

Natural Language Processing with Deep Learning

RUHR
UNIVERSITÄT
BOCHUM

RUB

Lecture 3 — Backpropagation and binary text classification

Prof. Dr. Ivan Habernal

October 30, 2025

www.trusthlt.org

Trustworthy Human Language Technologies Group (TrustHLT)

Ruhr University Bochum & Research Center Trustworthy Data Science and Security



CENTER FOR TRUSTWORTHY
DATA SCIENCE AND SECURITY

This lecture

- Recap: computational graph
- Backpropagation
- Binary text classification
- Log-linear models, Cross-entropy loss, Stochastic gradient descent, etc.

Computational graph

- DAG — directed acyclic graph (not necessarily a tree!)
- Each node — a differentiable function with arguments
- Leaves — variables (e.g., a , b) or constants
- Arrows — Function composition

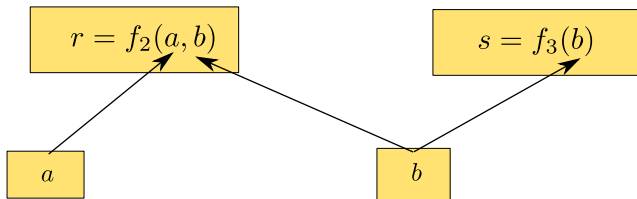
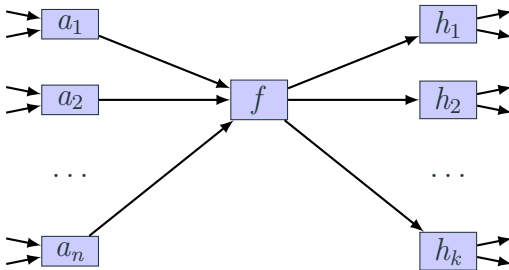


Figure 1: r , s are parents of b ; a , b are children (arguments) of r

Generic node in a computational graph



Adapted from J. Kun (2020). **A Programmer's Introduction to Mathematics**. 2nd ed., p. 265

Figure 2: A generic node of a computation graph. Node f has many inputs, its output feeds into many nodes, and each of its inputs and outputs may also have many inputs and outputs.

Backpropagation

- 1 Backpropagation
- 2 Text classification
- 3 Numerical representation of natural language text
- 4 Binary text classification

Working example

$$e = (a + b)(b + 1)$$

Compute gradient wrt. a and b

Working example

$$e = (a + b)(b + 1)$$

Compute gradient wrt. a and b

This one is easy by hand, but that's not the point

$$e = (a + b)(b + 1) = ab + a + b^2 + b$$

$$\frac{\partial e}{\partial a} = b + 1 \quad \frac{\partial e}{\partial b} = a + 2b + 1$$

Add some intermediate variables and function names

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$

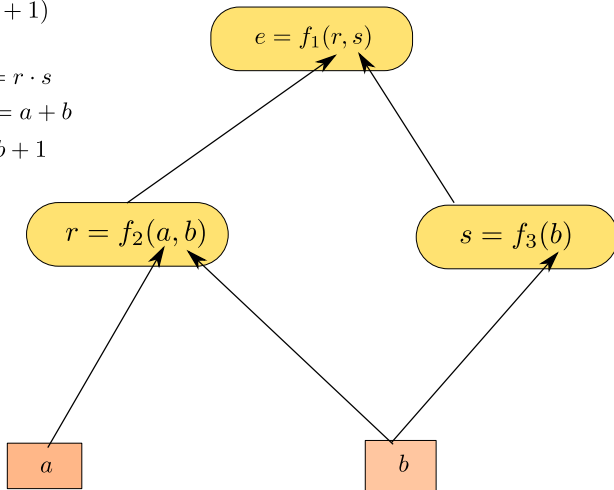
Build computational graph and evaluate (forward step)

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



Important: a, b will be some concrete real numbers, therefore r, s, e will be concrete real numbers too!

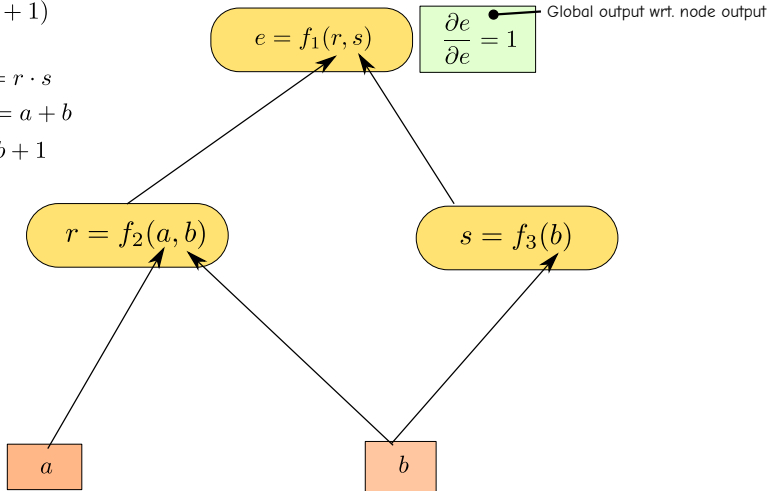
Goal: $\frac{\partial e}{\partial a}$ and $\frac{\partial e}{\partial b}$ (gradient), but let's do $\frac{\partial e}{\partial \star}$ for every node

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



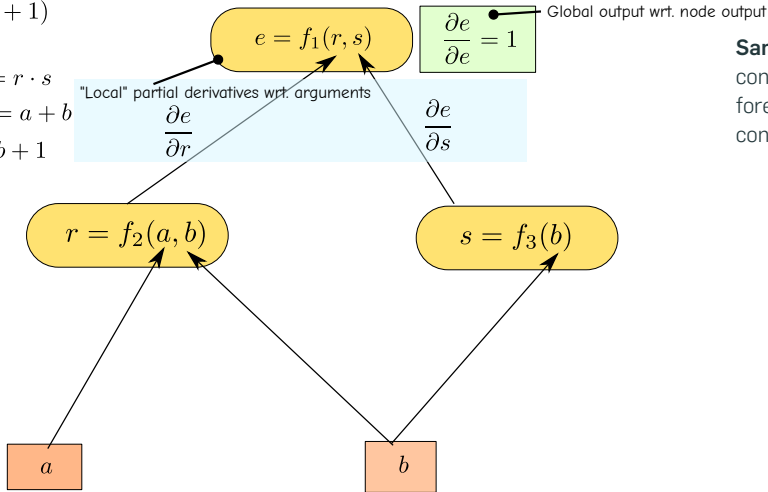
Since $e = r \cdot s$, partial derivatives are easy: $\frac{\partial e}{\partial r} = s$ and $\frac{\partial e}{\partial s} = r$

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



Sanity check: r, s are some concrete real numbers, therefore $\frac{\partial e}{\partial r}$ and $\frac{\partial e}{\partial s}$ will be concrete real numbers too!

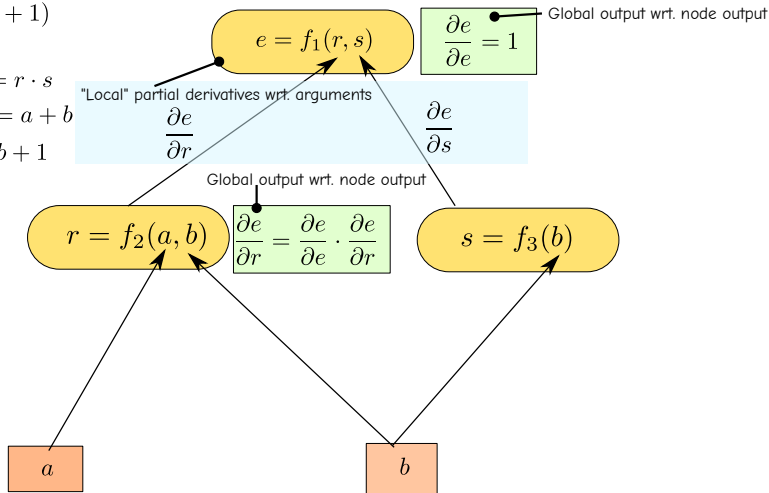
Proceed to next child r and compute $\frac{\partial e}{\partial r}$ – use chain rule!

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



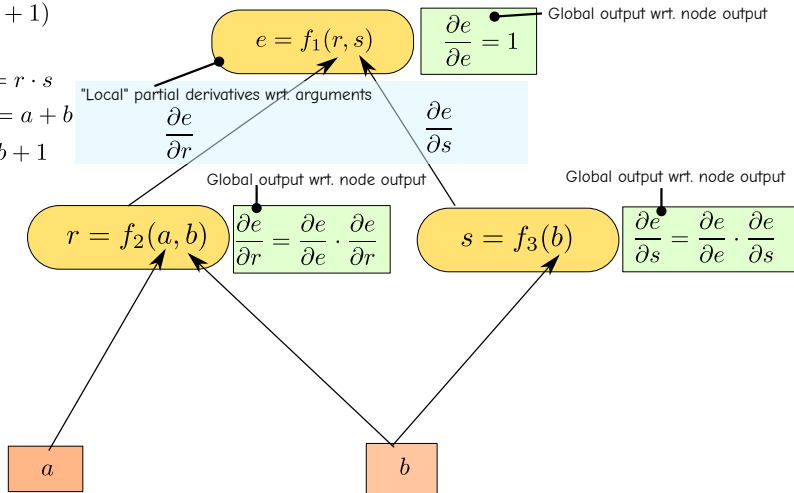
Proceed to next child s and compute $\frac{\partial e}{\partial s}$ – use chain rule!

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



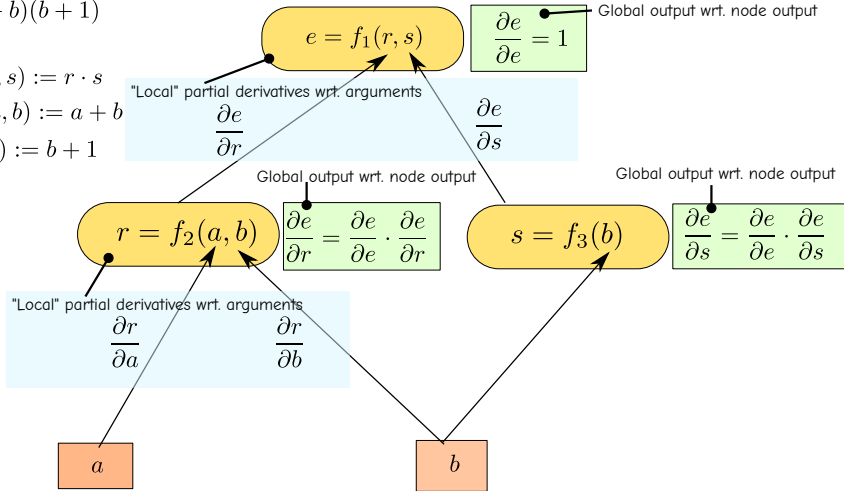
Since $r = a + b$, partial derivatives are easy: $\frac{\partial r}{\partial a} = 1$ and $\frac{\partial r}{\partial b} = 1$

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



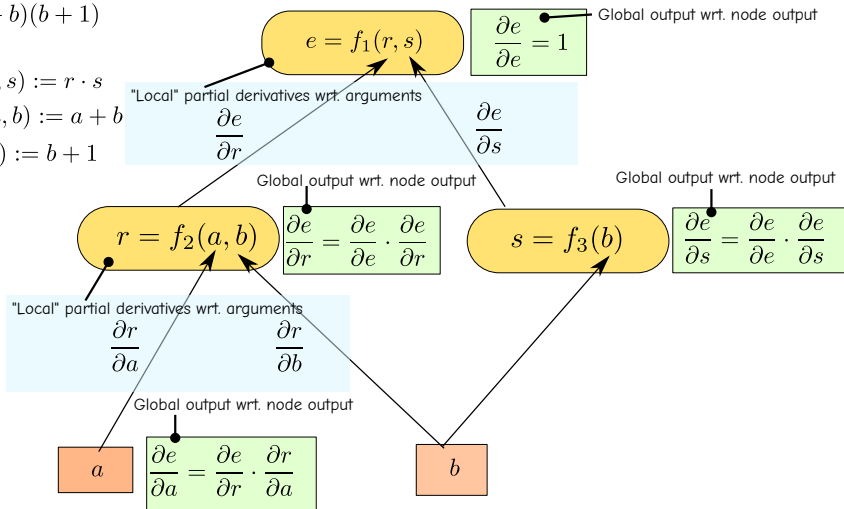
Proceed to next child a and compute $\frac{\partial e}{\partial a}$ – use chain rule!

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



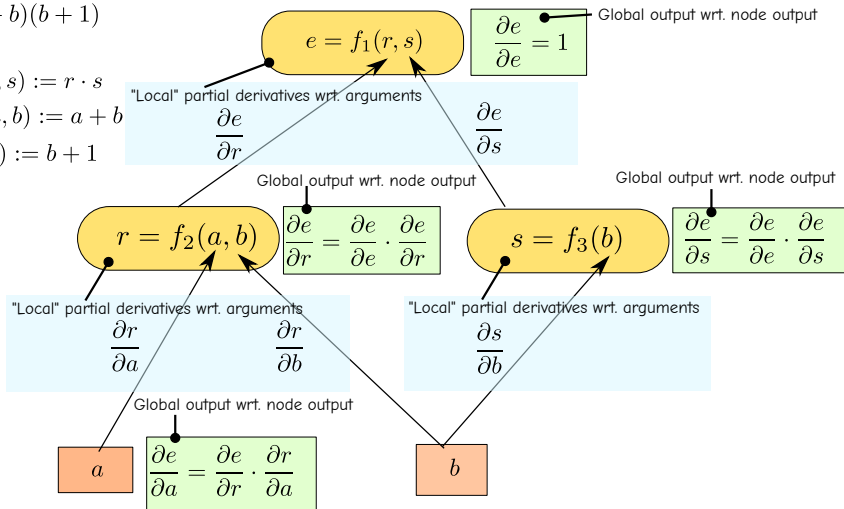
Since $s = b + 1$, partial derivatives are easy: $\frac{\partial s}{\partial b} = 1$

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



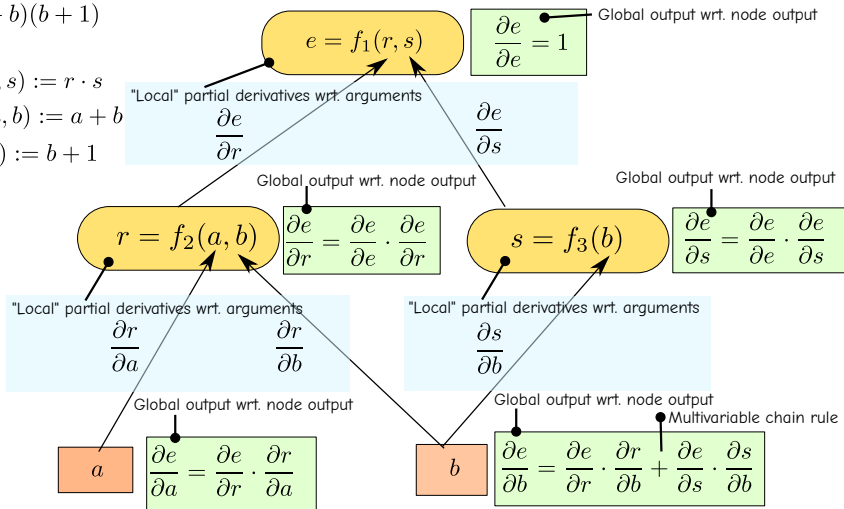
Proceed to b and compute $\frac{\partial e}{\partial b}$ – use multivariate chain rule!

$$e = (a + b)(b + 1)$$

$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



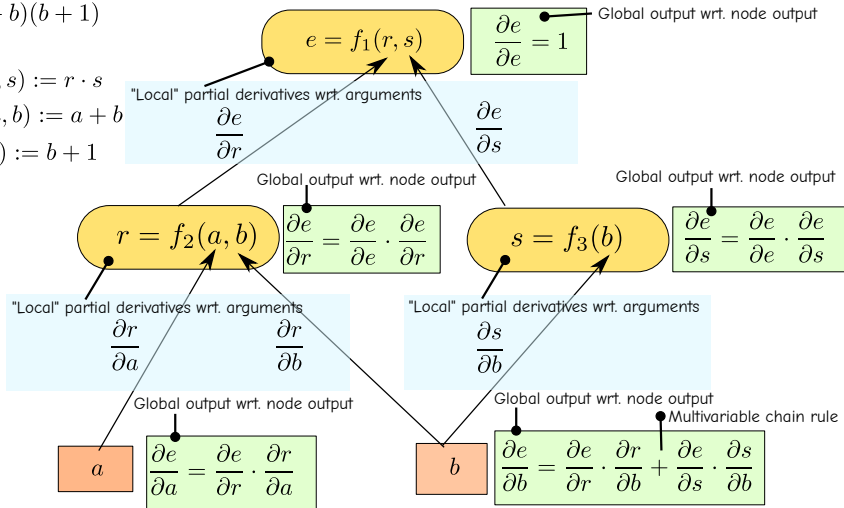
Goal: $\nabla e = \left(\frac{\partial e}{\partial a}; \frac{\partial e}{\partial b} \right)$ — we computed it for concrete a and b !

$$e = (a + b)(b + 1)$$

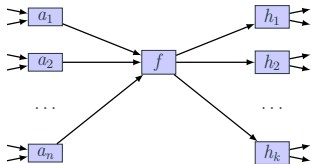
$$e = f_1(r, s) := r \cdot s$$

$$r = f_2(a, b) := a + b$$

$$s = f_3(b) := b + 1$$



Generic node in a computational graph $f(a_1, \dots, a_n)$



Assuming the graph is a function $e = g(\dots)$, we compute

$$\frac{\partial e}{\partial f} = \sum_{i=1}^k \frac{\partial e}{\partial h_i} \cdot \frac{\partial h_i}{\partial f}$$

and

$$\frac{\partial f}{\partial a_i} \quad \text{for } a_i, \dots, a_n$$

What each node must implement?

For example a function $s = f(a, b, c, d)$

- How to compute the output value s (given the parameters a, b, c, d)
- How to compute partial derivatives wrt. the parameters, i.e. $\frac{\partial s}{\partial a}, \frac{\partial s}{\partial b}, \frac{\partial s}{\partial c}, \frac{\partial s}{\partial d}$

Backpropagation

- Forward computation: Compute all nodes' output (and cache it)
- Backward computation (Backprop): Compute the overall function's partial derivative with respect to each node

Ordering of the computations? Recursively or build a graph's topology upfront and iterate

Backpropagation: Recap

- We can express any arbitrarily complicated function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ as a computational graph
- For computing the gradient ∇f at a concrete point (x_1, x_2, \dots, x_n) we run the forward pass and backprop
- When caching each node's intermediate output and partial derivatives, we avoid repeating computations \rightarrow efficient algorithm

Text classification

- 1 Backpropagation
- 2 Text classification**
- 3 Numerical representation of natural language text
- 4 Binary text classification

What are we going to achieve

Example task: Binary sentiment classification into positive and negative

What are we going to achieve

Example task: Binary sentiment classification into positive and negative

Recall the IMDB dataset

We will learn a simple yet powerful supervised machine learning model

- Known as logistic regression, maximum entropy classifier
- In fact, it is a single-layer neural network
- An essential important building block of deep neural networks

High-dimensional linear functions

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

$$f(\mathbf{x}) \text{ or } f(\mathbf{x}; \underbrace{\mathbf{W}, \mathbf{b}}_{\text{Explicit parameters}}) = \mathbf{x}\mathbf{W} + \mathbf{b}$$

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

Vector \mathbf{x} is the **input**, matrix \mathbf{W} and vector \mathbf{b} are the **parameters** — typically denoted $\Theta = \mathbf{W}, \mathbf{b}$

High-dimensional linear functions

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

$$f(\mathbf{x}) \text{ or } f(\mathbf{x}; \underbrace{\mathbf{W}, \mathbf{b}}_{\text{Explicit parameters}}) = \mathbf{x}\mathbf{W} + \mathbf{b}$$

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

Vector \mathbf{x} is the **input**, matrix \mathbf{W} and vector \mathbf{b} are the **parameters** — typically denoted $\Theta = \mathbf{W}, \mathbf{b}$

Goal of learning

Find \mathbf{W} and \mathbf{b} such that the function behaves as intended on a collection of input values $\mathbf{x}_{1:k} = \mathbf{x}_1, \dots, \mathbf{x}_k$ and the corresponding desired outputs $\mathbf{y}_{1:k} = \mathbf{y}_1, \dots, \mathbf{y}_k$

Linear real-valued function in binary classification

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

$$f(\mathbf{x}) \text{ or } f(\mathbf{x}; \underbrace{\mathbf{w}, b}) = \mathbf{x} \cdot \mathbf{w} + b$$

Explicit parameters

However, for binary text classification

- Our input is in the form of a natural language text
- Our labels are two categories, e.g., positive and negative

Linear real-valued function in binary classification

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

$$f(\mathbf{x}) \text{ or } f(\mathbf{x}; \underbrace{\mathbf{w}, b}) = \mathbf{x} \cdot \mathbf{w} + b$$

Explicit parameters

However, for binary text classification

- Our input is in the form of a natural language text
- Our labels are two categories, e.g., positive and negative

Let's start with the labels

Very easy: Just arbitrarily map the categories into 0 and 1 (e.g., negative = 0, positive = 1)

Numerical representation of natural language text

- 1 Backpropagation
- 2 Text classification
- 3 Numerical representation of natural language text**
- 4 Binary text classification

Goal: Transform text into a fixed-size vector of real numbers

What's our setup:

$$f(\mathbf{x}) : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R} \quad f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w} + b$$

What we need:

$$\mathbf{x} \in \mathbb{R}^{d_{in}}$$

Goal: Transform text into a fixed-size vector of real numbers

What's our setup:

$$f(\mathbf{x}) : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R} \quad f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w} + b$$

What we need:

$$\mathbf{x} \in \mathbb{R}^{d_{in}}$$

What we have:

One of my favorite movies ever, The Shawshank Redemption is a modern day classic as it tells the story of two inmates who become friends and find solace over the years in which this movie takes place. Based on a Stephen King novel, ...

What is a “word”?

A matter of debate among linguists

Very simplistic definition: words are sequences of letters separated by whitespace

But: `dog`, `dog?`, `dog.`, and `dog)` would be different words

Better: words separated by whitespace or punctuation

A process called **tokenization** splits text into tokens based on whitespace and punctuation

- English: the job of the tokenizer is quite simple
- Hebrew, Arabic: sometimes without whitespace
- Chinese: no whitespaces at all

Y. Goldberg (2017). **Neural Network Methods for Natural Language Processing**. Morgan & Claypool

Tokens

Symbols `cat` and `Cat` have the same meaning, but are they the same word?

Something like `New York`, is it two words, or one?

- We distinguish between words and tokens
- We refer to the output of a tokenizer as a token, and to the meaning-bearing units as words

Keep in mind

We use the term **word** very loosely, and take it to be interchangeable with **token**.

In reality, the story is more complex than that.

Vocabulary

We build a fix-sized static **vocabulary** (e.g., by tokenizing training data)

- Typical sizes: 20,000 – 100,000 words

Each word has a unique fixed index

$$V = \left(a_1 \quad \text{abandon}_2 \quad \dots \quad \text{cat}_{852} \quad \dots \quad \text{zone}_{2,999} \quad \text{zoo}_{3,000} \right)$$

(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

Example: a cat sat \rightarrow a, cat, sat

$$V = \begin{pmatrix} a_1 & abandon_2 & \dots & cat_{852} & \dots & zone_{2,999} & zoo_{3,000} \end{pmatrix}$$

$$a = \mathbf{x}^{D[1]} = \begin{pmatrix} 1_1 & 0_2 & 0_3 & \dots & 0_{2,999} & 0_{3,000} \end{pmatrix}$$

$$cat = \mathbf{x}^{D[2]} = \begin{pmatrix} 0_1 & \dots & 1_{852} & \dots & 0_{2,999} & 0_{3,000} \end{pmatrix}$$

$$sat = \mathbf{x}^{D[3]} = \begin{pmatrix} 0_1 & \dots & 1_{2,179} & \dots & 0_{2,999} & 0_{3,000} \end{pmatrix}$$

Averaged bag-of-words example: $x \in \mathbb{R}^{3,000}$

Example: a cat sat \rightarrow a, cat, sat

$$a = \mathbf{x}^{D_{[1]}} = \begin{pmatrix} 1_1 & 0_2 & 0_3 & \dots & 0_{2,999} & 0_{3,000} \end{pmatrix}$$

$$\text{cat} = \mathbf{x}^{D_{[2]}} = \begin{pmatrix} 0_1 & \dots & 1_{852} & \dots & 0_{2,999} & 0_{3,000} \end{pmatrix}$$

$$\text{sat} = \mathbf{x}^{D_{[3]}} = \begin{pmatrix} 0_1 & \dots & 1_{2,179} & \dots & 0_{2,999} & 0_{3,000} \end{pmatrix}$$

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

$$= \begin{pmatrix} 0.33_1 & 0_2 & \dots & 0_{851} & 0.33_{852} & 0_{853} & \dots & 0.33_{2,179} & \dots & 0_{3,000} \end{pmatrix}$$

Out-of-vocabulary (UNK) tokens

Words in a language are very unevenly distributed

- There is always a large 'tail' of rare words

When building the vocabulary, use the most frequent words, all others represented by an unknown token (UNK or OOV)

Example vocabulary, most common 3,000 words and UNK

$$V = (a_1 \quad \text{abandon}_2 \quad \dots \quad \text{zone}_{2,999} \quad \text{zoo}_{3,000} \quad \text{UNK}_{3,001})$$

- In machine translation, how to translate the UNK word?

P. Koehn (2020). **Neural Machine Translation.** (not freely available). Cambridge University Press

Subword units: Byte-pair encoding

- 1 The words in the corpus are split into characters (marking original spaces with a special space character) — this is the initial vocabulary V
- 2 The most frequent pair of characters is merged and added to V
- 3 Repeat 2 for a fixed given number of times
- 4 Each of these steps increases V by one, beyond the original inventory of single characters

When done over large corpora with multiple languages and writing systems, BPE prevents OOV!

Byte-pair encoding example on a toy corpus (part 1)

t h i s _ f a t _ c a t _ w i t h _ t h e _ h a
t _ i s _ i n _ t h e _ c a v e _ o f _ t h e _
t h i n _ b a t

Most frequent: t h (6 times), merge into a single token

th i s _ f a t _ c a t _ w i t h _ t h e _ h a t _
i s _ i n _ t h e _ c a v e _ o f _ t h e _ t h i n _
_ b a t

Most frequent: a t (4 times), merge into a single token

th i s _ f a t _ c a t _ w i t h _ t h e _ h a t _ i
s _ i n _ t h e _ c a v e _ o f _ t h e _ t h i n _
b a t

Byte-pair encoding example on a toy corpus (part 2)

th i s _ f at _ c at _ w i th _ th e _ h at _ i
s _ i n _ th e _ c a v e _ o f _ th e _ th i n _
b at

At the end of this process,
the most frequent words
will emerge as single tokens,
while rare words consist of
still unmerged subwords

Most frequent: th e (3 times), merge into a single token

th i s _ f at _ c at _ w i th _ the _ h at _ i s
_ i n _ the _ c a v e _ o f _ the _ th i n _ b
at

$V =$

{t, h, i, s, _, f, a, c, w, e, n, v, o, f, b, th, at, the}

SentencePiece: A variant of byte pair encoding

Byte-pair example. Word splits indicated with @@.

```
[the] [relationship] [between] [Obama]  
[and] [Net@@] [any@@] [ahu] [is] [not]  
[exactly] [friendly] [.]
```

SentencePiece escapes the whitespace with _ and tokenizes the input into an arbitrary subword sequence

SentencePiece example of "Hello world."

```
[Hello] [_wor] [ld] [.]
```

Lossless tokenization — all the information to reproduce the normalized text is preserved

T. Kudo and J. Richardson (2018).
"SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing". In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Brussels, Belgium: Association for Computational Linguistics, pp. 66–71

Recap: Transform text into a fixed-size vector of real numbers

What's our setup:

$$f(\mathbf{x}) : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R} \quad f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w} + b$$

What we need:

$$\mathbf{x} \in \mathbb{R}^{d_{in}}$$

What we have:

*One of my favorite movies ever, The Shawshank Redemption
is a modern day classic ...*

Simple solution:

- Bag-of-words (tokenized), $d_{in} = |V|$

Binary text classification

- 1 Backpropagation
- 2 Text classification
- 3 Numerical representation of natural language text
- 4 Binary text classification**

Binary text classification

Binary classification as a function

Linear function and its derivatives

We have this linear function

$$f(\mathbf{x}) : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R} \quad f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w} + b = x_{[1]} w_{[1]} + \dots + x_{[d_{in}]} w_{[d_{in}]} + b$$

Linear function and its derivatives

We have this linear function

$$f(\mathbf{x}) : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R} \quad f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w} + b = \mathbf{x}_{[1]} \mathbf{w}_{[1]} + \dots + \mathbf{x}_{[d_{in}]} \mathbf{w}_{[d_{in}]} + b$$

Derivatives wrt. parameters w and b

$$\frac{df}{d\mathbf{w}_{[i]}} = \mathbf{x}_{[i]} \quad \frac{df}{db} = 1$$

Non-linear mapping to $[0, 1]$

We have this linear function

$$f(\mathbf{x}) : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R} \quad f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w} + b$$

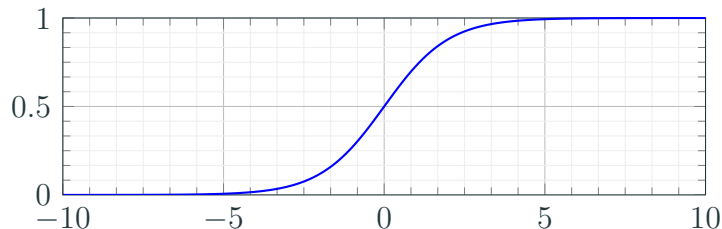
which has an unbounded range $(-\infty, +\infty)$

However, each example's label is $y \in \{0, 1\}$

Sigmoid (logistic) function

Sigmoid function $\sigma(t) : \mathbb{R} \rightarrow \mathbb{R}$

$$\sigma(t) = \frac{\exp(t)}{\exp(t) + 1} = \frac{1}{1 + \exp(-t)}$$



Symmetric function, range of $\sigma(t) \in [0, 1]$,

Sigmoid $\sigma(t) = \frac{1}{1+\exp(-t)}$

Derivative of sigmoid wrt. its input

$$\begin{aligned}\frac{d\sigma}{dt} &= \frac{\exp(t) \cdot (1 + \exp(t)) - \exp(t) \cdot \exp(t)}{(1 + \exp(t))^2} \\ &= \dots \\ &= \sigma(t) \cdot (1 - \sigma(t))\end{aligned}$$

Our binary text classification function

Linear function through sigmoid — log-linear model

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

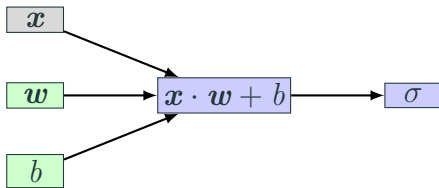


Figure 3: Computational graph; green nodes are trainable parameters, gray are inputs

Decision rule of log-linear model

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

- Prediction = 1 if $\hat{y} > 0.5$
- Prediction = 0 if $\hat{y} < 0.5$

Natural interpretation: Conditional probability of prediction
= 1 given the input \mathbf{x}

$$\sigma(f(\mathbf{x})) = \Pr(\text{prediction} = 1 | \mathbf{x})$$

$$1 - \sigma(f(\mathbf{x})) = \Pr(\text{prediction} = 0 | \mathbf{x})$$

Binary text classification

Finding the best model's parameters

Binary cross-entropy loss (logistic loss)

$$L_{\text{logistic}} = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

Partial derivative wrt. input \hat{y}

$$\frac{dL_{\text{Logistic}}}{d\hat{y}} = - \left(\frac{y}{\hat{y}} - \frac{1 - y}{1 - \hat{y}} \right) = - \frac{y - \hat{y}}{\hat{y}(1 - \hat{y})}$$

Full computational graph

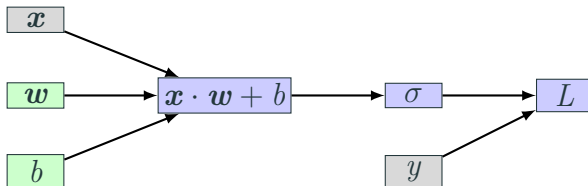


Figure 4: Computational graph; green nodes are trainable parameters, gray are constant inputs

How can we minimize this loss function wrt. w and b ?

Full computational graph

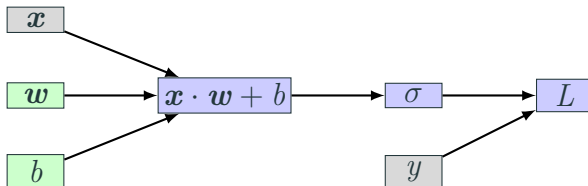


Figure 4: Computational graph; green nodes are trainable parameters, gray are constant inputs

How can we minimize this loss function wrt. w and b ?

- Recall: (a) Gradient descent and (b) backpropagation

(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

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

Recap

- 1 Backpropagation
- 2 Text classification
- 3 Numerical representation of natural language text
- 4 Binary text classification

Take aways

- Tokenization is tricky
- Simplest representation of text as bag-of-word features
- 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

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>

Appendix

Mathematical notation

Scalars, vectors, matrices

Lowercase letters represent scalars: x, y, b

Bold lowercase letters represent vectors: $\mathbf{w}, \mathbf{x}, \mathbf{b}$

Bold uppercase letters represent matrices: \mathbf{W}, \mathbf{X}

Indexing

$[\cdot]$ as the index operator of vectors and matrices

$\mathbf{b}_{[i]}$ is the i -th element of vector \mathbf{b}

$\mathbf{W}_{[i,j]}$ is the i -th row, j -th column of matrix \mathbf{W}

Notation

Sequences

$\mathbf{x}_{1:n}$ is a sequence of vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$

$[\mathbf{v}_1; \mathbf{v}_2]$ is vector concatenation

Note! We use vectors as row vectors

$$\mathbf{x} \in \mathbb{R}^d$$

Example $d = 5$:

$$\mathbf{x} = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \end{pmatrix}$$

which is simply a list (1-d array) of numbers (1, 2, 3, 4, 5)

Multiplication example

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

Example: $\mathbf{y} = \mathbf{x}\mathbf{W} + \mathbf{b}$, $d_{in} = 3$, $d_{out} = 2$

$$\begin{pmatrix} \mathbf{x}_{[1]} & \mathbf{x}_{[2]} & \mathbf{x}_{[3]} \end{pmatrix} \begin{pmatrix} \mathbf{W}_{[1,1]} & \mathbf{W}_{[1,2]} \\ \mathbf{W}_{[2,1]} & \mathbf{W}_{[2,2]} \\ \mathbf{W}_{[3,1]} & \mathbf{W}_{[3,2]} \end{pmatrix} + \begin{pmatrix} \mathbf{b}_{[1]} & \mathbf{b}_{[2]} \end{pmatrix} = \begin{pmatrix} \mathbf{y}_{[1]} & \mathbf{y}_{[2]} \end{pmatrix}$$

Mult. simplified with dot product $u \cdot v = \sum_i u_{[i]} v_{[i]}$

Example: $y = xW + b$, $d_{in} = 3$, $d_{out} = 1$

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

$$\begin{pmatrix} \mathbf{x}_{[1]} & \mathbf{x}_{[2]} & \mathbf{x}_{[3]} \end{pmatrix} \begin{pmatrix} \mathbf{W}_{[1,1]} \\ \mathbf{W}_{[2,1]} \\ \mathbf{W}_{[3,1]} \end{pmatrix} + b = y$$

Equivalent dot product: $y = \mathbf{x} \cdot \mathbf{w} + b$, $d_{in} = 3$, $d_{out} = 1$

$$\mathbf{x} \in \mathbb{R}^{d_{in}} \quad \mathbf{w} \in \mathbb{R}^{d_{out}} \quad b \in \mathbb{R}$$

$$\begin{pmatrix} \mathbf{x}_{[1]} & \mathbf{x}_{[2]} & \mathbf{x}_{[3]} \end{pmatrix} \cdot \begin{pmatrix} \mathbf{w}_{[1]} & \mathbf{w}_{[2]} & \mathbf{w}_{[3]} \end{pmatrix} + b = y$$

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>