

Text Classification with Multi-Layer Perceptrons

2025-26

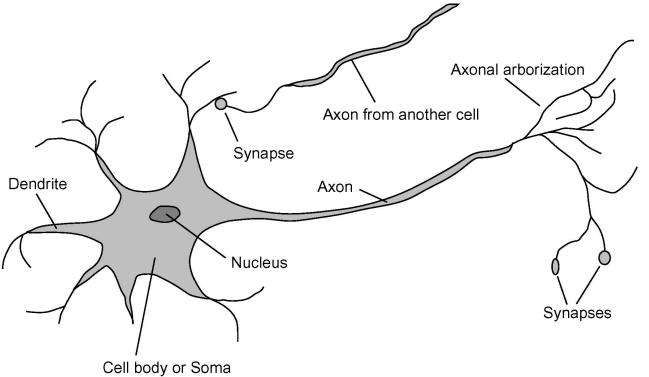
Ion Androutsopoulos

http://www.aueb.gr/users/ion/

Contents

- Natural and artificial neural networks (NNs).
- Multi-Layer Perceptrons (MLPs) and backpropagation.
- MLPs for text classification, regression, window-based token classification (e.g., for POS tagging, NER).
- Dropout, batch and layer normalization.
- Pre-training word embeddings with Word2Vec.
- Advice for training large neural networks.

Natural neural networks

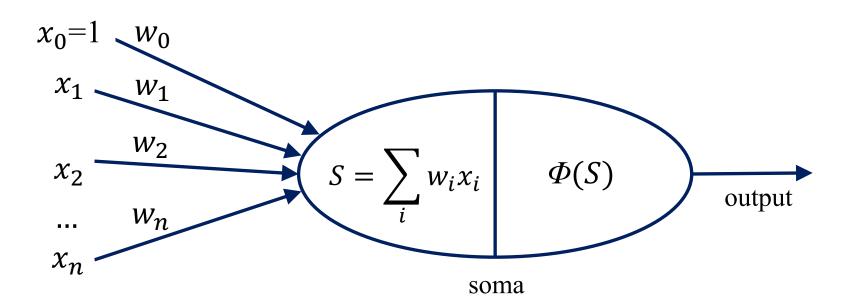


- Neuron: cell of the brain.
 - Cell body or soma: the main part, includes the nucleus.
 - **Dendrites** receive signals from other neurons.
 - Axon: transmits a single output to other neurons. Often much longer than the diameter of the soma.
 - **Synapses**: axon-dendrite interfaces, whose conductivities vary.
- Neural network: network of many neurons.

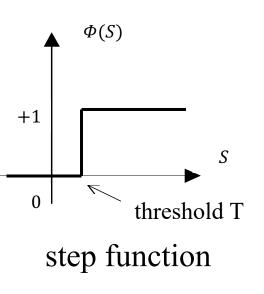
Artificial neural networks

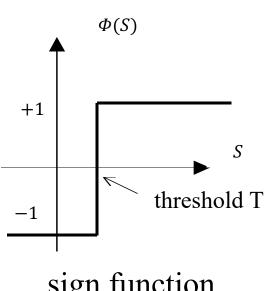
• Artificial neuron:

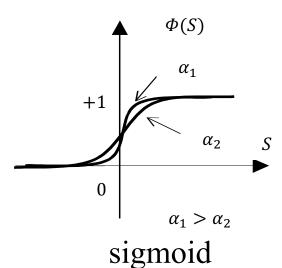
- Input: real variables.
- Input weights: real variables (roughly synapses).
- Soma: computes the weighted sum of the inputs, then applies an activation function to the sum.



Activation functions

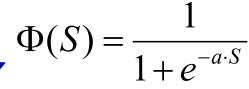






sign function

(logistic function)

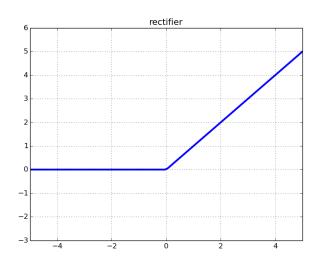


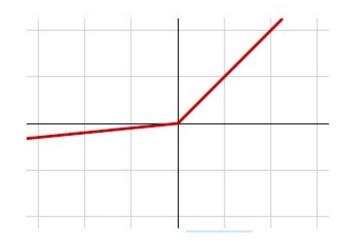
The sigmoid is differentiable.

+1tanh

The hyperbolic tangent (tanh) is very similar to the sigmoid, but with values from -1 to +1. Usually better, unless we really want values in (0, 1).

Activation functions – continued





Rectified Linear Unit (ReLU)

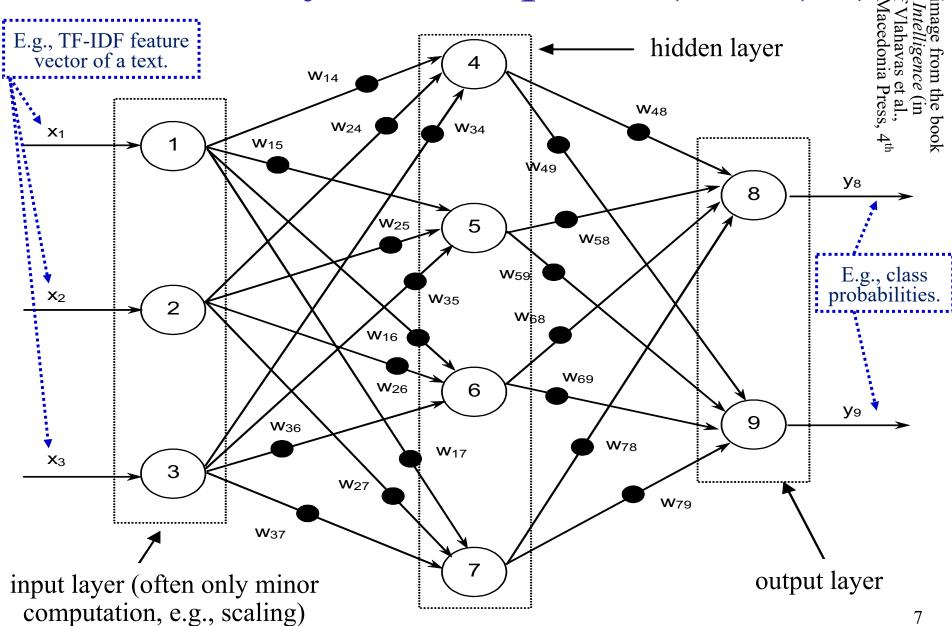
$$relu(x) = \begin{cases} 0 & \text{if } x < 0, \\ x & \text{otherwise} \end{cases}$$

$$leakyrelu(x) = \begin{cases} ax \text{ if } x < 0, \\ x \text{ otherwise} \end{cases}$$

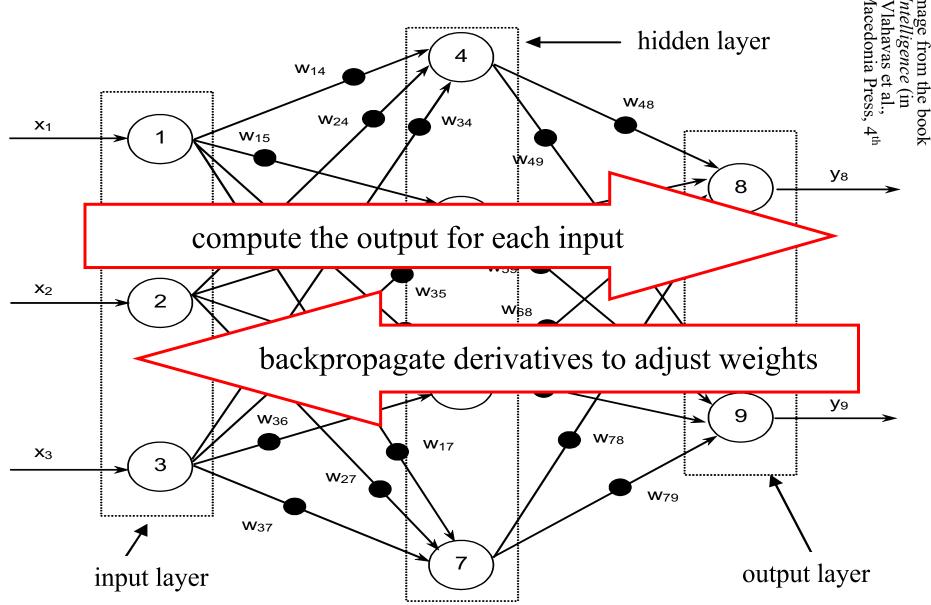
 α is a small positive.

ReLU and variants are popular choices.

Multi-layer Perceptron (MLP)



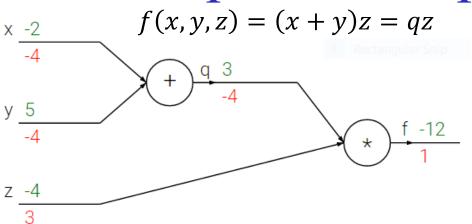
Backpropagation



Backpropagation

- Initialize all the weights to small random values.
 - \circ E.g., sample from a **zero-centered Gaussian** with **small** σ .
 - o Better initializations exist (see DL course).
 - o **Normalize** the **features** too (see "Important tricks" of Part 2).
- In each epoch, for each training example (or mini-batch):
 - o Compute the **output** $< o_1, o_2, ... >$ for the training example.
 - o For each weight w_{ij} , compute $\frac{\partial E}{\partial w_{ij}}$, where E the loss on the training example. We compute derivatives right to left.
 - 0 Update each weight as: $w_{ij} \leftarrow w_{ij} \eta \cdot \frac{\partial E}{\partial w_{ij}}$, i.e., for all the weights together: $W \leftarrow W \eta \cdot \nabla_W E$
 - o Hence, we use **SGD** (or variants). **No guarantee** SGD will find the **best solution**, but it (often) works in practice!

Example of computation graph

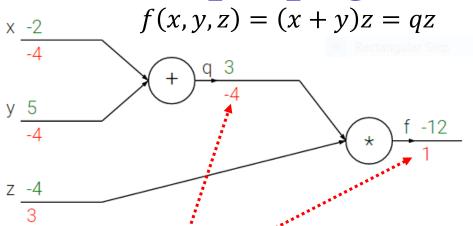


Example and figure from Stanford's "CNNs for Visual Recognition" (2016, F.-F. Li, A. Karpathy, J. Johnson) http://cs231n.github.io/optimization-2/

- Forward pass: $\langle x, y, z \rangle = \langle -2, 5, -4 \rangle$, q = 3, f = -12
- Imagine we wish to minimize f using SGD.
 - o In a more realistic scenario, f would be a **loss function**, and $\langle x, y, z \rangle$ the weights vector.

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} \leftarrow \begin{bmatrix} x \\ y \\ z \end{bmatrix} - \eta \nabla f(x, y, z) = \begin{bmatrix} x \\ y \\ z \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} & \frac{\partial f}{\partial z} \\ \frac{\partial f}{\partial y} & \frac{\partial f}{\partial z} \end{bmatrix}$$
We need:
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

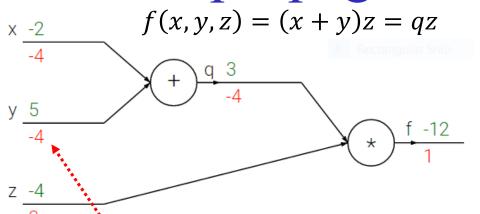
Backpropagation in the graph



Example and figure from Stanford's "CNNs for Visual Recognition" (2016, F.-F. Li, A. Karpathy, J. Johnson) http://cs231n.github.io/optimization-2/

- Backpropagation: We compute derivatives right to left.
 - $\circ \frac{\partial f}{\partial f} = 1 \text{ by definition.}$
 - $\circ \frac{\partial f}{\partial a} = z$. And for this $\langle x, y, z \rangle$ input, z = -4.
 - O During the **forward pass**, we need to **save the outputs of all the nodes** (e.g., here we need the value of z).

Backpropagation in the graph

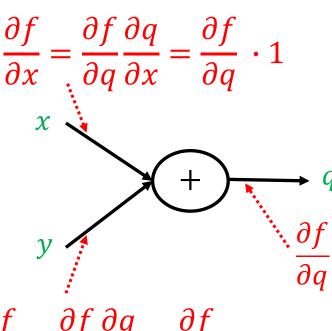


Example and figure from Stanford's "CNNs for Visual Recognition" (2016, F.-F. Li, A. Karpathy, J. Johnson) http://cs231n.github.io/optimization-2/

- **Backpropagation**: We compute derivatives right to left.
 - 1 by definition.
 - $\circ \frac{\partial f}{\partial g} = z$. And for this $\langle x, y, z \rangle$ input, z = -4.
 - $\circ \frac{\partial f}{\partial z} = q$. And for this $\langle x, y, z \rangle$ input, q = 3.

 - $\circ \frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = \frac{\partial f}{\partial q} \cdot 1. \text{ And here } \frac{\partial f}{\partial a} \text{ is } -4.$

Plug-and-play gates



$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y} = \frac{\partial f}{\partial q} \cdot 1$$

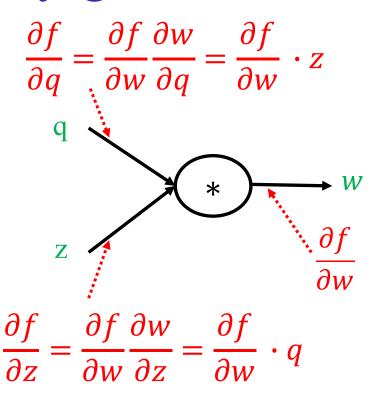
class PlusGate:

forward(x, y):

return x+y

backward($\frac{\partial f}{\partial q}$):

return
$$<\frac{\partial f}{\partial q}, \frac{\partial f}{\partial q}>$$



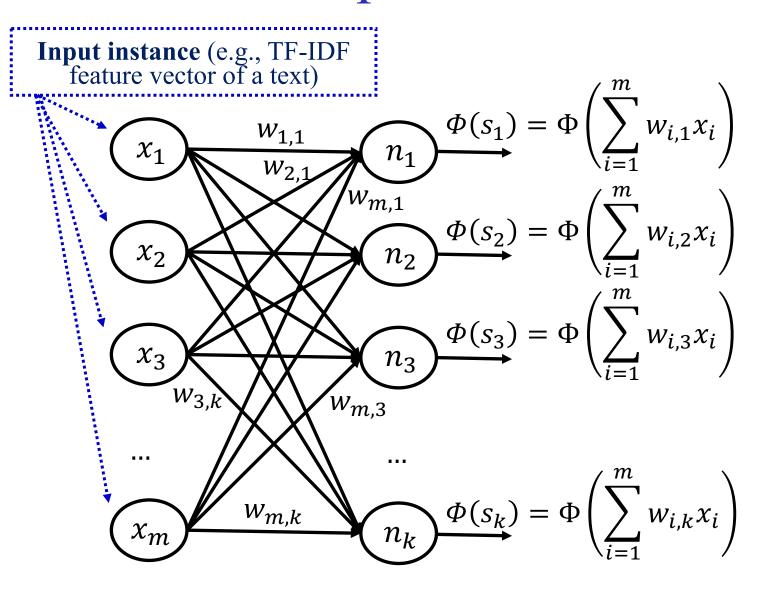
class StarGate:

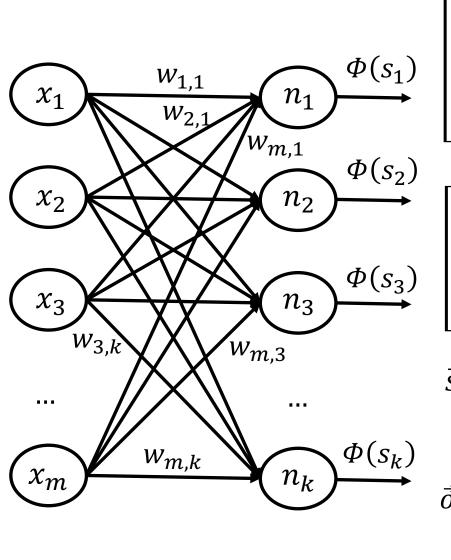
forward(q, z):

return q * z

backward($\frac{\partial f}{\partial w}$):

return $< \frac{\partial f}{\partial w} \cdot z, \frac{\partial f}{\partial w} \cdot q >$





$$\begin{bmatrix} S_1 \\ S_2 \\ S_3 \\ \dots \\ S_k \end{bmatrix} = \begin{bmatrix} w_{1,1}x_1 + w_{2,1}x_2 + \dots + w_{m,1}x_m \\ w_{1,2}x_1 + w_{2,2}x_2 + \dots + w_{m,2}x_m \\ w_{1,3}x_1 + w_{2,3}x_2 + \dots + w_{m,3}x_m \\ \dots \\ w_{1,k}x_1 + w_{2,k}x_2 + \dots + w_{m,k}x_m \end{bmatrix}$$

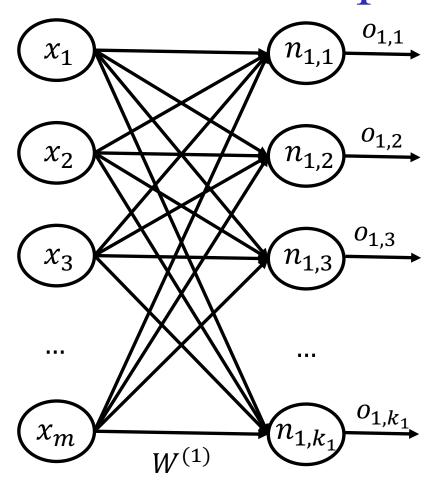
$$\begin{bmatrix} S_1 \\ S_2 \\ S_3 \\ \dots \\ S_k \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{2,1} & \dots & w_{m,1} \\ w_{1,2} & w_{2,2} & \dots & w_{m,2} \\ w_{1,3} & w_{2,3} & \dots & w_{m,3} \\ \dots & \dots & \dots & \dots \\ w_{1,k} & w_{2,k} & \dots & w_{m,k} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_m \end{bmatrix}$$

$$\vec{s} = W\vec{x}$$

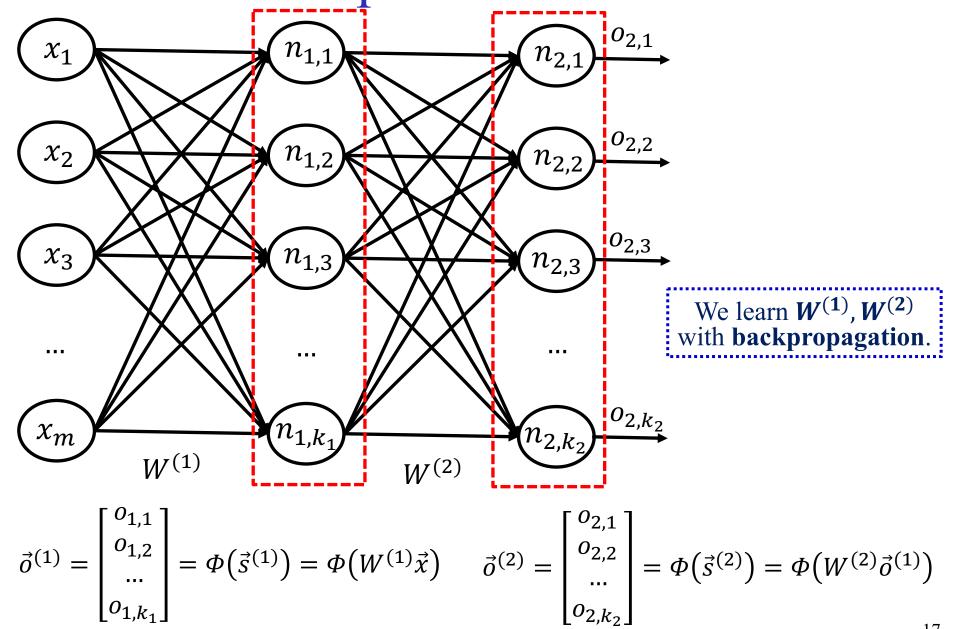
$$\vec{o} = \begin{bmatrix} o_1 \\ o_2 \\ o_3 \\ \dots \\ o_k \end{bmatrix} = \begin{bmatrix} \phi(s_1) \\ \phi(s_2) \\ \phi(s_3) \\ \dots \\ \phi(s_K) \end{bmatrix}$$

We learn **W** with **backpropagation**.

$$= \Phi(\vec{s}) = \Phi(W\vec{x})$$

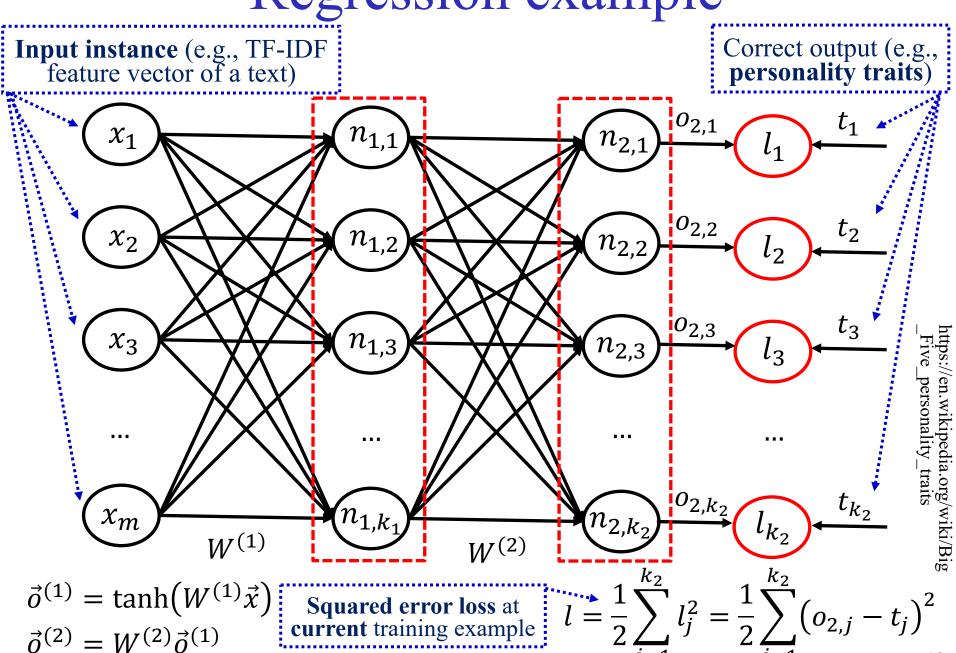


$$\vec{o}^{(1)} = \begin{bmatrix} o_{1,1} \\ o_{1,2} \\ \dots \\ o_{1,k_1} \end{bmatrix} = \Phi(\vec{s}^{(1)}) = \Phi(W^{(1)}\vec{x})$$



17

Regression example



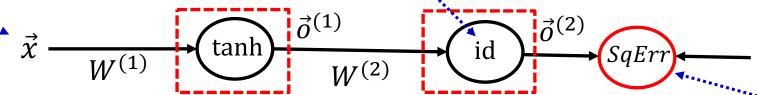
18

Regression example – more compact

Input instance (e.g., TF-IDF feature vector of a text)

$$\Phi(s) = s$$
 (identity)

Correct output (e.g., personality traits)

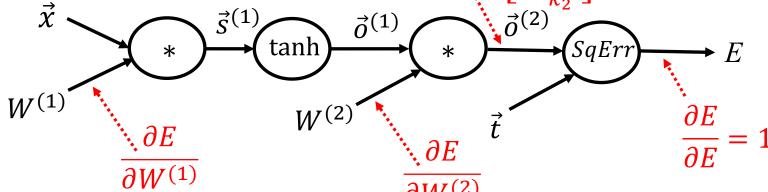


$$\vec{o}^{(1)} = \tanh(W^{(1)}\vec{x}) \quad \vec{o}^{(2)} = W^{(2)}\vec{o}^{(1)}$$

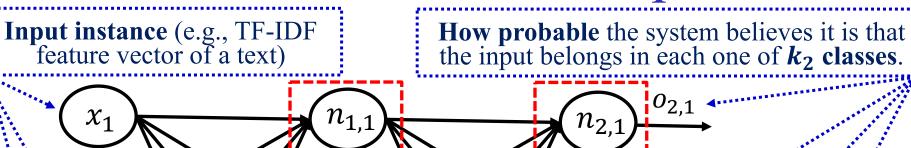
Squared error loss for the **current training instance**.

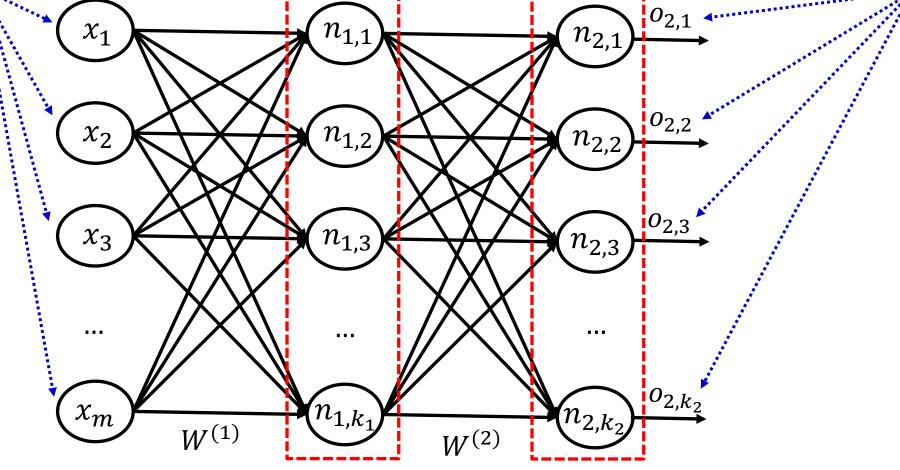
Or as a **computation graph**:

$$\frac{\partial E}{\partial \vec{o}^{(2)}} = \begin{vmatrix} \overline{\partial o_1^{(2)}} \\ \cdots \\ \overline{\partial E} \\ \overline{\partial o_{k_2}^{(2)}} \end{vmatrix} = \begin{bmatrix} o_1^{(2)} - i \\ \cdots \\ o_{k_2}^{(2)} - t \end{bmatrix}$$



Classification example





$$\vec{o}^{(1)} = \tanh(W^{(1)}\vec{x})$$

$$\vec{o}^{(2)} = \text{softmax}(W^{(2)}\vec{o}^{(1)})$$

Softmax

$$W^{(2)}\vec{o}^{(1)} = \vec{s}^{(2)} = \begin{bmatrix} s_{2,1} \\ s_{2,2} \\ ... \\ s_{2,k} \end{bmatrix}$$
 Output of the last layer, without an activation function. Confidence scores, one for each class. We want to convert them to probabilities summing to 1.

$$\operatorname{softmax}(W^{(2)}\vec{o}^{(1)}) = \operatorname{softmax}(\vec{s}^{(2)}) = \operatorname{softmax}\left(\begin{bmatrix} s_{2,1} \\ s_{2,2} \\ \dots \\ s_{2,k_2} \end{bmatrix}\right)$$

$$= \begin{bmatrix} \frac{\exp(s_{2,1})}{\sum_{j=1}^{k_2} \exp(s_{2,j})} \\ \frac{\exp(s_{2,2})}{\sum_{j=1}^{k_2} \exp(s_{2,j})} \\ \dots \\ \frac{\exp(s_{2,k_2})}{\sum_{j=1}^{k_2} \exp(s_{2,j})} \end{bmatrix}$$

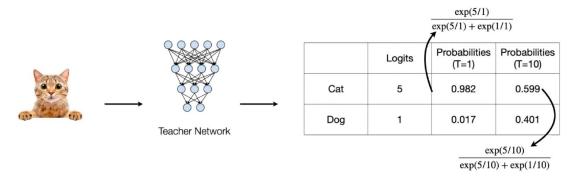
Softmax also moves the largest of its inputs towards 1 and the other inputs towards 0. Intuitively a **soft argmax!**

Softmax with temperature

The key insights (2): Temperature T

Softmax function with temperature to smooth the output

$$\frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \longrightarrow q_i = \frac{exp(z_i/T)}{\sum_j exp(z_j/T)}$$

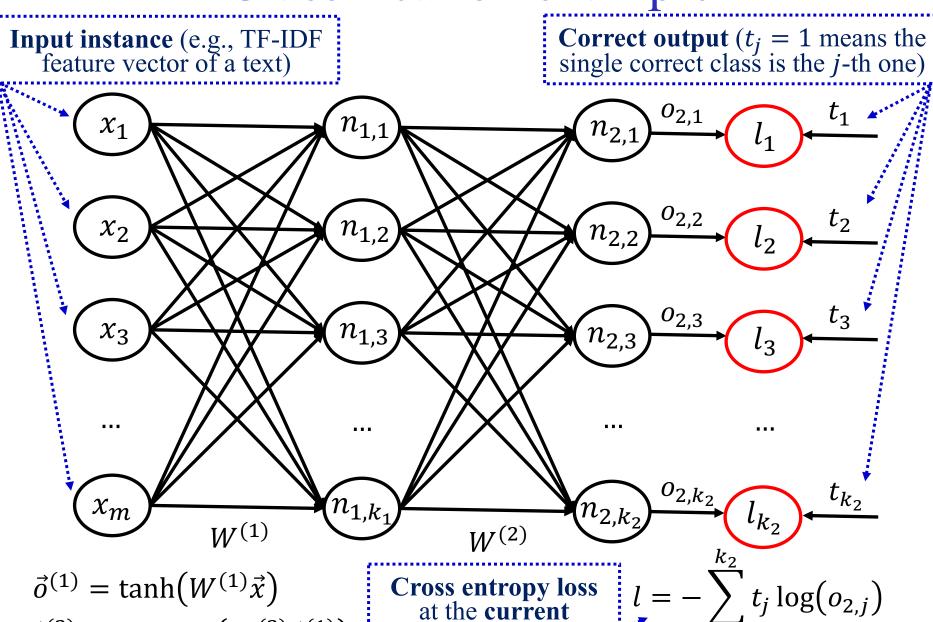


A larger temperature smooths the output probability distribution.

Google

Slide from the talk of Burak Gokturk at Berkeley's LLM Agents course (2024), https://rdi.berkeley.edu/llm-agents/f24.

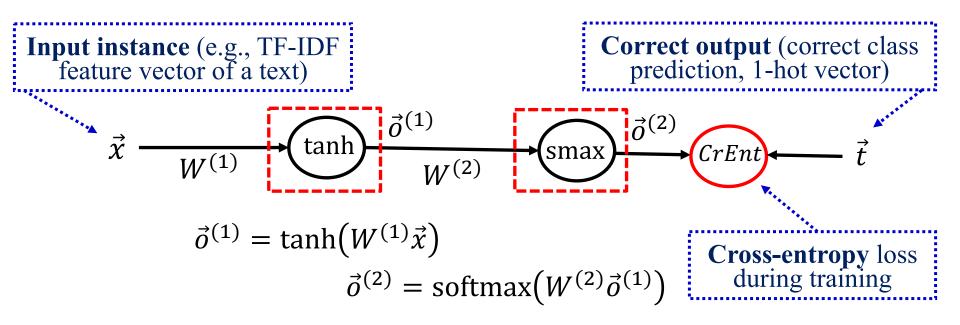
Classification example



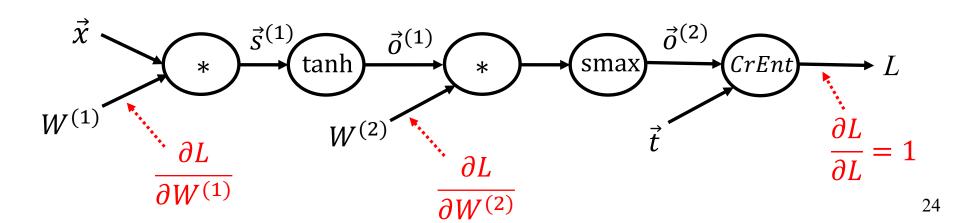
training instance

 $\vec{o}^{(2)} = \text{softmax}(W^{(2)}\vec{o}^{(1)})$

Classification example – more compact



Or as a **computation graph**:



Extracting contract elements

THIS AGREEMENT is made the 15th day of October 2009 (The "Effective Date") BETWEEN:

- (1) Sugar 13 Inc., a corporation whose office is at James House, 42-50 Bond Street, London, EW2H TL ("Sugar");
- (2) **E2 UK Limited**, a limited company whose registered office is at 260 Bathurst Road, Yorkshire, SL3 4SA ("**Provider**").

RECITALS:

A. The Parties wish to enter into a framework agreement which will enable Sugar, from time to time, to [...]

B. [...]

NO THEREFORE IT IS AGREED AS FOLLOWS:

ARTICLE I - DEFINITIONS

"Sugar" shall mean: Sugar 13 Inc.

"Provider" shall mean: E2 UK Limited

"1933 Act" shall mean: Securities Act of 1933

ARTICLE II - TERMINATION

The Service Period will be for five (5) years from the Effective Date (The "Initial Term"). The agreement is considered to be terminated in October 16, 2014.

ARTICLE III - PAYMENT - FEES

During the service period monthly payments should occur. The estimated fees for the Initial Term are £154,800.

ARTICLE IV - GOVERNING LAW

This agreement shall be governed and construed in accordance with the Laws of England & Wales. Each party hereby irrevocably submits to the exclusive jurisdiction of the courts sitting in Northern London.

IN WITNESS WHEREOF, the parties have caused their respective duly authorized officers to execute this Agreement.

BY: George Fake Authorized Officer Sugar 13 Inc.

BY: Olivier Giroux CEO E2 UK LIMITED

Identify start/end dates, duration, contractors, amount, legislations refs, jurisdiction etc. Similar to Named Entity Recognition (NER).

25

I. Chalkidis, I. Androutsopoulos, A. Michos, "Extracting Contract Elements", ICAIL 2017, http://nlp.cs.aueb.gr/pubs/icail2017.pdf.

Window-based NER example

i-th word of the text being classified

3-word **window** (often larger)

yesterday language (tech) announced that...

$$\vec{x}_{i-1} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ ... \\ 0 \end{bmatrix} \vec{x}_i = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ ... \\ 0 \end{bmatrix} \vec{x}_{i+1} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ ... \\ 1 \end{bmatrix}$$
1-hot vectors (|V|×1) of the words in the window. (|V| is the vocabulary size).

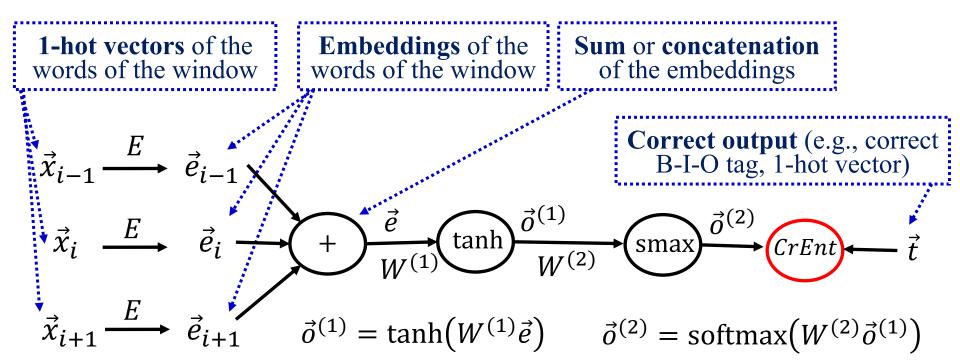
Embeddings (d×1) of the words in the window. (d is the dimensionality of the dimensional dimen

words in the window. (d is the dimensionality of

$$\vec{e}_{i-1} = \begin{bmatrix} 1.8 \\ 2.3 \\ -1.4 \\ 3.7 \\ ... \\ -1.1 \end{bmatrix} \vec{e}_i = \begin{bmatrix} 2.4 \\ -3 \\ 9.3 \\ 5.1 \\ ... \\ 3.9 \end{bmatrix} \vec{e}_{i+1} = \begin{bmatrix} 2.2 \\ 3.8 \\ 1.2 \\ -6.4 \\ ... \\ 7.1 \end{bmatrix}$$
is the **dimensionality** of the word embeddings).

Let **E** be a matrix $(d \times |V|)$ that **contains** all the **embeddings** of the **vocabulary** as **columns**. Then: $\vec{e}_{i-1} = E\vec{x}_{i-1}, \ \vec{e}_i = E\vec{x}_i, \dots$

Window-based NER example



We learn $W^{(1)}$, $W^{(2)}$ with **backpropagation**. We can also learn (or modify) the **word embeddings** E during **backpropagation**! But when we don't have large training datasets (e.g., corpus manually annotated with B-I-O tags), it may be better to use **pre-trained embeddings**, which can be obtained from large non-annotated corpora (e.g., via Word2Vec, GloVe, to be discussed).

We can use the same window-based approach for **POS-tagging**, **chunking**, ...

Cross-entropy loss

Word being classified. 3-word window (often larger).

yesterday language (tech) announced that...

$$\vec{o} = \begin{bmatrix} P_m(C = c_1) \\ P_m(C = c_2) \\ P_m(C = c_3) \\ \dots \\ P_m(C = c_k) \end{bmatrix} = \begin{bmatrix} 0.05 \\ 0.12 \\ 0.08 \\ \dots \\ 0.14 \end{bmatrix}$$

Probability estimates produced by the classifier for the class of the word "tech".

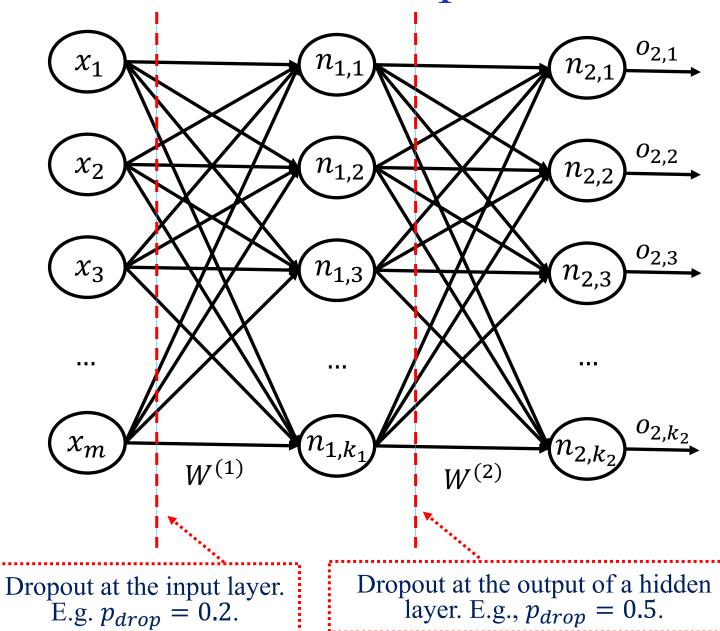
The correct "probabilities" for the class of "tech". A 1-hot vector.

$$\vec{t} = \begin{bmatrix} P(C = c_1) \\ P(C = c_2) \\ P(C = c_3) \\ \dots \\ P(C = c_k) \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

The **log-likelihood** of the correct class according to the classifier (with a minus sign).

$$H_{P_m}(C) = -\sum_{i=1}^k P(C = c_i) \log_2 P_m(C = c_i) = -\log_2 P_m(C = c_2)$$

Dropout

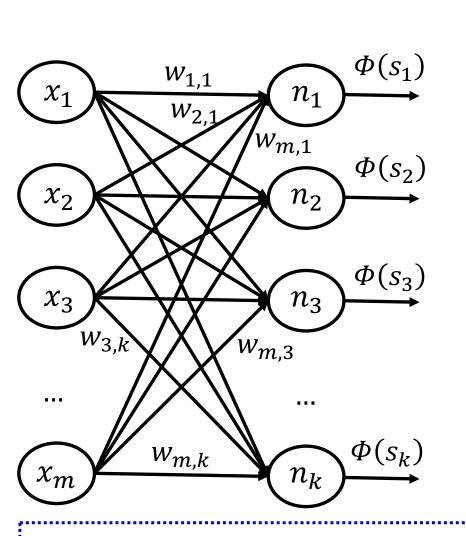


29

Dropout

- For each training instance (or mini-batch), we drop (remove) each neuron of the layer where dropout is applied with **probability** $p_{drop} = 1 - p_{keep}$.
 - Helps the neural net avoid relying too much on particular **neurons** (or inputs). A form of **regularization**. Works well!
 - o **Gradients** also **do not flow** through dropped neurons.
 - o Alternative explanation: we train an ensemble of networks, containing all the pruned networks that dropout creates.
- During testing, we multiply the output of each neuron (of the layer where dropout was applied) by p_{keep} , so that the neuron's expected output value will be as in training.
 - \circ Or we divide the output by p_{keep} during training instead.
 - We don't drop neurons during testing (only during training).

Batch normalization



At each layer, instead of:

$$s_j = \sum_{i=1}^n w_{i,j} x_i$$

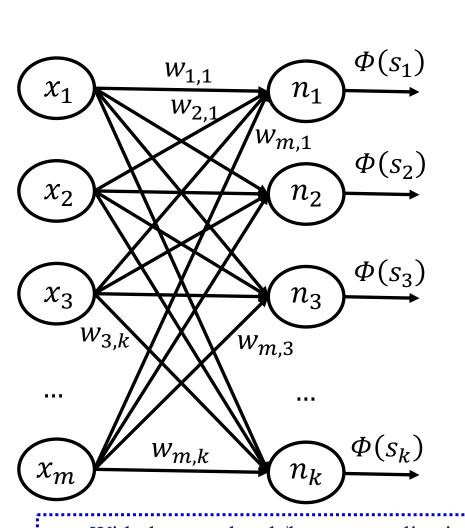
we use:

$$\bar{s}_j = \frac{g_j}{\sigma_i} (s_j - \mu_j) + b_j$$

- μ_j , σ_j are the **mean** and **std. dev. of** s_i in the **mini-batch**.
- g_j , b_j are **learned** parameters (constant after training).
- Φ now applied to \bar{s}_i .

See https://arxiv.org/pdf/1607.06450.pdf for batch vs. layer normalization. The latter is better for RNNs (next part), where layers are time-steps.

Layer normalization



At each layer, instead of:

$$s_j = \sum_{i=1}^m w_{i,j} x_i$$

we use:

$$\bar{s}_j = \frac{g_j}{\sigma} (s_j - \mu) + b_j$$

- μ , σ are the **mean** and **std**. **dev. of s₁, ..., s_k** in the <u>layer</u>.
- g_j , b_j are **learned** parameters (constant after training).
- Φ applied to \bar{s}_i .

With dropout, batch/layer normalization, residuals (to be discussed) and other additions, strictly speaking we no longer have an "MLP". Some people prefer "Feed Forward Neural Network" (FFNN), but "MLP" still often used as synonym.

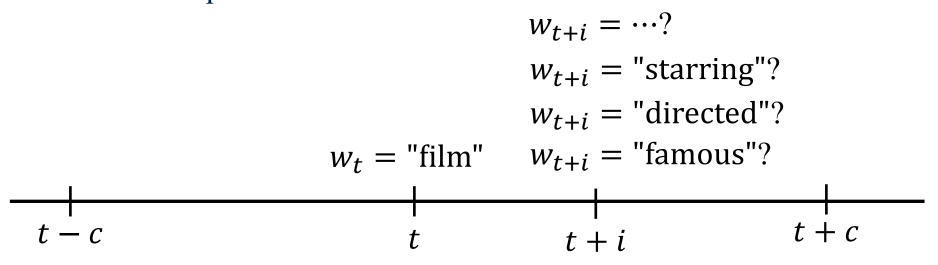
Pretraining word embeddings with Word2Vec

(skip-gram version)

• Every word w of the vocabulary has two vectors:

$$\overrightarrow{w}^{(in)}, \overrightarrow{w}^{(out)}$$

- The **vectors** are **randomly initialized**. We **learn** them.
- For every **token** w_t at **position** t of a **corpus** and **every position** t + i ($i \neq 0$) within a **window** [t c, t + c] around position t:



We want to be able to predict which vocabulary word occurs at position t + i by multiplying $\overrightarrow{w}_t^{(in)}$ and $\overrightarrow{w}_{t+i}^{(out)}$.

Word2Vec (skip-gram version)

In the **skip-gram** version of Word2Vec, the **central word** of each window "**predicts**" the **other words** of the **window**. In the **CBOW version**, the **context** (sum of the embeddings of the other words of the window) "**predicts**" the **central word**.

 $w_{t+i} = \cdots?$ $w_{t+i} = \text{"starring"?}$ $w_{t+i} = \text{"directed"?}$ $w_{t+i} = \text{"famous"?}$

$$w_t = \text{"film"}$$

$$t-c$$
 t
 $t+i$
 $t+c$

prediction

$$P(w_{t+i}|w_t) = \operatorname{softmax}\left(\overrightarrow{w}_{t+i}^{(out)} \cdot \overrightarrow{w}_t^{(in)}\right)$$

$$= \frac{\exp\left(\overrightarrow{w}_{t+i}^{(out)} \cdot \overrightarrow{w}_{t}^{(in)}\right)}{\sum_{w \in V} \exp\left(\overrightarrow{w}^{(out)} \cdot \overrightarrow{w}_{t}^{(in)}\right)}$$

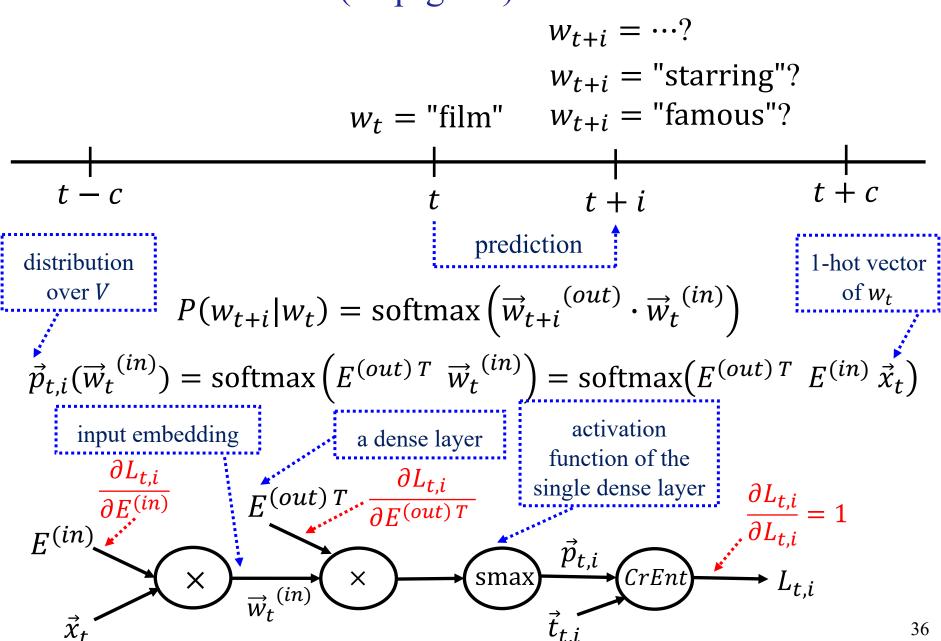
Word2Vec (skip-gram version)

• We learn the $\overrightarrow{w}^{(in)}$, $\overrightarrow{w}^{(out)}$ by maximizing the probability assigned to the w_{t+i} that actually occurs at each position t+i, i.e., we maximize the likelihood of the correct predictions:

$$\langle \hat{E}^{(in)}, \hat{E}^{(out)} \rangle = \underset{\langle E^{(in)}, E^{(out)} \rangle}{\operatorname{argmax}} \sum_{t=1+c}^{T-c} \sum_{-c \le i \le c, i \ne 0} \ln P(w_{t+i}|w_t)$$

- O Matrices $E^{(in)}$, $E^{(out)}$ contain in their columns all the *in* and *out* vectors (word embeddings) of all vocabulary words.
- \circ T is the corpus size, t is the center of the sliding window.
- o For each t value, we get 2c training examples.
- For batch gradient ascent, we would do steps: $\langle \hat{E}^{(in)}, \hat{E}^{(out)} \rangle \leftarrow \langle \hat{E}^{(in)}, \hat{E}^{(out)} \rangle + \eta \nabla L$
- o In practice, we use SGD (or variants), i.e., we use the likelihood L_i of a mini-batch of training examples (e.g., all 2c of a window): $\langle \hat{E}^{(in)}, \hat{E}^{(out)} \rangle \leftarrow \langle \hat{E}^{(in)}, \hat{E}^{(out)} \rangle + \eta \nabla L_i$

Word2Vec (skip-gram) as a shallow NN



Word2Vec (skip-gram with negative sampling)

When the **vocabulary** *V* is **large**, computing the **softmax** is very **time-consuming**. A more efficient alternative is **negative sampling**, a kind of **contrastive learning**.

We construct positive (+) and negative (-) examples, using the word w_{t+i} that actually occurs at position t+i, and random words r_{t+i} that do not actually occur at position t+i.

$$r_{t+i}$$
 = "medical" (random, negative)
 w_t = "film" w_{t+i} = "famous" (true, positive)

$$t-c$$
 t
 $t+i$
 $t+c$

$$\max_{\langle E^{(in)}, E^{(out)} \rangle} \sum_{t=1+c}^{T-c} \sum_{-c \le i \le c, i \ne 0} \ln P(+|w_{t+i}, w_t) + \ln P(-|r_{t+i}, w_t)$$

We try to learn to assign **high probabilities** to the **correct classes**. In practice, we use **multiple random words** r_{t+i} at each position t + i.

Word2Vec (skip-gram with negative sampling)

$$w_{t} = \text{"film"} \quad w_{t+i} = \text{"medical" (random, negative)}$$

$$v_{t} = \text{"film"} \quad w_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad w_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{"famous" (true, positive)}$$

$$v_{t} = \text{"film"} \quad v_{t+i} = \text{$$

We **no longer** try to produce a **probability distribution over the vocabulary** for the words w_{t+i} that may occur at t + i.

We are **given** w_t and a **particular** w_{t+i} or r_{t+i} and we need to decide if it is a positive or negative case.

Loss as a function of epochs

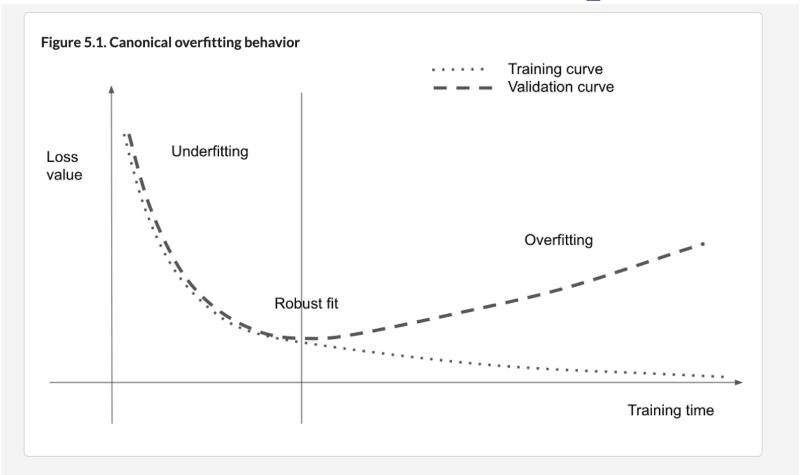


Figure from the **recommended book "Deep Learning with Python"** by F. Chollet, Manning Publications, 2nd edition. The 1st edition is freely available.

https://www.manning.com/books/deep-learning-with-python-second-edition

Practical advice for training deep NNs

- Check simple baselines: (e.g., majority, random, ...)
 - o If you can't beat them, you may have bugs, data problems, ...
 - o Look at how data are tokenized, pre-processed, ...
 - o Examine misclassification errors (e.g., extreme/frequent cases).
- Get the training and validation loss to start falling:
 - Tune the **learning rate** and the **mini-batch size**.
 - Use appropriate **models** (e.g., for sequences, images, ...).
- Reach the overfitting behavior of the previous slide.
 - o The training and validation loss (or metric) both fall up to a point, then the training loss continues to improve ideally reaching near zero, the validation loss deteriorates.
 - o Increase capacity (e.g., layers, neurons per layer), ...
- Then dropout, early stopping, batch/layer norm, ...
- Check Chollet's Chapter 5 for more advice...

Regularizing a high-capacity model

Figure 5.21 Effect of dropout on validation loss

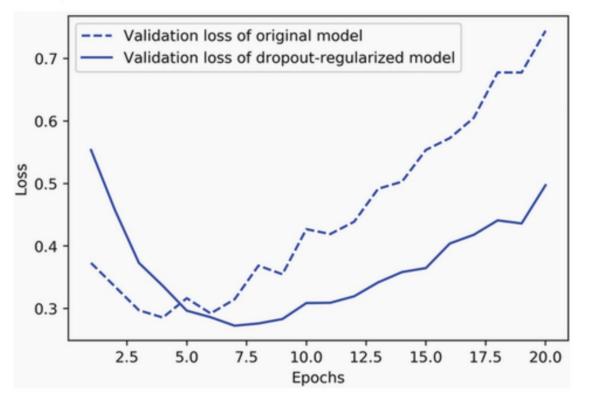


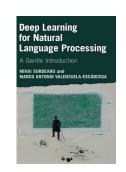
Figure from the **recommended book "Deep Learning with Python"** by F. Chollet, Manning Publications, 2nd edition. The 1st edition is freely available.

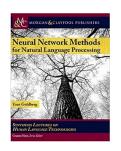
https://www.manning.com/books/deep-learning-with-python-second-edition

Additional optional reading slides.

Recommended reading

- M. Surdeanu and M.A. Valenzuela-Escarcega, *Deep Learning for Natural Language Processing: A Gentle Introduction*, Cambridge Univ. Press, 2024.
 - Chapters 5–9.
 - https://clulab.org/gentlenlp/text.html
 - Also available at AUEB's library.
- Y. Goldberg, Neural Network Models for Natural Language Processing, Morgan & Claypool Publishers, 2017.
 - o Mostly chapters 3–5 and 10. Available at AUEB's library.
- Jurafsky & Martin's, *Speech and Language Processing* is being revised (3rd ed.) to include Deep Learning methods.
 - http://web.stanford.edu/~jurafsky/slp3/







Other recommended resources

- Useful maths background: T. Parr και J. Howard, *The Matrix Calculus You Need for Deep Learning*.
 - https://explained.ai/matrix-calculus
- PyTorch tutorials: https://pytorch.org/tutorials/
- Stanford's NLP with Deep Learning course.
 - http://web.stanford.edu/class/cs224n/. Videos on YouTube.