

# Modern Regression

Daniel Lawson University of Bristol

Lecture 02.1.2 (v1.0.1)

# Signposting

- ▶ The previous section 02.1.1 is about interpretation of Regression in general.
- ▶ This lecture contains the mathematical content for Modern Regression (Matrix representation).

# Linear algebra view of covariance

- ▶ The **covariance matrix** of a random variable  $X$
- ▶ Where  $X$  is an  $n \times 1$  matrix, i.e. a column vector,
- ▶ has entries:

$$\text{Cov}(X)_{ij} = \text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)].$$

- ▶ The matrix form for this is:

$$\Sigma = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T],$$

- ▶ Where  $\mu = \mathbb{E}[X]$ .

# Linear algebra view of correlation

- ▶ Division by standard deviations is required to correctly generalise the **scalar correlation**:

$$\text{Corr}(X, Y) = \frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}.$$

- ▶ The **matrix form** for correlation is:

$$\text{Corr}(X) = (\text{diag}(\Sigma))^{-1/2} \Sigma (\text{diag}(\Sigma))^{-1/2}$$

- ▶ The matrix inversion is not computationally challenging because it is for a **diagonal matrix**.

# Regression is analogous to linear algebra with noise

- ▶ Most problems in Linear Algebra can be seen as **solving a system of linear equations**:

$$Ax + b = 0.$$

- ▶ However, data are not usually generated from a linear model.
- ▶ We therefore typically seek the least-bad fit that we can:

$$\min ||Ax + b||_2^2 = \min \sum_{i=1}^N (Ax_i - b)^2$$

- ▶ i.e. we find  $A$  and  $b$  such that they minimise the distance (in the squared  $L_2$  norm)
- ▶ Linear Algebra is therefore a very powerful way to view regression.

# Matrix form of least squares

- ▶ Consider data  $X'$  with  $p'$  **features** (columns) and  $n$  observations.
- ▶ Given the regression problem:

$$\mathbf{y} = X'\beta' + \mathbf{b} + \mathbf{e}$$

- ▶ to find  $\beta'$  (a matrix dimension  $p' \times 1$ )
- ▶ and  $b$  to minimise 'error':
- ▶ in  $e^2 = \sum_{i=1}^n \epsilon_i^2$

# Matrix form of least squares

- ▶ We construct a simpler representation by adding a constant feature:

$$X = \begin{bmatrix} 1 & X_{11} & \cdots & X_{1p'} \\ & & \cdots & \\ 1 & X_{n1} & \cdots & X_{np'} \end{bmatrix}$$

- ▶ which has  $p = p' + 1$  features.
- ▶ We now solve the analogous equation:

$$\mathbf{y} = X\boldsymbol{\beta} + \mathbf{e}$$

- ▶ which has the same solution but is in a more convenient form.

# Mean Squared Error (MSE)

- ▶ The **prediction error** is:

$$\mathbf{e}(\beta) = \mathbf{y} - \mathbf{X}\beta$$

- ▶ And the **estimation error** can be written:

$$\text{MSE}(\beta) = \frac{1}{n} \mathbf{e}^T \mathbf{e}$$



# Minimising MSE

- ▶ Taking (vector) derivatives with respect to  $\beta$ :

$$\nabla \text{MSE}(\beta) = \frac{1}{n}(\nabla \mathbf{y}^T \mathbf{y} - 2\nabla \beta^T x^T \mathbf{y} + \nabla \beta^T \mathbf{X}^T \mathbf{X} \beta) \quad (1)$$

$$= \frac{1}{n}(0 - 2x^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \beta) \quad (2)$$

- ▶ which is zero at the optimum  $\hat{\beta}$ :

$$\mathbf{X}^T \mathbf{X} \hat{\beta} - \mathbf{X}^T \mathbf{y} = 0$$

- ▶ with the solution:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

- ▶ Exercise: For the case  $p' = 1$ , check that this solution is the same as you can find in regular linear algebra textbooks.

# The Hat Matrix

- ▶ There is an important and **response independent** quantity hidden in the prediction:

$$H = X(X^T X)^{-1} X^T$$

- ▶ The fitted values are:

$$\hat{y} = X\hat{\beta} = X(X^T X)^{-1} X^T y = Hy$$

- ▶  $H$  is dimension  $N \times N$
- ▶  $H$  “projects”  $y$  into the fitted value space  $\hat{y}$

# Properties of the Hat Matrix

- ▶ **Influence:**  $\frac{\partial \hat{y}_i}{\partial y_j} = H_{ij}$ . So  $H$  controls how much a change in one observation changes the estimates of each other point.
- ▶ **symmetry:**  $H^T = H$ . So influence is symmetric.
- ▶ **Idempotency:**  $H^2 = H$ . So the predicted value for any projected point is the predicted value itself.
- ▶ You should read up on these and other vector algebra properties.

# Residuals and the Hat Matrix

- ▶ The residuals can be written:

$$e = y - Hy = (I - H)y$$

- ▶  $I - H$  is also symmetric and idempotent, and can also be interpreted in terms of Influence.
- ▶ Because of this,

$$\text{MSE}(\hat{\beta}) = \frac{1}{n} \mathbf{y}^T (I - H)^T (I - H) \mathbf{y} = \frac{1}{n} \mathbf{y}^T (I - H) \mathbf{y}$$

# Expectations

- ▶ If the data were generated by our model(!) then they are described by a random variable  $\mathbf{Y}$ :

$$\mathbf{Y} = \mathbf{x}\beta + \epsilon$$

- ▶ where  $\epsilon$  is an  $n \times 1$  matrix of RVs with mean  $\mathbf{0}$  and covariance  $\sigma^2 \mathbf{I}$ .
- ▶ From this it is straightforward to show that the **fitted values are unbiased**:

$$\mathbb{E}[\hat{y}] = \mathbb{E}[\mathbf{H}\mathbf{Y}] = \mathbf{x}\beta$$

- ▶ using the properties of Expectations with the symmetry and idempotency of  $\mathbf{H}$ .

# Covariance

- ▶ Similarly, it is straightforward to show that

$$\text{Var}[\hat{y}] = \sigma^2 H$$

using the properties of Variances with the symmetry and idempotency of  $H$ .

# Reflection

- ▶ By the end of the course, you should:
  - ▶ Be able to define **correlation** and **regression** in multivariate context
  - ▶ Be able to perform basic calculations using these concepts
  - ▶ Be able to extend intuition about their application.
- ▶ This is something worth reading up on
  - ▶ You should really understand univariate regression

# Signposting

- ▶ The mathematics behind Modern Regression is entirely analogous to the mathematics of a huge range of advanced, scalable Machine Learning tools.
- ▶ This is one of the places you should do your homework.



# References

There is a lot more technical detail in Cosma Shalizi's course on **Modern Regression**:

Modern Regression, by Cosma Shalizi

<http://www.stat.cmu.edu/~cshalizi/mreg/15/lectures/>