

Formula Sheet

Master Data Science

Tim Downie

Machine Learning 1

Summer Semester 2019

Last edit: June 17, 2019

K-Means Clustering Minimise $W(C_1, \dots, C_k) = \sum_{k=1}^K W(C_k) = \sum_{k=1}^K \sum_{i \in C_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2$,

$\boldsymbol{\mu}_k = (\mu_{k,1}, \dots, \mu_{k,p})$, with $\mu_{k,j} = \frac{1}{|C_k|} \sum_{i \in C_k} x_{i,j}$ and

$\|\mathbf{x} - \mathbf{y}\|^2 = \sum_{j=1}^p (x_j - y_j)^2$ is the squared *euclidean distance* between the vectors \mathbf{x} and \mathbf{y}

For K -Medians use Manhattan distance $\|\mathbf{x} - \mathbf{y}\| = \sum_{j=1}^p |x_j - y_j|$.

Mean Square Error (MSE)

$$\text{MSE}(f, \mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 = \frac{1}{n} \text{RSS} \text{ (Residual sum of squares)}$$

Regression

Linear Models A linear model has the form $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$.

In vector and matrix form: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$. \mathbf{X} is called the design matrix and $\boldsymbol{\beta}$ is the vector of parameters.

The linear model estimates are the least squares estimates $\hat{\boldsymbol{\beta}}$, which minimises the residual sum of squares.

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}.$$

$$\text{Fitted values: } \hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$$

$$\text{Residuals: } \hat{\boldsymbol{\epsilon}} = \mathbf{y} - \hat{\mathbf{y}}$$

An **Anova model** is a linear model where the explanatory variables are discrete or nominal. Example one factorial Anova model, an overall parameter β_0 is fitted and then one parameter is fitted for all but the first level.

$$y_i = \beta_0 + \beta_1 \mathcal{I}(x_i \in A_2) + \dots + \beta_{K-1} \mathcal{I}(x_i \in A_K) + \epsilon_i$$

$\mathcal{I}(x_i \in A_k)$ is equal to 1 if x_i is in the group A_k and zero if not. A_1 is called the base line group.

Ridge Regression minimises

$$\text{RSS} + \lambda \sum_{j=1}^p \hat{\beta}_j^2 = \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \hat{\beta}_j^2$$

The Lasso minimises

$$\text{RSS} + \lambda \sum_{j=1}^p |\hat{\beta}_j| = \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\hat{\beta}_j|$$

Classification

Binary classification Predict $Y=0$ if $P(Y=1|x_1, \dots, x_p) < \alpha$ else predict $Y=1$.

Logistic regression

Logit function: $\text{logit}(s) = \log\left(\frac{s}{1-s}\right)$

Logistic regression linear predictor $\text{logit}(f(x_1, \dots, x_p)) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

$$\Rightarrow f(x_1, \dots, x_p) = P(Y=1|x_1, \dots, x_p) = \frac{\exp\{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p\}}{1 + \exp\{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p\}}$$

The **sensitivity** of a classifier is the true positive rate. The **specificity** of a classifier is the true negative rate.

The conditional probability of A given that B occurs is $P(A|B) = \frac{P(A \cap B)}{P(B)}$

Bayes Thm 2nd version $P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\bar{A})P(\bar{A})}$

Bayes Classifier Choose the group y_j which maximises $P(Y = y_k|x_1, \dots, x_p)$, for all $k = 1, \dots, K$

Bayes Error Rate measures how often we can expect to make a false classification.

$$1 - E(\max_k P((Y = y_k|X)))$$

Linear discriminant analysis (LDA) The LDA classifier rule is to choose the value of k giving the largest $P(Y = y_k|x)$, the posterior probability our outcome variable is in group y_k given x

$$P(Y = y_k|x) = \frac{f_{X|Y=y_k}(x)P(Y = y_k)}{\sum_{j=1}^K f_{X|Y=y_j}(x)P(Y = y_j)}$$

$P(Y = y_k)$ is the prior probability that our outcome variable is in group y_k .

The density $f_{X|Y=y_k}(x)$ in each group is modelled using a multinomial distribution:

$$(X_1, X_2, \dots, X_p)|Y = y_k \sim N(\boldsymbol{\mu}_k, \Sigma^2)$$

$\boldsymbol{\mu}_k$ is a vector of length p . Σ^2 is the common variance-covariance matrix.

Quadratic discriminant analysis (QDA) As LDA but use a different covariance in each group.

$$(X_1, X_2, \dots, X_p)|Y = y_k \sim N(\boldsymbol{\mu}_k, \Sigma_k^2)$$

Σ_k^2 is the variance-covariance matrix in group k .

Tree Methods

Regression Trees Minimise $RSS = \sum_{j \in [T]} \sum_{i \in R_j} (y_i - \hat{y}_j)^2$, where the R_j s are terminal nodes in the

Tree T and $\hat{y}_j = \frac{1}{|R_j|} \sum_{i \in R_j} y_i$ mean of Y in the region R_j .

Pruning The full tree T_0 is pruned back to the subtree $T_\alpha \subset T_0$ by minimising

$$PLS(\alpha) = \sum_{j \in |T|} \sum_{i \in R_j} (y_i - \hat{y}_j)^2 + \alpha |T|,$$

with complexity parameter $\alpha \geq 0$, and each R_j the terminal nodes in T .

Classification trees

Classification error rate $E = \sum_{m=1}^M 1 - \max_k \{\hat{p}_{mk}\}$

Gini Index $G = \sum_{m=1}^M \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk})$

Entropy or Information index $D = - \sum_{m=1}^M \sum_{k=1}^K \hat{p}_{mk} \log(\hat{p}_{mk})$