

# Natural Language Processing with Deep Learning

## Lecture 4 — Text classification 2: Deep neural networks

---

Prof. Dr. Ivan Habernal

November 3, 2023

Natural Language Processing Group  
Paderborn University

We focus on Trustworthy Human Language Technologies



[www.trusthlt.org](http://www.trusthlt.org)

# Where we finished last time

---

Where we finished last time

Log-linear multi-class classification

Representations

From multi-dimensional linear transformation to probabilities

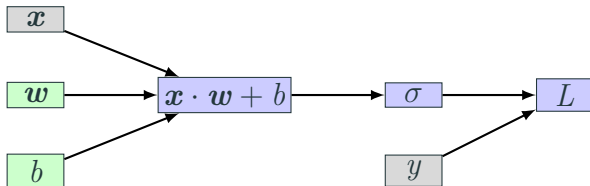
Loss function for softmax

Stacking transformations and non-linearity

# Our binary text classification function

Linear function through sigmoid — log-linear model

$$\hat{y} = \sigma(f(\mathbf{x})) = \frac{1}{1 + \exp(-(\mathbf{x} \cdot \mathbf{w} + b))}$$



**Figure 1:** Computational graph; green nodes are trainable parameters, gray are constant inputs

How can we minimize this function?

## (Online) Stochastic Gradient Descent

- 1: **function**  $\text{SGD}(f(\mathbf{x}; \Theta), (\mathbf{x}_1, \dots, \mathbf{x}_n), (\mathbf{y}_1, \dots, \mathbf{y}_n), L)$
- 2:     **while** stopping criteria not met **do**
- 3:         Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
- 4:         Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
- 5:          $\hat{\mathbf{g}} \leftarrow$  gradient of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  wrt.  $\Theta$
- 6:          $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
- 7:     **return**  $\Theta$

## (Online) Stochastic Gradient Descent

- 1: **function** SGD( $f(\mathbf{x}; \Theta)$ ,  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ ,  $(\mathbf{y}_1, \dots, \mathbf{y}_n)$ ,  $L$ )
- 2:     **while** stopping criteria not met **do**
- 3:         Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
- 4:         Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
- 5:          $\hat{\mathbf{g}} \leftarrow$  gradient of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  wrt.  $\Theta$
- 6:          $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
- 7:     **return**  $\Theta$

Loss in line 4 is based on a **single training example**  $\rightarrow$  a rough estimate of the corpus loss  $\mathcal{L}$  we aim to minimize

## (Online) Stochastic Gradient Descent

- 1: **function**  $\text{SGD}(f(\mathbf{x}; \Theta), (\mathbf{x}_1, \dots, \mathbf{x}_n), (\mathbf{y}_1, \dots, \mathbf{y}_n), L)$
- 2:     **while** stopping criteria not met **do**
- 3:         Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
- 4:         Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
- 5:          $\hat{\mathbf{g}} \leftarrow$  gradient of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  wrt.  $\Theta$
- 6:          $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
- 7:     **return**  $\Theta$

Loss in line 4 is based on a **single training example**  $\rightarrow$  a rough estimate of the corpus loss  $\mathcal{L}$  we aim to minimize

The noise in the loss computation may result in inaccurate gradients

# Minibatch Stochastic Gradient Descent

```
1: function MBSGD( $f(\mathbf{x}; \Theta)$ ,  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ ,  $(\mathbf{y}_1, \dots, \mathbf{y}_n)$ ,  $L$ )
2:   while stopping criteria not met do
3:     Sample  $m$  examples  $\{(\mathbf{x}_1, \mathbf{y}_1), \dots (\mathbf{x}_m, \mathbf{y}_m)\}$ 
4:      $\hat{\mathbf{g}} \leftarrow 0$ 
5:     for  $i = 1$  to  $m$  do
6:       Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$ 
7:        $\hat{\mathbf{g}} \leftarrow \hat{\mathbf{g}} + \text{gradient of } \frac{1}{m}L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i) \text{ wrt. } \Theta$ 
8:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$ 
9:   return  $\Theta$ 
```

# Properties of Minibatch Stochastic Gradient Descent

The minibatch size can vary in size from  $m = 1$  to  $m = n$

Higher values provide better estimates of the corpus-wide gradients, while smaller values allow more updates and in turn faster convergence

Lines 6+7: May be easily parallelized



# Log-linear multi-class classification

---

Where we finished last time

Log-linear multi-class classification

- Representations

- From multi-dimensional linear transformation to probabilities

Loss function for softmax

Stacking transformations and non-linearity

## From binary to multi-class labels

So far we mapped our gold label  $y \in \{0, 1\}$

What if we classify into distinct categorical classes?

- Categorical: There is no 'ordering'
- Example: Classify the language of a document into 6 languages (En, Fr, De, It, Es, Other)

## From binary to multi-class labels

So far we mapped our gold label  $y \in \{0, 1\}$

What if we classify into distinct categorical classes?

- Categorical: There is no 'ordering'
- Example: Classify the language of a document into 6 languages (En, Fr, De, It, Es, Other)

### One-hot encoding of labels

$$\text{En} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \quad \text{Fr} = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\text{De} = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 \end{pmatrix} \quad \dots$$

$$\mathbf{y} \in \mathbb{R}^{d_{out}} \quad \text{where } d_{out} \text{ is the number of classes}$$

## Possible solution: Six weight vectors and biases

Consider for each language  $\ell \in \{\text{En, Fr, De, It, Es, Other}\}$

- Weight vector  $w^\ell$  (e.g.,  $w^{\text{Fr}}$ )
- Bias  $b^\ell$  (e.g.,  $b^{\text{Fr}}$ )

## Possible solution: Six weight vectors and biases

Consider for each language  $\ell \in \{\text{En, Fr, De, It, Es, Other}\}$

- Weight vector  $\mathbf{w}^\ell$  (e.g.,  $\mathbf{w}^{\text{Fr}}$ )
- Bias  $b^\ell$  (e.g.,  $b^{\text{Fr}}$ )

We can predict the language resulting in the highest score

$$\hat{y} = f(\mathbf{x}) = \underset{\ell \in \{\text{En, Fr, De, It, Es, Other}\}}{\operatorname{argmax}} \quad \mathbf{x} \cdot \mathbf{w}^\ell + b^\ell$$

## Possible solution: Six weight vectors and biases

Consider for each language  $\ell \in \{\text{En, Fr, De, It, Es, Other}\}$

- Weight vector  $\mathbf{w}^\ell$  (e.g.,  $\mathbf{w}^{\text{Fr}}$ )
- Bias  $b^\ell$  (e.g.,  $b^{\text{Fr}}$ )

We can predict the language resulting in the highest score

$$\hat{y} = f(\mathbf{x}) = \underset{\ell \in \{\text{En, Fr, De, It, Es, Other}\}}{\operatorname{argmax}} \quad \mathbf{x} \cdot \mathbf{w}^\ell + b^\ell$$

But we can re-arrange the  $\mathbf{w} \in \mathbb{R}^{d_{in}}$  vectors into columns of a matrix  $\mathbf{W} \in \mathbb{R}^{d_{in} \times 6}$  and  $\mathbf{b} \in \mathbb{R}^6$ , to get

$$f(\mathbf{x}) = \mathbf{x}\mathbf{W} + \mathbf{b}$$

**Projecting input vector to output vector  $f(\boldsymbol{x}) : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^{d_{out}}$**

# Projecting input vector to output vector $f(\mathbf{x}) : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^{d_{out}}$

## Recall from lecture 3: High-dimensional linear functions

Function  $f(\mathbf{x}) : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^{d_{out}}$

$$f(\mathbf{x}) = \mathbf{x}\mathbf{W} + \mathbf{b}$$

where  $\mathbf{x} \in \mathbb{R}^{d_{in}}$        $\mathbf{W} \in \mathbb{R}^{d_{in} \times d_{out}}$        $\mathbf{b} \in \mathbb{R}^{d_{out}}$

The simplest neural network — a perceptron (simply a linear model)

- How to find the prediction  $\hat{y}$ ?



# Prediction of multi-class classifier

Project the input  $x$  to an output  $y$

$$\hat{y} = f(x) = xW + b$$

and pick the element of  $\hat{y}$  with the highest value

$$\text{prediction} = \hat{y} = \underset{i}{\operatorname{argmax}} \hat{y}_{[i]}$$

## Sanity check

What is  $\hat{y}$ ?

# Prediction of multi-class classifier

Project the input  $x$  to an output  $y$

$$\hat{y} = f(x) = xW + b$$

and pick the element of  $\hat{y}$  with the highest value

$$\text{prediction} = \hat{y} = \underset{i}{\operatorname{argmax}} \hat{y}_{[i]}$$

## Sanity check

What is  $\hat{y}$ ?

Index of 1 in the one-hot

For example, if  $\hat{y} = 3$ , then the document is in German

$$\text{De} = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 \end{pmatrix}$$

# Log-linear multi-class classification

---

## Representations

## Two representations of the input document

$$\hat{y} = xW + b$$

Vector  $x$  is a document representation

- Bag of words, for example ( $d_{in} = |V|$  dimensions, sparse)

Vector  $\hat{y}$  is **also** a document representation

- More compact (only 6 dimensions)
- More specialized for the language prediction task

## Matrix $W$ as learned representation — columns

$\hat{y} = xW + b \rightarrow$  two views of  $W$ , as rows or as columns

	En	Fr	De	It	Es	Ot
a	•	•	•	•	•	•
at	•	•	•	•	•	•
...						
zoo	•	•	•	•	•	•

## Matrix $W$ as learned representation — columns

$\hat{y} = xW + b \rightarrow$  two views of  $W$ , as rows or as columns

	En	Fr	De	It	Es	Ot
a	•	•	•	•	•	•
at	•	•	•	•	•	•
...						
zoo	•	•	•	•	•	•

Each of the 6 columns (corresponding to a language) is a  $d_{in}$ -dimensional vector representation of this language in terms of its characteristic word unigram patterns (e.g., we can then cluster the 6 language vectors according to their similarity)

## Matrix $W$ as learned representation — rows

$$\hat{y} = xW + b$$

	En	Fr	De	It	Es	Ot
a	•	•	•	•	•	•
at	•	•	•	•	•	•
...						
zoo	•	•	•	•	•	•

Each of the  $d_{in}$  rows corresponds to a particular unigram, and provides a 6-dimensional vector representation of that unigram in terms of the languages it prompts

# From bag-of-words to continuous bag-of-words

## Recall from lecture 3 — Averaged bag of words

$$\mathbf{x} = \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbf{x}^{D[i]}$$

$D[i]$  — word in doc  $D$  at position  $i$ ,  $\mathbf{x}^{D[i]}$  — one-hot vector

$$\hat{\mathbf{y}} = \mathbf{x} \mathbf{W} =$$



# From bag-of-words to continuous bag-of-words

## Recall from lecture 3 — Averaged bag of words

$$\mathbf{x} = \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbf{x}^{D[i]}$$

$D[i]$  — word in doc  $D$  at position  $i$ ,  $\mathbf{x}^{D[i]}$  — one-hot vector

$$\hat{\mathbf{y}} = \mathbf{x} \mathbf{W} = \left( \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbf{x}^{D[i]} \right) \mathbf{W}$$

# From bag-of-words to continuous bag-of-words

## Recall from lecture 3 — Averaged bag of words

$$\mathbf{x} = \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbf{x}^{D[i]}$$

$D[i]$  — word in doc  $D$  at position  $i$ ,  $\mathbf{x}^{D[i]}$  — one-hot vector

$$\hat{\mathbf{y}} = \mathbf{x} \mathbf{W} = \left( \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbf{x}^{D[i]} \right) \mathbf{W} = \frac{1}{|D|} \sum_{i=1}^{|D|} (\mathbf{x}^{D[i]} \mathbf{W})$$

=

# From bag-of-words to continuous bag-of-words

## Recall from lecture 3 — Averaged bag of words

$$\mathbf{x} = \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbf{x}^{D[i]}$$

$D[i]$  — word in doc  $D$  at position  $i$ ,  $\mathbf{x}^{D[i]}$  — one-hot vector

$$\begin{aligned}\hat{\mathbf{y}} = \mathbf{x} \mathbf{W} &= \left( \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbf{x}^{D[i]} \right) \mathbf{W} = \frac{1}{|D|} \sum_{i=1}^{|D|} (\mathbf{x}^{D[i]} \mathbf{W}) \\ &= \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbf{w}^{D[i]}\end{aligned}$$

(we ignore the bias  $\mathbf{b}$  here)

# From bag-of-words to continuous bag-of-words (CBOW)

Two equivalent views;  $W^{D[i]}$  is the  $D[i]$ -th row of matrix  $W$

$$\hat{y} = \frac{1}{|D|} \sum_{i=1}^{|D|} W^{D[i]} \quad \hat{y} = \left( \frac{1}{|D|} \sum_{i=1}^{|D|} x^{D[i]} \right) W$$

# From bag-of-words to continuous bag-of-words (CBOW)

Two equivalent views;  $W^{D[i]}$  is the  $D[i]$ -th row of matrix  $W$

$$\hat{y} = \frac{1}{|D|} \sum_{i=1}^{|D|} W^{D[i]} \quad \hat{y} = \left( \frac{1}{|D|} \sum_{i=1}^{|D|} x^{D[i]} \right) W$$

The continuous-bag-of-words (CBOW) representation

## From bag-of-words to continuous bag-of-words (CBOW)

Two equivalent views;  $W^{D[i]}$  is the  $D[i]$ -th row of matrix  $W$

$$\hat{y} = \frac{1}{|D|} \sum_{i=1}^{|D|} W^{D[i]} \quad \hat{y} = \left( \frac{1}{|D|} \sum_{i=1}^{|D|} x^{D[i]} \right) W$$

The continuous-bag-of-words (CBOW) representation

- Either by summing word-representation vectors

# From bag-of-words to continuous bag-of-words (CBOW)

Two equivalent views;  $W^{D[i]}$  is the  $D[i]$ -th row of matrix  $W$

$$\hat{y} = \frac{1}{|D|} \sum_{i=1}^{|D|} W^{D[i]} \quad \hat{y} = \left( \frac{1}{|D|} \sum_{i=1}^{|D|} x^{D[i]} \right) W$$

The continuous-bag-of-words (CBOW) representation

- Either by summing word-representation vectors
- Or by multiplying a bag-of-words vector by a matrix in which each row corresponds to a dense word representation (also called **embedding matrix**)

# Learned representations — central to deep learning

Representations are central to deep learning

One could argue that the main power of deep-learning is the ability to learn good representations



# **Log-linear multi-class classification**

---

**From multi-dimensional linear  
transformation to probabilities**

# Turning output vector into probabilities of classes

## Recap: Categorical probability distribution

Categorical random variable  $X$  is defined over  $K$  categories, typically mapped to natural numbers  $1, 2, \dots, K$ , for example  $\text{En} = 1$ ,  $\text{De} = 2$ ,  $\dots$

# Turning output vector into probabilities of classes

## Recap: Categorical probability distribution

Categorical random variable  $X$  is defined over  $K$  categories, typically mapped to natural numbers  $1, 2, \dots, K$ , for example En = 1, De = 2, ...

Each category parametrized with probability

$$\Pr(X = k) = p_k$$

# Turning output vector into probabilities of classes

## Recap: Categorical probability distribution

Categorical random variable  $X$  is defined over  $K$  categories, typically mapped to natural numbers  $1, 2, \dots, K$ , for example En = 1, De = 2, ...

Each category parametrized with probability

$$\Pr(X = k) = p_k$$

Must be valid probability distribution:  $\sum_{i=1}^K \Pr(X = i) = 1$

# Turning output vector into probabilities of classes

## Recap: Categorical probability distribution

Categorical random variable  $X$  is defined over  $K$  categories, typically mapped to natural numbers  $1, 2, \dots, K$ , for example En = 1, De = 2, ...

Each category parametrized with probability

$$\Pr(X = k) = p_k$$

Must be valid probability distribution:  $\sum_{i=1}^K \Pr(X = i) = 1$

How to turn an **unbounded** vector in  $\mathbb{R}^K$  into a categorical probability distribution?

# The softmax function $\text{softmax}(\mathbf{x}) : \mathbb{R}^K \rightarrow \mathbb{R}^K$

## Softmax

Applied element-wise, for each element  $\mathbf{x}_{[i]}$  we have

$$\text{softmax}(\mathbf{x}_{[i]}) = \frac{\exp(\mathbf{x}_{[i]})}{\sum_{k=1}^K \exp(\mathbf{x}_{[k]})}$$

# The softmax function $\text{softmax}(\mathbf{x}) : \mathbb{R}^K \rightarrow \mathbb{R}^K$

## Softmax

Applied element-wise, for each element  $\mathbf{x}_{[i]}$  we have

$$\text{softmax}(\mathbf{x}_{[i]}) = \frac{\exp(\mathbf{x}_{[i]})}{\sum_{k=1}^K \exp(\mathbf{x}_{[k]})}$$

- Nominator: Non-linear bijection from  $\mathbb{R}$  to  $(0; \infty)$
- Denominator: Normalizing constant to ensure

$$\sum_{j=1}^K \text{softmax}(\mathbf{x}_{[j]}) = 1$$

# The softmax function $\text{softmax}(\mathbf{x}) : \mathbb{R}^K \rightarrow \mathbb{R}^K$

## Softmax

Applied element-wise, for each element  $\mathbf{x}_{[i]}$  we have

$$\text{softmax}(\mathbf{x}_{[i]}) = \frac{\exp(\mathbf{x}_{[i]})}{\sum_{k=1}^K \exp(\mathbf{x}_{[k]})}$$

- Nominator: Non-linear bijection from  $\mathbb{R}$  to  $(0; \infty)$
- Denominator: Normalizing constant to ensure

$$\sum_{j=1}^K \text{softmax}(\mathbf{x}_{[j]}) = 1$$

We also need to know how to compute the partial derivative of  $\text{softmax}(\mathbf{x}_{[i]})$  wrt. each argument  $\mathbf{x}_{[k]}$ :  $\frac{\partial \text{softmax}(\mathbf{x}_{[i]})}{\partial \mathbf{x}_{[k]}}$



## Softmax can be smoothed with a 'temperature' $T$

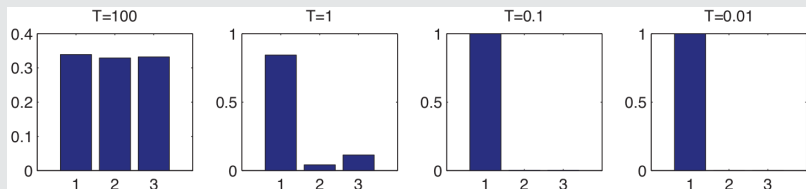
$$\text{softmax}(\mathbf{x}_{[i]}; T) = \frac{\exp(\frac{x_{[i]}}{T})}{\sum_{k=1}^K \exp(\frac{x_{[k]}}{T})}$$

# Softmax can be smoothed with a 'temperature' $T$

$$\text{softmax}(\mathbf{x}_{[i]}; T) = \frac{\exp(\frac{x_{[i]}}{T})}{\sum_{k=1}^K \exp(\frac{x_{[k]}}{T})}$$

Figure from K. Murphy (2012). ***Machine Learning: a Probabilistic Perspective***.  
MIT Press

## Example: Softmax of $\mathbf{x} = (3, 0, 1)$ at different $T$



High temperature  $\rightarrow$  uniform distribution

Low temperature  $\rightarrow$  'spiky' distribution, all mass on the largest element

# Loss function for softmax

---

Where we finished last time

Log-linear multi-class classification

Representations

From multi-dimensional linear transformation to probabilities

Loss function for softmax

Stacking transformations and non-linearity

# Categorical cross-entropy loss (aka. negative log likelihood)

Vector representing the gold-standard categorical distribution over the classes/labels  $1, \dots, K$ :

$$\mathbf{y} = (\mathbf{y}_{[1]}, \mathbf{y}_{[2]}, \dots, \mathbf{y}_{[K]})$$

Output from softmax:

$$\hat{\mathbf{y}} = (\hat{\mathbf{y}}_{[1]}, \hat{\mathbf{y}}_{[2]}, \dots, \hat{\mathbf{y}}_{[K]})$$

which is in fact  $\hat{\mathbf{y}}_{[i]} = \Pr(y = i | \mathbf{x})$

## Cross entropy loss

$$L_{\text{cross-entropy}}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{k=1}^K \mathbf{y}_{[k]} \log(\hat{\mathbf{y}}_{[k]})$$

## Background: K-L divergence (also known as *relative entropy*)

Let  $Y$  and  $\hat{Y}$  be categorical random variables over same categories, with probability distributions  $P(Y)$  and  $Q(\hat{Y})$

$$\begin{aligned}\mathbb{D}(P(Y) || Q(\hat{Y})) &= \mathbb{E}_{P(Y)} \left[ \log \frac{P(Y)}{Q(\hat{Y})} \right] \\ &= \mathbb{E}_{P(Y)} \left[ \log P(Y) - \log Q(\hat{Y}) \right] \\ &= \mathbb{E}_{P(Y)} [\log P(Y)] - \mathbb{E}_{P(Y)} [\log Q(\hat{Y})] \\ &= -\mathbb{E}_{P(Y)} \left[ \log \frac{1}{P(Y)} \right] - \mathbb{E}_{P(Y)} [\log Q(\hat{Y})] \\ &= -\mathbb{H}_P(Y) - \mathbb{E}_{P(Y)} [\log Q(\hat{Y})]\end{aligned}$$

# Stacking transformations and non-linearity

---

Where we finished last time

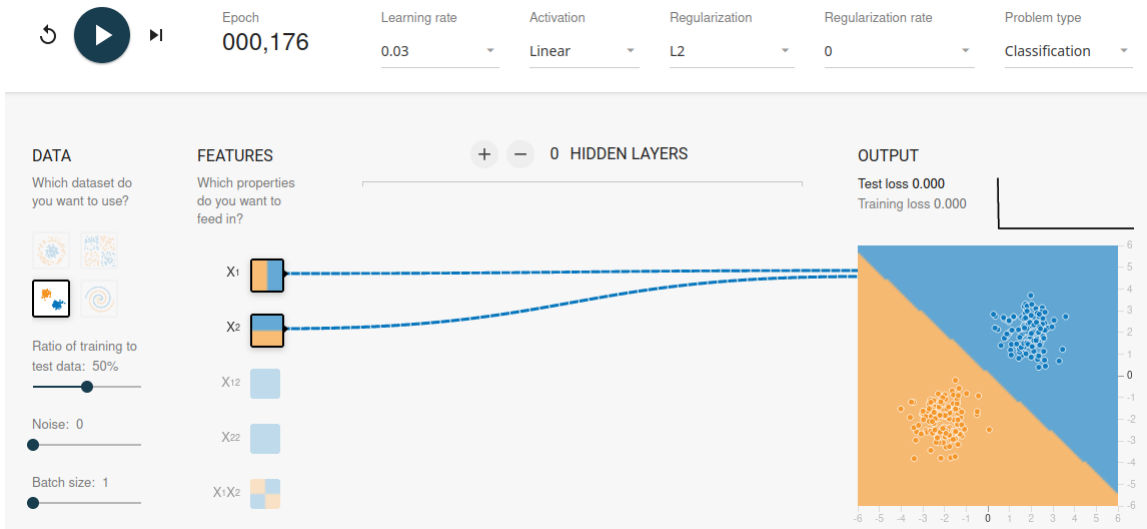
Log-linear multi-class classification

Representations

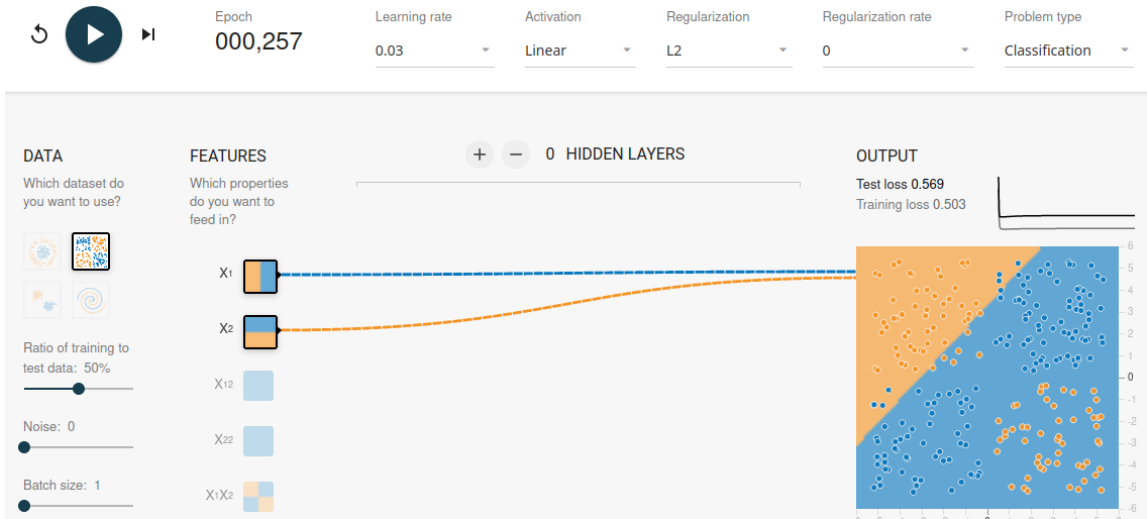
From multi-dimensional linear transformation to probabilities

Loss function for softmax

Stacking transformations and non-linearity



**Figure 2:** Linear model can tackle only linearly-separable problems (<http://playground.tensorflow.org>)



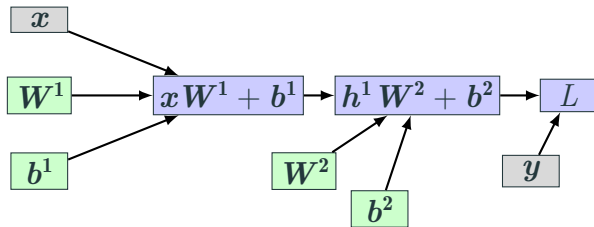
**Figure 3:** Linear model can tackle only linearly-separable problems (<http://playground.tensorflow.org>)



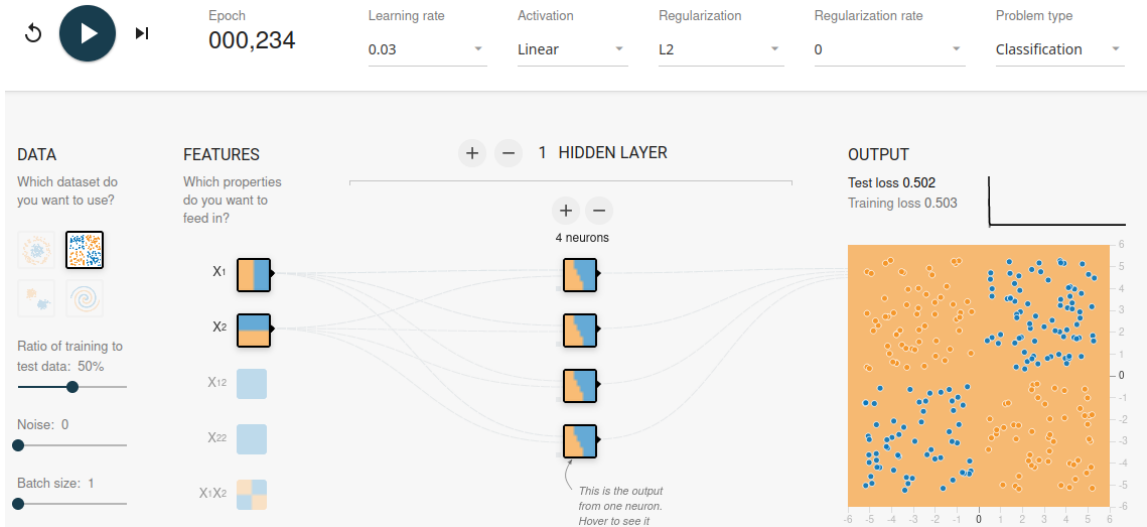
## Stacking linear layers on top of each other — still linear!

$$\mathbf{x} \in \mathbb{R}^{d_{in}} \quad \mathbf{W}^1 \in \mathbb{R}^{d_{in} \times d_1} \quad \mathbf{b}^1 \in \mathbb{R}^{d_1} \quad \mathbf{W}^2 \in \mathbb{R}^{d_1 \times d_{out}} \quad \mathbf{b}^2 \in \mathbb{R}^{d_{out}}$$

$$f(\mathbf{x}) = (\mathbf{x}\mathbf{W}^1 + \mathbf{b}^1)\mathbf{W}^2 + \mathbf{b}^2$$



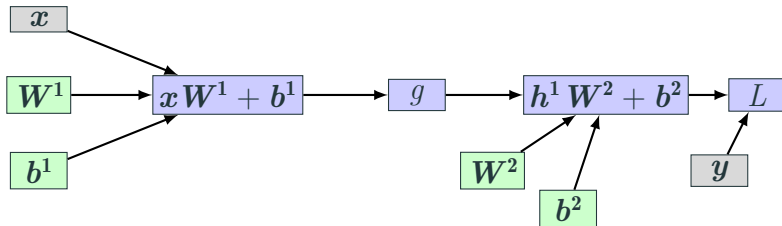
**Figure 4:** Computational graph; green circles are trainable parameters, gray are constant inputs



**Figure 5:** Linear hidden layers do not help  
(<http://playground.tensorflow.org>)

## Adding non-linear function $g : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_1}$

$$f(x) = g(xW^1 + b^1)W^2 + b^2$$

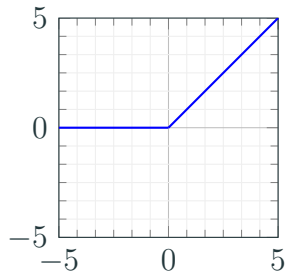


**Figure 6:** Computational graph; green circles are trainable parameters, gray are constant inputs

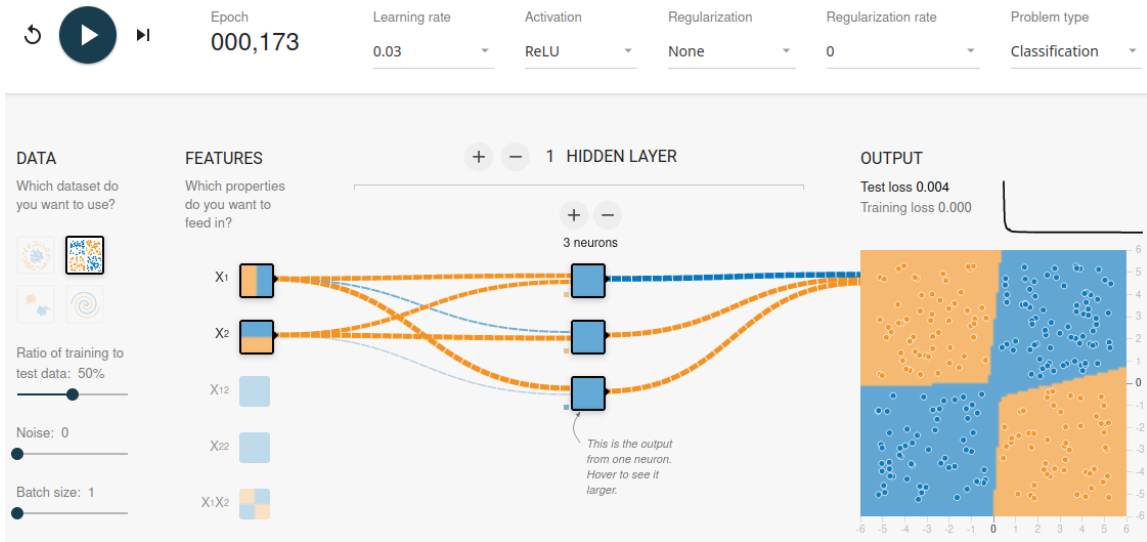
## Non-linear function $g$ : Rectified linear unit (ReLU) activation

$$\text{ReLU}(z) = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

or  $\text{ReLU}(z) = \max(0, x)$

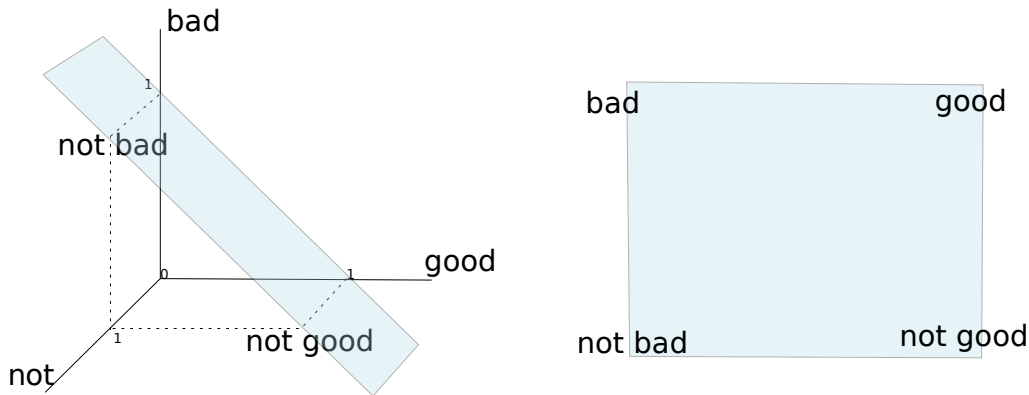


**Figure 7:** ReLU function



**Figure 8:** XOR solvable with, e.g., ReLU  
(<http://playground.tensorflow.org>)

# XOR example in super-simplified sentiment classification



**Figure 9:**  $V = \{\text{not}, \text{bad}, \text{good}\}$ , binary features  $\in \{0, 1\}$

# Recap

---

Where we finished last time

Log-linear multi-class classification

Representations

From multi-dimensional linear transformation to probabilities

Loss function for softmax

Stacking transformations and non-linearity

## Take aways

- Binary classification as a linear function of words and a sigmoid
- Binary cross-entropy (logistic) loss
- Training as minimizing the loss using minibatch SGD and backpropagation
- Stacking layers and non-linear functions: MLP (Multi-Layer Perceptron)
- ReLU as a go-to activation function in NLP



# License and credits

Licensed under Creative Commons  
Attribution-ShareAlike 4.0 International  
(CC BY-SA 4.0)



## Credits

Ivan Habernal

Content from ACL Anthology papers licensed under CC-BY  
<https://www.aclweb.org/anthology>