

STAT 543/641 – WINTER 2019 – HOMEWORK #2

DUE MARCH ??, 2019

Let X_1, \dots, X_m and Y_1, \dots, Y_n be random samples from populations with means μ_X and μ_Y and variances σ_X^2 and σ_Y^2 , respectively. Let S_X^2 and S_Y^2 be the standard unbiased estimators of σ_X^2 and σ_Y^2 , respectively.

- (1) Suppose $\sigma_X^2 = \sigma_Y^2$ and write σ^2 for this common value.

$$S := \frac{(m-1)S_X^2 + (n-1)S_Y^2}{m+n-2}$$

is an unbiased estimator of σ^2 . It's called the *pooled variance estimator*.

- (2) Suppose, in addition to having common variance, that the X_i are independent of the Y_i . What is the distribution of S_X^2 ? What is its variance? Compare the mean squared errors S_X^2 , S_Y^2 , and S^2 .
- (3) Generalize these results from the case of $K = 2$ populations to that of an arbitrary K . Compare with equation (4.15) in [2].
- (4) **[Bonus]** Can you prove analogous results with covariance matrices in place of variances?

[1, Exercise 12.16] This exercise examines an extreme case in which the likelihood equations for logistic regression have no solution.

Consider the following 20-point data set:

$$(0, 1), (0, 1), (0, 1), (0, 1), (0, 1), (0, 1), (0, 1), (0, 1), (0, 1), (0, 1) \\ (1, 0), (1, 0), (1, 0), (1, 0), (1, 0), (1, 0), (1, 1), (1, 1), (1, 1), (1, 1)$$

- (1) Observe that, empirically, $\text{Prob}(Y = 1|X = 0) = 1$ and $\text{Prob}(Y = 1|X = 1) = 0.5$. Let $\sigma(t) = (1 + e^{-t})^{-1}$ be the sigmoid function. Are there a and b such that $\sigma(a + b \cdot 0) = 1$ and $\sigma(a + b \cdot 1) = 0.5$?
- (2) Let $\mathcal{L}(a, b)$ be the likelihood function associated to fitting a logistic regression model to this data set. Show that

$$L := \lim_{b \rightarrow \infty} \mathcal{L}(-b, b) = \sup_{(a, b) \in \mathbb{R}^2} \mathcal{L}(a, b) < \infty$$

and that $\mathcal{L}(a, b) \neq L$ for any $(a, b) \in \mathbb{R}^2$. What are

$$\lim_{b \rightarrow \infty} \sigma(-b + b \cdot 0) \quad \text{and} \quad \lim_{b \rightarrow \infty} \sigma(-b + b \cdot 1)?$$

Let (\mathbf{X}, Y) be jointly distributed, where \mathbf{X} is a p -dimensional random vector and Y takes values in $\{1, \dots, K\}$. Suppose that, for each k , $\mathbf{X}|Y = k$ has Gaussian distribution with mean $\boldsymbol{\mu}_k$ and variance Σ , with the latter independent of k .

Consider a data set $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})$, where $\mathbf{x}^{(i)} \in \mathbb{R}^{1 \times p}$ and $y^{(i)} \in \{1, \dots, K\}$. For $1 \leq k \leq K$, let

$$I_k = \{i : y^{(i)} = k\}, \quad n_k = |I_k|, \quad \hat{\pi}_k = \frac{n_k}{n}.$$

Define sample means $\boldsymbol{\mu}_k$ and a pooled sample covariance Σ by

$$\begin{aligned} \hat{\boldsymbol{\mu}}_k &= \hat{\boldsymbol{\mu}}_{k,\mathbf{x}} = \frac{1}{n_k} \sum_{i \in I_k} \mathbf{x}^{(i)} \in \mathbb{R}^{p \times 1}, \\ \hat{\Sigma} &= \hat{\Sigma}_{\mathbf{x}} = \frac{1}{n - K} \sum_{k=1}^K \sum_{i \in I_k} (\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}}_k)^T (\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}}_k) \in \mathbb{R}^{p \times p}. \end{aligned}$$

Define linear discriminant functions, $\delta_k = \delta_{k,\mathbf{x}}$, by

$$\delta_k(\mathbf{v}) = \delta_{k,\mathbf{x}}(\mathbf{v}) = \mathbf{v} \hat{\Sigma} \hat{\boldsymbol{\mu}}_k^T - \frac{1}{2} \hat{\boldsymbol{\mu}}_k^T \hat{\Sigma} \hat{\boldsymbol{\mu}}_k + \log \hat{\pi}_k, \quad \mathbf{v} \in \mathbb{R}^{p \times 1}.$$

Let $\mathbf{a} \in \mathbb{R}^{p \times 1}$ and let

$$\begin{aligned} \mathbf{w}^{(i)} &= \mathbf{x}^{(i)} - \mathbf{a}. \\ \hat{\boldsymbol{\mu}}_{k,\mathbf{w}} &= \hat{\boldsymbol{\mu}}_{k,\mathbf{x}} - \mathbf{a}, \quad \Sigma_{\mathbf{w}} = \Sigma_{\mathbf{x}} \\ \delta_{k_1,\mathbf{w}}(v - \mathbf{a}) - \delta_{k_2,\mathbf{w}}(v - \mathbf{a}) &= \delta_{k_1,\mathbf{x}}(v) - \delta_{k_2,\mathbf{x}}(v) \end{aligned}$$

Let $U \in \mathbb{R}^{p \times p}$ be an orthogonal matrix and let $\mathbf{w}^{(i)} = U \mathbf{x}^{(i)}$. Then

$$\begin{aligned} \delta_{k,U\mathbf{x}}(Uv) &= \delta_{k,\mathbf{x}}(v). \\ \sum_{k=1}^K \pi_k \mu_k &= \mu \end{aligned}$$

Logistic regression (with and without ridge regularization, with and without PCA), LDA, Gaussian naïve Bayes, for breast cancer data set. Plot in 2d with decision boundary. Optional: Lasso

Document classification with multinomial naïve Bayes

Ridge regression via constrained optimization.

REFERENCES

- [1] Casella, Bergger, *Statistical Inference (2nd ed.)*, Duxbury, 2002.
- [2]