

Classification *

Logistic regression

• logistic function $p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$ and we fit with maximum likelihood
 $\Rightarrow \logit \Rightarrow \log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x$ (this is linear in x).
 $\log odds$ $\hookrightarrow R(\beta_0, \beta_1) = \prod_{i: y_i=1} p(x_i) \prod_{i: y_i=0} (1-p(x_i))$

• like linear regression, we can measure the accuracy of the coefficient with the z -statistic of β_1 is $\beta_1 / SE(\beta_1)$ and larger value (absolute) against the null $H_0: \beta_1 = 0$ which implies $p(x) = \frac{e^{\beta_0}}{1 + e^{\beta_0}}$
 • and multiple logistic regression $\log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ and $p(x) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}$
 when there is correlation among the predictors, the result obtained by one predictor will be different from using multiple predictors (confounding).

Logistic for classification

• when we wish to classify more than two classes \Rightarrow multinomial logistic regression.
 For k classes, we select the k^{th} as baseline and

$$P_k(Y=k | X=x) = \frac{e^{\beta_{k0} + \beta_{k1} x_1 + \dots + \beta_{kp} x_p}}{1 + \sum_{l=1}^{k-1} e^{\beta_{l0} + \beta_{l1} x_1 + \dots + \beta_{lp} x_p}} \quad \text{for } k=1, \dots, k-1 \text{ and } P_k(Y=k | X=x) = \frac{1}{1 + \sum_{l=1}^{k-1} e^{\beta_{l0} + \dots + \beta_{lp} x_p}}$$

$$\text{and } \log\left(\frac{P_k(Y=k | X=x)}{P_k(Y=k | X=x)}\right) = \beta_{k0} + \beta_{k1} x_1 + \dots + \beta_{kp} x_p$$

\Rightarrow but for multinomial we can use an alternative \Rightarrow the softmax rather than selecting a baseline each class is treated symmetrically

$$P_k(Y=k | X=x) = \frac{e^{\beta_{k0} + \beta_{k1} x_1 + \dots + \beta_{kp} x_p}}{\sum_{l=1}^k e^{\beta_{l0} + \beta_{l1} x_1 + \dots + \beta_{lp} x_p}} \quad \text{and } \log\left(\frac{P_k(Y=k | X=x)}{P_k(Y=k' | X=x)}\right)$$

Generative model classification (Umbrella LDA, QDA, NB)

$$= (\beta_{k0} - \beta_{l0}) + (\beta_{k1} - \beta_{l1}) x_1 + \dots$$

in logistic we model $P_k(Y=k | X=x)$ using logistic function \rightarrow alternative \rightarrow model the distribution of the predictors x separately in each of the response classes. Then use Bayes' theorem to flip those around into estimates for $P_k(Y=k | X=x) \Rightarrow$ why do that \Rightarrow when substantial separation between classes the parameters of logistic regression are unstable.

• Let $\hat{\pi}_k$ the probability that a random observation comes from the k^{th} class

$$\bullet f_k(x) = P_k(X | Y=k) \quad \text{and then } P_k(Y=k | X=x) = \frac{\hat{\pi}_k f_k(x)}{\sum_{l=1}^k \hat{\pi}_l f_l(x)}$$

to approximate $f_k(x)$ we can use linear discriminant, quadratic discriminant and naive Bayes

Linear discriminant

P=1

$f_h(x)$ is assume normal, with shared variance among classes $\sigma_1^2 = \dots = \sigma_h^2$

$$f_h(x) = \frac{1}{\sqrt{2\pi}\sigma_h} \exp\left(-\frac{1}{2\sigma_h^2}(x-\mu_h)^2\right)$$

The Bayes classifier involves assigning an observation $x=x$ to the class for which P is largest. It amounts to assigning the observation, to the class for which

$$f_h(x) = x \cdot \frac{\mu_h}{\sigma^2} - \frac{\mu_h^2}{2\sigma^2} + \log(\pi_h)$$

is largest

$$\hat{\mu}_h = \frac{1}{n_h} \sum_{i: y_i=h} x_i \quad \hat{\sigma}_h^2 = \frac{1}{n_h-1} \sum_{i: y_i=h} (x_i - \hat{\mu}_h)^2$$

$$\hat{\pi}_h = n_h/n$$

• LDA outperform logistic regression when the normality assumption hold.

• to test the model we use a confusion matrix. The Bayes classifier (and by extension LDA) use a threshold of 50% for cutoff. we can change the threshold to get the best model. The AOC displays the two types of error for all possible threshold. To compare classifier use the AUC (best possible is 1, random classifier is 0.5).

quadratic discriminant

Unlike LDA, $x \sim N(\mu_h, \Sigma_h)$ where Σ_h is a covariance matrix for the h class.

$$f_h(x) = -\frac{1}{2} x^T \Sigma_h^{-1} x + x^T \Sigma_h^{-1} \mu_h - \frac{1}{2} \mu_h^T \Sigma_h^{-1} \mu_h - \frac{1}{2} \log |\Sigma_h| + \log \pi_h$$

LDA is less flexible classifier, so lower variance, but if the assumption of common variance is off, there will be a higher bias. \rightarrow few training \rightarrow LDA

\rightarrow many training or \neq matrix \rightarrow QDA

Naive Bayes

The naive Bayes takes a different tack for estimating $f_1(x), \dots, f_p(x)$. Instead of assuming a distribution, we make the assumption, that within the h^{th} class, the p predictors are independent

$f_h(x) = f_{h1}(x_1) \times f_{h2}(x_2) \times \dots \times f_{hp}(x_p)$ where f_{hj} is the density function of the j th predictor among observations in the h^{th} class. It introduces some bias but reduces variance.

$$P(Y=h|X=x) = \frac{\pi_h \times f_{h1}(x_1) \times \dots \times f_{hp}(x_p)}{\sum_{k=1}^K \pi_k \times f_{k1}(x_1) \times \dots \times f_{kp}(x_p)} \quad \text{for } h=1, \dots, K$$

to estimate f_{hj} \rightarrow x_j is quantitative $\rightarrow x_j|Y=h \sim N(\mu_{jh}, \sigma_{jh}^2)$
 \rightarrow as histogram, as kernel density estimator
 \rightarrow x_j is qualitative \rightarrow count per classes.

NB good when p is large or N is small and we need to reduce the variance.

P>1

multivariate gaussian, $N(\mu_h, \Sigma)$ with common cov matrix Σ

$$f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

$$f_h(x) = x^T \Sigma^{-1} \mu_h - \frac{1}{2} \mu_h^T \Sigma^{-1} \mu_h + \log \pi_h$$