# Theory Recap:

**① Conditional Probability + Bayes (Short)**

- A is an event, B is another event

① $P(A), P(B)$.     $P(A|B) = \dfrac{P(A \cap B)}{P(B)}$

$\dfrac{1/6}{1/2} = 1/3$

② $P(A|B) P(B) = P(A \cap B) = P(B|A) P(A)$

(Never bother to learn Bayes Rule!)

$P(A|B) = \dfrac{P(B|A) P(A)}{P(B)}$ } Bayes Rule

③
- why? Interesting? <u>Data Science Interview</u> (3 simple steps to there)
- ① Define A and B. A is cloudy, B rain.
- ② Write down exactly what you have $P(A), P(B), P(A|B)$ }
- ③ Plug it all into Bayes Rule. Simple but really fiddly!

~ pause ~

---

|   | α | R | 1/6 |
|---|---|---|---|
| A |   | /////// |  |
| B |   |  |  |
| C |   |  |  |

} 1/3

- Rare events question
  ↓
  Unintuitive answer!

---

**② Random Variables and Distributions** → "take different values due to randomness"

① • What is a random variable? An object or quantity depending on randomness.
- ① X number when you roll a dice. Can be $\{1, 2, \ldots, 6\}$
② ② X the height of an individual picked at random
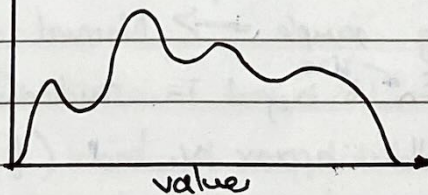- Random variables can be <u>discrete</u> or <u>continuous</u> → Someone to define each
  Difference between them? ⊗ a realisation of X is either from a
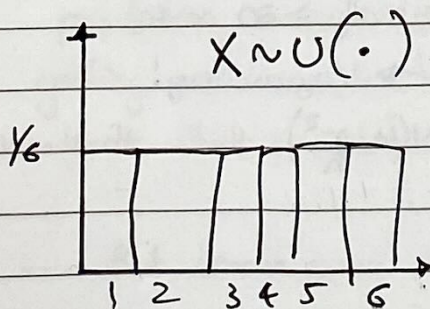  predefined lot or not

  $X, \; x = 3$

  ③ → $P(X = x)$ Notation. $\sum_{i=1}^{n} P(X = x_i) = 1$
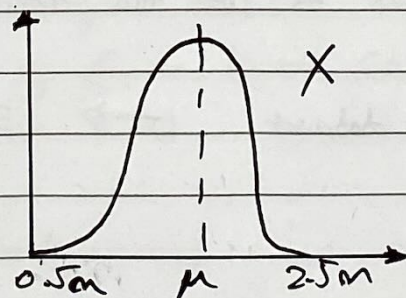  notation

prob
density

④
  X follows distribution →

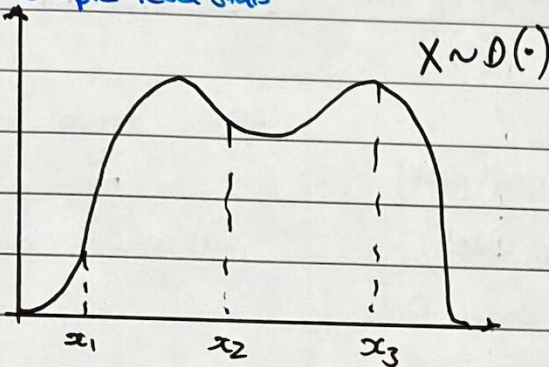<span style="writing-mode: vertical">Distributions</span>

$X \sim U(\cdot)$

value

1 2 3 4 5 6

0.5m   μ   2.5m

X

$x = 1, 2 \ldots$

◉ Distributions about Random Variables have properties.

- $E(X) = \mu$, $\text{Var}(X) = \sigma^2$  } Population Moments, Population mean

  {Never memorise!} $= E[(X-\mu)^2] = E(X^2) - E(X)^2$    Population sample

<u>Samples</u> are drawn from distribution.

↳ Ogive! Sample level stats



$X \sim D(\cdot)$

$x_1, x_2, x_3$

$(x_i \quad i = 1...n)$

$x_1, x_2, ..., x_n \quad \longrightarrow \quad \underbrace{X \sim D(\circ)}_{?}$

{ Statistics is about sample to distribution!

⟹ Claim about the world

<u>Problems</u>  ① Don't know the distribution

② Most distributions are messy so <s>difficult</s> to do inference

↳ Hard to do sample → population!

③ <u>Central Limit Theorem</u>:

⊛ Arguably the most beautiful thing in mathematics/statistics... ⊛

⚠ — <u>whatever</u> the underlying distribution, you can still do inference because the sample mean will be normally distributed as $n$ gets large

— Gateway to all statistics

"Approximately distributed"

- Sample: $x_1, x_2, ..., x_{100}$  { $\bar{x} \overset{a}{\sim} N(\mu, \frac{\sigma^2}{n})$ } Explain this visually

- Lots and Lots of sample ⟶ Normal dist

- Sample mean $\bar{x}_n$ ← size is in fact a random variable with its own distribution

- If $n$ is small then approx v. bad (if $n$ is 1 then just the underlying dist)

- $n$ must be "large" to use the approx (typically $> 30$ or $40$)

↳ Very arbitrary!

$z_1, z_2, ..., z_n$ { Your dataset  ⟶  $\bar{x}_n \overset{a}{\sim} N(\mu, \frac{\sigma^2}{n})$ think it should be

<span style="color:red">actually collected</span>

## One example

$$\bar{x}_n \overset{a}{\sim} N\left(\mu, \frac{\sigma^2}{n}\right)$$



Theoretical Mean
of pop.  10

← Distribution of $\bar{x}_n$

$5 \qquad \mu=10$

- $x_1, x_2, ..., x_{1000} \rightarrow \bar{x}_{1000}$.  5 . 5 is different to 10
- Can we say our sample is different?
- CLT allows us to work out the probability of us getting $\boxed{\bar{x}=5}$ given we think $\mu=10$. If its really low ⇒ our sample different...
  Hypothesis test which you'll see next time

Build the intuition for hypothesis testing! I think the population mean is 10, I collect loads of data (big sample) and the sample mean is 5. This seems lower than 10 so I suspect the original hypothesis that the population mean is 10 might be wrong. But how can I conclude this (inferential statistics...)

... well given the Central Limit Theorem I know the distribution of $\bar{x}$, the sample mean, and I know it is distributed around the population mean. I can therefore calculate the probability of getting a sample mean of 5 or lower, GIVEN, my hypothesis that the population mean is 10

... But more on this next week!