

Basic methods for data analysis

Philipp Renz

January 3, 2022

Contents

1	Introduction	i
2	Summarizing univariate data	i
3	Linear algebra: Basics	i
3.1	Matrix multiplication	ii
4	Summarizing multivariate data	ii
4.1	Bivariate data	ii
4.2	More features	ii
4.3	Principal component analysis	ii
5	Linear algebra: Singular value decomposition	ii
5.1	Another look at PCA	ii
6	Multiple linear regression	ii
7	Probability theory	ii

1 Introduction

2 Summarizing univariate data

3 Linear algebra: Basics

In data analysis we will often represent our data as matrices. Consequently many algorithms are described using the concept of matrices and their properties which could be broadly called the study of linear algebra. Here I don't want to give an axiomatic introduction but rather provide some basic results with the prerequisite that the reader already knows some basics about vectors and matrices. We will denote vectors as bold lowercase letters

$$\mathbf{a} = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix}, \quad (1)$$

where a_i is the i^{th} entry of the vector \mathbf{a} .

Matrices are denoted by bold uppercase letters

$$\mathbf{A} = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix}. \quad (2)$$

Depending on the context it is useful to view a matrix as a concatenation of either row or column vectors.

$$\mathbf{A} = \begin{pmatrix} - & \mathbf{a}_1^r & - \\ & \vdots & \\ - & \mathbf{a}_m^r & - \end{pmatrix} = \begin{pmatrix} | & & | \\ \mathbf{a}_1^c & \dots & \mathbf{a}_n^c \\ | & & | \end{pmatrix}. \quad (3)$$

We made an explicit distinction between row- and column vectors using superscripts here but these will often be left out for better readability.

3.1 Matrix multiplication

Given two matrices

$$\mathbf{A} = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix} \in R^{m \times n}, \quad \mathbf{B} = \begin{pmatrix} b_{11} & \dots & b_{1p} \\ \vdots & \ddots & \vdots \\ b_{n1} & \dots & b_{np} \end{pmatrix} \in R^{n \times p} \quad (4)$$

their product $\mathbf{C} = \mathbf{AB}$ is of shape $(m \times p)$ and its entries are

$$c_{ik} = \sum_{j=1}^n a_{ij} b_{jk} \quad (5)$$

4 Summarizing multivariate data

4.1 Bivariate data

4.2 More features

4.3 Principal component analysis

5 Linear algebra: Singular value decomposition

5.1 Another look at PCA

6 Multiple linear regression

7 Probability theory