Statistical Theory and Modeling (ST2601) Lecture 8 - Linear regression in vector form

Mattias Villani

Department of Statistics Stockholm University













Vectors and matrices - minimal intro to linear algebra

Linear regression in vector form

Multivariate normal distribution

Goals of the lecture

Linear regression in vector form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \qquad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

Least squares estimate of regression coefficients

$$\hat{\boldsymbol{eta}} = (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} \boldsymbol{y}$$

Multivariate normal distribution $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with pdf $f(\mathbf{x}) = |2\pi \boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$

What's the deal with all the bold letters? Linear algebra.

Worth the trip. Very useful for linear regression and more.

Vectors

Linear algebra: a vector is an object containing real numbers

$$\boldsymbol{a} = \begin{pmatrix} 1\\ 3\\ 5\\ 3 \end{pmatrix}$$

Common default: a vector is a column vector.
The transpose of a vector is a row vector

$$\boldsymbol{a}^{\top} = \left(\begin{array}{cccc} 1 & 3 & 5 & 3 \end{array} \right)$$

R:

> a = c(1,2,5,3)
> t(a) # transpose

Visualizing vectors in 2D

2D vector. Directed line (arrow) in \mathbb{R}^2 .



Visualizing vectors in 3D

3D vector. Directed line (arrow) in \mathbb{R}^3 .

$$\mathbf{a} = \begin{pmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \mathbf{a}_3 \end{pmatrix} = \begin{pmatrix} 3 \\ 2 \\ 2 \end{pmatrix}$$

$$\mathbf{a}=\left(3,2,2
ight)^{ op}$$



Vector addition and subtraction

Adding two vectors with the same number of elements

$$oldsymbol{a} = \left(egin{array}{c} oldsymbol{a}_1\ oldsymbol{a}_2\ oldsymbol{a}_3\end{array}
ight), \quad oldsymbol{b} = \left(egin{array}{c} b_1\ b_2\ b_2\ b_3\end{array}
ight) \quad oldsymbol{a} + oldsymbol{b} = \left(egin{array}{c} oldsymbol{a}_1 + b_1\ oldsymbol{a}_2 + b_2\ oldsymbol{a}_3 + b_3\end{array}
ight)$$

Substracting a vector from another vector

$$oldsymbol{a}-oldsymbol{b}=\left(egin{array}{cc} oldsymbol{a}_1-oldsymbol{b}_1\ oldsymbol{a}_2-oldsymbol{b}_2\ oldsymbol{a}_3-oldsymbol{b}_3\end{array}
ight)$$

Both these operations can be visualized geometrically.

Vector multiplication

In R a*b will do elementwise multiplication

$$oldsymbol{a} * oldsymbol{b} = \left(egin{array}{c} a_1 b_1 \ a_2 b_2 \ a_3 b_3 \end{array}
ight)$$

In a%*%b will compute the dot product

$$oldsymbol{a}\cdotoldsymbol{b}=oldsymbol{a}^{ op}oldsymbol{b}=oldsymbol{(} egin{array}{c} b_1\ b_2\ b_2\ b_3\end{array}ig)=oldsymbol{a}_1b_1+oldsymbol{a}_2b_2+oldsymbol{a}_3b_3$$

In general the dot product is

$$\boldsymbol{a}\cdot\boldsymbol{b}=\sum_{i=1}^na_ib_i$$

Orthogonal vectors

Two vectors are orthogonal if their dot product is zero



Matrices

A matrix is like a table, it has rows and columns

$$\mathbf{X} = \left(\begin{array}{rrr} 2 & 3 & 1 \\ 3 & 2 & 0 \end{array}\right)$$

This is a 2 × 3 matrix since it has 2 rows and 3 columns.
View a p × q matrix as q column vector stacked horizontally

$$\boldsymbol{X} = \left(\begin{array}{cccc} | & | & | \\ \boldsymbol{x}_1 & \boldsymbol{x}_2 & \cdots & \boldsymbol{x}_q \\ | & | & | \end{array}\right)$$

Example: the following three vectors give the matrix above

$$\mathbf{x}_1 = \begin{pmatrix} 2\\3 \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} 3\\2 \end{pmatrix}, \quad \mathbf{x}_3 = \begin{pmatrix} 1\\0 \end{pmatrix}$$

> x1 = c(2,3); x2 = c(3,2); x3=c(1,0);

> cbind(x1,x2,x3) # column bind. Also rbind exists

Matrix-Vector multiplication

- **A** is an $m \times n$ matrix **A**
- **b** is an *n*-element vector
- Matrix-vector product: dot product of each row in A with b

$$\mathbf{A}_{(m \times n)} = \begin{pmatrix} - & \mathbf{a}_{1}^{\top} & - \\ - & \mathbf{a}_{2}^{\top} & - \\ \vdots & \\ - & \mathbf{a}_{m}^{\top} & - \end{pmatrix} \qquad \mathbf{b}_{(n \times 1)} = \begin{pmatrix} b_{1} \\ \vdots \\ b_{n} \end{pmatrix}$$
$$\mathbf{Ab} = \begin{pmatrix} \mathbf{a}_{1}^{\top} \mathbf{b} \\ \mathbf{a}_{2}^{\top} \mathbf{b} \\ \vdots \\ \mathbf{a}_{m}^{\top} \mathbf{b} \end{pmatrix}$$

Matrix-Matrix multiplication

Matrix product of A and B: pairwise dot product of a row in A and a column in B

$$\mathbf{A}_{(m \times n)} = \begin{pmatrix} - & \mathbf{a}_{1}^{\top} & - \\ - & \mathbf{a}_{2}^{\top} & - \\ \vdots & \\ - & \mathbf{a}_{m}^{\top} & - \end{pmatrix} \qquad \mathbf{B}_{(n \times r)} = \begin{pmatrix} | & | & | & | \\ \mathbf{b}_{1} & \mathbf{b}_{2} & \cdots & \mathbf{b}_{r} \\ | & | & | & | \end{pmatrix}$$

$$\boldsymbol{A}\boldsymbol{B} = \begin{pmatrix} \boldsymbol{a}_1^{\top}\boldsymbol{b}_1 & \boldsymbol{a}_1^{\top}\boldsymbol{b}_2 & \cdots & \boldsymbol{a}_1^{\top}\boldsymbol{b}_r \\ \boldsymbol{a}_2^{\top}\boldsymbol{b}_1 & \boldsymbol{a}_2^{\top}\boldsymbol{b}_2 & \cdots & \boldsymbol{a}_2^{\top}\boldsymbol{b}_r \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{a}_m^{\top}\boldsymbol{b}_1 & \boldsymbol{a}_m^{\top}\boldsymbol{b}_2 & \cdots & \boldsymbol{a}_m^{\top}\boldsymbol{b}_r \end{pmatrix}$$

Matrix-Matrix multiplication

Example $\mathbf{A} = \begin{pmatrix} 2 & 3 \\ 3 & 2 \end{pmatrix}, \qquad \mathbf{B} = \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}$ $\mathbf{AB} = \begin{pmatrix} (2 & 3) \begin{pmatrix} 1 \\ 0 \end{pmatrix} & (2 & 3) \begin{pmatrix} 2 \\ 1 \end{pmatrix} \\ (3 & 2) \begin{pmatrix} 1 \\ 0 \end{pmatrix} & (3 & 2) \begin{pmatrix} 2 \\ 1 \end{pmatrix} \end{pmatrix}$ $= \begin{pmatrix} 2 \cdot 1 + 3 \cdot 0 = 2 & 2 \cdot 2 + 3 \cdot 1 = 7 \\ 3 \cdot 1 + 2 \cdot 0 = 3 & 3 \cdot 2 + 2 \cdot 1 = 8 \end{pmatrix}$

> A = matrix(c(2,3,3,2), 2, 2, byrow = TRUE) > B = matrix(c(1,2,0,1), 2, 2, byrow = TRUE) > A%*%B # A*B would do elementwise multiplication

Linear regression - one observation

One observation

$$y = \beta_1 x_1 + \ldots + \beta_p x_p + \varepsilon$$

In vector form

$$y = \underbrace{\left(\begin{array}{ccc} x_1 & \cdots & x_p\end{array}\right)}_{\mathbf{x}^\top} \underbrace{\left(\begin{array}{c} \beta_1 \\ \vdots \\ \beta_p\end{array}\right)}_{\boldsymbol{\beta}} + \varepsilon = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon$$

Add a one for the **intercept**

$$\begin{pmatrix} 1 & x_1 & \cdots & x_p \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

The ith observation

$$y_i = \boldsymbol{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$$

Linear regression - all observations

The ith observation

$$y_i = \boldsymbol{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$$

All i = 1, ..., n observations stacked under each other

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^\top \boldsymbol{\beta} \\ \mathbf{x}_2^\top \boldsymbol{\beta} \\ \vdots \\ \mathbf{x}_n^\top \boldsymbol{\beta} \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

With matrix-vector multiplication

$$\begin{pmatrix} \mathbf{x}_{1}^{\top}\boldsymbol{\beta} \\ \mathbf{x}_{2}^{\top}\boldsymbol{\beta} \\ \vdots \\ \mathbf{x}_{n}^{\top}\boldsymbol{\beta} \end{pmatrix} = \underbrace{\begin{pmatrix} \mathbf{x}_{1}^{\top} \\ \mathbf{x}_{2}^{\top} \\ \vdots \\ \mathbf{x}_{n}^{\top} \end{pmatrix}}_{\mathbf{X}} \boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta}$$

X is the $n \times p$ covariate matrix with *n* observations as rows.

Linear regression

Linear regression in vector form

$$oldsymbol{y} = oldsymbol{X}oldsymbol{eta} + arepsilon$$

Least squares estimate = maximum likelihood estimate

$$\hat{\boldsymbol{eta}} = (\boldsymbol{X}^{ op} \boldsymbol{X})^{-1} \boldsymbol{X}^{ op} \boldsymbol{y}$$

We now understand that

$$\mathbf{X}^{\top}\mathbf{X} = \begin{pmatrix} \sum_{i=1}^{n} x_{1i}^{2} & \sum_{i=1}^{n} x_{1i}x_{2i} & \cdots & \sum_{i=1}^{n} x_{1i}x_{pi} \\ \sum_{i=1}^{n} x_{1i}x_{2i} & \sum_{i=1}^{n} x_{2i}^{2} & \cdots & \sum_{i=1}^{n} x_{2i}x_{pi} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^{n} x_{1i}x_{pi} & \sum_{i=1}^{n} x_{2i}x_{pi} & \cdots & \sum_{i=1}^{n} x_{pi}^{2} \end{pmatrix}$$
$$\mathbf{X}^{\top}\mathbf{y} = \begin{pmatrix} \sum_{i=1}^{n} x_{1i}y_{i} \\ \sum_{i=1}^{n} x_{2i}y_{i} \\ \vdots \\ \sum_{i=1}^{n} x_{pi}y_{i} \end{pmatrix}$$

But what does $(\mathbf{X}^{\top}\mathbf{X})^{-1}$ mean? Inverse of a matrix? 😻

Matrix inverse

The inverse of regular number x is x^{-1} which is defined by

$$x^{-1}x = xx^{-1} = \frac{x}{x} = 1$$

Inverse of $p \times p$ matrix **A** is denoted by \mathbf{A}^{-1} and defined by

$$\boldsymbol{A}^{-1}\boldsymbol{A}=\boldsymbol{A}\boldsymbol{A}^{-1}=\boldsymbol{I}_{p}$$

where I_p is the $p \times p$ identity matrix

$$\boldsymbol{I}_{\boldsymbol{p}} = \left(\begin{array}{cccc} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{array} \right)$$

> A = matrix(c(2,3,3,2), 2, 2, byrow = TRUE) > invA = solve(A) > invA %*% A # returns the identity matrix

Least squares estimate

Least squares minimizes the sum of squared residuals

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Find minimum of $Q(\beta_0, \beta_1)$ by solving system of equations

$$\frac{\partial \mathbf{Q}(\beta_0, \beta_1)}{\partial \beta_0} = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i)(-1) = 0$$
$$\frac{\partial \mathbf{Q}(\beta_0, \beta_1)}{\partial \beta_1} = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i)(-x_i) = 0$$

gives the so called normal equations

$$n\bar{y} = n\beta_0 + \beta_1 n\bar{x}$$
$$\sum_{i=1}^n x_i y_i = \beta_0 n\bar{x} + \beta_1 \sum_{i=1}^n x_i^2$$

With solution

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \qquad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Least squares estimate - vector form

Sum of squared residuals in vector notation

$$Q(\boldsymbol{\beta}) = \sum_{i=1}^{n} (y_i - x_i^{\top} \boldsymbol{\beta})^2 = (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta})^{\top} (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta})$$

Set gradient vector equal to zero

$$\frac{\partial}{\partial \beta} Q(\beta) = -2 \boldsymbol{X}^{\top} (\boldsymbol{y} - \boldsymbol{X} \beta) = \boldsymbol{0}$$

gives the normal equations

$$\boldsymbol{X}^{\top} \boldsymbol{X} \boldsymbol{\beta} = \boldsymbol{X}^{\top} \boldsymbol{y}$$

Multiply both sides with the matrix inverse of $\mathbf{X}^{\top}\mathbf{X}$

$$\left(\boldsymbol{X}^{\top} \boldsymbol{X} \right)^{-1} \boldsymbol{X}^{\top} \boldsymbol{X} \boldsymbol{\beta} = \left(\boldsymbol{X}^{\top} \boldsymbol{X} \right)^{-1} \boldsymbol{X}^{\top} \boldsymbol{y}$$

gives the least squares solution

$$\hat{\boldsymbol{eta}} = (\boldsymbol{X}^{ op} \boldsymbol{X})^{-1} \boldsymbol{X}^{ op} \boldsymbol{y}$$

Gradients

Bivariate function z = f(x, y).

Partial derivative in *x*: change in *x*, holding *y* constant

$$f_x(x,y) = \frac{\partial}{\partial x} f(x,y)$$

Partial derivative in *y*: change in *y*, holding *x* constant

$$f_y(x,y) = \frac{\partial}{\partial y}f(x,y)$$

Gradient is the vector of partial derivatives

$$\left(\begin{array}{c}f_x(x,y)\\f_y(x,y)\end{array}\right)$$

General $f(x_1, \ldots, x_p)$ or $f(\mathbf{x})$. Gradient is *p*-dim vector

$$\frac{\partial}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial}{\partial x_1} f(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_p} f(\mathbf{x}) \end{pmatrix}$$

Gradients



Determinant of a square matrix

Let **A** be a 2 × 2 matrix

$$\mathbf{A} = \left(\begin{array}{cc} \mathbf{a}_{11} & \mathbf{a}_{12} \\ \mathbf{a}_{21} & \mathbf{a}_{22} \end{array}\right)$$

The determinant is the number

$$|\mathbf{A}| = a_{11}a_{22} - a_{12}a_{21}$$

Better intuition soon on why the determinant is important.

Formulas for larger matrices are complicated. Use a computer.

$$> A = matrix(c(2,3,3,2), 2, 2)$$

> det(A) # returns -5

Bivariate normal distribution

X and Y follow a bivariate normal distribution

$$(X, Y) \sim N(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$$

with joint pdf

$$f(x,y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} \\ \times \exp\left(-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x-\mu_x}{\sigma_x}\right)^2 + \left(\frac{y-\mu_y}{\sigma_y}\right)^2 - 2\rho\left(\frac{x-\mu_x}{\sigma_x}\right)\left(\frac{y-\mu_y}{\sigma_y}\right)\right]\right)$$

Parameters:

- \blacktriangleright μ_x the mean of X
- \blacktriangleright μ_y the mean of Y
- σ_x the standard deviation of X
- \triangleright σ_y the standard deviation of Y
- \triangleright ρ the correlation between X and Y

Bivariate normal distribution



Bivariate normal - widget



Multivariate normal distribution

x = $(X_1, X_2, ..., X_p)^{\top}$ and follows a multivariate normal distribution

 $m{x} \sim m{N}(m{\mu}, m{\Sigma})$

with joint pdf

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{p/2}} |\mathbf{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\top} \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

Clash in notation: *small* bold letters for random vectors.
Parameters when p = 2:

Mean vector

$$oldsymbol{\mu} = \left(egin{array}{c} \mu_1 \ \mu_2 \end{array}
ight)$$

Covariance matrix

$$\boldsymbol{\Sigma} = \left(\begin{array}{cc} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{array} \right)$$

Multivariate normal distribution

Determinant measures total variance

$$|\mathbf{\Sigma}| = \sigma_1^2 \sigma_2^2 (1 - \rho^2)$$

No correlation: |Σ| = σ₁²σ₂²
 Strong positive correlation: |Σ| small
 Strong negative correlation: |Σ| small
 The quadratic form

$$(\boldsymbol{x} - \boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1} (\boldsymbol{x} - \boldsymbol{\mu})$$

is the vector version of a squared standardized variable

$$\left(\frac{X-\mu}{\sigma}\right)^2$$