

$$f(\underline{a}, \underline{b}) = \frac{1}{1 + \exp(-\underline{a}^T \underline{b})}$$

$$\underline{a} \in \mathbb{R}^n$$

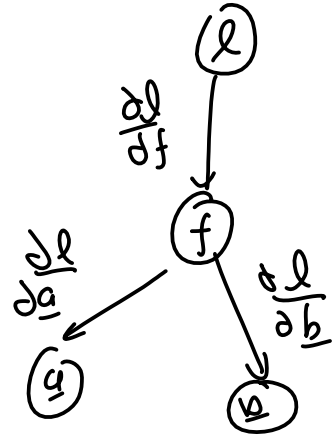
$$\underline{b} \in \mathbb{R}^n$$

$$f \in \mathbb{R}$$

Vector Jacobian Product.

$$\frac{\partial \ell}{\partial \underline{a}} =$$

$$\frac{\partial \ell}{\partial \underline{b}} =$$



$$\frac{d}{dx} \frac{1}{x} = ? \quad \Rightarrow \frac{dx^{-1}}{dx} = -1(x)^{-1-1} = -x^{-2} = -\frac{1}{x^2}$$

$$\frac{d}{dx} \exp(x) = ? = \exp(x)$$

$$\frac{\partial}{\partial \underline{a}} \underline{a}^T \underline{b} = \underline{b}^T \quad \text{and} \quad \frac{\partial}{\partial \underline{b}} \underline{a}^T \underline{b} = \underline{a}^T$$

$$\frac{\partial \ell}{\partial \underline{a}} = \frac{\partial \ell}{\partial f} \left(\frac{-1}{(1 + \exp(-\underline{a}^T \underline{b}))^2} \right) \left(\exp(-\underline{a}^T \underline{b}) \right) \frac{\partial (-\underline{a}^T \underline{b})}{\partial \underline{a}}$$

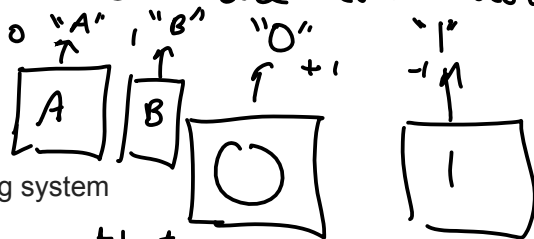
$$= \frac{\partial \ell}{\partial f} \frac{\underline{b}^T \exp(-\underline{a}^T \underline{b})}{(1 + \exp(-\underline{a}^T \underline{b}))^2}$$

$$\frac{\partial \ell}{\partial \underline{b}}$$

ML Data, Models and Learning

ML
 { Supervised
 { Classification = Output labels are discrete
 { Regression = Output labels are continuous
 { Unsupervised

Ref: MML Book (Chapter 8) <https://mml-book.github.io/>



There are three major components of a machine learning system

1. Data

Input data \rightarrow output labels
 SUPERVISED Learning

- A. Example: Handwritten digit images and corresponding labels
- B. Example: Road density and Salt concentration
- C. Example: X-Y coordinate of point cloud and corresponding Z coordinate

2. Models

Input data \rightarrow Predicted label
 function \rightarrow Model

- A. Example: Linear model: Equation of line or plane
- B. Example: Multi Layer perceptron: Two Linear models sandwiching a non-linear activation function.

3. Learning

- A. Example: Least square solution.
- B. Example: Gradient descent.

$$\frac{\partial}{\partial \underline{x}} l(\underline{x}) \Big|_{\underline{x}^*} = 0$$

$$\underline{x}_{t+1} = \underline{x}_t - \alpha \frac{\partial}{\partial \underline{x}} l(\underline{x})$$

Data as Vectors

Let us consider the problem of identifying the digit from handwritten images based on data. This is called a supervised learning problem, where we have a label y_i (the digit) associated with each example \underline{x}_i (the handwritten image). The label y_i has various other names, including target, response variable and annotation. A dataset is written as a set of example-label pairs $\{(\underline{x}_1, y_1), \dots, (\underline{x}_i, y_i), \dots, (\underline{x}_n, y_n)\}$. The features $\{\underline{x}_1, \dots, \underline{x}_n\}$ are often concatenated and written as a big matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ and the labels $\mathbf{y} \in \mathbb{R}^n$.

Models as functions

Once we have data in an appropriate vector representation, we can get to the business of constructing a predictive function (known as a predictor).

There are two major approaches:

- 1. a predictor as a function, $f: \mathbb{R}^d \mapsto \mathbb{R}$

A. Example: Linear Model: $f(\underline{x}) = \underline{w}^T \underline{x} + w_0$

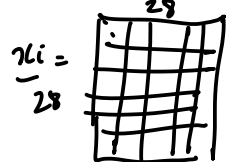
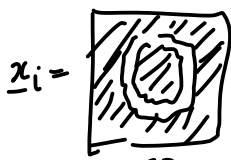
$$= w_1 x_1 + w_2 x_2 + \dots + w_n x_n + w_0$$

$$f(\underline{x}) = y \in \mathbb{R} \quad \underline{x} \in \mathbb{R}^d$$

$$f(x) = mx + c$$

Predicted label $\hat{y} = m x + c$
 Input

Predicted label $\hat{y} = \text{sign}(\hat{y} - (m x + c))$
 Input



Raw input

$\underline{x}_i = \begin{cases} \# \text{ of white pixels} \\ \text{spread of white pixels on y-axis} \end{cases}$
 feature vector

Input Output label

Training data

$$\mathbf{X} = \begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}$$

activation function

$$f(x) = \underline{w}_2^T \sigma(\underline{w}_1 x + \underline{w}_0) + w_{20} v_2 \in \mathbb{R}^m \quad \begin{matrix} w_1 \in \mathbb{R}^{m \times d} \\ w_0 \in \mathbb{R}^m \end{matrix}$$

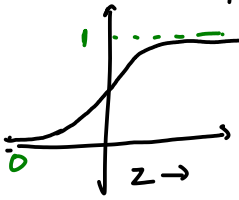
$$\hookrightarrow \text{ReLU}(z) = \begin{cases} z & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

B. Example: MLP: $f(x) = \underline{w}_2^T \sigma(\underline{w}_1 x + \underline{w}_0)$, where $\sigma : \mathbb{R}^m \mapsto \mathbb{R}^m$ is some non-linear activation function like ReLU, sigmoid or tanh.

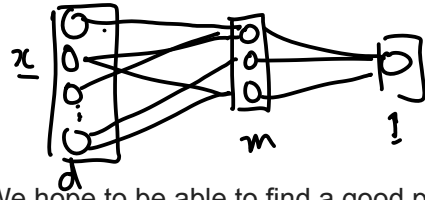
2. a predictor as a probabilistic model. Later

σ is any activation function

$$\text{sigmoid}(z) = \frac{1}{1 + \exp(-\alpha z)}$$



Hypothesis class of functions



$f(x)$

MLP

A predictor $f : \mathbb{R}^d \mapsto \mathbb{R}$, parametrized by θ . We hope to be able to find a good parameter θ^* such that we fit the data well, that is, $f(x_i, \theta^*) \cong y_i$ for all $i = 1, \dots, n$.

We use the notation $\hat{y}_i = f(x_i, \theta^*)$ to represent the output of the predictor.

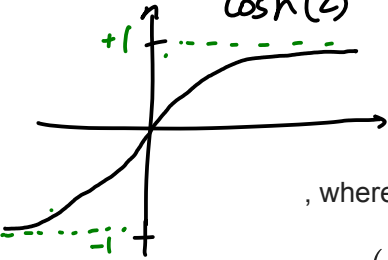
$\tanh(z)$ = Hyperbolic tangent

$$f(x) = \underline{w}^T x + w_0$$

Loss Function for Training

$$\begin{cases} f(x) = mx + c \\ f_1(x) = 2x + 1 \\ f_2(x) = 3x + 0 \end{cases}$$

$$\tanh(z) = \frac{\sinh(z)}{\cosh(z)}$$



$$r_{\text{emp}}(f, \mathbf{X}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n l(y_i, \hat{y}_i)$$

, where $\hat{y}_i = f(x_i, \theta^*)$.

$r_{\text{emp}}(f, \mathbf{X}, \mathbf{y})$ is called the empirical risk.

$f(x; m, c)$ ← parameters

Homework 5

MLP

$$f(\underline{x}) = \underline{w}_2^T a \left(\underline{w}_1 \underline{x} + \underline{w}_0 \right) + w_{20}$$

activation function

Parameters in MLP = $\{ \underline{w}_2, w_{20}, \underline{w}_1, \underline{w}_0 \}$

Linear model

$$f(x) = \underline{w}^T x + w_0$$

Parameter = $\{ \underline{w}, w_0 \}$

3-Layer Perceptron or 3-Layer NN

$$f(\underline{x}) = \underline{w}_3^T a_2 \left(w_2 a_1 (\underline{w}_1 \underline{x} + \underline{w}_0) + \underline{w}_2 \right) + w_{30}$$

Parameters = $\{ \underline{w}_3, w_2, w_1, \underline{w}_1, \underline{w}_2, w_{30} \}$

Learning:

while (not converged):

$$\text{Parameters}_{t+1} = \text{Parameters}_t - \alpha \frac{\partial l}{\partial \text{Parameter}}$$

$$f(\underbrace{\underline{x}}_{\text{input}}; \underbrace{\underline{\theta}}_{\text{Parameters (weights and biases)}})$$

$f_{\underline{\theta}}(\underline{x})$

① Data $\{ (x_i, y_i) \dots (x_n, y_n) \}$

② Model -

③ Learning

↳ Gradient Descent

↳ Loss function

Model: $f(x; \theta)$

→ $\hat{y}_i = f(x_i; \theta)$

$$\hat{y}_i = f(x_i; m, c)$$

$$\hat{y}_i = mx_i + c$$

common
notation
for
predicted
label

$$\hat{y}_i \approx y_i$$

Least square loss function

$$\text{Loss function } \ell(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2 = (y_i - \underbrace{mx_i + c}_{\hat{y}_i})^2$$

Thresholded L1 loss

$$\hat{y}_i = f(x_i; \underline{w}) = \text{sign}(\underline{w}^T x_i + w_0)$$

$$\hat{y}_i \in \{-1, +1\}$$

$$\ell(y_i, \hat{y}_i) = \begin{cases} 0 & \text{if } y_i = \hat{y}_i \\ \underbrace{|\underline{w}^T x_i + w_0|}_{\text{abs}} & \text{if } y_i \neq \hat{y}_i \end{cases}$$

The learning problem in general is formulated as an optimization problem

$$\theta^* = \arg \min_{\theta} \pi_{\text{emp}} \left(f, \{(\underline{x}_1, y_1), \dots, (\underline{x}_n, y_n)\} \right)$$

$$\underbrace{\pi_{\text{emp}}(f, \dots)} = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{y}_i)$$

} Data
(\underline{x}_i, y_i)
being
identically
independently
distributed.
(i.i.d.)