

Chapter 5. Vector random variables

- A **vector random variable** $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a collection of random numbers with probabilities assigned to outcomes.
- \mathbf{X} can also be called a **multivariate random variable**.
- The case with $n = 2$ we call a **bivariate random variable**.
- Saying X and Y are **jointly distributed random variables** is equivalent to saying (X, Y) is a bivariate random variable.
- Vector random variables let us model relationships between quantities.

Example: midterm and final scores

- We will look at the anonymized test scores for a previous course.

```
download.file(destfile="course_progress.txt",  
url="https://ionides.github.io/401f18/01/course_progress.txt")
```

```
# Anonymized scores for a random subset of 50 students
```

```
"final" "quiz" "hw" "midterm"
```

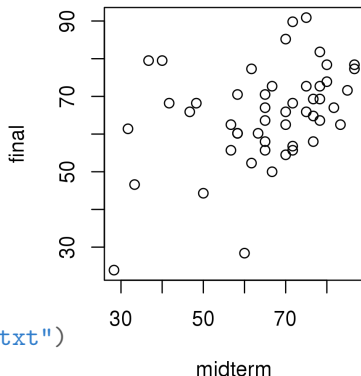
```
"1" 52.3 76.7 91 61.7
```

```
"2" 68.2 65.4 94.5 48.3
```

```
"3" 78.4 91.2 95.5 80
```

- A probability model lets us answer a question like, “What is the probability that someone gets at least 70% in both the midterm and the final”

```
x <- read.table("course_progress.txt")  
plot(final~midterm,data=x)
```



The bivariate normal distribution and covariance

- Let $X \sim \text{normal}(\mu_X, \sigma_X)$ and $Y \sim \text{normal}(\mu_Y, \sigma_Y)$.
- If X and Y are bivariate random variables we need another parameter to describe their dependence. If X is big, does Y tend to be big, or small, or does the value of X make no difference to the outcome of Y ?
- This parameter is the **covariance**, defined to be

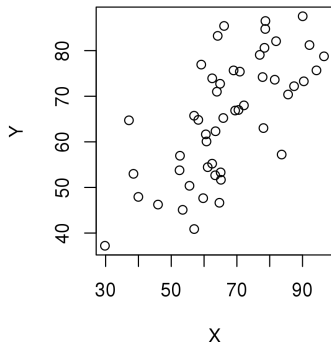
$$\text{Cov}(X, Y) = \text{E}[(X - \text{E}[X]) (Y - \text{E}[Y])]$$

- The parameters of the bivariate normal distribution in matrix form are the **mean vector** $\mu = (\mu_X, \mu_Y)$ and the **variance/covariance matrix**,

$$\mathbb{V} = \begin{bmatrix} \text{Var}(X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & \text{Var}(Y) \end{bmatrix}$$

- In R, the `mvtnorm` package lets us simulate the bivariate and multivariate normal distribution. It uses the vector and matrix form for the parameters.

Experimenting with the bivariate normal distribution



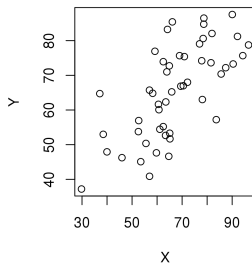
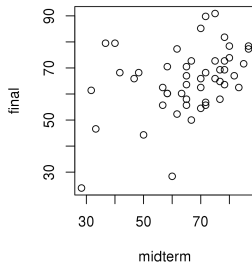
```
library(mvtnorm)
mvn <- rmvnorm(n=50,
  mean=c(X=65,Y=65),
  sigma=matrix(
    c(200,100,100,150),
    2,2)
)
plot(Y~X,data=mvn)
```

- We write $(X, Y) \sim \text{MVN}(\boldsymbol{\mu}, \mathbb{V})$, where MVN is read "multivariate normal".

Question 5.1. What are μ_X , μ_Y , $\text{Var}(X)$, $\text{Var}(Y)$, and $\text{Cov}(X, Y)$ for this simulation?

The bivariate normal as a model for exam scores

Question 5.2. Compare the data on midterm and final scores with the simulation. Does a normal model seem to fit? Would you expect it to? Why, and why not?



More on covariance

- Covariance is **symmetric**: we see from the definition

$$\begin{aligned}\text{Cov}(X, Y) &= \mathbb{E}\left[(X - \mathbb{E}[X]) (Y - \mathbb{E}[Y])\right] \\ &= \mathbb{E}\left[(Y - \mathbb{E}[Y]) (X - \mathbb{E}[X])\right] = \text{Cov}(Y, X)\end{aligned}$$

- Also, we see from the definition that $\text{Cov}(X, X) = \text{Var}(X)$.
- The **sample covariance** of n pairs of measurements $(x_1, y_1), \dots, (x_n, y_n)$ is

$$\text{cov}(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

where \bar{x} and \bar{y} are the sample means of $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{y} = (y_1, \dots, y_n)$.

Scaling covariance to give correlation

- The standard deviation of a random variable is interpretable as its scale.
- Variance is interpretable as the square of standard deviation

```
var(x$midterm)
## [1] 218.2155
var(x$final)
## [1] 169.7518
cov(x$midterm,x$final)
## [1] 75.61269
```

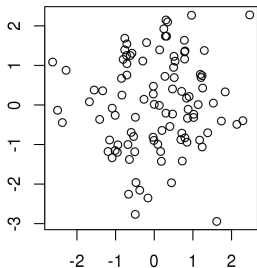
- Covariance is interpretable when scaled to give the **correlation**

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

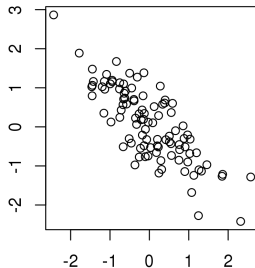
$$\text{cor}(\mathbf{x}, \mathbf{y}) = \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{\sqrt{\text{var}(\mathbf{x})\text{var}(\mathbf{y})}}$$

```
cor(x$midterm,x$final)
## [1] 0.3928662
```

```
rho <- 0
```

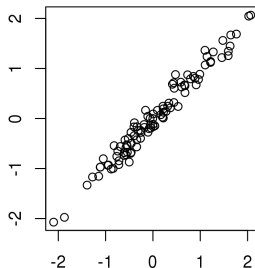


```
rho <- -0.8
```



```
library(mvtnorm)
mvn <- rmvnorm(n=100,
  mean=c(X=0,Y=0),
  sigma=matrix(
    c(1,rho,rho,1),
    2,2)
)
```

```
rho <- 0.95
```



More on interpreting correlation

- Random variables with a correlation of ± 1 (or data with a sample correlation of ± 1) are **linearly dependent**.
- Random variables with a correlation of 0 (or data with a sample correlation of 0) are **uncorrelated**.
- Random variables with a covariance of 0 are also uncorrelated!

Question 5.3. Suppose two data vectors $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{y} = (y_1, \dots, y_n)$ have been **standardized**. That is, each data point has had the sample mean subtracted and then been divided by the sample standard deviation. You calculate $\text{cov}(\mathbf{x}, \mathbf{y}) = 0.8$. What is the sample correlation, $\text{cor}(\mathbf{x}, \mathbf{y})$?

The variance of a sum

- A basic property of covariance is

(Eq. C1)
$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$$

- Sample covariance has the same formula,

(Eq. C2)
$$\text{var}(\mathbf{x} + \mathbf{y}) = \text{var}(\mathbf{x}) + \text{var}(\mathbf{y}) + 2 \text{cov}(\mathbf{x}, \mathbf{y})$$

- These nice formulas mean it can be easier to calculate using variances and covariances rather than standard deviations and correlations.

Question 5.4. Rewrite (Eq. C1) to give a formula for $\text{SD}(X + Y)$ in terms of $\text{SD}(X)$, $\text{SD}(Y)$ and $\text{Cor}(X, Y)$.

More properties of covariance

- Covariance is not affected by adding constants to either variable

(Eq. C3)
$$\text{Cov}(X + a, Y + b) = \text{Cov}(X, Y)$$

- Recall the definition $\text{Cov}(X, Y) = E[(X - E[X]) (Y - E[Y])]$. In words, covariance is the mean product of deviations from average. These deviations are unchanged when we add a constant to the variable.

- Covariance scales **bilinearly** with each variable

(Eq. C3)
$$\text{Cov}(aX, bY) = ab \text{Cov}(X, Y)$$

- Covariance distributes across sums

(Eq. C4)
$$\text{Cov}(X, Y + Z) = \text{Cov}(X, Y) + \text{Cov}(X, Z)$$

- Sample covariances also have these properties. You can test them in R using bivariate normal random variables, constructed as previously using `'rmvnorm()'`.

The variance/covariance matrix of vector random variables

- Let $\mathbf{X} = (X_1, \dots, X_p)$ be a vector random variable. For any pair of elements, say X_i and X_j , we can compute the usual scalar covariance, $v_{ij} = \text{Cov}(X_i, X_j)$.
- The variance/covariance matrix $\mathbb{V} = [v_{ij}]_{p \times p}$ collects together all these covariances.

$$\mathbb{V} = \text{Var}(\mathbf{X}) = \begin{bmatrix} \text{Cov}(X_1, X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \text{Cov}(X_2, X_2) & & \text{Cov}(X_2, X_p) \\ \vdots & & \ddots & \vdots \\ \text{Cov}(X_p, X_1) & \text{Cov}(X_p, X_2) & \dots & \text{Cov}(X_p, X_p) \end{bmatrix}$$

- The diagonal entries of \mathbb{V} are $v_{ii} = \text{Cov}(X_i, X_i) = \text{Var}(X_i)$ for $i = 1, \dots, p$ so the variance/covariance matrix can be written as

$$\mathbb{V} = \text{Var}(\mathbf{X}) = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & & \text{Cov}(X_2, X_p) \\ \vdots & & \ddots & \vdots \\ \text{Cov}(X_p, X_1) & \text{Cov}(X_p, X_2) & \dots & \text{Var}(X_p) \end{bmatrix}$$

The correlation matrix

- Covariance is harder to interpret than correlation, but easier for calculations.
- We can put together all the correlations into a correlation matrix, using the fact that $\text{Cor}(X_i, X_i) = 1$.

$$\text{Cor}(\mathbf{X}) = \begin{bmatrix} 1 & \text{Cor}(X_1, X_2) & \dots & \text{Cor}(X_1, X_p) \\ \text{Cor}(X_2, X_1) & 1 & & \text{Cor}(X_2, X_p) \\ \vdots & & \ddots & \vdots \\ \text{Cor}(X_p, X_1) & \text{Cor}(X_p, X_2) & \dots & 1 \end{bmatrix}$$

- Multivariate distributions can be very complicated.
- The variance/covariance and correlation matrices deal only with **pairwise** relationships between variables.
- Pairwise relationships can be graphed.

The sample variance/covariance matrix

- The **sample variance/covariance matrix** places all the sample variances and covariances in a matrix.
- Let $\mathbb{X} = [x_{ij}]_{n \times p}$ be a data matrix made up of p data vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$ each of length n .

$$\text{var}(\mathbb{X}) = \begin{bmatrix} \text{var}(\mathbf{x}_1) & \text{cov}(\mathbf{x}_1, \mathbf{x}_2) & \dots & \text{cov}(\mathbf{x}_1, \mathbf{x}_p) \\ \text{cov}(\mathbf{x}_2, \mathbf{x}_1) & \text{var}(\mathbf{x}_2) & & \text{cov}(\mathbf{x}_2, \mathbf{x}_p) \\ \vdots & & \ddots & \vdots \\ \text{cov}(\mathbf{x}_p, \mathbf{x}_1) & \text{cov}(\mathbf{x}_p, \mathbf{x}_2) & \dots & \text{var}(\mathbf{x}_p) \end{bmatrix}$$

- R uses the same notation. If x is a matrix or dataframe, $\text{var}(x)$ returns the sample variance/covariance matrix.

```
var(x)
```

```
##           final      quiz      hw      midterm
## final  169.75184  78.14294  51.27143  75.61269
## quiz   78.14294  224.39664 103.57755 107.32550
## hw     51.27143  103.57755 120.13265  61.44694
## midterm 75.61269 107.32550  61.44694 218.21553
```

The sample correlation matrix

- The **sample correlation matrix** places all the sample correlations in a matrix.
- Let $\mathbb{X} = [x_{ij}]_{n \times p}$ be a data matrix made up of p data vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$ each of length n .

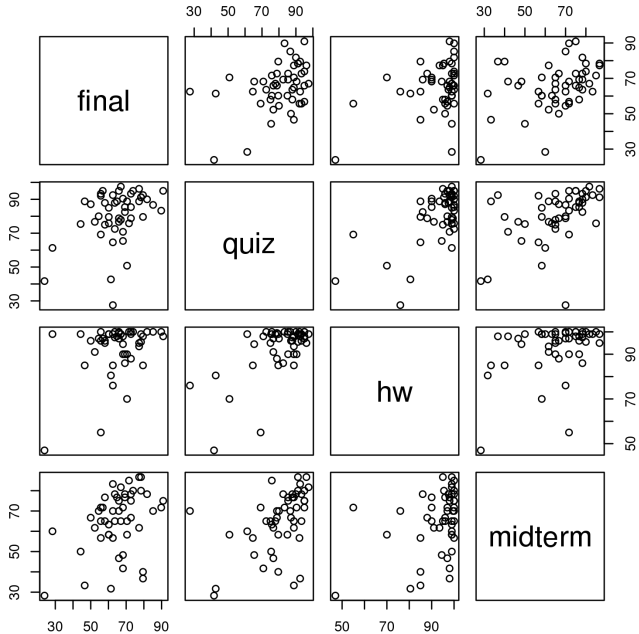
$$\text{cor}(\mathbb{X}) = \begin{bmatrix} 1 & \text{cor}(\mathbf{x}_1, \mathbf{x}_2) & \dots & \text{cor}(\mathbf{x}_1, \mathbf{x}_p) \\ \text{cor}(\mathbf{x}_2, \mathbf{x}_1) & 1 & & \text{cor}(\mathbf{x}_2, \mathbf{x}_p) \\ \vdots & & \ddots & \vdots \\ \text{cor}(\mathbf{x}_p, \mathbf{x}_1) & \text{cor}(\mathbf{x}_p, \mathbf{x}_2) & \dots & 1 \end{bmatrix}$$

- R uses the same notation. If x is a matrix or dataframe, $\text{cor}(x)$ returns the sample correlation matrix.

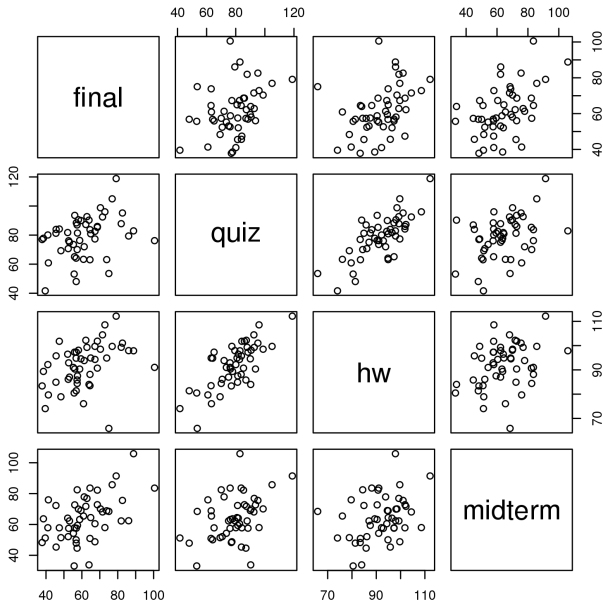
```
cor(x)
```

```
##           final      quiz      hw      midterm
## final      1.0000000  0.4003818  0.3590357  0.3928662
## quiz       0.4003818  1.0000000  0.6308512  0.4850114
## hw         0.3590357  0.6308512  1.0000000  0.3795132
## midterm    0.3928662  0.4850114  0.3795132  1.0000000
```

`pairs(x)`




```
mvn <- rmvnorm(50,mean=apply(x,2,mean),sigma=var(x))  
pairs(mvn)
```



Question 5.5. From looking at the scatterplots, what are the strengths and weaknesses of a multivariate normal model for test scores in this course?

Question 5.6. To what extent is it appropriate to summarize the data by the mean and variance/covariance matrix (or correlation matrix) when the normal approximation is dubious?

Linear combinations

- Let $\mathbf{X} = (X_1, \dots, X_p)$ be a vector random variable with mean vector $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)$ and $p \times p$ variance/covariance matrix \mathbb{V} .
- Let \mathbb{X} be a $n \times p$ data matrix.
- Let \mathbb{A} be a $q \times p$ matrix.
- $\mathbf{Z} = \mathbb{A}\mathbf{X}$ is a collection of q linear combinations of the p random variables in the vector \mathbf{X} , viewed as a **column** vector.
- $\mathbb{Z} = \mathbb{X}\mathbb{A}^T$ is an $n \times q$ collection of linear combinations of the p data points in each **row** of \mathbb{X} .
- Mental gymnastics are required: vectors are often interpreted as **column vectors** (e.g., $p \times 1$ matrices) but the vector of measurements for each unit is a **row vector** when considered as a row of an $n \times p$ data matrix.

Question 5.7. How would you construct a simulated data matrix \mathbb{Z}_{sim} from n realizations $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ of the random column vector $\mathbf{Z} = \mathbb{A}\mathbf{X}$? Be careful with transposes and keep track of dimensions.

Solution:

- There is a useful matrix variance/covariance formula for a linear combination.

$$\text{Var}(\mathbb{A} \mathbf{X}) = \mathbb{A} \text{Var}(\mathbf{X}) \mathbb{A}^T$$

$$\text{var}(\mathbb{X} \mathbb{A}^T) = \mathbb{A} \text{var}(\mathbb{X}) \mathbb{A}^T$$

Question 5.8. Add dimensions to each term in these equations to check they make sense.

Testing the variance/covariance formula

- Suppose that the overall course score is weighted 40% on the final and 20% on each of the midterm, homework and quiz.
- We can find the sample variance of the overall score two different ways.
 - (i) Directly computing the overall score for each student.

```
weights <- c(final=0.4,quiz=0.2,hw=0.2,midterm=0.2)
overall <- as.matrix(x) %*% weights
var(overall)
```

```
##           [,1]
## [1,] 104.2624
```

- (ii) Using $\text{var}(\mathbb{X} \mathbb{A}^T) = \mathbb{A} \text{var}(\mathbb{X}) \mathbb{A}^T$.

```
weights %*% var(x) %*% weights
```

```
##           [,1]
## [1,] 104.2624
```

- R interprets the vector 'weights' as a row or column vector as necessary.