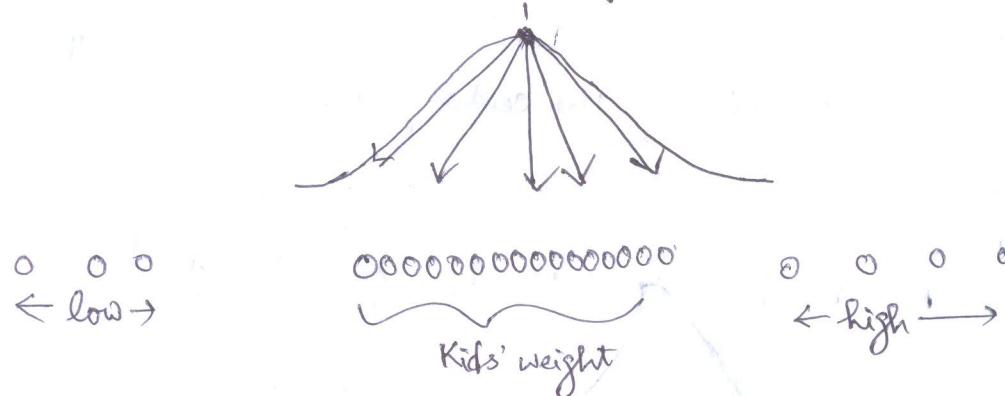


(5)

## MLE & MAP (ML Lecture 11)

(11)

Experiment: Weighing a bunch of nursery kids

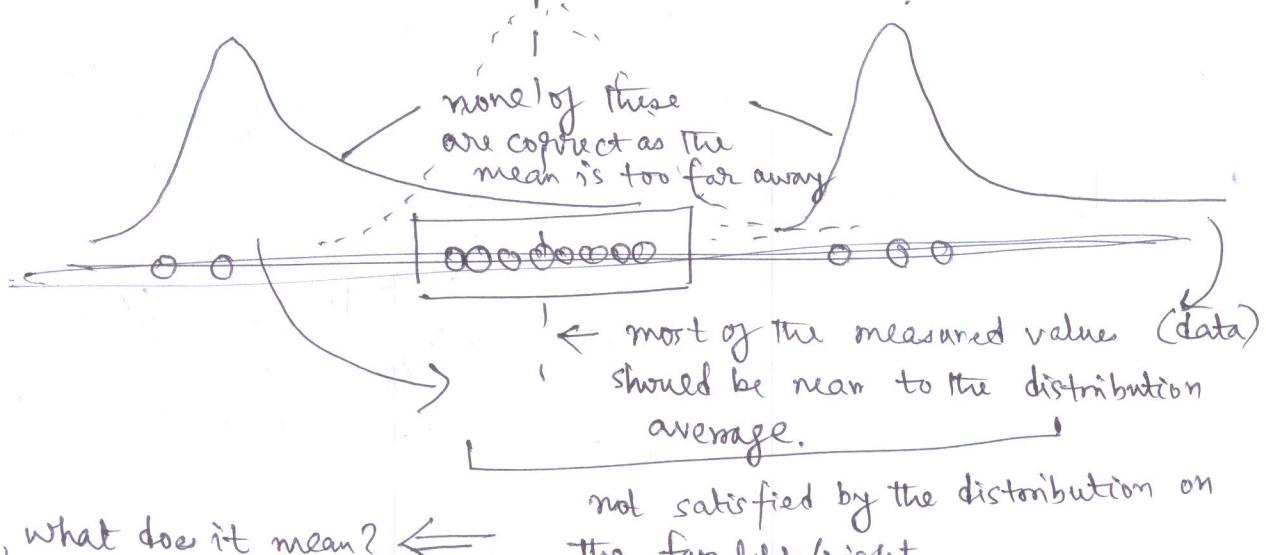


Goal: MLE fits distribution to data, in a maximum likely way. There exist diff distributions (normal, exponential, gamma & many more...)  
 ⇒ Convenient!

In this case, we assume the weights to be normally distributed. This implies

- most of the measurements to be close to ~~zero~~ the mean
- most of the measurements relatively symmetrical to the mean.

Once we settle on the shape, we have to settle the location of the distribution. Is one location better than the other?

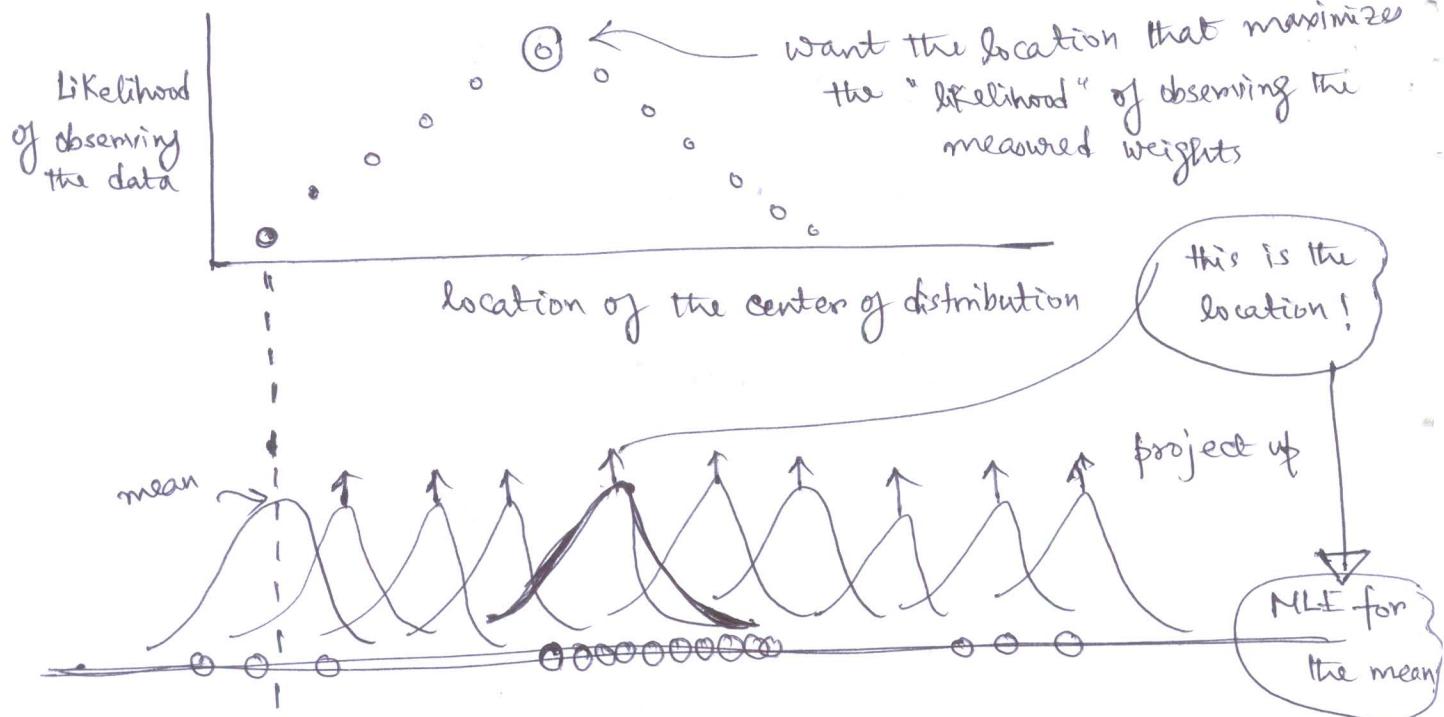


So, what does it mean? ↪

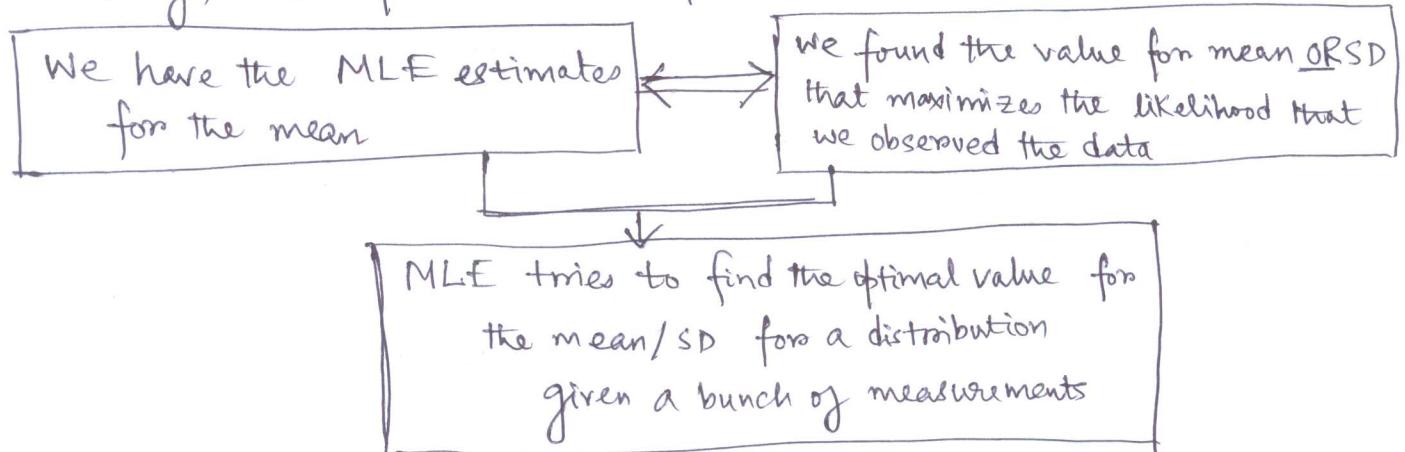
not satisfied by the distribution on the far left/right.



the likelihood of observing all the weights is low; however for the distribution at the center, the likelihood of observing most measurement is high.



Similarly, we compute the MLE for the standard deviation.



MLE for the Normal distribution:

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

location of  $\mu$       location of  $\sigma^2$

likelihood function

$$P(D|\theta) \Leftrightarrow P(\theta|D)$$

$$P(D|\theta) \Leftrightarrow P(\theta|D/x)$$

Location of  $\mu$ :

smaller  $\mu$  value  $\rightarrow$  moves the distribution to the left  
 larger  $\mu$  value  $\rightarrow$  - - - - - right

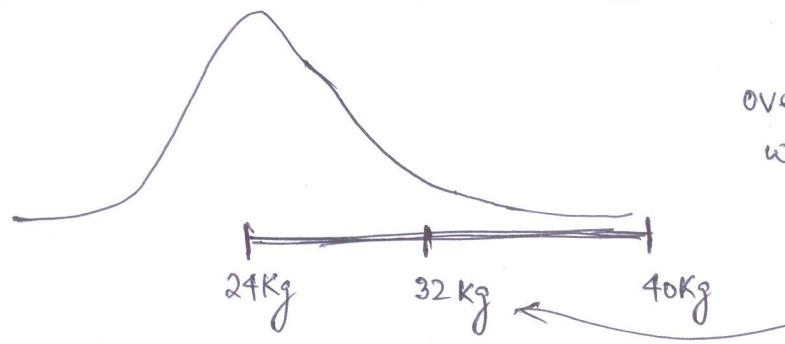
Location of  $\sigma$ :

Determines the width of the distribution

smaller  $\sigma$   
 distribution taller & narrower  
 (tall)

larger  $\sigma$   
 distribution fatter & wider  
 (bored)

(2)



overlay a normal distribution w/  $\mu = 28$  &  $\sigma = 2$  onto the data

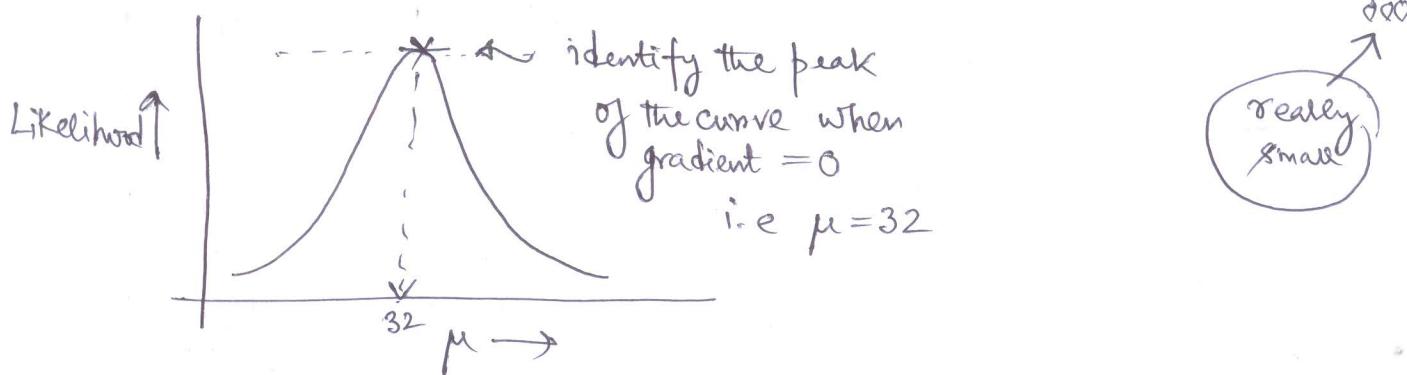
single measurement of weight

To determine the likelihood of the data given this distribution, we can plug the numbers in the likelihood function:  $\rightarrow L(\mu = 28, \sigma = 2 | x = 32)$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2} = \frac{1}{\sqrt{2\pi \cdot 2^2}} e^{-(32-28)^2/2 \cdot 2^2} = 0.03$$

Now, shift the distribution to the right i.e  $\mu = 30 \rightarrow L = 0.12$

Do this for a bunch of different distributions ( $\mu = 20 \rightarrow L() = 0.0000000003$ )



Fix 32 for  $\mu$  & vary  $\sigma$  for the appropriate distribution  $\Rightarrow$  need to do this for multiple data points. Therefore,

MLE for  $\mu$ :  $\sigma$  is constant, find the  $\mu$  for which slope = 0

MLE for  $\sigma$ :  $\mu$  is constant, —  $\sigma$  — slope = 0

More data:

$$\begin{cases} x_1 = 32 \text{ Kg} \\ x_2 = 34 \text{ Kg} \end{cases}$$

$$L(\mu = 28, \sigma = 2 | x_1 = 32) ; L(\mu = 28, \sigma = 2 | x_2 = 34)$$

$$x_1, x_2 \in D$$

But what's the likelihood when  $x_1 = 32$  &  $x_2 = 34$  i.e

$$L(\mu = 28, \sigma = 2 | x_1 = 32 \text{ &} x_2 = 34) = L_1 * L_2$$

$$L(D|D)$$

Should be written as  $L(D|\theta)$

$\Rightarrow$  IID (measurements are independent, weighing  $x_1$  doesn't have an effect on weighing  $x_2$ )

Keep adding data points ; i.e  $D = \{x_1, \dots, x_n\}$  then

$$L(\theta | D) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\underbrace{L(\theta | D)}_{L(\theta | D)} = L(\mu, \sigma | x_1, \dots, x_n) = L(\mu, \sigma | x_1) L(\mu, \sigma | x_2) \dots L(\mu, \sigma | x_n)$$

MLE estimation

derivative w.r.t  $\mu$  ( $\sigma$  constant)

derivative w.r.t  $\sigma$  ( $\mu$  constant)

$$\ln L(\theta | D) = \ln \left( \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right)$$

$$= \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_1 - \mu)^2}{2\sigma^2}} \right) + \dots + \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_n - \mu)^2}{2\sigma^2}} \right)$$

$$= -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x_1 - \mu)^2}{2\sigma^2}$$

$$= -\frac{1}{2} \ln(2\pi) - \frac{1}{2} * 2 \ln(\sigma) - \frac{(x_1 - \mu)^2}{2\sigma^2} + \dots$$

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

$$\frac{\partial L}{\partial \mu} = 0 \Rightarrow \sum_{i=1}^n \frac{(x_i - \mu)}{\sigma^2} \Rightarrow \sum_{i=1}^n \frac{x_i - \mu}{\sigma^2} = 0 \Rightarrow \mu = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$$

↑  
Sample mean

$$\frac{\partial L}{\partial \sigma} = 0 \Rightarrow -\frac{n}{\sigma} + \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^3} (\cancel{\sigma}) \sigma^{-3} = 0$$

$$\Rightarrow -\frac{n}{\sigma} + \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^3} = 0$$

$$\Rightarrow \sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \rightarrow \text{SD of measurements}$$

$$\Rightarrow \sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}}$$

MAP: Given,  $D = \{x_1, \dots, x_n | x_i \in \mathbb{R}^n\}$ ; assume a joint distribution  $p(D|\theta)$   
 $\theta$  is R.V; Goal: choose a "good"  $\theta \sim \text{for } D$  (hypothesis class)

$$\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} p(\theta|D) \rightarrow \text{posterior density}$$

Compare w/  
MLE

$$\theta_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} p(D|\theta) \rightarrow \text{likelihood function}$$

Pros: (i) easy to compute & interpret;  $p(D|\theta) = \underbrace{p(D|\theta)}_{\text{MLE}} \underbrace{p(\theta)}_{\text{prior}}$   
 MAP interpolates b/w  
 MLE & prior

(ii) avoids overfitting  $\rightarrow$  MLE does



Cons: (i) point estimate  $\Rightarrow$  no representation of uncertainty  
 (ii) must assume prior on  $\theta$



MAP for univariate Gaussian:

$$D = \{x_1, \dots, x_n\}, \theta \rightarrow \text{R.V.} \sim N(\mu, \sigma^2)$$

$\{x_1, \dots, x_n\}$  IID

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

$$\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

MLE for  $\theta \sim N(\mu, \sigma^2)$

Bayes Rule

$$= \underset{\theta}{\operatorname{argmax}} (p(D|\theta)p(\theta))$$

$$= \underset{\theta}{\operatorname{argmax}} \ln(p(D|\theta)p(\theta))$$

$$= \underset{\theta}{\operatorname{argmax}} (\ln p(D|\theta) + \ln p(\theta))$$

MLE for  $\mu$ :

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

$$\Rightarrow \frac{1}{\sigma^2} \sum_{i=1}^n x_i - n\mu = 0$$

$$\frac{\partial}{\partial \theta} (\ln p(D|\theta) + \ln p(\theta))$$

$$\Rightarrow \frac{1}{\sigma^2} \left( \sum_{i=1}^n x_i - n\theta \right) = 0$$

$$\begin{aligned} & \frac{\partial}{\partial \theta} \ln p(\theta) \\ &= \frac{\partial}{\partial \theta} \ln \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\theta-\mu)^2} \right) \\ &= \frac{\partial}{\partial \theta} \left( -\frac{1}{2} \ln 2\pi - \frac{1}{2}(\theta-\mu)^2 \right) \\ &= (\cancel{\theta}) \mu - \theta \end{aligned}$$

$$\therefore \frac{\partial}{\partial \theta} \left( \text{function} \right) = \frac{1}{\sigma^2} \left( \sum_{i=1}^n x_i - n\theta \right) + \cancel{\left( \frac{(\mu-\theta)}{\sigma^2} \right)} = 0$$

$$\Rightarrow (n+1)\theta = \frac{1}{\sigma^2} \sum_{i=1}^n x_i + \mu$$

$$\Rightarrow \theta = \frac{\frac{1}{\sigma^2} \sum_{i=1}^n x_i + \mu}{\frac{n}{\sigma^2} + 1} = \frac{n \bar{x} + \mu \sigma^2}{n + \sigma^2}$$

$$\boxed{\theta_{MAP} = \frac{n}{n+\sigma^2} \bar{x} + \frac{\sigma^2}{n+\sigma^2} \mu}$$

Convex combination of sample mean  $\mu$  & prior mean,  $\bar{x}$  (from prior data)

Convex combination :  $\alpha x + (1-\alpha) y ; \alpha \in (0,1)$ ,

$x \neq y$

$$\xrightarrow[\substack{\text{for all} \\ x \neq y \\ \alpha \in (0,1)}]{\quad} \begin{array}{ll} \alpha = 0, y \\ \alpha = 1, x \end{array}$$

$\therefore \theta_{MAP} \sim \bar{x}$  when  $\sigma^2 = 0$  (No variance)  
 $(x_{MLE}) \Rightarrow MLE$  is a special case of MAP

(i) As  $n$  &  $\sigma^2$  vary, we obtain a range of values b/w sample mean & prior mean ( $\mu$ )  
 $(\bar{x})$

(ii)  $n=0$  (No data),  $\theta_{MAP} \sim \mu$  (prior mean)

(iii)  $n \rightarrow \infty$ ,  $\theta_{MAP} \sim \bar{x}$  (sample mean)

Large # of data, MAP estimator is the sample mean  
 $\Leftrightarrow MLE$

(4)

MAP: An illustration

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

likelihood      prior  
↓                  ↓  
posterior

Recall the "cancer test" problem. What is MAP hypothesis?

Recall;

- (i) cancer  $\rightarrow$  + class ;  $\neg$  cancer  $\rightarrow$  - class ; Assume  $\theta \equiv h$
- (ii) Confusion matrix :

Actual	$h$	Prediction	
		-	+
-	0.97	0.03	
+	0.02	0.98	

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

(iii)

$$P(h-) = 0.992$$

$$P(h+) = 0.008$$

$$P(-|\neg \text{cancer}) = 0.97$$

$$P(-|\text{cancer}) = 0.02$$

$$P(+|\neg \text{cancer}) = 0.03$$

$$P(+|\text{cancer}) = 0.98$$



Now,

$$h_{\text{MAP}} = \underset{h}{\operatorname{argmax}} \quad P(D|h) P(h)$$

prob(cancer / tested positive)

$$P(+|\text{cancer}) P(\text{cancer}) = 0.98 \times 0.008 = 0.00784 = P(h+|D)$$

$$P(+|\neg \text{cancer}) P(\neg \text{cancer}) = 0.03 \times 0.992 = 0.02976 = P(h-|D)$$

prob( $\neg$ cancer / tested positive)

$$h_{\text{MAP}} = 0.02976$$

Given data, D; choose a good hypothesis / class label for D.

$h+ \equiv \text{Cancer}$  } two hypotheses ;  
 $h- \equiv \neg \text{cancer}$

For each hypothesis in  $H$ , compute the posterior probability

$$P(h+|D) = \frac{P(D|h+) P(h+)}{P(D)} ; \quad P(h-|D) = \frac{P(D|h-) P(h-)}{P(D)}$$

$$= 0.00784 \qquad \qquad \qquad = 0.02976$$

Output the hypothesis  $h_{\text{MAP}}$  with the highest posterior probability

$$h_{\text{MAP}} = \underset{h \in H}{\operatorname{argmax}} \{P(h+|D), P(h-|D)\} = P(h-|D)$$

## Note on MLE & MAP:

1. MAP can be interpreted as forming & changing opinion as an evolutionary process. You have an opinion belief that decides in forming your initial opinion until that gets modified by new evidence you gather from interactions (likelihood on Data) from data.

Ex: choosing a candidate (swing vote)

if you are an ideologue, there is no likelihood estimation interacting with your past beliefs to change your posterior prediction. There, the likelihood function plays "no role" (i.e no data i.e  $\bar{x} = 0$ ) i.e

$$\theta_{MAP} \sim \mu \text{ (prior mean)}$$

2. In the absence of prior, MLE is used
3. In the case of flat prior,  $MAP \sim MLE$
4. Use MLE if you have lot of data!
5. In the absence of data, just use prior mean.

Talking points:

(i) MLE of mean of a Gaussian :  $\mu_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$   
 MLE of std. dev. :  $\sigma_{MLE} = \left( \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \right)^{1/2}$

(ii) MLE on multiple linear regression (refer to Raj Jain's slides)

$$y = Xb + \varepsilon; \quad \begin{array}{l} \text{regression} \\ \downarrow \\ \text{Coefficients} \\ \text{to determine} \end{array} \quad \begin{array}{l} \text{(Gaussian Noise)} \\ \text{Datamatrix} \end{array}$$

$$b = (X^T X)^{-1} X^T y \\ = \text{pseudo inverse}$$

OK, let's look at the general LS problem:  $\min \sum_{i=1}^n (y_i - a_0 x_i - a_1)^2$

$\underbrace{\varepsilon}_{\text{residuals}}$

Write the residuals as n-vectors

$$\begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} - \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix}$$

$$\Rightarrow \varepsilon = y - Xb; \quad \text{in vectorial representation}$$

Therefore, the multiple regression problem i.e. LS minimization of the residuals become

$$\|\varepsilon\|^2 = \varepsilon^T \varepsilon = (y - Xb)^T (y - Xb)$$

Objective: Min  $(y - Xb)^T (y - Xb)$  wrt "b"

$$\begin{aligned} \frac{\partial}{\partial b} \{ (y - Xb)^T (y - Xb) \} &= X^T y - 2X^T Xb + y^T X \\ &= \frac{\partial}{\partial b} \{ y^T y - 2b^T X^T y + b^T X^T Xb \} = 2X^T (y - Xb) = 0 \\ &\Rightarrow -2X^T y + 2X^T Xb = 0 \\ &\Rightarrow y - Xb = X^T Xb = X^T y \\ &\Rightarrow b = (X^T X)^{-1} X^T y \end{aligned}$$

$$\begin{aligned} (y - Xb)^T (y - Xb) &= ((Xb)^T + y^T)(y - Xb) \\ &= b^T X^T y + b^T X^T Xb + y^T y - y^T Xb \\ &= y^T y - 2b^T X^T y + b^T X^T Xb \end{aligned}$$

MLE for Gaussian noise  
 $b = (X^T X)^{-1} X^T y$

regression coefficients

Ex:2: Given, samples  $\{0, 1, 0, 0, 1, 0\}$  from a binomial distribution  
 $p(x=0) = 1-\mu$ ,  $p(x=1) = \mu$ ;  $\hat{\mu}_{MLE} = ?$

$$L(\mu) = p(x=0)p(x=1)p(x=0)p(x=1)p(x=0)$$
$$= (1-\mu)^4 \mu^2$$

$$\log L(\mu) = 4 \ln(1-\mu) + 2 \ln \mu$$

$$\frac{1}{\mu} \frac{\partial L}{\partial \mu} = -\frac{4}{1-\mu} + \frac{2}{\mu}$$

$$\Rightarrow \frac{\partial L}{\partial \mu} = \mu \left( -\frac{4}{1-\mu} + \frac{2}{\mu} \right) = 0 \Rightarrow \frac{2}{\mu} = \frac{4}{1-\mu}$$

$$\Rightarrow 1-\mu = 2\mu$$

$$\Rightarrow \boxed{\mu = 1/3}$$

## Naive Bayes

Goal: To build a classifier that says whether a text is about sports or not. Training data has 5 sentences:

Data	Label
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

Test: A very close game     $\left[ \begin{array}{l} \text{Sports} \\ \text{Not sports} \end{array} \right]$  (?)

Notes - Naive Bayes is a probabilistic classifier and is a special case of Bayes' Rule. We would like to compute  $P(\text{Sports} | \text{a very close game})$  &  $P(\text{Not sports} | \text{a very close game})$  & pick the one w/ the larger value.

Step I: Feature Engineering  $\rightarrow$  extract useful features from data for use in the classification problem & remove the ones w/ no bearing on the model (such as identifier information) Text as data  $\Rightarrow$  "word frequencies"  $\Rightarrow$  ignore word orders & sentence construction, treating every document as a set of the words. So, features  $\equiv$  counts of these words.

Step II: Rewrite Bayes Rule in the context of the problem:

$$P(\text{Sports} | \text{a very close game}) = \frac{P(\text{a very close game} | \text{Sports}) \times P(\text{Sports})}{P(\text{a very close game})}$$

$P_1$

$$\& P(\text{Not sports} | \text{a very close game}) = \frac{P(\text{a very close game} | \text{Not sports}) \times P(\text{Not sports})}{P(\text{a very close game})}$$

$P_2$

i.e. discard the denominators as they are same! We need to compute

$$P_1 \sim P(\text{a very close game} | \text{sports}) \times P(\text{sports}) \quad \&$$

$$P_2 \sim P(\text{a very close game} | \neg \text{sports}) \times P(\neg \text{sports})$$

NOTE:

"a very close game" doesn't appear in the training data.

↳ Naïve handling of the above problem

$$P(\text{a very close game}) = P(a) \times P(\text{very}) \times P(\text{close}) \times P(\text{game})$$

$\underbrace{\hspace{10em}}$  strong assumption

Therefore,

$$\Rightarrow P(\text{a very close game}) = P(a) * P(\text{a} | \text{sports}) \times P(\text{very} | \text{sports})$$

$\underbrace{\hspace{10em}}$  \*  $P(\text{close} | \text{sports}) * P(\text{game} | \text{sports})$

Similarly,

$$\Rightarrow P(\text{a very close game} | \neg \text{sports}) = P(a | \neg \text{sports}) * P(\text{very} | \neg \text{sports})$$

$\underbrace{\hspace{10em}}$  \*  $P(\text{close} | \neg \text{sports}) * P(\text{game} | \neg \text{sports})$

Compute These

Step III:  $P(\text{sports}) = 3/5$ ,  $P(\neg \text{sports}) = 2/5$ ; a priori

$$P(\text{game} | \text{sports}) = \frac{\text{word frequency of "game" in sports}}{\text{total # of words in sports}} = \frac{2}{11}$$

$$P(\text{close} | \text{sports}) = 0 ! \rightarrow \text{how do we handle this} \Rightarrow$$

add 1 to every count so that it's never zero,  $\rightarrow$  to balance this we add the # of possible words to the divisor, so that the ratio  $\leq 1$ . In this case, # of possible words = 14

(2)

$$P(a| \text{sports}) = \frac{2+1}{11+14} ; \quad P(a| \text{7sports}) = \frac{1+1}{9+14}$$

$$P(\text{very}/\text{sports}) = \frac{1+1}{11+14} ; \quad P(\text{very}/\text{7sports}) = \frac{0+1}{11+14}$$

$$P(\text{close}/\text{sports}) = \frac{0+1}{11+14} ; \quad P(\text{close}/\text{7sports}) = \frac{1+1}{9+14}$$

$$P(\text{game}/\text{sports}) = \frac{2+1}{11+14} ; \quad P(\text{game}/\text{7sports}) = \frac{0+1}{9+14}$$

$$\therefore p_1 = 0.0000276 ; \quad p_2 = 0.00000572$$

↙  $p(\text{sports}/\text{a very close game}) > p(\text{7sports}/\text{a very close game})$

↓

Test data Label : sports
--------------------------