

$Y \sim \text{Bernoulli}(\pi)$ - indiv. data

$$E\{Y\} = \pi, \quad \text{Var}\{Y\} = \pi(1-\pi)$$



$$\pi = E\{Y\} = g(\delta_0 + \delta_1 x) = \frac{e^{\delta_0 + \delta_1 x}}{1 + e^{\delta_0 + \delta_1 x}} = \frac{1}{1 + e^{-(\delta_0 + \delta_1 x)}} \quad \begin{matrix} \text{- CDF of} \\ \text{Logistic distr.} \end{matrix}$$

$\hookrightarrow \text{not a linear}$

$$\text{Likelihood} = \prod_{i=1}^n \pi_i^{y_i} (1-\pi_i)^{1-y_i} \quad \text{model}$$

Logistic regression for classification:

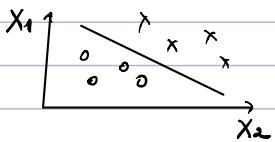
$$\hat{Y} = \begin{cases} 1, & \text{if } E\{Y|X\} \geq c \\ 0, & \text{if } E\{Y|X\} < c \end{cases}$$



→ confusion matrix
ROC

The classification cutoff separates points by a hyperplane:

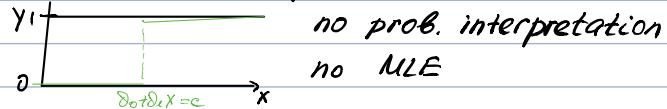
$$\frac{1}{1 + e^{-(\delta_0 + \delta_1 x_1 + \delta_2 x_2)}} = \frac{1}{2} \Rightarrow \delta_0 + \delta_1 x_1 + \delta_2 x_2 = 0 \Rightarrow x_1 = -\frac{\delta_0}{\delta_1} + \frac{\delta_2}{\delta_1} x_2$$



Comments:

- Linear separability: no unique MLE estimates $\text{Var}(\delta) \rightarrow \infty$
- Alternative ways to define Linear decision boundary:

$$f(x) = \begin{cases} 1, & \text{if } x^\top \delta \geq c \\ 0, & \text{if } x^\top \delta < c \end{cases} \quad \begin{matrix} g(x): \\ \text{step function} \end{matrix}$$



- Variable selection

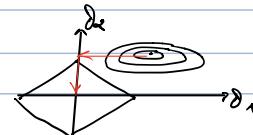
$$AIC = -2\log L + 2p \quad \begin{matrix} \text{no residuals} \Rightarrow \text{no } R^2, \text{ no SSE or MSE} \\ \text{search for subsets of predictors} \end{matrix}$$

$$BIC = -2\log L + \log N \cdot p \quad \begin{matrix} \text{that minimize AIC or BIC} \\ \text{train / dev / test sets or cross validation} \end{matrix}$$

- Regularization: $\min\{-\log L + \lambda \|\delta\|_2^2\}$

Components:

(1) objective function, (2) optimization algo



possible classes

Multi-class classification

		Y
		1 2 ... K
X		$n_{11} n_{12} \dots n_{1K}$
X ₁		$n_{11} n_{12} \dots n_{1K}$
X ₂		$n_{21} n_{22} \dots n_{2K}$
...		\dots
X _N		$n_{N1} n_{N2} \dots n_{NK}$
grouped data		
n _{ijk} ∈ {0, 1} if indiv. data		
multinomial		

$$\pi_k(x) = P\{Y=k | X=x\}, \quad \sum_{k=1}^K \pi_k(x) = 1$$

Multinomial Logistic regression:

$$\log \frac{\pi_k(x)}{\pi_1(x)} = \delta_k + \beta_k x, \quad k=2, \dots, K$$

arbitrary category as baseline,
to satisfy $\sum_{k=1}^K \pi_k = 1$

useful if interested
in param. interpretation

separate params
for each category

assumption of multi-class logistic regression

Solving for π_k (in the case of individual data, using $\sum_{k=1}^K \pi_k = 1$ and $y_{ik} = 1 - \sum_{k=2}^K y_{ik}$) (2)

$$\text{MLE: log Likelihood}_i = \log \left[\prod_{k=1}^K \pi_k(x)^{y_{ik}} \right] = \sum_{k=2}^K y_{ik} \log \pi_k(x) + \left(1 - \sum_{k=2}^K y_{ik} \right) \log \pi_1(x)$$

for one data point } for full likelihood, take product over all:

Plug into the

Likelihood the assumed π_i : $\pi_1(x) = \frac{1}{1 + \sum_{k=1}^K e^{d_k + \beta_k x}}$, $\pi_k = \frac{e^{d_k + \beta_k x}}{1 + \sum_{k=1}^K e^{d_k + \beta_k x}}$

Alternative parametrization: softmax regression

$$\pi_k(x) = \frac{e^{d_k + \beta_k x}}{\sum_{k=1}^K e^{d_k + \beta_k x}}$$

more convenient
if goal is prediction
because symmetric

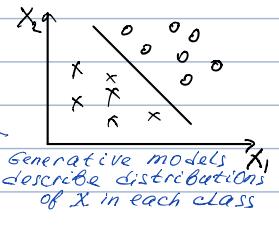
- only $K-1$ free parameter sets
- set $d_1 = \beta_1 = 0$ to have the previous case

Generative models

Bayes rule: $p(Y|X) = \frac{p(X|Y) p(Y)}{p(X)}$

Linear regression
Logistic regression

\Rightarrow discriminative
models



In classification, Y is discrete $P\{Y=k|X\} = \frac{f_k(x) \cdot \pi_k}{\sum_{l=1}^K f_l(x) \pi_l}$

most probable
class given the
data

distribution of X \rightarrow Likelihood

$f_k(x) \cdot \pi_k$ \leftarrow class prior
 $\sum_{l=1}^K f_l(x) \pi_l$ \leftarrow normalizing constant

Goal:

find class k maximizing
the posterior probability

MAP (maximum a posteriori learning)

More formally: minimize the 0-1 loss function

$$J(Y, \hat{Y}(x)) = \begin{cases} 0, & \hat{Y}(x) = Y \text{ — correct prediction} \\ 1, & \hat{Y}(x) \neq Y \text{ — wrong prediction} \end{cases}$$

In Linear reg.,
 $J(Y, \hat{Y}) = \sum_{i=1}^N (Y_i - \hat{Y}_i)^2$

Bayes risk: Loss function, averaged over $f_k(x)$ and K :

$$R(\hat{Y}(x)) = E \left\{ E \left\{ J(Y, \hat{Y}(x)) \right\} \right\}$$

over $f_k(x|Y=k)$ — all data of a class
over π_k — all classes

The MAP estimator minimizing Bayes risk under 0-1 Loss is

$$\hat{Y} = \arg \max_k P\{Y=k|X\} = \arg \max_k f_k(x) \pi_k - \text{ignore the denominator}$$

Specification of $f_k(x)$ define the classifier: variance-covariance

- $f_k(x)$ multivariate Gaussian, same Σ per class \rightarrow LDA
 - $f_k(x)$ multivariate Gaussian, different $\Sigma_k \rightarrow$ QDA
- $f(x) = \text{mixture of Gaussians}$
- $f_k(x)$ - component of the mixture
- $f_k(x) = \prod_{p=1}^P f_k(x_p)$ - naïve Bayes
- in each class, orthogonal dimensions
- typically (but not always) $f_k(x_p)$ is non-parametric (e.g., kernel density)

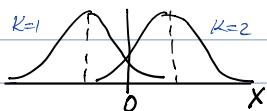
quadratic discriminant analysis

Linear discriminant analysis

Example: one predictor

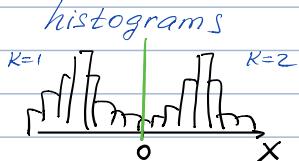
Population:

$N(\mu_k, \sigma^2)$, same



$$f_k(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu_k)^2}$$

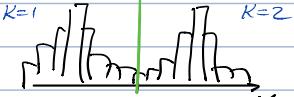
Data:



Q: what is the optimal decision boundary?

$$P\{Y=k|X\} = \frac{\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_k)^2\right) \cdot \pi_k}{\sum_{l=1}^2 \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_l)^2\right) \cdot \pi_l}$$

same variance



Predict class k :

$$\hat{y}(x) = \arg \max_k P\{Y=k|X\} = \arg \max_k \log P\{Y=k|X\}$$

$$= \arg \max_k \left\{ -\frac{1}{2\sigma^2}(x^2 - 2x\mu_k + \mu_k^2) + \log \pi_k \right\}$$

does not depend on k

$$= \arg \max_k \left\{ \underbrace{\frac{x\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log \pi_k}_{{\delta}_k(x)} \right\} = \arg \max_k {\delta}_k(x)$$

$\delta_k(x)$, linear in $x \rightarrow$ linear discriminant function

Assume $\pi_1 = \pi_2 = \frac{1}{2}$

$\hat{y}(x) = \arg \max_k \left\{ x\mu_k - \frac{1}{2}\mu_k^2 \right\}$. The boundary is $\{x\}$ where both classes have same posterior prob.

$$x\mu_1 - \frac{1}{2}\mu_1^2 = x\mu_2 - \frac{1}{2}\mu_2^2$$

$$x = \frac{\mu_1^2 - \mu_2^2}{2(\mu_1 - \mu_2)} = \frac{\mu_1 + \mu_2}{2}$$

In practice, estimate μ_k , σ and π_k from the data

Multiple classes and predictors:

Population:



- multivariate Gaussian $N(\mu_k, \Sigma)$
in this case: 2×2 2×2

$$\hat{y}(x) = \arg \max_k P\{Y=k|X\}$$

$$= \arg \max_k \frac{\int_{\mathbb{R}^P} f_k(x) \pi_k}{\sum_{l=1}^L f_l(x) \pi_l}$$

$$= \arg \max_k f_k(x) \pi_k$$

$$= \arg \max_k [\log f_k(x) + \log \pi_k]$$

Next, substitute multivariate Normal distribution:

$$\hat{Y}(X) = \arg \max_k \left\{ -\frac{1}{2} \log(2\pi)^{\frac{P}{2}} |\Sigma|^{\frac{1}{2}} - \frac{1}{2} (X - \mu_k)^T \Sigma^{-1} (X - \mu_k) + \log \pi_k \right\}$$

does not depend
on k

$$\stackrel{\text{open}}{=} \arg \max_k \left\{ -\frac{1}{2} X^T \Sigma^{-1} X + X^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k \right\}$$

Linear expression in $X \rightarrow \delta_k(X)$ - linear discriminant function

Decision boundary between classes K and L :

$$\left\{ X : P\{Y=K|X\} = P\{Y=L|X\} \right\} = \left\{ X : \delta_K(X) = \delta_L(X) \right\} = \left\{ X : \delta_K(X) - \delta_L(X) = 0 \right\}$$

MVN $\Rightarrow \left\{ X : X^T \Sigma^{-1} (\mu_K - \mu_L) - \frac{1}{2} (\mu_K - \mu_L)^T \Sigma^{-1} (\mu_K - \mu_L) + \frac{\log \pi_K}{\log \pi_L} = 0 \right\}$

For example, $K=2$:

$$\log \frac{\pi_1}{\pi_2} - \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2) = a_0$$

new notation

$$\underbrace{\Sigma^{-1} (\mu_1 - \mu_2)}_{\substack{P \times P \\ P \times 1}} = (a_1 \ a_2 \dots a_P)^T \Rightarrow \text{classify } \hat{Y}=1 \text{ if } a_0 + \sum_{j=1}^P a_j x_j > 0$$

Linear combination of all predictors (i.e., no feature selection!)

$$\text{Report posterior probability } P\{Y=K|X\} = \frac{f_K(x)\pi_K}{\sum_{k=1}^K f_k(x)\pi_k} = \frac{\exp\{-\frac{1}{2}\delta_K(X)\}}{\sum_{k=1}^K \exp\{-\frac{1}{2}\delta_k(X)\}}$$

Non-linear decision boundaries: create higher order features x_1^2, x_2^2, x_1x_2 etc

Summary :- decision boundary is a linear combination of all features

- no feature selection
- multivariate Gaussian (or at least continuous predictors)