# Word Embeddings

Luís Filipe Cunha
lfc@di.uminho.pt
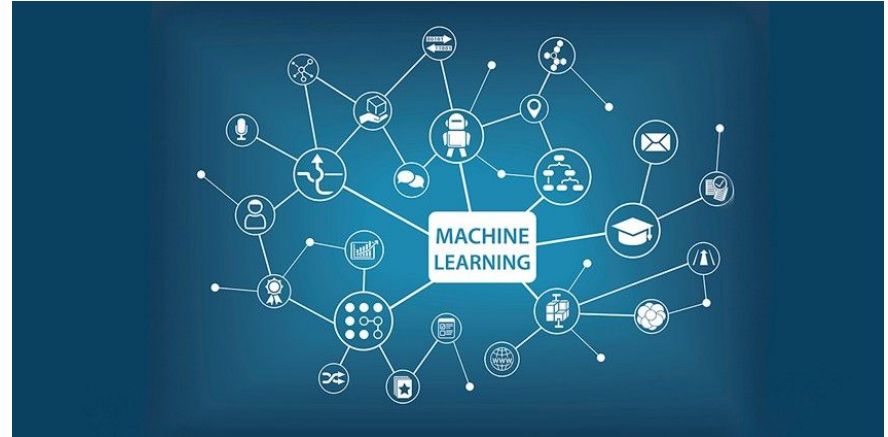
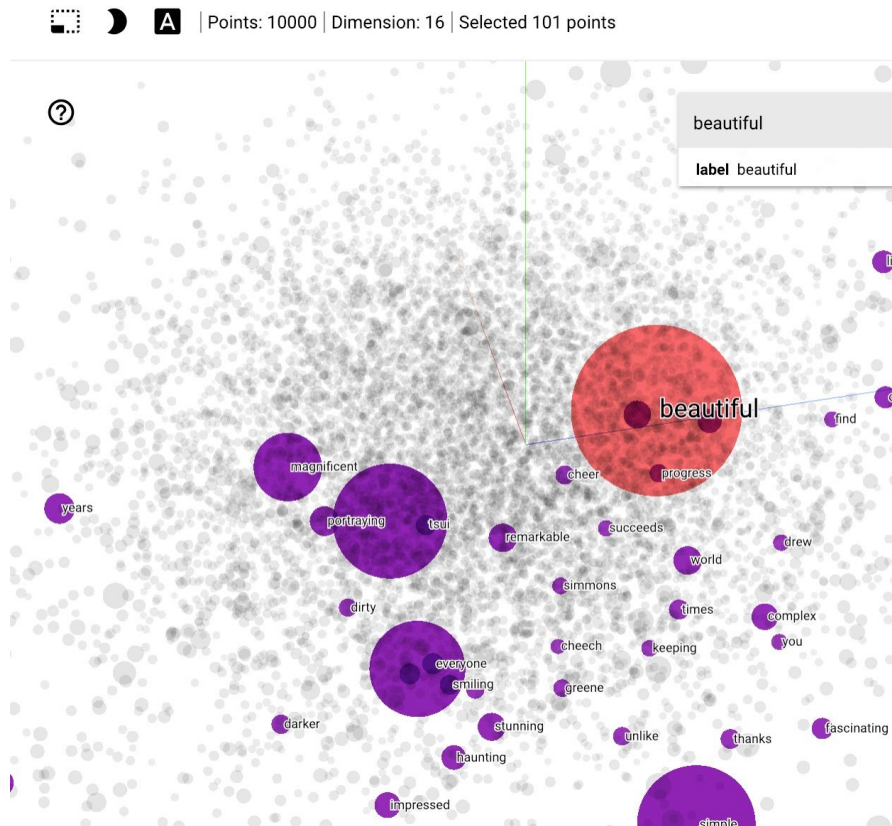José João Almeida
jj@di.uminho.pt

# Natural Language Processing

- Manual feature identification
  - Rule-based approaches
  - statistical models
- Deep Learning
  - Just feed the input data
  - Automatic feature learning

# Words Representations

- ML algorithms prefer well defined fixed-length inputs and outputs

- ML algorithms cannot work with raw text directly

- Numeric Vocabulary

- Bag of Words

- Word Embeddings

# Bag of Words (BOW)

Review 1: Game of Thrones is an amazing tv series!

Review 2: Game of Thrones is the best tv series!

Review 3: Game of Thrones is so great

- Tokenization
- Stop words
- Punctuation
- Ignore case
- Reducing words to their lemma
  - (e.g. "play" from "playing")

| | amazing | an | best | game | great | is | of | series | so | the | thrones | tv |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 |
| 2 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 |

# Limitations

- **Vocabulary:** Vector Length N (100k)
- **Sparsity:** Sparse Vectors
  - [0, 0, 0, 1, 0, …. 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
  - Large memory usage and expensive computation
- **Unknown words:** Words outside of vocabulary are ignored

| | amazing | an | best | game | great | is | of | series | so | the | thrones | tv |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 |
| 2 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 |

# Bag of Words (BOW)

- Sequence order is lost
  - Trabalhar para viver
  - Viver para trabalhar
- N-grams . Vector Dimensionality = V^N
- Vocabulary trigrams = 100k^3
- 1.000,000,000,000,000
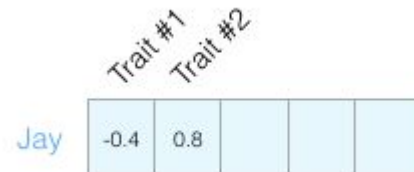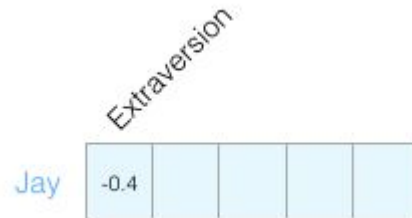- Semantic Meaning of the words lost
- Context is lost

| | amazing | an | best | game | great | is | of | series | so | the | thrones | tv |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 |
| 2 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 |

| | amazing tv | best tv | game thrones | thrones amazing | thrones best | thrones great | tv series |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |
| 2 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |

# Word Embeddings

| | | | |
|---|---|---|---|
| Openness to experience | ···· | 79 out of 100 |
| Agreeableness | ··················· | 75 out of 100 |
| Conscientiousness | ··············· | 42 out of 100 |
| Negative emotionality | ······· | 50 out of 100 |
| Extraversion | ························ | 58 out of 100 |

Extraversion

Jay | -0.4 | | | |

Trait #1  Trait #2

Jay | -0.4 | 0.8 | | | |

# Word Embeddings

- Dense
- Multidimensional
- length (50-1000)
- Words with similar meaning have similar numeric representation

## A 4-dimensional embedding

| cat => | 1.2 | -0.1 | 4.3 | 3.2 |
|--------|-----|------|-----|-----|
| mat => | 0.4 | 2.5 | -0.9 | 0.5 |
| on => | 2.1 | 0.3 | 0.1 | 0.4 |



Source: ■ Maria ■ Cão ■ Rei and Tenente

"In practice, short dense vectors work better"

# Embedding Layer

- Tokenization

- Create numeric vocabulary (*N* size)

- Create data batches

- Truncate and Padding

2 {'Data': 1, 'Local': 2, 'O': 3, 'Organizacao': 4, 'Pessoa': 5, 'Profissao': 6}

1 {'de': 1,            'Natural': 13,        'Meringolo': 9177,    'Adelina': 9189,
2 'e': 2,              'Filiação': 14,       'Pardo': 9178,        'Lbânia': 9190,
3 'do': 3,             'distrito': 15,       '2633': 9179,         'Rufino': 9191,
4 'ou': 4,             'º': 16,              '2016': 9180,         'Espírito': 9192,
5 'em': 5,             'o': 17,              'Atente': 9181,       'Prazeres': 9193,
6 'a': 6,              'n': 18,       (...)  'Joanesburgo': 9182,  'Etelvina': 9194,
7 'da': 7,             'que': 19,            'Gavela': 9183,       '1933': 9195,
8 'Maria': 8,          'Registo': 20,        'Calanga': 9184,      '1988': 9196,
9 'concelho': 9,       'Manuel': 21,         'Mambiça': 9185,      'Jesuína': 9197,
10 'país': 10,         'Pai': 22,            'Sotero': 9186,       'Sara': 9198,
11 'actual': 11,       'Mãe': 23,            '1951': 9187,         'Libânia': 9199
12 'residente': 12,    'para': 24,           'Bairros': 9188,      'terceiras': 9200}

9 words = [[2125, 1, 1482, 2, 2126, 695, 426, 1, 165, 1, 560, 1, 2755, 271, 1038, 347, 2, 225, 8,
      357, 2, 958, 106, 2, (...), 0, 0, 0, 0, 0], (...)]
10
11 labels = [[3, 3, 3, 3, 3, 3, 3, 3, 1, 1, 1, 1, 1, 3, 5, 5, 3, 3, 5, 5, 3, 5, 5, 3, (...), 0, 0,
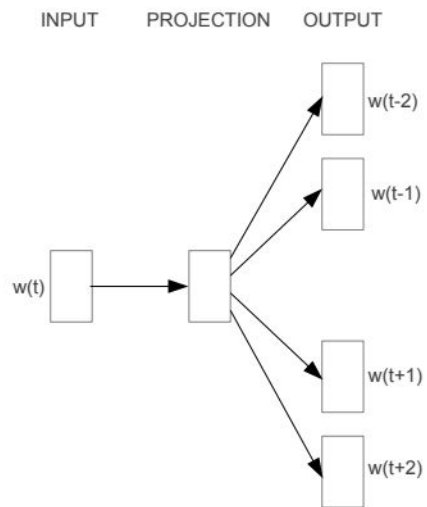      0, 0, 0], (...)]

# Word2Vec

- Trained to predict if a word belongs to the context
- "You shall know a word by the company it keeps" - John Rupert Firth
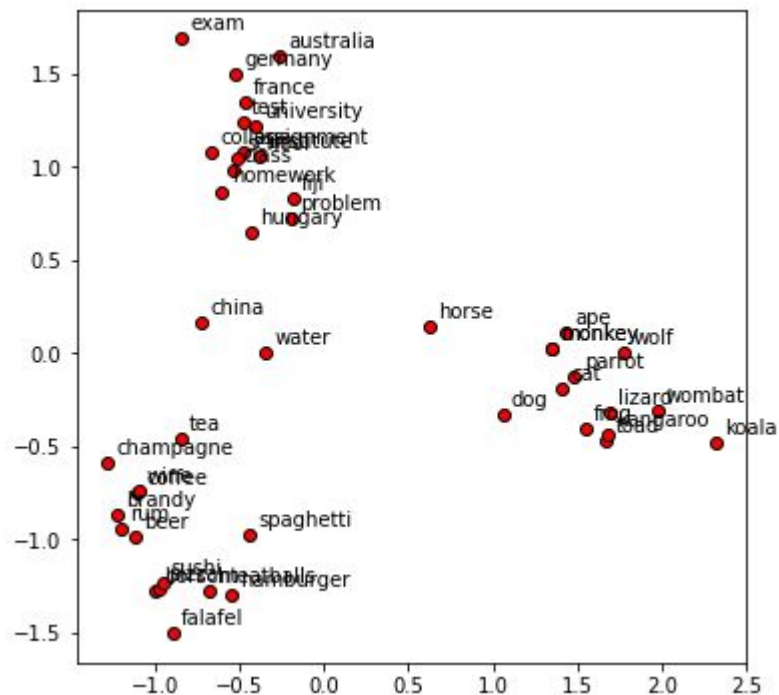- Milk is a likely word given "The cat was drinking"

# Word2Vec



INPUT    PROJECTION    OUTPUT

w(t-2)

w(t-1)

SUM

w(t+1)    w(t)

w(t+2)

**CBOW**

INPUT    PROJECTION    OUTPUT

w(t)    w(t-2)

w(t-1)

w(t+1)

w(t+2)

**Skip-gram**

# Similarity

# Analogies

king – man + woman ~= queen

# Limitations

- One vector per word (even if the word has multiple senses)

- Inability to handle unknown or OOV

- Scaling to new languages requires new embedding matrices

- Embeddings reflect cultural bias implicit in training text

# Reusing Word Embeddings (Transfer Learning)

- Train embeddings
- Use pre-trained word Embeddings
    - Glove
    - Word2vec

```
Corpora  >    Train Word
              Embeddings  >
```

```
Classifier Task1  >
Classifier Task2  >
Classifier Task3  >
```

# BIAS

- Ask "Paris : France :: Tokyo : x"
  - x = Japan
- Ask "father : doctor :: mother : x"
  - x = nurse
- Ask "man : computer programmer :: woman : x"
  - x = homemaker

# GPT-3 BIAS

- GPT-3 model presented biases towards gender, race, and religion (Brown et. al., 2020)
- Words such as "Islam" are associated with "terrorism".
- The word "female" word was usually associated with "naughty" or "beautiful"
- The "male" word is associated with "large", and "lazy".

# GPT3-Chat bot

# GPT3-Chat bot

# Data Visualization

# Dimension Reduction

- **PCA: Principal Component Analysis**

- **t-SNE: t-Distributed Stochastic Neighbor Embedding**



3D → 2D

# Principal Component Analysis (PCA)

- Dimensionality-reduction method

- Identifying patterns

- Trade a little accuracy for simplicity

- Preserving as much information as possible

1. Standardize the Dataset
2. Calculate the covariance matrix
3. Calculate the eigenvectors and eigenvalues
4. Choose Principal Components
5. Deriving the new data set (reorient the data)

# Standardize the Dataset

$$z = \frac{value - mean}{standard\ deviation}$$

|   |   | f1 | f2 | f3 | f4 |
|---|---|---|---|---|---|
| μ | = | 4 | 3 | 3 | 3.4 |
| σ | = | 3 | 1.58114 | 1.73205 | 2.30217 |

| f1 | f2 | f3 | f4 |
|---|---|---|---|
| 1 | 2 | 3 | 4 |
| 5 | 5 | 6 | 7 |
| 1 | 4 | 2 | 3 |
| 5 | 3 | 2 | 1 |
| 8 | 1 | 2 | 2 |

| f1 | f2 | f3 | f4 |
|---|---|---|---|
| -1 | -0.63246 | 0 | 0.26062 |
| 0.33333 | 1.26491 | 1.73205 | 1.56374 |
| -1 | 0.63246 | -0.57735 | -0.17375 |
| 0.33333 | 0 | -0.57735 | -1.04249 |
| 1.33333 | -1.26491 | -0.57735 | -0.60812 |

# Calculate the covariance matrix

Understand how the variables of the input data set are varying from the mean

Variables highly correlated can contain redundant information

p × p symmetric matrix (where p is the number of dimensions) that has as entries the covariances associated with all possible pairs of the initial variables

(Cov(a,a)=Var(a)),  (Cov(a,b)=Cov(b,a))

$$var(X) = \frac{\sum_{i=1}^{n}(X_i - \bar{X})(X_i - \bar{X})}{(n-1)}$$

$$cov(X,Y) = \frac{\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y})}{(n-1)}$$

$$\begin{bmatrix} Cov(x,x) & Cov(x,y) & Cov(x,z) \\ Cov(y,x) & Cov(y,y) & Cov(y,z) \\ Cov(z,x) & Cov(z,y) & Cov(z,z) \end{bmatrix}$$

# Covariance and Correlation

- if positive then : the two variables increase or decrease together (correlated)

- if negative then : One increases when the other decreases (Inversely correlated)

- covariance matrix summaries the correlations between all the possible pairs of variables.

$$cov = \begin{pmatrix} .616555556 & .615444444 \\ .615444444 & .716555556 \end{pmatrix}$$

# Calculate the eigenvectors and eigenvalues

$$Av = \lambda v$$

- Eigenvectors of the Covariance matrix are the directions of the axes where there is the most variance (information)

- Eigenvalues give the amount of variance

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \begin{pmatrix} 1 \\ 3 \end{pmatrix} = \begin{pmatrix} 11 \\ 5 \end{pmatrix}$$

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 12 \\ 8 \end{pmatrix} = 4 \times \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

# Calculate the eigenvectors and eigenvalues

- Provide information about patterns in the data.

- Line that best fits the data.

- Allow to create lines that characterise the data.

# Calculate the eigenvectors and eigenvalues

- 1st Eigenvector shows how the two sets of points are related along the line.

- 2nd Eigenvector shows that the points are off to the side of the main line by some amount (less important).

- Eigenvectors are pendicular to each other. (non correlated)



Mean adjusted data with eigenvectors overlayed

# Principal Components

$$eigenvalues = \begin{pmatrix} .0490833989 \\ 1.28402771 \end{pmatrix}$$

- Order them by eigenvalue, highest to lowest.

$$eigenvectors = \begin{pmatrix} -.735178656 & -.677873399 \\ .677873399 & -.735178656 \end{pmatrix}$$

- Components in order of significance.

- Information loss, however, if the eigenvalues are small, we don't lose much.

# Deriving the new data set

- Reorient the data from the original axes to the ones represented by the principal components

Original data restored using only a single eigenvector



| $x$ |
| --- |
| -.827970186 |
| 1.77758033 |
| -.992197494 |
| -.274210416 |
| -1.67580142 |
| -.912949103 |
| .0991094375 |
| 1.14457216 |
| .438046137 |
| 1.22382056 |

$$FinalDataSet = FeatureVector^T * StandardizedOriginalDataSet^T$$

# PCA

# t-Distributed Stochastic Neighbor Embedding

- Discover natural clusters

- Preserve the neighborhood

- Distant points correspond to dissimilar

  objects

1. Calculate similarity of points in High Dimension
2. Project all the points in the low dim space randomly
3. Calculate similarity of points in Low Dimension
4. Cost Function and gradient descendant

# t-Distributed Stochastic Neighbor Embedding

# Calculate similarity of points in High Dimension

- Calculate similarity of points in High Dimension

- Calculate similarity of points in Low Dimension

$$p_{ij} = \frac{exp(-||x_i - x_j||^2/2\sigma^2)}{\sum_{k \neq l} exp(-||x_l - x_k||^2/2\sigma^2)}$$

$$q_{ij} = \frac{(1 + ||y_i - y_j||^2)^{-1}}{\sum_{k \neq l}(1 + ||y_k - y_l||^2)^{-1}}$$

# Cost Function

KL(P || Q)

**Kullback Leibler Divergence**
Given two probabilities P and Q the KL divergence measures the how much does P as a distribution diverges from Q

Large Pij modeled by small qij: Large penalty

Small pij modeled by large qij: Small penalty

Minimization of the cost function
Gradient descent

$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

$$\frac{\delta C}{\delta y_i} = 2 \sum_j (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})(y_i - y_j).$$

# The crowding problem



Similarity in high dimension

Similarity in low dimension

There is much more space in high dimensions.

Blue = Gaussian
Red = Student's T

# The crowding problem



Student-t distribution has heavier tails.

# The crowding problem



(a) Gradient of SNE.     (b) Gradient of UNI-SNE.     (c) Gradient of t-SNE.

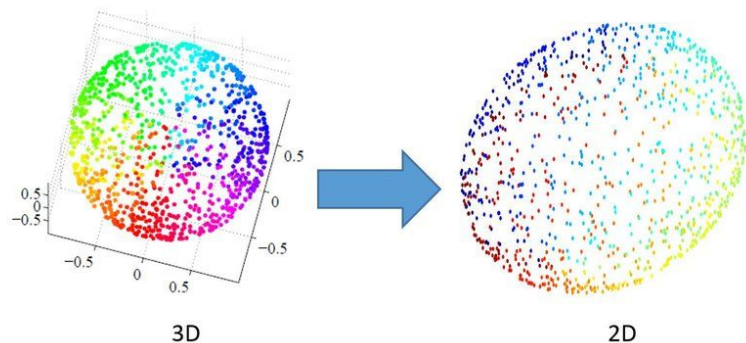t-SNE introduces strong repulsions between dissimilar datapoints that are

https://distill.pub/2016/misread-tsne/

# TSNE

# Dimension Reduction



3D            2D

**PCA: Principal Component Analysis**

- Preserve the global structure of the data

- Deterministic

- Preserves variance
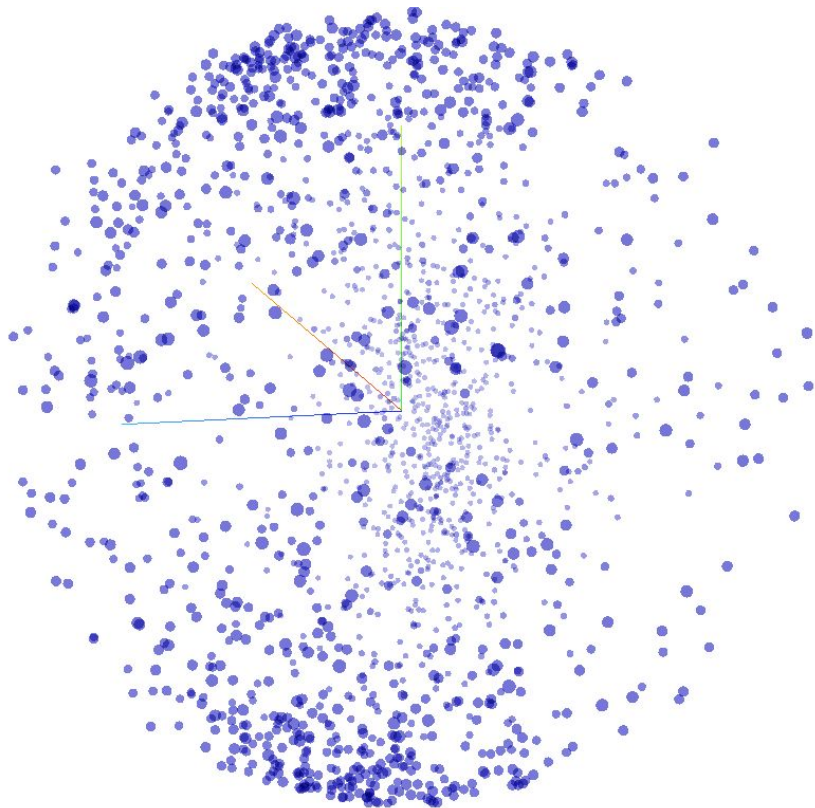
**t-SNE: t-Distributed Stochastic Neighbor Embedding**

- Preserve the local structure of data.
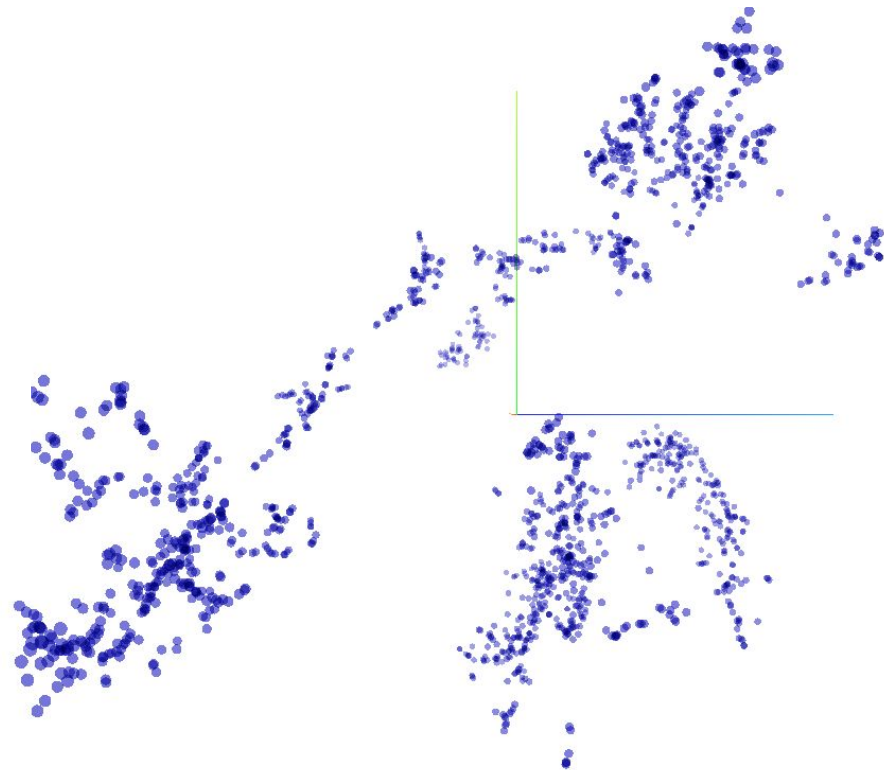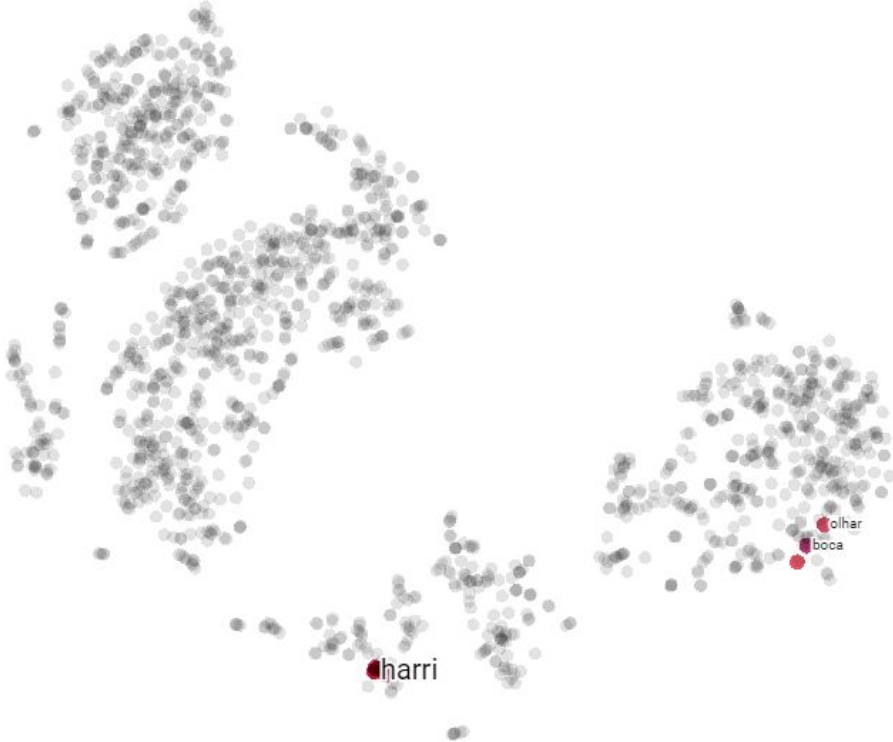
- Non-deterministic

- Preserves distance

# Data Visualization

# Word Embeddings

Luís Filipe Cunha
lfc@di.uminho.pt

José João Almeida
jj@di.uminho.pt