

## Natural Language Processing

2.2 Scaling word2vec and introduction to Neural Nets for NLP

Prof. Iacopo Masi and Prof. Stefano Faralli

```
import matplotlib.pyplot as plt
import scipy
import random
import numpy as np
import pandas as pd
pd.set_option('display.colheader_justify', 'center')
```

```
import matplotlib
import matplotlib.pyplot as plt
import numpy as np
matplotlib_inline
plt.style.use('seaborn-whitegrid')
font = {'family': 'Times',
        'weight': 'bold',
        'size' : 12}
matplotlib.rc('font', **font)

# Aux functions
def plot_grid(Xs, Ys, axis=None):
    """ Aux function to plot a grid"""
    t = np.arange(Xs.size) # define progression of int for indexing colormap
    if axis is None:
        ax=plt(0, 0, marker="*", color="r", linestyle="none") #plot origin
        ax.scatter(Xs,Ys, c=t, cmap='jet', marker="*") # scatter x vs y
        ax.axis('scaled') # axis scaled
    else:
        plt.plot(0, 0, marker="*", color="r", color="r", linestyle="none") #plot origin
        plt.scatter(Xs,Ys, c=t, cmap='jet', marker="*") # scatter x vs y
        plt.axis('scaled') # axis scaled

def linear_map4(Xs, Ys):
    """Map src points with A*"""
    # (N0,N0) -> (N0x2) # add 3-rd axis, like adding another layer
    src = np.stack((Xs,Ys), axis=Xs.ndim)
    g = np.prod(src.shape[1:-1]) # first two dimension
    # (N0)x2
    src_r = src.reshape(-1,src.shape[-1]) # ask reshape to keep last dimension and adjust the rest
    g = g * src_r.T # dot product
    dst = A @ src_r.T # ZxN
    # (N0)x2 and then reshape as NxNc2
    dst = (dst.T).reshape(src.shape)
    # Access X and Y
    return dst[:,0], dst[:,1]
```

My own latex definitions

## Today's lecture

- Skip-Gram and Continuous Bag of Words (CBOW)
- How to scale word2vec
- Negative Sampling
- Hierarchical Softmax
- Introduction to Neural Nets for NLP

This lecture material is taken from

- [Chapter 6 Jurafsky Book](#)
  - [Chapter 14.5 Eisenstein Book](#)
  - [Stanford Slide Word2Vec](#)
  - [Stanford Lecture Word2Vec](#)
  - [Stanford Notes on Word2Vec](#)
- Another good yet short resource is [\[d2l.ai\] Word embedding](#)
- [Research papers on word2vec and hierarchical softmax:](#)
    - First paper: word2vec + hierarchical softmax
    - Negative Sampling paper
    - A Scalable Hierarchical Distributed Language Model

word2vec is a generic framework

- `word2vec` presents two algorithms:
1. [Skip-Gram](#)
  2. [Continuous Bag-of-Word \(CBOW\) \(we see it today!\)](#)

Also it offers different training methods:

- [with naive softmax](#)
- [negative sampling](#) from [Mikolov et al. 2013]
- [hierarchical softmax](#) (we see them today!)

## word2vec: Skip-Gram Self-Supervision

With strong naive conditional independence assumption

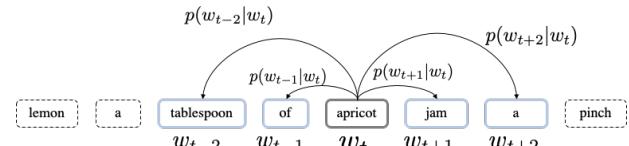
$$p(w_{t-2}, w_{t-1}, w_t, w_{t+1}, w_{t+2} | w_t, \theta) \approx \prod_{-m \leq j \leq m} p(G_{t+j} | w_t, \theta)$$

lemon	a	[tablespoon of apricot jam a] pinch
w_{t-2}	w_{t-1}	w <sub>t</sub> w_{t+1} w_{t+2}

## word2vec: Skip-Gram Self-Supervision

Given the example below, we have to compute:

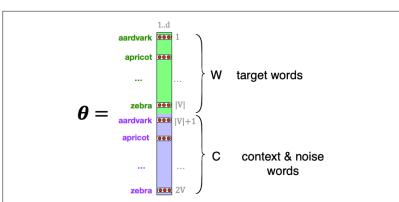
$$p(w_{t-2} | w_t) \cdot p(w_{t-1} | w_t) \cdot p(w_{t+1} | w_t) \cdot p(w_{t+2} | w_t)$$



## word2vec: Skip-Gram with softmax

Parameters to learn:

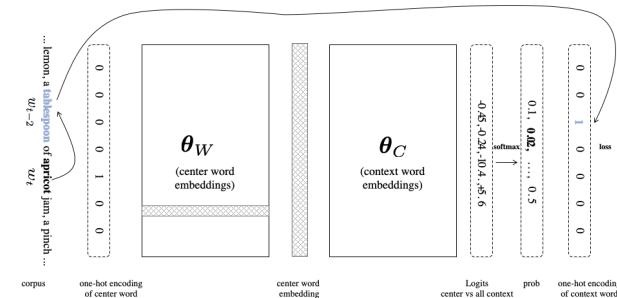
$$\theta = [\theta_W; \theta_C]$$



**Figure 6.13** The embeddings learned by the skipgram model. The algorithm stores two embeddings for each word, the target embedding (sometimes called the input embedding) and the context embedding (sometimes called the output embedding). The parameter  $\theta$  that the algorithm learns is thus a matrix of  $2|V|$  vectors, each of dimension  $d$ , formed by concatenating two matrices, the target embeddings  $\mathbf{W}$  and the context+noise embeddings  $\mathbf{C}$ .

## word2vec with Skip-Gram at a glance

... and why it can be seen as a tiny neural net.



### Loss function: compare two discrete distributions

$$p = \text{softmax}(\theta_C \cdot \theta_W[i]^T)$$

You can think  $p$  as of this form:

	lemon	tablespoon	gelato	...	jam
$p$ (word2vec prediction)	0.001	0.1	0.03	...	0.15

Let's consider the label as a one-hot encoding vector where 1 is over the ground-truth word given by the text.

	lemon	tablespoon	gelato	...	jam
$y$ label (ground-truth $w_{t-1}$ )	0	1	0	...	0

We want to adjust the weights  $\theta$  so that  $p$  matches the label!

### One-hot encoding is a selector!

- $\gamma = [0, 0, 0, 1, 0] \rightarrow$  works as a selector of probability of the actual ground-truth word that we removed!
- Of the probabilities returned by word2vec select that for which the index  $gr$  corresponds to the  $1$  in the label  $\gamma$
- in our case,  $gr$  is the index of 'tablespoon' in  $|V|$ .

$$L(w_{t-1}, w_t; \theta) = -\log[p(w_{t-1}|w_t)]|gr|$$

### Loss function simplified

We can select immediately  $|gr|$  in the numerator.

$$L(w_{t-1}, w_t; \theta) = -\log \left( \frac{\exp(\theta_C[gr] \cdot \theta_W[t]^T)}{\sum_{i=1}^{|V|} \exp(\theta_C[i] \cdot \theta_W[t]^T)} \right)$$

Sometimes is shown as:

$$L(w_{t-1}, w_t; \theta) = -\theta_C[gr] \cdot \theta_W[t]^T + \log \left( \sum_{i=1}^{|V|} \exp(\theta_C[i] \cdot \theta_W[t]^T) \right)$$

similarity center vs context      make sure it is a probability

### word2vec: Continuous Bag-of-Word (C-BOW)

#### Skip-Gram

$$p(w_t | w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2}; \theta)$$

#### C-BOW

$$p(w_t | w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2}; \theta)$$

### word2vec: C-BOW

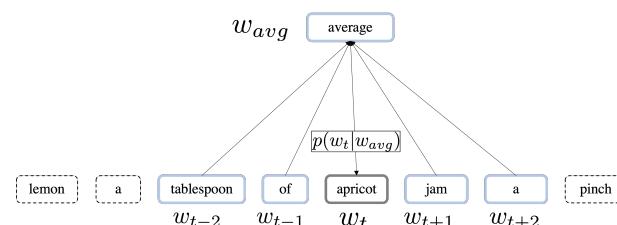
Given the context  $w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2} \rightarrow p(w_t)$

$$\begin{array}{ccccccccc} \text{lemon, } & \text{a } & \text{tablespoon } & \text{of } & \dots & \text{jam } & \text{a } & \text{pinch } \\ \text{---} & \text{---} \\ w_{t-2} & w_{t-1} & w_t & w_{t+1} & w_{t+2} & \theta_C & \theta_W & \theta_W \end{array}$$

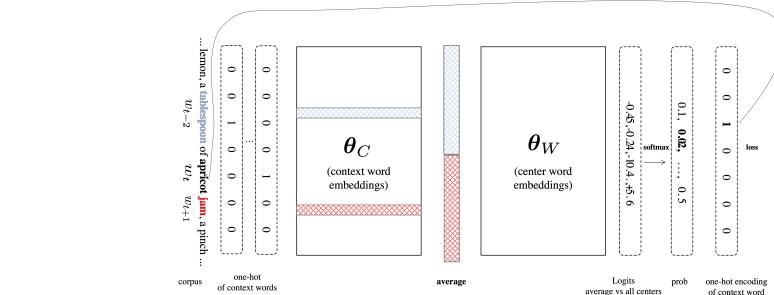
### word2vec: Continuous Bag-of-Word (CBOW)

$$w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2} \rightarrow p(w_t)$$

$$\theta_{C\_avg} = \sum_{m \leq j \leq m, j \neq 0} \theta_C[j+1]$$



## word2vec: Continuous Bag-of-Word (CBOW)



How to compute  $p(w_{t-1} | w_{avg})$ ?

1.  $\theta_{C_{avg}}$  is the average embedding from the context parameters  $\theta_C$ . The average is computed once defined the context window size.

$$\frac{\sum}{|V| \times 1} = \frac{\theta_C}{|V| \times D} \cdot \frac{\theta_W^T}{D \times 1}$$

1.  $x$  is logits and encodes the similarity via dot product of the average context word embedding  $\theta_{C_{avg}}$  against all vocabulary words taken as center  $\theta_C$ .

2. We pass  $x$  through softmax operator to get a distribution over  $|V|$  as:

$$p = \text{softmax}(x)$$

You can think  $p$  as of this form:

	lemon	tablespoon	gelato	...	jam
$p$	0.001	0.1	0.03	...	0.15

## Skipgram vs CBOW at a glance

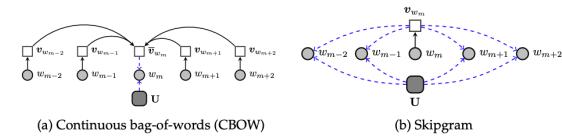


Figure 14.3: The CBOW and skipgram variants of WORD2VEC. The parameter  $U$  is the matrix of word embeddings, and each  $v_m$  is the context embedding for word  $w_m$ .

## Computational Complexity

- word2vec computes a normalized probability over word tokens (over the vocabulary  $V$ )
- A naive implementation of this probability requires summing over the entire vocabulary  $V$

$$L(w_{t-1}, w_t; \theta) = -\log \left( \frac{\exp(\theta_c[g] \cdot \theta_g[l]^T)}{\sum_{v \in V} \exp(\theta_c[v] \cdot \theta_g[l]^T)} \right)$$

The normalization in the denominator is of the order of:

$O(|V| \times D)$

where:

- $|V|$  is the vocabulary size and e.g.  $|V| = 3M$
- $D$  is the dimension of the embeddings e.g.  $D = 300$ .

## Two solutions to approximate the denominator

1. Negative sampling (Contrastive method)
2. Hierarchical Softmax (Tree-based solution)

## Scaling word2vec with Negative Sampling

## Skip-gram with Negative Sampling (SGNS)

Instead of doing:

1. Center word vs ground-truth context embedding  $\rightarrow \theta_c[g] \cdot \theta_g[l]^T$
2. normalize as a distribution: all context vs center word  $\rightarrow \sum_{v \in V} \exp(\theta_c[v] \cdot \theta_g[l]^T)$

$$L(w_{t-1}, w_t; \theta) = -\theta_c[g] \cdot \theta_g[l]^T + \log \left( \sum_{v=1}^{|V|} \exp(\theta_c[v] \cdot \theta_g[l]^T) \right)$$

similarity center vs context      make sure it is a probability

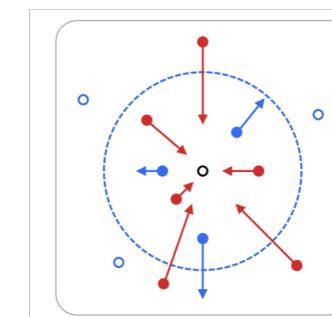
## Skip-gram with Negative Sampling (SGNS)

We "relax" the denominator and do:

1. [Same as before] Center word vs ground-truth context embedding  $\rightarrow \theta_c[g] \cdot \theta_g[l]^T$
2. Approximate the denominator as:
  - Sample  $k$  context words from  $V$ .
  - Compare each  $i$ -th sampled word vs center word

## Negative Sampling is a form of Contrastive Method

Attractive Forces      Repulsive Forces



## Negative Sampling is a form of Contrastive Method

- Positive Constraint:
  - Pull  $\theta_c[g]$  and  $\theta_g$  to be "close" in space (give high dot-product score)
  - Center and ground-truth context need to have high similarity
  - In doing so we may make other words unrelated to be close in the space, so we have to push them away
- Negative Constraint:
  - We have to push away  $\theta_g$  for any other remaining  $\theta_g[i]$   $\forall i \in V, i \neq g$
  - Before, this was done by forcing the total "mass" over the vocabulary distribution to sum to 1

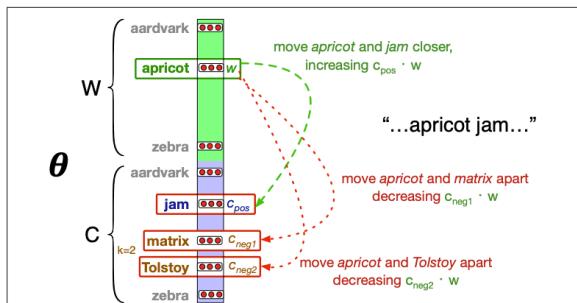


## Loss with negative sampling

If we minimize, we have to invert the sign:

$$\min_{\theta} -\log \left( \theta_{c(g)}^T \theta_H(l) \right) - \sum_{k=1}^K \log \left[ e^{-\theta_c(k)^T \theta_H(l)} \right]$$

## Visualization



**Figure 6.14** Intuition of one step of gradient descent. The skip-gram model tries to shift embeddings so the target embeddings (here for *apricot*) are closer to (have a higher dot product with) context embeddings for nearby words (here *jam*) and further from (lower dot product with) context embeddings for noise words that don't occur nearby (here *Tolstoy* and *matrix*).

## How to sample the negatives

Which is, how to select the indexes  $k \in \{1 \dots K\}, k \neq g$  in  $\sum_{i=1}^K \log [e^{-\theta_c(k)^T \theta_H(i)}] \text{?}$

We sample from an **Unigram model** defined over the corpus  $I$  as:

$$P(v)_a = \frac{\text{count}(v)^a}{\sum_v \text{count}(v)^a}$$

and fixing  $a = \frac{3}{4} = 0.75$

Setting  $a = 0.75$  gives better performance because **gives rare noise words slightly higher probability**. For rare words,  $P(v)_a > P(v)$  while **common words are more or less kept the same**.

## How to sample the negatives

$$P(v)_a = \frac{\text{count}(v)^a}{\sum_v \text{count}(v)^a}$$

and fixing  $a = \frac{3}{4} = 0.75$

Example:

$$\text{is: } 0.9^{.75} = 0.92 \text{ Constitution: } 0.09^{.75} = 0.16 \text{ bombastic: } 0.01^{.75} = 0.032$$

## Sparse Gradients

- We iteratively take gradients at each window for SGD
- In each window, we only have at most  $2w + 1$  words plus  $2wm$  negative words with negative sampling, so the gradient over a window  $\nabla_{\theta} L(\theta)$  is very sparse!
- Computationally, it is important to not have to send gigantic updates around.

$$\nabla_{\theta} L(\theta) = \begin{pmatrix} \theta \\ \vdots \\ \nabla_{\theta_{\text{vocab}}} \\ \vdots \\ \theta \\ \vdots \\ \nabla_{\theta_{\text{dim}}} \\ \vdots \\ \nabla_{\theta_{\text{learning}}} \\ \vdots \end{pmatrix} \in \mathbb{R}^{2wF}$$

## Scaling word2vec with Hierarchical Softmax

## Hierarchical Softmax in NLP

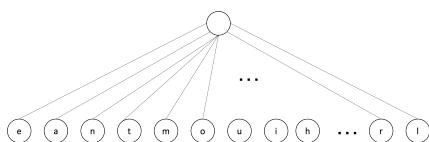
It is an alternative to **Negative Sampling**. We do not use a contrastive method yet we approximate the normalization over a large vocabulary with a **balanced binary tree structure**.

Computational cost reduces from  $O(|V|)$  to  $O(\log_2(|V|))$  in the best case.

### Regular Softmax is a degenerate tree!

We change point of view: softmax is a tree of depth 1 and  $|V|$  children that are leaf too!

Corpus this is an example of a huffman tree and assume \*\*word tokens are characters\*\* to simplify.



## Hierarchical Softmax in NLP

1. Given a vocabulary of word token  $I$  how to construct the tree (there are multiple ways of doing it)

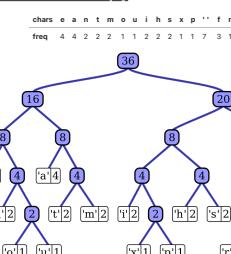
2. How to train with the tree

3. How to perform inference

## Huffman tree

Given a vocabulary of word token  $I$  how to construct the tree? We use **Huffman trees**

The corpus is `this is an example of a huffman tree` and assume \*\*word tokens are characters\*\* to simplify.

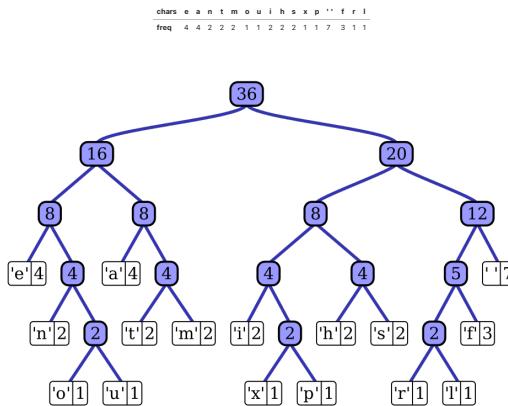


picture from Wikipedia

## Huffman tree

- vocabulary  $I \rightarrow$  word frequencies with **Unigram model**.
- More frequent word tokens are placed closer to the root; rare words are at deeper layers.
  - If you think Information Theory, we "spend" less in encoding frequent words and each word has a variable code length.
  - We want to encode "e" with a bit string. Convention: left is **0** and right is **1**
  - Thus encoding of "e" is left->left->left which is **000** 3 bits
  - Yet the encoding of "p" is right->left->left->right->right which is **10011** 5 bits
- Each node has always two children

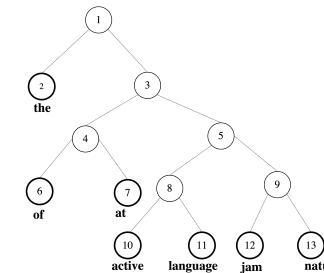
## Huffman tree



picture from Wikipedia

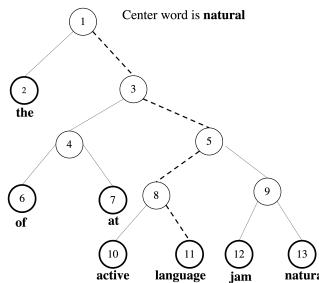
## Training with Hierarchical Softmax

We do not model anymore  $\#$  as the number of words in  $I'$ , yet we model  $\#$  as the number of internal nodes in the tree  $I'-1$ .  
We have a feature vector to be learned at each node  $i$  of the tree for a total of  $2I'-1$  vectors to be learned (context and center).  
We assume Skip-gram and we want to compute  $p(w_{t-1} | w_t = \text{natural})$ ; we also assume the ground-truth  $w_{t-1}$  is language.



## Training with Hierarchical Softmax

Assume Skip-gram and we want to compute  $p(w_{t-1} | w_t = \text{natural})$  and we assume the ground-truth  $w_{t-1}$  is language.  
So we "bypass" all paths except the one that from root leads to `language`.



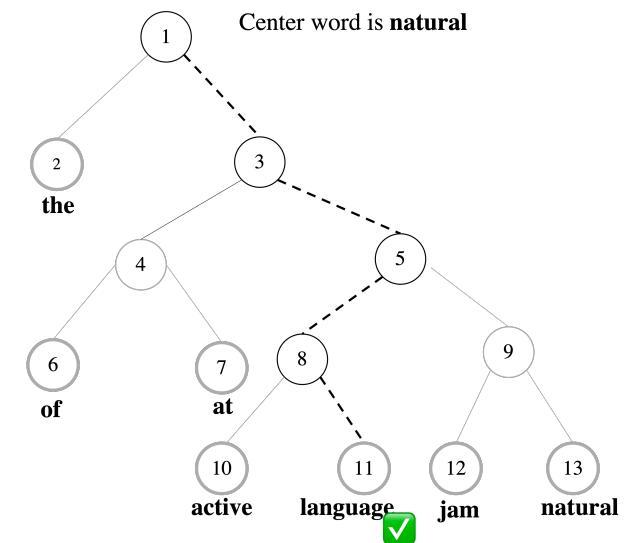
## Training with Hierarchical Softmax

$p(w_{t-1} | w_t = \text{natural})$  and we assume the ground-truth  $w_{t-1}$  is "language".

$$p(w_{t-1} | w_t = \text{natural}) = \prod_{n \in \text{path}(\text{root} \rightarrow \text{language})} p_{\text{branch}}(n, i)$$

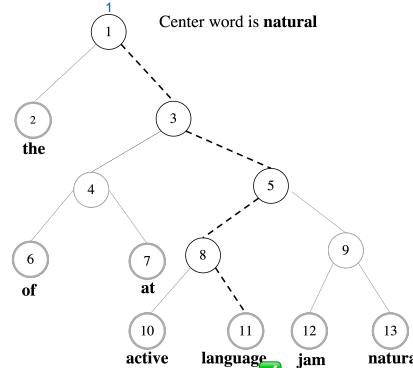
where:

$$p_{\text{branch}}(n, i) = \begin{cases} \sigma(\theta_i(a)^T \theta_B(l)), & \text{if left} \\ 1 - \sigma(\theta_i(a)^T \theta_B(l)), & \text{if right} \end{cases}$$

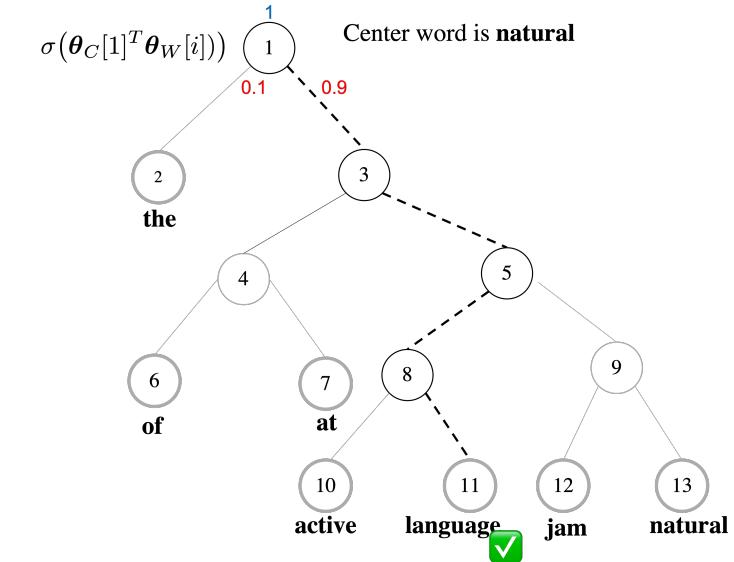


### Exercise: compute the loss with Hierarchical Softmax

$P(w_{t-1} | w_t = \text{natural})$  and we assume the ground-truth  $w_{t-1}$  is language.

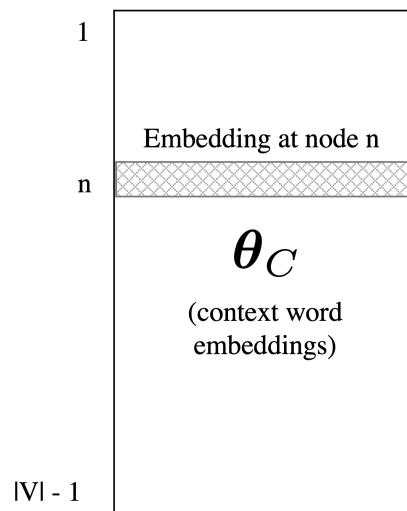


### Exercise: compute the loss with Hierarchical Softmax



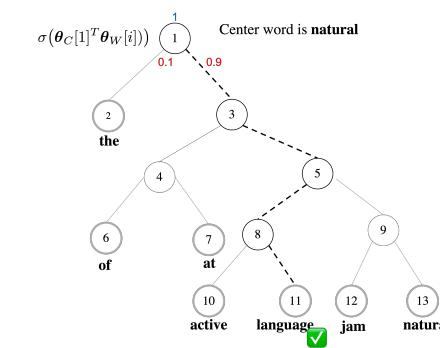
### Embeddings as a Matrix

The embedding  $\theta_C$  are still in a matrix, where each row of the matrix is indexed by the node  $n$  index.



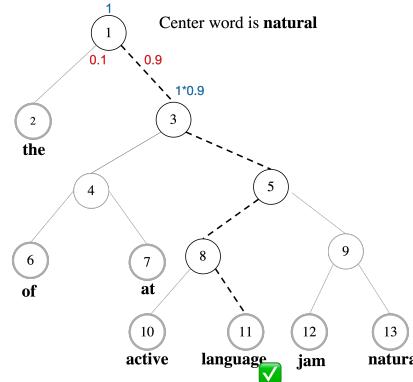
### Exercise: compute the loss with Hierarchical Softmax

$P(w_{t-1} | w_t = \text{natural})$  and we assume the ground-truth  $w_{t-1}$  is language.



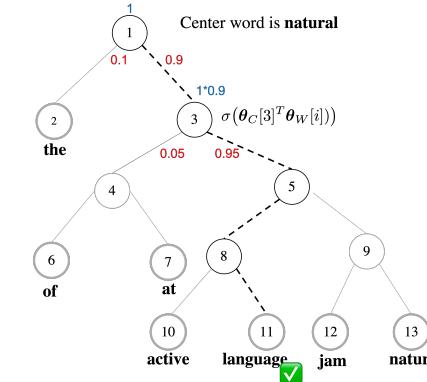
### Exercise: compute the loss with Hierarchical Softmax

$P(w_{t-1}|w_t = \text{natural})$  and we assume the ground-truth  $w_{t-1}$  is language.



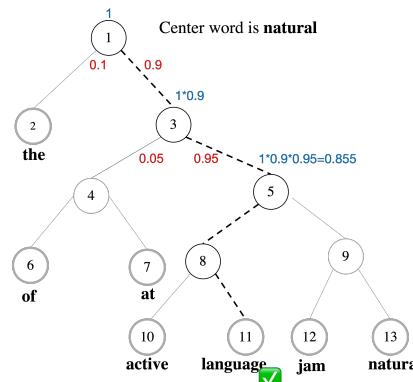
### Exercise: compute the loss with Hierarchical Softmax

$P(w_{t-1}|w_t = \text{natural})$  and we assume the ground-truth  $w_{t-1}$  is language.



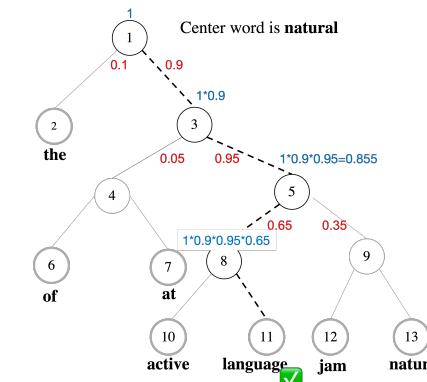
### Exercise: compute the loss with Hierarchical Softmax

$P(w_{t-1}|w_t = \text{natural})$  and we assume the ground-truth  $w_{t-1}$  is language.



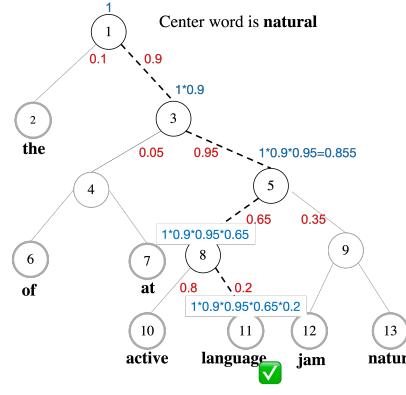
### Exercise: compute the loss with Hierarchical Softmax

$P(w_{t-1}|w_t = \text{natural})$  and we assume the ground-truth  $w_{t-1}$  is language.



### Exercise: compute the loss with Hierarchical Softmax

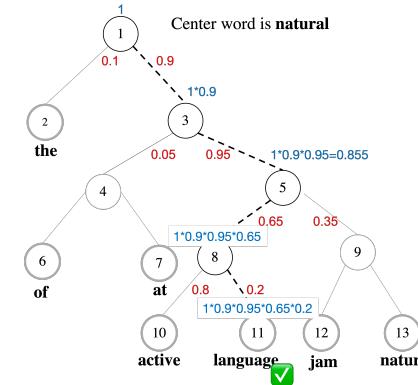
$P(w_{t-1} | w_t = \text{natural})$  and we assume the ground-truth  $w_{t-1}$  is language.



### Exercise: compute the loss with Hierarchical Softmax

$P(w_{t-1} | w_t = \text{natural}) = 1 \cdot 0.9 \cdot 0.95 \cdot 0.65 \cdot 0.2 = 0.11115$

Loss is  $-\log(P(w_{t-1} | w_t = \text{natural})) = -\log(0.11115)$



### Inference with Hierarchical Softmax

Important: In inference with do not have the label!

1. Exhaustive search [too complex]
2. Greedy search (at each branch take the branch at maximum probability) [too greedy]
3. Beam search (we will cover later on)

### Google Allo project used Beam search/Hierarchical Softmax

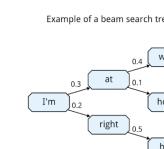
Google (now dead 🛡) project Allo used a Hierarchical tree to speed up inference.

Taken from <https://ai.googleblog.com/2016/05/chat-smarter-with-allo.html>

### Google Allo project used Beam search/Hierarchical Softmax

Google (now dead 🛡) project Allo used a Hierarchical tree to speed up inference.

A beam search is used to efficiently select the top N highest scoring responses from among the very large set of possible messages that a LSTM can generate. A snippet of the search space explored by such a beam-search technique is shown below.



As with any large-scale product, there were several engineering challenges we had to solve in generating a set of high-quality responses efficiently. For example, in spite of the two staged architecture, our first few networks were very slow and required about half a second to generate a response. This was obviously a deal breaker when we are talking about real time communication app! So we had to evolve our neural network architecture further to reduce the latency to less than 200ms. We moved from using a softmax layer to a hierarchical softmax layer which traverses a tree of words instead of traversing a list of words thus making it more efficient.

### Google Allo project used Hierarchical Softmax

As with any large-scale product, there were several engineering challenges we had to solve in generating a set of high-quality responses efficiently. For example, in spite of the two staged architecture, our first few networks were very slow and required about half a second to generate a response. This was obviously a deal breaker when we are talking about real time communication app!

So we had to evolve our neural network architecture further to reduce the latency to less than 200ms. We moved from using a softmax layer to a hierarchical softmax layer which traverses a tree of words instead of traversing a list of words thus making it more efficient.

Taken from <https://ai.googleblog.com/2016/05/chat-smarter-with-allo.html>

## Hints on fastText

The rest of fastText is the same as the skip-gram model. Compared with the skip-gram model, the vocabulary in fastText is larger, resulting in more model parameters. Besides, to calculate the representation of a word, all its subword vectors have to be summed, leading to higher computational complexity. However, thanks to shared parameters from subwords among words with similar structures, rare words and even out-of-vocabulary words may obtain better vector representations in fastText.



- Morphological structure of a word carries important information about the meaning of the word.
- fastText attempts to solve this by treating each word as the aggregation of its subwords.
- Morphologically rich languages (German, Turkish) in which a single word can have a large number of morphological forms, each of which might occur rarely, thus making it hard to train good word embeddings.

## Subword Embedding

In English, words such as

"helps", "helped", and "helping" are

inflected forms of the same word "help". The relationship between "dog" and "dogs" is the same as that between "cat" and "cats", and the relationship between "boy" and "boyfriend" is the same as that between "girl" and "girlfriend".

In other languages such as French and Spanish, many verbs have over 40 inflected forms, while in Finnish, a noun may have up to 15 cases.

In linguistics, morphology studies word formation and word relationships. However, the internal structure of words was neither explored in word2vec nor in GloVe.

## The fastText Model

To use morphological information, the fastText model proposed a subword embedding approach, where a **subword is a character n-gram**.

Instead of learning word-level vector representations, fastText can be considered as the **subword-level skip-gram**, where each center word is represented by the sum of its subword vectors.

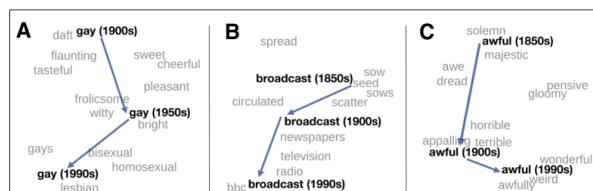
## The fastText Model

Let's illustrate how to obtain subwords for each center word in fastText using the word "where". First, add special characters "<" and ">" at the beginning and end of the word to distinguish prefixes and suffixes from other subwords. Then, extract character n-grams from the word. For example, when  $n = 3$ , we obtain all subwords of length 3: "<wh", "wh<", "<he", "he<", "<re", "re>", and the special subword "<where>".

In fastText, for any word  $v$ , denote by  $G_v$  the union of all its subwords of length between 3 and 6 and its special subword. The vocabulary is the union of the subwords of all words. Letting  $\mathbf{z}_g$  be the vector of subword  $g$  in the dictionary, the vector  $\mathbf{v}_v$  for word  $v$  as a center word in the skip-gram model is the sum of its subword vectors:

$$\mathbf{v}_v = \sum_{g \in G_v} \mathbf{z}_g$$

## Embedding and Historical Semantics



**Figure 6.17** A t-SNE visualization of the semantic change of 3 words in English using word2vec vectors. The modern sense of each word, and the grey context words, are computed from the most recent (modern) time-point embedding space. Earlier points are computed from earlier historical embedding spaces. The visualizations show the changes in the word *gay* from meanings related to "cheerful" or "frolicsome" to referring to homosexuality, the development of the modern "transmission" sense of *broadcast* from its original sense of sowing seeds, and the pejoration of the word *awful* as it shifted from meaning "full of awe" to meaning "terrible or appalling" (Hamilton et al., 2016b).

## Bias and Embeddings

In addition to their ability to learn word meaning from text, embeddings, alas, also **reproduce the implicit biases and stereotypes** that were latent in the text

Besides the notorious

man : king = woman : queen

the same embeddings analogies also exhibit gender bias

man : computer programmer = woman : homemaker

father : doctor = mother : nurse

## Evaluating Word Embeddings

Related to general evaluation in NLP/Machine Learning: **Intrinsic vs. Extrinsic**

- Intrinsic:**
  - Evaluation on a specific/intermediate subtask
  - Fast to compute
  - Helps to understand that system
  - Not clear if really helpful unless correlation to real task is established
  - Often involve **correlation with human judgments**

- Extrinsic:**
  - Evaluation on a real downstream task
  - Can take a long time to compute accuracy
  - Unclear if the subsystem is the problem or its interaction or other subsystems
  - If replacing exactly one subsystem with another improves accuracy → Winner! (**Ablation study**)
  - Always perform validation of the hyper-parameter on a validation (or dev) set; when you are sure then test once.

Taken from cs224n Stanford NLP with Deep learning

## Bias and Embeddings

Bias in the embeddings can cause **allocational harm**

when a system allocates resources (jobs or credit) unfairly to different groups. For example algorithms that use embeddings as part of a search for hiring potential programmers or doctors might thus incorrectly downweight documents with women's names

## Bias and Embeddings

Bias in the embeddings can cause **representational harm** as in bias towards the ethnicity groups

Using such methods, people in the United States have been shown to associate African-American names with unpleasant words (more than European-American names), male names more with mathematics and female names with the arts, and old people's names with unpleasant words

## Debiasing in NLP/Machine Learning is an Open-Problem

### Word Embedding: Topics we did NOT cover

- Global Vector - GloVe Model: taking the best of both word (SVD-based and iterative, word2vec based). Invented by Stanford.
- We only skimmed through fastText: extension of word2vec to deal with the **Out of Vocabulary (OOV)** problem
- FastText demo with Gensim

## Natural Language Processing with Deep Learning

### Brief recap of Neural Networks

#### Today's lecture

##### Supervised, **Parametric Models**

Propaedeutic part for Deep Learning

##### 0) Optimization in Deep Learning

##### 1) Network Structure: Multi-Layer Perceptron (MLP) is a Fully-Connected Neural Net

##### 2) Backpropagation

This lecture material is taken from

#### Chapter 7 Jurafsky Book

Another good yet short resource is [d2l.ai] Word embedding

- d2l.ai - Multi Variable Calculus
- Karpathy (Tesla Machine Learning Directory) Lecture on Backprop
- Stanford Neural Nets and Backprop lecture
- Stanford ML notes on Neural Nets
- Stanford ML notes on Backprop
- Animation from jermwatt.github.io

## Deep Learning

### 0) Quick Intro to Optimization in Deep Learning

#### 1) What is a Neural Net (just Multi-Layer Perceptron)

#### 2) How to obtain gradients on the weights

#### Gradient Descent or Batch GD

- Compute the gradient of the loss wrt to params for all n training samples
- $\theta \leftarrow \theta - \sum_{i=1}^n \nabla_{\theta} l(\theta, x_i, y_i)$

#### Stochastic Gradient Descent or SGD

- Compute the gradient of the loss wrt to params for a single random training samples
- $\theta \leftarrow \theta - \gamma \nabla_{\theta} l(\theta, x_j, y_j)$

#### How to optimize a Neural Net - SGD over mini-batches

1. In-between Batch GD and SGD with a single sample
2. We load randomly k samples over the n; usually k is a power of 2.
3. mini batch of 32, 64, 128 but could also be 100
4.  $\theta \leftarrow \sum_{i=1}^k \nabla_{\theta} l(\theta, x_i, y_i)$
5. Practically you take your training set X and you **shuffle** it, then go over it k by k. Simulate uniform random sampling without replacement.
6. When the list is over, re-start and shuffle again.
5. When you have performed a full pass on the shuffled data, this is called an **EPOCH**
6. You can train NN over iterations or over EPOCHS

## Training scheme Pseudo-code

```
from random import shuffle
training = list(range(1,11)) # each index points to a training sample, could be a matrix xmtx0x3, label y
shuffle(training)
converge, it, max_it, epoch = False, 0, 100, 3, 0
while not converge and it < max_it: # if you training convergence scheme
    print(f'[Epoch {epoch}]')
    for b in range(0, len(training), k): # Data Loader gives you a batch k x matrices
        mini_batch = training[b:b+k] # so mini-batch is a tensor HxWx3k
        if len(mini_batch) != k: # a possible way of handling the offset
            continue
        print('SGD step taken over', mini_batch) # compute the loss/gradients and update your model
        loss.backward() # get the gradients
        optimizer.step() # Incorporate the gradients
        # check convergence and set it to True
        if it == max_it:
            converge = True
        epoch += 1 # an epoch is done, we reshuffle the training set
        shuffle(training)

# Original unshuffle training set [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
[Epoch 0]
SGD step taken over [10, 8, 1, 2, 6, 4, 9, 7, 5, 3]
SGD step taken over [2, 6, 4]
SGD step taken over [9, 7, 5]
[Training set [10, 8, 1, 2, 6, 4, 9, 7, 5, 3]
[Epoch 1]
SGD step taken over [1, 10, 6]
SGD step taken over [9, 3, 7]
SGD step taken over [8, 4, 5]
[Training set [6, 3, 10, 6, 9, 3, 7, 8, 4, 5, 2]
[Epoch 2]
SGD step taken over [6, 3, 10]
SGD step taken over [5, 9, 8]
SGD step taken over [4, 7, 2]
[Training set [1, 2, 5, 10, 6, 7, 9, 8, 3, 4]
[Epoch 3]
SGD step taken over [1, 2, 5]
SGD step taken over [10, 6, 7]
SGD step taken over [9, 8, 3]
[Training set [2, 3, 1, 9, 6, 8, 4, 10, 7, 5]
[Epoch 4]
```

### 1) SGD over mini-batches

#### Initialization - Very Important if the function is not strictly convex

$\theta \sim N(\cdot, \sigma)$  omit details for now

With N(0, 1) random initialization from a distribution (There are different methods). We do not set them all to zero

#### 2. Repeat until convergence:

- Compute the gradient of the loss wrt to the parameters  $\theta$  given the mini-batch
- Take a small step in the opposite direction of steepest ascent (so steepest descent).

$$\theta \leftarrow \theta - \sum_{i=1}^k \nabla_{\theta} l(\theta, x_i, y_i)$$

3. When convergence is reached, you final estimate is in  $\theta$

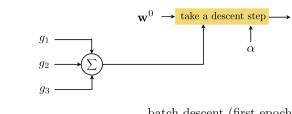
### Change of vocabulary - A bunch of training samples is a mini-batch

- We train NN **Stochastic Gradient Descent** over mini-batches with momentum (or variations thereof)
- When you train NN you 'sample' a mini-batch  $X_k$  from your big dataset  $X$ .

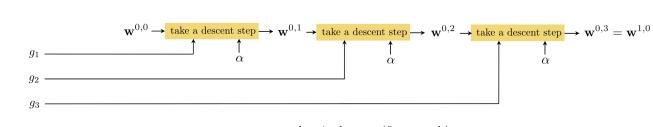
Below this holds for the final linear layer:

$$\frac{\mathbf{y}}{\mathbf{g}^{L+1}} = \frac{\mathbf{W}}{\mathbf{g}^{L+1}} \frac{\mathbf{x}_k + \mathbf{b}}{\mathbf{g}^L}$$

### Mini-Batch, Visually



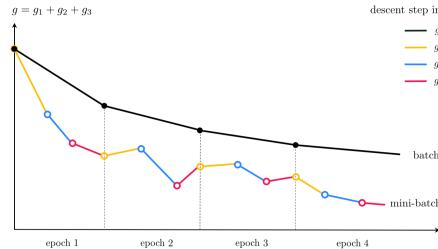
batch descent (first epoch)



stochastic descent (first epoch)

## Mini-Batch SGD vs [Batch] GD

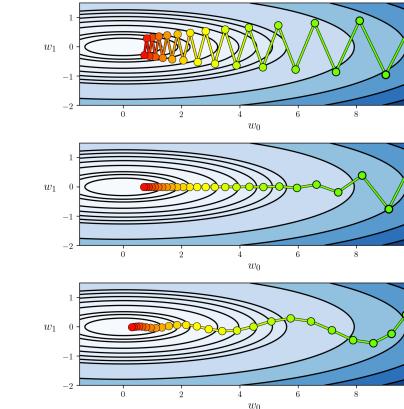
Loss in NN in **non-convex** with lots of local-minima so stochasticity adds noise that let the optimization escape from local minima.



Mini-batch is a sort of smoothing of the single point SGD

There is another smoothing technique: Momentum

Top: SGD; Bottom: SGD with momentum increasing memory of previous steps



## SGD over mini-batches with Momentum

1. Initialization - **Very Important** if the function is not strictly convex

$\theta \sim N$  omit details

With NN random initialization from a distribution (There are different methods). We do not set them all to zero

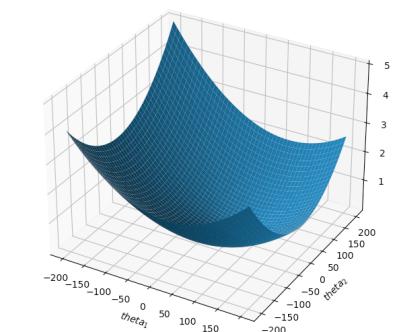
2. Repeat until **convergence**:

- Compute the gradient of the loss wrt to the parameters  $\theta$  given the **mini-batch**
- Take a small step in the opposite direction of steepest ascent (**so steepest descent**).

$$\begin{aligned} \Delta_{t+1} &= \alpha \Delta_t + (1 - \alpha) \sum_{i=1}^k \nabla \varphi(\theta; \mathbf{x}_i, y_i) \\ &\text{new update} \\ \theta &\leftarrow \theta - \gamma \Delta_{t+1} \end{aligned}$$

3. When convergence is reached (or **EARLY STOPPING**), you final estimate is in  $\theta$

## Loss Surface for Linear Regression $\ell_2^2$ loss with $d = 2$ parameters in $\theta$



With Deep Learning optimization is highly non-convex and #params explode!

## Loss Surface for ResNet-20 with no skip connection on ImageNet

ResNet-20, number of parameters  $\theta$  of the order of.....millions!

GPT-3 (LM behind chatGPT) has 150 billions parameters

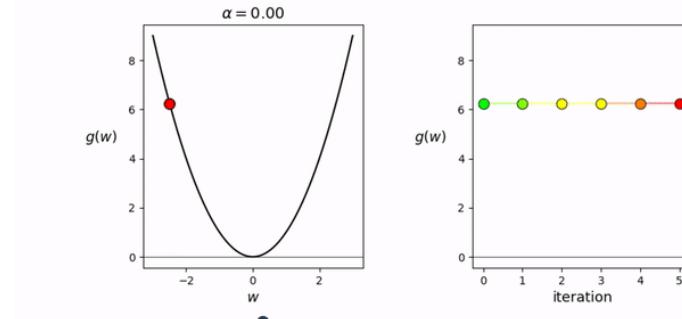
Visualization of mode connectivity for ResNet-20 with no skip connections on ImageNet dataset. The visualization by Javier Ideami



Taken from [https://zmailovpavel.github.io/curves\\_blogpost/](https://zmailovpavel.github.io/curves_blogpost/)

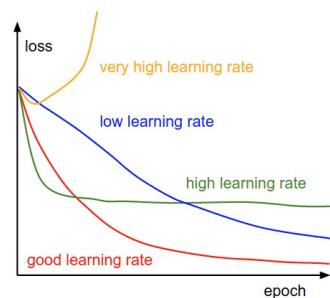
Video for the curious student

Learning rate is very important

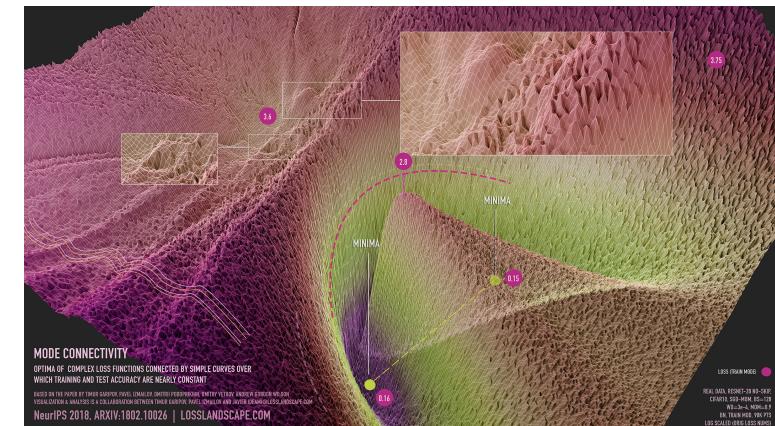


## Babysitting the training process

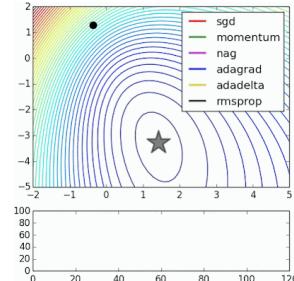
Loss in function of epochs



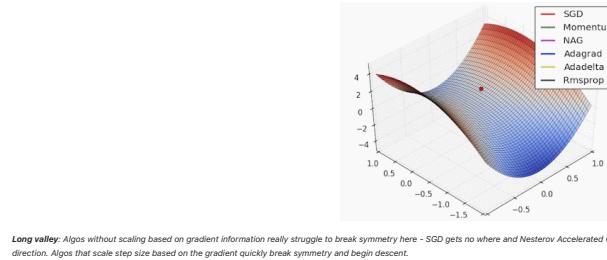
## Valleys, Hills, Noisy Surface



## Dynamics of Training

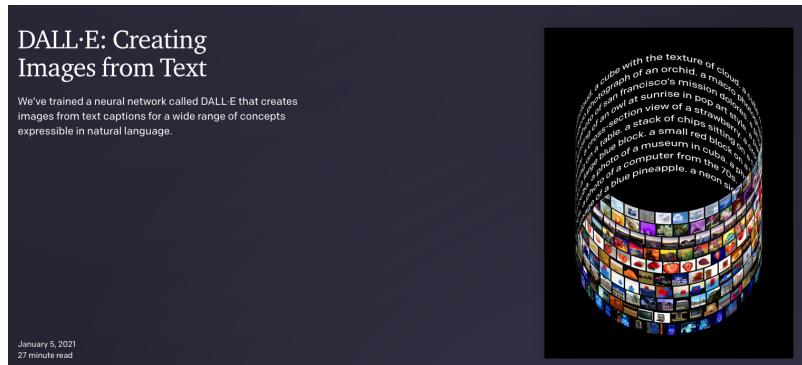


## Dynamics of Training



Just to give you an hint on where the community is headed with Deep Learning

## DALL-E OpenAI (January 2021)



## DALL-E OpenAI

TEXT PROMPT

an illustration of a baby daikon radish in a tutu walking a dog

AI-GENERATED IMAGES

TEXT PROMPT

an armchair in the shape of an avocado....

AI-GENERATED IMAGES

Edit prompt or view more images+

Edit prompt or view more images+

## OpenAI DALL-E - 12-billion parameters trained with self-supervision

Yikes!  $12 \times 10^9$  floating points parameters to train!

DALL-E is a **12-billion parameter** version of GPT-3 trained to generate images from text descriptions, using a dataset of text-image pairs. We've found that it has a diverse set of capabilities, including creating anthropomorphized versions of animals and objects, combining unrelated concepts in plausible ways, rendering text, and applying transformations to existing images.

## 0) Quick Intro to Optimization in Deep Learning

1) Network Structure: Multi-Layer Perceptron (MLP) is a Fully-Connected Neural Net

2) Backpropagation

1) Network Structure: Multi-Layer Perceptron (MLP)

is a Fully-Connected Neural Net

Networks and Topics that we do NOT cover

You will meet them at Deep Learning course

- Convolutional Neural Nets (good for images or any matrix data like as input)
- Generative Adversarial Networks (GAN) and adversarial training
- AutoEncoders or Variational Autoencoders
- Adversarial Attacks to NN

Networks and Topics we will cover

- Brief Recap on Feedforward NN
- Recurrent Neural Nets (RNN such as GRU, LSTM)
- Transformer Networks
- BERT, GPT-3

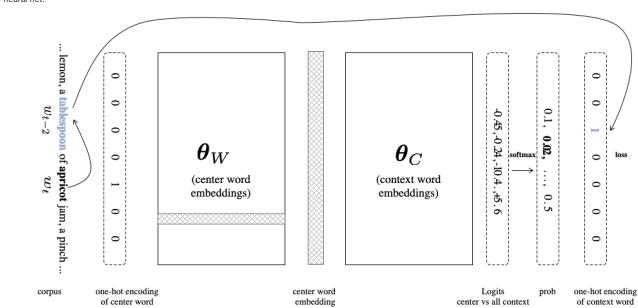
Let's go back to single layer, linear soft-max regression or "linear" neural network

Let's recall last classification layer of a neural net as pipeline

$$\mathbf{x} \implies \mathbf{z} = \mathbf{Wx} + \mathbf{b} \implies e^{\mathbf{z}} \implies \mathbf{p} = \frac{e^{\mathbf{z}}}{\sum_i e^{\mathbf{z}}} \implies -\ln(\mathbf{p}_y)$$

## word2vec with Skip-Gram at a glance

... and why it can be seen as a tiny neural net.



## Representation of a Single Layer

Let's consider our linear softmax regressor

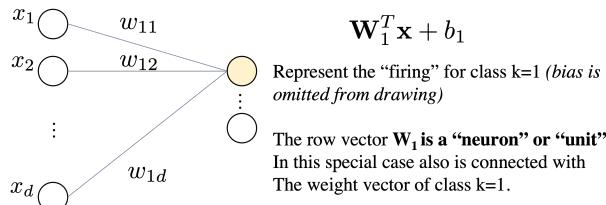
$$\frac{\mathbf{z}}{g^{(k)}} = \frac{\mathbf{W} \cdot \mathbf{x} + \mathbf{b}}{g^{(k)} \cdot g^{(k-1)} \cdot g^k}$$

We interpret as Linear Layer  $\mathbf{Wx} + \mathbf{b}$  followed by Non-Linear Activation function  $\sigma$

$$\sigma(\mathbf{Wx} + \mathbf{b}) = \sigma \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ \vdots & \cdots & \ddots & \vdots \\ w_{k1} & w_{k2} & \cdots & w_{kd} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{pmatrix} = \sigma \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \end{pmatrix}$$

## Representation of a Single Layer

$$\sigma(\mathbf{Wx} + \mathbf{b}) = \sigma \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ \vdots & \cdots & \ddots & \vdots \\ w_{k1} & w_{k2} & \cdots & w_{kd} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{pmatrix} = \sigma \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \end{pmatrix}$$

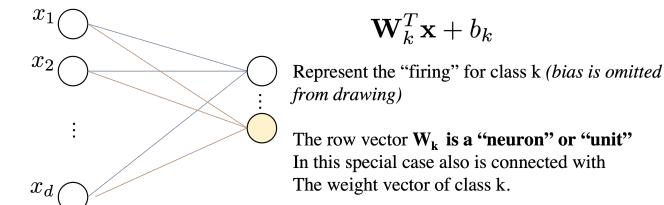


Represent the "firing" for class  $k=1$  (bias is omitted from drawing)

The row vector  $\mathbf{W}_1$  is a "neuron" or "unit"  
In this special case also is connected with  
The weight vector of class  $k=1$ .

## Representation of a Single Layer

$$\sigma(\mathbf{Wx} + \mathbf{b}) = \sigma \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ \vdots & \cdots & \ddots & \vdots \\ w_{k1} & w_{k2} & \cdots & w_{kd} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{pmatrix} = \sigma \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \end{pmatrix}$$

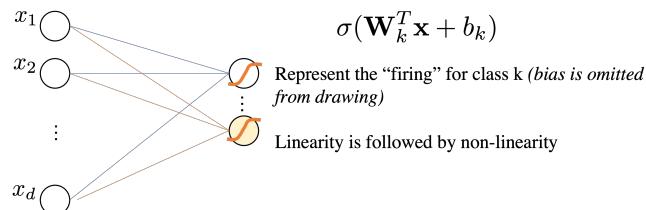


Represent the "firing" for class  $k$  (bias is omitted from drawing)

The row vector  $\mathbf{W}_k$  is a "neuron" or "unit"  
In this special case also is connected with  
The weight vector of class  $k$ .

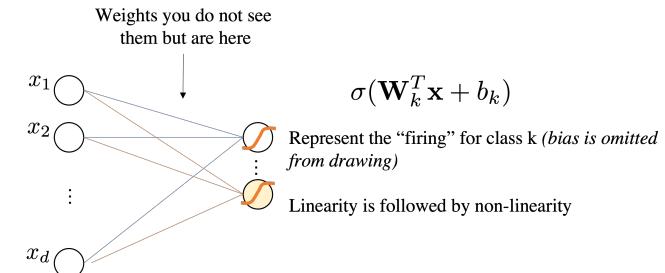
## Representation of a Single Layer: Linear plus non-Linear

$$\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) = \sigma * \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k1} & w_{k2} & \cdots & w_{kd} \\ \vdots & \vdots & \ddots & \vdots \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{pmatrix} = \sigma * \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \end{pmatrix}$$

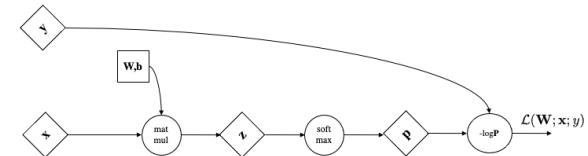


## Representation of a Single Layer: Linear plus non-Linear

$$\mathbf{w}_k = \begin{pmatrix} \text{unit} \\ \vdots \\ \text{unit} \\ \hline \mathbf{x} \\ \hline \text{unit} \end{pmatrix}$$

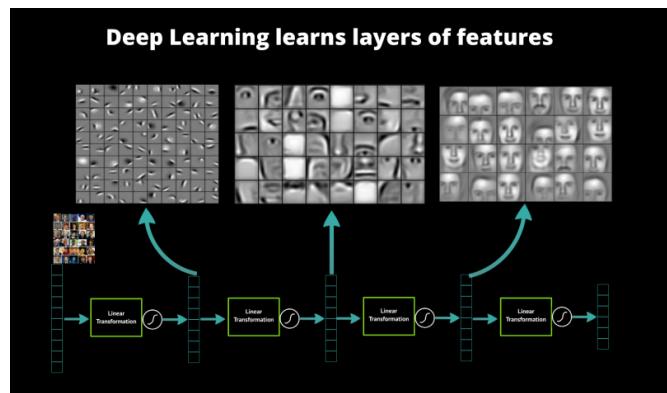


## Representation as a computational graph

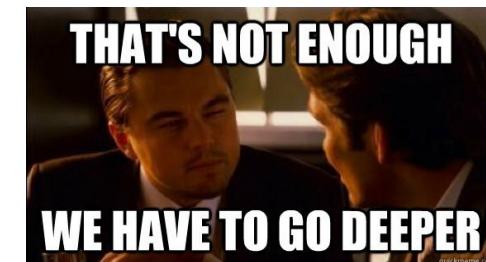


Damn, until now is all linear. So now the "Deep"!

- Damn, until now is all linear.
- Our beloved SoftMax+CE linear layer is there in the end (classifier).



A single linear layer is not enough for highly non-linear problems



## Adding another non-linear layer before the classifier

- We improve the expressiveness of our learned function by adding another linear layer **before** the classification layer.
- Think this new layer as a feature map  $\mathbf{x} \rightarrow \mathbf{v} \rightarrow \mathbf{y}$ : it maps our attribute to a feature space
- Now the classifier does not classify anymore directly  $\mathbf{v}$  but the feature  $\phi(\mathbf{v})$ .
- Sorry, notation becomes complex. Upper script means layer index; lower-script selects the unit
- $\mathbf{W}^1 \in \mathbb{R}^{d' \times d}, \mathbf{b}^1 \in \mathbb{R}^{d'}$  so then  $\mathbf{W}^2 \in \mathbb{R}^{m \times d'}, \mathbf{b}^2 \in \mathbb{R}^m$

$$p = \sigma(\mathbf{W}^2 \underbrace{\left( \sigma(\mathbf{W}^1 \mathbf{x} + \mathbf{b}^1) + \mathbf{b}^2 \right)}_{\phi(\mathbf{x})})$$

dim. analysis:  $d \mapsto p \mapsto k$

$\mathbf{W}^1 \in \mathbb{R}^{d \times p}$  is an Hidden Layer

Because it maps the original attribute in  $x$  from an dimensionality  $p$  and then  $p$  is used for classifying.

A priori you do not know what  $\mathbf{W}^1$  may learn.

$$p = \sigma(\mathbf{W}^2 \left( \sigma(\mathbf{W}^1 x + b^1) \right) + b^2)$$

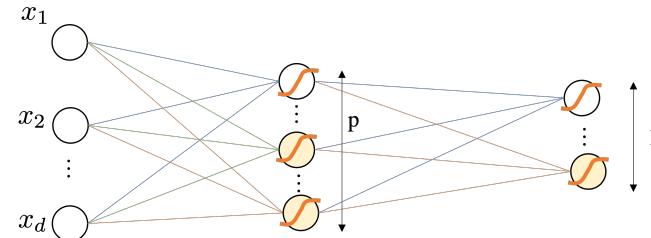
$$\phi(z)$$

dim analysis:  $d \mapsto p \mapsto k$

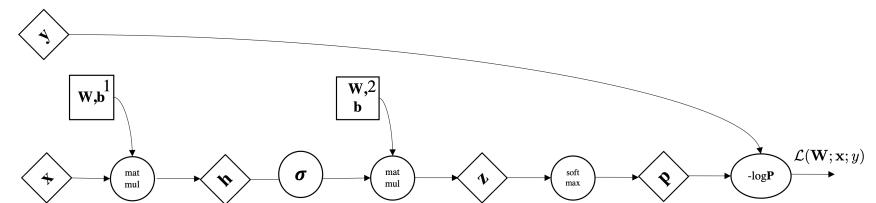
Let's update our visualizations

Multi-Layer Perceptron (MLP) with one hidden layer

Given the nature of these layers, they're called Fully-Connected NN



Multi-Layer Perceptron with one hidden layer

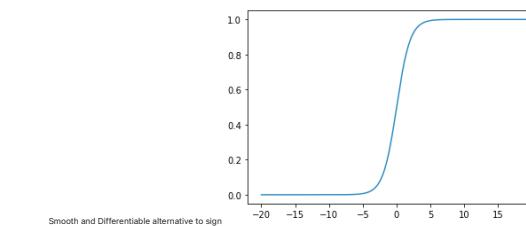


Non-linear activation functions: Sigmoid

Very important: Activation Functions are computed element-wise.

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

sigmoid or logistic function

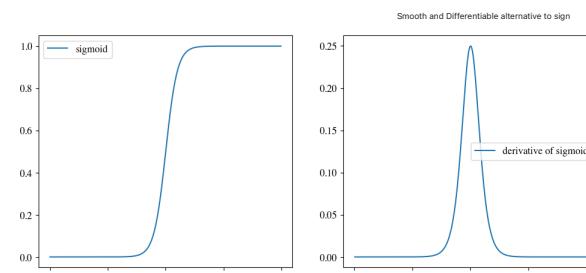


Non-linear activation functions: Sigmoid

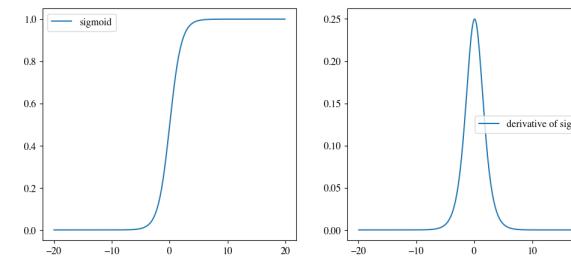
Very important: Activation Functions are computed element-wise.

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

sigmoid or logistic function



```
import numpy as np;
import matplotlib.pyplot as plt; step=0.1;
x = np.arange(-20.0, 20.0, step);
y = np.zeros(len(x));
for i in range(1, len(x), 1):
    y[i] = 1/(1+np.exp(-x[i]));
dy = np.diff(y);
dx = [step]*len(dy);
axes[0].legend(['sigmoid']);
axes[1].plot(x[1:],dy/step);
axes[1].legend(['derivative of sigmoid']);
```

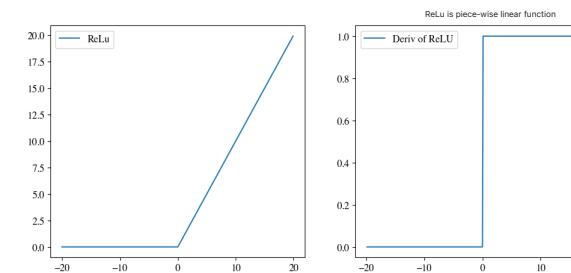


Non-linear activation functions: ReLU - Rectified Linear Unit

Very important: Activation Functions are computed element-wise.

$$\sigma(z) = \max(0, z)$$

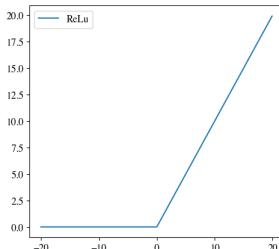
ReLU



```

In [14]:
import numpy as np;
import matplotlib.pyplot as plt; step=0.1;
x = np.arange(-20.0, 20.0, step);
fig, axes = plt.subplots(1,2,figsize=(12,5));
y = np.sin(x);
dy = np.diff(y);
axes[0].plot(x,y)
axes[0].set_title('ReLU')
axes[0].legend(['ReLU']);
axes[1].plot(dy)
axes[1].set_title('Deriv of ReLU')
axes[1].legend(['Deriv of ReLU']);

```



### Sigmoid

- Used to model output probability
- Nowdays not used in middle layers
- Have to compute  $\exp()$
- Vanishing gradients for large input magnitude

### ReLU

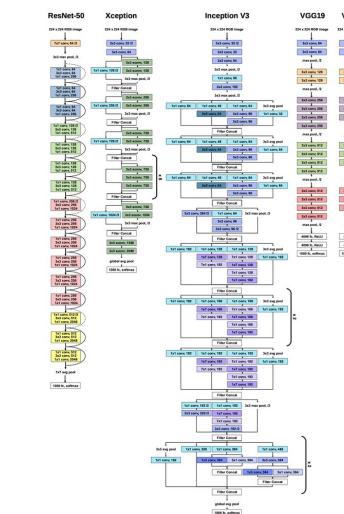
- Computationally efficient (no  $\exp()$ )
- No vanishing gradients but do not pass gradients for negative values
- Converge much faster than sigmoid (6x)
- Not differentiable in zero (subgradients)

There are other activation functions we do not cover

TanH, Leaky ReLU, parametrized ReLU, ELU

### Backpropagation and Differential Programming

NN can be huge composition of functions! 🍷



### Three ways of computing the gradients $\nabla_w L(x, y; w)$

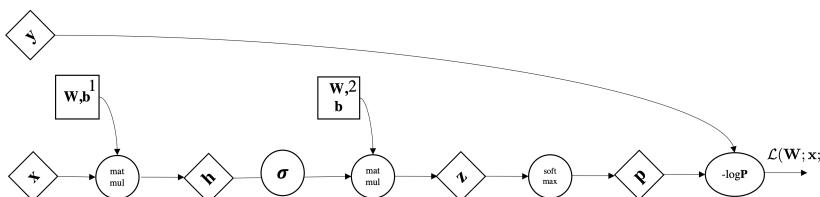
- Manually** (if we change the network, we have to adjust it for a 100 layer neural net) maybe not a good idea, does not scale, even if we use symbolic derivation tools such as Mathematica
- Finite Difference** good to check the gradients once you have an automatic way of computing it; **very slow, unfeasible in training!**
- Backpropagation:** application of chain rule of calculus to tensors with a computational graph with caching (**differential programming with automatic differentiation**)

### Backpropagation

Let's be clear on what we need to compute

$\forall l \in [1, \dots, L]$ :

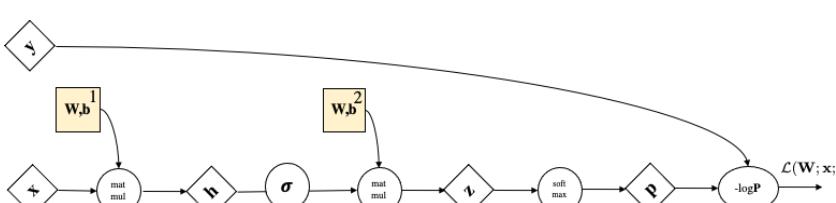
- $\nabla_{\mathbf{W}} L(x, y; (\mathbf{W}, b))$
- $\nabla_{\mathbf{b}} L(x, y; (\mathbf{W}, b))$



Once you have gradients on ALL weights  $\implies$  We can update

$\forall l \in [1, \dots, L]$ :

- $\mathbf{W}^l \leftarrow \mathbf{W}^l - \gamma \nabla_{\mathbf{W}^l} L(x, y; (\mathbf{W}, b))$
- $\mathbf{b}^l \leftarrow \mathbf{b}^l - \gamma \nabla_{\mathbf{b}^l} L(x, y; (\mathbf{W}, b))$



### How do we get all the weight updates?

Mostly taken from here

#### Chain Rule

Returning to functions of a single variable, suppose that  $y = f(g(x))$  and that the underlying functions  $y = f(u)$  and  $u = g(x)$  are both differentiable. The chain rule states that

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}$$

What is the derivative of loss wrt to x in the equation below?

$$y = \text{loss}(g(h(x)))$$

$$\frac{\partial \text{loss}}{\partial x} = \frac{\partial \text{loss}}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial x}$$

#### Chain Rule on Directed Acyclic Graph (DAG)

Automate the computation of derivatives with computer science

$$L(x, y, z) = (x + y)z$$

$$x = -2; y = 5; z = -4;$$

#### Now with Chain Rule

$$L(x, y, z) = (x + y)z$$

$$x = -2; y = 5; z = -4;$$

but we can re-write it with the chain rule:

$$\frac{\partial L(x, y, z)}{\partial x} = ((x + y)z) = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$$

#### Now with Chain Rule

$$L(x, y, z) = (x + y)z$$

$$x = -2; y = 5; z = -4;$$

but we can re-write it with the chain rule:

$$\frac{\partial L(x, y, z)}{\partial x} = ((x + y)z)' = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

## Now with Chain Rule

$$L(x, y, z) = (x + y)z$$

$$x = -2; y = 5; z = -4;$$

but we can re-write it with the chain rule:

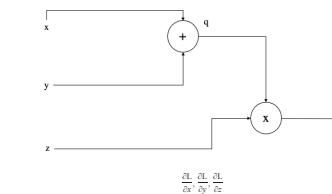
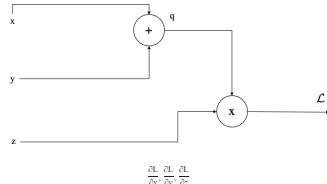
$$\frac{\partial L(x, y, z)}{\partial x} = ((x + y)z) \cdot \frac{\partial (x + y)z}{\partial x} = z \cdot 1 = z$$

## Chain Rule on Directed Acyclic Graph (DAG)

$$L(x, y, z) = (x + y)z$$

$$x = -2; y = 5; z = -4;$$

The computer science way:

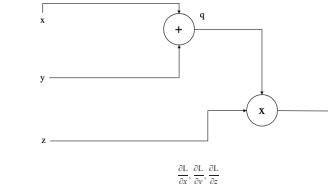


Who is input of  $L$ ?  
 $q$  and  $z$  are input of  $L$ .

$$\frac{\partial L}{\partial q}, \frac{\partial L}{\partial z}$$

## Chain Rule on Directed Acyclic Graph (DAG)

Even if the problem is very small, we break it down to subproblem so that we can **automate** it:

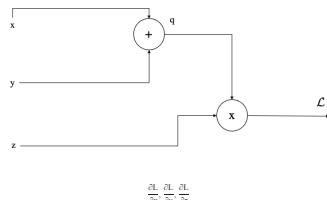


What is the derivate? (Just check the operation at the gate)

$$\frac{\partial L}{\partial q}$$

## Chain Rule on Directed Acyclic Graph (DAG)

Even if the problem is very small, we break it down to subproblem so that we can **automate** it:



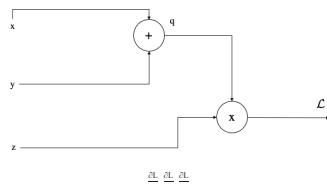
>

What is the derivate? (Just check the operation at the gate)?

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial z}{\partial z} = q$$

## Chain Rule on Directed Acyclic Graph (DAG)

Even if the problem is very small, we break it down to subproblem so that we can **automate** it:

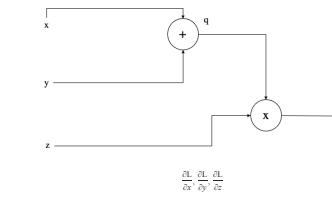


What is the derivate? (Just check the operation at the gate)?

$$\frac{\partial q}{\partial z} = 1, \frac{\partial q}{\partial y} = 1$$

## Chain Rule on Directed Acyclic Graph (DAG)

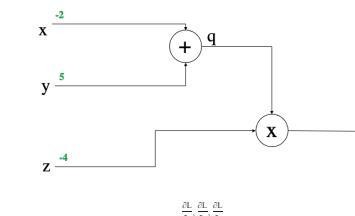
Even if the problem is very small, we break it down to subproblem so that we can **automate** it:



OK now we have all the **analytical "local"** partial derivatives, we can compute something

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial z}{\partial z} = q \cdot \frac{\partial z}{\partial z} = 1, \frac{\partial z}{\partial y} = 1$$

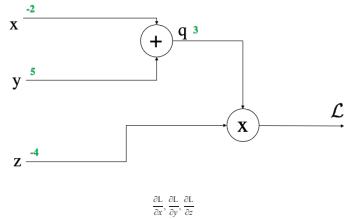
## Forward Pass



OK now we have all the **analytical** partial derivatives, we can compute something

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial z}{\partial z} = q \cdot \frac{\partial z}{\partial z} = 1, \frac{\partial z}{\partial y} = 1$$

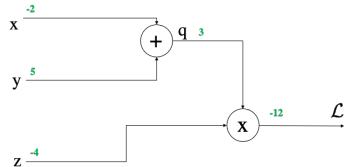
## Forward Pass



OK now we have all the analytical partial derivatives, we can compute something

$$\frac{\partial L}{\partial q} = z \frac{\partial L}{\partial z} + q \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Forward Pass



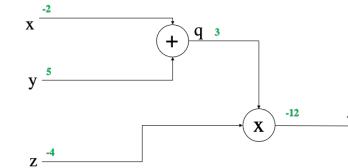
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \frac{\partial L}{\partial z} + q \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



Act as a base case for the recursion:

$$\frac{\partial L}{\partial L} = 1$$

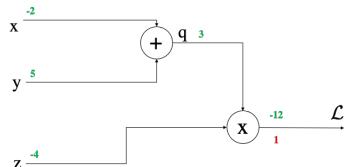
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \frac{\partial L}{\partial z} + q \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



Act as a base case for the recursion:

$$\frac{\partial L}{\partial L} = 1$$

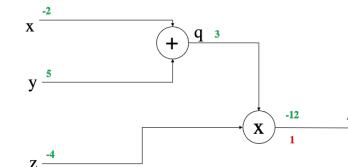
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \frac{\partial L}{\partial z} + q \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of L on z?:

$$\frac{\partial L}{\partial z}$$

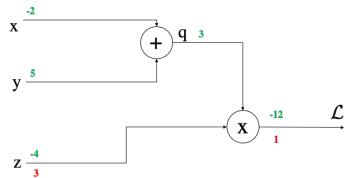
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \frac{\partial L}{\partial z} + q \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $L$  on  $z$ ?

$$\frac{\partial L}{\partial z} = q = 3$$

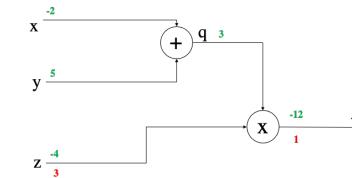
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $L$  on  $q$ ?

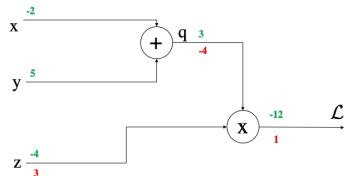
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $L$  on  $q$ ?

$$\frac{\partial L}{\partial q} = z = -4$$

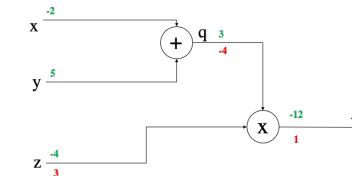
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $L$  on  $y$ ?

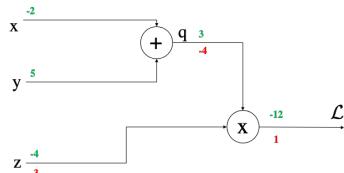
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $L$  on  $y$ ?

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial q} \frac{\partial q}{\partial y} = z \cdot 1 = -4$$

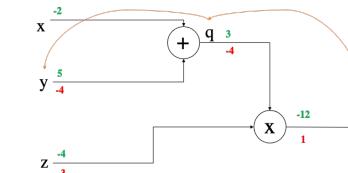
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z = \frac{\partial L}{\partial z} = q = 3, \frac{\partial L}{\partial y} = 1, \frac{\partial L}{\partial x} = -4$$

## Backward Pass



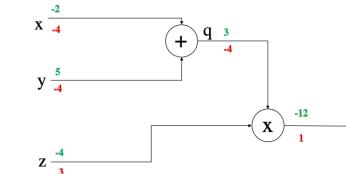
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = \frac{\partial L}{\partial z} = q = 3, \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

Check with our manual derivation



The high school way (as we did until now):

$$\begin{aligned}\frac{\partial L(x,y,z)}{\partial x} &= (xz + yz)' = (xz)' + (yz)' = z = -4 \\ \frac{\partial L(x,y,z)}{\partial y} &= (xz + yz)' = (xz)' + (yz)' = z = -4 \\ \frac{\partial L(x,y,z)}{\partial z} &= x + y = +3\end{aligned}$$

You know what? I do not trust math, I want to verify with a machine

## Pytorch check

```
from torch import tensor
def neural_net(x,y,z):
    return (x*y)*z
x, y, z = tensor(-2., requires_grad=True), tensor(5., requires_grad=True), tensor(-4., requires_grad=True)
loss = neural_net(x,y,z) # forward pass
loss.backward() # backward (after this I can check the gradients)
for el in [x,y,z]:
    print(el.grad)
```

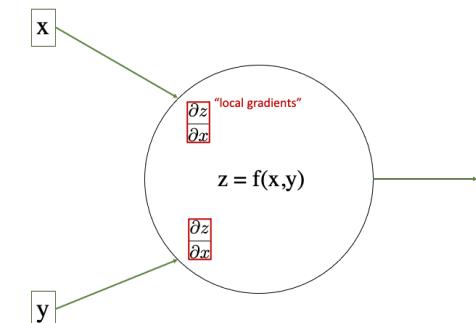
```
from torch import tensor
def neural_net(x,y,z):
    loss = (x*y)*z
    return loss
x, y, z = tensor(-2.), tensor(5., requires_grad=True), tensor(-4., requires_grad=True)
loss = neural_net(x,y,z) # forward pass
loss.backward() #backward (ok now I can check the gradients)
for el in [x,y,z]:
    print(el.grad)
```

```
from torch import tensor
def neural_net(x,y,z):
    loss = (x*y)*z
    return loss
x, y, z = tensor(-2.), tensor(5., requires_grad=True), tensor(-4., requires_grad=True)
loss = neural_net(x,y,z)
loss.backward()
for el in [x,y,z]:
    print(el.grad)
```

```
File "/var/folders/r7/lgn41t1489270pr_18qn1_c9000gp/T/ipykernel_72216/537812238.py", line 7
x, y, z = tensor(-2.), tensor(5., requires_grad=True), tensor(-4., requires_grad=True)
SyntaxError: invalid syntax
```

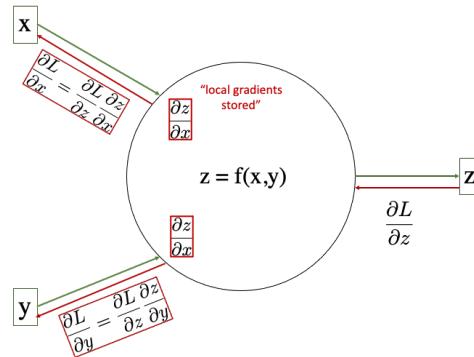
## General Recipe for Chain Rule over DAGs [Forward]

Just remember what you have to do at a generic gate:



## General Recipe for Chain Rule over DAGs [Backward]

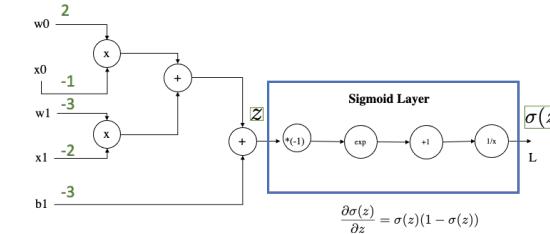
Multiply the gradient that you receive with your local gradient



## Logistic Regression Computational Graph could be simplified

$$f(w, x) = \frac{1}{1 + e^{-w_0 x_0 + w_1 x_1 + b}}$$

This is what implement the Sigmoid Layer in Pytorch:



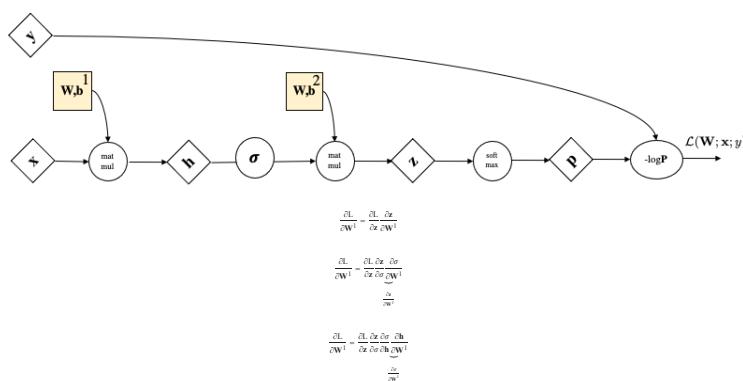
Are we done with training neural nets?

Not completely: till now scalars, but we have matrices and vectors!

Now this looks more familiar

$\forall i \in [1, \dots, L]:$

1.  $\mathbf{W}^i \leftarrow \mathbf{W}^i - \gamma \nabla_{\mathbf{W}^i} \mathcal{L}(\mathbf{x}, \mathbf{y}; (\mathbf{W}, \mathbf{b}))$
2.  $\mathbf{b}^i \leftarrow \mathbf{b}^i - \gamma \nabla_{\mathbf{b}^i} \mathcal{L}(\mathbf{x}, \mathbf{y}; (\mathbf{W}, \mathbf{b}))$



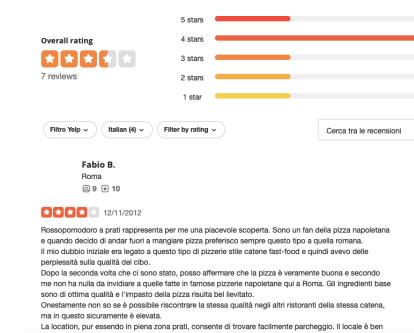
## Sentiment Analysis

- We focus on one common text categorization task: **sentiment analysis**.

- The extraction of sentiment, the positive or negative orientation that a writer expresses toward some object

A review of a movie, book, or product on the web expresses the author's sentiment toward the product, while an editorial or political text expresses sentiment toward a candidate or political action. Extracting consumer or public sentiment is thus relevant for fields from marketing to politics.

## Sentiment Analysis



## Natural Language Processing with Deep Learning

### NLP Task

#### Sentiment Analysis as a classification problem

## Sentiment Analysis: what we want

- + ...zany characters and richly applied satire, and some great plot twists
- It was pathetic. The worst part about it was the boxing scenes...
- + ...awesome caramel sauce and sweet toasty almonds. I love this place!
- ...awful pizza and ridiculously overpriced...

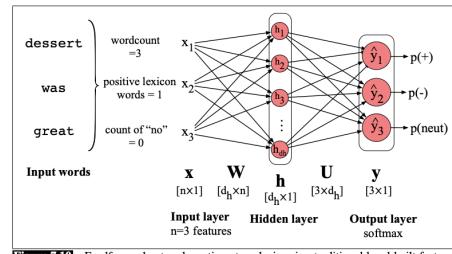
## Sentiment Analysis: what we have

	Cat	Documents
Training	-	just plain boring
	-	entirely predictable and lacks energy
	-	no surprises and very few laughs
	+	very powerful
Test	+	the most fun film of the summer
	?	predictable with no fun

## Feedforward Neural Net for Sentiment Analysis

w/ "hand-designed features"

$$\begin{aligned} \mathbf{x} &= [x_1, x_2, \dots, x_n] \quad (\text{each } x_i \text{ is a hand-designed feature}) \\ \mathbf{h} &= \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) \quad (\text{hidden layer}) \\ \mathbf{z} &= \mathbf{U}\mathbf{h} \quad (\text{classification layer}) \\ \hat{\mathbf{y}} &= \text{softmax}(\mathbf{z}) \end{aligned}$$



## Feedforward Neural Net for Sentiment Analysis

w/ word embeddings

We pool all the embeddings across the  $N$  word tokens we have, using arithmetic mean (early fusion).

$$\begin{aligned} \mathbf{x} &= \text{mean}(\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_N) \quad (\text{each } \mathbf{e}_i \text{ is a word2vec embedding}) \\ \mathbf{h} &= \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) \quad (\text{hidden layer}) \\ \mathbf{z} &= \mathbf{U}\mathbf{h} \quad (\text{classification layer}) \\ \hat{\mathbf{y}} &= \text{softmax}(\mathbf{z}) \end{aligned}$$

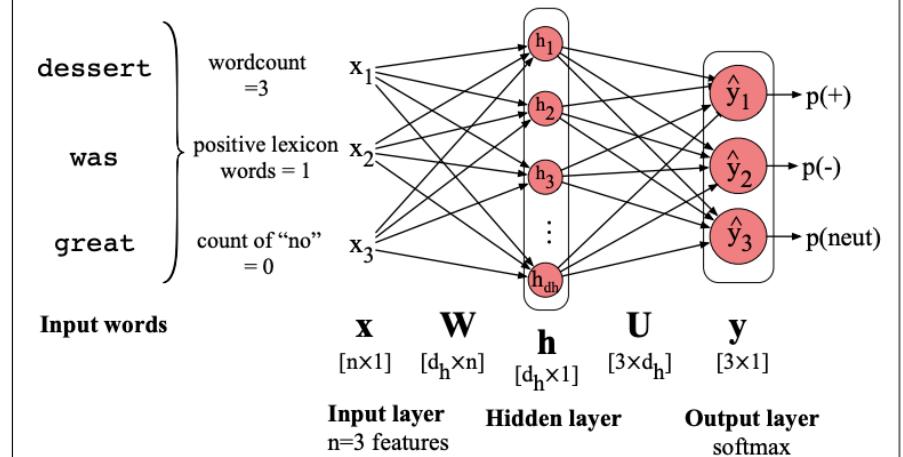


Figure 7.10 Feedforward network sentiment analysis using traditional hand-built features of the input text.

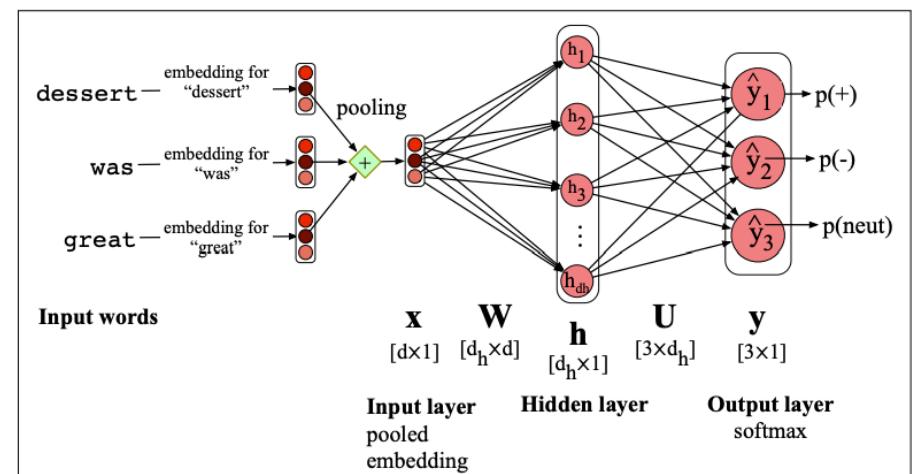
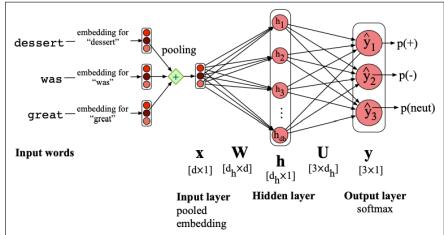


Figure 7.11 Feedforward sentiment analysis using a pooled embedding of the input words.

## Pretraining

Learning a representation on top of another is a very important concept which is related to **pretraining** in Deep Learning. You can think of two alternatives:

- Pretrain word2vec (**offline**): then tune a 2 layer NN on top of them
- Pretrain word2vec (**offline**) but also tune the word2vec embeddings as you train the 2 layer NN



**Figure 7.11** Feedforward sentiment analysis using a pooled embedding of the input words.

## NLP Task

### Neural Language Model with Feedforward Neural Net

Bengio et al., 2003

### Neural Language Model (Pros)

Neural language models can handle:

- much longer histories, can generalize better over contexts of similar words and;
- are more accurate at word-prediction.

### Neural Language Model (Cons)

On the other hand, neural net language models are:

- much more complex, are slower and need more energy to train
- are less interpretable than n-gram models

## Language Model (LM)

Language Modeling (LM) is the task of predicting what word comes next:

| the students opened their ...

books, laptops, exams, minds?

## Language Model (LM)

More formally, given a sequence of words  $w_1, \dots, w_r$  compute the probability distributions over the text words  $w_{r+1}$  where the support of the probability is a vocabulary  $V$ :

$$p(w_{r+1} | w_r, \dots, w_1)$$

where  $w_{r+1}$  can be any word in  $V = \{w_1, \dots, w_r\}$

A system that does this is called a **Language Model**

## Language Model (LM)

You can also think a LM as a machinery that assign a probability to a piece of text of T words:

$$p(w_r, \dots, w_1) = p(w_r, \dots, w_2 | w_1) \cdot p(w_1)$$

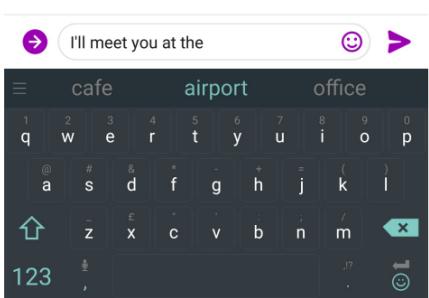
recursion

base case

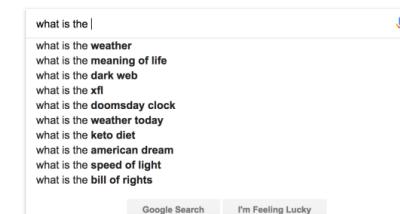
$$= p(w_r, \dots, w_2 | w_1) \cdot p(w_1) = p(w_r, \dots, w_3 | w_1, w_2) \cdot p(w_1, w_2) \cdot p(w_3 | w_2, w_1) \cdot p(w_1)$$

$$= p(w_r, \dots, w_3 | w_1, w_2) \cdot p(w_3 | w_2, w_1) \cdot p(w_r, \dots, w_4 | w_1, w_2, w_3) \cdot p(w_1, w_2, w_3) \cdot p(w_4 | w_1, w_2, w_3, w_4)$$

We use LM everyday!



We use LM everyday!



Language Model: the order matters!

| the cat is small

$$\frac{0 \quad 1 \quad 2 \quad 3}{\text{the cat is small}}$$

$p(\text{the, cat, is, small}) = 0.35$

| small the is cat

$$\frac{0 \quad 1 \quad 2 \quad 3}{\text{small the is cat}}$$

$p(\text{small, the, cat, is}) = 0.0034$

## Language Model: N-gram

- the students opened their ...
- unigrams: "the", "students", "opened", "their"
- bigrams: "the students", "students opened", "opened their"
- trigrams: "the students opened", "students opened their"
- four-grams: "the students opened their"

$$p(\text{the} | \text{its water is so transparent that}) = \frac{\text{count}(\text{its water is so transparent that the})}{\text{count}(\text{its water is so transparent that})}$$

## Language Model: N-gram problems

As you increase the "window" of your context:

1. Sparsity issue increases
2. Storage problem increases

Neural LM took over N-gram

## A fixed-window Neural LM

- Like in N-gram, we consider only a portion of text / N window or context size, not the entire text.

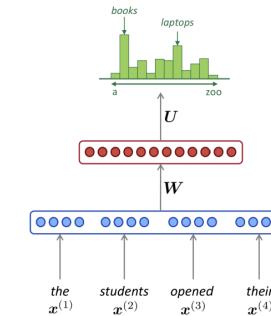
$$p(w_{t+1} | w_p, \dots, w_t) \approx p(w_{t+1} | w_p, \dots, w_{t-N})$$

- Given the context we use a Feedforward NN to predict the next word (similarly to word2vec).

As the prector started the class, the students opened their ...

- Unlike N-gram, it uses **distributed representations**

## A fixed-window Neural LM



## A fixed-window Neural LM

### output distribution

$$\hat{y} = \text{softmax}(U\mathbf{h} + b_2) \in \mathbb{R}^{|V|}$$

### hidden layer

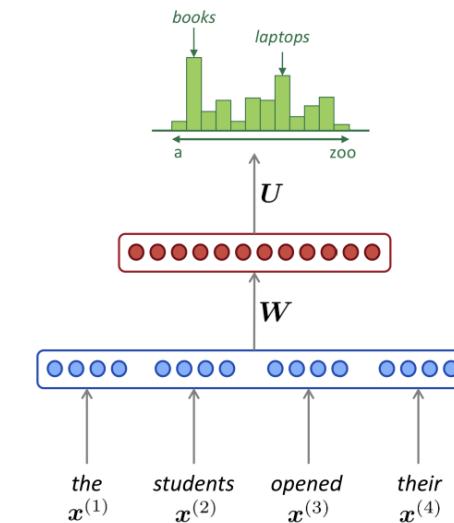
$$\mathbf{h} = f(W\mathbf{e} + b_1)$$

### concatenated word embeddings

$$\mathbf{e} = [\mathbf{e}^{(1)}; \mathbf{e}^{(2)}; \mathbf{e}^{(3)}; \mathbf{e}^{(4)}]$$

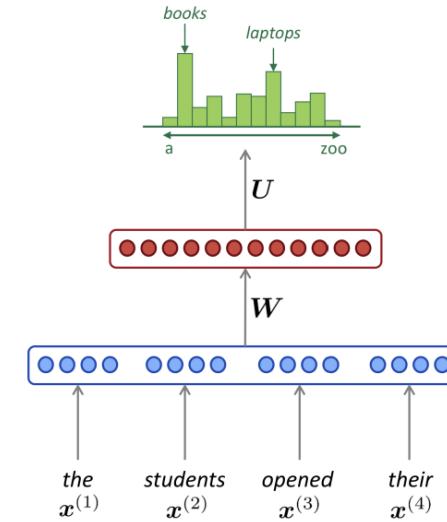
### words / one-hot vectors

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$$



## A fixed-window Neural LM

- Improvements over n-gram LM
  - No sparsity problem
  - No need to store all observed n-grams!
- Remaining problems:
- Fixed window is too small
  - Enlarging window enlarges  $W$
  - Window can never be large enough!
  - No symmetry in how the inputs are processed.



## A fixed-window Neural LM:

Journal of Machine Learning Research 3 (2003) 1137–1155

Submitted 4/02; Published 2/03

**A Neural Probabilistic Language Model**

Yoshua Bengio  
Réjean Ducharme  
Pascal Vincent  
Christian Jauvin  
Département d'Informatique et Recherche Opérationnelle  
Centre de Recherche Mathématiques  
Université de Montréal, Montréal, Québec, Canada

Editors: Jaz Kandola, Thomas Hofmann, Tomaso Poggio and John Shawe-Taylor

### Abstract

A goal of statistical language modeling is to learn the joint probability function of sequences of words in a language. This is intrinsically difficult because of the **curse of dimensionality**: a word sequence on which the model will be tested is likely to be different from all the word sequences seen during training. Traditional but very successful approaches based on n-grams obtain generalization by concatenating very short overlapping sequences seen in the training set. We propose to fight the curse of dimensionality by **learning a distributed representation for words** which allows each training sentence to inform the model about an exponential number of semantically neighboring sentences. The model learns simultaneously (1) a distributed representation for each word along with (2) the probability function for word sequences, expressed in terms of these representations.

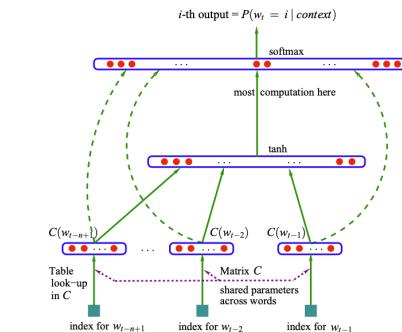
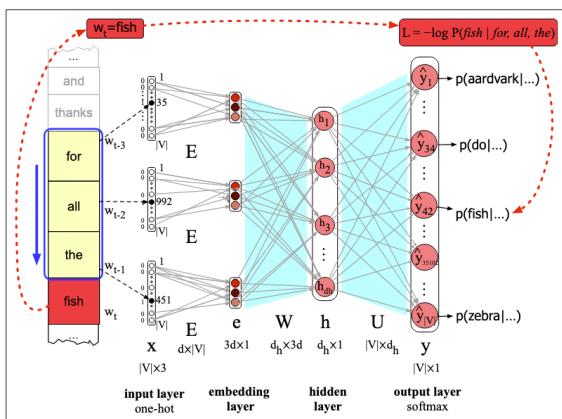


Figure 1: Neural architecture:  $f(i, w_{t-1}, \dots, w_{t-n+1}) = g(i, C(w_{t-1}), \dots, C(w_{t-n+1}))$  where  $g$  is the neural network and  $C(i)$  is the  $i$ -th word feature vector.

## Training the Neural LM: sliding window



**Figure 7.18** Learning all the way back to embeddings. Again, the embedding matrix  $E$  is shared among the 3 context words.

## A fixed-window Neural LM: What we need to move forward

- An architecture that can process **any arbitrary length** of context
- Handle word order in a **better way** than concatenating embeddings.

## Appendix on Deep Learning

```
> Original unshuffled training set [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
> Training set [10, 8, 1, 2, 6, 4, 9, 7, 5, 3]
[Epoch 0]
SGD step taken over [10, 8, 1]
SGD step taken over [2, 6, 4]
SGD step taken over [9, 7, 5]
[Training set [10, 6, 9, 3, 7, 8, 4, 5, 2]
[Epoch 1]
SGD step taken over [1, 10, 6]
SGD step taken over [9, 3, 7]
SGD step taken over [8, 4, 5]
[Training set [6, 3, 10, 5, 9, 8, 4, 7, 2, 1]
[Epoch 2]
SGD step taken over [6, 3, 10]
SGD step taken over [5, 9, 8]
SGD step taken over [4, 7, 2]
[Training set [1, 2, 5, 10, 6, 7, 9, 8, 3, 4]
[Epoch 3]
SGD step taken over [1, 2, 5]
SGD step taken over [10, 6, 7]
SGD step taken over [9, 8, 3]
[Training set [2, 3, 1, 9, 6, 8, 4, 10, 7, 5]
[Epoch 4]
```

## Images - Mini-batch is a tensor $H \times W \times 3 \times k$

as an example with RGB images of size  $H \times W$ , you have a tensor that contains  $k$  images in the mini-batch

## Video Frames - Mini-batch is a tensor $H \times W \times 3 \times t \times k$

as set of frames from a video, you have a tensor that contains  $t$  frames over  $k$  time instants in the mini-batch

```
from random import shuffle
training = list(range(1,11))
shuffle(training)
print("Original unshuffled training set", training)
converge, it, max_it, k, epoch = False, 0, 100, 3, 0
while not converge and it < max_it: # you training converge scheme
    print("Training set", training)
    for b in range(len(training), k): # Data Loader
        mini_batch = training[b:k]
        if len(mini_batch) != k:
            continue
        print("SGD step taken over", mini_batch) # compute the loss and update your model
        it += 1
    epoch += 1
    if epoch == 10: # epoch is done we reshuffle the training set
        shuffle(training)
print("> Training set", training)
```

## Gradient Descent or Batch GD

- Compute the gradient of the loss wrt to params for **all  $n$  training samples**
- $\theta \leftarrow \theta - \gamma \sum_{i=1}^n \nabla_{\theta} l(\theta, x_i, y_i)$

## Stochastic Gradient Descent or SGD

- Compute the gradient of the loss wrt to params for a **single random training samples**
- $\theta \leftarrow \theta - \gamma \nabla_{\theta} l(\theta, x_j, y_j)$

## How to optimize a Neural Net - SGD over mini-batches

1. In-between Batch GD and SGD with a single sample
2. We load randomly  $k$  samples over the  $n$ ; usually  $k$  is a power of 2.
  - mini batch of 32, 64, 128 but could also be 1000
3.  $\theta \leftarrow \gamma \sum_{i=1}^k \nabla_{\theta} l(\theta, x_i, y_i)$
4. Practically you take your training set  $X$  and you **shuffle** it, then go over it  $k$  by  $k$ . Simulate uniform random sampling without replacement.
  - When the list is over, re-shuffle and shuffle again.
5. When you have performed a full pass on the shuffled data, this is called an **EPOCH**
6. You can train NN over iterations or over **EPOCHS**

## Training scheme Pseudo-code

```
from random import shuffle
training = list(range(1,11)) # each index points to a training sample, could be a matrix x:HxWx3, label y
while not converge and it < max_it: # you training convergence scheme
    for b in range(len(training), k): # Data Loader gives you a batch  $k$  x matrices
        mini_batch = training[b:b+k] # so mini-batch is a tensor HxWx3k
        if len(mini_batch) != k: # a possible way of handling the offset
            continue
        print("SGD step taken over", mini_batch) # compute the loss/gradients and update your model
        loss.backward() # get the gradients
        optimizer.step() # Incorporate in the model
        # check convergence and set it to True
        it += 1
    epoch += 1 # an epoch is done, we reshuffle the training set
    shuffle(training)
```

## 1) SGD over mini-batches

### 1. Initialization - Very Important if the function is not strictly convex

$\theta \sim N(\cdot)$  omit details for now

With NN random initialization from a distribution (There are different methods). We do not set them all to zero

### 2. Repeat until convergence:

- Compute the gradient of the loss wrt to the parameters  $\theta$  given the **mini-batch**
- Take a small step in the opposite direction of steepest ascent (**so steepest descent**).

$$\theta \leftarrow \theta - \gamma \sum_{i=1}^k \nabla_{\theta} l(\theta, x_i, y_i)$$

3. When convergence is reached, you final estimate is in  $\theta$

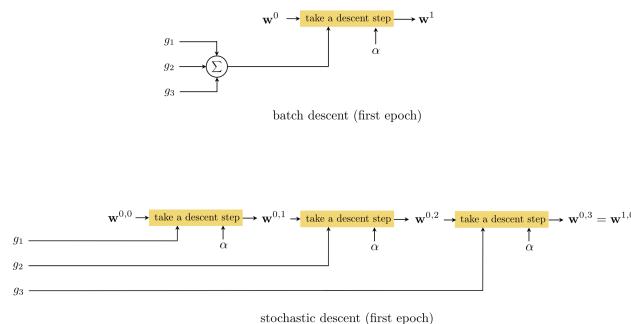
## Change of vocabulary - A bunch of training samples is a mini-batch

- We train NN **Stochastic Gradient Descent** over mini-batches with momentum (or variations thereof)
- When you train NN you 'sample' a mini-batch  $X_b$  from your big dataset  $X$ .

Below this holds for the final linear layer:

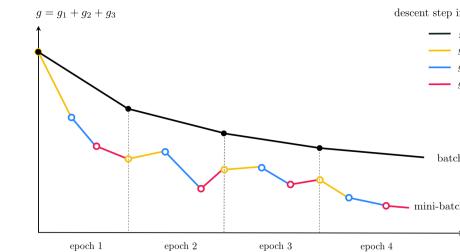
$$\frac{\mathbf{y}}{\mathbf{R}^{k \times n}} = \frac{\mathbf{W}^T \mathbf{x}_b + \mathbf{b}}{\mathbf{R}^{k \times d} \mathbf{R}^{d \times n} \mathbf{R}^k}$$

## Mini-Batch, Visually



## Mini-Batch SGD vs [Batch] GD

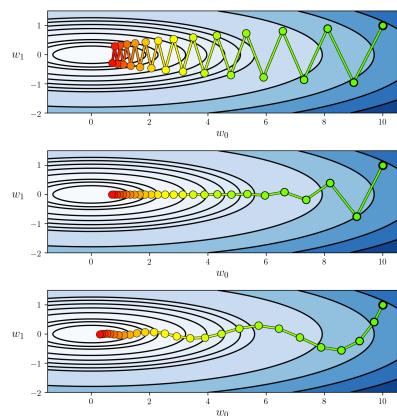
Loss in NN in **non-convex** with lots of local-minima so stochasticity adds noise that let the optimization escape from local minima.



Mini-batch is a sort of smoothing of the single point SGD

There is another smoothing technique: Momentum

Top: SGD; Bottom: SGD with momentum increasing memory of previous steps



## SGD over mini-batches with Momentum

1. Initialization - Very Important if the function is not strictly convex

$\theta \sim N$  omit details

With NN random initialization from a distribution (There are different methods). We do not set them all to zero

2. Repeat until convergence:

- Compute the gradient of the loss wrt to the parameters  $\theta$  given the **mini-batch**
- Take a small step in the opposite direction of steepest ascent (**so steepest descent**).

$$\begin{aligned}\Delta_{t+1} &= \alpha \Delta_t + (1-\alpha) \sum_{i=1}^k \nabla \varphi(\theta; \mathbf{x}_i, \mathbf{y}_i) \\ &\text{new update} \\ \theta &\leftarrow \theta - \gamma \Delta_{t+1}\end{aligned}$$

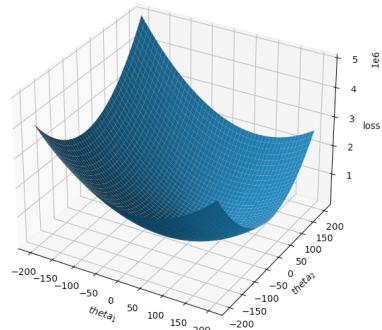
3. When convergence is reached (or **EARLY STOPPING**), your final estimate is in  $\theta$

## SGD over mini-batches with Momentum

- We introduce an additional term to remember what happened to the gradient in the previous iteration.
- This memory dampens oscillations and smoothes out the gradient updates.
- The memory is implemented with an **exponential moving average**
- Usually  $\alpha$  (the memory param) is set to 0.9, a good value.

$$\begin{aligned}\Delta_t &= \alpha \Delta_{t-1} + (1-\alpha) \sum_{i=1}^k \nabla \varphi(\theta; \mathbf{x}_i, \mathbf{y}_i) \\ &\text{new update} \\ \theta &\leftarrow \theta - \gamma \Delta_t\end{aligned}$$

Loss Surface for Linear Regression  $\ell_2^2$  loss with  $d = 2$  parameters in  $\theta$



Loss Surface for ResNet-20 with no skip connection on ImageNet

ResNet-20, number of parameters  $\theta$  of the order of.....millions!

GPT-3 (LM behind chatGPT) has 150 billions parameters

Visualization of mode connectivity for ResNet-20 with no skip connections on ImageNet dataset. The visualization by Javier Ideami

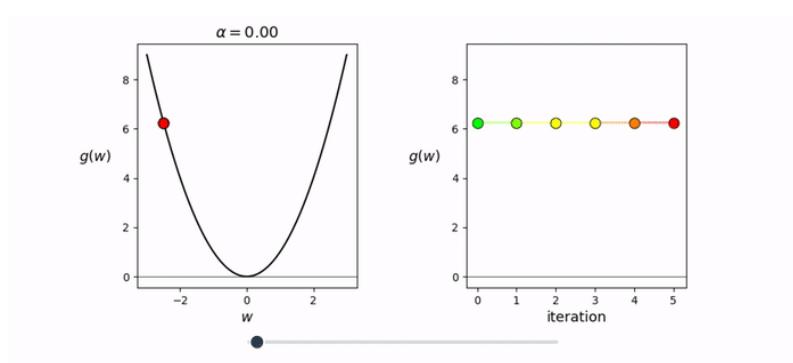


Taken from [https://zmailovpavel.github.io/curves\\_blogpost/](https://zmailovpavel.github.io/curves_blogpost/)

Video for the curious student

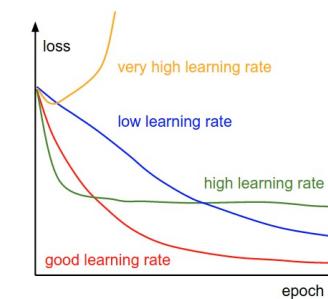
With Deep Learning optimization is highly non-convex and #params explode!

Learning rate is very important



Babysitting the training process

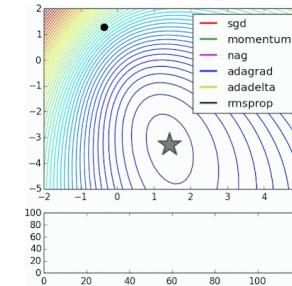
Loss in function of epochs



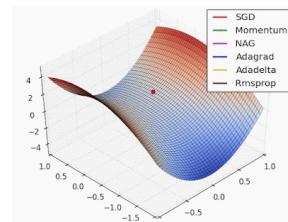
## Valleys, Hills, Noisy Surface



## Dynamics of Training



## Dynamics of Training



Long valley: Algos without scaling based on gradient information really struggle to break symmetry here - SGD gets no where and Nesterov Accelerated Gradient / Momentum exhibits oscillations until they build up velocity in the optimization direction. Algos that scale step size based on the gradient quickly break symmetry and begin descent.

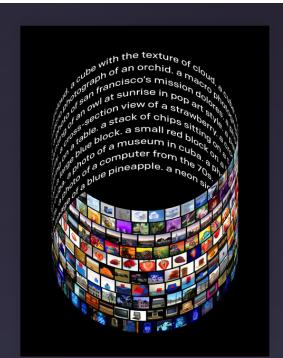
Just to give you an hint on where the community is headed with Deep Learning

## DALL-E OpenAI (January 2021)

### DALL·E: Creating Images from Text

We've trained a neural network called DALL·E that creates images from text captions for a wide range of concepts expressible in natural language.

January 5, 2021  
27 minute read





## OpenAI DALL-E - 12-billion parameters trained with self-supervision

Yikes!  $12 \times 10^9$  floating points parameters to train

DALLE is a **12-billion parameter** version of GPT-3 trained to generate images from text descriptions, using a dataset of text-image pairs. We've found that it has a diverse set of capabilities, including creating anthropomorphized versions of animals and objects, combining unrelated concepts in plausible ways, rendering text, and applying transformations to existing images.

## 0) Quick Intro to Optimization in Deep Learning

1) Network Structure: Multi-Layer Perceptron (MLP) is a Fully-Connected Neural Net

2) Backpropagation

## 1) Network Structure: Multi-Layer Perceptron (MLP)

is a Fully-Connected Neural Net

## Networks and Topics that we do NOT cover

You will meet them at Deep Learning course

- Convolutional Neural Nets (good for images or any matrix data like as input)
- Generative Adversarial Networks (GAN) and adversarial training
- AutoEncoders or Variational Autoencoders
- Adversarial Attacks to NN

## Networks and Topics we will cover

- Brief Recap on Feedforward NN
- Recurrent Neural Nets (RNN such as GRU, LSTM)
- Transformer Networks
- BERT, GPT-3

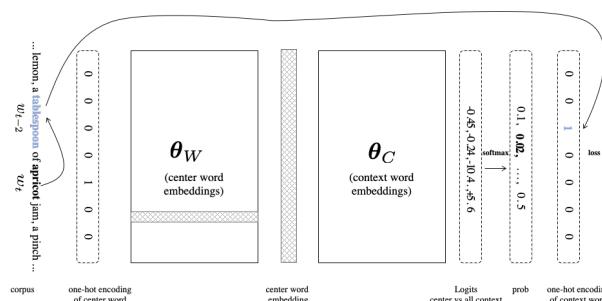
Let's go back to single layer, linear soft-max regression or "linear" neural network

Let's recall last classification layer of a neural net as pipeline

$$\mathbf{x} \implies \mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b} \implies e^{\mathbf{z}} \implies \mathbf{p} = \frac{e^{\mathbf{z}}}{\sum_i e^{\mathbf{z}}} \implies -\ln(p_y)$$

## word2vec with Skip-Gram at a glance

... and why it can be seen as a tiny neural net.



## Representation of a Single Layer

Let's consider our linear softmax regressor

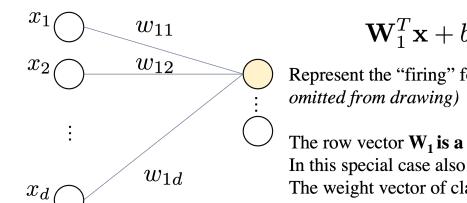
$$\begin{aligned} \mathbf{z} &= \mathbf{W} \mathbf{x} + \mathbf{b} \\ \mathbf{g}^{k+1} &= \mathbf{h}^k \mathbf{g}^{k+1} \end{aligned}$$

We interpret as Linear Layer  $\mathbf{Wx} + b$  followed by Non-Linear Activation function  $\sigma$ 

$$\sigma(\mathbf{Wx} + \mathbf{b}) = \sigma \cdot \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ \vdots & \cdots & \ddots & \vdots \\ w_{k1} & w_{k2} & \cdots & w_{kd} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{pmatrix} = \sigma \cdot \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \end{pmatrix}$$

## Representation of a Single Layer

$$\sigma(\mathbf{Wx} + \mathbf{b}) = \sigma \cdot \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ \vdots & \cdots & \ddots & \vdots \\ w_{k1} & w_{k2} & \cdots & w_{kd} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \end{pmatrix} = \sigma \cdot \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \end{pmatrix}$$

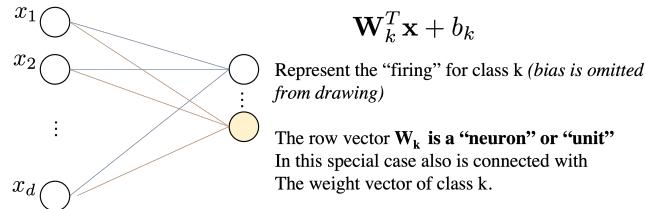


Represent the "firing" for class  $k=1$  (*bias is omitted from drawing*)

The row vector  $\mathbf{W}_1$  is a "neuron" or "unit". In this special case also is connected with The weight vector of class  $k=1$ .

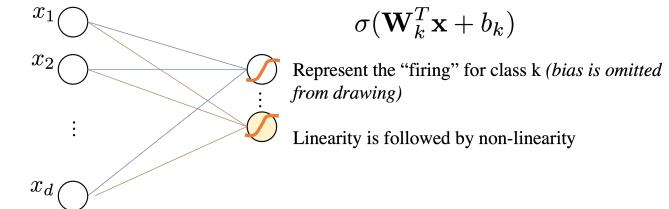
## Representation of a Single Layer

$$\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) = \sigma * \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k1} & w_{k2} & \cdots & w_{kd} \\ \vdots & \vdots & \ddots & \vdots \\ w_{d1} & w_{d2} & \cdots & w_{dd} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_d \end{pmatrix} = \sigma * \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \\ \vdots \\ z_d \end{pmatrix}$$



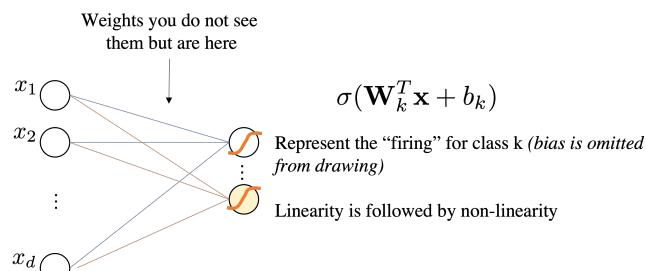
## Representation of a Single Layer: Linear plus non-Linear

$$\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) = \sigma * \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k1} & w_{k2} & \cdots & w_{kd} \\ \vdots & \vdots & \ddots & \vdots \\ w_{d1} & w_{d2} & \cdots & w_{dd} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_d \end{pmatrix} = \sigma * \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \\ \vdots \\ z_d \end{pmatrix}$$

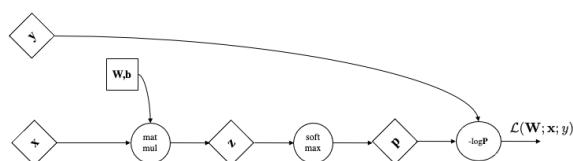


## Representation of a Single Layer: Linear plus non-Linear

$$\mathbf{w}\mathbf{x} = \begin{pmatrix} -\text{unit} \\ \vdots \\ -\text{unit} \end{pmatrix} \begin{pmatrix} 1 \\ \mathbf{x} \\ 1 \end{pmatrix}$$

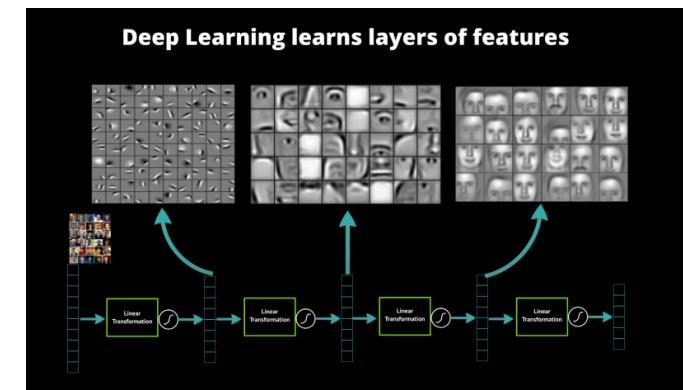


## Representation as a computational graph



## Damn, until now is all linear. So now the "Deep"!

- Damn, until now is all linear.
- Our beloved SoftMax+CE linear layer is there in the end (classifier).



A single linear layer is not enough for highly non-linear problems



### Adding another non-linear layer before the classifier

- We improve the expressiveness of our learned function by adding another linear layer **before** the classification layer.
- Think this new layer as a feature map  $x \rightarrow \phi(x)$ : it maps our attribute to a feature space
- Now the classifier does not classify anymore directly  $x$  but the feature  $\phi(x)$ .
- Sorry, notation becomes complex. Upper script means layer index; lower-script selects the unit
- $W^1 \in \mathbb{R}^{d \times p}$ ,  $b^1 \in \mathbb{R}^p$  so then  $W^2 \in \mathbb{R}^{p \times k}$ ,  $b^2 \in \mathbb{R}^k$

$$\begin{aligned} p &= \sigma(W^2 \underbrace{\left(\sigma(W^1 x + b^1)\right)}_{\phi(x)} + b^2) \\ &\quad \downarrow \\ \text{dim analysis: } d &\mapsto p \mapsto k \end{aligned}$$

$W^1 \in \mathbb{R}^{d \times p}$  is an Hidden Layer

Because it maps the original attribute in  $x$  from an dimensionality  $d$  and then  $p$  is used for classifying.

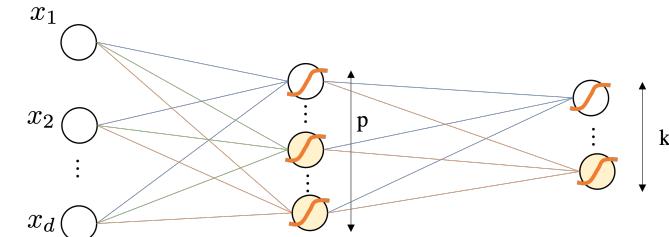
A priori you do not know what  $W^1$  may learn.

$$\begin{aligned} p &= \sigma(W^2 \underbrace{\left(\sigma(W^1 x + b^1)\right)}_{\phi(x)} + b^2) \\ &\quad \downarrow \\ \text{dim analysis: } d &\mapsto p \mapsto k \end{aligned}$$

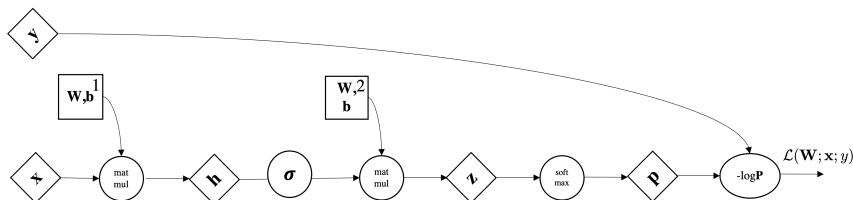
Let's update our visualizations

Multi-Layer Perceptron (MLP) with one hidden layer

Given the nature of these layers, they're called Fully-Connected NN



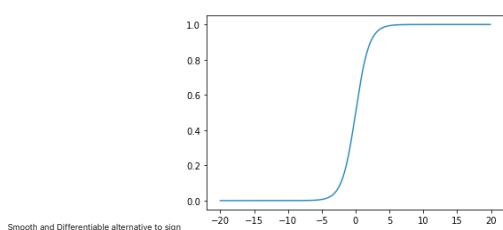
### Multi-Layer Perceptron with one hidden layer



### Non-linear activation functions: Sigmoid

Very important: Activation Functions are computed element-wise.

$$\sigma(z) = \frac{1}{1 + \exp^{-z}} \quad \text{sigmoid or logistic function}$$



### Non-linear activation functions: Sigmoid

Very important: Activation Functions are computed element-wise.

$$\sigma(z) = \frac{1}{1 + \exp^{-z}} \quad \text{sigmoid or logistic function}$$

Smooth and Differentiable alternative to sign

```
{import numpy as np; import matplotlib.pyplot as plt; step=0.1; x = np.arange(-20.0, 20.0, step); fig, axes = plt.subplots(1,2,figsize=(12,5)); y = 1/(1+np.exp(-x)); dy = np.diff(y); axes[0].plot(x[1:],dy/step); axes[0].legend(['derivative of sigmoid']);}
```



### Non-linear activation functions: ReLU - Rectified Linear Unit

Very important: Activation Functions are computed element-wise.

$$\sigma(z) = \max(0, z) \quad \text{ReLU}$$

ReLU is piece-wise linear function

```
{import numpy as np; import matplotlib.pyplot as plt; step=0.1; x = np.arange(-20.0, 20.0, step); fig, axes = plt.subplots(1,2,figsize=(12,5)); y = np.maximum(0,x); dy = np.diff(y); axes[0].plot(x,y); axes[1].plot(x[1:],dy/step); axes[0].legend(['ReLU']); axes[1].legend(['Deriv of ReLU']);}
```



### Sigmoid

- Used to model output probability
- Nowdays not used in middle layers
- Have to compute  $\exp(z)$
- Vanishing gradients** for large input magnitude

## ReLU

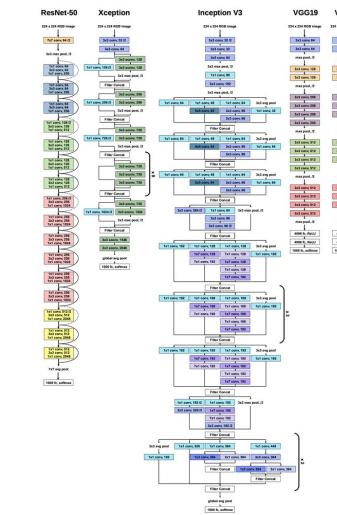
- Computationally efficient (no exp)
- No vanishing gradients but do not let pass gradients for negative values
- Converge much faster than sigmoid ( $\text{f}_x$ )
- Not differentiable in zero (subgradients)

There are other activation functions we do not cover

TanH, Leaky ReLU, parametrized ReLU, ELU

Backpropagation and Differential Programming

NN can be huge composition of functions! 🍷



Three ways of computing the gradients  $\nabla_w \mathcal{L}(x, y; w)$

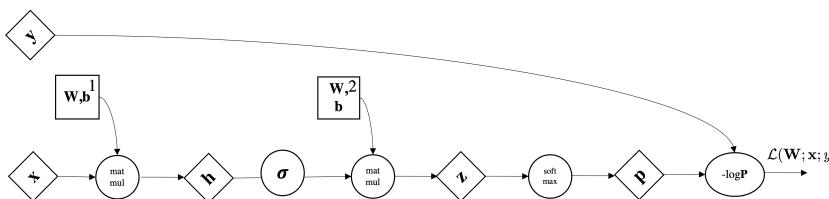
1. **Manually** (if we change the network, we have to adjust it for a 100 layer neural net) maybe not a good idea, does not scale, even if we use symbolic derivation tools such as Mathematica 📈
2. **Finite Difference** good to check the gradients once you have an automatic way of computing it; **very slow, unfeasible in training!** 🚫
3. **Backpropagation**: application of chain rule of calculus to tensors with a computational graph with caching (differential programming with automatic differentiation) 🚫

## Backpropagation

Let's be clear on what we need to compute

$\forall l \in [1, \dots, L]$ :

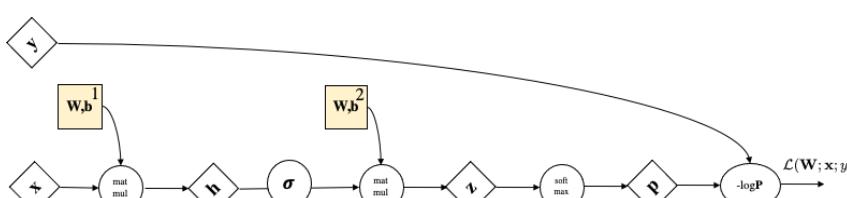
1.  $\nabla_{\mathbf{W}} \mathcal{L}(x, y; (\mathbf{W}, b))$
2.  $\nabla_b \mathcal{L}(x, y; (\mathbf{W}, b))$



Once you have gradients on ALL weights  $\implies$  We can update

$\forall l \in [1, \dots, L]$ :

1.  $\mathbf{W}^l \leftarrow \mathbf{W}^l - \gamma \nabla_{\mathbf{W}} \mathcal{L}(x, y; (\mathbf{W}, b))$
2.  $\mathbf{b}^l \leftarrow \mathbf{b}^l - \gamma \nabla_b \mathcal{L}(x, y; (\mathbf{W}, b))$



How do we get all the weight updates?

Mostly taken from here

Chain Rule

Returning to functions of a single variable, suppose that  $y = f(g(x))$  and that the underlying functions  $y = f(u)$  and  $u = g(x)$  are both differentiable. The chain rule states that

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}$$

What is the derivative of loss wrt to x in the equation below?

$$y = \text{loss}(g(h(x)))$$

$$\frac{\partial \text{loss}}{\partial x} = \frac{\partial \text{loss}}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial x}$$

Chain Rule on Directed Acyclic Graph (DAG)

Automate the computation of derivatives with computer science

The high school way

$$\frac{\partial L(x, y, z)}{\partial x} = ?$$

$$\frac{\partial L(x, y, z)}{\partial y} = ?$$

$$\begin{aligned} \frac{\partial L(x, y, z)}{\partial z} &= (xz + yz)' = (xz)' + (yz)' = z \\ \frac{\partial L(x, y, z)}{\partial y} &= (xz + yz)' = (xz)' + (yz)' = z \\ \frac{\partial L(x, y, z)}{\partial x} &= x + y \end{aligned}$$

The high school way (as we did until now):

## Now with Chain Rule

$$L(x, y, z) = (x + y)z$$

$$x = -2; y = 5; z = -4;$$

but we can re-write it with the chain rule:

$$\frac{\partial L(x, y, z)}{\partial x} = ((x + y)z)' = \frac{\partial L}{\partial q} \frac{\partial q}{\partial x}$$

## Now with Chain Rule

$$L(x, y, z) = (x + y)z$$

$$x = -2; y = 5; z = -4;$$

but we can re-write it with the chain rule:

$$\frac{\partial L(x, y, z)}{\partial x} = ((x + y)z)' = \frac{\partial L}{\partial q} \frac{\partial q}{\partial x} = \frac{\partial q}{\partial x}$$

## Now with Chain Rule

$$L(x, y, z) = (x + y)z$$

$$x = -2; y = 5; z = -4;$$

but we can re-write it with the chain rule:

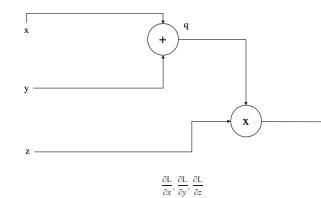
$$\frac{\partial L(x, y, z)}{\partial x} = ((x + y)z)' = \frac{\partial L}{\partial q} \frac{\partial q}{\partial x} = z \cdot 1 = z$$

## Chain Rule on Directed Acyclic Graph (DAG)

$$L(x, y, z) = (x + y)z$$

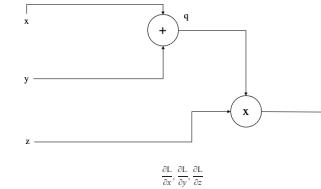
$$x = -2; y = 5; z = -4;$$

The computer science way:



## Chain Rule on Directed Acyclic Graph (DAG)

Even if the problem is very small, we break it down to subproblem so that we can **automate** it:



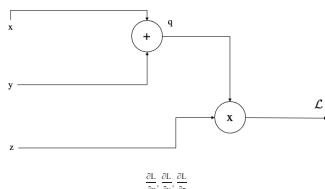
Who is input of L?

q and z are input of L.

$$\frac{\partial L}{\partial q}, \frac{\partial L}{\partial z}$$

## Chain Rule on Directed Acyclic Graph (DAG)

Even if the problem is very small, we break it down to subproblem so that we can **automate** it:

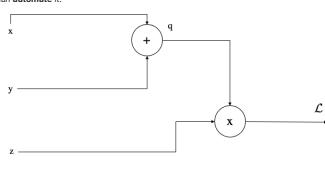


What is the derivate? (Just check the operation at the gate)

$$\frac{\partial L}{\partial q}, \frac{\partial L}{\partial z}$$

## Chain Rule on Directed Acyclic Graph (DAG)

Even if the problem is very small, we break it down to subproblem so that we can **automate** it:



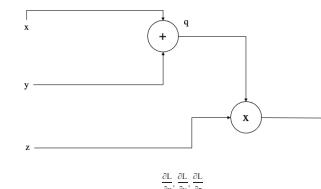
>

What is the derivate? (Just check the operation at the gate)?

$$\frac{\partial L}{\partial q} = z, \frac{\partial L}{\partial z} = q$$

## Chain Rule on Directed Acyclic Graph (DAG)

Even if the problem is very small, we break it down to subproblem so that we can **automate** it:

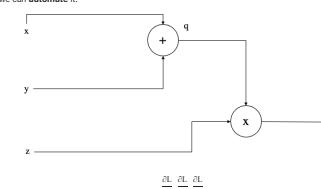


What is the derivate? (Just check the operation at the gate)?

$$\frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

## Chain Rule on Directed Acyclic Graph (DAG)

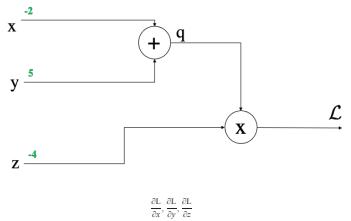
Even if the problem is very small, we break it down to subproblem so that we can **automate** it:



Ok now we have all the **analytical "local"** partial derivatives, we can compute something

$$\frac{\partial L}{\partial q} = z, \frac{\partial L}{\partial z} = q; \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

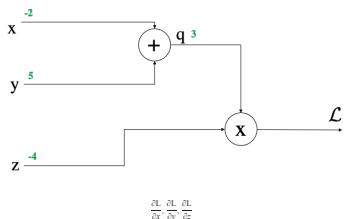
## Forward Pass



Ok now we have all the analytical partial derivatives, we can compute something

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

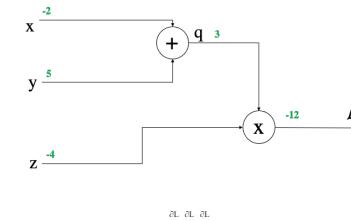
## Forward Pass



Ok now we have all the analytical partial derivatives, we can compute something

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Forward Pass



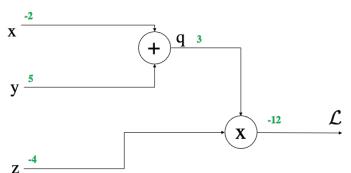
This is what we wanted:

$$\frac{\partial L}{\partial q} = \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



Act as a base case for the recursion:

$$\frac{\partial L}{\partial z} = ?$$

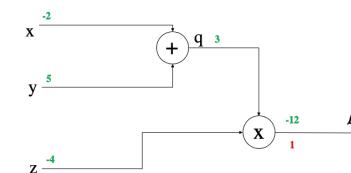
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



Act as a base case for the recursion:

$$\frac{\partial L}{\partial z} = 1$$

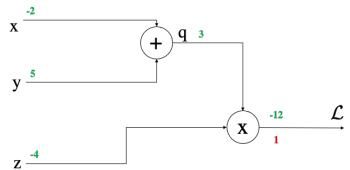
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $L$  on  $z$ ?

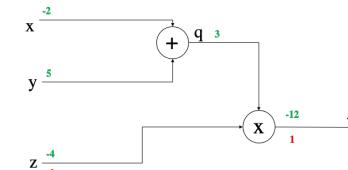
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} + 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $L$  on  $z$ ?

$$\frac{\partial L}{\partial z} = q = 3$$

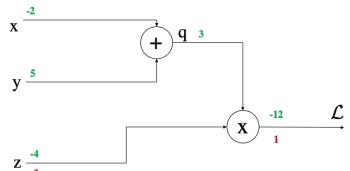
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} + 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $L$  on  $q$ ?

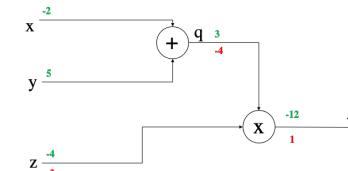
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} + 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $L$  on  $q$ ?

$$\frac{\partial L}{\partial q} = z = -4$$

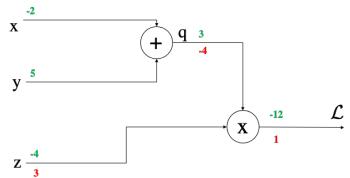
This is what we wanted:

$$\frac{\partial L}{\partial x}, \frac{\partial L}{\partial y}, \frac{\partial L}{\partial z}$$

This is what we have

$$\frac{\partial L}{\partial q} = z \cdot \frac{\partial L}{\partial z} = q \cdot \frac{\partial q}{\partial z} + 1 \cdot \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $\mathcal{L}$  on  $y$ ?

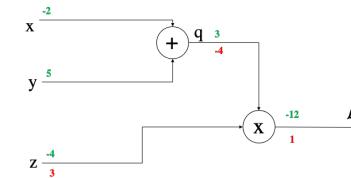
This is what we wanted:

$$\frac{\partial \mathcal{L}}{\partial x}, \frac{\partial \mathcal{L}}{\partial y}, \frac{\partial \mathcal{L}}{\partial z}$$

This is what we have

$$\frac{\partial \mathcal{L}}{\partial y} = z \cdot \frac{\partial \mathcal{L}}{\partial z} + q \cdot \frac{\partial q}{\partial z} = 1 \cdot \frac{\partial q}{\partial z} - 1, \frac{\partial q}{\partial z} = 1$$

## Backward Pass



What is the value of the gradient of  $\mathcal{L}$  on  $y$ ?

$$\frac{\partial \mathcal{L}}{\partial y} = \frac{\partial \mathcal{L}}{\partial q} \cdot \frac{\partial q}{\partial y} = z \cdot 1 = -4$$

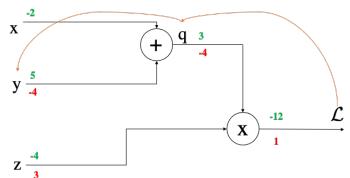
This is what we wanted:

$$\frac{\partial \mathcal{L}}{\partial x}, \frac{\partial \mathcal{L}}{\partial y}, \frac{\partial \mathcal{L}}{\partial z}$$

This is what we have

$$\frac{\partial \mathcal{L}}{\partial y} = z \cdot \frac{\partial \mathcal{L}}{\partial z} + q \cdot \frac{\partial q}{\partial z} - 1, \frac{\partial q}{\partial z} = 1$$

## Backward Pass



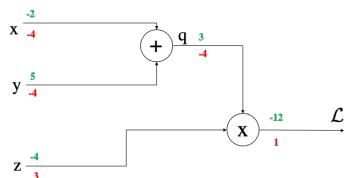
This is what we wanted:

$$\frac{\partial \mathcal{L}}{\partial x}, \frac{\partial \mathcal{L}}{\partial y}, \frac{\partial \mathcal{L}}{\partial z}$$

This is what we have

$$\frac{\partial \mathcal{L}}{\partial y} = z \cdot \frac{\partial \mathcal{L}}{\partial z} + q \cdot \frac{\partial q}{\partial z} - 1, \frac{\partial q}{\partial z} = 1$$

Check with our manual derivation



The high school way (as we did until now):

$$\begin{aligned}\frac{\partial \mathcal{L}(x,y,z)}{\partial z} &= (xz + yz)' = (xz)' + (yz)' = z = -4 \\ \frac{\partial \mathcal{L}(x,y,z)}{\partial y} &= (xz + yz)' = (yz)' + (yz)' = z = -4 \\ \frac{\partial \mathcal{L}(x,y,z)}{\partial z} &= x + y = +3\end{aligned}$$

You know what? I do not trust math, I want to verify with a machine

### Pytorch check

```
from torch import tensor
def neural_net(x,y,z):
    return (x*y)*z
x, y, z = tensor(-2., requires_grad=True), tensor(5.,requires_grad=True), tensor(-4., requires_grad=True)
loss = neural_net(x,y,z) # forward pass
loss.backward() # backward (after this I can check the gradients)
for el in [x,y,z]:
    print(el.grad)
```

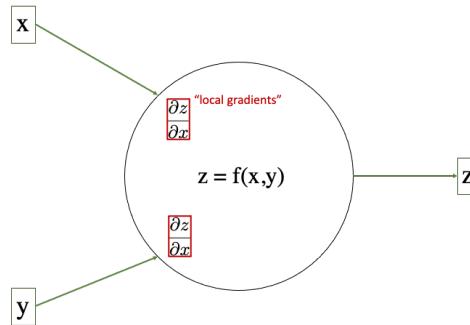
```
tensor(-4.)
tensor(-4.)
tensor(3.)
```

```
from torch import tensor
def neural_net(x,y,z):
    loss = (x*y)*z
    return loss
x, y, z = tensor(-2.), tensor(5.,requires_grad=True), tensor(-4., requires_grad=True)
loss = neural_net(x,y,z) # forward pass
loss.backward() #backward (now I can check the gradients)
for el in [x,y,z]:
    print(el.grad)
```

```
from torch import tensor
def neural_net(x,y,z):
    loss = (x*y)*z
    return loss
x, y, z = tensor(-2.), tensor(5.,requires_grad=True), tensor(-4., requires_grad=True)
loss = neural_net(x,y,z)
loss.backward()
for el in [x,y,z]:
    print(el.grad)
```

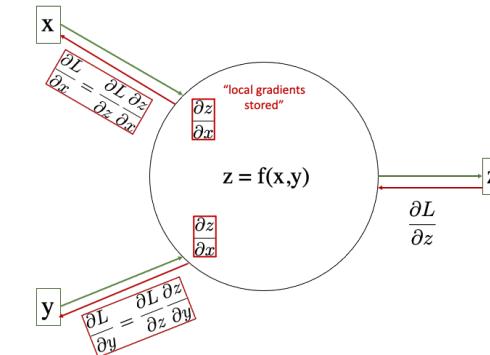
### General Recipe for Chain Rule over DAGs [Forward]

Just remember what you have to do at a generic gate:



### General Recipe for Chain Rule over DAGs [Backward]

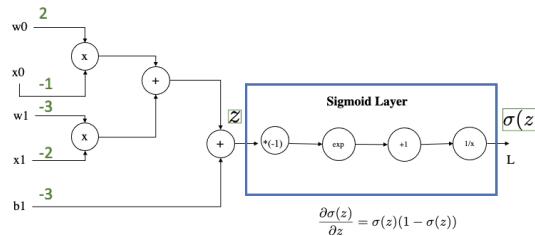
Multiply the gradient that you receive with your local gradient



### Logistic Regression Computational Graph could be simplified

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + b)}}$$

This is what implement the Sigmoid Layer in Pytorch:



Are we done with training neural nets?

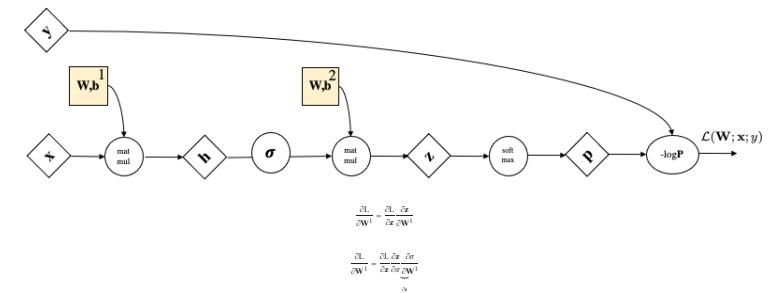
Not completely: till now scalars, but we have matrices and vectors!

### Now this looks more familiar

$\forall l \in [1, \dots, L]:$

$$1. \mathbf{W}^l \leftarrow \mathbf{W}^l - \gamma \nabla_{\mathbf{W}^l} \mathcal{L}(\mathbf{x}, \mathbf{y}; (\mathbf{W}, \mathbf{b}))$$

$$2. \mathbf{b}^l \leftarrow \mathbf{b}^l - \gamma \nabla_{\mathbf{b}^l} \mathcal{L}(\mathbf{x}, \mathbf{y}; (\mathbf{W}, \mathbf{b}))$$



$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^1} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^1}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^1} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^1} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^1}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^1} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^1} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^1} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^1}$$