CS 505: Introduction to Natural Language Processing Wayne Snyder Boston University

Lecture 10: Deep Learning: Artificial Neural Networks: Design and Implementation; Workflow for Classification Tasks; Evaluating Classifiers



Deep Learning refers to Supervised Learning using an Artificial Neural Network, which has the following features:

- It is a network/graph of small computation units called artificial neurons, loosely modeled on the neurons in our brains, which send signals to each other. The signals are floating-point numbers.
- The network is typically organized in layers: the first layer is the input layer, the last is the output layer, and others are called hidden layers.
- The input layer is array/vector of floats, and the output layer produces an array of floats. Thus, the network computes a function from vectors to vectors.







Each neuron:

- Is connected to each neuron in the previous layer, and each connection has a weight or parameter which determines the strength of the signal (importance of this input to the neuron);
- Performs logistic regression, using the sigmoid or other non-linear activation function.

Gradient descent is used to learn the weights to minimize some cost function on the outputs.





Features of artificial neural networks:

- Additional layers may perform data aggregation (e.g., convolution and pooling), dropout, or other kinds of data manipulation (e.g., softmax = transforming the output into a probability distribution).
- In a feedforward network, the network transforms an array of floats through the layers into another array of floats; in a sequence model, the inputs and outputs are sequences of vectors; and recurrent layers have cyclical connections which act as memory.
- BERT, GPT, and other large networks learn to pay attention to complex patterns in the input sequence (e.g., words in a sentence).







A neuron is a higher-dimensional version of our logistic regression algorithm:



Inputs



p Learning

is the sigmoid from logistic



But non-linear activation functions besides sigmoid are more often used!



Hyperbolic Tangent (Tanh)



Rectified Linear Unit (Relu)

A small amount of Linear Algebra can be used to compactly specify the logistic regression performed by a neuron.

Inputs and weights/parameters are just vectors:



often written as tuples:

$$\mathbf{x} = (x_1, x_2, \dots, x_m).$$
 $\mathbf{w} = (w_1, w_2, \dots, w_m).$

Then the neuron performs logistic regression by performing a dot product of inputs \mathbf{X} and weights \mathbf{W} ,

$$\mathbf{w} \cdot \mathbf{x} = (w_1 \cdot x_1) + \dots + (w_k \cdot x_k) = \mathbf{w}^{\mathrm{T}} \mathbf{x}$$

and then applying the activation function f:

$$\sigma(m * x + b) \qquad a = f(\mathbf{w}^T \mathbf{x}).$$
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \qquad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix}$$

A small but important detail: Each neuron has a bias term, because they are simply doing logistic regression!

This connection has a weight (the value b) but is not connected to the inputs; it serves to scale the output, just as b does in linear regression.

A convenient way to encode this is to assume that all input vectors to all neurons contain a constant 1.0 value: $\sigma(m * x + b * 1.0)$









The simplest network is simply a row of such neurons, where

- o Each has its own weight/parameter vector, but
- Shares the same input vector

Activations (outputs)

 a_1 a_2 a_3 a_{Δ} Neuron 1 Neuron 3 Neuron 4 Neuron 2 $\mathbf{w}_2^T \rightarrow \mathbf{x}$ $\mathbf{w}_{3}^{T} \overset{\downarrow}{\rightarrow} \times$ $\mathbf{w}_4^T \rightarrow \mathbf{x}$ $\mathbf{w}_1^T - \mathbf{x}$ $w_{1.1}$ $w_{4.5}$ $w_{1,2}$ b_4 1.0 \mathcal{X}_{1} χ_4 χ_5 χ_2 X_3 Inputs

The inputs and outputs for a layer are vectors:

$$\mathbf{a} = (a_1, a_2, a_3, a_4)$$

The weights/parameters for a layer form a matrix:

$$W \mathbf{x} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & w_{1,4} & w_{1,5} & b_1 \\ w_{2,1} & w_{2,2} & w_{2,3} & w_{2,4} & w_{2,5} & b_2 \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} & w_{3,5} & b_3 \\ w_{4,1} & w_{4,2} & w_{4,3} & w_{4,4} & w_{4,5} & b_4 \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_2 \\ x_4 \\ x_5 \\ 1.0 \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x} \\ \mathbf{w}_2^T \mathbf{x} \\ \vdots \\ \mathbf{w}_n^T \mathbf{x} \end{bmatrix}.$$

 $\mathbf{x} = (x_1, x_2, x_3, x_4, x_5, 1.0)$



Again, a little bit of linear algebra can make this a whole lot simpler:

$$W \mathbf{x} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & w_{1,4} & w_{1,5} & b_1 \\ w_{2,1} & w_{2,2} & w_{2,3} & w_{2,4} & w_{2,5} & b_2 \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} & w_{3,5} & b_3 \\ w_{4,1} & w_{4,2} & w_{4,3} & w_{4,4} & w_{4,5} & b_4 \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_2 \\ x_4 \\ x_5 \\ 1.0 \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x} \\ \mathbf{w}_2^T \mathbf{x} \\ \vdots \\ \mathbf{w}_n^T \mathbf{x} \end{bmatrix}$$

$$\mathbf{a} = (a_1, a_2, a_3, a_4)$$
$$\mathbf{f} \begin{pmatrix} \mathbf{w}_1^T \mathbf{x} \\ \mathbf{w}_2^T \mathbf{x} \\ \mathbf{w}_3^T \mathbf{x} \\ \mathbf{w}_4^T \mathbf{x} \end{pmatrix} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \mathbf{a}_3 \\ \mathbf{a}_4 \end{bmatrix}$$

Inputs

Our FFNN:





 $a^2 = f(W_2(f(W_2\mathbf{x}) + [1.0]))$

Classification Methods: Supervised ML

Input:

- a fixed set of classes $C = \{c_1, c_2, ..., c_J\}$
- a randomly-permuted set of labeled documents
 (d₁, c₁),...,(d_n, c_n) split into
 - a training set (d₁, c₁),...,(d_m, c)
 - a testing set d_{m+1}, \dots, d_n (labels withheld)
- Output:
 - A classifier $\gamma: d \rightarrow c$ trained the training set
 - The testing set with labels calculated by y
 - Test results (confusion matrix, metrics, etc.)

Classification: Binary, Multiclass, and Multilabel

In Binary Classification, we have 2 labels, and we must choose one; often this is phrased as "something" or "not something" (spam or ham, misinformation or not, etc.)

In Multiclass Classification, we have more than 2 labels, and our task is to assign a single label to each sample:

In Multilabel Classification, we have more than 2 labels, and our task is to assign any appropriate labels (not just one):





Multilabel Classification



Classification Methods: Supervised ML

- There are many different kinds of classifiers for labeled data
 - Naïve Bayes
 - Logistic regression
 - Neural networks

Classification with Deep Learning

All of these types of classification are typically implemented by having the network output a probability distribution on all the labels.

To convert the output values into a distribution, we used a generalization of the sigmoid called the softmax:



Softmax = a generalization of sigmoid which scales k numbers into a probability distribution.

• For a vector *z* of dimensionality *k*, the softmax is:

softmax(z) =
$$\left[\frac{\exp(z_1)}{\sum_{i=1}^{k} \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^{k} \exp(z_i)}, \dots, \frac{\exp(z_k)}{\sum_{i=1}^{k} \exp(z_i)}\right]$$

Examples:

5 A = [1.2, 0.5, 3.2, 2.2, 1.9]6 softmax(A)

array([0.07343397, 0.03646623, 0.54260772, 0.19961422, 0.14787785])

 $s(z) = \frac{e^z}{e^z + 1}$

```
: 1 A = [1.2, 2.3]
2 softmax(A)
```

: array([0.24973989, 0.75026011])

Data Format for Deep Learning Classifier

The training set consists of a matrix of features X and a vector of labels Y.

For each input observation $x^{(i)}$, we have a vector of features $[x_1, x_2, ..., x_n]$.

Feature *j* for input $x^{(i)}$ is x_i , more precisely $x_i^{(i)}$.



Multiclass Classification

We will consider the MNIST database of handwritten digits: this is the "Hello World" of deep learning.



Corinna Cortes Google Labe New York

Multiclass Classification

The MNIST digit database consists of 70,000 28x28 BW pixel images, stored as 28*28=784 floating-point numbers in the range [0..1]; labels are integers 0..9:



Recall: Supervised Machine Learning Workflow



Training involves multiple phases of evaluation with a validation set to find optimal values for the hyperparameters.

Cost/Loss Function for Classification

One more detail! What is the cost (loss) function used with the output of a classifier?



Image Features

Cost/Loss Function for Classification

We need to compare two probability distributions

Label:
$$[0, 0, 0, 0, 0, 0, 0, 0, 0, 0]$$
 \mathbf{y} Network output: $[p_0, p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8, p_9]$ $\mathbf{\hat{y}}$

to measure how different they are. The function used to do this is called the Cross-Entropy Loss:

$$\hat{\mathbf{y}} = softmax(\mathbf{a})$$
 $CE(\hat{\mathbf{y}}, \mathbf{y}) = -\frac{1}{n} \sum_{i=1}^{n} y_i \log(\hat{y}_i)$

One great feature of this cost function is that the derivative of the softmax plus CE Loss function is absurdly simple: it is just the difference of the two distributions:

$$\frac{\partial CE(\hat{y}_i, y_i)}{\partial z_i} = \hat{y}_i - y_i$$

Recall: Supervised Machine Learning Workflow



Training involves multiple phases of evaluation with a validation set to find optimal values for the hyperparameters.