

Towards Modern Regression

Daniel Lawson University of Bristol

Lecture 02.1 (v2.0.0)

Signposting

- ▶ This lecture covers:
 - ▶ Classical regression
 - ▶ Towards Modern Regression - the vectorised version, which uses Matrix algebra.
 - ▶ Leave-one-out Cross Validation
- ▶ The maths here underpins almost all modern data science.

Correlation and Covariance

- ▶ **Correlation** and **Covariance** are quantifications of a relationship between x and y .
- ▶ They quantify the **linear relationship**.
- ▶ They ask, “How does **variation** in x and y associate?”
- ▶ Consequently, they are purely descriptive and do not attempt to establish any cause and effect.
- ▶ Covariance is a generalisation of variance; it summarises the 2-D marginals of high dimensional data.

Covariance

- ▶ A reminder: covariance is simply the second (central) moment:

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

- ▶ it is straightforward to show that

$$\text{cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X] \mathbb{E}[Y].$$

- ▶ Recall that we typically use **unbiased** estimators which often slightly differ from natural theoretical analogue. The **sample covariance** is:

$$\text{cov}(X, Y) = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})$$

Correlation

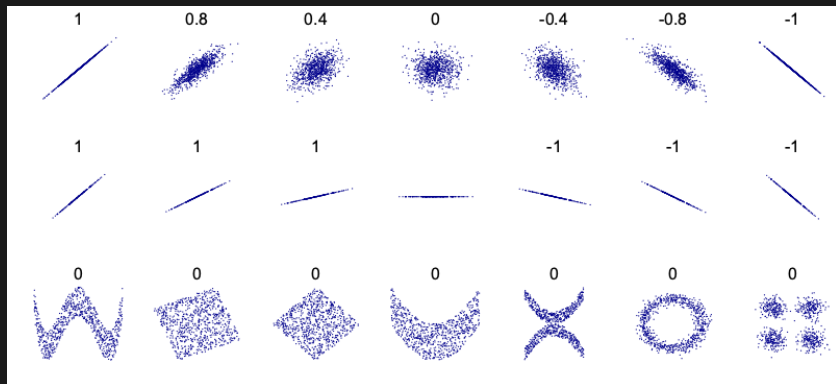
- ▶ Correlation is simply a **normalised measure** of covariance.

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

- ▶ It takes values between -1 and 1.
- ▶ Sample correlation uses the unbiased estimator of covariance, to account for the number of degrees of freedom in the data.
- ▶ What should we take the correlation of?
 - ▶ See rank correlation, canonical correlation, etc.

Examples

From Wikipedia: Correlation_and_dependence



Regression

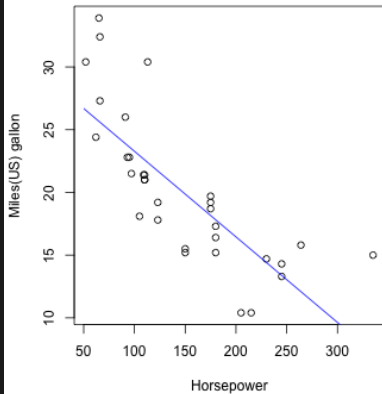
- ▶ **Regression**, considers the relationship of a response variable as determined by one or more explanatory variables.
 - ▶ Regression is designed to help **make predictions** of y when we observe x .
 - ▶ It is **not** a joint model of x and y .
 - ▶ It predicts the *best guess*.
 - ▶ There is a probabilistic interpretation based on Normal Distributions.
- ▶ Regression is often used as a tool to establish causality. . .
 - ▶ A and B share a causal relationship if a regression for B given A, conditional on C ($C=\text{everything else}$), has an association
 - ▶ This does not resolve whether A causes B, or B causes A
 - ▶ Since we don't measure **everything else**, regression rarely establishes causality!
 - ▶ Assumptions are needed to make a causal connection. This is known as **causal inference** and there are frameworks to establish causality.

Discrete predictors

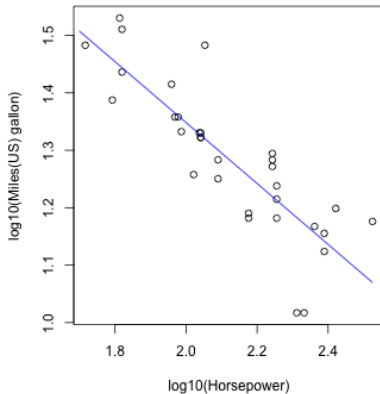
- ▶ If you include categorical/factor predictors, each **level** or unique value is used as a binary predictor.
- ▶ This is called **One Hot Encoding**.

Regression example

a) Scatter plot (mtcars MPG-HP)



b) Log-scale Scatter plot (mtcars MPG-HP)



Multiple Regression example

```
> lm(mpg ~ cyl + hp + wt, data=mtcars) %>% summary
```

Call:

```
lm(formula = mpg ~ cyl + hp + wt, data = mtcars)
```

Residuals:

Min	1Q	Median	3Q	Max
-3.9290	-1.5598	-0.5311	1.1850	5.8986

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	38.75179	1.78686	21.687	< 2e-16	***
cyl	-0.94162	0.55092	-1.709	0.098480	.
hp	-0.01804	0.01188	-1.519	0.140015	
wt	-3.16697	0.74058	-4.276	0.000199	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Important measures of regression

- ▶ **R squared** (and adjusted R squared): variance explained/total variance. This tells us how predictable y is.
- ▶ The coefficients β_i .
 - ▶ These should be compared to their error $\hat{\sigma}_i$.
 - ▶ The ratio is a t-value ($t_i = \beta_i / \hat{\sigma}_i$) from which a p-value can be calculated.
- ▶ F statistic and F test p-value:
 - ▶ F is the ratio of the explained to unexplained variance, accounting for the **degrees of freedom**.
 - ▶ The full model compared to a null in which there are no explanatory variables.
 - ▶ Used in variable selection, ANOVA, etc.

Vector Notation

- ▶ There are several choices of convention that we have to make
- ▶ Vectors of length k are also matrices, but are they $k \times 1$ or $1 \times k$?
- ▶ We use $k \times 1$, i.e. column vectors
- ▶ Similarly there are choices about matrix derivatives
- ▶ We use derivative with respect to a column vector as a row vector
- ▶ Some resources will have everything transposed as a consequence

Linear algebra view of covariance

- ▶ The **covariance matrix** of a random variable X
- ▶ Where X is a vector-valued RV with length k ,
- ▶ 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**:

$$XA + \mathbf{b} = \mathbf{0}.$$

- ▶ Where X is an n by p matrix of data,
- ▶ A is an p by 1 matrix of coefficients,
- ▶ and $-\mathbf{b}$ is a n -vector of target values.
- ▶ However, data are not *usually* generated from a linear model.
- ▶ We therefore typically seek the least-bad fit that we can:

$$\operatorname{argmin} \|XA + \mathbf{b}\|_2^2 = \sum_{i=1}^N (\mathbf{x}_i A + b_i)^2$$

- ▶ i.e. we find A and \mathbf{b} such that they minimise the distance (in the squared L_2 norm)
- ▶ Linear algebra allows this very effectively!
- ▶ 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:
 - ▶ β' (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:

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

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

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

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

Mean Squared Error (MSE)

- ▶ The **prediction error** is:

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

- ▶ Using the notation that \mathbf{e} is a p by 1 matrix
- ▶ The **estimation error** is written in matrix form:

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

- ▶ Why? $\mathbf{e}^T \mathbf{e} = \sum_{i=1}^n e_i^2$
 - ▶ Hence $\text{MSE}(\beta)$ is a 1×1 matrix, i.e. a scalar, and $|\text{MSE}(\beta)| = \text{MSE}(\beta)$.
 - ▶ Noticing this sort of thing makes the matrix algebra easier.
- ▶ We want to minimise this MSE with respect to the parameters β .

How to do the Matrix Algebra

Lecture 13 of Cosma Shalizi's notes is a really helpful reminder!

- ▶ Look at the Matrix Algebra Cheat Sheet - specifically:
 - ▶ How does a transpose work?
 - ▶ How do you re-order elements?
 - ▶ How does a gradient work in linear and quadratic forms?

Minimising MSE

- ▶ Taking (vector) derivatives with respect to β :

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

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

- ▶ 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:

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

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

$$\mathbf{e} = \mathbf{y} - \mathbf{H}\mathbf{y} = (\mathbf{I} - \mathbf{H})\mathbf{y}$$

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

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

Expectations

- ▶ If the data were generated by our model(!) then they are described by an RV \mathbf{Y} (an n -vector):

$$\mathbf{Y}_i = \mathbf{x}_i\beta + \epsilon_i$$

- ▶ \mathbf{x}_i is still a vector but *not* a Random Variable!
- ▶ ϵ is an $n \times 1$ matrix of RVs with mean $\mathbf{0}$ and covariance $\sigma_s^2 \mathbf{I}$.
- ▶ From this it is straightforward to show that the **fitted values are unbiased**:

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

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

Covariance

- ▶ Similarly, it is straightforward to show that

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

using the properties of Variances with the symmetry and idempotency of \mathbf{H} .

- ▶ In other words, the covariance of the fitted values is determined entirely by the structure of the covariates, via the Hat matrix.

Motivation: Residuals

- ▶ The **residual sum of squares** for n predictions of a univariate y :

$$R^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- ▶ The expected value of the prediction error $\mathbb{E}(e^2) = R^2/n$.
- ▶ What happens if **compare two models** M_1 and M_2 , where M_1 is a subset of M_2 ?

Linear Models - Model selection

- ▶ For illustration, consider

$$Y = \mathbf{x}_1 A_1 + \epsilon_1$$

- ▶ and

$$Y = \mathbf{x}_1 A_1 + \mathbf{x}_2 A_2 + \epsilon_2.$$

- ▶ Unless $\mathbf{x}_2 = 0$ or $\mathbf{x}_2 \equiv \mathbf{x}_1$, then ϵ_2^2 will be smaller than ϵ_1^2 .
 - ▶ This is an example of a more general rule: **larger models always have better predictions.**
- ▶ So prediction error is OK to use to fit models with the same dimension, but is incomplete for **model selection.**

Cross-Validation Motivation

- ▶ Usually we are not interested in properties of **our sample**.
- ▶ We instead wish to know how our inference will generalise to **new samples**.
- ▶ The most straight forward way to predict how a model generalises is to test in **held-out data**.
- ▶ **Cross Validation** is a procedure to leave-out some data for testing.
- ▶ How much data?
 - ▶ **Leave-one-out Cross-Validation** (LOOCV) leaves out one datapoint at a time for testing.
 - ▶ **k-Fold Cross Validation** (k-fold CV) keeps a fraction $(k - 1)/k$ of the data for learning parameters and $1/k$ for testing.

Prediction accuracy in linear regression

- In linear regression, the errors are

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\beta = \mathbf{y} - \mathbf{H}\mathbf{y} = \mathbf{y} - \hat{\mathbf{y}}$$

- Recall the \mathbf{H} matrix describes the influence of y_i on \hat{y}_j , i.e. that y_i and \hat{y}_j covary.
- We show in Worksheet 2.2A that the expected MSE for the i -th datapoint is:

$$\mathbb{E}(e_i^2) = \mathbb{E} \left[(y_i - \hat{y}_i)^T (y_i - \hat{y}_i) \right] = \mathbb{E} \left[(y_i - \hat{y}_i)^2 \right] \quad (3)$$

$$= \text{Var}[y_i] + \text{Var}[\hat{y}_i] - 2\text{Cov}[y_i, \hat{y}_i] + [\mathbb{E}(y_i) - \mathbb{E}(\hat{y}_i)]^2 \quad (4)$$

- This is shown by rearranging the formula for $\text{Var}[y_i - \hat{y}_i]$

Out-of-sample prediction accuracy in linear regression

- ▶ We can write the same thing when predicting an **out-of-sample** y'_i :

$$\mathbb{E}(e_i'^2) = \mathbb{E} \left[(y'_i - \hat{y}_i)^T (y'_i - \hat{y}_i) \right] \quad (5)$$

$$= \text{Var}[y'_i] + \text{Var}[\hat{y}_i] - 2\text{Cov}[y'_i, \hat{y}_i] + [\mathbb{E}(y'_i) - \mathbb{E}(\hat{y}_i)]^2 \quad (6)$$

- ▶ But out-of-sample, $\text{Cov}[y'_i, \hat{y}_i] = 0$ whereas within-sample, $\text{Cov}[y_i, \hat{y}_i] \neq 0$.
- ▶ Therefore:

$$\mathbb{E}(e_i'^2) = \mathbb{E}(e_i^2) + 2\text{Cov}[y_i, \hat{y}_i]$$

Quantifying Out-of-sample prediction accuracy

- ▶ Fortunately we already did the work required to describe this:

$$\text{Cov}[y_i, \hat{y}_i] = \sigma^2 H_{ii}$$

- ▶ The mean out-of-sample prediction error is

$$\mathbb{E}(e'^2) = n^{-1} \sum_{i=1}^n e_i'^2 = n^{-1} \sum_{i=1}^n e_i^2 + 2n^{-1} \text{tr}(H)$$

- ▶ We show in Worksheet 2.2A that $\text{tr}(H) = \sigma^2 p$ where p =number of predictors.
- ▶ The **optimism** is defined as $2n^{-1}\sigma^2 p$.
- ▶ The optimism grows with σ^2 and p but shrinks with n . It is used to define the **model selection criteria** ΔC_p which is minimised:

$$\Delta C_p = MSE_1 - MSE_2 + \frac{2}{n} \hat{\sigma}^2 (p_1 - p_2)$$

Linear model optimism and AIC

- ▶ Minimising **Akaike's Information Criterion**:

$$AIC = -2\mathbb{L}(\hat{\theta}) + 2\text{Dim}(\theta)$$

- ▶ reduces to maximising ΔC_p when the Likelihood \mathbb{L} is a Normal distribution.
- ▶ There are many other Information Criteria...

LOOCV

- ▶ We write a statistic \hat{s} based on all data $\{y\}$ except i as $\hat{s}^{(-i)}$ and the data is $\{y\}^{(-i)}$.
- ▶ For a general **loss function** we can write:

$$LOOCV = \frac{1}{n} \sum_{i=1}^n \text{Loss} \left(y_i; \hat{\theta} | y^{(-i)} \right)$$

- ▶ i.e. we evaluate the loss function for each datapoint using the estimate from the remaining data.
- ▶ NB A loss function is something that we choose the parameters θ to minimise. It can be:
 - ▶ the MSE,
 - ▶ the (negative log) likelihood,
 - ▶ a penalised version of these,
 - ▶ or any other convenient quantity.

LOOCV for linear models

- ▶ For the MSE of a linear model we can write:

$$LOOCV = \frac{1}{n} \sum_{i=1}^n \left(y_i - \hat{y}_i^{(-i)} \right)^2$$

- ▶ It is not particularly straightforward¹ to show that:

$$LOOCV = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - H_{ii}} \right)^2$$

- ▶ This is a very important quantity, often called the Studentized residual
- ▶ i.e. the LOOCV can be directly computed from a regression containing all data, by “downweighting” low-leverage data and upweighting high-leverage (hard to predict) data.

¹Our references avoid proving this, but do discuss the motivation. Proofs are available but beyond scope.

Leave-one-out Cross-Validation

- ▶ Leaving out a single datapoint is going to be insufficient unless the **data are independent**.
- ▶ The real world is rarely completely independent.
- ▶ However, there is often a computationally convenient way to compute LOOCV, and it is still better than leaving nothing out. It converges to C_p for large n .
- ▶ Analogous tricks work for:
 - ▶ **Linear models** including **Best Linear Unbiased Predictors** (BLUPs)
 - ▶ **Kernel methods**
 - ▶ **Nearest neighbour** methods
 - ▶ And others

Asymptotics

- ▶ Here are some facts about the asymptotic behaviour of LOOCV:
 - ▶ As $n \rightarrow \infty$, the expected out-of-sample MSE of the model picked by LOO cross-validation is **close to that of the best model** considered.
 - ▶ As $n \rightarrow \infty$, if the true model is among those being compared, LOOCV tends to pick a **strictly larger model** than the truth.
- ▶ LOOCV is not the right tool for choosing the **right model**.
- ▶ It is however an excellent tool for choosing the model with the best out-of-sample **predictive power**.
- ▶ ...when the data to be predicted come from the **same distribution as the data!**

Implications

- ▶ Matrix form is a massive simplification of complex algebra
- ▶ It is easy to check that e.g. dimensions make sense
- ▶ These vector calculations are repeated in many machine-learning methods
- ▶ The details change but the principle remains
- ▶ Linear-Algebra loss minimisation techniques are extremely important
- ▶ They often sit inside a wider argument, e.g. updated conditional on some other parameters

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.
 - ▶ Be able to follow the reasoning in a paper where things get complicated.
- ▶ Matrix algebra is worth reading up on!
 - ▶ Describe it for example in your assessments' reflection.

Signposting

- ▶ Make sure to look at **02.1-Regression.R**
- ▶ The mathematics behind Modern Regression is analogous to the mathematics underpinning scalable Machine Learning. **It is very important.**
- ▶ For accessible material see Cosma Shalizi's Modern Regression Lectures (Lectures 13-14)
- ▶ Further reading in chapters 2.3 and 3.2 of The Elements of Statistical Learning: Data Mining, Inference, and Prediction (Friedman, Hastie and Tibshirani)
- ▶ Next up: 2.2 Statistical Testing