

# What is Relevant in a Text Document?

An Interpretable Machine Learning Approach

---

Leila Arras, Franziska Horn, Grégoire Montavon,  
Klaus-Robert Müller, Wojciech Samek

Christoph Schaller

# Word2Vec

---

- Is an approach to learn vector representations of words
- Using the context words to create the initial vectors

## Skipgram

- Is better to represent infrequent words
- Nearby context words have higher weight
- Trained by each context against the word

## CBOW

- Predicts a word given a window of context words
- Order of context words has no weight
- Trained by each word against its context

# Layer-Wise Relevance Propagation

---

# Identifying Relevant Words

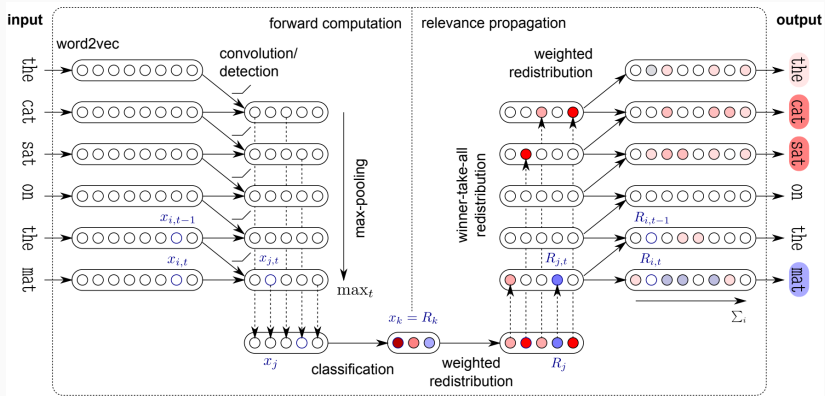


Figure 1: Diagram of a CNN-based interpretable machine learning system

# Identifying Relevant Words

Needs a vector-based word representation and a neural network

## Step One

Compute input representation of text document

## Step Two

Forward-propagate input representation

## Step Three

Backward-propagate using the layer-wise relevance propagation

## Step Four

Pool the relevance score onto the input neurons

## Computing the input representation of a text document

Words	Vectors
the	[0.035, -0.631, ...
cat	[0.751, -0.047, ...
sat	[0.491, 0.002, ...
on	[-0.181, -0.086, ...
the	[0.035, -0.631, ...

**Table 1:** CBOW vector example

### Forward-propagate the input representation until the output is reached

- We begin with our  $D \times L$  matrix-representation of the document
  - $D$  is the embedding dimension
  - $L$  is the document size
- The convolutional layer produces a new representation of  $F$  features maps of length  $L - H + 1$
- ReLU is applied element wise
- Features maps are pooled by computing the maximum over the text sequence of the document
- The maxpooled features are fed into a logistic classifier



## Step Two

ML Model	Test Accuracy (%)
CNN1 (H=1, F=600)	79.79
CNN2 (H=2, F=800)	<b>80.19</b>
CNN3 (H=3, F=600 )	79.75

**Table 2:** Performance of different CNN models

### Backward-propagate using the layer-wise relevance propagation

- Delivers one scalar relevance value per input variable, input data point and possible target class
- Redistributes the score of a predicted class back to the input space
- The Neuron that had the maximum value in the pool is granted all the relevance

### Pool the relevance score onto the convolutional layer

- $R_{(i,t-\tau)\leftarrow(j,t)} = \frac{z_{i,j,\tau}}{\sum_{i,\tau} z_{i,j,\tau}}$
- Similar to the Equation used for LRP
- More complex due to the convolutional structure of the layer

### Pool the relevance score onto the input neurons

- $R_{i,t} = \sum_{j,\tau} R_{(i,t)\leftarrow(j,t+\tau)}$
- The Word that had the maximum value in the pool is granted all the relevance

# Identifying Relevant Words

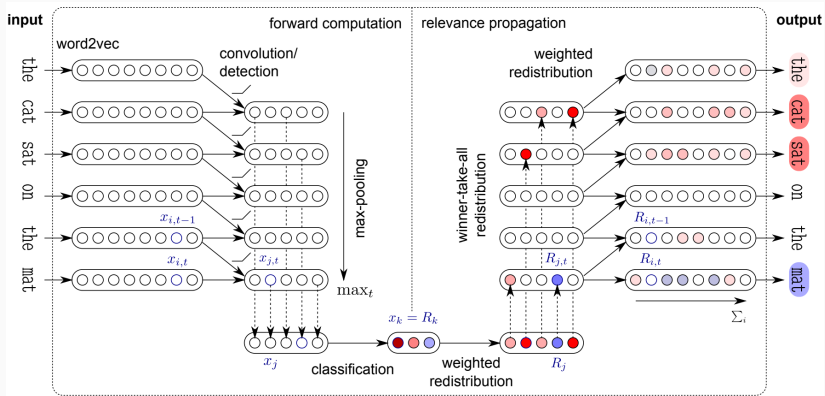


Figure 2: Diagram of a CNN-based interpretable machine learning system

## Obtaining relevance over all dimensions of word2vec

- $R_t = \sum_i R_{i,t}$   
pool relevances over all dimensions
- $\forall_i : d_i = \sum_t R_t \cdot x_{i,t}$   
condense semantic information
- $\forall_i : d_i = \sum_t R_{i,t} \cdot x_{i,t}$   
build document summary vector without pooling

## Bag of Words

- Documents are represented as vectors
- Each entry is TFIDF of a word in the training vocabulary

## Support Vector Machine

- Hyperplanes are learned to separate classes
- Linear prediction scores for each class are obtained
- $s_c = w_c^T x + b_c$
- $w_c$  are class specific weights
- $b_c$  is class specific bias

# Performance Comparison

ML Model	Test Accuracy (%)
BoW/SVM (V=70631 words)	80.10
CNN1 (H=1, F=600)	79.79
CNN2 (H=2, F=800)	<b>80.19</b>
CNN3 (H=3, F=600 )	79.75

**Table 3:** Performance of different ML Models

## LRP Decomposition

- $R_i = (w_c)_i \cdot x_i + b_c/D$
- $D$  is the number of non-zero entries of  $x$

## Vector Document Representation

- $d$  is built component-wise
- $\forall_j : d_j = R_j \cdot \tilde{x}_j$
- Replacing  $R_j$  with a TFIDF score allows comparability

## Why is this Approach the Baseline?

- Relies on word frequencies
- All words in the embeddings are equidistant



# Quality of Word References

---

# Comparing Relevance Scores

## How to compare relevance scores assigned by algorithms?

### Intrinsic Validation

Counting Words	Deleting Words
Creating a list of the most relevant words for a category across all documents	Removing words and measuring the decrease of the classification score

The second approach grants an objective estimation to compare relevance decomposition methods

## How to compare the explanatory power of ML models?

### Extrinsic Validation

#### Problems

- Need for common evaluation basis
- Classifiers differ in their reaction to removed words

#### Approach

Comparing models by how 'semantic extractive' their word relevances are

# Measuring Model Explanatory Power

## How to compare the explanatory power of ML models?

### Step One

Compute document summary vectors for all test set documents

### Step Two

- Normalize document summary vectors to euclidean norm
- perform K-nearest neighbor classification

### Step Three

- Repeat Step Two over ten random data splits
- Average KNN classification accuracies

The maximum KNN accuracy is used as explanatory power index

# Results

---

# Identification of Relevant Words

## CNN2

Yes, **weightlessness** does feel like falling. It may feel strange at first, but the body does adjust. The feeling is not too different from that of sky diving.

sci.space (8.1)

>And what is the motion sickness  
>that some **astronauts** occasionally experience?

It is the body's reaction to a strange environment. It appears to be induced partly to physical discomfort and part to mental distress. Some people are more prone to it than others, like some people are more prone to get sick on a roller coaster **ride** than others. The mental part is usually induced by a lack of clear indication of which way is up or down, ie: the Shuttle is normally oriented with its cargo bay pointed towards **Earth**, so the Earth (or ground) is "above" the head of the **astronauts**. About 50% of the **astronauts** experience some form of motion sickness, and **NASA** has done numerous tests in **space** to try to see how to keep the number of occurrences down.

Yes, weightlessness does feel like falling. It may feel strange at first, but the body does adjust. The feeling is not too different from that of sky diving.

sci.med (4.1)

>And what is the motion **sickness**  
>that some astronauts occasionally experience?

It is the **body's** reaction to a strange environment. It appears to be induced partly to physical **discomfort** and part to mental distress. Some people are more prone to it than others, like some people are more prone to get sick on a roller coaster **ride** than others. The mental part is usually induced by a lack of clear indication of which way is up or down, ie: the Shuttle is normally oriented with its cargo bay pointed towards **Earth**, so the Earth (or ground) is "above" the head of the **astronauts**. About 50% of the **astronauts** experience some form of motion **sickness**, and **NASA** has done numerous tests in **space** to try to see how to keep the number of occurrences down.

## SVM

Yes, weightlessness **does** feel like falling. It may feel strange at first, but the body **does** adjust. The feeling **is** not too different from that of **sky** diving.

sci.space (0.3)

>And what **is** the motion sickness  
>that some **astronauts** occasionally experience?

It **is** the body's reaction to a strange environment. It appears to be induced partly to physical discomfort and part to mental distress. Some **people** are more prone to it than others, like some **people** are more prone to get sick on a roller coaster **ride** than others. The mental part **is** usually induced by a lack of clear indication of which way is up or down, ie: the **Shuttle** **is** normally oriented with its cargo bay pointed towards **Earth**, so the **Earth** (or ground) **is** "above" the head of the **astronauts**. About 50% of the **astronauts** experience some form of **motion** sickness, and **NASA** has done numerous tests in **space** to try to see how to keep the number of occurrences down.

Yes, weightlessness does feel like falling. **It** may feel strange at first, but **the** **body** does adjust. **The** feeling **is** not too different from that **of** **sky** diving.

sci.med (-0.6)

>And what **is** **the** motion sickness  
>that **some** astronauts occasionally experience?

It **is** **the** **body's** reaction to a strange environment. It appears to be induced partly to physical **discomfort** and **part** **is** mental distress. Some people are more prone to it than others, like some people are more prone to get sick on a roller coaster **ride** than others. **The** mental part **is** usually induced by a lack **of** clear indication **of** which way is up or down, ie: **the** Shuttle **is** normally oriented with its cargo bay pointed towards **Earth**, so **the** **Earth** (or ground) **is** "above" **the** head **of** **the** **astronauts**. About 50% **of** **the** **astronauts** experience some form **of** **motion** sickness, and **NASA** has done numerous tests in **space** to try to see how to keep **the** number **of** occurrences **down**.

Figure 3: LRP heatmaps, positive is red, negative is blue

# Identification of Relevant Words

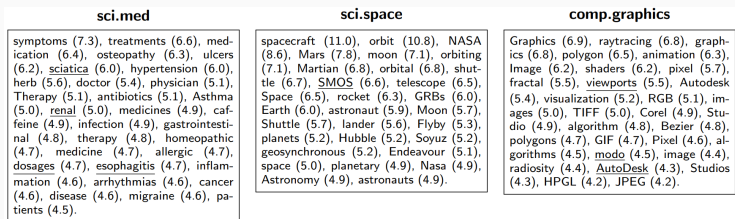


Figure 4: The 30 most relevant words for CNN2

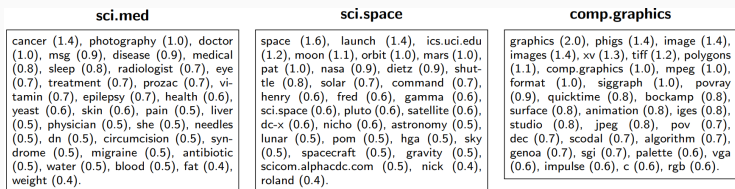


Figure 5: The 30 most relevant words for Bow/SVM

# Document Summary Vectors

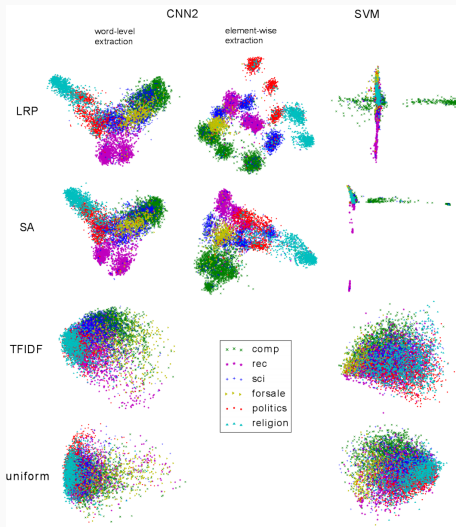


Figure 6: The 30 most relevant words for Bow/SVM



## How good is LRP in identifying relevant words?

- Delete Sequence of words from document
- Classify document again
- Report as function of accuracy and number of missing words

## Three different approaches

1.
  - Start with correctly classified documents
  - Delete words in decreasing order of their relevance
2.
  - Start with falsely classified documents
  - Delete words in increasing order of their relevance
3.
  - Start with falsely classified documents
  - Delete words in decreasing order of their score

# Evaluating LRP

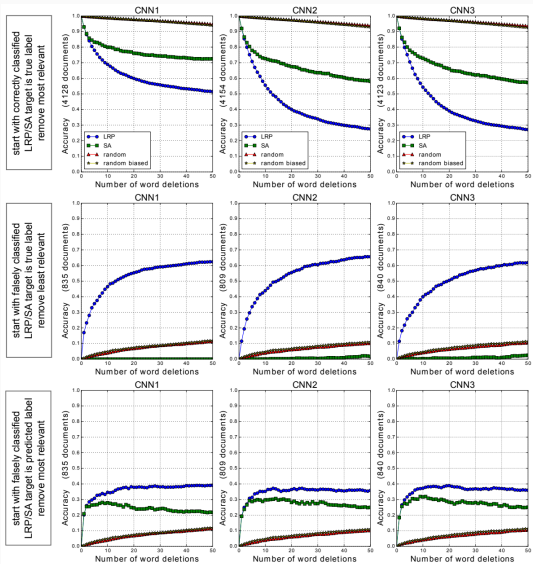


Figure 7: Word deletion experiments

# Quantifying Explanatory Power

Semantic Extraction		Explanatory Power Index (EPI)	KNN parameter
word2vec/CNN1	LRP (ew)	0.8045 ( $\pm$ 0.0044)	K = 10
	SA (ew)	0.7924 ( $\pm$ 0.0052)	K = 9
	LRP	0.7792 ( $\pm$ 0.0047)	K = 8
	SA	0.7773 ( $\pm$ 0.0041)	K = 6
word2vec/CNN2	LRP (ew)	0.8076 ( $\pm$ 0.0041)	K = 10
	SA (ew)	0.7993 ( $\pm$ 0.0045)	K = 9
	LRP	0.7847 ( $\pm$ 0.0043)	K = 8
	SA	0.7767 ( $\pm$ 0.0053)	K = 8
word2vec/CNN3	LRP (ew)	0.8034 ( $\pm$ 0.0039)	K = 13
	SA (ew)	0.7931 ( $\pm$ 0.0048)	K = 10
	LRP	0.7793 ( $\pm$ 0.0037)	K = 7
	SA	0.7739 ( $\pm$ 0.0054)	K = 6
word2vec	TFIDF	0.6816 ( $\pm$ 0.0044)	K = 1
	uniform	0.6208 ( $\pm$ 0.0052)	K = 1
BoW/SVM	LRP	0.7978 ( $\pm$ 0.0048)	K = 14
	SA	0.7837 ( $\pm$ 0.0047)	K = 17
BoW	TFIDF	0.7592 ( $\pm$ 0.0039)	K = 1
	uniform	0.6669 ( $\pm$ 0.0061)	K = 1

Table 4: Results over 10 random data splits

# Quantifying Explanatory Power

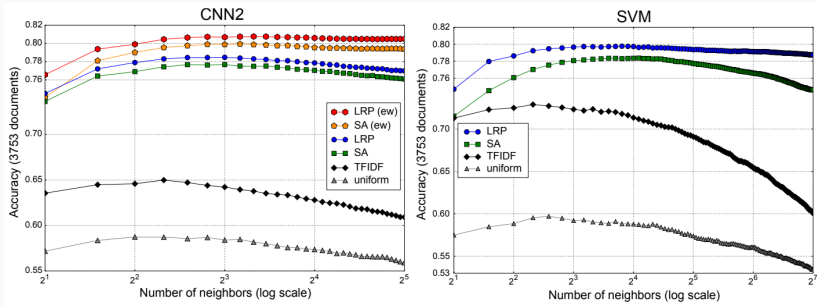


Figure 8: Word deletion experiments

Questions?