Why probability?

Phenomenon
↓
mathematical model

Deterministic → Probabilistic
- parameters DETERMINE
- fixed outcome

- parameters DESCRIBE
- output is variable

## Definitions

Sample space: set of all possible outcomes $\Omega$

EVENT: subset of $\Omega$, $A \in \Omega$

PROBABILITY $P$: is a function that assigns likelihood to EVENTS in $\Omega$ occurring
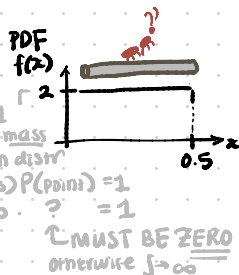
must satisfy the axioms

(1) $\sum_{A_i \in \Omega} P(A_i) = 1$ — something must happen

(2) $P(\phi) = 0$ — nothing never happens

(3) $P(A) + P(\bar{A}) = 1$ — EITHER Raining or not Raining

(4) $P(A_i) \in [0,1]$ — mass is positive

finite $\Omega \to$ mass
infinite $\Omega \to$ density

(5) If $A_1 \cdots A_n$ disjoint,

$A_i \cap A_j = \phi \quad \forall (i,j)$

$\sum P(\cup A_i) = \sum P(A_i)$

$\cup$: union (in either A or B)
$\cap$: intersection (in Both A and B)
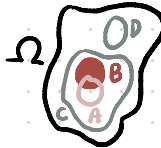
(6) $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ — union

A ∪ B

PDF $f(x)$

≥1
- density, not mass
- for uniform distr'n
  (# of points) P(point) = 1
  ∞ · ? = 1
  └ MUST BE ZERO
  otherwise $\int \to \infty$

2

0.5

INDEPENDENCE: $P(A,B) = P(A)P(B)$ iff A & B are INDEPENDENT $A \perp\!\!\!\perp B$

Both events happening

arises from assumption or by construction
↓ BEB1103          ↓
coin flips      6-dice: A: {2,4,6}, B: {1,2,3,4}

CONDITIONAL PROBABILITY: $P(A,B) = P(A|B)P(B)$

SUPPOSE B has OCCURRED.
This information changes probabilities of other events

$P(D|B) = 0$   (D, B disjoint: $D \cap B = \phi$)
$P(C|B) = 1$   ($C \supset B$)
$P(A|B) = ???$  (not disjoint: $A \cap B \neq \phi$)

└ given B has happened, what is prob. of A happening?
when do B & A OVERLAP?

$$P(A|B) = \frac{P(A,B)}{P(B)}$$

What is $P(A|B)$ when $A \perp\!\!\!\perp B$ (they are independent?)

$P(A|B)$ Shows up in

BAYES' LAW:
$P(A,B) = P(B,A)$
$P(A|B)P(B) = P(B|A)P(A)$

$$\boxed{P(A|B) = \frac{P(B|A)P(A)}{P(B)}}$$

- Bayes + billiards
- indep. thought of Laplace:
  IP of cause given an event is proportional to
  IP of event given the cause
- updates!

Marginalization:
1. partition $\Omega$ into $B_i$'s: $B_1 \cdots B_n$

$B_2 \quad B_3 \quad B_4$

find $P(A)$ in terms of $B_i$'s

$B_i$'s disjoint $\cup P(AB) = \sum P(AB)$

$$P(A) = \sum_i P(A|B_i) P(B_i)$$ → continuous analog

2.

marginal distr's

|     | | |
|-----|------|------|
| 0.1 | 0.2  | 0.3  |
| 0.3 | 0.4  | 0.7  |
| [0.4] | [0.6] | |

$P(\cdot | x) = \frac{P(\cap x)}{P(x)} = \frac{0.2}{0.6}$

$P(x,y)$

3. given joint distribution
$$P(x) = \int P(x|y) P(y) \, dy$$

$P(x,y)$

# Random Variable: RV $X: \Omega \to \mathbb{R}$ <span style="color:gray">a very unfortunate name...</span>

DETERMINISTIC function that **maps outcomes in sample space** to **real numbers**.

So if $\Omega$ has a probability distribution over it, this induces a **PROBABILITY distribution of RV on $\mathbb{R}$**

$$
\begin{array}{c|c|c}
\Omega & \mathbb{R} & P(\ ) \\
\text{coins} & H \to 1 & 0.4 \\
& T \to 0 & 0.6
\end{array}
\qquad \text{indicator R.V.}
$$

$$
\begin{array}{c|c|c}
\Omega & \mathbb{R} & P(\ ) \\
\{1,1\} \to \text{sum}(1,1) & & 1/36 \\
& & \\
\{6,6\} \to \text{sum}(6,6) & & 1/36
\end{array}
\qquad
\begin{array}{l}
S: \{i,j\} \\
\text{RV: } i+j
\end{array}
$$

DIST. of sum of 2 DICE

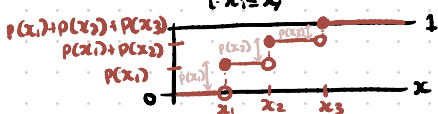**Random variables : MAPPING**

<span style="color:gray">Brownian motion</span>

for $s \in \Omega$, $x \in \mathbb{R}$ RV: $s \to x$, can be

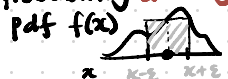**DISCRETE**: $x$ can take on finite # values

probability **mass** function: pmf: $p(x_i) = P(X = x_i)$

CDF: $F(x) = P(X \le x) = \sum_{i: x_i \le x} p(x_i)$



$$P(x_1) + P(x_2) + P(x_3)$$
$$P(x_1) + P(x_2)$$
$$P(x_1)$$

**continuous**: $x$ can take on infinite # values

probability **density** function

pdf $f(x)$



1. $f(x) \ge 0$
2. $\int f(x) = 1$
3. $P(a \le x \le b) = \int_a^b f(x) dx$

$P(x = a) = 0 \quad \forall a \in \mathbb{R}$

$T_1$ height $= (T_1)$
$T_2$ height $= (T_2)$
$\mathbb{R}$ ∟ infinite Range

$F(x) = \int_{-\infty} f(x) dx$

$f(x) = \dfrac{dF(x)}{dx}$

$P(x-\epsilon \le X \le x+\epsilon) = \int_{x-\epsilon}^{x+\epsilon} f(x) dx \approx 2\epsilon f(x)$

---

# CUMULATIVE DISTRIBUTION function of RV $X$

$$\text{CDF}(x) = F_X(x) = \Pr(X \le x)$$

$\lim_{x \to -\infty} F_X(x) = 0$

$\lim_{x \to +\infty} F_X(x) = 1$

# Expectation of RV $X$, $\exists$ function $g$ $g(X)$

|  | DISCRETE $X$ — pmf | continuous $X$ — pdf |
|---|---|---|
| $E[X] =$ | $\sum_i x_i \, p(x_i)$ | $\int x \, f(x) \, dx$ |
| $E[g(X)] =$ | $\sum_i g(x_i) \, p(x_i)$ | $\int g(x) \, f(x) \, dx$ |

<span style="color:gray">$g(x)$: new random variable, want expectation.</span>

**\* statistical interpretation**: $E[X] \approx \dfrac{1}{n} \sum_{i=1}^{n} X_i$

<span style="color:red">Law of Large Numbers</span> ↑

<span style="color:gray">REPEATED sampling $X_i = X(s_i)$</span>

**\* linearity**: 2 rvs $X$ & $Y$: $E(aX + bY) = aE(X) + bE(Y)$

<span style="color:gray">Question: can 10 dots be covered with 10 identical coins with NO overlap???</span>



area of plane covered by ◯ (tile w/ no gaps using hexagon) ⬡
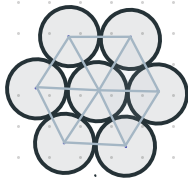
$\dfrac{\text{area of ◯ in } \square}{\text{area of } \square} : \dfrac{\pi r^2 + 6(\frac{1}{3}\pi r^2)}{\frac{3\sqrt{3}}{2}(2r)^2} = \dfrac{\pi}{2\sqrt{3}} = \boxed{0.9069}$

By linearity of expectation $10 \times 0.9069 \sim 9.069$

$E\left(\begin{smallmatrix}\text{\# dots you}\\\text{cover}\end{smallmatrix}\right) = 9.069$

<span style="color:gray">with this tiling</span>

Since $E > 9$, there must be a tiling that can cover $\boxed{\text{ALL TEN}}$ points

<span style="color:gray">Note that this approach does not work with 11...</span>
<span style="color:gray">$11(0.9069) = 9.98...$ so we are only promised coverage of 10/11 points by the above logic</span>

# moment generating functions (MGF): $n^{th}$ moment of RV $X$

$$E[X^n] = \begin{cases} \sum_i x_i^n \, p(x_i) & \text{discrete} \\ \int x^n f(x) \, dx & \text{continuous} \end{cases}$$

<span style="color:gray">useful for studying $\Sigma$ RV's</span>

| moment | UNCENTERED | CENTERED | |
|---|---|---|---|
| 1st | $E[X] = \mu$ | | avg |
| 2nd | $E[X^2]$ | $E[(X-\mu)^2]$ | spread |
| 3rd | $E[X^3]$ | $E[(X-\mu)^3]$ | asymm. |
| 4th | $E[X^4]$ | $E[(X-\mu)^4]$ | tails |

**\*\*\* Both CDF's & MGFs UNIQUELY define a probability distribution \*\*\***

# Storied Distributions

## 1. Bernoulli RV: two outcomes: success & failure
probability

$$X = \begin{cases} 1 & \text{success} \quad \theta \\ 0 & \text{failure} \quad 1-\theta \end{cases}$$

└ Distributed according to.

$$X \sim \text{Bernoulli}(\theta)$$

$$P(x): \begin{cases} \theta & \text{if } x=1 \\ 1-\theta & \text{if } x=0 \end{cases}$$

## 2. Binomial RV: independent Bernoullis w/ p success $\theta$

$N = \text{\# trials}$
$X = \text{\# successes in trials} \sim \{0,1,2 \dots N\}$

$$X \sim \text{Binomial}(N,\theta)$$

$$P(x): \binom{N}{x}\theta^x(1-\theta)^{N-x}$$

$$\binom{N}{x} = \frac{N!}{x!\,(N-x)!}$$

ABCDE

$$\left[ \frac{-\,-\,-}{5\;4\;3} \right]$$

\# arrangements order matters $= 5 \times 4 \times 3 = \dfrac{5!}{2!}$

But what if order doesn't matter? $\left(\dfrac{5!}{2!}\right)\dfrac{1}{3!}$
↓
Shuffling chosen

## 3. Geometric RV: independent Bernoullis

$X = \text{\# failures before success}$

$$X \sim \text{Geom}(\theta)$$

$$p(x): (1-\theta)^x\theta$$

## 4. Poisson RV: essentially interested in \# arrivals given rate of arrivals for MEMORYLESS process

mathematically it is a Binomial with $\lambda$ finite successes
But $N$ \# trials $\to \infty$
$\theta$ p success $\to 0$   $N\theta = \text{constant} = \lambda$
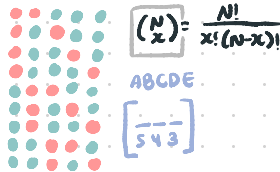
$X = \text{\# arrivals}$
$\lambda = \text{rate of arrivals}$

$$X \sim \text{Poisson}(\lambda)$$

$$P(x) = \frac{\lambda^x}{x!}e^{-\lambda}$$

*** DERIVATION ***

$$\binom{N}{x}\theta^x(1-\theta)^{N-x}$$

$e^{-\lambda}$

$$\lim_{\substack{N\to\infty \\ \theta\to 0}}\left[\frac{N!}{x!(N-x)!}\left(\frac{\lambda}{N}\right)^x\left(\frac{N-\lambda}{N}\right)^{N-x}\right] = \underbrace{\frac{n(n-1)\cdots(n-x+1)}{n^x}}_{1}\,\frac{\left(1-\frac{\lambda}{n}\right)^n}{\left(1-\frac{\lambda}{n}\right)^x}\,\frac{\lambda^x}{x!}$$

$e^{-\lambda} \to e^{-\lambda} = \lim_{t\to\infty}\left(1+\frac{-\lambda}{n}\right)^n$

$$= \boxed{\frac{\lambda^x}{x!}e^{-\lambda}}$$

mailman week by week

A Poisson distribution & a Poisson process are not the same
└ see LHS of next page

## 1. Uniform RV

$$X \sim \text{unif}(a,b)$$

$$f(x) = \begin{cases} \dfrac{1}{b-a} & \text{if } x \in [a,b] \\ 0 & \text{otherwise} \end{cases}$$



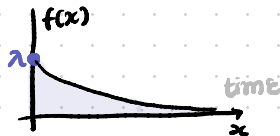## 2. Exponential RV:

interarrival time of Poisson processes
$\lambda$: arrival rate (same $\lambda$ in Poisson distr).

$$X \sim \text{Expon}(\lambda)$$

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$



DERIVATION: during interarrival, nothing arrives → Poisson w/ 0 arrivals in one unit of time

So $P(0) = \dfrac{\lambda^0}{0!}e^{-\lambda} = e^{-\lambda}$

$$P\binom{\text{no arrivals}}{\text{in time } t \text{ units}} = P\binom{\text{no arrivals}}{\text{in time } 0\text{-}1 \text{ unit}}\,P\binom{\text{no arrivals}}{\text{in time } 1\text{-}2 \text{ units}}\cdots$$

$$= e^{-\lambda}\,e^{-\lambda}\cdots e^{-\lambda}$$

So in $x$ time units,

Recognize CDF!!

$$P(X > x) = e^{-\lambda x}$$
$$1 - P(X \leq x) = e^{-\lambda x}$$
$$F(x) = 1 - e^{-\lambda x}$$
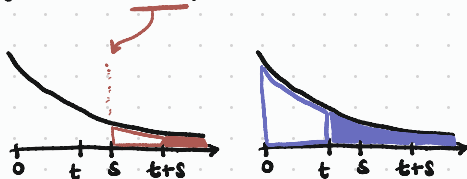$$f(x) = \frac{d}{dx}F(x) = \boxed{\lambda e^{-\lambda x}} \quad \square$$

The exponential distribution is MEMORYLESS
BOUNDARY for LIGHT/HEAVY tails.

# Mathematical Encoding of Memorylessness

- Poisson processes are MEMORYLESS
- The Exponential Distr. is the only memoryless continuous distr.
- A probability distribution is memoryless if

$$Pr(x > t+s \mid x > s) = Pr(x > t)$$



Shifting by 's' makes no difference

$$\boxed{P(x > t) = e^{-\lambda t}}$$ Recall from exp distr derivation

conditional probability definition

$$P(x > t+s \mid x > s)\, P(x > s) = P(x > t+s)$$
$$\underbrace{\quad}_{A}\ \underbrace{\quad}_{B}\qquad \underbrace{\quad}_{B}\qquad \underbrace{\quad}_{A}$$
$$\qquad\qquad e^{-\lambda s}\qquad e^{-\lambda(t+s)}$$

$$\boxed{P(x > t+s \mid x > s) = e^{-\lambda t}}$$

## Poisson process

# of ● → ∞
λ rate ●/●



$a_1 \quad a_2\ a_3 \qquad a_4$ ...

# of a's : Poisson Distr.
↔ of ● Btw ●'s : Expo. Distr.
↔ of ● Btw α ●'s : Gamma Distr, $\alpha = \mathbb{Z}^+$

<u>Heavy tail</u>: tails HEAVIER than the Exponential

light tails ↔ finite MGF ∀ orders

---

## 3. Gamma R.V.

waiting time for α arrivals of Poisson process

λ : rate of arrivals
α : # of arrivals

$$f(x; \lambda, \alpha) = \frac{1}{\Gamma(\alpha)} \frac{(\lambda x)^{\alpha}}{x} e^{-\lambda x}$$

DERIVATION       ┌─ SAME AS BEFORE BUT     k=1
                              cannot arrive           k=2
$$P(x \le x) = 1 - P(x > x)$$                          k=k

addition for mutually exclusive events

$$F(x) = 1 - \sum_{k=0}^{\alpha-1} \frac{(\lambda x)^k e^{-\lambda x}}{k!} \;\rightarrow\; \frac{dF(x)}{dx} = \frac{\lambda^{\alpha} x^{\alpha-1} e^{-\lambda x}}{\Gamma(\alpha)} \quad\square$$
         └ # of events

## 4. Gaussian Distribution RV (aka Normal)

$$x \sim \mathcal{N}(\mu, \sigma^2)$$
                                    ┌ very light tails
$$f(x) = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

central       Many quantities modeled by sums of R.V.'s
limit theorem
$$\approx \sum_{i=1}^{n} x_i \dot\sim \mathcal{N}$$
                          └ CLT



If both 1st & 2nd moment defined,
normalized sum of independent random variables w/
<u>ANY</u> underlying distribution approaches Normal. what does this mean?

Pick any $a, b \in \mathbb{R}$. $\bar{y} = \frac{1}{n}\sum_{i=1}^{n} x_i, \; x_i \sim \mathcal{N}(\mu, \sigma)$

concentration     $$\lim_{n\to\infty} Pr\left[a\frac{\sigma}{\sqrt{n}} \le \bar{y} - \mu \le b\frac{\sigma}{\sqrt{n}}\right] = \frac{1}{\sqrt{2\pi}}\int_a^b e^{-\frac{1}{2}t^2}\, dt : \bar{y} \sim \mathcal{N}\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$
bound
                                                                    └ PDF

# Parameter Estimation

Statistical functional $T: F \to \mathbb{R}$

can define $\Theta$, parameter of distribution as a functional $\Theta = T(F)$
EX: means variances medians

approximate CDF with ECDF

$\qquad \downarrow d/dx \qquad \downarrow d/dx$

$\qquad f(x) \qquad \hat{f}(x) \quad \frac{1}{n}\sum_{i=1}^{n}\delta(x-x_i)$

Estimate $\hat{\Theta} = T(\hat{F})$ ← plug in empirical data
$\hat{\Theta}$ is called a plug-in estimate

$\hat{\Theta}$'s have biases: $\langle\hat{\Theta}\rangle - \Theta = \int \hat{\Theta} f(x) \, dx - T(F)$
How off is this estimate on average?

**Question:** Given a set of data, we can estimate parameters of INTEREST via plug-in estimates, but how do we account for sampling variation?

**Solution:** Bootstrapping! → confidence intervals!

sampling your dataset     95% of the time, a 95% interval
with Replacement          of $\hat{\Theta}$ will contain $\Theta$

---

### Maximum Likelihood Estimate

Likelihood: $L(\Theta; \vec{y}) = f(\vec{y}; \Theta)$ $\qquad \vec{y} = \{y_1, y_2, y_3, \dots y_n\}$

for iid, $L(\Theta; \vec{y}) = \prod_{i=1}^{n} f(y_i)$  *sums are nicer than products*

Log-Likelihood
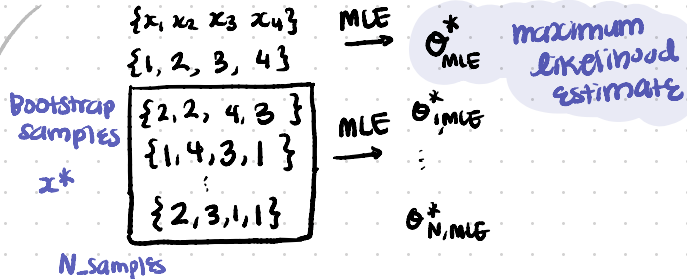$\quad \ell(\Theta; \vec{y}) = \log L(\Theta; \vec{y}) = \sum_i \log f(y_i)$

$\Theta_{MLE} = \text{argmax}_\Theta \, \ell(\Theta; \vec{y})$  set $\frac{\partial \ell}{\partial \Theta} = 0$   *logarithm is monotonic
ORDER IS PRESERVED*

---

## (A) Constructing conf ints non-parametrically

$\{x_1, x_2, x_3, x_4\} \quad \Theta$
$\{1, 2, 3, 4\} \longrightarrow \hat{\Theta}$ plug in estimate

Bootstrap samples $x^*$

$\begin{cases} \{2, 2, 4, 3\} \\ \{1, 4, 3, 1\} \\ \vdots \\ \{2, 3, 1, 1\} \end{cases} \longrightarrow$ $\hat{\Theta}_1^*$  $\hat{\Theta}^*$ Bootstrap Replicates

$\qquad \qquad \qquad \hat{\Theta}_{N\_samples}^*$

N_samples

for αth conf ints, RETRIEVE PERCENTILES of $\hat{\Theta}^*$

## (B) Constructing conf. ints parametrically

$\{x_1, x_2, x_3, x_4\} \quad$ MLE
$\{1, 2, 3, 4\} \longrightarrow \Theta_{MLE}^*$  maximum likelihood estimate

Bootstrap samples $x^*$

$\begin{cases} \{2, 2, 4, 3\} \\ \{1, 4, 3, 1\} \\ \vdots \\ \{2, 3, 1, 1\} \end{cases}$ MLE $\Theta_{1,MLE}^*$

$\qquad \qquad \qquad \Theta_{N,MLE}^*$

N_samples

for αth conf ints, RETRIEVE PERCENTILES of $\hat{\Theta}^*$

• Sometimes MLE = plug-in, but this is not always the case
• parametric inference assumes the model you have is true, and then we optimize under that assumption

1. Prof. Leonard Schulman's [CS 150 2018 Lecture Notes](#)
2. Prof. Kostia Zuev's [ACM 116 Course Webpage](#)