

ML as a random experiment

- Let's make prediction whether a fish is a Salmon or a Mackerel by its size



$$X = \mathbb{R}$$

$$Y = \left\{ \begin{array}{c} 1 \\ \uparrow \\ \text{Mackerel} \end{array} , \begin{array}{c} 2 \\ \uparrow \\ \text{salmon} \end{array} \right\}$$

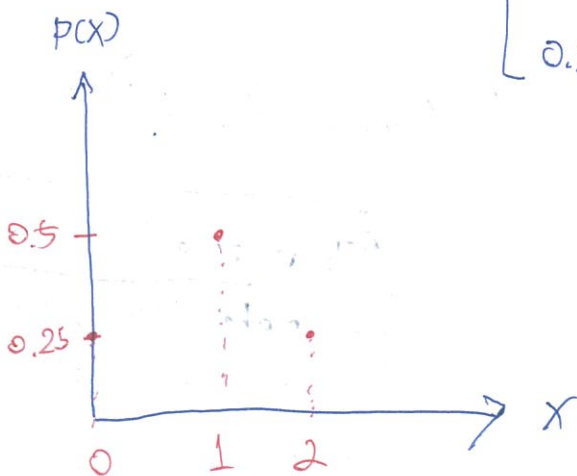
- Each sample  $(x_i, y_i) \sim P(\tilde{X}, \tilde{Y})$

Joint probability distribution  
which we have no access to

- 
- Probability Distribution: A function that gives the probability of occurrence of different outcomes

EX: 
$$P(X) = \begin{cases} 0.5 & \text{if } x = 1 \\ 0.25 & \text{if } x = 0 \\ 0.25 & \text{if } x = 2 \end{cases}$$

Here,  $X$  is a discrete R.V. denoting the number of heads observed in a double coin toss

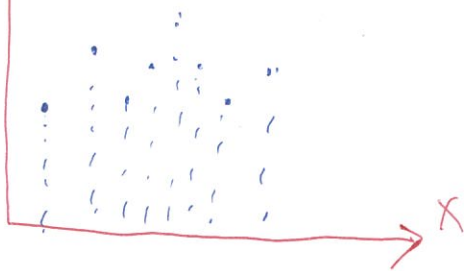


- Probability mass function: A function that gives the probability that a discrete R.V. is equal to some value

Probability density function: A function whose value at any given sample is a likelihood that the value of r.v. equal that sample

P.M.F.

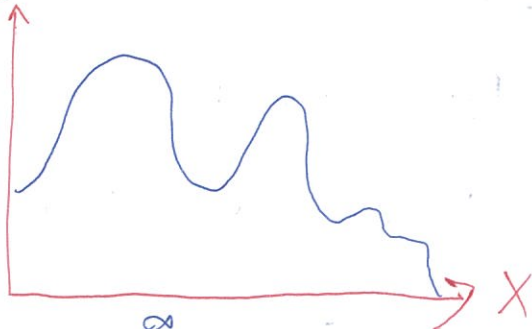
$f(x)$



$$\sum_i f(x_i) = 1$$

p.d.f.

$f(x)$



$$\int_{-\infty}^{\infty} f(x) dx = 1$$

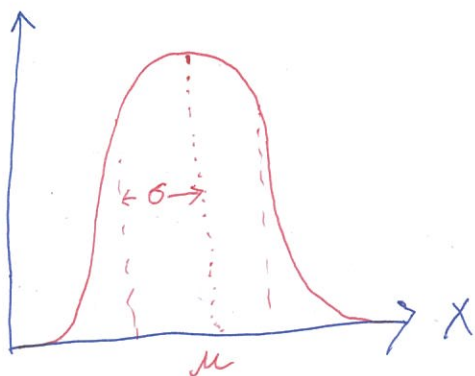
Our distribution  $P(X)$ :

- Consider one type of fishes,  $x_i \sim P(X)$

R.V.

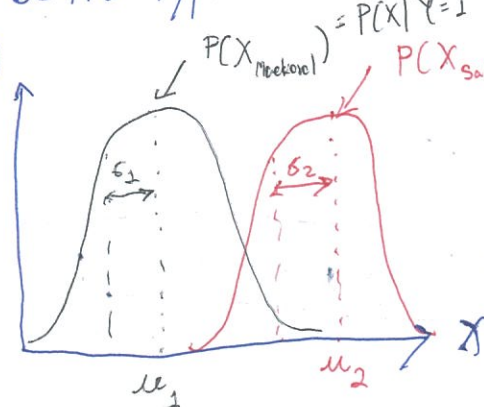
We shall assume a normal/Gaussian distribution

$P(X)$



- For both types of fishes, we would have

$P(X)$



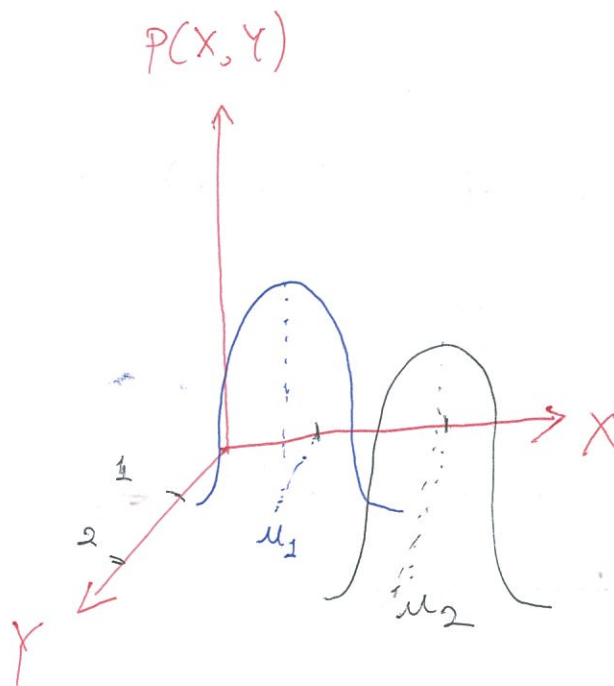
$P(X_{\text{Mackinac}}) = P(X|Y=1)$

$P(X_{\text{Salmon}}) = P(X|Y=2)$

$\Leftarrow P(X_{\text{Mackinac}}, X_{\text{Salmon}})$

- How to think of  $P(X, Y)$ ?

Joint probability distribution  
A probability distribution of at least two jointly r.v.

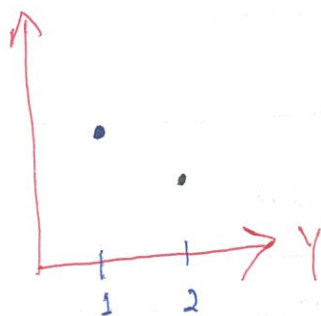


- Let's say if we have ~~no~~ access to  $P(X, Y)$ , then we could build an optimal classifier that make prediction by

$$h(x) = \underset{y}{\operatorname{argmax}} P(Y|x)$$

Bayes classifiers

Ex:  $P(Y|x)$  "score  $x$ "



$\Rightarrow$  optimal prediction:  $h(x) = 1$

- Conditional probability distribution:

Given two jointly r.v.,  $P(Y|x)$  is the distribution of  $Y$  when  $x$  is known to be a particular value

- Huge observation: If accessing  $P(X, Y)$  was possible, we then would have an optimal classifier.

## Two ML approaches:

① Approximate / estimate  $P(Y|x)$  directly

Discriminative learning

(E.g. k-NN, Perceptrons)

② Approximate / estimate  $P(X, Y)$ .

Then, apply Bayes classifier

Generative learning

Big- Q: How to learn the joint distribution from data?

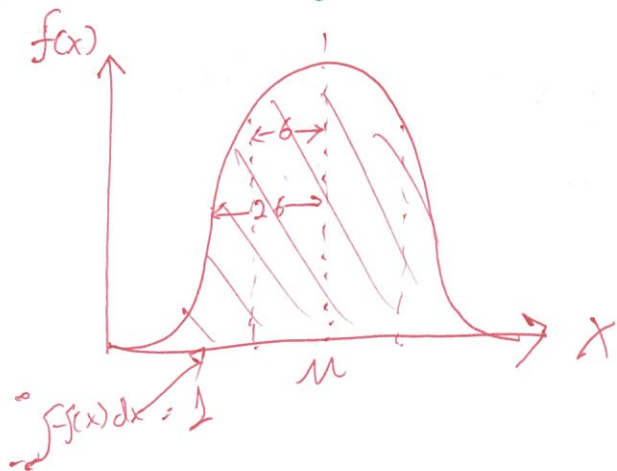
Chain rule for R.V.:

$$\begin{aligned} P(X, Y) &= P(X|Y) \cdot P(Y) \\ &= P(Y|X) \cdot P(X) \end{aligned}$$

## Estimating Probabilities from Data

- Goal: Build a distribution that models the real distribution (approximate)

- Ex. The modelling distribution is a Gaussian



$$\text{p.d.f. } f(x | (\mu, \sigma)) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

parameter vector of the modeling distribution



- For a Gaussian p.d.f.  $f(x|\mu, \sigma)$ , <sup>the</sup> ~~each~~ value of the function at any  $x$  is the likelihood (not probability, but can be viewed alike).

---

### Simple scenario I: Coin Toss

- Experiment: Toss a coin  $n$  times. How would we estimate  $P(H)$ ?

$$\text{Ex } D = \{H, T, T, H, T, H, H, T, T, T\}$$

$n = 10$

What is  $P(H)$ ?

- Intuitively,  $P(H) \approx \frac{n_H}{n_T + n_H} = \frac{4}{10} = 0.4$

- Let's try to formalize what we just did here

---

### Maximum Likelihood Estimation:

- ① Make an explicit modelling assumption about what type of distribution your data was sampled from,
- ② Set the parameters of this distribution so that the data you observed is as likely as possible.

## Coin Toss (cont.)

- Let  $\theta$  be a parameter of the modelling distribution

- For our coin toss scenario,

$$D = \{ \underset{x_1}{H, T, T}, \underset{x_2}{H, T, T}, \underset{x_3}{H, T, T, T} \}$$

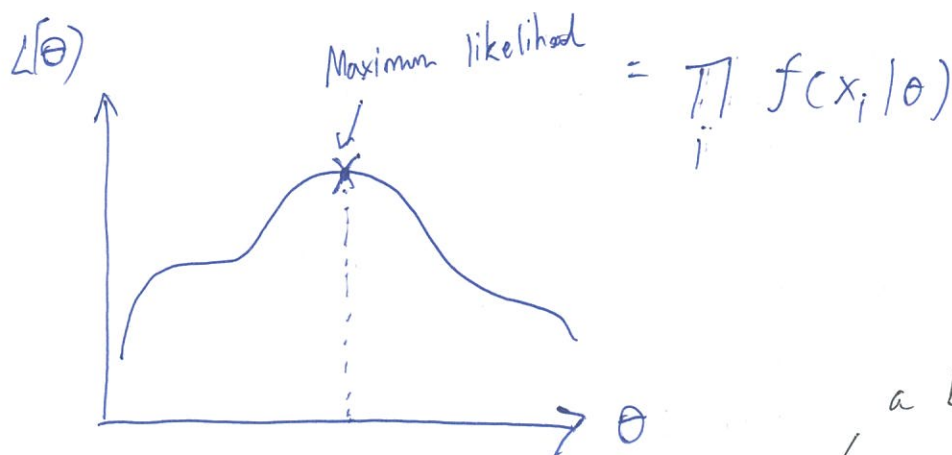
$x_1$  through  $x_n$  are a random sample from  $f(x|\theta)$   
i.i.d

- To measure the likelihood, we define the likelihood function

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

- ~~Therefore~~

$$= f(x_1 | \theta) \dots f(x_n | \theta)$$



a binomial distribution

- Let's assume

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

- Our goal is to find  $\theta_{ML} = \arg \max_{\theta} L(\theta, x)$