

# Linear Classifiers: Discriminants

Prof. Rafal Kustra

Public Health Sciences

October 12, 2016

# Binary classification

- Simplest but very popular case: each object comes from one of two classes
- For convenience, let's code them  $y_i \in \{0, 1\}$
- Remember that general regression problem was a conditional expectation:

$$E(y|\mathbf{x})$$

which here becomes:

## Posterior probability

$$E(Y|\mathbf{x}) = P(Y = 1|\mathbf{x})$$

- In multiclass problems, many classifiers return an estimate of  $P(y = k|\mathbf{x})$  rather than just class  $y$ .
- These posterior probabilities are central from decision theory point of view

- As in linear regression we start with the linear predictor in  $\mathbf{x}$ :

$$\theta(\mathbf{x}) = \sum_j \beta_j x_j$$

- In regression case this is used to approximate the conditional expectation directly, i.e.

$$\hat{y}(\mathbf{x}) = E(y|\mathbf{x}) = \theta(\mathbf{x})$$

- At least in binary case, this could also work in classification and would approximate posterior probability
- But direct LS approach is inconvenient:
  - ①  $\theta(\mathbf{x})$  can easily be outside of  $[0, 1]$  interval
  - ② minimizing square distance between  $\theta(x)$  and  $y = (0, 1)$  means far away points even on a good side can badly influence  $\hat{\beta}$  estimate.

- For two classes, one could code response  $y$  as 0/1 and run OLS
- This would estimate  $P(Y = 1|\mathbf{x})$ . It is also equivalent to setting up *two* responses,  $y_1, y_2$  (two columns of now-matrix,  $Y$ ) each being an indicator for it's class. Solution  $\beta$  (a  $p \times 2$  matrix) is:

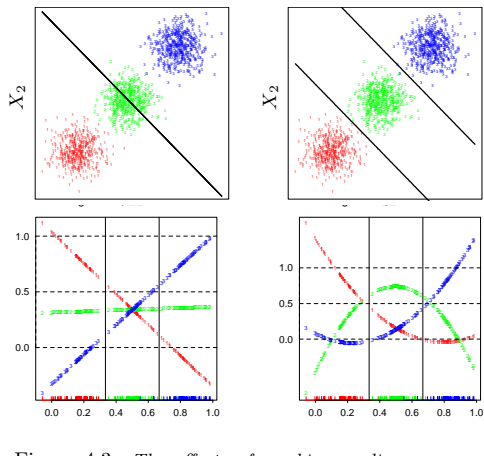
$$\hat{\beta} = (X'X)^{-1}X'Y$$

with corresponding two columns of  $\hat{Y}$

- The two columns of  $\hat{Y}$  are just reflections of each other (check!). Remember contrasts.
- One could extend this approach to  $K$ -class, coding each as binary column of  $Y$ . This gives you  $K$  columns of  $\hat{Y}$ , each as an estimate of  $P(Y = k|\mathbf{x})$
- For  $K > 2$  we have a potential of *masking* as a problem with this approach (in addition to improper range of  $\hat{Y}$ )

- In  $p$ -dim space of predictors, a set of points ('curve') where classifier switches from predicting class A to class B
- For Bayes classifier: a set of points where posterior probabilities for two (or more) classes are equal: predicting either of these classes leads to same expected cost
- For OLS classifier above: set of points where  $\hat{y}_1 = \hat{y}_2$  for class 1 and 2 (i.e., fitted values corresponding to column 1 or 2 of response *matrix*  $Y$ )

# Masking Geometry

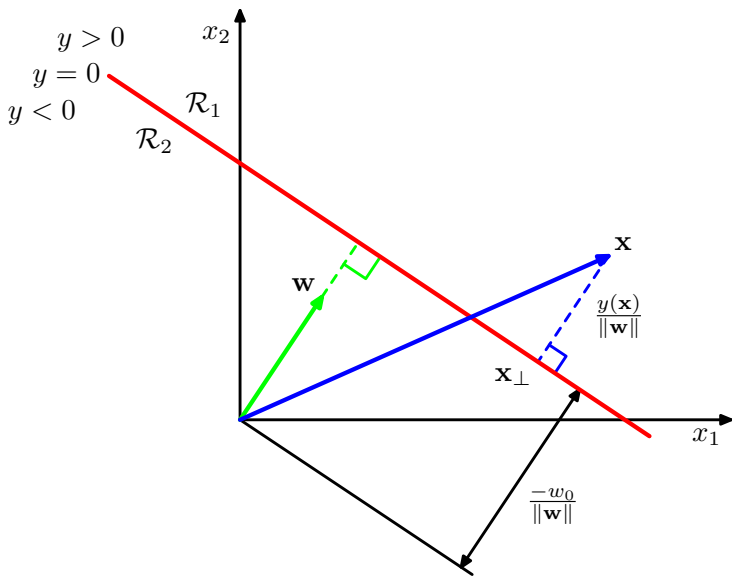


# Linear Discriminants

- One way to deploy linear functions of  $\mathbf{x}$  is through *discriminant* functions
- We will seek  $\hat{\beta}$  that predicts the class of object  $\mathbf{x}$  by the sign of  $\theta(\mathbf{x})$ , i.e. in binary case, we predict class  $y = 1$ , if  $\theta(\mathbf{x}) > 0$  and class  $y = 0$  otherwise
- This creates a linear boundary,  $\mathcal{D}$  characterized by all points  $\mathbf{x}$  such that  $\theta(\mathbf{x}) = 0$
- As it stands boundary  $\mathcal{D}$  passes through the origin. Adding an intercept term we remedy that:

$$\theta(\mathbf{x}) = \alpha + \beta^T \mathbf{x}$$

- In this formulation one can see that  $\alpha$  locates a surface  $\alpha/||\beta||$  units away from origin
- Vector  $\beta$  is a normal to the decision hyperplane: we can think of classification procedure as projecting points  $\mathbf{x}$  on (extension of)  $\beta$  and checking if they lie to the left or right of  $\alpha/||\beta||$





# Fisher Discriminant Analysis for 2 classes

- In 1936, Fisher proposed finding a discriminant direction,  $\beta$  such that when points from two classes were projected onto the line two things happen:
  - 1 The class means would be well separated
  - 2 The spread within classes was small
- This leads to finding vector  $\beta$  that maximizes the spread but also minimizes the within-class variance. For  $K$  classes spread of means can be measured by *between-class* variance:

$$B = \sum_{k=1}^2 (\bar{\mathbf{x}}_k - \bar{\mathbf{x}})(\bar{\mathbf{x}}_k - \bar{\mathbf{x}})^T = \frac{1}{2}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T$$

- The common within-class covariance matrix can be estimated by:

$$W = (N - K)^{-1} \sum_i (\mathbf{x}_i - \bar{\mathbf{x}}_{k(i)}) (\mathbf{x}_i - \bar{\mathbf{x}}_{k(i)})^T$$

- The direction of the line onto which projected data would maximize the spread of class means while minimizing the within-class spread is  $\mathbf{a}$ :

$$\mathbf{a} = \operatorname{argmax}_{\alpha} \frac{\alpha^T B \alpha}{\alpha^T W \alpha}$$

- This direction is one which is a compromise between the line that is perpendicular to the line joining the two class means, and the line which is perpendicular to the principal axis of the *pooled* covariance structure
- Another way to look at it is to assume that the points in different classes have different means but the same covariance matrix,  $\Sigma$
- Then the variance of an observation,  $\mathbf{x}_i$  projected on any vector,  $\alpha$ , is:

$$\operatorname{Var}(\mathbf{x}_i^T \alpha) = \alpha^T \Sigma \alpha$$

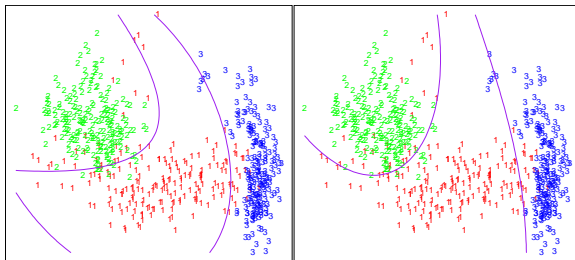


Figure 4.6: *Two methods for fitting quadratic boundaries. The left plot shows the quadratic decision bound-*

- The within-class covariance matrix,  $W$ , is a standard estimate of the common covariance of points with different means. Hence the denominator of the Fisher criterion is an estimate of the variance of an observation projected on vector,  $\alpha$
- It is easy to check the the solution to FDA criterion is:

$$\mathbf{a} \propto W^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$$

and is usually normalized so that  $\mathbf{a}^T W \mathbf{a} = 1$

- Classification is done by measuring the distance of a new point,  $\mathbf{x}_0$ , to the class means *in the projected space* (here a line) and choosing the class with the closest mean, ie:

$$C(\mathbf{x}_0) = \underset{k}{\operatorname{argmin}} \mathbf{x}_0^T \mathbf{a} - \bar{\mathbf{x}}_k^T \mathbf{a} = \underset{k}{\operatorname{argmin}} \mathbf{a}^T (\mathbf{x}_0 - \bar{\mathbf{x}}_k)$$

# LDA for more classes

- With two classes there is one hyperplane separating the classes in a linear classifier. With  $K$  classes one needs  $K - 1$  hyperplanes
- Another way to look at this: the Fisher criterion:

$$\frac{\alpha^T B \alpha}{\alpha^T W \alpha}$$

has rank  $K - 1$ , because of the numerator (unless the dimensionality of the data is smaller than number of classes). JN Rao in 1956 extended the FDA to LDA for multiple classes, by finding up to  $K - 1$  directions,  $\mathbf{a}_k$ , starting with the maximum (Fisher) one,  $\mathbf{a}_1$  and subsequent ones,  $\mathbf{a}_k$ , for  $k = 2, \dots, K - 1$  subject to:

$$\mathbf{a}_k^T W \mathbf{a}_j, \quad j = 1, \dots, k - 1$$

This is a kind-of orthogonality constraint in the metric of pooled covariance structure

- Classification of new point,  $\mathbf{x}_0$  again happens by first projecting the new point and class means onto the space spanned by discriminant directions,  $\mathbf{a}_k$ , and finding the closest class mean to the new point using Euclidean distance. Let  $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_{K-1}]$  be a matrix with  $K - 1$  columns being discriminant directions. Then:

$$C(\mathbf{x}_0) = \underset{k}{\operatorname{argmin}} \ ||A(\mathbf{x}_0 - \bar{\mathbf{x}}_k)||$$

# Classification via Maximum Likelihood

- Lets enter a Gaussian world:

$$f_k(\mathbf{x}) = N_p(\mu_k, \Sigma)$$

so that class-densities are Gaussian with *the same* covariance matrix,  $\Sigma$  and class-specific means

- The joint log-likelihood of data separates:

$$-1/2p \log[(2\pi)|\Sigma|] \sum_i \left\{ (\mathbf{x}_i - \mu_{k(i)})^T \Sigma^{-1} (\mathbf{x}_i - \mu_{k(i)}) \right\}$$

- The ML estimator (or Bayes estimator, under the same prior probabilities) assigns  $\mathbf{x}_0$  to that class,  $k$ , that has maximum likelihood for  $\mathbf{x}_0$ , or equivalently *minimum Mahalanobis distance* to class mean:

$$(\mathbf{x}_0 - \mu_k)^T \Sigma^{-1} (\mathbf{x}_0 - \mu_k)$$

- One inserts the MLE estimates of  $\mu_k, \Sigma$  in above
- It is an easy theorem to show that minimum Mahalanobis distance is same as Euclidean distance in the space of LDA directions

- In LDA with Gaussian assumption we assumed a distribution for the  $\mathbf{X}$  vector of inputs.
- In Logistic regression we go straight after the posterior,  $P(y|\mathbf{x})$ . For *binary classification* code two classes of  $y$  as  $(0, 1)$ . We need a link between  $E(y|x) = P(y = 1|x)$  and  $x$ .
- Logit link is most popular:

$$P(y = 1|\mathbf{x}; \beta) = \frac{\exp(\mathbf{x}^T \beta)}{1 + \exp(\mathbf{x}^T \beta)}$$

- This leads to the inverse function:

$$\log \left( \frac{P(y = 1|\mathbf{x}; \beta)}{1 - P(y = 1|\mathbf{x}; \beta)} \right) = \mathbf{x}^T \beta = \sum_j \beta_j x_j$$

the LHS is frequently called *log-odds* at  $\mathbf{x}$



# Parameter estimation

- We have a model for predicting posterior probabilities given any  $\mathbf{x}$  input that depends on parameter vector,  $\beta$
- Maximum-likelihood estimate of  $\beta$  starts with forming a log-likelihood function
- The model is conditional on  $\mathbf{x}$ : the only randomness is in  $y$  which is binary
- Given a training set with  $y_i$ ,  $i = 1, \dots, N$  where each  $y_i \in \{0, 1\}$  we have  $N$  Bernoulli trials, each with probability,  $p(\mathbf{x}_i)$ . We can write down a likelihood as follows:

$$\prod_i p(\mathbf{x}_i)^{y_i} (1 - p(\mathbf{x}_i))^{1-y_i}$$

which leads to log-likelihood:

$$\mathcal{L}(\beta) = \sum_i y_i \log(p(\mathbf{x}_i)) (1 - y_i) \log(1 - p(\mathbf{x}_i))$$

- or:

$$\sum_i y_i \log \left( \frac{p(\mathbf{x}_i)}{1 - p(\mathbf{x}_i)} \right) + \log(1 - p(\mathbf{x}_i))$$

- In terms of  $\beta$  this becomes:

$$\sum_i \sum_j y_i x_{ij} \beta_j - \sum_i \log \left( 1 + \exp \left( \sum_j x_{ij} \beta_j \right) \right)$$

- MLE estimates of  $\beta$  are of course obtained by maximizing the log-likelihood. As this is a non-linear function, iterative process is used. Normally a slightly modified version of well-known Newton-Raphson function minimizer, called Iteratively Reweighted Least Squares is used. Book has the details in Ch. 4.3.3.

# Boundary regions for Logistic regression

- So far we have a model for posterior probabilities. To classify as usual we predict the class with maximum posterior probability
- Even though the probabilities are not linear in  $\mathbf{x}$  the boundary between the two classes is.
- The boundary is defined as a set of  $\mathbf{x}$  points where posterior probabilities are equal: here 0.5 each
- This is equivalent to the odds ratio being 1:

$$p(\mathbf{x}) = 1 - p(\mathbf{x}) \iff \frac{p(\mathbf{x})}{1 - p(\mathbf{x})} = 1 \iff \text{logit}(p(\mathbf{x})) = 0$$

which in logistic regression means set of  $\mathbf{x}$  such that:

$$\mathbf{x}^T \hat{\beta} = 0$$

which of course is linear in  $\mathbf{x}$

# Comparison with Gaussian LDA

- LDA starts with a Gaussian assumption on class-densities,  $f_k$  (for binary case,  $k = 1, 2$ ):

$$f_k(\mathbf{x}) \equiv f_{X|k}(\mathbf{x}; \Sigma, \mu_k) \propto |\Sigma|^{-p/2} \exp \left\{ -(\mathbf{x} - \mu_k)^T \Sigma^{-1} (\mathbf{x} - \mu_k) \right\}$$

- With 0-1 loss, we classify by highest posterior probability, which in this case is:

$$P(y = k|\mathbf{x}) = \frac{f_k(\mathbf{x}; \Sigma, \mu_k) \pi_k}{f_X(\mathbf{x})}$$

- Since denominator does not depend on class  $k$  it is not relevant for classification. Hence for binary case the boundary region becomes:

$$\begin{aligned} \log \frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})} &= (\mathbf{x} - \mu_0)^T \Sigma^{-1} (\mathbf{x} - \mu_0) \\ &\quad - (\mathbf{x} - \mu_1)^T \Sigma^{-1} (\mathbf{x} - \mu_1) + \log \frac{\pi_1}{\pi_0} \end{aligned}$$

which is also linear in  $\mathbf{x}$  since  $\mathbf{x}^T \Sigma^{-1} \mathbf{x}$  term cancels out

- So we have two binary procedures with quite different assumptions producing linear classification boundaries
- One crucial difference is that LDA is forced to model the whole distribution of the  $x$  space. For high-dimension problems this can be very difficult and wasteful
  - ① Normality assumption may be quite wrong
  - ② Even if normality is plausible, equal-covariance assumption may not be
  - ③ Number of parameters to estimate could be enormous:  $O(p^2)$  rather than  $O(p)$  that logistic regression requires
- Logistic on the other hand goes straight after posterior and can a much better choice for most problems. However:
  - ① Having distributional framework for  $X$  space can be beneficial when extending the model
  - ② IF we can get the distribution of  $X$  right, we can have better prediction in  $X$  different from ones in the training set

# Polychotomous Logistic regression

- Quite easy to extend logistic regression to  $K$  classes
- First need to choose a reference class. This is arbitrary in a sense that model predictions will not change (but interpretation of parameters will).
- With  $K$  classes, let's code  $y$  as:

$$y_i \in (0, \dots, K - 1)$$

- As before, we will model log-odds ratio as a linear function of  $\mathbf{x}$ . Reference class used for denominator. So, for class  $k$ ,  $k = 1, \dots, K - 1$  we have:

$$\log \frac{P(y = k|\mathbf{x})}{P(y = 0|\mathbf{x})} = \mathbf{x}^T \beta_k$$

- Of course all those probabilities have to sum up to one across classes, which leads to inverse functions:

$$P(y = k|\mathbf{x}; \beta) = \frac{\exp(\mathbf{x}^T \beta_k)}{1 + \sum_{\kappa=1}^{K-1} \exp(\mathbf{x}^T \beta_{\kappa})}; \quad k = 1, \dots, K - 1$$

- and for the reference class ( $k = 0$ ):

$$P(y = k|\mathbf{x}; \beta) = \frac{1}{1 + \sum_{\kappa=1}^{K-1} \exp(\mathbf{x}^T \beta_{\kappa})}$$

- Hence for  $K$  classes there are  $K - 1$  vectors of parameters,  $\beta_k$  defining  $K - 1$  boundaries.
- Estimation is similar as in binary case for logistic regression. The Hessian matrix is block-diagonal but the weight matrix is not which leads to slightly more complicated algorithm.