

# Notes on Deep Learning

Qiang Hu

January 11, 2018

## Abstract

In this note, the concepts and knowledge on Deep Learning are recorded.

## 1 Introduction to Deep Learning

- Broadly speaking, deep learning is a more approachable name for an artificial neural network. The *deep* in deep learning refers to the depth of the network. An artificial neural network can be very shallow.
- The basic level of a neural network is the **perceptron**, the mathematical representation of a biological neuron. There are several layers of interconnected perceptrons in a neural network.
- **Deep learning** focuses on using **neural networks** for complex practical problems.
- Deep neural networks are used for *object recognition* and *image analysis*, for various modules of self-driving cars, for chatbots and natural language understanding problems.
- Linear models and stochastic optimization methods are crucial for training deep neural networks.
- Popular building blocks of neural networks includes **fully connected layers**, *convolutional* and *recurrent layers*. These blocks are used to define complex modern architectures in **TensorFlow** and **Keras** frameworks.
- Prerequisite knowledge:
  - **Machine learning** involves computer intelligence that doesn't know the answers up front. Instead, the program will run against training data, verify the success of its attempts, and modify its approach accordingly. Machine learning typically requires a sophisticated education, spanning software engineering and computer science to statistical methods and linear algebra.

- With **unsupervised learning**, there aren't any predefined or corresponding answers. The goal is to figure out the *hidden patterns* in the data. It's usually used for **clustering** and **associative tasks**, like grouping customers by behavior. Amazon's "customers who also bought..." recommendations are a type of associative task.
- Deep learning has proven to be an effective unsupervised learning technique. *Computer vision* has been a main beneficiary of deep learning. Speech recognition is another area that is felt deep learning's impact, because spoken languages are so vast and ambiguous. Through deep learning, generalizing the two languages (e.g., English and Mandarin) didn't require much additional design effort (i.e., design very different features).
- **Linear Regression**: a supervised machine learning algorithm where the predicted output is continuous and has a constant slope.
- **Mean Squared Error (MSE)**: as a cost function, it is used to optimize the weights. MSE measures the average squared difference between an observation's actual and predicted values. The output is a single number representing the cost, or score, associated with our current set of weights.

$$MSE = \frac{1}{N} \sum_{i=1}^n (y_i - y_i')^2 \quad (1)$$

$y_i$  is the actual value of an observation, and  $y_i'$  is the prediction.

- **Gradient descent**: a method to minimize MSE. Calculates the gradient of the cost function. To capture the impact of each weight on the final prediction, we use **partial derivatives**. To find the partial derivatives, we use the **Chain rule**. E.g.,

$$f'(m, b) = \left[ \begin{array}{c} \frac{df}{dp} \\ \frac{df}{db} \end{array} \right] \quad (2)$$

To minimize the cost function, we iteratively move in the direction of steepest descent as defined by the negative of the gradient.

- **Learning rate**: it is size of the update on the weights. The calculated gradient tells us the slope of our cost function at our current position (i.e., weights) and the direction we should update to reduce our cost function. We should move in the direction opposite the gradient.
- **Training**: training a model is the process of iteratively improving your prediction equation by looping through the dataset multiple times, each time updating the weights values in the direction indicated by the slope of the cost function. Training is completed when we reach an acceptable error threshold, or when subsequent training

iterations fail to reduce our cost. Before training, we need to initialize our weights (set default values), set our **hyperparameters** (i.e., learning rate and number of iterations).

The optimization problem can be formed as,

$$\min_w L(w) = \min_w \sum_{i=1}^l L(w; x_i, y_i) \quad (3)$$

$L(w)$  is the loss (cost) function.

$w^0$  is the initial weights.

while True,

$$w^t = w^{t-1} - \eta_t \nabla L(w^{t-1}) \quad (4)$$

if  $\|w^t - w^{t-1}\| < \epsilon$ , break

Note that,  $l$  gradients should be computed on each step. If the dataset doesn't fit in memory, it should be read from the disk on every GD step.

- **Normalization:** As the number of features grows, calculating gradient takes longer to compute. We can speed this up by “normalizing” our input data to ensure all values are within the same range. This is especially important for datasets with high standard deviations or differences in the ranges of the attributes. Our goal now will be to normalize our features so they are all in the range -1 to 1. First, subtract the mean of the column (mean normalization); second, divide by the range of the column (feature scaling).
- **Logistic regression:** a classification algorithm used to assign observations to a discrete set of classes. Logistic regression transforms its output using the logistic sigmoid function to return a *probability* value which can then be mapped to two or more discrete classes.
- **Sigmoid activation:** it is used to map predicted values to probabilities. The function maps any real value into another value between 0 and 1.

$$S(z) = \frac{1}{1 + e^{-z}} \quad (5)$$

It's easy to calculate the derivative of the sigmoid function.

$$s'(z) = s(z)(1 - s(z)) \quad (6)$$

- **Decision boundary:** the threshold value or tipping point used to map the probability value to a discrete class.
- **Prediction function:** in logistic regression, a prediction function returns the probability of our observation being positive, True, or “Yes”. We call this class 1 and its notation is  $P(class = 1)$ . As the probability gets closer to 1, our model is more confident that the observation is in class 1.

- **Cross-Entropy**: a.k.a., **Log Loss**. It can be divided into two separate cost functions: one for  $y = 1$  and one for  $y = 0$ .

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) \\ \text{Cost}(h_{\theta}(x), y) &= -\log h_{\theta}(x) \quad \text{if } y = 1 \\ \text{Cost}(h_{\theta}(x), y) &= -\log (1 - h_{\theta}(x)) \quad \text{if } y = 0 \end{aligned} \quad (7)$$

The key thing to note is the cost function penalizes confident and wrong predictions more than it rewards confident and right predictions!

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \quad (8)$$

And the derivative of the cost function is:

$$C' = x \cdot (s(z) - y) \quad (9)$$

- **Multiclass logistic regression**: When  $y = 0, 1, 2, \dots, n$ ,
  1. Divide the problem into  $n + 1$  binary classification problems;
  2. For each class ...
  3. Predict the probability the observations are in that single class.
  4. prediction = max(probability of the classes)
- **Softmax transform**: transfer the score to probability.

$$z = (w_1^T x, \dots, w_K^T x) \Rightarrow (e^{z_1}, \dots, e^{z_K}) \quad (10)$$

$$\sigma(z) = \left( \frac{e^{z_1}}{\sum_{k=1}^K e^{z_k}}, \dots, \frac{e^{z_K}}{\sum_{k=1}^K e^{z_k}} \right) \quad (11)$$

- **Model regularization**: add a regularizer  $R(w)$  to the original loss function  $L(w)$ . The regularizer penalizes the model for large weights.  $\lambda$  is the regularization strength which controls the model quality on a training set and model complexity.

$$L_{\text{reg}}(w) = L(w) + \lambda R(w) \Rightarrow \min_w L_{\text{reg}}(w) \quad (12)$$

To minimize the regularized loss function  $L_{\text{reg}}(w)$ , it equals to the optimization problem,

$$\begin{cases} \min_w L(w) \\ \text{s.t. } R(w) \leq C \end{cases} \quad (13)$$

If the L1 penalty  $\|w\| = \sum |w_i|$  is used,

- \* Drives some weights exactly to zero;
- \* Learns sparse models;

- \* Cannot be optimized with simple gradient methods.

If the L2 penalty  $\|w\|^2 = \sum w_j^2$  is used,

- \* Drives all weights closer to zero;
- \* Can be optimized with gradient methods.

Other regularization techniques:

- \* Dimensionality reduction, e.g., remove features, principal component analysis;
- \* Data augmentation, e.g., upon image data, flip, rotate the image to get more data;
- \* Dropout;
- \* Early stopping;
- \* Collect more data.

- **Stochastic gradient descent:** Similar to the regular GD, the training starts with an initial  $w_0$ , but in every step, we randomly choose an example with index  $i$  between 1 and  $l$  to update the weight. This process leads to very noisy approximations. But if we make enough numbers of iterations, it converges to some minimum.

- \* Noisy updates lead to fluctuations;
- \* Needs only one example on each step;
- \* Can be used in online setting;
- \* Learning rate  $\eta_t$  should be chosen very carefully.

- **Mini-batch gradient descent:** Similar to the stochastic GD, but in every step, we randomly choose  $m$  examples from the dataset to calculate the gradient.

- \* Still can be used in online setting;
- \* Reduces the variance of gradient approximations;
- \* Learning rate  $\eta_t$  should be chosen very carefully.

- **Gradient descent extensions:** to optimize GD for difficult functions with complex level sets, we can define the **momentum** function  $h_t = \alpha h_{t-1} + \eta_t g_t$ , where  $g_t$  is the gradient calculated. Then we update the weight using  $h_t$ , i.e.,  $w^t = w^{t-1} - h_t$ .

- \* It tends to move in the same direction as on previous steps;
- \*  $h_t$  accumulates values along dimensions where gradients have the same sign;
- \* Usually,  $\alpha = 0.9$ .

**Nesterov momentum:**  $h_t = \alpha h_{t-1} + \eta_t \nabla L(w^{t-1} - \alpha h_{t-1})$

- **AdaGrad:** choose learning rate adaptively.

$$\begin{aligned} G_j^t &= G_j^{t-1} + g_{t,j}^2 \\ w_j^t &= w_j^{t-1} - \frac{\eta_t g_{t,j}}{\sqrt{G_j^t + \epsilon}} \end{aligned} \tag{14}$$

where  $g_{t,j}$  is the gradient with respect to  $j$ -th parameter.

- \* Separate learning rates for each dimension;
- \* Suits for sparse data
- \* Learning rate can be fixed as  $\eta_t = 0.01$ ;
- \*  $G_j^t$  always increases, leads to early stops.

– **RMSprop:**

$$\begin{aligned} G_j^t &= \alpha G_j^{t-1} + (1 - \alpha) g_{t,j}^2 \\ w_j^t &= w_j^{t-1} - \frac{\eta_t g_{t,j}}{\sqrt{G_j^t + \epsilon}} \end{aligned} \quad (15)$$

- \*  $\alpha$  is about 0.9;
- \* Learning rate adapts to latest gradient steps.

– **Adam:**

$$\begin{aligned} m_j^t &= \frac{\beta_1 m_j^{t-1} + (1 - \beta_1) g_{t,j}}{1 - \beta_1^t} \\ v_j^t &= \frac{\beta_2 v_j^{t-1} + (1 - \beta_2) g_{t,j}^2}{1 - \beta_2^t} \\ w_j^t &= w_j^{t-1} - \frac{\eta_t m_j^t}{\sqrt{v_j^t + \epsilon}} \end{aligned} \quad (16)$$

Note in the first step, the normalization is very large, but as the iteration goes on, the denominator is close to 1. This is used to remove the bias of  $v$  closing to 0 at the beginning. It also combines momentum and individual learning rates.

## 2 Introduction to neural networks

- In image processing, linear models are not helpful. Because pixel representations are dependencies, pixel representation in the class is very complex and very valuable.
- Transform the objects into some space in which the problem is easy.
- The old method to do image processing is that human experts construct features from lines and edges, which are combined into simple parts of image, into features, which are again, combined into discrete choices.
- Directly stacking multiple layers of linear models does not help the non-linear classification, because it still is linear models in use.
- Applying the output of the previous linear model to a non-linear function introduces the non-linearity. Some examples of non-linear functions:

$$f(a) = \frac{1}{1 + e^a} \quad (17)$$

$$f(a) = \max(0, a) \quad (18)$$

$$f(a) = \tanh(a) \quad (19)$$

$$f(a) = \log(1 + e^a) \quad (20)$$

Most of them increase monotonically. And most of them are constant when the argument approaches  $-\infty$ . The choice usually depends on the network architecture.

- **Layer** is a building block for neural network:

- Input layer: input data;
- Dense layer:  $f(x) = Wx + b$  (linear regressors);
- Non-linearity layer:  $f(x) = \sigma(x)$ ;
- A few more: ...
- Output layer

All layers except input and output layers are *hidden layers*. There are generally no connections between nodes in the same layer.

- **Activation** is the output of a layer. (i.e., some intermediate signal in the neural network) The neurons in the hidden layers are only activated by nodes in the previous layer.
- Deep model is just a (composite) function.
  - $f_1(f_2(f_3(x, w_3), w_2), w_1) \rightarrow \text{loss}$
  - $x$  is feature vector
  - $w_{1,2,3}$  are model parameters.
- The differentiation of the above deep model function is a *Jacobian matrix* when  $X = [x_1, x_2, \dots, x_n]$ ,  $f = [f_1, f_2, \dots, f_m]$ :

$$\frac{df}{dx} = \begin{bmatrix} \frac{df}{dx_1} & \dots & \frac{df}{dx_n} \end{bmatrix} = \begin{bmatrix} \frac{df_1}{dx_1} & \dots & \frac{df_1}{dx_n} \\ \dots & \dots & \dots \\ \frac{df_m}{dx_1} & \dots & \frac{df_m}{dx_n} \end{bmatrix} \quad (21)$$

- When  $F : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{p \times q}$ ,  $X$  is a  $m \times n$  matrix,  $F$  is also a  $p \times q$  matrix. Thus,  $F^i$ 
  - a matrix of matrices
  - fourth-rank tensor
  - four-dimensional array

$$\frac{\partial F}{\partial X} = \begin{bmatrix} \frac{\partial F}{\partial X_{11}} & \dots & \frac{\partial F}{\partial X_{1n}} \\ \dots & \dots & \dots \\ \frac{\partial F}{\partial X_{m1}} & \dots & \frac{\partial F}{\partial X_{mn}} \end{bmatrix} \quad (22)$$

$$\frac{\partial F}{\partial X_{ij}} = \begin{bmatrix} \frac{\partial F_{11}}{\partial X_{ij}} & \dots & \frac{\partial F_{1q}}{\partial X_{ij}} \\ \dots & \dots & \dots \\ \frac{\partial F_{p1}}{\partial X_{ij}} & \dots & \frac{\partial F_{pq}}{\partial X_{ij}} \end{bmatrix} \quad (23)$$