

Fall 2019 CSCI 662

NON LINEAR CLASSIFIERS

Neural Networks / Deep Learning
"Introduction".

Thamme Gowda "TG"
<tg@isi.edu> ತಮ್ಮಗೌಡ

Based on chapter 3 of Jacob Eisenstein's
NLP book draft.



OVERVIEW

- Logistic Regression
 - What's deep learning is offering
 - Matrix and Vectors as transformation functions
 - Feed forward Networks
 - Activation Functions
 - Back propagation
- Computation Graph

LOGISTIC REGRESSION

→ Discriminative ; Probabilistic

$$p(y | x; \theta) = \frac{\exp(\theta \cdot f(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(\theta \cdot f(x, y'))}$$

↑ For each y ↑ same θ ↑ different featurizer

$$\theta \in \mathbb{R}^d, f(x, y) \in \mathbb{R}^d$$

y is a discrete random variable

Two key components:

- ① Dot product $\theta \cdot f(x, y)$
- ② Softmax or logistic function

VISUALISING

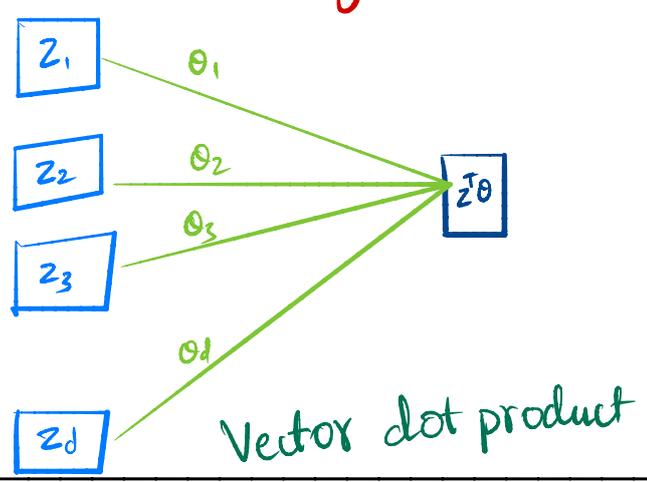
DOT PRODUCT

Model computes
say, $z = f(x, y)$

$f(x, y) \cdot \theta$
↑ features $\in \mathbb{R}^d$ ← parameters $\in \mathbb{R}^d$
a kind of transformation function
dotproduct: $\mathbb{R}^d \rightarrow \mathbb{R}$
↑ weighted sum

In general $z, \theta \in \mathbb{R}^d$
 $z \cdot \theta = \theta \cdot z = z^T \theta = \theta^T z = \sum_{i=1}^d z_i \theta_i$

Computation graph



- ① Dot product is special case of matrix multiplication
 $z^T \theta = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}_d \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}_d \quad [1 \times d][d \times 1] = [1 \times 1]$
- ② Matrix multiplication is special case of Tensor product
⇒ Lets visualize Matrix multiplication

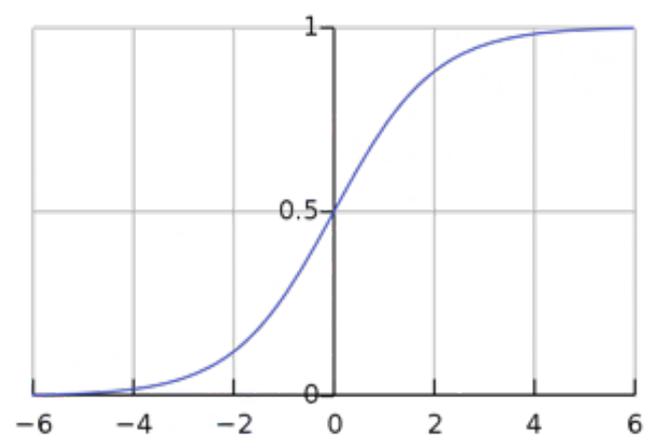
SIGMOID FUNCTION aka Std. Logistic

⇒ Maps $x \in \mathbb{R} \rightarrow$ probability scale

$$S(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$$

Good for Binary classes

For Multiclass ?



SOFTMAX FUNCTION

$$z \in \mathbb{R}^d$$

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_{j=1}^d e^{z_j}}$$

SOFTMAX

```
>>> b = np.array([25, 25, 25, 25])
>>> b / b.sum()
array([0.25, 0.25, 0.25, 0.25])
>>> be = np.exp(b)
>>> be / be.sum()
array([0.25, 0.25, 0.25, 0.25])
>>>
>>> a = np.array([25, 25, 24, 26])
>>> a / a.sum()
array([0.25, 0.25, 0.24, 0.26])
>>> ae = np.exp(a)
>>> ae / ae.sum()
array([0.2, 0.2, 0.07, 0.53])
```

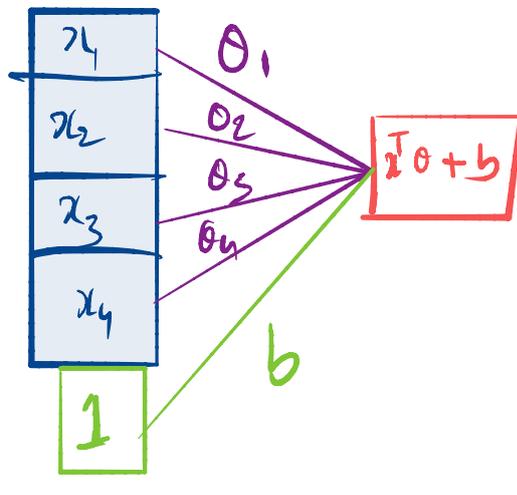
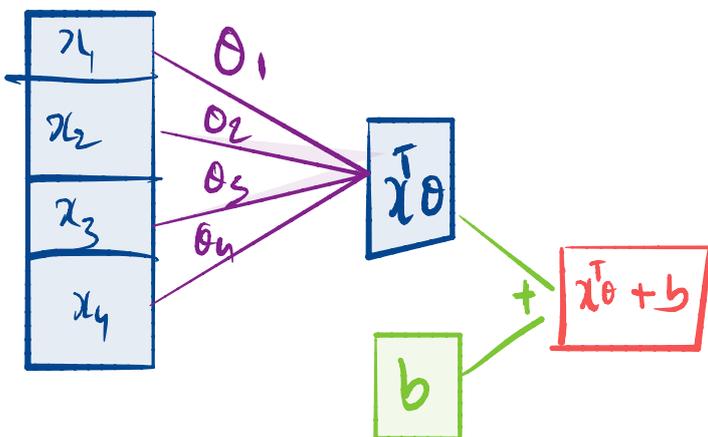
Preference
Amplification

Huge
punishment
or reward

BIAS TERM

$$y = x^T \theta + b$$

$$y = [x; 1]^T [\theta; b]$$



→ Don't be surprised if bias term is dropped from proofs

LOGISTIC REGRESSION - Parameter Estimation

Maximum Likelihood of parameter θ on a dataset $(x^{1:N}, y^{1:N}) \Leftrightarrow$ IID assumption

Maximize $p(y^{1:N} | x^{1:N}; \theta) = \prod_{i=1}^N p(y^{(i)} | x^{(i)}; \theta)$ ← Numerical instability

$\log p(y^{1:N} | x^{1:N}; \theta) = \sum_{i=1}^N \log p(y^{(i)} | x^{(i)}; \theta)$ ← Stable

→ Minimize Negative Log Likelihood = $-\sum_{i=1}^N \log p(y^{(i)} | x^{(i)}; \theta)$

Remember, $p(y^{(i)} | x^{(i)}) = \frac{e^{f(x^{(i)}, y^{(i)}) \cdot \theta}}{\sum_{y' \in Y} e^{f(x^{(i)}, y') \cdot \theta}}$

FEATURE ENGINEERING vs POWERFUL MODELING

$$\hat{y} = \operatorname{argmax}_{y \in Y} p(y | x; \theta)$$

$$x \rightarrow \begin{matrix} f(x, y=y_1) \cdot \theta \\ f(x, y=y_2) \cdot \theta \\ \vdots \\ f(x, y=y_k) \cdot \theta \end{matrix} \Rightarrow \text{Softmax} \Rightarrow \operatorname{Argmax} \Rightarrow \hat{y}$$

← aka feature engineering

→ Too much effort in building $f(x, y) \forall y \in Y$

↳ Model is simple: dot product + softmax



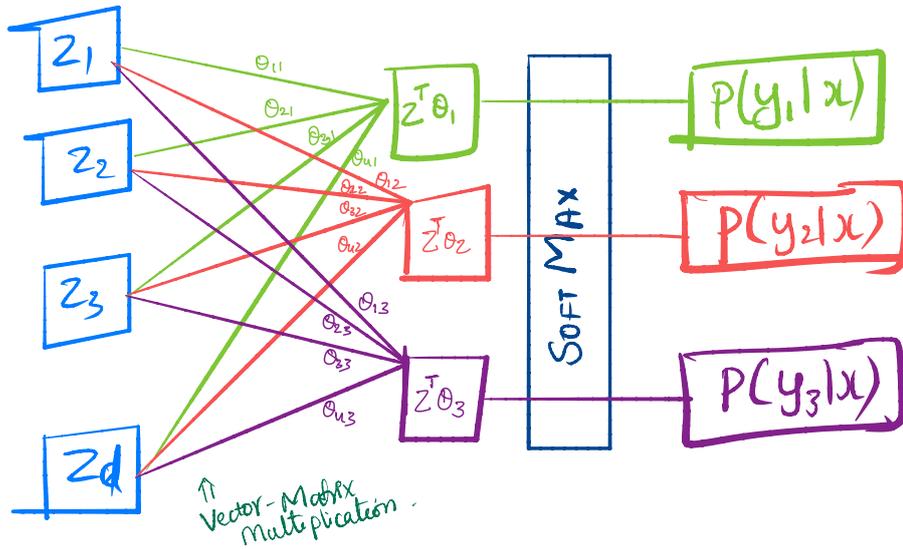
→ good $f(x, y)$ is time consuming, requires domain expertise

⇒ Would you like to spend less time in feature engineering but more time in designing models beyond dot product?

⇒ WELCOME TO DEEP LEARNING.

LOGISTIC REGRESSION (2.0)

- Previous model had same θ for all $y \in Y$, $z = f(x, y)$ but different
 - Let's have same $f(x)$ but different θ for each y .
 - Let's place all those θ s in a matrix
- For feature \mathbb{R}^d and classes k , $\theta \in \mathbb{R}^{d \times k}$

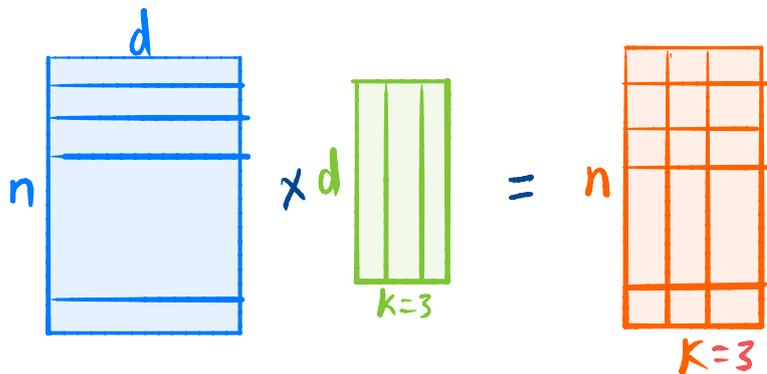


z_1	z_2	z_3	z_4
-------	-------	-------	-------

θ_{11}	θ_{12}	θ_{13}
θ_{21}	θ_{22}	θ_{23}
θ_{31}	θ_{32}	θ_{33}
θ_{41}	θ_{42}	θ_{43}

VISUALISING MATRIX MULTIPLICATION

$$x \in \mathbb{R}^{n \times d}; \quad w \in \mathbb{R}^{d \times k}; \quad xw \in \mathbb{R}^{n \times k}$$



GPUs are great for fast matrix multiplications.

- n examples in minibatch
- d dimension of model
- "hidden dimensions"
- k classes or hidden dimension

Look at rows in $x \Rightarrow n$ examples with \mathbb{R}^d representations
 Look at columns in $w \Rightarrow k$ \mathbb{R}^d vectors for dot product.

$w \in \mathbb{R}^{d \times k} : \mathbb{R}^d \rightarrow \mathbb{R}^k$: Transformation operation "or functions"

WHY DEEP LEARNING

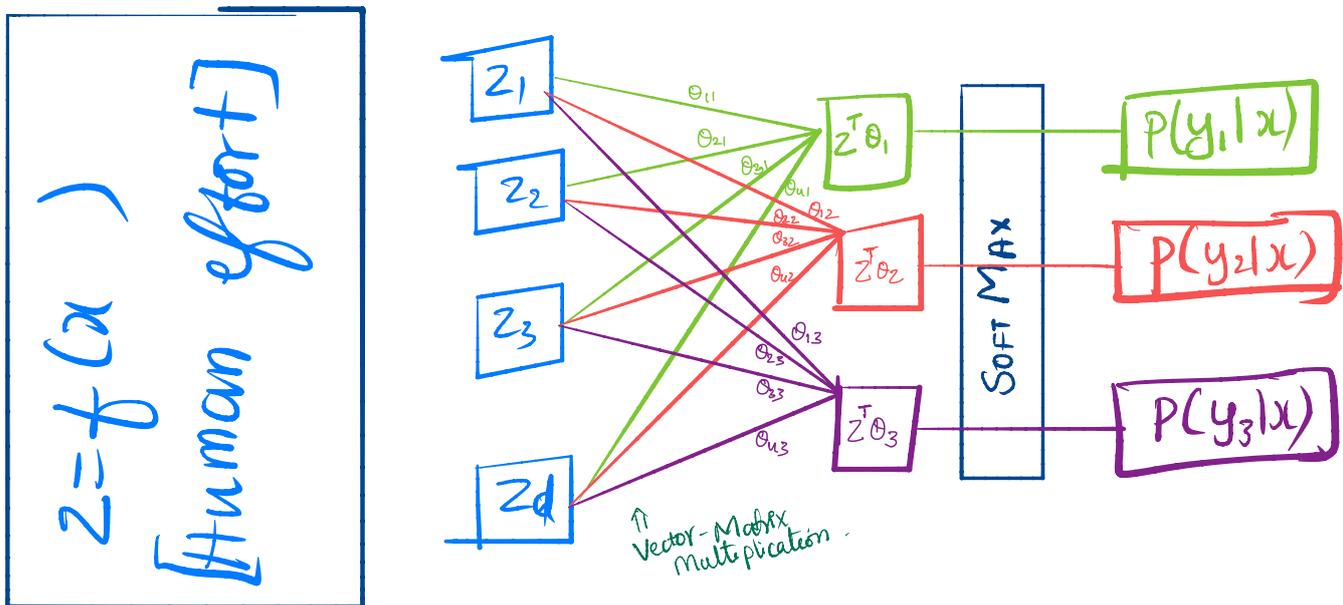
- Go beyond linear \leftrightarrow More power to models
- No need for extensive feature engineering
 - ↳ Models can learn useful features
- Optimize end-to-end in a single network.
 - ↳ NLP ~~is~~ ^{was} "pipeline" of subtasks

BUT...

- * Computationally expensive → Look at GPUs and TPUs
- * Lots of parameters → need lots of data
 - ↳ We may have it already
 - ↳ We may be able to get it.
 - ↳ That's a good research @'
 - ↳ Need efficient estimation
- * Not readily interpretable \Rightarrow Hmm! 😬

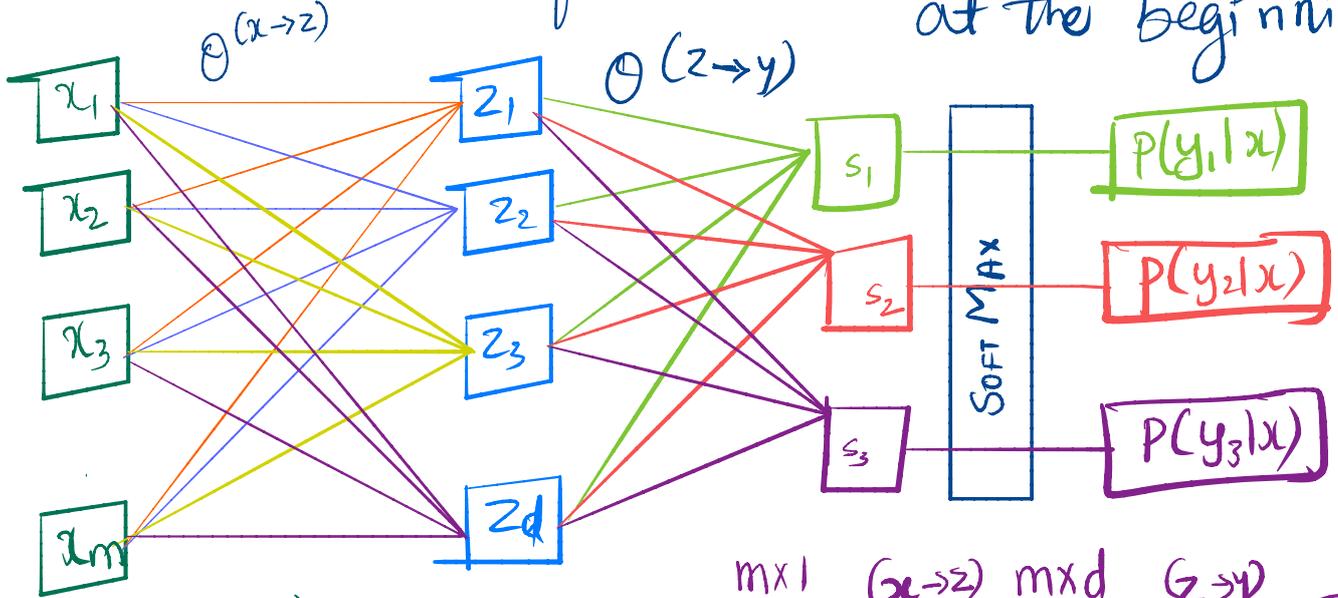
LOGISTIC REGRESSION → FEED FORWARD NN

Goal: Learn the features → Add another layer at the beginning



LOGISTIC REGRESSION → FEED FORWARD NN

Goal: Learn the features → Add another layer at the beginning



$$s = (x^T \theta^{(x \rightarrow z)}) \theta^{(z \rightarrow y)}$$

$$\hat{y} = \operatorname{argmax} \operatorname{Softmax}(s)$$

$x \in \mathbb{R}^{m \times 1}$ $\theta^{(x \rightarrow z)} \in \mathbb{R}^{m \times d}$ $\theta^{(z \rightarrow y)} \in \mathbb{R}^{d \times k}$
 What iff $W \in \mathbb{R}^{m \times k}$ such that $W = \theta^{(x \rightarrow z)} \theta^{(z \rightarrow y)}$

FEED FORWARD NETWORKS

$$s = (x^T \theta^{(x \rightarrow z)}) \theta^{(z \rightarrow y)} = x^T (\theta^{(x \rightarrow z)} \theta^{(z \rightarrow y)}) = x^T W$$

Two or more linear transformations \Leftrightarrow One Layer

Non Linear functions \rightarrow Violates the above \rightarrow Extra power

$$s = (A(x^T \theta^{(x \rightarrow z)})) \theta^{(z \rightarrow y)} \quad | \quad x \in \mathbb{R}^{m \times 1}; \quad \theta^{(x \rightarrow z)} \in \mathbb{R}^{m \times d}; \quad \theta^{(z \rightarrow y)} \in \mathbb{R}^{d \times k}$$

A is some elementwise non-linear function
 A is commonly called as **ACTIVATION FUNCTION**

Example: Sigmoid, Tanh, ReLU, GELU

ACTIVATION FUNCTIONS

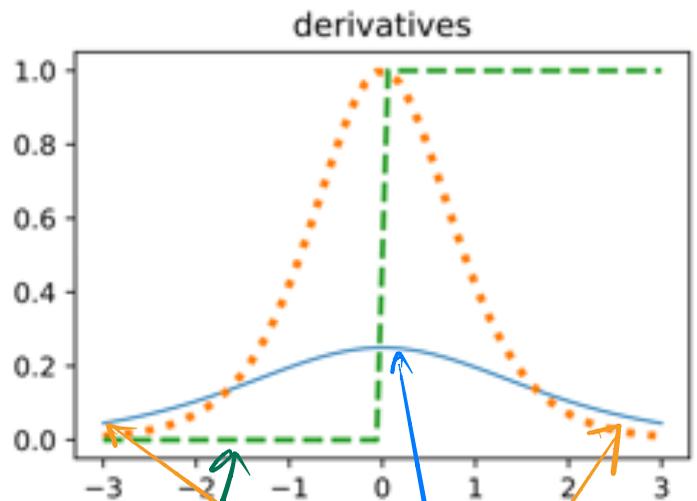
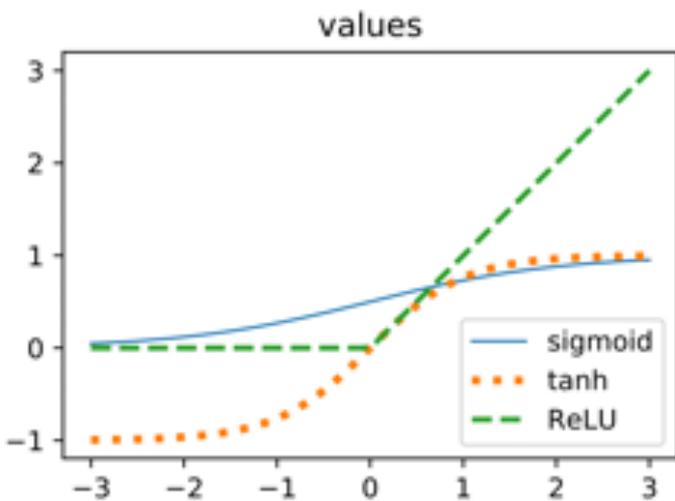
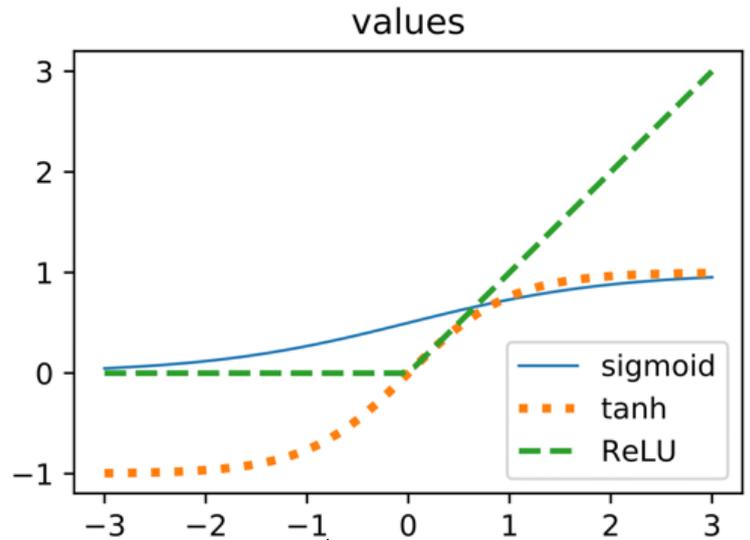
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\text{ReLU}(x) = x^+ = \max(0, x)$$

Consider these:-

- Computational needs for calculating $f(x)$ and $f'(x)$
- Vanishing Gradients problem
- Dead neurons



Problems with

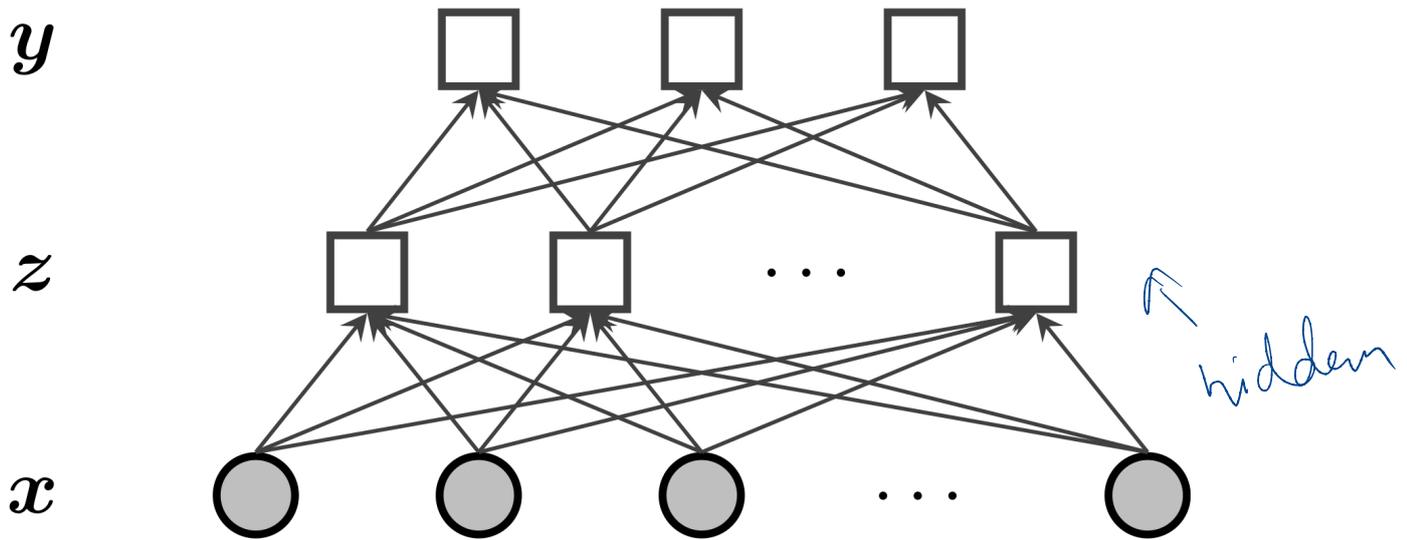
- Sigmoid
- Tanh
- ReLU

deadzone

Too small
Many multiplications
yield underflow

FEED FORWARD NETWORKS

From Textbook

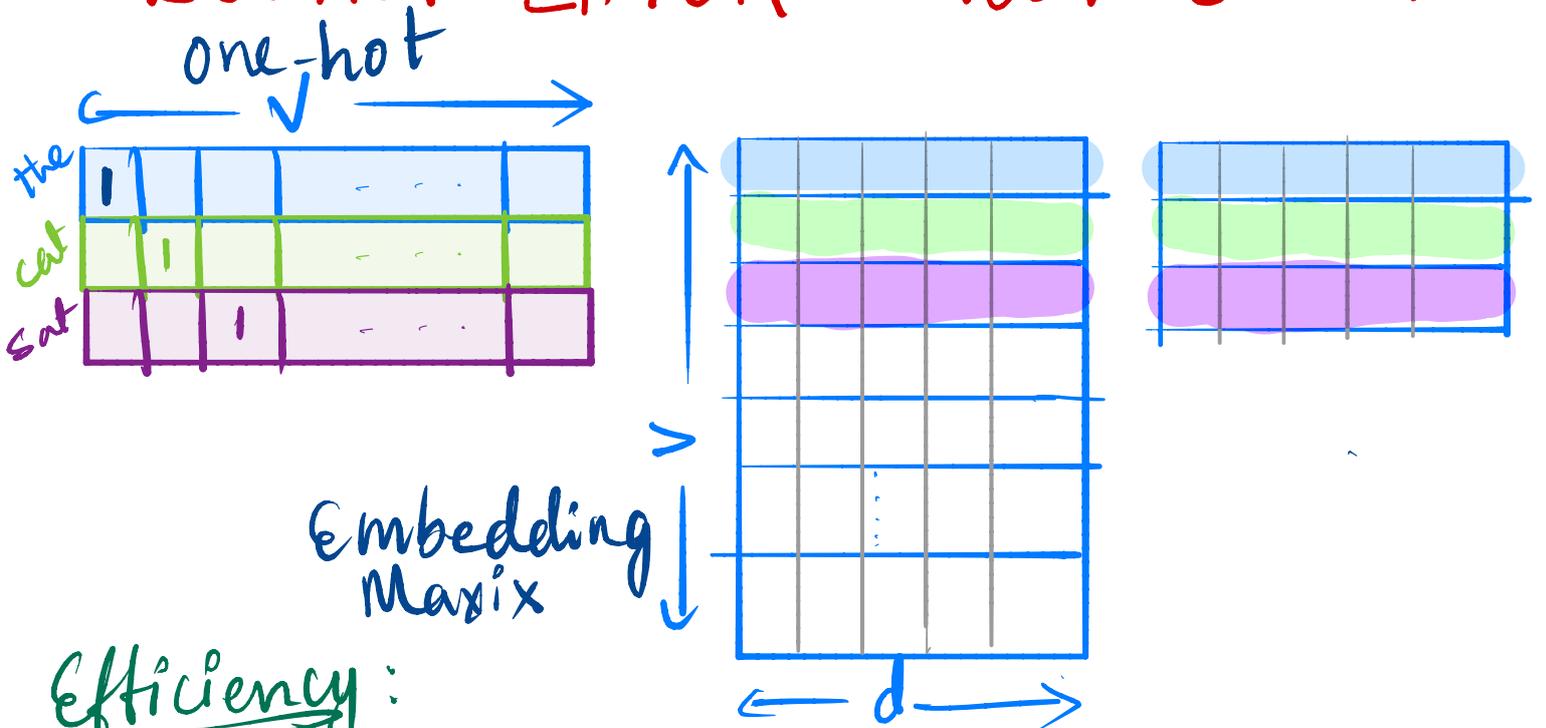


$$p(z | x; \Theta^{(x \rightarrow z)}) = \sigma(\Theta^{(x \rightarrow z)} x)$$

$$p(y | z; \Theta^{(z \rightarrow y)}, b) = \text{SoftMax}(\Theta^{(z \rightarrow y)} z + b),$$

LOOKUP LAYER

WORD EMBEDDINGS



Efficiency:

⇒ Instead of multiplying matrix of one-hot vectors, lookup vectors for words directly

PARAMETER ESTIMATION

Use an optimizer to approximate model's parameters that best fit a dataset

Loss function \uparrow
Training data \uparrow

Optimizer \rightarrow batch \rightarrow mini-batch
 \rightarrow online

LOSS FUNCTION

Conditional Log Likelihood

$$-\mathcal{L} = - \sum_{i=1}^N \log p(y^{(i)} | \mathbf{x}^{(i)}; \Theta)$$

close to 1
close to 0

alternatively as follows:

$$\tilde{y}_j \triangleq \Pr(y = j | \mathbf{x}^{(i)}; \Theta)$$
$$-\mathcal{L} = - \sum_{i=1}^N \mathbf{e}_{y^{(i)}} \cdot \log \tilde{\mathbf{y}}$$

one-hot vector

CALCULUS REVISION [DIFFERENTIAL ONLY!]

Constant	$\frac{d c}{d x} = 0$	$\frac{d}{d x} a^x = a^x \log_e a$
power	$\frac{d}{d x} x^n = n x^{n-1}$	$\frac{d}{d x} \log_a x = \frac{1}{x \log_e a}$ \Rightarrow If $a=e$, $\log_e e = 1$
Sum or Difference	$\frac{d}{d x} [f(x) \pm g(x)] = f'(x) \pm g'(x)$	
Product Rule	$\frac{d}{d x} [f(x) g(x)] = f(x) g'(x) + g(x) f'(x)$ \rightarrow Similarly Quotient Rule	
Chain Rule	$\frac{d}{d x} f(g(x)) = f'(g(x)) \cdot g'(x)$	

Back propagation is an efficient algorithm to compute derivatives of such complex functions.

PARTIAL DERIVATIVES

$f(x)$ \leftarrow function of single variable \Rightarrow Total $= \frac{d f(x)}{d x}$

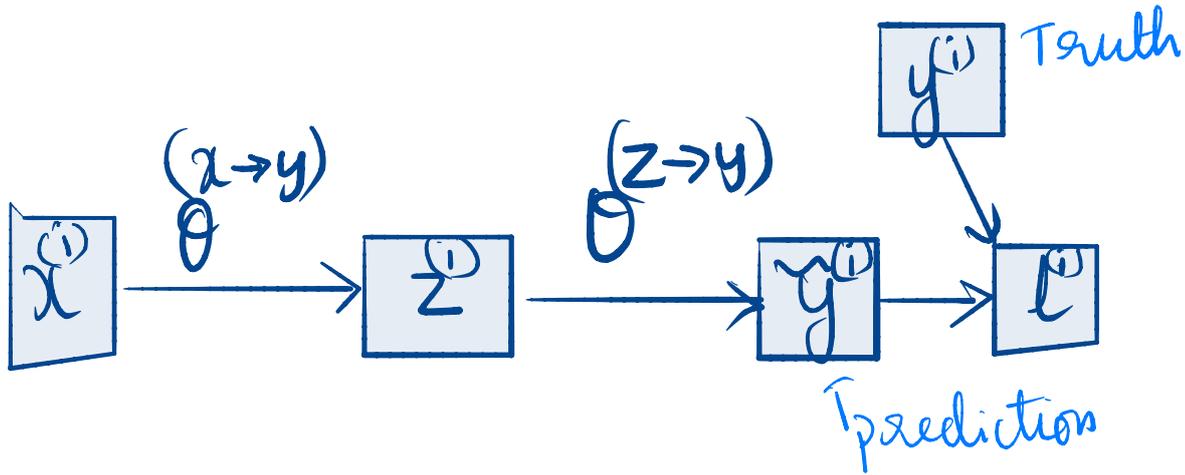
$f(x, y)$ \leftarrow function of two or more vars \Rightarrow Partial $\begin{cases} \frac{\partial f(x, y)}{\partial x} \\ \frac{\partial f(x, y)}{\partial y} \end{cases}$

\rightarrow NOTE: x and y are variables in calculus
Here we have parameters in Tensors

$\rightarrow z = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{bmatrix}$

∇_z

6 GRADIENT UPDATE RULE



$$\theta_k^{(z \rightarrow y)} \leftarrow \theta_k^{(z \rightarrow y)} - \eta^{(t)} \nabla_{\theta_k^{(z \rightarrow y)}} l^{(i)}$$

$$\theta_n^{(x \rightarrow z)} \leftarrow \theta_n^{(x \rightarrow z)} - \eta^{(t)} \nabla_{\theta_n^{(x \rightarrow z)}} l^{(i)}$$

$$z^{(i)} = f(x^{(i)}; \theta^{(x \rightarrow z)})$$

$$\hat{y}^{(i)} = g(z^{(i)}; \theta^{(z \rightarrow y)})$$

$$L = l(\hat{y}^{(i)}, y^{(i)})$$

↓ FORWARD
 ↑ BACKWARD

How to calculate $\frac{\partial L}{\partial \theta^{(x \rightarrow z)}}$ Chain Rule

Back propagation is an efficient algorithm

BACKPROPAGATION

Algorithm 6 General backpropagation algorithm. In the computation graph G , every node contains a function f_t and a set of parent nodes π_t ; the inputs to the graph are $\mathbf{x}^{(i)}$.

```
1: procedure BACKPROP( $G = \{f_t, \pi_t\}_{t=1}^T, \mathbf{x}^{(i)}$ )
2:    $v_{t(n)} \leftarrow x_n^{(i)}$  for all  $n$  and associated computation nodes  $t(n)$ .
3:   for  $t \in \text{TOPOLOGICALSORT}(G)$  do  $\triangleright$  Forward pass: compute value at each node
4:     if  $|\pi_t| > 0$  then
5:        $v_t \leftarrow f_t(v_{\pi_{t,1}}, v_{\pi_{t,2}}, \dots, v_{\pi_{t,N_t}})$ 
6:    $g_{\text{objective}} = 1$   $\triangleright$  Backward pass: compute gradients at each node
7:   for  $t \in \text{REVERSE}(\text{TOPOLOGICALSORT}(G))$  do
8:      $g_t \leftarrow \sum_{t': t \in \pi_{t'}} g_{t'} \times \nabla_{v_t} v_{t'}$   $\triangleright$  Sum over all  $t'$  that are children of  $t$ , propagating
       the gradient  $g_{t'}$ , scaled by the local gradient  $\nabla_{v_t} v_{t'}$ 
9:   return  $\{g_1, g_2, \dots, g_T\}$ 
```

Demo:

<https://google-developers.appspot.com/machine-learning/crash-course/backprop-scroll/>

Further Reading:

Yes you should understand backprop by Andrej Karpathy

<https://medium.com/@karpathy/yes-you-should-understand-backprop-e2f06eab496b>

REGULARIZATION AND DROPOUT

→ All models overfit, non-linear models are no exception

→ Regularization → avoid overfitting

$$L = \sum_{i=1}^N L^{(i)} + \lambda \|\theta\|_2^2$$

→ Dropout: Randomly set some nodes to zero.

