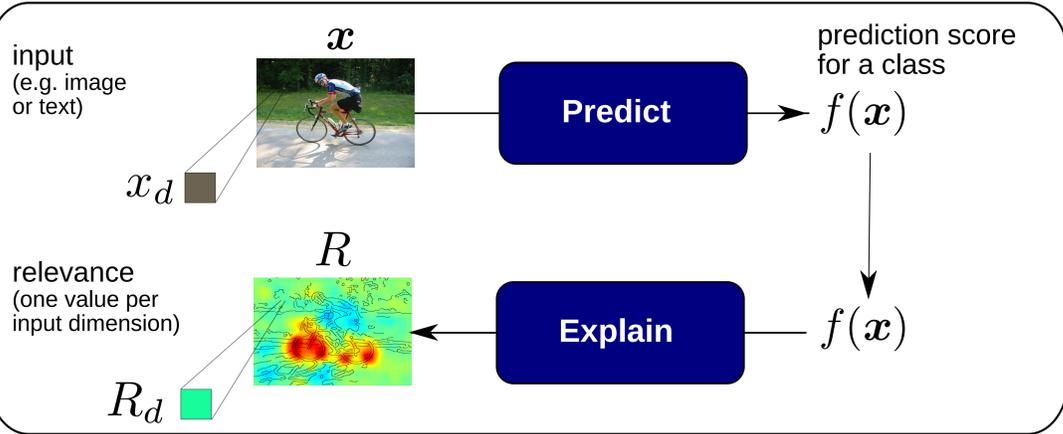


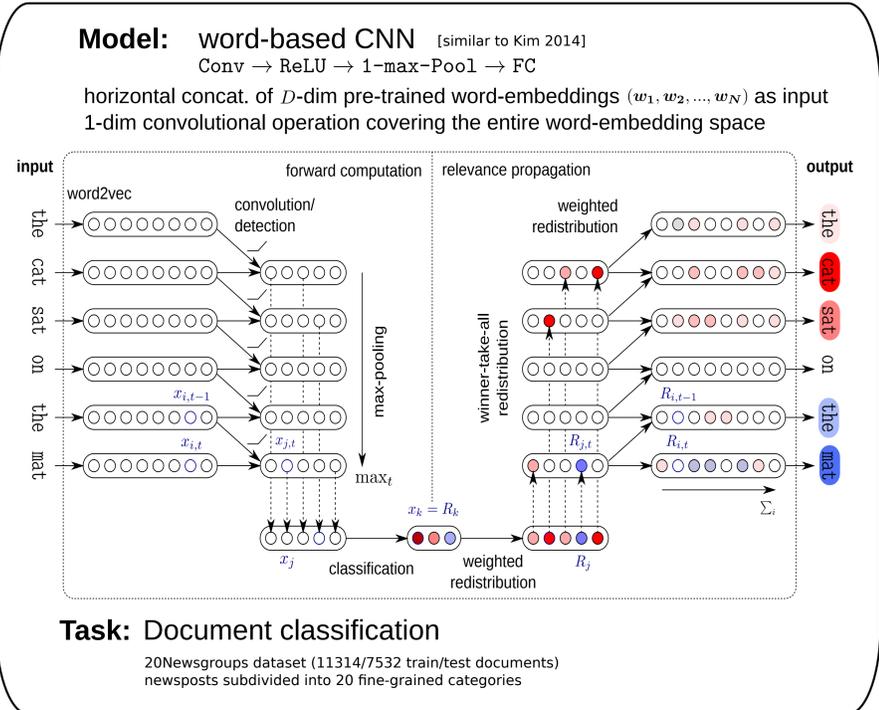
# Explaining Predictions of Non-Linear Classifiers in NLP

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

## Explaining Classifier Predictions



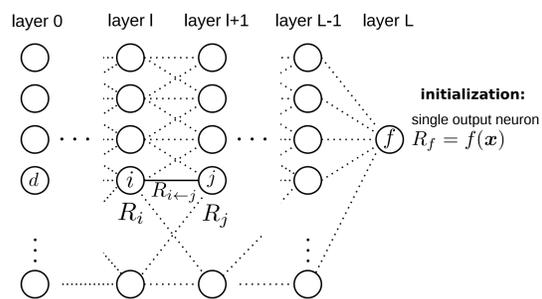
## Experimental Setup



## Explaining Neural Network Predictions

### LRP

Layer-wise Relevance Propagation  
[recently proposed by Bach et al. 2015]



Linear layer  $z_{ij} = x_i w_{ij} + \frac{b_j}{n}$   $n$ : receptive field of neuron  $j$   
 $z_j = \sum_i z_{ij}$   
 $s(z_j) = \epsilon \cdot (1_{z_j \geq 0} - 1_{z_j < 0})$  stabilizing term  
message proportional to neuron contribution in forward pass  $R_{i \leftarrow j} = \frac{s(z_j) + \frac{z_j}{n}}{\sum_i z_{ij} + s(z_j)} R_j$  message  
 $R_i = \sum_j R_{i \leftarrow j}$  sum of incoming messages

Activation layer  $R_i = R_j$  non-linear activation is taken into account by next linear layer

Max-pooling layer  $R_i = \begin{cases} R_j & i = \arg \max_l x_l \\ 0 & \text{else} \end{cases}$

### Properties:

- layer-wise conservation of score:  
 $\sum_d R_d = \dots = \sum_i R_i^{(l)} = \sum_j R_j^{(l+1)} = \dots = f(\mathbf{x})$
- signed relevance: indicate input regions that support or inhibit a specific classification decision
- relevance indicates contribution of input dim to actual classification decision (static)

### SA

Sensitivity Analysis or Gradient-Magnitude  
[e.g. Dimopoulos et al. 1995, Denil et al. 2014]

$$R_d = \left( \frac{\partial f}{\partial x_d} \right)^2 \quad \text{Backpropagation}$$

Relevance aggregation: for both LRP and SA by summation

Word-level Relevance w.r.t. our CNN Model:

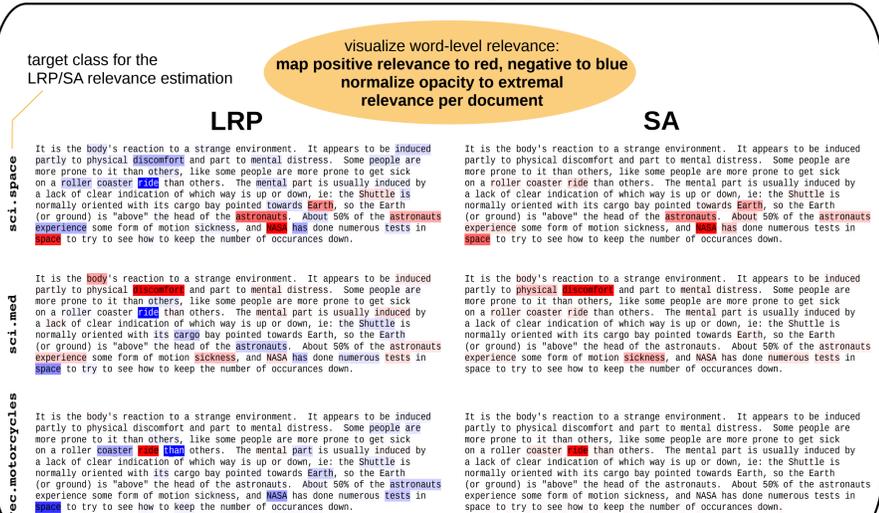
$$R(\mathbf{w}_t) = \sum_{i=1}^D R_{i,t}$$

$\mathbf{w}_t$   $t^{\text{th}}$  input word  
 $D$  word-embedding dimension

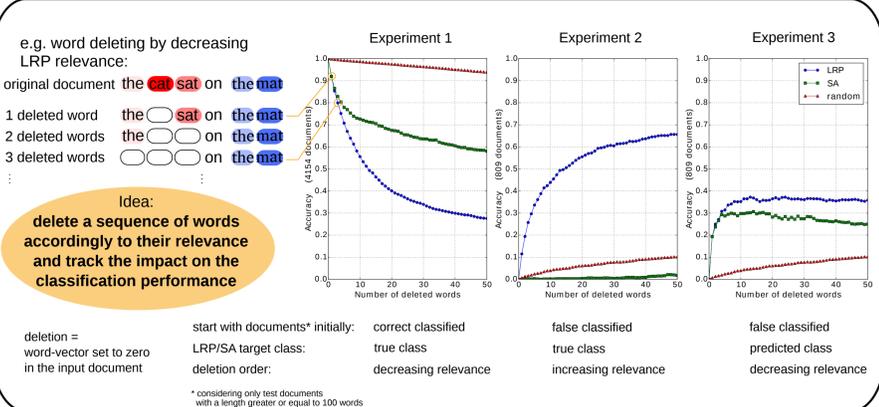
### Properties:

- conservation of squared gradient norm:  
 $\sum_d R_d = \|\nabla_{\mathbf{x}} f(\mathbf{x})\|_2^2$
- positive relevance: does not discriminate between input regions that support or inhibit a specific classification decision
- relevance indicates sensitivity of classification decision to changes in the input dimensions (dynamic)

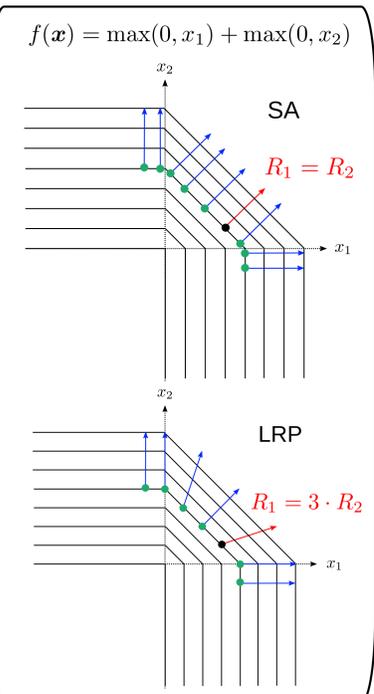
## Document Highlighting



## Quantitative Eval.: Word Deleting



## Intuition



## From word-vectors to document-vectors

