

Introduction to Deep Learning

Machine Learning

Linear Regression

$$\hat{y}_i = \sum_{j=1}^d x_{ij} \theta_j$$

d = input dimension

x = input data

θ = weights

Loss function:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

in matrix notation:

$$J(\theta) = (\mathbf{X}\theta - \mathbf{y})^T (\mathbf{X}\theta - \mathbf{y})$$

differentiate:

$$\frac{\partial J(\theta)}{\partial \theta} = 2\mathbf{X}^T \mathbf{X} \theta - 2\mathbf{X}^T \mathbf{y} = 0$$

$$\theta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Analytical solution for a convex problem.

Logistic Regression

$$\hat{\mathbf{y}} = p(\mathbf{y} = 1 | \mathbf{X}, \theta) = \prod_{i=1}^n p(y_i = 1 | \mathbf{x}_i, \theta)$$

where

$$\hat{y}_i = \sigma(\mathbf{x}_i \theta)$$

σ is the sigmoid/logistic function.

Loss function (binary cross entropy loss):

$$L(\hat{y}_i, y_i) = -[y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)]$$

No analytical solution. Iterative method needed (gradient descent).

Loss Functions

Measure the goodness of the prediction.

Large loss indicates bad performance.

Choice of loss function depends on problem.

Regression Loss

L1 Loss:

$$L(y, \hat{y}; \theta) = \frac{1}{n} \sum_i^n ||y_i - \hat{y}_i||_1$$

MSE Loss:

$$L(y, \hat{y}; \theta) = \frac{1}{n} \sum_i^n ||y_i - \hat{y}_i||_2^2$$

where

L1 norm (Manhattan distance):

$$||x||_1 = \sum |x_i|$$

- Promotes sparsity
- More robust to outliers

L2 norm (Euclidean distance):

$$||x||_2 = \sqrt{\sum x_i^2}$$

- Differentiable
- Penalizes large errors more heavily

Squared L2 norm:

$$||x||_2^2 = \sum x_i^2$$

- Differentiable

- Computationally efficient

Classification Loss

Binary Cross Entropy loss (binary classification):

$$L(y, \hat{y}; \theta) = -\frac{1}{n} \sum_i^n [y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)]$$

y is a probability output ($[0, 1]$) of the first class. Therefore, $1 - y$ is the probability of the second class. Works well with a sigmoid activation in the last layer.

Log penalizes bad predictions: $\log 0 = -\infty$, $\log 1 = 0$ (sort of inverts the prediction \hat{y}).

Encourages the model to make confident predictions close to 0 or 1.

Cross Entropy loss (multi-class classification):

$$L(y, \hat{y}; \theta) = - \sum_{i=1}^n \sum_{k=1}^k (y_{ik} \cdot \log \hat{y}_{ik})$$

The second sum is a generalization of the binary cross entropy case for multiple classes.

Backpropagation

Backpropagation aims to minimize the loss function by adjusting network weights and biases. The level of adjustment is determined by the gradients of the loss function with respect to those parameters.

Neural network is a chain of functions: linear transformations and activation functions. During backpropagation, the gradients of each linear transformation are calculated using the chain rule, starting from the derivative of the loss function and going backwards through the network.

Chain rule:

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \frac{\partial u}{\partial x}$$

or

$$h'(x) = f'(g(x))g'(x)$$

Calculate the gradients of the last layer:

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}}$$

where z are the outputs of the layer before the activation function.

For any weight update in the hidden layers, sum all gradients of all neurons that are influenced by that weight.

Useful Derivatives

ReLU:

$$\frac{\partial ReLU}{\partial x} = \text{if } x < 0 \text{ then } 0; \text{ if } x > 0 \text{ then } 1$$

Sigmoid:

$$\frac{\partial \sigma}{\partial x} = \sigma(x)(1 - \sigma(x))$$

Gradient descent

Gradient update steps in direction of negative gradient:

$$W' = W - \alpha \nabla f(W)$$

where α is the learning rate.

Regularization

Prevent overfitting with regularization.

Add regularization to the loss:

$$L'(y, \hat{y}; \theta) = L(y, \hat{y}; \theta) + \lambda R(\theta)$$

L2 Regularization:

$$R(\theta) = \sum_{i=1}^n \theta_i^2$$

- Is lower for smaller, distributed values.
- Enforces that weights have similar values.
- Model will take all information into account to make decisions.

L1 Regularization:

$$R(\theta) = \sum_{i=1}^n |\theta_i|$$

- Is overall lower than L2 with some higher weight values.
- Enforces sparsity.
- Model will focus attention on a few key features.

Regularization coefficient λ sets the level of regularization (higher λ means more regularization).