

CS 581 – ADVANCED ARTIFICIAL INTELLIGENCE

TOPIC: PROBABILITY THEORY



Mustafa Bilgic



<http://www.cs.iit.edu/~mbilgic>

MOTIVATION

- The agent needs reason in an uncertain world
- Uncertainty can be due to
 - Noisy sensors (e.g., temperature, GPS, camera, etc.)
 - Imperfect data (e.g., low resolution image)
 - Missing data (e.g., lab tests)
 - Imperfect knowledge (e.g., medical diagnosis)
 - Exceptions (e.g., all birds fly except ostriches, penguins, birds with injured wings, dead birds, ...)
 - Changing data (e.g., flu seasons, traffic conditions during rush hour, etc.)
 - ...
- The agent still must act (e.g., step on the breaks, diagnose a patient, order a lab test, ...)

TENTATIVE PLAN

- Probability background
- Classification
 - Naïve Bayes, logistic regression, neural networks
 - Maximum likelihood estimation, Bayesian estimation, gradient optimization, backpropagation
- Decision making
 - Episodic decision making, Markov decision processes, multi-armed bandits
 - Value of information, Bellman equations, value iteration, policy iteration, UCB1, ϵ -greedy
- Reinforcement learning
 - Prediction, control, Monte-Carlo methods, temporal difference learning, Sarsa, Q-learning

SOME EXERCISES

- In a class, 70% of the hardworking students got an A. John got an A. What is the probability that John is a hardworking student?
- You design a Covid test with the following behavior
 - $P(+ | covid) = 0.95$; $P(- | covid) = 0.05$
 - $P(+ | \sim covid) = 0.10$; $P(- | \sim covid) = 0.90$
 - John takes the test, and the result is +. What is the probability that John has covid?
- In a town, 70% of the hospitalized are vaccinated. Do the vaccines provide any protection against hospitalization?
- $P(toothache | cavity) = 0.75$. $P(cavity | toothache) = ?$

RANDOM VARIABLES

- Pick variables of interest
 - Medical diagnosis
 - Age, gender, weight, temperature, LT1, LT2, ...
 - Loan application
 - Income, savings, payment history, ...
 - Earlier examples
 - Grad student, Grade, Covid, Test result, Ache, X-Ray
- Every variable has a domain
 - Binary (e.g., True/False)
 - Categorical (e.g., Red/Green/Blue)
 - Real-valued (e.g., 97.8)
- **Possible world**
 - An assignment to all variables of interest

PROBABILITY MODEL

- A **probability model** associates a numerical probability $P(w)$ with each possible world w
 - $P(w)$ sums to 1 over all possible worlds
- An **event** is the set of possible worlds where a given predicate is true
 - Roll two dice
 - The possible worlds are (1,1), (1,2), ..., (6,6); 36 possible worlds
 - Predicate = two dice sum to 10
 - Event = {(4,6), (5,5), (6,4)}
 - Toothache and cavity
 - Four possible worlds: $(t, c), (t, \sim c), (\sim t, c), (\sim t, \sim c)$
 - Some worlds are more likely than others
 - Predicate can be anything about these variables: $t \wedge c, t, t \vee \sim c,$

AXIOMS OF PROBABILITY

1. The probability $P(a)$ of a proposition a is a real number between 0 and 1
2. $P(\text{true}) = 1$, $P(\text{false}) = 0$
3. $P(a \vee b) = P(a) + P(b) - P(a \wedge b)$

$P(\neg a)$

- $P(a \vee \neg a) = P(a) + P(\neg a) - P(a \wedge \neg a)$
- $P(\text{true}) = P(a) + P(\neg a) - P(\text{false})$
- $1 = P(a) + P(\neg a) - 0$
- $P(\neg a) = 1 - P(a)$
- Intuitive explanation:
 - The probability of all possible worlds is 1
 - Either a or $\neg a$ holds in one world
 - The worlds that a holds and the worlds that $\neg a$ holds are mutually exclusive and exhaustive

RANDOM VARIABLES – NOTATION

- Capital: X : a variable
- Lowercase: x : a particular value of X
- $\text{Val}(X)$: the set of values X can take
- Bold Capital: \mathbf{X} : a set of variables
- Bold lowercase: \mathbf{x} : an assignment to all variables in \mathbf{X}
- $P(X=x)$ will be shortened as $P(x)$
- $P(X=x \cap Y=y)$ will be shortened as $P(x,y)$

JOINT DISTRIBUTION

- We have n random variables, V_1, V_2, \dots, V_n
- We are interested in the probability of a possible world, where
 - $V_1 = v_1, V_2 = v_2, \dots, V_n = v_n$
- $P(V_1, V_2, \dots, V_n)$ associates a probability for each possible world \equiv the **joint distribution**
 - How many entries are there, if we assume the variables are all binary?

TOOTHACHE EXAMPLE

Ache	X-Ray	P(A, X)
toothache	cavity	0.15
toothache	\neg cavity	0.10
\neg toothache	cavity	0.05
\neg toothache	\neg cavity	0.70

PRIOR AND POSTERIOR

- Prior probability
 - Probability of a proposition in the absence of any other information
 - E.g., $P(V_1, V_3, V_5)$
- Conditional/posterior probability
 - Probability of a proposition given another piece of information
 - E.g., $P(V_2, V_3 \mid V_5 = T, V_7 = F)$
 - $P(A \mid B) = P(A \wedge B) / P(B)$

MARGINALIZATION

- Given a distribution over n variables, you can calculate the distribution over any subset of the variables by summing out the irrelevant ones
- For example
 - Given $P(A, B, C, D)$
 - Calculate
 - $P(A)$
 - $P(A, C)$
 - ... (any subset)

LET'S ANSWER A FEW QUERIES

Ache	X-Ray	P(A, X)
toothache	cavity	0.15
toothache	\neg cavity	0.10
\neg toothache	cavity	0.05
\neg toothache	\neg cavity	0.70

- $P(\text{cavity}) = ?$ 0.20
- $P(\neg \text{cavity}) = ?$ 0.80
- $P(\text{toothache}) = ?$ 0.25
- $P(\neg \text{toothache}) = ?$ 0.75

CONDITIONAL DISTRIBUTION

- $P(A, B, C \mid D, E, F, G) = \frac{P(A, B, C, D, E, F, G)}{P(D, E, F, G)}$

LET'S ANSWER A FEW QUERIES

Ache	X-Ray	P(A, X)
toothache	cavity	0.15
toothache	¬cavity	0.10
¬toothache	cavity	0.05
¬toothache	¬cavity	0.70

- $P(\text{cavity} \mid \text{toothache}) = ?$ ← $\frac{P(C, t)}{P(t)} = \frac{0.15}{0.25} = 0.6$
- $P(\text{cavity} \mid \neg\text{toothache}) = ?$
- $P(\neg\text{cavity} \mid \text{toothache}) = ?$
- $P(\neg\text{cavity} \mid \neg\text{toothache}) = ?$
- $P(\text{toothache} \mid \text{cavity}) = ?$
- $P(\neg\text{toothache} \mid \text{cavity}) = ?$
- $P(\text{toothache} \mid \neg\text{cavity}) = ?$
- $P(\neg\text{toothache} \mid \neg\text{cavity}) = ?$

BAYES' RULE

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

BAYES' RULE

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

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

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

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

BAYES' RULE

- $P(B|A) = \frac{P(A|B)*P(B)}{P(A)}$
- Example use
 - $P(\text{cause} | \text{effect}) = P(\text{effect} | \text{cause}) * P(\text{cause}) / P(\text{effect})$
- Why is this useful?
 - Because in practice it is easier to get probabilities for $P(\text{effect} | \text{cause})$ and $P(\text{cause})$ than for $P(\text{cause} | \text{effect})$
 - E.g., $P(\text{disease} | \text{symptoms}) = P(\text{symptoms} | \text{disease}) * P(\text{disease}) / P(\text{symptoms})$
 - It is easier to know what symptoms diseases cause. It is harder to diagnose a disease given symptoms

BAYES RULE

- Can we compute $P(\alpha|\beta)$ from $P(\beta|\alpha)$?

CLASS EXAMPLE

- In a class, 70% of the hardworking students got an A. John got an A. What is the probability that John is a hardworking student?
- Possible worlds: 4
 - $\langle h, a \rangle, \langle h, \sim a \rangle, \langle \sim h, a \rangle, \langle \sim h, \sim a \rangle$
- Let's say there are 100 students in a class
- Let's say 10 of them work hard (h), 90 do not ($\sim h$)
- Probability of a randomly picked student being hardworking
 - $P(h) = 0.1$
- We are told that 70% of the hardworking students got an A.
 - $P(a | h) = 0.7$
 - 7 hardworking students got an A; 3 did not get an A.
- What is $P(h|a) = ?$

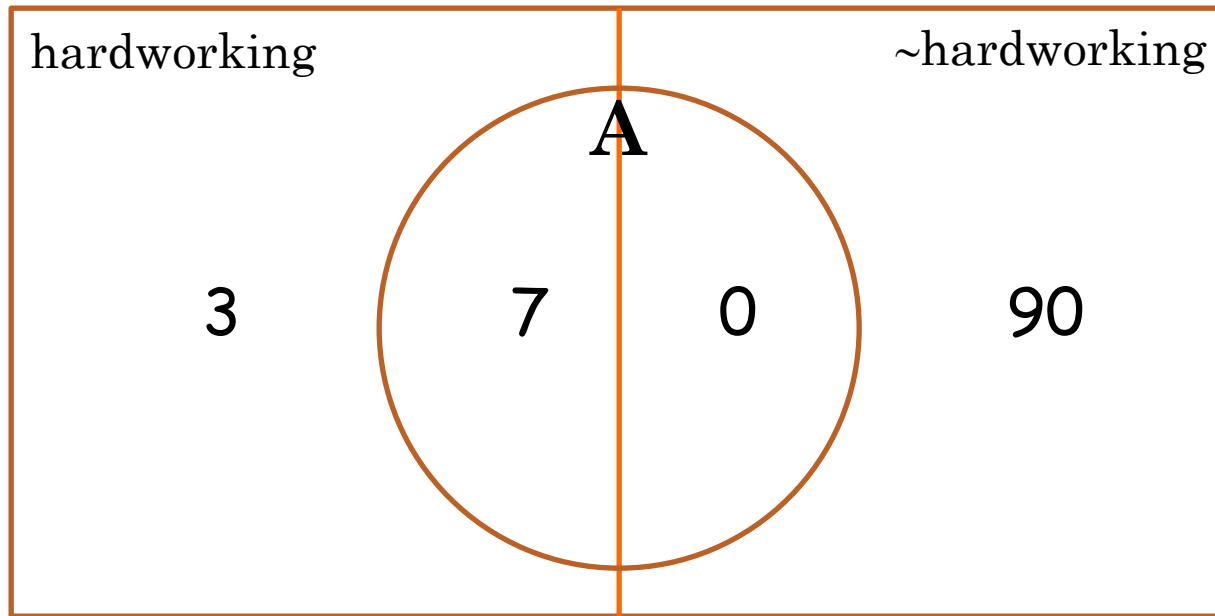
H, A

$$P(a|h) = 0.7$$

CLASS EXAMPLE

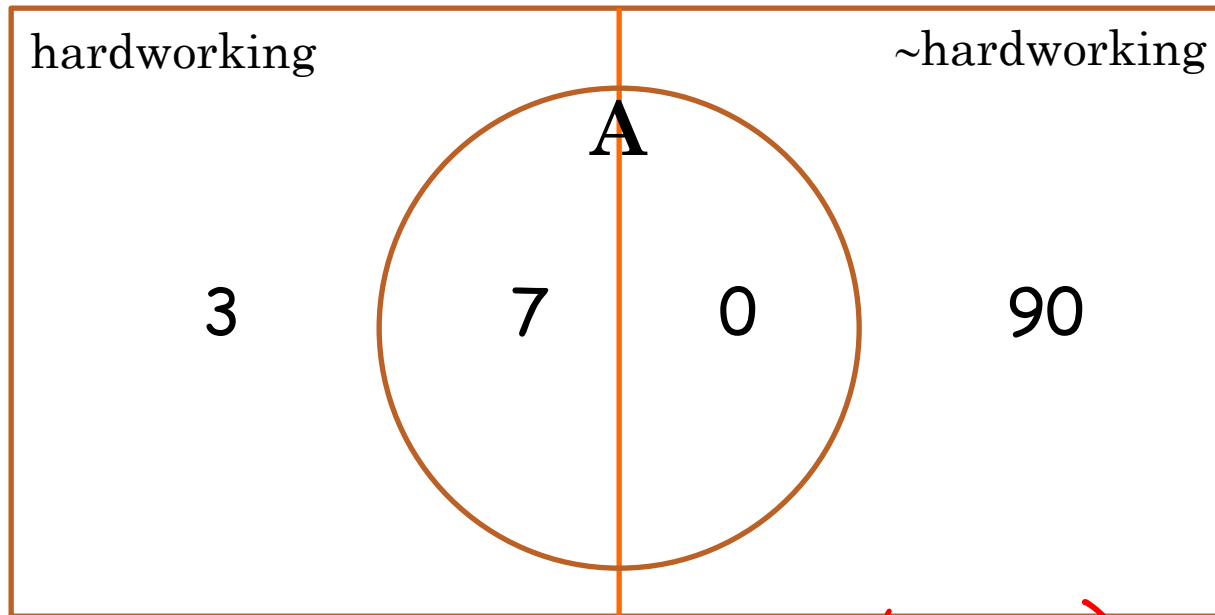
- In a class, 70% of the hardworking students got an A. John got an A. What is the probability that John is a hardworking student? $P(h|a)$
- Possible worlds: 4
 - $\langle h, a \rangle, \langle h, \sim a \rangle, \langle \sim h, a \rangle, \langle \sim h, \sim a \rangle$
- Let's say there are 100 students in a class
- Let's say 10 of them work hard (h), 90 do not ($\sim h$)
- Probability of a randomly picked student being hardworking
 - $P(h) = 0.1$
- We are told that 70% of the hardworking students got an A.
 - $P(a|h) = 0.7$
 - 7 hardworking students got an A; 3 did not get an A.
- What is $P(h|a) = ?$

VERY DIFFICULT CLASS



$$P(h | a) = ?$$

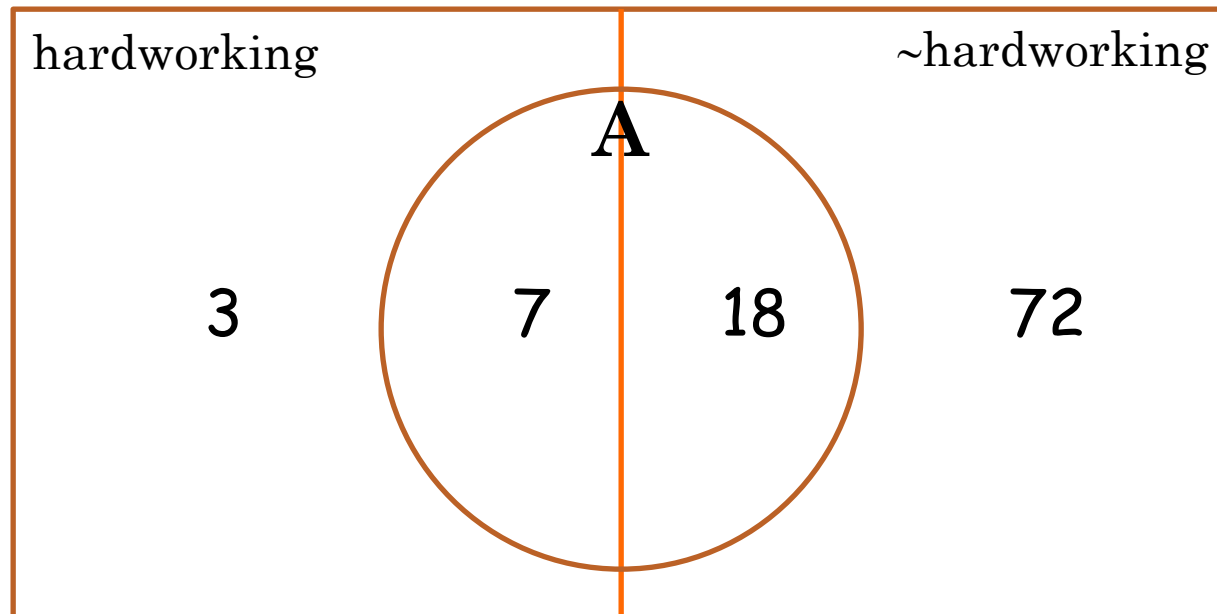
VERY DIFFICULT CLASS



$$P(h | a) = ?$$

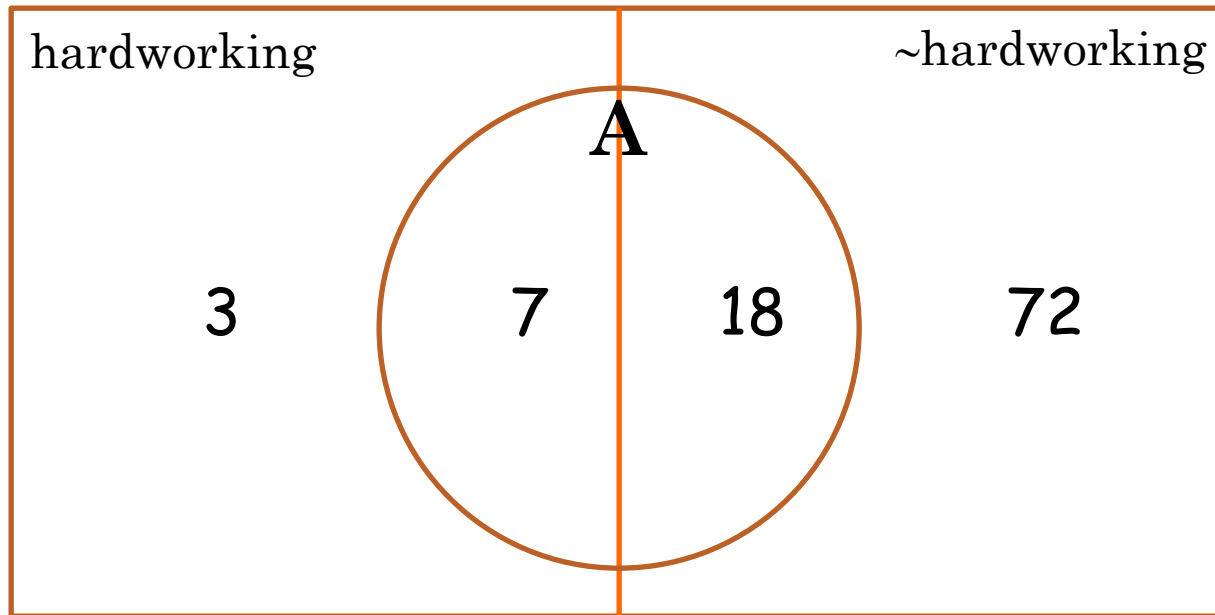
$$\frac{P(h, a)}{P(a)} = \frac{7}{7}$$

MEDIUM DIFFICULT CLASS



$$P(h | a) = ?$$

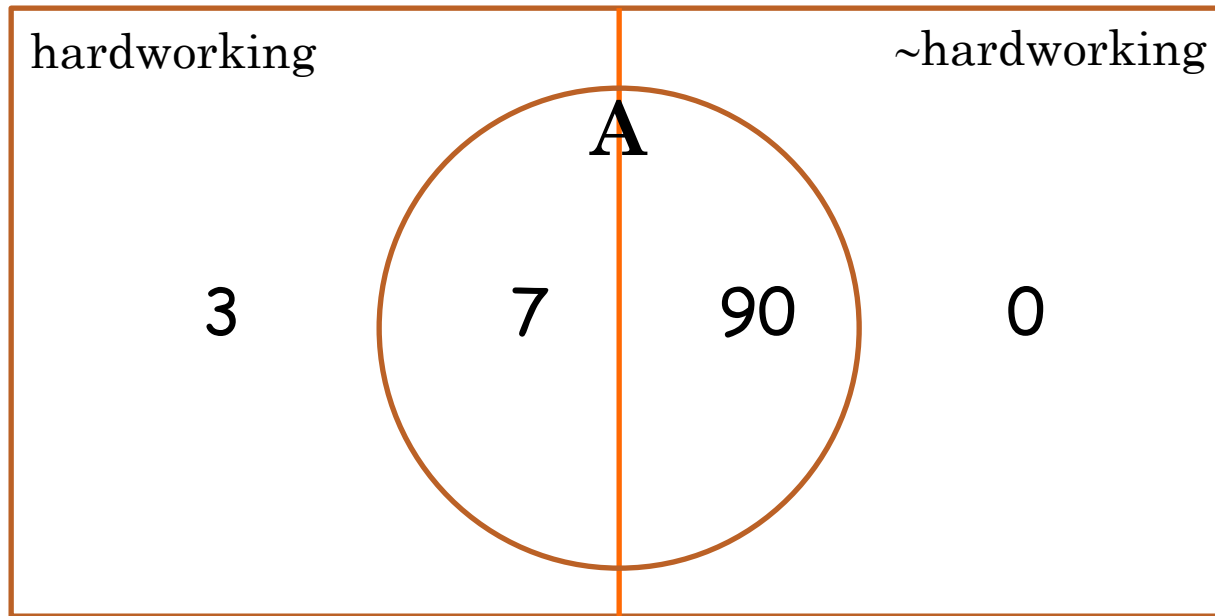
MEDIUM DIFFICULT CLASS



$$P(h | a) = ?$$

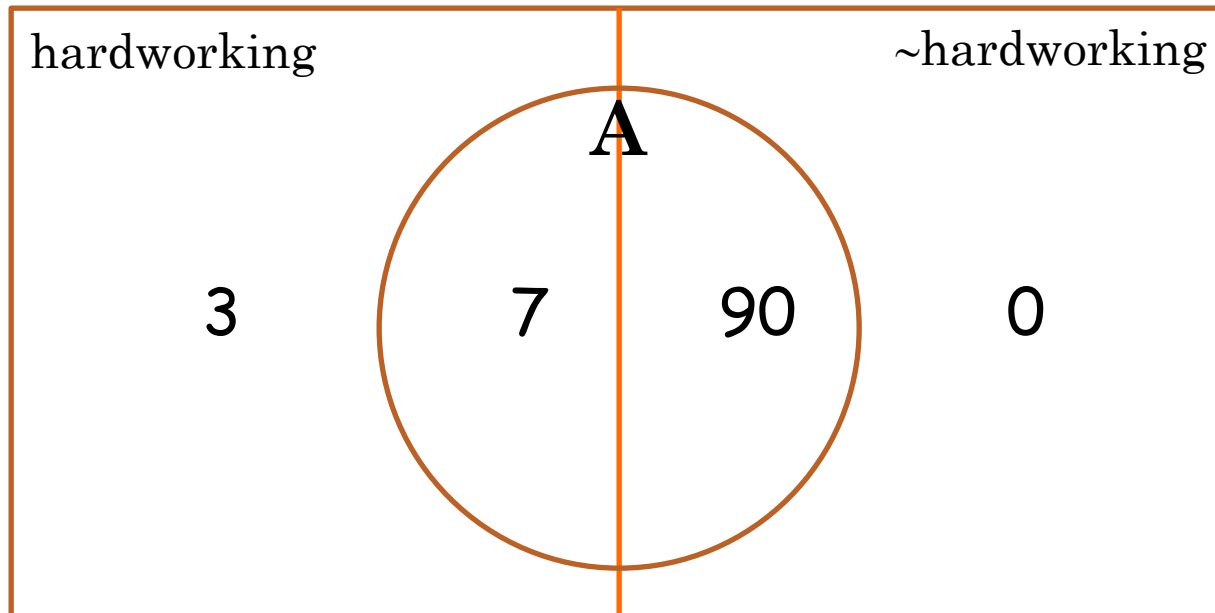
$$\frac{7}{25} = 0.28$$

WEIRD CLASS



$$P(h | a) = ?$$

WEIRD CLASS



$$P(h | a) = ?$$

$\frac{7}{97} \approx 0.07$

CHAIN RULE

- $P(X_1, X_2, X_3, \dots, X_k) =$
 - $P(X_1) P(X_2 | X_1) P(X_3 | X_1, X_2) \dots P(X_k | X_1, X_2, X_3, \dots, X_{k-1})$
 - or
 - $P(X_2) P(X_1 | X_2) P(X_3 | X_1, X_2) \dots P(X_k | X_1, X_2, X_3, \dots, X_{k-1})$
 - or
 - $P(X_2) P(X_3 | X_2) P(X_1 | X_3, X_2) \dots P(X_k | X_1, X_2, X_3, \dots, X_{k-1})$
 - or
 - Pick an order, then
 - $P(\text{first})P(\text{second} | \text{first})P(\text{third} | \text{first}, \text{second}) \dots P(\text{last} | \text{all_previous})$

$$P(A, B, C, D) = P(A) P(B|A) P(C|A, B) P(D|A, B, C)$$

A B C D

B D A C

$$= P(B) P(D|B)$$

$$P(A|B, D)$$

$$P(C|B, D, A)$$

MARGINAL INDEPENDENCE

- An event α is **independent** of event β in P , denoted as $P \models \alpha \perp \beta$, if
 - $P(\alpha | \beta) = P(\alpha)$, or
 - $P(\beta) = 0$
- Proposition: A distribution P satisfies $\alpha \perp \beta$ if and only if
 - $P(\alpha, \beta) = P(\alpha) P(\beta)$
 - Can you prove it?
- Corollary: $\alpha \perp \beta$ implies $\beta \perp \alpha$

$$P(\alpha, \beta) = P(\beta) P(\alpha | \beta)$$
$$P(\beta) P(\alpha)$$

MARGINAL INDEPENDENCE

$$P(X) \cdot P(Y)$$

X	Y	P(X, Y)
t	t	0.18
t	f	0.42
f	t	0.12
f	f	0.28

$$0.6 \times 0.3 = 0.18$$

$$0.6 \times 0.7 = 0.42$$

$$0.4 \times 0.3 = 0.12$$

$$0.4 \times 0.7 = 0.28$$

Is $X \perp Y$?

$$P(X, Y) = P(X) P(Y)$$

$$\langle 0.6, 0.4 \rangle \cdot \langle 0.3, 0.7 \rangle$$

$$P(X|Y) \stackrel{?}{=} P(X) \quad P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

MARGINAL INDEPENDENCE

$$\frac{0.18}{0.3} \stackrel{?}{=} 0.6$$

$$\frac{0.42}{0.7} \stackrel{?}{=} 0.6$$

$$P(X) \cdot P(Y)$$

X	Y	P(X, Y)
t	t	0.18
t	f	0.42
f	t	0.12
f	f	0.28

$$0.6 \times 0.3 = 0.18$$

$$0.6 \times 0.7 = 0.42$$

$$0.4 \times 0.3 = 0.12$$

$$0.4 \times 0.7 = 0.28$$

Is $X \perp Y$?

$$P(X, Y) = P(X) P(Y)$$

$$\langle 0.6, 0.4 \rangle \cdot \langle 0.3, 0.7 \rangle$$

J M T

~~J I M~~
J I M | T

CONDITIONAL INDEPENDENCE

- Two events are independent given another event
- An event α is **independent** of event β given event γ in P , denoted as $P \models (\alpha \perp \beta \mid \gamma)$, if
 - $P(\alpha \mid \beta, \gamma) = P(\alpha \mid \gamma)$, or
 - $P(\beta, \gamma) = 0$
- Proposition: A distribution P satisfies $\alpha \perp \beta \mid \gamma$ if and only if
 - $P(\alpha, \beta \mid \gamma) = P(\alpha \mid \gamma) P(\beta \mid \gamma)$

~~H A K~~

Height Knowledge
 $H \perp K \mid A$

NUMBER OF PARAMETERS

- Assuming everything is binary
- $P(V_1)$ requires
 - 1 independent parameter
- $P(V_1, V_2, \dots, V_n)$ requires
 - $2^n - 1$ independent parameters
- $P(V_1 | V_2)$ requires
 - 2 independent parameters
- $P(V_1, V_2, \dots, V_n | V_{n+1}, V_{n+2}, \dots, V_{n+m})$ requires
 - $2^m \times (2^n - 1)$ independent parameters

NUMBER OF PARAMETERS

- Assuming everything is binary
- $P(V_1)$ requires
 - 1 independent parameter
- $P(V_1, V_2, \dots, V_n)$ requires
 - $2^n - 1$ independent parameters
- $P(V_1 | V_2)$ requires
 - 2 independent parameters
- $P(V_1, V_2, \dots, V_n | V_{n+1}, V_{n+2}, \dots, V_{n+m})$ requires
 - $2^m \times (2^n - 1)$ independent parameters

CONTINUOUS SPACES

- Assume X is continuous and $\text{Val}(X) = [0,1]$
- If you would like to assign the same probability to all real numbers in $[0, 1]$, what is, for e.g., $P(X=0.5) = ?$

PROBABILITY DENSITY FUNCTION

- We define **probability density function**, $p(x)$, a non-negative integrable function, such that $\int_{Val(X)} p(x)dx = 1$

$$P(X \leq a) = \int_{-\infty}^a p(x)dx$$

$$P(a \leq X \leq b) = \int_a^b p(x)dx$$

UNIFORM DISTRIBUTION

- A variable X has a uniform distribution over $[a,b]$ if it has the PDF

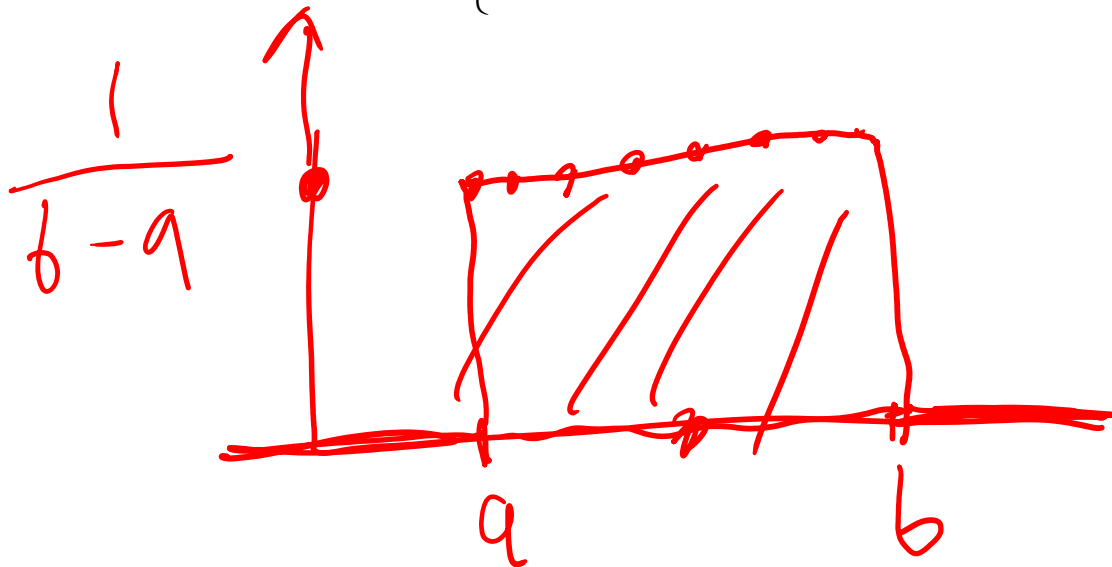
$$p(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

Check and make sure that $p(x)$ integrates to 1.

UNIFORM DISTRIBUTION

- A variable X has a uniform distribution over $[a,b]$ if it has the PDF

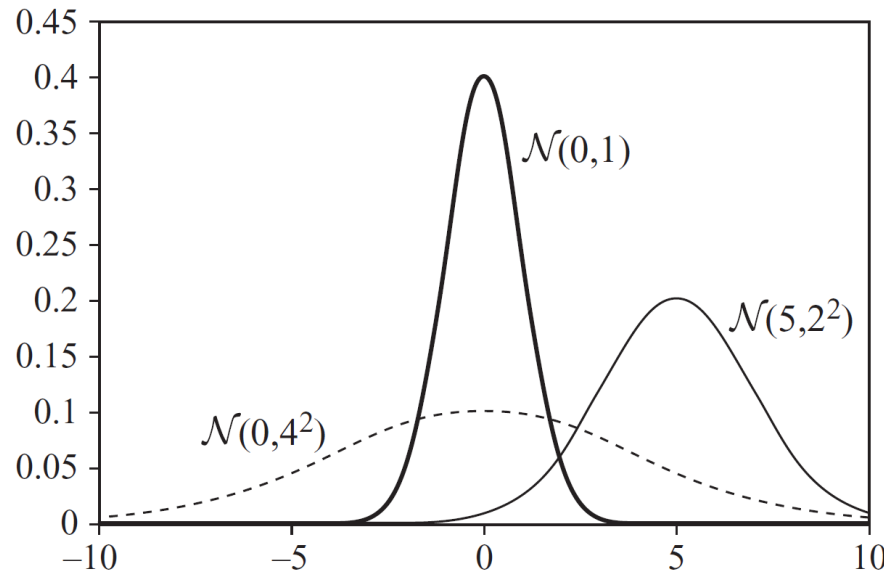
$$p(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$



GAUSSIAN DISTRIBUTION

- A variable X has a Gaussian distribution with mean μ and variance σ^2 , if it has the PDF

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Can $p(x)$ be ever greater than 1?

Yes!

CONDITIONAL PROBABILITY

- We want $P(Y | X=x)$ where X is continuous, Y is discrete
- $P(Y | X=x) = P(Y, X=x) / P(X=x)$
 - What's wrong with this expression?
- Instead, we use the following expression

$$P(Y | X = x) = \lim_{\varepsilon \rightarrow 0} P(Y | x - \varepsilon \leq X \leq x + \varepsilon)$$

CONDITIONAL PROBABILITY

- We want $p(Y | X)$ where X is discrete, Y is continuous
- How would you represent it?

$$\begin{array}{ll} p(Y | X=a) & N(\mu_a, \sigma_a^2) \\ p(Y | X=b) & N(\mu_b, \sigma_b^2) \\ \vdots & \vdots \\ p(Y | X=z) & \end{array}$$

EXPECTATION

$$E_P[X] = \sum_x xP(x)$$

$$E_P[X] = \int_x xp(x)dx$$

$$E_P[aX + b] = aE_P[X] + b$$

$$E_P[X + Y] = E_P[X] + E_P[Y]$$

$$E_P[X | y] = \sum_x xP(x | y)$$

What about $E[X*Y]$?

EXPECTATION

$$E_P[X] = \sum_x xP(x)$$

$$E_P[X] = \int_x xp(x)dx$$

$$E_P[aX + b] = aE_P[X] + b$$

$$E_P[X + Y] = E_P[X] + E_P[Y]$$

$$E_P[X | y] = \sum_x xP(\underline{x} | y)$$

What about $E[X*Y]$? $\neq E[X] * E[Y]$

$$\begin{array}{ccc} 3 & 5 & 10 \\ 0.1 & 0.2 & 0.7 \\ 0.3 + 1 + 7 = 8.3 \end{array}$$

$$\begin{array}{ccc} 10 & 14 & 24 \\ 0.1 & 0.2 & 0.7 \end{array}$$

$x^2 + 4$

$$\begin{array}{ccc} 2 & 15 & 16 \\ 0.3 & 0.4 & 0.3 \end{array}$$

VARIANCE

$$\text{Var}_P[X] = E_P \left[\left(X - E_P[X] \right)^2 \right]$$

$$\text{Var}_P[X] = E_P[X^2] - \left(E_P[X] \right)^2$$

Can you derive the second expression using the first expression?

$$\text{Var}_P[aX + b] = a^2 \text{Var}_P[X]$$

What is $\text{Var}[X+Y]$?

3 5 10
0.1 0.2 0.7

VARIANCE

$$\text{Var}_P[X] = E_P \left[(X - E_P[X])^2 \right] \leftarrow \left[X^2 - 2 \cdot X \cdot \bar{X} + \bar{X}^2 \right]$$

$$\text{Var}_P[X] = E_P[X^2] - (E_P[X])^2 \leftarrow$$

Can you derive the second expression using the first expression?

$$\text{Var}_P[aX + b] = a^2 \text{Var}_P[X]$$

What is $\text{Var}[X+Y]$?

UNIFORM AND GAUSSIAN DISTRIBUTION

- If $X \sim N(\mu, \sigma^2)$, then $E[X] = \mu$, $\text{Var}[X] = \sigma^2$
- What about the expectation and variance of a uniform distribution?