

# Machine Learning 2020

Volker Roth

Department of Mathematics & Computer Science  
University of Basel

16th March 2020

## Section 3

# Classification

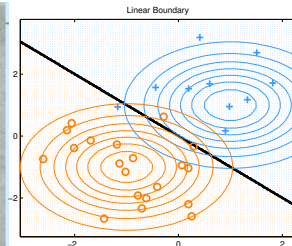
# Classification

## Example

Sorting fish according to species using optical sensing

Features:

- Length
- Brightness
- Width
- Shape of fins



# Bayesian Decision Theory

- Assign observed  $\mathbf{x} \in \mathbb{R}^d$  into one of  $k$  classes. A **classifier** is a mapping that assigns labels to observations

$$f_{\alpha} : \mathbf{x} \rightarrow \{1, \dots, k\}.$$

- For any observation  $\mathbf{x}$  there exists a set of  $k$  **possible actions**  $\alpha_i$ , i.e.  $k$  different assignments of labels.
- The **loss**  $L$  incurred for taking action  $\alpha_i$  when the true label is  $j$  is denoted by a loss matrix  $L_{ij} = L(\alpha_i | c = j)$ .
- “Natural” 0 – 1 loss function can be defined by simply **counting misclassifications**:  $L_{ij} = 1 - \delta_{ij}$ , where

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{otherwise.} \end{cases}$$

# Bayesian Decision Theory (cont'd)

- A classifier is trained on a set of **observed pairs**  $\{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, c) = p(c|\mathbf{x})p(\mathbf{x})$
- The probability that a given  $\mathbf{x}$  is member of class  $c_j$ , i.e. the posterior probability of membership in class  $j$ , is obtained via the **Bayes rule**:

$$P(c_j|\mathbf{x}) = \frac{\text{Given the label, observation is generated } p(\mathbf{x}|c=j)}{p(\mathbf{x})} \quad \text{Nature picks a label first } P(c=j),$$

where

$$p(\mathbf{x}) = \sum_{j=1}^k p(\mathbf{x}|c=j)P(c=j).$$

- Given an observation  $\mathbf{x}$ , the expected loss associated with choosing action  $\alpha_i$  (the **conditional risk** or **posterior expected loss**) is

$$R(f_{\alpha_i}|\mathbf{x}) = \sum_{j=1}^k L_{ij}P(c_j|\mathbf{x}) \stackrel{(\text{if } L_{ij}=1-\delta_{ij})}{=} \sum_{j \neq i} P(c_j|\mathbf{x}) = 1 - P(c_i|\mathbf{x}).$$

# Bayesian Decision Theory (cont'd)

- **Goal:** minimize the **overall risk** of the classifier  $f_\alpha$ :

$$R(f_\alpha) = \int_{\mathbf{R}^d} R(f_\alpha(\mathbf{x})|\mathbf{x})p(\mathbf{x}) d\mathbf{x}.$$

- If  $f_\alpha(\mathbf{x})$  minimizes the conditional risk  $R(f_\alpha(\mathbf{x})|\mathbf{x})$  for every  $\mathbf{x}$ , the overall risk will be minimized as well.
- This is achieved by the **Bayes optimal classifier** which chooses the mapping

$$f(\mathbf{x}) = \operatorname{argmin}_i \sum_{j=1}^k L_{ij} p(c = j|\mathbf{x}).$$

- For 0 – 1 loss this reduces to classifying  $\mathbf{x}$  to the class **with highest posterior probability**:

$$f(\mathbf{x}) = \operatorname{argmax}_i p(c = i|\mathbf{x}).$$

# Bayesian Decision Theory (cont'd)

- **Simplification:** only 2 classes:  $c$  is Bernoulli RV.
- Bayes optimal classifier is defined by the **zero crossings** of the **Bayes optimal discriminant function**

$$G(\mathbf{x}) = P(c_1|\mathbf{x}) - P(c_2|\mathbf{x}), \text{ or } g(\mathbf{x}) = \log \frac{P(c_1|\mathbf{x})}{P(c_2|\mathbf{x})}.$$

- Link to regression: use **encoding**  $\{+1, -1\}$  for the two possible states  $c_{1,2}$  of  $c$ . The conditional expectation of  $c|\mathbf{x}$  equals the Bayes discriminant function:

$$E[c|\mathbf{x}] = \sum_{c \in \{+1, -1\}} cP(c|\mathbf{x}) = P(c_1|\mathbf{x}) - P(c_2|\mathbf{x}) = G(\mathbf{x}).$$

- Classification can be viewed as a (local) approximation of  $G(\mathbf{x}) = E[c|\mathbf{x}]$  near its zero crossings.

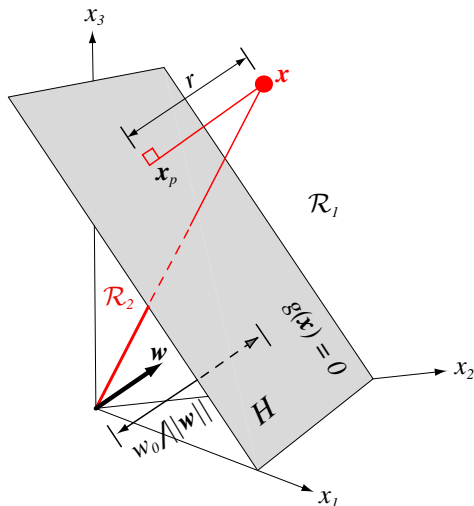
# Linear Discriminant Functions

- Problem: direct approximation of  $G$  would require the knowledge of the Bayes optimal discriminant.
- One approach: Define a **parametrized family of classifiers**  $\mathcal{F}_{\mathbf{w}}$  from which we can choose one (or more) function(s) by some **inference mechanism**.
- One such family is the set of linear discriminant functions
$$g(\mathbf{x}; \mathbf{w}) = w_0 + \mathbf{w}^t \mathbf{x}.$$
- Two-category case: Decide  $c_1$  if  $g(\mathbf{x}; \mathbf{w}) > 0$  and  $c_2$  if  $g(\mathbf{x}; \mathbf{w}) < 0$ .
- Equation  $g(\mathbf{x}; \mathbf{w}) = 0$  defines the **decision surface**.  
Linearity of  $g(\mathbf{x}; \mathbf{w}) \rightsquigarrow$  **hyperplane**  $\rightsquigarrow$   $\mathbf{w}$  is orthogonal to any vector lying in the plane.
- The hyperplane divides the feature space into half-spaces  $\mathcal{R}_1$  (“positive side”) and  $\mathcal{R}_2$  (“negative side”).



# Decision Hyperplanes

- $g(\mathbf{x}; \mathbf{w})$  defines distance  $r$  from  $\mathbf{x}$  to the hyperplane:  $\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$ .
- $g(\mathbf{x}_p) = 0 \Rightarrow g(\mathbf{x}) = r\|\mathbf{w}\| \Leftrightarrow r = g(\mathbf{x})/\|\mathbf{w}\|$ .



# Generalized Linear Discriminant Functions

Use basis functions  $\{b_1(\mathbf{x}), \dots, b_m(\mathbf{x})\}$ , where each  $b_i(\mathbf{x}) : \mathbb{R}^d \mapsto \mathbb{R}$ , and  $g(\mathbf{x}; \mathbf{w}) = w_0 + w_1 b_1(\mathbf{x}) + \dots + w_m b_m(\mathbf{x}) =: \mathbf{w}^t \mathbf{y}$  (note that we have redefined  $\mathbf{y}$  here in order to be consistent with the following figure)

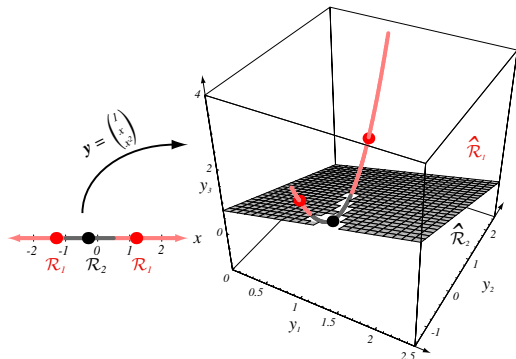


Fig 5.5 in (Duda& Hart)

# Generalized Linear Discriminant Functions

Use basis functions  $\{b_1(\mathbf{x}), \dots, b_m(\mathbf{x})\}$ , where each  $b_i(\mathbf{x}) : \mathbb{R}^d \mapsto \mathbb{R}$ , and

$$g(\mathbf{x}; \mathbf{w}) = w_0 + w_1 b_1(\mathbf{x}) + \dots + w_m b_m(\mathbf{x}) =: \mathbf{w}^t \mathbf{y}.$$

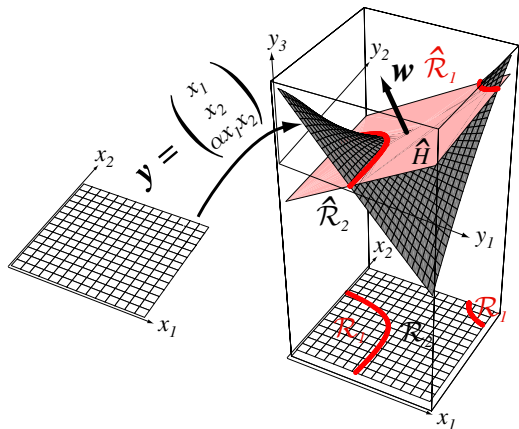


Fig 5.6 in (Duda& Hart)

# Separable Case

- Consider sample  $\{\mathbf{y}_i, c_i\}_{i=1}^n$ . If there exists  $f(\mathbf{y}; \mathbf{w}) = \mathbf{y}^t \mathbf{w}$  which is positive for all examples in class 1 and negative for all examples in class 2, we say that the sample is **linearly separable**.
- “Normalization”: replace all samples labeled  $c_2$  by their negatives  $\rightsquigarrow$  simply write  $\mathbf{y}^t \mathbf{w} > 0$  for all samples.
- Each sample places a constraint on the possible location of  $\mathbf{w}$   $\rightsquigarrow$  **solution region**.

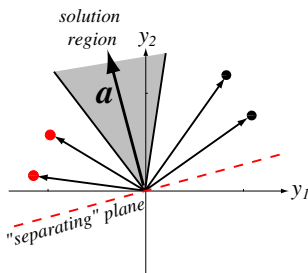
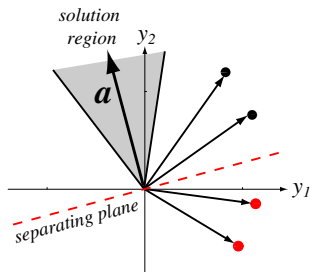


Fig 5.8 in (Duda& Hart)

## Separable Case: margin

- Different solution vectors may have different **margins**  $b : \mathbf{y}^t \mathbf{w} \geq b > 0$ .
- Intuitively, **large margins are good**.

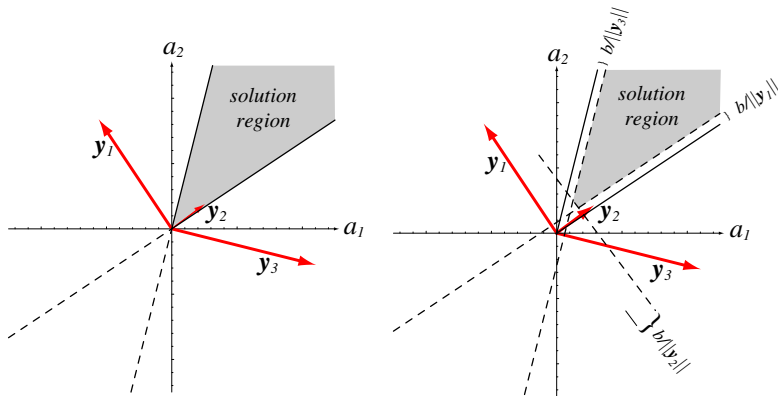


Fig 5.9 in (Duda & Hart)

# Gradient Descent

- Solve  $\mathbf{y}^t \mathbf{w} > 0$  by defining  $J(\mathbf{w})$  such that a minimizer of  $J$  is a solution.
- Start with initial  $\mathbf{w}(1)$ , and choose next value by moving in the direction of steepest gradient:  $\mathbf{w}(k+1) = \mathbf{w}(k) - \eta(k) \nabla J(\mathbf{w}(k))$ .
- Alternatively, use second order expansion (Newton):  
 $\mathbf{w}(k+1) = \mathbf{w}(k) - H^{-1} \nabla J(\mathbf{w}(k))$ .

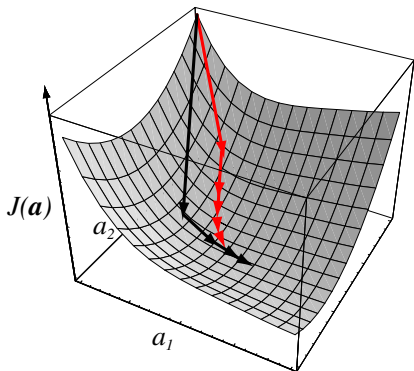


Fig 5.10 in (Duda&Hart)

# Minimizing the Perceptron Criterion

- Solve  $\mathbf{y}^t \mathbf{w} > 0$  by defining  $J(\mathbf{w})$  such that a minimizer of  $J$  is a solution.
- Most obvious: number of misclassifications, but **not differentiable**.
- Alternative choice:  $J_p(\mathbf{w}) = \sum_{\mathbf{y} \in \mathcal{M}} -\mathbf{y}^t \mathbf{w}$ , where  $\mathcal{M}(\mathbf{w})$  is the set of samples **misclassified** by  $\mathbf{w}$ .
- Since  $\mathbf{y}^t \mathbf{w} < 0 \forall \mathbf{y} \in \mathcal{M}$ ,  $J_p$  is non-negative, and zero only if  $\mathbf{w}$  is a solution.
- Gradient:  $\nabla J(\mathbf{w}) = - \sum_{\mathbf{y} \in \mathcal{M}} \mathbf{y}$   
$$\mathbf{w}(k+1) = \mathbf{w}(k) + \eta(k) \sum_{\mathbf{y} \in \mathcal{M}} \mathbf{y}.$$

This defines the **Batch Perceptron algorithm**.

## Minimizing the Perceptron Criterion (2)

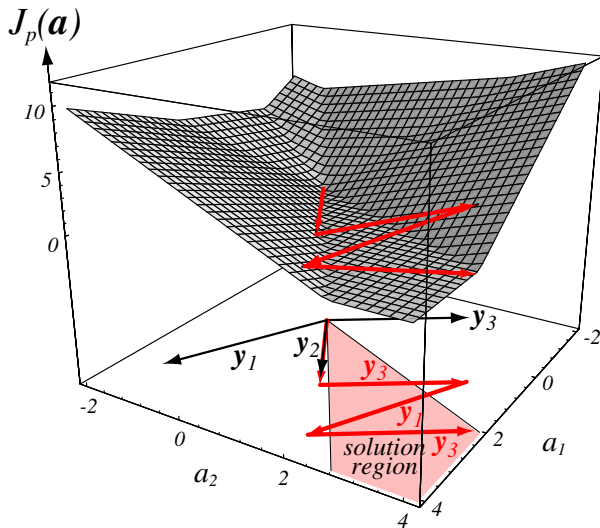


Fig 5.12 in (Duda& Hart)



# Fixed-Increment Single Sample Perceptron

- Fix learning rate  $\eta(k) = 1$ .
- Sequential single-sample updates: use superscripts  $\mathbf{y}^1, \mathbf{y}^2, \dots$  for misclassified samples  $\mathbf{y} \in \mathcal{M}$ . Ordering is irrelevant.
- Simple algorithm:

$\mathbf{w}(1)$  arbitrary

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \mathbf{y}^k, \quad k \geq 1$$

## Perceptron Convergence Theorem

If the samples are linearly separable, the sequence of weight vectors given by the Fixed-Increment Single Sample Perceptron algorithm will terminate at a solution vector.

**Proof:** exercises.

A number of problems with the perceptron algorithm:

- When the data are separable, there are **many solutions**, and which one is found **depends on the starting values**.
- In particular, no separation margin can be guaranteed (however, there exist modified versions...)
- The number of steps can be **very** large.
- When the data are **not separable**, the algorithm will **not necessarily converge**, and cycles may occur. The cycles can be long and therefore hard to detect.
- Method “technical” in nature, no (obvious) probabilistic interpretation (but we will see that there is one).

But the perceptron algorithm is historically important (1957, one of the first ML algorithms!), was even implemented in analog hardware(!)

# Generative (or Informative) vs Discriminative

- Notation: For the following discussion it is more convenient to go back to the original  $\mathbf{x}$ -vectors (potentially after some basis expansion) instead of using the “normalized” representation  $\mathbf{y}$ .
- Two main strategies:
  - ▶ **Generative:** Generative classifiers specify how to generate data using the **class densities**. Likelihood/posterior of each class is examined and classification is done by assigning to the most likely class.
  - ▶ **Discriminative:** These classifiers focus on modeling the **class boundaries** or the class membership probabilities directly. No attempt is made to model the underlying class conditional densities.

# Generative Classifiers

- Central idea: model the conditional class densities  $p(\mathbf{x}|c)$ .
- Assuming a parametrized class conditional density  $p_{\mathbf{w}_j}(\mathbf{x}|c = j)$  and collecting all model parameters in a vector  $\mathbf{w}$ , a typical (Frequentist) approach now proceeds by maximizing the log likelihood

$$\hat{\mathbf{w}}_{MLE} = \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^n \log p_{\mathbf{w}}(\mathbf{x}_i|c_i)$$

- The resulting estimate  $\hat{\mathbf{w}}_{MLE}$  might then be plugged into Bayes rule to compute the posteriors:

$$P(c_j|\mathbf{x}) = \frac{p_{\hat{\mathbf{w}}_{MLE}}(\mathbf{x}|c = j)}{p(\mathbf{x})} P(c = j).$$

# Generative Classifiers: LDA

- In *Linear Discriminant Analysis* (LDA), a Gaussian model is used where all classes share a common covariance matrix  $\Sigma$ :

$$p_{\mathbf{w}}(\mathbf{x}|c = j) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \Sigma).$$

- The resulting discriminant functions are linear:

$$\begin{aligned} g(\mathbf{x}) &= \log \frac{P(c_1|\mathbf{x})}{P(c_2|\mathbf{x})} = \log \frac{P(c_1)\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma)}{P(c_2)\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \Sigma)} \\ &= \underbrace{\log \frac{P(c_1)}{P(c_2)} - \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)^t \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}_{w_0} \\ &\quad + \underbrace{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \Sigma^{-1} \mathbf{x}}_{\mathbf{w}^t \mathbf{x}} \\ &= w_0 + \mathbf{w}^t \mathbf{x}, \quad \text{with } \mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2). \end{aligned}$$

# LDA algorithm

- Let  $\hat{\Sigma}$  be an estimate of the shared covariance matrix  $\Sigma$ :

$$\Sigma_c = \frac{1}{n_c} \sum_{\mathbf{x} \in \mathcal{X}_c} (\mathbf{x} - \mathbf{m}_c)(\mathbf{x} - \mathbf{m}_c)^t, \quad c \in \{c_1, c_2\}$$

$$\hat{\Sigma} = \frac{1}{2}(\Sigma_1 + \Sigma_2).$$

- Let  $\mathbf{m}_j$  an estimate of  $\mu_j$ :

$$\mathbf{m}_c = \frac{1}{n_c} \sum_{\mathbf{x} \in \mathcal{X}_c} \mathbf{x}, \quad n_c = |\mathcal{X}_c|.$$

- Fisher's LDA finds the weight vector

$$\mathbf{w}^F = \hat{\Sigma}^{-1}(\mathbf{m}_1 - \mathbf{m}_2).$$

This  $\mathbf{w}^F$  asymptotically coincides with the Bayes-optimal  $\mathbf{w}$   
if Gaussian model is correct.

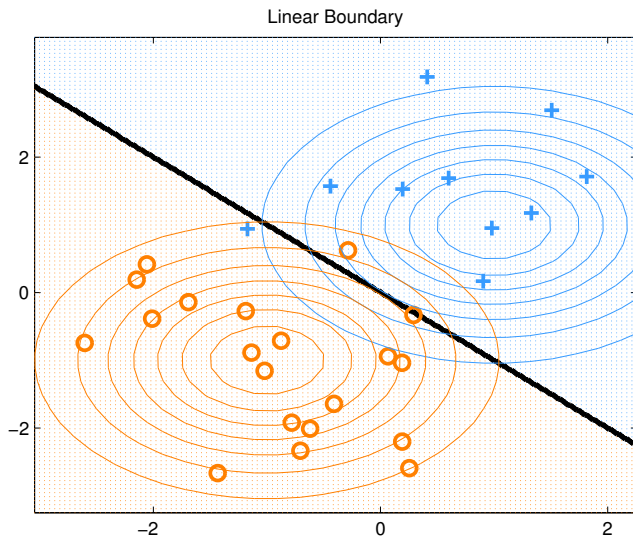


Fig 4.5 in K. Murphy

# Fishers discriminant and least squares

**Remark:** The Fisher vector  $\hat{\mathbf{w}}^F = \Sigma_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$  coincides with the solution of the LS problem  $\hat{\mathbf{w}}^{LS} = \arg \min_{\mathbf{w}} \|\mathbf{A}\mathbf{w} - \mathbf{b}\|^2$  if

$n_1$  = # samples in class 1

$n_2$  = # samples in class 2

$$\mathbf{b} = \begin{pmatrix} +1/n_1 \\ \cdot \\ +1/n_1 \\ -1/n_2 \\ \cdot \\ -1/n_2 \end{pmatrix}, \quad \mathbf{A} = \begin{pmatrix} \mathbf{x}_1^t \\ \cdot \\ \mathbf{x}_{n_1}^t \\ \mathbf{x}_{n_1+1}^t \\ \cdot \\ \mathbf{x}_{n_1+n_2}^t \end{pmatrix},$$

with  $\mathbf{m} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$  (i.e. origin in sample mean).



# Fishers discriminant and least squares (cont'd)

## Proofsketch:

- Shared covariance matrix also called “within class covariance”  
 $\Sigma_W \propto \sum_{\mathbf{x} \in \mathcal{X}_c} (\mathbf{x} - \mathbf{m}_c)(\mathbf{x} - \mathbf{m}_c)^t, \quad c = c_1, \text{ or } c = c_2.$
- Its counterpart is the “between class covariance”  
 $\Sigma_B \propto (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^t$
- The sum of both is the “total covariance”  $\Sigma_B + \Sigma_W = \Sigma_T$   
 $\Sigma_T \propto \sum_i (\mathbf{x}_i - \mathbf{m})(\mathbf{x}_i - \mathbf{m})^t = A^t A.$
- We know that  $\mathbf{w}^F \propto \Sigma_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2) \rightsquigarrow \Sigma_W \mathbf{w}^F \propto (\mathbf{m}_1 - \mathbf{m}_2).$
- Now  $\Sigma_B \mathbf{w}^F = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^t \mathbf{w}^F \rightsquigarrow \Sigma_B \mathbf{w}^F \propto (\mathbf{m}_1 - \mathbf{m}_2).$
- $\Sigma_T \mathbf{w}^F = (\Sigma_B + \Sigma_W) \mathbf{w}^F \rightsquigarrow \Sigma_T \mathbf{w}^F \propto (\mathbf{m}_1 - \mathbf{m}_2).$
- With  $A^t A = \Sigma_T$ ,  $A^t \mathbf{b} = \mathbf{m}_1 - \mathbf{m}_2$ , we arrive at  
 $\mathbf{w}^F \propto \Sigma_T^{-1}(\mathbf{m}_1 - \mathbf{m}_2) = \Sigma_T^{-1} A^t \mathbf{b} = (A^t A)^{-1} A^t \mathbf{b} = \mathbf{w}^{LS}.$

## Fishers discriminant and least squares (cont'd)

- Focus on last equation. For notational simplicity, denote the least-squares estimate  $\mathbf{w}^{LS}$  by  $\mathbf{w}$ .
- Introducing the “residual sum of squares” as the least-squares cost function, the equation follows from:

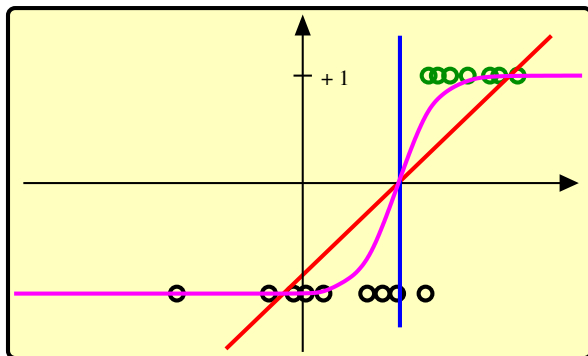
$$\begin{aligned}RSS(\mathbf{w}) &= \sum_{i=1}^n [b_i - \mathbf{w}^t \mathbf{x}_i]^2 \\ \frac{\partial RSS(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} [\mathbf{b}^t \mathbf{b} - 2\mathbf{b}^t A\mathbf{w} + \mathbf{w}^t A^t A\mathbf{w}] \\ &= -2A^t \mathbf{b} + 2A^t A\mathbf{w} \stackrel{!}{=} \mathbf{0} \Rightarrow \mathbf{w} = (A^t A)^{-1} A^t \mathbf{b}.\end{aligned}$$

$A^t \mathbf{b} = A^t A\mathbf{w}$  are called the **normal equations**.

- We have used the following results from matrix calculus:

$$\begin{aligned}\frac{\partial}{\partial \mathbf{x}} \mathbf{y}^t \mathbf{x} &= \mathbf{y} \\ \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^t M \mathbf{x} &= 2M\mathbf{x}, \text{ if } M \text{ is symmetric}\end{aligned}$$

## Fishers discriminant and least squares (cont'd)



- Two-class LDA solution viewed as indicator regression.
- Magenta curve: Bayes-optimal discriminant function  
 $G(\mathbf{x}) = P(c = +1|\mathbf{x}) - P(c = -1|\mathbf{x})$
- Red line: Regression fit  $\rightsquigarrow$  zero crossing determines the separating hyperplane (vertical blue line).

# Discriminative classifiers

- Discriminative classifiers focus directly on the discriminant function.
- In general, they are more flexible with regard to the class conditional densities they are capable of modeling.
- Notation: Can use any class encoding scheme. Here:  $c \in \{0, 1\}$ .
- Bayes formula:

$$\begin{aligned} g(\mathbf{x}) &= \log \frac{P(c = 1|\mathbf{x})}{P(c = 0|\mathbf{x})} \\ &= \log \frac{p(\mathbf{x}|c = 1)P(c = 1)}{p(\mathbf{x}|c = 0)P(c = 0)}, \end{aligned}$$

- Can model any conditional probabilities that are exponential “tilts” of each other:

$$p(\mathbf{x}|c = 1) = e^{g(\mathbf{x})} p(\mathbf{x}|c = 0) \frac{P(c = 0)}{P(c = 1)}$$

# Logistic Regression (LOGREG)

- **Logistic regression** uses a **linear discriminant function**, i.e.  $g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$ .
- For the special case  $p(\mathbf{x}|c) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{0,1}, \Sigma)$ , same as LDA:

$$p(\mathbf{x}|c=1) = \mathcal{N}(\mathbf{x}, \boldsymbol{\mu}_1, \Sigma) = e^{g(\mathbf{x})} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_0, \Sigma) \frac{P(c=0)}{P(c=1)}$$

$$\Rightarrow g(\mathbf{x}) = w_0 + \mathbf{w}^t \mathbf{x} = \log \frac{P(c=1)\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma)}{P(c=0)\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_0, \Sigma)}$$

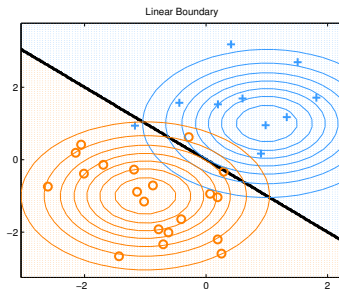
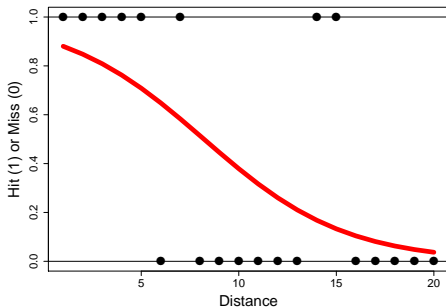


Fig 4.5 in K. Murphy

# Logistic Regression (LOGREG)

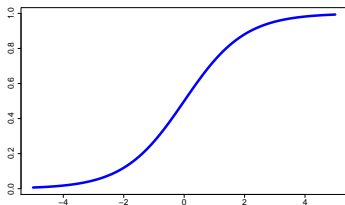
- Two-class problem with Bernoulli RV  $c$  taking values in  $\{0, 1\}$   
 $\rightsquigarrow$  sufficient to represent  $P(1|\mathbf{x})$ , since  $P(0|\mathbf{x}) = 1 - P(1|\mathbf{x})$ .
- “Success probability” of the Bernoulli RV:  $\pi(\mathbf{x}) := P(1|\mathbf{x})$ .
- Probability of miss ( $c = 0$ ) or hit ( $c = 1$ ) as a function of  $\mathbf{x}$ :  
$$p(c|\mathbf{x}) = \pi(\mathbf{x})^c (1 - \pi(\mathbf{x}))^{1-c}, \quad \pi(\mathbf{x}) = P(1|\mathbf{x}) = E[c|\mathbf{x}].$$
- Basketball example:



Adapted from Fig. 7.5.1 in (B. Flury)

# Logistic Regression (LOGREG)

- LOGREG:  $g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 = \log \frac{P(c=1|\mathbf{x})}{P(c=0|\mathbf{x})} = \log \frac{\pi(\mathbf{x})}{1-\pi(\mathbf{x})}$
- This implies  $\frac{\pi(\mathbf{x})}{1-\pi(\mathbf{x})} = \exp\{g(\mathbf{x})\}$   
 $\Rightarrow \pi(\mathbf{x}) = P(c = 1|\mathbf{x}) = \frac{\exp\{g(\mathbf{x})\}}{1+\exp\{g(\mathbf{x})\}} =: \sigma(g(\mathbf{x})).$
- **Sigmoid** or **logistic** “squashing function”  $\sigma(z) = \frac{e^z}{1+e^z} = \frac{1}{1+e^{-z}}$   
turns linear predictions into probabilities



- Simple extension for  $K$  classes: the **softmax** function:  
$$P(c = k|\mathbf{x}) = \frac{\exp\{g_k(\mathbf{x})\}}{\sum_{m=1}^K \exp\{g_m(\mathbf{x})\}}.$$

# Logistic Regression (LOGREG)

- Assume that  $w_0$  is “absorbed” in  $\mathbf{w}$  using  $\mathbf{x} \leftarrow (1, \mathbf{x})$ . Estimate  $\mathbf{w}$  by maximizing the conditional likelihood

$$\hat{\mathbf{w}}_{DISCR} = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n (\pi(\mathbf{x}_i; \mathbf{w}))^{c_i} (1 - \pi(\mathbf{x}_i; \mathbf{w}))^{1-c_i},$$

or by maximizing the corresponding log likelihood  $l$ :

$$l(\mathbf{w}) = \sum_{i=1}^n [c_i \log \pi(\mathbf{x}_i; \mathbf{w}) + (1 - c_i) \log(1 - \pi(\mathbf{x}_i; \mathbf{w}))].$$

- The **score functions** are defined as the gradient of  $l$ :

$$\mathbf{s}(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} l(\mathbf{w}) = \sum_{i=1}^n \mathbf{x}_i (c_i - \pi_i).$$



# Logistic Regression (LOGREG)

- $\pi_i$  depends non-linearly on  $\mathbf{w}$ 
  - $\rightsquigarrow$  equation system  $\mathbf{s}(\mathbf{w}) = \mathbf{0}$  cannot be solved analytically
  - $\rightsquigarrow$  iterative techniques needed.
- Newton's method: Update  $\mathbf{w}$  at the  $r$ -th step as

$$\mathbf{w}^{(r+1)} = \mathbf{w}^{(r)} + \{H^{(r)}\}^{-1} \mathbf{s}^{(r)},$$

where  $H^{(r)}$  is the Hessian of  $l$ , evaluated at  $\mathbf{w}^{(r)}$ :

$$\begin{aligned} H^{(r)} &= \left( \frac{\partial^2 l}{\partial \mathbf{w} \partial \mathbf{w}^t} \right) \Big|_{\mathbf{w}=\mathbf{w}^{(r)}} \\ &= \sum_{i=1}^n \pi_i^{(r)} (1 - \pi_i^{(r)}) \mathbf{x}_i \mathbf{x}_i^t. \end{aligned}$$

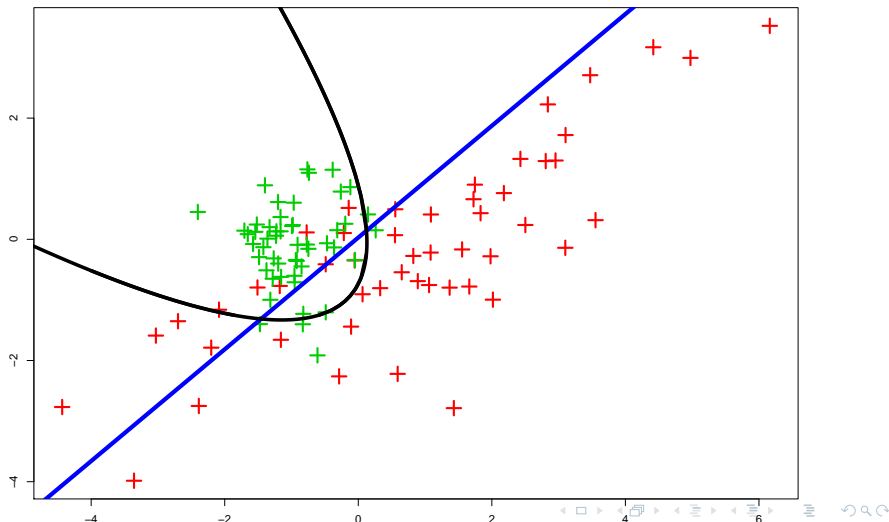
# Logistic Regression (LOGREG)

Newton updates  $\rightsquigarrow$  **Iterated Re-weighted Least Squares (IRLS):**

- The **Hessian**  $H^{(r)}$  is equal to  $(X^t W^{(r)} X)$ , with
$$W = \text{diag} \{ \pi_1(1 - \pi_1), \dots, \pi_n(1 - \pi_n) \}.$$
- **Score functions**:  $\mathbf{s}^{(r)} = X^t W^{(r)} \mathbf{e}^{(r)}$ , where  $\mathbf{e}$  is a vector with entries
$$e_j = (c_j - \pi_j) / W_{jj}.$$
- With  $\mathbf{q}^{(r)} := X \mathbf{w}^{(r)} + \mathbf{e}^{(r)}$ , the updates read
$$\begin{aligned} H^{(r)} \mathbf{w}^{(r+1)} &= H^{(r)} \mathbf{w}^{(r)} + \mathbf{s}^{(r)} \\ (X^t W^{(r)} X) \mathbf{w}^{(r+1)} &= X^t W^{(r)} \mathbf{q}^{(r)}. \end{aligned}$$
- These are the **normal equations of a LS problem**  $\|A\mathbf{w} - \mathbf{b}\|^2$  with input matrix  $A = (W^{(r)})^{1/2} X$  and r.h.s.  $\mathbf{b} = (W^{(r)})^{1/2} \mathbf{q}^{(r)}$ .
- The values  $W_{jj}$  are functions of  $\mathbf{w}$   $\rightsquigarrow$  iteration is needed.

# Logistic Regression (LOGREG)

Simple binary classification problem in  $\mathbb{R}^2$ . Solved with LOGREG using polynomial basis functions.



# Loss functions

- LOGREG maximizes log likelihood

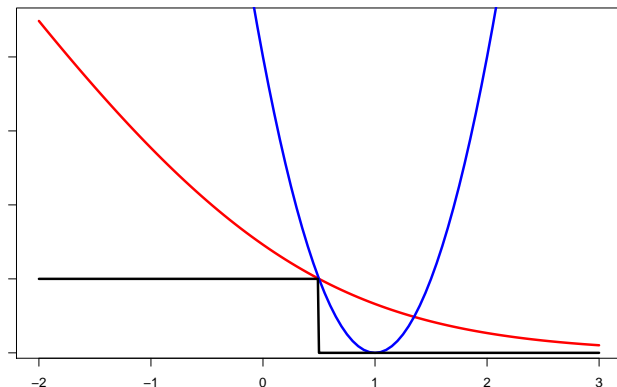
$$l(\mathbf{w}) = \sum_{i=1}^n [c_i \log \pi(\mathbf{x}_i; \mathbf{w}) + (1 - c_i) \log(1 - \pi(\mathbf{x}_i; \mathbf{w}))],$$

where  $z = \mathbf{w}^t \mathbf{x}$ ,  $\pi = \frac{1}{1+e^{-z}}$ ,  $1 - \pi = \frac{e^{-z}}{1+e^{-z}}$ .

- This is the same as minimizing

$$\begin{aligned} -l(\mathbf{w}) &= \sum_{i=1}^n [-c_i \log \pi - (1 - c_i) \log(1 - \pi)] \\ &=: \sum_{i=1}^n \text{Loss}(c_i, z_i). \end{aligned}$$

# Loss functions



Using  $\{0, 1\}$  encoding of the two classes, and approximating a target with  $c = +1$ . Black: 0/1-loss, red: logistic loss, blue: quadratic loss (LDA).

# LOGREG and Perceptron

- Gradient of negative log-likelihood:

$$\nabla_{\mathbf{w}^{(r)}} = \frac{\partial}{\partial \mathbf{w}} -l(\mathbf{w})|_{\mathbf{w}^{(r)}} = \sum_{i=1}^n \mathbf{x}_i(\pi_i - c_i).$$

- Gradient descent:  $\mathbf{w}^{(r+1)} = \mathbf{w}^{(r)} - \eta \nabla_{\mathbf{w}^{(r)}}$ .
- Assume stream of data  $\rightsquigarrow$  online update for new observation  $\mathbf{x}_i$ :  
 $\mathbf{w}^{(r+1)} = \mathbf{w}^{(r)} - \eta(\pi_i - c_i)\mathbf{x}_i$ , with  $\pi_i = P(c = 1|\mathbf{x}_i, \mathbf{w}^{(r)})$ .
- Now consider approximation: define **most probable label**  
 $\hat{c}_i = \arg \max_{c \in \{0,1\}} P(c|\mathbf{x}_i, \mathbf{w}^{(r)})$  and replace  $\pi_i$  with  $\hat{c}_i$ .
- If we predicted correctly, then  $\hat{c}_i = c_i \rightsquigarrow$  approximate gradient is zero  
 $\rightsquigarrow$  update has no effect.
- If  $\hat{c}_i = 0$  but  $c_i = 1$ :  $\mathbf{w}^{(r+1)} = \mathbf{w}^{(r)} - \eta(\hat{c}_i - c_i)\mathbf{x}_i = \mathbf{w}^{(r)} + \eta\mathbf{x}_i$ .
- Note that this is again the **perceptron algorithm**.
- Simple solution to most problems of the perceptron: use exact gradient instead of approximation based on most probable labels.