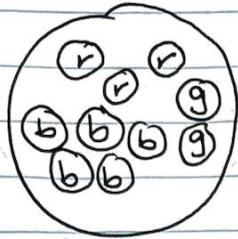


Bayesian ML - Homework 1 - Probability and Priors

Problem 1: Basic Probability



10 balls total
 $\frac{5}{10}$ blue
 $\frac{3}{10}$ red
 $\frac{2}{10}$ green

(a) The probability of drawing a red ball from the bag is $\frac{3}{10}$

(b) ~~$\frac{4}{9}$ blue balls~~ Probability that

$\frac{5}{10}$ probability that the first ball is blue

$\frac{4}{9}$ probability that the next ball is also blue

$\frac{5}{10} \times \frac{4}{9} = \frac{20}{90} = \frac{2}{9}$, probability that the next ball is also blue

Problem 2: Independent Events

$P(\text{Server being down}) = 0.05$, server failures are independent.

Probability of 2 independent events $P(A \text{ and } B) = P(A) \times P(B)$

$P(\text{server 1 being down AND server 2 being down}) = P(0.05) \times P(0.05) = 0.0025$

(a) 0.0025

Probability of at least 1 being down

$P(\text{server 1 being down OR server 2 being down}) =$

$P(0.05) + P(0.05) = 0.1$

(b) 0.1

Courtney
Hodge

Problem 3: Conditional Probability

	DS	Non-DS	
30%	♀ ♂	♀ ♂	70%
Non-PhD /	\ PhD	Non-PhD /	\ PhD
60%	40%	90%	10%
♂	♀	♀	♂

(a) If an employee is chosen randomly, what is the probability that the employee is a Data Scientist?

$$P(\text{DS}) \times P(\text{PhD} | \text{DS}) = \left(\frac{3}{10}\right) \times \left(\frac{4}{10}\right) = \boxed{\frac{12}{100}}$$

(b) Given that an employee has a PhD, what is the probability that the employee is a Data Scientist?

Bayes' Theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \rightarrow P(\text{DS}|\text{PhD}) = \frac{P(\text{PhD}|\text{DS})P(\text{DS})}{P(\text{PhD})}$$

#1 Calculate PhD using the law of total probability:

$$P(A) = \sum_{i=1}^n P(A|B_i)P(B_i)$$

$$\begin{aligned} P(\text{PhD}) &= P(\text{PhD} | \text{DS})P(\text{DS}) + P(\text{PhD} | \text{Non-DS})P(\text{Non-DS}) \\ &= \left(\frac{4}{10}\right)\left(\frac{3}{10}\right) + \left(\frac{1}{10}\right)\left(\frac{7}{10}\right) \\ &= \frac{12}{100} + \frac{7}{100} = \frac{19}{100} \end{aligned}$$

→ #2 Solve: $P(\text{DS}|\text{PhD}) = \frac{\left(\frac{4}{10}\right)\left(\frac{3}{10}\right)}{\left(\frac{19}{100}\right)} = 0.63$

Problem 4: Law of Total Probability

Law of Total Probability

$$P(A) = \sum_{i=1}^n P(A|B_i) P(B_i)$$

$$P(\text{Tested Positive} | \text{Has Disease}) = 0.95$$

$$P(\text{Tested Positive} | \text{Doesn't Have Disease}) = 0.1$$

$$P(\text{Has Disease}) = 0.5\% = 0.005$$

$$P(\text{Doesn't Have Disease}) = 0.995$$

(a) If a person tested positive in the test, what is the probability that the person actually has the disease?

$$P(\text{Has Disease} | \text{Tested Positive}) = \frac{P(\text{Has Disease}) P(\text{Tested Positive} | \text{Has Disease})}{P(\text{Tested Positive})}$$

$$P(A \cap B) \quad P(A) \quad \xrightarrow{\quad P(B|A) \quad}$$

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

$$P(A|B) = \frac{0.005(0.95)}{P(B)} = \frac{0.005(0.95)}{0.10425} = 0.0456$$

$$\begin{aligned} \rightarrow P(B) &= P(\text{Tested Positive}) = P(\text{Tested Pos} | \text{Has Disease}) P(\text{Has Disease}) \\ &\quad + P(\text{Tested Pos} | \text{Doesn't Have Disease}) * \\ &\quad P(\text{Doesn't Have Disease}) \end{aligned}$$

$$\begin{aligned} P(\text{Tested Positive}) &= 0.95(0.005) + (0.1)0.995 \\ &= 0.10425 \end{aligned}$$

(b) What is the total probability of a person testing positive?

$$P(\text{Tested Positive}) = 0.10425$$

Courtney
Hodge

Problem 5

Bayes' Theorem

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

$$P(\text{Spam}) = .9$$

$$P(\text{Not-Spam}) = 0.1 \quad P(\text{Identified Not-Spam} | \text{Spam}) = 0.05$$

$$P(\text{Identifies Spam} | \text{Spam}) = 0.95$$

$$P(\text{Identifies Not-Spam} | \text{Not-Spam}) = 0.85$$

$$P(\text{Identified Spam} | \text{Not-Spam}) = \boxed{0.15}$$

(a) If an email is picked at random, and the filter classifies it as spam, what is the probability that it is actually spam?

$$P(\text{Spam}) \times P(\text{Identifies Spam} | \text{Spam}) = 0.9 \times 0.95 = 0.855$$

$$\rightarrow P(\text{Spam} | \text{Identified Spam}) = \frac{P(\text{Identified Spam} | \text{Spam})P(\text{Spam})}{P(\text{Identified Spam})}$$

$$P(A \downarrow, B \downarrow)$$

$$P(\text{Identified Spam})$$

$$P(A|B) = \frac{0.95(0.9)}{P(B)} = \frac{0.95(0.9)}{0.81} = 0.983$$

Law of Total Probability

$$P(A) = \sum_{i=1}^n P(A|B_i)P(B_i)$$

$$\begin{aligned} P(\text{Identified Spam}) &= P(\text{Identified Spam} | \text{Spam})P(\text{Spam}) + \\ &\quad P(\text{Identified spam} | \text{Not-Spam})P(\text{Not-Spam}) \\ &= 0.95(.9) + (\boxed{0.15})0.1 \\ &= 0.81 \end{aligned}$$

$$P(\text{Identified Spam} | \text{Not-Spam}) = 0.15$$

Courtney
Hodge

Problem 5

(b) If an email is classified as "not spam", what is the probability that it is actually spam?

$P(\text{Spam} \mid \text{Identified Not-Spam})$

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

$$= \frac{P(\text{Spam}) P(\text{Identified Not-Spam} \mid \text{Spam})}{P(\text{Identified Not-Spam})}$$

$$= \frac{0.9 (0.05)}{P(\text{Identified Not-Spam})} = \frac{0.9(0.05)}{0.13} = 0.346$$

Law of Total Probability

$$P(A) = \sum_{i=1}^n P(A \mid B_i) P(B_i)$$

$$\begin{aligned} P(\text{Identified}^{\text{^}} \text{Spam}) &= P(\text{Identified}^{\text{^}} \text{Spam} \mid \text{Spam}) P(\text{Spam}) \\ &\quad + P(\text{Identified Not Spam} \mid \text{Not-Spam}) * \\ &\quad P(\text{Not-Spam}) \\ &= 0.05 (0.9) + (0.85) 0.1 \\ &= 0.13 \end{aligned}$$

Contrey
Hodge

Problem 6: Expectation of a Discrete Random Variable

(a) Define the random variable X that models this game.

$$\frac{1}{6} = 6, \text{ win \$10}$$

$$\frac{5}{6} = 1, 2, 3, 4 \text{ or } 5, \text{ lose \$2}$$

$$\cancel{P\left(\frac{1}{6}\right) + P\left(\frac{5}{6}\right)}$$

$X = 10$, if you roll a 6

$X = -2$, if you roll 1 through 5

$$X = \begin{cases} 10, & \text{if a 6 appears} \\ -2, & \text{if any number 1 to 5 appears} \end{cases}$$

(b) Compute the expected value of X

$$P(\text{if a 6 appears})X + P(\text{if you roll 1 to 5})X$$

expected

To compute the value of X , we use the expected value of a discrete random variable

$E(X) = \sum_i P(X_i) \cdot X_i$, where X_i are the possible outcomes and $P(X_i)$ are the corresponding probabilities.

Prob of rolling a 6 is $P(X=10) = \frac{1}{6}$

Prob of rolling anything else $P(X=-2) = \frac{5}{6}$

$$E(X) = P(X=10)(\cancel{10}) + P(X=-2)(\cancel{-2}) - 2$$

$$\cancel{\frac{1}{6}} \cdot \cancel{10} - \frac{10}{6} = 0$$

$$\left(\frac{1}{6}\right)(10) + \left(\frac{5}{6}\right)(-2) = \frac{10}{6} + \left(-\frac{10}{6}\right) = 0 \rightarrow$$

$$E(X) = 0$$

Problem 7: Expectation of a Continuous Random Variable

(a) Compute the expected value $E[X]$ of X

Probability Density Function

$$f(x) = 2e^{-2x} \text{ for } x \geq 0 \quad \text{PDF} = f(x)$$

Integration By Parts

$$u = x \quad du = dx$$

$$v = -e^{-2x} \quad dv = 2e^{-2x}$$

$$\int 2e^{-2x} dx \quad u = -2x \quad du = -2dx$$

$$\rightarrow -e^u \quad \rightarrow \int u dv = uv - \int v du \quad (\text{take neg out})$$

$$= -e^u \quad = -e^{-2x} \quad E(X) = \left[-xe^{-2x} \right]_0^\infty + \int_0^\infty e^{-2x} dx$$

Solving, we get $\underset{A}{\therefore} x=0, -xe^{-2x}=0$

$\underset{A}{\therefore} x=\infty, -xe^{-2x} \rightarrow 0 = 0$

$$\text{Solving } \int_0^\infty e^{-2x} dx = \left[-\frac{1}{2}e^{-2x} \right]_0^\infty = (0 - (-\frac{1}{2})) = \frac{1}{2}$$

$$\text{Therefore, } \boxed{E(X) = 0 + \frac{1}{2} = \frac{1}{2}}$$

(b) Compute the variance $\text{Var}[X]$ of X

Computing $E(X^2)$: $E(X^2) = \int_0^\infty x^2 \cdot f(x) dx$ Variance of a continuous random variable X is

Since we know that the PDF of an exponential distribution is $f(x) = \lambda e^{-\lambda x}$, then... $E(X^2) = \int_0^\infty x^2 / (2e^{-2x}) dx$ $\text{Var}(X) = E(X^2) - [E(X)]^2$

$$E(X) = 1/\lambda$$

$$E(X^2) = 2/\lambda^2 \dots \text{with } \lambda=2$$

$$E(X^2) = \frac{2}{\lambda^2} = \frac{2}{2^2} = \frac{2}{4} = \frac{1}{2} \rightarrow \text{Now we can compute the Variance}$$

$$V(X) = E(X^2) - [E(X)]^2 = 0.5 - (0.5)^2 = 0.5 - 0.25 = 0.25$$

$$V(X) = 0.25 \text{ yrs}^2$$

(c) Interpret your findings for parts a and b in the context of the server's processing time.

$E(X) = 0.5$ means that on average, it will take servers 30 minutes to process a query of this type.

$V(X) = 0.25$ is the spread of variability of the server's processing times around the average, meaning processing times can take significantly longer or shorter times to process based on the query.

$$V(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n (x_i - 0.5)^2 = \frac{1}{n} \sum_{i=1}^n (x_i^2 - 1) = \text{minutes}$$

$$\frac{1}{n} \sum_{i=1}^n x_i^2 = \bar{x} + V(X) = (\bar{x})^2$$

$\bar{x} \neq [x]$ not same since not squared (a)

new rows is to deviate at (x_1, x_2, \dots, x_n) $\rightarrow (x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x})$

$$E(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) = (\bar{x}) - \bar{x} = 0$$

$$E(X^2) = (\bar{x})^2$$

$$(S = \sum_{i=1}^n x_i^2) \rightarrow (S - n\bar{x}^2) = \sum_{i=1}^n (x_i - \bar{x})^2$$

deviate into original row so we have $S = \sum_{i=1}^n x_i^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n 2\bar{x}(x_i - \bar{x}) + \sum_{i=1}^n \bar{x}^2 = S - n\bar{x}^2 + n\bar{x}^2 = S$

$$E(S) = E(S - n\bar{x}^2 + n\bar{x}^2) = E(S) = E(\sum_{i=1}^n (x_i - \bar{x})^2) = (E(x - \bar{x}))^2 = 0$$

$$E(S) = (E(x))^2$$