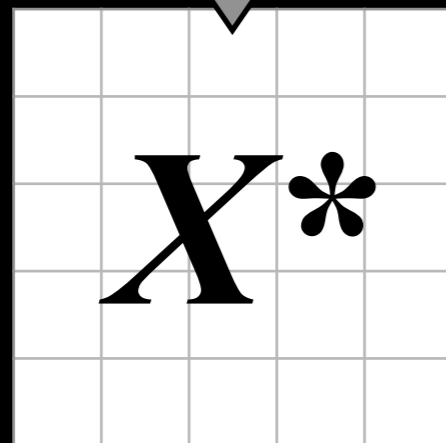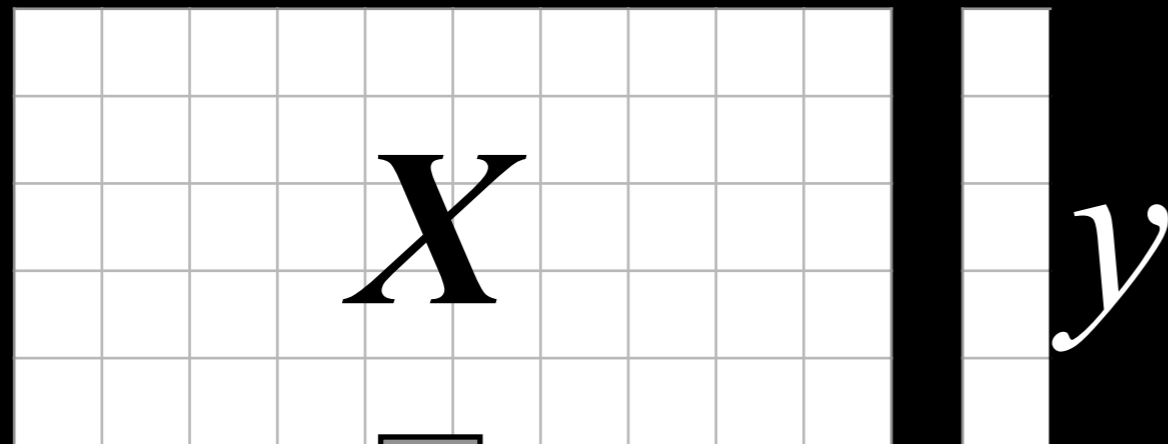
# Regularization Norms

# Feature Selection

# Chi-Squared Selection



$X$

$y$

$X^*$

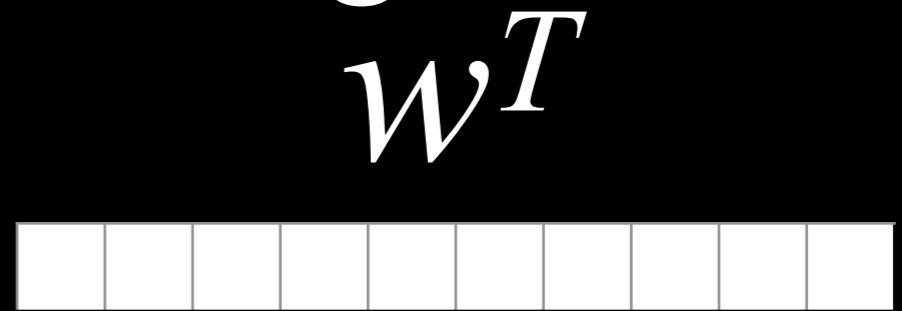Measure Chi2 value (correlation) for each Feature with target, select top K by cutoff

Bocconi

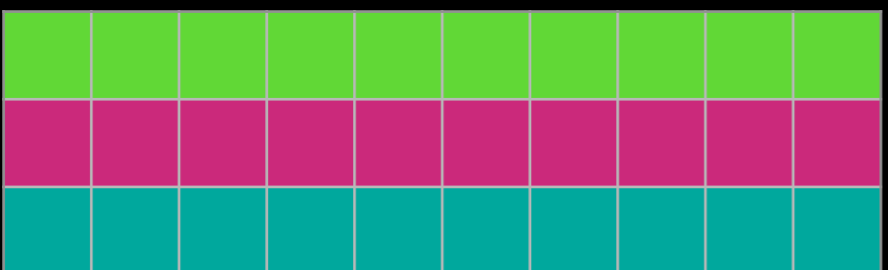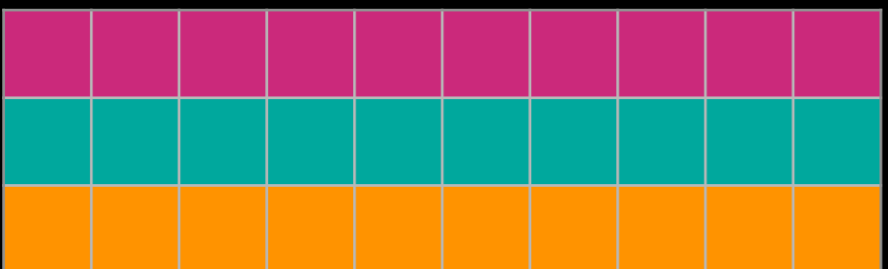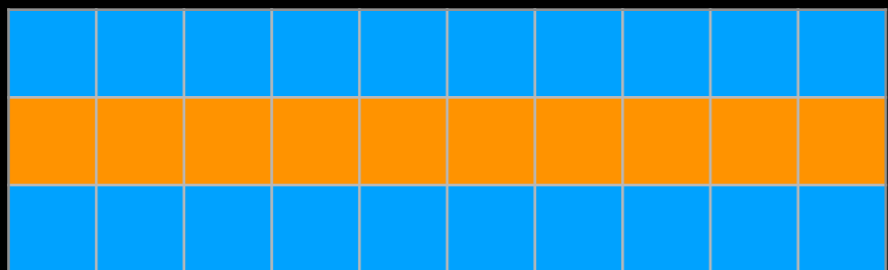# Dimensionality Reduction



$X$

$y$

$X*$

*Reduce Dimensionality to prevent spurious correlations with target, bring out latent dimensions*

Bocconi

# Randomized Logistic Regression

$$w^T$$

Fit N models wi L1 norm on subsets

AVERAGE | 1 | .3 | .6 | 0 | 1 | .6 | .3 | 0 | 1 | .3 |

# Wrapping Up

# Take Home Points

- **Preprocessing** removes noise and unwanted variation

- Words and texts can be represented as:

  - **Sparse, discrete** feature vectors (counts/TFIDF)

  - **Dense, continuous embedding** vectors

- Choose the appropriate performance **metric**

- Choose an informative **baseline**

- **Regularize, regularize, regularize**

- **Feature selection** can improve performance and provide insights

Bocconi