Marginal Distributions $\quad p(A), p(B)$

Joint Distributions $\quad p(A,B) = p(A)p(B) \quad$ if $A, B$ are independent events.

Conditional Distributions $\quad p(A|B), p(B|A)$

$$p(A) = \sum_B p(A,B), \quad (B) = \sum_A p(A)$$

$$p(A|B) = \frac{p(A,B)}{p(B)} = \frac{p(A,B)}{\sum_A p(A,B)}, \quad (B|A) = \frac{p(A,B)}{p(A)} = \frac{p(A,B)}{\sum_B p(A,B)}$$

Given $p(A|B), p(B)$, we can also find $P(B|A)$.

$$p(A|B) = \frac{p(A,B)}{p(B)} \implies p(A,B) = p(A|B)p(B)$$

$$p(B|A) = \frac{p(A,B)}{p(A)} = \frac{p(A,B)}{\sum_B p(A,B)} = \frac{p(A|B)p(B)}{\sum_B p(A|B)p(B)}$$

Joint Probability Density $\quad p(x,y)$

Marginal Probability Density $\quad p(x) = \int p(x,y)\,dy \quad, \quad p(y) = \int p(x,y)\,dx$

$$p(x|y) = \frac{p(x,y)}{p(y)} = \frac{p(x,y)}{\int p(x,y)\,dx}. \qquad p(y|x) = \frac{p(x,y)}{p(x)} = \frac{p(x,y)}{\int p(x,y)\,dy}$$

$$= \frac{p(y|x)p(x)}{\int p(y|x)p(x)\,dx} \qquad\qquad = \frac{p(x|y)p(y)}{\int p(x|y)p(y)\,dy}.$$

$\longrightarrow$ we can find $p(x|y)$
given that we know $p(y|x)$ and $p(x)$ !

# Maximum Likelihood Estimation

← parameter that we are trying to solve for our distribution.

Let $X$ be a discrete r.v. with pmf $p$ depending on parameter $\theta$.

$L(\theta | x) = P_\theta(x) = P_\theta(X = x)$ is the likelihood function, given the outcome $x$ of the r.v. $X$.

**For Bernoulli r.v. $X$,**

$$L(\theta | x) = P_\theta(x) = \theta^x (1-\theta)^{1-x} \quad \text{s.t. } P_\theta(X=1) = \theta \quad ; \quad P_\theta(X=0) = (1-\theta)$$
$$= 1 - P_\theta(X=1)$$

Suppose we have collected some data ... (bernoulli coin tosses)

i.i.d.

$$\text{data} = \{x_1, x_2, x_3, \dots, x_N\}$$

Likelihood is : $\quad L(\theta | \text{data}) = P_\theta(\text{data}) = \prod_{i=1}^{N} P_\theta(x_i) = \prod_{i=1}^{N} \theta^{x_i}(1-\theta)^{1-x_i}$

↳ Probability of observing the data that we observe assuming each coin toss was iid

$$\text{e.g. } L(\theta | x_1=1, x_2=0, x_3=1) = \theta^1 (1-\theta)^1 \theta^1$$

We want to find the value of $\theta$ that maximizes $L(\theta | \text{data})$!

i.e. the value of $\theta$ that makes the data that we collect most probable.

To Maximize $L(\theta | \text{data})$, we want to take its derivative w.r.t. $\theta$. $\quad$ s.t. $\frac{dL}{d\theta} = 0$

$$\hat{\theta} = \arg\max_{\theta} L(\theta).$$

Note that argmax is the operation that finds the argument that gives the max value from function $L(\theta)$.

Most of the time it is better to take the log of the likelihood before differentiating,

· and it usually leads to a simpler expression for the derivative that is easier to set to 0 and solve.

· since log is a monotonous function, i.e. whatever maximizes $L$ also maximizes $\log L$.

$$\ell(\theta) = L(\theta) = \log \prod_{i=1}^{N} \theta^{x_i}(1-\theta)^{1-x_i} = \sum_{i=1}^{N} \{x_i \log\theta + (1-x_i)\log(1-\theta)\}.$$
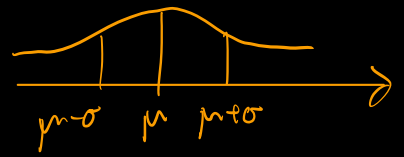
Set $\quad \frac{d\ell}{d\theta} = \frac{1}{\theta}\sum_{i=1}^{N} x_i - \frac{1}{1-\theta}\sum_{i=1}^{N}(1-x_i) = 0$

Solve for $\theta$
$$(1-\theta)\sum_{i=1}^{N} x_i = \theta \sum_{i=1}^{N}(1-x_i).$$
$$\sum_{i=1}^{N} x_i - \theta\sum_{i=1}^{N} x_i = N\theta - \theta\sum_{i=1}^{N} x_i$$
$$\Rightarrow \theta = \frac{1}{N}\sum_{i=1}^{N} x_i$$

Note that for Binomial dist.,
$$P(x=k) = \binom{N}{k}\theta^k (1-\theta)^{N-k}$$

For Gaussian (Normal) r.v.;

$$\text{data} = \{x_1, x_2, \ldots, x_N\} \quad \text{iid.}$$

$$L(\theta \mid \text{data}) = \prod_{i=1}^{N} p_\theta(x_i), \quad \text{where } \theta = \{\mu, \sigma^2\}$$

$$= \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x_i - \mu}{\sigma}\right)^2\right]$$

Goal :  Find parameters that best describe data collected.
These values will maximizes the likelihood

$$\hat{\mu}, \hat{\sigma}^2 = \arg\max_{\mu, \sigma^2} L(\mu, \sigma^2 \mid \text{data})$$

$$\ell(\mu, \sigma^2) = \log L(\mu, \sigma^2) = \log \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x_i - \mu}{\sigma}\right)^2\right]$$

$$= \sum_{i=1}^{N} \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x_i - \mu}{\sigma}\right)^2\right] = \sum_{i=1}^{N}\left[\log\frac{1}{\sqrt{2\pi\sigma^2}} + \log e^{-\frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2}\right]$$

$$= \sum_{i=1}^{N}\left[-\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2}\left(\frac{x_i-\mu}{\sigma}\right)^2\right] = \underbrace{-\frac{1}{2}\log(2\pi\sigma^2)\sum_{i=1}^{N}(1)}_{\text{Constant.}} - \frac{1}{2}\sum_{i=1}^{N}\left(\frac{x_i-\mu}{\sigma}\right)^2$$

$$\frac{\partial \ell}{\partial \mu} = \sum_{i=1}^{N}\left(\frac{x_i - \mu}{\sigma}\right)\frac{1}{\sigma}$$

Set $\frac{\partial L}{\partial \mu} = 0$, $\quad \frac{1}{\sigma^2}\sum_{i=1}^{N}(x_i - \mu) = 0$.

$$\sum_{i=1}^{N} x_i - N\mu = 0.$$

$$\hat{\mu} = \frac{1}{N}\sum_{i=1}^{N} x_i = f(x_1, x_2, \ldots, x_N).$$

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

$$\ell(\mu, v) = -\frac{N}{2} \log(2\pi v) - \frac{1}{2} \frac{1}{v} \sum_{i=1}^{N} (x_i - \mu)^2, \quad \text{where } v = \sigma^2$$

$$\frac{\partial \ell}{\partial v} = -\frac{N}{2} \frac{1}{v} - \frac{1}{2}(-1) \frac{1}{v^2} \sum_{i=1}^{N} (x_i - \mu)^2$$

Set $\frac{\partial \ell}{\partial v} = 0$, $\quad N = \frac{1}{v} \sum_{i=1}^{N} (x_i - \mu)^2$

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

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^{N} (x_i - \hat{\mu})^2 = g(x_1, x_2, \ldots, x_N).$$

Our MLE estimates are also random variables!

So, we can ask:
- What is their distribution?
- What is their expectation?

Can use characteristics to prove.

$$\hat{\mu} \sim N\left(\mu, \frac{\sigma^2}{N}\right)$$

$$E(\hat{\mu}) = \mu$$

$$E(\hat{\sigma}^2) \neq \sigma^2$$

$$E(\hat{\sigma}^2) = \frac{N-1}{N} \sigma^2 \to \sigma^2 \text{ as } N \to \infty.$$