

THEORY OF PROBABILITY TO BAYESIAN INFERENCE

2022/02/18

SIGMA ALGEBRA AND MEASURE SPACES

[From: SCHULLER QM LECTURE 5]

DEFN: LET M BE A NON-EMPTY SET. THEN A COLLECTION OF SUBSETS $\sigma \subseteq P(M)$ [WHERE $P(M)$ IS THE POWER SET OF M , I.E., THE SET OF ALL SUBSETS OF M] IS CALLED A σ -ALGEBRA OF M IF:

$$(i) M \in \sigma$$

" M WITHOUT A "
I.E., THE COMPLEMENT OF A IN M , A^c

$$(ii) A \in M \Rightarrow M \setminus A \in \sigma$$

$$(iii) \bigcup_{n \in N} A_n \in \sigma \quad \text{FOR ANY SEQUENCE } \{A_n\}_{n \in N} \text{ WITH } A_n \in \sigma \forall n.$$

■ $A \in \sigma$ IS CALLED A MEASURABLE SET IN M

■ THE PAIR (M, σ) IS A MEASURABLE SPACE.

DEFN: GIVEN A MEASURABLE SPACE (M, σ) , A

MEASURE $\mu: \sigma \rightarrow (\bar{\mathbb{R}})^+$ IS A MAP SATISFYING:

$$(i) \mu(\emptyset) = 0$$

$$\bar{\mathbb{R}} = (\mathbb{R} \cup \{-\infty, \infty\})$$

"EXTENDED REAL NUMBER LINE"

(ii) FOR A SEQUENCE OF PAIRWISE DISJOINT MEASURABLE SETS $A_1, A_2, \dots \in \sigma$ (I.E., $(A_i \cap A_j = \emptyset \forall i, j)$):

$$\mu\left(\bigcup_{n \geq 1} A_n\right) = \sum_{n \geq 1} \mu(A_n) \quad \text{"}\sigma\text{-ADDITIONITY"}$$

■ THE TRIPLE (M, σ, μ) IS A MEASURE SPACE

■ PROPERTIES OF A MEASURE:

(i) MONOTONY: $A_1 \subseteq A_2 \Rightarrow \mu(A_1) \leq \mu(A_2)$

(ii) SUB-ADDITIONITY:

$$\mu\left(\bigcup_{n \geq 1} A_n\right) \leq \sum_{n \geq 1} \mu(A_n) \quad \text{(NO PAIRWISE-DISJOINTNESS)}$$

(iii) CONTINUITY FROM BELOW:

$$A_1 \subseteq A_2 \subseteq \dots \text{ w/ } A := \bigcup_{n \geq 1} A_n \Rightarrow \lim_{n \rightarrow \infty} \mu(A_n) = \mu(A)$$

(iv) CONTINUITY FROM ABOVE

$$A_1 \supseteq A_2 \supseteq \dots \text{ w/ } A := \bigcap_{n \geq 1} A_n \text{ AND } \mu(A_1) < \infty$$

Axioms of Probability and Their Consequences

- DEFN: A MEASURE SPACE (M, \mathcal{E}, P) IS CALLED A PROBABILITY SPACE WITH M THE SAMPLE SPACE, \mathcal{E} THE EVENT SPACE, AND $P(E)$ THE PROBABILITY OF EVENT $E \in \mathcal{E}$, WHEN THE PROBABILITY MEASURE, $P: \mathcal{E} \rightarrow (\bar{\mathbb{R}})^+$ SATISFIES:

$$P(M) = 1$$

[THEREFORE, APART FROM THE DEFINITION OF A MEASURE SPACE, THE ONLY AXIOM IS "TOTAL MEASURE ONE"]

- THE SO-CALLED AXIOMS OF PROBABILITY FOLLOW FROM THE ABOVE-DEFINED PROPERTIES OF A MEASURE SPACE:

$$\textcircled{1} \quad P(E) \geq 0 \quad \forall E \in \mathcal{E} \quad (\text{DEFN OF A MEASURE})$$

$$\textcircled{2} \quad P(M) = 1 \quad (\text{THE ONLY REAL AXIOM HERE})$$

$$\textcircled{3} \quad P\left(\bigcup_{i=1}^{\infty} E_i\right) = \sum_{i=1}^{\infty} P(E_i), \text{ FOR PAIRWISE DISJOINT } \{E_1, \dots, E_n\} \quad (\sigma\text{-ADDITIVITY OF A MEASURE})$$

Consequences

■ MONOTONICITY: $A \subseteq B \Rightarrow P(A) \leq P(B)$

$$\begin{aligned} B &= B \cap M \\ &= B \cap (A \cup A^c) \\ &= (B \cap A) \cup (B \cap A^c) \quad \text{Now A} \\ &\quad \text{DISJOINT UNION!} \\ &= A \cup (B \cap A^c) \quad \text{) } \sigma\text{-ADDITIVITY} \end{aligned}$$

$$\Rightarrow P(B) = P(A) + P(B \cap A^c) \geq P(A) \checkmark$$

■ PROBABILITY OF THE COMPLEMENT: $P(E^c) = 1 - P(E)$

$$M = E \cup E^c \quad \text{DISJOINT UNION}$$

$$P(M) = P(E) + P(E^c) \Rightarrow 1 = P(E) + P(E^c) \quad \checkmark$$

■ PROBABILITY OF THE EMPTY SET: $P(\emptyset) = P(M^c) = 1 - P(M) = 0$

■ TARGET OF THE PROBABILITY MEASURE: $P: \mathcal{E} \rightarrow [0, 1]$

$$P(E^c) = 1 - P(E) \geq 0 \Rightarrow P(E) \leq 1$$

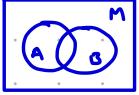
$$\Rightarrow 0 \leq P(E) \leq 1 \quad \forall E \in \mathcal{E}$$

EMPTY SET NEED
NOT BE THE ONLY
EVENT w/ $P(E) = 0$

④ "ADDITION LAW OF PROBABILITY"

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

WE CAN WRITE



$$\rightarrow (A \cup B) = (A \cap B) \cup (B^c \cap A) \cup (A^c \cap B)$$

$$A = (A \cap B) \cup (A \cap B^c)$$

$$B = (B \cap A) \cup (B \cap A^c)$$

} ALL DISJOINT UNIONS

$$P(A \cup B) = P(A \cap B) + P(B^c \cap A) + P(A^c \cap B)$$

$$P(A) = P(A \cap B) + P(A \cap B^c)$$

$$P(B) = P(B \cap A) + P(B \cap A^c)$$

$$\Rightarrow P(A \cup B) = P(A) + [P(B) - P(B \cap A)] \quad \checkmark$$

⑤ SCIENCE OF PROBABILITY: THE PROBABILITY SPACE IS ASSOCIATED TO A WELL-DEFINED "EXPERIMENT" WITH OEM EACH OF THE POSSIBLE ELEMENTARY OUTCOMES OF THAT EXPERIMENT. AN EVENT $E \in \mathcal{E}$ IS AN ELEMENTARY OR DERIVED/COMPOSITE OUTCOME, THE PROBABILITY OF AN EVENT OCCURRING IS $P(E)$.

EXAMPLES:

■ EXPERIMENT: FLIPPING A COIN TWICE

SAMPLE SPACE: $M = \{(H, H), (H, T), (T, H), (T, T)\}$

EXAMPLE EVENTS AND PROBABILITIES $= \{o_1, o_2, o_3, o_4\}$

• FIRST COIN IS HEADS: $E = \{o_1, o_2\}$

$$P(E) = \underbrace{P(o_1) + P(o_2)}_{\text{DISJOINT EVENTS}} = \underbrace{\frac{1}{4} + \frac{1}{4}}_{\text{ASSUMING FAIR & INDEP.}} = \frac{1}{2}$$

• ONE COIN IS TAILS: $E = \{o_2, o_3, o_4\}$, $P(E) = \frac{3}{4}$

■ EXPERIMENT: PLACING IN A HORSE RACE AMONG 7 HORSES

SAMPLE SPACE: $M = \{\text{THE } 7! \text{ PERMUTATIONS OF } 1234567\}$

EVENTS:

• HORSE 4 FINISHES 1ST: $E = \{6! \text{ PM OF } \boxed{4} \underline{123567}\}$, $P(E) = \frac{1}{7}$

Ross's FIRST
Course in
PROBABILITY,
And Kolmogorov's
Found of Probability

THESE DISCRETE EXAMPLES ASSUME THAT OEM $\Rightarrow o \in E$ I.E., THE ELEMENTARY OUTCOMES ARE MEASURABLE, AND THAT THEY ARE EQUAL-PROBABLE

- EXPERIMENT: ROLLING TWO DICE IN SEQUENCE
SAMPLE SPACE: $M = \{(1,1), (1,2) \dots \langle 36 \text{ ELEMENTS} \rangle\}$
EVENTS:
- A TOTAL OF 7: $E = \{(1,6), (6,1), (3,4), (4,3), (2,5), (5,2)\}$
 - $P(E) = 6/36 = 1/6$

THIS TYPE OF CONTINUOUS
EXAMPLE WAS UNIFIED w/
THE DISCRETE CASE BY
KOLMOGOROV. THERE IS
NO MEASURE ASSOCIATED TO
 $X \in M$, i.e., FOR $x \in \mathbb{R}$
 $x \notin E$. POINTS ARE NOT
MEASURABLE, BUT:
 $(a,b), (a,b], [a,b] \in E$

- EXPERIMENT: MEASURING THE LIFETIME OF A LIGHTBULB
SAMPLE SPACE: $M = \mathbb{R}_+^+$ (Hours)
EVENTS:
- BULB DOES NOT LAST 12h: $E = [0, 12)$, $P(E) = ?$ *

TERMINOLOGY:

Theory of Sets	Events in Probability Expt.
$A \cap B = \emptyset$ OR $A_1 \cap A_2 \cap \dots \cap A_n$	EVENTS ARE <u>INCOMPATIBLE</u> (MUTUALLY-EXCLUSIVE, INDEPENDENT?)
$E = A_1 \cap A_2 \cap \dots \cap A_n$	EVENT E IS THE <u>SIMULTANEOUS OCCURRENCE</u> OF THE EVENTS $\{A_i\}$
$E = A_1 \cup A_2 \cup \dots \cup A_n$	EVENT E IS THE <u>OCCURRENCE</u> OF AT LEAST ONE OF $\{A_i\}$
COMPLEMENT OF A, A^c	THE <u>NON-OCCURRENCE</u> OF A
$A = \emptyset$	EVENT A IS <u>IMPOSSIBLE</u>
$A = M$	EVENT A MUST OCCUR
THE <u>PARTITION*</u> OF M BY $\Delta = \{A_1, A_2, \dots, A_n\} \subseteq \Sigma$	<u>ONE AND ONLY ONE</u> <u>EVENT</u> A_i WILL OCCUR
$B \subset A$	THE OCCURRENCE OF B IMPLIES THE OCCURRENCE OF A.

* - THE SETS ARE
PAIRWISE DISJOINT,
 $A_i \cap A_j = \emptyset \forall i \neq j$
AND $\bigcup A_i = M$.

ALGEBRA OF SETS (FOR $A, B, C \subseteq M$, WITH \emptyset THE EMPTY SET)

SINCE A PROBABILITY SPACE HAS A SET AS ITS BASE STRUCTURE...
LET US REVIEW.

(WIKIPEDIA)

WARNING:

THIS "MAKES SENSE" w/
 $a \times (b+c) = ab+ac$
But This Does NOT!
 $a+(b+c) \neq (a+b)+(a+c)$

■ UNION AND INTERSECTION ARE COMMUTATIVE:

$$A \cup B = B \cup A \text{ AND } A \cap B = B \cap A$$

■ UNION AND INTERSECTION ARE ASSOCIATIVE

$$A \cup (B \cup C) = (A \cup B) \cup C \text{ AND } (\text{same for } \cap)$$

■ DISTRIBUTIVE PROPERTIES

$$A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$$

$$A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$$

■ IDENTITIES

$$A \cup \emptyset = A \text{ AND } A \cap M = A$$

■ COMPLEMENTS

$$A \cup A^c = M \text{ AND } A \cap A^c = \emptyset$$

WE COULD ALSO PROVE DEMORGAN'S LAWS:

$$\left(\bigcup_{i=1}^n A_i \right)^c = \bigcap_{i=1}^n A_i^c \quad \text{AND} \quad \left(\bigcap_{i=1}^n A_i \right)^c = \bigcup_{i=1}^n A_i^c$$

CONDITIONAL PROBABILITIES AND BAYES' THEOREM

2022/02/19

■ DEFN: THE CONDITIONAL PROBABILITY OF EVENT A UNDER CONDITION B IS DEFINED (WHEN $P(B) > 0$):

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

KOLMOGOROV USES:
 $P_B(A) = P(A|B)$
 WHICH IS Nicer BECAUSE $P_B(A)$ IS THE "PROBABILITY OF A".

IN EXPERIMENTAL TERMS, THIS IS THE PROBABILITY THAT EVENT A OCCURS UNDER THE CONDITION THAT EVENT B ALSO OCCURS.

■ THE MULTIPLICATION THEOREM FOLLOWS FROM INDUCTION

OF $P(A \cap B) = P(A) \cdot P(B|A)$:

$$P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1) P(A_2 | A_1) P(A_3 | A_1 \cap A_2) \dots \times P(A_n | A_1 \cap A_2 \cap \dots \cap A_{n-1})$$

■ FOR FIXED EVENT B, $(M, E, P(\cdot|B))$ IS A PROBABILITY SPACE SINCE:

$$P(M|B) = P(M \cap B) / P(B) = P(B) / P(B) = 1$$

AND FOR DISJOINT $\{A_1, \dots, A_n\}$

$$P\left(\bigcup_{i=1}^n A_i | B\right) = \frac{P((\cup A_i) \cap B)}{P(B)} = \frac{P((A_1 \cap B) \cup \dots \cup (A_n \cap B))}{P(B)}$$

$$= \sum \frac{P(A_i \cap B)}{P(B)} = \sum_{i=1}^n P(A_i | B)$$

Note, e.g.,

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

$$P(A \cap B \cap C) = P(A \cap B) \cap C$$

$$= P(A \cap B) \cdot P(C|A \cap B)$$

$$= P(A) \cdot P(B|A) \cdot P(C|A \cap B)$$

$P(\cdot|B)$ SATISFIES

σ -ADDITIVITY
AND

$$P(M|B) = 1$$

✓

THEOREM OF BAYES : THE DEFINITIONS OF $P(A|B)$ AND $P(B|A)$ ALLOW US TO WRITE (USING $P(A \cap B) = P(B \cap A)$) :

$$P(A|B) P(B) = P(B|A) P(A)$$

OR

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

IF WE THEN CONSIDER A PARTITION OF M , $A = \{A_1, \dots, A_n\}$, WE CAN WRITE:

$$\begin{aligned} B &= B \cap M = B \cap \left(\bigcup_{i=1}^n A_i \right) \\ &= (B \cap A_1) \cup (B \cap A_2) \cup \dots \cup (B \cap A_n) \end{aligned}$$

AND THEREFORE:

$$P(B) = \sum_{i=1}^n P(B \cap A_i) = \sum_{i=1}^n P(B|A_i) P(A_i)$$

AND SO, FOR ANY A_i WE CAN WRITE:

$$P(A_i|B) = \frac{P(B|A_i) P(A_i)}{P(B|A_1) P(A_1) + \dots + P(B|A_n) P(A_n)}$$

[FOR THE PARTITION OF M , $A = (A_1, \dots, A_n)$]

WHICH IS CALLED BAYES' THEOREM.

THE EVENTS $A_i \in A$ ARE USUALLY CALLED HYPOTHESES, AND THE ABOVE EQUATION GIVES THE PROBABILITY OF THE A_i HYPOTHESIS GIVEN THAT EVENT B HAS OCCURRED.

MORE TERMINOLOGY:

- ① $P(A_i)$ IS THE PRIOR PROBABILITY OF HYPOTHESIS A_i
- ② $P(A_i|B)$ IS THE POSTERIOR PROBABILITY (ONCE B HAS OCCURRED)
- ③ DENOMINATOR $\sum_i P(B|A_i) P(A_i)$ IS A NORMALIZING CONST

EXAMPLE: PROBABILITY OF HAVING DISEASE (D) GIVEN A POSITIVE TEST (+). $A = \{D, D^c\}$, $M = D \cup D^c = \{\text{+}\} \cup \{\text{-}\}$

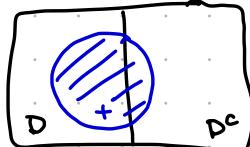
$$P(D|+) = \frac{P(+|D) \cdot P(D)}{P(+|D) P(D) + P(+|D^c) (1 - P(D))}$$

WHERE $P(D)$ IS THE PREVALENCE OF DISEASE IN THE COMMUNITY, $P(+|D)$ IS THE SENSITIVITY OF THE TEST AND $(1 - P(+|D^c))$ IS THE SPECIFICITY OF THE TEST.

WE COULD USE A
PARTITION OF B
HERE, RIGHT?
 $B = B \cap B$

(INITIAL ASSUMPTIONS)
ARE MODIFIED
GIVEN NEW EVIDENCE

(SUM OF THE TERMS
IN THE NUMERATOR)



AS A CONNECTION TO APPLICATION, CALL THE PARTITION THE HYPOTHESES $\mathcal{H} = \{H_1, H_2, \dots, H_3\}$ CALL AN OBSERVED EVENT EVIDENCE, E. GIVEN A PRIOR PROBABILITY OF HYPOTHESIS H_i , $P(H_i)$, WE FIND THE UPDATED POSTERIOR PROBABILITY OF H_i , GIVEN NEW EVIDENCE,

$$P(H_i | E) = \frac{P(E | H_i) \cdot P(H_i)}{\sum_j P(E | H_j) P(H_j)}$$

WHERE $P(E | H_i)$ IS SOMETIMES CALLED THE LIKELIHOOD OF THE HYPOTHESIS.

EXAMPLE : A DETECTIVE IS 50% CONFIDENT OF A PARTICULAR SUSPECT'S GUILT ($P(H) = 0.5$). NEW EVIDENCE SHOWS THAT THE PERPETRATOR IS LEFT-HANDED ($P(E | H) = 1$). IF THE PARTICULAR SUSPECT IS LEFT-HANDED (LIKE 19% OF THE POPULATION, $P(E | H^c) = 0.19$):

$$P(G | L) = \frac{P(L | G) \cdot P(G)}{P(L | G) P(G) + P(L | H^c) P(H^c)} = \frac{1 \cdot 0.5}{1 \cdot 0.5 + 0.19 \cdot 0.4} = 87\%$$

EXAMPLE : 2-CHILD FAMILY MOVES INTO NEIGHBORHOOD. YOU MEET MOTHER & DAUGHTER. WHAT IS PROBABILITY OTHER CHILD IS A GIRL?

$$P(GG | g) = \frac{P(g | GG) \cdot P(GG)}{P(g | GG) \cdot P(GG) + P(g | GB) \cdot P(GB) + P(g | BB) \cdot P(BB)}$$

ASSUMING THAT IT IS EQUALLY LIKELY FOR THE MOTHER TO BE w/ ANY CHILD AND THAT THE COUPLE'S SEQUENCE OF OFFSPRING IS EQUALM LIKELY TO BE $\{GG, GB, BG, BB\}$:

$$P(GG | g) = \frac{1 \cdot \frac{1}{4}}{\frac{1}{4} + \frac{1}{2} \cdot \frac{1}{2} + 0 \cdot \frac{1}{4}} = \frac{\frac{1}{4}}{\frac{1}{4} + \frac{1}{4}} = \frac{1}{2}$$

EXAMPLE : 3 TWO-SIDED CARDS w/ RR (RED-RED), RB (RED-BLACK) BB (BLACK-BLACK) ON FRONT-BACK. GIVEN A TIE IS SHOWN, WHAT CHANCE IS IT RR?

$$P(RR | R) = \frac{P(R | RR) \cdot P(RR)}{P(R | RR) \cdot P(R) + P(R | RS) P(RB) + P(R | BS) \cdot P(BB)} = \frac{\frac{1}{3} \cdot \frac{1}{2}}{\frac{1}{3} + \frac{1}{2} \cdot \frac{1}{3}} = \frac{2}{3}$$

AS WE WILL SEE LATER,
 $P(E | H)$ IS ALSO CALLED
THE LIKELIHOOD FUNCTION,
I.E., WHEN VIEWED AS A
FUNCTION OF H, WITH FIXED
E:
 $L(H) = P(E | H)$

H = GUILTY (G)

H^c = NOT GUILTY (NG)

E = LEFT-HANDED (L)

EXPERIMENT HERE
INVOLVES THE
PROPERTIES OF THE
SUSPECT, SO:

$P(L) = \underbrace{P(L | NG)}_{\text{LEFT-HANDED}} \text{ OF POPULATION}$

M = {GG, GB, BB}

$g = \text{MET ONE}$
GIRL

THERE ARE 6 EQUALLY LIKELY OUTCOMES HERE:

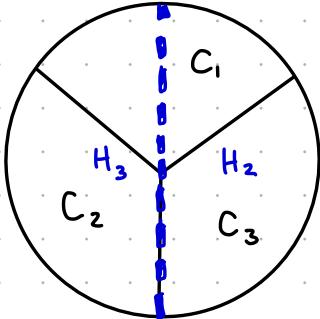
R_F(RR), R_B(RR), R_B(RR)

B_F(RB), B_F(BS), B_B(BS)

SO 2 OF 5 "R" OPTIONS ARE FROM THE RR CARD.

This problem became famous after appearing in the "Ask Marilyn" advice column (Parade Magazine, 1990/9/9). The solution given was, essentially:

The probability of your selection did not change (from $\frac{1}{3}$) just because the host opens a door, but the probability of the other two doors ($\frac{2}{3}$) also didn't change... it is just all in the unopened.



EXAMPLE (MONTY HALL): BEHIND THREE DOORS ARE A CAR AND THREE GOATS. YOU CHOOSE DOOR 1. THE GAME SHOW HOST OPENS ANOTHER DOOR, SHOWS YOU A GOAT AND OFFERS YOU THE CHANCE TO CHANGE YOUR CHOICE. SHOULD YOU?

$$\begin{aligned} P(C_1 \text{ IN } | \text{Host opens } H_2) &= P(C_1 | H_2) \\ &= \frac{P(H_2 | C_1) \cdot P(C_1)}{P(H_2 | C_1) P(C_1) + P(H_2 | C_2) P(C_2) + P(H_2 | C_3) P(C_3)} \\ &= \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2} \cdot \frac{1}{3} + 0 \cdot \frac{1}{3} + 1 \cdot \frac{1}{3}} \\ &= \frac{\frac{1}{6}}{\frac{1}{6} + \frac{2}{6}} = \frac{1}{3} \end{aligned} \quad (C_1 \cup C_2 \cup C_3 = M)$$

YES, YOU SHOULD CHANGE! CONSIDER ALSO:

$$\begin{aligned} P(C_3 | H_2) &= \frac{P(H_2 | C_3) P(C_3)}{P(H_2 | C_3) P(C_3) + P(H_2 | C_2) P(C_2) + P(H_2 | C_1) P(C_1)} \\ &= \frac{1 \cdot \frac{1}{3}}{\frac{1}{3} + 0 + \frac{1}{2} \cdot \frac{1}{3}} = \frac{2}{3} \end{aligned}$$

AND:

$$\begin{aligned} P(C_1 | H_3) &= \frac{P(H_3 | C_1) P(C_1)}{P(H_3 | C_1) P(C_1) + P(H_3 | C_2) P(C_2) + 0} \\ &= \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2} \cdot \frac{1}{3} + 1 \cdot \frac{1}{3}} = \frac{1}{3} \end{aligned}$$

So, whichever door the host opens, your current door only has $\frac{1}{3}$ probability of containing the car, while the other closed door has $\frac{2}{3}$ prob.

DEFINITION: THE ODDS RATIO OF A PARTICULAR HYPOTHESIS H IS:

$$OR(H) := \frac{P(H)}{P(H^c)} = \frac{P(H)}{1 - P(H)}$$

IF, FOR EXAMPLE $P(H) = \frac{3}{4}$ THEN $OR(H) = 3$. ONE SAYS "THE ODDS ARE 3 TO 1 IN FAVOR".

DEFN THE INTRODUCTION OF NEW EVIDENCE CAN CHANGE THE ODDS RATIO VIA BAYES THEOREM:

$$\text{OR}(H|E) = \frac{P(H|E)}{P(H^c|E)} = \frac{P(E|H)}{P(E|H^c)} \left(\frac{P(H)}{P(H^c)} \right)$$

OR(H)

This will be called
THE BAYES FACTOR

"ODDS ARE 1 TO 2 AGAINST"
FINDING THE CAR IN "1"
(ALTHOUGH SAME FOR ANY DOOR)

SO THE RATIO OF THE LIKELIHOODS OF THE EVIDENCE GIVEN THE HYPOTHESIS OR ITS COMPLEMENT YIELDS THE FACTOR BY WHICH THE ODDS RATIO CHANGES.

- Ex: IN THE Monty Hall Problem:

$$\text{OR}(\text{car in } 1) = \frac{P(1)}{P(1^c)} = \frac{1/3}{2/3} = \frac{1}{2}$$

BUT

$$\text{OR}(\text{car in } 1 \mid \text{host opens door 2}) = \frac{P(H_2|1)}{P(H_2|1^c)} \cdot \text{OR}(1)$$

$$\text{OR}(1|H_2) = \frac{1/2}{1} \cdot \text{OR}(1) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

SO THE ODDS HAVE DECREASED!

- DEFN** A MODEL OF A SYSTEM, m , COULD BE CONSIDERED A HYPOTHESIS, AND DATA ACCUMULATED, D , COULD BE CONSIDERED EVIDENCE IN FAVOR (OR AGAINST) SUCH A MODEL:

$$\xrightarrow{\text{Posterior}} P(m|D) = \frac{P(D|m) \cdot P(m)}{P(D)}$$

Prior

WHERE WE AGAIN CALL $P(D|m)$ THE LIKELIHOOD OF MODEL GIVEN THE DATA. IF THE MODEL DEPENDS ON PARAMETERS, $\vec{\theta}$, WE CAN INCLUDE DEPENDENCE ON THE VALUES OF THOSE PARAMETERS!

$$P(m \cap (\vec{\theta} = \vec{\theta}^*) | D) = \frac{P(D|m \cap (\vec{\theta} = \vec{\theta}^*)) \cdot P(m \cap (\vec{\theta} = \vec{\theta}^*))}{P(D)}$$

- DEFN**: IN THE COMPARISON OF TWO MODELS, THE BAYES FACTOR IS THE RATIO OF THEIR LIKELIHOODS:

$$K(m_1, m_2) := \frac{P(D|m_1)}{P(D|m_2)} = \frac{P(m_1|D)}{P(m_2|D)} \cdot \frac{P(m_2)}{P(m_1)}$$

WHICH IS THE RATIO OF THEIR POSTERIOR PROBABILITIES IF THE MODELS ARE INITIALLY EQUALLY LIKELY.

Interpretation of Bayes Factor (Jeffreys, 1961)

K	SIGNIFICANCE
$< 10^0$	NEGATIVE (Supports m_2)
$10^0 - 10^1 (1-5)$	BARELY ANYTHING
$10^1 - 10^2 (5-10)$	SUBSTANTIAL
$10^2 - 10^3 (10-10)$	STRONG
$10^3 - 10^4 (10-20)$	VERY STRONG
$> 10^4 (> 20)$	DECISIVE

IF $P(D|m \cap (\vec{\theta} = \hat{\vec{\theta}}))$ IS USED, WHERE $\vec{\theta}$ IS THE MAXIMUM LIKELIHOOD ESTIMATE OF THE PARAMETER, THEN THIS BECOMES THE LIKELIHOOD TEST OF CLASSIC (FREQUENTIST) STATISTICS. BUT FOR THE BAYES FACTOR, IT IS ASSUMED THAT, FOR A MODEL THAT DEPENDS ON PARAMETERS:

$$P(D|m) = \sum P(D|m \cap \theta_i) \cdot P(\theta_i|m)$$

WHERE THE EVENT $\theta_i := (\vec{\theta} = \vec{\theta}_i)$ AND THE SUM IS OVER ALL POSSIBLE VALUES OF THE PARAMETER VECTOR (WHICH IS A PARTITION).

INDEPENDENT TRIALS AND INDEPENDENT EVENTS

3/10-18/2022

- DEF'N: GIVEN A PROBABILITY SPACE (M, \mathcal{E}, P) , A PARTITION $A = \{A_i\}$ OF M INTO DISJOINT EVENTS $\{A_i\}$ IS CALLED A TRIAL. THIS CAN GENERALLY BE THOUGHT OF AS A SPECIFIC QUESTION FOR WHICH THE DISJOINT EVENTS PROVIDE ALL POSSIBLE ANSWERS.

EXAMPLES

- TWO FAIR COINS TOSSED: FIRST COIN HEADS?

$A = \text{FIRST COIN IS HEADS}$

$$A = \{A, A^c\}$$

- TWO FAIR COINS TOSSED: HOW MANY HEADS?

$A_i = i \text{ HEADS APPEARED}$

$$A = \{A_0, A_1, A_2\}$$

- FIVE CARDS DEALT: WHAT KIND OF POKER HAND?

$$A = \{(HIGH CARD), (ONE PAIR), (TWO PAIR), \dots, (TSF)\}$$

FOR A SET OF n TRIALS, $(A^{(1)}, A^{(2)}, \dots, A^{(n)})$, WITH

$$A^{(i)} = \{A_1^{(i)}, A_2^{(i)}, \dots, A_{r_i}^{(i)}\}$$

WE KNOW:

$$P(\vec{\theta}) = \sum P(\vec{\theta}|A_i) P(A_i)$$

FOR A PARTITION $\{A_i\}$, BUT $P(D|m)$

FOR FIXED m IS A PROBABILITY ITSELF... CALL IT (A LA KOLMOGOROV):

$$P_m(D) := P(D|m)$$

THEN:

$$P_m(D) = \sum P_m(D|A_i) P_m(A_i)$$

WHICH IS, UNWRAPPING...

$$P(D|m) = \sum P(D|m \cap A_i) \cdot P(A_i|m)$$

I'M CALLING KOLMOGOROV'S "EXPERIMENTS" "TRAILS". SEEMS THAT THE MODERN EQUIVALENT IS "CUSES" SEE math.stackexchange.com/questions/4402669

THERE ARE $\Gamma = \Gamma_1, \Gamma_2, \dots, \Gamma_n$ PROBABILITIES FOR THE OUTCOME OF THOSE TRIALS (FOR THE n ANSWERS TO THE n QUESTIONS ASKED):

$$P_{g_1, \dots, g_n} := P\left(A_{g_1}^{(1)} \cap A_{g_2}^{(2)} \cap \dots \cap A_{g_n}^{(n)}\right) ; \sum_{\vec{g}} P_{\vec{g}} = 1$$

- DEF'N : A SET OF n TRIALS $(A^{(1)}, \dots, A^{(n)})$ ARE MUTUALLY INDEPENDENT TRIALS IF FOR ANY ALLOWED (g_1, \dots, g_n) THE FOLLOWING EQUATION HOLDS TRUE:

$$P_{g_1, g_2, \dots, g_n} = P(A_{g_1}^{(1)}) \cdot P(A_{g_2}^{(2)}) \cdot \dots \cdot P(A_{g_n}^{(n)}),$$

AMONG THESE $\Gamma = \Gamma_1 \cdot \Gamma_2 \cdot \dots \cdot \Gamma_n$ EQUATIONS THERE ARE ONLY $(\Gamma - \sum_i \Gamma_i + (n-1))$ INDEPENDENT EQUATIONS.

- THEOREM : IF n TRIALS ARE INDEPENDENT, THEN ANY $m < n$ OF THEM ARE ALSO INDEPENDENT.

- DEF'N : A SET OF n EVENTS (A_1, \dots, A_n) ARE MUTUALLY INDEPENDENT EVENTS IF THE TRIALS:

$$\mathcal{A}^{(i)} := \{A_i, A_i^c\} \quad (i=1, \dots, n)$$

ARE MUTUALLY INDEPENDENT TRIALS.

- From THE ABOVE WE SEE THAT $\Gamma_i = 2 \Rightarrow \Gamma = 2^n$ AND THEREFORE THERE ARE $2^n - 2n + n - 1 = 2^n - n - 1$ INDEPENDENT EQUATIONS WHICH CAN BE WRITTEN:

$$P(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_m}) = P(A_{i_1}) \cdot P(A_{i_2}) \cdot \dots \cdot P(A_{i_m})$$

$$\text{w/ } 1 \leq i_1 < i_2 < \dots < i_m \leq n \quad (m = 1, 2, \dots, n)$$

- For TWO INDEPENDENT EVENTS THERE IS ONE EQN:

$$P(A_1 \cap A_2) = P(A_1) \cdot P(A_2)$$

WHICH IS AN INDEPENDENT EQN OF THE FOUR:

I DON'T SEE THIS
BUT HAVEN'T
THOUGHT TOO LONG...

I'M NOT SURE HOW TO
COUNT THESE TO GET

$2^n - n - 1$
BUT HAVEN'T THOUGHT MUCH

$$A_1 = \{A_1, A_1^c\}$$

$$P(A_1 \cap A_2) = P(A_1) \cdot P(A_2)$$

$$A_2 = \{A_2, A_2^c\}$$

$$P(A_1^c \cap A_2) = P(A_1^c) \cdot P(A_2)$$

$$P(A_1 \cap A_2^c) = P(A_1) \cdot P(A_2^c)$$

$$P(A_1^c \cap A_2^c) = P(A_1^c) \cdot P(A_2^c)$$

CONSIDER, FOR EXAMPLE, THE LHS OF THE 3RD EQN:

$$\begin{aligned} P(A_1 \cap A_2^c) &= P(A_1 \cap (\underbrace{A_2 \cup A_2^c}_m)) - P(A_1 \cap A_2) \\ &= P(A_1) - P(A_1 \cap A_2) \\ &= P(A_1) - P(A_1) \cdot P(A_2) \quad \text{USING 1ST EQN} \\ &= P(A_1) [1 - P(A_2)] \\ &= P(A_1) P(A_2^c) \end{aligned}$$

SO WE SEE THAT EQN (1) IMPLIES THE OTHERS.

② RELATIONSHIP TO CONDITIONAL PROBABILITY: THE BASIC EQUATION REGARDING INDEPENDENCE OF n EXPERIMENTS, AND, IN PARTICULAR, ITS RESTRICTION TO ANY SUBSET OF $m \leq n$ DISTINCT EXPERIMENTS:

(SEE THEOREM AFTER DEFN)

$$P(A_{g_1}^{(i_1)} \cap \dots \cap A_{g_m}^{(i_m)}) = P(A_{g_1}^{(i_1)}) \cdot \dots \cdot P(A_{g_m}^{(i_m)}),$$

TOGETHER W/ THE "MULTIPLICATION THEOREM," IMPLIES THAT:

$$P(A_{g_m}^{(i_m)} | A_{g_1}^{(i_1)} \cap \dots \cap A_{g_{m-1}}^{(i_{m-1})}) = P(A_{g_m}^{(i_m)})$$

I.E., CONDITIONING A TRIAL'S EVENT ON THE OCCURANCE OF EVENTS IN OTHER, INDEPENDENT, TRIALS DOES NOT CHANGE THAT EVENT'S PROBABILITY. FOR TWO INDEPENDENT EVENTS:

OBVIOUS FROM
DEFN OF
CONDITIONAL PROBS

$$P(A \cap B) = P(A) \cdot P(B) \Rightarrow P(A|B) = P(A)$$

(AND $P(B|A) = P(B)$)

■ THEOREM: A NECESSARY AND SUFFICIENT CONDITION FOR THE INDEPENDENCE OF TRIALS $A^{(1)}, A^{(2)}, \dots, A^{(n)}$ IS THAT THE CONDITIONAL PROBABILITY OF SOME $A_g^{(i)}$ UNDER THE ASSUMPTION THAT SEVERAL OTHER TRIALS, $A^{(i_1)}, A^{(i_2)}, \dots, A^{(i_k)}$ HAVE HAS DEFINITE RESULTS, $A_{g_1}^{(i_1)}, A_{g_2}^{(i_2)}, \dots, A_{g_k}^{(i_k)}$ IS EQUAL TO THE (UNCONDITIONAL / ABSOLUTE) PROBABILITY OF $A_g^{(i)}$ (WHEN $P(A_g^{(i)}) > 0$)

■ THEOREM: SAME AS ABOVE BUT FOR EVENTS.

④ DEFN: EVENTS A AND B ARE CONDITIONALLY INDEPENDENT GIVEN C IF

$$P(A \cap B | C) = P(A | C) \cdot P(B | C)$$

④ EXAMPLES

■ (6 ON EACH ROLL OF DIE) VS. (6 ON FIRST, SUM IS 8)

■ (DRAWING TWO CARDS w/ REPLACEMENT) VS. (w/o REPLACE.)

■ A MARKOV CHAIN (AKA, MARKOV PROCESS)

DESCRIBES A SEQUENCE OF EVENTS FOR WHICH THE PROBABILITY OF THE NEXT EVENT DEPENDS ONLY ON THE CURRENT STATE. THUS, CONDITIONED ON THE CURRENT STATE, THE FUTURE AND PAST STATES ARE INDEPENDENT



RANDOM VARIABLES

2022/02/26

④ RECALL THAT A MEASURABLE MAP FROM ONE MEASURABLE SPACE (M, \mathcal{E}) TO ANOTHER (N, \mathcal{F}) IS A MAP, $g: M \rightarrow N$, FOR WHICH $\text{PREIM}_g(f \in \mathcal{F}) \in \mathcal{E}$. IF (M, \mathcal{E}, P) IS A MEASURE SPACE, THEN ONE CAN DEFINE A MEASURE ON (N, \mathcal{F}) AS THE PUSH-FORWARD OF P:

$$g_* P: \mathcal{F} \rightarrow \overline{\mathbb{R}}^+, (g_* P)(f) := P(\text{PREIM}_g(f))$$

SCHÜLER
QM OF
(measures)

MEASURABLE
MAPS ARE (WITH
MINOR ADDITIONS)
THOSE THAT ARE
INTEGRABLE.

DEFINITION: Given a probability space (M, \mathcal{E}, P) and a measurable space (N, \mathcal{F}) , a measurable map, $g: M \rightarrow N$, is called an (N, \mathcal{F}) -valued random variable.

RECALL THAT, given a topological space (M, Θ_M) , the Borel σ -algebra of (M, Θ_M) is:

$$\sigma(\Theta_M) = \{A \in \mathcal{M} \mid \forall \sigma\text{-algebras } \mathcal{T} \text{ of } M \text{ with } \Theta_M \subseteq \mathcal{T} : A \in \mathcal{T}\}$$

i.e., it is the σ -algebra generated by the topology.

For $(\mathbb{R}^d, \Theta_{std})$ we call $\sigma(\Theta_{std})$ the Borel σ -algebra of \mathbb{R}^d .

PROPOSITION: Given a probability space (M, \mathcal{E}, P) , a map $X: M \rightarrow \mathbb{R}$ is a real-valued random variable if

$$\text{preim}_X((-\infty, a)) \in \mathcal{E} \quad \forall a \in \mathbb{R}$$

PROOF: WE ASSUME THE TARGET IS THE MEASURE-ASUE SPACE $(\mathbb{R}, \sigma(\Theta_{std}))$. THE ABOVE DEFINITION GUARANTEES THAT X IS A MEASURABLE MAP FOR A TARGET $(\mathbb{R}, \mathcal{E}_{loc})$ WHERE

$$\mathcal{E}_{loc} = \sigma(E_{loc}) \text{ w/ } E_{loc} = \{(-\infty, a)\}$$

IS THE σ -ALGEBRA GENERATED BY THE SET OF "LEFT-OPEN RAYS". THUS WE MUST PROVE THAT

$$\sigma(E_{loc}) = \sigma(\Theta_{std})$$

WE WILL USE THE FACT THAT, FOR A SUBSET OF A σ -ALGEBRA: $A \subseteq \sigma \Rightarrow \sigma(A) \subseteq \sigma$. NOW, KNOWING THAT AN OPEN INTERVAL IS OPEN, WE CAN WRITE A HALF-CLOSED INTERVAL AS THE COUNTABLE INTERSECTION OF OPEN INTERVALS:

GENERATOR:
 $\sigma(E)$ FOR
 $E \subseteq P(M)$ IS
 $\sigma(E) = \bigcap_{\text{all } \sigma\text{-measurable } E \text{ containing } E} \sigma$
AND IS THE
SMALLEST σ -
ALGEBRA CONTAIN-
ING E .

(Simpler sans:)

LEMMA: To show
that $g: M \rightarrow N$ is
measurable for (M, \mathcal{E}) ,
 $(N, \sigma(E))$, it
suffices to check the
elements of the
generating set:
 $\text{preim}_g(E \in \mathcal{E}) \in \mathcal{E}$

SEE Prop. 1.9 of
people.clas.utl.edu/
pascocj/files/
6616notes01dec2017.pdf

SEE:
[mpaldrige.github.io/
teaching/
math0042-notes-02.pdf](https://mpaldrige.github.io/teaching/math0042-notes-02.pdf)

$$[a, b) = \bigcap_{n=1}^{\infty} (a - \frac{1}{n}, b)$$

FURTHER, WE CAN WRITE: $\in E_{\text{OI}}$

$$(-\infty, a) = \bigcup_{n \in \mathbb{N}} [n, a)$$

Thus, AS SETS,

$$E_{\text{LOR}} \subseteq E_{\text{HCl}} \subseteq E_{\text{OI}} \subseteq \Theta_{\text{STD}}$$

AND, THEREFORE

$$\sigma(E_{\text{LOR}}) \subseteq \sigma(\Theta_{\text{STD}})$$

SO THEN IT REMAINS TO SHOW:

$$\sigma(\Theta_{\text{STD}}) \subseteq \sigma(E_{\text{LOR}})$$

WHICH COULD BE DONE BY SHOWING ... ???

- DEF'N: THE PROBABILITY FUNCTION OF A REAL-VALUED RANDOM VARIABLE ON (M, \mathcal{E}, P) IS THE PUSH-FORWARD OF P :

$$P(X): \sigma(\Theta_{\text{STD}}) \rightarrow \overline{\mathbb{R}}^+, P^{(X)}(u) := P(\text{Pacm}_X(u))$$

This allows us to view $(\mathbb{R}, \sigma(\Theta_{\text{STD}}), P^{(X)})$ as a PROBABILITY SPACE.

- DEF'N: THE DISTRIBUTION FUNCTION (AKA THE CUMULATIVE DISTRIBUTION FUNCTION) OF A RANDOM VARIABLE X IS THE FUNCTION

$$F_X(a) = P^{(X)}((-\infty, a]) = "P(X < a)"$$

And $F_X(a)$ can be stated, as written above, as "THE PROBABILITY THAT THE RANDOM VARIABLE IS LESS THAN a ".

• CLEARLY,

$$F_X(-\infty) = 0 \quad \text{and} \quad F_X(+\infty) = 1$$

• AND

$$P^{(X)}([a, b)) = P^{(X)}((-\infty, a]^c \cap (-\infty, b))$$

$\sigma(\Theta_{\text{STD}})$ THE
BOREL σ -ALGEBRA
OF \mathbb{R} .

I.E., A MAP:
 $F: \mathbb{R} \rightarrow [0, 1]$

* math.stackexchange.com/questions/2142162

Let's allow ourselves to write:

$$P^{(X)}((-\infty, a]) = "P(X \leq a)"$$

WHERE

$$"P(X \leq a)" = P("X \leq a")$$

AND

$$"X \leq a" = \{\omega \in \Omega | X(\omega) \leq a\}$$

Recall the Addition Law:

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

$$= 1 - P^X((-\infty, a)) + P^X((-\infty, b))$$

$$- P^X((-\infty, a)^c \cup (-\infty, b))$$

↑

$$= 1 - F_X(a) + F_X(b) - P^X((-\infty, +\infty))$$

$$= F_X(b) - F_X(a)$$

② Which implies that

$$F_X(a) \leq F_X(b) \text{ for } a \leq b$$

Showing that $F_X(x)$ is a nondecreasing fcn.

■ DEF'N: WHEN THE DISTRIBUTION FUNCTION $F^X(a)$ IS DIFFERENTIABLE WE DEFINE THE PROBABILITY DENSITY OF X AT a :

$$f_X(a) = \frac{d}{dx} F_X(x) \Big|_{x=a}$$

AND IF THE FOLLOWING HOLDS $\forall x \in \mathbb{R}$ THEN THE DISTRIBUTION IS CALLED CONTINUOUS:

$$F_X(x) = \int_{-\infty}^x f_X(a) da$$

AND THEN THE PROBABILITY OF ANY BOREL SET $A \subseteq \mathbb{R}$ IS

$$P^X(A) = \int_A f_X(a) da$$

■ THE DISTRIBUTION AND DENSITY FUNCTIONS CAN BE GENERALIZED FOR CONDITIONAL PROBABILITIES. DEFINE:

$$P^X(A | B) := "P(X \in A | B)"$$

$$\stackrel{\text{(NOTE AGAIN: DIFFERENT P MAPS)}}{\Rightarrow} = P(\{E \in \mathcal{E} | X(E) \in A\} | B)$$

THEN:

$$F_X(a | B) = P(X < a | B)$$

AND

$$f_X(a | B) = \frac{d}{dx} F_X(x | B) \Big|_{x=a}$$

MULTI-DIMENSIONAL DISTRIBUTION FUNCTIONS : GIVEN N RANDOM VARIABLES $X = (X_1, X_2, \dots, X_N)$ WE HAVE

(AS ABOVE, ONE CAN EXTEND THE TARGET OF X FROM INTERVALS TO THE ENTIRE BOREL σ -ALGEBRA.)

$$X : M \rightarrow \mathbb{R}^N$$

AND A PROBABILITY FUNCTION P BOREL σ -ALGEBRA OF \mathbb{R}^N

$$P(X_1, X_2, \dots, X_N) : \sigma(\Omega_{\text{sto}}) \rightarrow \bar{\mathbb{R}}^+$$

DEFINED IN THE SAME MANNER AS $P(x)$. A DISTRIBUTION FUNCTION IS DEFINED

$$F_{X_1, X_2, \dots, X_N}(a_1, a_2, \dots, a_N) := P^{(x_1, \dots, x_N)}(L_{a_1, a_2, \dots, a_N})$$

WHERE $L_{a_1, a_2, \dots, a_N} = \{x \in \mathbb{R}^N \mid x_1 < a_1, x_2 < a_2, \dots\}$

ALTHOUGH F GIVES DIRECTLY ONLY THE VALUES $P^{(x)}$ ON THE "L" SETS, $P^{(x_1, \dots, x_N)}$ IS UNIQUELY DETERMINED FOR ALL $A \in \sigma(\Omega_{\text{sto}})$ BY KNOWLEDGE OF F_{X_1, X_2, \dots, X_N} .

AND ... IF THE DERIVATIVE EXISTS, WE HAVE THE N-DIMENSIONAL PROBABILITY DENSITY OF $X = (X_1, X_2, \dots, X_N)$:

OR, MORE PRECISELY:
 $f_{X_1, X_2, \dots, X_N}(a_1, \dots, a_N) = \frac{d^{(N)}}{dx_1 dx_2 \dots dx_N} F_{X_1, \dots, X_N}(x_1, \dots, x_N) \Big|_{\vec{a}}$

WHICH FOR CONTINUOUS F_X (SEE ABOVE) IMPLIES:

$$P^{(x)}(A) = \int \dots \int_A f(a_1, \dots, a_N) d^n a$$

FOR ALL BOREL SETS $A \subseteq \mathbb{R}^N$.

EXPECTATION OF DISCRETE RANDOM VARIABLES : CONSIDER A RANDOM VARIABLE $X : M \rightarrow \mathbb{R}$ THAT TAKES ONLY FINITELY MANY VALUES $\{a_n\}$ AND THIS PARTITIONS THE SPACE M INTO $\{A_n = \text{PREIM}_X(a_n)\}$. THE MATHEMATICAL EXPECTATION OF X IS:

$$E[X] = \sum_n a_n P(A_n) = \sum_n a_n P^{(x)}(a_n) =: \sum_n a_n p_X(a_n)$$

WHERE p_X IS THE PUSH-FORWARD OF P RESTRICTED TO THE SET $\{a_n\}$, AND IS CALLED THE PROBABILITY MASS FUNCTION.

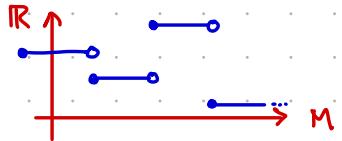
AKA, A DISCRETE Random Var.

RECALL: $P(\text{PREIM}_X(a_n))$

LEBESGUE INTEGRATION

DEFN: GIVEN THE MEASURABLE SPACES (M, σ) AND $(\mathbb{R}, \sigma(\text{Borel}))$, A MEASURABLE FUNCTION $S: M \rightarrow \mathbb{R}^+$ IS CALLED A SIMPLE FUNCTION (OR, STEP MAP) IF IT TAKES M INTO ONLY FINITELY-MANY VALUES, I.E.,

$$S(M) = \{s_1, s_2, \dots, s_N\}$$



THE FUNCTION THUS PARTITIONS M INTO THE FINITE PARTITION $\{A_n\}$ AND CAN THEREFORE BE EXPRESSED AS

$$S = \sum_{n=1}^N s_n \chi_{A_n}; \quad \chi_A: M \rightarrow \mathbb{R}, m \mapsto \begin{cases} 1 & \text{MEA} \\ 0 & \text{m} \notin A \end{cases}$$

WHERE χ_A IS AN INDICATOR FUNCTION AND $A_n = \text{PREIM}_S(s_n)$.

DEFN: LET (M, σ, μ) BE A MEASURE SPACE AND $(\bar{\mathbb{R}}, \bar{\sigma})$ A MEASURABLE SPACE. THE INTEGRAL OF A NON-NEGATIVE MEASURABLE MAP, $f: M \rightarrow \bar{\mathbb{R}}$ ($f(m) \geq 0 \forall m \in M$), IS:

INDICATOR Fcn
= CHARACTERISTIC Fcn

$$\int f d\mu := \sup \left\{ \sum_{z \in S(M)} z \cdot \mu(\text{PREIM}_S(z)) \mid \begin{array}{l} S \text{ IS A} \\ \text{SIMPLE Fcn} \\ \text{w/ } S \leq f \end{array} \right\}$$

HEIGHT WIDTH

FOR THE INTEGRAL OVER A SUBSET OF THE DOMAIN, $(A \subseteq M) \in \sigma$, WE DEFINE THE INTEGRAL OF A NON-NEGATIVE MEASURABLE MAP f OVER A AS:

$$\int_A f d\mu := \int f \cdot \chi_A d\mu$$

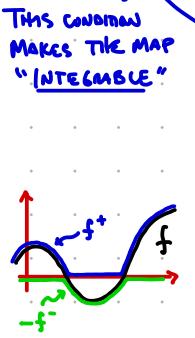
RECALL: PRODUCT OF MAPS
 $(f \cdot \chi)(m) := f(m) \cdot \chi(m)$

AN ALTERNATIVE NOTATION IS TO EXPRESS THE INTEGRAL IN TERMS OF THE VALUES OF THE FUNCTION (AKA, THE DUMMY VARIABLE(S) OF INTEGRATION). E.G., FOR $\text{mult}: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, WE WRITE:

$$\int_A \text{mult} d\mu = \int_A \text{mult}(x, y) d\mu(x, y) = \int_A x \cdot y d\mu(x, y) = \int_A x \cdot y \mu(dx, dy)$$

COMMON IN
PROBABILITY TH.

DEFN: THE INTEGRAL OF A MEASUREABLE MAP (NOT NECESSARILY POSITIVE), $f: M \rightarrow \bar{\mathbb{R}}$, FOR WHICH $\int |f| d\mu < \infty$, IS DEFINED AS



THIS CONDITION MAKES THE MAP "INTEGRABLE"

$$\int f d\mu := \int f^+ d\mu - \int f^- d\mu$$

WHERE

$$f^+ := \max(f, 0) \quad \text{AND} \quad f^- = \max(-f, 0)$$

CHANGE OF VARIABLE: LET $T: M \rightarrow M'$ BE A MEASUREABLE MAP FROM (M, \mathcal{F}, μ) TO (M', \mathcal{F}') . WE CAN DEFINE A MEASURE ON (M', \mathcal{F}') AS USUAL VIA THE PUSHFORWARD:

$$M'_T: \mathcal{F}' \rightarrow \bar{\mathbb{R}}_0^+,$$

$$A' \mapsto M'_T(A') := \mu(P_{\text{preim}_T}(A'))$$

IF $f: M' \rightarrow \mathbb{R}$ IS A MEASUREABLE REAL FUNCTION ON (M', \mathcal{F}') THEN $(f \circ T): M \rightarrow \mathbb{R}$ IS A MEASUREABLE REAL FUNCTION ON (M, \mathcal{F}) .

Theorem: IF f IS NONNEGATIVE, THEN

$$(*) \quad \int_M (f \circ T) d\mu = \int_{M'} f d\mu'_T$$

AND

$$(**) \quad \int_{\text{preim}_T(A')} (f \circ T) d\mu = \int_{A'} f d\mu'_T$$

FURTHER, A FUNCTION f (NOT NECESSARILY NONNEGATIVE) IS INTEGRABLE WITH RESPECT TO μ_T IF AND ONLY IF $(f \circ T)$ IS INTEGRABLE WITH RESPECT TO μ .

2022/03/11

From Billingsley's
PROBABILITY &
MEASURE (2nd Ed)
Section 16

PROOF: IF $f = \chi_{A'}$ (INDICATOR FUNCTION) THEN:

$$f \circ T = \chi_{A'} \circ T = \chi_{\text{PREIM}_+(A')}$$

SO EQN (*) BECOMES

$$\int_M \chi_{\text{PREIM}_+(A')} d\mu = \int_{M'} \chi_{A'} d\mu'$$

$$M(\text{PREIM}_+(A')) = M'_+(A')$$

INDICATOR FUNC IS A 2-VALUED SIMPLE FUNCTION w/ VALUE 1 OVER A' AND 0 ELSEWHERE. SO THE "SUP" IS CANCELLED TRIVIALLY.

WHICH IS SIMPLY THE DEFINITION OF M'_+ . BY THE LINEARITY OF THE INTEGRAL, EQN (*) THEN HOLDS FOR ALL NONNEGATIVE SIMPLE FUNCTIONS. AND, IF $\{f_n\}$ IS A SEQUENCE OF SIMPLE FUNCTIONS $f_n \rightarrow f$, THEN $(f_n \circ T) \rightarrow (f \circ T)$ AND (*) FOLLOWS BY THE MONOTONE CONVERGENCE THEOREM, APPLICATION OF (*) TO $|f|$ "ESTABLISHES THE ASSERTION ABOUT INTEGRABILITY", AND (*) FOLLOWS FOR GENERAL f BY DECOMPOSITION INTO f^+ & f^- . FINALLY, (***) FOLLOWS FROM (*) VIA $f \rightarrow f \cdot \chi_{A'}$ (MULTIPLY NOT COMPOSITION).

DEFN: GIVEN A MEASURABLE MAP $f: [a, b] \rightarrow \mathbb{R}$, AND A FUNCTION $g: [a, b] \rightarrow \mathbb{R}$ OF BOUNDED VARIATION, THE LEBESGUE-STIELTJES INTEGRAL IS DEFINED AS THE FOLLOWING LEBESGUE INTEGRAL:

$$\int_a^b f(x) dg(x) := \int_{[a, b]} f d\mu_g$$

WHERE μ_g IS THE UNIQUELY-DEFINED LEBESGUE-STIELTJES MEASURE ON $[a, b]$ THAT HAS VALUE ON ANY INTERVAL:

$$M((s, +]) = g(+)-g(s)$$

THE IMPORTANCE OF THIS REWRITING OF A LEBESGUE INTEGRAL IS THAT, FOR CONTINUOUS f , THE L.S. INTEGRAL IS RIEMANNIAN.

NOT SURE BUT
DON'T CARE
RIGHT NOW.

THERE IS A
GENERALIZATION
TO:
 $f, g: \mathbb{R}^d \rightarrow \mathbb{R}$

SEE:
BROWNSLEY
END OF SECTION 17

MATHEMATICAL EXPECTATION

- DEFN: GIVEN A RANDOM VARIABLE $X: M \rightarrow \mathbb{R}$ ON THE PROBABILITY SPACE (M, \mathcal{F}, P) , THE MATHEMATICAL EXPECTATION OF X IS THE LEBESGUE INTEGRAL:

$$E[X] := \int_M X \, dP = \int_M X(m) \, P(dm)$$

↑
JUST
NOTATION

NOTE THAT OUR PRIOR DEFINITION OF THE SAME QUANTITY — VALID FOR A DISCRETE, FINITELY-VALUED RANDOM VARIABLE — COINCIDES WITH THIS DEFINITION SINCE SUCH A RANDOM VARIABLE IS A SIMPLE FUNCTION.

- RECALL THAT WE DEFINED THE PROBABILITY FUNCTION OF X AS:

$$P(x): \mathcal{F}(\Omega_{\text{end}}) \rightarrow \overline{\mathbb{R}}^+, \quad n \mapsto P^{(x)}(n) := P(\text{preim}_X(n))$$

THIS IS PRECISELY THE MEASURE M'_T DEFINED IN THE CHANGE OF VARIABLE FORMULA, WITH $M = P$, $T = X$, AND $M'_T = P^{(x)}$. THEREFORE WE CAN WRITE

$$\begin{aligned} E[X] &= \int_M X \, dP \\ &= \int_M \text{Id}_{\mathbb{R}} \circ X \, dP \\ &= \int_{\mathbb{R}} \text{Id}_{\mathbb{R}} \, dP^{(x)} \\ &= \int_{\mathbb{R}} x \, dP^{(x)}(x) = \int_{\mathbb{R}} x \, P^{(x)}(dx) \end{aligned}$$

↘ $T = X$
 ↘ $f = \text{Id}_{\mathbb{R}}$ (IDENTITY
FUNCTION)
 ↘ CHANGE OF VARIABLES
THEOREM
 ↘ ALTERNATIVE
NOTATION

- USING THE LEBESGUE-STIELTJES NOTATION, WE CAN WRITE

Corresponds to L-S def'n
 $f = \text{Id}_{\mathbb{R}}$
 $g = F_X$
 $\mu_g = P^{(x)}$

$$\int \text{Id}_{\mathbb{R}} \, dP^{(x)} = : \int_{-\infty}^{\infty} x \, dF_X(x)$$

WHERE WE HAVE IDENTIFIED $P^{(x)}$ AS THE UNIQUELY-DEFINED MEASURE ASSOCIATED TO $F_x: \mathbb{R} \rightarrow \mathbb{R}$ THAT, FOR ANY INTERVAL $(a, b] \subseteq \mathbb{R}$, HAS VALUE:

$$M_{F_x}((a, b]) = F_x(b) - F_x(a) := P^{(x)}((a, b]),$$

* - IS IT? THE LAST EQUALITY IS (EQUIVALENT TO \star) HOW WE DEFINED THE DISTRIBUTION FUNCTION ORIGINALLY. FOR CONTINUOUS $F_x(x)$ THIS STIELTJES INTEGRAL IS RIEMANNIAN. AND, IF $F_x(x)$ IS DIFFERENTIABLE, THEN

$$\begin{aligned} \int_a^b x dF(x) &= \lim_{N \rightarrow \infty} \sum_{n=0}^{\infty} x_n \cdot (F(x_{n+1}) - F(x_n)) \\ &= \lim_{N \rightarrow \infty} \sum_{n=0}^{\infty} x_n \frac{F(x_n + \Delta x) - F(x_n)}{\Delta x} \Delta x \\ &= \int_a^b x f_x(x) dx \end{aligned}$$

RECALL:
 $f(x)|_a = \frac{dF(x)}{dx}|_a$

Thus, we can further write the mathematical expectation (when F_x and f_x are well-behaved) as:

$$E[X] = \int_{-\infty}^{\infty} x f_x(x) dx$$

WHERE f_x IS THE PROBABILITY DENSITY OF X .

❷ CONDITIONAL EXPECTATION OF A RV WITH RESPECT TO AN EVENT
 WE HAVE ALREADY SEEN THAT A CONDITIONAL PROBABILITY, $P(A|B)$, IS A PROBABILITY MEASURE FOR THE SPACE $(M, \mathcal{E}, P(\cdot|B))$. WE COULD WRITE THIS PROBABILITY USING KOLMOGOROV'S NOTATION

$$P_B(A) := P(A|B)$$

WE CAN THEN DEFINE THE CONDITIONAL EXPECTATION OF A RANDOM VARIABLE X WITH RESPECT TO THE

(SHORTHAND
RIEMANNIAN
INTEGRATION)

EVENT B:

$$E_B[X] := \int_M X dP_B \quad \left(P_B(\cdot) = P(\cdot | B) \right)$$

\downarrow

$$= \int_B X dP_B$$

\downarrow

$$E_B[X] = \frac{1}{P(B)} \int_B X dP$$

SINCE
 $P_B(A) = 0$ FOR
 $A \cap B = \emptyset$

SINCE
 $P_B(A) = \frac{P(A)}{P(B)}$
 WHEN $A \subset B$
 (SINCE $A \cap B = A$)

WE COULD CONSIDER TWO CONDITIONAL PROBABILITIES

$$P(A) E_A[X] + P(B) E_B[X] = \int_A X dP + \int_B X dP$$

$$= \int_{A+B} X dP$$

SUCH THAT

$$E_{A+B}[X] = \frac{P(A) E_A[X] + P(B) E_B[X]}{P(A+B)}$$

AND FOR $A+B=M$:

$$E[X] = P(A) E_A[X] + P(A^c) E_{A^c}[X]$$

More generally, for a partition $M = A_1 \cup A_2 \cup \dots \cup A_n$,
 $E[X] = \sum_{i=1}^n P(A_i) E_{A_i}[X]$

CONDITIONAL PROBABILITIES WRT TRIALS AND RANDOM VARIABLES

DEF'N: THE CONDITIONAL PROBABILITY OF AN EVENT, B , AFTER A TRIAL Ω IS A RANDOM VARIABLE

THAT, FOR EVERY EVENT $A_i \in \Omega$, TAKES THE VALUE $P(B | A_i)$. WE CAN WRITE THIS AS

These are Kolmogorov's words but they don't make sense b.c.
 $X: M \rightarrow \mathbb{R}$
 It's not known that $A_i \in M$

2022/03/13 - 20

math.stackexchange.com/
 questions/4402669

DEFN: For a trial A with countable - many outcomes $\{A_i\}$, we can associate an integer-valued random variable:

$$X_A : M \rightarrow N, m \mapsto i \text{ when } m \in A_i$$

so the preim $X_A^{-1}(i) = A_i$ ($i=1, 2, \dots, k$). We then define the conditional probability of the event B after the trial A as the random variable:

JUST NOTATION:
 $B \in \mathcal{E}$: EVENT
 A : Partition of M
 $[P(B|A)]$: Random Variable

$$[P(B|A)] : M \rightarrow \mathbb{R}$$

$$m \mapsto [P(B|A)](m) = P(B | A_{X_A(m)})$$

NOTATIONALLY, WE COULD WRITE

$$P(B | A = A_{X_A(\cdot)}) := [P(B|A)](\cdot)$$

OR, SINCE THE VALUES OF X_A ARE 1-1 WITH THE ELEMENTS OF THE PARTITION, WE COULD JUST USE:

$$P(B | X_A = X_A(\cdot)) := [P(B|A)](\cdot)$$

THE EXPECTATION OF THIS DISCRETE RANDOM VARIABLE IS:

$$\begin{aligned} E[[P(B|A)]] &= \sum_{i=1}^k P(B | A = A_i) P(A_i) \\ &= \sum_{i=1}^k P(B | A_i) P(A_i) \end{aligned}$$

AND THE CONDITIONAL EXPECTATION OF $[P(B|A)]$ WITH RESPECT TO (FIXED EVENT) C IS:

$$\begin{aligned} E_C[[P(B|A)]] &= \frac{1}{P(C)} \int_C [P(B|A)] dP \\ &= \frac{1}{P(C)} \sum_{\{i | A_i \subseteq C\}} P(B | A_i) P(A_i) \end{aligned}$$

Given multiple trials $A^{(1)}, A^{(2)}, \dots, A^{(n)}$, we can define

$$[P(B | A^{(1)} A^{(2)} \dots A^{(n)})] : M \rightarrow \mathbb{R}$$

as the random variable that takes the value:

$$P(B | A_{i_1}^{(1)} \cap A_{i_2}^{(2)} \cap \dots \cap A_{i_n}^{(n)})$$

When acting on an element, m , for which $m \in A_{i_1}^{(1)} \cap A_{i_2}^{(2)} \cap \dots \cap A_{i_n}^{(n)}$. The condition for independent trials can then be written as:

$$[P(A_g^{(k)} | A^{(1)} A^{(2)} \dots A^{(k-1)})] = P(A_g^{(k)})$$

for arbitrary k, g ,

④ Conditional Probability with respect to a Random Variable

A random variable $X : M \rightarrow \mathbb{R}$ can be thought of as partitioning M through its preimage, and thus X can be considered a trial. In the case of a discrete random variable, the prior discussions regarding trials and conditioning with respect to trials can be immediately applied to random variables

(For a discrete random variable)

$$[P(B | X)] : M \rightarrow \mathbb{R}$$

$$m \mapsto P(B | \text{Preim}_X(X(m)))$$

where we could use the notation above to write:

$$P(B | X = X(m)) := [P(B | X)](m)$$

For general random variables, which can take an uncountable number of outcomes, the set

$$"X = X(m)" = \text{Preim}_X(X(m))$$

GIVEN TWO PARTITIONS OF M , $P^{(1)}$ & $P^{(2)}$, LET THEIR PRODUCT (OR MEET), $P^{(1)}P^{(2)}$, BE THE NEW PARTITION THAT IS THE SET OF ALL NONEMPTY INTERSECTIONS:

$$\{P_{i_1}^{(1)} \cap P_{i_2}^{(2)}\} \quad (P_{i_k}^{(1)} \in P^{(1)})$$

$$\begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \\ \text{D} \end{array} \quad \begin{array}{c} \text{E} \\ \text{F} \\ \text{G} \\ \text{H} \end{array} = \begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \\ \text{D} \\ \text{E} \\ \text{F} \\ \text{G} \\ \text{H} \end{array}$$

(6 of the 12 intersections are \emptyset)

SEE THEOREM NEAR
END OF "INDEPENDENT
TRAILS & EVENTS"

KEY POINT / MOTIVATION OF THIS ENTIRE ANALYSIS

OFTEN HAS PROBABILITY ZERO. THEREFORE THE ABOVE RESULT CANNOT BE USED SINCE $P(B|A) = P(B \cap A)/P(A)$ ($P(A) > 0$). WE MUST THEREFORE DEFINE THE CONDITIONAL PROBABILITY WITH RESPECT TO A RANDOM VARIABLE INDIRECTLY.

DEF'N: For a generic random variable, $X: M \rightarrow \mathbb{R}$, from (M, \mathcal{E}, P) to $(\mathbb{R}, \sigma(\Omega_{std}), P^{(x)})$, the conditional probability of an event $B \in \mathcal{E}$ with respect to X is the unique random variable $[P(B|X)]$ satisfying:

$$P(B | \text{Preim}_X(u)) = E_{\text{Preim}_X(u)} [P(B|X)]$$

For all $u \in \sigma(\Omega_{std})$ with $P^{(x)}(u) > 0$.

Proof: We must show: ① that this agrees with the above definition for discrete random variables and, ② that it uniquely* defines $[P(B|X)]$ for generic random variables.

① For a discrete X with target values

$A = \{A_i\}$ $\{x_1, x_2, \dots, x_n\}$, let $A_i = \text{Preim}_X(x_i)$, and let $A_u = \text{Preim}_X(u)$ for $u \in \{x_1, \dots, x_n\}$:

CONDITIONAL EXPECTATION WAS DEFINED AS A LEBESGUE INTEGRAL, BUT THAT WORKS FOR A DISCRETE RV (I.E., A STEP FN)

$$E_{A_u} [P(B|X)] := \frac{1}{P(A_u)} \sum_{\{i | A_i \subseteq A_u\}} [P(B|X)]_i P(A_i)$$

If we let $[P(B|X)]_i = P(B|X=x_i) = P(B|A_i)$, as defined above for the discrete case, then:

$$\begin{aligned} \frac{1}{P(A_u)} \sum_{\{i | A_i \subseteq A_u\}} P(B|A_i) P(A_i) &= \frac{1}{P(A_u)} \sum_{\{i | A_i \subseteq A_u\}} P(B \cap A_i) \\ &= \frac{1}{P(A_u)} P\left(B \cap \bigcup_{\{i | A_i \subseteq A_u\}} A_i\right) = \frac{P(B \cap A_u)}{P(A_u)} = P(B|A_u) \end{aligned}$$

WHERE WE HAVE USED THE FACT THAT SOME SUBSET OF THE $A = \{A_i\}$ PARTITIONS A_u SINCE $A_u \subseteq A$.

② WE WILL PROVE UNIQUENESS. For the existence proof, SEE CHAPTER V IN KOLMOGOROV (USING NIKODYM'S THM.).

STARTING FROM OUR DEFINING EQUATION FOR $[P(B|X)]$:

$$P(B | \text{Preim}_X(u)) = E_{\text{Preim}_X(u)} [P(B|X)]$$

MULTIPLY BOTH SIDES BY $P(\text{Preim}_X(u))$. FIRST THE LHS:

$$P(\text{Preim}_X(u)) \cdot P(B | \text{Preim}_X(u)) = P(B \cap \text{Preim}_X(u))$$

AND THEN THE RHS:

$$\begin{aligned} & P(\text{Preim}_X(u)) \cdot \left[\frac{1}{P(\text{Preim}_X(u))} \int_{\text{Preim}_X(u)} [P(B|X)] dP \right] \\ &= \int_{\text{Preim}_X(u)} [P(B|X)] dP \end{aligned}$$

RECALL: For $T: M \rightarrow M'$

$$\int_M (f \circ T) d\mu = \int_{M'} f d\mu_T$$

WHERE

$$\mu_T(m) = \mu(\text{Preim}_T(m))$$

NOW WE WISH TO USE THE CHANGE OF VARIABLES THEOREM USING $T = X: M \rightarrow \mathbb{R}$ AND DEFINING $f: \mathbb{R} \rightarrow \mathbb{R}$ VIA:

$$f \circ X: M \rightarrow \mathbb{R}$$

$$m \mapsto (f \circ X)(m) = [P(B|X)](m)$$

Thus, the function f maps the target of the random variable X into the target of the random variable $[P(B|X)]$ for the same $m \in M$:

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

$$x \mapsto f(x) = [P(B|X)](\text{Preim}_X(x))$$

NOTATIONALLY, IT MAKES SENSE TO WRITE THIS AS

$$P(B|X=\cdot) := [P(B|X)](\text{Preim}_X(\cdot))$$

SO, WE FIND:

$$\int_{\text{Preim}_X(u)} [P(B|X)] dP = \int_{\text{Preim}_X(u)} P(B|X=\cdot) \circ X dP$$

$$= \int_u P(B|X=\cdot) dP^{(X)}$$

(NOW AN INTEGRAL
OVER $u \subseteq \mathbb{R}$)

AND SETTING LHS EQUAL TO RHS:

$$P(B \cap \text{Preim}_X(u)) = \int_u P(B|X=\cdot) dP^{(x)}$$

ONE OF THE FUNDAMENTAL PROPERTIES OF THE LEbesgue INTEGRAL (PROPERTY IX, SECTION 1, CHAPTER IV) IS THAT IF, FOR A PROBABILITY SPACE (M, \mathcal{E}, P) ,

$$\int_A X dP = \int_A Y dP \quad \forall A \in \mathcal{E}$$

THEN X AND Y ARE EQUIVALENT RANDOM VARIABLES ($P(\{x \neq y\}) = 0$). THEREFORE THE RHS ABOVE UNIQUELY SPECIFIES $P(B|X=\cdot)$ ON $(\mathbb{R}, \mathcal{O}(\mathbb{R}), P^{(x)})$.

■ NOTATION: First is KOLMOGOROV. Second is MINE.

$$P_X(a; B) = P(B|X=a) = [P(B|X)](\text{Preim}_X(a))$$

■ TWO RANDOM VARIABLES X AND Y INDUCE THE SAME PARTITIONING OF M IF THERE EXISTS A BIJECTIVE MAP FROM THE TARGET OF X , $U \subseteq \mathbb{R}$, TO THE TARGET OF Y , $V \subseteq \mathbb{R}$, $g: U \rightarrow V$, SUCH THAT

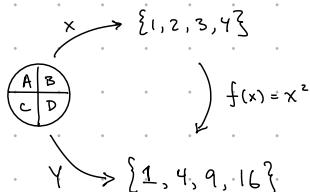
$$g \circ X : M \rightarrow V$$

$$m \mapsto (g \circ X)(m) = Y(m)$$

I.E., $Y = g \circ X$. IN THAT CASE, THE RANDOM VARIABLES FOR CONDITIONAL PROBABILITY, AS DEFINED ABOVE ARE THE SAME*!

$$[P(B|X)] = [P(B|Y)]$$

THUS WE SAY THAT $[P(B|X)]$ IS DEFINED UNIQUELY UP TO EQUIVALENCE.



* I REALLY DON'T SEE THE EQUALITY OF $[P(B|X)]$ AND $[P(B|Y)]$.

CONSTRUCTION OF CONDITIONAL PROBABILITY RV : ALTHOUGH WE HAVE SHOWN IT EXISTS AND IS UNIQUE (UP TO SOME EQUIVALENCE), WE HAVE NOT YET WRITTEN DOWN AN EXPRESSION FOR $[P(B|X)]$. WE DID TRANSFORM THE DEFINING EQUATION FROM

$$P(B | \text{Preim}_X(u)) = E_{\text{Preim}_X(u)} [P(B|X)]$$

(FOR ALL $u \in \sigma(\Theta_{\text{std}})$ FOR WHICH $P(\text{Preim}_X(u)) > 0$)
TO:

$$P(B \cap \text{Preim}_X(u)) = \int_u P(B|X=\cdot) dP^{(x)}$$

(FOR ALL $u \in \sigma(\Theta_{\text{std}})$ w/ $P(\text{Preim}_X(u)) > 0$), WHERE:
WE DEFINED $P(B|X=\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ VIA:

$$[P(B|X)] = P(B|X=\cdot) \circ X$$

Now, WE CAN WRITE THIS INTEGRAL IN TERMS OF THE DISTRIBUTION F_X USING THE LEbesgue-STIELTJES Formalism, WITH $g = F_X$ AND THEREFORE $\mu_g = P^{(x)}$.
RECALL:

$$\int_a^b f(x) dg(x) := \int_{(-\infty, b]} f(x) d\mu_g$$

$$\mu_g((s, t]) = g(t) - g(s)$$

THUS, LETTING $u = (-\infty, a]$, WE HAVE:

$$P(B \cap \text{Preim}_X((-\infty, a])) = \int_{-\infty}^a P(B|X=x) dF_X(x)$$

KOLMOGOROV NOTES (SECTION 3, CHAPTER IV) THAT A THEOREM OF LEbesgue PROVIDES A SOLUTION TO THE INTEGRAND:

$$P(B|X=a) = P(B) \left[\lim_{h \rightarrow 0} \frac{F_X(a+h|B) - F_X(a|B)}{F_X(a+h) - F_X(a)} \right]$$

WHERE WE DEFINE THE CONDITIONAL DISTRIBUTION, $F_X(x|B)$,

VIA:

This is where Kolmogorov's notation
 $P_B(A) = P(A|B)$
 MAKES MORE SENSE.
 $P_B^{(x)}$ IS A PROBABILITY MEASURE ON $(\mathbb{R}, \mathcal{F}(\mathbb{R}))$
 THAT HAPPENS TO BE CONSTRUCTED FROM A CONDITIONAL PROBABILITY P_B ON (M, \mathcal{E}) .
 MY NOTATION CARRIES WITH IT SOME INTUITION THAT THE A & B IN $P(A|B)$ ARE EVENTS IN THE SAME PROBABILITY SPACE BUT THEY NEED NOT BE.

$$P^{(x)}(u) := P(P_{\text{REM}_X}(u))$$

$$\Rightarrow P^{(x)}(u | B) := P(P_{\text{REM}_X}(u) | B)$$

$$\Rightarrow F_x(x | B) := P^{(x)}((-\infty, x) | B)$$

$$= P(P_{\text{REM}_X}(-\infty, x) | B)$$

Now, IF THE DERIVATIVES OF $F_x(x)$ AND $F_x(x | B)$ EXIST, AND IF $f_x(a) = \frac{dF_x(x)}{dx} \Big|_{x=a} > 0$, WE CAN FURTHER WRITE:

$$P(B | X = a) = P(B) \frac{f_x(a | B)}{f_x(a)}$$

WHERE

$$f_x(a | B) = \frac{dF_x(x | B)}{dx} \Big|_{x=a}$$

Thus, we reduced the question of the random variable $[P(B | X)]$ to the function $P(B | X = \cdot) : \mathbb{R} \rightarrow \mathbb{R}$ and have now defined it in terms of densities of the original random variable X .

- ② MARGINALIZING OVER CONDITIONAL PROBABILITIES: RETURNING TO OUR LAST INTEGRAL DERIVED FROM THE DEFINING EQUATION, BUT LETTING $u \rightarrow \mathbb{R}$, WE FIND

$$P(B) = \int_{-\infty}^{+\infty} P(B | X = x) dF_x(x)$$

$$P(B) = \int_{-\infty}^{\infty} P(B | X = x) f_x(x) dx$$

) DEFN ABOVE

$$P(B) = \int_{-\infty}^{\infty} P(B) f_x(x | B) dx \Rightarrow \int f_x(x | B) dx = 1$$

WE CAN REARRANGE THE EQUATION WE DERIVED FOR $P(B|X=.)$ AND UTILIZE THE ABOVE MARGINALIZATION TO WRITE:

$$f_x(a|B) \cdot P(B) = P(B|X=a) \cdot f_x(a)$$

$$f_x(a|B) = \frac{P(B|X=a) \cdot f_x(a)}{P(B)}$$

This seems to
be a tautology
based on the
definition of
 $P(B|X=x)$...
WHAT IS GAINED HERE?

$$f_x(a|B) = \frac{P(B|X=a) \cdot f_x(a)}{\int_{-\infty}^{\infty} P(B|X=x) f_x(x) dx}$$

WHICH KOLMOGOROV CALLS BAYES THEOREM FOR CONTINUOUS DISTRIBUTIONS. I GUESS THE IDEA HERE IS: IF WE CAN DEDUCE $P(B|X=.)$ OTHERWISE (OR MEASURE IT IN A TRIAL/EXPERIMENT) THEN ITS RELATIONSHIP TO THESE OTHER PROBABILITY DENSITIES IS THIS.

BAYES THEOREM FOR PROBABILITY DENSITIES

2022/03/23

IN THE PREVIOUS SECTION, WE MENTIONED IN PASSING THE CONDITIONAL PROBABILITY DENSITY BUT LET US SPECIFY IT EXPLICITLY NOW:

DEF'N: THE PROBABILITY DENSITY OF A RANDOM VARIABLE, X , CONDITIONED ON THE EVENT B , FOR WHICH $P(B) > 0$, IS:

$$f_x(x|B) = \frac{d}{dx} [F_x(x|B)]$$

$$= \frac{d}{dx} [P^{(x)}(-\infty, x|B)]$$

$$= \frac{d}{dx} [P(\text{PREIM}_X(-\infty, x)|B)]$$

$$= \frac{d}{dx} [P(\text{PREIM}_X(-\infty, x) \cap B) / P(B)]$$

WE ARE ACTUALLY DEFINING BOTH
 $f_x(x|B)$

AND
 $F_x(x|B)$

HERE... BUT THE DEFINITIONS ARE COMPLETELY NATURAL

AND IN KOLMOGOROV'S NOTATION
 $f_B^{(x)}(x) \leftarrow F_B^{(x)}(x) \leftarrow P_B^{(x)}$
NOT REALLY EVEN NEW QUANTITIES.

WE ARE INTERESTED IN THE PROBABILITY DENSITY OF ONE RANDOM VARIABLE CONDITIONED ON THE VALUE OF A SECOND RANDOM VARIABLE, I.E., CONDITIONED ON THE EVENT THAT THE 2ND RANDOM VARIABLE TAKE SOME VALUE!

$$f_X(x \mid \{Y=y\})$$

THIS IS PERFECTLY WELL DEFINED FOR A DISCRETE RANDOM VARIABLE Y , SO LONG AS $P(\{Y=y\}) > 0$.

BUT FOR A CONTINUOUS RANDOM VARIABLE IT IS OFTEN THE CASE THAT $P(\{Y=y\}) = 0 \forall y$! Thus, we CANNOT DIRECTLY APPLY THE DEFINITION GIVEN ABOVE. HOWEVER, GIVEN THE RESULTS OF THE PREVIOUS SECTION WE CAN MAKE THE FOLLOWING DEFINITION.

DEF'N: THE CONDITIONAL PROBABILITY DENSITY OF A RANDOM VARIABLE X WITH RESPECT TO THE EVENT THAT A SECOND RANDOM VARIABLE, Y , HAS TAKEN THE VALUE y IS:

$$f_X(x \mid \{Y=y\}) := \frac{d}{dx} \left[P(\text{PREIM}_X(-\infty, x) \mid \{Y=y\}) \right]$$

WHERE WE USE THE RANDOM VARIABLE NOTATION $[P(B \mid \{Y=y\})]$ OF THE PREVIOUS SECTION.

GIVEN THE DERIVED EXPRESSION IN THE PREVIOUS SECTION, WE CAN FURTHER CALCULATE:

$$\begin{aligned} f_X(x \mid \{Y=y\}) &= \frac{d}{dx} \left[P(\text{PREIM}_X(-\infty, x)) \right. \\ &\quad \times \left. \frac{f_Y(y \mid \text{PREIM}_X(-\infty, x))}{f_Y(y)} \right] \end{aligned}$$

$$= \frac{\frac{d}{dx} \left[P(\text{PREIM}_X(-\infty, x)) \cdot \frac{d}{dy} \left[P(\text{PREIM}_Y(-\infty, y) \mid \text{PREIM}_X(-\infty, x)) \right] \right]}{f_Y(y)}$$

NOTE: THE COMMONLY-USED NOTATION:

$f_{XY}(x|y)$
IS BAD! WE ARE
TALKING ABOUT THE
DENSITY OF X WITH
RESPECT TO A FIXED
VALUE OF Y .

$$\begin{aligned}
 &= \frac{\frac{\partial^2}{\partial x \partial y} \left[P(\text{PrElm}_X(-\infty, x)) \cdot P(\text{PrElm}_Y(-\infty, y) \mid \text{PrElm}_X(-\infty, x)) \right]}{f_Y(y)} \\
 &= \frac{\frac{\partial^2}{\partial x \partial y} \left[P(\text{PrElm}_X(-\infty, x) \cap \text{PrElm}_Y(-\infty, y)) \right]}{f_Y(y)}
 \end{aligned}$$

This density is commonly written in inference applications as:

" $f(x|y)$ "

which is terrible because it is missing so much information:

- conditioning w/r/t an event $\{Y=y\}$
- distribution of the RV "X", $f_X(x|\dots)$
- So X is a domain var.; y is fixed

$$f_X(x \mid \{Y=y\}) = \frac{f_{XY}(x, y)}{f_Y(y)}$$

② BAYES THEOREM FOR DENSITIES: we can similarly find:

$$f_Y(y \mid \{X=x\}) = \frac{f_{XY}(x, y)}{f_X(x)}$$

Therefore:

$$f_X(x \mid \{Y=y\}) \cdot f_Y(y) = f_Y(y \mid \{X=x\}) \cdot f_X(x)$$

OR:

$$f_X(x \mid \{Y=y\}) = \frac{f_Y(y \mid \{X=x\}) \cdot f_X(x)}{f_Y(y)}$$

OR, ALTERNATIVELY:

$$f_X(x \mid \{Y=y\}) = \frac{f_Y(y \mid \{X=x\}) \cdot f_X(x)}{\int_{-\infty}^{\infty} f_Y(y \mid \{X=x'\}) \cdot f_X(x') dx'}$$

$\underbrace{f_{XY}(x, y)}$

BAYES'
THEOREM
FOR
DENSITIES

OVERVIEW OF BAYESIAN INFERENCE

2022/03/25

[Following D. Mackay's LECTURES
9-10, Inf.Th., Prob., Rec., Neural]

➊ Toy Problem: WE ARE GIVEN A FEW REAL NUMBERS $\{x_i\}_{i=1}^N$ AND WE ARE TOLD (OR WE ASSUME, AS A HYPOTHESIS MODEL, H_N) THAT THEY ARE SAMPLES FROM A NORMAL DISTRIBUTION. WHAT CAN WE INFER ABOUT THAT DISTRIBUTION FROM THE DATA?

- ➋ UNDER THE HYPOTHESIS, H_N , OUR UNDERLYING PROBABILITY SPACE, M , IS THE SET OF ALL POSSIBLE N-ELEMENT DATA SETS GENERATED FROM ALL POSSIBLE NORMAL DISTRIBUTIONS (SPECIFIED BY THE MEAN AND STANDARD DEVIATION, (μ, σ)). ONE COULD IMAGINE EXPANDING THIS SPACE BY REMOVING THE NORMAL DISTRIBUTION HYPOTHESIS, H_N , AND CONSIDERING, E.G., A MULTITUDE OF OTHER DATA-GENERATING DISTRIBUTIONS. IN THAT MANNER ONE COULD DO MODEL-HYPOTHESIS COMPARISON. WE COULD ALSO EXPAND TO DATA SETS OF ANY SIZE.
- ➋ THE PARTICULAR DATA SET THAT WE HAVE OBSERVED COULD BE WRITTEN AS THE EVENT, D . BECAUSE THERE ARE UNCOUNTABLY-MANY POSSIBILITIES, THIS IS A PROBABILITY-ZERO EVENT AND SHOULD BE HANDLED WITH A DENSITY, I.E., DEFINE A VECTOR OF RANDOM VARIABLES $\vec{D} = (D_1, D_2, \dots, D_N)$ WITH JOINT DENSITY $f_D(\vec{d})$, AND WRITE THE EVENT AS $D = \{\vec{D} = \vec{d}\}$ WHERE $\vec{d} = \{x_i\}$ IS OUR OBSERVED DATA SET.
- ➋ THE PARAMETERS OF THE NORMAL DISTRIBUTION ARE ALSO RANDOM VARIABLES, WITH JOINT DISTRIBUTION $f_{M\Sigma}(\mu, \sigma)$. ANOTHER PART OF OUR HYPOTHESIS IS ITS SPECIFICATION (THE PRIOR) ON M .

INTERESTING: AS SOON AS YOU CALL THE PARAMETERS RANDOM VARIABLES IT REQUIRES YOU TO DEFINE THE PRIOR, I.E., THE JOINT DISTRIBUTION.

ANSWER : THE BAYESIAN ANSWER TO THE QUESTION (WHAT CAN WE INFERENCE ABOUT $N(\mu, \sigma)$ FROM THE DATA?) IS THE ("POSTERIOR") JOINT PROBABILITY DENSITY OF (μ, Σ) CONDITIONED ON THE DATA EVENT, $D = \{\vec{D} = \vec{d}\}$, I.E.,

$$f_{M\Sigma}(\mu, \sigma \mid \{\vec{D} = \vec{d}\}) = \frac{f_{\vec{D}}(\vec{d} \mid \{(\mu, \Sigma) = (\mu, \sigma)\}) \cdot f_{M\Sigma}(\mu, \sigma)}{f_{\vec{D}}(\vec{d})}$$

OR, THE
MARGINAL
LIKELIHOOD OF H

WHERE THE MARGINAL DENSITY OF THE DATA — WHICH INCORPORATES OUR HYPOTHESIS OF THE NORMAL DISTRIBUTION, AND OUR PRIOR DISTRIBUTION OF (μ, Σ) , AND IS THEREFORE ALSO CALLED THE EVIDENCE FOR THE MODEL GIVEN THE DATA — IS WRITTEN:

$$f_{\vec{D}}(\vec{d}) = \iint f_{\vec{D}}(\vec{d} \mid \{(\mu, \Sigma) = (\mu', \sigma')\}) f_{M\Sigma}(\mu', \sigma') d\mu' d\sigma'$$

AND THE CONDITIONAL DENSITY OF \vec{D} — CONDITIONED ON THE VALUES OF (μ, Σ) , AND EVALUATED AT THE OBSERVED DATA SET, \vec{d} — BUT VIEWED AS A FUNCTION OF THE PARAMETERS (μ, σ) , IS CALLED THE LIKELIHOOD FUNCTION (OF THE PARAMETERS, GIVEN THE OBSERVED DATA):

$$L(\mu, \sigma; \vec{d}) = f_{\vec{D}}(\vec{d} \mid \{(\mu, \Sigma) = (\mu, \sigma)\})$$

IN THE BOXED SOLUTION ABOVE, WE ARE ASSUMING A FIXED VALUE OF \vec{d} (THE OBSERVED DATA) BUT THE POSTERIOR DENSITY OF (μ, Σ) IS A FUNCTION OF (μ, σ) , WHOSE VALUE FOR A PARTICULAR (μ, σ) IS FOUND BY EVALUATING THE RHS AT (μ, σ) .

(PRIMES TO INDICATE THESE ARE (UMMY) INTEGRATION VARIABLES NOT THE (μ, σ) IN THE DOMAIN OF:
 $f_{M\Sigma}(\mu, \sigma)$)

GENERAL FORMALISM: TO MAKE EXPLICIT THAT WE ALWAYS HAVE SOME ASSUMPTIONS ABOUT THE SYSTEM IN QUESTION — E.G., DATA FROM A CERTAIN TYPE OF DISTRIBUTION, NUMBER OF DATA POINTS — WE WILL CONDITION DISTRIBUTIONS ON A HYPOTHESIS-MODEL EVENT, H . WE WILL WRITE THE SET OF ALL PARAMETERS AS A VECTOR OF RANDOM VARIABLES, $\vec{\theta}$, WITH (PRIOR) JOINT DISTRIBUTION $f_{\vec{\theta}}(\vec{\theta})$. THE OBSERVED DATA, \vec{d} , WILL BE ASSUMED TO BE A SAMPLE FROM THE JOINT DISTRIBUTION OF A RANDOM VARIABLE VECTOR \vec{D} . THUS BAYES THEOREM IS WRITTEN:

$$f_{\vec{\theta}}(\vec{\theta} \mid \{\vec{D} = \vec{d}\} \cap H) = \frac{f_{\vec{D}}(\vec{d} \mid \{\vec{\theta} = \vec{\theta}\} \cap H) \cdot f_{\vec{\theta}}(\vec{\theta} \mid H)}{\int f_{\vec{D}}(\vec{d} \mid \{\vec{\theta}' = \vec{\theta}'\} \cap H) \cdot f_{\vec{\theta}}(\vec{\theta}' \mid H) d\vec{\theta}'}$$

↑
Posterior DENSITY

↑
Likelihood Function, $L(\vec{\theta})$

↑
Prior DENSITY

↑
 $f_{\vec{D}}(\vec{d} \mid H)$
Marginal distribution of the data ... more common: Marginal Likelihood of (or Evidence for) H

WE DON'T USUALLY MAKE ASSUMPTIONS ABOUT THE FORM OF $f_{\vec{\theta}}(\vec{\theta} \mid \dots)$, BUT:

- CHOOSE A BROAD ENOUGH PRIOR $f_{\vec{\theta}}(\vec{\theta} \mid H)$ TO ALLOW FOR ALL VARIANCES OF INTEREST, E.G., UNIFORM DISTRIBUTION ON $U \subseteq \mathbb{R}^d$
- THE ASSUMPTIONS ABOUT THE DISTRIBUTION OF \vec{D} (THE LIKELIHOOD FUNCTION) WILL SHAPE $f_{\vec{\theta}}(\dots)$

EXAMPLE: CONSIDER 6 DATA POINTS DRAWN FROM A NORMAL DISTRIBUTION.

[I DREW THESE FROM $N(\mu, \Sigma)$]

$$\vec{d} = [3.81, 6.96, 7.02, -0.86, 5.15, 3.87]$$

WHAT CAN WE DETERMINE ABOUT THAT DISTRIBUTION? WE WILL ASSUME THAT THE ELEMENTS OF \vec{D} ARE INDEPENDENT

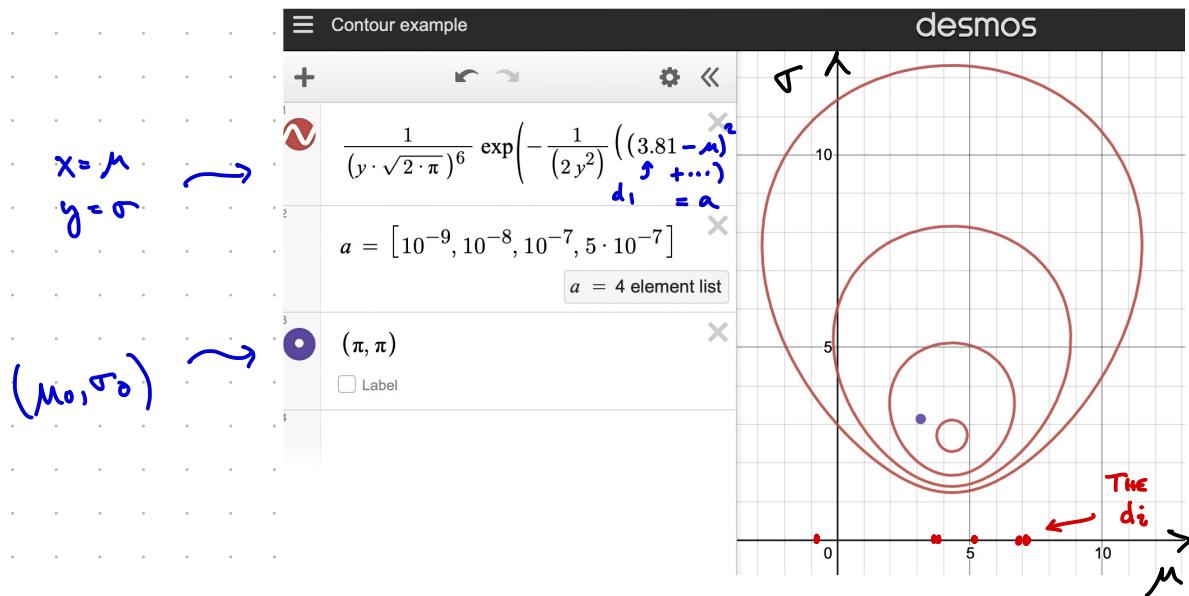
AND IDENTICALLY DISTRIBUTED RANDOM VARIABLES, SO:

$$f_{\vec{D}}(\vec{d}) = f_{D_1}(d_1) \cdot f_{D_2}(d_2) \cdot \dots \cdot f_{D_6}(d_6)$$

WITH $f_{D_i} \sim N(\mu, \sigma^2)$ WHERE (THE DISTRIBUTIONS OF) μ AND σ ARE TO BE INFERRED FROM THE DATA. THE LIKELIHOOD FUNCTION IS THEREFORE:

$$\begin{aligned} L(\mu, \sigma) &= f_{\vec{D}}(\vec{d} \mid \{\mu, \sigma\} = (\mu, \sigma)) \sim H \\ &= \prod_{i=1}^6 \frac{1}{\sigma \sqrt{2\pi}} \exp \left[-\frac{(d_i - \mu)^2}{2\sigma^2} \right] \\ &= \left(\frac{1}{\sigma \sqrt{2\pi}} \right)^6 \exp \left[-\sum_{i=1}^6 \frac{(d_i - \mu)^2}{2\sigma^2} \right] \end{aligned}$$

THE CONTOUR PLOT OF $L(\mu, \sigma)$ LOOKS LIKE THIS:



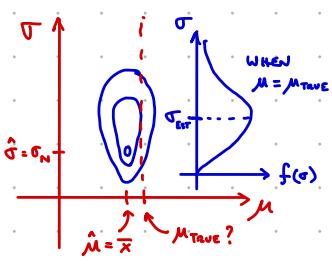
WHERE WE SEE THAT THE ACTUAL $(\mu_0, \sigma_0) = (\pi, \pi)$ WITH WHICH I GENERATED THE 6 DATA POINTS LIES VERY CLOSE TO THE PEAK OF THE LIKELIHOOD. IF WE MULTIPLY THE LIKELIHOOD BY, E.G., A UNIFORM DIST OVER $-10 \leq \mu \leq 20$ AND $0 \leq \sigma \leq 20$, THE QUALITATIVE PICTURE DOES NOT CHANGE. THEN DIVISION BY THE CONSTANT $f_{\vec{D}}(\vec{d})$ WOULD ALSO NOT CHANGE THE PICTURE. SO THE LIKELIHOOD BASICALLY TELLS THE WHOLE STORY!

THE PEAK OF THE LIKELIHOOD FUNCTION WILL GIVE YOU (OBVIOUSLY) THE VALUES FOUND BY MAXIMUM LIKELIHOOD. THE VALUES AT THE PEAK ARE

$$\hat{\mu} = \bar{x} := \frac{1}{N} \sum_{i=1}^N d_i$$

$$\hat{\sigma}^2 = \sigma_N^2 := \frac{1}{N} \sum_{i=1}^N (d_i - \bar{x})^2$$

THE SHAPE OF THE LIKELIHOOD FUNCTION IS NOT SYMMETRIC FOR SECTIONS OF CONSTANT Σ . IT IS EASY TO SEE WHY: IF THE ACTUAL DISTRIBUTION MEAN IS OFF SET FROM $\hat{\mu}$ THEN THE OBSERVED DATA ARE FURTHER FROM μ_{true} , ON AVERAGE, THAN FROM $\hat{\mu}$, AND THEM IMPLY A LARGER σ^2 .



ONE COULD ASK EITHER:

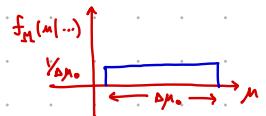
(i) I KNOW WHAT σ IS EXACTLY. BUT WHAT DOES THE DATA TELL ME ABOUT μ ?

(ii) I DON'T KNOW μ AND DON'T CARE... TELL ME ABOUT THE NOISE. WHAT DOES THE DATA SAY ABOUT Σ ?

THE FIRST QUESTION IS ANSWERED BY:

$\{\Sigma = \sigma^2\}$

$$f_M(\mu | D \cap \{\Sigma = \sigma^2\} \cap H) = \frac{f_D(\bar{d} | \{(M, \Sigma) = (\mu, \sigma^2)\} \cap H) \cdot f_M(\mu | \{\Sigma = \sigma^2\} \cap H)}{f_D(\bar{d} | \{\Sigma = \sigma^2\} \cap H)}$$



CHOOSING SOME BROAD PRIOR $f_M(\mu | \dots)$, AND CALCULATING THE NORMALIZING CONSTANT:

$$f_D(\bar{d} | \{\Sigma = \sigma^2\} \cap H) = \int f_D(\bar{d} | \{(M, \Sigma) = (\mu, \sigma^2)\} \cap H) \times f_M(\mu | \{\Sigma = \sigma^2\} \cap H) d\mu$$

AND THE SECOND QUESTION IS ANSWERED BY:

$$f_{\Sigma}(\sigma | D \cap H) = \frac{f_{\bar{D}}(\bar{d} | \{\Sigma = \sigma\} \cap H) \cdot f_{\Sigma}(\sigma | H)}{f_{\bar{D}}(\bar{d} | H)}$$

NOTICE IN THIS SECOND QUESTION THAT THE LIKELIHOOD IS EXACTLY THE NORMALIZING CONSTANT OF THE FIRST QUESTION (THAT "WE DON'T CARE ABOUT") IF YOU DO THE MARGINALIZATION INTEGRATION OVER μ TO OBTAIN $f_{\bar{D}}(\bar{d} | \{\Sigma = \sigma\} \cap H)$ THEN THE PEAK OF THIS MARGINAL LIKELIHOOD FOR σ IS:

$$\max(f(\sigma | D \cap H)) = \sigma_{N-1} := \sqrt{\frac{\sum_{i=1}^N (d_i - \mu)^2}{N-1}}$$

LASTLY: THE NORMALIZING CONSTANT IN OUR SECOND QUESTION IS CALLED THE MARGINAL LIKELIHOOD OF H (THE HYPOTHESIS-MODEL):

$$f_{\bar{D}}(\bar{d} | H) := \iint \left[f_D(\bar{d} | \{M, \Sigma = \mu, \sigma\} \cap H) \times f_M(\mu | \{\Sigma = \sigma\} \cap H) \cdot f_{\Sigma}(\sigma | H) \right] d\mu d\sigma$$

OR: $f_{M\Sigma}(\mu, \sigma | H)$

IT IS ALSO SOMETIMES CALLED THE "EVIDENCE FOR H" (GIVEN THE OBSERVED DATA).

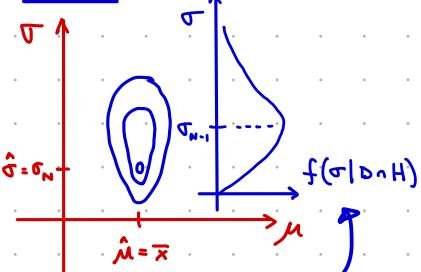
NON STANDARD

EXAMPLE: A SOURCE EMITS PARTICLES THAT DECAY AT POSITIONS THAT ARE EXPONENTIALLY DISTRIBUTED WITH LENGTH SCALE λ BUT ARE ONLY OBSERVED IN THE WINDOW:

$$(x_{\min}, x_{\max}) = (a, b) = (1, 20)$$

GIVEN N OBSERVATIONS $\{x_1, x_2, \dots, x_N\}$ WHAT IS λ ?

MARGINALIZED DISTRIBUTION OF σ :



IT IS THE MARGINAL DISTRIBUTION AFTER MULTIPLYING BY PRIOR AND NORMALIZING. BECAUSE THAT IT IS THE MARGINAL LIKELIHOOD.

STANDARD TERM

WE START BY WRITING DOWN OUR GIVEN ASSUMPTION (THE ASSUMPTIONS OF OUR HYPOTHESIS MODEL). FOR A SINGLE OBSERVED DECM, THE POSITION IS DRAWN FROM:

$$f_x(x | \{\lambda = \lambda\} \cap H) := \begin{cases} \frac{1}{Z(\lambda)} \exp(-\frac{x}{\lambda}) & (x \in (a, b)) \\ 0 & \text{otherwise} \end{cases}$$

WHERE THE NORMALIZING CONSTANT IS:

$$Z(\lambda) = \int_a^b \exp\left(-\frac{x}{\lambda}\right) dx = \left[-\lambda e^{-\frac{x}{\lambda}}\right]_a^b = \lambda [e^{-\lambda x} - e^{-b/\lambda}]$$

LET $\vec{x} = [x_1, x_2, \dots, x_N]$ BE THE DATA, AND ASSUME THAT THE RANDOM VARIABLES FOR EACH ARE INDEPENDENT AND IDENTICALLY DISTRIBUTED, SUCH THAT $f_{\vec{x}} = (f_x)^N$, I.E.

$$f_{\vec{x}}(\vec{x} | \{\lambda = \lambda\} \cap H) = \frac{1}{[Z(\lambda)]^N} \exp\left[-\sum_{i=1}^N \frac{x_i}{\lambda}\right]$$

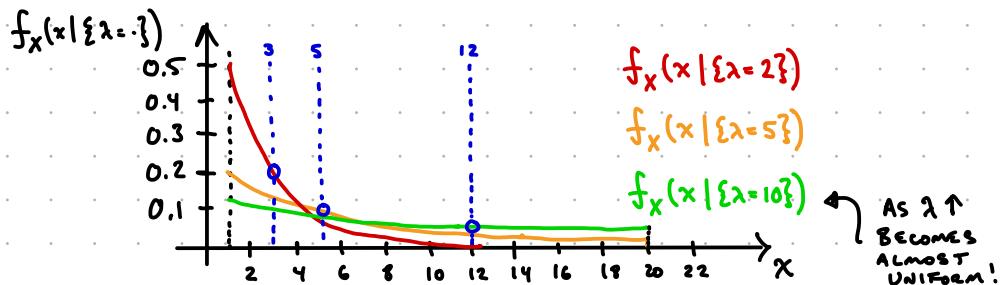
WE CAN THEN WRITE DOWN BAYES THEOREM:

$$f_{\lambda}(\lambda | \{\vec{x} = \vec{x}\} \cap H) = \frac{f_{\vec{x}}(\vec{x} | \{\lambda = \lambda\} \cap H) \cdot f_{\lambda}(\lambda | H)}{\int_0^{\infty} f_{\vec{x}}(\vec{x} | \{\lambda = \lambda\} \cap H) \cdot f_{\lambda}(\lambda | H) d\lambda}$$

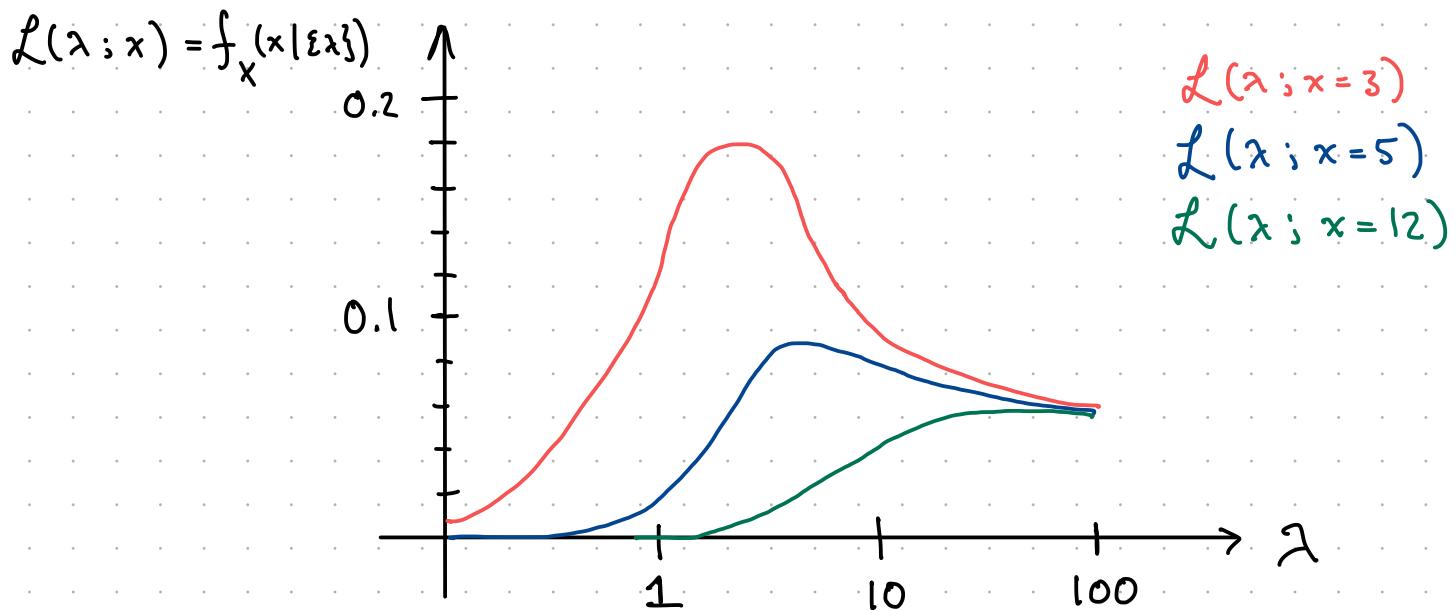
THIS CAN VERY EASILY BE TYPED INTO A COMPUTER (SPECIFICALLY JUST THE LIKELIHOOD FUNCTION, WHICH GIVES ALL INTERESTING INFORMATION FOR $f_{\lambda}(\lambda | \dots)$ GIVEN A UNIFORM PRIOR). BUT LET'S GET SOME INTUITION FOR THE LIKELIHOOD FOR A FEW DATA POINTS, E.G.,

$$N=3, \vec{x} = (3, 5, 12)$$

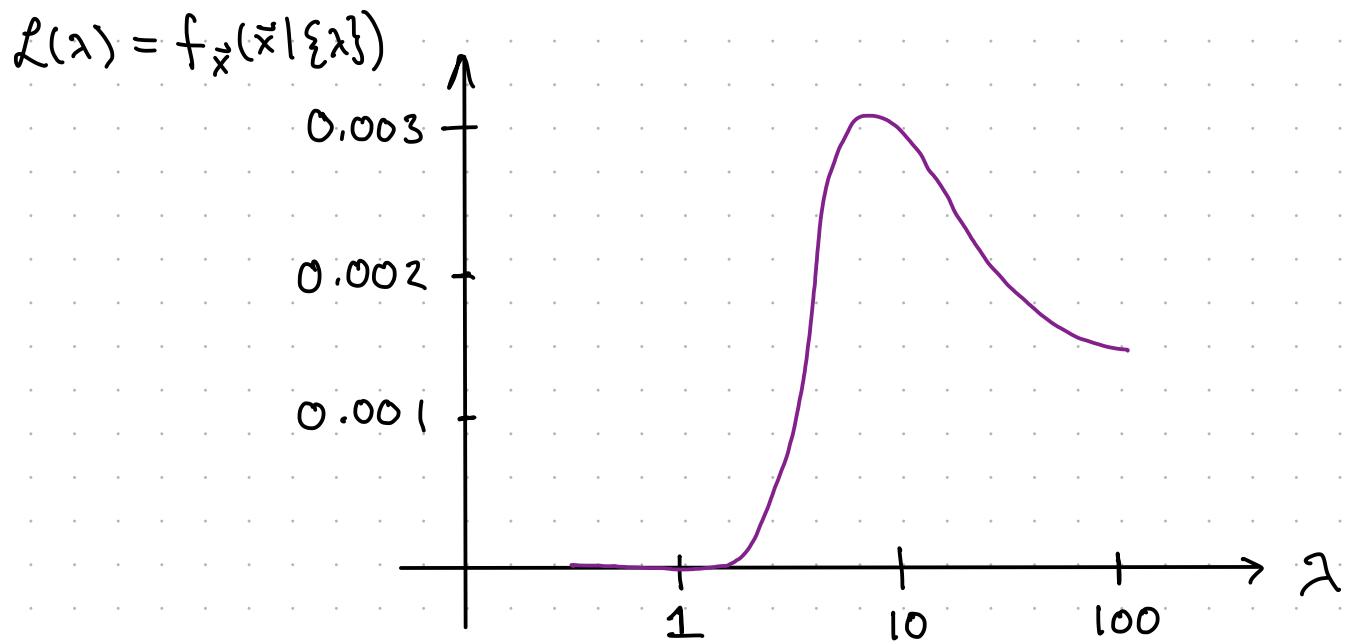
IF WE PLOT SOME TRIAL EXPONENTIALS, THERE ARE SOME λ VALUES THAT ARE GOOD FOR SOME DATA POINTS



NOW, THE LIKELIHOOD IS JUST THE EXPONENTIAL FUNCTION, BUT VIEWED AS A FUNCTION OF λ (NOTE, THOUGH, THE DEPENDENCE OF Z ON λ AS WELL). WE CAN PLOT THE LIKELIHOODS FOR EACH DATA POINT INDIVIDUALLY:



AMAZINGLY, THE SINGLE DATA POINT $x=3$ PREDICTS A WELL-DEFINED PEAK. WE CAN THEN PLOT THEIR PRODUCT, THE OVERALL $N=3$ LIKELIHOOD:



WHICH GIVES US THE MAXIMUM AMOUNT OF INFORMATION ABOUT "WHAT IS λ ?" THAT CAN BE HAD FROM THESE THREE DATA POINTS. (JUST GIVING $\hat{\lambda}$ IS TOO LITTLE INFO!)

⌚ MODEL COMPARISON:

[FINISH WRITING UP MACKAY LECTURE 10]

② INDEPENDENT RANDOM VARIABLES: TWO RANDOM VARIABLES X AND Y ARE SAID TO BE INDEPENDENT IF, FOR ANY TWO BOREL SETS $A, B \subseteq \mathbb{R}$:

$$\text{" } P(X \in A \cap Y \in B) = P(X \in A) \cdot P(Y \in B) \text{"}$$

IN OTHER WORDS, TWO RANDOM VARIABLES ARE INDEPENDENT IF, FOR ALL A AND B , THE EVENTS

$$E_A = \{e \in \mathcal{E} \mid X(e) \in A\}$$

$$E_B = \{f \in \mathcal{E} \mid Y(f) \in B\}$$

ARE INDEPENDENT.

■ IF X AND Y ARE INDEPENDENT RANDOM VARIABLES W/ PROBABILITY DENSITIES $f_X(x)$ AND $f_Y(y)$, THEN:

$$f_{XY}(x,y) = f_X(x) \cdot f_Y(y)$$

■ RANDOM VARIABLES THAT ARE NOT INDEPENDENT ARE CALLED DEPENDENT.