

Maximum Likelihood Estimation

Nipun Batra

August 13, 2023

IIT Gandhinagar

Agenda

Revision - Prior, Posterior, MLE, MAP

Distributions, IID

MLE

MLE for Bernoulli Distribution

MLE for Univariate Normal Distribution

MLE for Multivariate Normal Distribution

MLE for Gaussian Linear Regression

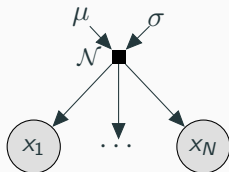
Revision - Prior, Posterior, MLE, MAP

Distributions, IID

Notebook

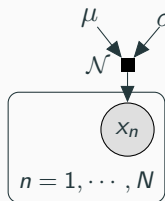
Graphical model

Assume model parameters are θ and data is D . We can write the joint probability distribution as:



Graphical model

Assume model parameters are θ and data is D . We can write the joint probability distribution as:



$$\begin{aligned}P(D|\theta) &= P(x_1, x_2, \dots, x_n|\theta) \\&= P(x_1|\theta) \cdot P(x_2|\theta) \cdot \dots \cdot P(x_n|\theta)\end{aligned}$$

MLE

We have three courses: C1, C2, C3. Assume no student takes more than one course. The scores of students in these courses are normally distributed with the following parameters:

- C1: $\mu_1 = 80, \sigma_1 = 10$
- C2: $\mu_2 = 70, \sigma_2 = 10$
- C3: $\mu_3 = 90, \sigma_3 = 5$

Pop Quiz

We have three courses: C1, C2, C3. Assume no student takes more than one course. The scores of students in these courses are normally distributed with the following parameters:

- C1: $\mu_1 = 80, \sigma_1 = 10$
- C2: $\mu_2 = 70, \sigma_2 = 10$
- C3: $\mu_3 = 90, \sigma_3 = 5$

I randomly pick up a student and ask them their marks. They say 82. Which course do you think they are from? To keep things simple, for now assume that all three courses have equal number of students.

Pop Quiz

We have three courses: C1, C2, C3. Assume no student takes more than one course. The scores of students in these courses are normally distributed with the following parameters:

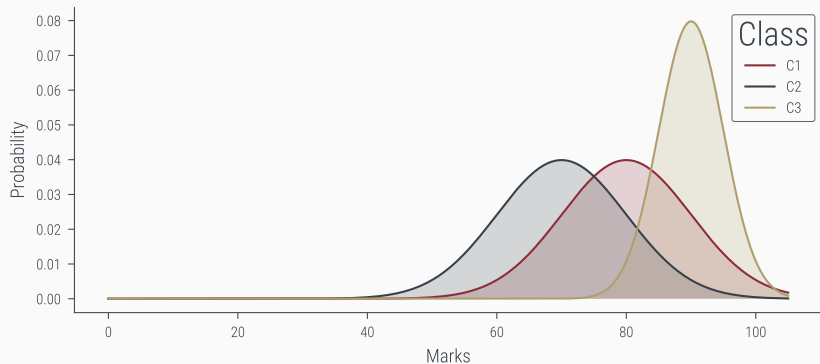
- C1: $\mu_1 = 80, \sigma_1 = 10$
- C2: $\mu_2 = 70, \sigma_2 = 10$
- C3: $\mu_3 = 90, \sigma_3 = 5$

I randomly pick up a student and ask them their marks. They say 82. Which course do you think they are from?

Most likely C1. But why?

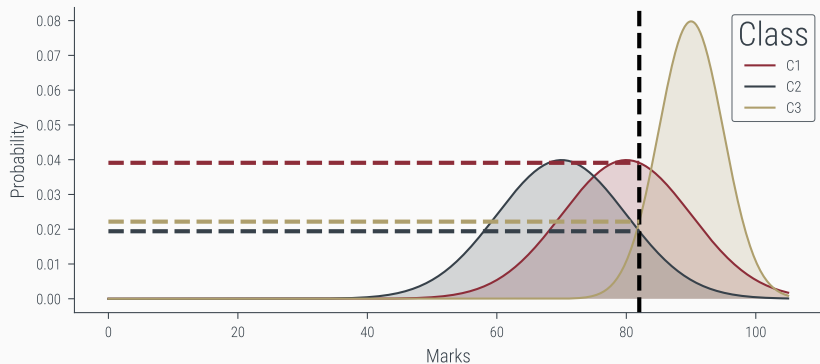
Pop Quiz

Let us plot the probability density functions of the three courses.



Pop Quiz

Let us plot the probability density functions of the three courses.



Notebook

Pop Quiz 2

Let us say we observed a value of 20. We know it came from a normal distribution with $\sigma = 1$. What is the most likely value of μ ?

Pop Quiz 2

Let us say we observed a value of 20. We know it came from a normal distribution with $\sigma = 1$. What is the most likely value of μ ?

20. But why?

Pop Quiz 2

Let us say we observed a value of 20. We know it came from a normal distribution with $\sigma = 1$. What is the most likely value of μ ?

20. But why?

Let us evaluate probability density function at 20 for different values of μ for $\sigma = 1$, i.e., $f(x = 20|\mu, \sigma = 1)$.

Pop Quiz 2

Let us say we observed a value of 20. We know it came from a normal distribution with $\sigma = 1$. What is the most likely value of μ ?

20. But why?

Let us evaluate probability density function at 20 for different values of μ for $\sigma = 1$, i.e., $f(x = 20|\mu, \sigma = 1)$.

Importantly, this is a function of μ and not x (which is fixed at 20).

Notebook

Pop Quiz 3

Let us now go back to our original problem. We have three courses: C1, C2, C3. Assume no student takes more than one course.

We ask two students their marks. The first student says 82 and the second student says 72. Which course do you think they are from? Assumption: Both are from the same course.

Let us create a table of probabilities for each course:

MLE for Bernoulli Distribution

MLE for Bernoulli Distribution

The probability mass function of a bernoulli distribution is given by:

$$f(x|\theta) = \theta^x(1 - \theta)^{(1-x)} \quad (1)$$

Let us assume we have a dataset $D = \{x_1, x_2, \dots, x_n\}$, where each x_i is an independent sample from the above distribution and $x_i \in \{0, 1\}$. We want to estimate the parameter θ from the data.

Our likelihood function is given by:

$$P(D|\theta) = \mathcal{L}(\theta) = \prod_{i=1}^n f(x_i|\theta) \quad (2)$$

Log Likelihood Function

Log-likelihood function:

$$\log \mathcal{L}(\theta) = \sum_{i=1}^n \log f(x_i|\theta) \quad (3)$$

Simplifying the above equation, we get:

$$\begin{aligned} \log \mathcal{L}(\theta) &= \sum_{i=1}^n \log f(x_i|\theta) \\ &= \sum_{i=1}^n \log \left(\theta^{x_i} (1 - \theta)^{(1-x_i)} \right) \\ &= \sum_{i=1}^n \left(\log (\theta^{x_i}) + \log \left((1 - \theta)^{(1-x_i)} \right) \right) \end{aligned}$$

$$\log \mathcal{L}(\theta) = \sum_{i=1}^n (x_i \log(\theta) + (1 - x_i) \log(1 - \theta))$$

Log Likelihood Function for Bernoulli Distribution

Log-likelihood function for bernoulli distributed data is:

$$\log \mathcal{L}(\theta) = \sum_{i=1}^n (x_i \log(\theta) + (1 - x_i) \log(1 - \theta))$$

Maximum Likelihood Estimate for θ

To find the MLE for θ , we differentiate the log-likelihood function with respect to θ and set it to zero:

$$\begin{aligned}\frac{\partial \log \mathcal{L}(\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \left(\sum_{i=1}^n (x_i \log(\theta) + (1 - x_i) \log(1 - \theta)) \right) \\&= \sum_{i=1}^n \left(\frac{\partial}{\partial \theta} (x_i \log(\theta)) + \frac{\partial}{\partial \theta} ((1 - x_i) \log(1 - \theta)) \right) \\&= \sum_{i=1}^n \left(x_i \frac{\partial}{\partial \theta} \log(\theta) + (1 - x_i) \frac{\partial}{\partial \theta} \log(1 - \theta) \right) \\&= \sum_{i=1}^n \left(\frac{x_i}{\theta} - \frac{(1 - x_i)}{1 - \theta} \right) = 0\end{aligned}$$

$$\begin{aligned}
\frac{\partial \log \mathcal{L}(\theta)}{\partial \theta} &= \sum_{i=1}^n \left(\frac{x_i(1 - \theta) - \theta(1 - x_i)}{\theta(1 - \theta)} \right) = 0 \\
&= \sum_{i=1}^n \left(\frac{x_i - x_i\theta - \theta + \theta x_i}{\theta(1 - \theta)} \right) \\
&= \sum_{i=1}^n \left(\frac{x_i - \theta}{\theta(1 - \theta)} \right) \\
&= \sum_{i=1}^n (x_i - \theta) = 0 \\
&= \sum_{i=1}^n x_i - \sum_{i=1}^n \theta = 0 \\
&= \sum_{i=1}^n x_i - n\theta = 0
\end{aligned}$$

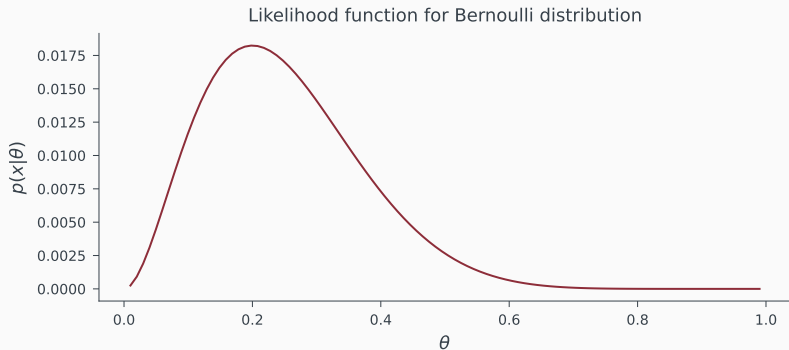
$$\theta = \frac{\sum_{i=1}^n x_i}{n}$$

Maximum Likelihood Estimate for θ

MLE of θ , denoted as $\hat{\theta}_{\text{MLE}}$, is given by:

$$\hat{\theta}_{\text{MLE}} = \frac{\sum_{i=1}^n x_i}{n}$$

For example if we have a Bernoulli Distribution with $\theta = 0.2$, the likelihood, $P(D|\theta)$ is given below:



MLE for Univariate Normal Distribution

Univariate Normal Distribution

The probability density function of a univariate normal distribution is given by:

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \quad (4)$$

Let us assume we have a dataset $D = \{x_1, x_2, \dots, x_n\}$, where each x_i is an independent sample from the above distribution. We want to estimate the parameters $\theta = \{\mu, \sigma\}$ from the data.

Our likelihood function is given by:

$$P(D|\theta) = \mathcal{L}(\mu, \sigma^2) = \prod_{i=1}^n f(x_i|\mu, \sigma^2) \quad (5)$$

Log Likelihood Function

Log-likelihood function:

$$\log \mathcal{L}(\mu, \sigma^2) = \sum_{i=1}^n \log f(x_i | \mu, \sigma^2) \quad (6)$$

Simplifying the above equation, we get:

$$\begin{aligned} \log \mathcal{L}(\mu, \sigma^2) &= \sum_{i=1}^n \log f(x_i | \mu, \sigma^2) \\ &= \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(x_i - \mu)^2}{2\sigma^2} \right) \right) \\ &= \sum_{i=1}^n \left(\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \log \left(\exp \left(-\frac{(x_i - \mu)^2}{2\sigma^2} \right) \right) \right) \end{aligned}$$

$$\begin{aligned}
 \log \mathcal{L}(\mu, \sigma^2) &= \sum_{i=1}^n \left(\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{(x_i - \mu)^2}{2\sigma^2} \right) \\
 &= \sum_{i=1}^n \left(-\frac{1}{2} \log(2\pi\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2} \right) \\
 &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2
 \end{aligned}$$

Log Likelihood Function for Univariate Normal Distribution

Log-likelihood function for normally distributed data is:

$$\log \mathcal{L}(\mu, \sigma^2) = -\frac{n}{2} \log(2\pi) - n \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

Maximum Likelihood Estimate for μ

To find the MLE for μ , we differentiate the log-likelihood function with respect to μ and set it to zero:

$$\frac{\partial \log \mathcal{L}(\mu, \sigma^2)}{\partial \mu} = \frac{\partial}{\partial \mu} \left(-\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right) = 0$$
$$\frac{\partial}{\partial \mu} \left(\sum_{i=1}^n (x_i - \mu)^2 \right) = 0$$

Maximum Likelihood Estimate for μ

MLE of μ , denoted as $\hat{\mu}_{\text{MLE}}$, is given by:

$$\hat{\mu}_{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n x_i$$

MLE for σ for normally distributed data

Recall that the log-likelihood function is given by:

$$\log \mathcal{L}(\mu, \sigma^2) = \sum_{i=1}^n \log f(x_i | \mu, \sigma^2) \quad (7)$$

Let us find the maximum likelihood estimate of σ^2 now. We can do this by taking the derivative of the log-likelihood function with respect to σ^2 and equating it to zero.

$$\frac{\partial \log \mathcal{L}(\mu, \sigma^2)}{\partial \sigma^2} = \sum_{i=1}^n \frac{\partial \log f(x_i | \mu, \sigma^2)}{\partial \sigma^2} = 0 \quad (8)$$

MLE for σ for normally distributed data

Log Likelihood Function for Univariate Normal Distribution

Log-likelihood function for normally distributed data is:

$$\log \mathcal{L}(\mu, \sigma^2) = -\frac{n}{2} \log(2\pi) - n \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

Now, we can differentiate the log-likelihood function with respect to σ and equate it to zero.

MLE for σ for normally distributed data

$$\frac{\partial}{\partial \sigma} \log \mathcal{L}(\mu, \sigma^2) = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

Multiplying through by σ^3 , we have:

$$-n\sigma^2 + \sum_{i=1}^n (x_i - \mu)^2 = 0$$

Maximum Likelihood Estimate for σ^2

MLE of σ^2 , denoted as $\hat{\sigma}_{\text{MLE}}^2$, is given by:

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

MLE for Multivariate Normal Distribution

MLE for Multivariate Normal Distribution

The probability density function of a multivariate normal distribution is given by:

$$f(x|\mu, \Sigma) = (2\pi)^{-\frac{k}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} \quad (9)$$

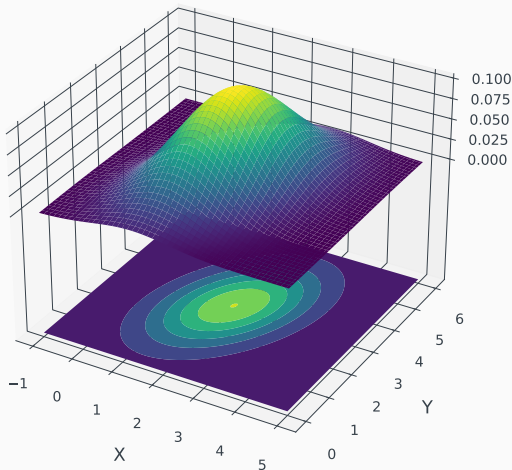
Let us assume we have a dataset $D = \{x_1, x_2, \dots, x_n\}$, where each x_i is an independent sample from the above distribution. We want to estimate the parameters $\theta = \mu, \sigma$ from the data.

Our likelihood function is given by:

$$P(D|\theta) = \mathcal{L}(\mu, \Sigma) = \prod_{i=1}^n f(x_i|\mu, \Sigma) \quad (10)$$

For example: A bivariate Normal distribution can be visualized as given below:

Covariance Matrix:

$$\begin{bmatrix} 1. & 0.5 \\ 0.5 & 2. \end{bmatrix}$$


Log Likelihood Function

Log-likelihood function:

$$\log \mathcal{L}(\mu, \Sigma) = \sum_{i=1}^n \log f(x_i | \mu, \Sigma) \quad (11)$$

Simplifying the above equation, we get:

$$\begin{aligned} \log \mathcal{L}(\mu, \Sigma) &= \sum_{i=1}^n \log f(x_i | \mu, \Sigma) \\ &= \sum_{i=1}^n \log \left((2\pi)^{-\frac{k}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp^{-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1}(x_i - \mu)} \right) \\ &= \sum_{i=1}^n \log((2\pi)^{-\frac{k}{2}}) + \sum_{i=1}^n \log(\det(\Sigma)^{-\frac{1}{2}}) + \\ &\quad \sum_{i=1}^n \log(\exp^{-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1}(x_i - \mu)}) \end{aligned}$$

Continuing, we get:

$$= -\frac{kn}{2} \log(2\pi) - \frac{n}{2} \log(\Sigma) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

Log Likelihood Function for Multivariate Normal Distribution

Log-likelihood function for multivariate normally distributed data is:

$$-\frac{kn}{2} \log(2\pi) - \frac{n}{2} \log(\Sigma) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

Maximum Likelihood Estimate for μ

To find the MLE for μ , we differentiate the log-likelihood function with respect to μ and set it to zero:

$$\begin{aligned} &= \frac{\partial}{\partial \mu} \left(-\frac{kn}{2} \log(2\pi) - \frac{n}{2} \log(\Sigma) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) \\ &= \frac{\partial}{\partial \mu} \left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) \\ &= -\frac{1}{2} \sum_{i=1}^n \left(\Sigma^{-1} (x_i - \mu) + (x_i - \mu)^T \Sigma^{-1} \right) = 0 \\ &= -\frac{1}{2} \sum_{i=1}^n 2 \Sigma^{-1} (x_i - \mu) = 0 \\ &\quad \text{as } (x_i - \mu)^T \Sigma^{-1} = \Sigma^{-1} (x_i - \mu) \end{aligned}$$

$$\begin{aligned} &= \Sigma^{-1} \sum_{i=1}^n (x_i - \mu) = 0 \\ &= \sum_{i=1}^n (x_i) - n\mu = 0 \\ \mu &= \frac{\sum_{i=1}^n x_i}{n} \end{aligned}$$

Maximum Likelihood Estimate for μ

MLE of μ , denoted as $\hat{\mu}_{\text{MLE}}$, is given by:

$$\hat{\mu}_{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n x_i$$

MLE for Σ for multivariate normally distributed data

Recall that the log-likelihood function is given by:

$$\log \mathcal{L}(\mu, \Sigma) = \sum_{i=1}^n \log f(x_i | \mu, \Sigma) \quad (12)$$

Let us find the maximum likelihood estimate of Σ now. We can do this by taking the derivative of the log-likelihood function with respect to Σ and equating it to zero.

$$\frac{\partial \log \mathcal{L}(\mu, \Sigma)}{\partial \Sigma} = \sum_{i=1}^n \frac{\partial \log f(x_i | \mu, \Sigma)}{\partial \Sigma} = 0 \quad (13)$$

After differentiating and simplifying, we get:

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

Maximum Likelihood Estimate for Σ

MLE of Σ , denoted as $\hat{\Sigma}_{\text{MLE}}$, is given by:

$$\hat{\Sigma}_{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

MLE for Gaussian Linear Regression

MLE for Gaussian Linear Regression

Let us assume we have a dataset

$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, where $x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$. Also

$Y \sim N(w^T x, \sigma^2)$, Y_1, \dots, Y_n are independent i.e.

$Y_i \sim N(w^T x_i, \sigma^2)$. Assume σ^2 is known. We want to estimate the parameters $\theta = W$ from the data.

Our likelihood function (Normal distribution) is given by:

$$P(D|\theta) = P(y_1, \dots, y_n | x_1, \dots, x_n, \theta) = \prod_{i=1}^n p(y_i | x_i, \theta) \quad (14)$$

$$P(D|\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp^{-\frac{1}{2\sigma^2} (y_i - w^T x_i)^2} \quad (15)$$

Continuing, we get:

$$\begin{aligned} P(D|\theta) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp^{-\frac{1}{2\sigma^2}(y_i - w^T x_i)^2} \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2} \end{aligned}$$

Since $w^T x_i$ is a scalar, $w^T x_i = x_i^T w$ (transpose of a scalar is the scalar itself):

$$\sum_{i=1}^n (y_i - w^T x_i)^2 = (y - Xw)^T (y - Xw) = \|y - Xw\|^2$$

$$P(D|\theta) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp^{-\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)}$$

To maximize $P(D|\theta)$, we need to minimize $(y - Xw)^T(y - Xw)$.

$$\begin{aligned}\mathcal{L}(w) &= (y - Xw)^T(y - Xw) \\ &= y^T y - 2y^T Xw + w^T X^T Xw\end{aligned}$$

Since w is a matrix, we follow these two rules:

$$\frac{\partial(A\theta)}{\partial\theta} = A^T \tag{16}$$

$$Z = X^T X, \frac{\partial(\theta^T Z \theta)}{\partial\theta} = 2Z^T \theta \tag{17}$$

Differentiating with respect to w , we get:

$$-2X^T y + 2X^T Xw = 0$$

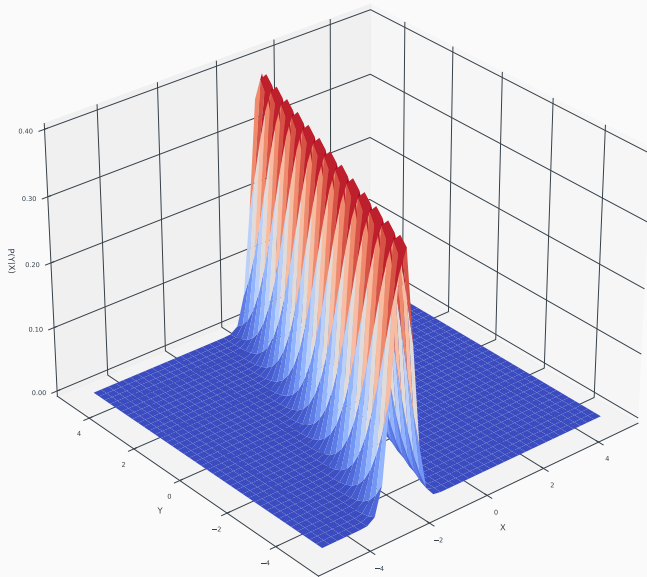
Thus, we get:

$$w = (X^T X)^{-1} X^T y \quad (18)$$

Maximum Likelihood Estimate for w

MLE of w , denoted as \hat{w}_{MLE} , is given by:

$$\hat{w}_{\text{MLE}} = (X^T X)^{-1} X^T y$$



Random variable and Random sample

Random variable: $X : \Omega \rightarrow \mathbb{R}$ is a function from the sample space to the real line.

Random sample: Collection of n independent and identically distributed (i.i.d.) random variables $X_1, X_2, X_3, \dots, X_n$. A group of experiments constitutes a sample.

For example:

Random variable: Y (possible outcomes 1 to 6)

Random sample: 4,2,6 (outcomes of three consecutive die tosses)

Bias of an Estimator

The bias of an estimator $\hat{\theta}$ of a parameter θ is defined as:

$$\text{Bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta$$

where $\mathbb{E}(\hat{\theta})$ is the expected value of the estimator $\hat{\theta}$.

- An estimator is said to be unbiased if $\text{Bias}(\hat{\theta}) = 0$.
- An estimator is said to be biased if $\text{Bias}(\hat{\theta}) \neq 0$.

Question: What is the expectation of $\hat{\mu}_{MLE}$ calculated over? What is the source of randomness?

If X_i 's are normally distributed random variables with mean μ and variance σ^2 respectively, then $E(X_i) = \mu$ and $Var(X_i) = \sigma^2$.

Recall that if an estimator $\hat{\theta}$ of a parameter θ is unbiased then:
 $\mathbb{E}(\hat{\theta}) = \theta$.

$$\begin{aligned}\mathbb{E}(\hat{\mu}_{MLE}) &= \mathbb{E}(\bar{X}) \\ &= \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) \\ &= \frac{1}{n}(n\mu) = \mu\end{aligned}$$

Estimator $\hat{\mu}_{MLE}$ is unbiased

$$\mathbb{E}(\hat{\mu}_{MLE}) = \mu$$

The MLE of σ^2 is given by

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x}x_i + \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x} \frac{1}{n} \sum_{i=1}^n x_i + \bar{x}^2 \\ &= \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2\end{aligned}$$

$$\begin{aligned}\mathbb{E}(\hat{\sigma}^2) &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 \right] = \left[\frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i^2) \right] - \mathbb{E}(\bar{X}^2) \\ &= \frac{1}{n} \sum_{i=1}^n (\sigma^2 + \mu^2) - \left(\frac{\sigma^2}{n} + \mu^2 \right) \\ &= \frac{1}{n} (n\sigma^2 + n\mu^2) - \frac{\sigma^2}{n} - \mu^2 \\ &= \sigma^2 - \frac{\sigma^2}{n} = \frac{n\sigma^2 - \sigma^2}{n} = \frac{(n-1)\sigma^2}{n}\end{aligned}$$

Estimator $\hat{\sigma}_{MLE}$ is biased

$$\mathbb{E}(\hat{\sigma}_{MLE}) = \frac{(n-1)\sigma^2}{n}$$

