

Scientific Computing for Biologists

Linear Algebra Review II & Multiple Regression

Instructor: Paul M. Magwene

17 September 2013

Overview of Lecture

■ More Linear Algebra

- Linear combinations and Spanning Spaces
- Subspaces
- Basis vectors
- Dimension
- Rank

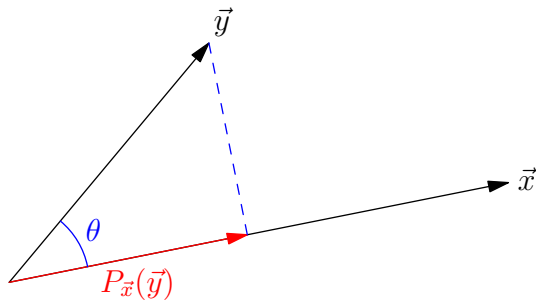
■ More on Regression

- Multiple regression
- Curvilinear regression
- Major axis regression

Hands-on Session

- Multiple regression

A bit of review



$$\cos \theta = ?$$

$$P_{\vec{x}}(\vec{y}) = ?$$

$$C_{\vec{x}}(\vec{y}) = ?$$

Space Spanned by a List of Vectors

Definition

Let X be a finite list of n -vectors. The **space spanned** by X is the set of all vectors that can be written as linear combinations of the vectors in X .

A space spanned includes the zero vector and is closed under addition and multiplication by a scalar.

Remember that a *linear combination* of vectors is an equation of the form $z = b_1x_1 + b_2x_2 + \cdots + b_px_p$

Subspaces

\mathbb{R}^n denotes the set of real n -vectors - the set of all $n \times 1$ matrices with entries from the set \mathbb{R} of real numbers.

Definition

A **subspace** of \mathbb{R}^n is a subset S of \mathbb{R}^n with the following properties:

- 1 $0 \in S$
- 2 If $u \in S$ then $ku \in S$ for all real numbers k
- 3 If $u \in S$ and $v \in S$ then $u + v \in S$

Examples of subspaces of \mathbb{R}^n :

- any space spanned by a list of vectors in \mathbb{R}^n
- the set of all solutions to an equation $Ax = 0$ where A is a $p \times n$ matrix, for any number p .

Basis

Let S be a subspace of \mathbb{R}^n . Then there is a finite list, X , of vectors from S such that S is the space spanned by X .

Let S be a subspace of \mathbb{R}^n spanned by the list (u_1, u_2, \dots, u_n) . Then there is a linearly independent sublist of (u_1, u_2, \dots, u_n) that also spans S .

Definition

A list X is a **basis** for S if:

- X is linearly independent
- S is the subspace spanned by X

Dimension

Let S be a subspace of \mathbb{R}^n .

Definition

The **dimension** of S is the number of elements in a basis for S .

Rank of a Matrix

Let A be an $n \times p$ matrix.

Definition

The **rank** of A is equal to the dimension of the row space of A which is equal to the dimension of the column space of A .

Where the row space of A is the space spanned by the list of rows of A and the column space of A is defined similarly.

Regression Models

Variable space view of multiple regression

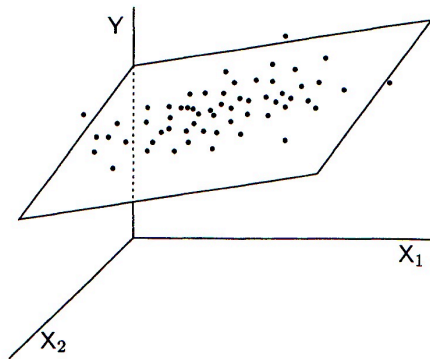


Figure 4.1: *The regression of Y onto X_1 and X_2 as a scatterplot in variable space.*

Subject Space Geometry of Multiple Regression

Multiple Regression

Let Y be a vector of values for the outcome variable. Let X_i be explanatory variables and let x_i be the mean-centered explanatory variables.

$$Y = \hat{Y} + e$$

where –

Uncentered version:

$$\hat{Y} = a1 + b_1X_1 + b_2X_2 + \dots + b_pX_p$$

Centered version:

$$\hat{y} = b_1x_1 + b_2x_2 + \dots + b_px_p$$

Statistical Model for Multiple Regression

In matrix form:

$$y = Xb + e$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} ; X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} ;$$

$$b = \begin{bmatrix} a \\ b_1 \\ b_2 \\ \vdots \\ b_p \end{bmatrix} ; e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

Estimating the Coefficients for Multiple Regression

$$y = Xb + e$$

Estimate b as:

$$b = (X^T X)^{-1} X^T y$$

Multiple Regression Loadings

The regression **loadings** should be examined as well as the regression coefficients.

Loadings are given by:

$$\cos \theta_{\vec{x}_j, \hat{\vec{y}}} = \frac{\vec{x}_j \cdot \hat{\vec{y}}}{|\vec{x}_j| |\hat{\vec{y}}|}$$

Multiple regression: Cautions and Tips

- Comparing the size of regression coefficients only makes sense if all the predictor variables have the same scale
- The predictor variables (columns of X) must be linearly independent; when they're not the variables are **multicollinear**
- Predictor variables that are **nearly multicollinear** are, perhaps, even more difficult to deal with

Why is near multicollinearity of the predictors a problem?

(a)

Nearly
collinear
pre-
dic-
tors

Figure: When predictors are nearly collinear, small differences in the vectors can result in large differences in the estimated regression.

What can I do if my predictors are (nearly) collinear?

- Drop some of the linearly dependent sets of predictors.
- Replace the linearly dependent predictors with a combined variable.
- Define orthogonal predictors, via linear combinations of the original variables (PC regression approach)
- 'Tweak' the predictor variables so that they're no longer multicollinear (Ridge regression).

Curvilinear Regression

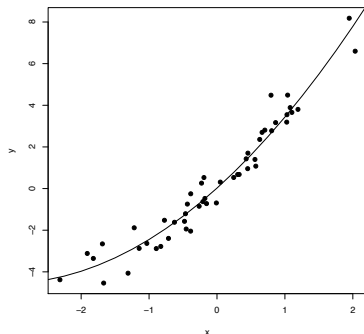
Curvilinear regression using **polynomial models** is simply multiple regression with the x_i replace by powers of x .

$$\hat{y} = b_1x + b_2x^2 + \dots + b_px^n$$

Note:

- this is still a *linear* regression (linear in the coefficients)
- best applied when a specific hypothesis justifies there use
- generally not higher than quadratic or cubic

Example of Curvilinear Regression



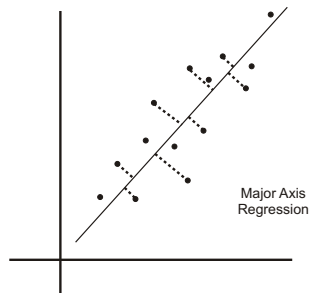
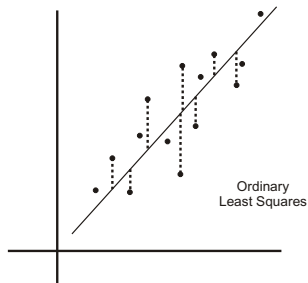
$$y = 3x + 0.5x^2 + e$$

```
lm(formula = y ~ x + I(x^2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	0.02229	0.11651	0.191	0.849	
x	2.94001	0.09693	30.331	< 2e-16	***
I(x^2)	0.47146	0.07685	6.135	1.68e-07	***

Least Squares Regression vs. Major Axis Regression



Vector Geometry of Major Axis Regression

(a)

~~OLS~~

Major

Axis

Re-

gres-

sion

Figure: Vector geometry of ordinary least-squares and major axis regression.