# Probability Theory

## Applications for Data Science
## Module 5: Expectation, Variance, Covariance, and Correlation
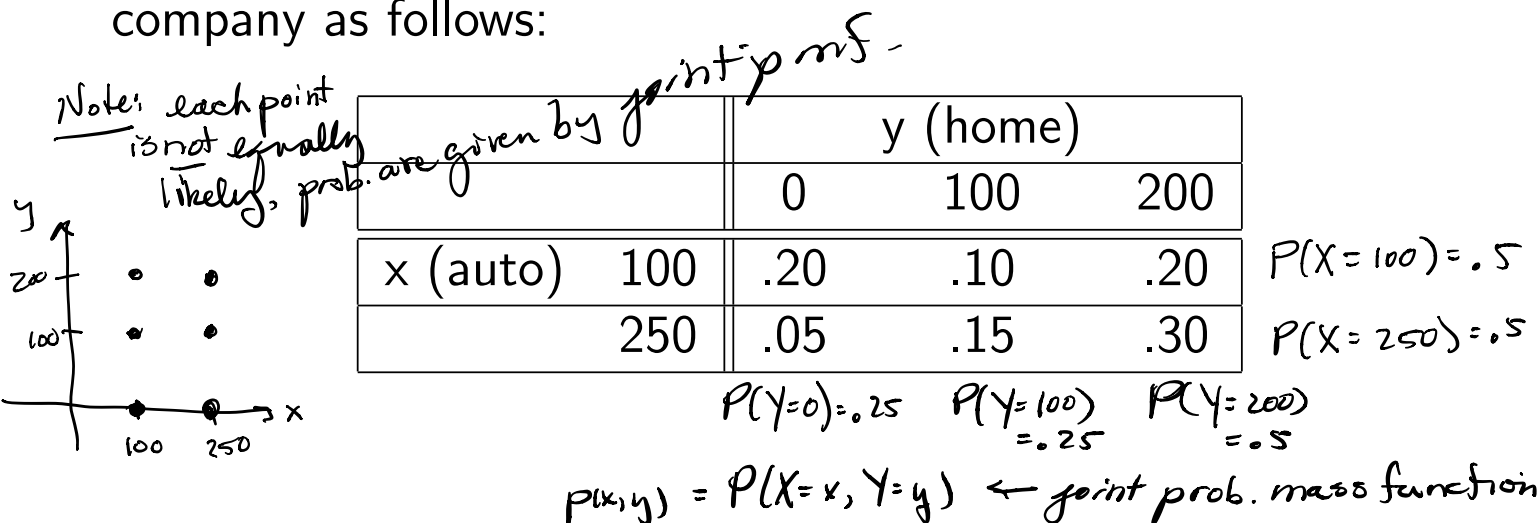
Anne Dougherty

March 20, 2021

# TABLE OF CONTENTS

# Expectation, Variance, Covariance, and Correlation

At the end of this module, students should be able to

▶ Compute the mean, variance, and standard deviation of a function of a random variable (i.e. $g(X)$).

▶ Explain the concept of jointly distributed random variables, for two random variables $X$ and $Y$.

▶ **Define, compute, and interpret the covariance between two random variables $X$ and $Y$.**

▶ **Define, compute, and interpret the correlation between two random variables $X$ and $Y$.**

Example: An insurance agency services customers who have both a homeowner's policy and an automobile policy. For each type of policy, a deductible amount must be specified. For an automobile policy, the choices are $100 or $250 and for the homeowner's policy, the choices are $0, $100, or $200.

Suppose the **joint probability table** is given by the insurance company as follows:

*joint p mf.*

| | | y (home) 0 | 100 | 200 |
|---|---|---|---|---|
| x (auto) | 100 | .20 | .10 | .20 |
| | 250 | .05 | .15 | .30 |

Note: each point is not equally likely, prob. are given by joint p mf.

$P(X=100) = .5$
$P(X=250) = .5$

$P(Y=0) = .25$  $P(Y=100) = .25$  $P(Y=200) = .5$

$p(x,y) = P(X=x, Y=y)$ ← joint prob. mass function

y
20
100
100   250
x

When two random variables, $X$ and $Y$, are not independent, it is frequently of interest to assess how strongly they are related to each other.

Definition: The **covariance** between two rv's, $X$ and $Y$, is defined as:

$$\text{Cov}(X, Y) = E\left[(X - E(x))(Y - E(Y))\right]$$

$$= E\left[(X - \mu_x)(Y - \mu_y)\right]$$

expectation
need to sum
over all possible
$x$ & $y$ values

$$= \begin{cases} \sum_x \sum_y (x - \mu_x)(y - \mu_y) P(X = x, Y = y), & (\text{discrete}) \\[2em] \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu_x)(y - \mu_y) f(x,y) dx\, dy, & (\text{cont}) \end{cases}$$

$X$ & $Y$
discrete
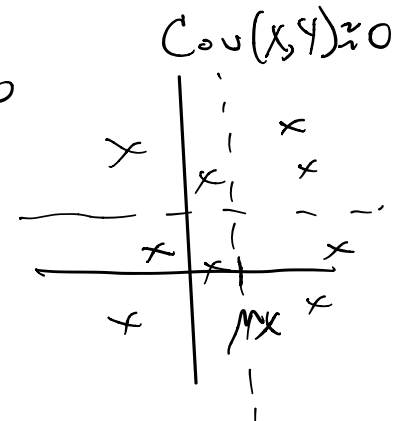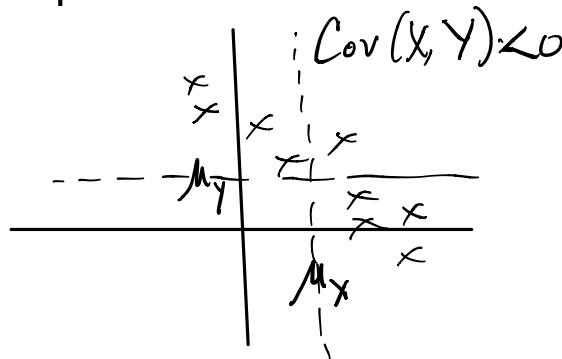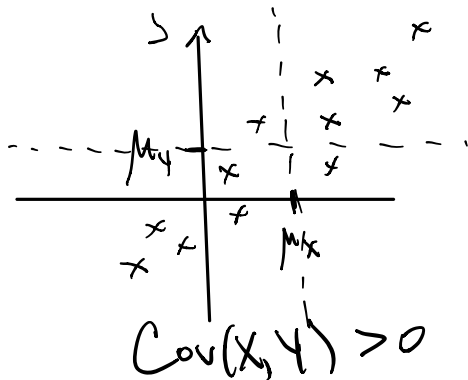
Definition: Covariance of $X$ and $Y$ is given by
$$Cov(X, Y) = E[(X - E(X))(Y - E(Y))]$$

To calculate covariance:

$$Cov(X, Y) = \begin{cases} \sum_x \sum_y (x - \mu_X)(y - \mu_Y)P(X = x, Y = y) \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu_X)(y - \mu_Y)f(x, y) \, dx \, dy \end{cases}$$

The covariance depends on both the set of possible pairs and the probabilities for those pairs.



$Cov(X, Y) > 0$

$Cov(X, Y) < 0$

$Cov(X, Y) \approx 0$

$Cov(X, Y) = E[(X - E(X))(Y - E(Y))]$

▶ If both variables tend to deviate in the same direction (both go above their means or below their means at the same time), then the covariance will be positive.

▶ If the opposite is true, the covariance will be negative.

▶ If $X$ and $Y$ are not strongly (linearly) related, the covariance will be near 0.

Aside: It is possible to have a strong relationship between $X$ & $Y$ and still have $Cov(X, Y) \approx 0$

e.g.

$$Cov(X,Y) = \sum_x \sum_y (x-\mu_X)(y-\mu_Y) \, P\{X=x, Y=y\}$$

Covariance example calculation:

|  |  | y (home) | | |
|---|---|---|---|---|
|  |  | 0 | 100 | 200 |
| x (auto) | 100 | .20 | .10 | .20 |
|  | 250 | .05 | .15 | .30 |

$$\mu_X = \sum_x x\,P(X=x) = 100\,(.5) + 250\,(.5) = 175$$

$$\mu_Y = \sum_y y\,P(Y=y) = 0\,(.25) + 100\,(.25) + 200\,(.5) = 125$$

| x | y | $x-\mu_X$ | $y-\mu_Y$ | $P(X=x, Y=y)$ |
|---|---|---|---|---|
| 100 | 0 | −75 | −125 | .2 |
| 250 | 0 | 75 | −125 | .05 |
| 100 | 100 | −75 | −25 | .1 |
| 250 | 100 | 75 | −25 | .15 |
| 100 | 200 | −75 | 75 | .2 |
| 250 | 200 | 75 | 75 | .3 |

$$Cov(X,Y) = 1875$$

Is this a strong relationship between X & Y? It seems like a "big" number, but it's hard to say. The correlation coef. will help.

But, before we get to correlation, there's a
few more ideas related to covariance we need
to discuss.

Computational formula for covariance:

$$Cov(X, Y) = E(XY) - E(X)E(Y)$$

$$\left( \text{Recall: } V(X) = E\left((X - E(X))^2\right) = E(X^2) - (E(X))^2 \right.$$

$$\rightarrow Cov(X, Y) = E\left((X - E(X))(Y - E(Y))\right)$$

$$= E\{ XY - X E(Y) - Y E(X) + E(X)E(Y) \}$$

$$= E(XY) - E(\underbrace{X E(Y)}_{\text{constant}}) - E(\underbrace{Y E(X)}) + E(\underbrace{E(X)E(Y)}_{\text{constant}})$$

$$= \underline{E(XY) - E(X)E(Y)}$$

What if $X$ and $Y$ are independent?

If $X \& Y$ are indep., $P(X=x, Y=y) = P(X=x)P(Y=y)$ for all $x, y$.
$\quad \& $ discrete

$$Cov(X,Y) = \sum_x \sum_y (x-\mu_x)(y-\mu_y) P(X=x, Y=y)$$

$$\overset{X \& Y \text{ indep}}{=} \sum_x \sum_y (x-\mu_x)(y-\mu_y) P(X=x) P(Y=y)$$

$$= \left[\sum_x (x-\mu_x) P(X=x)\right]\left[\sum_y (y-\mu_y) P(Y=y)\right\}$$

$$= \left[\underbrace{\sum_x x P(X=x)}_{E(X)} - \mu_x \underbrace{\sum_x P(X=x)}_{1}\right]\left[\quad \text{similar} \quad \right\}$$

$$= 0$$

So, if $X \& Y$ are indep, $Cov(X, Y) = 0$    $X \& Y$ may
* It does not go the other way. If $Cov(X, Y) = 0$   still be *dependent*.

Useful formulas for random variables $X$ and $Y$ and real numbers $a$ and $b$:

- $E(aX + bY) = aE(X) + bE(Y)$

- $V(aX + bY) = a^2 V(X) + b^2 V(Y) + 2abCov(X, Y)$

$$V(aX+bY) = E\left[\left(aX+bY - E(aX+bY)\right)^2\right]$$
$$= E\left[\left(a(X - E(X)) + b(Y - E(Y))\right)^2\right]$$
$$= a^2 E\left[(X - E(X))^2\right] + b^2 E\left[(Y - E(Y))^2\right]$$
$$+ 2ab \, E\left[(X - E(X))(Y - E(Y))\right]$$
$$= a^2 V(X) + b^2 V(Y) + 2ab \, Cov(X, Y)$$

Definition: The **correlation coefficient** of $X$ and $Y$, denoted by $Cor(X, Y)$ or just $\rho_{XY}$, is defined by

$$\rho_{X,Y} = \frac{Cov\,(X,Y)}{\sigma_X\,\sigma_Y}$$

It represents a "scaled" covariance. The correlation is always between -1 and 1.

Two ~~extreme example~~s: *special cases*

- ▶ What if $X$ and $Y$ are independent?

$$Cov(X,Y)=0 \quad so \quad \rho_{X,Y} = \frac{Cov(X,Y)}{\sigma_X \sigma_Y} = 0$$

- ▶ What if $Y = aX + b$? $\quad \rho_{X,Y} = \frac{Cov(X,Y)}{\sigma_X \sigma_Y}$

① $Cov(X,Y) = Cov(X, aX+b)$
$$= E[(X-E(X))(aX+b-E(aX+b))]$$
$$= a E[(X-E(X))^2] = a V(X) = a\sigma_X^2$$

② Also $V(Y) = E[(Y-E(Y))^2]$
$$= E\{(aX+b-E(aX+b)\} = a^2 \cdot V(X)$$

③ $\rho_{X,Y} = \dfrac{a \sigma_X^2}{\sigma_X \cdot |a|\sigma_X} \quad \overset{\sigma_Y=|a|\sigma_X}{=} \quad \dfrac{a}{|a|} = \begin{cases} 1 & \text{if } a > 0 \\ -1 & \text{if } a < 0 \end{cases}$

|  | y (home) | | |
|---|---|---|---|
|  | 0 | 100 | 200 |
| x (auto)   100 | .20 $\frac{1}{6}$ | .10 $\frac{1}{6}$ | .20 $\frac{1}{6}$  $p$  $\frac{1}{2}$ |
| 250 | .05 $\frac{1}{6}$ | .15 $\frac{1}{6}$ | .30 $\frac{1}{6}$  $\frac{1}{2}$ |
|  | $\frac{1}{3}$ | $\frac{1}{3}$ | $\frac{1}{3}$ |

Find $\rho_{XY}$

Earlier: $\text{Cov}(X,Y) = 1875$, $E(X) = 175$, $E(Y) = 125$

You should verify $V(X) = 75^2$ & $V(Y) = 6875 = \sigma_y^2$

$$\rho_{X,Y} = \frac{1875}{75\sqrt{6875}} \approx .3$$

Conclusions: Correlation measures the strength of the linear relationship between $X$ & $Y$. If $X$ & $Y$ are indep, $\rho_{X,Y} = 0$. But if you compute $\rho_{X,Y} = 0$ cannot conclude independence.

In the next module we'll transition to multiple random variables.