

* Probability Distributions:-

* Binomial (Discrete)

$$PMF = {}^n C_x p^x (1-p)^{n-x}$$

n, p are the parameters of distribution

Mean and Variance:-

$$\mu = E(x) = np$$

$$\sigma^2 = V(x) = np(1-p)$$

* Poisson Dist:-

Let λ be the rate at which the event occurs.

t be the length of time interval
 x be the total no. of events in that time interval.

$$\mu = \lambda t$$

$$P(x=x) = e^{-\mu} \frac{\mu^x}{x!}$$

μ is parameter of the distribution

$$E(x) = \mu$$

$$V(x) = \mu$$

Continuous dist:-

* Uniform dist:-

If Random variable is defined by

$$f(x) = \frac{1}{b-a}, \quad -\infty < a \leq x \leq b < \infty$$

$$E(x) = \frac{a+b}{2}$$

$$V(x) = \frac{(b-a)^2}{12}$$

* Normal Dist:-

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

$$E(x) = \mu, \quad V(x) = \sigma^2$$

* Exponential Dist:-

$$f(x) = \lambda e^{-\lambda x}, \quad \forall x \geq 0$$

λ - Rate (Instantaneous failure rate)

$$E(x) = \frac{1}{\lambda}, \quad V(x) = \left(\frac{1}{\lambda}\right)^2$$

$$P\{x \leq x\} = (1 - e^{-\lambda x})$$

$$P\{x > x\} = e^{-\lambda x}$$

$$P\{x_1 < x \leq x_2\} = e^{-\lambda x_1} - e^{-\lambda x_2}$$

MLE of Poisson Dist:-

Given x_1, \dots, x_n are independent random Poisson variables with rate (mean) λ and,

$$PMF (x_i = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

for $k = 0, 1, 2, \dots$

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

$$\begin{aligned} \ln L(\lambda | x_1, x_2, \dots, x_n) &= \sum_{i=1}^n \log \left(\frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \right) \\ &= \sum_{i=1}^n \left[-\lambda + x_i \log \lambda - \log(x_i!) \right] \end{aligned}$$

To find maximum likelihood we diff w.r.t λ :-

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

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

MLE of normal dist:-

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

$$L(\mu, \sigma | x_1, \dots, x_n) = \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2 \right) \right]$$

$$\ln L(\mu, \sigma | x_1, \dots, x_n) =$$

$$\sum_{i=1}^n \left[-\frac{1}{2} \log(2\pi) - \log(\sigma) - \frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2 \right]$$

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

for μ we diff w.r.t μ and for σ we diff w.r.t σ .

$$\hat{\mu} = \frac{\sum_{i=1}^n x_i}{n} \quad \left| \quad \hat{\sigma} = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}}$$

MLE for Bernoulli:-

$$f(x) = p^x (1-p)^{1-x}$$

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

$$\ln L(p) = \sum_{i=1}^n \left[x_i \log p + (1-x_i) \log (1-p) \right]$$

Diff w.r.t p :-

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

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

*MLE for Beta dist:-

It is a continuous dist. on $[0, 1]$ parameterised by 2 +ve shape parameters denoted as α and β which represent prior knowledge about dist.

$$L(\alpha, \beta | x_1, \dots, x_n) = \prod_{i=1}^n \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \Gamma(\beta)} x_i^{\alpha-1} (1-x_i)^{\beta-1}$$

Γ is gamma function

$$\begin{aligned} \ln L(\alpha, \beta) &= \sum_{i=1}^n (\alpha-1) \log(x_i) + (\beta-1) \log(1-x_i) \\ &\quad - \log(\Gamma(\alpha)) - \log(\Gamma(\beta)) \\ &\quad + \log(\Gamma(\alpha + \beta)) \end{aligned}$$

We will diff and will get:-

$$\hat{\alpha} = \frac{\sum_{i=1}^n x_i (\alpha-1) + \sum_{i=1}^n (1-x_i) (\beta-1)}{\alpha + \beta - 2n}$$

$$\hat{\beta} = \frac{n(\alpha + \beta)}{\alpha + \beta - 2n}$$

Bayes Thm:-

$$P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)}$$

* Maximum A-Posteriori Estimation:-

From Bayes:-

$$P(\theta|D) = \frac{\overbrace{P(D|\theta)}^{\text{Likelihood}} \overbrace{P(\theta)}^{\text{Prior}}}{P(D)}$$

$$\theta_{ML} = \arg \max_{\theta} \sum_i \log P(x_i|\theta)$$

$$\theta_{MAP} = \arg \max_{\theta} \left[\sum_i \log P(x_i|\theta) + \sum_i \log P(\theta) \right]$$

* Linear Regression:-

$$\hat{y} = \omega_0 + \omega_1 x_i$$

$$= x^T \omega$$

Learn a function which passes through as many lines as possible: Minimizing Least squares Error:-

$$E(\omega) = \frac{1}{2} \sum_i (y_i - \hat{y})^2$$

diff wrt ω

$$\nabla E(\omega) = \nabla \left[\frac{1}{2} \sum_i \| (y_i - x_i^T \omega) \|^2 \right]$$

$$= x y - x x^T \omega = 0$$

$$\omega_{ML} = (x x^T)^{-1} x y$$

* While calculating linear regression:-

$$m = \frac{(n \sum xy - \sum y \sum x)}{n \sum x^2 - (\sum x)^2}$$

$$c = \frac{(\sum y - m \sum x)}{n}$$

* Polynomial Regression:-

$$y(x, \omega) = \omega_0 + \omega_1 x + \omega_2 x^2 + \dots + \omega_{n-1} x^{n-1}$$

$$= \sum_{j=0}^{n-1} \omega_j x^j$$

Coefficients can be learned by minimizing Error function:-

$$E(\omega) = \frac{1}{2} \sum_{n=1}^n (y(x_n, \omega) - t_n)^2$$

t_n - actual data-point

$$\phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{n-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{n-1}(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_n) & \phi_1(x_n) & \dots & \phi_{n-1}(x_n) \end{bmatrix}$$

$N \times n$

$$\omega = (\phi^T \phi)^{-1} \phi^T t$$

> Adding Regularization coefficient to Error function in order to minimise over-fitting:-

$$\hat{E}(\omega) = \frac{1}{2} \sum_{n=1}^N \{ y(x_n, \omega) - t_n \}^2 + \frac{\lambda}{2} \|\omega\|^2$$

Regularized Least Squares:-

$$E(\omega) = \frac{1}{2} \sum_{n=1}^N \{ t_n - \omega^T \phi(x_n) \}^2 + \sum_{j=1}^M \sum_i |\omega_j|^q$$

If $q=2 \rightarrow$ Ridge Regression

If $q=1 \rightarrow$ Lasso Regression

Generally the family of regression containing this regularized least squares error function with different values of $q \rightarrow$ Elastic net regularization

The regularized least squares solution is:-

$$\omega = (\lambda I + \phi^T \phi)^{-1} \phi^T t$$

* Least Squares = MLE

$$\arg \min_{\omega} E = \arg \max_{\omega} L$$

* Regularized Least Squares = MAP

* Bias Variance decomposition:-

$$\begin{aligned} & E_0 [\{y(x;0) - h(x)\}^2] \\ &= \underbrace{\{E_0[y(x;0) - h(x)]\}^2}_{\text{bias}^2} + \underbrace{E_0 \{y(x;0) - E_0[y(x;0)]\}^2}_{\text{Variance}} \end{aligned}$$

→ Overfitting - High variance

→ Underfitting - High bias

* Linear Discriminant Analysis:-

$$m_1 = \frac{1}{N} \sum_{i \in c_1} x_i, \quad m_2 = \frac{1}{N} \sum_{i \in c_2} x_i$$

$$\text{Maximise } m_2 - m_1 = \omega^T (m_2 - m_1)$$

Fisher's Discriminant Analysis:-

Improvement of LDA.

Maximise a function that will give large separation b/w projected classes while giving small variance within each class.

∴ find ω that maximises

$$J(\omega) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

$$m_1 = \frac{\sum_t \omega^T x^t}{\sum_t x^t}, \quad s_1^2 = \sum_t (\omega^T x^t - m_1)^2$$

$$\text{FLO solution} \rightarrow \omega = C \cdot S^{-1} (m_1 - m_2) \\ S_\omega = S_1 + S_2$$

* WLS:-

$$Y = \beta x + \epsilon \rightarrow \text{error}$$

↑
Params
↓
dependent term

Here,
 ϵ assumed to be (multi-variate) normally distributed with mean vector 0 and non-constant variance-covariance matrix.

$$\begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix}$$

If we define the reciprocal of each variance σ_i^2 , as the weight $w_i = 1/\sigma_i^2$, then let matrix ω be a diagonal matrix containing these weights:-

$$\omega = \begin{bmatrix} w_1 & 0 & \dots & 0 \\ 0 & w_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & w_n \end{bmatrix}$$

The weighted least squares estimate is then:-

$$\begin{aligned} \hat{\beta}_{WLS} &= \arg \min_{\beta} \sum_{i=1}^n (\epsilon_i)^2 \\ &= (X^T \omega X)^{-1} X^T \omega Y \end{aligned}$$

In OLS:-

$$\text{Least squares} = \sum_{i=1}^n (y_i - x_i \beta)^2$$

$$WLS = \sum_{i=1}^n w_i (y_i - x_i \beta)^2$$

* Polynomial Regression:-

$$y = b_0 + b_1x + b_2x^2$$

$$\begin{bmatrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \\ \sum x_i^2 y_i \end{bmatrix}$$

* Mid Sem revision:-

* Supervised Learning - Regression and Classification

- Linear regression
- Polynomial "
- Logistic regression
- Poisson Regression
- Naive Bayes
- KNN
- Decision Trees
- Ensemble methods
- Model Selection and regularization
- Bias and Variance

* Unsupervised Learning:-

- Dimensionality Reduction
- LDA, FDA

→ Bin Prediction:-

Assuming Bernoulli dist:-

$$L(P|x) = (p)^{N_0} (1-p)^{N_1}$$

$$l(P|x) = N_0 \log p + N_1 \log (1-p)$$

$$MLE = \frac{d}{dp} (l(P|x)) = \frac{N_0}{p} - \frac{N_1}{1-p} = 0$$

$$\therefore p = \frac{N_0}{N_0 + N_1}$$

Assuming a Beta prior:-

$$MAP = \frac{N_0 + \alpha - 1}{N_0 + N_1 + \alpha + \beta - 2}$$