

Parameter Estimation

Given samples $D = \{x^1, \dots, x^N\}$, form of the distribution $p(x) = \underline{f(x; \mathbf{w})}$
estimate values of the parameters w .

i.i.d assumption: each instance is independently and identically distributed.

Two methods:

1. Maximum likelihood estimation
2. Bayesian estimation

Aug 31 | Sep 1 2 3 4 5 $N=5$
0 1 1 0 0 $\equiv D$

$Y \sim -100 -90 -75 -100 -95$
 $p(x) \sim \text{Bernoulli}(x; w)$

$p(\text{rain on Friday, 8th Sep 2023}) \equiv 0.5$

$Y \equiv$ amount of rainfall in 1st
10 days of September.

$$P(Y) \sim N(\mu; \sigma^2)$$

$$f(Y; w) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(Y-\mu)^2}{2\sigma^2}}$$

$$\hat{\mu}_D = -\frac{463}{5}$$

$$\hat{\sigma}_D^2 = \frac{10}{10}$$

$$f(x; w) \equiv w^x (1-w)^{1-x}$$

$$\hat{w} = \frac{2}{5}$$

Maximum Likelihood Estimation (MLE)

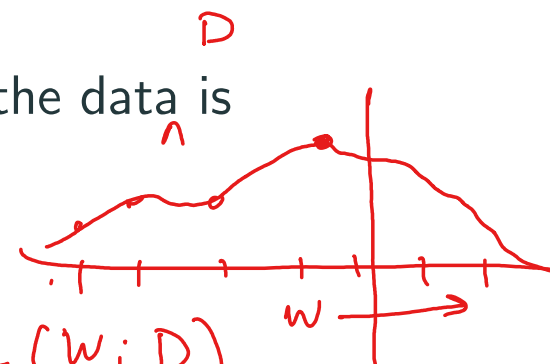
- Find the value \hat{w} for which the probability (likelihood) of the data is maximized.

- Likelihood of data $L(w; D)$

$$L(w; D) = \prod_{i=1}^N f(x^i; w)$$

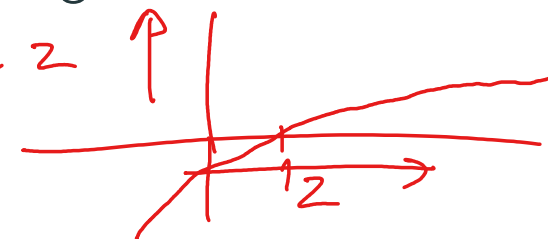
probability function given to us.

$$\hat{w} \equiv \underset{w}{\operatorname{argmax}} L(w; D)$$



- Maximizing log-likelihood of data is equivalent to maximizing likelihood.

$$\hat{w} \equiv \underset{w}{\operatorname{argmax}} \log \prod_{i=1}^N f(x^i; w) = \sum_{i=1}^N \log f(x^i; w)$$



Solving the MLE objective

$$\max_{\mathbf{w}} \underline{LL(\mathbf{w}; D)}$$

Apply numerical optimization algorithms...e.g. stochastic gradient ascent

$\mathbf{w}^0 \equiv$ initial value of parameters

for $t = 1$ to \dots

$x^i \sim$ sample one example from D .

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t + \eta \nabla_{\mathbf{w}} \log f(x^i; \mathbf{w})$$

\nearrow learning rate

For many simple distributions, the objective is concave in \mathbf{w} . Maxima of \mathbf{w} iff gradients w.r.t \mathbf{w} is zero. Eg: $\log f(x^i; \mathbf{w})$ is concave for Bernoulli, Gaussian, ...

If $\nabla_{\mathbf{w}} LL(\mathbf{w}^*; D) = 0$ then \mathbf{w}^* is global maxima:

$$\Leftrightarrow \nabla_{\mathbf{w}} \sum_{i=1}^N \log f(x^i; \mathbf{w}) = 0$$

Parameter estimation for Bernoulli distribution

MLE for Bernoulli: $f(x; w) = w^x (1-w)^{1-x}$ $x^i \in \{0, 1\}$

$$LL(w; D \equiv \{x^1, x^2, \dots, x^N\}) = \sum_{i=1}^N \log w^{x^i} (1-w)^{1-x^i}$$
$$= \sum_{i=1}^N x^i \log w + (1-x^i) \log (1-w)$$

$$\nabla_w LL(w; D) = \frac{\partial}{\partial w} \sum_{i=1}^N x^i \log w + (1-x^i) \log (1-w)$$
$$= \sum_{i=1}^N x^i \frac{1}{w} + (1-x^i) \cdot \frac{1(-1)}{1-w} = \frac{\sum_{i=1}^N x^i}{w} - \frac{\sum_{i=1}^N (1-x^i)}{1-w} = 0$$

$$\Rightarrow (1-w)n_1(D) - w n_0(D) = 0$$

$$\Rightarrow n_1(D) = w(n_0(D) + n_1(D)) = wN \Rightarrow \hat{w} = \frac{n_1(D)}{N}$$

MLE for Gaussian distribution

$$f(x^i; \mathbf{w} \equiv (\mu, \sigma^2)) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x^i - \mu)^2}{2\sigma^2}} \quad D \equiv \{x^1, x^2, \dots, x^N\} \quad x^i \in \mathbb{R}$$

$$LL(\mathbf{w}; D) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x^i - \mu)^2}{2\sigma^2}}$$

$$= \sum_{i=1}^N -\frac{(x^i - \mu)^2}{2\sigma^2} - \log \sigma \quad (-\log \sqrt{2\pi}) \text{ \textit{constant - ignore.}}$$

$$\nabla_{\mathbf{w}} LL(\mathbf{w}; D) = 0 = \begin{bmatrix} \frac{\partial}{\partial \mu} LL(\mathbf{w}; D) \\ \frac{\partial}{\partial \sigma} LL(\mathbf{w}; D) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\frac{\partial}{\partial \mu} \left[-\sum_{i=1}^N \frac{(x^i - \mu)^2}{2\sigma^2} - \log \sigma \right] \Rightarrow -\sum_{i=1}^N \frac{2(x^i - \mu)}{2\sigma^2} (-1) = 0 \Rightarrow \sum_{i=1}^N x^i = \sum_{i=1}^N \mu = N\mu$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^N (x^i - \hat{\mu})^2}{N}$$

MLE estimate -
was derived using $\hat{\mu}$ estimated from same data

As a result $\hat{\sigma}$ is known to have some "bias".

Statistical correction for the bias is obtained

using
$$\tilde{\sigma}^2 = \frac{\sum_{i=1}^N (x^i - \hat{\mu})^2}{N-1}.$$

MLE for Multinomial distribution

$$x \in \{1, 2, \dots, K\} \quad K \geq 2$$

x is also denoted as a "one-hot" vector

$$x^i = \begin{bmatrix} x_1^i \\ x_2^i \\ \vdots \\ x_K^i \end{bmatrix}$$

$$x = 3 \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Lagrangian multiplier

$$f(x; w = \{p_1, p_2, \dots, p_K\}) = \prod_{j=1}^K p_j^{x_j}$$

$$\sum_{j=1}^K p_j = 1 \quad p_j \geq 0$$

$$D = \{x^1, x^2, \dots, x^N\} \quad x^i \in \{0, 1\}^K$$

$$\max_{w = \{p_1, \dots, p_K\}} \sum_{i=1}^N \log \prod_{j=1}^K p_j^{x_j^i}$$

using Lagrangian multiplier:

$$\max_{p_1, p_2, \dots, p_K} \sum_{i=1}^N \sum_{j=1}^K x_j^i \log p_j + \lambda \left(\sum_{j=1}^K p_j - 1 \right)$$

$$\text{s.t. } p_1 + p_2 + \dots + p_K = 1$$

$$\nabla_{p_j} F(w; \lambda) = \sum_{i=1}^N x_j^i / p_j + \lambda \cdot 1 = 0 \Rightarrow p_j = \frac{\sum_{i=1}^N x_j^i}{N} = n_j(D)$$

$$\because \sum_j p_j = 1 \text{ we get } \sum_{j=1}^K \frac{n_j(D)}{(-\lambda)} = 1 \Rightarrow (-\lambda) = \sum_{j=1}^K n_j(D) = N$$

$$\hat{p}_j = \frac{n_j(D)}{N}$$

Example: states from which students in a class came
 $K = 30$; $\{1 = \text{Gujarat}, 2 = \text{TN}; 3 = \text{Mah}, 4 = \text{Odisha}, \dots\}$

We took three samples.

$$X^1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$X^2 = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$X^3 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

MLE estimates

$$\hat{p}_1 = \frac{1}{3} ; \hat{p}_2 = \frac{1}{3} ; \hat{p}_3 = 0 ; \hat{p}_4 = \frac{1}{3} , \hat{p}_5 = \hat{p}_6 = \dots = \hat{p}_{30} = 0$$