

Given data $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}^N$

Find a function

$f: \mathbb{R}^D \rightarrow \mathbb{R}^N$ so that

[1] $f(x_i) \approx y_i \quad 1 \leq i \leq n$

[2] For each new data point (x, y) , $f(x) \approx f(y)$



Train Accuracy : 100%

Test Accuracy: 10%

x_i = input data / attributes

y_i = labels / output variables

f_θ = model / depending on θ

Machine learning - goal is to learn θ

Algorithm 3.1

1. Select a model
2. Split data into training + test data
3. Learn parameters θ
4. Validation of θ , choose of parameter

Validation options

~ Choose parameter θ^* which minimizes

ERM • Empirical risk $\sim R(\theta) = \frac{1}{n} \sum_{i=1}^n l(y_i, f_\theta(x_i))$

$l: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ loss function such as
distance, squared distance, etc.

MLE • likelihood function $\sim \prod_{i=1}^n P_\theta(y_i | x_i)$

log-likelihood \sim

$$l(\theta) = \sum_{i=1}^n \log P_\theta(y_i | x_i)$$

MAP: posteriori functions $x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^{n \times 1}$ $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^{n \times N}$

$$\alpha(\theta) = P(\theta | X, Y)$$

When ERM = linear Regression

$$f_\theta(x) = x^T \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_D \end{pmatrix} + \theta_0 \quad l(y, \hat{y}) = (y - \hat{y})^2$$

MLE and MAP - an overview

We will first focus on Linear MLE, then MAP.

For nonlinear versions see Thm 3.13, Prop 3.14.

Select a model

$$(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}$$

Deterministic

• linear model $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}$

$$\theta = (\theta_0, \dots, \theta_D) \in \mathbb{R}^{D+1}$$

(w/ quadratic loss called linear regression)

$$f_\theta(x) = \theta_0 + x^T \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_D \end{pmatrix}$$

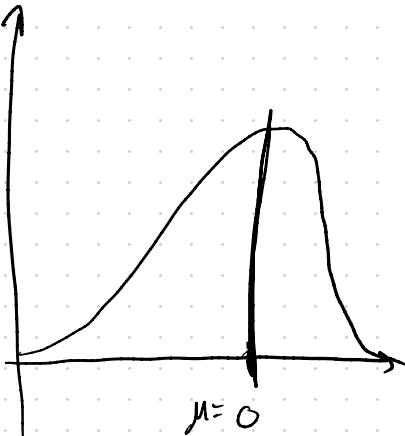
• linear probabilistic model

mean variance

$$P_\theta(y|x) = \Phi(y | f_\theta(x), \sigma^2)$$

density function of normal distribution

$$= \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(f_\theta(x) - y)^2}{2\sigma^2}}$$



] linear model + random noise

$$\Rightarrow y = f_{\theta}(x) + \varepsilon \quad \text{for } \varepsilon \sim N(0, \sigma^2)$$

Julia notation

$$f(x::\text{Number}, \theta) = \theta[z] * x + \theta[I]$$

$$f(X::\text{Vector}, \theta) = [f(x_i|\theta) \text{ for } x \in X]$$

Types must be specified in first input because function f has 2 definitions & Julia needs to know which one to use.

- Given a model, we need to "learn" parameter θ , i.e. choose optimal value of θ with some measure

Deterministic \circ Linear regression, Empirical Risk minimization version

$$\theta_* = \underset{\theta}{\operatorname{argmin}} R(\theta)$$

$$R(\theta) = \frac{1}{n} \sum_{i=1}^n l(y_i, f_{\theta}(x_i))$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$

$$\Sigma = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \in \mathbb{R}^{n \times (D+1)}$$

feature matrix

Theorem 3.8

When $\text{rank}(\Sigma) = r(\Sigma) = D+1$
then $\underset{\theta}{\operatorname{argmin}} R(\theta)$ has a unique solution

$$\theta^* = \Sigma^+ y.$$

$$\Sigma^+ = (\Sigma^T \Sigma)^{-1} \Sigma^T$$

Julia Note ~ in example we had

$$\theta_{\text{ERM}} = \Sigma \setminus y$$

\uparrow left-division

By Prop 1.8 (see HW 1 exercise)

- Since $\Sigma \in \mathbb{R}^{n \times (D+1)}$ has rank $D+1$, it is left-invertible, so can compute left division.

Thm 3.8 pt

$$R(\theta) = \frac{1}{n} \sum_{i=1}^n (f_\theta(x_i) - y_i)^2 = \frac{1}{n} \sum_{i=1}^n (x_i^T \theta + \theta_0 - y_i)^2$$

$$= \frac{1}{n} \sum_{i=1}^n (x_i^T \theta + \theta_0 - y_i)(x_i^T \theta + \theta_0 - y_i)$$

Notice

$$\Sigma \theta = \begin{bmatrix} & x_1 \\ & \vdots \\ & x_n \\ \hline n \times D+1 & \end{bmatrix} \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_D \\ \hline (D+1) \times 1 \end{bmatrix} = \theta_0 + x_i^T \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_D \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n \times 1$$

$$= \frac{1}{n} (\Sigma \theta - y)^T (\Sigma \theta - y)$$

$$(n \times n) \times (n \times 1)$$

Recall Pseudoinverse

$$\Sigma \theta = y \text{ solved by } \Sigma \theta^+ y = \tilde{\theta}$$

$$\tilde{\theta} = \text{Im}(\Sigma^T), \quad \Sigma \tilde{\theta} = y_0 \text{ and}$$

$$y_0 = \underset{w \in \text{Im}(\Sigma)}{\text{argmin}} \quad \boxed{\| w - y \|_2^2}$$

$$(w - y)^T (w - y)$$

Regression for (linear) probabilistic models

Same idea as before w/ randomized term,
no normalization by \sqrt{n} .

- Likelihood function (error from yesterday)

$$L(\theta) = \prod_{i=1}^n p_\theta(y_i | x_i)$$

- Log-likelihood function $l(\theta) = \sum_{i=1}^n \log(p_\theta(y_i | x_i))$

$$\log(L(\theta)) = l(\theta), \text{ minimizing } l(\theta) + L(\theta) \text{ equivalent since logarithm is non-decreasing.}$$

Fix σ^2 , $p_\theta(y_i | x_i) = \mathbb{P}(y_i | f_\theta(x_i), \sigma^2)$

Julia note $\mathbb{P}(x := \text{Number}, \theta) = \sigma * \text{randn}() + f(x, \sigma) \sim N(\mu, 1)$

Theorem
(3.12)

When $r(\Sigma) = D+1$, MLE is solved by $\hat{\theta}_{ML} = \Sigma^+ Y$

$$\boxed{\text{PF}} \quad \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(f_\theta(x)-y)^2}{2\sigma^2}} = p_\theta(y | f_\theta(x), \sigma^2)$$

$$\log p_\theta(y_i | f_\theta(x_i), \sigma^2) = \log \left(\frac{e^{-\frac{(f_\theta(x_i)-y_i)^2}{2\sigma^2}}}{\sqrt{2\pi} \sigma} \right)$$

$$> \log \left(e^{-\frac{(f_\theta(x_i)-y_i)^2}{2\sigma^2}} \right) - \log (\sqrt{2\pi} \sigma)$$

$$= -\frac{(f_\theta(x_i)-y_i)^2}{2\sigma^2} - \frac{1}{2} \log (2\pi \sigma^2)$$

So log-likelihood function is

$$l(\theta) = \sum_{i=1}^n \log (p_\theta(y_i | f_\theta(x_i), \sigma^2))$$

$$= -\sum_{i=1}^n \frac{1}{2} \log (2\pi \sigma^2) - \boxed{\sum_{i=1}^n \left(x_i^T \begin{pmatrix} \theta_0 \\ \theta_D \end{pmatrix} + \theta_0 - y_i \right)^2}$$

From linear term

$$= \underbrace{-\frac{n}{2} \log (2\pi \sigma^2)}_{\text{independent of } \theta} - \underbrace{\frac{1}{2\sigma^2} \| \Sigma \theta - y \|^2}_{\text{Unique Solution follows}}$$

from linear case.

Regularization \rightarrow used to prevent overfitting.

ERM tends towards overfitting
Deterministic linear \rightarrow Tikhonov regularization
 $\lambda \in \mathbb{R}$

$$l(y_i, f_\theta(x_i)) = (f_\theta(x_i) - y_i)^2 + \lambda \|\theta\|^2$$

$$\|\theta\|^2 = \theta_0^2 + \theta_1^2 + \dots + \theta_d^2$$

Corresponding regression problem called ridge regression

* See Theorem 3.11

Probabilistic linear \sim MLE tends towards overfitting.

Instead of adding $\lambda \|\theta\|^2$, we put a distribution on θ called a prior.

posterior function is then $P(\theta | X, Y)$

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^{n \times d} \quad Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$$

Want to view $X+Y$ as random variables.

[Bayes' Theorem] $P(\theta | X, Y) = P(\theta) \frac{P(Y | X, \theta)}{P(Y | X)}$

Taking logarithm

constant w.r.t. θ .

$$\log [P(\theta | x, y)] = \underbrace{\log P(\theta)}_{\text{distribution of } \theta} + \underbrace{\log P(y | x, \theta)}_{\text{MLE case}} - \underbrace{\log (y(x))}_{\text{"regularizing" term}}$$

* Special case when $D=1$ (Thm 3.15 for general)

$$\Theta = \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \underbrace{\begin{bmatrix} \sigma & \sigma \\ \sigma & \sigma \end{bmatrix}}_{\Sigma} \right)$$

$$\Rightarrow P(\theta) = \frac{1}{\sqrt{(2\pi)^2(\sigma^2)}} \exp \left(-\frac{1}{2\sigma^2} \|\theta\|^2 \right)$$

Also assume ϵ as in MLE case w/ independent variables

$$P(y | x, \theta) = \prod_{i=1}^n \phi(y_i | f_\theta(x_i), \sigma^2)$$

$$= \frac{1}{\sqrt{(2\sigma^2)^n}} \exp \left(-\frac{1}{2\sigma^2} \|y - \sum \theta_i x_i\|^2 \right)$$

\ominus_{MAP} maximizes $\log(P(\theta|x,y))$
 ~ Maximize by taking derivative w.r.t. θ .

$$\log[P(\theta|x,y)] = \log P(\theta) + \log P(y|x,\theta) - \log(Y/x)$$

$$= -\frac{1}{2\sigma^2} \|\theta\|^2 - \log(\sqrt{\pi\sigma^2})$$

$$-\frac{1}{2\sigma^2} \|y - \theta x\|^2 - \log((\sqrt{2\sigma^2})^n)$$

+ constant w.r.t. θ

Taking derivative, green goes away
 we have

$$\frac{d}{d\theta} \log(P(\theta|x,y)) = -\frac{1}{\sigma^2} \theta - \frac{1}{\sigma^2} \sum (y_i - \theta x_i)$$

$$\begin{aligned} \|\theta\|^2 &= \theta_0^2 + \theta_1^2 + \dots + \theta_n^2 \\ \frac{d}{d\theta} &= [2\theta_0, 2\theta_1, \dots, 2\theta_n] = 2\theta \end{aligned}$$

$$\begin{aligned} (y_1 - (\theta x_1))^2 \\ 2[y_1 - (\theta x_1)][1 \cdot x_1] \end{aligned}$$

Set derivative = 0

$$\theta = -\Sigma^T (\Sigma \theta - y)$$

$$(\Sigma^T \Sigma - \text{Id}) \theta = \Sigma^T y$$

This is
Theorem 3.5
in a special case

$\mu = 0$, $\Sigma = \sigma \text{Id}_2$

$\theta_{\text{MAP}} = (\Sigma^T \Sigma - \text{Id}_2)^{-1} \Sigma^T y$