



SWAYAM & NPTEL COURSE ON

Mathematics for Machine Learning

by

Prof. Debjani Chakraborty

DEPARTMENT OF MATHEMATICS
IIT Kharagpur

Module 8.1

Lecture 36: MLE for Discrete Probability Distribution

NPTEL



Concepts Covered

Maximum Likelihood Estimation for some discrete probability distributions

NPTEL



MLE for Binomial distribution parameter

Suppose that the random variables $X = \{X_1, X_2, \dots, X_n\}$ form a *i.i.d.* random sample from a population of Binomial distribution $\text{Bin}(k|n, \theta)$. We would like to find θ_{MLE} for θ .

$$\begin{aligned}\text{Bin}(k|n, \theta) \\ = \binom{n}{k} \theta^k (1 - \theta)^{n-k}\end{aligned}$$

if X is discrete random variable, $f(x|\theta)$ is point mass function, then for every observed random sample x_1, x_2, \dots, x_n , we define joint probability or the **Likelihood Function**

$$f(x_1, x_2, \dots, x_n | \theta) = f(x_1 | \theta) \cdots f(x_n | \theta)$$

The likelihood function is

$$L(\theta) = f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n \binom{n}{x_i} \theta^{x_i} (1 - \theta)^{(n-x_i)}$$



MLE for Binomial distribution parameter

Suppose that the random variables $X = \{X_1, X_2, \dots, X_n\}$ form a *i.i.d.* random sample from a population of Binomial distribution $Bin(k|n, \theta)$. We would like to find θ_{MLE} for θ .

The likelihood function is

$$L(\theta) = f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n \binom{n}{x_i} \theta^{x_i} (1 - \theta)^{(n-x_i)}$$

The log-likelihood function is

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



MLE for Binomial distribution parameter

Suppose that the random variables $X = \{X_1, X_2, \dots, X_n\}$ form a *i.i.d.* random sample from a population of Binomial distribution $Bin(k|n, \theta)$. We would like to find θ_{MLE} for θ .

The log-likelihood function is:

$$\log \sum_{i=1}^n \left(\log \binom{n}{x_i} + x_i \log \theta + (1-x_i) \log (1-\theta) \right)$$

Taking derivative of log-likelihood function and equating to zero –

$$\frac{d L(\theta)}{d\theta} = \sum_{i=1}^n \left(\frac{x_i}{\theta} - \frac{1-x_i}{1-\theta} \right) = 0$$

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



MLE for Poisson distribution parameter

Suppose that the random variables $X = \{X_1, X_2, \dots, X_n\}$ form a *i.i.d.* random sample from a population of Poisson distribution $\text{Poisson}(k|\lambda)$. We would like to find λ_{MLE} for λ .

$$\text{Poisson}(k|\lambda) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad \lambda > 0$$

if X is discrete random variable, $f(x|\lambda)$ is point mass function, then for every observed random sample x_1, x_2, \dots, x_n , we define joint probability or the **Likelihood Function**

$$f(x_1, x_2, \dots, x_n | \lambda) = f(x_1 | \lambda) \cdots f(x_n | \lambda)$$

The likelihood function is

$$L(\lambda) = f(x_1, x_2, \dots, x_n | \lambda) = \prod_{i=1}^n \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$



MLE for Poisson distribution parameter

Suppose that the random variables $X = \{X_1, X_2, \dots, X_n\}$ form a *i.i.d.* random sample from a population of Binomial distribution $\text{Bin}(k|n, \theta)$. We would like to find θ_{MLE} for θ .

The likelihood function is

$$L(\lambda) = f(x_1, x_2, \dots, x_n | \lambda) = \prod_{i=1}^n \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} = \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}$$

The log-likelihood function is

$$\begin{aligned} \log L(\lambda) &= \log \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!} \\ &= -n\lambda + \log \lambda \sum_{i=1}^n x_i - \sum_{i=1}^n \log x_i! \end{aligned}$$



MLE for Binomial distribution parameter

Suppose that the random variables $X = \{X_1, X_2, \dots, X_n\}$ form a *i.i.d.* random sample from a population of Binomial distribution $Bin(k|n, \theta)$. We would like to find θ_{MLE} for θ .

The log-likelihood function is:

$$-n\lambda + \log \lambda \sum_{i=1}^n x_i - \sum_{i=1}^n \log x_i!$$

Taking derivative of log-likelihood function and equating to zero –

$$\frac{d L(\lambda)}{d\lambda} = -n + \frac{1}{\lambda} \sum_{i=1}^n x_i = 0$$

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



Example: Let a population has been considered with the following probability distribution

$$f(x|\theta) = \begin{cases} \frac{1}{4} & \text{if } x = 0 \\ \frac{1 + 2\theta}{4} & \text{if } x = 1 \\ \frac{1}{4} & \text{if } x = 2 \\ \frac{1 - 2\theta}{4} & \text{if } x = 3 \end{cases}$$

A sample with 10 observations has been considered, which is as follows:

(3, 0, 2, 1, 3, 2, 1, 0, 2, 1)

Find θ_{MLE} for θ .



Example: Let a population has been considered with the following probability distribution

Sample: (3, 0, 2, 1, 3, 2, 1, 0, 2, 1)

The likelihood function is

$$L(\theta) = P(x = 3)P(x = 0)P(x = 2)P(x = 1)P(x = 3)P(x = 2)P(x = 1)P(x = 0) \\ P(x = 2)P(x = 1)$$

$$= \left(\frac{1}{4}\right)^2 \left(\frac{1+2\theta}{4}\right)^3 \left(\frac{1}{4}\right)^3 \left(\frac{1-2\theta}{4}\right)^2$$

$$f(x|\theta) = \begin{cases} \frac{1}{4} & \text{if } x = 0 \\ \frac{1+2\theta}{4} & \text{if } x = 1 \\ \frac{1}{4} & \text{if } x = 2 \\ \frac{1-2\theta}{4} & \text{if } x = 3 \end{cases}$$

The log-likelihood function is:

$$\log L(\theta) = \{2 \log 1 - 2 \log 4\} + \{3 \log (1+2\theta) - 3 \log 4\} + \{3 \log 1 - 3 \log 4\} + \\ \{2 \log (1-2\theta) - 2 \log 4\}$$



Example: Let a population has been considered with the following probability distribution

Sample: (3, 0, 2, 1, 3, 2, 1, 0, 2, 1)

The log-likelihood function is:

$\text{Log } L(\theta)$

$$\begin{aligned} &= \{2 \log 1 - 2 \log 4\} + \{3 \log (1 + 2\theta) - 3 \log 4\} + \{3 \log 1 \\ &\quad - 3 \log 4\} + \{2 \log (1 - 2\theta) - 2 \log 4\} \end{aligned}$$

Taking derivative of log-likelihood function and equating to zero –

$$\frac{d L(\theta)}{d\theta} = \frac{6}{1 + 2\theta} - \frac{4}{1 - 2\theta} = 0$$

$$\theta_{MLE} = .1$$

$$f(x|\theta) = \begin{cases} \frac{1}{4} & \text{if } x = 0 \\ \frac{1 + 2\theta}{4} & \text{if } x = 1 \\ \frac{1}{4} & \text{if } x = 2 \\ \frac{1 - 2\theta}{4} & \text{if } x = 3 \end{cases}$$



Example: Let a population has been considered with the following probability distribution

$$f(x|\theta) = \begin{cases} \frac{1+2\theta}{4} & \text{if } x = 0 \\ \frac{1-\theta}{4} & \text{if } x = 1 \\ \frac{1+\theta}{4} & \text{if } x = 2 \\ \frac{1-2\theta}{4} & \text{if } x = 3 \end{cases}$$

A sample with 10 observations has been considered, which is as follows:

(3, 0, 2, 1, 3, 2, 1, 0, 2, 1)

Find θ_{MLE} for θ .



Example: Let's assume our observed data points are: $x = [1,3,2,4,5,2,1,3,2,4]$ from a Poisson distribution.

The likelihood function for the Poisson distribution is:

$$L(\lambda) = \prod_{i=1}^N \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

Taking the logarithm of the likelihood function, we get the log-likelihood function:

$$\log L(\lambda) = -N\lambda + \sum_{i=1}^N x_i \log \lambda - \sum_{i=1}^N \log(x_i!)$$

To find the maximum likelihood estimate for λ , we differentiate the log-likelihood function with respect to λ , set it equal to zero, and solve for λ .

$$\begin{aligned}\frac{d}{d\lambda} \log L(\lambda) &= -N + \frac{1}{\lambda} \sum_{i=1}^N x_i = 0 \\ \Rightarrow \quad \lambda &= \frac{1}{N} \sum_{i=1}^N x_i\end{aligned}$$

Now, let's calculate the MLE of using our sample data:

$$\begin{aligned}\lambda_{MLE} \\ &= \frac{1}{10} \sum_{i=1}^{10} x_i = \frac{1}{10} (1 + 3 + 2 + 4 + 5 + 2 + 1 + 3 + 2 + 4) = \frac{27}{10} = 2.7\end{aligned}$$



- ✓ Many practical problems we examine appropriate conditions and select the probability distribution
- ✓ A range of discrete probability distributions are considered here

Different continuous probability distributions are in the next.....

References

- ✓ Linear Algebra and Learning from Data (2019), Gilbert Strang, Wellesley Cambridge Press
- ✓ Machine Learning: A Probabilistic Perspective, Kevin P. Murphy (MIT Press), 2021 edition
- ✓ Deisenroth MP, Faisal AA, Ong CS. Mathematics for machine learning. Cambridge University Press; 2020 Apr 23.





NPTEL ONLINE CERTIFICATION COURSES
IIT KHARGPUR



SWAYAM & NPTEL COURSE ON

Mathematics for Machine Learning

by

Prof. Debjani Chakraborty

DEPARTMENT OF MATHEMATICS
IIT Kharagpur

Module 8.2

Lecture 37: MLE for Continuous Probability Distribution

NPTEL



Concepts Covered

Maximum Likelihood Estimation for some continuous probability distributions

NPTEL



MLE for Exponential distribution parameter

Maximum likelihood estimation can be applied to a vector valued parameter. For a simple random sample of n exponential random variables, we would like to find λ_{MLE} for λ .

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

if X is continuous random variable, $f(x|\lambda)$ is probability density function, then for every observed random sample x_1, x_2, \dots, x_n , we define joint probability or the **Likelihood Function**

$$f(x_1, x_2, \dots, x_n | \lambda) = f(x_1 | \lambda) \cdots f(x_n | \lambda)$$

The likelihood function is

$$L(\lambda) = f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n \lambda e^{-\lambda x_i}$$

$$\begin{aligned}\hat{\lambda}_{MLE} &= \underset{\lambda}{\operatorname{argmax}} L(\lambda) \\ &= \underset{\lambda}{\operatorname{argmax}} \log L(\lambda)\end{aligned}$$



MLE for Exponential distribution parameter

Maximum likelihood estimation can be applied to a vector valued parameter. For a simple random sample of n exponential random variables, we would like to find λ_{MLE} for λ .

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

The likelihood function is

$$L(\lambda) = f(x_1, x_2, \dots, x_n | \lambda) = \prod_{i=1}^n \lambda e^{-\lambda x_i}$$

The log-likelihood function is

$$\log L(\lambda) = \underline{n \log \lambda} - \lambda \sum_{i=1}^n x_i$$



MLE for Exponential distribution parameter

Maximum likelihood estimation can be applied to a vector valued parameter. For a simple random sample of n exponential random variables, we would like to find λ_{MLE} for λ .

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

The log-likelihood function is

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

Taking derivative of log-likelihood function w.r.t. μ and equating to zero –

$$\frac{d \log L(\lambda)}{d \lambda} = \frac{n}{\lambda} - \sum_{i=1}^n x_i = 0 \quad \frac{n}{\lambda} = \sum_{i=1}^n x_i$$

$$\lambda_{MLE} = \frac{1}{\frac{1}{n} \sum_{i=1}^n x_i} = \frac{1}{\bar{x}}$$

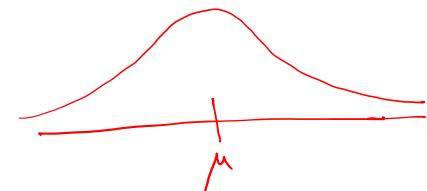


MLE for Normal distribution parameters

Maximum likelihood estimation can be applied to a vector valued parameter. For a simple random sample of n normal random variables We would like to find μ_{MLE} for μ and σ_{MLE} for σ .

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

where $-\infty < x, \mu < \infty, \sigma > 0$



if X is continuous random variable, $f(x|\mu, \sigma)$ is probability density function, then for every observed random sample x_1, x_2, \dots, x_n , we define joint probability or the **Likelihood Function**

$$f(x_1, x_2, \dots, x_n | \mu, \sigma) = f(x_1 | \mu, \sigma) \cdots f(x_n | \mu, \sigma)$$

$$\mu_{MLE} = \underset{\mu}{\operatorname{argmax}} L(\mu, \sigma)$$
$$\sigma_{MLE}$$

The likelihood function is

$$L(\theta) = f(x_1, x_2, \dots, x_n | \mu, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i - \mu)^2 / 2\sigma^2}$$



MLE for Normal distribution parameters

Maximum likelihood estimation can be applied to a vector valued parameter. For a simple random sample of n normal random variables We would like to find μ_{MLE} for μ and σ_{MLE} for σ .

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

where $-\infty < x, \mu < \infty, \sigma > 0$

The likelihood function is

$$L(\mu, \sigma) = f(x_1, x_2, \dots, x_n | \mu, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i-\mu)^2/2\sigma^2}$$

The log-likelihood function is

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



MLE for Normal distribution parameters

Maximum likelihood estimation can be applied to a vector valued parameter. For a simple random sample of n normal random variables We would like to find μ_{MLE} for μ and σ_{MLE} for σ .

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

where $-\infty < x, \mu < \infty, \sigma > 0$

The log-likelihood function is

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

Taking derivative of log-likelihood function w.r.t. μ and equating to zero –

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

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

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



MLE for Normal distribution parameters

Maximum likelihood estimation can be applied to a vector valued parameter. For a simple random sample of n normal random variables We would like to find μ_{MLE} for μ and σ_{MLE} for σ .

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

where $-\infty < x, \mu < \infty, \sigma > 0$

The log-likelihood function is

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

partial

Taking derivative of log-likelihood function w.r.t. σ^2 and equating to zero –

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

$$\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$



Example: Now, let's take a numerical example with some sample data. Suppose we have the following 10 observations drawn from a Normal distribution: $X=\{2.5, 3.1, 4.0, 2.8, 3.7, 8.9, 4.3, 10.3, 1.1, 5.6\}$. Find MLE of μ and σ^2 .

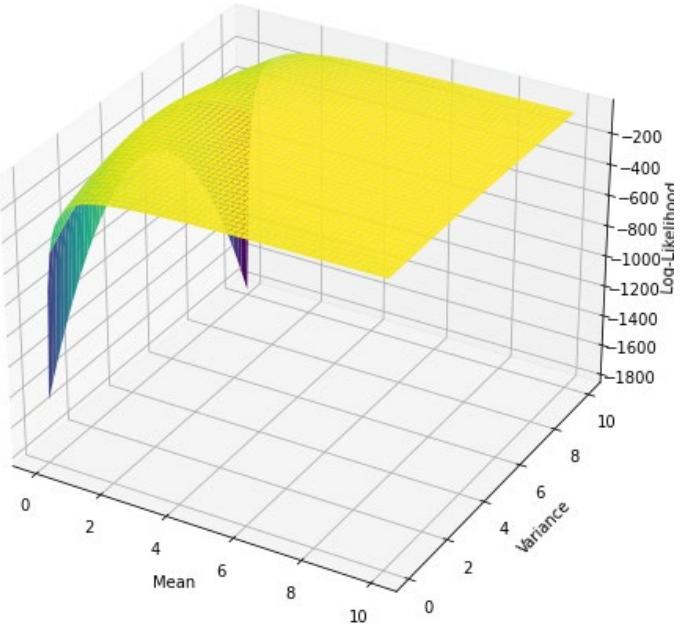
$$\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

The log-likelihood function is

$$\log L(\mu, \sigma) = -\frac{10}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

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

Likelihood Function for Normal Distribution



NPTEL

$$\mu_{MLE} = 4.63.$$

$$\sigma_{MLE}^2 = 7.54.$$



Example: Now, let's take a numerical example with some sample data. Suppose we have the following 10 observations drawn from the following distribution: $X=\{2.5, 3.1, 4.0, 2.8, 3.7, 8.9, 4.3, 10.3, 1.1, 5.6\}$. Find MLE of θ . Find Mean and standard deviation also.

$$f(x|\theta) = (1 - 2\theta) + 4\theta x, \quad 0 < x < 1, \quad -1 < \theta < 1$$

$$\begin{aligned} E(x) &= \int_0^1 x f(x) dx \\ &= \int_0^1 \{(1 - 2\theta)x + 4\theta x^2\} dx \\ &= \left[\frac{(1 - 2\theta)x^2}{2} \right]_0^1 + \left[\frac{4\theta x^3}{3} \right]_0^1 \\ &= \frac{1 - 2\theta}{2} + \frac{4\theta}{3} = \left(\frac{1}{2} + \frac{\theta}{3} \right) \\ \text{Var}(x) &= \int_0^1 \{(x - E(x))^2 f(x) dx \end{aligned}$$

$$\begin{aligned} \int_0^1 f(x) dx &= \int_0^1 \{(1 - 2\theta)x + 4\theta x^2\} dx \\ &= \left[(1 - 2\theta)x \right]_0^1 + \left[2\theta x^2 \right]_0^1 \\ &= 1 - 2\theta + 2\theta = 1 \end{aligned}$$



Example: Now, let's take a numerical example with some sample data. Suppose we have the following 10 observations drawn from the following distribution: $X=\{2.5, 3.1, 4.0, 2.8, 3.7, 8.9, 4.3, 10.3, 1.1, 5.6\}$.
Find MLE of θ . Find Mean and standard deviation also.

$$f(x|\theta) = (1 - 2\theta) + 4\theta x, \quad 0 < x < 1, \quad -1 < \theta < 1$$

$$\begin{aligned} \text{var}(x) &= \int_0^1 \left\{ x - \left(\frac{1}{2} + \frac{\theta}{3}\right) \right\}^2 \{(1-2\theta) + 4\theta x\} dx \\ &= \int_0^1 \left[x^2 + \left(\frac{1}{2} + \frac{\theta}{3}\right)^2 - 2x\left(\frac{1}{2} + \frac{\theta}{3}\right) \right] [(1-2\theta) + 4\theta x] dx \end{aligned}$$

$$\begin{aligned} \log L(\theta) &= \log \prod_{i=1}^{10} f(x_i|\theta) = \log \prod_{i=1}^{10} \{(1-2\theta) + 4\theta x_i\} \\ &= \log g(\theta) \end{aligned}$$

$$\frac{d \log L(\theta)}{d\theta} = 0, \quad \theta_{MLE} = ?$$



Example: Now, let's take a numerical example with some sample data. Suppose we have the following 10 observations drawn from the following distribution: $X=\{2.5, 3.1, 4.0, 2.8, 3.7, 8.9, 4.3, 10.3, 1.1, 5.6\}$. Find MLE of θ .

$$f(x|\theta) = \frac{\theta}{x^{\theta+1}}$$

$$L(\theta) = \prod_{i=1}^{10} \frac{\theta}{x_i^{\theta+1}}, \quad \log L(\theta) = 10 \log \theta - (\theta+1) \sum_{i=1}^{10} \log x_i$$

$x_i \in X$

$$\frac{d \log L(\theta)}{d\theta} = \frac{10}{\theta} - \cancel{10} \cancel{\log} \sum \log x_i = 0$$

$$\theta_{MLE} = \frac{10}{\sum_{i=1}^n \log x_i}$$



Example: Now, let's take a numerical example with some sample data. Suppose we have the following 10 observations drawn from the following distribution: $X=\{2.5, 3.1, 4.0, 2.8, 3.7, 8.9, 4.3, 10.3, 1.1, 5.6\}$. Find MLE of θ .

$$f(x|\theta) = \frac{\theta}{x^{\theta+1}}$$

```
import math

def MLE_theta(sample_points):
    # This code computes the MLE estimate of alpha
    log_sum = 0
    for x in sample_points:
        log_sum += math.log(x)
    n = len(sample_points)
    return n / log_sum

def main():
    sample_points = [2.5, 3.1, 4.0, 2.8, 3.7, 8.9, 4.3, 10.3, 1.1, 5.6]
    theta = MLE_theta(sample_points)
    print(theta)

if __name__ == '__main__':
    main()
```

$$\theta_{MLE} = 0.7370920578760686$$

$$\theta_{MLE} = \frac{10}{\sum_{i=1}^{10} \log x_i}$$



```

import math
import numpy as np
import matplotlib.pyplot as plt

def MLE_theta(sample_points):
    # This code computes the MLE estimate of theta
    log_sum = 0
    for x in sample_points:
        log_sum += math.log(x)
    n = len(sample_points)
    return n / log_sum

def show_pdf(x, theta):
    # distribution PDF
    return theta / x**theta

def main():
    sample_points = [2.5, 3.1, 4.0, 2.8, 3.7, 8.9, 4.3, 10.3, 1.1, 5.6]
    theta = MLE_theta(sample_points)
    print(theta)

    # Plot the histogram of the sample data
    plt.hist(sample_points, bins=20, density=True,
            color='lightblue', label='Sample Data')

    # Generate points for the fitted distribution
    x_values = np.linspace(min(sample_points), max(sample_points), 1000)
    y_values = [show_pdf(x, theta) for x in x_values]

    # Plot the fitted distribution
    plt.plot(x_values, y_values, color='red', label='Distribution Fit')

    plt.xlabel('Value')
    plt.ylabel('Probability Density')
    plt.title('Distribution Fit to Sample Data')
    plt.legend()
    plt.show()

if __name__ == '__main__':
    main()

```



```

import math
import numpy as np
import matplotlib.pyplot as plt

def MLE_theta(sample_points):
    # This code computes the MLE estimate of theta
    log_sum = 0
    for x in sample_points:
        log_sum += math.log(x)
    n = len(sample_points)
    return n / log_sum

def show_pdf(x, theta):
    # distribution PDF
    return theta / x***(theta + 1)

def main():
    sample_points = [2.5, 3.1, 4.0, 2.8, 3.7, 8.9, 4.3, 10.3, 1.1, 5.6]
    theta = MLE_theta(sample_points)
    print(theta)

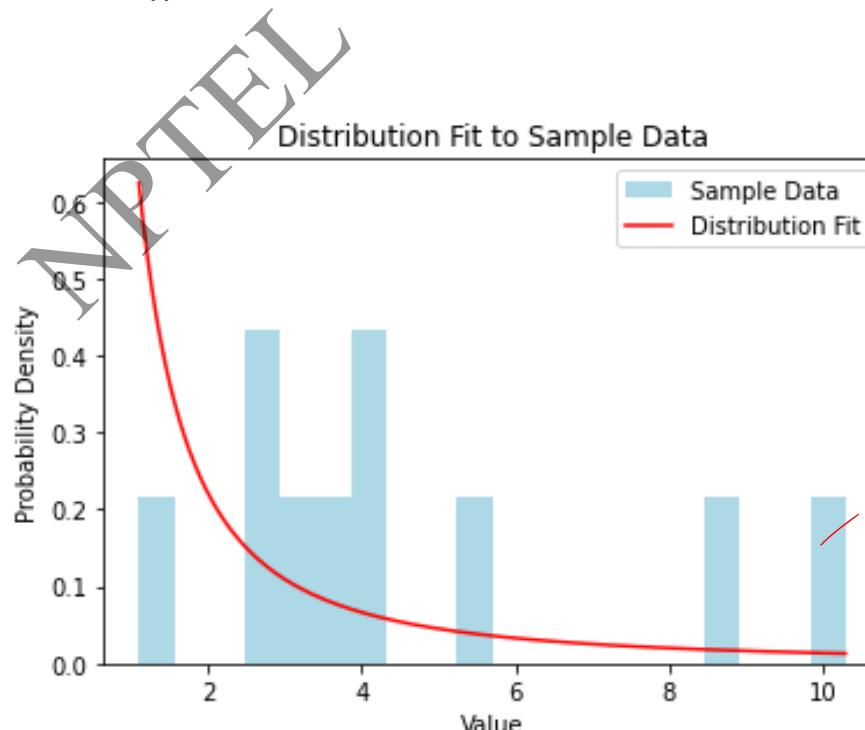
    # Plot the histogram of the sample data ✅
    plt.hist(sample_points, bins=20, density=True,
            color='lightblue', label='Sample Data')

    # Generate points for the fitted distribution
    x_values = np.linspace(min(sample_points), max(sample_points), 1000)
    y_values = [show_pdf(x, theta) for x in x_values]

    # Plot the fitted distribution
    plt.plot(x_values, y_values, color='red', label='Distribution Fit')
    plt.xlabel('Value')
    plt.ylabel('Probability Density')
    plt.title('Distribution Fit to Sample Data')
    plt.legend()
    plt.show()

if __name__ == '__main__':
    main()

```



EXAMPLE:

Let us consider a random variable $X \sim N(\mu, \sigma^2)$. If we consider a i.i.d. sample of size 3 and consider these two estimators for μ

$$\text{Estimator 1: } \hat{\theta}_1 = x_1 + x_2 - x_3 \quad \checkmark$$

$$\text{Estimator 2: } \hat{\theta}_2 = x_1 + x_2 - 2x_3$$

$$\mu_{MLE} = \bar{x} = \frac{x_1 + x_2 + x_3}{3}$$

$$\text{Estimator 1: } E(\hat{\theta}_1) = E[x_1 + x_2 - x_3] = \mu \quad \checkmark$$

$$\text{Estimator 2: } E(\hat{\theta}_2) = E(x_1 + x_2 - 2x_3) = \mu \quad \checkmark$$

But,

$$\text{Estimator 1: } var(\hat{\theta}_1) = var[x_1 + x_2 - x_3] = 3\sigma^2, \quad \checkmark$$

$$\text{Estimator 2: } var(\hat{\theta}_2) = var(x_1 + x_2 - 2x_3) = 6\sigma^2$$

What about MSE?



Conclusion

- ✓ MLE for some continuous probability distributions

To do: Properties of MLE



References

- ✓ Linear Algebra and Learning from Data (2019), Gilbert Strang, Wellesley Cambridge Press
- ✓ Machine Learning: A Probabilistic Perspective, Kevin P. Murphy (MIT Press), 2021 edition
- ✓ Deisenroth MP, Faisal AA, Ong CS. Mathematics for machine learning. Cambridge University Press; 2020 Apr 23.





Thank You...

NPTEL ONLINE CERTIFICATION COURSES
IIT KHARGPUR



SWAYAM & NPTEL COURSE ON

Mathematics for Machine Learning

by

Prof. Debjani Chakraborty

DEPARTMENT OF MATHEMATICS
IIT Kharagpur

Module 8.3

Lecture 38: Properties of Estimators

NPTEL



Concepts Covered

- 1. Unbiased Estimator**
- 2. Consistent Estimators**
- 3. Minimum Variance Estimators**
- 4. Bias vs Variance**

NPTEL



Estimators: $\hat{\theta} = f(X_1, X_2 \dots X_n)$

Estimate: Value of the Estimator

Bias: Expected difference between the value of the Estimator and True value of Estimator

$$Bias = E(\hat{\theta} - \theta) = E_{X_1, X_2 \dots X_n} [argmax_{\theta} f(X_1, X_2 \dots X_n) - \theta]$$

- Even the sample is large, due to bias estimator does not converge to true estimator
- That is, it is not expressive enough to approximate the true value arbitrarily well
- If Bias is zero it is called unbiased Estimator
- To converge, how many sample points we need!!!!



Estimators: $\hat{\theta} = f(X_1, X_2 \dots X_n)$

Estimate: Value of the Estimator

Variance: Expected difference between the value of the Estimator and True value of Estimator

$$\begin{aligned} \text{Variance} &= E \left((\hat{\theta} - E(\hat{\theta}))^2 \right) \\ &= E_{X_1, X_2 \dots X_n} ((\operatorname{argmax}_{\theta} f(X_1, X_2 \dots X_n)) - E(\operatorname{argmax}_{\theta} f(X_1, X_2 \dots X_n)))^2 \end{aligned}$$

- If more training examples are taken then variance decreases
- If variance of an estimator is minimum compared to any other estimator of the same parameter and it is unbiased then it is called **Minimum Variance Unbiased Estimator (MVUE)**



Sample Mean is unbiased estimator of population Mean for $N(\mu, \sigma^2)$

$$\underline{E(\bar{x})} = E\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n} \sum_{i=1}^n E(x_i) = \frac{1}{n} \cdot n\mu = \mu \quad \boxed{\bar{x}}$$

$$\text{Bias}(\bar{x}) = E(\bar{x} - \mu) = 0$$

$$\begin{aligned}\underline{\text{Var}(\bar{x})} &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n x_i\right) \\ &= \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n x_i\right) \\ &= \frac{1}{n^2} \cdot \sum_{i=1}^n \underline{\text{Var}(x_i)} \\ &= \frac{1}{n^2} \cdot n \sigma^2 = \frac{\sigma^2}{n}\end{aligned}$$



$$\text{variance}(\bar{X}) = \frac{\sigma^2}{n}$$

NPTEL



Sample Variance is not unbiased estimator of population Variance for $N(\mu, \sigma^2)$

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

$\times S_{\text{var}} = \frac{1}{n} \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right) \leftarrow$
 $\sqrt{S_{\text{var}}^2} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 - \frac{1}{n} \sum_{i=1}^n (\bar{x} - \mu)^2 \\ E(S_{\text{var}}^2) &= \frac{1}{n} \sum_{i=1}^n E(x_i - \mu)^2 - \frac{1}{n} \sum_{i=1}^n E(\bar{x} - \mu)^2 \\ &= \frac{1}{n} \sum_{i=1}^n \text{Var}(x_i) - \text{Var}(\bar{x}) \\ &= \sigma^2 - \frac{\sigma^2}{n} = \frac{n-1}{n} \sigma^2 \end{aligned}$$



Sample Variance is not unbiased estimator of population Variance for $N(\mu, \sigma^2)$

$$E(S_{\text{var}}^2) = \frac{n-1}{n} \sigma^2$$

$$\text{Bias}(S_{\text{var}}^2) = \frac{n-1}{n} \sigma^2 - \sigma^2 = -\frac{\sigma^2}{n}$$

$= -\text{var}(\bar{x})$

$$\boxed{S_{\text{var}}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$E\left(\frac{1}{n-1} S_{\text{var}}^2\right) = \frac{n}{n-1} \cdot \frac{n-1}{n} \sigma^2 = \sigma^2$$

\uparrow
 S_{var}^2

$$\boxed{E\left(\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2\right) = \sigma^2}$$



Consistent Estimators

If $\hat{\theta}_n$ is estimator of θ for $f(X_1, X_2 \cdots X_n | \theta)$ based on a sample of size n , if
 $\hat{\theta}_n \rightarrow \theta$ in probability when $n \rightarrow \infty$ then $\hat{\theta}_n$ is **consistent** estimator of θ .

$$P \left\{ \left| \hat{\theta}_n - \theta \right| \rightarrow 0 \right\} \rightarrow 1 \quad \text{as } n \rightarrow \infty$$

Efficient Estimators

If $\hat{\theta}_1$ and $\hat{\theta}_2$ are two unbiased estimators of θ for $f(X_1, X_2 \cdots X_n | \theta)$ based on same sample of size n , then $\hat{\theta}_1$ is more efficient than $\hat{\theta}_2$ when
 $var(\hat{\theta}_1) < var(\hat{\theta}_2)$.



Let $X_1, X_2 \dots X_n$ be a random sample from $X \sim \text{Exponential}(\lambda)$. Consider the following

estimators $\hat{\theta}_1 = \frac{1}{n} \sum_{i=1}^n X_i$ and $\hat{\theta}_2 = \frac{1}{n+1} \sum_{i=1}^n X_i$ of the population mean.

1. Are these unbiased estimators? ✓
2. Find the Mean Square Errors. ✓
3. Which estimator is better and why? ✓

$$\begin{aligned} E(\hat{\theta}_1) &= \frac{1}{n} \sum E(X_i) \\ &= \frac{\theta}{n} \cdot n = \theta. \end{aligned}$$

$$\theta = \frac{n}{n+1}$$

$$E(\hat{\theta}_1) - \theta = 0 \quad \checkmark$$

$$\begin{aligned} E(\hat{\theta}_2) - \theta &= \frac{1}{n+1} \sum_{i=1}^n E(X_i) - \theta \\ &= \frac{\theta \cdot n}{n+1} - \theta = -\frac{\theta}{n+1} \quad \checkmark \end{aligned}$$

$$\text{Var}(\hat{\theta}_1) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \cdot \frac{n}{\lambda^2} = \frac{\theta^2}{n}$$



$$\hat{\theta}_1 = \frac{1}{n} \sum_{i=1}^n X_i \text{ and } \hat{\theta}_2 = \frac{1}{n+1} \sum_{i=1}^n X_i$$

1. Are these unbiased estimators?
2. Find the Mean Square Errors.
3. Which estimator is better and why?

$$\text{Var}(\hat{\theta}_2) = \frac{1}{(n+1)} \sum_{i=1}^n \text{Var}(X_i) = \frac{n}{(n+1)} \cdot \frac{1}{\lambda^2} = \frac{n\theta^2}{(n+1)^2} \quad \text{Var}(\hat{\theta}_1) = \frac{\theta^2}{n}$$

$$\text{MSE}(\hat{\theta}_1) = \text{Bias} + \text{Variance} = 0 + \frac{\theta^2}{n} = \frac{\theta^2}{n}$$

$$\text{MSE}(\hat{\theta}_2) = \left(\frac{\theta}{n+1}\right)^2 + \frac{n\theta}{(n+1)^2} = \frac{\theta^2 + n\theta}{(n+1)^2} = \frac{\theta}{n+1}$$

$$\text{Bias}(\hat{\theta}_1) < \text{Bias}(\hat{\theta}_2)$$

$$\text{MSE}(\hat{\theta}_1) > \text{MSE}(\hat{\theta}_2)$$



$$\widehat{\theta}_1 = \frac{1}{n} \sum_{i=1}^n X_i \text{ and } \widehat{\theta}_2 = \frac{1}{n+1} \sum_{i=1}^n X_i$$

1. Are these unbiased estimators?
2. Find the Mean Square Errors.
3. Which estimator is better and why?

NPTEL



Conclusion

- ✓ Criteria to check efficiency of different estimators

NPTEL



References

- ✓ Linear Algebra and Learning from Data (2019), Gilbert Strang, Wellesley Cambridge Press
- ✓ Machine Learning: A Probabilistic Perspective, Kevin P. Murphy (MIT Press), 2021 edition
- ✓ Deisenroth MP, Faisal AA, Ong CS. Mathematics for machine learning. Cambridge University Press; 2020 Apr 23.





Thank You...

NPTEL ONLINE CERTIFICATION COURSES
IIT KHARGPUR



SWAYAM & NPTEL COURSE ON

Mathematics for Machine Learning

by

Prof. Debjani Chakraborty

DEPARTMENT OF MATHEMATICS
IIT Kharagpur

Module 8.4

Lecture 39: MLE for Linear and Logistic Regression

NPTEL



Concepts covered

- ✓ MLE for Linear Regression
- ✓ MLE for logistic Regression

NPTEL



The Method of Maximum Likelihood Estimation

If n random samples X_1, X_2, \dots, X_n are drawn from a population with probability function $f(X|\theta)$, then joint probability function L of X_1, X_2, \dots, X_n is,

$$L(\theta) = f(x_1|\theta)f(x_2|\theta) \cdots f(x_n|\theta) = \prod_{i=1}^n f(x_i|\theta)$$

L: Likelihood function.

Since, X_1, X_2, \dots, X_n are independently and identically distributed.

Objective: To find θ , such that Likelihood function is maximum
i.e. $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta)$



Linear Regression (Parameter Fitting with Maximum Likelihood)

- X = arbitrary distribution
- If $X = x$ then $Y = w_0 + w_1x + \beta$, for some constants w_0, w_1 , and some random noise variable β .
- $\beta \sim N(0, \sigma^2)$ which is independent of X .
- Here, we are assuming noise variable β has particular Gaussian distribution.

Conditional probability density function of Y , $p(y|X = x; w_0, w_1, \sigma^2)$ for each x .

$$\beta = Y - w_0 - w_1 x \sim N(0, \sigma^2)$$

The probability density under the model for the given data:

$$\prod_{i=1}^n p(y_i|x_i; w_0, w_1, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - (w_0 + w_1 x_i))^2}{2\sigma^2}}$$

Training Data

x_i	y_i	Prob.
x_1	y_1	p_1
x_2	y_2	p_2
...
x_n	y_n	p_n

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$



Linear Regression (Parameter Fitting with Maximum Likelihood)

When we visualize data, we don't know the true parameters, consider $(\theta_0, \theta_1, s^2)$ gives the probability density:

$$\prod_{i=1}^n p(y_i|x_i; \theta_0, \theta_1, s^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi s^2}} e^{-\frac{(y_i - (\theta_0 + \theta_1 x_i))^2}{2s^2}}$$

The above equation is the likelihood, a function of the parameter values.

It is easier to work with **log-likelihood**,

$$\begin{aligned} L(\theta_0, \theta_1, s^2) &= \log \prod_{i=1}^n p(y_i|x_i; \theta_0, \theta_1, s^2) \\ &= \sum_{i=1}^n \log p(y_i|x_i; \theta_0, \theta_1, s^2) = -\frac{n}{2} \log 2\pi - n \log s - \frac{1}{2s^2} \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))^2 \end{aligned}$$



Linear Regression (Parameter Fitting with Maximum Likelihood)

$$L(\theta_0, \theta_1, s^2) = -\frac{n}{2} \log 2\pi - n \log s - \frac{1}{2s^2} \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))^2 \quad \text{-----(1)}$$

To find the parameters θ_0 and θ_1 , first we will partial differentiate equation (1) with respect to θ_0 ,

$$\frac{\partial L}{\partial \theta_0} = 0$$

$$\frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$$

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

$$\frac{1}{2s^2} \times 2 \sum_{i=1}^n \{y_i - (\theta_0 + \theta_1 x_i)\} = 0 \Rightarrow \sum_{i=1}^n \{y_i - (\theta_0 + \theta_1 x_i)\} = 0$$

$$\sum_{i=1}^n y_i - \sum_{i=1}^n \theta_0 - \sum_{i=1}^n \theta_1 x_i = 0$$



Linear Regression (Parameter Fitting with Maximum Likelihood)

$$L(\theta_0, \theta_1, s^2) = -\frac{n}{2} \log 2\pi - n \log s - \frac{1}{2s^2} \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))^2 \quad \text{-----(1)}$$

To find the parameters θ_0 and θ_1 , first we will partial differentiate equation (1) with respect to θ_0 ,

$$\begin{aligned} \frac{\partial L}{\partial \theta_0} &= 0 \\ \sum_{i=1}^n y_i - \sum_{i=1}^n \theta_0 - \sum_{i=1}^n \theta_1 x_i &= 0 \\ n\bar{y} - n\theta_0 - n\theta_1 &\times \left(\frac{1}{n} \sum_{i=1}^n x_i \right) = 0 \end{aligned}$$

$$\underline{n\bar{y} - n\theta_0 - n\theta_1 \bar{x}} = 0 \Rightarrow \cancel{\bar{y}} - \theta_0 - \theta_1 \bar{x} = 0$$

$$\underline{\theta_0^{\text{MLE}}} = \bar{y} - \underline{\theta_1^{\text{MLE}}} \bar{x} \quad \text{-----(2)}$$



Linear Regression (Parameter Fitting with Maximum Likelihood)

$$L(\theta_0, \theta_1, s^2) = -\frac{n}{2} \log 2\pi - n \log s - \frac{1}{2s^2} \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))^2 \quad \dots \quad (1)$$

partial differentiate equation (1) with respect to θ_1 , $\frac{\partial L}{\partial \theta_1} = 0$

$$\frac{1}{2s^2} \times 2 \sum_{i=1}^n \{y_i - (\theta_0 + \theta_1 x_i)\}x_i = 0 \Rightarrow \sum_{i=1}^n \{y_i - (\theta_0 + \theta_1 x_i)\}x_i = 0$$

$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n \theta_0 x_i - \sum_{i=1}^n \theta_1 x_i^2 = 0 \Rightarrow \sum_{i=1}^n x_i y_i - n \times \theta_0 \times \frac{1}{n} \sum_{i=1}^n x_i - \theta_1 \sum_{i=1}^n x_i^2 = 0$$

$$\sum_{i=1}^n x_i y_i - n\theta_0 \bar{x} - \theta_1 \sum_{i=1}^n x_i^2 = 0 \quad \text{-----(3) } \checkmark$$



Linear Regression (Parameter Fitting with Maximum Likelihood)

$$L(\theta_0, \theta_1, s^2) = -\frac{n}{2} \log 2\pi - n \log s - \frac{1}{2s^2} \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 x_i))^2 \quad \dots \dots \dots (1)$$

$$\theta_0 = \bar{y} - \underline{\theta_1 \bar{x}} \quad \dots \dots \dots (2)$$

$$\frac{\partial L}{\partial \theta_0} = 0$$

$$\sum_{i=1}^n x_i y_i - \underline{n \theta_0 \bar{x}} - \theta_1 \sum_{i=1}^n x_i^2 = 0 \quad \dots \dots \dots (3)$$

$$\frac{\partial L}{\partial \theta_1} = 0$$

Using equation (2) in equation (3),

$$\sum_{i=1}^n x_i y_i - \theta_1 \sum_{i=1}^n x_i^2 - n(\bar{y} - \underline{\theta_1 \bar{x}})\bar{x} = 0 \Rightarrow \sum_{i=1}^n x_i y_i - \theta_1 \sum_{i=1}^n x_i^2 - n\bar{x}\bar{y} + n\theta_1 \bar{x}^2 = 0$$

$$\begin{aligned} \theta_1 \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) &= \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \\ \Rightarrow \theta_1 &= \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} \end{aligned}$$



Linear Regression (Parameter Fitting with Maximum Likelihood)

$$\theta_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}$$

$$\theta_1 = \frac{\sum_{i=1}^n x_i y_i - \bar{y} \sum_{i=1}^n x_i}{\sum_{i=1}^n x_i^2 - n \bar{x}^2 + n \bar{x}^2 - n \bar{x}^2} = \frac{\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \bar{y}}{\sum_{i=1}^n x_i^2 - 2n \bar{x}^2 + n \bar{x}^2}$$

$$\theta_1 = \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i^2 - \sum_{i=1}^n 2 \bar{x} x_i + \sum_{i=1}^n \bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x} + \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i^2 - 2x_i \bar{x} + \bar{x}^2)}$$

$$\theta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + \sum_{i=1}^n \bar{x}(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + n \bar{x} \bar{y} - n \bar{x} \bar{y}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\theta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$



Linear Regression (Parameter Fitting with Maximum Likelihood)

In the maximum likelihood method, the following estimators we will obtain:

$$\theta_1^{\text{MLE}} \approx w_1^{\text{MLE}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$w_0^{\text{MLE}} = \bar{y} - w_1^{\text{MLE}} \bar{x}$$

$$\sigma_{\text{MLE}}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0^{\text{MLE}} + w_1^{\text{MLE}} x_i))^2$$

$$E(\theta_1^{\text{MLE}}) = w_1 ? \\ \text{Var}(\theta_1^{\text{MLE}}) = ?$$

- The above estimators for the slope and intercept is exactly similar to the least square estimators.
- This is the crucial property of considering independent Gaussian noise.



Linear Regression (Parameter Fitting with Maximum Likelihood)

In the maximum likelihood method, the following estimators we will obtain:

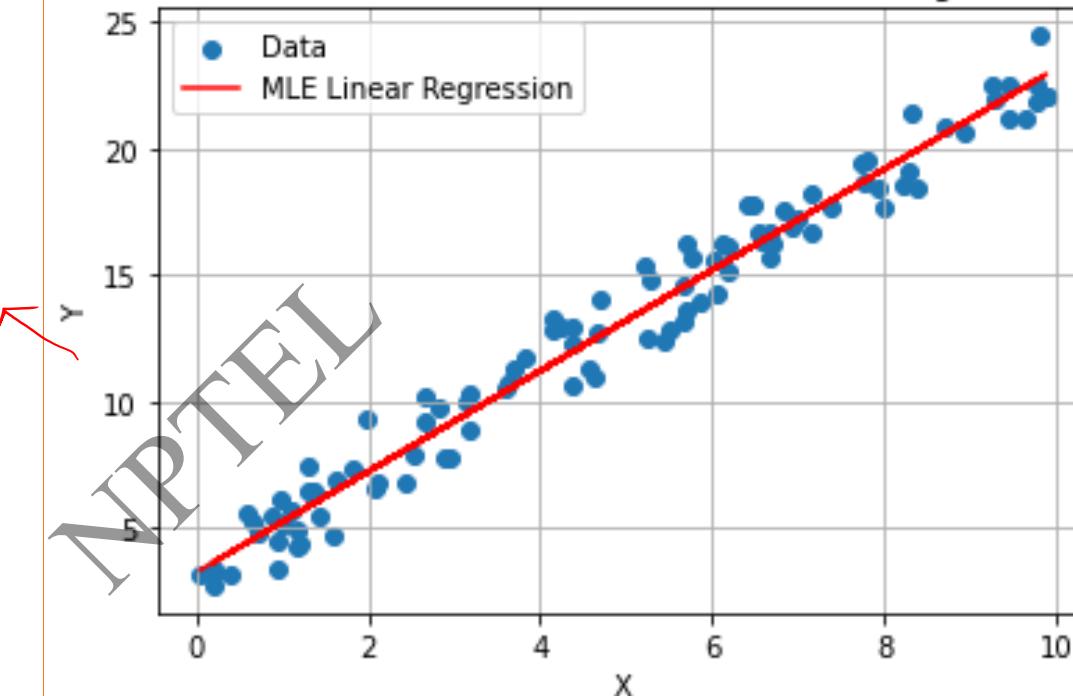
$$w_1^{MLE} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$w_0^{MLE} = \bar{y} - w_1^{MLE} \bar{x}$$

$$\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n \left(y_i - (w_0^{MLE} + w_1^{MLE} x_i) \right)^2$$

- The above estimators for the slope and intercept is exactly similar to the least square estimators.
- This is the crucial property of considering independent Gaussian noise. ✎

Maximum Likelihood Estimation for Linear Regression



Logistic Regression (Parameter Fitting with Maximum Likelihood)

$$\sigma(x) = h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}} \quad \checkmark$$

Where $0 \leq \sigma(x) \leq 1$ is the sigmoid function used in logistic regression.

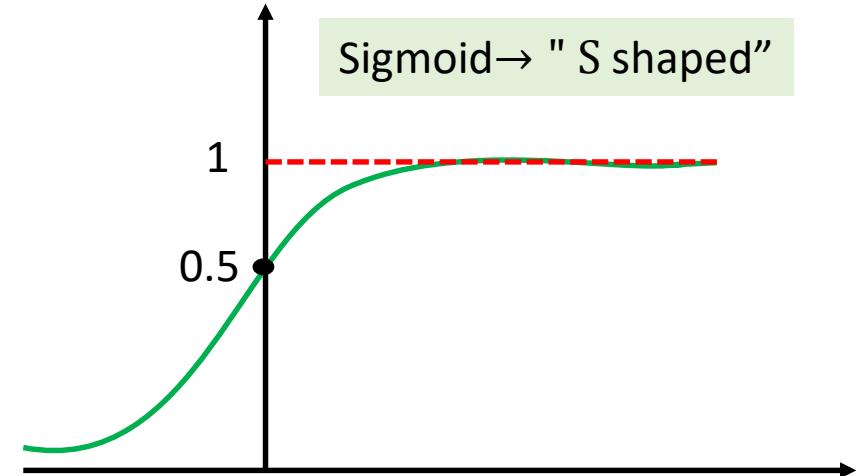
$\sigma(x)$ Gives the probability

$$P(y=1|x=x_i) = p(x_i) = \sigma(x_i) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_i)}} \quad \checkmark$$

Rearrange: $z = \theta_0 + \theta_1 x_i = \ln \left[\frac{p(x_i)}{1-p(x_i)} \right]$

$$1 - p(x_i) = 1 - \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_i)}}$$

$$= \frac{e^{-(\theta_0 + \theta_1 x_i)}}{1+e^{-(\theta_0 + \theta_1 x_i)}} \quad \checkmark$$



Logistic Regression (Parameter Fitting with Maximum Likelihood)

Probability for being in category 1,

$$p(y = 1|x = x_i) = p(x_i) = \sigma(x_i) = \frac{1}{1+e^{-(\theta_0+\theta_1x_i)}} \quad \text{--- (1)}$$

Probability for being in category 0,

$$p(y = 0|x = x_i) = 1 - p(x_i) \quad \text{--- (2)}$$

By equation (1) and (2),

$$P(y = y_i|x = x_i) = p(x_i)^{y_i} [1 - p(x_i)]^{1-y_i} \text{ where } y_i \in \{0,1\}$$

$$\text{when } y_i = 1, p(y = y_i|x = x_i) = p(x_i) \quad \text{--- (3)}$$

$$\text{when } y_i = 0, p(y = y_i|x = x_i) = 1 - p(x_i) \quad \text{--- (4)}$$

Fitting using maximum likelihood,

$$p(x_i) = \frac{1}{1 + e^{-(\theta_0+\theta_1x_i)}}$$

$$\begin{aligned} y_i &= 0 \\ y_i &= 1 \end{aligned}$$

$$P(y = y_i|x = x_i) = p(x_i)^{y_i} [1 - p(x_i)]^{1-y_i}$$



Logistic Regression (Parameter Fitting with Maximum Likelihood)

For n independent training example, we can write

$$= P_1 \times P_2 \times \cdots \times P_n \quad \checkmark$$

$$= \prod_{i=1}^n p(x_i)^{y_i} [1 - p(x_i)]^{1-y_i}$$

For n data points the likelihood function.

Find θ_0 and $\theta_i, i = 1, 2 \dots n$ that maximize L or $\log(L)$.

$$\sigma(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n)}}$$

$$\log \left[\frac{p(x_i)}{1 - p(x_i)} \right] = \theta_0 + \theta_1 x_1 + \cdots + \theta_n x_n$$

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

It does not have a closed- form solution so we cannot write down a simple expression of the estimator

~~UNPUBLISHED~~

x_i $y_i \in \{0, 1\}$

Training Data

x_i	y_i (1 or 0)	Prob.
x_1	y_1	P_1
x_2	y_2	P_2
...
x_n	y_n	P_n



Logistic Regression (Parameter Fitting with Maximum Likelihood)

To find the parameters maximize $\log(L)$ equivalently we can minimize $\log(\text{likelihood})$ (NLL) to maintain uniformity with other algorithm:

$$\begin{aligned} J(w) &= -\log(L) = -\sum_{i=1}^n \{y^{(i)} \log p(x_i) + (1 - y^{(i)}) \log(1 - p(x_i))\} \\ &= -\sum_{i=1}^n \{y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})\} \quad \times \end{aligned}$$



References

- ✓ Linear Algebra and Learning from Data (2019), Gilbert Strang, Wellesley Cambridge Press
- ✓ Machine Learning: A Probabilistic Perspective, Kevin P. Murphy (MIT Press), 2021 edition
- ✓ Deisenroth MP, Faisal AA, Ong CS. Mathematics for machine learning. Cambridge University Press; 2020 Apr 23.





SWAYAM & NPTEL COURSE ON

Mathematics for Machine Learning

by

Prof. Debjani Chakraborty

DEPARTMENT OF MATHEMATICS
IIT Kharagpur

Module 8.5

Lecture 40: MLE for Naïve Bayes Model

NPTEL



Concepts covered

- ✓ MLE for Naïve Bayes Model

NPTEL



The Method of Maximum Likelihood Estimation

If k random samples X_1, X_2, \dots, X_k are drawn from a population with probability function $f(X|\theta)$, then joint probability function L of X_1, X_2, \dots, X_k is,

$$L(\theta) = f(x_1|\theta)f(x_2|\theta) \cdots f(x_k|\theta) = \prod_{i=1}^k f(x_i|\theta)$$

L: Likelihood function.

Since, X_1, X_2, \dots, X_k are independently and identically distributed.

Objective: To find θ , such that Likelihood function is maximum
i.e. $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta)$



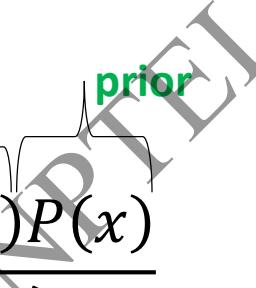
Assume,

1. we have some prior knowledge $p(x)$
2. a random variable y , assumes value from random sample
3. We can update the knowledge $p(x)$ after observing y , $P(x/y)$

$$P(x/y) = \frac{P(y/x)P(x)}{P(y)}$$

Diagram illustrating the components of Bayes' Theorem:

- posterior**: The final result, $P(x/y)$.
- likelihood**: $P(y/x)$
- prior**: $P(x)$
- evidence**: $P(y)$



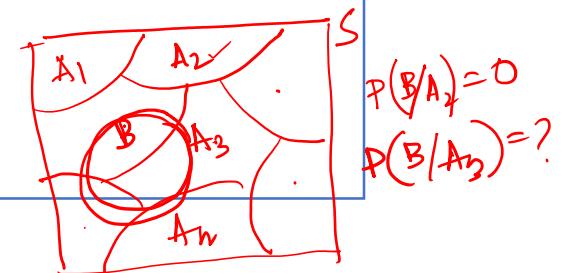


Bayes Rule :

Hypothesis: (Equal division of ignorance): If nothing is known about the prior probabilities $P(A_1) \cdots P(A_n)$ then they are all equal.

Let $A_1, A_2 \cdots A_n$ denote a disjoint partition of a outcome set S and let B be any event. Let $P(A_i) \neq 0, i = 1, 2 \cdots n$ and $P(B) \neq 0$, then for $i = 1, 2 \cdots n$

$$\underline{P(A_i/B)} = \frac{P(B/A_i)P(B)}{\sum_{i=1}^n P(B/A_i)P(B)}$$



Result: Let $A_1, A_2 \cdots A_n$ denote a disjoint partition of a outcome set S and let B be any event. Let $P(A_i) \neq 0, i = 1, 2 \cdots n$ and $P(B) \neq 0$, then for $i = 1, 2 \cdots n$

$$\underline{P(B)} = \sum_{i=1}^n P(A_i)P(B/A_i)$$



Naïve Bayes Model

Example: Classifying documents using bag of words

Document classification is the problem of classifying text documents into different categories.

One simple approach is to represent each document as a binary vector,

each word is present or not, so $x_{ij} = 1$ iff word j occurs in document i , otherwise $\underline{x_{ij}} = 0$

Email Subject	Label
Offer: Double your income	Spam
Free gift	spam
You have been selected for Award	spam
Open for a free webinar	Not Spam
Meeting scheduled for tomorrow	Not Spam
Reminder: Review meeting agenda	Not spam

In Naive Bayes, we make a "naive" assumption that all features are conditionally independent given the class label.



Example: Classifying documents using bag of words

$$P(\text{Spam} | \text{"offer", "free"}) = \frac{P(\text{Spam}) \cdot P(\text{"offer"}|\text{Spam}) \cdot P(\text{"free"}|\text{Spam})}{P(\text{"offer", "free"})}$$

$$\underline{P(\text{spam}) = .5}, \underline{P(\text{not spam}) = .5}$$

$$\begin{aligned} & P(\text{"offer"}|\text{"spam"}) \\ &= \frac{\text{number of occurrences of feature} + \alpha}{\text{total number of instances of the given class} + \alpha * \text{Dimension of the data}} \\ &= \frac{3 + 1}{3 + 2} = \frac{4}{5} = .8 \end{aligned}$$

$\alpha = 1$

$$P(\text{"offer"}|\text{"Not spam"}) = .2$$

Email Subject	Label
Offer: Double your income	Spam
Free gift	spam
You have been selected for Award	spam
Open for a free webinar	Not Spam
Meeting scheduled for tomorrow	Not Spam
Reminder: Review meeting agenda	Not spam

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$



Email Subject	Label
Double your income	Spam
Offer: Free gift	spam
You have been selected for Award	spam
Open for a free webinar	Not Spam
Meeting scheduled for tomorrow	Not Spam
Reminder: Review meeting agenda	Not spam

Naïve Bayes Model

Classifying documents using bag of words

$$P(\text{Spam} \mid \text{"offer", "free"})$$

$$= \frac{P(\text{Spam}) \cdot P(\text{"offer"}|\text{Spam}) \cdot P(\text{"free"}|\text{Spam})}{P(\text{"offer", "free"})}$$

$$P(\text{spam}) = .5, P(\text{not spam}) = .5$$

$$P(\text{"free"}|\text{"spam"})$$

$$= \frac{\text{number of occurrences of feature} + \alpha}{\text{total number of instances of the given class} + \alpha * \text{Dimension of the data}}$$

$$= \frac{1 + 1}{3 + 2} = \frac{2}{5} = .4 \quad \times$$

$$\alpha = 1$$

$$P(\text{"free"}|\text{"Not spam"}) = .6 \quad \checkmark$$



Email Subject	Label
Offer: Double your income	Spam
Free gift	spam
You have been selected for Award	spam
Open for a free webinar	Not Spam
Meeting scheduled for tomorrow	Not Spam
Reminder: Review meeting agenda	Not spam

Naïve Bayes Model

Classifying documents using bag of words

$$P(\text{Spam} \mid \text{"offer", "free"})$$

$$= \frac{P(\text{Spam}) \cdot P(\text{"offer"}|\text{Spam}) \cdot P(\text{"free"}|\text{Spam})}{P(\text{"offer", "free"})}$$

$$\underline{P(\text{spam}) = .5, P(\text{not spam}) = .5}$$

$$\frac{P(\text{Spam}) \cdot P(\text{"offer"}|\text{Spam}) \cdot P(\text{"free"}|\text{Spam})}{P(\text{"offer", "free"})} \text{ equivalent to } .5 \cdot .8 \cdot .4 = .16$$



Naïve Bayes Model

✓

Document classification is the problem of classifying text documents into different categories.
One simple approach is to represent each document as a binary vector,

each word is present or not, so $x_{ij} = 1$ iff word j occurs in document i , otherwise $x_{ij} = 0$

We can then use the following class conditional density:

$$P(\text{Spam} \mid \text{"offer", "free"}) \\ = \frac{P(\text{Spam}) \cdot P(\text{"offer"}|\text{Spam}) \cdot P(\text{"free"}|\text{Spam})}{P(\text{"offer", "free"})}$$

APTEL

$$P(\text{class}_i / \text{feature}_{n \in [1..N]}) \\ = P(\text{class}_i) \prod_{n=1}^N P(\text{feature}_n / c_i)$$



Fitting a naive Bayes model with maximum likelihood

$$P(\text{class}_i/\text{feature}_{n \in [1..N]}) = P(\text{class}_i) \prod_{n=1}^N P(\text{feature}_n/c_i)$$

We use the maximum likelihood method in finding parameters that maximize the likelihood of the observed w_n data set of size D , $\{w_n, c_i\}_{i=1..C, n=1..N}$ in class c_i

Then likelihood under a certain model with parameters $f(X|\theta_i)$ of class c_i

$$L(D|\theta_{i \in [1..C]}, f(X|\theta_i)) = \prod_{d=1}^D \left(P(c_i) \prod_{n=1}^N P(w_n/c_i) \right)$$

Taking log-likelihood

$$\log L = \sum_{d=1}^D \sum_{i=1}^C \log P(c_i) + \sum_{d=1}^D \sum_{c=1}^C \sum_{n=1}^N \sum_{v=1}^V \log P(w_n|c_i),$$



Fitting a naive Bayes model with maximum likelihood

$$P(\text{class}_i / \text{feature}_{n \in [1..N]}) \\ = P(\text{class}_i) \prod_{n=1}^N P(\text{feature}_n / c_i)$$

Taking log-likelihood

$$\log L = \sum_{d=1}^D \sum_{i=1}^C \log P(c_i) + \sum_{d=1}^D \sum_{c=1}^C \sum_{n=1}^N \sum_{v=1}^V \log P(w_n | c_i),$$

e.g. D: Number of emails C: Number of classes
V: Total number of words in the vocabulary
N: Number of words / features we are looking for

$$P(c_i) = \frac{\text{Total number of time we see the class } c_i}{\text{Number of data}}$$

$$\sum_{i=1}^C P(c_i) = 1 \quad \checkmark \\ \sum_{n=1}^N \sum_{i=1}^C P(w_n | c_i) = 1 \quad \checkmark$$



Fitting a naive Bayes model with maximum likelihood

$$P(\text{class}_i / \text{feature}_{n \in [1..N]}) \\ = P(\text{class}_i) \prod_{n=1}^N P(\text{feature}_n / c_i)$$

Taking log-likelihood

$$\log L = \sum_{d=1}^D \sum_{i=1}^C \log P(c_i) + \sum_{d=1}^D \sum_{c=1}^C \sum_{n=1}^N \sum_{v=1}^V \log P(w_n | c_i),$$

e.g. D: Number of emails C: Number of classes
V: Total number of words in the vocabulary
N: Number of words / features we are looking for

$$P(c_i) = \frac{\text{Total number of time we see the class } c_i}{\text{Number of data}} \quad \checkmark$$

Each c_i has its own probability of success θ_i in the probability distribution

$$\theta_i^{MLE} = \frac{\text{number of time we see the word in the class } c_i}{\text{Total number of times we see the word in all classes}} \quad \checkmark$$



Conclusion

Maximum likelihood estimation based on training data

NPTEL



References

- ✓ Linear Algebra and Learning from Data (2019), Gilbert Strang, Wellesley Cambridge Press
- ✓ Machine Learning: A Probabilistic Perspective, Kevin P. Murphy (MIT Press), 2021 edition
- ✓ Deisenroth MP, Faisal AA, Ong CS. Mathematics for machine learning. Cambridge University Press; 2020 Apr 23.





Thank You...

NPTEL ONLINE CERTIFICATION COURSES
IIT KHARGPUR