Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
00000000

**Imperial College London**

# Lecture 11: Introduction to Linear Models
## Statistical Modelling I

Dr. Riccardo Passeggeri

Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
00000000

## Last time

**Lectures 1**-**10**: focus on methods for inference in samples that are iid

Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
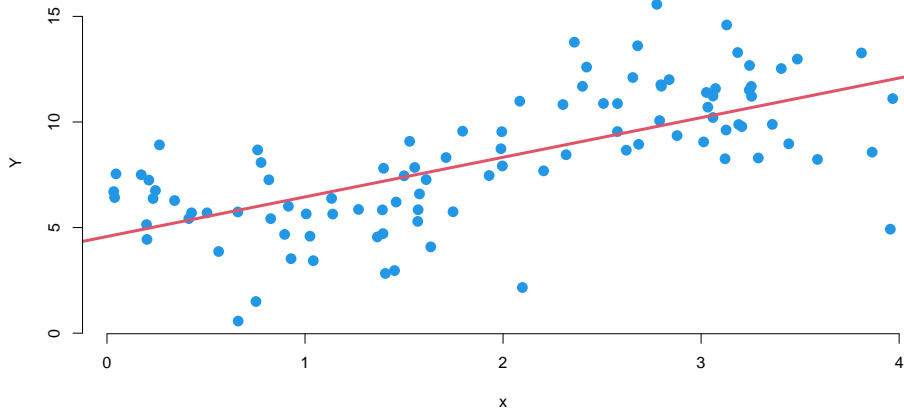00000000

## Outline

1. Introduction

2. Matrix Algebra

3. Expectations of Random Vectors

# Introduction

## Why linear models?

## Definition: Simple Linear Model

$$Y_i = \beta_1 + x_i\beta_2 + \epsilon_i, \quad i = 1, \ldots, n$$

- ► $Y_i$ "outcome", "response"; observable random variable.
- ► $x_i$ "covariate"; observable constant.
- ► $\beta_1$, $\beta_2$ unknown parameters.
- ► Error $\epsilon_1, \ldots, \epsilon_n$ iid, $E(\epsilon_i) = 0$, $Var(\epsilon_i) = \sigma^2$ for $i = 1, \ldots, n$.
- ► $\sigma^2 > 0$ is another unknown parameter.
- ► The errors $\epsilon_1, \ldots, \epsilon_n$ are not observable.

## Least squares estimators

The *least squares estimators* $\hat{\beta}_1$, $\hat{\beta}_2$ of $\beta_1$ and $\beta_2$ are defined as the minimisers of

$$S(\beta_1, \beta_2) = \sum_{i=1}^{n}(y_i - \beta_1 - x_i\beta_2)^2.$$

Note that:

- $e_i = y_i - \hat{\beta}_1 - x_i\hat{\beta}_2$, the so-called residuals, are observable. They are not iid, as dependence is introduced via $\hat{\beta}_1$, $\hat{\beta}_2$.
- The unknown parameters are $\beta_1$, $\beta_2$ and $\sigma^2$.
- In linear regression models $Y_1, \ldots, Y_n$ are generally not iid observations. Independence will still hold if the errors $\epsilon_1, \ldots, \epsilon_n$ are independent. However, the $Y_i$ do not have the same distribution; the distribution of $Y_i$ depends on the covariate $x_i$.

Introduction
○○○○

Matrix Algebra
●○○○○

Expectations of Random Vectors
○○○○○○○○

# Matrix Algebra

Introduction
0000

Matrix Algebra
0●000

Expectations of Random Vectors
00000000

## A toolkit for linear algebra

Linear regression naturally leads to a connection between statistics and linear algebra

This lecture, we highlight some useful results about matrices.

$A^T$ denotes the transpose of a matrix. I will use the terms "invertible" and "non-singular" synonymously.

**Matrix transposition, multiplication and inversion:**
- ▶ Let $A \in \mathbb{R}^{n \times m}, B \in \mathbb{R}^{m \times n}$. Then $(AB)^T = B^T A^T$
- ▶ Let $A \in \mathbb{R}^{n \times n}$ be non-singular. Then $(A^{-1})^T = (A^T)^{-1}$.

Introduction
0000

Matrix Algebra
00●00

Expectations of Random Vectors
00000000

## Transpose and trace

**(Trace)** Let $A = (A_{ij}) \in \mathbb{R}^{n \times n}$. Then

$$\text{trace}(A) = \sum_{i=1}^{n} A_{ii}$$

**Lemma.** $\text{trace}(AB) = \text{trace}(BA)$.

**Proof.** Recall that $AB = (\sum_j A_{ij} B_{jk})_{i,k}$. Thus, we have that

$$\text{trace}(AB) = \sum_i \sum_j A_{ij} B_{ji} = \sum_j \sum_i B_{ji} A_{ij} = \text{trace}(BA).$$

**Example.** Let $A = (1, 1)$, $B = (1, 1)^T$. Then $AB = 2 \neq \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} = BA$, but $\text{trace}(AB) = 2 = \text{trace}(BA)$.

Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
00000000

# Rank of $X^T X$

Let $X$ be an $n \times p$ matrix. Then $\operatorname{rank}(X^T X) = \operatorname{rank}(X)$.

**Proof.** Let $\operatorname{kern}(X) = \{x \in \mathbb{R}^p : Xx = 0\}$. Then $p = \operatorname{rank} X + \dim \operatorname{kern}(X)$. Similarly, $p = \operatorname{rank} X^T X + \dim \operatorname{kern}(X^T X)$

It suffices to show: $\operatorname{kern}(X) = \operatorname{kern}(X^T X)$.

If $x \in \operatorname{kern}(X)$ then $0 = Xx$ and hence $0 = X^T Xx$ which shows $x \in \operatorname{kern}(X^T X) = \{y : X^T Xy = 0\}$.

If $x \in \operatorname{kern}(X^T X)$ then $0 = X^T Xx$ and thus

$$0 = x^T X^T Xx = (Xx)^T Xx = \|Xx\|^2$$

which shows $Xx = 0$, i.e. $x \in \operatorname{kern}(X)$.

Introduction
0000

Matrix Algebra
0000●

Expectations of Random Vectors
00000000

A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is **positive definite** if

$$\forall \mathbf{x} \in \mathbb{R}^n \setminus \{0\} : \mathbf{x}^T A \mathbf{x} > 0.$$

**Lemma.** $A \in \mathbb{R}^{n \times n}$ is symmetric $\implies$ $\exists$ orthogonal matrix $P$ (i.e. $P^T P = I$) s.t. $P^T A P$ is diagonal (with diagonal entries equal to the eigenvalues of $A$).

$A$ an $n \times n$ positive definite symmetric matrix $\implies$ $\exists$ non-singular matrix $Q$ s.t. $Q^T A Q = I_n$.

### Proof.

First part is a standard linear algebra result.
The second result can be derived from it: $A$ p.d. $\implies$ its eigenvalues are $> 0$.
Hence, $P^T A P = D = \text{diag}(\lambda_1, \ldots, \lambda_n)$ where $\lambda_i > 0 \, \forall i$.
Let $E = D^{\frac{1}{2}} = \text{diag}(\sqrt{\lambda_1}, \ldots, \sqrt{\lambda_n})$ and define $Q = PE^{-1}$. Then

$$Q^T A Q = (PE^{-1})^T APE^{-1} = (E^{-1})^T P^T APE^{-1} = (E^{-1})^T EEE^{-1} = I.$$

$\square$

Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
●0000000

# Expectations of Random Vectors

Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
0●000000

## Why do we need expectations of random vectors?

Linear regression models describe the relationship between $Y$ and $x$ based on $E(Y \mid x)$.

The parameter *vector* $(\beta_0, \beta_1)$ suggests there may be correlation between least squares estimators.

Let $\boldsymbol{X} = (X_1, \ldots, X_n)^T$ be a random vector.

Then

$$E(\boldsymbol{X}) = (E\,X_1, \ldots, E\,X_n)^T,$$

i.e. the expectation is defined componentwise. For random matrices the expectation is also defined componentwise.

Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
00●00000

## Lemma: Linearity of expectations

Let $\boldsymbol{X}$ and $\boldsymbol{Y}$ be n-variate random vectors. Then the following hold:

- $E(\boldsymbol{X} + \boldsymbol{Y}) = E\,\boldsymbol{X} + E\,\boldsymbol{Y}$.
- Let $a \in \mathbb{R}$ then $E(a\boldsymbol{X}) = a\,E(\boldsymbol{X})$
- Let $A$, $B$ be deterministic matrices of "suitable dimensions" (deterministic means that they are not random). Then $E(A\boldsymbol{X}) = A\,E(\boldsymbol{X})$ and $E(\boldsymbol{X}^T B) = E(\boldsymbol{X})^T B$.

**Proof.** Use properties of one-dimensional random variables, for example

$$E(A\boldsymbol{X}) = (E(\sum_j A_{ij} X_j))_i = (\sum_j A_{ij}\,E(X_j))_i = A\,E(\boldsymbol{X}).$$

Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
000●0000

## Covariance of random vectors

If $\boldsymbol{X}$, $\boldsymbol{Y}$ are random vectors then

$$\begin{aligned}
\text{cov}(\boldsymbol{X}, \boldsymbol{Y}) :=& (\text{cov}(X_i, Y_j))_{i,j} \\
=& \, \mathsf{E}[(\boldsymbol{X} - \mathsf{E}(\boldsymbol{X}))(\boldsymbol{Y} - \mathsf{E}(\boldsymbol{Y}))^T] = \mathsf{E}[\boldsymbol{X}\boldsymbol{Y}^T] - \mathsf{E}(\boldsymbol{X})\,\mathsf{E}(\boldsymbol{Y})^T.
\end{aligned}$$

Furthermore $\text{cov}(\boldsymbol{X}) := \text{cov}(\boldsymbol{X}, \boldsymbol{X})$.

Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
0000●000

## Lemma: Covariance properties

If $\boldsymbol{X}$, $\boldsymbol{Y}$ and $\boldsymbol{Z}$ are random vectors, $A$, $B$ are deterministic matrices and $a, b \in \mathbb{R}$ are constants then (assuming appropriate dimensions)

- $\operatorname{cov}(\boldsymbol{X}, \boldsymbol{Y}) = \operatorname{cov}(\boldsymbol{Y}, \boldsymbol{X})^T$
- $\operatorname{cov}(a\boldsymbol{X} + b\boldsymbol{Y}, \boldsymbol{Z}) = a\operatorname{cov}(\boldsymbol{X}, \boldsymbol{Z}) + b\operatorname{cov}(\boldsymbol{Y}, \boldsymbol{Z})$
- $\operatorname{cov}(A\boldsymbol{X}, B\boldsymbol{Y}) = A\operatorname{cov}(\boldsymbol{X}, \boldsymbol{Y})B^T$
- $\operatorname{cov}(A\boldsymbol{X}) = A\operatorname{cov}(\boldsymbol{X})A^T$
- $\operatorname{cov}(\boldsymbol{X})$ is positive semidefinite and symmetric,
  i.e. $\boldsymbol{c}^T \operatorname{cov}(\boldsymbol{X})\boldsymbol{c} \geq 0$ for all vectors $\boldsymbol{c}$, or, equivalently, all eigenvalues of $\operatorname{cov}(\boldsymbol{X})$ are nonnegative.
- If $\boldsymbol{X}$ and $\boldsymbol{Y}$ are independent then $\operatorname{cov}(\boldsymbol{X}, \boldsymbol{Y}) = 0$.

**Proof.** Work from properties of one-dimensional covariance or work with one of the vector definitions of the covariance.

Introduction
0000

Matrix Algebra
00000

Expectations of Random Vectors
00000●00

## Examples 1 and 2

Let $X \sim Binomial(17, 0.4)$. Then

$cov(X) =$

If $Y_1, \ldots, Y_n$ are independent then

$$cov\left(\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}\right) =$$

Introduction
○○○○

Matrix Algebra
○○○○○

Expectations of Random Vectors
○○○○○○○●○

## Example 3

Let $X, Y$ be independent r.v. with $X \sim N(5, 2)$ and $Y \sim$ Binomial$(10, 0.5)$. Then

$$\text{cov}\left(\begin{pmatrix} X \\ -X \end{pmatrix}\right) =$$

$$\text{cov}\left(\begin{pmatrix} X \\ X + Y \end{pmatrix}\right) =$$

$$\text{cov}\left(X, \begin{pmatrix} 2X \\ X - Y \end{pmatrix}\right) =$$

## Looking ahead

In the next lecture we discuss how to use these concepts to specify and work with general linear models (with multiple predictors)