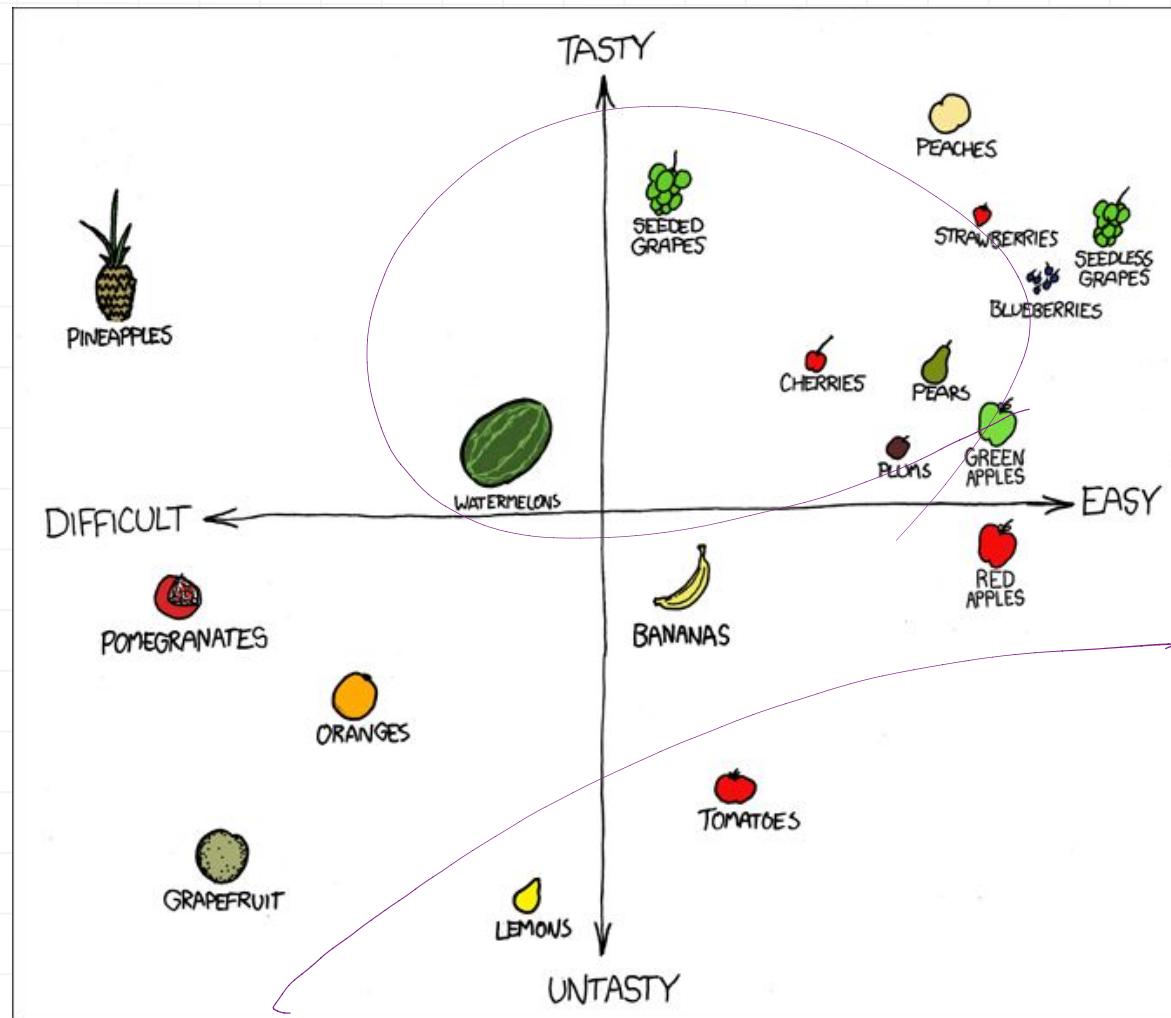


<https://xkcd.com/388/>



# Announcements

In person lecture: Wed 16 March, PHYS T (will try *simulcast* in Teams, with us luck!)

<https://studentvip.com.au/anu/main/maps/142757>

Quiz 1 next week, due Thu (releasing by Mon)

# Linear models for classification

Classification problems

A glimpse of decision theory (Sec 1.5)

Discriminant functions - why least squares doesn't work here

The perceptron algorithm

Probabilistic generative models - origin of the logistic function

Probabilistic discriminative models

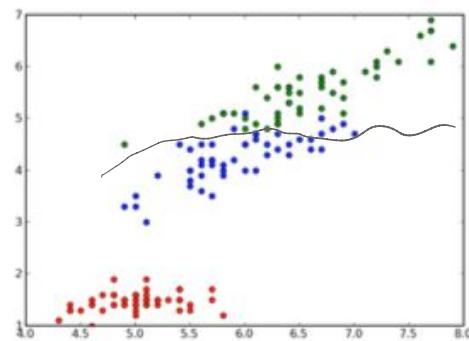
Logistic regression

Laplace approximation (Bayesian logistic regression) - later

Origin of the logistic function, logistic regression and how it connects to perceptron

# Classification

- Goal : Given input data  $\mathbf{x}$ , assign it to one of  $K$  discrete classes  $\mathcal{C}_k$  where  $k = 1, \dots, K$ .
- Divide the input space into different regions.
- Equivalently: map each point to a categorical label.



Length of petal [in cm] vs sepal [cm] for three types of flowers  
(Iris Setosa, Iris Versicolor, Iris Virginica).

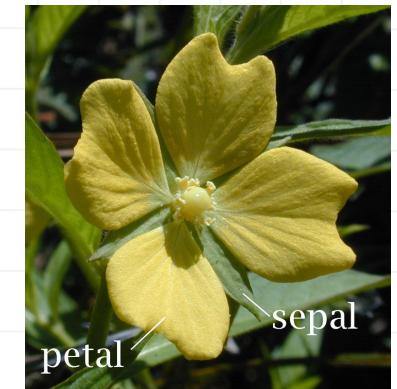


Image from wikipedia

# Representing classes

binary  $t \in \{0, 1\}$

Meeting in person/remote, pass/fail, benign/malignant, should be given credit (Y/N), ...

one-hot  $\underbrace{\mathbf{t} = (0, 1, 0, 0, 0)^T}_{\sum_k t_k = 1}$  (4.1)

Also denote  $t_k \longrightarrow \mathcal{C}_k$

Iris flowers, object classes in photos, natural language, ...

Other representations abound, e.g. structured (time series, sequences – PRML Chapter 13, graphs),  
Sec 4.1.5 (an output representation that connects Fisher Discriminant to least squares),  
Sec 4.1.7 perceptrons, SVMs  $\{-1, +1\}$

# A primer on decision theory (Sec 1.5)

Inference: compute  $p(\mathbf{x}, \mathbf{t})$ , e.g. from a set of training data.

Decision theory: take a specific action based on our understanding of the values  $\mathbf{t}$  is likely to take.

$$\underbrace{p(\mathcal{C}_k | \mathbf{x})}_{\text{---}} = \frac{p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)}{p(\mathbf{x})}. \quad (1.77)$$



(end of page 38) "We shall see that the decision stage is generally very simple, even trivial, once we have solved the inference problem."

- many counter-examples! e.g. multiple medical tests, test + quarantine for covid
- frontiers of theoretical and practical ML, e.g. active/online learning, ML and economics

# Minimise error

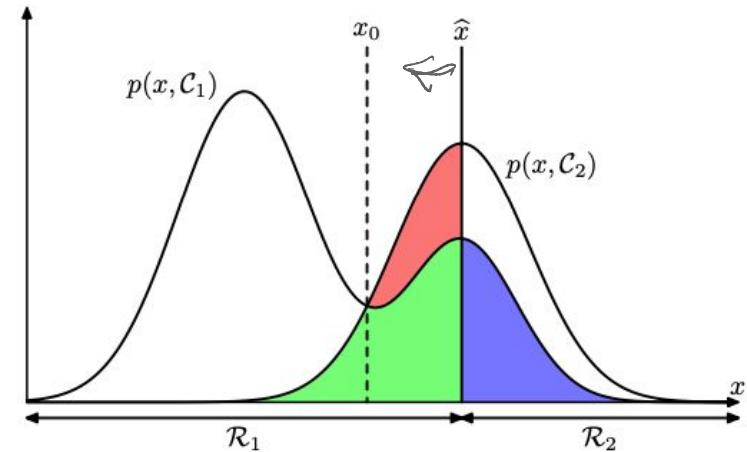
$$p(\text{correct}) = \sum_{k=1}^K p(\mathbf{x} \in \mathcal{R}_k, \mathcal{C}_k)$$

$$= \sum_{k=1}^K \int_{\mathcal{R}_k} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}$$

*red+green*

$$p(\text{mistake}) = p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1)$$

$$= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}. \quad (1.78)$$



**Figure 1.24** Schematic illustration of the joint probabilities  $p(x, \mathcal{C}_k)$  for each of two classes plotted against  $x$ , together with the decision boundary  $x = \hat{x}$ . Values of  $x \geq \hat{x}$  are classified as class  $\mathcal{C}_2$  and hence belong to decision region  $\mathcal{R}_2$ , whereas points  $x < \hat{x}$  are classified as  $\mathcal{C}_1$  and belong to  $\mathcal{R}_1$ . Errors arise from the blue, green, and red regions, so that for  $x < \hat{x}$  the errors are due to points from class  $\mathcal{C}_2$  being misclassified as  $\mathcal{C}_1$  (represented by the sum of the red and green regions), and conversely for points in the region  $x \geq \hat{x}$  the errors are due to points from class  $\mathcal{C}_1$  being misclassified as  $\mathcal{C}_2$  (represented by the blue region). As we vary the location  $\hat{x}$  of the decision boundary, the combined areas of the blue and green regions remains constant, whereas the size of the red region varies. The optimal choice for  $\hat{x}$  is where the curves for  $p(x, \mathcal{C}_1)$  and  $p(x, \mathcal{C}_2)$  cross, corresponding to  $\hat{x} = x_0$ , because in this case the red region disappears. This is equivalent to the minimum misclassification rate decision rule, which assigns each value of  $x$  to the class having the higher posterior probability  $p(\mathcal{C}_k|x)$ .

# What are the potential problems for minimising classification error?

$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}. \end{aligned} \quad (1.78)$$

Example scenario: medical diagnosis

# Minimise loss, balanced metric, having a reject option

**Figure 1.25** An example of a loss matrix with elements  $L_{kj}$  for the cancer treatment problem. The rows correspond to the true class, whereas the columns correspond to the assignment of class made by our decision criterion.

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}. \quad (1.80)$$

$$\text{Decision: } \operatorname{argmin}_j \sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x}) \quad (1.81)$$

		cancer	normal
True	cancer	( 0      1000 )	
	normal	( 1      0 )	

		Predicted 1-Cancer	Predicted 0-No Cancer
Actual	1-Cancer	90 [TP]	4 [FN]
	0-No Cancer	1 [FP]	5 [TN]

$$\frac{TP}{P} = \frac{90}{94}$$

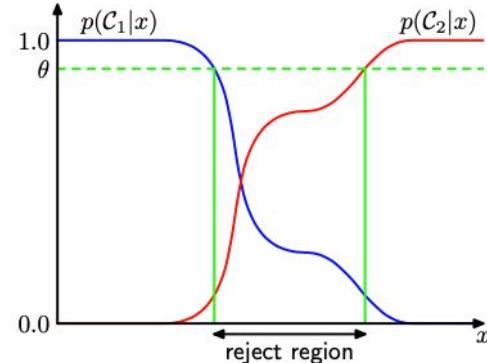
→ Cancer Records [90+4=94]

→ Non Cancer Records [5+1=6]

<https://towardsdatascience.com/confusion-matrix-clearly-explained-fee63614dc7>

$$\frac{TN}{N} = \frac{5}{6}$$

**Figure 1.26** Illustration of the reject option. Inputs  $x$  such that the larger of the two posterior probabilities is less than or equal to some threshold  $\theta$  will be rejected.



Brodersen, K. H.; Ong, C. S.; Stephan, K. E.; Buhmann, J. M., 2010.

The balanced accuracy and its posterior distribution.

International Conference on Pattern Recognition.

$$\frac{1}{2} (\overline{TP/P} + \overline{TN/N})$$

(many other metrics also address this problem)

# Three types of models for decision problems

①

- Find a **discriminant function**  $f(\mathbf{x})$  which maps each input directly onto a class label.

②

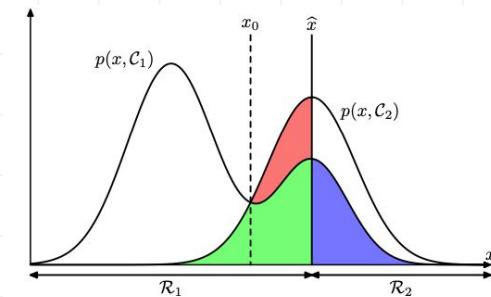
- Discriminative Models

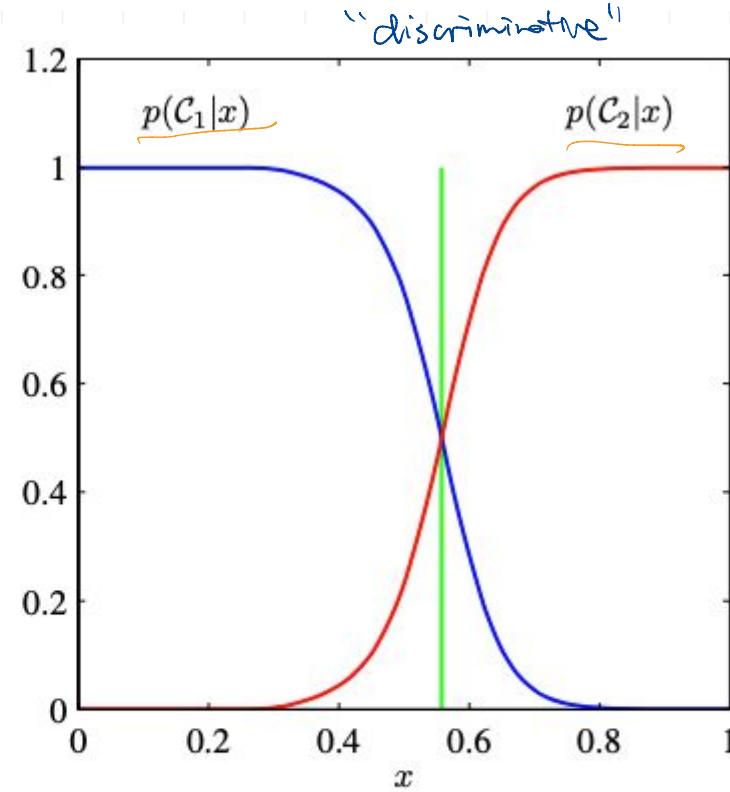
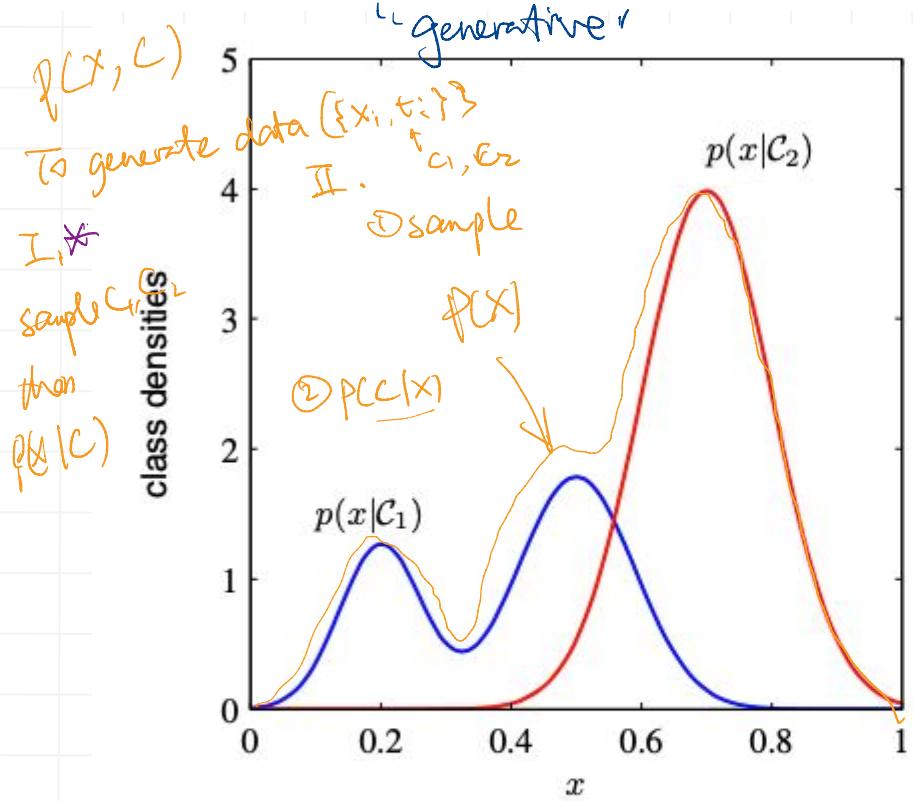
- ➊ Solve the inference problem of determining the posterior class probabilities  $p(C_k | \mathbf{x})$ .
- ➋ Use decision theory to assign each new  $\mathbf{x}$  to one of the classes.

③

- Generative Models

- ➊ Solve the inference problem of determining the class-conditional probabilities  $p(\mathbf{x} | C_k)$ .
- ➋ Also, infer the prior class probabilities  $p(C_k)$ .
- ➌ Use Bayes' theorem to find the posterior  $p(C_k | \mathbf{x})$ .
- ➍ Alternatively, model the joint distribution  $p(\mathbf{x}, C_k)$  directly.
- ➎ Use decision theory to assign each new  $\mathbf{x}$  to one of the classes.

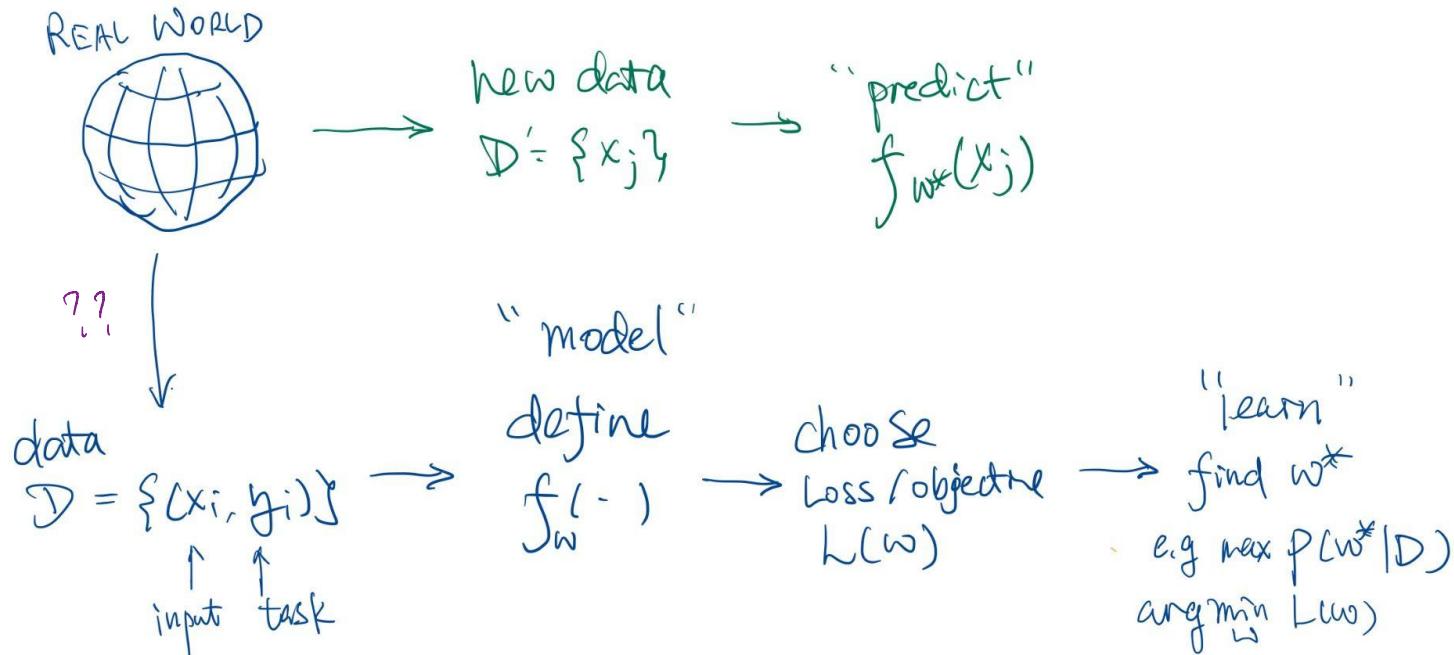




**Figure 1.27** Example of the class-conditional densities for two classes having a single input variable  $x$  (left plot) together with the corresponding posterior probabilities (right plot). Note that the left-hand mode of the class-conditional density  $p(x|C_1)$ , shown in blue on the left plot, has no effect on the posterior probabilities. The vertical green line in the right plot shows the decision boundary in  $x$  that gives the minimum misclassification rate.

high level goal of "learning": the non-trivial dependency between  $X, C$

# A modular view of (supervised) ML



# What we covered so far

Classification problems

A glimpse of decision theory (Sec 1.5)

Discriminant functions - why least squares doesn't work here

The perceptron algorithm

Probabilistic generative models - origin of the logistic function

Probabilistic discriminative models

Logistic regression

(Laplace approximation, Bayesian logistic regression)

# Discriminant function

## Definition

A **discriminant** is a function that maps from an input vector  $\mathbf{x}$  to one of  $K$  classes, denoted by  $\mathcal{C}_k$ .

- Consider first two classes ( $K = 2$ ).
- Construct a linear function of the inputs  $\mathbf{x}$

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (4.4)$$

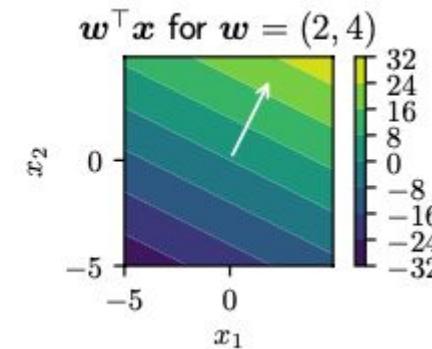
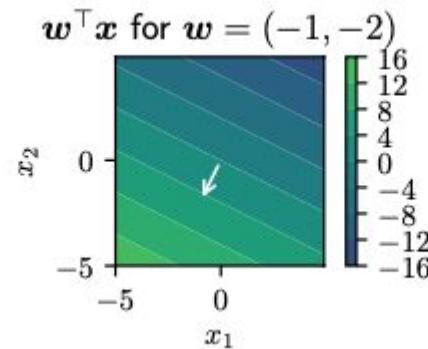
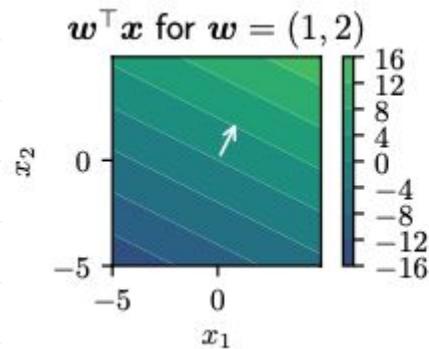
such that  $\mathbf{x}$  being assigned to class  $\mathcal{C}_1$  if  $y(\mathbf{x}) \geq 0$ , and to class  $\mathcal{C}_2$  otherwise.

- **weight vector**  $\mathbf{w}$
- **bias**  $w_0$  ( sometimes  $-w_0$  called **threshold** )

# Linear discriminant functions

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (4.4)$$

- Decision boundary  $y(\mathbf{x})=0$  -- is a hyperplane of D-1 dimensions (in a D-dimensional input space).
  - $\mathbf{w}$  is orthogonal to any vector lying in the decision surface.
- 
- Contour lines of  $y(\mathbf{x})$  are hyperplanes too
  - $y(\mathbf{x})$  is the signed distance of point  $\mathbf{x}$  from the hyper-plane (see next page)



# Discriminant for two classes

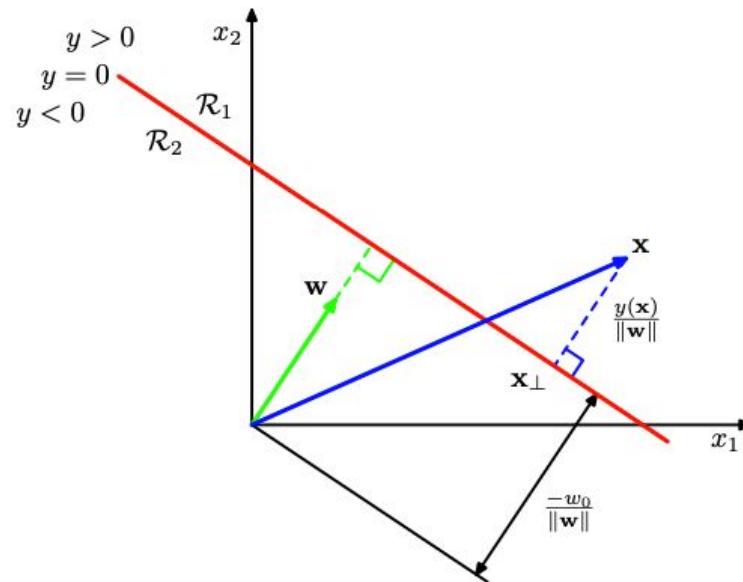
$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

(4.4)

$y(\mathbf{x})$  gives a **signed** measure of the perpendicular distance  $r$  from the decision surface to  $\mathbf{x}$ , that is  $r = y(\mathbf{x})/\|\mathbf{w}\|$ .

$$y(\mathbf{x}) = \mathbf{w}^T \overbrace{\left( \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right)}^{\mathbf{x}} + w_0 = r \frac{\mathbf{w}^\top \mathbf{w}}{\|\mathbf{w}\|} + \overbrace{\mathbf{w}^\top \mathbf{x}_\perp + w_0}^0 = r \|\mathbf{w}\|$$

**Figure 4.1** Illustration of the geometry of a linear discriminant function in two dimensions. The decision surface, shown in red, is perpendicular to  $\mathbf{w}$ , and its displacement from the origin is controlled by the bias parameter  $w_0$ . Also, the signed orthogonal distance of a general point  $\mathbf{x}$  from the decision surface is given by  $y(\mathbf{x})/\|\mathbf{w}\|$ .

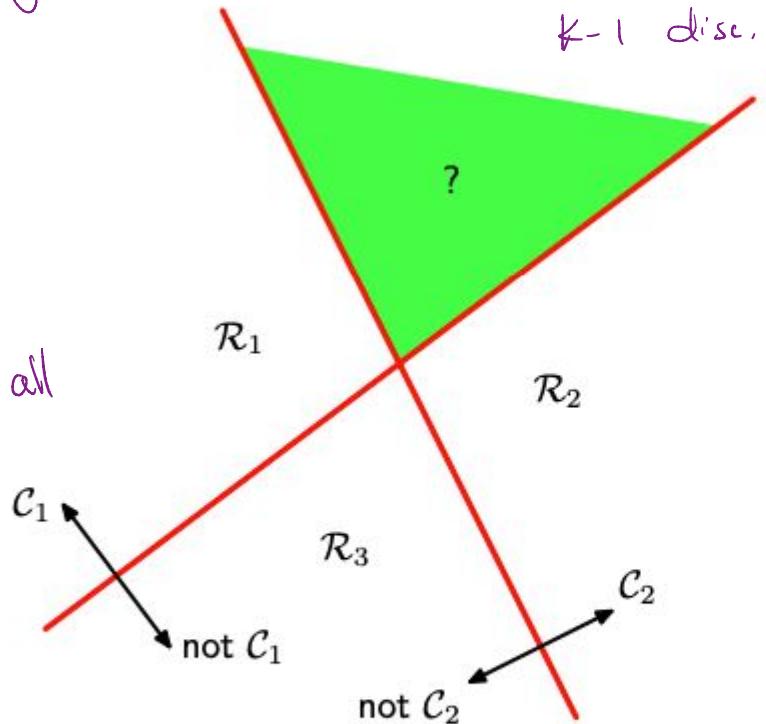


$$y(x) = w^T x$$

*k* classes

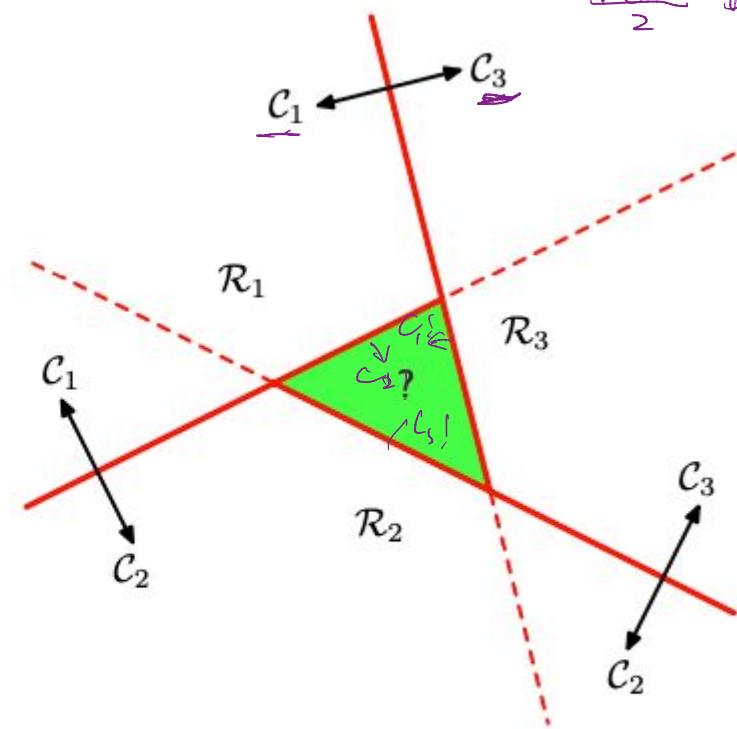
*k*-1 disc. fn.

One vs all



One-vs-one

$$\frac{k(k-1)}{2} \# \text{ of } y_s$$



**Figure 4.2** Attempting to construct a  $K$  class discriminant from a set of two class discriminants leads to ambiguous regions, shown in green. On the left is an example involving the use of two discriminants designed to distinguish points in class  $C_k$  from points not in class  $C_k$ . On the right is an example involving three discriminant functions each of which is used to separate a pair of classes  $C_k$  and  $C_j$ .

# Discriminant function for K>=2 classes

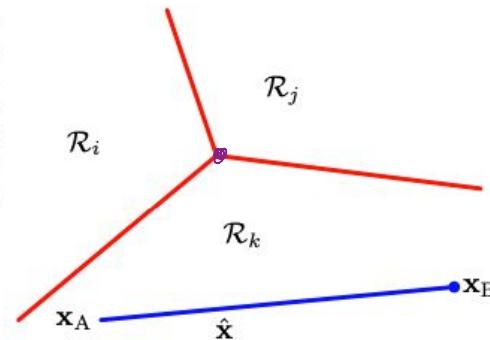
Learn a series of functions  $k=1, \dots, K$ , assigning a point  $\mathbf{x}$  to class  $\mathcal{C}_k$  if  $y_k(\mathbf{x}) > y_j(\mathbf{x})$  for all  $j \neq k$ .

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (4.9)$$

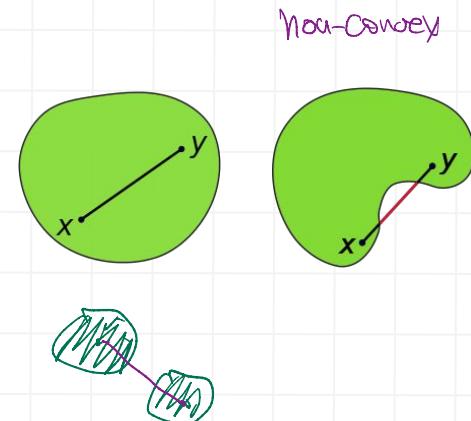
$$k^* \leftarrow \operatorname{argmax}_k y_k(\mathbf{x})$$

Claim: The decision regions of such a discriminant are always singly connected and convex.

**Figure 4.3** Illustration of the decision regions for a multiclass linear discriminant, with the decision boundaries shown in red. If two points  $\mathbf{x}_A$  and  $\mathbf{x}_B$  both lie inside the same decision region  $\mathcal{R}_k$ , then any point  $\hat{\mathbf{x}}$  that lies on the line connecting these two points must also lie in  $\mathcal{R}_k$ , and hence the decision region must be singly connected and convex.



$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B).$$



Non-Convex

# Can we use least squares?

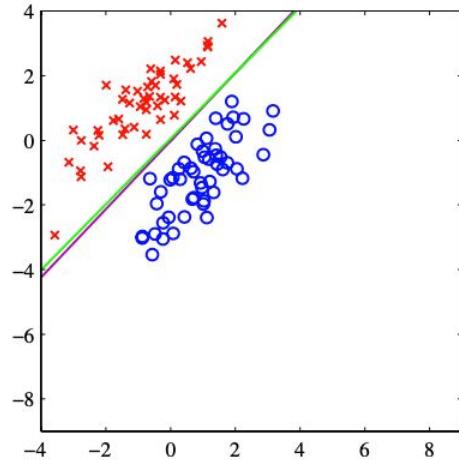
$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (4.13)$$

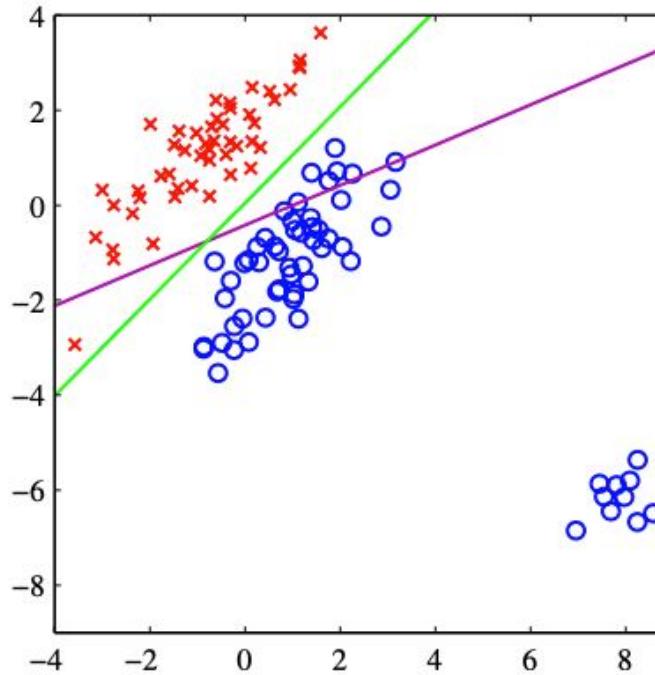
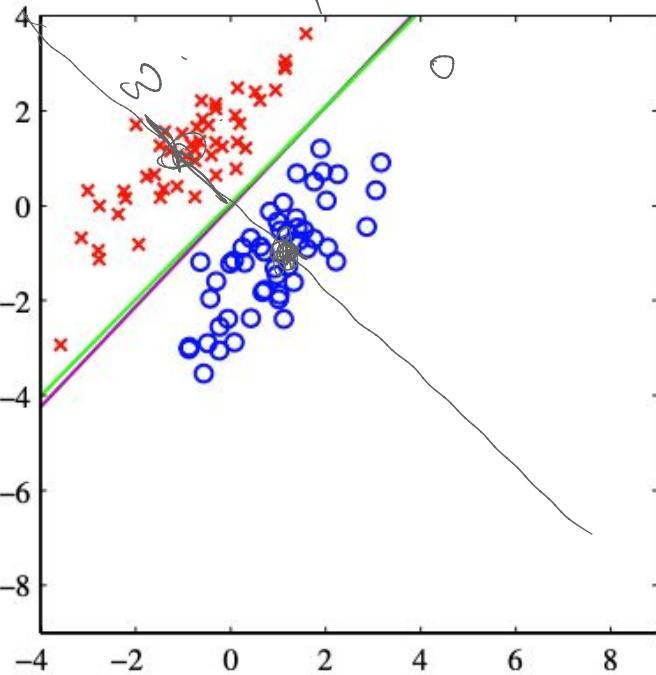
$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}} \quad (4.14)$$

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T})^T (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T}) \right\}. \quad (4.15)$$

*t<sub>u</sub> ∈ {0, 1}*

$$\widetilde{\mathbf{W}} = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{T} = \widetilde{\mathbf{X}}^\dagger \mathbf{T} \quad (4.16)$$





**Figure 4.4** The left plot shows data from two classes, denoted by red crosses and blue circles, together with the decision boundary found by least squares (magenta curve) and also by the logistic regression model (green curve), which is discussed later in Section 4.3.2. The right-hand plot shows the corresponding results obtained when extra data points are added at the bottom left of the diagram, showing that least squares is highly sensitive to outliers, unlike logistic regression.

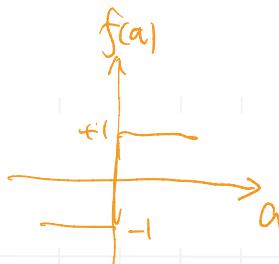
Mismatch between model requirement and model specification.

# The (historically significant) perceptron [Rosenblatt 1962]

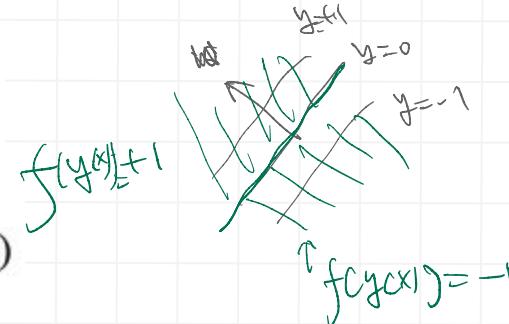
$$y(\mathbf{x}) = f(\underline{\mathbf{w}^T \phi(\mathbf{x})})$$

(4.52)

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0. \end{cases}$$



(4.53)



$$t \in \{-1, +1\}$$

$\min$ ,

$$E_P(\mathbf{w}) = - \sum_{n \in M} \mathbf{w}^T \underline{\phi_n t_n}$$

(4.54)

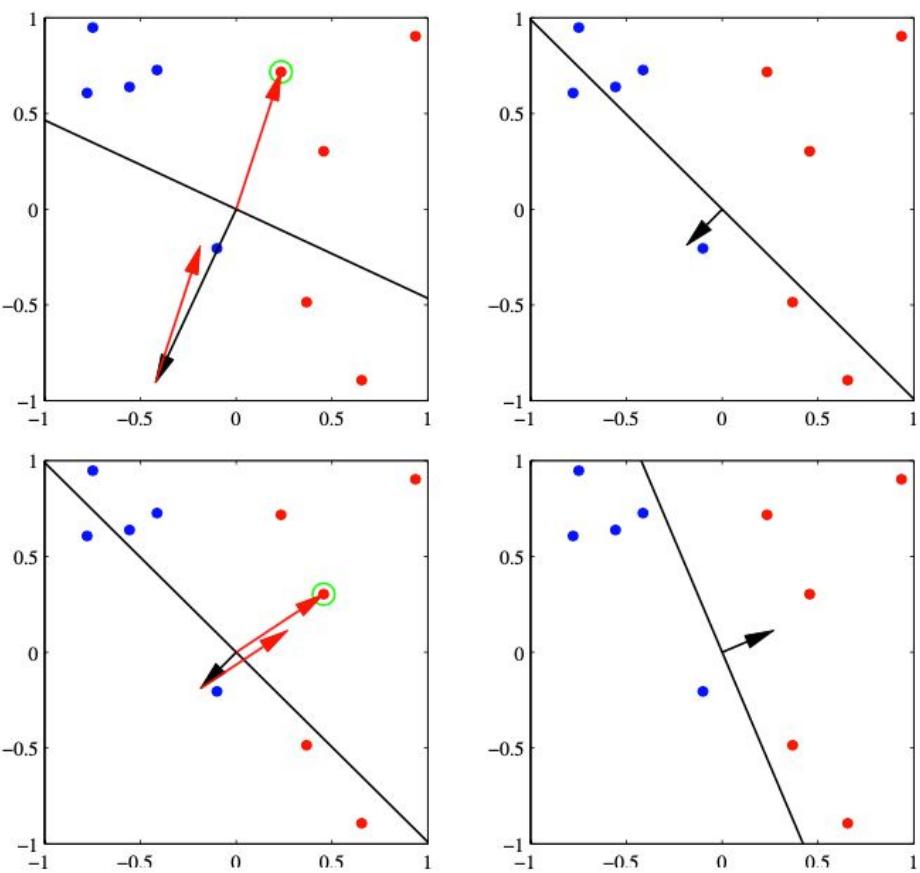
learning rate  
↓

Stochastic gradient descent

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \underline{\phi_n t_n}$$

(4.55)

"The *perceptron convergence theorem* states that if there exists an exact solution (in other words, if the training data set is linearly separable), then the perceptron learning algorithm is guaranteed to find an exact solution in a finite number of steps ... however, until convergence is achieved, we will not be able to distinguish between a nonseparable problem and one that is simply slow to converge."



**Figure 4.7** Illustration of the convergence of the perceptron learning algorithm, showing data points from two classes (red and blue) in a two-dimensional feature space ( $\phi_1, \phi_2$ ). The top left plot shows the initial parameter vector  $w$  shown as a black arrow together with the corresponding decision boundary (black line), in which the arrow points towards the decision region which classified as belonging to the red class. The data point circled in green is misclassified and so its feature vector is added to the current weight vector, giving the new decision boundary shown in the top right plot. The bottom left plot shows the next misclassified point to be considered, indicated by the green circle, and its feature vector is again added to the weight vector giving the decision boundary shown in the bottom right for which all data points are correctly classified.



**Figure 4.8** Illustration of the Mark 1 perceptron hardware. The photograph on the left shows how the inputs were obtained using a simple camera system in which an input scene, in this case a printed character, was illuminated by powerful lights, and an image focussed onto a  $20 \times 20$  array of cadmium sulphide photocells, giving a primitive 400 pixel image. The perceptron also had a patch board, shown in the middle photograph, which allowed different configurations of input features to be tried. Often these were wired up at random to demonstrate the ability of the perceptron to learn without the need for precise wiring, in contrast to a modern digital computer. The photograph on the right shows one of the racks of adaptive weights. Each weight was implemented using a rotary variable resistor, also called a potentiometer, driven by an electric motor thereby allowing the value of the weight to be adjusted automatically by the learning algorithm.

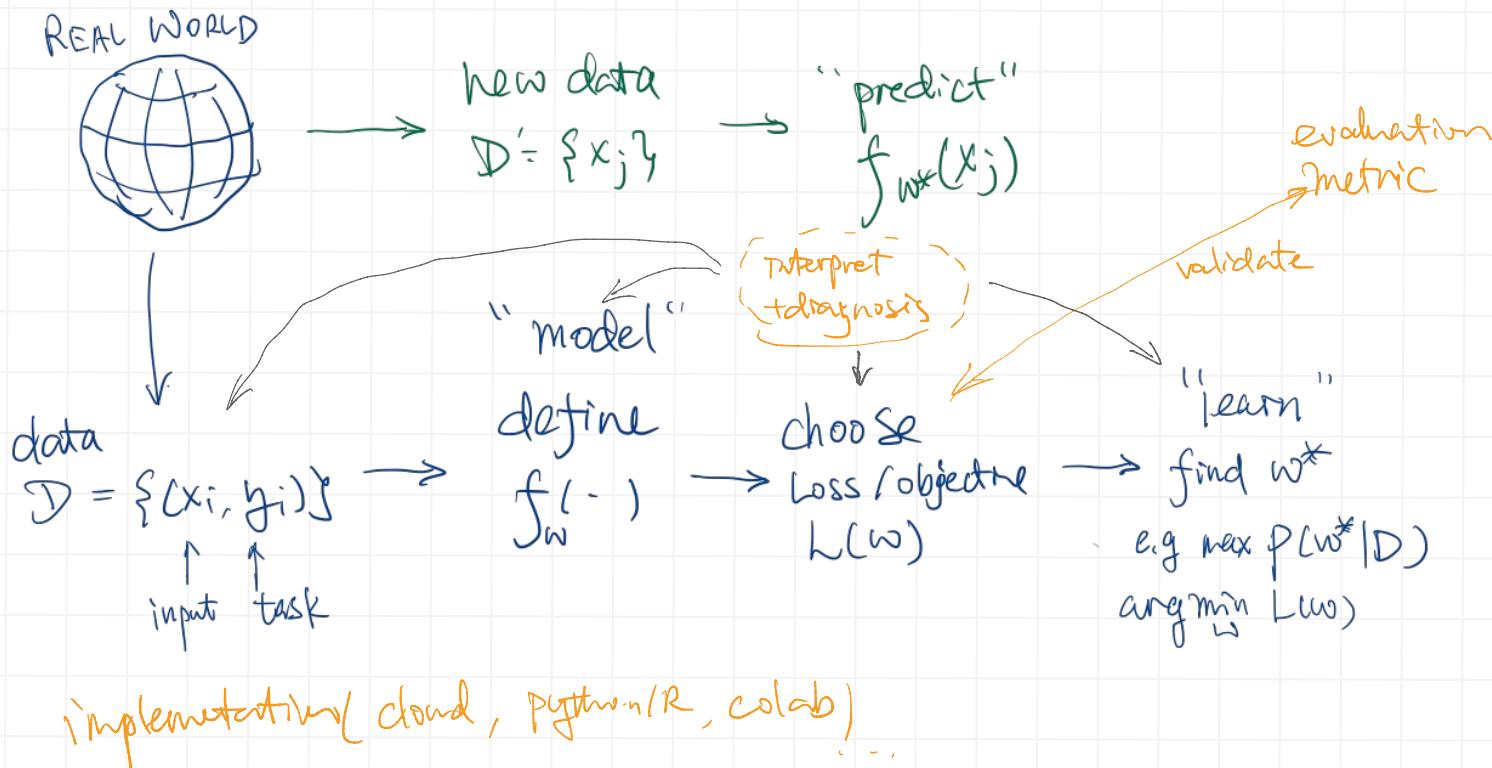
Frank Rosenblatt  
1928–1969



Rosenblatt's perceptron played an important role in the history of machine learning. Initially, Rosenblatt simulated the perceptron on an IBM 704 computer at Cornell in 1957, but by the early 1960s he had built special-purpose hardware that provided a direct, parallel implementation of perceptron learning. Many of his ideas were encapsulated in "Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms" published in 1962. Rosenblatt's work was criticized by Marvin Minsky, whose objections were published in the book "Perceptrons", co-authored with

Papert. This book was widely misinterpreted at the time as showing that neural networks were fatally flawed and could only learn solutions for linearly separable problems. In fact, it only proved such limitations in the case of single-layer networks such as the perceptron and merely conjectured (incorrectly) that they applied to more general network models. Unfortunately, however, this book contributed to the substantial decline in research funding for neural computing, a situation that was not reversed until the mid-1980s. Today, there are many hundreds, if not thousands, of applications of neural networks in widespread use, with examples in areas such as handwriting recognition and information retrieval being used routinely by millions of people.

# Another look at the modular view of supervised ML



# What we covered so far

Classification problems

A glimpse of decision theory (Sec 1.5)

Discriminant functions - why least squares doesn't work here

The perceptron algorithm

Probabilistic generative models - origin of the logistic function

Probabilistic discriminative models

Logistic regression

(Laplace approximation, Bayesian logistic regression)

## Bayes theorem, again

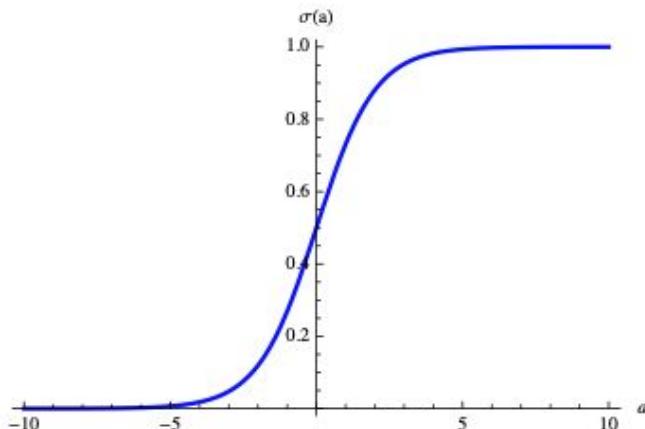
$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \Leftarrow p(\mathbf{x})$$

$$a(\mathbf{x}) = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} = \ln \frac{p(\mathbf{x}, \mathcal{C}_1)}{p(\mathbf{x}, \mathcal{C}_2)}$$

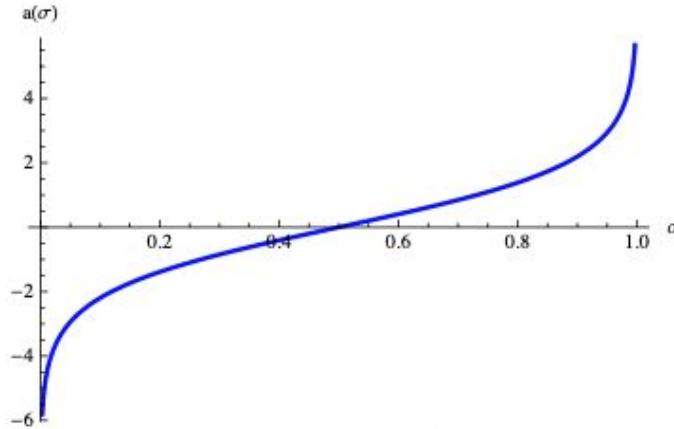
$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

logistic sigmoid function

$$= \frac{1}{1 + \exp(-a)} = \sigma(a) \quad (4.57)$$



**Sigmoid**  $\sigma(a) = \frac{1}{1+\exp(-a)}$



**Logit**  $a(\sigma) = \ln\left(\frac{\sigma}{1-\sigma}\right)$

$\ln [p(\mathcal{C}_1|\mathbf{x})/p(\mathcal{C}_2|\mathbf{x})]$   
a.k.a log odds

**Symmetry:**  $\sigma(-a) = 1 - \sigma(a)$

**Derivative:**  $\frac{d}{da} \sigma(a) = \sigma(a) \sigma(-a) = \sigma(a) (1 - \sigma(a))$

# Probabilistic Generative Models - Multiclass

- The **normalised exponential** is given by

$$p(\mathcal{C}_k \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \mathcal{C}_k) p(\mathcal{C}_k)}{\sum_j p(\mathbf{x} \mid \mathcal{C}_j) p(\mathcal{C}_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

where

$$a_k = \ln(p(\mathbf{x} \mid \mathcal{C}_k) p(\mathcal{C}_k)).$$

- Usually called the **softmax function** as it is a smoothed version of the **arg max** function, in particular:

$$a_k \gg a_j \quad \forall j \neq k \Rightarrow \left( p(\mathcal{C}_k \mid \mathbf{x}) \approx 1 \wedge p(\mathcal{C}_j \mid \mathbf{x}) \approx 0 \right)$$

- So, **softargmax** is a more descriptive though less common name.

# Class-conditional distribution being Gaussians

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}. \quad (4.64)$$

cancel      cancel  
 ↓

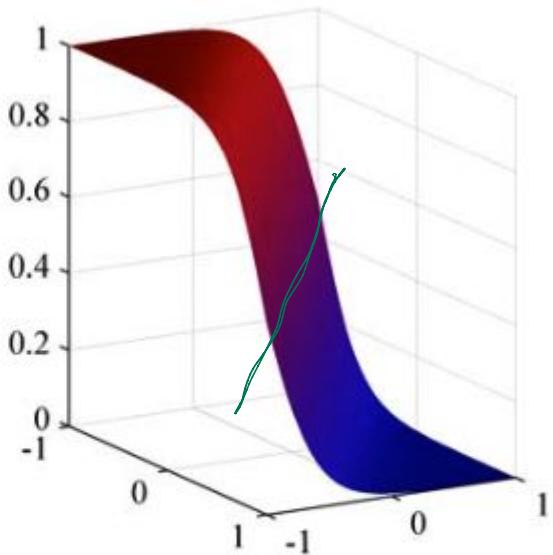
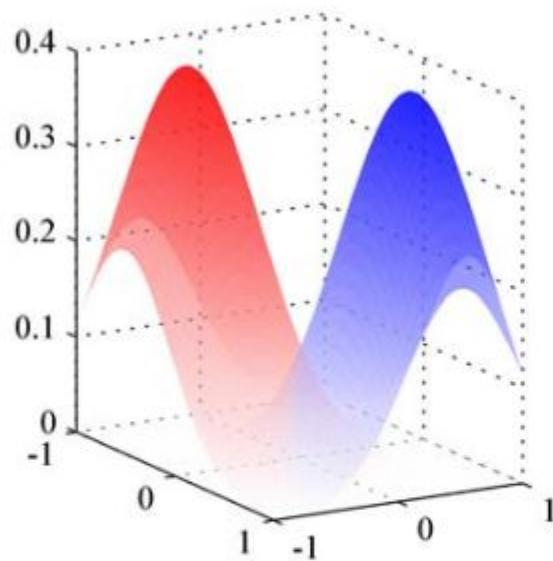
two classes  $\mathcal{C}_1, \mathcal{C}_2$   
 class-means  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2$   
 share cov.  $\Sigma$   
 $\underline{x^T \Sigma^{-1} x} \rightarrow \text{cancel}$

$$a(\mathbf{x}) = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) \quad (4.65)$$

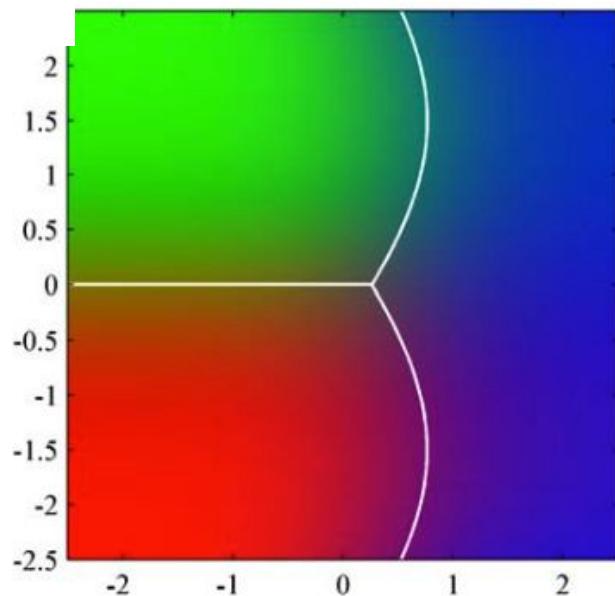
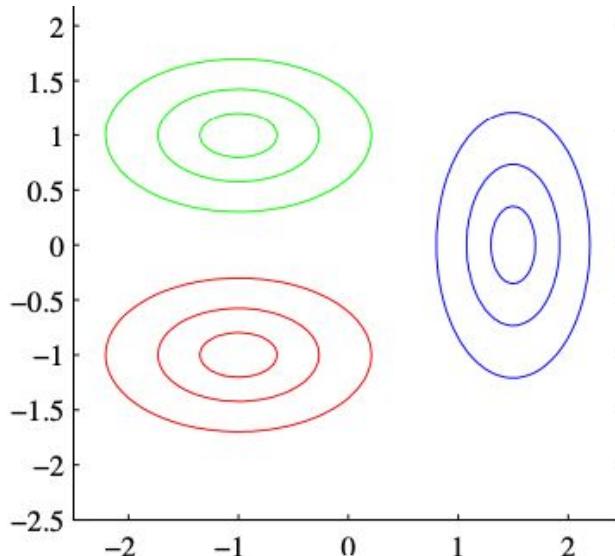
$$\mathbf{w} = \underbrace{\Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}_{\Downarrow} \quad (4.66)$$

$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}. \quad (4.67)$$



**Figure 4.10** The left-hand plot shows the class-conditional densities for two classes, denoted red and blue. On the right is the corresponding posterior probability  $p(C_1|x)$ , which is given by a logistic sigmoid of a linear function of  $x$ . The surface in the right-hand plot is coloured using a proportion of red ink given by  $\underline{p(C_1|x)}$  and a proportion of blue ink given by  $p(C_2|x) = 1 - p(C_1|x)$ .

- If the class-conditional distributions have different covariances, the quadratic terms  $-\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x}$  do not cancel out.
- We get a quadratic discriminant.



**Figure 4.11** The left-hand plot shows the class-conditional densities for three classes each having a Gaussian distribution, coloured red, green, and blue, in which the red and green classes have the same covariance matrix. The right-hand plot shows the corresponding posterior probabilities, in which the RGB colour vector represents the posterior probabilities for the respective three classes. The decision boundaries are also shown. Notice that the boundary between the red and green classes, which have the same covariance matrix, is linear, whereas those between the other pairs of classes are quadratic.

# Maximum likelihood estimates

$$\pi^{P(C_1)} \quad P(C_2) = 1 - \pi \\ \{ \pi, \mu_1, \mu_2, \Sigma \}$$

Given the functional form of the class-conditional densities  $p(\mathbf{x} | C_k)$ , how can we determine the parameters  $\mu$  and  $\Sigma$  and the class prior?

PLA

$$p(\mathbf{x}_n, C_1) = p(C_1)p(\mathbf{x}_n | C_1) = \pi \mathcal{N}(\mathbf{x}_n | \mu_1, \Sigma)$$

$$p(\mathbf{x}_n, C_2) = p(C_2)p(\mathbf{x}_n | C_2) = (1 - \pi) \mathcal{N}(\mathbf{x}_n | \mu_2, \Sigma)$$

$$\{C_1, C_2\} \{0, 1\}$$

$$p(\mathbf{t}, \mathbf{X} | \pi, \mu_1, \mu_2, \Sigma)$$

$$= \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \mu_1, \Sigma)]^{t_n} \times [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \mu_2, \Sigma)]^{1-t_n}$$

$t_n = C_1 \text{ OR } t_n = C_2$

$$L = \log \left\{ \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} \times [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n} \right\}$$

$\sum_{n=1}^N t_n [\log \pi + \log \mathcal{N}(\cdot | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})] + (1-t_n) [\log (1-\pi) + \log \mathcal{N}(\cdot | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]$

The term depending on  $\pi$  is

$$\frac{\partial L}{\partial \pi} = \frac{\partial}{\partial \pi} \sum_{n=1}^N (t_n \ln \pi + (1-t_n) \ln(1-\pi))$$

$$\Rightarrow \sum_{n=1}^N t_n \frac{1}{\pi} + (1-t_n) \frac{-1}{1-\pi}$$

The maximum is at:

$$\pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$$

# MLE solution for means and covariance

{  $\pi$ ,  $\mu_1$ ,  $\mu_2$ ,  $\Sigma$  }  
↑ ↑ ↑ ↑  
I D D OUT!

$$\boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n \quad (4.75)$$

$$\boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n \quad (4.76)$$

$$\mathbf{S} = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2 \quad (4.78)$$

$$\mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T \quad (4.79)$$

$$\mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T. \quad (4.80)$$

O(D<sup>2</sup>) parameters!

Left as reading:  
Sec 4.2.3 naive bayes  
Sec 4.2.4 exponential family

## Discrete features - Naive Bayes

- Assume the input space consists of discrete features, in the simplest case  $x_i \in \{0, 1\}$ .
- For a  $D$ -dimensional input space, a general distribution would be represented by a table with  $2^D$  entries.
- Together with the normalisation constraint, this are  $2^D - 1$  independent variables.
- Grows exponentially with the number of features.
- The **Naïve Bayes** assumption is that, given the class  $\mathcal{C}_k$ , the features are independent of each other:

$$\begin{aligned} p(\mathbf{x} \mid \mathcal{C}_k) &= \prod_{i=1}^D p(x_i \mid \mathcal{C}) \\ &= \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i} \end{aligned}$$

# Naive Bayes classifier

$$p(\mathbf{x}|\mathcal{C}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i} \quad (4.81)$$

$$p(\mathcal{C}_k | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x} | \mathcal{C}_j)p(\mathcal{C}_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$a_k(\mathbf{x}) = \sum_{i=1}^D \{x_i \ln \mu_{ki} + (1 - x_i) \ln(1 - \mu_{ki})\} + \ln p(\mathcal{C}_k) \quad (4.82)$$

Linear functions of input  $\mathbf{x}$

# What we covered so far

Classification problems

A glimpse of decision theory (Sec 1.5)

Discriminant functions - why least squares doesn't work here  
(Fisher discriminant)

The perceptron algorithm

Probabilistic generative models - origin of the logistic function

Probabilistic discriminative models

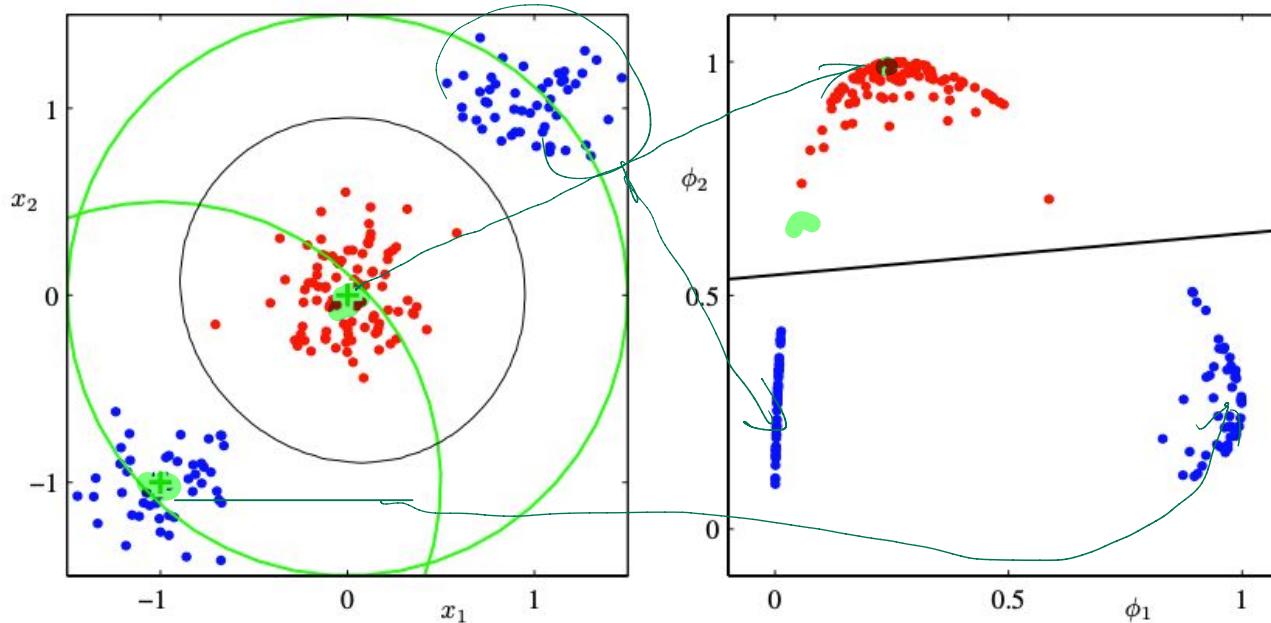
Logistic regression

(Laplace approximation, Bayesian logistic regression)

# Probabilistic Discriminative Models

- **Discriminative** training: learn only to discriminate between the classes.
- Maximise a likelihood function defined through the conditional distribution  $p(C_k | \mathbf{x})$  directly.
- Typically fewer parameters to be determined.
- As we learn the posterior  $p(C_k | \mathbf{x})$  directly, prediction may be better than with a generative model where the class-conditional density assumptions  $p(\mathbf{x} | C_k)$  poorly approximate the true distributions.
- But: discriminative models can not create synthetic data, as  $p(\mathbf{x})$  is not modelled.
- As an aside: *certain theoretical analyses show that generative models converge faster to their — albeit worse — asymptotic classification performance and are superior in some regimes.*

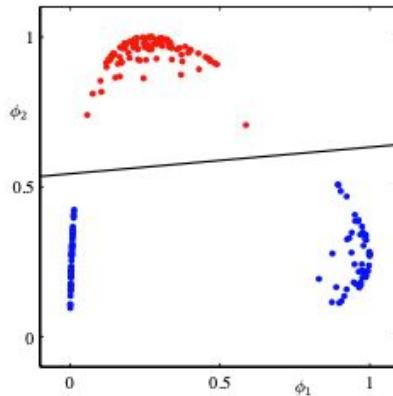
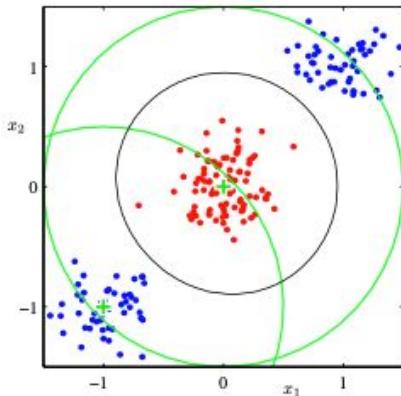
# Fixed basis functions



**Figure 4.12** Illustration of the role of nonlinear basis functions in linear classification models. The left plot shows the original input space  $(x_1, x_2)$  together with data points from two classes labelled red and blue. Two 'Gaussian' basis functions  $\phi_1(x)$  and  $\phi_2(x)$  are defined in this space with centres shown by the green crosses and with contours shown by the green circles. The right-hand plot shows the corresponding feature space  $(\phi_1, \phi_2)$  together with the linear decision boundary obtained given by a logistic regression model of the form discussed in Section 4.3.2. This corresponds to a nonlinear decision boundary in the original input space, shown by the black curve in the left-hand plot.

# Original input vs feature space

- Linear decision boundaries in the feature space generally correspond to nonlinear boundaries in the input space.
- Classes which are NOT linearly separable in the input space may become linearly separable in the feature space:



- If classes overlap in input space, they will also overlap in feature space — nonlinear features  $\phi(\mathbf{x})$  can not remove the overlap; but they may increase it.

# Logistic regression

$$p(C_1|\phi) = y(\phi) = \sigma(w^T \phi)$$

(4.87)

→ ignore  $P(C_i)$   
ignore shape of  
 $P(x|C_i)$

The name came from statistics. It's a model for classification rather than regression!

Compare to generative model of one Gaussian per class with shared variance

$$p(C_1|x) = \sigma(w^T x + w_0)$$

(4.65)

model

$P(x|C_1)$  ~ Gaussian

$P(C_1)$  ~ number -

Same expression for posterior estimates, but linear number of parameters rather than quadratic.

D in millions / billions . sparse .  
↑

# Maximum likelihood for logistic regression

- Determine the parameter via maximum likelihood for data  $(\phi_n, t_n), n = 1, \dots, N$ , where  $\phi_n = \phi(\mathbf{x}_n)$ . The class membership is coded as  $t_n \in \{0, 1\}$ .
- Likelihood function

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$

$y_n = p(C_1 | \phi_n)$

- Error function : negative log likelihood resulting in the **cross-entropy** error function

$$y_n = p(C_1 | \phi_n) = \sigma(\mathbf{w}^T \phi_n)$$

$$\min_{\mathbf{w}} E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

↳ Assumes every data point counts as the same .

Gradient of Logistic regression

$$E(\omega) = -\ln p(t|\hat{\omega}) = - \sum_{n=1}^N \{ t_n \ln y_n + (1-t_n) \ln(1-y_n) \}$$
$$= - \sum_{n=1}^N \{ t_n \ln \sigma(\omega^T \phi_n) + (1-t_n) \ln(1-\sigma(\omega^T \phi_n)) \}$$

$$y_n = p(C_1 | \phi_n) = \sigma(\mathbf{w}^T \phi_n)$$

$$\frac{d\sigma}{dx} = \sigma(1-\sigma)$$

$$\nabla \ln \sigma(\omega^T \phi_n) =$$

$$\nabla \ln(1-\sigma(\omega^T \phi_n)) =$$

$$\nabla E(\omega) =$$

## Gradient of Logistic regression

$$E(\omega) = -\ln p(\hat{t}|\hat{\omega}) = - \sum_{n=1}^N \left\{ t_n \ln y_n + (1-t_n) \ln(1-y_n) \right\}$$

$$\rightarrow \frac{d\sigma}{dx} = \sigma(1-\sigma)$$

$$= - \sum_{n=1}^N \left\{ t_n \ln \sigma(\omega^\top \phi_n) + (1-t_n) \ln(1-\sigma(\omega^\top \phi_n)) \right\}$$

$$\nabla_{\omega} \ln \sigma(\omega^\top \phi_n) = \frac{1}{\sigma(\omega^\top \phi_n)} \sigma(\omega^\top \phi_n) (1-\sigma(\omega^\top \phi_n)) \phi_n$$

$\nabla_{\omega} E(\omega) = 0$   
 → maybe a local maximum

$$\omega^{(t+1)} = \omega^{(t)} + \eta \nabla E(\omega)$$

$$\begin{aligned} \nabla \ln(1-\sigma(\omega^\top \phi_n)) &= \frac{1}{1-\sigma(\omega^\top \phi_n)} [-\sigma(\omega^\top \phi_n)(1-\sigma(\omega^\top \phi_n)) \phi_n] \\ &= -\sigma(\omega^\top \phi_n) \phi_n \end{aligned}$$

$$\begin{aligned} \nabla E(\omega) &= - \sum_{n=1}^N \left\{ \underbrace{t_n(1-\sigma(\omega^\top \phi_n)) \phi_n}_{\text{scalar}} + (1-t_n)(-\sigma(\omega^\top \phi_n) \phi_n) \right\} \\ &= - \sum_{n=1}^N \left\{ t_n \phi_n + t_n \sigma(\omega^\top \phi_n) \phi_n - \sigma(\omega^\top \phi_n) \phi_n + t_n \sigma(\omega^\top \phi_n) \phi_n \right\} \\ \boxed{\nabla E(\omega)} &= \sum_{n=1}^N (\gamma_n - t_n) \phi_n \end{aligned}$$

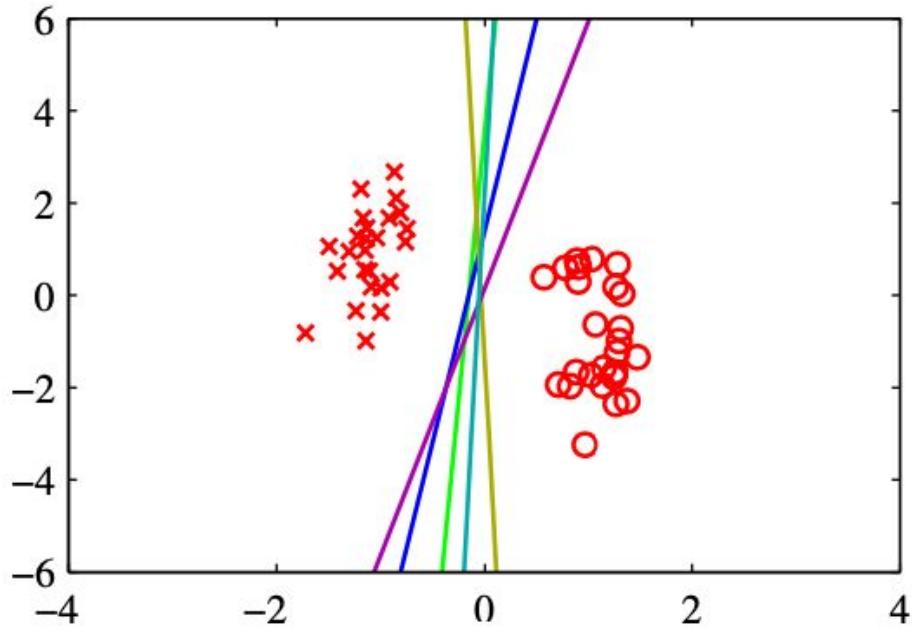
stochastic gradient

$$\nabla E_n(\omega) = \frac{(\gamma_n - t_n) \phi_n}{\text{error}}$$

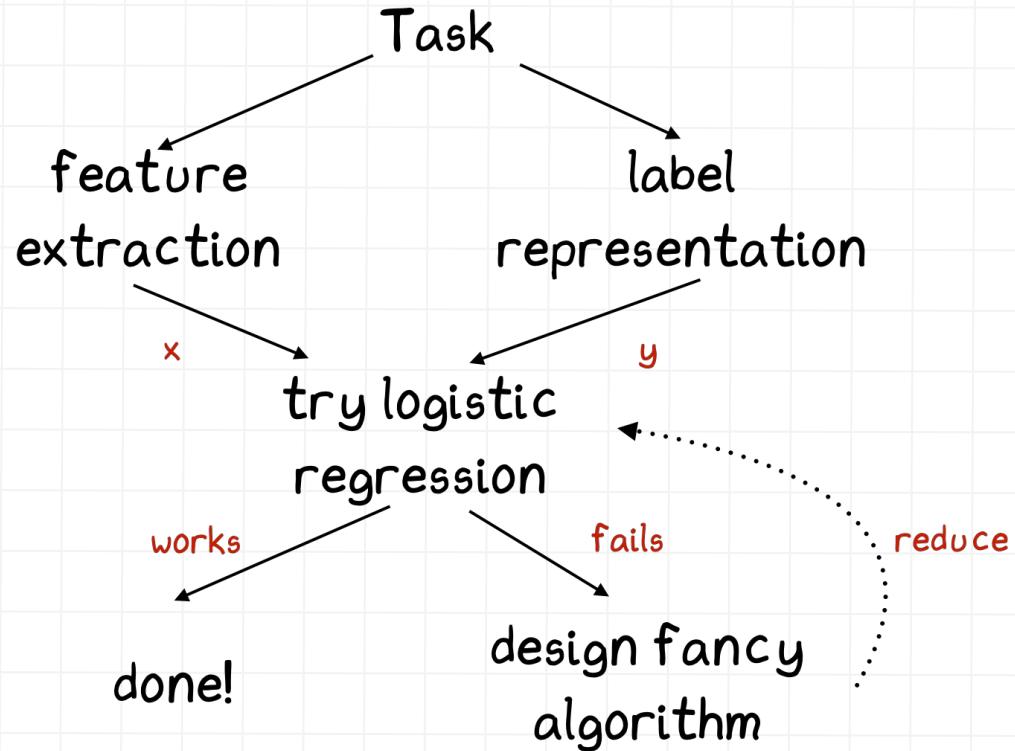
# A critique of logistic regression

$$P(C_1(x) = 1 | \underline{w}^T x)$$

PRML Fig 10.13



# the machine learning “master algorithm”\*



# Summary of Linear Models for classification

Classification problems

A glimpse of decision theory (Sec 1.5)

Discriminant functions - why least squares doesn't work here

The perceptron algorithm

Probabilistic generative models - origin of the logistic function

Probabilistic discriminative models

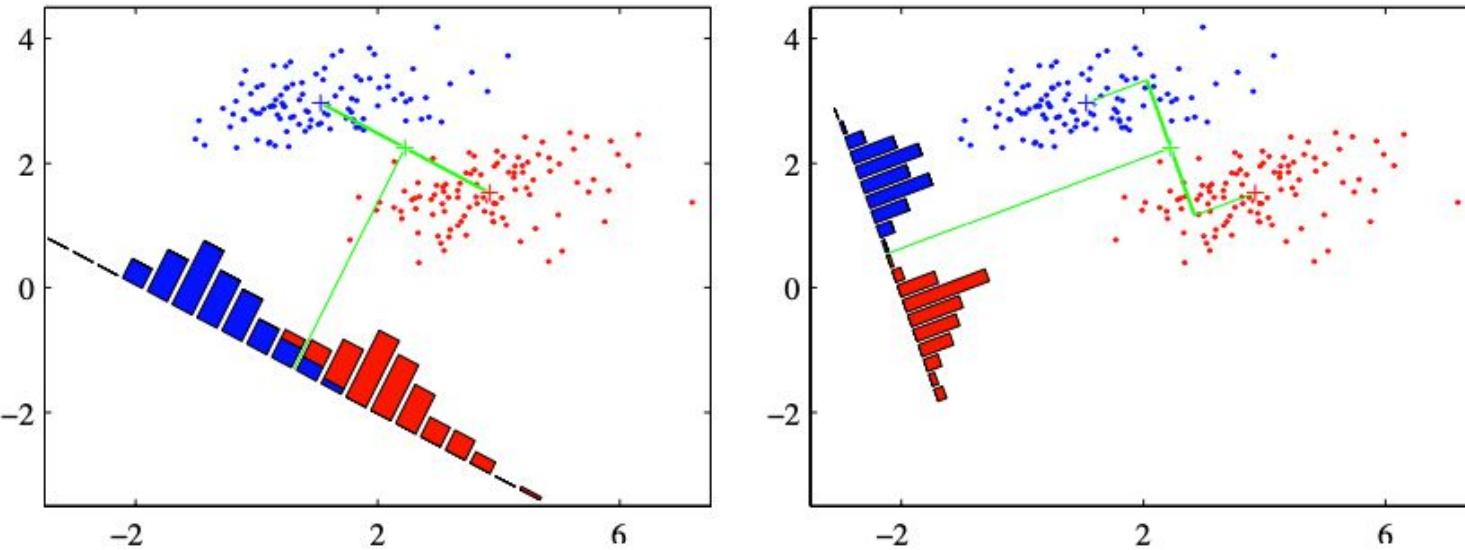
Logistic regression

(Laplace approximation, Bayesian logistic regression) – later in the semester

Origin of the logistic function, logistic regression and how it connects to perceptron

# Laplace approximation

## Other linear discriminants (e.g. Fisher)



**Figure 4.6** The left plot shows samples from two classes (depicted in red and blue) along with the histograms resulting from projection onto the line joining the class means. Note that there is considerable class overlap in the projected space. The right plot shows the corresponding projection based on the Fisher linear discriminant, showing the greatly improved class separation.

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad (4.26)$$

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1). \quad (4.30)$$