

# 1. Introduction to High Dimensional Data Analysis

Lieven Clement

statOmics, Ghent University (<https://statomics.github.io>)

## Contents

|                     |   |           |
|---------------------|---|-----------|
| <b>1</b>            | <b>Introduction</b>   | <b>1</b>  |
| 1.1                 | What are high dimensional data? . . . . .   | 2         |
| <b>2</b>            | <b>Important tasks in high dimensional data analysis?</b>   | <b>2</b>  |
| 2.1                 | Example: Kang et al. (2018)'s droplet-based scRNA-seq data of PBMCs cells from 8 lupus patients measured before and after 6h-treatment with INF- $\beta$ (16 samples in total). . . . . | 2         |
| 2.2                 | Data exploration and dimensionality reduction . . . . .   | 3         |
| 2.3                 | Prediction . . . . .  | 5         |
| 2.4                 | Large scale hypothesis testing . . . . .  | 6         |
| <b>3</b>            | <b>Linear regression</b>  | <b>8</b>  |
| 3.1                 | Scalar form . . . . .   | 8         |
| 3.2                 | Vector/Matrix form . . . . .  | 8         |
| 3.3                 | Interpretation . . . . .  | 9         |
| 3.4                 | Least Squares (LS) . . . . .  | 9         |
| 3.5                 | Variance Estimator? . . . . .   | 16        |
| 3.6                 | Prediction error . . . . .  | 16        |
| <b>Session info</b> |   | <b>17</b> |

## 1 Introduction

- We live in a big data era
- Massive Datasets on our location, surfing behavior, consumer data, social media, ...
- Life Sciences: advent of high throughput technologies has enabled us to measure brain activity (brain images) as well as the expression of thousands of genes, proteins, ... for each subject or even single cell.
- Industry: process control with many sensors that follow process in real time...
- Data driven journalism
- ...

Challenge: We have to learn from high dimensional data!

## 1.1 What are high dimensional data?

- We typically observe multiple variables/features ( $p$ ) for each subject/experimental unit  $i = 1, \dots, n$  i.e.

$$\mathbf{x}_i^T = [x_{i1}, \dots, x_{ip}]$$

- Multivariate statistics have a long history, but were designed for the  $n >> p$  case,
- Nowadays many high throughput technologies generate multivariate data with many variables (large  $p$ ) as compared to the number of independent replicates or samples (sample size  $n$ ), resulting in **{high-dimensional data}**, which is characterised by

$$p >> n.$$

- New statistical methods for dealing with the  $p >> n$  case have been developed in the last 20 years. Some of them are adaptations of multivariate methods.

---

Issues with high-dimensional data:

- *Computational problems*: large matrices, numerical accuracy of computers become an issue
- *Classical asymptotic theory does not hold* for  $p \rightarrow \infty$  as  $n \rightarrow \infty$
- *Model (or feature) selection* requires specialised methods to deal with the enormous number of possible models. (e.g. in linear regression with  $p$  potential predictors:  $2^p$  possible models)
- Models that can potentially depend on large- $p$  predictors are vulnerable to *huge overfitting*.
- In searching for associations between an outcome and large  $p$  potential exploratory variables, we are at risk to make *many false discoveries*
- The *Curse of Dimensionality* may cause a prediction model to become useless. ( $n$  data points become sparse in large- $p$  dimensional space)

---

## 2 Important tasks in high dimensional data analysis?

### 2.1 Example: Kang et al. (2018)'s droplet-based scRNA-seq data of PBMCs cells from 8 lupus patients measured before and after 6h-treatment with INF- $\beta$ (16 samples in total).

```
library(tidyverse)
## Packages to load and visualize the single-cell data
library(ExperimentHub)
library(scater)
```

```

## Load Kang single-cell data from ExperimentHub
eh <- ExperimentHub()
sce <- eh[["EH2259"]]

sce

class: SingleCellExperiment
dim: 35635 29065
metadata(0):
assays(1): counts
rownames(35635): MIR1302-10 FAM138A ... MT-ND6 MT-CYB
rowData names(2): ENSEMBL SYMBOL
colnames(29065): AACATACAATGCC-1 AACATACATTCC-1 ... TTTGCATGGTTGG-1
    TTTGCATGTCTTAC-1
colData names(5): ind stim cluster cell multiplets
reducedDimNames(1): TSNE
mainExpName: NULL
altExpNames(0):

```

- Data on gene expression of 35635 genes of 29065 cells.

## 2.2 Data exploration and dimensionality reduction

- Visualisation is a first essential step to learn from data

```
counts(sce)[18990:19000, 9:20]
```

```
11 x 12 sparse Matrix of class "dgCMatrix"
```

|                | RP11-467L20.10 | MYRF | TMEM258 | FEN1 | FADS2 | FADS1 | FADS3 | RAB3IL1 | BEST1 | FTH1 | AP003733.1 |     |
|----------------|----------------|------|---------|------|-------|-------|-------|---------|-------|------|------------|-----|
| RP11-467L20.10 | .              | .    | .       | .    | .     | .     | .     | .       | .     | .    | .          |     |
| MYRF           | .              | .    | .       | .    | .     | .     | .     | .       | .     | .    | .          |     |
| TMEM258        | .              | .    | 3       | .    | .     | 1     | 1     | .       | .     | .    | .          |     |
| FEN1           | .              | .    | .       | .    | .     | .     | .     | .       | .     | .    | .          |     |
| FADS2          | .              | .    | .       | .    | .     | .     | .     | .       | .     | .    | .          |     |
| FADS1          | .              | .    | .       | .    | .     | .     | .     | .       | .     | .    | .          |     |
| FADS3          | .              | .    | .       | .    | .     | .     | .     | .       | .     | .    | .          |     |
| RAB3IL1        | .              | .    | .       | .    | .     | .     | .     | .       | .     | .    | .          |     |
| BEST1          | 1              | .    | .       | .    | .     | 1     | .     | .       | .     | 1    | 1          |     |
| FTH1           | 163            | 8    | 271     | 4    | 5     | 174   | 38    | 58      | 31    | 3    | 663        | 612 |
| AP003733.1     | .              | .    | .       | .    | .     | .     | .     | .       | .     | .    | .          |     |

- It is impossible to learn about the structure in the data by staring at the data matrix.
- We should be able to explore the data in a low dimensional projection

```
plotReducedDim(sce, dimred="TSNE", colour_by="cell")
```



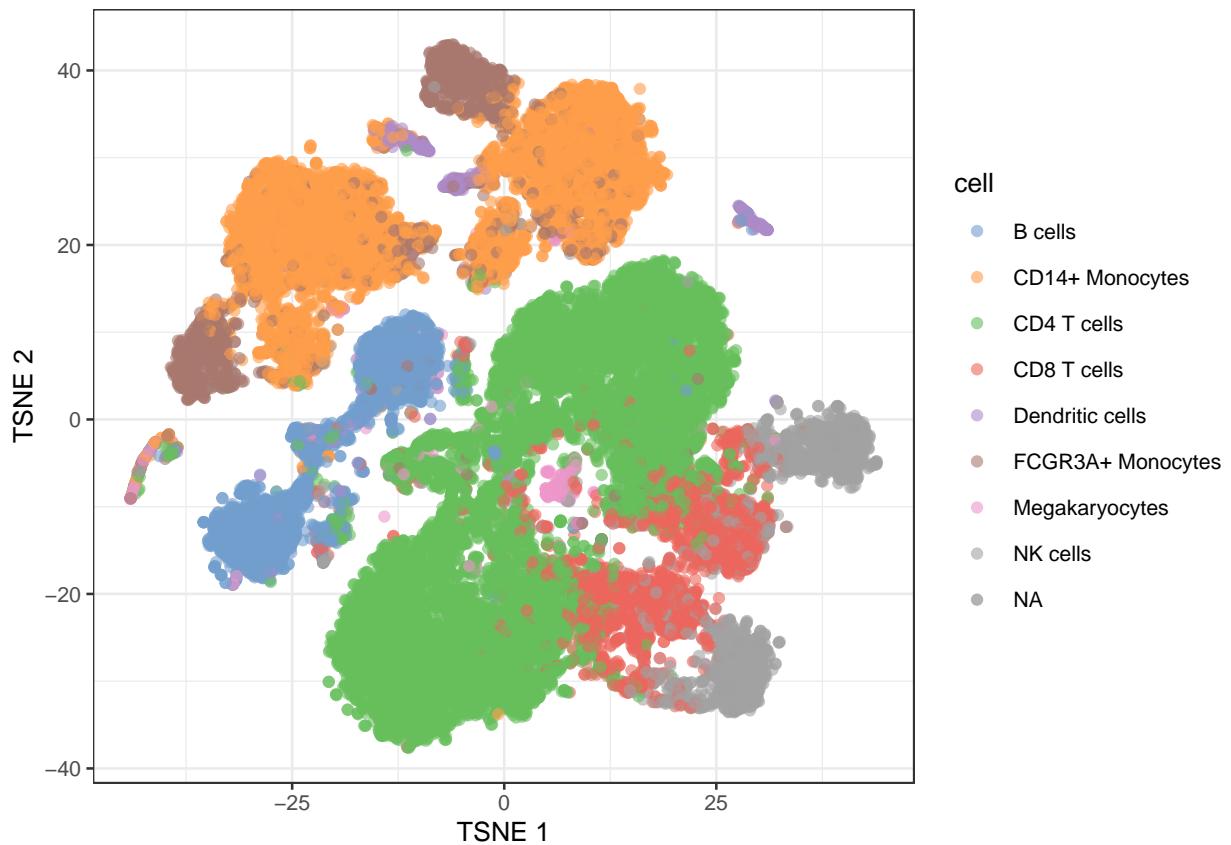
```
plotReducedDim(sce, dimred="TSNE", colour_by="stim")
```



- Note, that we see huge effect of treatment. If I see this I am always on my guard!
- We contacted the authors and learned that all control cells were sequenced in a first run and all stimulated cells were on a second sequencing run. So the large effect might be an effect of batch!

## 2.3 Prediction

```
plotReducedDim(sce, dimred="TSNE", colour_by="cell")
```



- In single cell analysis it is key to identify cell types in scRNA-seq data sets before in-depth investigations of their functional and pathological roles.
- Use models that were build based on reference data sets to predict cell types in new data sets based on their gene expression pattern
- Problem: we have 35635 genes at our disposal to build this prediction model!
- Other examples
  - Prediction of risk on mortality is used on a daily basis in intensive care units to prioritise patient care.
  - Facebook: predict the identity of the faces of people that are on a new image that is uploaded.
  - Netflix: Suggest movies that you would like

## 2.4 Large scale hypothesis testing

```
plotReducedDim(sce, dimred="TSNE", colour_by="cell", shape_by="stim")
```



- Which genes are differentially expressed between control and stimulated treatment (assess in each cell type)
- For which genes we see that the differential expression according to treatment is changing according to the cell type (interaction)
- We have to model the gene expression for each gene

$$\begin{cases} y_{ig} & \sim NB(\mu_{ig}, \phi_g) \\ E[y_{ig}] & = \mu \\ \log(\mu_{ig}) & = \eta_{ig} \\ \eta_{ig} & = \beta_{0,g} + \sum_{c=1}^C \beta_{c,g} X_{ic} + \beta_{s,g} X_{is} + \sum_{c=1}^C \beta_{c:s,g} X_{ic} X_{is} + \alpha_{ig} \end{cases}$$

With

- Cell type  $c = 2 \dots C$  and  $X_{ic}$  is a dummy variable that  $X_{ic} = 1$  if cell  $i$  is of cell type  $c$  and  $X_{ic} = 0$  otherwise. Note that cell type  $c = 1$  is the reference cell type
- Indicator  $X_{is}$  indicates if cell  $i$  was stimulated  $X_{is} = 1$  or not  $X_{is} = 0$ . So the control treatment is the reference treatment.

Suppose we want to test if the effect of the stimulus (average difference in expression between stimulated and non stimulated cells) is different in cell type  $c$  than in the reference cell type 1?

- $H_0 : \beta_{c:s,g} = 0$

- $H_1 : \beta_{c:s,g} \neq 0$
- We have to assess this for 35635 genes!
- If we assess each test at the 5% level we can expect  $0.05 * 35635 = 1782$  false positives.

→ massive multiple testing problem!

Note, that we cannot differentiate between batch and treatment because of the flaw in the experimental design!

Other examples

- Find regions (voxels) in the brain (on brain image) that are associated with a certain condition/treatment
  - Evaluation of trading rules
- 

### 3 Linear regression

- Linear regression is a very important statistical tool to study the association between variables and to build prediction models.

#### 3.1 Scalar form

- Consider a vector of predictors  $\mathbf{x} = (x_1, \dots, x_p)$  and
- a real-valued response  $Y$
- then the linear regression model can be written as

$$Y = f(\mathbf{x}) + \epsilon = \beta_0 + \sum_{j=1}^p x_j \beta_j + \epsilon$$

with i.i.d.  $\epsilon \sim N(0, \sigma^2)$

- We will typically work on mean centered variables:  $Y - \bar{Y}$  and  $X - \bar{X}$ , then the intercept is dropped from the model:

$$Y = f(\mathbf{x}) + \epsilon = \sum_{j=1}^p x_j \beta_j + \epsilon$$

#### 3.2 Vector/Matrix form

- $n$  observations  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  with  $\mathbf{x}_1^T = [1 \ x_1 \ \dots \ x_p]$
- Regression in matrix notation

$$\mathbf{Y} = \mathbf{X} +$$

$$\text{with } \mathbf{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix} \text{ or } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix} \text{ and } \epsilon = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Note, that upon centering of  $\mathbf{X}$  and  $\mathbf{Y}$  the 1 is dropped from each  $\mathbf{x}_i$  and thus the column of 1 is dropped in  $\mathbf{X}$

---

### 3.3 Interpretation

From the linear regression we get

$$E[Y | \mathbf{x}] = \mathbf{x}^T \boldsymbol{\beta}.$$

- Hence, the  $\beta$  parameters relate the regressor  $\mathbf{x}$  to the mean outcome.
- If we know the covariate pattern  $\mathbf{x}$  we can use the model to predict  $Y$ .

Upon centering we get the following for a model with a single regressor

$$E[Y | x] = \beta_1 x$$

and

$$\beta_1 = E[Y | x + 1] - E[Y | x].$$

Hence, the  $\beta_1$  parameter has an interpretation as the mean increase of the outcome as the regressor increases with one unit. The parameter does not say very about individual outcomes. It is the residual variance  $\sigma^2$  that determines how much individual outcomes vary about the mean outcome.

The  $\beta$  parameters are used to measure association, but a  $\beta \neq 0$  does not necessarily mean that the model will give good predictions.

- In a later chapter we will discuss the problem of large scale hypothesis testing: testing many hypotheses in a single study (ten to hundred thousands of hypotheses).
- A statistical test is constructed to control the type I error rate at the significance level. However, when many hypotheses are to be tested in a single study, the probability to find false associations is no longer controlled if p-values are compared to the significance level.
- We will later introduce the concept of false discovery rates to overcome the problem.

### 3.4 Least Squares (LS)

- Minimize the residual sum of squares

$$\begin{aligned} RSS(\boldsymbol{\beta}) &= \sum_{i=1}^n e_i^2 \\ &= \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \end{aligned}$$

- or in matrix notation

$$\begin{aligned} RSS(\boldsymbol{\beta}) &= (\mathbf{Y} - \mathbf{X})^T (\mathbf{Y} - \mathbf{X}) \\ &= \|\mathbf{Y} - \mathbf{X}\|_2^2 \end{aligned}$$

with the  $L_2$ -norm of a  $p$ -dim. vector  $v$   $\|\mathbf{v}\|_2 = \sqrt{v_1^2 + \dots + v_p^2} \rightarrow \hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\boldsymbol{\beta}} \|\mathbf{Y} - \mathbf{X}\|^2$

### 3.4.1 Minimize RSS

$$\begin{aligned}
 \frac{\partial RSS}{\partial \beta} &= \mathbf{0} \\
 \frac{(\mathbf{Y} - \mathbf{X})^T (\mathbf{Y} - \mathbf{X}\beta)}{\partial \beta} &= \mathbf{0} \\
 -2\mathbf{X}^T (\mathbf{Y} - \mathbf{X}\beta) &= \mathbf{0} \\
 \mathbf{X}^T \mathbf{X} &= \mathbf{X}^T \mathbf{Y} \\
 \hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}
 \end{aligned}$$

It can be shown that the estimator is unbiased:

$$E[\hat{\beta}] = \beta$$


---

### 3.4.2 Geometric Interpretation

Toy Example: fit without intercept

- n=3 and p=2

$$\mathbf{X} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \\ 0 & 0 \end{bmatrix}$$

```
set.seed(4)
x1 <- c(2, 0, 0)
x2 <- c(0, 2, 0)
y <- x1 * 0.5 + x2 * 0.5 + rnorm(3, 2)
fit <- lm(y ~ -1 + x1 + x2)
```

- $\mathbf{Y}$  is now given by

$$\mathbf{Y} = \begin{bmatrix} 1 + \epsilon_1 \\ 1 + \epsilon_2 \\ 0 + \epsilon_3 \end{bmatrix}$$

with the  $\epsilon_i$ s some random noise.

- The actual values of  $\mathbf{Y}$  are:

```
## True Y
y
```

[1] 3.216755 2.457507 2.891145

```
## Fitted Y
fit$fitted.values
```

```
1         2         3
3.216755 2.457507 0.000000
```

```
# predict values on regular xy grid
x1pred <- seq(-1, 3, length.out = 10)
x2pred <- seq(-1, 3, length.out = 10)
xy <- expand.grid(x1 = x1pred,
x2 = x2pred)
ypred <- matrix (nrow = 30, ncol = 30,
data = predict(fit, newdata = data.frame(xy),
interval = "prediction"))

library(plot3D)

# fitted points for droplines to surface
th=20
ph=5
scatter3D(x1,
           x2,
           y,
           pch = 16,
           col="darkblue",
           cex = 1,
           theta = th,
           ticktype = "detailed",
           xlab = "x1",
           ylab = "x2",
           zlab = "y",
           colvar=FALSE,
           bty = "g",
           xlim=c(-1,3),
           ylim=c(-1,3),
           zlim=c(-2,4))

for (i in 1:3)
  lines3D(
    x = rep(x1[i],2),
    y = rep(x2[i],2),
    z = c(y[i],fit$fitted[i]),
    col="darkblue",
    add=TRUE,
    lty=2)

z.pred3D <- outer(
  x1pred,
  x2pred,
```

```

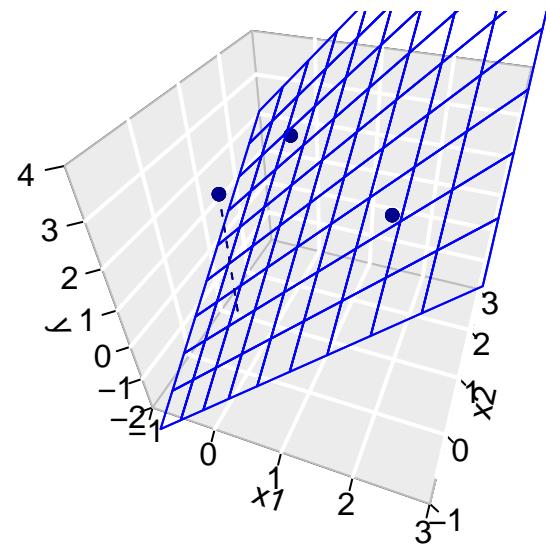
function(x1,x2)
{
  fit$coef[1]*x1+fit$coef[2]*x2
})

x.pred3D <- outer(
  x1pred,
  x2pred,
  function(x,y) x)

y.pred3D <- outer(
  x1pred,
  x2pred,
  function(x,y) y)

surf3D(
  x.pred3D,
  y.pred3D,
  z.pred3D,
  col="blue",
  facets=NA,
  add=TRUE)

```



### 3.4.2.1 Visualise fit

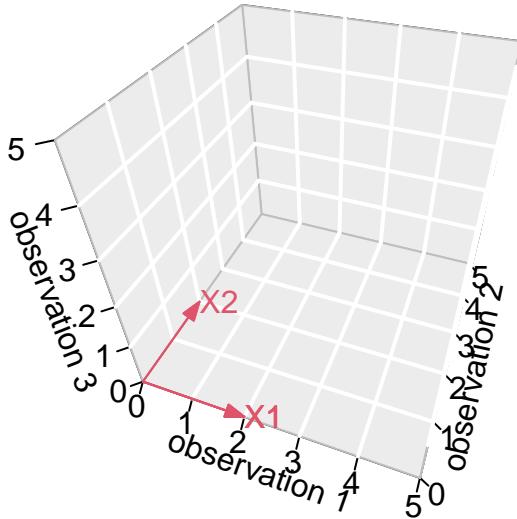
### 3.4.2.2 Projection

- We can also interpret the fit as the projection of the  $n \times 1$  vector  $\mathbf{Y}$  on the column space of the matrix  $\mathbf{X}$ .
- So each column in  $\mathbf{X}$  is also an  $n \times 1$  vector.
- For the toy example  $n=3$  and  $p=2$ . So the column space of  $\mathbf{X}$  is a plane in the three dimensional space.

$$\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

1. Plane spanned by column space:

```
arrows3D(
  0, 0, 0, x1[1], x1[2], x1[3],
  xlim = c(0, 5), ylim = c(0, 5), zlim = c(0, 5), bty = "g",
  theta = th, col = 2, ticktype = "detailed",
  xlab = "observation 1", ylab = "observation 2", zlab = "observation 3"
)
text3D(x1[1],x1[2],x1[3],labels="X1",col=2,add=TRUE)
arrows3D(0,0,0,x2[1],x2[2],x2[3],add=TRUE,col=2)
text3D(x2[1],x2[2],x2[3],labels="X2",col=2,add=TRUE)
```



2. Vector of  $\mathbf{Y}$ :

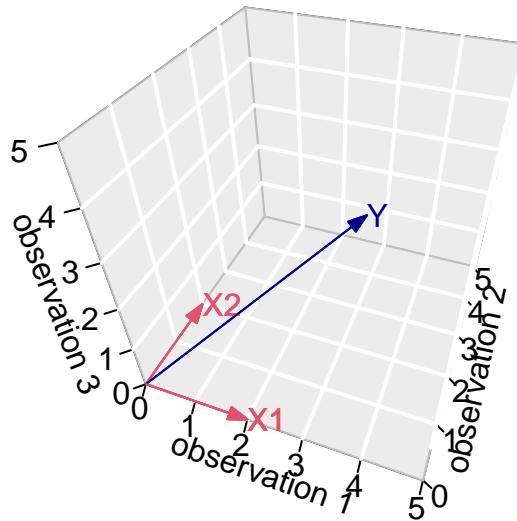
Actual values of  $\mathbf{Y}$ :

```
y
```

```
[1] 3.216755 2.457507 2.891145
```

$$\mathbf{Y} = \begin{bmatrix} 3.2167549 \\ 2.4575074 \\ 2.8911446 \end{bmatrix}$$

```
arrows3D(  
  0, 0, 0, x1[1], x1[2], x1[3],  
  xlim = c(0, 5), ylim = c(0, 5), zlim = c(0, 5), bty = "g",  
  theta = th, col = 2, ticktype = "detailed",  
  xlab = "observation 1", ylab = "observation 2", zlab = "observation 3"  
)  
text3D(x1[1],x1[2],x1[3],labels="X1",col=2,add=TRUE)  
arrows3D(0,0,0,x2[1],x2[2],x2[3],add=TRUE,col=2)  
text3D(x2[1],x2[2],x2[3],labels="X2",col=2,add=TRUE)  
arrows3D(0,0,0,y[1],y[2],y[3],add=TRUE,col="darkblue")  
text3D(y[1],y[2],y[3],labels="Y",col="darkblue",add=TRUE)
```



3. Projection of Y onto column space

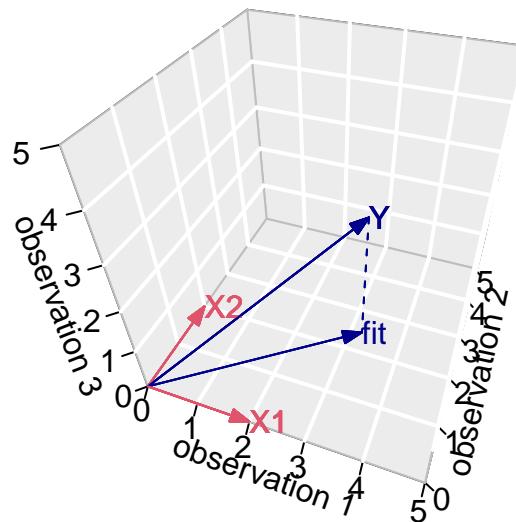
Actual values of fitted values  $\hat{\mathbf{Y}}$ :

```
fitted.values(fit)
```

```
1           2           3  
3.216755  2.457507  0.000000
```

$$\mathbf{Y} = \begin{bmatrix} 3.2167549 \\ 2.4575074 \\ 0 \end{bmatrix}$$

```
arrows3D(  
  0, 0, 0, x1[1], x1[2], x1[3],  
  xlim = c(0, 5), ylim = c(0, 5), zlim = c(0, 5), bty = "g",  
  theta = th, col = 2, ticktype = "detailed",  
  xlab = "observation 1", ylab = "observation 2", zlab = "observation 3"  
)  
text3D(x1[1],x1[2],x1[3],labels="X1",col=2,add=TRUE)  
arrows3D(0,0,0,x2[1],x2[2],x2[3],add=TRUE,col=2)  
text3D(x2[1],x2[2],x2[3],labels="X2",col=2,add=TRUE)  
arrows3D(0,0,0,y[1],y[2],y[3],add=TRUE,col="darkblue")  
text3D(y[1],y[2],y[3],labels="Y",col="darkblue",add=TRUE)  
arrows3D(0,0,0,fit$fitted[1],fit$fitted[2],fit$fitted[3],add=TRUE,col="darkblue")  
segments3D(y[1],y[2],y[3],fit$fitted[1],fit$fitted[2],fit$fitted[3],add=TRUE,lty=2,col="darkblue")  
text3D(fit$fitted[1],fit$fitted[2],fit$fitted[3],labels="fit",col="darkblue",add=TRUE)
```



- Note, that it is also clear from the equation in the derivation of the LS that the residual is orthogonal on the column space:

$$-2\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta) = 0$$

- Imagine what happens when  $p \gg n$ !!!

→ curse of dimensionality!

### 3.5 Variance Estimator?

$$\begin{aligned}\hat{\Sigma}_{\hat{\beta}} &= \text{var}[(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}] \\ &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\text{var}[\mathbf{Y}]\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} \\ &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T(\mathbf{I}\sigma^2)\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1} \\ &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{I}\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\sigma^2 \\ &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\sigma^2 \\ &= (\mathbf{X}^T\mathbf{X})^{-1}\sigma^2\end{aligned}$$

The fact that  $\hat{\beta}$  is unbiased and has a variance of  $(\mathbf{X}^T\mathbf{X})^{-1}\sigma^2$  will be important when assessing association!

---

### 3.6 Prediction error

Least squares estimators are unbiased and consistent, but these properties are not very important for prediction models.

A prediction model is considered good if it can predict well outcomes.

The **prediction error** for a prediction at predictor  $\mathbf{x}$  is given by

$$\hat{Y}(\mathbf{x}) - Y^*,$$

where

- $\hat{Y}(\mathbf{x}) = \mathbf{x}^T\hat{\beta}$  is the prediction at  $\mathbf{x}$
- $Y^*$  is an outcome at predictor  $\mathbf{x}$

Since prediction is typically used to predict an outcome before it is observed, the outcome  $Y^*$  is not observed yet. Hence, the prediction error cannot be computed.

---

The problem of unobservable prediction errors is partly solved by the **expected conditional test error** (sometimes referred to as the mean squared error, MSE)

$$\text{Err}(\mathbf{x}) = E[(\hat{Y}(\mathbf{x}) - Y^*)^2].$$

With (suppressing the dependence on  $\mathbf{x}$ )

$$\mu = E[\hat{Y}] \text{ and } \mu^* = E[Y^*]$$

the error can be expressed as

$$\begin{aligned}\text{Err} &= E\left\{\left[(\hat{Y} - \mu) - (Y - \mu^*) - (\mu^* - \mu)\right]^2\right\} \\ &= E[(\hat{Y} - \mu)^2] + E[(Y - \mu^*)^2] + E[(\mu^* - \mu)^2] \\ &= \text{var}[\hat{Y}] + \text{var}[Y] + \text{bias}^2\end{aligned}$$

The term  $\text{var}[Y]$  (irreducible error) does not depend on the model and may therefore be ignored when Err is used for comparing prediction models.

---

In this introductory chapter we only aim to give a rough discussion on prediction errors. Later definitions will be refined and the notation will be more accurate. Also a more detailed discussion on the bias-variance trade-off will follow. For the moment it is sufficient to vaguely know that:

- the expected conditional test error is introduced to circumvent the problem that the prediction error cannot be observed. In later chapters we will look at estimators of the expected error.
- the expected conditional test error is in some literature also known as the *mean squared error* (MSE), but we do not adopt this terminology because MSE is also commonly used to refer to SSE divided by the residual degrees of freedom in a linear regression model.
- The identity  $\text{Err} = \text{var}[\hat{Y}] + \text{var}[Y] + \text{bias}^2$  is known as the bias-variance trade-off. It shows that a good prediction model (i.e. a model resulting in a small Err), can be obtained by a model that shows a small bias as long as this bias is compensated with a large reduction of the variance or the predictions. A more detailed discussion will follow in later chapters.
- For prediction models with a large number of predictors we will therefore introduce penalized regression. This will induce some bias in the estimation, but will allow us to reduce the variance considerably.

## Session info

Session info

```
[1] "2022-01-24 11:17:15 UTC"
```

```
- Session info -----
  setting  value
  version  R version 4.1.2 (2021-11-01)
  os        macOS Big Sur 10.16
  system   x86_64, darwin17.0
  ui        X11
```

```

language (EN)
collate  en_US.UTF-8
ctype    en_US.UTF-8
tz       UTC
date    2022-01-24
pandoc  2.7.3 @ /usr/local/bin/ (via rmarkdown)

```

| - Packages ----- |                    |           |            |                    |
|------------------|--------------------|-----------|------------|--------------------|
|                  | ! package          | * version | date (UTC) | lib source         |
| P                | AnnotationDbi      | 1.54.1    | 2021-06-08 | [?] Bioconductor   |
| P                | AnnotationHub      | * 3.0.1   | 2021-06-20 | [?] Bioconductor   |
| P                | assertthat         | 0.2.1     | 2019-03-21 | [?] CRAN (R 4.1.0) |
| P                | backports          | 1.2.1     | 2020-12-09 | [?] CRAN (R 4.1.0) |
| P                | beachmat           | 2.8.1     | 2021-08-12 | [?] Bioconductor   |
| P                | beeswarm           | 0.4.0     | 2021-06-01 | [?] CRAN (R 4.1.0) |
| P                | Biobase            | * 2.52.0  | 2021-05-19 | [?] Bioconductor   |
| P                | BiocFileCache      | * 2.0.0   | 2021-05-19 | [?] Bioconductor   |
| P                | BiocGenerics       | * 0.38.0  | 2021-05-19 | [?] Bioconductor   |
| P                | BiocManager        | 1.30.16   | 2021-06-15 | [?] CRAN (R 4.1.0) |
| P                | BiocNeighbors      | 1.10.0    | 2021-05-19 | [?] Bioconductor   |
| P                | BiocParallel       | 1.26.2    | 2021-08-22 | [?] Bioconductor   |
| P                | BiocSingular       | 1.8.1     | 2021-06-08 | [?] Bioconductor   |
| P                | BiocVersion        | 3.13.1    | 2021-03-09 | [?] Bioconductor   |
| P                | Biostrings         | 2.60.2    | 2021-08-05 | [?] Bioconductor   |
| P                | bit                | 4.0.4     | 2020-08-04 | [?] CRAN (R 4.1.0) |
| P                | bit64              | 4.0.5     | 2020-08-30 | [?] CRAN (R 4.1.0) |
| P                | bitops             | 1.0-7     | 2021-04-24 | [?] CRAN (R 4.1.0) |
| P                | blob               | 1.2.2     | 2021-07-23 | [?] CRAN (R 4.1.0) |
| P                | broom              | 0.7.9     | 2021-07-27 | [?] CRAN (R 4.1.0) |
| P                | bslib              | 0.3.0     | 2021-09-02 | [?] CRAN (R 4.1.0) |
| P                | cachem             | 1.0.6     | 2021-08-19 | [?] CRAN (R 4.1.0) |
| P                | cellranger         | 1.1.0     | 2016-07-27 | [?] CRAN (R 4.1.0) |
| P                | cli                | 3.1.1     | 2022-01-20 | [?] CRAN (R 4.1.2) |
| P                | colorspace         | 2.0-2     | 2021-06-24 | [?] CRAN (R 4.1.0) |
| P                | crayon             | 1.4.1     | 2021-02-08 | [?] CRAN (R 4.1.0) |
| P                | curl               | 4.3.2     | 2021-06-23 | [?] CRAN (R 4.1.0) |
| P                | DBI                | 1.1.1     | 2021-01-15 | [?] CRAN (R 4.1.0) |
| P                | dbplyr             | * 2.1.1   | 2021-04-06 | [?] CRAN (R 4.1.0) |
| P                | DelayedArray       | 0.18.0    | 2021-05-19 | [?] Bioconductor   |
| P                | DelayedMatrixStats | 1.14.3    | 2021-08-26 | [?] Bioconductor   |
| P                | digest             | 0.6.28    | 2021-09-23 | [?] CRAN (R 4.1.0) |
| P                | dplyr              | * 1.0.7   | 2021-06-18 | [?] CRAN (R 4.1.0) |
| P                | ellipsis           | 0.3.2     | 2021-04-29 | [?] CRAN (R 4.1.0) |
| P                | evaluate           | 0.14      | 2019-05-28 | [?] CRAN (R 4.1.0) |
| P                | ExperimentHub      | * 2.0.0   | 2021-05-19 | [?] Bioconductor   |
| P                | fansi              | 0.5.0     | 2021-05-25 | [?] CRAN (R 4.1.0) |
| P                | farver             | 2.1.0     | 2021-02-28 | [?] CRAN (R 4.1.0) |
| P                | fastmap            | 1.1.0     | 2021-01-25 | [?] CRAN (R 4.1.0) |
| P                | filelock           | 1.0.2     | 2018-10-05 | [?] CRAN (R 4.1.0) |
| P                | forcats            | * 0.5.1   | 2021-01-27 | [?] CRAN (R 4.1.0) |
| P                | fs                 | 1.5.0     | 2020-07-31 | [?] CRAN (R 4.1.0) |
| P                | generics           | 0.1.0     | 2020-10-31 | [?] CRAN (R 4.1.0) |
| P                | GenomeInfoDb       | * 1.28.4  | 2021-09-09 | [?] Bioconductor   |
| P                | GenomeInfoDbData   | 1.2.6     | 2022-01-24 | [?] Bioconductor   |

|                          |          |            |     |                                  |
|--------------------------|----------|------------|-----|----------------------------------|
| P GenomicRanges          | * 1.44.0 | 2021-05-19 | [?] | Bioconductor                     |
| P ggbeeswarm             | 0.6.0    | 2017-08-07 | [?] | CRAN (R 4.1.0)                   |
| P ggplot2                | * 3.3.5  | 2021-06-25 | [?] | CRAN (R 4.1.0)                   |
| P glue                   | 1.4.2    | 2020-08-27 | [?] | CRAN (R 4.1.0)                   |
| P gridExtra              | 2.3      | 2017-09-09 | [?] | CRAN (R 4.1.0)                   |
| P gtable                 | 0.3.0    | 2019-03-25 | [?] | CRAN (R 4.1.0)                   |
| P haven                  | 2.4.3    | 2021-08-04 | [?] | CRAN (R 4.1.0)                   |
| P highr                  | 0.9      | 2021-04-16 | [?] | CRAN (R 4.1.0)                   |
| P hms                    | 1.1.1    | 2021-09-26 | [?] | CRAN (R 4.1.0)                   |
| P htmltools              | 0.5.2    | 2021-08-25 | [?] | CRAN (R 4.1.0)                   |
| P httpuv                 | 1.6.3    | 2021-09-09 | [?] | CRAN (R 4.1.0)                   |
| P httr                   | 1.4.2    | 2020-07-20 | [?] | CRAN (R 4.1.0)                   |
| P interactiveDisplayBase | 1.30.0   | 2021-05-19 | [?] | Bioconductor                     |
| P IRanges                | * 2.26.0 | 2021-05-19 | [?] | Bioconductor                     |
| P irlba                  | 2.3.3    | 2019-02-05 | [?] | CRAN (R 4.1.0)                   |
| P jquerylib              | 0.1.4    | 2021-04-26 | [?] | CRAN (R 4.1.0)                   |
| P jsonlite               | 1.7.2    | 2020-12-09 | [?] | CRAN (R 4.1.0)                   |
| P KEGGREST               | 1.32.0   | 2021-05-19 | [?] | Bioconductor                     |
| P knitr                  | 1.33     | 2021-04-24 | [?] | CRAN (R 4.1.0)                   |
| P labeling               | 0.4.2    | 2020-10-20 | [?] | CRAN (R 4.1.0)                   |
| P later                  | 1.3.0    | 2021-08-18 | [?] | CRAN (R 4.1.0)                   |
| P lattice                | 0.20-45  | 2021-09-22 | [?] | CRAN (R 4.1.2)                   |
| P lifecycle              | 1.0.1    | 2021-09-24 | [?] | CRAN (R 4.1.0)                   |
| P lubridate              | 1.7.10   | 2021-02-26 | [?] | CRAN (R 4.1.0)                   |
| P magrittr               | 2.0.1    | 2020-11-17 | [?] | CRAN (R 4.1.0)                   |
| P Matrix                 | 1.3-4    | 2021-06-01 | [?] | CRAN (R 4.1.2)                   |
| P MatrixGenerics         | * 1.4.3  | 2021-08-26 | [?] | Bioconductor                     |
| P matrixStats            | * 0.60.1 | 2021-08-23 | [?] | CRAN (R 4.1.0)                   |
| P memoise                | 2.0.0    | 2021-01-26 | [?] | CRAN (R 4.1.0)                   |
| P mime                   | 0.12     | 2021-09-28 | [?] | CRAN (R 4.1.0)                   |
| P misc3d                 | 0.9-0    | 2020-09-06 | [?] | CRAN (R 4.1.0)                   |
| P modelr                 | 0.1.8    | 2020-05-19 | [?] | CRAN (R 4.1.0)                   |
| P munsell                | 0.5.0    | 2018-06-12 | [?] | CRAN (R 4.1.0)                   |
| P muscData               | * 1.6.0  | 2021-05-20 | [1] | Bioconductor                     |
| P pillar                 | 1.6.3    | 2021-09-26 | [?] | CRAN (R 4.1.0)                   |
| P pkgconfig              | 2.0.3    | 2022-01-24 | [?] | Github (r-lib/pkgconfig@b81ae03) |
| P plot3D                 | * 1.4    | 2021-05-22 | [?] | CRAN (R 4.1.0)                   |
| P png                    | 0.1-7    | 2013-12-03 | [?] | CRAN (R 4.1.0)                   |
| P promises               | 1.2.0.1  | 2021-02-11 | [?] | CRAN (R 4.1.0)                   |
| P purrr                  | * 0.3.4  | 2020-04-17 | [?] | CRAN (R 4.1.0)                   |
| P R6                     | 2.5.1    | 2021-08-19 | [?] | CRAN (R 4.1.0)                   |
| P rappdirs               | 0.3.3    | 2021-01-31 | [?] | CRAN (R 4.1.0)                   |
| P Rcpp                   | 1.0.7    | 2021-07-07 | [?] | CRAN (R 4.1.0)                   |
| P RCurl                  | 1.98-1.4 | 2021-08-17 | [?] | CRAN (R 4.1.0)                   |
| P readr                  | * 2.0.1  | 2021-08-10 | [?] | CRAN (R 4.1.0)                   |
| P readxl                 | 1.3.1    | 2019-03-13 | [?] | CRAN (R 4.1.0)                   |
| P renv                   | 0.14.0   | 2021-07-21 | [?] | CRAN (R 4.1.0)                   |
| P reprex                 | 2.0.1    | 2021-08-05 | [?] | CRAN (R 4.1.0)                   |
| P rlang                  | 0.4.11   | 2021-04-30 | [?] | CRAN (R 4.1.0)                   |
| P rmarkdown              | 2.10     | 2021-08-06 | [?] | CRAN (R 4.1.0)                   |
| P RSQLite                | 2.2.8    | 2021-08-21 | [?] | CRAN (R 4.1.0)                   |
| P rstudioapi             | 0.13     | 2020-11-12 | [?] | CRAN (R 4.1.0)                   |
| P rsvd                   | 1.0.5    | 2021-04-16 | [?] | CRAN (R 4.1.0)                   |
| P rvest                  | 1.0.1    | 2021-07-26 | [?] | CRAN (R 4.1.0)                   |

|                        |          |                |                |
|------------------------|----------|----------------|----------------|
| P S4Vectors            | * 0.30.0 | 2021-05-19 [?] | Bioconductor   |
| P sass                 | 0.4.0    | 2021-05-12 [?] | CRAN (R 4.1.0) |
| P ScaledMatrix         | 1.0.0    | 2021-05-19 [?] | Bioconductor   |
| P scales               | 1.1.1    | 2020-05-11 [?] | CRAN (R 4.1.0) |
| P scater               | * 1.20.1 | 2021-06-15 [?] | Bioconductor   |
| P scuttle              | * 1.2.1  | 2021-08-05 [?] | Bioconductor   |
| P sessioninfo          | 1.2.2    | 2021-12-06 [?] | CRAN (R 4.1.0) |
| P shiny                | 1.7.0    | 2021-09-22 [?] | CRAN (R 4.1.0) |
| P SingleCellExperiment | * 1.14.1 | 2021-05-21 [?] | Bioconductor   |
| P sparseMatrixStats    | 1.4.2    | 2021-08-08 [?] | Bioconductor   |
| P stringi              | 1.7.4    | 2021-08-25 [?] | CRAN (R 4.1.0) |
| P stringr              | * 1.4.0  | 2019-02-10 [?] | CRAN (R 4.1.0) |
| P SummarizedExperiment | * 1.22.0 | 2021-05-19 [?] | Bioconductor   |
| P tibble               | * 3.1.5  | 2021-09-30 [?] | CRAN (R 4.1.0) |
| P tidyverse            | * 1.1.4  | 2021-09-27 [?] | CRAN (R 4.1.0) |
| P tidyselect           | 1.1.1    | 2021-04-30 [?] | CRAN (R 4.1.0) |
| P tzdb                 | 0.1.2    | 2021-07-20 [?] | CRAN (R 4.1.0) |
| P utf8                 | 1.2.2    | 2021-07-24 [?] | CRAN (R 4.1.0) |
| P vctrs                | 0.3.8    | 2021-04-29 [?] | CRAN (R 4.1.0) |
| P vipor                | 0.4.5    | 2017-03-22 [?] | CRAN (R 4.1.0) |
| P viridis              | 0.6.1    | 2021-05-11 [?] | CRAN (R 4.1.0) |
| P viridisLite          | 0.4.0    | 2021-04-13 [?] | CRAN (R 4.1.0) |
| P withr                | 2.4.2    | 2021-04-18 [?] | CRAN (R 4.1.0) |
| P xfun                 | 0.25     | 2021-08-06 [?] | CRAN (R 4.1.0) |
| P xml2                 | 1.3.2    | 2020-04-23 [?] | CRAN (R 4.1.0) |
| P xtable               | 1.8-4    | 2019-04-21 [?] | CRAN (R 4.1.0) |
| P XVector              | 0.32.0   | 2021-05-19 [?] | Bioconductor   |
| P yaml                 | 2.2.1    | 2020-02-01 [?] | CRAN (R 4.1.0) |
| P zlibbioc             | 1.38.0   | 2021-05-19 [?] | Bioconductor   |

```
[1] /Users/runner/work/HDDA21/HDDA21/renv/library/R-4.1/x86_64-apple-darwin17.0
[2] /private/var/folders/24/8k48j16d249_n_qfxwsl6xvm0000gn/T/RtmpgQEB16/renv-system-library
```

P -- Loaded and on-disk path mismatch.

---