# Assumptions of Linear Models

# Linear Model Assumptions
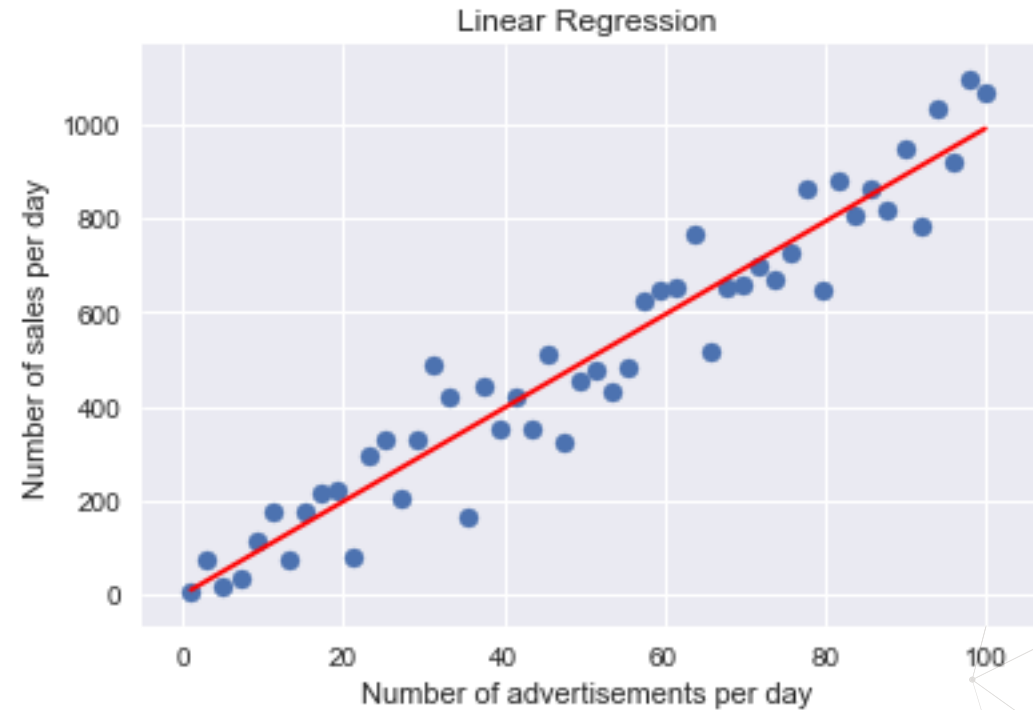
Linear models make the following assumptions about the independent variables (Xs)

- Linear relationship between the variables and the target

- Multivariate normality

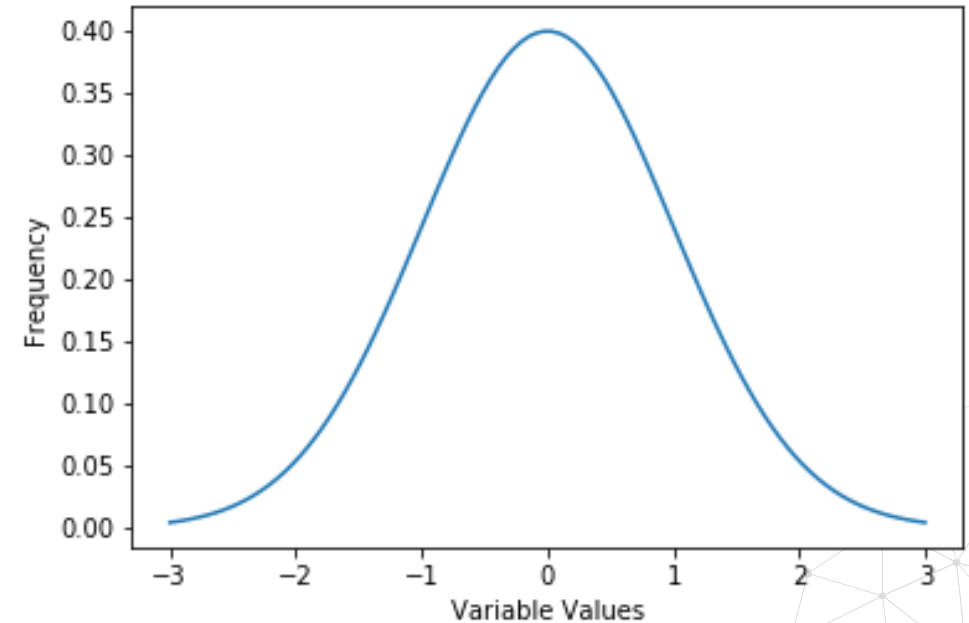- No or little co-linearity

- Homoscedasticity

# Linear Relationship

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_n X_n$$

- Linear relationship can be assessed with scatter plots

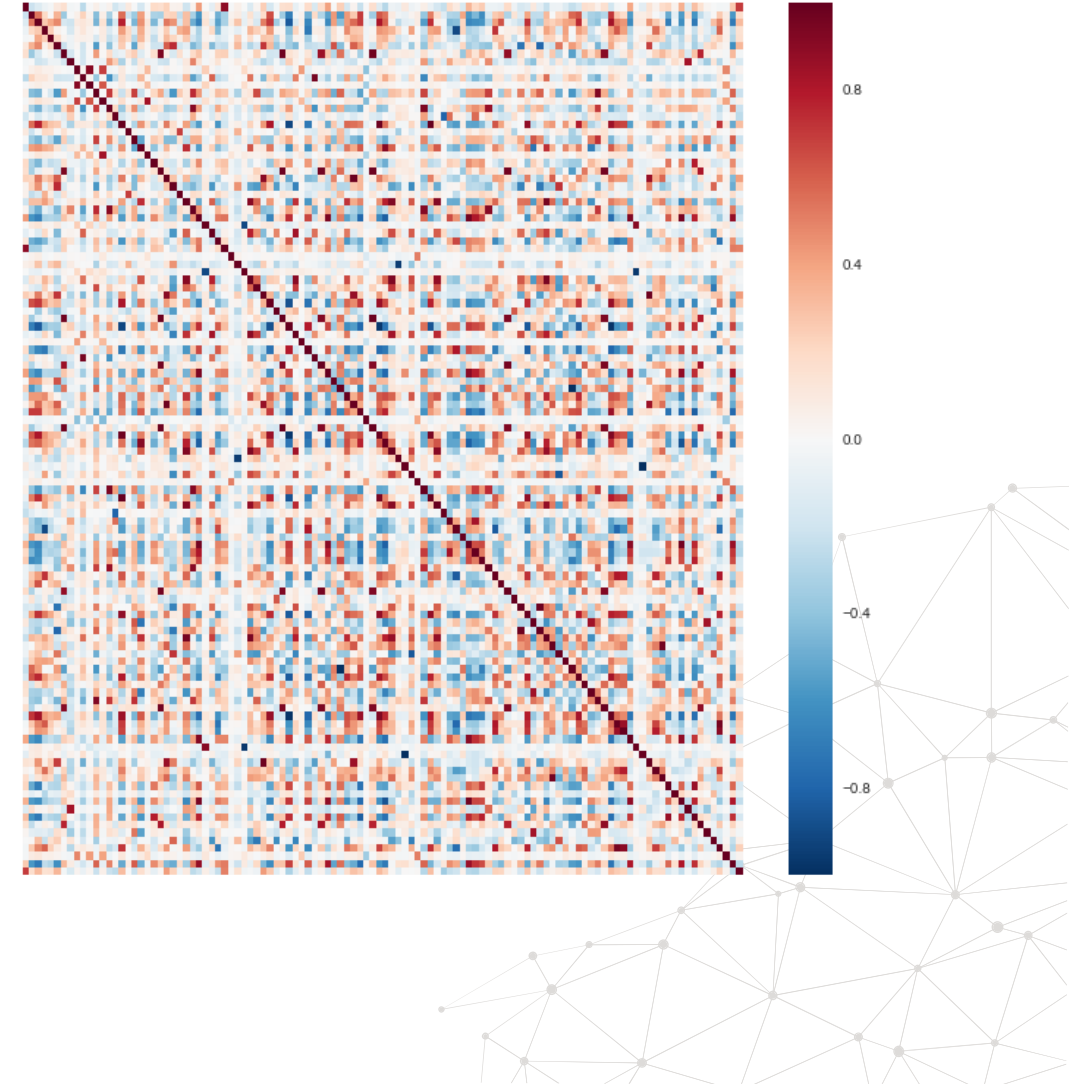- Sometimes non-linear transformations of the variables and the target improve the linear relationship

# Normality

- Variables follow a Gaussian Distribution

- Normality can be assessed with histograms and Q-Q plots

- Normality can be statistically tested, for example with the Kolmogorov-Smirnov test.

- When the variable is not normally distributed a non-linear transformation (e.g., logarithm-transformation) may fix this issue.

# No co-linearity

- Multicollinearity occurs when the independent variables are correlated with each other

- Multicollinearity can be assessed with a correlation matrix or the variance inflation factor (VIF)
  - Outside of the scope of this course
  - Check the course Feature Selection for Machine Learning

# Homoscedasticity

- The independent variables have the same finite variance.

- Also known as homogeneity of variance.

- There are tests and plots to determine homoscedasticity.
  - Residuals plot
  - Levene's test
  - Barlett's test
  - Goldfeld-Quandt Test

- Non-linear transformations and feature scaling can help improve homogeneity of variance
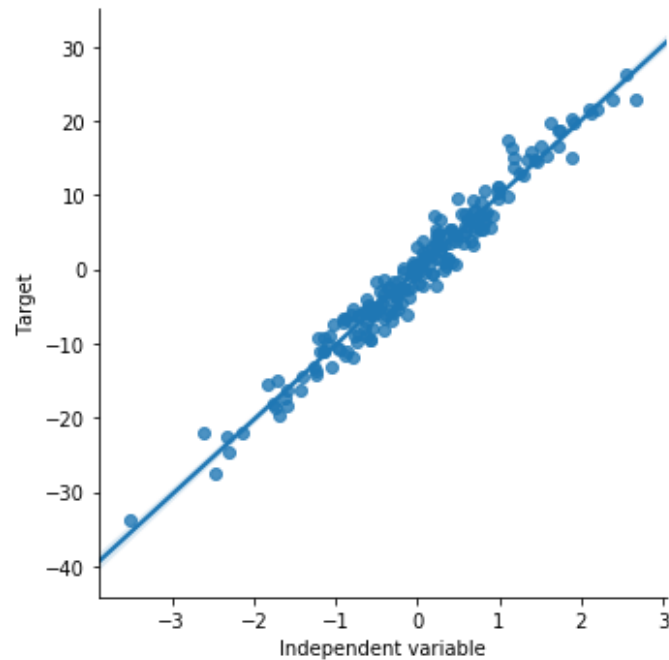
# **Evaluate model assumptions**

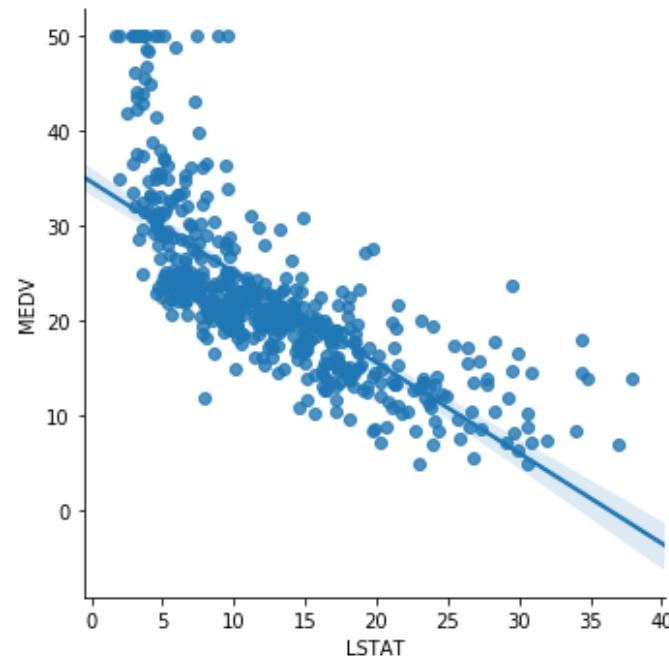Compare model assumptions in simulated and real data
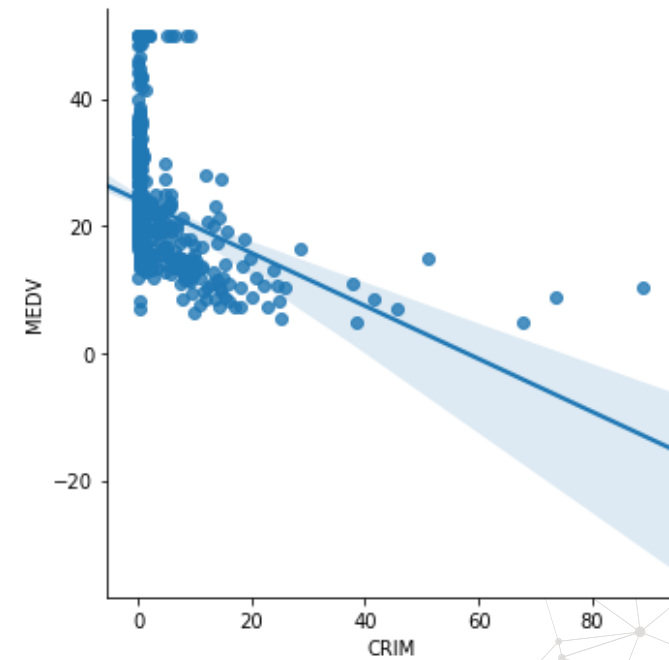
# Linear Relationship – Scatter plots



Expected – Simulated data
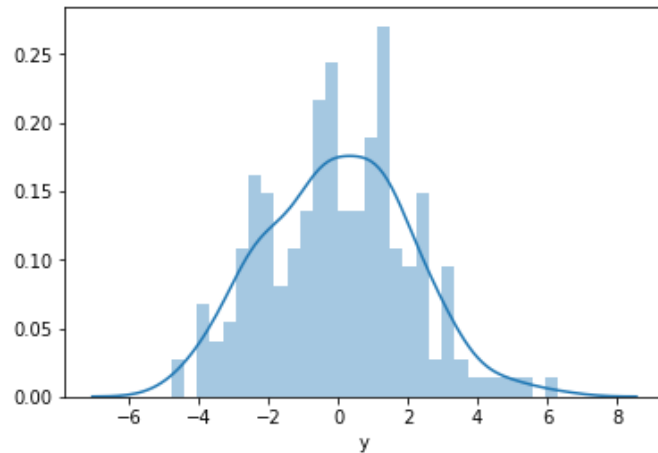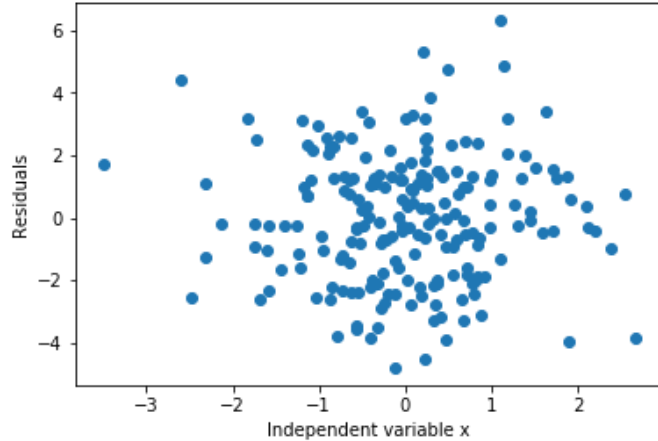
Somewhat linear relationship
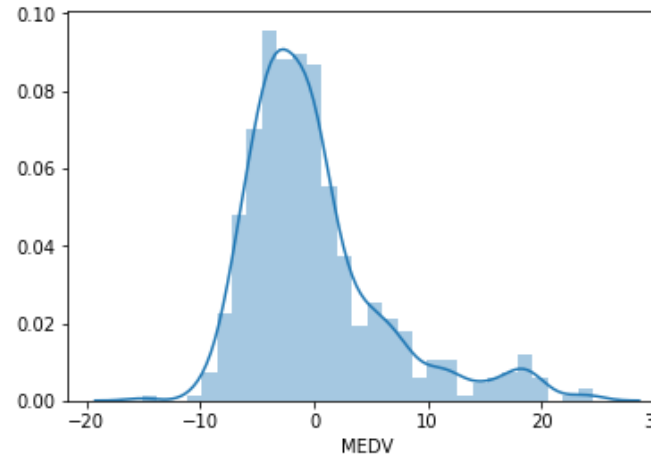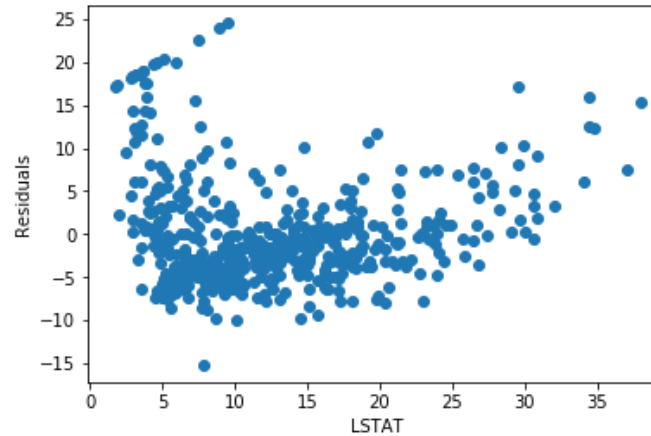
Non-linear relationship

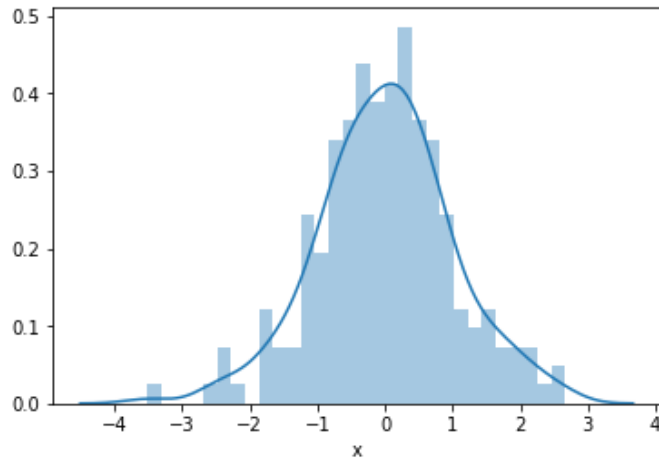# Linear Relationship – Residual plots



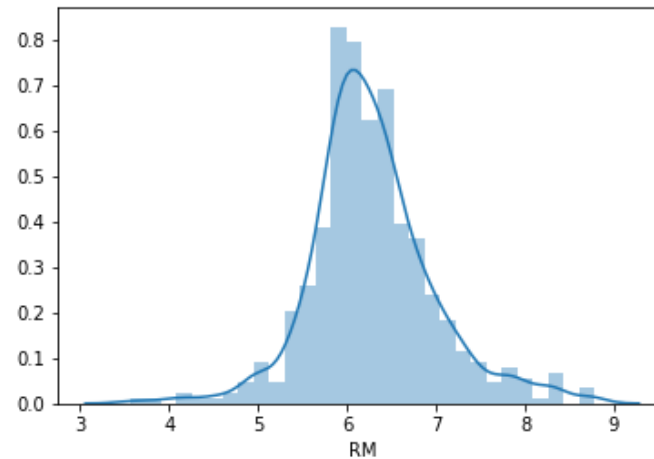Expected – Simulated data



Somewhat linear relationship

- If relationship between X and y is linear, residuals should be normally distributed and centred around 0

- Residuals are the difference between the predictions and the real value y.

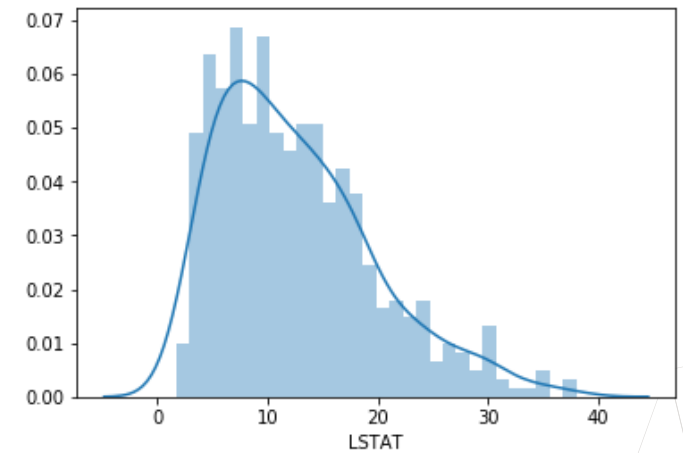# Normality – Histograms

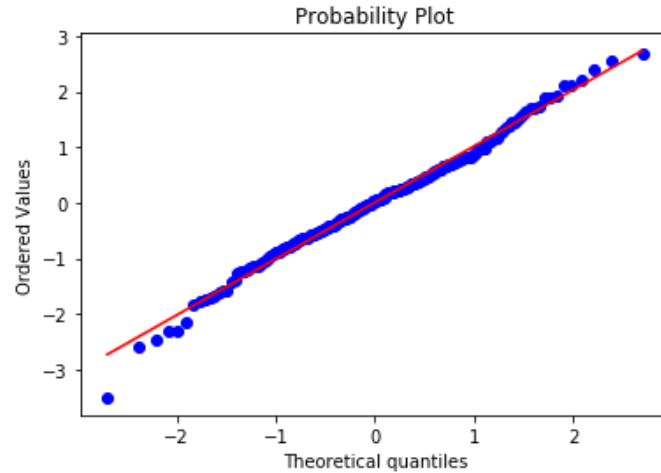Expected – Simulated data    Somewhat linear relationship (RM)    Non-linear relationship (LSTAT)
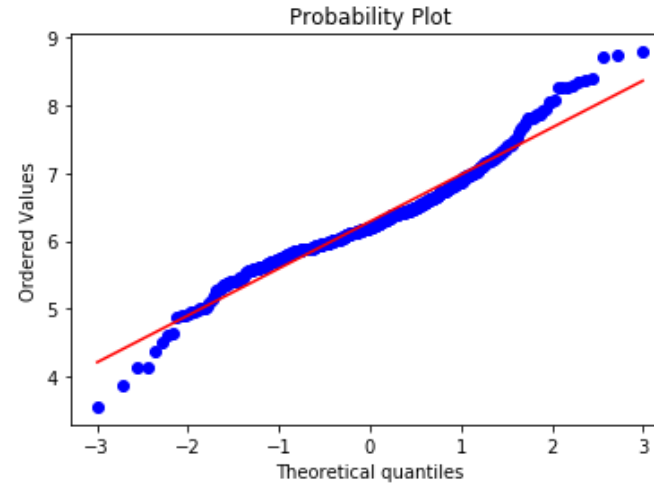


Gaussian distributions adopt a bell shape
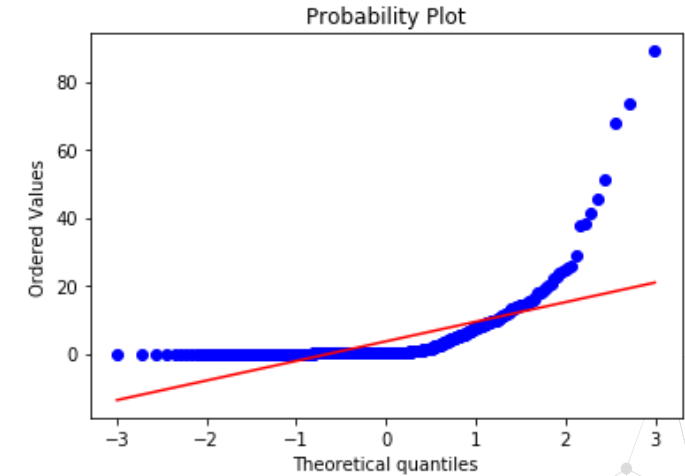
# Normality – Q-Q plots

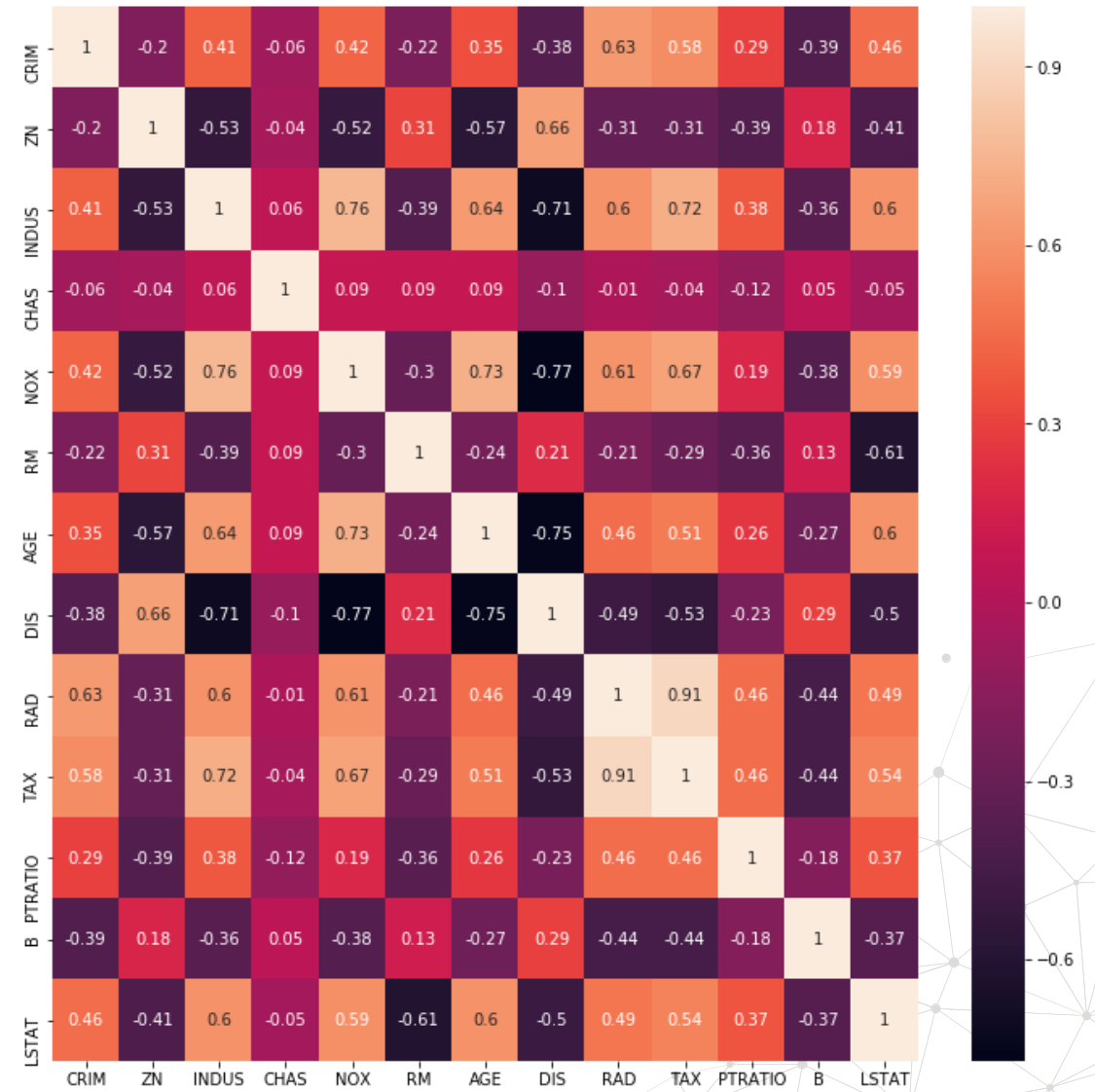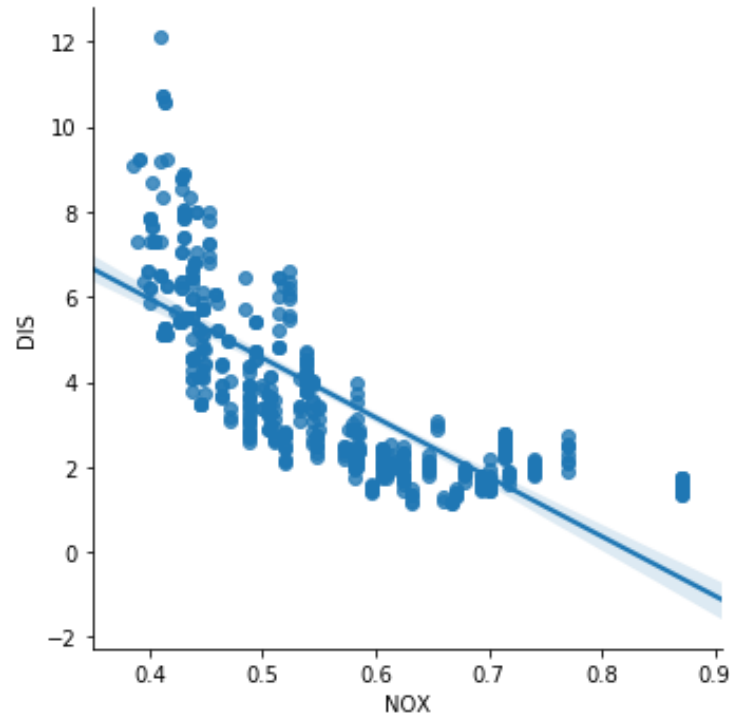Expected – Simulated data          Somewhat linear relationship          Non-linear relationship
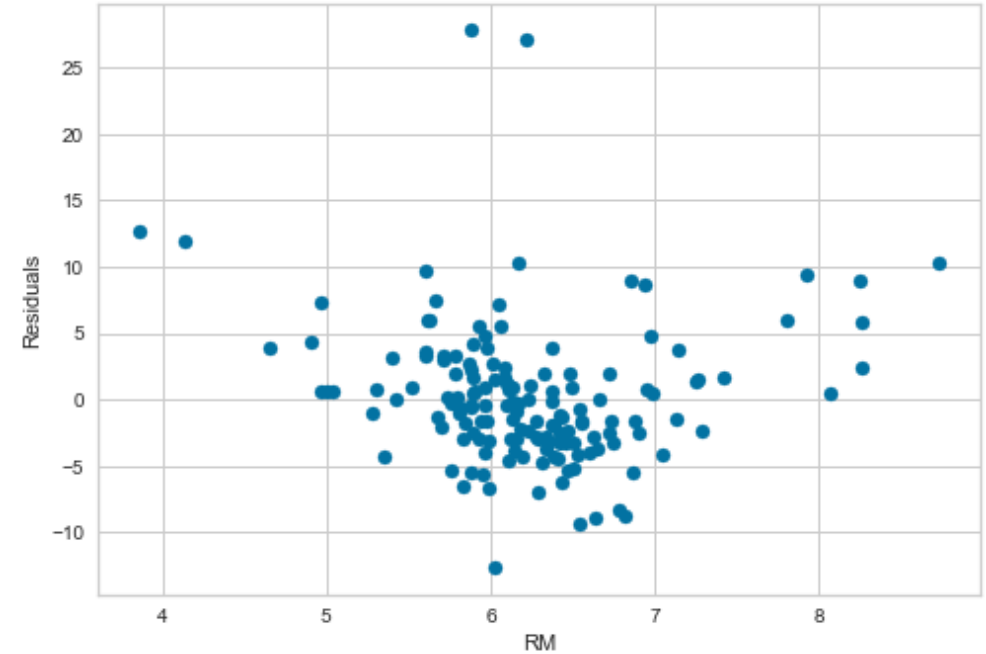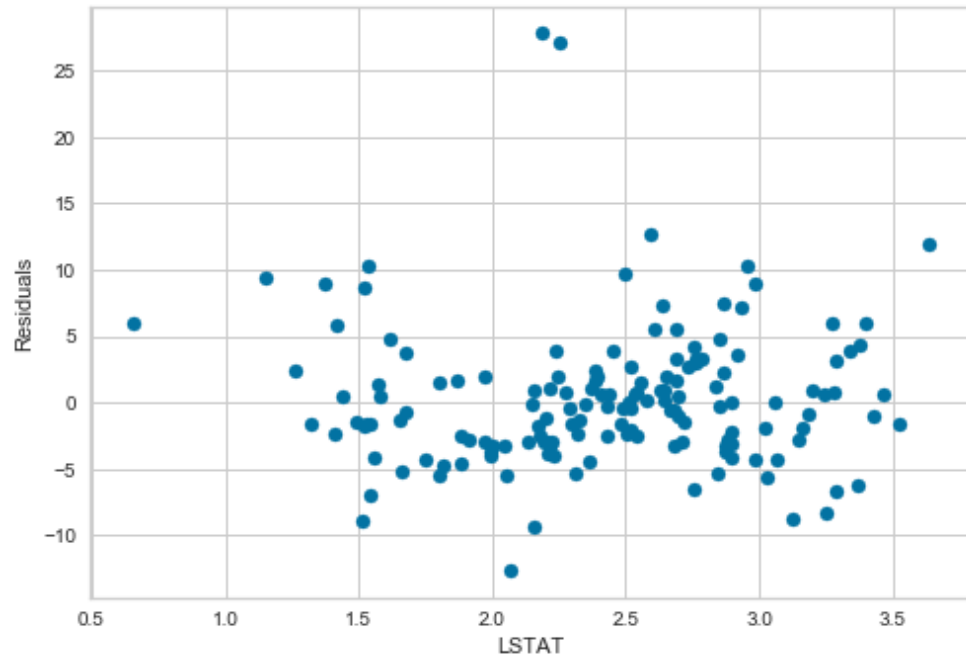


- Q-Q plots plot the variable quantiles in the y-axis and the expected quantiles of the normal distribution on the x-axis.

- If variable is normally distributed, the blue dots should fall on a 45 degree line

# Multi Co-linearity

Evaluated by correlation

# Homoscedasticity



**Homoscedasticity**: the error term (that is, the "noise" in the relationship between the independent variables X and the dependent variable Y) is the same across all the independent variables.

To identify homoscedasticity we need to plot the residuals vs each of the independent variables.

The distributions should be similar.

# Accompanying Jupyter Notebook



- Read the accompanying Jupyter Notebook

- Full demonstration of the linear assumptions and the influence of non-linear transformations

# THANK YOU

www.trainindata.com