# Regression Modelling

Maxwell Munford
PhD Biomechanics, Mechanical Engineering

# Objectives

- ▷ Correlation
- ▷ Linear Regression
- ▷ Multi-Variable Regression
- ▷ Interpretation and Application
- ▷ NOT focussing on theory

# What Is Regression?

… and why do we use it?

4

# What Is Regression?

# What Is Regression?

▹ A statistical analysis that attempts to predict the effects of one or more variables on another variable

▹ Correlation is a mutual relationship between 2 or more variables

# The Process

- ▷ Scatter plots
- ▷ Measure the degree of linearity between two variables
- ▷ Quantify this relationship

# What do we assume?

▷ Assume a causal relationship

▷ Equation of line of best fit

▷ Test the significance

▷ Analyse residuals

▷ Predication

Analyse model

# 8

# Why Do We Start With Scatter Plots?

54.26
X Mean

47.83
Y Mean

16.76
X Standard Deviation
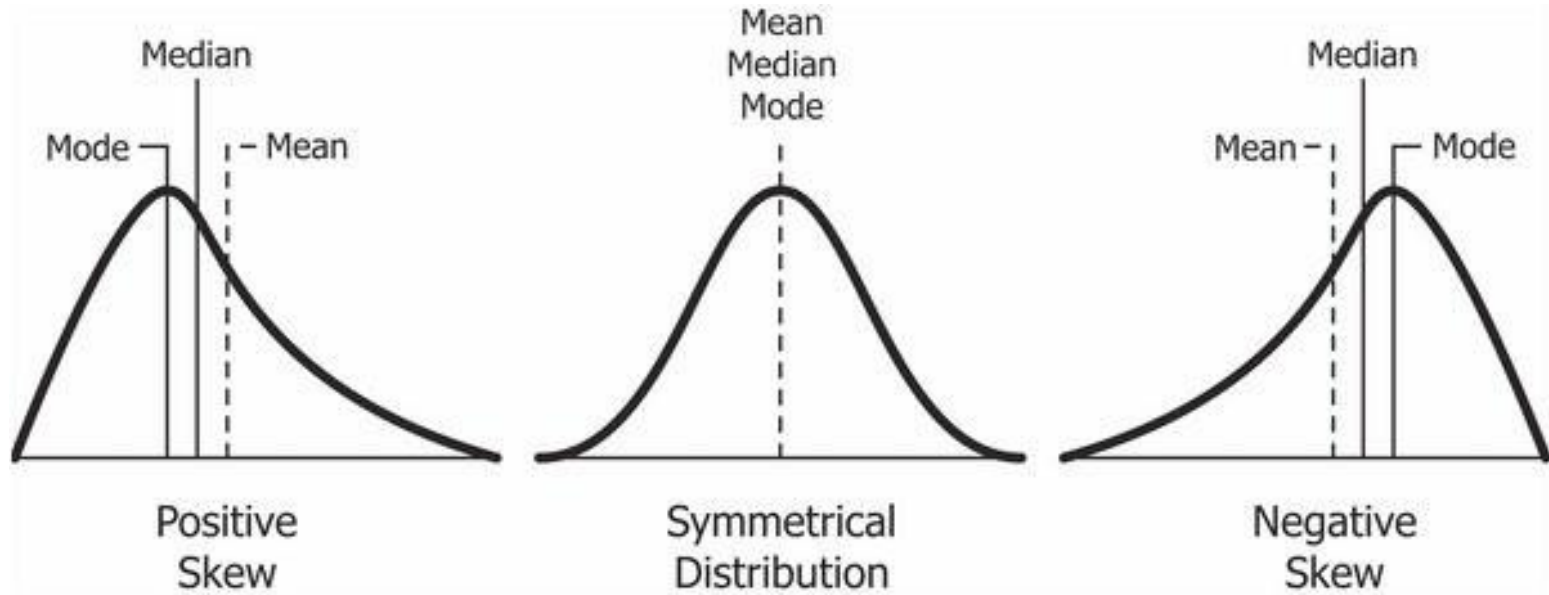
26.93
Y Standard Deviation

−0.06
Correlation Coefficient

# Scatter Plots

▷ Plot a scatter diagram and look for evidence of linear trend

# Skew

▹ Skewness, in statistics, is the degree of distortion from the symmetrical bell curve in a probability distribution.

▹ Can be positive, negative or zero, to a varying degree.

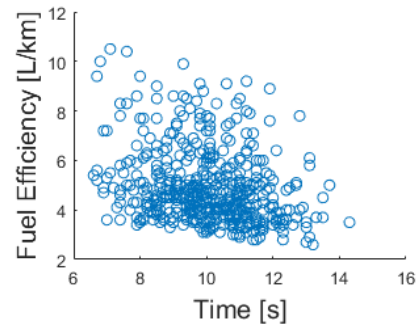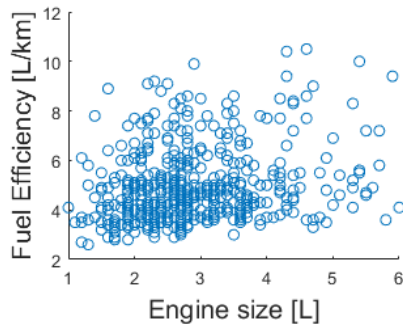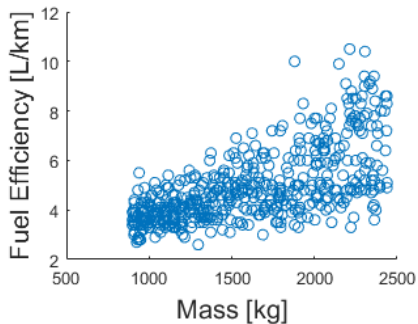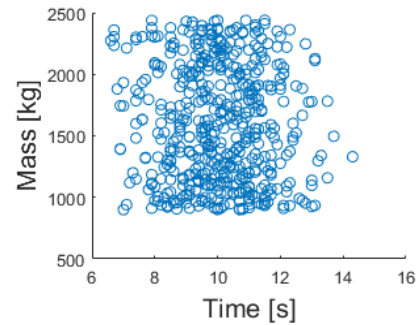▹ Helps to consider the extremes of data, not just the averages, by standardising data.
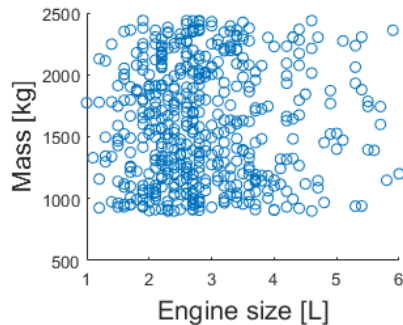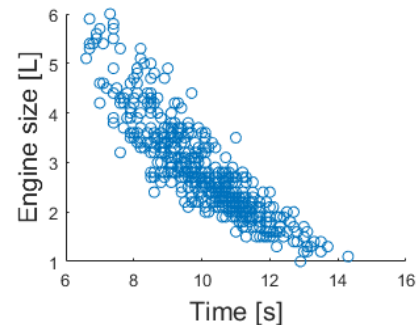
# Skew

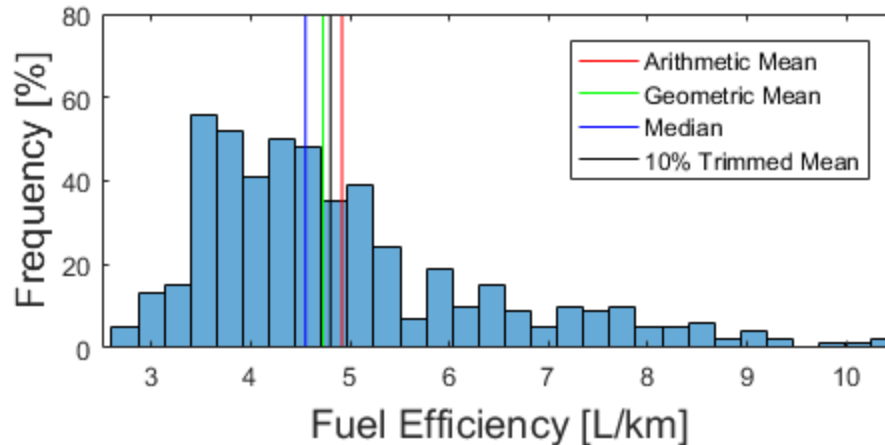# Car Emissions Data Task

## Regression Analysis

Generate and evaluate a model for fuel efficiency in terms of time, vehicle mass, engine size, fuel type and colour

## Scatter Plots

▷ Looking for possible correlation

Histogram

▹ Looking for skew in the response variable
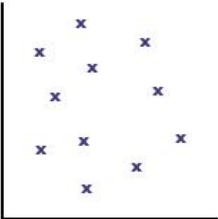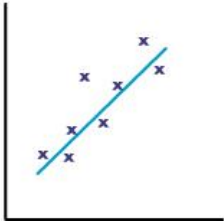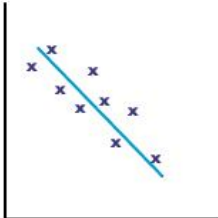
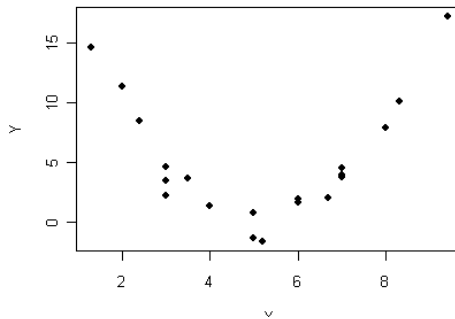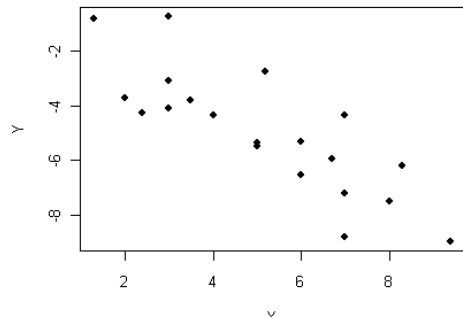▹ Positive skew        =>        mode < median < mean

# Initial Thoughts...

▷ Correlation exists between fuel efficiency, vehicle mass, acceleration time and engine size and engine size and acceleration time.

▷ Weaker correlation between fuel efficiency and engine size

▷ Positive skew, so we should standardise by the geometric mean, $x' = \frac{x - \bar{x}_{Geometric}}{\sigma_{Geometric}}$

## Correlation

- Having decided that $x_i$ and $y_i$ are paired
- We want to study the relationship between them

- The Correlation coefficient measures the degree of linear relation between x and y

# Correlation

| R | Results | Correlation | Example |
|---|---------|-------------|---------|
| 0 | Uncorrelated | No correlation |  |
| >0 | Positively correlated | $y = \beta x + \alpha$ |  |
| <0 | Negatively corerlated | $y = -\beta x + \alpha$ |  |

# Correlation

# Correlation



correlation = 0.78

# Correlation

# Correlation

# Correlation

# Correlation

# Correlation

# Correlation

- ▹ Not necessarily interested in the correlation coefficient

- ▹ Where does weaker correlation exist

- ▹ Interaction term may be needed between engine size and time

Regression

- $y_i = \alpha + \beta x_i$
- Let a = estimated $\alpha$
  b = estimated β



- a = average fuel consumption at height = 0
- b = increase in fuel consumption for 1cm increase in height

▷ Our model is not perfect

▷ $y_i = \alpha + \beta x_i + \epsilon_i$

▷ Residuals are the error between our model predictions and the actual data

▷ We assume these are normally distributed

# Regression: Least Squares

▸ Find the line of best fit (estimate $\alpha$ and $\beta$) to minimise this error

▸ $\sum_{i=1}^{n} \epsilon_i^2 = \sum_{i=1}^{n} (y_i - \alpha - \beta x_i)^2$

▸ MATLAB, Python and excel all have libraries to solve this

# Regression: Least Squares

▷ $\sum_{i=1}^{n} \epsilon_i^2 = \sum_{i=1}^{n} (y_i - \alpha - \beta x_i)^2$

▷ This is an optimisation in terms of $\alpha$ and $\beta$

▷ The resulting model predicts $y$ values given $x$

▷ $y = a + bx$

# Regression: Confidence Interval

- How accurate are our estimated values of $a$ and $b$?
- Confidence intervals based on the standard error
- 95% CI in $b$ calculated from t scores
- $( b - t\,SE(b) \, , \, b + t\,SE(b) )$

| df/p | 0.40 | 0.25 | 0.10 | 0.05 | 0.025 | 0.01 | 0.005 | 0.0005 |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.324920 | 1.000000 | 3.077684 | 6.313752 | 12.70620 | 31.82052 | 63.65674 | 636.6192 |
| 2 | 0.288675 | 0.816497 | 1.885618 | 2.919986 | 4.30265 | 6.96456 | 9.92484 | 31.5991 |
| 3 | 0.276671 | 0.764892 | 1.637744 | 2.353363 | 3.18245 | 4.54070 | 5.84091 | 12.9240 |
| 4 | 0.270722 | 0.740697 | 1.533206 | 2.131847 | 2.77645 | 3.74695 | 4.60409 | 8.6103 |
| 5 | 0.267181 | 0.726687 | 1.475884 | 2.015048 | 2.57058 | 3.36493 | 4.03214 | 6.8688 |
| 6 | 0.264835 | 0.717558 | 1.439756 | 1.943180 | 2.44691 | 3.14267 | 3.70743 | 5.9588 |
| 7 | 0.263167 | 0.711142 | 1.414924 | 1.894579 | 2.36462 | 2.99795 | 3.49948 | 5.4079 |
| 8 | 0.261921 | 0.706387 | 1.396815 | 1.859548 | 2.30600 | 2.89646 | 3.35539 | 5.0413 |
| 9 | 0.260955 | 0.702722 | 1.383029 | 1.833113 | 2.26216 | 2.82144 | 3.24984 | 4.7809 |
| 10 | 0.260185 | 0.699812 | 1.372184 | 1.812461 | 2.22814 | 2.76377 | 3.16927 | 4.5869 |
| 11 | 0.259556 | 0.697445 | 1.363430 | 1.795885 | 2.20099 | 2.71808 | 3.10581 | 4.4370 |
| 12 | 0.259033 | 0.695483 | 1.356217 | 1.782288 | 2.17881 | 2.68100 | 3.05454 | 4.3178 |
| 13 | 0.258591 | 0.693829 | 1.350171 | 1.770933 | 2.16037 | 2.65031 | 3.01228 | 4.2208 |
| 14 | 0.258213 | 0.692417 | 1.345030 | 1.761310 | 2.14479 | 2.62449 | 2.97684 | 4.1405 |
| 15 | 0.257885 | 0.691197 | 1.340606 | 1.753050 | 2.13145 | 2.60248 | 2.94671 | 4.0728 |
| 16 | 0.257599 | 0.690132 | 1.336757 | 1.745884 | 2.11991 | 2.58349 | 2.92078 | 4.0150 |
| 17 | 0.257347 | 0.689195 | 1.333379 | 1.739607 | 2.10982 | 2.56693 | 2.89823 | 3.9651 |
| 18 | 0.257123 | 0.688364 | 1.330391 | 1.734064 | 2.10092 | 2.55238 | 2.87844 | 3.9216 |
| 19 | 0.256923 | 0.687621 | 1.327728 | 1.729133 | 2.09302 | 2.53948 | 2.86093 | 3.8834 |
| 20 | 0.256743 | 0.686954 | 1.325341 | 1.724718 | 2.08596 | 2.52798 | 2.84534 | 3.8495 |
| 21 | 0.256580 | 0.686352 | 1.323188 | 1.720743 | 2.07961 | 2.51765 | 2.83136 | 3.8193 |
| 22 | 0.256432 | 0.685805 | 1.321237 | 1.717144 | 2.07387 | 2.50832 | 2.81876 | 3.7921 |
| 23 | 0.256297 | 0.685306 | 1.319460 | 1.713872 | 2.06866 | 2.49987 | 2.80734 | 3.7676 |
| 24 | 0.256173 | 0.684850 | 1.317836 | 1.710882 | 2.06390 | 2.49216 | 2.79694 | 3.7454 |
| 25 | 0.256060 | 0.684430 | 1.316345 | 1.708141 | 2.05954 | 2.48511 | 2.78744 | 3.7251 |
| 26 | 0.255955 | 0.684043 | 1.314972 | 1.705618 | 2.05553 | 2.47863 | 2.77871 | 3.7066 |
| 27 | 0.255858 | 0.683685 | 1.313703 | 1.703288 | 2.05183 | 2.47266 | 2.77068 | 3.6896 |
| 28 | 0.255768 | 0.683353 | 1.312527 | 1.701131 | 2.04841 | 2.46714 | 2.76326 | 3.6739 |
| 29 | 0.255684 | 0.683044 | 1.311434 | 1.699127 | 2.04523 | 2.46202 | 2.75639 | 3.6594 |
| 30 | 0.255605 | 0.682756 | 1.310415 | 1.697261 | 2.04227 | 2.45726 | 2.75000 | 3.6460 |
| z | 0.253347 | 0.674490 | 1.281552 | 1.644854 | 1.95996 | 2.32635 | 2.57583 | 3.2905 |
| CI | ——— | ——— | 80% | 90% | 95% | 98% | 99% | 99.9% |

# Warning

▹ Do not use your model to predict data outside the range of values in the <mark>domain</mark> of $x$

▹ Be cautious of overfitting

Regression with a Single Predictor Variable

▸ Model with one predictor:

▸ $Efficiency = \alpha + \beta Mass$

# Regression with a Single Predictor Variable

▹ Model with one predictor:

▹ $Efficiency = \alpha + \beta Mass$

▹ Rsquared = 0.4157

▹ MSE = 0.7402

▹ AIC = 12705

# What do we think of this model?

- ▹ Rsquared = 0.4157

- ▹ MSE = 0.7402

- ▹ AIC = 1271

# Multiple Variable Linear Regression?

# Multiple Explanatory Variables

- ▹ Multiple variables which simultaneously affect output variable

- ▹ Interpretation can become increasingly difficult with more variables

- ▹ Prevent cofounding and reduce residual variation

# Multiple Explanatory Variables

- $y_i = \alpha + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \beta_2 x_{2,i} + \cdots + \epsilon_i$

- $i$ = number of observations
- $x_{1,i} = i^{th}$ observation of the 1st variable
- $x_{2,i} = i^{th}$ observation of the 2nd variable
- $\beta_1$ = the increase in $y$ for a unit increase in $x_1$

# Categorical Variables

▹ Takes one of a limited number of values

▹ Binary variables take values 0 or 1

## Binary Variables

▸ $y_i = \alpha + \beta_1 x_i + \epsilon_i$

▸ This model fits the mean of $y$ for each category of $x$

▸ $\alpha$ = mean $y_i$ in 1st group

▸ $\alpha + \beta$ = mean $y_i$ in 2nd group

▸ $\beta$ = difference in $y_i$ between groups

# Categorical Variables

▹ Box plots of our data against known categorical variables

# Categorical Variables

▸ Box plots of our data against known categorical variables

# Categorical Variables

- ▹ Fuel efficiency looks likely to depend on Fuel type to some degree

- ▹ Fuel type could be a suitable predictor variable in the model

# Regression: Least Squares with Multiple Variables

▸ $\sum_{i=1}^{n} \epsilon_i^2 = \sum_{i=1}^{n} (y_i - \alpha - \beta_1 x_{1,i} - \beta_2 x_{2,i} - \beta_2 x_{2,i})^2$

▸ This is an optimisation in terms of $\alpha$ and $\beta_{1,2,3}$ ...

▸ The resulting model predicts $y$ values given $x$

# Regression with Multiple Variables

▸ Regression coefficients ($b$) report the effect of each variable while holding all others at their average values

# Confounding

▸ Consider a researcher attempting to assess the effectiveness of drug X, from population data in which drug usage was a patient's choice.

▸ Data show that gender differences influence a patient's choice of drug as well as their chances of recovery (Y).

▸ In this scenario, gender z confounds the relationship between X and Y since Z is a cause of both X and Y.

# Model Interactions

▸ $y_i = \alpha + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \beta_3 x_{3,i} + \beta_4 (x_{2,i}\, x_{3,i}) + \cdots + \epsilon_i$

▸ $x_{3,i}$ is a categorical variable

▸ $\beta_2$ = increase in $y_i$ for unit increase in $x_2$

▸ $\beta_3 + \beta_4$ = increase in $y_i$ for unit increase in $x_2$ for an observation in the categorical variable $x_3$

# Model Interactions

▸ $ArmLength_i = \alpha + \beta_1 Age_i + \beta_2 Height_i + \beta_3 Gender_i$
$$+\beta_4(Height_i Gender_i) + \epsilon_i$$

▸ Gender is the categorical variable

▸ $\beta_2$ $\qquad$ = increase in $ArmLength$ for unit increase in $Height$

▸ $\beta_2 + \beta_3$ $\qquad$ = increase in $ArmLength$ for unit increase in $Height$
$\qquad\qquad$ for a $Female$

# Model Interactions

▹ <mark>If the interaction term is significant then we should include both the interaction and individual terms</mark>
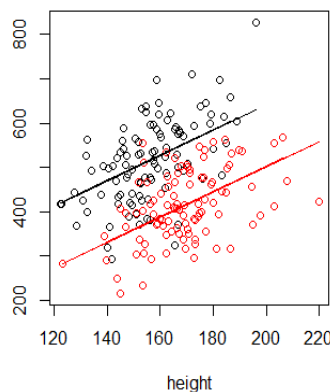
▹ The effect of gender is the different between males and females in this model on the value of height

▹ This must be found considering the interaction term

# Model Interactions

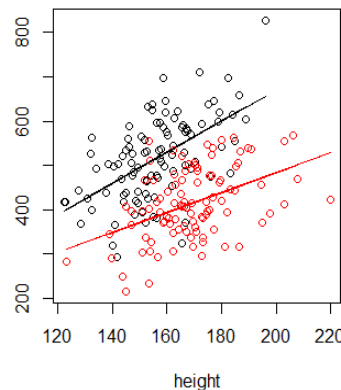▷ Assuming the effect of height is significant we can have three models depending on whether gender and the interaction are significant.



$$ArmLength_i = \alpha + \beta_1 Age_i$$
$$+\beta_2 Height_i$$

$$ArmLength_i = \alpha + \beta_1 Age_i$$
$$+\beta_2 Height_i + \beta_3 Gender_i$$

$$ArmLength_i = \alpha + \beta_1 Age_i$$
$$+\beta_2 Height_i + \beta_3 Gender_i$$
$$+\beta_4 (Height_i Gender_i)$$

# Improve on the Previous Model by Including more Terms

▸ Model with one predictor:

▸ $Efficiency = \alpha + \beta_1 Mass + \beta_2 FuelType$

▸ Rsquared = 0.6540
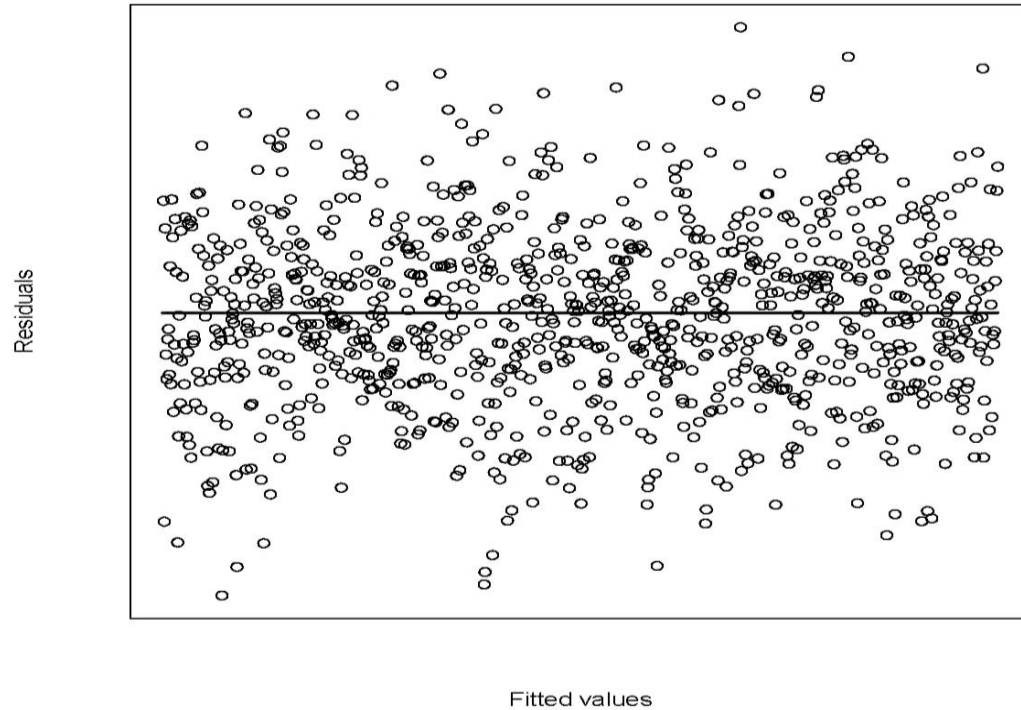
▸ MSE = 0.4383

▸ AIC = 1010

# Model Checking

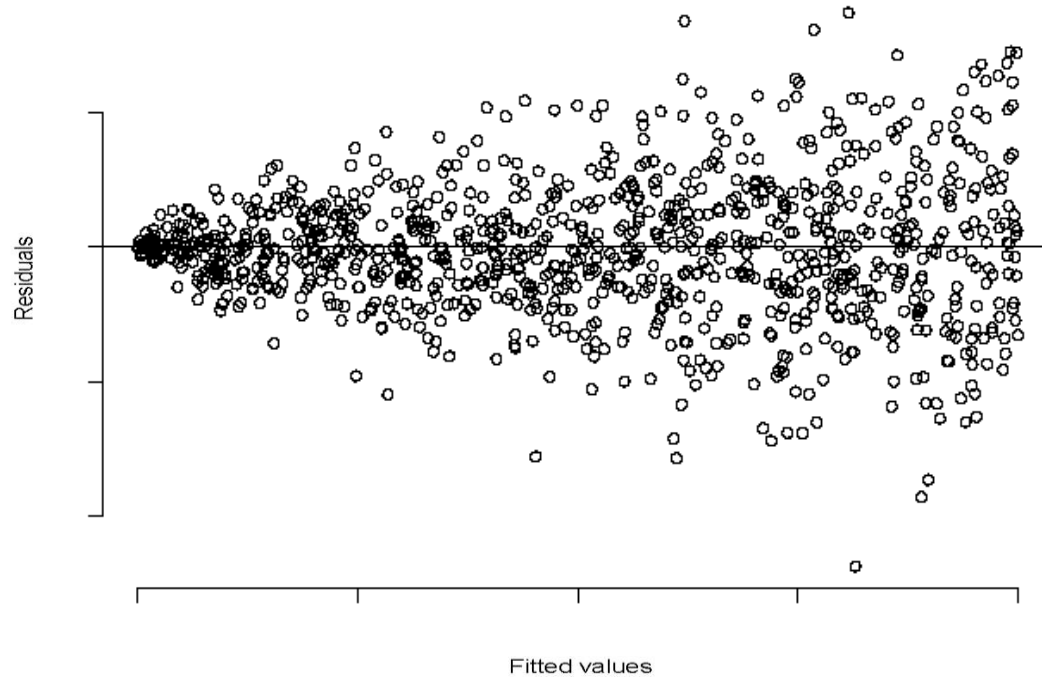Model checking should be performed to avoid erroneous extrapolation of data trends

## Residuals

▷ Residuals are the error between the observed and predicted data

▷ Look for trends or patterns in the residuals which indicate an assumption is not valid

# Residual Analysis

- ▹ Scatter plots of residuals against fitted values help identify:
    - ▸ Non-constant variance
    - ▸ Violation of the linearity assumption
    - ▸ Potential outliers

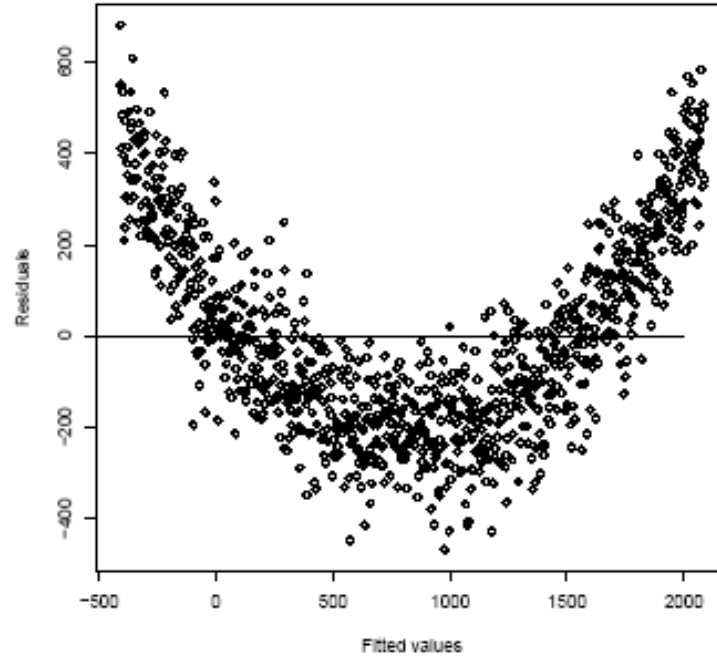- ▹ If these assumptions are valid then you will see no trend

# Residual Scatter Plot: Satisfactory

# Residual Scatter Plot: Non-constant Variance
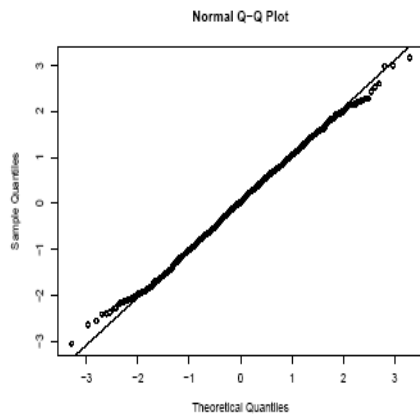
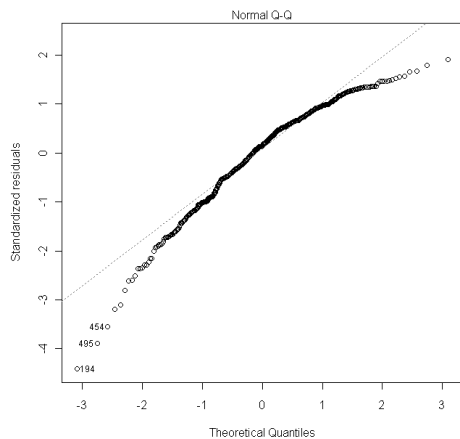# Residual Scatter Plot: Non-Linear Relationship

# Residual Analysis

▹ Standardise residuals by dividing by their standard deviation

▹ They should now be Normal with mean = 0 and variance = 1

▹ Box plots – symmetric?

▹ Proportion of standardised residuals inside percentiles

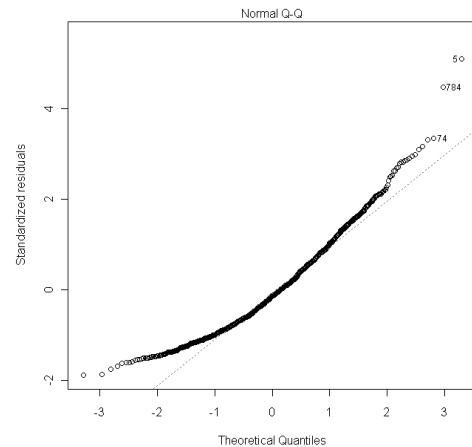▹ Q-Q plot – x = y graph should form

# Q–Q Plots and Skew

▷ Satisfactory

▷ Negative Skew

▷ Positive Skew

# Skew

- ▷ Median < Geometric Mean < Arithmetic Mean
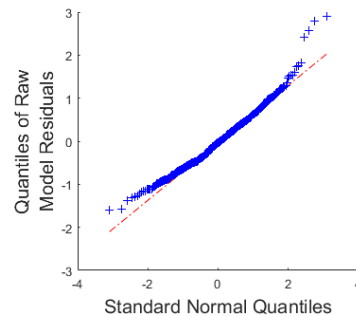- ▷ Positive skew => Standardise by the geometric mean

- ▷ Negative skew=> Standardise

# Identifying Outliers

▷ Any points which stand out as having larger residuals than other values should be checked

▷ Cook's Distance is given by most software and measures the influence of each individual point on the model

Residual Analysis

▷ $Efficiency = \alpha + \beta_1 Mass + \beta_2 FuelType$

▸ Possibly a non-linear relationship?



▸ Maybe very slightly positively skewed?

# Model Fitting

Model checking should be performed to avoid erroneous extrapolation of data trends

# Method

**Scatter Plots, Box Plots & Histograms**
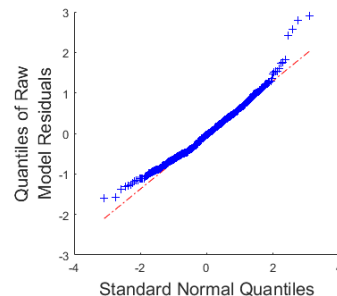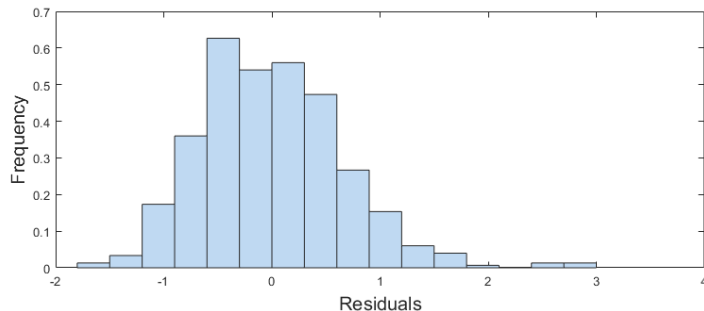Look for linearly related variables, skewness

**Interaction Terms**
These may improve the model of your data

**Adjust Model**
You may transform variables before fitting a model or include additional terms

**Fit Model**
Based on which variables you suspect are linearly related from scatter diagrams.

**Residual Analysis**
Look for non-constant variance, non-linearity and potential outliers

# MATLAB Script

67

Model Fitting

- Given data on car emissions in terms of vehicle mass, acceleration time, engine size, fuel type and colour.
- Fit models seeking to optimise for:
  - Error
  - AIC – relative loss of information

Search for Best Models

- ▸ Open file: Regression_Analysis_Car_Emissions
- ▸ Use the scatter and box plots to pick a starting model

# How to Use the Script

▷ Standardise by the arithmetic or geometric mean

```
94
95    %%%%%%%%%%%%%%%%%%%%%%
96    % 'G' for geometric standardisation or 'A' for arithmetic standardisation
97 -  standard = Standardise(cont,FieldNames,ArMean,ArStd,GeoMean,GeoStd,'A');
98    %%%%%%%%%%%%%%%%%%%%%%
99
```

▷ Edit line 115 with your trial model

```
113    %%%%%%%%%%%%%%%%%%%%%%%%%%%%
114    % Insert Model Here
115 -  model = ['Efficiency ~ Mass + FuelType'];
116
117    % Other Examples for models
118    % model = ['Efficiency ~ Mass + EngineSize^2'];
119    % AccelTime
120    % EngineSize
121    % Mass
122    % FuelType
123    % Mass:FuelType % Example of an interaction term
124    %%%%%%%%%%%%%%%%%%%%%%%%%%%%
```

# Search for Best Models

▹ Keep note of the MSE, Rsquared and AIC

▹ Seek to increase Rsquared

▹ Seek to reduce MSE and AIC

# Things to Try

▹ Increasing the number of terms

▹ Including categorical variables

▹ Using interaction terms (FuelType:AccelerationTime)

▹ Standardising Method

▹ Raising terms to a power

# How did you do?

Prizes for the 'best' model
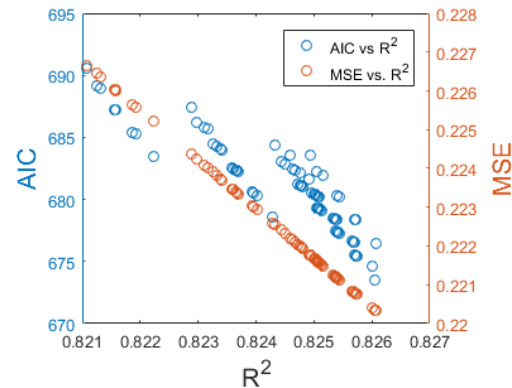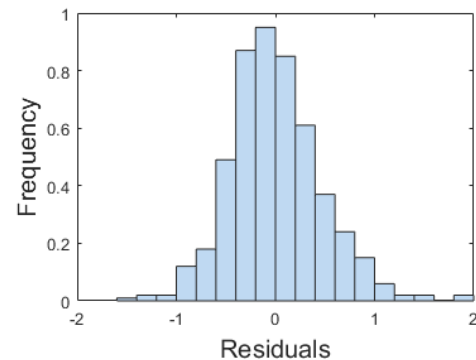
# An Optimum Model?

# Optimum Model

▷ Running the script:

```
125
126     % Make Model:
127     % [Rmax,mdl,Criteria] = MakeModel(tbl,model);
128 —   [CountOpt,Rmax,mdl,Criteria] = MakeOptimumModel(tbl,3)
129
```

▷ Normally distributed residuals
▷ AIC and MSE clearly decrease with Rsquared

# Model Simplicity

▷ Model 1

- ▸ Rsquared = 0.8261
- ▸ MSE = 0.220
- ▸ AIC = 673

$$Efficiency =$$
$$\alpha + \beta_1 Mass$$
$$+ \beta_2 EngineSize$$
$$+ \beta_3 AccelerationTime$$
$$+ \beta_4 MassFuelType^2$$

▷ Model 2

- ▸ Rsquared = 0.8261
- ▸ MSE = 0.220
- ▸ AIC = 676

$$Efficiency =$$
$$\alpha + \beta_1 Mass^3$$
$$+ \beta_2 EngineSiz$$
$$+ \beta_3 AccelerationTime^2$$
$$+ \beta_4 MassFuelType^2$$
$$+ \beta_4 EngineSize\ FuelType$$

# Further Reading

What are next steps

## Next Steps

- ▸ K-means cluster classification
- ▸ Bootstrapping (Confidence Intervals)
- ▸ Introductory courses to Machine Learning (Stanford, coursera)

# How did we do?

▹ Please let us know what you thought of this course and how we can improve it

▹ https://forms.gle/HAdcNGkFK5fkCimk6