

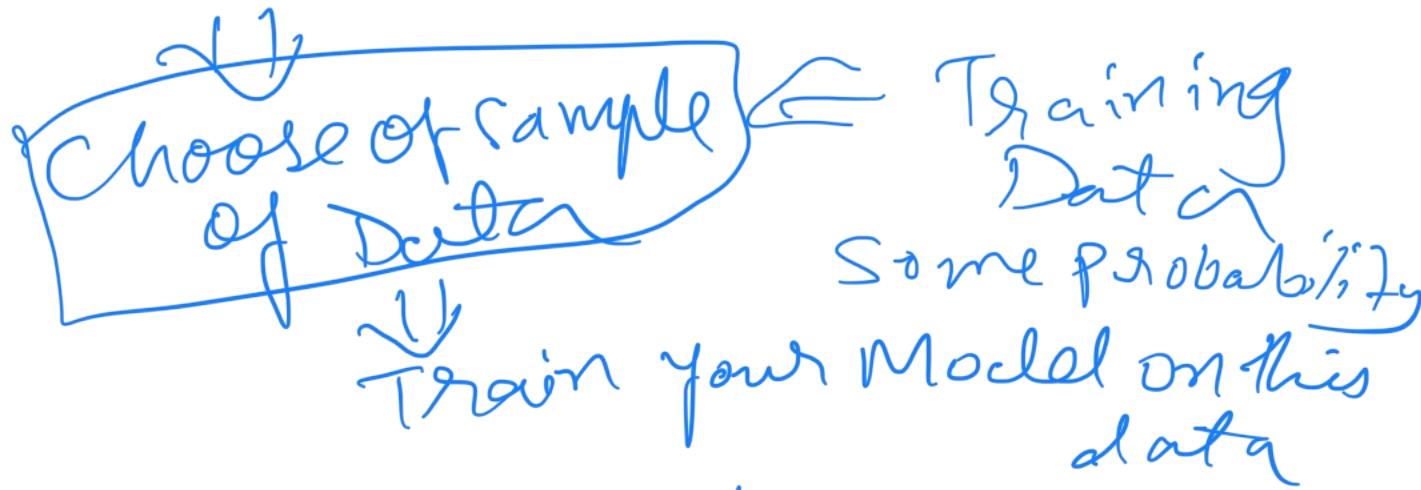
Why Probability in ML ?

→ Bayesian setup

- Entire Premise of ML is based on the notion of data coming from some distribution
- Classifiers Ptc. which Bay's rule

Infinite stream of data

$x_1 \ x_2 \ \dots \ x \ \dots$



Apply Model to make predictions
about unseen data
test data

Most material in these slides/notes is taken with permission from Prof.
Mukesh Tiwari , DA-IICT

What?

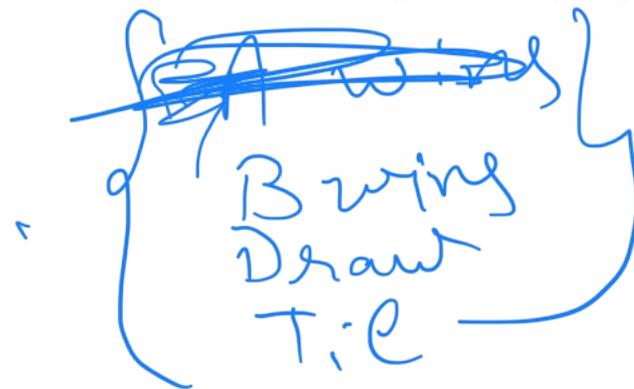
Quantifies the likelihood of the occurrence of an event.

Why?

- Model real world uncertainty in terms of simplified random processes.
- Make sense of data from real world: Family of models for many complex systems.

Sample Space and Events

- **Sample space (Ω):** set of all possible outcomes of a random experiment.
 - flipping of a coin $\Omega = \{H, T\}$
 - rolling of a dice $\Omega = \{1, 2, 3, 4, 5, 6\}$
- **Sample point ω :** Each element of Ω , e.g. $\Omega = \{H, T\}$, sample point H or T (one realization of the random experiment)
- **Event:** any subset E of Ω , e.g. $\Omega = \{H, T\}$, event: $\{H\}, \{T\}, \{H, T\}, \emptyset$
- **Probability:** For each subset $E \in \Omega$, there is a probability $P(E)$ (probability is defined on events).



Axioms of Probability

Properties of probability measure:

- $P(\Omega) = 1$
- $0 \leq P(E) \leq 1$.
- For any sequence of mutually exclusive events A, B ($A \cap B = \emptyset$)

$$A \cap B = \emptyset$$

Negation | Impossibility
 $P(A) = 1 \Rightarrow$ Certain

$$P(A \cup B) = P(A) + P(B)$$

\cup - oring
 \cap - anding

$$\begin{aligned} P(A \cap B) &= 0 \\ \xrightarrow{\text{A wins}} & \quad \xrightarrow{\text{B wins}} \\ P(A \cup B) &= P(A) + P(B) \end{aligned}$$

A large number of important mathematical results follow from these axioms:

- Since $E \cup E^c = \Omega$, and $E \cap E^c = \emptyset$, $P(E^c) = 1 - P(E)$

- $\underline{P(\emptyset) = 0}$

- $P(E \cup F) = \underline{P(E)} + \underline{P(F)} - \underline{P(E \cap F)}$

- If $E \subseteq F$, then $\underline{P(E)} \leq \underline{P(F)}$

$$P(E^c) = 1 - P(E)$$

$F \models P$ (Students > 50 marks)
 $E \models P$ (Student > 60 marks)

If it rains
 $P(\text{It rains}) = 0.7$

If does not rain
 $P(\text{It does not rain}) = 0.3$

Conditional Probability

$E = \text{Roads are wet}$

$F = \text{There is Dew}$

$$P(E|F) = \frac{P(E \cap F)}{P(F)}$$
$$P(E \cap F) = P(E|F)P(F)$$

F has already occurred

Law of total probability

F_1 Either you have B-Tech with S.J.Y.A
 F_2 BA in literature

Let $F_1 \dots F_n$ be disjoint events that form a partition of the sample space. Then for any event E

$$P(E) = \sum_{i=1}^n P(F_i)P(E|F_i)$$

(Total Probability Theorem)

$$P(E) = \sum_{i=1}^n P(F_i) P(E|F_i)$$

Baye's Formula

$$\underline{P(B|A)} = \frac{\underline{P(A|B)P(B)}}{\underline{P(A)}} = \frac{P(A|B)P(B)}{P(A|B)P(B) + P(A|B^c)P(B^c)}$$

B_1, B_2
 B, B^c

Notion of Posterior

θ - Model
 x - data

$$p(\theta|x) = \frac{p(x|\theta)}{p(x)} p(\theta)$$

Likelihood
Prior was used to generate data
P(Data | Model) . P(Model)

Posterior
overshading
MAP



Coin

whether it is a fair or
coin ~~you~~ not

$$P(H) = P(T) = 0.5$$

$$\theta = P(H) \quad P(H) = 0.5$$

$$x = \boxed{HTHTHTHHHTTTTH}$$

$$P(\theta|x) \propto P(x|\theta) - P(\theta)$$

$P(H) = 0.65$ fair coin

$\theta \rightarrow$ classifier \rightarrow Cars

$\theta_1, \theta_2, \theta_3, \dots, \theta_n \rightarrow$ Trucks

$X = [(\text{image}_1, \text{car}), \text{image}_2, \dots]$

$$P(\theta = \theta_1 | X) = P(X | \theta = \theta_1) \cdot P(\theta_1)$$

↑ Likelihood ↑ Prior

$$f_\theta(I) = \text{car}$$

Posterior \propto Likelihood \times Prior

Prior

$$P(\theta) = 0.5$$

$$P_\pi(H) = 0.5$$

[H H T T T H H H H]

EM Algorithm

Gradient Descent

$$P(\theta|D)$$

0.6

$$P = \frac{P(0.5)}{P(0.5) + P(0.1)}$$

$$P = \frac{0.5}{0.5 + 0.1} = 0.83$$

$$P(H) = 0.83$$

$$P(H=0.65) \times$$

$$P(O_1)$$

$$P(O_2)$$

X X X

O O O

Random Variable

- It frequently occurs that in performing an experiment we are interested in some function of the outcome rather than the outcome itself. These real valued functions defined on the sample space are known as *random variables*.
 - Discrete random variable: takes either finite or a countably infinite number of values.
 - Continuous random variable: continuum of possible values.

Example

- Outcome of tossing two coins ($\Omega = \{T, T\}, \{T, H\}, \{H, T\}, \{H, H\}\}$).
 - Define X as the number of heads: $X = 0, 1, 2$
 - Probability $P\{X = x\}$: $P\{X = 0\} = 1/4$, $P\{X = 1\} = 1/2$, $P\{X = 2\} = 1/4$
- Probability of the events are denoted as $P(\cdot)$
- Probability of the random variables are denoted as $P\{\cdot\}$

Convention: X, Y, Z - random variables, x, y, z - real numbers.

$X = \# \text{ of Heads}$
 in coin toss
 X can takes $\{0, 1, 2\}$

Expectation of a random variable

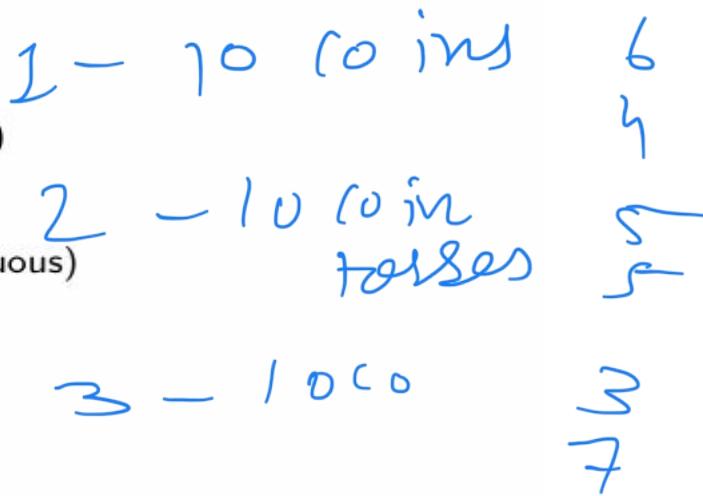
The expectation or the average value (mean, weighted average, center of mass) of a random variable is:

$$E[X] = \sum_x x p_X(x) \quad (\text{discrete})$$

$$E[X] = \int_{-\infty}^{\infty} dx \times f_X(x) \quad (\text{continuous})$$

Variance

$$\text{Var}[X] = \sigma_X^2 = E[(X - E[X])^2]$$



- Variance calculates how spread out a distribution is
- It puts special emphasis on the outliers
- standard deviation $\sigma = \sqrt{\text{Var}[X]}$
- Other (relevant) statistical quantities: Mode (most likely value), Median (middle of the distribution)

$$E[X] = \sum_{x_i} x_i p(x_i)$$

$$\begin{aligned} & \frac{1}{3}((4-5)^2 + (3.5-5)^2 + (6-5)^2) \\ &= 1 \cdot \left(\frac{1}{6}\right) + 2\left(\frac{1}{6}\right) + \dots + 6\left(\frac{1}{6}\right) \\ &= 4 \end{aligned}$$

$X = 0$

1

2

3

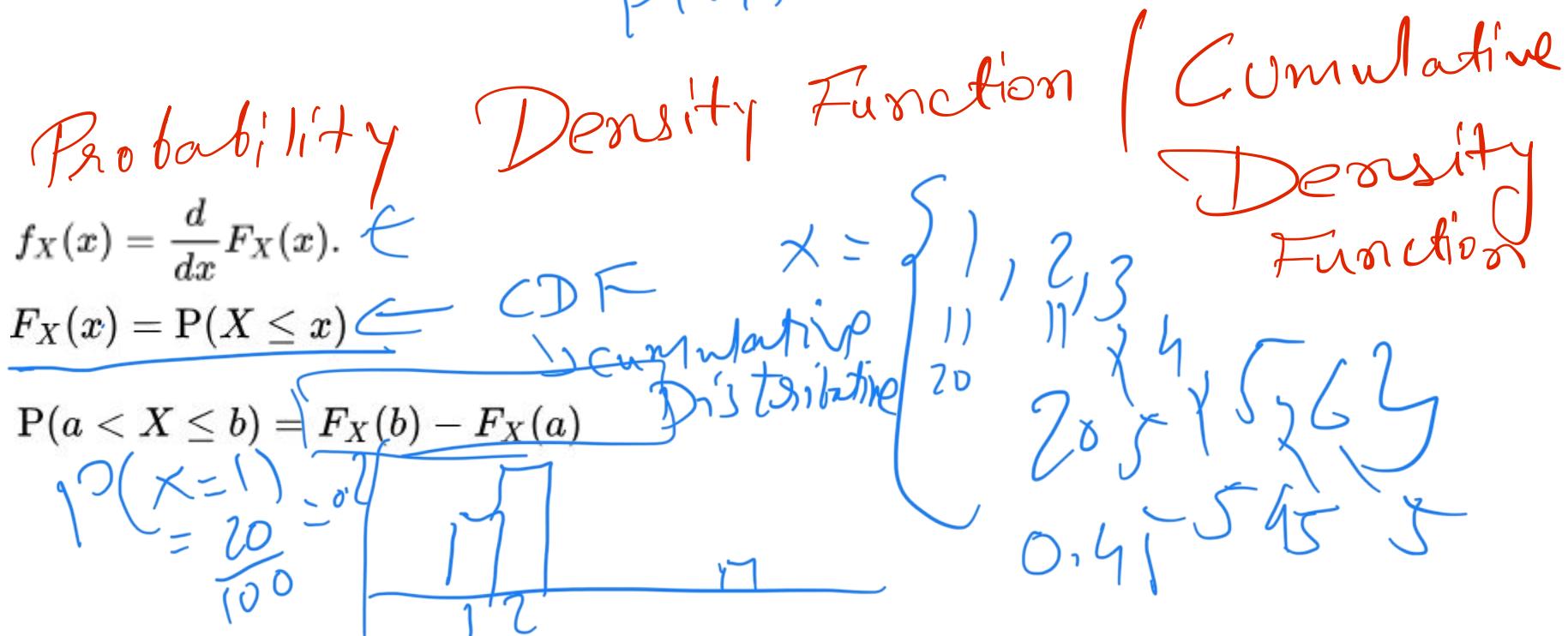


$t \rightarrow$
~~PDS~~ = Likelihood $\prod_{i=1}^n p_{\text{obs}}(x_i | \theta)$
~~PT~~ Prior(t) = $p(\theta|t)$



$P(X = 10.009)$, $P(X < 11)$
 $P(X = 10.010)$

PMF
 $p_X(x) = P(X = x)$ (Probability Mass Function)
 Discrete Random Variables
 PMF



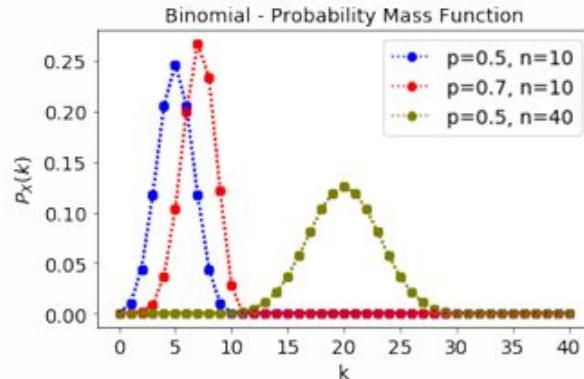
Bernoulli Distribution

Experiment where the outcome is either a *success* ($X = 1$) with probability p or a *failure* ($X = 0$) probability $1 - p$,

$$P\{X = 1\} = p, \quad P\{X = 0\} = 1 - p$$

Binomial Distributione $b(n, p)$

- X : number of successes in n (fixed) independent Bernoulli trials.
- $p_X(k) = \binom{n}{k} p^k (1-p)^{n-k}$, $k = 0, 1, 2, \dots, n$



Uniform Distribution

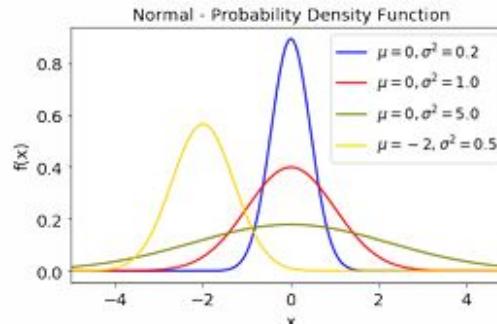
A random variable is said to be uniformly distributed over the interval (α, β) if its pdf is given by

$$f_X(x) = \begin{cases} \frac{1}{\beta - \alpha}, & \alpha < x < \beta \\ 0, & \text{otherwise} \end{cases}$$

Normal Distribution

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad -\infty < x < \infty$$

$$X \sim \mathcal{N}(\mu, \sigma^2)$$



Jointly distributed random variables

If we have to deal with two or more random variables (let us say X and Y)

- Joint **cumulative** distribution of X and Y

$$F_{X,Y}(x,y) = P\{X \leq x, Y \leq y\} = \int_{-\infty}^x \int_{-\infty}^y f_{X,Y}(x,y) dx dy$$

- Distribution (**marginal** distribution) of X ($F_X(x)$)

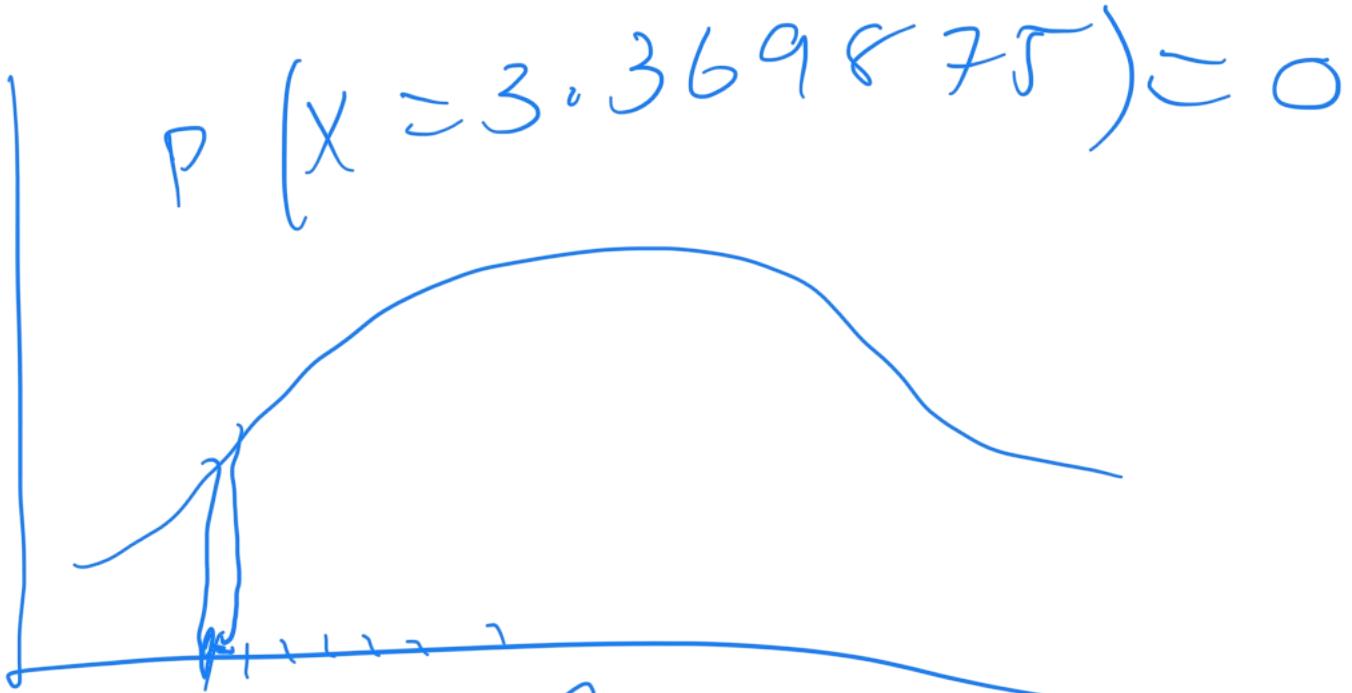
$$F_X(a) = P\{X \leq a\} = P\{X \leq a, Y < \infty\} = F(a, \infty)$$

- Joint probability **mass** function (convenient to have for discrete distribution)

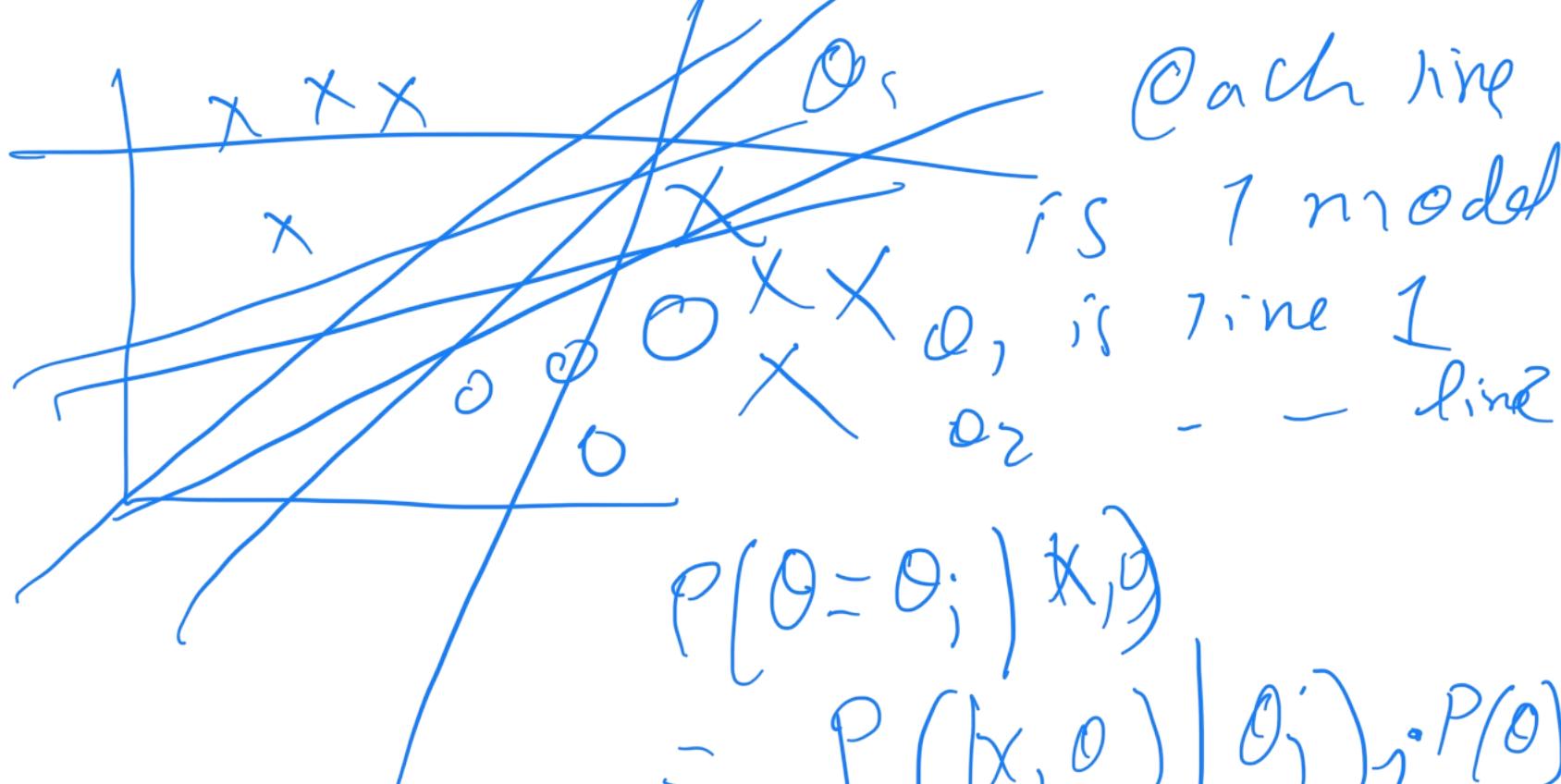
$$p(x,y) = P\{X = x, Y = y\}$$

- **Marginal PMFs:** Probability mass function of X ($p_X(x)$) or Y ($p_Y(y)$)

$$p_X(x) = \sum_{y: p(x,y) > 0} p(x,y)$$



$$(x, x + \delta x) \quad P(X \leq 0.3)$$



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

K L Divergence

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right).$$

