

Linear Regression



COGNIBOT
AI meets Industry

- Once we've acquired data with multiple variables, one very important question is how the variables are related.
- Regression is a set of techniques for estimating relationships

Fitting curves to bivariate data

- Bivariate data $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Model:
 - $Y_i = f(x_i) + E_i$ (where $f(x)$ is some function, E_i random error)

- Total squared error:

$$\sum_{i=1}^n E_i^2 = \sum_{i=1}^n (y_i - f(x_i))^2$$



- Model allows us to predict the value of y for any given value of x
 - X is called the independent variable
 - Y is the dependent variable

Examples of $f(x)$

- Lines: $y = ax + b + E$
- Polynomials: $y = ax^2 + bx + c + E$
- Other: $y = a\sin(x) + b + E$



Simple linear regression

Find the best-fitting line

- Fit a line to the data
 - $y_i = ax_i + b + E$ (E is of normal distribution)

- Total squared error:

$$\sum_{i=1}^n E_i^2 = \sum_{i=1}^n (y_i - ax_i - b)^2$$

- Goal: Find the values of a and b that give the best fit
- Best fit: minimizes the total squared error



Linear Regression

Find the best-fitting polynomial

- Fit a parabola to the data
 - $y_i = ax_i^2 + bx_i + c + E$ (E is of normal distribution)
- Total squared error:

$$\sum_{i=1}^n E_i^2 = \sum_{i=1}^n (y_i - ax_i^2 - bx_i - c)^2.$$

- Goal: Find the values of a, b, c that give the ‘best fitting parabola’.
- Best fit: minimizes the total squared error
- Can also fit higher order polynomials



- Parameters are linear:
 - $y = ax + b$
 - $y = ax^2 + bx + c$
- It is not because the curve being fit has to be a straight line—although this is the simplest and most common case.



Examples of $f(x)$

- Lines: $y = ax + b + E$
- Polynomials: $y = ax^2 + bx + c + E$
- Other: $y = a\sin(x) + b + E$



Measuring the fit

- $TSS = \sum (y_i - \bar{y})^2$ = total sum of squares = total variation.
- $RSS = \sum (y_i - \hat{y}_i)^2$ = residual sum of squares.
- RSS/TSS = unexplained fraction of the total error.
- $R^2 = 1 - RSS/TSS$ is measure of goodness-of-fit
- R^2 is the fraction of the variance of y explained by the model.

