

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/401f18/05/course_progress.txt")
```

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

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

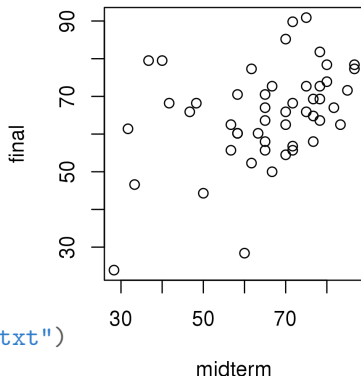
```
"1" 52.3 76.7 91 61.7
```

```
"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

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

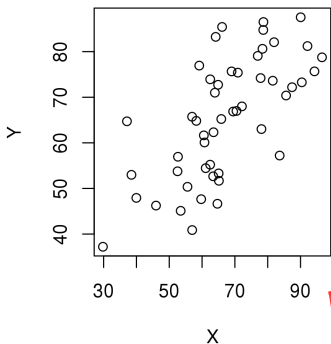
- The parameters of the bivariate normal distribution in matrix form are the **mean vector** $\mu = (\mu_X, \mu_Y)$ and the **variance/covariance matrix**,

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

install. packages ("mvtnorm") ; library("mvtnorm")

- In R, the mvtnorm package lets us simulate the bivariate and multivariate normal distribution via the `rmvnorm()` function. It has the mean vector and variance/covariance matrix as arguments.

Experimenting with the bivariate normal distribution



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

$$V = \begin{bmatrix} 200 & 100 \\ 100 & 150 \end{bmatrix} = \begin{bmatrix} \text{Var}(X) & \text{Cov}(X,Y) \\ \text{Cov}(Y,X) & \text{Var}(Y) \end{bmatrix}$$

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

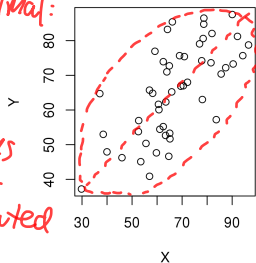
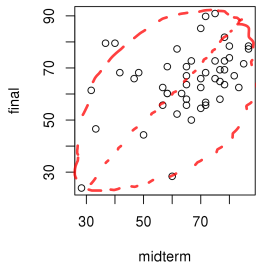
Question 5.1. What are μ_X , μ_Y , $\text{Var}(X)$, $\text{Var}(Y)$, and $\text{Cov}(X, Y)$ for this simulation?

$$\mu_X = 65, \mu_Y = 65, \text{Var}(X) = 200, \text{Var}(Y) = 150 \\ \text{Cov}(X, Y) = 100$$

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?

- Maybe the exam scores are clustered more in the top right?
- Maybe the exam scores are more variable; perhaps the distribution is longer tailed than normal?
- So far as the scatterplot does look normal: both are scattered around a line in something like a football shape.
- A central limit theorem here: both exams are the sum of many questions. The test-taking skill may not be normally distributed but the test uncertainty should be.



More on covariance

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

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

- 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

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

$\mathbf{X} = \begin{bmatrix} x \\ y \end{bmatrix}$
 $n \times 2$

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

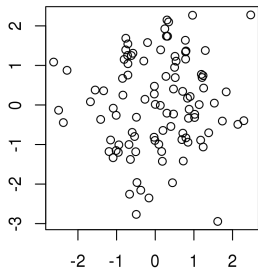
```
var(x$midterm)
## [1] 218.2155
var(x$final)
## [1] 169.7518
cov(x$midterm,x$final)
## [1] 75.61269
```

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

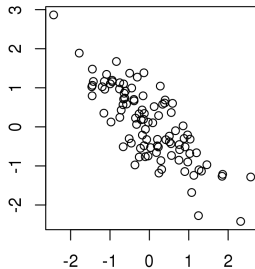
$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\text{SD}(X)\text{SD}(Y)} \text{cor}(\mathbf{x}, \mathbf{y}) = \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{\sqrt{\text{var}(\mathbf{x})\text{var}(\mathbf{y})}}$$

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

```
rho <- 0
```

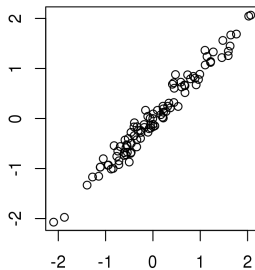


```
rho <- -0.8
```



```
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 called **linearly dependent** or **colinear**.
- 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, $\text{cor}(\mathbf{x}, \mathbf{y})$?

$$\text{cor}(\underline{x}, \underline{y}) = \frac{\text{cov}(\underline{x}, \underline{y})}{\underbrace{\text{sd}(\underline{x}) \text{sd}(\underline{y})}_{1 \times 1}} = \text{cov}(\underline{x}, \underline{y})$$

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 $\text{Cor}(X, Y)$.

$$\begin{aligned} [\text{SD}(X+Y)]^2 &= (\text{SD}(X))^2 + (\text{SD}(Y))^2 \\ &\quad + 2 \text{Cor}(X, Y) \text{SD}(X) \text{SD}(Y) \end{aligned}$$

$$\text{so, } \text{SD}(X+Y) = \sqrt{(\text{SD}(X))^2 + (\text{SD}(Y))^2 + 2 \text{Cor}(X, Y) \text{SD}(X) \text{SD}(Y)}$$

More properties of covariance

* i.e. for (X, Y, Z) a vector random variable.

- Covariance is not affected by adding constants to either variable

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

- 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

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

- Covariance distributes across sums * for X, Y, Z joint random variables

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

- Sample covariances also have these properties. You can test them in R using bivariate normal random variables, constructed as previously using 'rmvnorm()'. e.g. $\underline{x} = (x_1, \dots, x_n)$, $\underline{y} = (y_1, \dots, y_n)$, $\underline{z} = (z_1, \dots, z_n)$ three vectors with the same length n .

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, $v_{ij} = \text{Cov}(X_i, X_j)$. *sometimes called the variance matrix or the covariance matrix.*
- The variance/covariance matrix $\mathbb{V} = [v_{ij}]_{p \times p}$ collects together all these covariances.

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

$p \times p$ $p \times 1$

- 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 . *(check the case $p=2$ for intuition)*

(0 0)
not
(0 0)

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

n x p
p x p

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

`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 |

elements of a matrix are the (1,1), (2,2), ..., (p,p) entries which give $\text{var}(x_1), \dots, \text{var}(x_p)$.

The sample correlation matrix

- The **sample correlation matrix** places all the sample correlations in a matrix. $\text{cor}(X, X) = \frac{\text{cov}(X, X)}{\sqrt{\text{Var}(X) \cdot \text{Var}(X)}} = \frac{\text{Var}(X)}{\text{Var}(X)} = 1$. This is also true for $\text{cor}(\underline{x}, \underline{x})$
- 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)
```

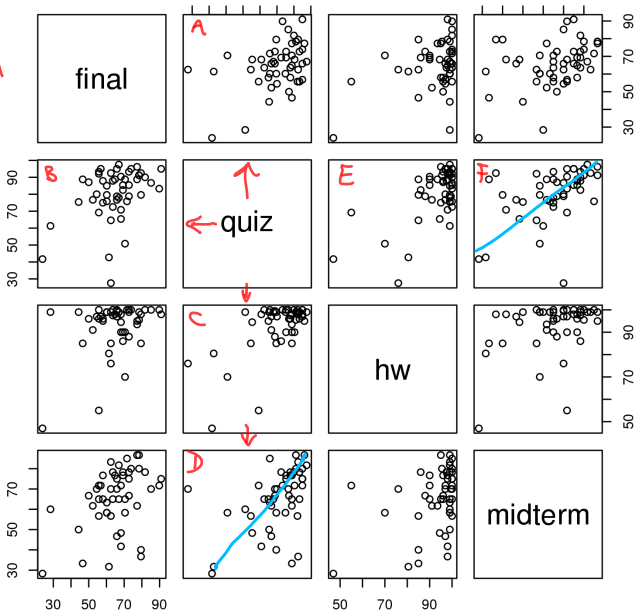
```
##           final      quiz      hw      midterm
## final    1.0000000  0.4003818  0.3590357  0.3928662
## quiz     0.4003818  1.0000000  0.6308512  0.4850114
## hw       0.3590357  0.6308512  1.0000000  0.3795132
## midterm  0.3928662  0.4850114  0.3795132  1.0000000
```

pairs(x)

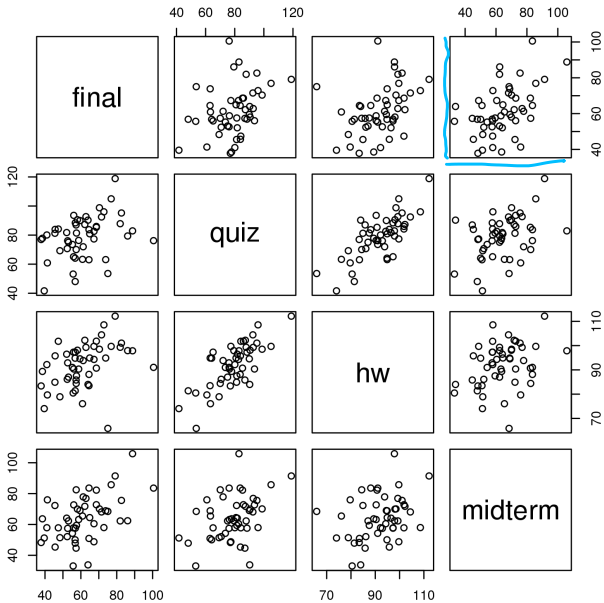
(A) has quiz on the x axis, final on the y axis.
(B) has quiz on the y axis, final on the x axis.

(A, C, D) all have quiz on the x-axis

(B, E, F) all have quiz on the y-axis




```
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? *Specifically, why does HW not follow a normal curve?*

Difficulty of the assignment? Office hrs?

Maybe HW tests attention to detail?

This means the central limit theorem is not holding.

Why not? Questions not independent - all too similar?

Each HW is all or nothing, so there are not a large number.

HW score is driven by a small number of rare events - missed HW.

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 ^{multivariate} approximation is dubious?

The normal distribution is fully determined by the mean & variance/covariance matrix. A football-shaped scatterplot is well described by the sample mean & sample variance/covariance matrix. If the scatterplot is not football shaped, there is information lost by summarizing it as mean & variance/covariance.

Linear combinations

here, $\underline{x} = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{bmatrix}$ is a column vector, written this way to match $\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$

- 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 \mathbf{V} .
- Let \mathbf{X} be a $n \times p$ data matrix.
- Let \mathbf{A} be a $q \times p$ matrix.
- $\mathbf{Z} = \mathbf{A}\mathbf{X}$ is a collection of q linear combinations of the p random variables in the vector \mathbf{X} , viewed as a **column** vector.
- $\mathbf{Z} = \mathbf{X}\mathbf{A}^T$ is an $n \times q$ collection of linear combinations of the p data points in each **row** of \mathbf{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.

one row of \mathbf{Z} is $\{x_{i1} \ x_{i2} \ \dots \ x_{ip}\} \mathbf{A}^T = \underline{x}^T \mathbf{A}^T$

transposing this gives a column vector $(\underline{x}^T \mathbf{A}^T)^T = \mathbf{A} \underline{x}$.

we are going to be interested in $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ for a probability model for the linear model.

this will be a fixed matrix \mathbf{A}

this will be our random variable \mathbf{X}

Variables of length p as column vectors or row vectors

Question 5.7. How would you construct a simulated data matrix \mathbf{Z}_{sim} from n realizations $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ of the random column vector $\mathbf{Z} = \mathbf{A}\mathbf{X}$? Hint: You are expected to write a matrix constructing \mathbf{Z}_{sim} by stacking together $\mathbf{Z}_1, \dots, \mathbf{Z}_n$. Be careful with transposes and keep track of dimensions. Recall that \mathbf{Z}_{sim} should be $n \times p$.

$$\mathbf{Z}_{\text{sim}} = \begin{bmatrix} \mathbf{Z}_1^T \\ \mathbf{Z}_2^T \\ \vdots \\ \mathbf{Z}_n^T \end{bmatrix}$$

Expectation and variance of linear combinations

- The expectation of a vector random variable is the vector of expectations of each element. If $\mathbf{X} = (X_1, \dots, X_n)$ then

$$E[\mathbf{X}] = (E[X_1], E[X_2], \dots, E[X_n])$$

- The expectation of a sum is the sum of the expectations.

$$E[X + Y] = E[X] + E[Y]$$

- This formula extends to any linear combination of n random variables. If $\mathbf{Z} = \mathbb{A}\mathbf{X}$ then $E[\mathbf{Z}] = \mathbb{A}E[\mathbf{X}]$. In other words,

$$E[\mathbb{A}\mathbf{X}] = \mathbb{A}E[\mathbf{X}]$$

- There is a useful matrix variance/covariance formula for a linear combination, which also works for sample variance/covariance.

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

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

*intuition:
expectation is*

*a long run
average.*

*averages have
this additive
property.*

Exercises with the matrix variance/covariance formula

Question 5.8. Add dimensions to each quantity in the equations

$$\text{Var}(\mathbf{A}\mathbf{X}) = \mathbf{A} \text{Var}(\mathbf{X}) \mathbf{A}^T \quad \text{and} \quad \text{var}(\mathbf{X}\mathbf{A}^T) = \mathbf{A} \text{var}(\mathbf{X}) \mathbf{A}^T.$$

$$\begin{aligned} \text{Var}(aX) &= a^2 \text{Var}(X), \\ \text{SD}(aX) &= a \text{SD}(X) \end{aligned}$$

$$\begin{array}{c} \begin{array}{|c|c|} \hline \overbrace{q \times p} & \overbrace{p \times p} \\ \hline \end{array} & \overbrace{p \times q} \\ \hline \underbrace{q \times p} & \underbrace{p \times p} \\ \hline \end{array} \quad \begin{array}{c} \overbrace{q \times p} \\ \overbrace{p \times p} \\ \hline \underbrace{q \times q} \end{array}$$

$$\begin{array}{c} \overbrace{q \times p} \\ \overbrace{p \times q} \\ \hline \underbrace{q \times q} \end{array} \quad \begin{array}{c} \overbrace{q \times p} \\ \overbrace{p \times p} \\ \overbrace{p \times q} \\ \hline \underbrace{q \times q} \end{array}$$

No. \tilde{X} should be a vector. Usually \tilde{X} is a column vector

Question 5.9. Let $\mathbf{A} = [1 \dots 1]$ be the $1 \times p$ row vector of ones. Let $\mathbf{X} = (X_1, \dots, X_p)$ be a vector random variable with variance/covariance matrix $\mathbf{V} = [V_{ij}]_{p \times p}$. Evaluate the variance/covariance formula in this case. Hence, find $\text{Var}(\bar{X})$ where $\bar{X} = (1/p) \sum_{i=1}^p X_i$.

What is $\mathbf{A}\tilde{X}$? $\mathbf{A}\tilde{X} = (1 \ 1 \ \dots \ 1) \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{pmatrix} = \sum_{i=1}^p X_i$

this multiplication sums over rows in each column
this multiplication sums columns in each row

What is $\text{Var}(\mathbf{A}\tilde{X})$? Using the variance/covariance formula,

$$\text{Var}(\mathbf{A}\tilde{X}) = (1 \ 1 \ \dots \ 1) \begin{pmatrix} V_{11} & V_{12} & \dots & V_{1p} \\ \vdots & \vdots & \dots & \vdots \\ V_{p1} & \dots & \dots & V_{pp} \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \sum_{i=1}^p \sum_{j=1}^p V_{ij}$$

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{p} \mathbf{A}\tilde{X}\right) = \frac{1}{p^2} \text{Var}(\mathbf{A}\tilde{X}) = \frac{1}{p^2} \sum_{i=1}^p \sum_{j=1}^p V_{ij}$$

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
```

50×4 4×1
 50×1

as.matrix turns the dataframe x into a matrix.

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

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

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

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

Independence

- Two events E and F are **independent** if

$$P(E \text{ and } F) = P(E) \times P(F)$$

Worked example 5.1. Suppose we have a red die and a blue die. They are ideal fair dice, so the values should be independent. What is the chance they both show a six?

- (a) Using the definition of independence.

$$\begin{aligned} & P(\{\text{red die shows } 6\} \text{ and } \{\text{blue die shows } 6\}) \\ &= P(\text{red die shows } 6) \times P(\text{blue die shows } 6) \quad : \text{ independence} \\ &= \frac{1}{6} \times \frac{1}{6} \quad \text{since each die is fair} \quad = \frac{1}{36} \end{aligned}$$

- (b) By considering equally likely outcomes, without using the definition.

There are 36 possible pairs of outcomes (red die, blue die). Assuming equally likely outcomes, each pair (red, blue) has probability $1/36$, so $P(\text{both dice show } 6) = \frac{1}{36}$.

- The multiplication rule agrees with an intuitive idea of independence.

Independence of random variables

- X and Y are **independent random variables** if, for any intervals $[a, b]$ and $[c, d]$,

$$P(a < X < b \text{ and } c < Y < d) = P(a < X < b) \times P(c < Y < d)$$

- This definition extends to vector random variables. $\mathbf{X} = (X_1, \dots, X_n)$ is a **vector of independent random variables** if for any collection of intervals $[a_i, b_i]$, $1 \leq i \leq n$,

$$P(a_1 < X_1 < b_1, \dots, a_n < X_n < b_n) = P(a_1 < X_1 < b_1) \times \dots \times P(a_n < X_n < b_n)$$

- $\mathbf{X} = (X_1, \dots, X_n)$ is a **vector of independent identically distributed (iid) random variables** if, in addition, each element of \mathbf{X} has the same distribution.
- “ X_1, \dots, X_n are n random variables with the **normal**(μ, σ) distribution” is written more formally as
“Let $X_1, \dots, X_n \sim \text{iid normal}(\mu, \sigma)$.”

Independent vs uncorrelated

a property of random variables, not data

a property of random variables, with a comparable sample version.

- If X and Y are independent they are uncorrelated.
- The converse is not necessarily true.
- **For normal random variables, the converse is true.**
- If X and Y are bivariate normal random variables, and $\text{Cov}(X, Y) = 0$, then X and Y are independent.
- The following slide demonstrated the possibility of being uncorrelated but not independent (for non-normal random variables).
- If the scatter plot of two variables looks normal and their sample correlation is small, the variables are appropriately modeled as independent.

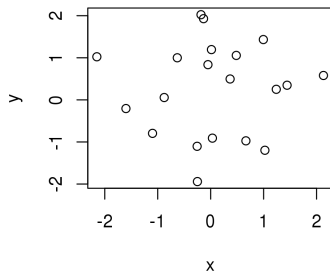
Zero correlation with and without independence

```
x <- rnorm(20)
y <- rnorm(20)
```

```
cor(x,y)
```

```
## [1] 0.01825057
```

```
plot(x,y)
```

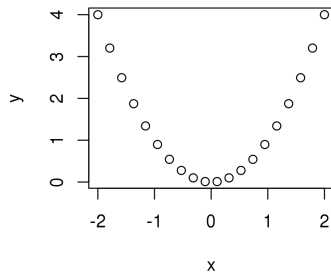


```
x <- seq(-2,2,length=20)
y <- x^2
```

```
cor(x,y)
```

```
## [1] -1.704156e-16
```

```
plot(x,y)
```



The measurement error model

- Let $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)$ be a vector consisting of n independent normal random variables, each with mean zero and variance σ^2 . (Cov(ϵ_i, ϵ_j) = 0 if $i \neq j$.)
- This is a common model for **measurement error** on n measurements.
- The mean vector and variance/covariance matrix are

$$\underline{E[\epsilon] = \mathbf{0}},$$

$$\underline{\text{Var}(\epsilon) = \sigma^2 \mathbb{I}}$$

$$= \begin{pmatrix} \sigma^2 & 0 & & 0 \\ 0 & \sigma^2 & & \\ & & \ddots & \\ 0 & & & \sigma^2 \end{pmatrix}$$

where $\mathbf{0} = (0, \dots, 0)$ and \mathbb{I} is the $n \times n$ identity matrix.

- For the measurement error model, we break our usual rule of using upper case letters for random variables.
- We can write $\epsilon \sim \text{MVN}(\mathbf{0}, \sigma^2 \mathbb{I})$ or $\epsilon \sim \text{iid normal}(0, \sigma)$.
- ϵ models **unbiased** measurements (meaning $E[\epsilon_i] = 0$) subject to **additive Gaussian error**.

Example: 5 repeated measurements x_1, \dots, x_5 of the speed of light could be modeled as $X_i = \mu + \epsilon_i$ for $i = 1, \dots, 5$, where μ is the unknown true value of this quantity.

A probability model for the linear model

- First recall the sample version of the linear model, which is

*Constant,
not random.*

$$\mathbf{y} = \mathbb{X} \mathbf{b} + \mathbf{e},$$

*vector random
variable*

where \mathbf{y} is the measured response, \mathbb{X} is an $n \times p$ matrix of explanatory variables, \mathbf{b} is chosen by least squares, and \mathbf{e} is the vector of residuals.

- We want to build a random vector \mathbf{Y} that provides a probability model for the data \mathbf{y} . We write this as

*vector
random
variable*

$$\mathbf{Y} = \mathbb{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}$$

*A probability model
in principle defines
the probability of
everything in the model.*

where \mathbb{X} is the explanatory matrix, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$ is an unknown coefficient vector (we don't know the true population coefficient!) and $\boldsymbol{\epsilon}$ is Gaussian measurement error with $E[\boldsymbol{\epsilon}] = \mathbf{0}$ and $\text{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbb{I}$.

- Our probability model asserts that the process which generated the response data \mathbf{y} was like drawing a random vector \mathbf{Y} constructed using a random measurement error model with known matrix \mathbb{X} for some fixed but unknown value of $\boldsymbol{\beta}$. *we treat \mathbb{X} as a constant matrix.*

A digression on “useful” models

“Now it would be very remarkable if any system existing in the real world could be exactly represented by any simple model. However, cunningly chosen parsimonious models often do provide remarkably useful approximations. For example, the law $PV = RT$ relating pressure P , volume V and temperature T of an *ideal* gas via a constant R is not exactly true for any real gas, but it frequently provides a useful approximation and furthermore its structure is informative since it springs from a physical view of the behavior of gas molecules. For such a model there is no need to ask the question ‘Is the model true?’. If *truth* is to be the *whole truth* the answer must be *No*. The only question of interest is ‘Is the model illuminating and useful.’ ” (Box, 1978)

“Essentially, all models are wrong, but some are useful.”

(Box and Draper, 1987)

- Perhaps the most useful statistical model ever is $\mathbf{Y} = \mathbb{X}\beta + \epsilon$.
- Anything so widely used is also widely abused. Our task is to understand $\mathbf{Y} = \mathbb{X}\beta + \epsilon$ so that we can use it wisely.

Expectation and variance/covariance of \mathbf{Y}

- 05 $E[\mathbf{X}\beta + \epsilon] = E[\mathbf{X}\beta] + E[\epsilon] = \mathbf{X}\beta + \mathbf{0} = \mathbf{X}\beta$
- Recall the linear model $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ where $\epsilon \sim \text{MVN}(\mathbf{0}, \sigma^2 \mathbf{I})$. $\beta = (\beta_1, \dots, \beta_p)$

Question 5.10. What is the expected value, $E[\mathbf{Y}]$?

$E[\mathbf{Y}]$ is a constant vector of length n . The i th element of $E[\mathbf{Y}]$ is $E[\sum_{j=1}^p x_{ij}\beta_j + \epsilon_i]$

since $Y_i = \sum_{j=1}^p x_{ij}\beta_j + \epsilon_i$ where $\mathbf{X} = [x_{ij}]$

$$\text{so, } E[Y_i] = \sum_{j=1}^p x_{ij}\beta_j + E[\epsilon_i]$$

$$\text{so, } E[Y_i] = \sum_{j=1}^p x_{ij}\beta_j$$

in matrix notation,

$$E[\mathbf{Y}] = \mathbf{X}\beta$$

since expectation of a sum is the sum of expectations, and expectation of a constant is constant

Question 5.11. What is the variance/covariance matrix, $\text{Var}(\mathbf{Y})$?

$$\text{Var}(\mathbf{Y}) = \text{Var}(\mathbf{X}\beta + \epsilon) = \text{Var}(\epsilon) = \sigma^2 \mathbf{I}$$

this is a constant vector, by assumption, and adding a constant doesn't change variance

by assumption of the measurement error model.

[note: assuming \mathbf{X} is constant is like conditioning on \mathbf{X} if \mathbf{X} is random - but we don't pursue this idea]

Expectation of the least squares coefficient

this is a random vector

Worked example 5.2. Let $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$. This is the probability model for the sample least squares coefficient $\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. Use linearity to calculate $E[\hat{\beta}]$.

Expectation is linear means that it preserves linear relationships. if

$$Y = aX + b,$$

$$E[Y] = a E[X] + b$$

Solution:

$$\begin{aligned} E[\hat{\beta}] &= E[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E[\mathbf{Y}] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta \\ &= \beta \end{aligned}$$

from previous slide

- Interpretation: If the data \mathbf{y} are well modeled as a draw from the probability model $\mathbf{Y} = \mathbf{X}\beta + \epsilon$, then the least squares estimate \mathbf{b} is well modeled by a random vector with mean β .

Variance/covariance matrix of the least squares coefficients

Question 5.12. Consider the linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ with $E[\boldsymbol{\epsilon}] = \mathbf{0}$ and $\text{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$. Use the formula $\text{Var}(\mathbf{A}\mathbf{Y}) = \mathbf{A}\text{Var}(\mathbf{Y})\mathbf{A}^T$ to find $\text{Var}(\hat{\boldsymbol{\beta}})$.

Hint: what should be our choice of \mathbf{A} so that $\boldsymbol{\beta} = \mathbf{A}\mathbf{Y}$?

scalar
the constant
 σ^2 comes out
of the product.
The matrix \mathbf{I} disappears
since $\mathbf{I}\mathbf{X} = \mathbf{X}$

$$\begin{aligned} & (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} \\ &= [(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})] (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \mathbf{I} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{\boldsymbol{\beta}}) &= \text{Var}(\boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}) \\ &= \text{Var}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned}$$

using $\mathbf{A} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
 $\mathbf{A}^T = \mathbf{X} [(\mathbf{X}^T \mathbf{X})^{-1}]^T$
 $= \mathbf{X} [(\mathbf{X}^T \mathbf{X})^{-1}]$
 $= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$

$\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}$ plus some error due to the measurement process. \rightarrow

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ \hat{\boldsymbol{\beta}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}) \\ &= \boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \end{aligned}$$

Standard errors for coefficients in the linear model

- The formula $\text{Var}(\hat{\beta}) = \sigma^2 (\mathbb{X}^T \mathbb{X})^{-1}$ needs extra work to be useful for data analysis.
- In practice, we know the model matrix \mathbb{X} but we don't know the measurement standard deviation σ .
- An estimate of the measurement standard deviation is the sample standard deviation of the residuals.

- For $\mathbf{y} = \mathbb{X}\mathbf{b} + \mathbf{e}$ with \mathbb{X} being $n \times p$, an estimate of σ is

$$s = \sqrt{\frac{1}{n-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \sqrt{\frac{1}{n-p} \sum_{i=1}^n (y_i - [\mathbb{X}\mathbf{b}]_i)^2}$$

- We will discuss later why we choose to divide by $n - p$.
- The **standard error** of b_k for $k = 1, \dots, p$ is

$$\text{SE}(b_k) = s \sqrt{[(\mathbb{X}^T \mathbb{X})^{-1}]_{kk}}$$

which is an estimate of $\sqrt{[\text{Var}(\hat{\beta})]_{kk}}$.

- These standard errors are calculated by `lm()` in R.

kth diagonal entry of $(\mathbb{X}^T \mathbb{X})^{-1}$
e.g. if $(\mathbb{X}^T \mathbb{X})^{-1} = [V_{ij}]_{p \times p}$, the k^{th} diagonal entry is V_{kk} .

```
lm1 <- lm(L_detrended~U_detrended) ; summary(lm1)
```

```
##
```

```
## Call:
```

```
## lm(formula = L_detrended ~ U_detrended)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -1.55654 -0.48641 -0.01867  0.40856  1.63118
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.28999    0.09343   3.104  0.00281 **
## U_detrended  0.13137    0.06322   2.078  0.04161 *
## ---
```

```
## Signif. codes:
```

```
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
```

```
## Residual standard error: 0.7705 on 66 degrees of freedom
## Multiple R-squared:  0.06141, Adjusted R-squared:  0.04718
## F-statistic: 4.318 on 1 and 66 DF, p-value: 0.04161
```

How does R obtain linear model standard errors?

- The previous slide shows output from our analysis of unemployment and mortality from Chapter 1.
- Let's first extract the estimates and their standard errors from R, a good step toward reproducible data analysis.

```
names(summary(lm1))
```

```
## [1] "call"          "terms"          "residuals"  
## [4] "coefficients"  "aliases"        "sigma"  
## [7] "df"           "r.squared"      "adj.r.squared"  
## [10] "fstatistic"    "cov.unscaled"
```

```
summary(lm1)$coefficients
```

```
##           Estimate Std. Error  t value    Pr(>|t|)  
## (Intercept) 0.2899928 0.09343146  3.103802 0.002812739  
## U_detrended 0.1313673 0.06321939  2.077959 0.041606370
```

Extracting the design matrix

```
X <- model.matrix(lm1)
head(X)

##      (Intercept) U_detrended
## 16              1 -1.0075234
## 17              1  1.1027941
## 18              1  0.4881116
## 19              1 -1.5349043
## 20              1 -1.8662535
## 21              1 -2.0059360
```

Computing the standard errors directly

```
s <- sqrt(sum(resid(lm1)^2)/(nrow(X)-ncol(X))) ; s  
  
## [1] 0.7704556  
  
V <- s^2 * solve(t(X)%*%X)  
sqrt(diag(V))  
  
## (Intercept) U_detrended  
## 0.09343146 0.06321939
```

- This matches the standard errors generated by `lm()`.

```
summary(lm1)$coefficients  
  
##           Estimate Std. Error  t value    Pr(>|t|)  
## (Intercept) 0.2899928 0.09343146  3.103802 0.002812739  
## U_detrended 0.1313673 0.06321939  2.077959 0.041606370
```

Extracting the coefficient variance/covariance matrix

- The fitted `lm` object in R stores the estimated variance/covariance matrix for the coefficients in the output of `summary()`.

```
s <- summary(lm1)$sigma
XX <- summary(lm1)$cov.unscaled
s^2 * XX

##                (Intercept) U_detrended
## (Intercept) 0.008729439 0.000000000
## U_detrended 0.000000000 0.003996692
```

- This matches what we get from calculating $s^2(\mathbf{X}^T\mathbf{X})^{-1}$ directly.

```
X <- model.matrix(lm1)
sum(resid(lm1)^2)/(nrow(X)-ncol(X)) * solve(t(X)%*%X)

##                (Intercept)  U_detrended
## (Intercept) 8.729439e-03 1.305064e-20
## U_detrended 1.305064e-20 3.996692e-03
```