

3.6) poisson pmf  $\text{poi}(x/\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$

To estimate the parameter  $\lambda$ , we do MLE by solving optimization problem

$$\hat{\lambda}_{MLE} = \underset{\lambda}{\operatorname{argmax}} \text{poi}(D/\lambda)$$

$$\approx \underset{\lambda}{\operatorname{argmax}} \log(\text{poi}(D/\lambda))$$

To solve this, we'll do differentiation ~~on both sides~~ on log-likelihood.

$$\frac{d}{d\lambda} (\log(\text{poi}(D/\lambda))) = 0$$

$$\frac{d}{d\lambda} \left( \sum_{i=1}^N \log e^{-\lambda} \cdot \frac{\lambda^{x_i}}{x_i!} \right) = 0$$

$$\frac{d}{d\lambda} \left( \sum_{i=1}^N \log e^{-\lambda} + \sum_{i=1}^N \log \lambda^{x_i} + \sum_{i=1}^N \log x_i! \right) = 0$$

$$\frac{d}{d\lambda} \left( -N\lambda + (\log \lambda) \sum_{i=1}^N x_i + \sum_{i=1}^N \log x_i! \right) = 0$$

$$-N + \frac{1}{\lambda} \sum_{i=1}^N x_i = 0 \Rightarrow \lambda = \frac{\sum_{i=1}^N x_i}{N}$$

MLE

So, the MLE estimate for parameter is empirical mean. When  $N$  is large, this coincides with the expected value of distribution for  $N$  large,  $E(x \sim \text{poi}(x/\lambda)) = \lambda = \lambda_{MLE}$

3.8)  
a) Uniform distribution centered on 0 with width  $2a$ .

$$x \sim \text{Unif}(0, 2a) = \frac{1}{2a} \mathbb{I}(x \in [-a, a])$$

To estimate the parameter  $a$ , we calculate MLE by solving optimization problems for likelihood.

$$\hat{a}_{\text{MLE}} = \underset{a}{\operatorname{argmax}} \text{Unif}(D/a) \quad D = \text{dataset}$$

$$= \underset{a}{\operatorname{argmax}} \prod_{i=1}^N \frac{1}{2a} \mathbb{I}(x_i \in [-a, a])$$

$$= \frac{1}{2a} \underset{a}{\operatorname{argmax}} \left( \frac{1}{2a} \right)^N \prod_{i=1}^N \mathbb{I}(x_i \in [-a, a])$$

In order to maximize this, 'a' value should be minimum and in between  $[-a, a]$

$$\text{So, } \hat{a}_{\text{MLE}} = \max(x_i \in [-a, a])$$

3.9)  
a) MLE estimation for exponential distribution

$$x \sim \exp(x/\theta) = \theta e^{-\theta x}$$

To estimate parameter  $\theta$ , we calculate MLE by solving optimization problem for likelihood.

$$\hat{\theta}_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \exp(D/\theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N \theta \cdot e^{-\theta x_i}$$

$$\sim \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log(\theta \cdot e^{-\theta x_i})$$



To get maximum for  $\theta$ , we do differentiation

$$\frac{d}{d\theta} \left( \sum_{i=1}^N \log \theta - e^{-\theta x_i} \right) = 0$$

$$\sum_{i=1}^N \frac{1}{\theta} - \sum_{i=1}^N x_i = 0$$

$$\frac{N}{\theta} = \sum_{i=1}^N x_i \Rightarrow \theta = \frac{1}{\frac{\sum_{i=1}^N x_i}{N}} = \frac{1}{\bar{x}}$$

If we say  $\bar{x}$  is empirical estimate for mean

then the MLE for parameter  $\theta$  is

$$\hat{\theta}_{MLE} = \frac{1}{\bar{x}}$$

and when  $N$  is large  $\bar{x} \approx E(X)$

So, MLE substantiates that exponential factor is  $1/E(X)$ .

3.11) According to above derivation

$$\hat{\theta}_{MLE} = \frac{1}{\bar{x}}$$

for  $\{x_1 = 5, x_2 = 6, x_3 = 4\}$  dataset

$$\bar{x} = \frac{5+6+4}{3} = 5$$

$$\hat{\theta}_{MLE} = \frac{1}{\bar{x}} = \frac{1}{5} = 0.2$$

4.1) It's given that  $X \sim U(-1, 1)$

$$Y = X^2$$

$$E(X) = \int_{-1}^1 x \cdot p(x) dx = \int_{-1}^1 x \cdot \frac{1}{2} dx = \left( \frac{x^2}{2} \right)_{-1}^1 = 0$$

$$E(Y) = \int_{-1}^1 x^2 \cdot p(x) dx = \int_{-1}^1 x^2 \cdot \frac{1}{2} dx = \frac{1}{2} \left( \frac{x^3}{3} \right)_{-1}^1$$

$$\text{Correlation Coefficient} = \frac{1}{3}$$

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

$$\text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$$

$$= E(XY) - E(X)E(Y)$$

$$E(XY) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy \cdot p(x)p(y) dx dy$$

$$= \int_{-\infty}^{\infty} x \cdot \left( \int_{-\infty}^{\infty} y \cdot p(y) dy \right) dp(x) dx$$

$$= \int_{-\infty}^{\infty} x \cdot \left( \int_{-1}^1 x^2 p(x) dx \right) dx \quad \left[ \text{change of variables formula} \right]$$

$$= \int_{-1}^1 x \cdot \left( \frac{1}{2} \cdot \frac{x^3}{3} \right)_{-1}^1 dx = \int_{-1}^1 \frac{1}{3} \cdot x dx = \left( \frac{x^2}{6} \right)_{-1}^1 = 0$$

$$\text{So, } \text{Cov}(X, Y) = 0 - 0 = 0.$$

From this we can say, even though  $X, Y$  are not independent, their  $\text{Cov}(X, Y)$  can be 0.

4.2)  $X \sim N(0, 1)$

a

and  $Y = WX$ ,  $W$  is discrete r.v.  $\begin{matrix} \nearrow 1 \\ \searrow -1 \end{matrix} \rightarrow$

Here  $W$  is discrete random variable with support  $\{-1, 1\}$  and

$$P(W) = 0.5$$

Now, we can see random variable

$$Y = \begin{cases} -X & \text{When } W = -1 \\ X & \text{When } W = 1 \end{cases}$$

As  $X$  follows Gaussian  $N(0, 1)$

$$P(X) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

If we put  $-x$  in place  $x$  we get same equation

So,  $P(Y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}}$ , so  $Y$  also follows

Gaussian with same parameters  $0, 1$ .

$$Y \sim N(0, 1)$$

b)  $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$

Since  $X, Y$  are both Gaussian  $E(X) = 0$   
 $E(Y) = 0$

$$E(XY) = E(E(XY|W))$$

$$= 0.5 E[XY|W=-1] + 0.5 E[XY|W=1]$$

$$= 0.5 E[-X^2] + 0.5 E[X^2]$$

$$= 0.5 E[X^2 - X^2] = 0$$



4.3) we can prove  $\rho(x, y)$  (correlation coefficient) is always between  $-1$  and  $1$ . by using Cauchy-Schwarz inequality

$$[\text{Cov}(x, y)]^2 \leq \text{Var}(x) \text{Var}(y)$$

$$\text{So, } |\text{Cov}(x, y)| \leq \sqrt{\text{Var}(x) \text{Var}(y)}$$

$$|\rho(x, y)| \leq 1$$

$$-1 \leq \rho(x, y) \leq 1$$

4.4) Correlation coefficient for linearly related variables is  $\pm 1$ .

let  $y$  is linearly dependent on  $x$

$$y = ax + b$$

$$\text{Var}(x) = E(x^2) - E^2(x) \quad \text{Var}(y) = a^2 [E(x^2) - E^2(x)]$$

$$\text{Cov}(x, y) = E(xy) - E(x)E(y)$$

$$= E[x(ax + b)] - [E(x)(aE(x) + b)]$$

$$= aE(x^2) + bE(x) - aE^2(x) - bE(x)$$

$$= aE(x^2) - aE^2(x)$$

$$\rho(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}} = \frac{a(E(x^2) - E^2(x))}{\sqrt{a^2 [E(x^2) - E^2(x)]}}$$

$$\text{So, when } a > 0 \Rightarrow \frac{a}{|a|} = 1 \quad = \frac{a}{|a|} = \pm 1.$$

$$\text{and } a < 0 \Rightarrow \frac{a}{|a|} = -1$$

\* we know that posterior

$$P(\theta/D) \propto P(D/\theta)P(\theta)$$

$\theta$  - parameters to be estimated from data

$D$  - Dataset

$P(D/\theta)$  - likelihood

$P(\theta)$  - prior for the parameters.

the Maximum A posteriori (MAP) estimate for  $\theta$

$$\hat{\theta}_{\text{MAP}} \Rightarrow \underset{\theta}{\operatorname{argmax}} P(D/\theta)P(\theta)$$

$$= \underset{\theta}{\operatorname{argmax}} [\log P(D/\theta) + \log P(\theta)]$$

In this equation prior is always constant, so we can write it

$$\hat{\theta}_{\text{MAP}} \approx \underset{\theta}{\operatorname{argmax}} \log P(D/\theta) = \hat{\theta}_{\text{MLE}}$$

This is MLE (Maximum likelihood estimate) for  $\theta$  and as the dataset size ( $N$ ) increases

$P(D/\theta)$  converges to single  $\theta$ . This is nothing

when data size increases it overwhelms the prior

This is analogous to central limit theorem, where

when we have enough sample data, parameters calculated on samples best fit the true underlying distribution.

In same way, as we are trying to optimize based on likelihood data  $\{ \arg \max_{\theta} p(D|\theta) \}$ , when we get enough data, this estimated parameters after solving optimization problem converges to true parameter.