

Introduction to Data Science (Lecture)

02/09

Def An experiment is one activity that produces distinct outcomes. (ϵ)

* The set of all outcomes from an experiment is called Ω (sample space)

* subsets of Ω are called events. An event occurs if the result of the experiment $w \in \Omega$ if $w \in A$

ex: Toss a dice $\Omega = \{1, 2, 3, 4, 5, 6\}$

$A = \{\text{"an odd number}\} = \{1, 3, 5\} \rightarrow$ Event

$w=1 \Rightarrow A \text{ occurs / happens}$ $w=2 \Rightarrow A \text{ does not happen}$

$w=3 \Rightarrow$

What's probability?

LTRF idea (Long Term Relative Frequency)

Let's say we repeat an experiment n times, we have an event A

$N(A, n) = \frac{\text{"number of times } A \text{ happens}}{n} \rightarrow \max. = 1, \min = 0$

ex: $A = \{H\} \quad \Omega = \{H, T\}$

Flip a coin 10 times and 5 heads $N(A, 10) = \frac{5}{10} = \frac{1}{2}$

LTRF $\lim_{n \rightarrow \infty} N(A, n) = p$ if you repeat the experiment a multiple of times
then it becomes a number (p) \rightarrow this is valid

if you assume that the events are independent from each other

Rules of N :

1) $0 \leq N(A, n) \leq 1, \quad N(\Omega, n) = 1$

2) Let A, B be two events $A \cap B = \emptyset$ (disjointed) then

$N(A \cup B, n) = N(A, n) + N(B, n)$ (addition rule)

3) Each experiment is independent of each other

single experiment \rightarrow Let's assume that all events are "observable" $A \in F$

Def: A function $P: F \rightarrow [0, 1]$ is a probability distribution if P satisfies rules 1,2

\rightarrow 2*) Let $A_1, A_2, \dots, A_i, \cap A_j = \emptyset$

$$P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$$

Rule 1,2 \Rightarrow a bunch of formulas

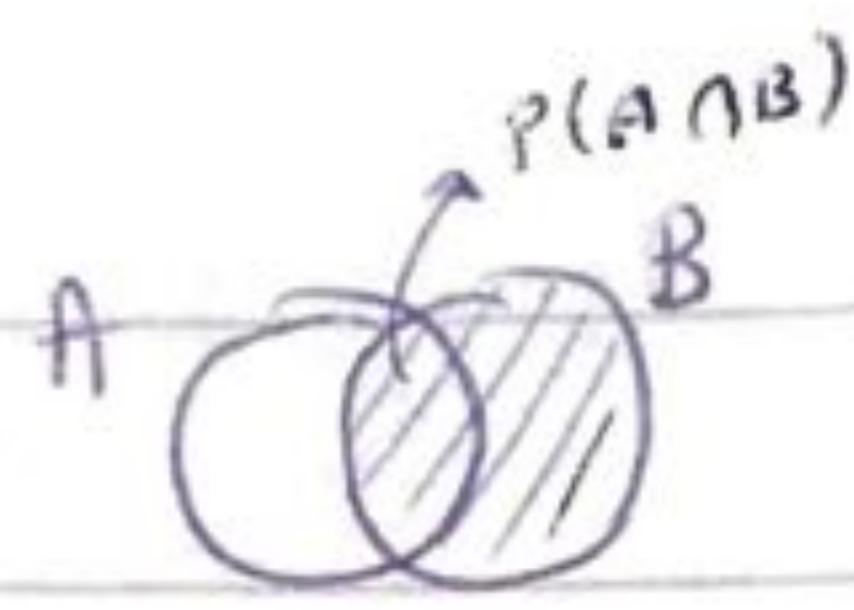
1) $P(A^c) = 1 - P(A) \quad A^c: \Omega \setminus A$

2) $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ Inclusion/Exclusion Principle

3) $P(A \cup B) \leq P(A) + P(B)$ Boole's Principle/Inequality

Conditional probability:

$$P(A|B) := \frac{P(A \cap B)}{P(B)} \quad \text{given } P(B) \neq 0$$



Independence: $P(A \cap B) = P(A) \cdot P(B)$

$\Downarrow \textcircled{1}$

$$P(A|B) = P(A) \text{ if } P(B) > 0$$

04/08

Vector Space:

{bread, butter(kg), eggs}

$$v = (\overset{\uparrow}{3}, \overset{\uparrow}{1}, \overset{\uparrow}{7})$$

Properties:

$$v = (v_1, v_2, v_3, \dots, v_n) \in V \quad \text{if } v+u = (v_1+u_1, v_2+u_2, v_3+u_3, \dots, v_n+u_n) = u+v$$

$$u = (u_1, u_2, u_3, \dots, u_n) \in V$$

$$\textcircled{2} \text{ if } v, u: \text{vectors} \quad v+u+w = (\overset{\uparrow}{v+u}) + \overset{\uparrow}{w} = z+w$$

$$\textcircled{3} \text{ } 0 \leftarrow \text{zero vector} = (0, 0, \dots, 0_n) \in V$$

$$\textcircled{4} \text{ } u+0=u \text{ identity element in } V$$

$$\textcircled{5} \text{ for every } u \in V \text{ (exists) } \exists -u \in V: u+(-u) = 0 \xrightarrow{\text{zero vector}} \text{the inverse of every vector should be included in the } V$$

$$R \setminus V \quad a \in R \quad u \in V \quad \text{e.g. } u = (1, 5, -2, \dots, 7)$$

$$a \cdot u \in V \quad \text{e.g. } a=2 \quad 2 \cdot u = 2(1, 5, -2, \dots, 7) = (2, 10, -4, \dots)$$

Warning! If $u \in V, v \in V$ we cannot write $u \cdot v$, when multiply 2 vectors, we don't get a vector but a scalar (dot product)

$$a, b \in R \quad u, v \in V \quad a(u+v) = (au+av) \in V \quad (a+b) \cdot u = a \cdot u + b \cdot u$$

$$a \cdot (bu) = (ab) \cdot u$$

\rightarrow not between two vectors though

* In R and V there are only two operations: multiplication and addition $\textcircled{1, 2}$

~~you cannot add $v+v$~~

f: $V \times V$ this means one element of vector 1 and another from vector 2 can be added.

$$u, v \xrightarrow{\oplus} u+v \in V$$

$$g: R \times V$$

$$a, v \xrightarrow{\odot} a \cdot v \in V$$

* If I have coordinates then the vector looks like: $V = \{(x, 0), (0, y)\}$

$$v = (2, 0) \in V \quad \text{if } u+v \in V? \quad u+v = (2, 0) + (0, 3) = (2, 3) \notin V$$

vector space if you combine them they are no longer a vector space

Intro to Data Science (Lecture)

c2. Def: A random variable as a "mapping" function $X : \Omega \rightarrow \mathbb{R}$ that assigns a real number to each outcome.

Ex.1 flip a coin (0) times $\Omega = \underbrace{\text{HHHTT...H}}_{10 \text{ times}}$

$X(\Omega)$ $X(w) = \text{"how many heads"}$

Ex.2 $\Omega = \{\text{all texts}\}$

$$x(w) = \begin{cases} 1, & \text{if text contains "free"} \\ 0, & \text{otherwise} \end{cases}$$

$x(w) = \text{"count how many times free appears"}$

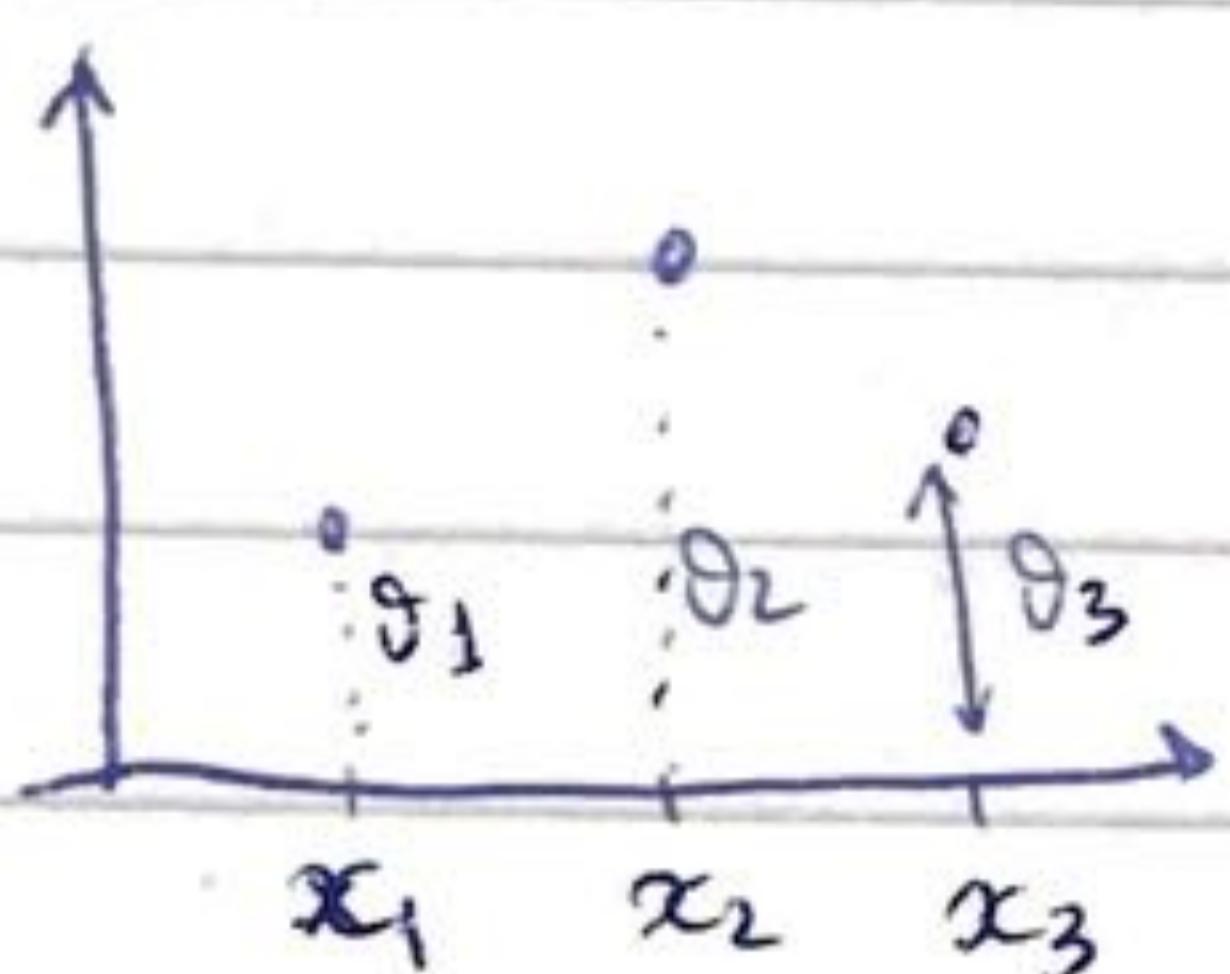
Distribution function:

$$X : \Omega \rightarrow \mathbb{R}$$

$$\mathbb{P}(X \leq x) = \mathbb{P}(\{w : X(w) \leq x\}) =: F_X(x)$$

Def: A discrete random variable is a random variable that takes discrete values, ex. $(0, 1, 2, 3, \dots)$. The PMF (Probability Mass Function).

$$f_X(x) = \mathbb{P}(X=x) = \begin{cases} \vartheta_i, & \text{if } x=x_i \\ 0, & \text{otherwise} \end{cases} \quad X \text{ takes values } (x_1, x_2, \dots)$$



$$\text{Formulas: } \begin{aligned} ① \quad F_X(x) &= \mathbb{P}(X \leq x) = \mathbb{P}(\{X=x_1\} \cup \{X=x_2\} \cup \dots \cup \{X=x_i\}) \\ &= \sum_{x_i \leq x} \mathbb{P}(X=x_i) = \sum_{x_i \leq x} f_X(x_i) \end{aligned}$$

$$② \quad F_X(b) - F_X(a), \quad b \geq a = \sum_{a < x_i \leq b} f_X(x_i)$$

$$③ \quad \sum_{x_i} f_X(x_i) = 1$$

Def: The Expected value of a discrete random variable X is $\mathbb{E}[X] := \sum_{x_i} x_i f_X(x_i)$ based on LTRF idea.

Let z_1, z_2, \dots, z_{100}

$$\sum_{i=1}^{100} z_i / 100 = (\sum_{z_i=1} z_i + \sum_{z_i=2} z_i + \dots + \sum_{z_i=6} z_i) / 100$$

$$= 1 \cdot N(z=1, 100) + 2 \cdot N(z=2, 100) + \dots + 6 \cdot N(z=6, 100)$$

Common expectations: $V[X] := \mathbb{E}[(X - \mathbb{E}[X])^2] = \sum_{x_i} (x_i - m)^2 \cdot f_X(x_i)$

k : th central moments

$$\mathbb{E}[(x - \mathbb{E}(x))^k]$$

standard deviation

$$\sqrt{V[x]} = \sigma(x)$$

$$\mathbb{E}[X] = \mu$$

$$\mathbb{E}\left[\left(\frac{x - \mathbb{E}[x]}{\sqrt{V[x]}}\right)^k\right] = \mathbb{E}\left[\left(\frac{x - \mu}{\sigma}\right)^k\right] \text{ k:th standardised moment}$$

$\kappa = 3$: skewness, $\kappa = 4$ = kurtosis

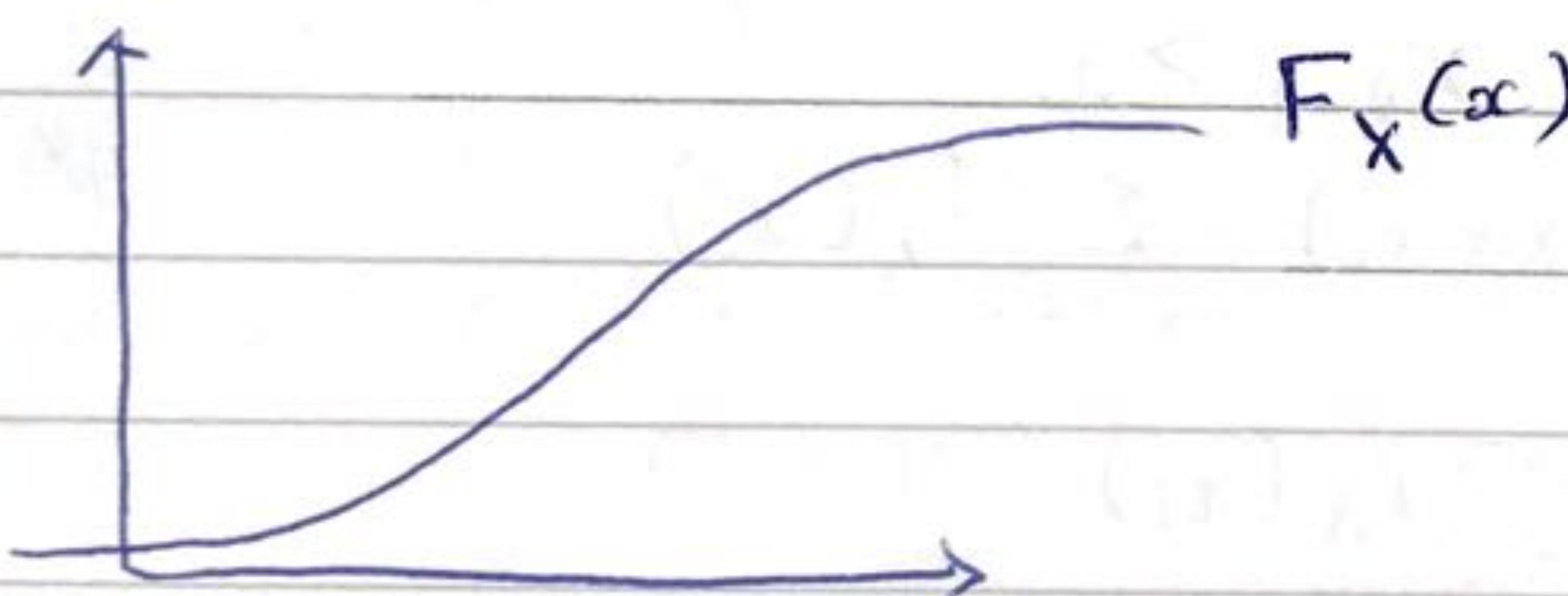
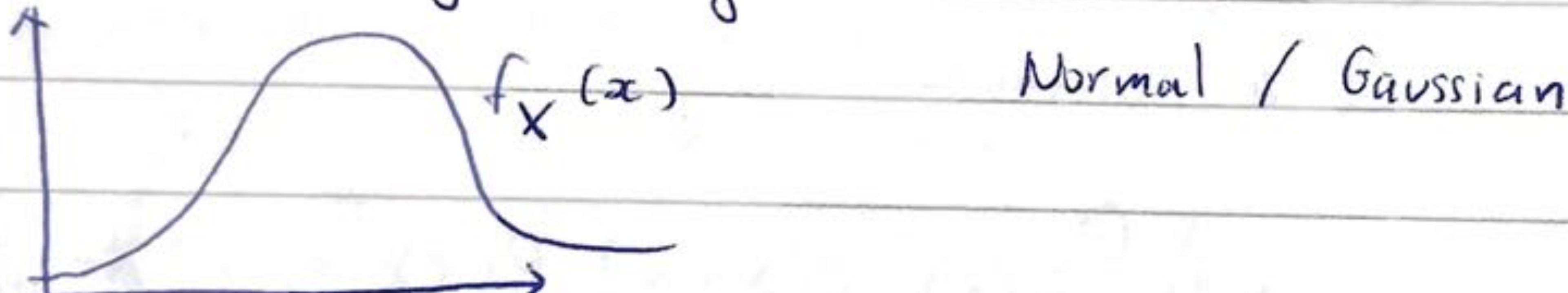
Def: We say that a random variable is continuous if there is a piecewise count function $f_X(x) : \mathbb{R} \rightarrow [0, +\infty]$ such that

$$F_X(x) = \int_{-\infty}^x f_X(s) ds = P(X \leq x)$$

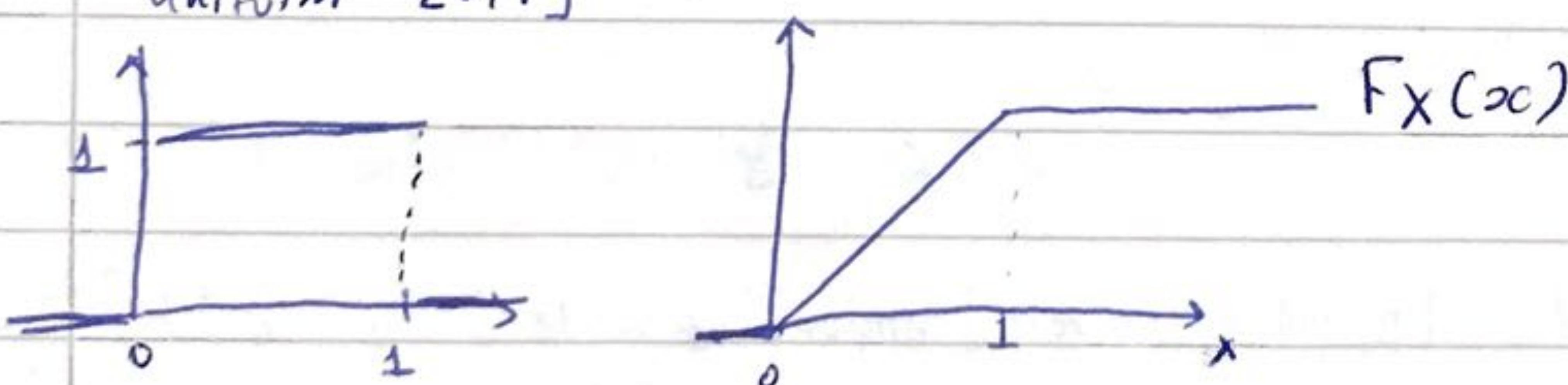
Warning: $P(X = x) = 0$

$$\begin{aligned} P(X = x) &= P(\{X \leq x\} \setminus (x < x)) \\ &= P(X \leq x) - P(X < x) \\ &= \int_{-\infty}^x f_X(s) ds - \int_{-\infty}^{x-} f_X(s) ds = 0 \end{aligned}$$

f_X : probability density



Uniform $[0, 1]$



Transformation: X is discrete Random Variable $g: \mathbb{R} \rightarrow \mathbb{R}$ $Y = g(x)$

$$\text{Ex. } g(x) = (x - m)^2$$

$F_Y(y) = P(Y \leq y) = P(g(X) \leq y)$ Let's say g is increasing ^{use the inverse} $= P(X \leq g^{-1}(y)) = P(X \leq g^{-1}(y)) = F_X(g^{-1}(y))$

Lecture 3

Def: A multivariate Random Variable \underline{X} is a mapping from Ω to \mathbb{R}^m
 $\underline{X}: \Omega \rightarrow \mathbb{R}^m$ That is \underline{X} represents m measurements
 $F_{\underline{X}}(x_1, x_2, \dots, x_m) = P(\{\underline{X} \leq \underline{x}\}) = P(\{x_1 \leq x_1\} \cap \{x_2 \leq x_2\} \cap \dots \cap \{x_m \leq x_m\})$

Joint Distribution function (JDF)

Def: Let $F_{\underline{X}}(x_1, x_2)$ a JDF for $\underline{X} = (x_1, x_2)$

$F_{x_1}(x_1) := F_{\underline{X}}(x_1, +\infty) = P(\{\underline{X} \leq \underline{x}\} \cap \{x_2 \leq \underbrace{\infty}_{\text{all } \Omega}\}) = P(x_1 \leq x_1)$
 Marginal Distribution

↳ it is the 03-Random Variables .pynb → def F_{-X-12} because y takes ∞ which is the largest value (like the ∞ in the JDF)

Def: A 2-dimensions Random Variable $\underline{Z} = (\underline{X}, Y)$ we can say \underline{X} and Y are independent if $F_{\underline{Z}}(x, y) = F_{\underline{X}}(x) \cdot F_Y(y)$ for all $(x, y) \in \mathbb{R}^2$

Think $P(A \cap B) = P(A) \cdot P(B)$

Def: Let $\underline{Z} = (\underline{X}, Y) \in \mathbb{R}^2$ then the conditional PMF/Density

$$f_{\underline{X}|Y}(x|y) = f_{\underline{X}Y}(x, y) / f_Y(y),$$

density (little f) of X conditioned on Y $f_Y(y) > 0$

e.g. discrete $F_{\underline{X}Y}(x, y) = P(\underbrace{\{\underline{X} \leq x\}}_A \cap \underbrace{\{Y = y\}}_B) = \frac{P(A \cap B)}{P(B)} = \frac{\sum_{x_i \leq x} P(\underline{X} = x_i, Y = y)}{P(Y = y)}$

$$= \sum_{x_i \leq x} \left(\frac{f_{\underline{X}Y}(x_i, y)}{f_Y(y)} \right) \rightsquigarrow \text{PMF}$$

Independence ($\Rightarrow f_{\underline{X}|Y} = f_X(x) f_Y(y) \geq 0$)

! head -n 5 data/corrs.csv

Sequences: We say that we have a sequence of Random Variable x_1, x_2, \dots, x_m , if for a fixed m $\underline{X} = (x_1, \dots, x_m)$ is a RV

* We say sequence is independent if $F_{x_1, x_2, \dots, x_m}(x_1, x_2, \dots, x_m) = F_{x_1}(x_1) \cdot F_{x_2}(x_2) \cdots F_{x_m}(x_m)$

* We say a sequence is identically distributed if $F_{x_i} = F_{x_j}$ for all i, j

* We say a sequence is iid if both independent and identically distributed

In regression we want $E[Y | \underline{X} = x] = r(x)$

$r(x)$
 regression
 function

17/09/24 Intro to Data Science; Lecture 4

Variance
R(·)

Learning from data? What is the avg weight of the population in Sweden?

1) Experiment is to select a random person $\omega \in \Omega = \{w_1, w_2, \dots, w_N\}$

weighting person $w \rightarrow x(w)$, $x(w) : \Omega \rightarrow \mathbb{R}^+$

what do we know about x ? x : continuous, $x \geq 0$ x : bounded $\leq M$

Assumptions about reality

2) Choose to select n people x_1, \dots, x_n (if there is replacement then iid)

x_1, \dots, x_n IID

knowing the weight of 1 person doesn't

$\frac{1}{n} \sum_{i=1}^n x_i$ LTRF idea

give you an indication of other person's weight

$$\mathbb{E}[x]$$

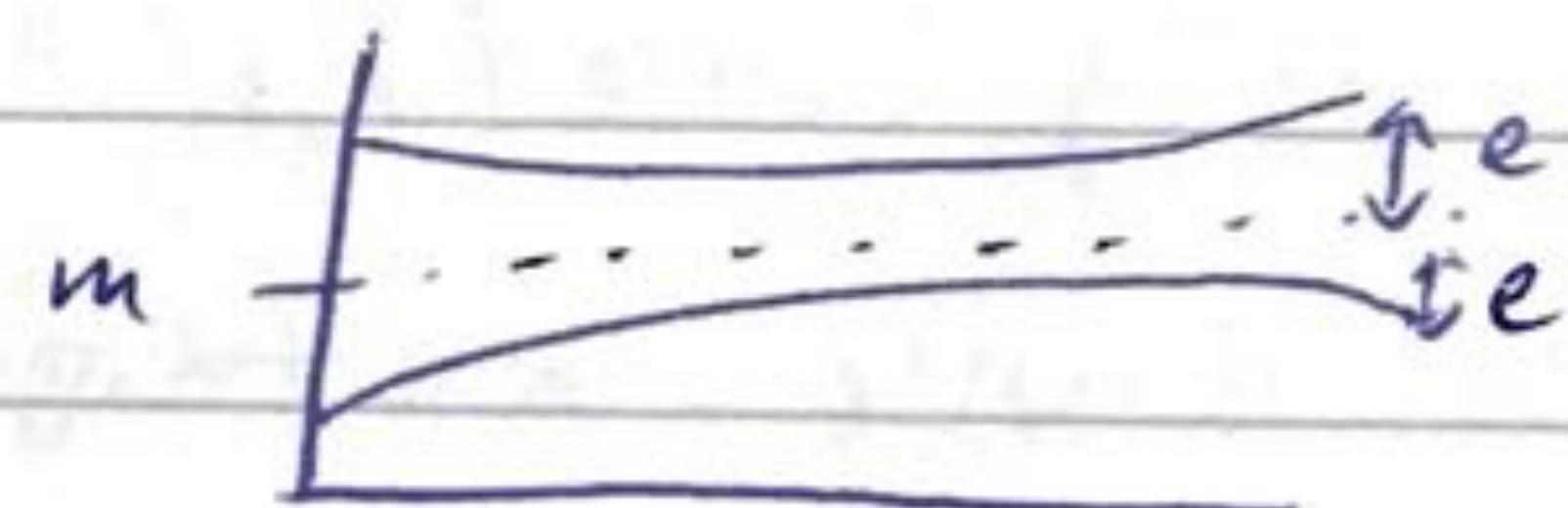
phenomenon of concentration (when you do an experiment a lot) as. coin toss

We started with LTRF we defined probability now we are back at CTRF, we gained that we can calculate things.

Fix ϵ (error) > 0 and define $m = \mathbb{E}[x_i]$, $\bar{x}_n := \frac{1}{n} \sum_{i=1}^n x_i$

$$\mathbb{P}(|\bar{x}_n - m| > \epsilon) \rightarrow 0$$

as $n \rightarrow \infty$



"convergence in probability"

* Remark: Sometimes we can calculate exactly $\mathbb{P}(|\bar{x}_n - m| > \epsilon)$ $x_i \stackrel{iid}{\sim} \text{Bernoulli}(p)$
then $n \cdot \bar{x}_n \sim \text{Binomial}(n, p)$

Probability inequalities: Chebychev: Let $X \in \mathbb{R}$ be a RV and let $\epsilon > 0$

then $\mathbb{P}(|X - \mathbb{E}[X]| > \epsilon) \leq \frac{\mathbb{V}[X]}{\epsilon^2}$, assuming $\mathbb{E}[X]$ exists and $\mathbb{V}[X]$ exists

Small variance \Rightarrow small probabilities

Large error \Rightarrow — — —

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \underbrace{\sum_{i=1}^n \mathbb{E}[x_i]}_{m} = m$$

take any one
bc they're identically
distributed

$$\mathbb{P}(|\bar{x}_n - m| > \epsilon) \leq \frac{\mathbb{V}[\bar{x}_n]}{\epsilon^2} \quad \text{Independence gives } \mathbb{V}[\bar{x}_n] = \frac{\mathbb{V}[x_i]}{n}$$

$$\begin{aligned} \mathbb{V}[\bar{x}_n] &= \mathbb{E}\left[(\bar{x}_n - \mathbb{E}[\bar{x}_n])^2\right] = \mathbb{E}\left[\left(\frac{1}{n} \cdot \sum_{i=1}^n (x_i - m)\right)^2\right] = \frac{1}{n^2} \mathbb{E}\left[\left(\sum_{i=1}^n y_i\right)^2\right] \\ &= \frac{1}{n^2} \mathbb{E}\left[\sum_{i,j} y_i y_j\right] \stackrel{\text{indep.}}{=} \frac{1}{n^2} \left(\sum_i \mathbb{E}[y_i] \cdot \underbrace{\mathbb{E}[y_j]}_{\text{indep.}} + \sum_{i \neq j} \mathbb{E}[y_i^2]\right) = \frac{\mathbb{V}[x_i]}{n} \end{aligned}$$

$$\mathbb{V}[x_i]$$

Variance of a sum = sum of the variances
 independence

$$P(|\bar{X}_n - m| > \epsilon) \leq \frac{\sqrt{V[\bar{X}_n]}}{\epsilon^2} = \frac{\sqrt{V[X_1]}}{n \cdot \epsilon^2}$$

$n \cdot \epsilon^2 \rightarrow 0$ as $n \rightarrow \infty$

Fix a number $a \in (0, 1)$ choose ϵ such that $\frac{\sqrt{V[X]}}{n \cdot \epsilon^2} = a$

$$\frac{\sqrt{V[X]}}{a} = \epsilon$$

$$P\left(|\bar{X}_n - m| > \frac{\sqrt{V[X]}}{\sqrt{n} \cdot \epsilon}\right) \leq a$$

(look at the prop) Hoeffding's inequality (will be used for the entire course - we need to learn it).

Let X_1, \dots, X_n be IID RVs such that $a \leq x_i \leq b$ and $\epsilon > 0$ then

$$P(|\bar{X}_n - m| > \epsilon) \leq 2e^{-\frac{2n\epsilon^2}{(b-a)^2}}$$

the bigger a then all the $e^{-\frac{2n\epsilon^2}{(b-a)^2}}$ becomes bigger

First thing: fix a solve $2e^{-\frac{2n\epsilon^2}{(b-a)^2}} = a$

$$\epsilon = \sqrt{\frac{(b-a)^2 \ln(\frac{2}{a})}{2n}}$$

then form an interval $I = (\bar{X}_n - \epsilon, \bar{X}_n + \epsilon)$

empirical mean

$$P(m \text{ is in } I) \geq 1-a \quad \text{let us make a decision based on data } (x_1, \dots, x_n)$$

Decision: I decided that the true mean is inside my interval I . (you are going to be correct $1-a$ % of times)

I : is for before you've seen the data then if you see the data is no longer a prob.

→ It is the fastest you learn as ~~data~~ ~~the~~ errors get smaller when sample size increases e.g. weight of a person (there is a max. and we know it is non-negative., coin toss)

Binomial is exact but it is difficult to calculate it for $n \geq 20$. but you can use Hoeffding to approximate.

If you know the variance you can do better: Bennett

$$2 \exp\left(-\frac{s^2}{\delta^2}\right) \dots \text{goes to zero when sample } \rightarrow \infty$$

25/09 Hoeffding's Inequality: true if $a \leq x \leq b$

Lecture 5: $E[e^{-s(x-\mu(x))}] \leq ? e^{s^2/\sigma^2}$ (assume this is true).

if we assume that is a valid assumption and as soon as true it is called sub-Gaussian. It's called sub-Gaussian because the tails are smaller than Gaussian

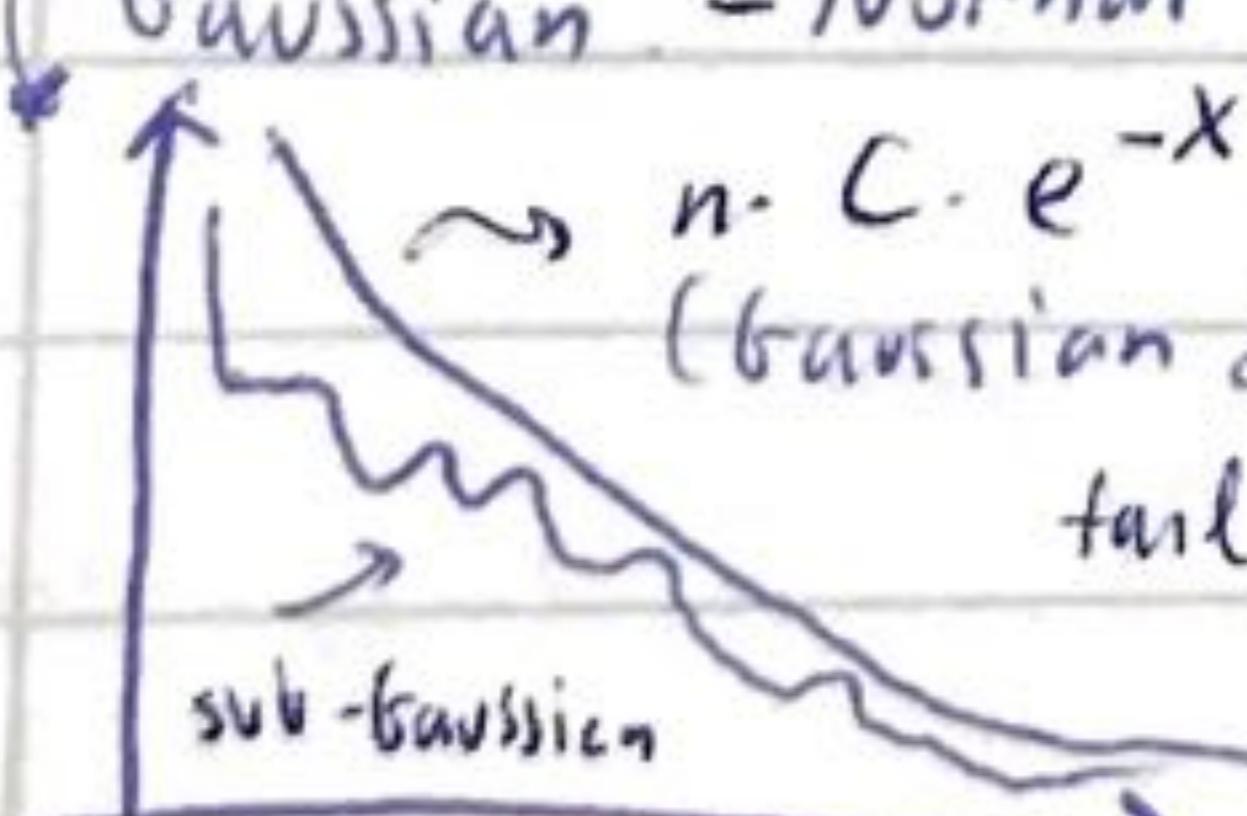
Gaussian - Normal

$$\text{density: } n \cdot C \cdot e^{-x^2/\sigma^2}, C > 1$$

(Gaussian distribution)

tail is smaller than Gaussian

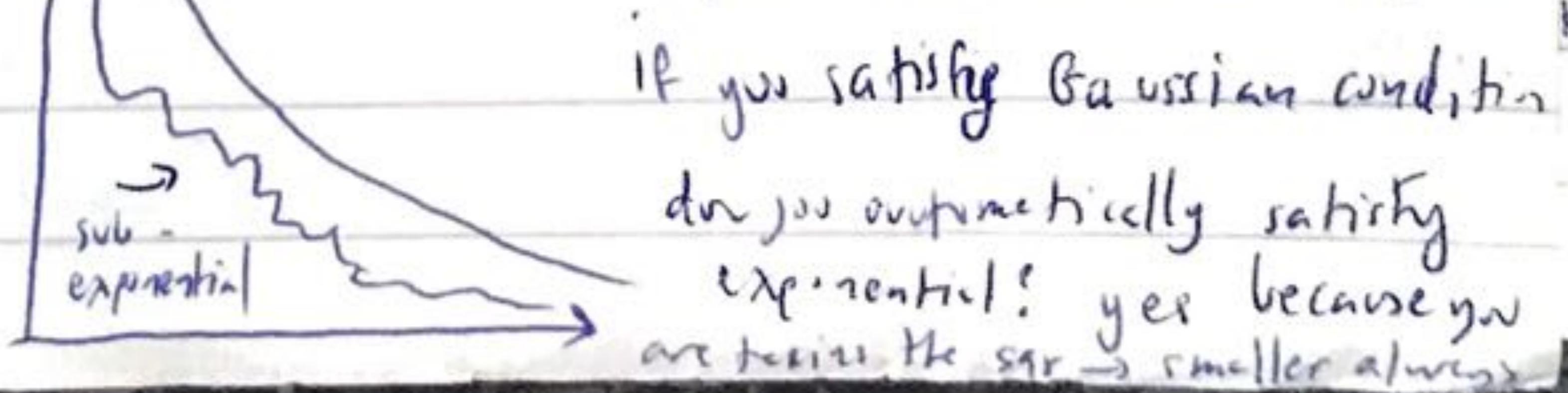
density:



$\rightarrow c \cdot e^{-x/\sigma}$ (exponential distribution)

if you satisfy Gaussian condition

do you automatically satisfy exponential? yes because you are tailing the same smaller always



Ex. 1: Normal is sub-Gaussian

Ex. 2: as $x \leq b$ also sub-Gaussian

Ex. 3: $X \sim N(0, 1) \rightarrow X^2$ is sub-exponential ($\chi^2 \sim \text{Chi-squared}$)
the tail gets an exponential distribution

if you want to estimate the variance not only the mean, then you are putting in more information using squared values

$\hat{V}_n := \sum_{i=1}^n (X_i - \mathbb{E}[X_i])^2$ the rate with respect to epsilon for sub-exponential is lower.

for sub-exponential: $P(|\bar{X}_n - \mathbb{E}[\bar{X}_n]| > \epsilon) \leq \max\left(2e^{-\frac{\epsilon^2 n}{C_1}}, 2e^{-\frac{\epsilon n}{C_2}}\right)$

as the error increases (for large ϵ) you get fatter tail → that is why there isn't Ω^2 .

Def: A statistical model is a set of distributions (densities/ regression functions)
↳ Wasserman.

$\mathcal{F} \quad S := \{F\}$ - set of distributions

These are our assumptions: We are implicitly assuming that our data comes from

$F^* \in S$

Ex. 1: $S := \left\{ \underbrace{\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}}_{\Phi_{\mu, \sigma}}, \mu \in \mathbb{R}, \sigma > 0 \right\}$ it is not a big set as I can describe it using only 2 numbers (mean & variance).

Then I will have a good approximation. This is called a parametric model

Non-parametric (⇒ not parametric)

Ex. $S := \{F_x : F_x(a) = 0, F_x(b) = 1\}$ loss all functions from a to b - you need a lot of parameters to describe it, you take a RV between a and b.

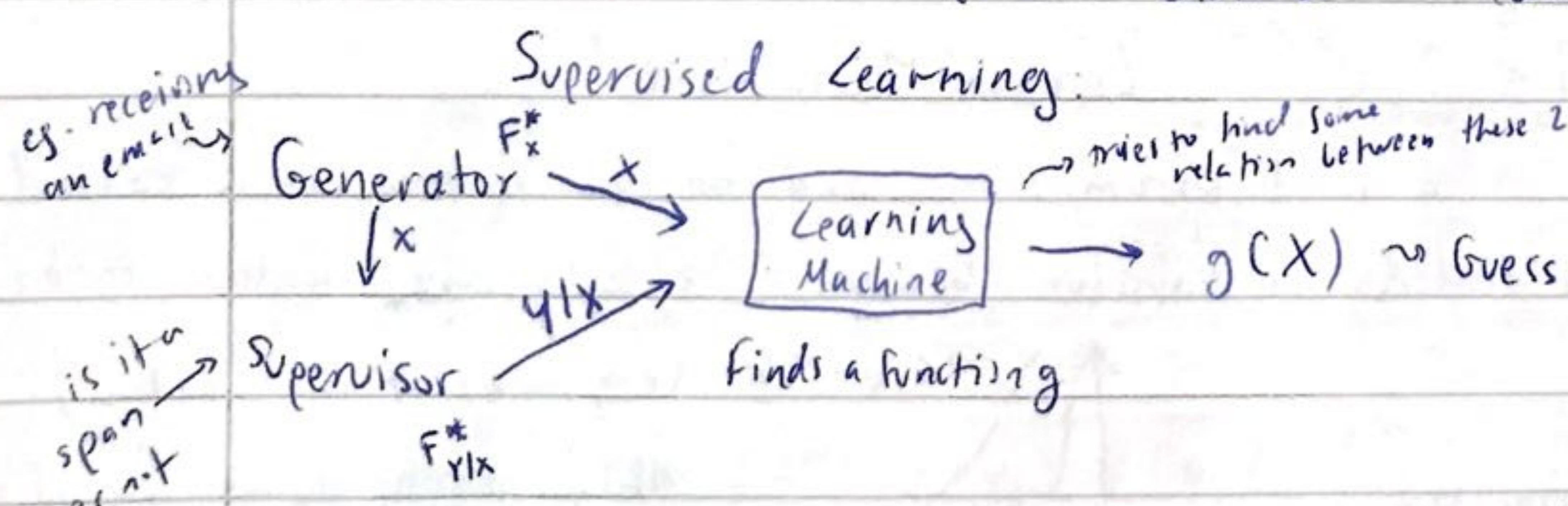
What can we learn?

* $\mathbb{E}[X] = \int x \cdot dF^*(x)$ - sum or integral based on if it is a discrete/continuous variable.

* learn $(f^*)' = f^*$ density - "Histogram" example of learning density

learn $F^*(x) : \hat{F}(x) = \frac{1}{n} \sum_{x_i \leq x} 1$ (how many sample is below this)

↳ Empirical Distribution function.



Def: Consider a function $L : \mathbb{R}^2 \rightarrow \mathbb{R}$ (loss function) Let $(X, Y) \sim F_{XY}^*$, then risk of a

function $g: \bar{X} \rightarrow \bar{Y}$ is $R(g) := \mathbb{E}[L(g(x), Y)]$ → it measures how close they are is \hat{Y} , it is a good guess (L) and \hat{s} : expected error

Ex. 1: $L(a, b) = (a - b)^2$ quadratic loss.

cares about large errors ↴

output is a number $R(g) = \mathbb{E}[(g(x) - Y)^2]$ least square problem if $g(x)$ = line (error very big when diff is big because 1^2)

Ex. 2: $L(a, b) = |a - b|$ absolute loss

$R(g) = \mathbb{E}[|g(x) - Y|]$ → measures everything the same (proportional).

Ex. 3: if the output $\bar{Y} = \{0, 1\}$:

$$L(a, b) = \begin{cases} 1, & \text{if } a \neq b \\ 0, & \text{otherwise} \end{cases}$$

$$R(g) = \mathbb{E}[L(g(x), Y)] = P(g(x) \neq Y)$$

↳ Pattern recognition

$1 - R(g)$ = "accuracy"

It is called 0-1 loss

Objective: The learning machine tries to minimise risk by searching a $g^* \in M$ [~model space] that solves $R(g) = \inf_{g \in M} R(g)$ → optimisation problem

Ex. 1: $M = \{kx + m, k \in \mathbb{R}, m \in \mathbb{R}\}$ → linear function

$$L(a, b) = (a - b)^2$$

Least Squares problem
↳ Ordinary Linear Regression

The LM doesn't have access to $R(g)$ because it doesn't know F^* . To solve this we use data (examples) which is a bunch of pairs $((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n))$.

Instead solves $\hat{R}(g^*) = \inf_{g \in M} \hat{R}(g)$

$$\text{where } \hat{R}(g) = \frac{1}{n} \sum_{i=1}^n L(g(x_i), y_i)$$

Once you have found \hat{g}^* (if you do the same thing again with new data then \hat{g}^* changes)
↳ \hat{g}^* is random (depends on size of M and # of data points)

How do we know its any good? Take new data (testing data)

$((x_{n+1}, y_{n+1}), \dots, (x_{n+m}, y_{n+m}))$ Then you calculate:

$$\hat{R}_m(\hat{g}^*) = \frac{1}{m} \sum_{j=1}^m L(\hat{g}(x_{n+j}), y_{n+j}) \leftarrow \text{testing error}$$

30/09 $\int = \{ \text{all admissible distributions} \}$

$\int : \text{statistical Model}$

$M = \{ \text{all guessing functions } g \text{ for the LM} \}$

Loss Function $L(a, b)$

Data: $((x_1, y_1), \dots, (x_n, y_n)) =: D_{Tr}$: "Training Data"

↳ IID - F_{XY}^{**} (unknown to us)

(if there is sorting in the data there is no longer IID → need scheduling)

$((x_{n+1}, y_{n+1}), \dots, (x_{n+m}, y_{n+m})) =: D_{Te}$: "Testing Data"

empirical risk

The LM minimises the "empirical risk" $\hat{g} = g \in M \min \frac{1}{n} \sum_{i=1}^n L(g(x_i), y_i)$ because they are actual data

We are thinking that the learning machine is approximately minimising the true risk

$$R(g) := \mathbb{E}[L(g(x), Y)]$$

Risk

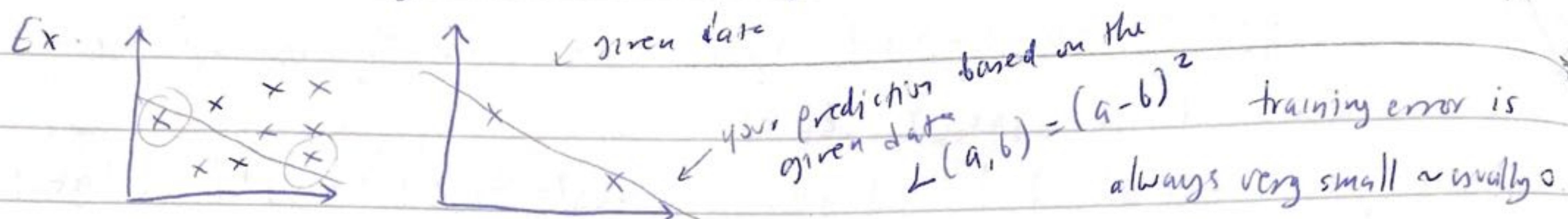
what is the approximation of this?
of the true risk

We test \hat{g} on D_{Te} by calculating $\frac{1}{m} \sum_{i=1}^m L(\hat{g}(x_{n+i}, y_{n+i}), E[R(\hat{g})])$
 \hat{g} : is not random anymore \rightarrow you guessed it $= E[R(\hat{g})] | D_{Tr}]$

actual function based on the training data and you want to fit it on the testing data

trained model
(posterior error)
 \hat{g} : is fixed since you know the training data.

for automatic fitting use fitHub



but when you applied it in ^{testing} data increases.

Def: Given a statistical model S , the data space $\bar{\mathbb{X}}_n$ is where $((x_i, y_i), \dots)$

$(x_n, y_n) \sim^{ID} F_x^*$ line $F_x^* \in S$. true distribution $\otimes x_i = \text{includes both } x \text{ and } y$
 Ex. $x_i \in \mathbb{R}$ $\bar{\mathbb{X}}_n = \mathbb{R}^{xn}$ but we don't know it.

$x_i \in \mathbb{R}^d$ $\bar{\mathbb{X}}_n = (\mathbb{R}^d)^{xn}$

Def: A statistic is a function on the data space $\bar{\mathbb{X}}_n$: $T: \bar{\mathbb{X}}_n \rightarrow \mathbb{T}$
 ↪ the possible values of T .

Ex: Supervised Learning $\hat{g} = \arg \min_{g \in M} \frac{1}{n} \sum_{i=1}^n L(g(x_i), y_i))$ \hat{g} : is statistic
 data space $\bar{\mathbb{X}}_n = (\mathbb{R}^2)^{xn}$ can be calculated based on data.

↪ all straight lines that $\hat{g} \in M$ are in the model space.

Ex: Testing error $\frac{1}{m} \sum_{i=1}^m L(\hat{g}(x_{n+i}), y_{n+i})$. this is statistic as we calculated it only on data.
 $=: T(D_{Te})$ a statistic $\bar{\mathbb{X}}_m$ \mathbb{T} : target space $= \mathbb{R}_+$

Properties: "An estimator" is a statistic T that approximates ϑ "

bias(T) = $E[T] - \vartheta$ if > 0 positively-biased, < 0 negatively-biased,
 ↑ "T estimates ϑ " if $= 0 \Rightarrow$ unbiased

Ex: $T(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i$ $\vartheta = E[x_i]$ bias(T) = 0

Ex 2: $T(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2$ $\vartheta = \text{Var}[x_i]$

bias(T) = $\frac{1}{n}$, if you have a lot of data, the bias is very small (bias goes to zero if you collect more data). If $\text{bias}(T) \rightarrow 0$ as $n \rightarrow +\infty$ we say T is asymptotically unbiased.

Ex. 3 $T(x_1, \dots, x_n) = x_1$, $\vartheta = E[x_i]$ what is the bias? $\text{Bias}(T) = 0$

Def: standard error $se(T) = \sqrt{V(T)}$

Ex. 1 $se(T) = \sqrt{\frac{V(x_1)}{n}}$ depends on n .

Ex. 3 $se(T) = \sqrt{V(x_1)}$ ↪ it doesn't matter how many data you get → will be the same

calculated } how far they are
from true

$$\text{Mean squared error (MSE)}(\tau) = \mathbb{E}[(\tau - \theta)^2] = \mathbb{E}[(\tau - \mathbb{E}[\tau])^2 + (\mathbb{E}[\tau] - \theta)^2]$$

$$= (\text{se}(\tau))^2 + (\text{bias}(\tau))^2$$

Bias - Variance decomposition

You can split the error into uncertainty (how far you are from the true value - if you push one to zero the other goes higher, the smaller the variance the higher the variance will be)

Bias can go to zero as you can find the true but \Rightarrow variance will go up.

4.10 Lab Estimation: x_1, \dots, x_n

- estimate parameters
for example: $x_1, \dots, x_n \sim N(\mu, \sigma^2) \rightarrow$ use parametric model

Point estimator: $\hat{\theta} = g(x_1, \dots, x_n)$

$\tau(x)$: estimator

Properties: 1) bias $E(\tau(x)) - \theta$ if $E(\tau(x)) = \theta \Rightarrow$ unbiased

2) $\text{se} = \sqrt{\text{var}(\tau(x))}$

MSE: Mean squared error:

$\text{Bias}^2(\hat{\theta}) + \text{Var}(\hat{\theta})$

Example: ① $x \sim Po(2)$ $\hat{\theta} = [P(x=0)]^2 = e^{-2}$

Choose unbiased estimator ($\tau(x)$)

$E(\tau(x)) = \theta = e^{-2}$

$$\sum_{x=0}^{\infty} \frac{e^{-2} 2^x}{x!} \tau(x) = e^{-2}$$

$$\tau(x) = -1^x$$

② AOS: Chapter 6: ① $x_1, \dots, x_n \sim Po(2)$

$$\hat{x} = \frac{\sum x_i}{n} \quad \text{find bias, se, MSE} \quad \text{1) expectation: } E\hat{x} = E\left[\frac{\sum x_i}{n}\right] = \frac{1}{n} \sum E x_i = \frac{n \cdot \mu}{n} = 2$$

sample mean is always unbiased bias = 0.

$$\text{MSE} = \text{Var}(\hat{x}) = \text{Var}\left(\frac{\sum x_i}{n}\right) = \frac{1}{n^2} \sum \text{Var}(x_i) = \frac{n \cdot \sigma^2}{n^2} = \frac{\sigma^2}{n}$$

$$\frac{\sum \text{Var}(x_i)}{n} = \frac{n \cdot \sigma^2}{n} = 2 \quad \Leftrightarrow \quad \text{se} = \frac{\sqrt{2}}{\sqrt{n}}$$

③ $x_1, \dots, x_n \sim Un(0, \sigma^2) \quad \hat{\theta} = 2\bar{x}_n \quad$ compute se, bias, MSE

$$\bar{x}_n = \frac{\sum x_i}{n} \quad E\hat{\theta} = 2 \frac{\sum x_i}{n} = 2 \frac{\sum (\epsilon x)}{n} = \frac{2 \cdot \frac{\sigma}{\sqrt{n}} \cdot n}{n} = \frac{2\sigma}{\sqrt{n}} = \sigma \quad \text{unbiased estimator} \quad \text{bias} = 0$$

$$\text{MSE} = \text{Var}(\hat{\theta}) = \text{Var}\left(\frac{2 \sum x_i}{n}\right) = \frac{4}{n^2} \sum \text{Var}(x) = \frac{4}{n^2} \cdot \frac{\sigma^2}{2} \cdot \frac{\sigma^2}{3n^2} = \frac{4\sigma^4}{6n^4} = \frac{2\sigma^4}{3n^2} = \frac{\sigma^4}{\frac{3n^2}{2}}$$

L7: Estimation Risk: std error if $\hat{f} \Rightarrow$ better

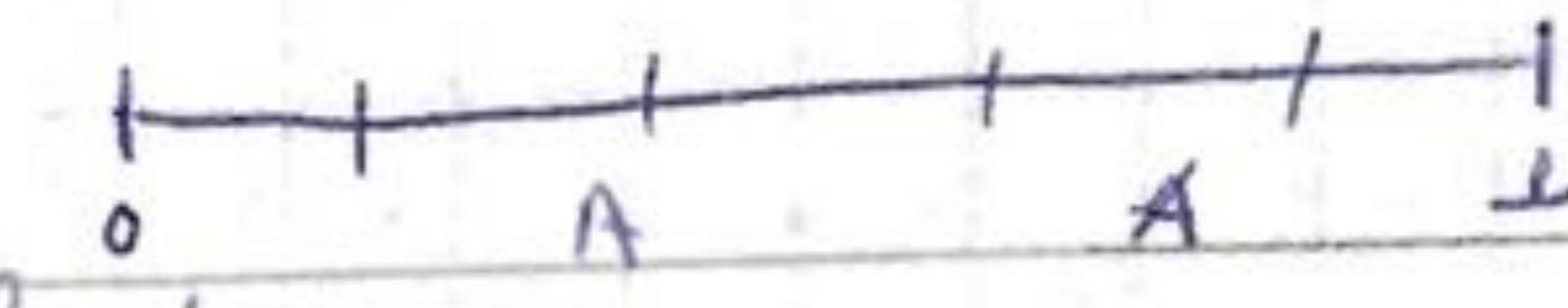
L8: Generating RVs Markov Chains: Pseudo $\xrightarrow{\text{behaves as it was random}}$ Random Number Generation (PRNG)

Minimum requirements (Naive): to Naive Uniform $[\underline{0}, \overline{1}]$

Create U_0, U_1, \dots, U_n for any set $A \subseteq [\underline{0}, \overline{1}]$ $\frac{N_n(A)}{n} \rightarrow \int_A dx$

$N_n(A)$ is the number of $U_i \in A$ $0 \leq i \leq n-1$

$$X \sim \text{unif}([\underline{0}, \overline{1}]) \quad P(X \in A) = \int_A dx.$$



To hard to do $\text{unif}([\underline{0}, \overline{1}])$ so let's do uniform $M = \{\underline{0}, \overline{1}, \dots, M-1\}$

Def: A sequence U_0, U_1, \dots is pseudorandom if for any number $a \in M$: $\frac{N_n(a)}{n} \xrightarrow{\text{Frequency}} \frac{1}{M}$

(every number should appear an equal amount of times).

e.g. $0, 1, 0, 1, \dots$ is pseudo random on $\{\underline{0}, \overline{1}\}$

$1, 1, 0, 0, 1, 1, 0, 0, \dots$ — //

Def: A congruential generator with parameters (a, b, M) : $D(x) = (ax + b) \bmod M$.
start at some U_0 -random seed $U_i = D(U_{i-1})$: $U_1 = D(U_0)$, $U_2 = D(U_1)$

Ex. 1: $a=2, b=1, M=3$ $U_0 = 1, U_1 = (2 \cdot 1 + 1) \bmod 3 = 0, U_2 = (2 \cdot 0 + 1) \bmod 3 = 1$

$U_0 = 2 : U_1 = (2 \cdot 2 + 1) \bmod 3 = 2$

Ex. 2: $a=3, b=1, M=7$ $U_0 = 1 : U_1 = (3 \cdot 1 + 1) \bmod 7 = 4$

$$U_1 = 4 : U_2 = (3 \cdot 4 + 1) \bmod 7 = 6$$

$$U_2 = 6 : U_3 = (3 \cdot 6 + 1) \bmod 7 = 5$$

$$U_3 = 5 : U_4 = (3 \cdot 5 + 1) \bmod 7 = 2$$

$$U_4 = 2 : U_5 = (3 \cdot 2 + 1) \bmod 7 = 1$$

$$U_5 = 1$$

$$U_6 = 3 : U_7 = (3 \cdot 3 + 1) \bmod 7 = 3$$

Def: The period of D starting at U_0 is the smallest number T_0 such that $U_{i+T_0} = U_i$, for all i .

Ex. 1: Period of D starting in 1 is 2 and if you start from $U_0 = 2$ is 1

Ex. 2: $U_0 = 1$, period 6, $U_0 = 3$ period 1.

Full period on $\{\underline{0}, \overline{1}, \dots, M-1\}$ is M

Start at $U_0 = 0$ then after M steps we get 0

$$\frac{N_n(o)}{n} \text{ or } \frac{N_{k \cdot M}(o)}{k \cdot M} = \frac{k}{k \cdot M} = \frac{1}{M} \quad \frac{N_1(o)}{n} \rightarrow \frac{1}{M}$$

choose a very large M and assume you have a full period V_0, V_1, \dots, V_n . Fix K such that K divides M and construct: $V_i = [V_i \cdot \frac{1}{M}] \rightarrow$ floor you take just the integer and leave behind the decimals (same as using $\text{int}(\frac{1}{M})$ in python). V_i is pseudorandom on $\{\underline{0}, \overline{1}, \dots, K-1\}$

The period of V_i is M (good candidate for randomness).

To create uniform $[\underline{0}, \overline{1}]$ we take $V_i = \frac{V_i}{M}$

Sampling: say we want to sample from F , calculate F^{-1}

Generate Uniform $([\underline{0}, \overline{1}])$ and consider $X = F^{-1}(V)$ then $X \sim F$

Lecture: Introduction to Data Science

Regression: Consider a pair $(x, y) \sim F_{xy}$

$$L(a, b) = (a - b)^2$$

Let's say we have a guess f : our guess

$$R(f) = \mathbb{E}[L(\bar{Y}, f(\bar{x}))] = \mathbb{E}[(\bar{Y} - f(\bar{x}))^2]$$

$$\text{Define } r(x) = \mathbb{E}[\bar{Y} | \bar{x} = x]$$

$$R(f) = \mathbb{E}[(Y - r(\bar{x})) + (r(\bar{x}) - f(\bar{x}))]^2$$

$$\begin{aligned} &= \mathbb{E}[(Y - r(\bar{x}))^2] + \mathbb{E}[(r(\bar{x}) - f(\bar{x}))^2] + 2\mathbb{E}[(Y - r(\bar{x}))(r(\bar{x}) - f(\bar{x}))] \\ &= \mathbb{E}[\mathbb{E}(Y | \bar{x})(r(\bar{x}) - f(\bar{x}) | \bar{x})] \\ &= \mathbb{E}[\mathbb{E}[Y | \bar{x}] - r(\bar{x})](r(\bar{x}) - f(\bar{x})) = 0 \end{aligned}$$

$$\text{Conclusion: } R(f) = \mathbb{E}[(Y - r(\bar{x}))^2] + \mathbb{E}[(r(\bar{x}) - f(\bar{x}))^2]$$

If you change f : nothing happens to the first term. If you choose $r = f$: the second term will be zero \Rightarrow will be the best that you can do.

$R(r)$ is the smallest risk.

you cannot do anything about it

①: calculates the measurement noise (how fat it is) "noise variance"

②: calculates the bias² (you can do something about it)

$r \rightarrow$ is called a regression function

Find "f" Now we have $\mathbb{E}[(Y - r(\bar{x}))^2] = 0$ this can happen when you know the real label.

good for image classification (e.g. NN) and text generation - sequence of words.

(Usually much easier)

Pattern Recognition: $Y \in \{0, 1\}$ $L(a, b) = \begin{cases} 1 & \text{if you make a mistake } a \neq b \\ 0 & \text{otherwise (no mistake)} \end{cases}$

$r(x) = \mathbb{E}[\bar{Y} | X = x] = P(Y=1 | X=x)$ \hookrightarrow class \hookrightarrow images \hookrightarrow true probability that class is in the image based on the image.

$$h^*(x) = \begin{cases} 1 & \text{if } r(x) > 1/2 \\ 0 & \text{if } r(x) \leq 1/2 \end{cases}$$

Bayes Classifying Rule

Property: $R(h^*) \leq R(h)$ for any guessing function h .

Maximum Likelihood Method:

Let $S = \{P_a(x) : a \in \mathbb{R}^d\}$ densities; it can be a Gaussian or any other distribution

a: parametric model

Loss: $L(x, a) = \ln P_a(x)$ log-loss

$$X \sim P_a \in S \quad R(a) = \mathbb{E}[L(X, a)] = \mathbb{E}[-\ln P_a(x)] = \int \ln P_a(x) \cdot P_{a^*}(x) dx$$

we cannot calculate the risk but we can calculate the empirical risk:

Empirical Risk: $x_1, \dots, x_n \stackrel{iid}{\sim} P_{\alpha^*}$

$$\hat{R}(\alpha) = \frac{1}{n} \sum_{i=1}^n -\ln P_{\alpha}(\bar{x}_i)$$

neg-log-likelihood

This is also called entropy - we are measuring $R(\alpha) - R(\alpha^*)$

Regression: If we assume $f_{XY} = f_{Y|X} \cdot f_X \Leftrightarrow f_{XY} = \ln f_{Y|X} + \ln f_X$

If we think that $f_{Y|X} = \text{Bern}(x|y) P_{\alpha^*}(y|x)$

Statistical Model is assumed to be $\mathcal{F} := \{P_{\alpha}(x,y) : P_{\alpha}(x) = P_0(x)$
 $P_{\alpha}(y|x) = f_{\alpha}(x,y)\}$

Ex. Linear regression

$$P_{\alpha}(y|x) = \frac{1}{\sqrt{2\pi \cdot \sigma^2}} \cdot e^{-\frac{(y-(ax+b))^2}{\sigma^2}}$$

$$a = (a, b, \sigma)$$

→ just a number

usually we do an assumption that this is poisson distribution
and plug it as below.

$$-\ln P_{\alpha}(y|x) = \ln \sqrt{2\pi} + \ln |\sigma| + \frac{(y - (ax+b))^2}{\sigma^2} \rightsquigarrow \text{only these part depends on the data.}$$

doesn't depend
on data (a)

$$\hat{R}(\alpha) = \ln |\sigma| + \frac{1}{n} \sum_{i=1}^n \frac{(\bar{y}_i - (a \cdot \bar{x}_i + b))^2}{\sigma^2} + \ln \sqrt{2\pi}$$

always same result when
using Gaussian

Logistic Regression: $P_{\alpha}(y|x) = \text{Bernoulli}(p)$

$$p = \frac{1}{1+e^{-(ax+b)}} = G(ax+b)$$

$$P_{\alpha}(y|x) = p^y (1-p)^{1-y}$$

$$G(x) = \frac{1}{1+e^{-x}} \rightsquigarrow \text{logistic function}$$

$$\ln P_{\alpha}(y|x) = y \ln p + (1-y) \cdot \ln (1-p)$$

when you guess the p then you plug it to the pattern recognition part and
you get 1 or 0 based on the result.

prob. that $X \leq x = P(\underline{X} \leq x) = P(F^{-1}(u) \leq x) = P(u \leq F(x)) = F(x)$

because $F_u(u) = u$

$$= 1 - e^{-\lambda x}$$

$$\Rightarrow y = 1 - e^{-\lambda x}$$

$$\Rightarrow e^{-\lambda x} = 1 - y \quad (\Rightarrow x = \frac{1}{\lambda} \ln(1-y) = F^{-1}(y))$$

Inversion Sampling

(works for discrete too - but use a different def for discr)

Accept-Reject Method: Input target density f

sampling density g

$$f(x) \leq Ng(x) \text{ for all } x.$$

Repeat: Draw $\underline{X} \sim g$ and calculate $r(\underline{X}) = f(\underline{X}) / Ng(\underline{X})$, Draw $U \sim \text{Uniform}([0, 1])$ and accept \underline{X} if $U \leq r(\underline{X})$ → accept or reject.

$$X \sim G \quad U \sim \text{unit } [0, 1]$$

$$I = \begin{cases} 1 & \text{if } U \leq r(x) \\ 0 & \text{otherwise} \end{cases} \quad \text{if we accept then } I=1 \quad \text{distribution that we are looking for: } P(X \leq x | I=1)$$

$$\text{Bayes' thm: } f_{\underline{X}/I} (x|i) = \frac{f_{I/\underline{X}}(i/x) \cdot f_{\underline{X}}(x)}{f_I(i)}$$

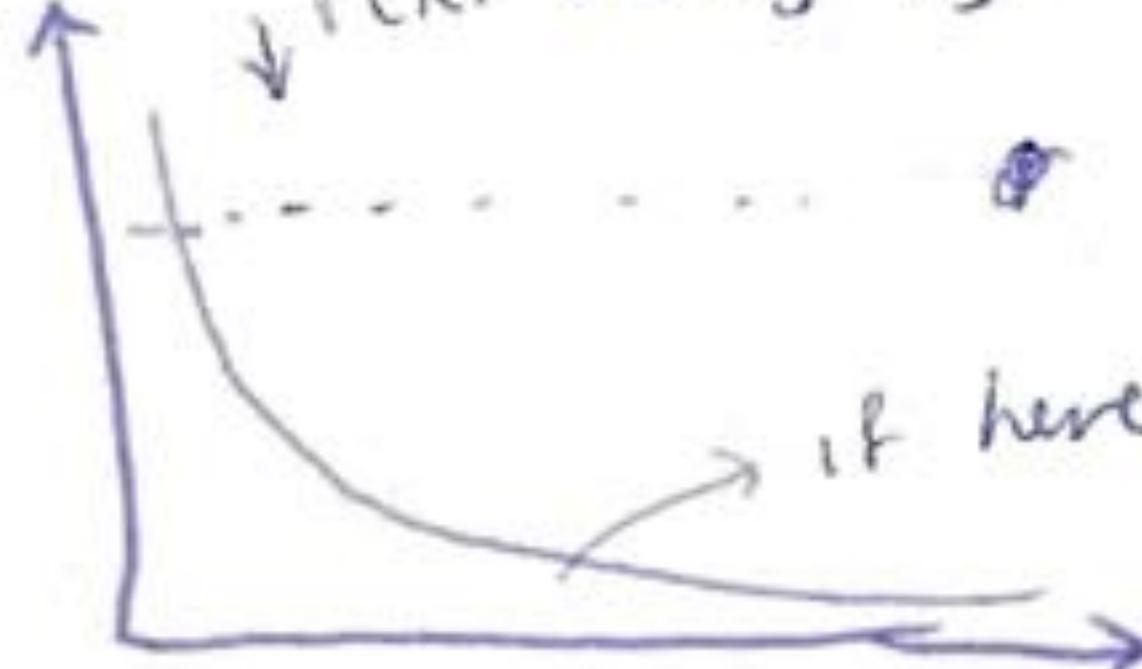
$$f_{I/\underline{X}}(i/x) = P(U \leq r(\underline{X}) | X=x) = r(x)$$

$$f_{\underline{X}}(x) = g(x)$$

$$f_I(i) = \text{using the law of total probability: } \int f_{I/\underline{X}}(i/x) \cdot f_{\underline{X}}(x) dx = \int_{\underline{X}=i}^{\underline{X}} \frac{r(x) \cdot g(x)}{N \cdot g(x)} dx = \frac{1}{N}$$

$$f_{\underline{X}/I}(x|i) = \frac{r(x) \cdot g(x)}{1/N} = f(x) \quad \text{That is } P(\underline{X} \leq x | I=i) = F(x)$$

If density is very very thin and sampling distribution is uniform,



Lab session: Random number generation:

1) inverse transform method: $u \sim U(0,1)$, $F^{-1}(u) = F \sim \text{target distribution}$

2) Box-Muller algorithm: $u_1, u_2 \sim U(0,1)$ - iid. $Z_0 = \sqrt{-2 \ln(u_1)} \cdot \cos(2\pi \cdot u_2)$

$$Z_1 = \sqrt{-2 \ln(u_1)} \cdot \sin(2\pi \cdot u_2)$$

$$Z_0, Z_1 \sim N(0,1)$$

$$Z = Z_0 \cdot \sigma + \mu$$

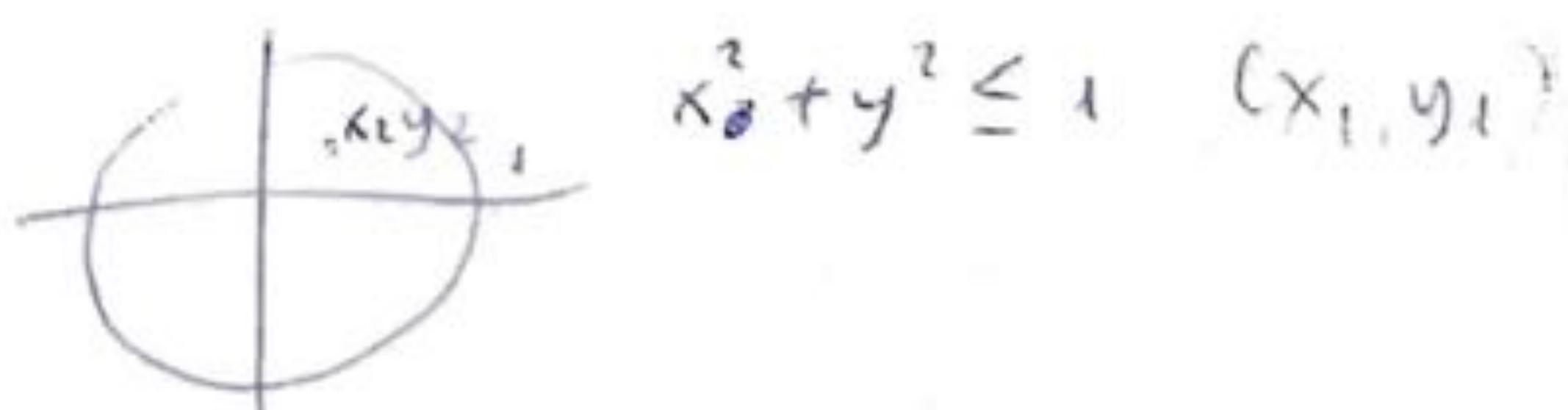
$$Z \sim N(\mu, \sigma^2)$$

3) Accept-reject: target distribution $f(x)$ & sampling distribution $g(x)$
 $f(x) \leq Cg(x)$

Method: generate x from $g(x)$ $r(x) = \frac{f(x)}{Cg(x)}$
 generate $u \sim U(0,1)$ if $u \leq r(x)$ $x^n = x^{\theta}$

else:

reject generation until $u \leq r(x)$



Markov's Chains:

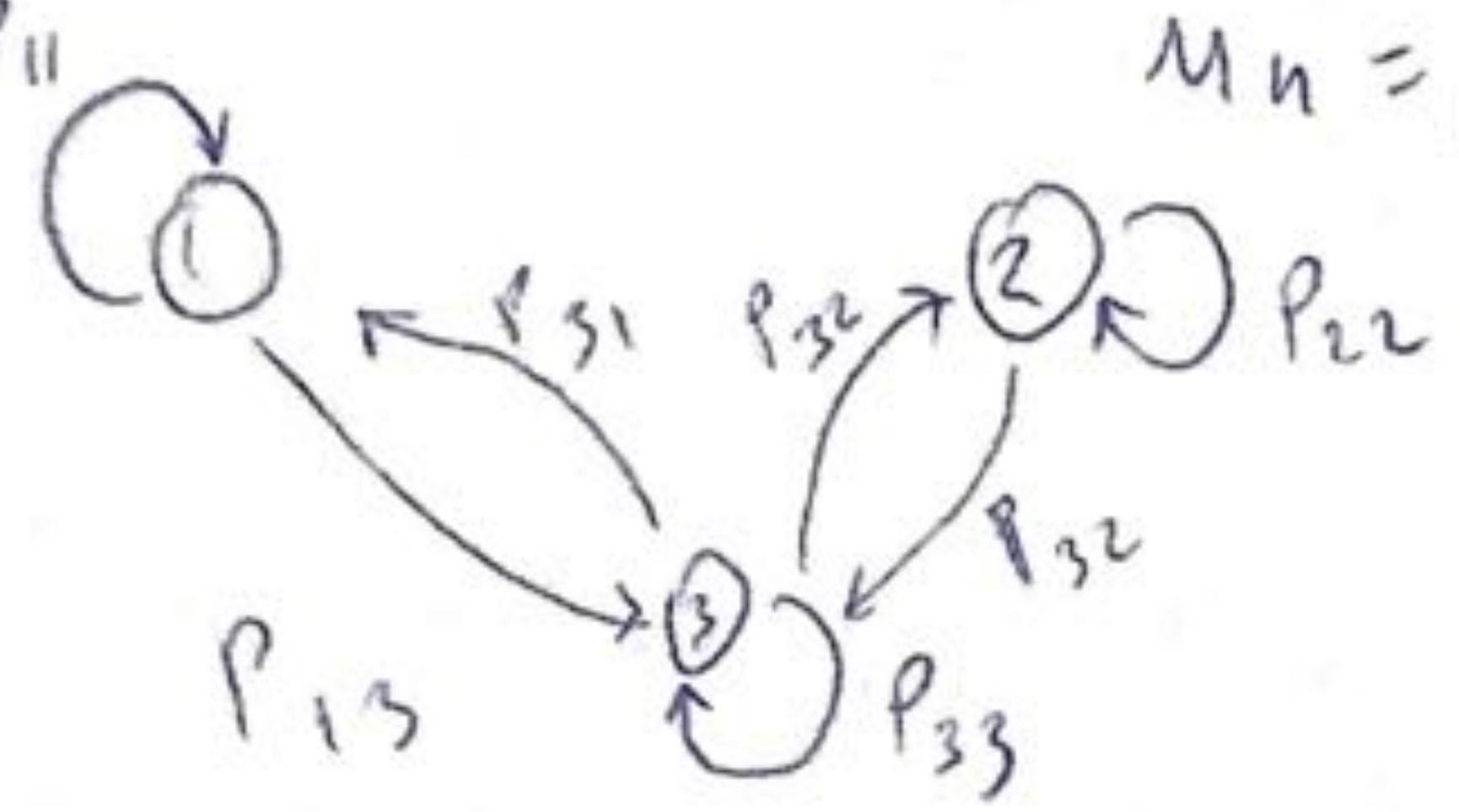
→ specialise of stochastic process.

finite MC : $X_n \sim$ stochastic probability

X_n is a MC if $P(X_{t+1} \leq x | f_t) = P(X_{t+1} \leq x)$

MC is defined with a transition matrix.

$$M_n = \{M_{ij}\}$$



Key Properties: $\sum_j P_{ij} = 1$ for any i

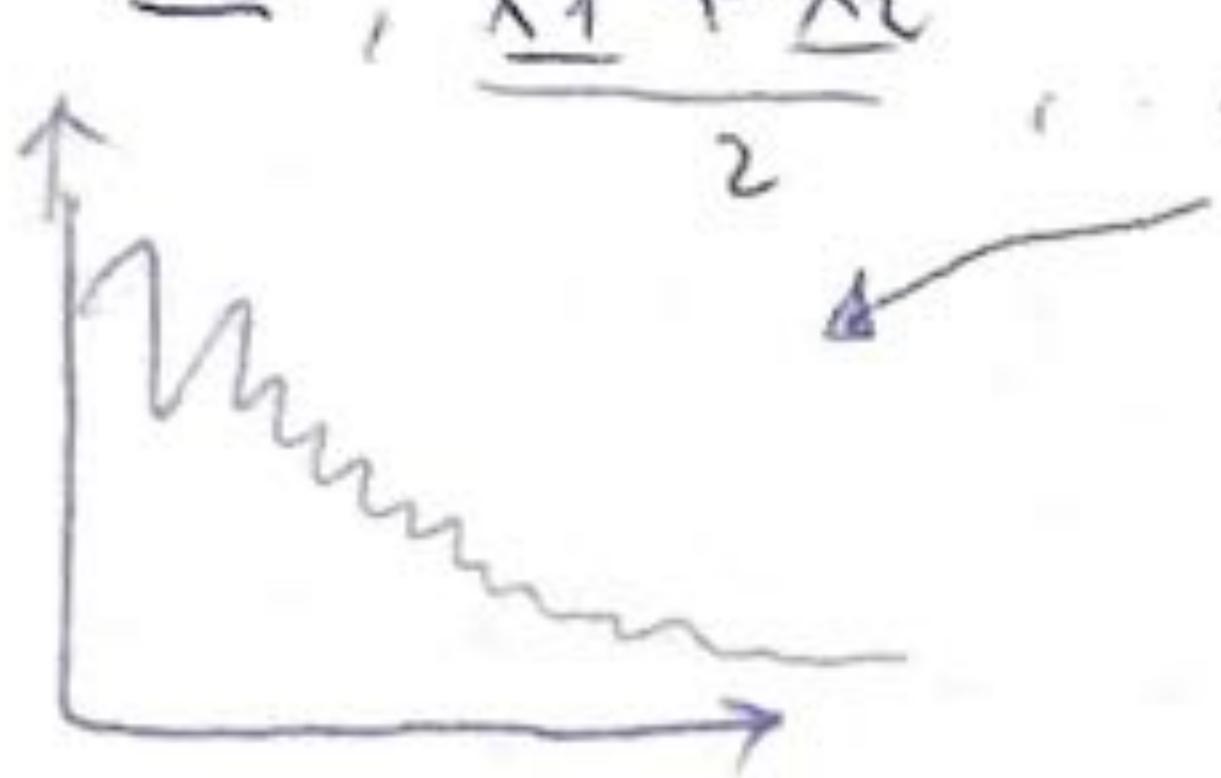
Definition: Lecture 9: ~~BBB~~ Markov's Chain

Def: A stochastic process is indexed set of random variables $\{X_i\}_{i \in \mathbb{I}} \rightarrow$ index can be a time or any number.

If $\mathbb{I} = \{1, 2, 3, \dots\}$ we say its discrete

Ex: $\underline{X_1}, \underline{X_2}, \dots, \underline{X_n} \sim F$ (taking independent samples) \rightarrow you cannot shuffle it as you will destroy it.

Ex: $\underline{X_1}, \underline{\frac{X_1 + X_2}{2}}, \dots, \frac{1}{n} \cdot \sum_{i=1}^n \underline{X_i}$ all connected as the next value depends on the previous.



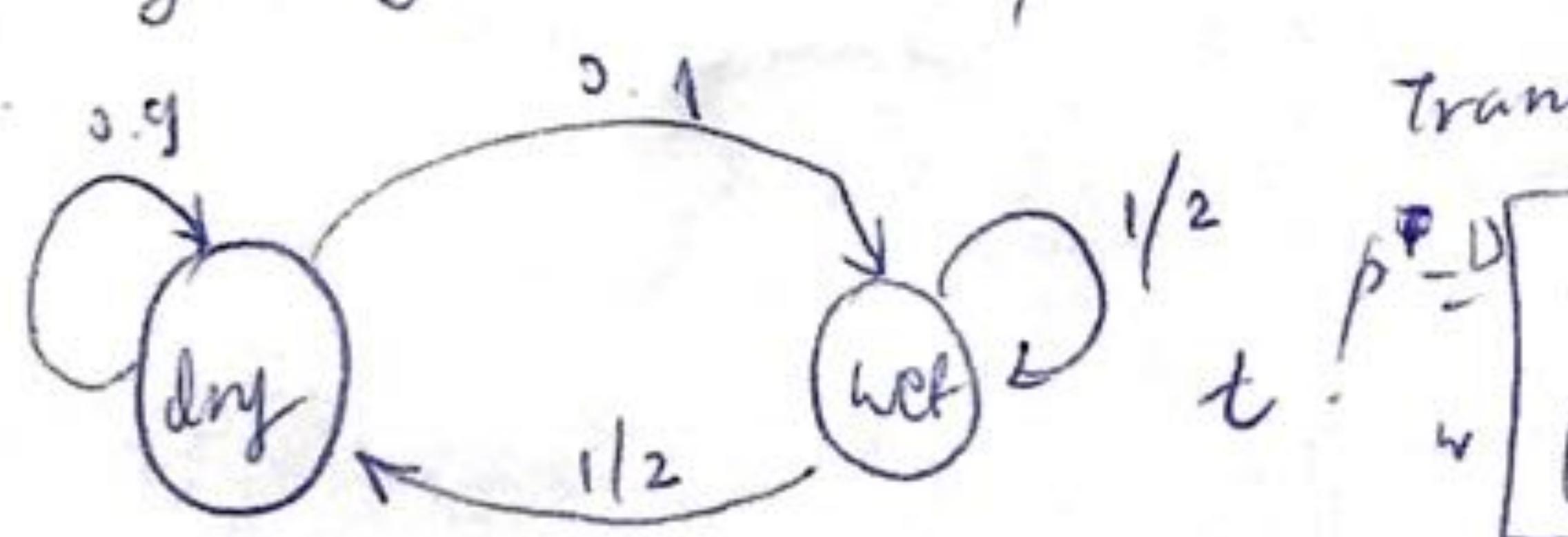
Def: Markov process: $\{X_i\}_{i \in \mathbb{N}}$, $\underline{X_i} \in \underline{\mathcal{S}}$ \rightarrow fix this to be the same in a Markov process

$$\underline{\mathcal{S}} = \{s_1, \dots, s_n\}$$

$P(X_t = y | X_1, X_2, \dots, X_{t-1}) = P(X_t = y | X_{t-1})$ as "Markov Property" (holds a memory of me)

"Memory only of one step"

Ex: 3.9



Transition diagram:

$$P = \begin{bmatrix} P(X_t = \text{dry} | X_{t-1} = \text{dry}) & P(X_t = \text{wet} | X_{t-1} = \text{dry}) \\ P(X_t = \text{dry} | X_{t-1} = \text{wet}) & P(X_t = \text{wet} | X_{t-1} = \text{wet}) \end{bmatrix}$$

$t-1$

$$\begin{bmatrix} P(X_t = \text{wet} | X_{t-1} = \text{dry}) \\ P(X_t = \text{wet} | X_{t-1} = \text{wet}) \end{bmatrix} = \begin{bmatrix} 0.9 & 0.5 \\ 0.1 & 0.5 \end{bmatrix}$$

Say that we know:

$$P_{t-1} = [P(X_{t-1} = \text{dry}), P(X_{t-1} = \text{wet})]$$

$P(X_t = \text{dry}) =$ if you are on the dry, you could have been previously to wet or dry and using the law of total probability $= P(X_t = \text{dry} | X_{t-1} = \text{dry}) \cdot P(X_t = \text{dry}) + P(X_t = \text{dry} | X_{t-1} = \text{wet}) P(X_{t-1} = \text{wet})$

$$P(X_t = \text{wet}) = P(X_t = \text{wet} | X_{t-1} = \text{dry}) \cdot P(X_{t-1} = \text{dry}) + P(X_t = \text{wet} | X_{t-1} = \text{wet}) P(X_{t-1} = \text{wet})$$

P_{t-1} = (1, 2) vector

$$\underline{P} = (2, 2) \quad P_{t-1} \cdot \underline{P} = (1, 2)$$

$$P_t = [P(X_t = \text{dry}), P(X_t = \text{wet})] = P_t = P_{t-1} \cdot \underline{P}$$

$$\text{Recursive formula: } P_t = P_{t-1} \cdot \underline{P} = (P_{t-2} \cdot \underline{P}) \cdot \underline{P} = P_0 \cdot \underline{P}^t$$

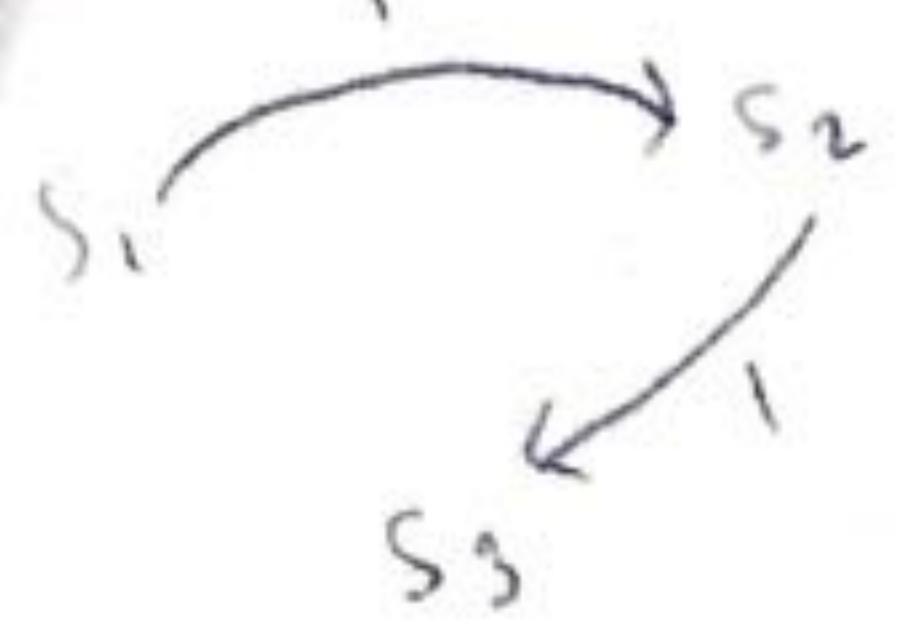
$$P_0 = [0, 1] \rightarrow \text{we start in the wet state}$$

$$P_0 = [0, 1] \cdot \underline{P}^0$$

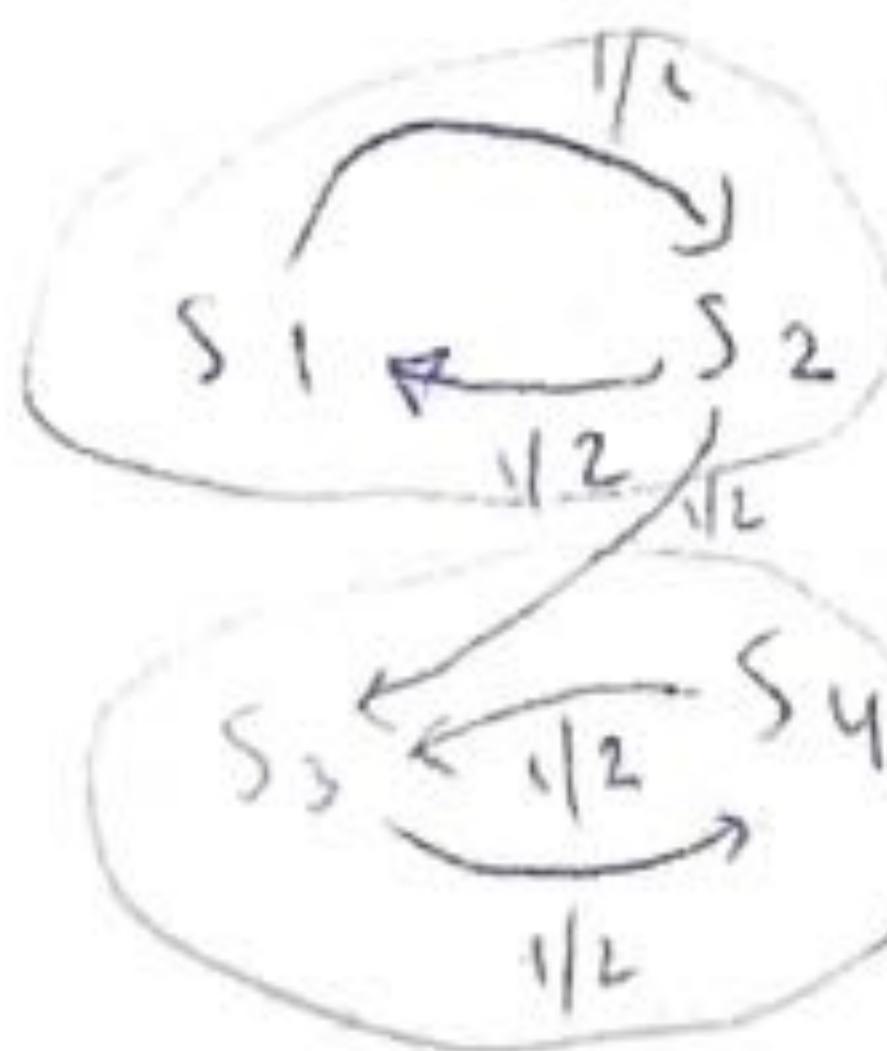
Homogeneous Markov chains we can do something as the transition matrix is the same over time.

$\{s_1, \dots, s_n\}$ we say $s_i \rightarrow s_j$
 s_i communicates with s_j if $P(X_t = s_j | X_0 = s_i)$ for some t .

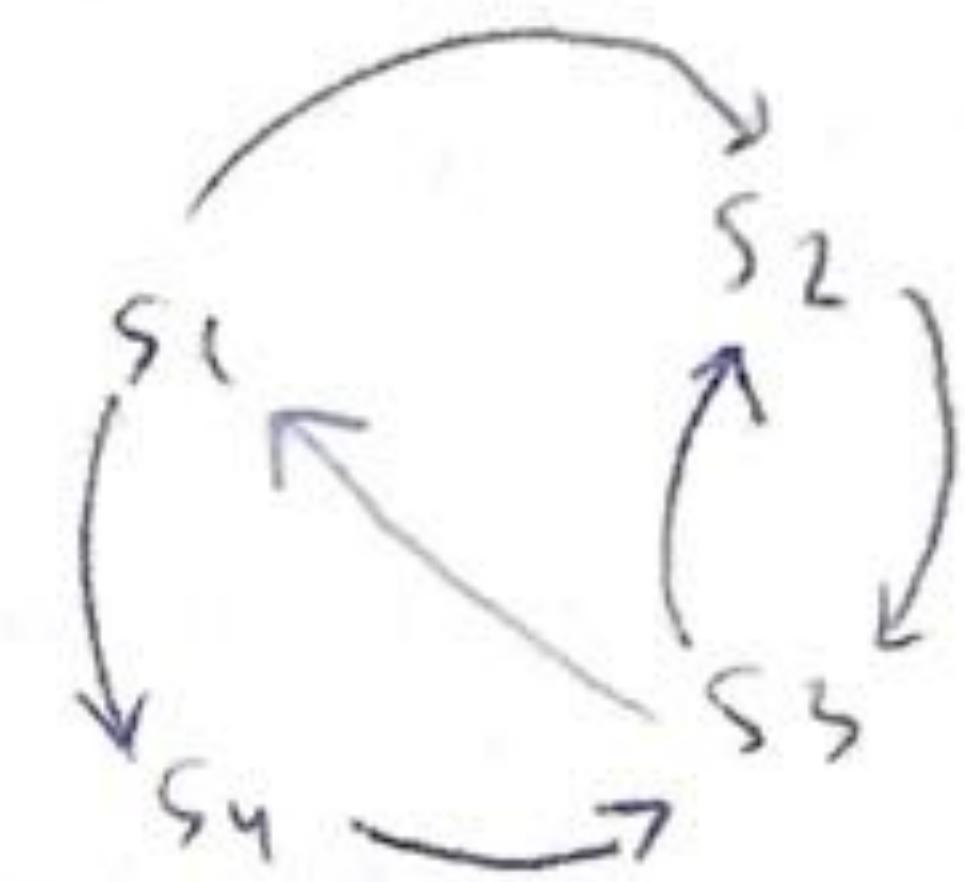
$s_i \rightarrow s_j$ and $s_j \rightarrow s_i$ intercommunicate



$$\begin{aligned} s_1 &\rightarrow s_2 \\ s_2 &\rightarrow s_1 \\ s_1 &\rightarrow s_3 \\ s_3 &\rightarrow s_1 \end{aligned}$$



"two clauses"
 check if there is
 a place you can
 get stuck



when everything communicates we say chain is irreducible.

you can stuck in a loop:

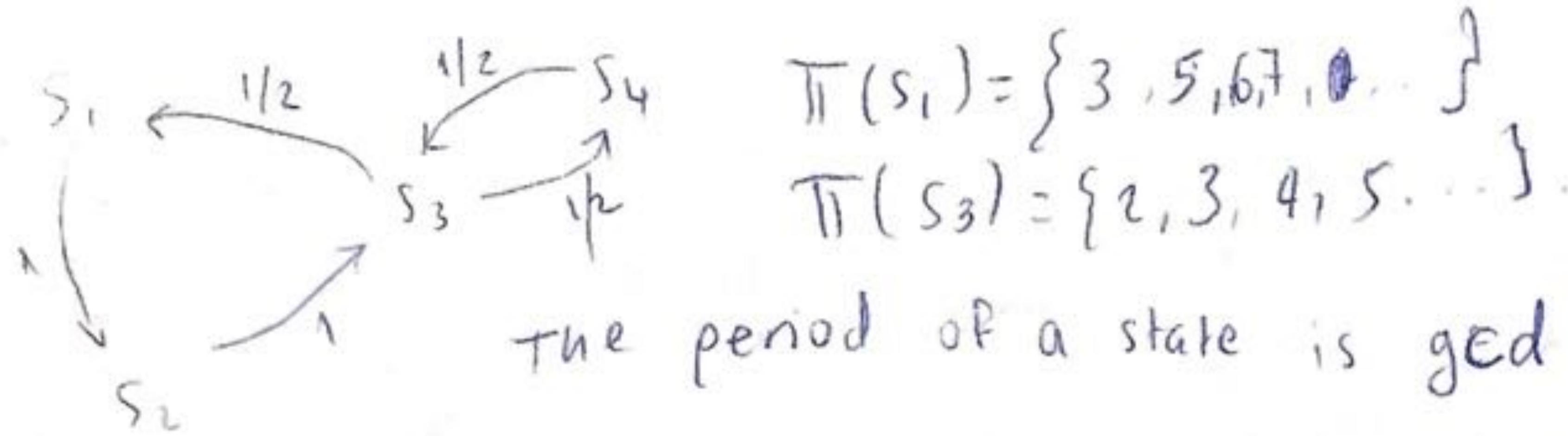


you will just going around. \rightarrow avoid having

returning times:

to mixing \rightarrow problematic

$$\Pi(x) = \{t \in \mathbb{N} : (P^t)_{xx} > 0\}$$



$$\begin{aligned} \Pi(s_1) &= \{3, 5, 6, 7, \dots\} \\ \Pi(s_3) &= \{2, 3, 4, 5, \dots\} \end{aligned}$$

The period of a state is $\text{gcd}(\Pi(x))$

we say that a chain is aperiodic if $\text{gcd}(\Pi(s_i)) = 1$ for all states

Theorem: If irreducible and aperiodic then $\underline{P}_0 P^t \xrightarrow[t \rightarrow \infty]{\text{mixing}} \Pi \text{ unique}$

also $\Pi = \Pi \cdot P$

\rightsquigarrow stationary distribution if this holds

$$\Pi^T = (\Pi \cdot P)^T = P^T \cdot \Pi^T \rightsquigarrow \text{eigenvector of } P^T$$

* remember to take P^T in the computer
 with an eigenvalue of 1.)

Perron-Frobenius theorem $\Pi_i > 0$ for all i

because the ~~row~~ was sum to 1.
 (columns if transpose)

how fast it goes to the stationary distribution?

$$(P_0 - \Pi) \cdot P = (\sum_i (P_0 - \Pi) \cdot v_i - \Pi) \cdot P = \sum_{i=1}^N (P_0 \cdot v_i) \lambda_i - \Pi = \sum_{i=1}^N \lambda_i \cdot (P_0 \cdot v_i)$$

$$\| (P_0 - \Pi) \cdot P \| \leq \lambda_2 \| (P_0 - \Pi) \| \rightsquigarrow \text{converges exponentially fast}$$

"Exponentially convergence to equilibrium"

Lab 3/11 Markov Chains:

1) Irreducibility: MC $x = \{s_0, \dots, s_n\}$ is irreducible if:

$$\begin{aligned} P(X_t = s_j | X_0 = s_i) &> 0 \\ P(X_t = s_i | X_0 = s_j) &> 0 \quad s_i, s_j \in \mathbb{R}^2 \end{aligned}$$

2) Aperiodicity: A period of the state - greatest common division of $T(x)$.

$$T(x) = \{t \in \mathbb{N} : p_t(x, x) > 0\}$$

Stationary Distribution: MC: $x = \{s_0, \dots, s_n\}, P = (P(x, y))$ has a stationary distribution $\Pi = \{\Pi(s_0), \Pi(\dots), \Pi(s_n)\}$

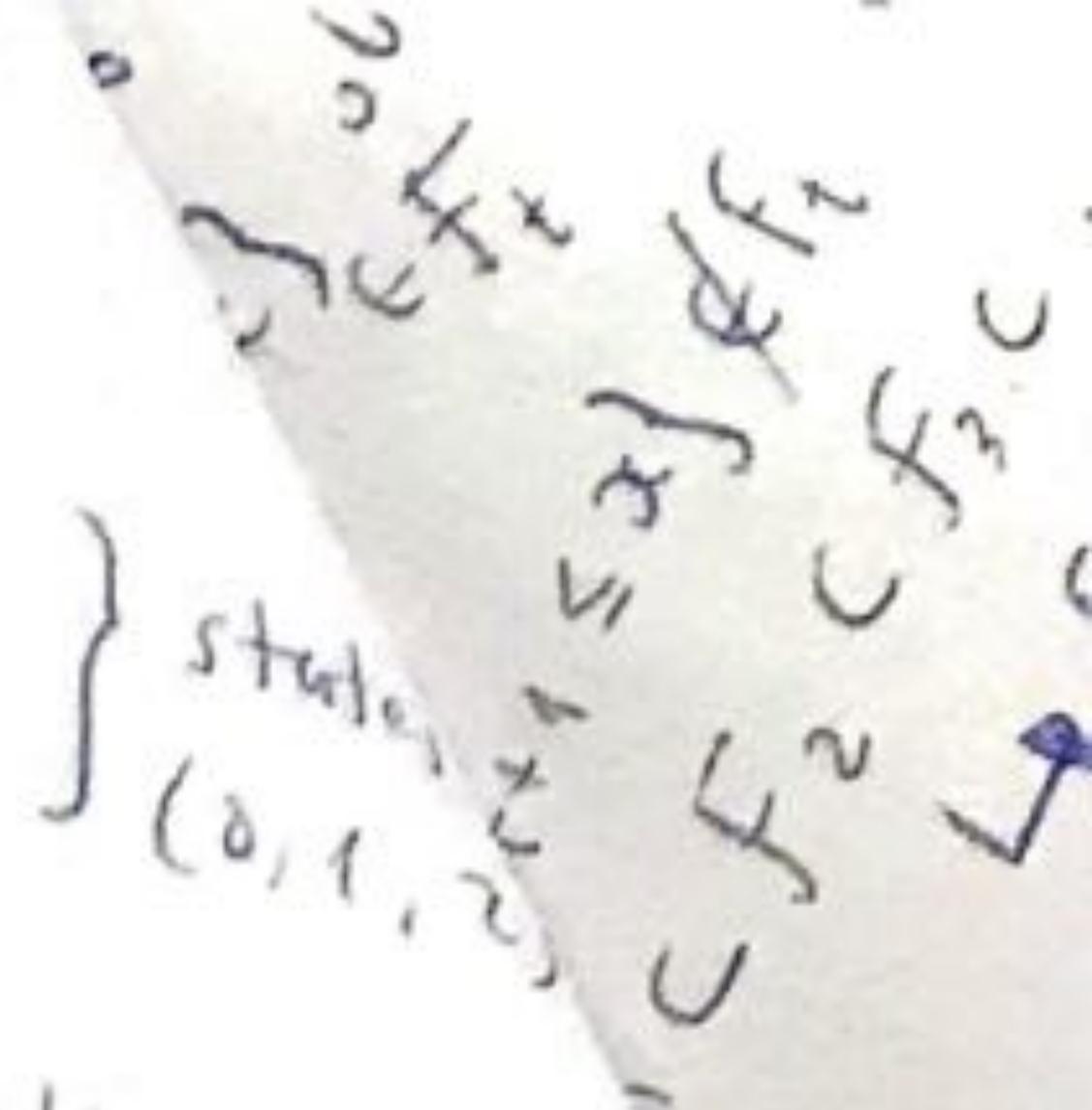
if 1) $\Pi(x) > 0$ for any x , $\sum \Pi(\lambda) = 1$

2) $\Pi \cdot P = \Pi$ (fixed point equation)

Homogeneity: $P_{ij}^n = P_{ij}^t$, any $n, t \in N$

exercises from lecture notes: T.35:

- splits in two
- stays the same
- dies



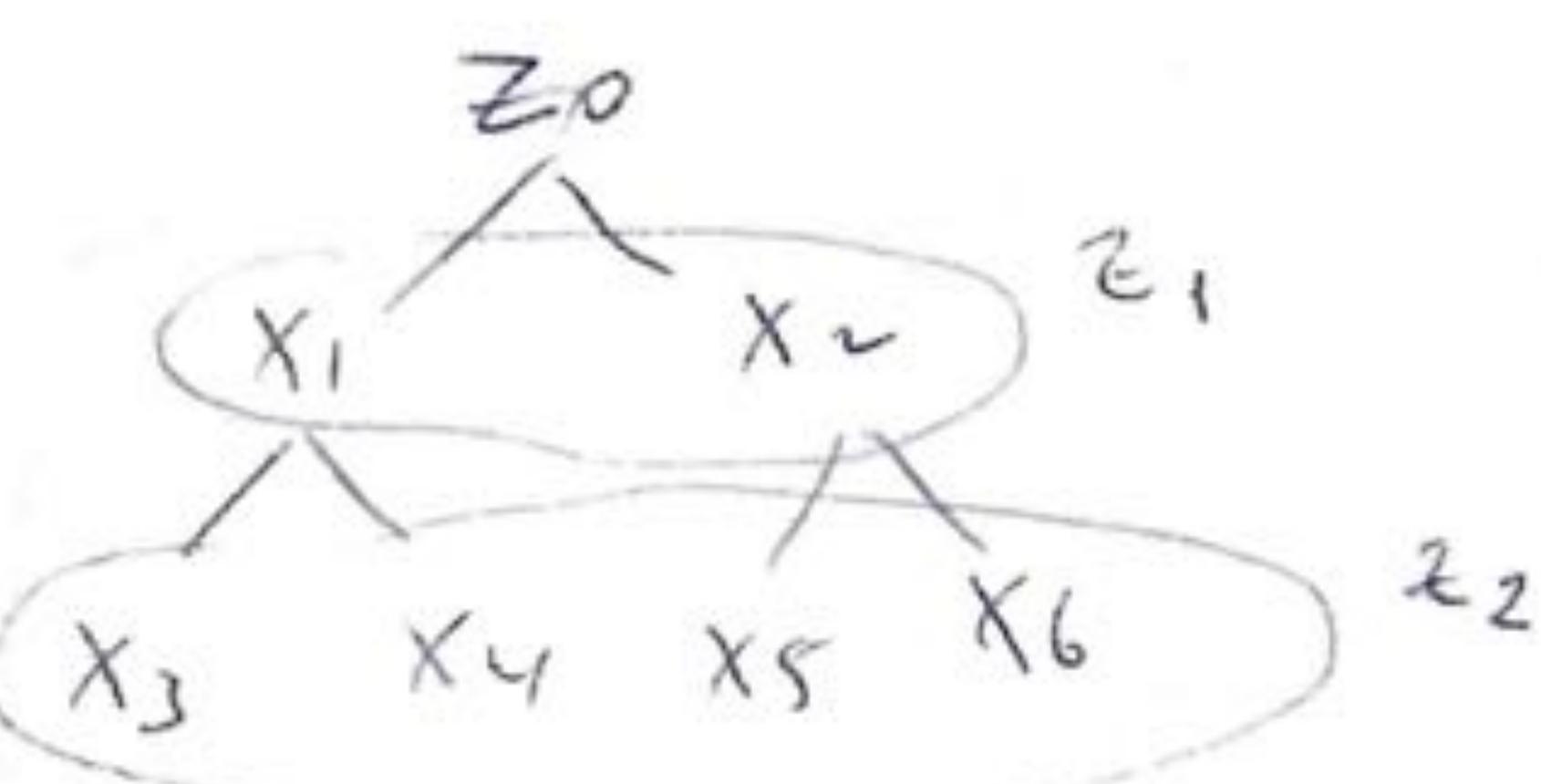
$X \sim$ offspring from each bacteria

$$Z_{n+1} = \sum_{i=1}^n \text{Avg } X_{n,i}$$

1) is Z_{n+1} a MC? 2) simulate the process. 3) Estimate the expectancy

it is a MC as it depends only on the previous state

Galton-Watson process



$X \sim$ # offspring each node generates

$$\bullet p = \{1/3, 1/3, 1/3\}$$

Extinction prob ($\mathbb{P}(X)$) = 1

$$\text{Ex. } < 1 \Rightarrow \lim_{n \rightarrow \infty} \mathbb{P}(Z_n = 0) = 1$$

$$\text{Ex. } > 1 : \lim_{n \rightarrow \infty} \mathbb{P}(Z_n = 0) < 1$$

13/11 Repetition Chapter 1-3

Experiment: A procedure that produces outcomes distinct.

Set of all outcomes: Ω "sample space" $\Omega = \{\omega\}$ outcomes from the experiment

Trial: One round of an experiment: \rightarrow n-fold trial (number of times you do the trial)

Events: subset of Ω , $A \subset \Omega$

A happens when $\omega \in A$

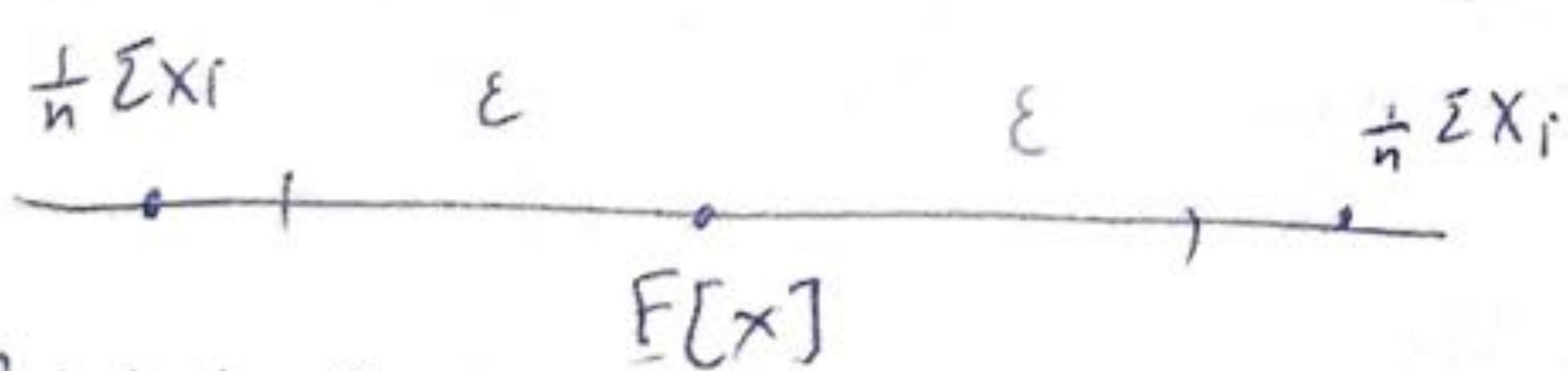
Ex. win loss $\Omega = \{0, 1\}$, $A = \{0\}$ = "we get a 0 in the trial"

Events A, B : $A \cup B = A$ or B , $A \cap B = "A \text{ and } B"$

$\mathbb{P}(\max\{X_1, X_2\} \leq x) = \mathbb{P}(\{X_1 \leq x\} \cap \{X_2 \leq x\})$

then is again OR. and if \geq then it is OR. and if \leq

$$\mathbb{P}\left(1 \frac{1}{n} \sum X_i - \mathbb{E}[X_i] / \epsilon > \epsilon\right) = \mathbb{P}\left(\frac{1}{n} \sum X_i - \mathbb{E}[X] / \epsilon > \epsilon\right) \cup \mathbb{P}\left(\mathbb{E}[X] - \frac{1}{n} \sum X_i > \epsilon\right)$$



→ continue this to find the Hoeffding inequality

Def: let Ω be a set then we call F (sigma(σ)-algebra), observables). F : set of events/subsets that are observable. To write it down you need 3 rules: 1) $\Omega \in F$ (check if the trial happened), 2) If $A \in F \Rightarrow A^c \in F$ (check if the trial did not happen), 3) If $A_1, A_2, \dots, A_n \in F \Rightarrow \bigcup A_i \in F$ (at least one happens)

x. MARKOV's chain: X_1, X_2, \dots, X_n .

$$F_i = \{\omega : \underline{X}_1 = x_1, \underline{X}_2 = x_2, \dots, \underline{X}_i = x_i \mid \text{for all } \underline{x}_{i+1}, \underline{x}_{i+2}, \dots, \underline{x}_n \in X\}$$

time t , that is F_t

observable

$\in F_t$

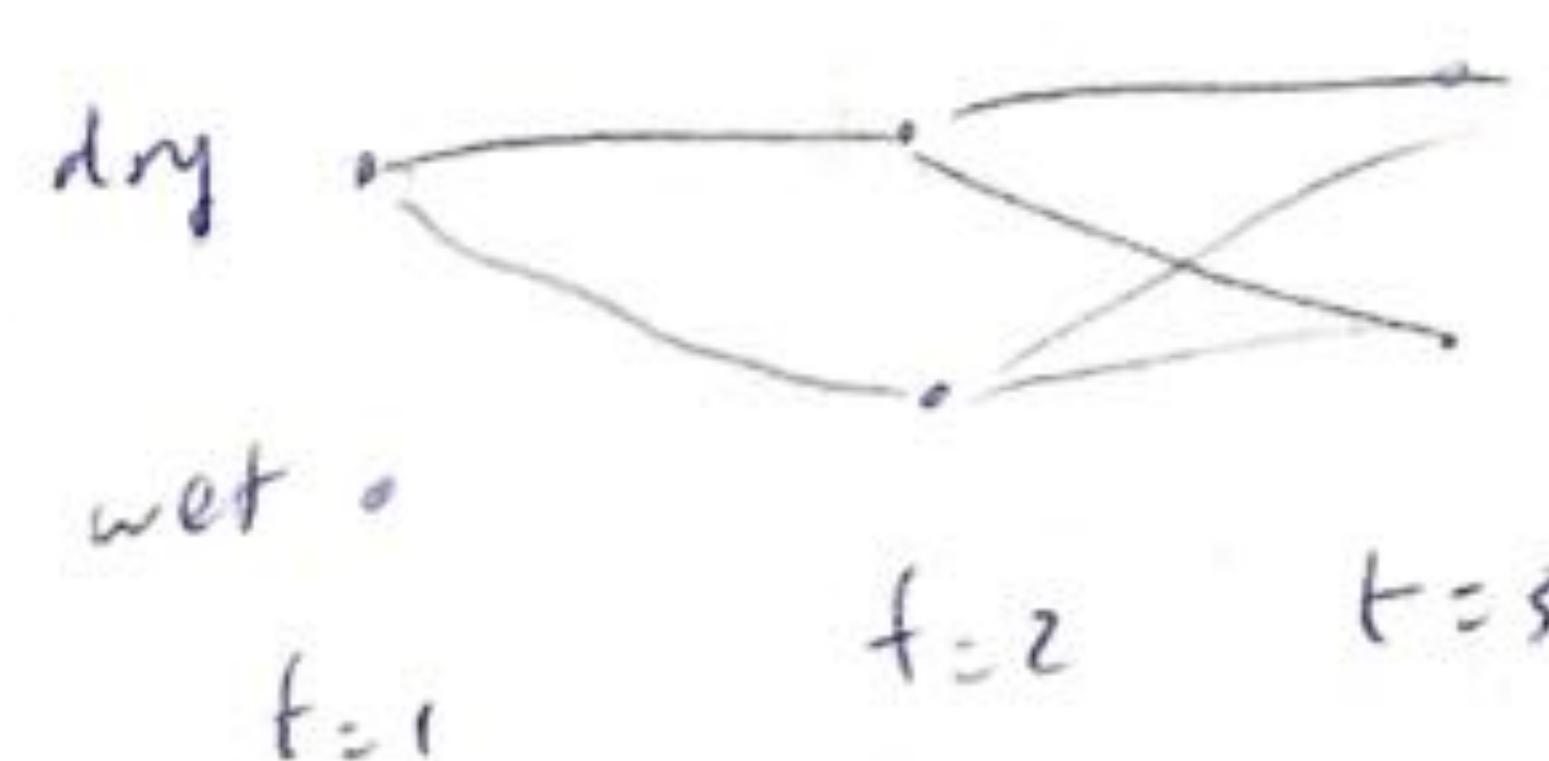
$\{X_i \leq x\} \notin F_t$ not observable

$\subset F_2 \subset F_3 \subset \dots$ (It is growing, more time the + you get more information)

↳ Filtration

$$F_1 = \{\{X_1 = 1\}, \{X_1 = 0\}, \{X_1 = 1 \cup X_1 = 0\}\}$$

$$\text{for ex. } \{X_1 = 1\} = \{w; X_1 = 1\} = \bigcup_{i=0}^{\infty} \{w; X_1 = 1, X_2 = i\} \in F_2$$



Def: Let Ω be a set and given an F (observables) then: $P: F \rightarrow [0, 1]$.

1) $P(\Omega) = 1$ 2) $A, B \in F$ $A \cap B = \emptyset$ then $P(A \cup B) = P(A) + P(B)$

3*) Rule 2: holds for sequences: $A_1, A_2, \dots \in F$ $A_i \cap A_j = \emptyset$ if $i \neq j$

$$P(\bigcup_i A_i) = \sum_i P(A_i)$$

* Once you have Ω you can construct F and then construct P

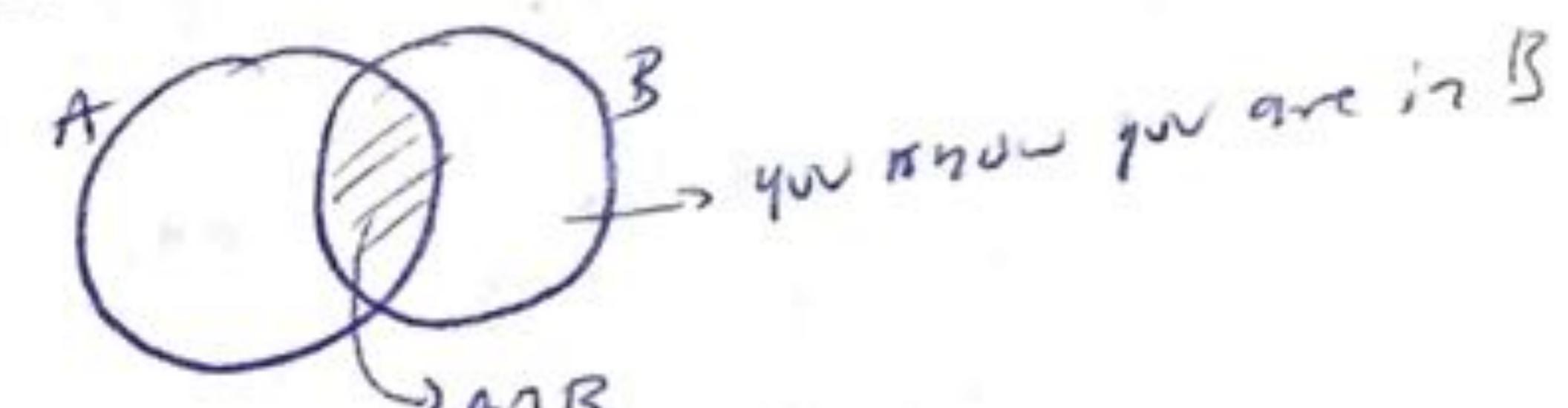
(Ω, F, P) : probability triple

Def: Conditional probability: $P(A|B) = \frac{P(A \cap B)}{P(B)}$

Law of Total Probability: $\forall A, B_1, \dots, B_n \in F \quad \bigcup B_i = \Omega$

$$P(A) = \sum_{i=1}^n P(A|B_i) \cdot P(B_i)$$

, given that A, B : observables $\in F$



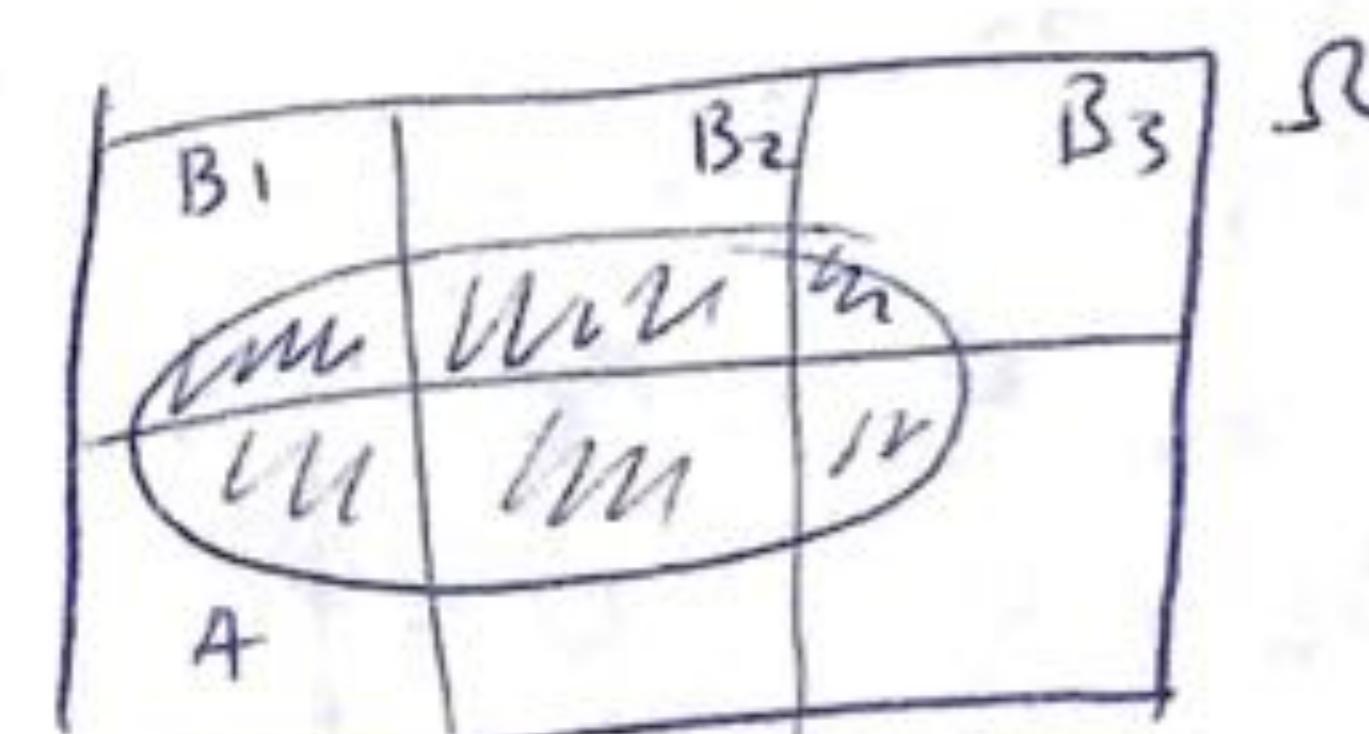
$$B_i \cap B_j = \emptyset \quad i \neq j$$

$$= \sum_{i=1}^n \frac{P(A \cap B_i)}{P(B_i)} P(B_i)$$

$$= \sum_{i=1}^n P(A \cap B_i)$$

$$(A \cap B_i) \cap (A \cap B_j) = \emptyset$$

$$= P(A)$$



Def: A RV is a function \bar{x} from Ω to \mathbb{R} $\bar{x}: \Omega \rightarrow \mathbb{R}$

2) $\{w: \bar{x} \leq x\} \in F$ (the little x it can observe), for all $x \in \mathbb{R}$

(Ω, F, P) $F_{\bar{x}}(x) = P(\bar{x} \leq x)$ CDF

defined for \bar{x} r.v.d.

$F_{\bar{x}_1, \dots, \bar{x}_n}(x_1, x_2, \dots, x_n) := P(\bar{x}_1 \leq x_1, \dots, \bar{x}_n \leq x_n)$

Let X_1, \dots, X_n of RV's Joint CDF $F_{X_1, \dots, X_n}(x_1, x_2, \dots, x_n)$

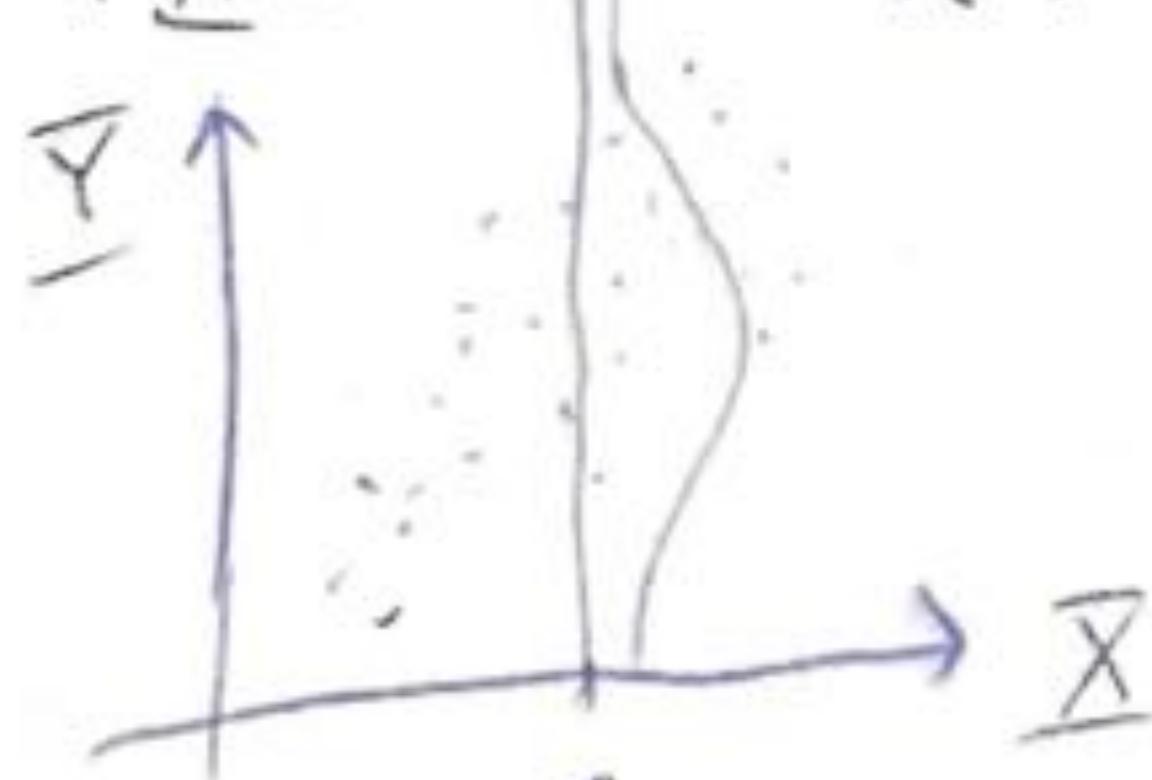
we say X and Y are independent if $F_{X,Y}(x,y) = F_X(x) \cdot F_Y(y)$

we can now write the conditional density (continuous RVs): $f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)}$

$f_{\bar{x}}(x) > 0$. $r(x) = E[Y|X=x] = \int y f_{Y|X}(y|x) dy$ ~ Regression function
prediction is always conditional.

$r(x) = "E[Y|X]"$ you calculate this for every x and then plug it in $r(x)$
over I will make.

$$E[(Y - r(\bar{x}))^2]$$



nb 14/11/2024

Markov Chains:

- We have some M: (with states $\{1, 2, 3, 4\}$)

$$P = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{3} & 0 & 0 & \frac{2}{3} \\ \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ \frac{1}{5} & \frac{2}{5} & \frac{2}{5} & 0 \end{bmatrix}$$

initial state $x_0=4$

- check if the rows sum to one
 a) Prob to return to x_0 in 2 steps?
 $P(X_2=4 | X_0=4) = \sum_{i=1}^4 P(X_2=4 | X_0=4, X_1=i) \cdot P(X_1=i | X_0=4)$
 $= \sum_{i=1}^4 P(X_2=4 | X_1=i) \cdot P(X_1=i | X_0=4) = 0 + \frac{2}{3} \cdot \frac{2}{5} + \frac{2}{5} \cdot \frac{1}{3} + 0 = \frac{3}{5}$

- Find stationary distribution:

$$\pi = P\pi$$

$$\pi_1 = \frac{1}{2}\pi_2 + \frac{1}{2}\pi_3$$

$$\pi_2 = \frac{1}{3}\pi_1 + \frac{2}{3}\pi_4$$

$$\pi_3 = \frac{1}{3}\pi_1 + \frac{1}{3}\pi_2 + \frac{1}{3}\pi_4$$

$$\pi_4 = \frac{1}{5}\pi_1 + \frac{2}{5}\pi_2 + \frac{2}{5}\pi_3$$

$$\pi_1 = \frac{46}{206}, \pi_2 = \frac{60}{206}, \pi_3 = \frac{45}{206}, \pi_4 = \frac{55}{206}$$

}

② Coin Flip : You flip a coin until you get 3 tails in the row
Define a MC with the states:

S - init (before flipping)

T = one tail

H = a head (which means your streak of T is over)

TT = 2 tails in a row

TTT = 3 tails in a row

Compute transition matrix

$$P = \begin{bmatrix} S & H & T & TT & TTT \\ S & 0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ H & 0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ T & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ TT & 0 & \frac{1}{2} & 0 & 0 & \frac{1}{2} \\ TTT & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$E[S] =$ expected number of steps to get TTT from state S. Compute $E(H), E(S), E(T), E(TT), E(TTT)$

$E(TTT) = 0$: you are already in this state - you don't need to do anything else

$E(S) = E(H)$: if you get head is at ~~is~~ you start again.

Intro to Data Science - Pattern Recognition

1. Supervised Learning

(x, y) $y \in \{0, 1, \dots, 1_{c-1}\}$: classes
 comes from the generator comes from the supervisor

$P_{\bar{Y}|\bar{X}}(y|x)$ conditional distribution given x

Model space $M = \{g_a(x) : g_a(x) \in \{0, 1, \dots, 1_{c-1}\}\}$

0-1 loss:

$$L(z, u) = \begin{cases} 1, & \text{if } u \neq z \\ 0, & \text{if } u = z \end{cases} \quad R(z) = E[L(g_z(x), Y)] \sim \text{Risk}$$

* Goal is always minimise Risk $z^* = \arg \min_z R(z)$

Linear Classifier

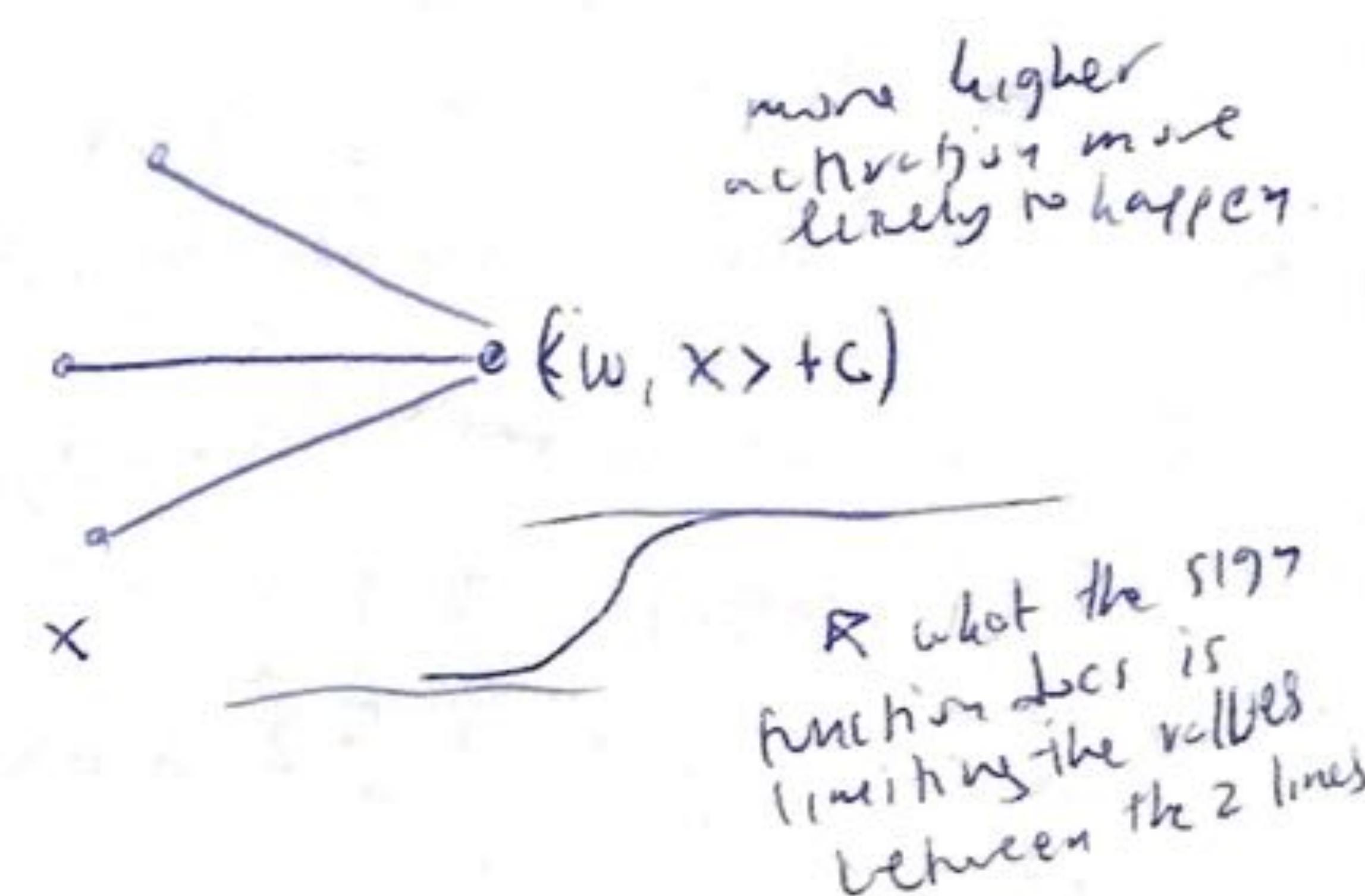
consider $x \in \mathbb{R}^d$, $y \in \{-1, 1\}$ and let $w \in \mathbb{R}^d$ and $c \in \mathbb{R}$

$$g_w(x) = \text{sign}(\langle w, x \rangle + c) \leftarrow (\text{McCullagh Pitts 4})$$

$$L(y, g_w(x)) = \text{sign}(\frac{\langle w, x \rangle + c}{2}) \cdot y + 1$$

$$y \neq g_w(x) \quad L(y, g_w(x)) = \frac{(-1) + 1}{2} = 1$$

$$y = g_w(x) \quad L(y, g_w(x)) = \frac{-(-1) + 1}{2} = 0$$

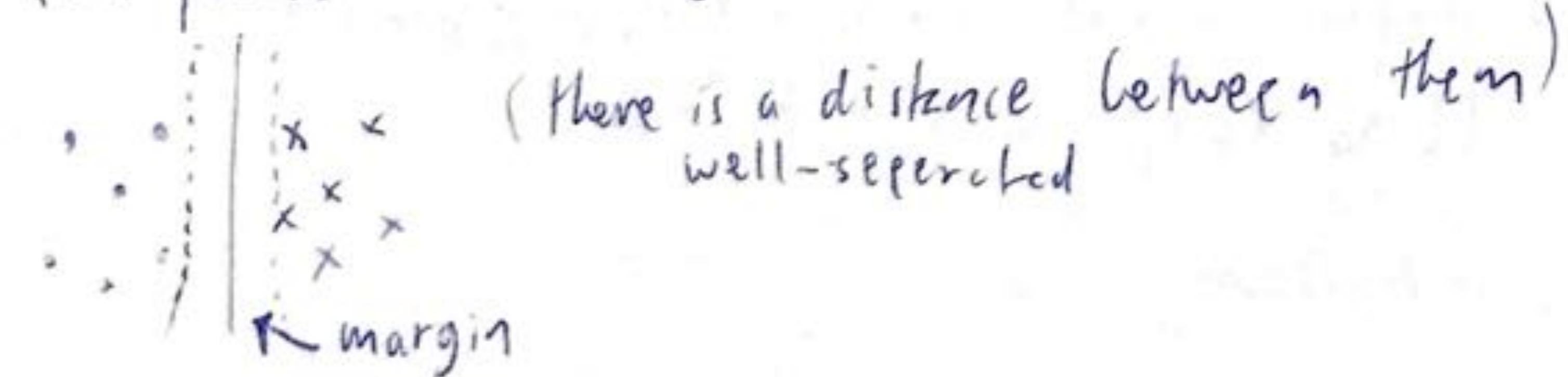


Dimension Trick:

Redefine $\tilde{w} = (w_1, \dots, w_d, c)$ $\tilde{x} = (x_1, \dots, x_d, 1)$ to remove the c (simplify the problem)

$$\langle x, w \rangle + c = \langle \tilde{x}, \tilde{w} \rangle$$

The problem to solve:



$\langle \tilde{w}, \tilde{x} \rangle = 0$: is a plane passing through the origin
 all vectors on this line are perpendicular to \tilde{w} \Rightarrow you get the zero vector

Perception Algorithm: (Rosenblatt 60)

① set $w = (0, 0, \dots, 0)$

② we have points (x_i, y_i) $i=1, 2, \dots, n$ if there is a point (x_i, y_i) such that $\langle w, x_i \rangle \cdot y_i < 0$ (classification was wrong) : $w \leftarrow w + x_i y_i \in \mathbb{R}^{d+1} \rightarrow \{-1, 1\}$

Repeat until all satisfy $\langle w, x_i \rangle \cdot y_i > 0$

Why this update rule?
we have $\langle w, x_i \rangle y_i < 0$, we want $\langle w, x_i \rangle y_i > 0$

This is a function of w , find which direction we need to change w to create the correct sign.

$$\nabla_w (\langle w, x_i \rangle y_i) =$$

$$\left(\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \dots, \frac{\partial f}{\partial w_{d+1}} \right) \rightarrow \text{direction of steepest (best) increase}$$

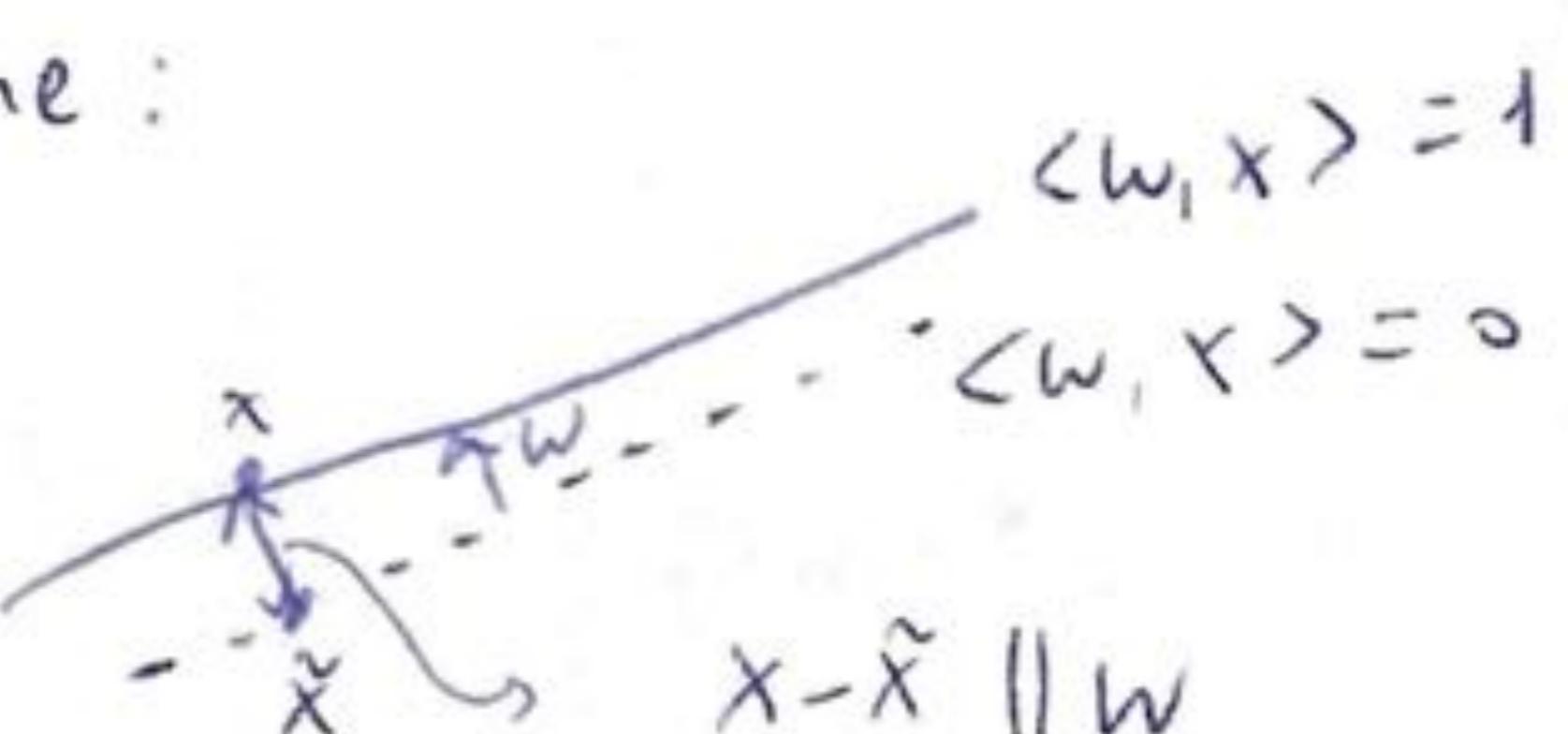
$$= \text{we have } \nabla_w (\langle w, x_i \rangle y_i) \rightarrow w_1(x_1)_i + w_2(x_2)_i + \dots + w_{d+1}(x_{d+1})_i \cdot d+1$$

If there are many solutions to a problem this algorithm will stop once it finds one solution.
(Hard Margin).

Support Vector Machine:

$$\text{Take } \|w\| = 1$$

$$\langle w, x \rangle = 1$$



$$|(x - \tilde{x})_w| = \underbrace{\|x - \tilde{x}\|}_{\perp} \underbrace{\|w\|}_{1} = \|x - \tilde{x}\|$$

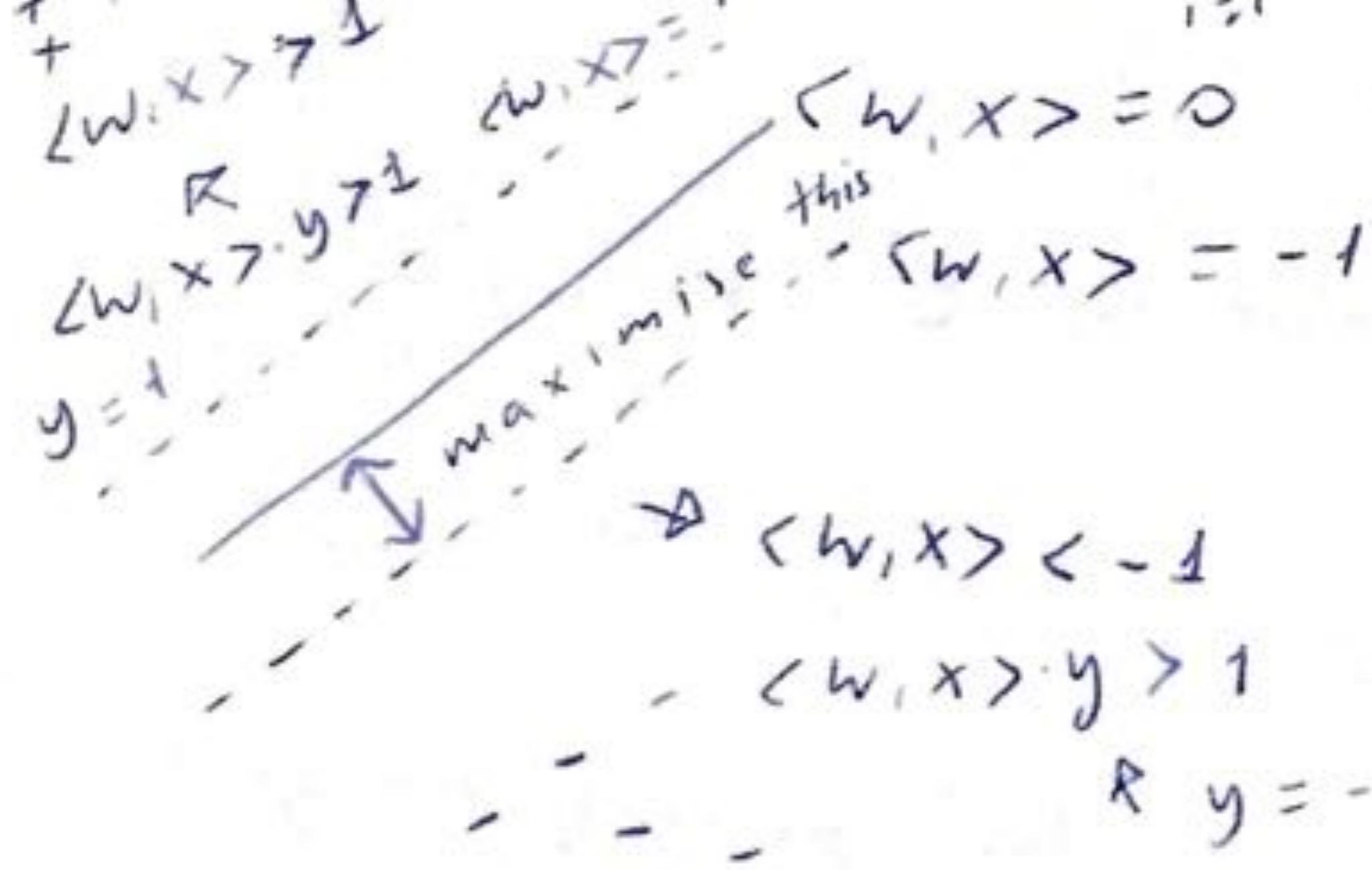
$$\text{if } \|w\| \neq 1 \text{ then } \|x - \tilde{x}\| = \frac{1}{\|w\|}$$

Consider the function: $\max\{0, 1 - \langle w, x \rangle \cdot y\}$ "Hinge loss"
minimise $\|w\|$ (you maximise the distance in other words)

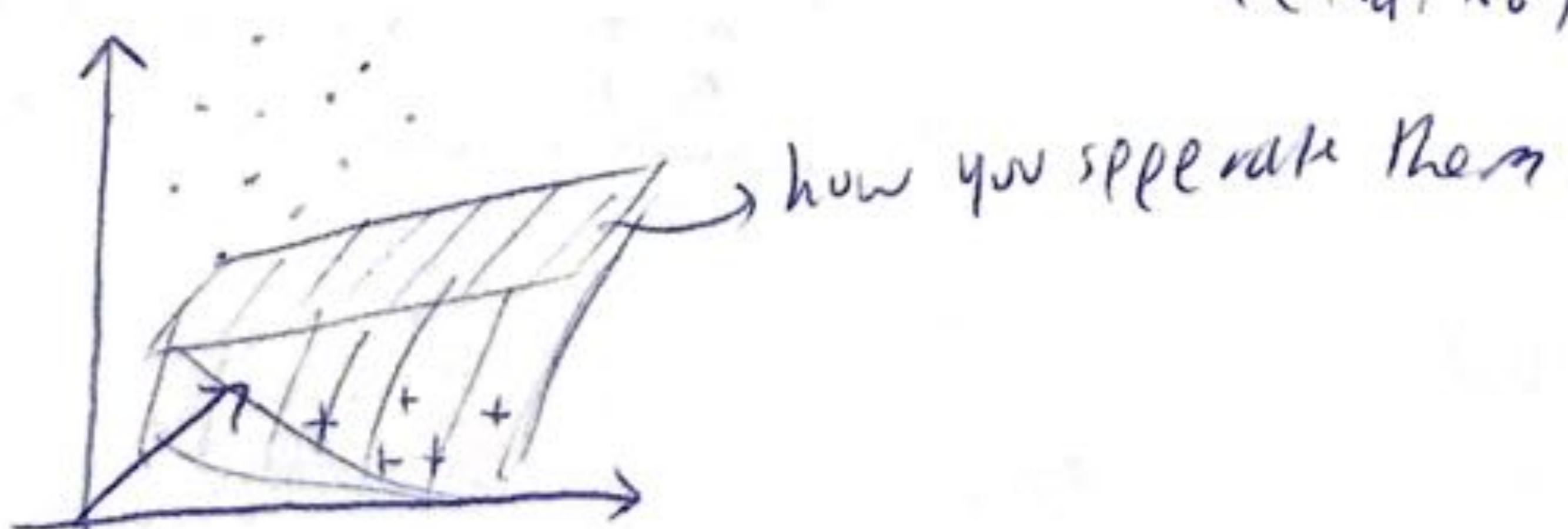
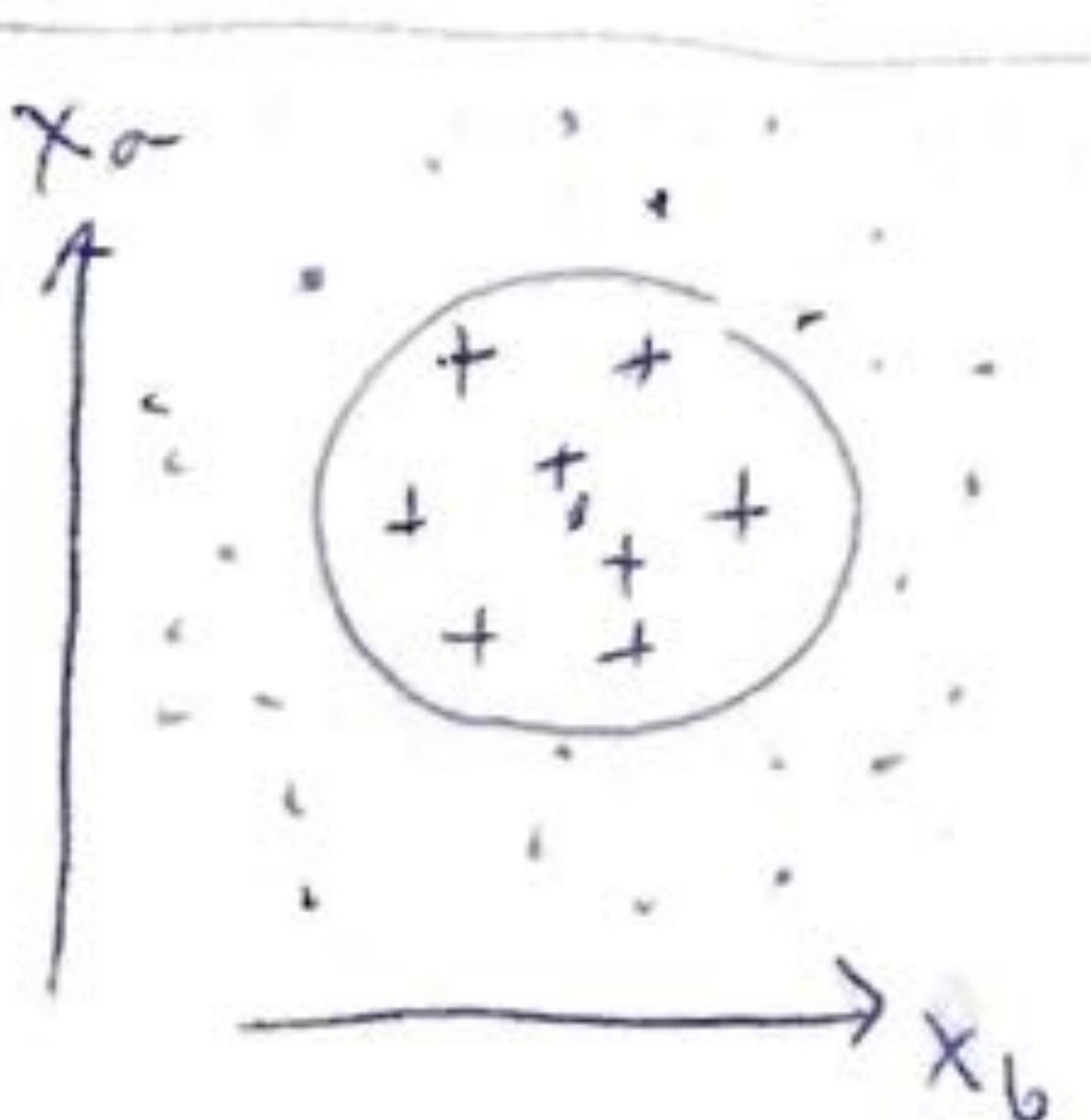
under the constraint $\langle w, x_i \rangle \cdot y_i \geq 1$ for all $i = 1, \dots, n$

Soft margin SVM constraint

$$\text{minimise } \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \max\{0, 1 - \langle w, x_i \rangle \cdot y_i\}$$



When data do not have a linear separator
 $\Phi(x_a, x_b) = (x_a, x_b, x_a^2 + x_b^2)$



The Kernel Trick: Perception update rule : initial $w = 0$ $+1 \quad -1$

$$w \leftarrow w + x_i y_i \quad \text{we will find } w \text{ of the form: } w = \sum_{i=1}^n c_i x_i$$

$$\langle w, x_k \rangle = \sum_{i=1}^n c_i \underbrace{\langle x_i, x_k \rangle}_{K_{ij}} \quad \text{Gram Matrix}$$

we have a ϕ

$$\phi(x_i)$$

$$\langle \phi(x_i), \phi(x_k) \rangle = \sum_{i=1}^n c_i \underbrace{\langle \phi(x_i), \phi(x_k) \rangle}_{K_{ik}}$$

If we have function $k(a, b)$

such that $K_{ik} = k(x_i, x_k)$ ~ we just need K and we don't need the ϕ anymore

Def: A kernel is a function $k(a, b)$ ^{extended} ^{symmetric}, if there is a function ϕ such that

$$\phi: \mathbb{R}^k \rightarrow \mathbb{R}^d : k(a, b) = \langle \phi(a), \phi(b) \rangle, a, b \in \mathbb{R}^k$$

Ex. 1) $k(a, b) = \langle a, b \rangle$ linear

2) $k(a, b) = (\gamma \langle a, b \rangle + m)^r$ polynomial

3) $k(a, b) = e^{-\|a-b\|^2}$ radial basis function (RBF)

4) $k(a, b) = \tanh(\gamma \langle a, b \rangle + m)$: sigmoidal \rightarrow infinite dimension \Rightarrow very complicated functions

RBF:

$$\langle w, x \rangle = 1$$

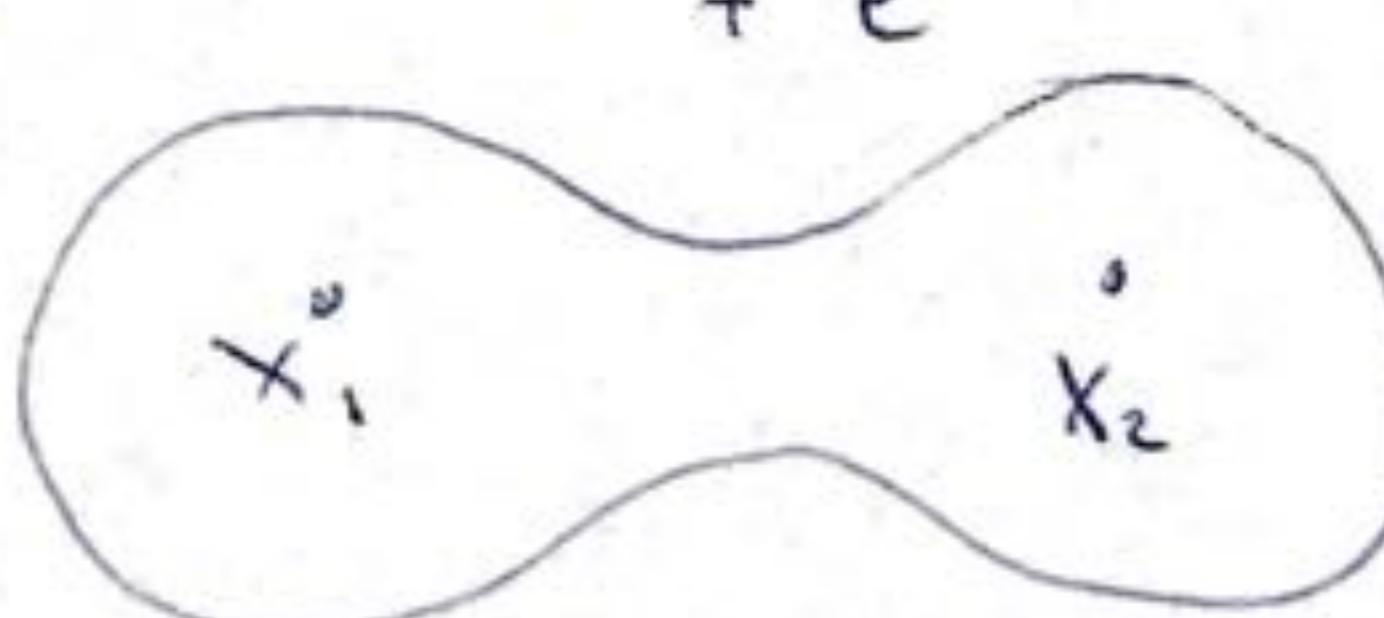
$1 = \sum_{i=1}^n c_i K(x_i, x) \sim$ find K such that this is true

1-point:

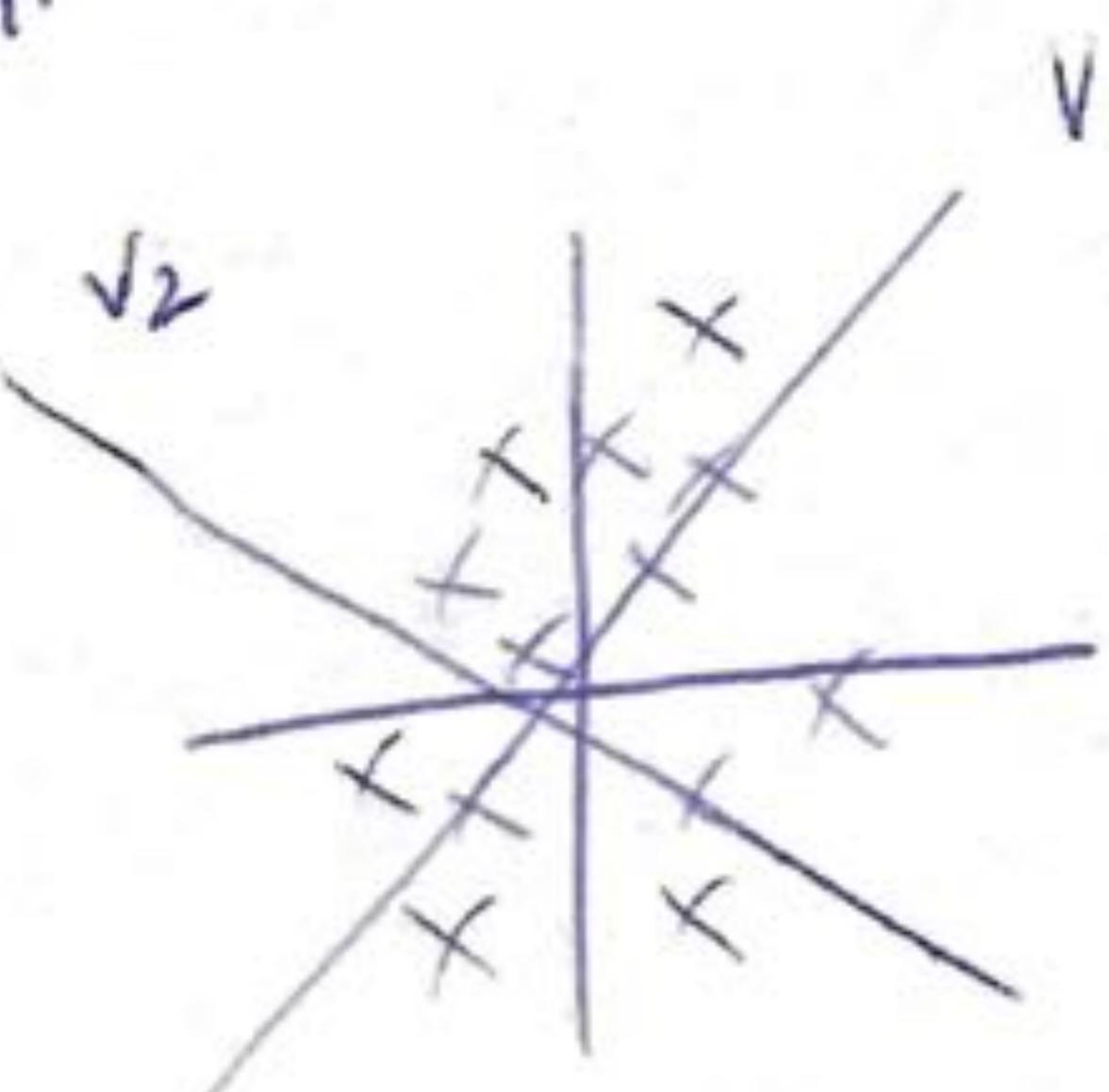


2-points:

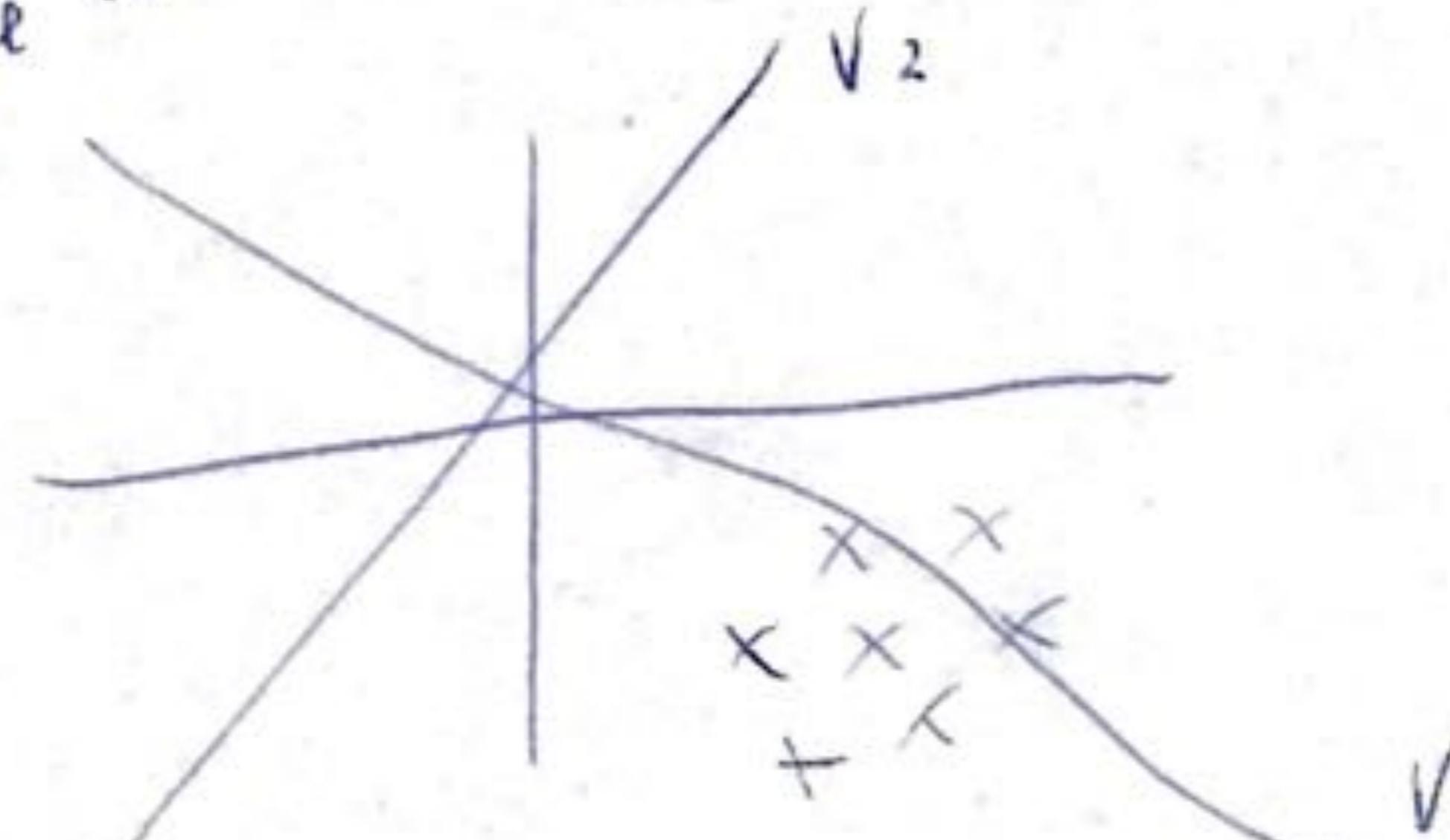
$$e^{-\|x_1 - x\|^2} + e^{-\|x_2 - x\|^2} = \frac{1}{2}$$



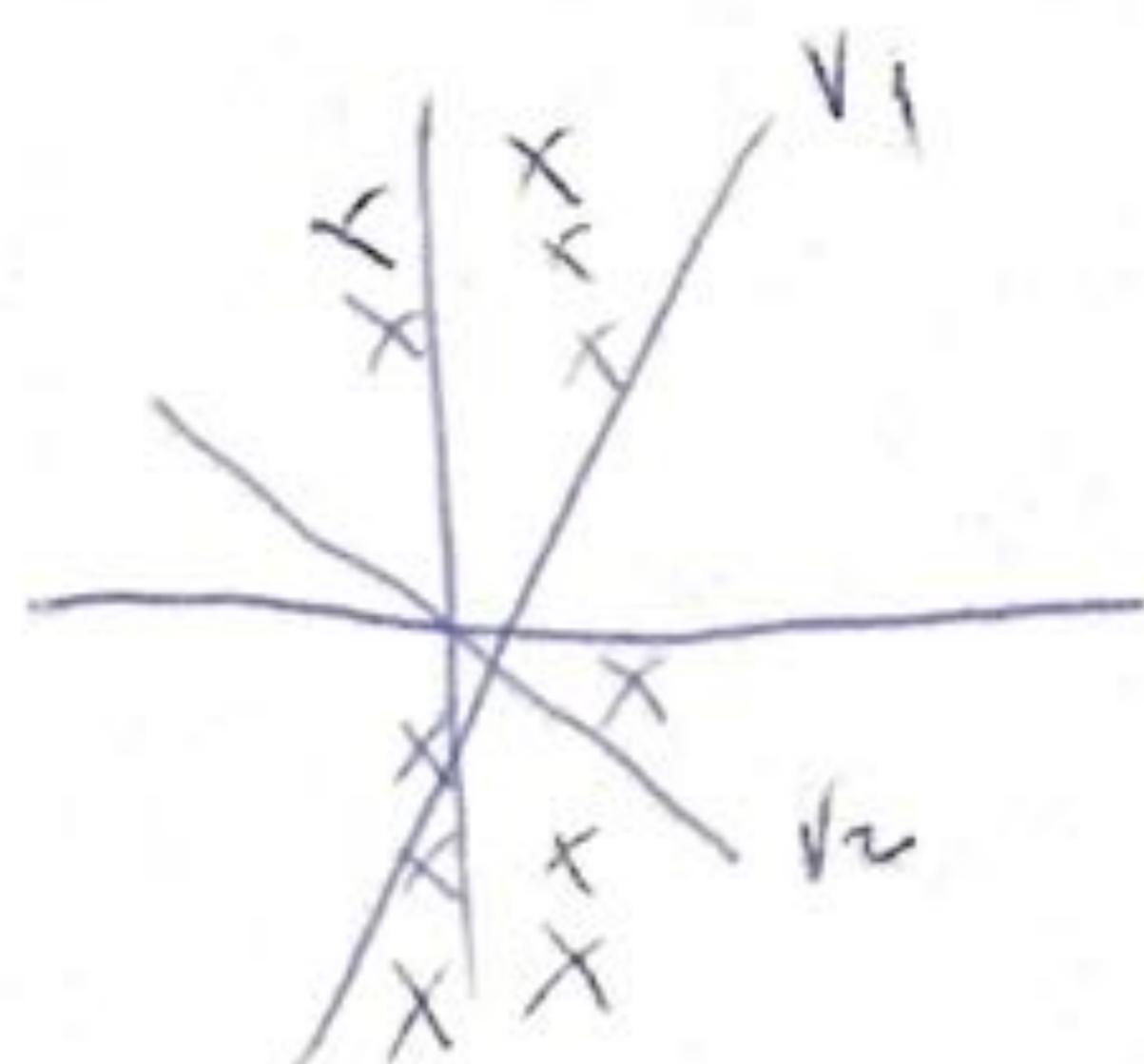
Q1/2



SVD
PCA "change of basis"
whitening in sklearn

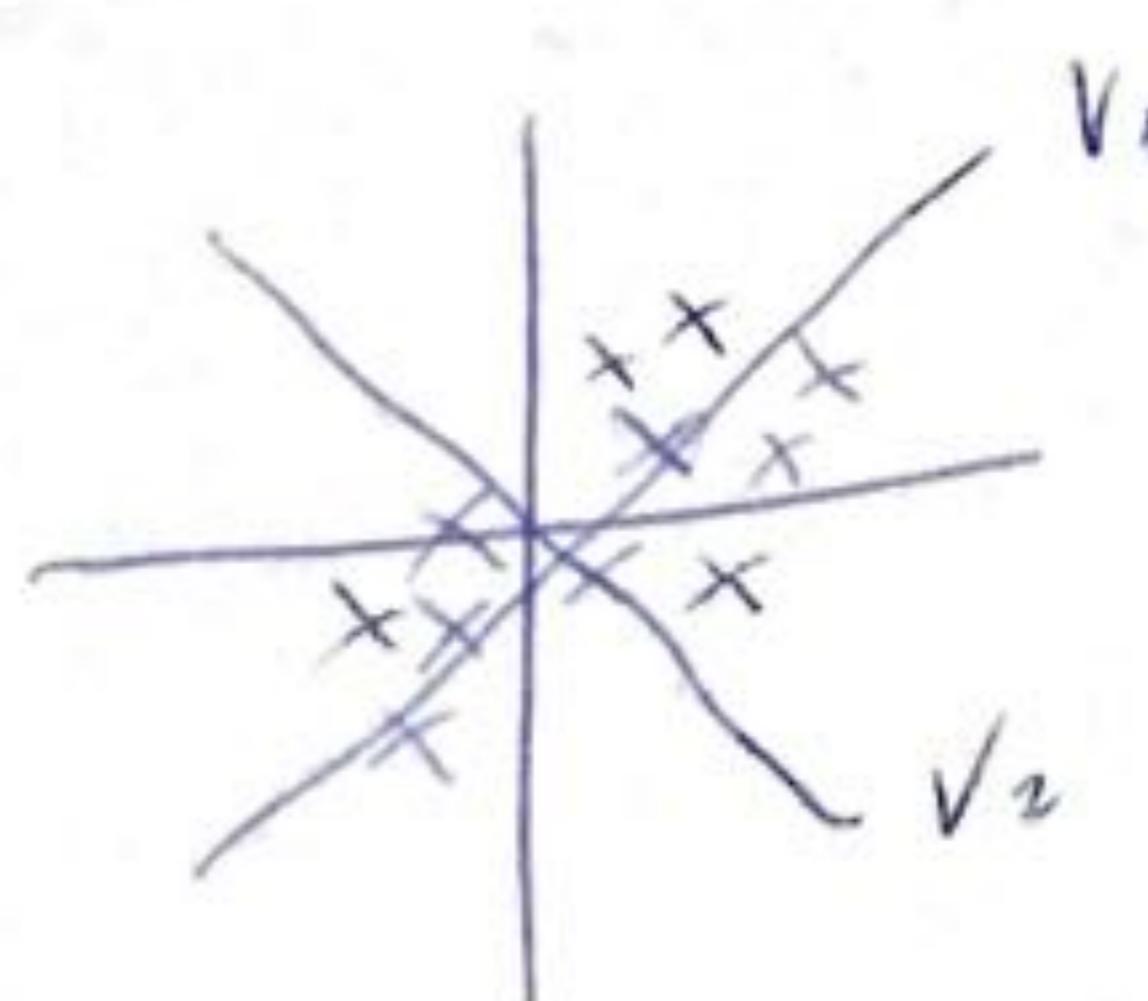


Centre your data!



standardised

if you standardise your data \Rightarrow

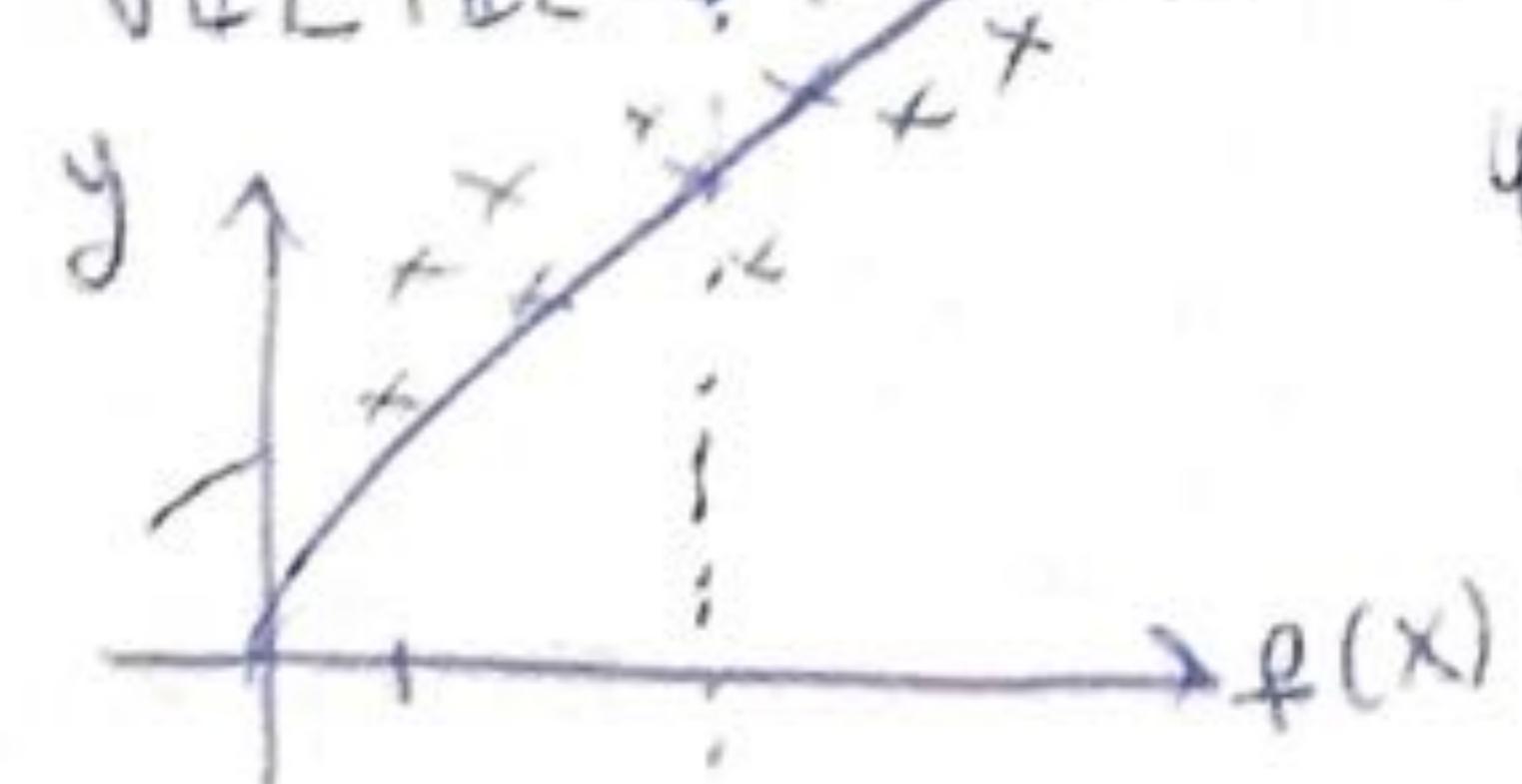


* does the data exploration in the train set only.

(calibration) error:

$$\sqrt{\mathbb{E}[|\mathbb{E}[Y|f(x)] - f(x)|^2]} =: C(f)$$

you want it to be 45° $y=x$ otherwise you are overpredicting or underpredicting



$y \in \{0, 1\}$

$$E[|P(Y=1|f(x)) - f(x)|^2] \rightarrow \text{usually in test set} \rightarrow \text{calibrated set}$$

predicted probability

Feature engineering

- * Construct features from data such that a "model" can understand it.
- * Construct more features to make a model more powerful.

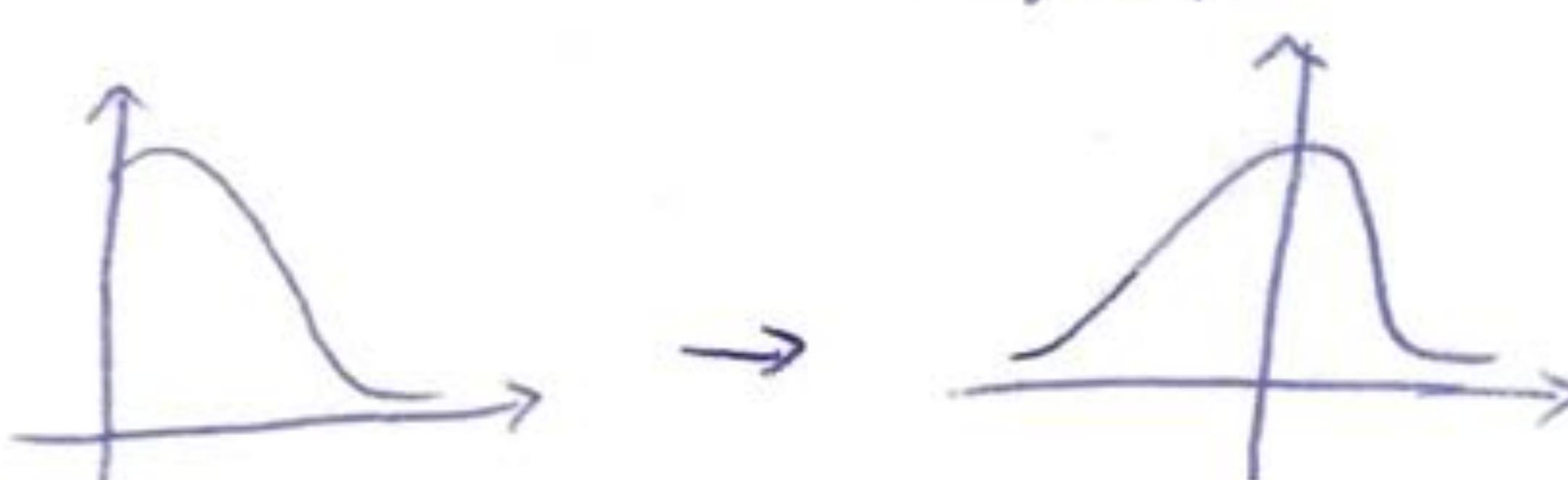
Ex: $f_{a_0, a_1}(x) = a_0 + a_1 \cdot x$

$$z = x^2$$

$$f_{a_0, a_1, a_2}(x, z) = a_0 + a_1 \cdot x + a_2 \cdot x^2$$

* Transformations of features

Ex: replace x with $\log(x)$



* Transformation of the target:

Ex: y is price

$$y' = \log(y)$$

$$\text{log scale } Y_1' - Y_2' = \log\left(\frac{Y_1}{Y_2}\right) : \text{proportional error}$$

Test data:

Dictionary: an enumeration of all possible words in your data

Bag of words model: ignore order/content \Rightarrow how many times does word x appears?

Term Frequency - Inverse Document Frequency (Tfidf Vectorizer in python)

\hookrightarrow is better than count: $f_{t,d} \leftarrow \text{document}$ $\hookrightarrow \text{term}$ Count of how many times t appears in d .

$$TF(t, d) = \frac{f_{t,d}}{\sum_{t \in d} f_{t,d}} \leftarrow \text{number of words in } d.$$

$$IDF(t) = \log\left(\frac{N}{|\{d : t \in d\}|}\right) \quad N: \text{number of docs}$$

$$TFIDF(t, d) = TF(t, d) \cdot IDF(t). \quad \hookrightarrow \text{in how many document does this word appear?}$$