

# UCSB, Physics 129L, Computational Physics

## Lecture notes, Week 9

Zihang Wang (UCSB), zihangwang@ucsb.edu

March 12, 2025

### Contents

<b>1</b>	<b>Multilayer neural network, optimal transport, and deep learning</b>	<b>1</b>
1.1	Regression Problem: maximum likelihood estimation . . . . .	2
1.2	Classification Problem: Cross entropy . . . . .	3
1.3	Transport cost function (loss function) . . . . .	4
1.4	Linear Regression as a Single-Layer Supervised Learning Algorithm	5
1.5	Layers in neural network . . . . .	5
1.6	Multilayer Perceptron (MLP) . . . . .	6

## 1 Multilayer neural network, optimal transport, and deep learning

Deep Learning is a subfield of machine learning that employs neural networks to model complex patterns in data. Unlike traditional machine learning approaches, deep learning models automatically learn hierarchical representations from raw data. In other words, we do not need to impose any pre-known models and assumptions.

When data is **labeled**, we call it **supervised deep learning**.  $\{X, y\}$  Both  $X$  and  $y$  can be scalar or vector. On the other hand, if the data is **unlabeled**,  $X$ , it is **unsupervised deep learning**.

In supervised deep learning, the task is to learn the relationship between  $y$  and  $X$ , such the prediction  $\hat{y}$  from the model  $\hat{y} = T(x)$  is as close as possible to true value.  $T$  is analogous to the transport map (plan) that represents large amount of hyperparameters that give the deep learning enough of degree's of freedom. We should note that the  $x$  and  $\hat{y}$  does not need to have the same dimension. This is similar to the **Kantorovich problem**. Number of parameters in state-of-the-art neural networks: **GPT-3: 175 billion GPT-4: 1.8 trillion**.

The performance of a neural network is measured by a **loss function**  $E$ . The loss function defined in the neural network shares the same idea as we have

in the **optimal transport**: we want to find the optimal way of designing a transport map  $T^*$  that minimizes the cost of transporting a source distribution  $\mu$  (model output) to a target distribution  $\nu$  (training set that contains true labels).

The question becomes, how to find the optimal transport plan  $T^*$  such that the predicted  $\hat{y} = T(x)$  can match the true value  $y$ : **The best map  $T^*$  is the ones that give the smallest prediction error on training data.**

In the perspective of optimal transport, the map is defined by a scalar potential gradient  $\nabla\varphi = T$  that minimizes the **transport cost**,

$$C(T) = \int_X c[x, T(\mathbf{x})] d\mu(\mathbf{x}) = \int_X c[\mathbf{x}, T(\mathbf{x})] \mu(\mathbf{x}) d\mathbf{x},$$

where in supervised learning, it is defined as the **loss function** quantifies the discrepancy between predicted values and actual values. The choice of loss function depends on the type of learning problem: **regression** and **classification**.

### 1.1 Regression Problem: maximum likelihood estimation

In regression tasks, the target variable  $y$  is continuous, meaning  $y \in \mathbb{R}$ . The most commonly used loss functions are: **Mean Squared Error** (MSE),

$$C(T) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

where  $\hat{y}_i = f_\theta(X_i)$  is the predicted value based on the transport plan  $T$ . MSE penalizes larger errors more heavily due to the squared term. There is also the **Mean Absolute Error** (MAE),

$$C(T) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|,$$

MAE gives equal weight to all errors, making it more robust to outliers compared to MSE.

Minimizing MSE corresponds to maximizing the likelihood under the assumption that errors are **normally distributed** (refer to lecture notes on stochastic calculus and log likelihood):

$$y_i = T(X_i) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2).$$

Recall in previous lecture notes, the MSE objective function is derived from the **negative log-likelihood** ( $\chi$ -square) of the Gaussian distribution (with **Frequentist statistics or Bayesian inference with uniform prior**),

$$C(T) = - \sum_{i=1}^n \log P(y_i | X_i, \theta) \sim \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

**This shows that minimizing MSE is equivalent to performing maximum likelihood estimation (MLE) under a Gaussian noise assumption.**

## 1.2 Classification Problem: Cross entropy

In **classification problems**, the target variable  $y$  is discrete and belongs to a finite set of classes. For a classification problem with  $K$  classes, the model assigns a raw score (**logit**)  $z_k$  to each class  $k$ . For a given input  $X_i$ , the logits are computed as,

$$z_k = w_k^T x_i + b_k,$$

where  $w_k$  is the weight for class  $k$ ,  $b_k$  is the bias for class  $k$ ,  $x_i$  is the input feature vector for sample  $i$ . The logits  $z_k$  do not represent probabilities. Instead, **they are arbitrary scores that can be positive or negative**. To convert the logits into probabilities, we apply the **softmax function**:

$$p_T(y_i = k|X_i) = \frac{\exp(z_k)}{\sum_{j=1}^K \exp(z_j)}$$

This ensures the probability  $p_T(y_i = k|X_i)$  is always in the range  $(0, 1)$ , and the probabilities sum to 1 across all classes. This is very similar to the Boltzmann weights we discussed previously.

Given a dataset with  $n$  samples, the cross-entropy loss for multi-class classification is:

$$C(T) = - \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}(y_i = k) \log p_T(y_i = k|X_i),$$

$\mathbb{1}(y_i = k)$  is an indicator function that equals 1 if the true class for  $X_i$  is  $k$ , otherwise 0 ( **one-hot encoded**). -  $p_T(y_i = k|X_i)$  is the predicted probability for class  $k$ . Since  $y_i$  is the true value, it only selects a single value out of this sum by the indicator function  $\mathbb{1}(y_i = k)$ .

Consider a classification problem with  $K = 3$  classes, where  $y \in \{0, 1, 2\}$ . Suppose for an input  $X_i$ , the model predicts the logits:

$$z_0 = 1.5, \quad z_1 = 2.0, \quad z_2 = 0.5.$$

Applying the **softmax function**:

$$p_\theta(y_i = 0|X_i) = \frac{e^{1.5}}{e^{1.5} + e^{2.0} + e^{0.5}} \approx 0.29$$

$$p_\theta(y_i = 1|X_i) = \frac{e^{2.0}}{e^{1.5} + e^{2.0} + e^{0.5}} \approx 0.49$$

$$p_\theta(y_i = 2|X_i) = \frac{e^{0.5}}{e^{1.5} + e^{2.0} + e^{0.5}} \approx 0.22$$

If the true label is  $y_i = 1$ , then the loss for this sample is:

$$L_1 = -\log p_\theta(y_i = 1|X_i) = -\log(0.49) \approx 0.71$$

If the true label was  $y_i = 2$ , the loss would be:

$$L_2 = -\log p_\theta(y_i = 2|X_i) = -\log(0.22) \approx 1.51$$

The cross-entropy loss directly corresponds to the negative log-likelihood function in maximum likelihood estimation,

$$C(T) = -\sum_{i=1}^n L_i.$$

Thus, **minimizing cross-entropy loss is equivalent to maximizing the likelihood of the observed data.**

In **unsupervised deep learning**, how well the model captures or reconstructs the inherent structure of the input data, rather than comparing its output to an externally provided label.

### 1.3 Transport cost function (loss function)

To update the weights in the network, we compute the gradient of the loss with respect to the network outputs.

**Cross entropy** measures the dissimilarity between two probability distributions. In **classification**, we often have a one-hot encoded target distribution (delta function) and a predicted probability distribution. **In regression, however, you're predicting continuous values rather than a distribution.**

**Regression** models typically output real numbers without the normalization constraint (i.e., they don't sum to 1 like probabilities). **Using cross entropy would require interpreting these outputs as probabilities, which usually doesn't align with the nature of regression tasks.**

In binary classification, even though the problem is about predicting a class (0 or 1), the model typically outputs a continuous value that represents a probability after applying a sigmoid function.

In binary classification, we use a single output neuron with a sigmoid activation function, which converts the logit  $z$  into a probability:

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}.$$

Since there are only two classes (0 and 1), one probability  $\hat{y}$  is enough because:

$$P(y = 1|X) = \hat{y}, \quad P(y = 0|X) = 1 - \hat{y}$$

The two probabilities are complementary and always sum to 1. However, **in multi class case, we do not have the complementary property.**

Therefore, if we want to interpret the output as probabilities, we need to use cross entropy method. On the other hand, if we look for prediction in values, we need to use mean-square.

## 1.4 Linear Regression as a Single-Layer Supervised Learning Algorithm

**Linear regression** can be viewed as a simple supervised learning algorithm with one layer,

$$\hat{y}_i = T(X_i) = \theta_0 + \theta_1 X_i,$$

where  $\theta = (\theta_0, \theta_1)$  are the model parameters that generate transport map from  $X_i$  to  $y_i$ . These parameters can be obtained via the stationary condition of the negative log-likelihood: It minimizes the total sum of squared residuals ( $\chi$  square),

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

This corresponds to minimizing the MSE loss.

For optimization, the gradient of the MSE loss function with respect to  $\theta$  is,

$$\frac{\partial C_{\theta}}{\partial \theta_1} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) X_i, \quad \frac{\partial C_{\theta}}{\partial \theta_0} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i).$$

which leads to updates in gradient descent:

$$\theta_{1,\text{new}}^{(l)} = \theta_{1,\text{old}}^{(l)} - \alpha_1 \frac{\partial C(T)}{\partial \theta_1^{(l)}}, \quad \theta_{0,\text{new}}^{(l)} = \theta_{0,\text{old}}^{(l)} - \alpha_0 \frac{\partial C(T)}{\partial \theta_0^{(l)}}$$

where  $\alpha$  is the learning rate.

## 1.5 Layers in neural network

A **neural network** comprises layers of interconnected nodes (neurons) that are usually fully-connected.

The main components include:

- **Input Layer:** Receives the raw data, denoted as  $\mathbf{x} \in \mathbb{R}^n$ .
- **Hidden Layers:** Intermediate layers that perform non-linear transformations. A network can have one or multiple hidden layers.
- **Output Layer:** Produces the final output, such as a class label or a regression value.

The term **deep** in deep learning refers to **the number of layers (or depth)** in the neural network. A neural network is considered deep when it has multiple hidden layers (beyond just an input and output layer). The **width of a neural network** refers to the **number of neurons per layer**.

Each neuron receives an  $n$  dimension data  $\mathbf{z}_1$  produced by the previous layer, and performs a linear transformation with **weight matrix**  $\theta_1^l$  and **bias**  $\theta_0^l$ , that

acts on previous layer data  $z^{(l-1)}$ . **In this case, we set the batch number to be 1.**

$$z^{(l)} = \theta_1^{(l)} a^{(l-1)} + \theta_0^{(l)}, \quad \theta_1^{(1)} \in \mathbb{R}^{m \times n}, \quad \theta_0^{(1)} \in \mathbb{R}^m \quad (1)$$

Since joint linear transformations can be written effectively as a single linear equation, we must introduce some non-linearity. To do that, we introduce the **activation function**,  $f()$ , that acts on the current layer output,

$$a^{(l)} = f(z^{(l)}). \quad (2)$$

Common activation functions  $f()$  include:

- **Sigmoid:**  $\sigma(z) = \frac{1}{1 + e^{-z}}$
- **ReLU:**  $\text{ReLU}(z) = \max(0, z)$
- **Tanh:**  $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

## 1.6 Multilayer Perceptron (MLP)

**Multilayer Perceptron (MLP)** is a simple example of the **feedforward network** with one hidden layer (no loop). Let the input vector be  $\mathbf{x} \in \mathbb{R}^n$ , the hidden layer contain  $m$  neurons, and the output be a scalar  $y$ . **We say they are fully connected: all neurons are connected to the neurons in previous and forward layers.**

The forward propagation is described as follows: We first have the **input layer** that takes  $n$  inputs and processed by  $n$  neurons. therefore, it outputs  $m$  dimensional vector,

$$z^{(1)} = \theta_1^{(1)} \mathbf{x} + \theta_0^{(1)}, \quad \theta_1^{(1)} \in \mathbb{R}^{m \times n}, \quad \theta_0^{(1)} \in \mathbb{R}^m \quad (3)$$

$$a^{(1)} = f(z^{(1)}). \quad (4)$$

We then pass the vector to the activation function. The result will be pass on to the **hidden layer**. The input of the hidden layer becomes  $m$  dimensional and the output becomes  $r$  dimensional,

$$z^{(2)} = \theta_1^{(2)} a^{(1)} + \theta_0^{(2)}, \quad \theta_1^{(2)} \in \mathbb{R}^{r \times m}, \quad \theta_0^{(2)} \in \mathbb{R}^r \quad (5)$$

$$a^{(2)} = f(z^{(2)}). \quad (6)$$

Then, we have the **output layer**:

$$z^{(3)} = \theta_1^{(3)} a^{(2)} + \theta_0^{(3)}, \quad \theta_1^{(3)} \in \mathbb{R}^{1 \times r}, \quad \theta_0^{(3)} \in \mathbb{R} \quad (7)$$

$$\hat{y} = a^{(3)} = \tilde{f}(z^{(3)}). \quad (8)$$

In this case,  $\tilde{f}$  is usually the softmax for classification problem and activation function if in regression problem.