Highlights for regression problems

# Data

| | TV | Radio | Newspaper | Sales |
|---|---|---|---|---|
| **1** | 230.1 | 37.8 | 69.2 | 22.1 |
| **2** | 44.5 | 39.3 | 45.1 | 10.4 |
| **3** | 17.2 | 45.9 | 69.3 | 9.3 |
| **4** | 151.5 | 41.3 | 58.5 | 18.5 |
| **5** | 180.8 | 10.8 | 58.4 | 12.9 |

# Data Exploration - outlier

# Feature Engineering

Log / Squared scaling to get linear connection

merge features by + / * -

# Correlation: BP, Age, Weight, BSA, Dur, Pulse, Stress

|        | BP    | Age   | Weight | BSA   | Dur   | Pulse |
|--------|-------|-------|--------|-------|-------|-------|
| Age    | 0.659 |       |        |       |       |       |
| Weight | 0.950 | 0.407 |        |       |       |       |
| BSA    | 0.866 | 0.378 | 0.875  |       |       |       |
| Dur    | 0.293 | 0.344 | 0.201  | 0.131 |       |       |
| Pulse  | 0.721 | 0.619 | 0.659  | 0.465 | 0.402 |       |
| Stress | 0.164 | 0.368 | 0.034  | 0.018 | 0.312 | 0.506 |

# Detecting collinearity

Correlation matrix

Transforming features to fit non-linear relationships

Simple linear regression can easily be extended to include multiple features. This is called **multiple linear regression**:

$$y = \beta_0 + \beta_1 x_1 + \ldots + \beta_n x_n$$

Each $x$ represents a different feature, and each feature has its own coefficient. In this case:

$$y = \beta_0 + \beta_1 \times TV + \beta_2 \times Radio + \beta_3 \times Newspaper$$

# Diagnosing model fit

# Model Evaluation Metrics for Regression

For classification problems, we have only used classification accuracy as our evaluation metric. What metrics can we used for regression problems?

**Mean Absolute Error** (MAE) is the mean of the absolute value of the errors:

$$\frac{1}{n}\sum_{i=1}^{n}|y_i - \hat{y}_i|$$

**Mean Squared Error** (MSE) is the mean of the squared errors:

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

**Root Mean Squared Error** (RMSE) is the square root of the mean of the squared errors:

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

Let's calculate these by hand, to get an intuitive sense for the results:

# sources

https://www.ritchieng.com/machine-learning-evaluate-linear-regression-model/