# Formula Sheet

## Data Mining

1. Equal-width binning

$$\texttt{width } w = \frac{\max - \min}{n}$$

2. Normalization Formulas

$$x_i' = \frac{(x_i - \min)}{\max - \min}(\max{}_{new} - \min{}_{new}) + \min{}_{new}$$

$$x_i' = \frac{x_i}{10^j}$$

3. Linear Regression

- Residuel:

$$e_i = |y_i - h_\theta(y_i)|$$

- Cost function:

$$J(\theta_0, \theta_1) = \frac{1}{2m}\Sigma_{i=1}^m (h_\theta(x_i) - y_i)^2$$

- Linear Regression on One-Dimensional Data:

$$\theta_1 = \frac{\Sigma_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\Sigma_{i=1}^m (x_i - \bar{x})^2}$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

4. Perceptron Algorithm

$$f(x) = \Sigma_{j=0}^n w_j x_j > 0, y = 1$$

$$f(x) = \Sigma_{j=0}^n w_j x_j < 0, y = -1$$

5. KNN

- Distance:

$$d(x, x') = \sqrt{\Sigma_{i=1}^n (x_i - x_i')^2}$$

6. KMean

- New Mean:

$$c_k' = \frac{1}{|C_k|}\Sigma_{x_i \in C_k} x_i$$

7. Hierarchical Clustering

- MAX, MIN, AVERAGE, CENTROID

# Statistical

Basic Statistics

- Mean

$$\bar{x} = \frac{1}{n}\Sigma_{i=1}^{n} x_i$$

- Median
- Range

$$R = \max - \min$$

- Population variance

$$\sigma^2 = \frac{1}{n}\Sigma_{i=1}^{n}(x_i - \bar{x})^2$$

- Sample variance

$$s^2 = \frac{1}{n-1}\Sigma_{i=1}^{n}(x_i - \bar{x})^2$$

- Population standard deviation

$$\sigma = \sqrt{\frac{1}{n}\Sigma_{i=1}^{n}(x_i - \bar{x})^2}$$

- Sample standard deviation

$$s = \sqrt{\frac{1}{n-1}\Sigma_{i=1}^{n}(x_i - \bar{x})^2}$$

- IQR

$$\texttt{IQR} = \texttt{Q3} - \texttt{Q1}$$

- Outliner

$$x_i < Q1 - 1.5 \cdot \texttt{IQR}$$

$$x_i > Q3 + 1.5 \cdot \texttt{IQR}$$

Inferential Statistics

- Standard Error

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

- Confidence Interval (95%)

$$\left[\bar{x} - 2\sigma_{\bar{x}}, \bar{x} + 2\sigma_{\bar{x}}\right]$$

$$\left[\bar{x} - 2\frac{\sigma}{\sqrt{n}}, \bar{x} + 2\frac{\sigma}{\sqrt{n}}\right]$$

- Single T-Test, Paired T-test

$$t = \frac{\bar{x} - \mu}{s_{\bar{x}}} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

- Indpendent Samples T-Test

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{S_{\text{Pooled}}^2\left(\frac{1}{m} + \frac{1}{n}\right)}}$$

$$S_{\text{Pooled}}^2 = \frac{\text{df}_{\text{x}}}{\text{df}_{\text{total}}}s_x^2 + \frac{\text{df}_{\text{y}}}{\text{df}_{\text{total}}}s_y^2$$
$$= \frac{(m-1)s_x^2 + (n-1)s_y^2}{m+n-2}$$

- One-Way ANOVA

$$\text{MST} = \frac{\text{SST}}{p-1} = \frac{\Sigma_{i=1}^p n_i(\bar{x}_i - \bar{x})^2}{p-1}$$

$$\text{MSE} = \frac{\text{SSE}}{n-p} = \frac{\Sigma_{i=1}^p \Sigma_{j=1}^{n_i}(x_{ij} - \bar{x}_i)^2}{n-p}$$
$$= \frac{1}{n-p}\left((Y_{11} - \bar{Y}_1)^2 + (Y_{21} - \bar{Y}_1)^2 + \cdots + (Y_{n_1 1} - \bar{Y}_1)^2 + \right.$$
$$(Y_{12} - \bar{Y}_2)^2 + (Y_{22} - \bar{Y}_2)^2 + \cdots + (Y_{n_2 2} - \bar{Y}_2)^2 +$$
$$\cdots +$$
$$\left. (Y_{1p} - \bar{Y}_p)^2 + (Y_{2p} - \bar{Y}_p)^2 + \cdots + (Y_{n_p p} - \bar{Y}_p)^2\right)$$

$$F = \frac{\text{MST}}{\text{MSE}}$$

- Post-Hoc Test

$$\text{Tukey's HSD} = \frac{\bar{Y}_i - \bar{Y}_j}{\sqrt{\frac{\text{MSE}}{n}}}$$

# Security and Privacy

- Prevalence by UCT

$$\textbf{Prevalence} = \textbf{Average of GroupA} - \textbf{Average of GroupB}$$

- Prevalence by NST

$$\textbf{Prevalence} = \frac{\Sigma \textbf{s}_\textbf{i}}{\Sigma \textbf{k}_\textbf{i}}$$

, where $s_i$ is the number of people he/she knows are engagingin sensitive activity and $k_i$ is total number of people he/she knows.

- Prevalence by NRRT

$$\textbf{Prevalence } s = (P - ct)/(1 - c)$$

, where $P$ is the proportion of people answer "Yes", $c$ is probability of answer "Yes" in first question (Coffee), $t$ is the number of people who answer "Yes" in alternative non-sensitive question (Taxi).

- Prevalence by RRT

$$\textbf{Prevalence } s = (1 - \theta - P)/(1 - 2\theta)$$

, where $P$ is the proportion of people answer "Yes", $\theta$ is the ratio of the positive question.