

# CS 6501 Natural Language Processing

## Feed-forward Neural Networks

---

Yangfeng Ji

October 3, 2018

Department of Computer Science  
University of Virginia



ENGINEERING

# Overview

1. Introduction
2. Feed-forward Neural Networks
3. Back Propagation
4. Further Comments

# Introduction

---

# Classification

Decision function

$$\Psi(x, y) = w_y^\top f(x, \theta) \quad (1)$$

- ▶  $x$ : data point
- ▶  $y$ : label
- ▶  $w_y$ : classification weights with respect to label  $y$
- ▶  $f(x, \theta)$ : feature function
- ▶  $\theta$ : parameter of feature function

# Example: Feature engineering

How to construct  $f(x, \theta)$ ?

## Example sentence

I love drinking coffee

- ▶ Unigram: I, love, drinking, coffee
- ▶ Bigram: I love, love drinking, ...
- ▶ POS tags:  $\langle I, IN \rangle, \dots$
- ▶ Production rules:  $S \rightarrow NP VP, \dots$
- ▶ ...

# Example: Feature engineering

How to construct  $f(x, \theta)$ ?

Example sentence

I love drinking coffee

Vocab	I	love	drinking	hate	coffee	tea
$x^\top$	[1	1	1	0	1	0]

$$f(x, \theta) = \mathbf{V}x$$

$$\Psi(x, y) = w_y^\top (\mathbf{V}x)$$

where  $\theta = \mathbf{V}$

# An Alternative View

Vocab	I	love	drinking	hate	coffee	tea
$x^\top$	[1	1	1	0	1	0]
$\mathbf{V}$	[ $v_I$	$v_{\text{love}}$	$v_{\text{drinking}}$	$v_{\text{hate}}$	$v_{\text{coffee}}$	$v_{\text{tea}}$ ]

# An Alternative View

Vocab	I	love	drinking	hate	coffee	tea
$x^\top$	[1	1	1	0	1	0]
$\mathbf{V}$	$[v_I$	$v_{\text{love}}$	$v_{\text{drinking}}$	$v_{\text{hate}}$	$v_{\text{coffee}}$	$v_{\text{tea}}]$

$$f(x, \theta) = v_I + v_{\text{love}} + v_{\text{drinking}} + v_{\text{coffee}} \quad (2)$$



# Linear Functions

Looking for a more powerful model then  $f(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{V}\mathbf{x}$ ?

How about

$$f(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{U}\mathbf{V}\mathbf{x}$$

# Linear Functions

Looking for a more powerful model then  $f(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{V}\mathbf{x}$ ?

How about

$$\begin{aligned}f(\mathbf{x}, \boldsymbol{\theta}) &= \mathbf{U}\mathbf{V}\mathbf{x} \\ &= (\mathbf{U}\mathbf{V})\mathbf{x}\end{aligned}$$

Not really, maybe a little. Essentially, it is still a linear function with a single matrix decomposed as  $\mathbf{U}\mathbf{V}$ .

# Nonlinearity

Add a nonlinear function  $h$

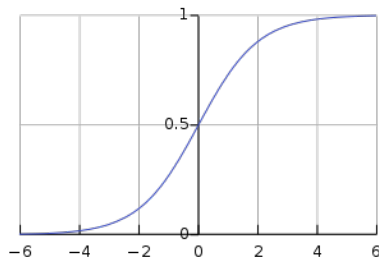
$$f(\mathbf{x}, \boldsymbol{\theta}) = h(\mathbf{V}\mathbf{x})$$

$$\Psi(\mathbf{x}, \mathbf{y}) = \mathbf{w}_y^\top h(\mathbf{V}\mathbf{x})$$

Now, it is a neural network!

Example: Sigmoid function

$$h(t) = \frac{1}{1 + e^{-t}}$$



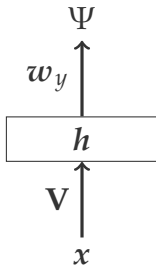
# Feed-forward Neural Networks

---

# A Simple Feed-forward Network

A fully-connected feed-forward neural network

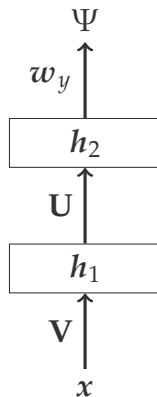
$$\Psi(x, y) = w_y^\top h(\mathbf{V}x) \quad (3)$$



# Another Feed-forward Network

$$\Psi(x, y) = w_y^\top \cdot \underbrace{h_2(\mathbf{U} \cdot h_1(\mathbf{V} \cdot x))}_{f(x, \theta)} \quad (4)$$

where  $h_1$  and  $h_2$  are nonlinear functions without parameters (hidden units).

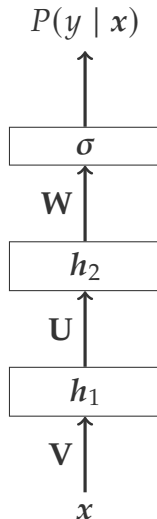


# Softmax Function

Normalize the score function to make a probability

$$\begin{aligned} P(y | x) &= \sigma(\Psi(x, y)) \\ &= \frac{\exp(\Psi(x, y))}{\sum_{y'} \exp(\Psi(x, y'))} \end{aligned} \quad (5)$$

- ▶ Main advantage is on **training**
- ▶ This is **not** a probabilistic model



# Loss Function: Cross entropy

Binary classification on a single data point  $\mathbf{x}$  with  $y \in \{0, 1\}$

$$\ell = -y \log P(y = 1 \mid \mathbf{x}) - (1 - y) \log(1 - P(y = 1 \mid \mathbf{x})) \quad (6)$$



# Loss Function: Cross entropy

Binary classification on a single data point  $\mathbf{x}$  with  $y \in \{0, 1\}$

$$\ell = -y \log P(y = 1 \mid \mathbf{x}) - (1 - y) \log(1 - P(y = 1 \mid \mathbf{x})) \quad (6)$$

► if  $y = 1$ :

$$\ell = -\log P(y = 1 \mid \mathbf{x})$$

► if  $y = 0$ :

$$\ell = -\log(1 - P(y = 1 \mid \mathbf{x})) = -\log P(y = 0 \mid \mathbf{x})$$

# Loss Function: Cross entropy

$K$ -class: convert label to  $K$ -dimensional one-hot vector with  $y_k = 1$ , if  $k$  is the label

$$\begin{aligned}\ell &= - \sum_{k=1}^K y_k \log P(y_k = 1 \mid \mathbf{x}) \\ &= -\log P(y_k = 1 \mid \mathbf{x})\end{aligned}\tag{7}$$

# Loss Function: Cross entropy

$K$ -class: convert label to  $K$ -dimensional one-hot vector with  $y_k = 1$ , if  $k$  is the label

$$\begin{aligned}\ell &= - \sum_{k=1}^K y_k \log P(y_k = 1 \mid \mathbf{x}) \\ &= -\log P(y_k = 1 \mid \mathbf{x})\end{aligned}\tag{7}$$

Essentially, it is the same as **negative log-likelihood** (NLL) in logistic regression.

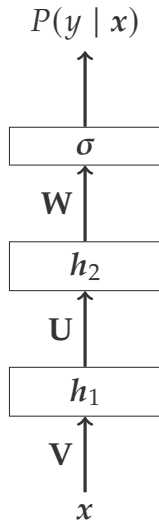
# Back Propagation

---

# Online Learning

Training NNs with **one** example at a time (redefine  $\theta$  as  $\{\mathbf{W}, \mathbf{U}, \mathbf{V}\}$ )

$$\ell(\theta) = -\log P(y_k = 1 \mid x) \quad (8)$$



# Gradient based Learning

Stochastic gradient descent

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \cdot \frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \quad (9)$$

For a subset of  $\boldsymbol{\theta}$ , e.g.,  $\boldsymbol{w}_k$

$$\boldsymbol{w}_k \leftarrow \boldsymbol{w}_k - \eta \cdot \frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{w}_k} \quad (10)$$

# Gradient based Learning (cont.)

Recall the definition of  $P(\mathbf{y} \mid \mathbf{x})$

$$\begin{aligned} \log P(\mathbf{y}_k \mid \mathbf{x}) = & \mathbf{w}_k^\top \cdot \mathbf{h}_2(\mathbf{U} \cdot \mathbf{h}_1(\mathbf{V} \cdot \mathbf{x})) \\ & - \log \sum_{k'} \exp(\mathbf{w}_{k'}^\top \cdot \mathbf{h}_2(\mathbf{U} \cdot \mathbf{h}_1(\mathbf{V} \cdot \mathbf{x}))) \end{aligned} \quad (11)$$

# Gradient based Learning (cont.)

Recall the definition of  $P(\mathbf{y} \mid \mathbf{x})$

$$\begin{aligned} \log P(\mathbf{y}_k \mid \mathbf{x}) &= \mathbf{w}_k^\top \cdot \mathbf{h}_2(\mathbf{U} \cdot \mathbf{h}_1(\mathbf{V} \cdot \mathbf{x})) \\ &\quad - \log \sum_{k'} \exp(\mathbf{w}_{k'}^\top \cdot \mathbf{h}_2(\mathbf{U} \cdot \mathbf{h}_1(\mathbf{V} \cdot \mathbf{x}))) \end{aligned} \quad (11)$$

Gradient wrt  $\mathbf{w}_k$

$$\begin{aligned} \frac{\partial \ell}{\partial \mathbf{w}_k} &= - \frac{\partial}{\partial \mathbf{w}_k} \log P(\mathbf{y} \mid \mathbf{x}) \\ &= - \mathbf{h}_2(\mathbf{U} \cdot \mathbf{h}_1(\mathbf{V} \cdot \mathbf{x})) \\ &\quad + P(\mathbf{y}_k \mid \mathbf{x}) \cdot \mathbf{h}_2(\mathbf{U} \cdot \mathbf{h}_1(\mathbf{V} \cdot \mathbf{x})) \end{aligned} \quad (12)$$



# Basic Derivatives

$$\blacktriangleright \frac{d(a \cdot z)}{dz} = a$$

$$\blacktriangleright \frac{d \log(z)}{dz} = \frac{1}{z}$$

$$\blacktriangleright \frac{de^z}{dz} = e^z$$

# Basic Derivatives

►  $\frac{d(a \cdot z)}{dz} = a$

►  $\frac{d \log(z)}{dz} = \frac{1}{z}$

►  $\frac{de^z}{dz} = e^z$

► Chain rule:  $\frac{df(g(z))}{dz} = \frac{df(g(z))}{dg(z)} \cdot \frac{dg(z)}{dz}$

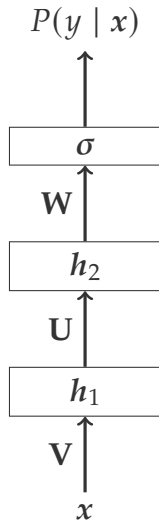
# One More Example

Given

$$\begin{aligned}\ell &= -\log P(y | x) \\ &= -\log \sigma(\mathbf{W} \cdot h_2(\mathbf{U} \cdot h_1(\mathbf{V} \cdot x)))\end{aligned}$$

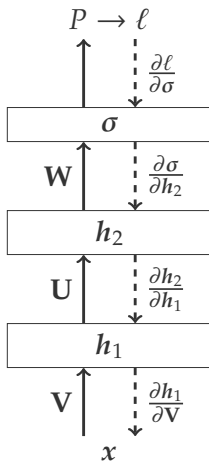
with the chain rule, we have

$$\frac{\partial \ell}{\partial \mathbf{V}} = \frac{\partial \ell}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial \mathbf{V}}$$



# Back Propagation

$$\frac{\partial \ell}{\partial \mathbf{V}} = \frac{\partial \ell}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial \mathbf{V}} \quad (13)$$

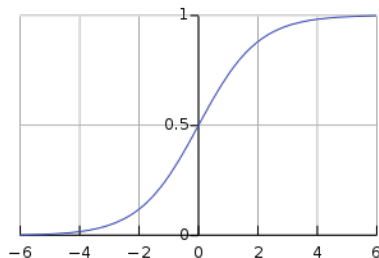


# Problems of Gradients

$$\frac{\partial \ell}{\partial \mathbf{V}} = \frac{\partial \ell}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial \mathbf{V}} \quad (14)$$

Vanishing gradients, if  
 $\|\frac{\partial \cdot}{\partial \cdot}\| \ll 1$

$$\|\frac{\partial \ell}{\partial \mathbf{V}}\| \rightarrow 0 \quad (15)$$



Solution: initialize the parameters carefully

## Problems of Gradients (cont.)

Exploding gradients, if  $\|\frac{\partial \cdot}{\partial \cdot}\| > 1$

$$\left\| \frac{\partial \ell}{\partial \mathbf{V}} \right\| > M \quad (16)$$

# Problems of Gradients (cont.)

Exploding gradients, if  $\|\frac{\partial \ell}{\partial \mathbf{z}}\| > 1$

$$\left\| \frac{\partial \ell}{\partial \mathbf{V}} \right\| > M \quad (16)$$

Solution: norm clipping [Pascanu et al., 2013]

$$\tilde{\mathbf{g}} \leftarrow \lambda \frac{\mathbf{g}}{\|\mathbf{g}\|} \quad (17)$$

where  $\mathbf{g} = \frac{\partial \ell}{\partial \mathbf{V}}$  and  $1 < \lambda \leq 5$ .

## Further Comments

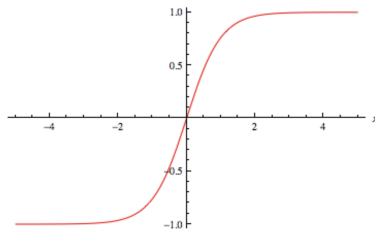
---



# Choices of Hidden Units

## Tanh function

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



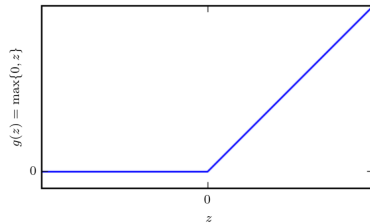
### Sigmoids

1. Symmetric sigmoids such as hyperbolic tangent often converge faster than the standard logistic function.
2. A recommended sigmoid [19] is:  $f(x) = 1.7159 \tanh(\frac{2}{3}x)$ . Since the tanh function is sometimes computationally expensive, an approximation of it by a ratio of polynomials can be used instead.
3. Sometimes it is helpful to add a small linear term, e.g.  $f(x) = \tanh(x) + ax$  so as to avoid flat spots.

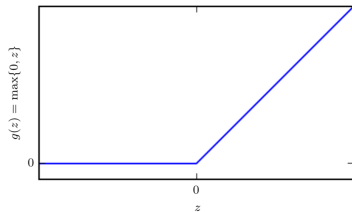
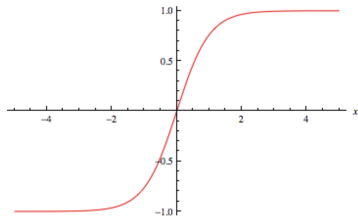
# Choices of Hidden Units (cont.)

Rectified linear unit

$$g(z) = \max\{0, z\}$$



# Comparison



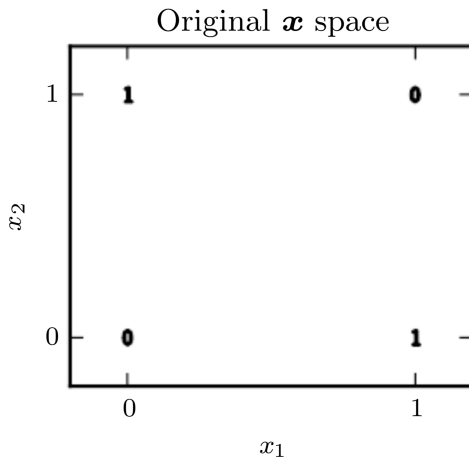
- ▶ Gradient
- ▶ Input range

Mini-batch size  $K$

$$\theta \leftarrow \theta - \eta \cdot \frac{1}{K} \sum_{k=1}^K \frac{\partial \ell_k(\theta)}{\partial \theta} \quad (18)$$

- ▶ Typically  $10 \leq K \leq 100$  [Hinton, 2012]
- ▶ Larger  $K$ 
  - ▶ gives more reliable estimate of gradient
  - ▶ takes advantage of matrix-vector multiplies on GPU
- ▶ Should work together with learning rate  $\eta$

# XOR Problem



# Summary

1. Introduction
2. Feed-forward Neural Networks
3. Back Propagation
4. Further Comments

# Reference



Hinton, G. E. (2012).  
A practical guide to training restricted boltzmann machines.  
In *Neural networks: Tricks of the trade*, pages 599–619. Springer.



LeCun, Y. A., Bottou, L., Orr, G. B., and Müller, K.-R. (2012).  
Efficient backprop.  
In *Neural networks: Tricks of the trade*, pages 9–48. Springer.



Pascanu, R., Mikolov, T., and Bengio, Y. (2013).  
On the difficulty of training recurrent neural networks.  
In *International Conference on Machine Learning*, pages 1310–1318.