

17/9/18

Lec-12

Parameter estimation:-

• classifiers were built on the assumption that we know the distro.

But in fact, we don't exactly know the distro.

and

In ~~regresssion~~ linear classifier / regression, we found how to update w (which is learning).

So, given a set of data points, how sure am I that the model works?

A Parameter Estimation.

Before learning, we assume that the distro is normal with elliptical contours.

If we want to fit a model to a data, we want to find parameters that will help me in predicting new points.

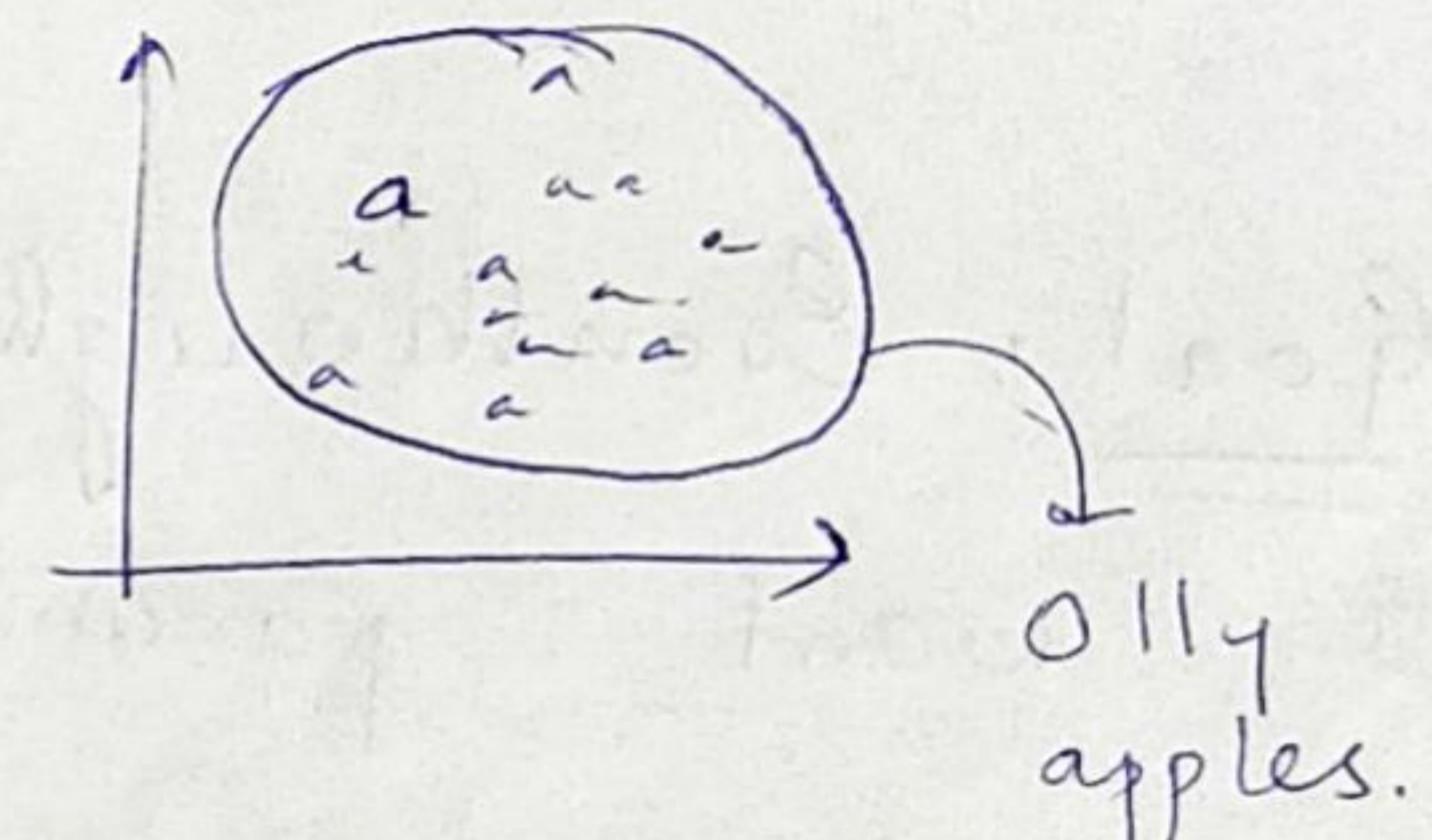
Some - terms

Population:

Every apple I could observe (if world = ∞)

Every apple in Gachibowli (if world = Gachibowli)

Prob: We can never know the whole population.

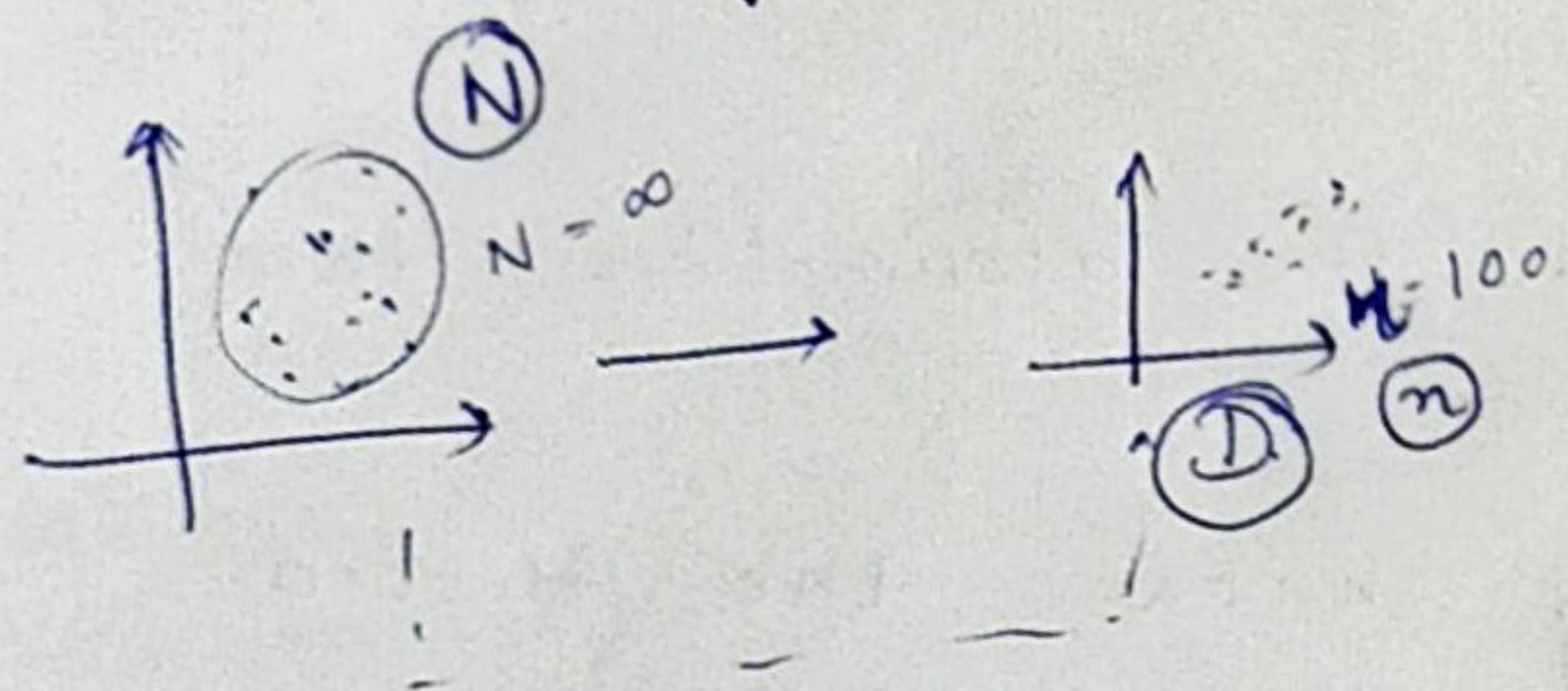


Olly apples.

Sampling:- Take a set of random points from the entire population. It is important, so that, we have no biases.

→ The population may ~~not~~ have any distribution

but sampling is uniform. So, the sampled distro also follows similar distro as original one



So sample + Params

→ Params

In normal distro:

$$\text{Params: } \mu, \Sigma \xrightarrow{\text{+ Sample}} \mu = \begin{bmatrix} a \\ b \end{bmatrix}, \Sigma = \begin{bmatrix} \sigma_1 & \sigma_{12} \\ \sigma_{21} & \sigma_2 \end{bmatrix}$$

Goal: Even though I am looking at a sample.
I want parameters over the sample

Two approaches:-

$$\Rightarrow \text{Compute } \theta : [\theta_1, \theta_2, \dots, \theta_n] \\ \omega : [\omega_1, \omega_2, \dots, \omega_n] \rightarrow \text{line}$$

Compute θ that explains D .

Maximum likelihood estimate: (MLE)

Find θ so, we observe max of D .

$$\hat{\theta}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} \ p(D|\theta)$$

If I predict some, μ and Σ , how good is the normal distro.

2) Bayes' style: [Maximum Posterior Estimate]
use Bayes theorem to find θ

$$p(\theta|D) = \frac{p(\theta|D) \times p(D)}{p(\theta)}$$

$$\Rightarrow \hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} [p(\theta|D) \times p(D)]$$

- If $p(\theta) = \text{constant} \Rightarrow$ No prior info about $\theta \Rightarrow \text{MLE} == \text{MAP}$
- Bayesian gives nothing with certainty
- So $\hat{\theta}$ is also a density fn. with some distro

MLE:

$$\underset{\theta}{\operatorname{argmax}} P(D|\theta) = \prod_{k=1}^n p(x_k|\theta)$$

If samples are independent

taking log won't affect

$$\Rightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log \left(\prod_{k=1}^n p(x_k|\theta) \right)$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_{k=1}^n \ln p(x_k|\theta)$$

To maximise, find $\frac{\partial L}{\partial \theta} = 0$

$$\frac{\partial L}{\partial \theta} = \sum_{k=1}^n \frac{\partial}{\partial \theta} \ln(p(x_k|\theta)) = 0$$

So each comp has to be zero.

Case 1: $N(\mu, \Sigma)$: Σ is known

$$\text{Likelihood: } \ln p(x_k|\mu) = \ln \left(\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x_k - \mu)^T \Sigma^{-1} (x_k - \mu)} \right)$$

$$= \ln(\text{stuff}) - \frac{1}{2} (x_k - \mu)^T \Sigma^{-1} (x_k - \mu)$$

$$\frac{\partial}{\partial \mu} \ln p(x_k|\mu) = \frac{\partial}{\partial \mu} \left(-\frac{1}{2} (x_k - \mu)^T \Sigma^{-1} (x_k - \mu) \right)$$

$\Sigma \rightarrow \text{Covariance}$

$$= -\frac{1}{2} \times 2 \Sigma^{-1} (x_k - \mu) \times -1$$

Substitute $\mu = \hat{\mu}_k$ as we know the best estimate $\hat{\mu}$

Satisfies $\underset{\mu}{\operatorname{argmax}} \ln p(x_k|\mu) = 0$

$\rightarrow \frac{\partial}{\partial \mu} \ln p(x_k|\mu) = 0$

Removing Σ

$$\frac{\partial}{\partial \mu} \ln p(x_k|\mu) =$$

$$\Rightarrow \sum_{k=1}^n (x_k - \hat{\mu}) = 0$$

$$\sum_{k=1}^n x_k - \sum_{k=1}^n \hat{\mu} = 0$$

$$\hat{\mu}_{\text{MLE}} = \frac{\sum_{k=1}^n x_k}{n}$$

Average of samples, is the best estimate of μ

Step 2

$$N(\mu, \sigma^2) \quad \boxed{N(\theta_1, \theta_2)} \\ \mu = \theta_1, \quad \sigma^2 = \theta_2$$

while ∇_{μ} , \Rightarrow same as Step 1: Assume const $\sigma^2 \Rightarrow \hat{\mu} =$

$$\nabla_{\sigma^2}$$

$$\nabla_{\theta_2} \ln(x_k/\theta) = \begin{bmatrix} \nabla \theta_1 \\ \nabla \theta_2 \end{bmatrix} \ln(x_k/\theta) = \begin{bmatrix} \text{same as} \\ \text{step 1} \\ \downarrow \\ \text{need to find} \end{bmatrix}$$

$$\nabla_{\theta_2} = \nabla \left(\ln \left[\frac{1}{\sqrt{2\pi\theta_2}} \right] - \frac{(x_k - \theta_1)^2}{2\theta_2} \right)$$

$$= \nabla \left[-\frac{1}{2\theta_2} - \frac{1}{2} \ln \theta_2 - \frac{(x_k - \theta_1)^2}{2\theta_2} \right]$$

$$\nabla_{\theta_2} = -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2}$$

$$\Rightarrow \sum_{k=1}^n \left[-\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} \right] = 0 \quad \text{at } \theta_2 = \hat{\theta}_2$$

$$\Rightarrow \sum_{k=1}^n \left[-\frac{1}{2\hat{\theta}_2} + \frac{(x_k - \theta_1)^2}{2\hat{\theta}_2^2} \right] = 0$$

$$\frac{-n}{2\hat{\theta}_2} + \sum_{k=1}^n \frac{(x_k - \theta_1)^2}{2\hat{\theta}_2^2} = 0$$

$$\Rightarrow \frac{1}{2\hat{\theta}_2^2} \sum_{k=1}^n (x_k - \theta_1)^2 = \frac{n}{2\hat{\theta}_2}$$

$$\Rightarrow \sum_{k=1}^n (x_k - \theta_1)^2 = \frac{n}{2} \times 2\hat{\theta}_2$$

$$\hat{\theta}_2 = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2 \Rightarrow \hat{\sigma}_2^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2$$

Hence, we proved

$$\begin{Bmatrix} \mu \\ \sigma^2 \end{Bmatrix}$$

of popul for a normal distro

$$\begin{Bmatrix} \mu \\ \sigma^2 \end{Bmatrix}$$

sample

11th for a multivariate gaussian $\Rightarrow \bar{\mu}, \Sigma$

As sampling changes $\Rightarrow \hat{\mu}$ changes.

What will happen if many $\hat{\mu}$ are averaged.

Will $\text{mean}(\hat{\mu}) = \mu_{\text{global}}$?

i.e. Will $\boxed{\mu_{\text{global}} = E(\hat{\mu})}$ unbiased

If small sample size $\Rightarrow \hat{\mu}$ jumps a lot.

If large $\Rightarrow \hat{\mu}$ remains same

and $\sigma_{\text{glob}}^2 / \text{Var}(\hat{\mu}) \propto \frac{1}{n}$

as $n \rightarrow \infty$, $\text{Var} \rightarrow 0$, so it peaks

Will $\hat{\sigma}_{\text{global}}^2 = \frac{\hat{\sigma}^2}{n} = E(\hat{\sigma}^2)$?

No, $\frac{\hat{\sigma}^2}{n}$ is not

In fact $E(\hat{\sigma}^2) = \frac{n-1}{n} \sigma_{\text{glob}}^2 \Rightarrow E(\hat{\sigma}^2) < \sigma_{\text{glob}}^2$ always Biased

So MLE $\hat{\sigma}^2$ is a biased estimate, [as we estimate smaller]

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \mu)^2$$

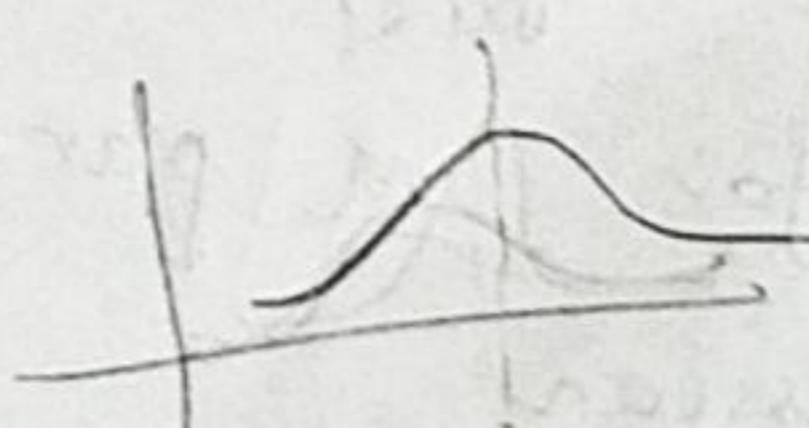
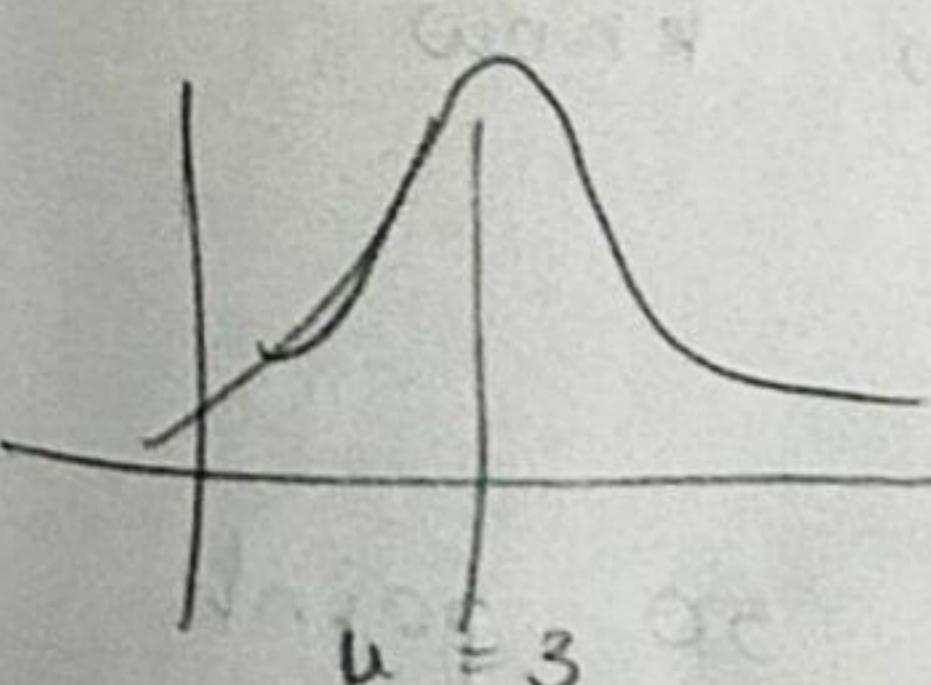
$$E(\hat{\sigma}^2) = \frac{n-1}{n} \times \frac{n-1}{n} \times \sigma^2$$

There are 3 D.H.S

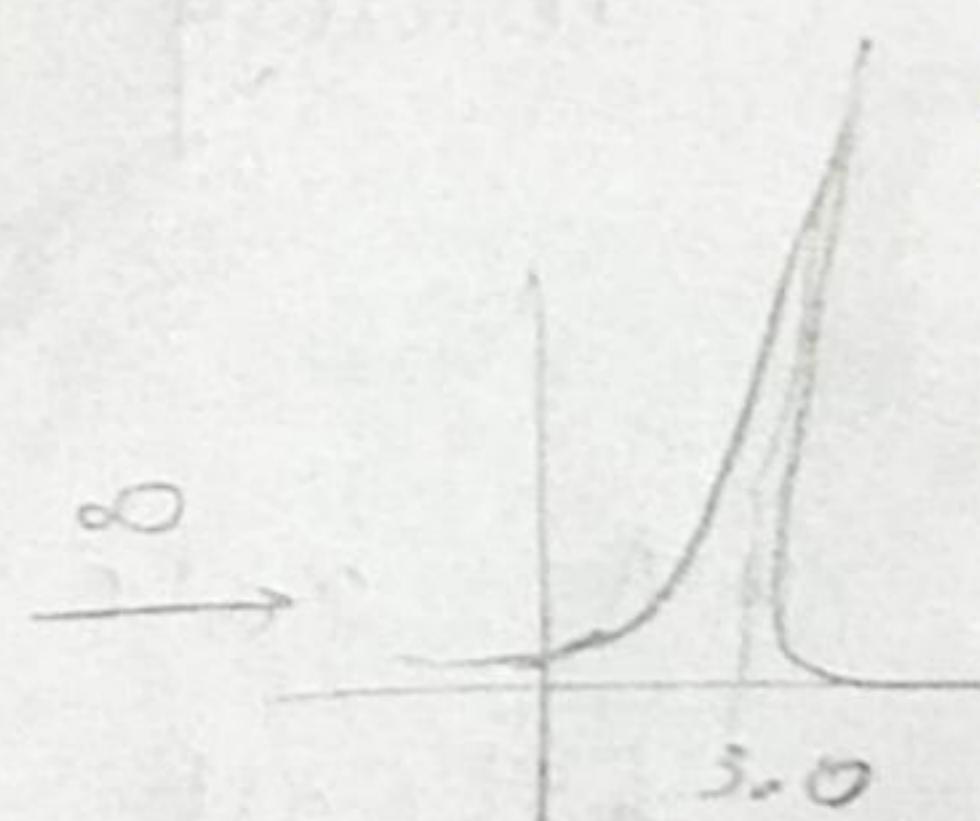
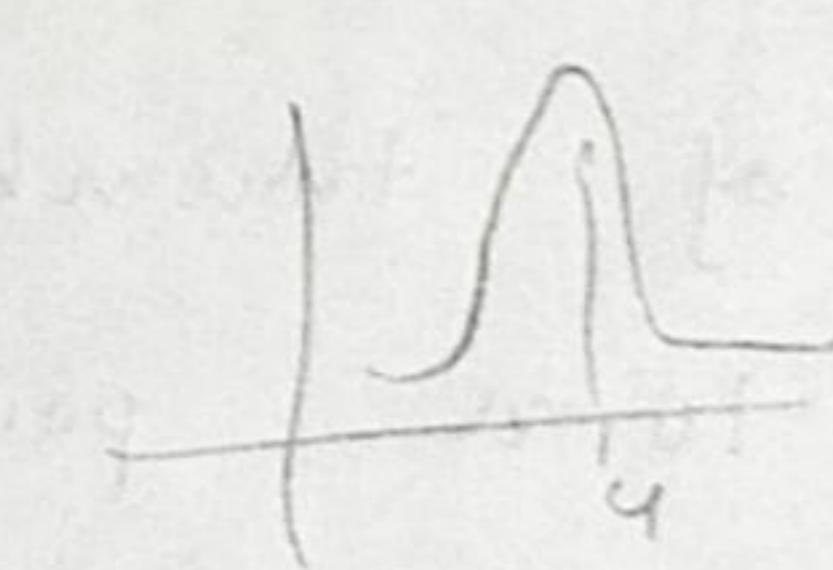
MAP estimate:

$p(\theta)$ is a distro and so, has μ_0, σ_0

Only one $\theta \Rightarrow \mu$



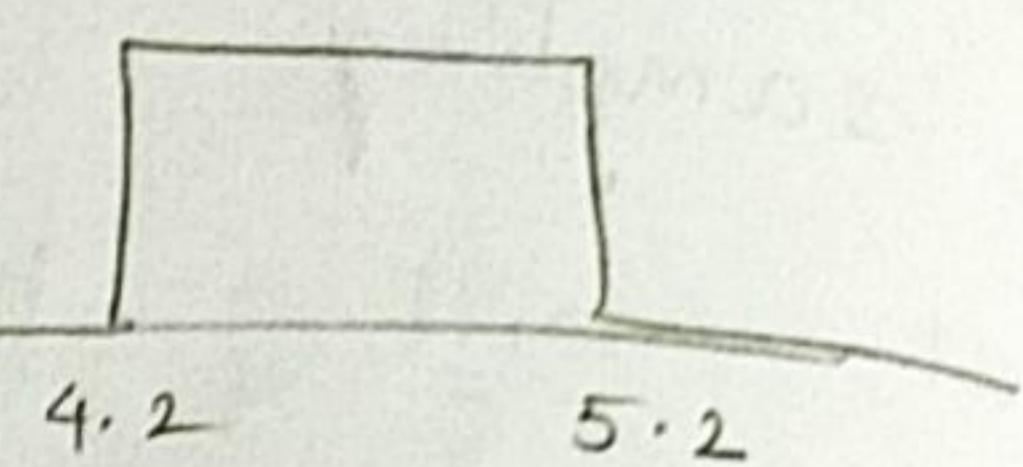
after many samples



Bayesian learning: Process of updating beliefs

What if $p(\theta) \rightarrow 4.2$ to $5.2 \Rightarrow$ uniform

If prior density is zero at a point, it will never be zero posterior



20/9/18

Lec-13

• Bayesian Parameter Estimation:-

$$p(w_i|x) = \frac{P(x|w_i)p(w_i)}{\int p(x)} \quad | \text{Bayes rule}$$

generative model

• Now, we need to know $P(x|w_i)$.

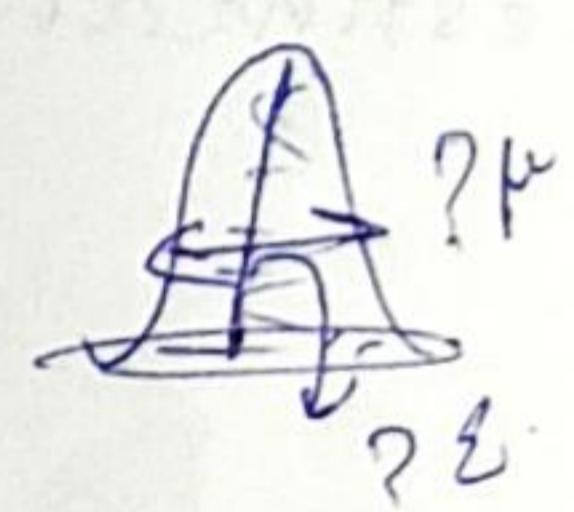
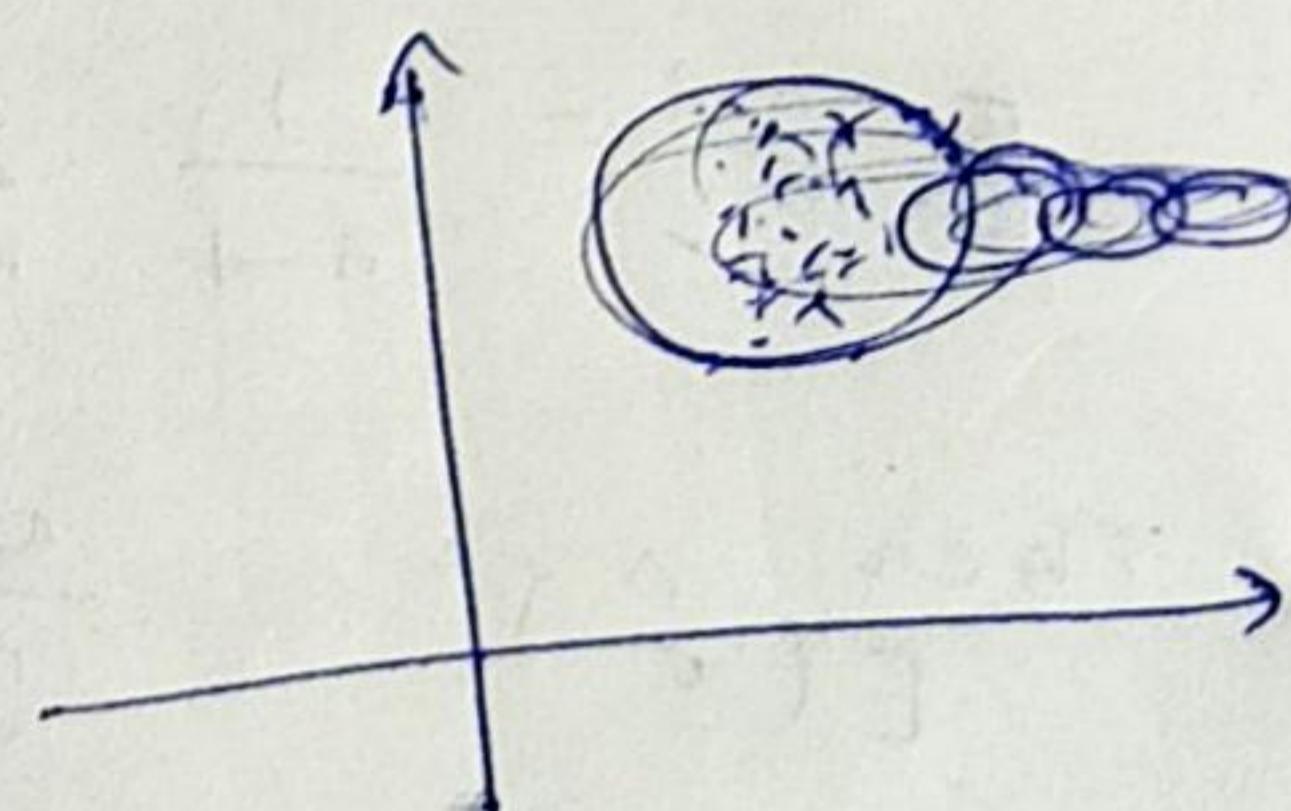
• If no assumptions are made on ~~of the~~ the probability distro, we can't use it.

• Many density estimate functions, need to know the form of the distro

• Here we are assuming the distro ~~too~~ is ~~the~~.

Normal distro.

[So, at first we needed ∞ params, (no idea of density at all) to normal distro



$$\begin{array}{c} \mu \\ \downarrow \\ \mu_1 \mu_2 \\ \downarrow \downarrow \\ \alpha_1 \alpha_2 = \alpha_1 \alpha_2 \end{array} \quad \underbrace{\quad}_{5 \text{ params}}$$

• A rule of thumb is you need to know atleast 10 n points for m params.

\Rightarrow 2D gaussian

\hookrightarrow Use atleast 50 points

- Classifier does not work properly
 - ↳ fn-form is not valid
 - ↳ fn-form is good, but samples are so less that improper params
- So when, we say Bayesian estimation is but, we ~~can~~ skip the above problems.

Let training data = $\{x_1, x_2, \dots, x_n\}$

$$p(x) \sim N(\mu, \sigma^2)$$

$$p(\theta) = \arg \max_{\theta} p(D|\theta)$$

$$= \arg \max_{\theta} \prod_{k=1}^n p(x_k|\theta)$$

$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{k=1}^n x_k, \quad \hat{\sigma}_{ME} = \sqrt{\frac{1}{n} \sum (x_k - \hat{\mu})^2}$$

MAP: estimate

$$\hat{\theta}_{map} = \arg \max_{\theta} p(D|\theta) \propto p(\theta) \quad \text{Is a density now.}$$

$$p(\theta|D) =$$

$$p_\theta(x|w_i) \sim N(\mu, \sigma^2)$$

$p(\mu)$ is a density fn.

μ is a random variable.

$$p(\mu) \sim N(\mu_0, \sigma_0^2)$$

$$p(\mu|D) = \frac{p(D|\mu) \times p(\mu)}{p(D)}$$

$$p(\mu) \sim N(\mu_0, \sigma_0^2) \rightarrow \text{Uncertainty in guess}$$

↳ Our best guess

μ_0, σ_0^2 won't affect, as we increase samples: $p(\mu)$ will go to

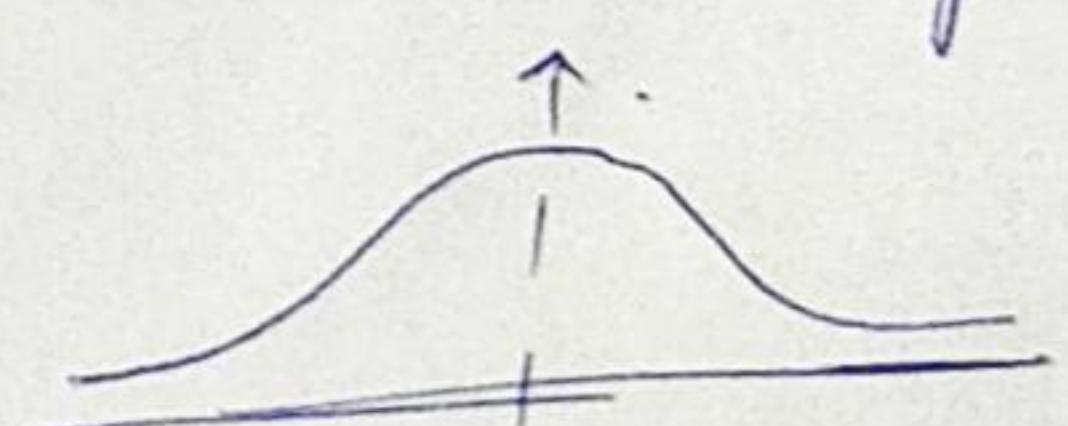
An example

$$p_\theta(x|w_i) \sim N(\mu, \sigma^2)$$

$\mu = 250, \sigma = 50$

$$p(\mu) \sim N(\mu_0, \sigma_0^2)$$

↓
Uncertainty
in
avg. wt //



$$p(\mu|D) = p(D|\mu) \times p(\mu)$$

correct μ

Case 1 $p(\mu|D)$, $p(x|\mu) \sim N(\mu, \sigma^2)$
 $p(\mu) \sim N(\mu_0, \sigma_0^2)$

$$p(\mu|D) = \frac{p(D|\mu) \cdot p(\mu)}{\int p(D|\mu) \times p(\mu) d\mu} \rightarrow \text{Const} = \alpha$$

$$p(\mu|D) = \alpha \prod_{k=1}^n g(x_k|\mu) \times \frac{p(\mu)}{\downarrow}$$

\downarrow A normal here \downarrow A normal here

Posterior is also a normal density

• Reproducing Density: If prior x $\xrightarrow{\text{posterior distro } a}$ posterior distro a

Both are same \Rightarrow

\Rightarrow The distro a is

$$p(\mu|D) = \alpha \prod_{k=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x_k-\mu}{\sigma})^2} \times \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{1}{2}(\frac{\mu-\mu_0}{\sigma_0})^2}$$
$$= \alpha \prod_{k=1}^n \frac{1}{\sqrt{2\pi}\sqrt{\sigma\sigma_0}} \times e^{-\frac{1}{2}\left\{(\frac{x_k-\mu}{\sigma})^2 - (\frac{\mu-\mu_0}{\sigma_0})^2\right\}}$$

$$p(\mu|D) = \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{1}{2}\left(\frac{\mu-\mu_n}{\sigma_n}\right)^2}$$

So, before observing we had $p(\mu) \sim N(\mu, \sigma^2)$; now
after (posterior) $p(\mu|D) \sim N(\mu_n, \sigma_n^2)$

where $\mu_n = \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$

Average of samples we observed

[MLE estimate]

$\mu_0 \rightarrow$ Before observing

$\hat{\mu}_n \rightarrow$ Only observing

Bayes says : use weighted μ_0 and $\hat{\mu}_n$

$\hat{\mu}_n$ as $n \rightarrow \infty$, $\hat{\mu}_n = \hat{\mu}_n$
 and if higher σ \Rightarrow higher n is needed
 meaning, if varied data \Rightarrow more samples
 are required

$$\sigma_n^2 = \frac{\sigma_0^2 + \sigma^2}{n\sigma_0^2 + \sigma^2} \quad \text{as } n \rightarrow \infty \quad \text{so, } \sigma_n^2 = 0.$$

, so the $\hat{\mu}_n$ value changes from $\hat{\mu}_n \rightarrow \hat{\mu}_n$ and $\hat{\sigma}_n \rightarrow 0$.

\Rightarrow Bayesian learning \uparrow =

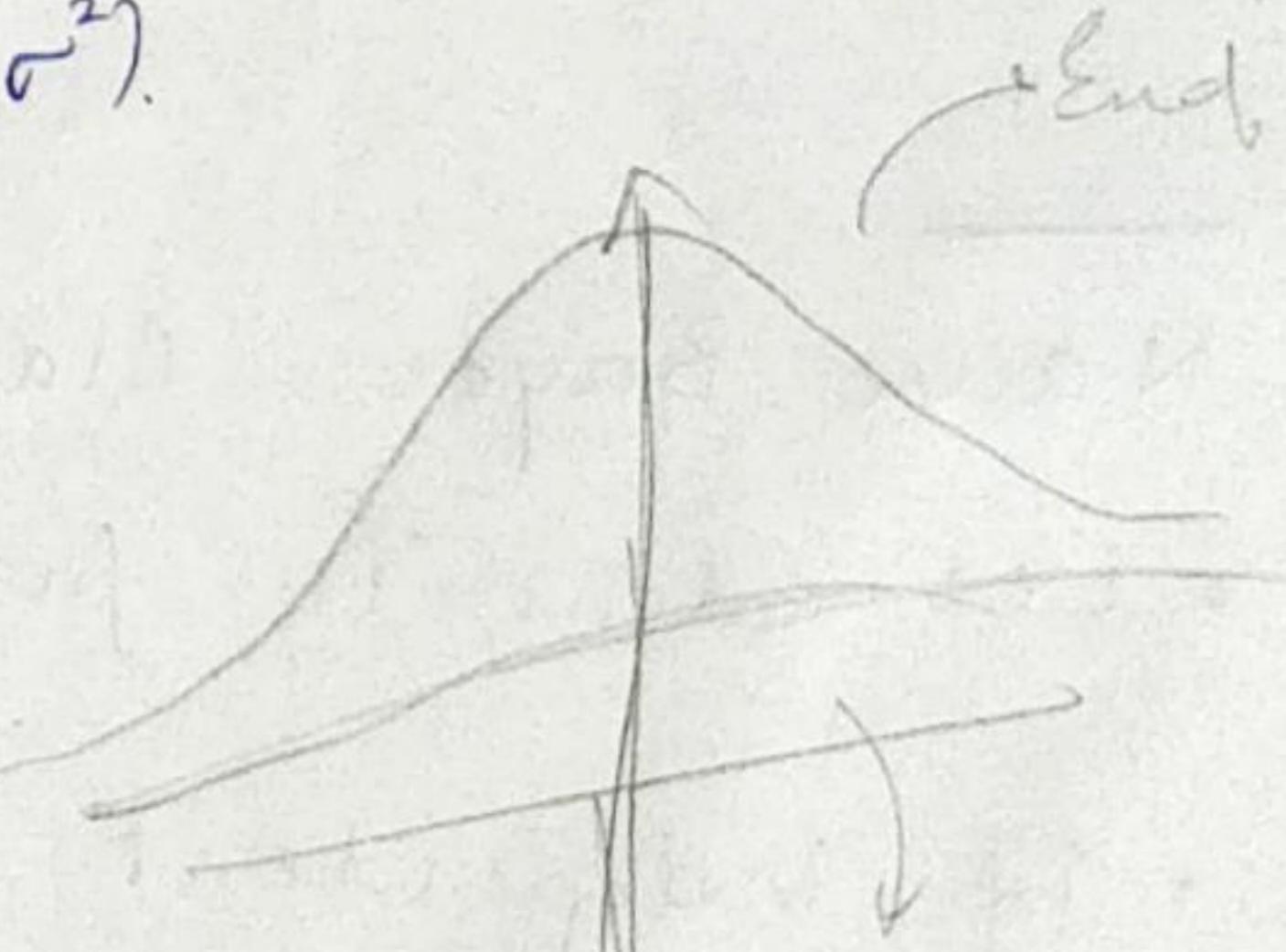
• Recursively: we can do the same, everytime
 we get a new sample.

$$p(D^n | \theta) = p(x_n | \theta) \times p(D^{n-1} | \theta)$$

$$p(x | D) = \int p(x | \mu) \times p(\mu | D) d\mu \quad \text{can be checked by integration}$$

$$\sim N(\hat{\mu}_n, \sigma^2 + \sigma_n^2)$$

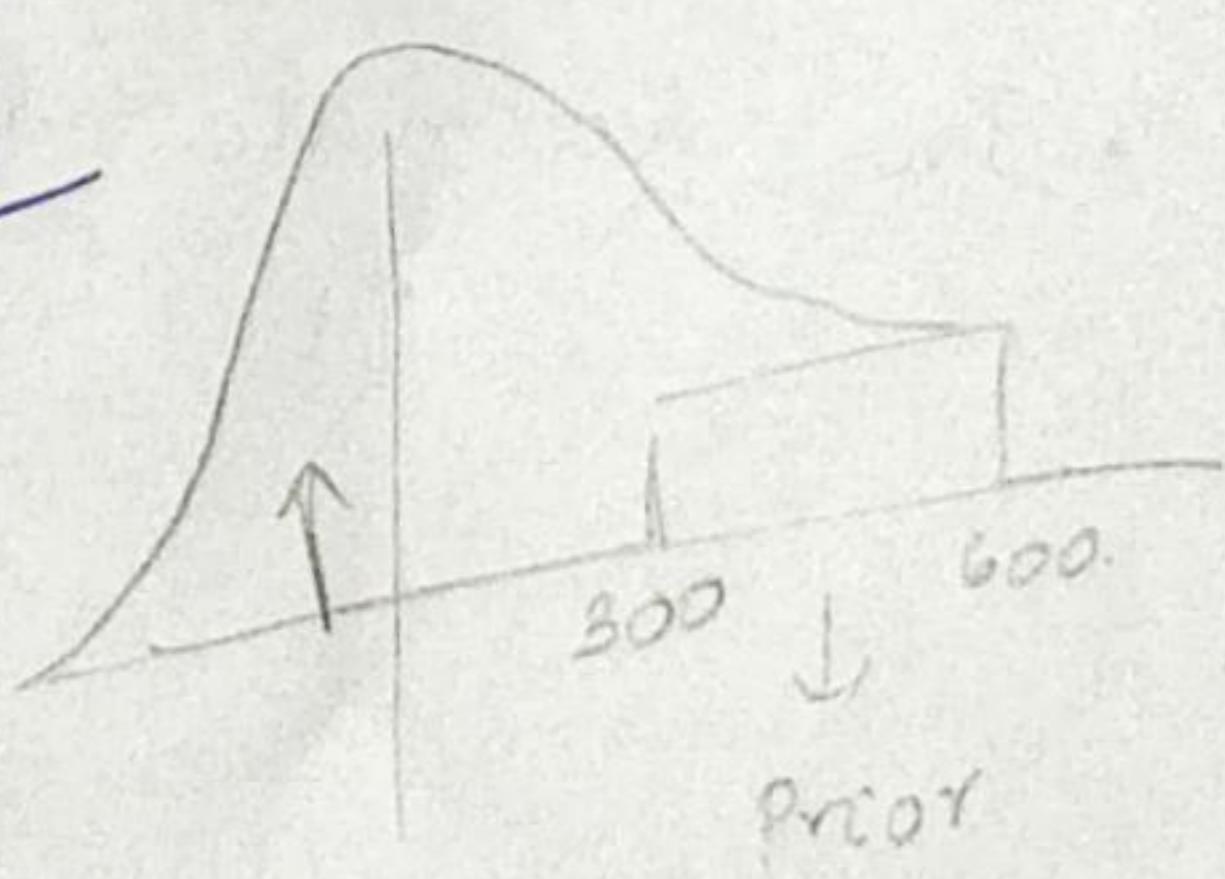
$$\text{as } n \rightarrow \infty, \quad p(x | D) \sim N(\hat{\mu}_n, \sigma^2).$$



$\hat{\mu}_n$ we think is impossible in
 whatever prior, \Rightarrow it will never

• If uniform density, $\hat{\mu}_{act}$ is outside,
 you can never reach $\hat{\mu}_{act}$

So, use some distro with range $(-\infty, \infty)$



$$D = \{4, 7, 8\}$$

TB example

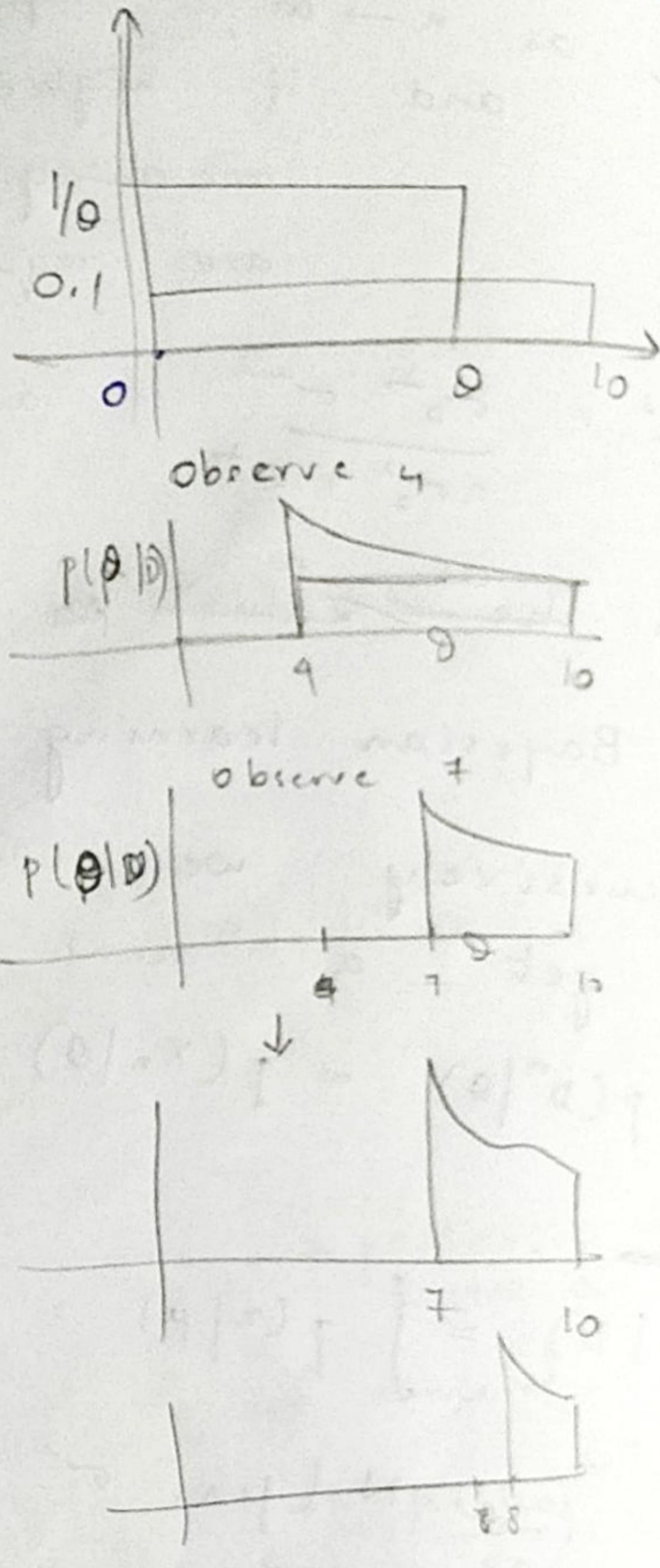
$$\text{If } \theta = 5$$

$$\Rightarrow P(D) = 0.$$

$$\theta \rightarrow \theta_{\max} \rightarrow 8.$$

$$\text{If } \theta = 10,$$

$$P(\theta \leq 4) = 0.$$



Naive Bayes Classifiers

- Let density f_x be normal.

- If independent features //

$$p(x) = p(x_1)p(x_2) \cdots p(x_d) \quad \mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix}, \quad \sigma^2 = \begin{bmatrix} \sigma_1^2 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \sigma_d^2 \end{bmatrix}_{d \times d}$$

- In normal distro:
It implies off diagonal