

Formula that you should remember for the exam

For calculation that involves formula outside this sheet will provide you the formula in the exam.

This is provided for you to study only. We **will NOT** provide you this sheet in the exam.

Data Mining

1. Equal-width binning

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

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

- Residual:

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

- Cost function:

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

- Linear Regression on One-Dimensional Data:

$$\theta_1 = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2}$$
$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

4. Perceptron Algorithm

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

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

5. KNN

- Distance:

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

6. KMean

- New Mean:

$$c'_k = \frac{1}{|C_k|} \sum_{x_i \in C_k} x_i$$

7. Hierarchical Clustering

- MAX, MIN, AVERAGE, CENTROID

Statistical

Basic Statistics

- Mean

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

- Median
- Range

$$R = \max - \min$$

- Population variance

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

- Sample variance

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

- Population standard deviation

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

- Sample standard deviation

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

- IQR

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

- Outliner

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

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

Inferential Statistics

- Standard Error

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

- Confidence Interval (95%)

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

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

- Single T-Test, Paired T-test

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

- Independent Samples T-Test

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

$$\begin{aligned} S_{\text{Pooled}}^2 &= \frac{\text{df}_x}{\text{df}_{\text{total}}} s_x^2 + \frac{\text{df}_y}{\text{df}_{\text{total}}} s_y^2 \\ &= \frac{(m-1)s_x^2 + (n-1)s_y^2}{m+n-2} \end{aligned}$$

- One-Way ANOVA

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

$$\begin{aligned} \text{MSE} &= \frac{\text{SSE}}{n-p} = \frac{\sum_{i=1}^p \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2}{n-p} \\ &= \frac{1}{n-p} ((Y_{11} - \bar{Y}_1)^2 + (Y_{21} - \bar{Y}_1)^2 + \dots + (Y_{n_1 1} - \bar{Y}_1)^2 + \\ &\quad (Y_{12} - \bar{Y}_2)^2 + (Y_{22} - \bar{Y}_2)^2 + \dots + (Y_{n_2 2} - \bar{Y}_2)^2 + \\ &\quad \dots + \\ &\quad (Y_{1p} - \bar{Y}_p)^2 + (Y_{2p} - \bar{Y}_p)^2 + \dots + (Y_{n_p p} - \bar{Y}_p)^2) \\ F &= \frac{\text{MST}}{\text{MSE}} \end{aligned}$$

- 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

$$\text{Prevalence} = \text{Average of Group A} - \text{Average of Group B}$$

- Prevalence by NST

$$\text{Prevalence} = \frac{\sum s_i}{\sum k_i}$$

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

- Prevalence by NRRT

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

, 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

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

, where P is the proportion of people answer "Yes", θ is the ratio of the positive question.