

CSE 176 Introduction to Machine Learning Lecture 9: Neural Network

Some slides from Simon Prince, Dan Jurafsky, Roni Sengupta, and Olga Veksler

Recap: Principal Component Analysis (PCA)





Recap: Principal Component Analysis (PCA)



Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error. Reduce from n-dimension to k-dimension: Find kvectors $u^{(1)}, u^{(2)}, \ldots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.



Principal Component Analysis

Goal: Find *r*-dim projection that best preserves variance

- 1. Compute mean vector μ and covariance matrix Σ of original points
- 2. Compute eigenvectors and eigenvalues of Σ
- 3. Select top r eigenvectors
- 4. Project points onto subspace spanned by them:

$$y = A(x - \mu)$$

where y is the new point, x is the old one, and the rows of A are the eigenvectors

Today's topic

Shallow Neural NetworkDeep Neural NetworkLosses





Shallow Neural Network

Linear Classification

For example: fish classification - salmon or sea bass?
 extract two features, fish length and fish brightness



□yⁱ is the output (label or target) for example xⁱ



Linear Classification (perceptron)

□For two class problem and 2-dimensional data (feature vectors)



Neural Unit

• Take weighted sum of inputs, plus a bias

$$z = b + w_i x_i$$

$$i$$

$$z = w \cdot x + b$$

• Instead of just using z, we'll apply a nonlinear activation function f:

$$y = a = f(z)$$



Non-linear Activation Function





Final function the unit is computing

$$y = s(w \cdot x + b) = \frac{1}{1 + exp(-(w \cdot x + b))}$$





 x_1

*x*₂

 x_3

Input layer



+1

An example

Suppose a unit has:

$$w = [0.2, 0.3, 0.9]$$

b = 0.5

What happens with input x:

$$x = [0.5, 0.6, 0.1]$$

$$y = s(w \cdot x + b) = \frac{1}{1 + e^{-(w \cdot x + b)}} = \frac{1}{1 + e^{-(w \cdot x + b)}} = \frac{1}{1 + e^{-0.87}} = .70$$



Other non-linear activation function





Perceptron

- A very simple neural unit
- Binary output (0 or 1)
- No non-linear activation function

$$y = \begin{cases} 0, & \text{if } w \cdot x + b \leq 0 \\ 1, & \text{if } w \cdot x + b > 0 \end{cases}$$



Perceptron from the 50's and 60's





https://www.youtube.com/watch?v=cNxadbrN_al&t=71s

The XOR problem

Minsky and Papert (1969)

□Can perceptron compute simple functions of input?

AND			OR			XOR		
xl	x2	У	x1	x2	У	x1	x2	У
0	0	0	0	0	0	0	0	0
0	1	0	0	1	1	0	1	1
1	0	0	1	0	1	1	0	1
1	1	1	1	1	1	1	1	0



Easy to build AND or OR with perceptron

$$y = \begin{cases} 0, & \text{if } w \cdot x + b \le 0\\ 1, & \text{if } w \cdot x + b > 0 \end{cases}$$



Is it possible to capture XOR with perceptrons?

□Pause the lecture and try for yourself!

□No!

Why? Perceptrons are linear classifiers



Decision boundaries



XOR is not a linearly separable function!



Solution to the XOR problem

XOR can't be calculated by a single perceptron
XOR can be calculated by a layered network of units.







a) The original *x* space

b) The new (linearly separable) h space

(With learning: hidden layers will learn to form useful representations)



Shallow Neural Network with Hidden Units

$$y = \phi_0 + \phi_1 \mathbf{a}[\theta_{10} + \theta_{11}x] + \phi_2 \mathbf{a}[\theta_{20} + \theta_{21}x] + \phi_3 \mathbf{a}[\theta_{30} + \theta_{31}x].$$

Break down into two parts:

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

where:

Hidden units
$$\begin{cases} h_1 = \mathbf{a}[\theta_{10} + \theta_{11}x] \\ h_2 = \mathbf{a}[\theta_{20} + \theta_{21}x] \\ h_3 = \mathbf{a}[\theta_{30} + \theta_{31}x] \end{cases}$$



1. Compute the linear function











$$y = \phi_0 + \phi_1 \mathbf{a}[\theta_{10} + \theta_{11}x] + \phi_2 \mathbf{a}[\theta_{20} + \theta_{21}x] + \phi_3 \mathbf{a}[\theta_{30} + \theta_{31}x].$$



Example shallow network = piecewise linear functions 1 "joint" per ReLU function

Depicting shallow neural networks

$$h_{1} = a[\theta_{10} + \theta_{11}x]$$

$$h_{2} = a[\theta_{20} + \theta_{21}x]$$

$$y = \phi_{0} + \phi_{1}h_{1} + \phi_{2}h_{2} + \phi_{3}h_{3}$$

$$h_{3} = a[\theta_{30} + \theta_{31}x]$$



With enough hidden units

Q... we can describe any 1D function to arbitrary accuracy





Universal approximation theorem

"a formal proof that, with enough hidden units, a shallow neural network can describe any continuous function on a compact subset of \mathbb{R}^D to arbitrary precision"



Universal approximation theorem

Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

Abstract. In this paper we demonstrate that finite linear combinations of compositions of a fixed, univariate function and a set of affine functionals can uniformly approximate any continuous function of n real variables with support in the unit hypercube; only mild conditions are imposed on the univariate function. Our results settle an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary decision regions can be arbitrarily well approximated by continuous feedforward neural networks with only a single internal, hidden layer and any continuous sigmoidal nonlinearity. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.

Key words. Neural networks, Approximation, Completeness.

Cybenko, George. "Approximation by superpositions of a sigmoidal function." *Mathematics of control, signals and systems* 2.4 (1989): 303-314.



Terminology



- Y-offsets = biases
- Slopes = weights
- Everything in one layer connected to everything in the next = fully connected network
- No loops = feedforward network
- Values after ReLU (activation functions) = activations
- Values before ReLU = pre-activations
- One hidden layer = shallow neural network
- More than one hidden layer = deep neural network
- Number of hidden units ≈ capacity





Deep Neural Network

Shallow network

1 input, 4 hidden units, 2 outputs

$$h_1 = \mathbf{a}[\theta_{10} + \theta_{11}x]$$
$$h_2 = \mathbf{a}[\theta_{20} + \theta_{21}x]$$
$$h_3 = \mathbf{a}[\theta_{30} + \theta_{31}x]$$
$$h_4 = \mathbf{a}[\theta_{40} + \theta_{41}x]$$





Network as composing function

$$\begin{aligned} h_1 &= \mathbf{a}[\theta_{10} + \theta_{11}x] \\ h_2 &= \mathbf{a}[\theta_{20} + \theta_{21}x] \\ h_3 &= \mathbf{a}[\theta_{30} + \theta_{31}x] \\ h_4 &= \mathbf{a}[\theta_{40} + \theta_{41}x] \end{aligned} \qquad \begin{aligned} h_1' &= \mathbf{a}[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3] \\ h_2' &= \mathbf{a}[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3] \\ h_3' &= \mathbf{a}[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3] \end{aligned}$$



Example of Multi Layer Perceptron (MLP)





Shallow vs deep networks

The best results are created by deep networks with many layers.

- □50-1000 layers for most applications
- Best results in
 - Computer vision
 - □ Natural language processing
 - Graph neural networks
 - Generative models
 - Reinforcement learning

All use deep networks. But why?

- Ability to approximate different functions?
- □Both obey the universal approximation theorem.
- □Argument: One layer is enough, and for deep networks could arrange for the other layers to compute the identity function.



Shallow vs deep networks

□Number of linear regions per parameter

Deep networks create many more regions per parameters



Shallow vs deep networks

□ Fitting and generalization

Figure 20.2 MNIST-1D training. Four fully connected networks were fit to 4000 MNIST-1D examples with random labels using full batch gradient descent, He initialization, no momentum or regularization, and learning rate 0.0025. Models with 1,2,3,4 layers had 298, 100, 75, and 63 hidden units per layer and 15208, 15210, 15235, and 15139 parameters, respectively. All models train successfully, but deeper models require fewer epochs.







Training Perceptron - First Attempt

 $[\mathbf{x}_2,...,\mathbf{x}_m]$

$$\Rightarrow L(W) = \sum_{i \in \text{train}} [\mathbf{y}^i \neq u(W^T X^i)]$$

classification error counts
since both $\mathbf{y}^i, u \in \{0,1\}$

extreme case of (so-called) vanishing gradients

Zero Gradients Problem

Classification error loss function *L*(W) is **piecewise constant**:



NOTE: in this case gradient ∇L is always either zero or does not exist "error count" loss function cannot be optimized via gradient descent

Work-around for Zero Gradients

Perceptron: $f(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$ approximate decision function *u* using its softer version (relaxation)



u(*t*) - *unit step* function (a.k.a. *Heaviside* function)



Work-around for Zero Gradients

Perceptron: $f(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$ approximate decision function *u* using its softer version (relaxation)



Relaxed predictions are often interpreted as prediction "probabilities"

$$\Pr(\mathbf{x}^{i} \in \text{Class1} | W) = \sigma(W^{T} X^{i})$$

$$\Pr(\mathbf{x}^{i} \in \text{Class0} | W) = 1 - \sigma(W^{T} X^{i}) \equiv \sigma(-W^{T} X^{i})$$

Training Perceptron - Second Attempt

Perceptron approximation: $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$



Classification error loss: now makes no sense at all

relaxed decision function (sigmoid)

never returns exactly 0 or 1



$$0 < \sigma(t) \equiv \frac{1}{1 + \exp(-t)} < 1$$

 $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$ Perceptron approximation: $L(\mathbf{y}, \sigma) = (\mathbf{y} - \sigma)^2$ Consider quadratic loss: NOTE: **6**(*t*) Loss $L(\mathbf{y}, \sigma(W^T X))$ is now differentiable with respect to Wbecause $L(\mathbf{y}, \sigma)$ is differentiable w.r.t. σ

Perceptron approximation: $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$ $L(\mathbf{y}, \sigma) = (\mathbf{y} - \sigma)^2$ Consider quadratic loss: **6**(*t*) 1 $(\mathbf{y}^i - \sigma(W^T X^i))$ $W^T_i X^i > 0$

misclassified example

Perceptron approximation: $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$ $L(\mathbf{y}, \sigma) = (\mathbf{y} - \sigma)^2$ Consider quadratic loss: **6**(*t*) $(\mathbf{y}^j - \sigma(W^T X^j))$ $W^T_X X^j < 0$ $\mathbf{v}^j = 1$ another misclassified example

 $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$ Perceptron approximation: $L(\mathbf{y}, \sigma) = (\mathbf{y} - \sigma)^2$ Consider quadratic loss: $\mathbf{b}(t)$ $(\mathbf{y}^j - \sigma(W^T X^j))$ $(\mathbf{y}^i - \sigma(W^T X^i))$ NOTE: loss function encourages W s.t. correctly classified points are moved further from the decision boundary, i.e. $W^T X^i \gg 0$ and $W^T X^j \ll 0$. correctly classified examples

Perceptron approximation: $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$

Consider quadratic loss:
$$L(\mathbf{y}, \sigma) = (\mathbf{y} - \sigma)^2$$



Cross-Entropy Loss (related to *logistic regression* loss)

Perceptron approximation: $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$

Consider two probability distributions over two classes (e.g. bass or salmon): $(\mathbf{y}, 1 - \mathbf{y})$ and $(\sigma, 1 - \sigma)$ bass salmon $\Pr(\mathbf{x}^{i} \in \text{Class1} | W) = \sigma(W^{T}X^{i})$ $\Pr(\mathbf{x}^{i} \in \text{Class0} | W) = 1 - \sigma(W^{T}X^{i})$

Distance between two distributions can be evaluated via **cross-entropy** (equivalent to *KL divergence* for fixed target)

$$H(\boldsymbol{p},\boldsymbol{q}) := -\sum_{k} p_k \, \ln q_k$$

Cross-Entropy Loss (related to *logistic regression* loss)

Perceptron approximation: $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$

Consider two probability distributions over two classes (e.g. bass or salmon): $(\mathbf{y}, 1 - \mathbf{y})$ and $(\sigma, 1 - \sigma)$



(binary) Cross-entropy loss:

$$L(\mathbf{y},\sigma) = -\mathbf{y}\ln\sigma - (1-\mathbf{y})\ln(1-\sigma)$$

Distance between two distributions can be evaluated via **cross-entropy** (equivalent to *KL divergence* for fixed target)

$$H(\boldsymbol{p},\boldsymbol{q}) := -\sum_{k} p_k \, \ln q_k$$

Cross-Entropy Loss (related to *logistic regression* loss)

Perceptron approximation: $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$

Consider two probability distributions over two classes (e.g. bass or salmon): $(\mathbf{y}, 1 - \mathbf{y})$ and $(\sigma, 1 - \sigma)$



(binary) **Cross-entropy loss:**

$$L(\mathbf{y},\sigma) = -\mathbf{y}\ln\sigma - (1-\mathbf{y})\ln(1-\sigma)$$

Each data label y provides "deterministic" distribution (y, 1 - y) that is either (1,0) or (0,1). This implies an equivalent alternative expression:

$$L(\mathbf{y},\sigma) = \begin{cases} -\ln\sigma & \text{if } \mathbf{y} = 1\\ -\ln(1-\sigma) & \text{if } \mathbf{y} = 0 \end{cases}$$

Cross-Entropy Loss (related to *logistic regression* loss)

Perceptron approximation: $\mathbf{f}(\mathbf{w}, \mathbf{x}^i) = u(W^T X^i) \approx \sigma(W^T X^i)$

Consider two probability distributions over two classes (e.g. bass or salmon): $(\mathbf{y}, 1 - \mathbf{y})$ and $(\sigma, 1 - \sigma)$



Total loss:
$$\sum_{i \in \text{train}} \left(-\mathbf{y}^i \ln \sigma(W^T X^i) - (1 - \mathbf{y}^i) \ln(1 - \sigma(W^T X^i)) \right)$$

$$\Rightarrow L(W) = -\sum_{\substack{i \in \text{train} \\ \mathbf{y}^i = 1}} \ln \sigma(W^T X^i) - \sum_{\substack{i \in \text{train} \\ \mathbf{y}^i = 0}} \ln(1 - \sigma(W^T X^i))$$

sum of Negative Log-Likelihoods (NLL)

Towards Multi-label Classification



Towards Multi-label Classification



Towards Multi-label Classification



Such "probability scores" $G_1, G_2, ..., G_K$ over K classes do not add up to 1

Common Approach: Soft-Max



Notation: K rows of matrix W are vectors W_k so that vector WX has elements $W_k^T X$

Soft-Max Function $\bar{\sigma} : \mathbb{R}^K \to \Delta_K$



Soft-max normalizes logits vector **a** converting it to distribution over classes

Soft-Max Function $\bar{\sigma} : \mathbb{R}^K \to \Delta_K$



Soft-max normalizes logits vector **a** converting it to distribution over classes

Soft-Max Function $\bar{\sigma} : \mathbb{R}^K \to \Delta_K$

NN Example: $\frac{\exp W_1^T X}{\sum_k \exp W_k^T X}$ NOTE: $W_1^T X$ $\exp W_2^T X$ soft-max generalizes sigmoid $\sum_{k} \exp W_{k}^{T} X$ to multi-class predictions. Indeed, $W_2^T X$ consider binary perceptron with scalar softmax linear discriminator $W^{T}X$ (e.g. for class 1) $\frac{\exp W_K^T X}{\sum_k \exp W_k^T X}$ $W_K^T X$ $\sigma(W^T X) = \frac{1}{1 + e^{-W^T X}}$ sigmoid $\equiv \frac{e^{\frac{1}{2}W^{T}X}}{e^{\frac{1}{2}W^{T}X} + e^{-\frac{1}{2}W^{T}X}} = \bar{\sigma}_{1} \begin{pmatrix} \frac{1}{2}W^{T}X \\ -\frac{1}{2}W^{T}X \end{pmatrix}$ class 1 output of soft-max for a combination of two linear predictors: $\frac{1}{2}W^{T}X$ for class 1 and - $\frac{1}{2}W^{T}X$ for class $\neg 1$ (class 0) $(\bar{\sigma}_1, \bar{\sigma}_2, \dots, \bar{\sigma}_K)$

Soft-max normalizes logits vector **a** converting it to distribution over classes

bass salmon sturgeon

(general multi-class case)

Cross-Entropy Loss

K-label perceptron's output: $\bar{\sigma}(\mathbf{W}X^i)$ for example X^i *k*-th index Multi-valued label $\mathbf{y}^i = k$ gives **one-hot** distribution $\bar{\mathbf{y}}^i = (0, 0, 1, 0, \dots, 0)$ Consider two probability distributions over K classes (e.g. bass, salmon, sturgeon): $\bar{\mathbf{y}}^i$ and $(\bar{\sigma}_1, \bar{\sigma}_2, \bar{\sigma}_3, ..., \bar{\sigma}_K)$ $\Pr(\mathbf{x}^i \in \operatorname{Class} k \,|\, W) = \bar{\sigma}_k(WX^i)$ bass salmon sturgeon cross entropy **Total loss:** $L(W) = \sum \left[\sum -\bar{\mathbf{y}}_k^i \ln \bar{\sigma}_k(WX^i) \right]$ $i \in \text{train} \quad k$ $L(W) = - \sum \ln \bar{\sigma}_{\mathbf{y}^i}(WX^i)$ $i \in \text{train}$

sum of Negative Log-Likelihoods (NLL)

soft-max vs arg-max Multi-label (linear) Classification

Define *K* linear transforms, from features *X* to K "*logits*" $logit_k(X) = W_k^T X$ for k = 1, 2, ... K

- arg-max assigns X to class k corresponding to the largest logit $\arg \max_k \{W_k^T X\}$
- Let **R**_k be decision region for class k all points X assigned to class k by *arg-max*

soft-max $\bar{\sigma}\{W_k^{\top}X\}$ softens hard **arg-max** predictions similarly to how sigmoid softens unit-step function



Summary

□Shallow neural network

□Universal function approximation theorem

Deep neural network

□ Multi layer perceptron

Loss

□Sigmoid, Softmax

Cross entropy loss, quadratic loss

