

# CMPT 210: Probability and Computing

## Lecture 21

---

Sharan Vaswani

November 21, 2024

- **Variance:** Standard way to measure the deviation from the mean. For r.v.  $X$ ,  $\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \sum_{x \in \text{Range}(X)} (x - \mu)^2 \Pr[X = x]$ , where  $\mu := \mathbb{E}[X]$ .
- **Alternate Definition:**  $\text{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$ .
- If  $X \sim \text{Ber}(p)$ ,  $\text{Var}[X] = p(1 - p) \leq \frac{1}{4}$ .
- **Standard Deviation:** For r.v.  $X$ , the standard deviation in  $X$  is defined as:  
 $\sigma := \sqrt{\text{Var}[X]} = \sqrt{\mathbb{E}[X^2] - (\mathbb{E}[X])^2}$ .
- For constants  $a, b$  and r.v.  $R$ ,  $\text{Var}[aR + b] = a^2 \text{Var}[R]$ .
- For pairwise independent random variables  $R_1, R_2, R_3, \dots, R_n$ ,  $\text{Var}[\sum_{i=1}^n R_i] = \sum_{i=1}^n \text{Var}[R_i]$ .

**Definition:** For two random variables  $R$  and  $S$ , the covariance between  $R$  and  $S$  is defined as:

$$\text{Cov}[R, S] := \mathbb{E}[(R - \mathbb{E}[R]) (S - \mathbb{E}[S])] = \mathbb{E}[RS] - \mathbb{E}[R] \mathbb{E}[S]$$

$$\text{Cov}[R, S] = \mathbb{E}[(R - \mathbb{E}[R]) (S - \mathbb{E}[S])]$$

$$= \mathbb{E}[RS - R \mathbb{E}[S] - S \mathbb{E}[R] + \mathbb{E}[R] \mathbb{E}[S]] \quad (\text{Expanding the terms})$$

$$= \mathbb{E}[RS] - \mathbb{E}[R \mathbb{E}[S]] - \mathbb{E}[S \mathbb{E}[R]] + \mathbb{E}[R] \mathbb{E}[S] \quad (\text{Linearity of Expectation})$$

$$\implies \text{Cov}[R, S] = \mathbb{E}[RS] - \mathbb{E}[R] \mathbb{E}[S] - \mathbb{E}[S] \mathbb{E}[R] + \mathbb{E}[R] \mathbb{E}[S] = \mathbb{E}[RS] - \mathbb{E}[R] \mathbb{E}[S]$$

- Covariance generalizes the notion of variance to multiple random variables.

$$\text{Cov}[R, R] = \mathbb{E}[R R] - \mathbb{E}[R] \mathbb{E}[R] = \text{Var}[R]$$

- If  $R$  and  $S$  are independent r.v's,  $\mathbb{E}[RS] = \mathbb{E}[R] \mathbb{E}[S]$  and  $\text{Cov}[R, S] = 0$ .
- The covariance between two r.v's is symmetric i.e.  $\text{Cov}[R, S] = \text{Cov}[S, R]$ .

# Covariance

- For two arbitrary (not necessarily independent) r.v's,  $R$  and  $S$ ,

$$\text{Var}[R + S] = \text{Var}[R] + \text{Var}[S] + 2 \text{Cov}[R, S]$$

- Recall from Slide 9 in Lecture 20, where we showed that,

$$\text{Var}[R + S] = \text{Var}[R] + \text{Var}[S] + 2(\mathbb{E}[RS] - \mathbb{E}[R] \mathbb{E}[S]) = \text{Var}[R] + \text{Var}[S] + 2 \text{Cov}[R, S].$$

If  $R$  and  $S$  are independent,  $\text{Cov}[R, S] = 0$  and we recover the formula for the sum of independent variables.

- If  $R = S$ ,  $\text{Var}[R + R] = \text{Var}[R] + \text{Var}[R] + 2\text{Cov}[R, R] = \text{Var}[R] + \text{Var}[R] + 2\text{Var}[R] = 4\text{Var}[R]$  which is consistent with our previous formula that  $\text{Var}[2R] = 2^2\text{Var}[R]$ .
- Generalization to multiple random variables  $R_1, R_2, \dots, R_n$  (Recall from Slide 10 in Lecture 20):

$$\text{Var} \left[ \sum_{i=1}^n R_i \right] = \sum_{i=1}^n \text{Var}[R_i] + 2 \sum_{1 \leq i < j \leq n} \text{Cov}[R_i, R_j]$$

## Covariance - Example

**Q:** If  $X$  and  $Y$  are indicator r.v.'s for events  $A$  and  $B$  respectively, calculate the covariance between  $X$  and  $Y$

We know that  $\text{Cov}[X, Y] = \mathbb{E}[XY] - \mathbb{E}[X] \mathbb{E}[Y]$ . Note that  $X = \mathcal{I}_A$  and  $Y = \mathcal{I}_B$ . We can conclude that  $XY = \mathcal{I}_{A \cap B}$  since  $XY = 1$  iff both events  $A$  and  $B$  happen.

$$\implies \mathbb{E}[X] = \Pr[A] ; \mathbb{E}[Y] = \Pr[B] ; \mathbb{E}[XY] = \Pr[A \cap B]$$

$$\implies \text{Cov}[X, Y] = \mathbb{E}[XY] - \mathbb{E}[X] \mathbb{E}[Y] = \Pr[A \cap B] - \Pr[A] \Pr[B]$$

If  $\text{Cov}[X, Y] > 0 \implies \Pr[A \cap B] > \Pr[A] \Pr[B]$ . Hence,

$$\Pr[A|B] = \frac{\Pr[A \cap B]}{\Pr[B]} > \frac{\Pr[A] \Pr[B]}{\Pr[B]} = \Pr[A]$$

If  $\text{Cov}[X, Y] > 0$ , it implies that  $\Pr[A|B] > \Pr[A]$  and hence, the probability that event  $A$  happens increases if  $B$  is going to happen/has happened. Similarly, if  $\text{Cov}[X, Y] < 0$ ,  $\Pr[A|B] < \Pr[A]$ . In this case, if  $B$  happens, then the probability of event  $A$  decreases.

# Correlation

**Definition:** The correlation between two r.v's  $R_1$  and  $R_2$  is defined as:

$$\text{Corr}[R_1, R_2] = \frac{\text{Cov}[R_1, R_2]}{\sqrt{\text{Var}[R_1] \text{Var}[R_2]}}$$

$\text{Corr}[R_1, R_2] \in [-1, 1]$  and indicates the strength of the relationship between  $R_1$  and  $R_2$ .

- If  $\text{Corr}[R_1, R_2] > 0$ , then  $R_1$  and  $R_2$  are said to be positively correlated, else if  $\text{Corr}[R_1, R_2] < 0$ , the r.v's are negatively correlated.

- If  $R_1 = R_2 = R$ , then,  $\text{Corr}[R, R] = \frac{\text{Cov}[R, R]}{\sqrt{\text{Var}[R] \text{Var}[R]}} = \frac{\text{Var}[R]}{\text{Var}[R]} = 1$ .

- If  $R_1$  and  $R_2$  are independent,  $\text{Cov}[R_1, R_2] = 0$  and  $\text{Corr}[R_1, R_2] = 0$ .

- If  $R_1 = -R_2 = R$ , then,

$$\begin{aligned}\text{Corr}[R, -R] &= \frac{\text{Cov}[R, -R]}{\sqrt{\text{Var}[R] \text{Var}[-R]}} = \frac{\text{Cov}[R, -R]}{\sqrt{\text{Var}[R] (-1)^2 \text{Var}[R]}} = \frac{\text{Cov}[R, -R]}{\text{Var}[R]} \\ &= \frac{\mathbb{E}[-R^2] - \mathbb{E}[R] \mathbb{E}[-R]}{\text{Var}[R]} = \frac{-\mathbb{E}[R^2] + \mathbb{E}[R] \mathbb{E}[R]}{\text{Var}[R]} = \frac{-\text{Var}[R]}{\text{Var}[R]} = -1\end{aligned}$$

# Matching Birthdays

**Q:** In a class of  $n$  students, what is the probability that two students share the same birthday? Assume that (i) each student is equally likely to be born on any day of the year, (ii) no leap years and (iii) student birthdays are independent of each other.

For  $d := 365$  (since no leap years),

$$\Pr[\text{two students share the same birthday}] = 1 - \frac{d \times (d-1) \times (d-2) \times \dots \times (d-(n-1))}{d^n}$$

**Q:** On average, how many pairs of students have matching birthdays?

Define  $M$  to be the number of pairs of students with matching birthdays. For a fixed ordering of the students, let  $X_{i,j}$  be the indicator r.v. corresponding to the event  $E_{i,j}$  that the birthdays of students  $i$  and  $j$  match. Hence,

$$M = \sum_{i,j|1 \leq i < j \leq n} X_{i,j} \implies \mathbb{E}[M] = \mathbb{E}\left[\sum_{i,j|1 \leq i < j \leq n} X_{i,j}\right] = \sum_{i,j|1 \leq i < j \leq n} \mathbb{E}[X_{i,j}] = \sum_{i,j|1 \leq i < j \leq n} \Pr[E_{i,j}]$$

(Linearity of expectation)

# Matching Birthdays

For a pair of students  $i, j$ , let  $B_i$  be the r.v. equal to the day of student  $i$ 's birthday.  $\text{Range}(B_i) = \{1, 2, \dots, d\}$ . For all  $k \in [d]$ ,  $\Pr[B_i = k] = 1/d$  (each student is equally likely to be born on any day of the year).

$$E_{i,j} = (B_i = 1 \cap B_j = 1) \cup (B_i = 2 \cap B_j = 2) \cup \dots$$

$$\Rightarrow \Pr[E_{i,j}] = \sum_{k=1}^d \Pr[B_i = k \cap B_j = k] = \sum_{k=1}^d \Pr[B_i = k] \Pr[B_j = k] = \sum_{k=1}^d \frac{1}{d^2} = \frac{1}{d}$$

(student birthdays are independent of each other)

$$\Rightarrow \mathbb{E}[M] = \sum_{i,j | 1 \leq i < j \leq n} \Pr[E_{i,j}] = \frac{1}{d} \sum_{i,j | 1 \leq i < j \leq n} (1) = \frac{1}{d} [(n-1) + (n-2) + \dots + 1] = \frac{n(n-1)}{2d}$$

Hence, in our class of 100 students, on average, there are  $\frac{(100)(50)}{365} = 13.7$  students with matching birthdays.



# Matching Birthdays

**Q:** Are the  $X_{i,j}$  r.v.'s mutually independent?

No, because if  $X_{i,k} = 1$  and  $X_{j,k} = 1$ , then,

$$\Pr[X_{i,j} = 1 | X_{j,k} = 1 \cap X_{i,k} = 1] = 1 \neq \frac{1}{d} = \Pr[X_{i,j} = 1].$$

**Q:** Are the  $X_{i,j}$  pairwise independent?

Yes, because for all  $i, j$  and  $i', j'$  (where  $i \neq i'$ ),  $\Pr[X_{i,j} = 1 | X_{i',j'} = 1] = \Pr[X_{i,j} = 1]$  because if students  $i'$  and  $j'$  have matching birthdays, it does not tell us anything about whether  $i$  and  $j$  have matching birthdays.

# Matching Birthdays

**Q:** If  $M$  is the random variable equal to the number of pairs of students with matching birthdays, calculate  $\text{Var}[M]$ .

$$\text{Var}[M] = \text{Var}\left[\sum_{i,j|1 \leq i < j \leq n} X_{i,j}\right]$$

Since  $X_{i,j}$  are pairwise independent, the variance of the sum is equal to the sum of the variance.

$$\begin{aligned} \Rightarrow \text{Var}[M] &= \sum_{i,j|1 \leq i < j \leq n} \text{Var}[X_{i,j}] = \sum_{i,j|1 \leq i < j \leq n} \frac{1}{d} \left(1 - \frac{1}{d}\right) = \frac{1}{d} \left(1 - \frac{1}{d}\right) \frac{n(n-1)}{2} \\ &\quad \text{(Since } X_{i,j} \text{ is an indicator (Bernoulli) r.v. and } \Pr[X_{i,j} = 1] = \frac{1}{d}\text{)} \end{aligned}$$

Hence, in our class of 75 students, the standard deviation for the matching birthdays is equal to  $\sqrt{\frac{(100)(50)}{365} \frac{364}{365}} \approx 3.7$ .

Questions?