

# M7. Classification: Linear Models (and Kernel Methods)

Manikandan Narayanan

Week 11 (Oct 6-, 2025)

PRML Jul-Nov 2025 (Grads Section)

# Acknowledgment of Sources

- Slides based on content from related
  - Courses:
    - IITM – Profs. Arun/Harish/Chandra’s PRML offerings (slides, quizzes, notes, etc.), Prof. Ravi’s “Intro to ML” slides – cited respectively as [AR], [HR], [CC], [BR] in the bottom right of a slide.
    - India – NPTEL PR course by IISc Prof. PS. Sastry (slides, etc.) – cited as [PSS] in the bottom right of a slide.
  - Books:
    - PRML by Bishop. (content, figures, slides, etc.) – cited as **[CMB]**
    - Pattern Classification by Duda, Hart and Stork. (content, figures, etc.) – [DHS]
    - Mathematics for ML by Deisenroth, Faisal and Ong. (content, figures, etc.) – [DFO]
    - Information Theory, Inference and Learning Algorithms by David JC MacKay – [DJM]

# Outline for Module M7

- M7. Classification (Linear models)
  - **M7.0 Introduction**
  - M7.1 Generative model ((Naïve) Bayes classifier) (brief as already seen)
  - M7.2 Discriminative model (Logistic regression)
  - M7.3 Discriminant function (Perceptron, *Linear/Fisher's Discriminant Analysis*) (very brief)
  - M7.4 Kernel methods
  - M7.5 Concluding thoughts

# Outline for Module M7 (detailed)

- M7. Classification (Linear models)
  - **M7.0 Introduction**
    - Motivating examples / Birds-eye view
  - M7.1 Generative model (brief as already seen)
    - Naïve Bayes classifier: Gaussian, identity covariance
    - Bayes classifier: Gaussian, general (shared) covariance
    - From generative model → Logistic fn., Linear decision surfaces
  - M7.2 Discriminative model
    - Logistic regression: MLE and Gradient-descent solution
    - Ridge Logistic regression: Bayesian MAP/regularization
  - M7.3 Discriminant function (very brief)
    - Geometry of decision surfaces
    - Perceptron objective function
    - Unifying view of objective functions
  - M7.4 Kernel methods
    - Representer theorem
    - Kernel logistic regression / etc. (also Kernel regression (as Assignment qn), Kernel PCA, leading to Kernel SVM, etc.)
  - M7.5 Concluding thoughts

# Recall: Classification – standard examples

- Identify a given image as mango or orange, given many example images of mango and orange.



M



M



O



M

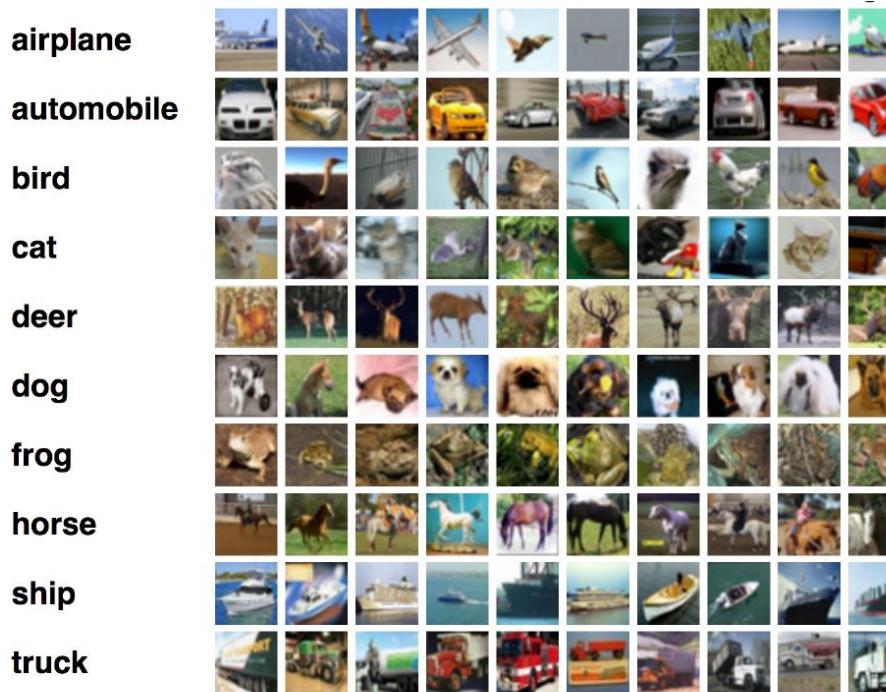


O



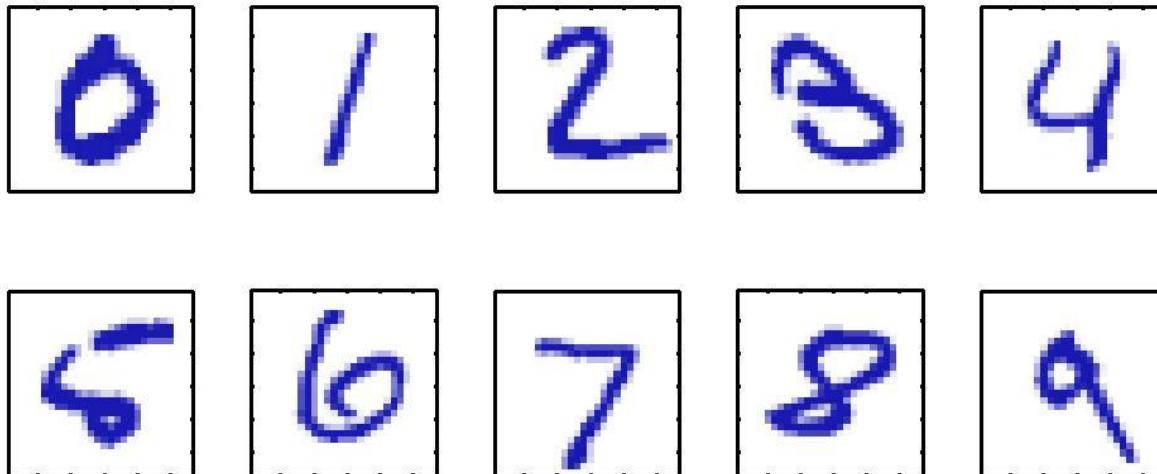
# Recall: Classification – popular examples

Classify an image into one of 10 classes,  
given a “training set” of images with classes.



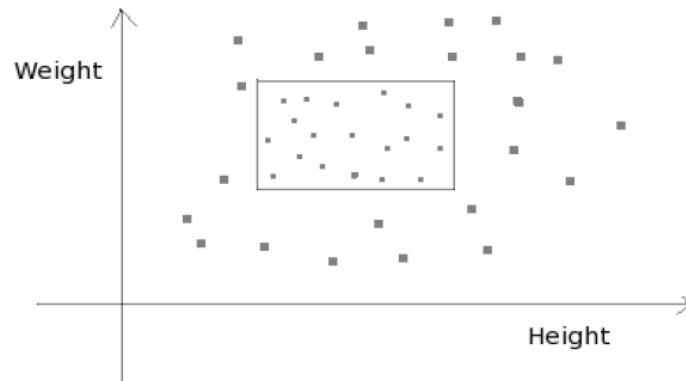
# Recall: Example classification problem

- Handwritten Digit Recognition:
  - Take a digit's  $28 \times 28$  pixel image, represented by a vector of 784 real numbers, as input feature vector  $\mathbf{x}$ , and output the identify of the digit 0,...,9.
  - Non-trivial problem due to the wide variability of handwriting.



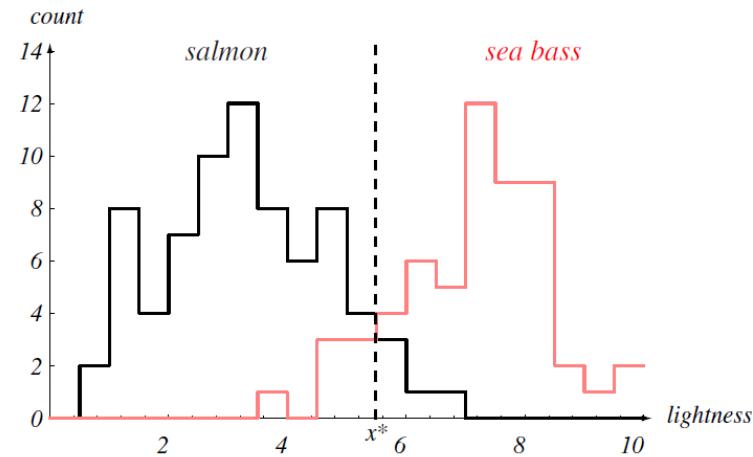
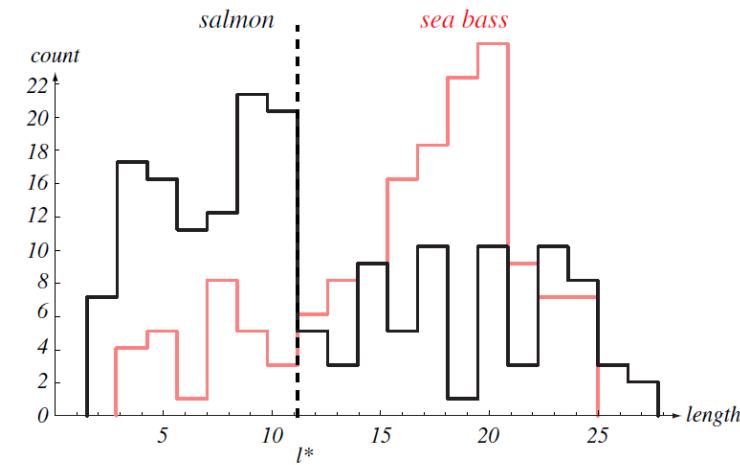
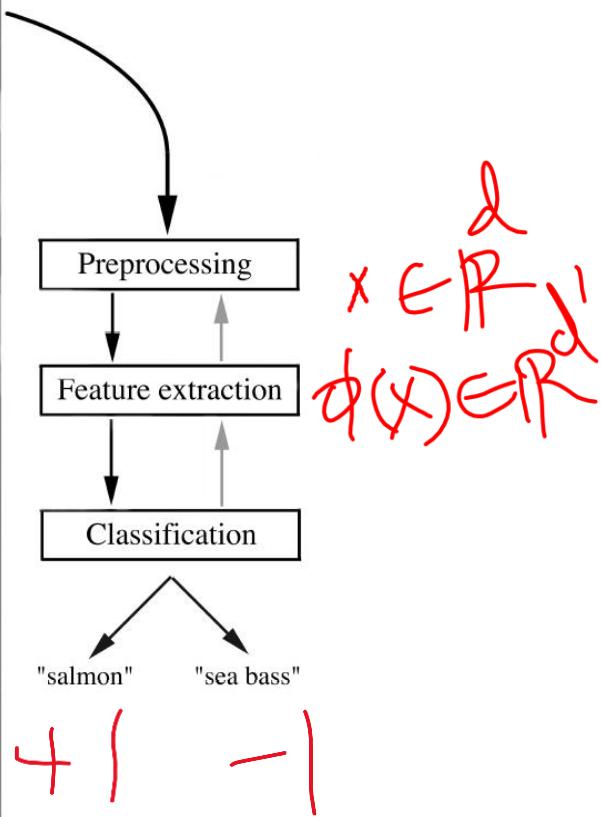
# Recall: Example (non-linear) classification problem

- Problem: recognize persons of ‘medium build’
- Features: Height and Weight



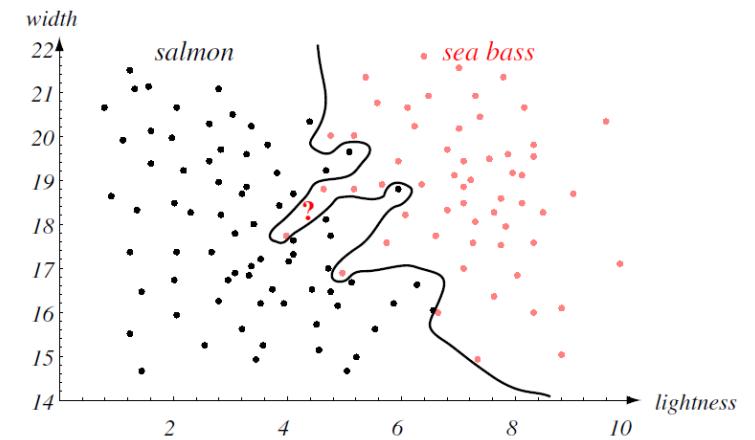
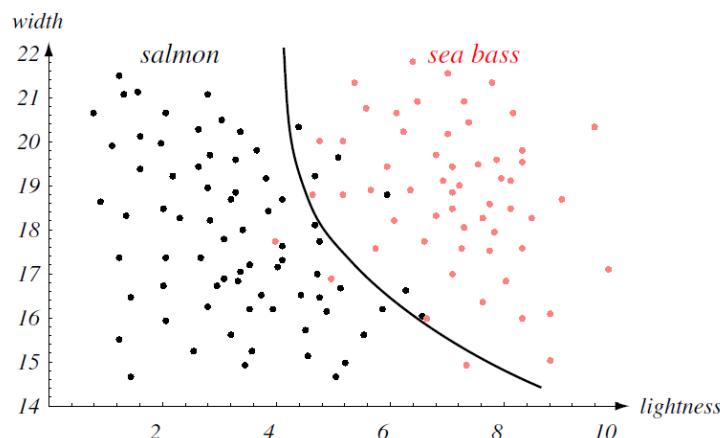
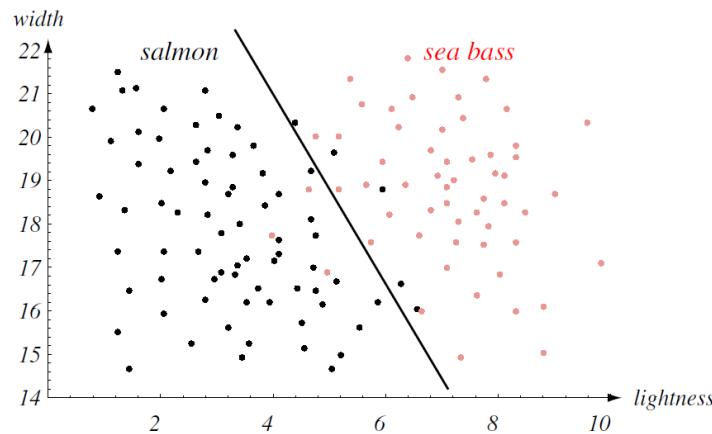
The classifier is ‘nonlinear’ here

# Recall: Classification – standard examples



[DHS]

# Recall: Decision boundary – from simple to complex



# Classification

- Learn a map from input variables (features  $x$ ) to a discrete/categorical output variable (target  $t \in [K] := \{1, 2, \dots, K\} := \{C_k\}_{k=1,\dots,K}$ ; **1-of-K (one-hot)** encoding of  $t$  (or) **-1/+1** for binary  $t$  is also popular).  
Given  $\{x_n, t_n\}_{n=1\dots N}$  pairs, we seek a function  $f_w(x)$  that approximates the map from  $x \rightarrow t$ .

Linear if  $f_w(x) := f(z(x, w)) = f(w^T x + w_0)$  where

- argument of  $f$  i.e.,  $z(x, w)$  is **linear** in adjustable parameters  $w$  and  $x$  (but non-linear basis fns.  $\phi(x)$  also possible).
- f is a non-linear activation function** (its inverse is called *link function* in generalized linear models in statistics)
- decision surfaces correspond to  $f_w(x) = \text{const.}$ , so that  $w^T x + w_0 = \text{const.}$  (so **decision surfaces are linear in  $x$** , even if  $f(\cdot)$  is non-linear, as long as its one-to-one)

- A basic supervised learning problem/algorithm like linear regression:

- Practical limitations for complex data, but basic linear classification methods (**Naïve Bayes classifier and Logistic regression**) is a good starting point for other sophisticated learning algorithms. Linear models being simpler are highly efficient to learn, highly interpretable, etc.
- Optimization problem doesn't result in a closed-form solution as in linear regression; but objective function is convex and gradient-descent methods work well.

- Our approach in this lecture: Mostly [Harish's logistic regression notes; CMB, Chapter 4].

IP:  $\rightarrow (x_1, \dots, x_N, \dots, x_N \in \mathbb{R}^D)$   
OP:  $\rightarrow \omega : f(x, w) \text{ s.t. } f(x, w) \text{ approx. } t_n$

# Example $f(\cdot)$ (activation functions)

- Binary (0-1) classifn.: Let  $z := w^T x + w_0$ . Then,

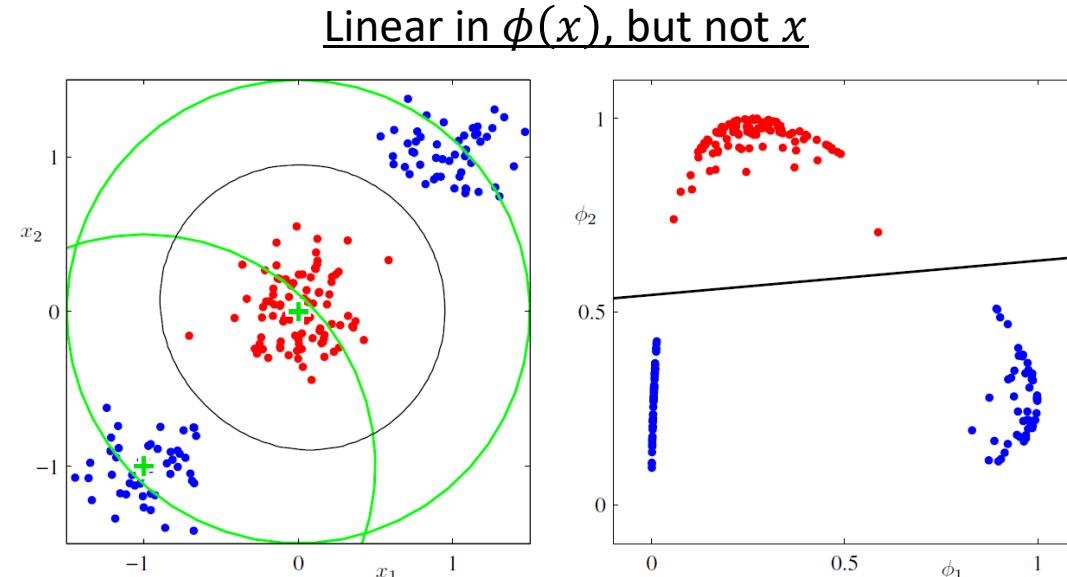
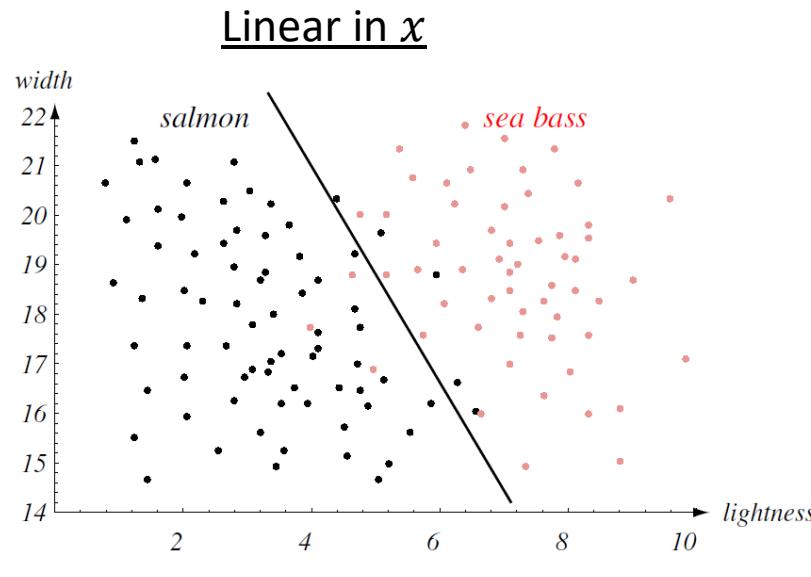
$$f(z) = \text{logistic}(z) := \sigma(z) := \frac{1}{1+\exp(-z)}.$$

- Mult-class (1-of-K one-hot) classifn.: Let  $z_k := w_k^T x + w_{k0}$ . Then,

$$f(z_1, z_2, \dots, z_K) = \text{soft(argmax}}(z_1, z_2, \dots, z_K) := \left\{ \frac{\exp(z_k)}{\sum_i \exp(z_i)} \right\}_{k=1 \text{ to } K}.$$

# A note on DB (Decision Boundary): linear in $x$ or $\phi(x)$ ?

If  $f$  one-to-one, then  $f(w^T x + w_0) = \text{cutoff}$  (say 0.5)  $\Rightarrow w^T x + w_0 = \text{const.}$  ( $x$  could've been  $\phi(x)$  too).



DB also linear in  $w$ , but of less use here than in linear regression (where obj. fn linear in  $w$  led to closed-form solution).

In linear models of classification, obj. fn. (empirical expected loss or empirical risk) is still convex in  $w$  to enable nice properties for optmzn. via gradient descent or similar methods (but closed-form solution for  $w$  not possible).

# Recall: Inference and decision: three approaches for classification

- Generative model approach:

- (I) Model  $p(x, C_k) = p(x|C_k)p(C_k)$

- (I) Use Bayes' theorem  $p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)}$

- (D) Apply optimal decision criteria

- Discriminative model approach:

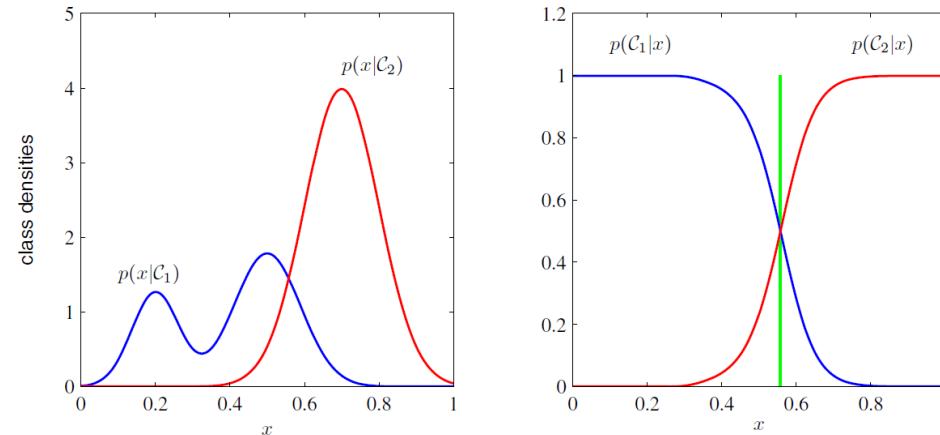
- (I) Model  $p(C_k|x)$  directly

- (D) Apply optimal decision criteria

- Discriminant function approach:

- (D) Learn a function that maps each  $x$  to a class label directly from training data

Note: No posterior probabilities!



# Outline for Module M7

- M7. Classification (Linear models)
  - M7.0 Introduction
  - **M7.1 Generative model**
    - Naïve Bayes classifier: Gaussian, identity covariance
    - Bayes classifier: Gaussian, general (shared) covariance
    - From generative model → Logistic fn., Linear decision surfaces
  - M7.2 Discriminative model (Logistic regression)
  - M7.3 Discriminant function (Perceptron) (very brief)
  - M7.4 Kernel methods
  - M7.5 Concluding thoughts

# K=2 Bayes classifier – Gaussian, identity covariance

Recall : Generative Models for classification

$$X, Y \in \mathbb{R}^d \times \{\pm 1\}$$

$$X | y=+1 \sim N(\mu_+, I) ; P(y=+1) = a$$

$$X | y=-1 \sim N(\mu_-, I) ; P(y=-1) = 1-a$$

$$\begin{aligned} P(y=1 | x=x) &= \frac{f_{x|y}(x|+1) P(y=+1)}{f_{x|y}(x|+1) P(y=+1) + f_{x|y}(x|-1) P(y=-1)} \\ &= \frac{a \cdot \exp(-\frac{1}{2} \|x - \mu_+\|^2)}{a \exp(-\frac{1}{2} \|x - \mu_+\|^2) + (1-a) \exp(-\frac{1}{2} \|x - \mu_-\|^2)} \end{aligned}$$

$$= \frac{1}{1 + \exp(x^T \mu_- - x^T \mu_+ - \frac{1}{2} \|\mu_-\|^2 + \frac{1}{2} \|\mu_+\|^2 + \ln(\frac{1-a}{a}))}$$

$$= \frac{1}{1 + \exp(-(w^T x + b))}$$

$$\text{where } w = \mu_+ - \mu_-$$

$$b = \frac{-1}{2} \|\mu_+\|^2 + \frac{1}{2} \|\mu_-\|^2 + \ln \frac{a}{1-a}$$

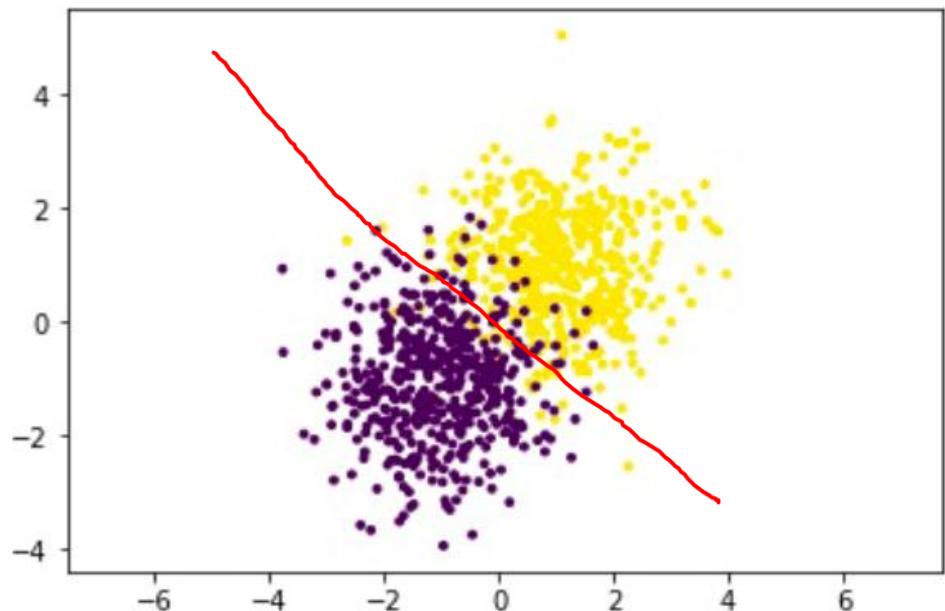
$$P(y=1|x) = \sigma(w^T x + b)$$

where

$$\sigma(u) = \frac{1}{1 + \exp(-u)} \Rightarrow$$



Example data 1



$$x | y = -1 \sim N \left( \begin{bmatrix} -1 \\ -1 \end{bmatrix}, I \right)$$

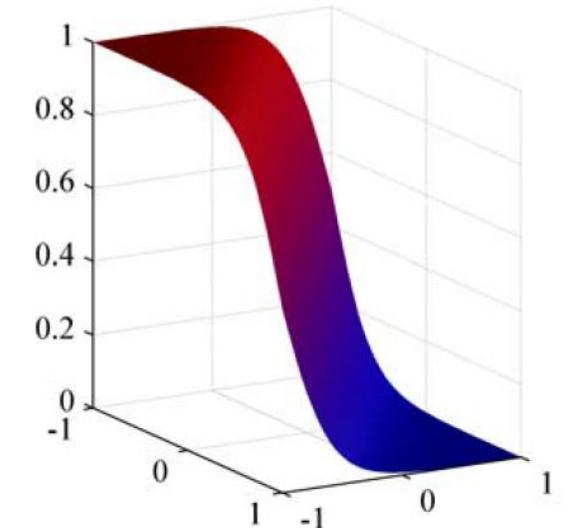
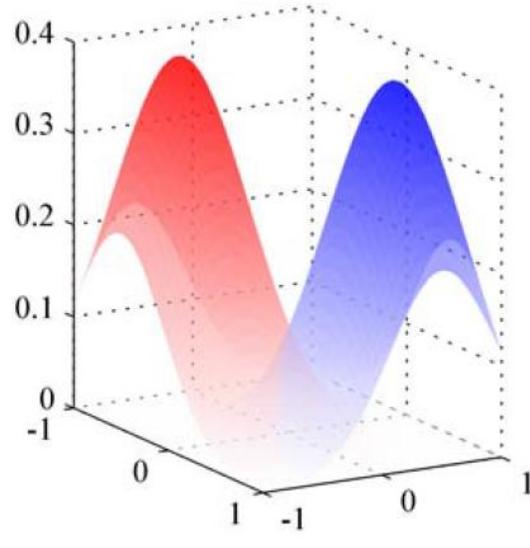
$$x | y = +1 \sim N \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix}, I \right)$$

$$P(y=1) = \frac{1}{2}$$

$$P(y=1 | x=x) = \sigma([2, 2]x)$$

# K=2 Bayes classifier – Gaussian, general (shared) covariance

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}. \quad p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$



$$\begin{aligned} \mathbf{w} &= \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ w_0 &= -\frac{1}{2}\boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}. \end{aligned}$$

[CMB]

# K>2 Bayes classifier – Gaussian, general (shared) covariance

- Normalized exponential or softmax:

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

- Gaussian:

(same as before)

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

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

$$\mathbf{w}_k = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k$$

$$w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \ln p(\mathcal{C}_k).$$

# Sanity check

That's all there is to generative models of classification...  
...or am I missing something?

How do we learn the generative model? Density estimation (MLE) of Gaussian ( $K=2$  shown;  $K>2$  is similar)

$$\hat{\pi}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{1}_{\{\mathbf{y}_n = C_1\}} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$$

$$\hat{\mu}_{\text{ML}} = \frac{1}{N_1} \sum_{n=1}^N \mathbf{1}_{\{\mathbf{y}_n = C_1\}} \mathbf{x}_n \quad \hat{\mu}_{\text{ML}} = \frac{1}{N_2} \sum_{n=1}^N \mathbf{1}_{\{\mathbf{y}_n = C_2\}} \mathbf{x}_n$$

$$\hat{\Sigma}_{\text{ML}} = \mathbf{S} = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2$$

$$\mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in C_1} (\mathbf{x}_n - \hat{\mu}_{\text{ML}})(\mathbf{x}_n - \hat{\mu}_{\text{ML}})^T$$

$$\mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in C_2} (\mathbf{x}_n - \hat{\mu}_{\text{ML}})(\mathbf{x}_n - \hat{\mu}_{\text{ML}})^T.$$

$$\begin{aligned} L(\theta) &= \prod_{n=1}^N P(x_n, y_n) \\ &\leftarrow (\pi, \mu_1, \mu_2, \Sigma) \\ &= \prod_{n=1}^N P(y_n | \pi) P(x_n | y_n, \mu_1, \mu_2, \Sigma) \\ LL(\theta) &= \sum_{n=1}^N \log P(y_n | \pi) \\ &+ \sum_{n: y_n = C_1} \log P(x_n | y_n, \mu_1, \Sigma) \\ &+ \sum_{n: y_n = C_2} \log P(x_n | y_n, \mu_2, \Sigma) \end{aligned}$$

# What if covariance matrix is not shared across classes?

- Then,  $LL(\theta)$  in last slide becomes:  
$$\begin{aligned} LL(\theta) = & \sum_n \log P(y_n | \pi) \\ & + \sum_{n: y_n = c_1} \log P(x_n | y_n, \mu_1) \\ & + \sum_{n: y_n = c_2} \log P(x_n | y_n, \mu_2) \end{aligned}$$
- and each of these three summation terms can be optimized separately!
- Therefore, all params' MLE remain the same, except for these params.:

$$\hat{\Sigma}_{1(m)} = S_1 \quad \text{and} \quad \hat{\Sigma}_{2(m)} = S_2 .$$

**Recall:** Inference & Decision (steps in detail for a generative model when covariance matrix is **not** shared):

Setup (training data):

Learning the Bayes classifier from data:

- $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  with  
 $x_i \in \mathcal{X}, y_i \in \{-1\}$

Inference (density est.):

- Assume  $P_{x|y}(x|1)$  and  $P_{x|y}(x|-1)$  are from  $P$ .

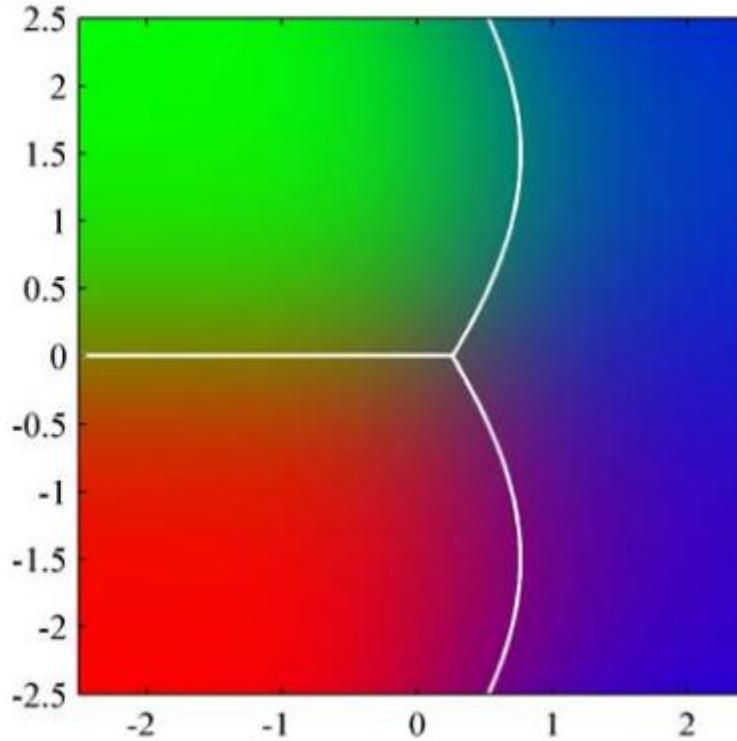
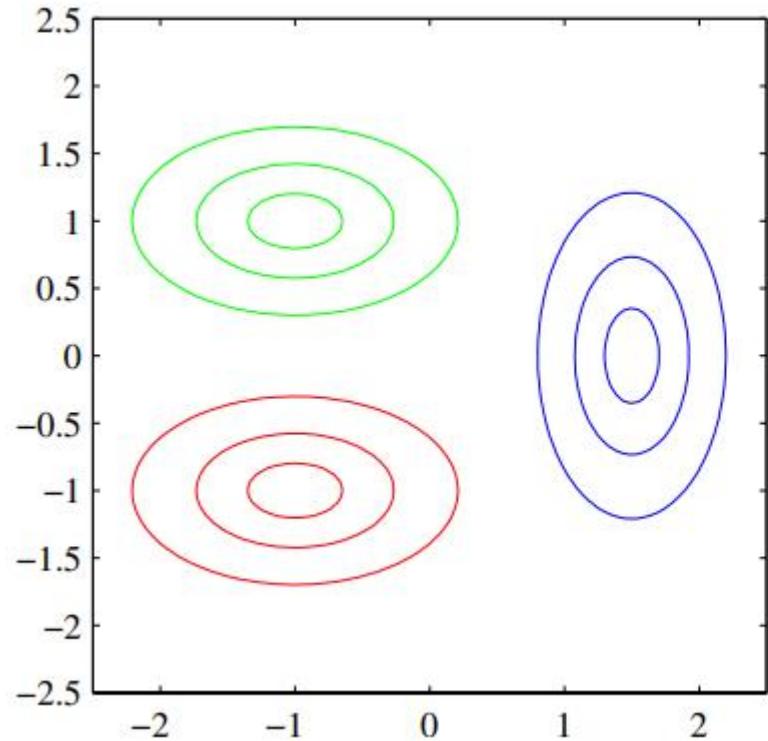
- Estimate  $P_{x|y}(x|1)$  from positively labelled samples  
and  $P_{x|y}(x|-1)$  from negatively labelled samples  
using Maximum Likelihood.
- Estimate  $P(y=1)$  &  $P(y=-1)$

- Use Bayes Rule to make an estimate of  
the posterior probability  $P(y=1|x=x) = \eta(x)$

Decision:

- Compute  $h^*(x) = \text{Sign}(2\eta(x)-1)$

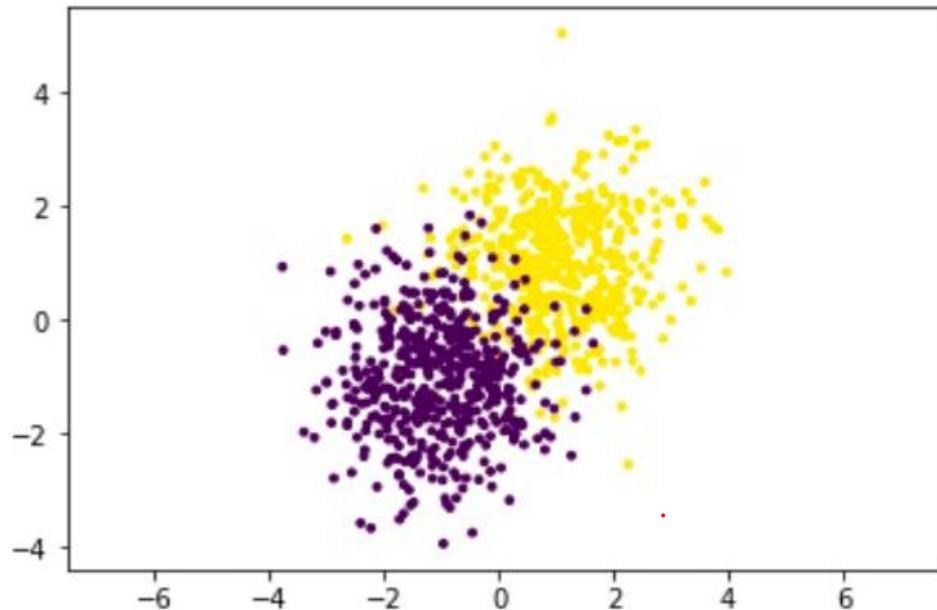
Shared vs. non-shared covariance matrix: linear vs. quadratic discriminant (DB – quadratic terms do not cancel out in non-shared case!)



# Sanity check

Now, we're really done with generative models of classification...  
...and let's motivate the discriminative model approach.

Example data 1



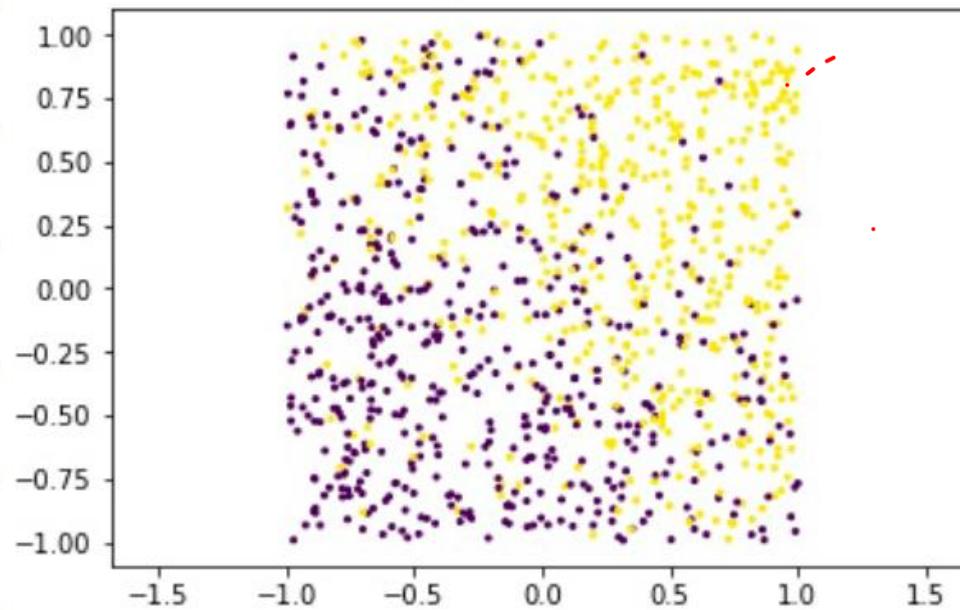
$$x | y = -1 \sim N \left( \begin{bmatrix} -1 \\ -1 \end{bmatrix}, I \right)$$

$$x | y = +1 \sim N \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix}, I \right)$$

$$P(y=1) = \frac{1}{2}$$

$$P(y=1 | x=x) = \sigma([2, 2]x)$$

Example data 2



$$P(y=1 | x=x) = \sigma([2, 2]x)$$

Compared to generative models:

Discriminative (logistic regn.): less restrictive (with enough data)

Discriminant (linear decision surface): even less restrictive (but no probabs.)

# Outline for Module M7

- M7. Classification (Linear models)
  - M7.0 Introduction
  - M7.1 Generative model ((Naïve) Bayes classifier)
  - **M7.2 Discriminative model**
    - Logistic regression: MLE and Gradient-descent solution
    - Ridge Logistic regression: Bayesian MAP/regularization
  - M7.3 Discriminant function (Perceptron) (very brief)
  - M7.4 Kernel methods
  - M7.5 Concluding thoughts

## Classification :

Discriminative models :

$$x, y \in \mathbb{R}^d \times \{-1, 1\} \sim D$$
$$x \sim D_x$$

$$P(y=1|x=x) = \sigma(w^T x);$$

$$P(y=-1|x=x) = 1 - \sigma(w^T x) = \sigma(-w^T x)$$

$$\left. \begin{array}{l} P(y=y|x=x) = \sigma(y w^T x) \end{array} \right\}$$

where  $w$  is some unknown constant.

1. Given  $(x_1, y_1), \dots, (x_m, y_m)$  i.i.d. samples drawn from  $D$ . Can  $w$  be discovered?
2. If  $w$  is known, how to predict for a given  $x$ ?

$$P(y_1, y_2, \dots, y_m | x_1, \dots, x_m, w) = \mathcal{L}(w)$$

$$= \prod_{i=1}^m P(y_i = y_i | x_i = x_i, w)$$

$$= \prod_{i=1}^m (\sigma(w^T x_i))^{I(y_i=1)} (\sigma(-w^T x_i))^{I(y_i=-1)} = \boxed{\prod_{i=1}^m \sigma(y_i w^T x_i)}$$

$$\mathcal{L}_d(w) = \sum_{i=1}^m \log \sigma(y_i w^T x_i)$$

$$= \sum_{i=1}^m \log \frac{1}{1 + e^{-y_i w^T x_i}}$$

$$-\mathcal{L}(w) = \sum_{i=1}^m \log (1 + \exp(-y_i w^T x_i)) = \hat{\mathcal{R}}(w)$$

↓

Empirical logistic loss

$$\nabla \hat{\mathcal{R}}(w) = \sum_{i=1}^m \frac{1}{1 + \exp(-y_i w^T x_i)} (\exp(-y_i w^T x_i)) (-y_i x_i)$$

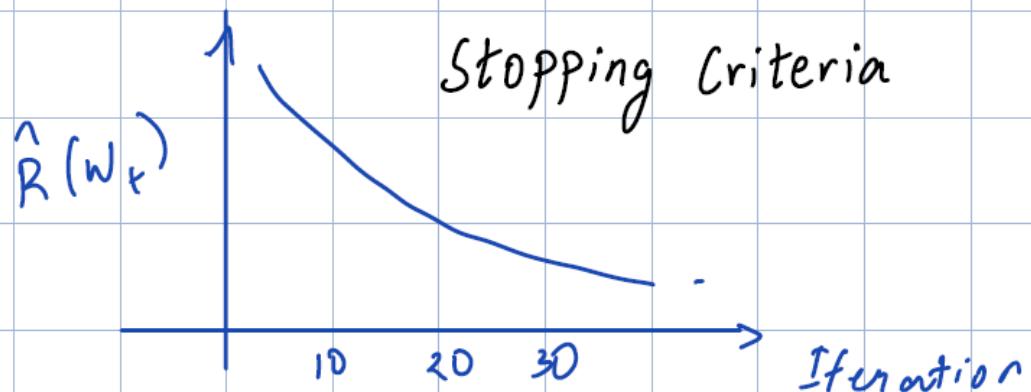
$$= \sum_{i=1}^m \sigma(-y_i w^T x_i) (-y_i x_i)$$

Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla \hat{R}(w_t)$$

Apply Gradient descent to Empirical Logistic Loss:

$$w \mapsto w + \eta \sum_{i=1}^m \sigma(-y_i w^\top x_i) y_i x_i$$



NOTE: Iterated reweighted least squares method  
(Newton-Raphson instead of gradient-descent) also  
possible and common with logistic regn.

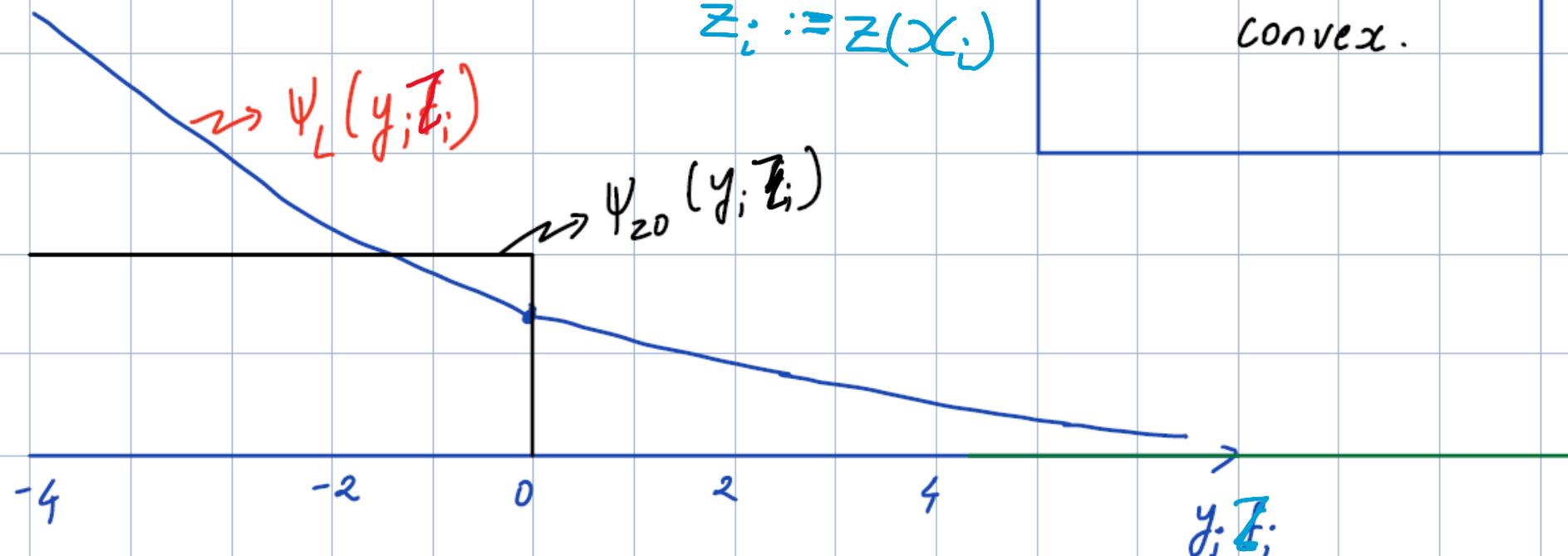
## Empirical Logistic loss Minimisation .

$$\hat{R}(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

$$= \sum_{i=1}^n \psi_L(y_i; w^T x_i)$$

$\underbrace{\zeta_i := z(x_i)}$

P.T  $\psi_L: \mathbb{R} \rightarrow \mathbb{R}$  is  
convex.



## Regularised/MAP Logistic Regression

$$w \sim N(0, \tau^2 I) \quad x \sim D_x$$

$$P(Y = +1 | x) = \sigma(w^T x)$$

$$\mathcal{L}^P(w) = P(w | x_1, \dots, x_n, y_1, \dots, y_n)$$

$$= \frac{P(y_1, \dots, y_n | x_1, \dots, x_n, w) P(w | x_1, \dots, x_n)}{P(y_1, \dots, y_n | x_1, \dots, x_n)}$$

Steps here

$$-\mathcal{L}^P(w) = -\sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \frac{\lambda}{2} \|w\|^2$$

(Subject to additive & multiplicative constants)

# Outline for Module M7

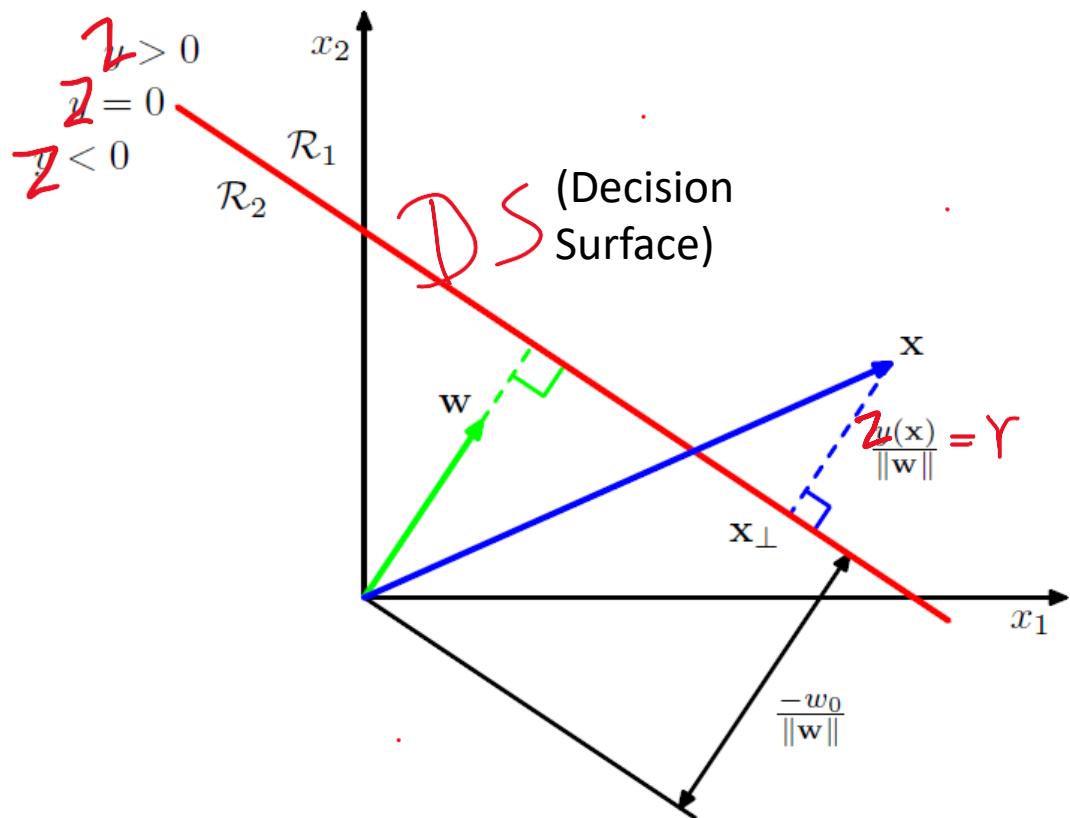
- M7. Classification (Linear models)
  - M7.0 Introduction
  - M7.1 Generative model ((Naïve) Bayes classifier)
  - M7.2 Discriminative model (Logistic regression)
  - **M7.3 Discriminant function (very brief)**
    - **Geometry of decision surfaces**
    - Perceptron objective function
    - Unifying view of objective functions
  - M7.4 Kernel methods
  - M7.5 Concluding thoughts

# Discriminant

- Discriminant is a function that takes an input vector  $x \in \mathbb{R}^d$  and assigns it to one of the  $K$  classes
  - we will assume  $K=2$  henceforth for simplicity!
- We focus only on **linear discriminants** (i.e., those for which DB is a hyperplane wrt  $x$  (or  $\phi(x)$ )).
  - $z(x) = w^T x + w_0$  (or  $w^T \phi(x)$ )
  - DB:  $z(x) = 0$  (hyperplane)
  - Prediction:  $f(z(x)) = \text{sign}(z(x))$  (i.e., Predict  $C_1$  if  $z(x) \geq 0$ , &  $C_2$  if  $z(x) < 0$ )

Recall defn. of hyperplane  $\{x \in \mathbb{R}^d: w^T x = b\}$ , which is a  $(d - 1)$ -dimensional (affine) subspace of a  $d$ -dim. vector space.

# Geometry of decision surfaces: 2-way or binary classification



DB is  $w^T x + w_0 = \text{const.}$ , but const. can be absorbed into  $w_0$  to get DB as  $w^T x + w_0 = 0$ .  
 Let  $z(w, x) = w^T x + w_0$  with the const. absorbed.

Signed dist.  $r$  of  $x$  to DB:  
 Express  $x = x_{\perp} + r \frac{w}{\|w\|}$

Mult. by  $w^T$  & add  
 $w_0$  on both sides:

$$\Rightarrow z(x) = 0 + r \|w\|$$

$$\Rightarrow r = \frac{z(x)}{\|w\|}$$

# Outline for Module M7

- M7. Classification (Linear models)
  - M7.0 Introduction
  - M7.1 Generative model ((Naïve) Bayes classifier)
  - M7.2 Discriminative model (Logistic regression)
  - **M7.3 Discriminant function (very brief)**
    - Geometry of decision surfaces
    - **Perceptron objective function**
    - Unifying view of objective functions
  - M7.4 Kernel methods
  - M7.5 Concluding thoughts

# Perceptron

(artificial neuron - elementary unit in ANNs)

- General form of linear models of classification:

$$f_w(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0), \text{ or}$$
$$f_w(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$$

- Choice of function  $f$  for perceptron (or any linear discriminant for that matter):

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

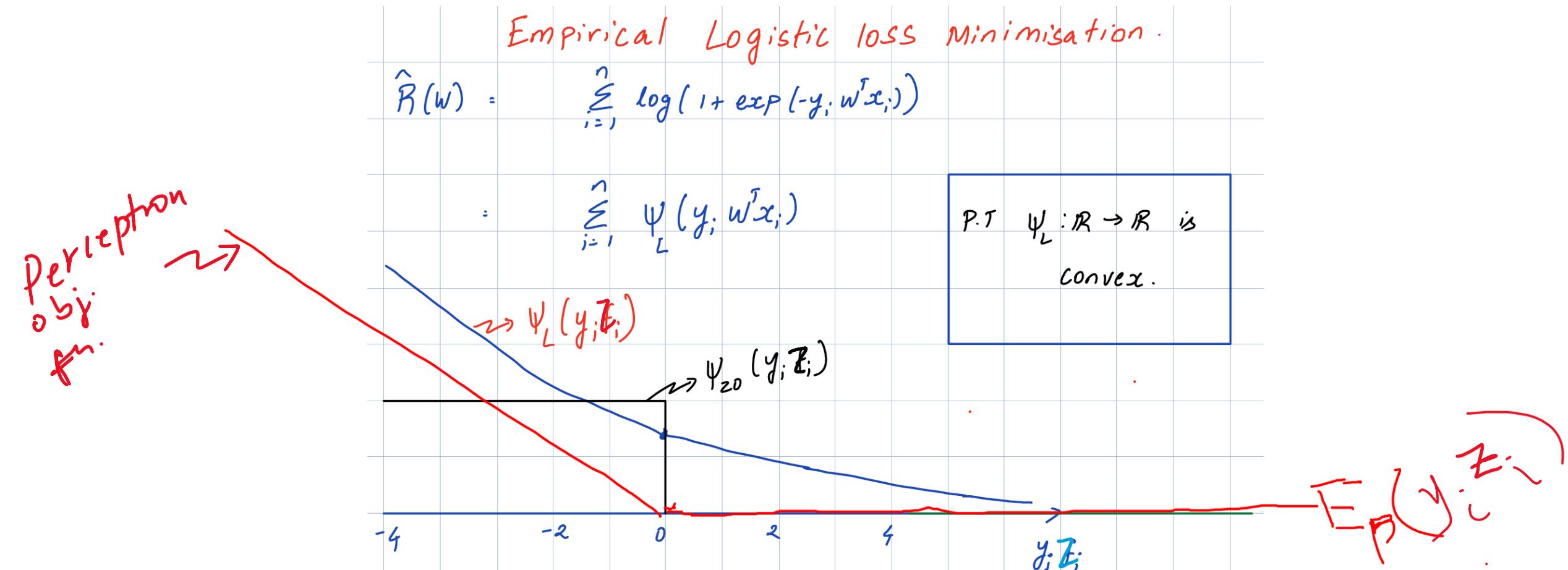
- Choice of objective function **specific to perceptron** to learn  $w$ :

$$E_P(w) = \sum_{i \in \mathcal{M}} (-y_i \underbrace{\mathbf{w}^T \phi(x_i)}_{\text{misclassified}})$$

misl|assif.

*distance from*  
 *$\phi(x_i)$*   
 *$D_S \cup D_B$*

# Recall: Logistic regression – direct (non-probabilistic) motivation using obj. fn.



Perceptron – simple (stochastic grad. descent) training algo. exists, but only of historical importance, as it suffers from many issues (doesn't converge for non-linearly separable dataset, multiple decision surfaces possible for linearly separable datasets, convergence in practice may take many steps,...)

# Outline for Module M7

- M7. Classification (Linear models)
  - M7.0 Introduction
  - M7.1 Generative model ((Naïve) Bayes classifier)
  - M7.2 Discriminative model
  - **M7.3 Discriminant function (very brief)**
    - Geometry of decision surfaces
    - Perceptron objective function
    - **Unifying view of objective functions**
  - M7.4 Kernel methods
  - M7.5 Concluding thoughts

The larger context: loss fn view offers a unified motivation of many classifn. methods (thereby allowing us to condense a “laundry” list of methods into a single framework)

## 1 LOSS FUNCTIONS FOR BINARY CLASSIFICATION

Given data :  $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \{\pm 1\}$

Most algorithms we see are “Surrogate risk” minimisation algorithms.

i.e) They want to find an  $\mathcal{L}: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\text{Sign } (\mathcal{L}(x_i)) = y_i$$

or equivalently minimise

$$z_0[f] = \sum_{i=1}^n \psi_{z_0}(y_i \mathcal{L}(x_i))$$

$$\begin{aligned} \text{where } \psi_{z_0}(u) &= 1 \quad \text{if } u \leq 0 \\ &= 0 \quad \text{otherwise.} \end{aligned}$$

# Surrogate loss fns

But due to the difficulty of minimising the zero one error we minimise one of these:

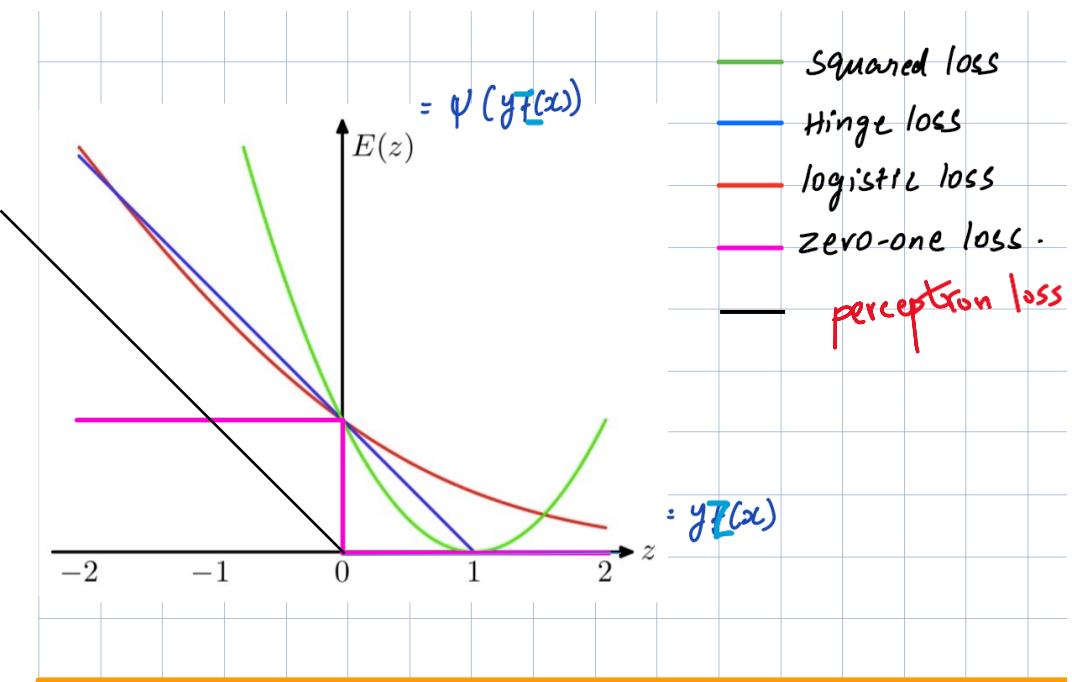
$$sq[\mathcal{E}] = \sum_{i=1}^n \psi_q(y_i; \mathcal{E}(x_i)) = \sum_{i=1}^n (1 - y_i \mathcal{E}(x_i))^2 = \sum_{i=1}^n (y_i - \mathcal{E}(x_i))^2$$

$$\text{logistic } [\mathcal{E}] = \sum_{i=1}^n \psi_L(y_i; \mathcal{E}(x_i)) = \sum_{i=1}^n \log(1 + \exp(-y_i \mathcal{E}(x_i)))$$

$$\text{Hinge } [\mathcal{E}] = \sum_{i=1}^n \psi_H(y_i; \mathcal{E}(x_i)) = \sum_{i=1}^n \max(0, 1 - y_i \mathcal{E}(x_i))$$

$$\text{exp } [\mathcal{E}] = \sum_{i=1}^n \psi_E(y_i; \mathcal{E}(x_i)) = \sum_{i=1}^n \exp(-y_i \mathcal{E}(x_i))$$

$$\text{perc}[z] = \sum_{i=1}^n \psi_{\text{perc}}(y_i; \mathcal{E}(x_i)) = \sum_{i=1}^n \max(0, -y_i \mathcal{E}(x_i))$$

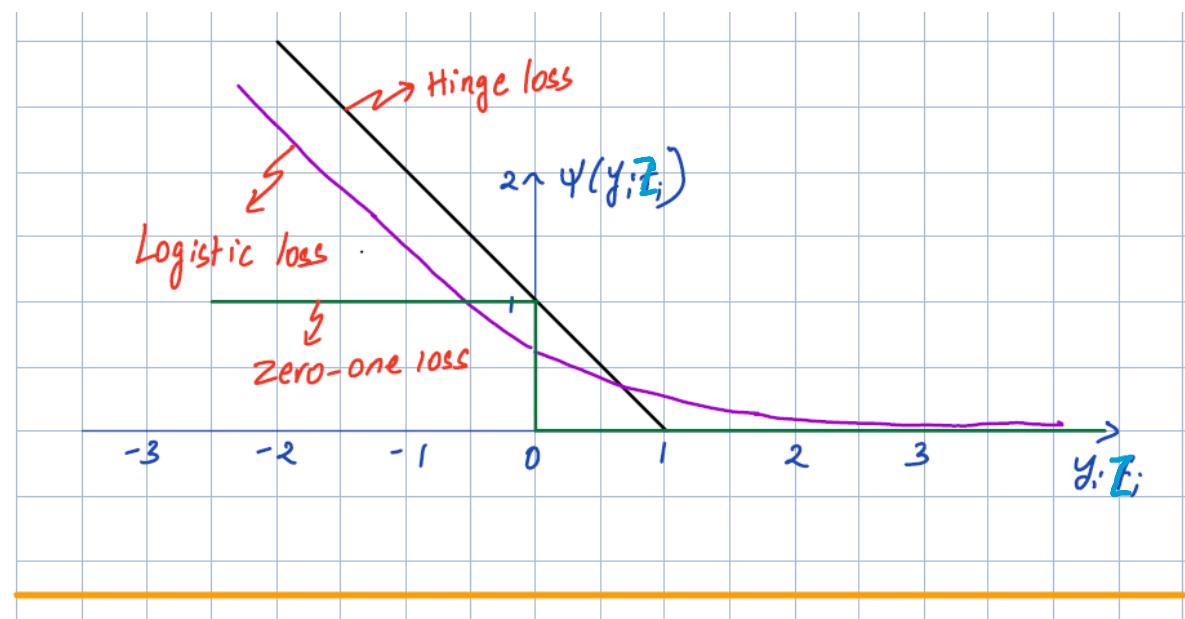


# More on Hinge loss and (soft-margin) SVM

$$\begin{aligned} \min_{w,b,\xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i (w^T x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned} \quad \left\{ \begin{array}{l} \Rightarrow \\ \min_{w,b} C \sum_{i=1}^n \psi_H(y_i(w^T x_i + b)) \\ + \frac{1}{2} \|w\|^2 \end{array} \right.$$

Express  $\xi_i$  as a function of  $w, b$ :

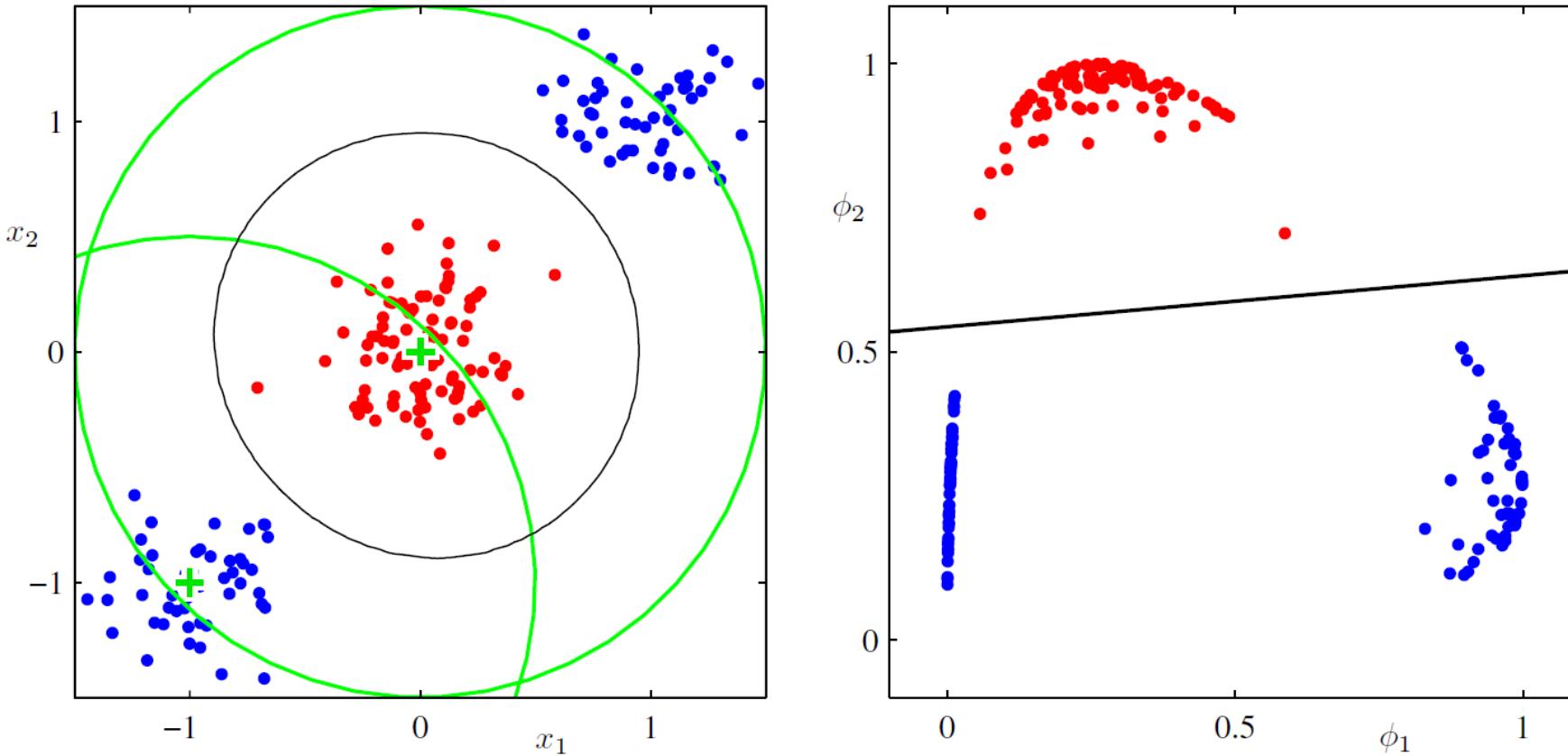
$$\begin{aligned} \xi_i &\geq 1 - y_i (w^T x_i + b) ; \quad \xi_i \geq 0 \\ \therefore \quad \xi_i &= \max(1 - y_i (w^T x_i + b), 0) \\ &= \max(1 - y_i z_i, 0) \quad \text{where } z_i = w^T x_i + b \\ &= \psi_H(y_i z_i) \end{aligned}$$



$$\psi_H(u) = \max(1 - u, 0)$$

# Outline for Module M7

- M7. Classification (Linear models)
  - M7.0 Introduction
  - M7.1 Generative model ((Naïve) Bayes classifier)
  - M7.2 Discriminative model (Logistic regression)
  - M7.3 Discriminant function (Perceptron) (very brief)
  - **M7.4 Kernel methods**
    - Representer theorem
    - Kernel logistic regression / Kernel regression (as Asst. qn.) / etc. (leading to Kernel SVM, etc.)
  - M7.5 Concluding thoughts



**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  $\phi_1(\mathbf{x})$  and  $\phi_2(\mathbf{x})$  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.

$\phi: \mathbb{R}^n \rightarrow \mathbb{R}^{d'}$  Kernel Logistic Regression

$$\min_{w \in \mathbb{R}^{d'}} \sum_{i=1}^n \log \left( 1 + \exp(-y_i w^T \phi(x_i)) \right) + \frac{\lambda}{2} \|w\|^2$$

$\uparrow\downarrow$  RT

$$\Phi = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix}_{n \times d'}$$

$$w = \Phi^T \alpha$$

$$K = \Phi \Phi^T$$

nxn

$$\min_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \log \left( 1 + \exp(-y_i \alpha^T \Phi \phi(x_i)) \right) + \frac{\lambda}{2} \alpha^T \Phi \Phi^T \alpha$$

$\uparrow\downarrow$  subst.

The  $i^{\text{th}}$  column of  $K$ .

$$\min_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \log \left( 1 + \exp(-y_i \alpha^T K_{:,i}) \right) + \frac{\lambda}{2} \alpha^T K \alpha$$

# Representer theorem

$$\min_{w \in \mathbb{R}^d} f(\Phi w) + \frac{\gamma}{2} \|w\|^2$$



$$w = u + v \quad u \in \text{NS}(\phi), v \in \text{RS}(\phi)$$

$$\min_{\substack{u \in \text{NS}(\phi) \\ v \in \text{RS}(\phi)}} f(\Phi(u+v)) + \frac{\gamma}{2} \|u+v\|^2$$



$$\min_{v \in \text{RS}(\phi)} f(\Phi v) + \frac{\gamma}{2} \|u\|^2 + \frac{\gamma}{2} \|v\|^2$$



# Representer theorem (contd.)

$$\min_{v \in \text{RS}(\Phi)} f(\bar{\Phi}v) + \frac{\lambda}{2} \|v\|^2$$

$$\Downarrow$$
$$\min_{\alpha \in \mathbb{R}^n} f(\bar{\Phi} \bar{\Phi}^T \alpha) + \frac{\lambda}{2} \|\bar{\Phi}^T \alpha\|^2$$

$$\Downarrow$$
$$\min_{\alpha \in \mathbb{R}^n} f(K \alpha) + \frac{\lambda}{2} \alpha^T K \alpha$$

## Gradient descent for kernel Logistic Regression:

$$R(\alpha) = \sum_{i=1}^n \log(1 + \exp(-y_i \alpha^T K_{:,i})) + \frac{\lambda}{2} \alpha^T K \alpha$$

$$\begin{aligned}\nabla R(\alpha) &= \sum_{i=1}^n \frac{\exp(-y_i \alpha^T K_{:,i})}{1 + \exp(-y_i \alpha^T K_{:,i})} \cdot (-y_i; K_{:,i}) \\ &\quad + \lambda K \alpha\end{aligned}$$

$$:= \sum_{i=1}^n \sigma(-y_i \alpha^T K_{:,i}) (-y_i; K_{:,i}) + \lambda K \alpha$$

$$\alpha^* = \min_{\alpha} R(\alpha)$$

### Gradient Descent:

1. Initialise  $\alpha \in \mathbb{R}^n$ , choose  $\eta > 0$ .
2. Repeat:
  - (a)  $\alpha \leftarrow \alpha - \eta \nabla R(\alpha)$
  - (b) Stop if  $R(\alpha)$  has not decreased much in the last few iterations.
3. Return  $\alpha$ .

[HR]

Predicting the class at a given  $x$ :

$$\text{Predicted value of } P(y=1 | x=x) = \sigma(w^* \phi(x))$$

$$\begin{aligned} \sigma(w^* \phi(x)) &= \sigma(\alpha^* \bar{\Phi} \phi(x)) \\ &= \sigma\left(\sum_{i=1}^n \alpha_i^* \kappa(x_i, x)\right) \end{aligned}$$

$$\text{Class Predicted : sign } (\sigma(w^* \phi(x)) - 0.5)$$

$$= \text{Sign}(w^* \phi(x))$$

# Brief Mention: Kernel (linear) regression

Representer Theorem:

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|\Phi w - y\|^2 + \frac{\lambda}{2} \|w\|^2$$

$\Phi \in \mathbb{R}^{n \times d}$   
 $\Phi^\top \alpha = \sum_{i=1}^n \alpha_i \phi(x_i)$

$\Updownarrow$  RT

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{2} \|\Phi \Phi^\top \alpha - y\|^2 + \frac{\lambda}{2} \|\Phi^\top \alpha\|^2$$

Why? Any  $w \in \mathbb{R}^d$  =  $v + u$  for some  $u \in \text{Null}(\Phi)$   
 $v \in \text{Row}(\Phi)$

$$\text{and } \|w\|^2 = \|v\|^2 + \|u\|^2$$

$$\text{Let } v = \Phi^\top \alpha$$

(will be explored further in Assignment)

[HR]

# Kernel functions and examples

$k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is called Kernel function,  
if  $k(u, v)$  represents  $\phi(u)^T \phi(v)$  for some  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$

Example  
Kernel :

$$K(u, v) = \exp\left(-\frac{\|u-v\|^2}{\sigma^2}\right)$$

$\begin{matrix} u & v \\ \mathbb{R}^d & \mathbb{R}^d \end{matrix}$

$$\exists \phi: \mathbb{R}^d \rightarrow \mathcal{H}$$

$$K(u, v) = \langle \phi(u), \phi(v) \rangle$$

Another  
Example Kernel:

$$K(u, v) = (u^T v)^7$$

## Illustration of Kernel Function:

To show:  $K(u, v) = (u^T v)^3$

is a valid kernel.

$$X = \begin{bmatrix} t_1 & t_2 \\ u_1 & u_2 \\ v_1 & v_2 \end{bmatrix} ; \Phi = \begin{bmatrix} t_1^3 & \sqrt{3}t_1^2 t_2 & \sqrt{3}t_1 t_2^2 & t_2^3 \\ u_1^3 & \sqrt{3}u_1^2 u_2 & \sqrt{3}u_1 u_2^2 & u_2^3 \\ v_1^3 & \sqrt{3}v_1^2 v_2 & \sqrt{3}v_1 v_2^2 & v_2^3 \end{bmatrix} = \begin{bmatrix} \phi(t)^T \\ \phi(u)^T \\ \phi(v)^T \end{bmatrix}$$

2d data & homogenous 3<sup>rd</sup> degree Polynomial regression

3 data points .

$$\phi(u) = [u^3, \sqrt{3}u^2, u^2, \sqrt{3}u, u^2, u^3]$$

$\phi$   
x ...

$$\Phi \Phi^T = \begin{bmatrix} \phi(t)^T \phi(t) & \phi(t)^T \phi(u) & \phi(t)^T \phi(v) \\ \phi(u)^T \phi(t) & \phi(u)^T \phi(u) & \phi(u)^T \phi(v) \\ \phi(v)^T \phi(t) & \phi(v)^T \phi(u) & \phi(v)^T \phi(v) \end{bmatrix}$$

$$\phi(u)^T \phi(v) = u_1^3 v_1^3 + 3u_1^2 u_2 v_1^2 v_2 + 3u_1 u_2^2 v_1 v_2^2 + u_2^3 v_2^3$$

$$= (u_1 v_1 + u_2 v_2)^3$$

$$= (u^T v)^3 = k(u, v)$$

↳

Homogenous 3<sup>rd</sup> degree Polynomial  
Kernel.

(i.e.)

you can compute  $\phi(u)^T \phi(v)$   
without computing  $\phi(u)$  &  $\phi(v)$

Qn: What functions  $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  are valid Kernel functions?

Exercise: Identify if following fns are kernels:

i) These kernels are for  $d=1$

$$(a) K(x, y) = x + y$$

$$(b) K(x, y) = x - y$$

$$(c) K(x, y) = xy$$

$$(d) K(x, y) = x^2$$

$$(e) K(x, y) = x^2 + y^2$$

2) These kernels are for  $d>1$

$$(a) K(x, y) = x^T y$$

$$(b) K(x, y) = \sum_{i=1}^d x_i^2 y_i^2$$

$$(c) K(x, y) = \exp(x^T y) \quad (e) K(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$$

$$(d) K(x, y) = (1+x^T y)^p$$

Solution:

i) Find  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^d$ , s.t  $K(x, y) = \phi(x)^T \phi(y)$

to show  $K$  is a valid kernel.

ii) For all sets of data points  $x_1, \dots, x_n$  with  $x_i \in \mathbb{R}^d$ ,  
the  $n \times n$  Kernel matrix must be symmetric and PSD for  $K$  to be a valid kernel.

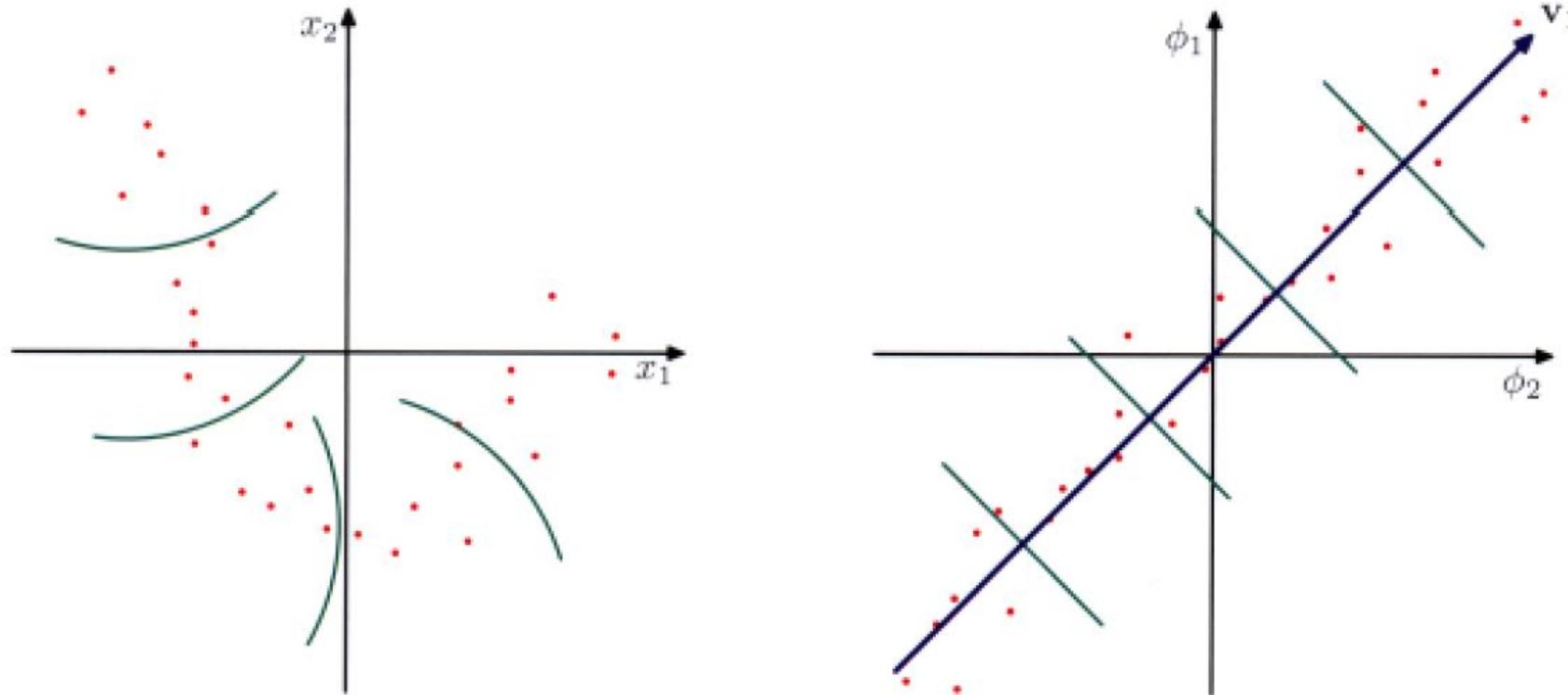
# Concluding thoughts

- Linear models for classification seen from all three approaches - generative/discriminative/direct.
  - Has similar pros/cons as linear regression - non-linear basis functions, along with the kernel method/trick, can allow for non-linear boundaries between classes, but still suffers from fundamental limitation of fixed basis functions (fixed before seeing the training data).
- Next steps:
  - Non-linear models of classification: Use training data that starts with a set of non-linear basis functions, and either selects a subset of it (kernel SVM; convex objective fn.) or adapts their parameters (multi-layer perceptron or neural network; non-convex obj.).

Thank you! (Appendix follows)

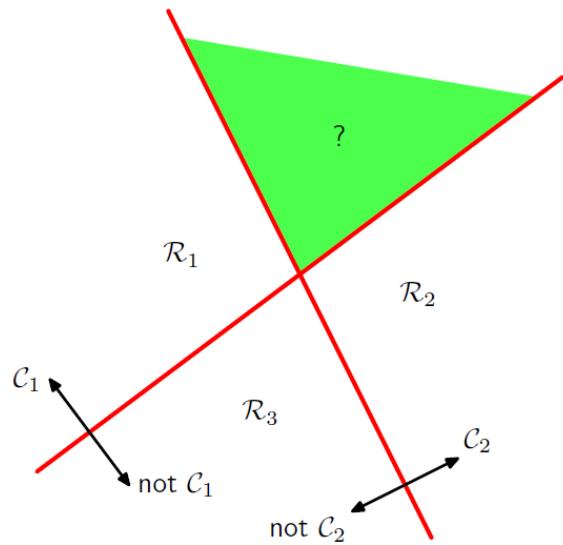
# Appendix

# Kernel PCA

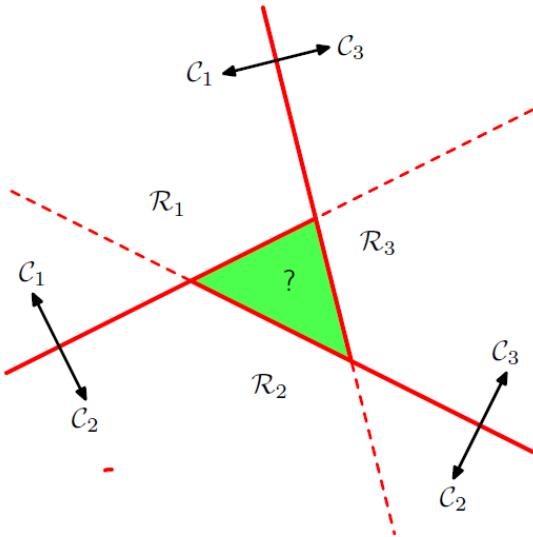


**Figure 12.16** Schematic illustration of kernel PCA. A data set in the original data space (left-hand plot) is projected by a nonlinear transformation  $\phi(x)$  into a feature space (right-hand plot). By performing PCA in the feature space, we obtain the principal components, of which the first is shown in blue and is denoted by the vector  $v_1$ . The green lines in feature space indicate the linear projections onto the first principal component, which correspond to nonlinear projections in the original data space. Note that in general it is not possible to represent the nonlinear principal component by a vector in  $x$  space.

Discriminant (DBs) for  $K > 2$  classes: Simple 2-way  $\rightarrow$  multi-way constructions don't work! However, multi-way  $\rightarrow$  2-way works!

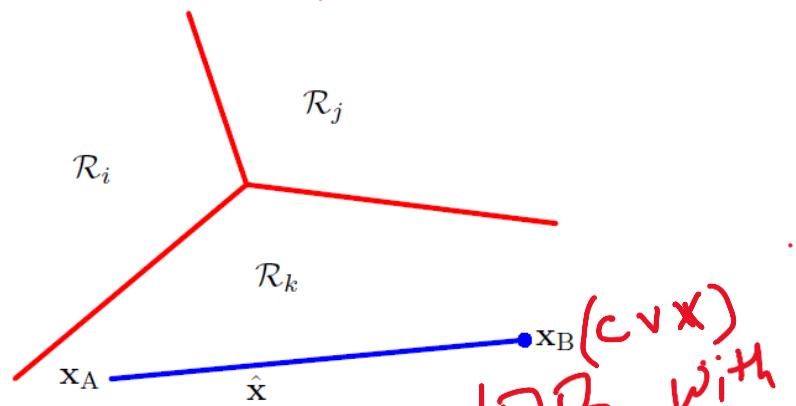


1-vs-rest



$c_i \text{ vs } c_j$

$$\begin{aligned} z_k(\mathbf{x}) &= \mathbf{w}_k^T \mathbf{x} + w_{k0} \\ k(x) &= \arg \max_{k \in [K]} z_k(\mathbf{x}) \end{aligned}$$

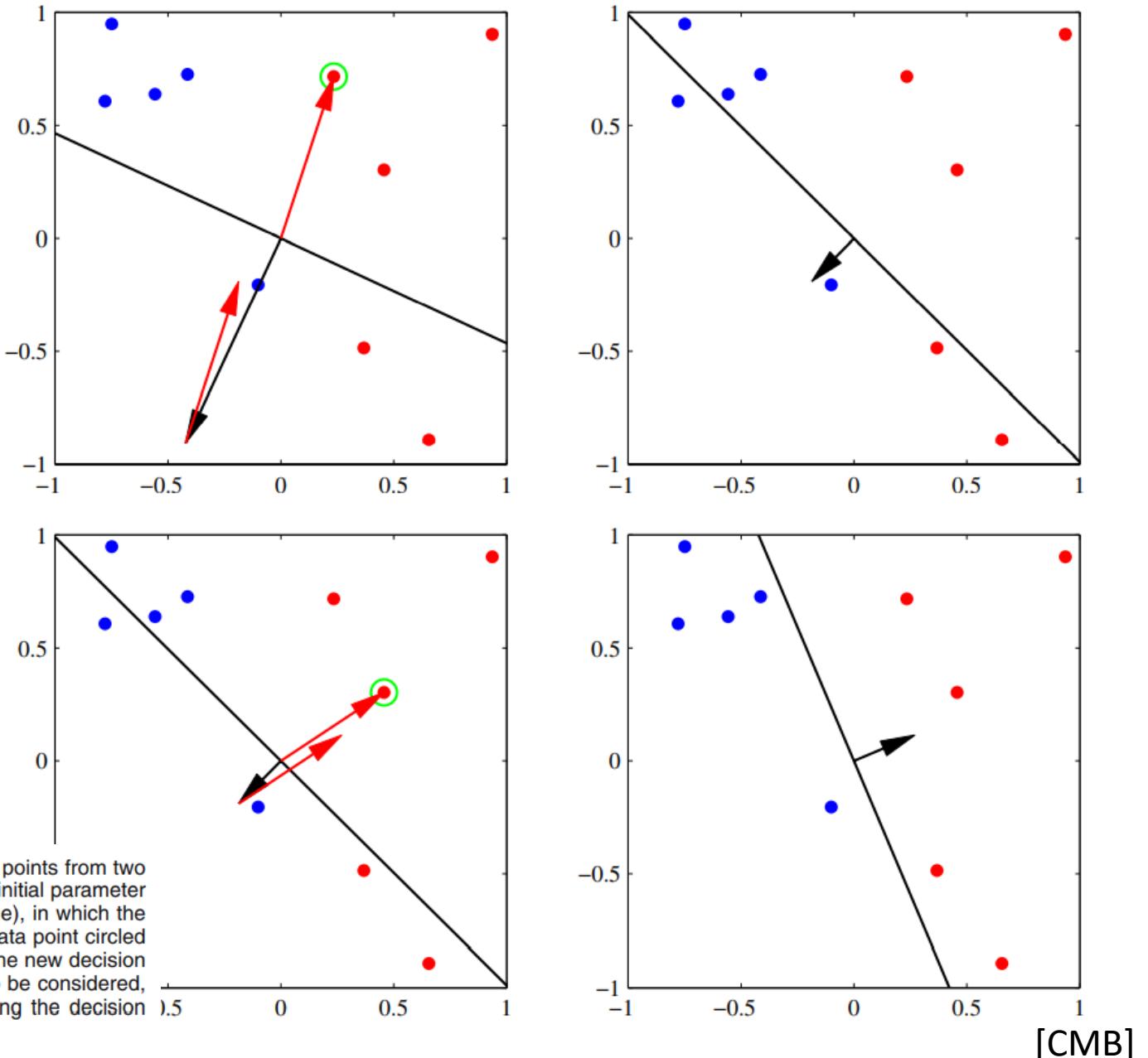


$z_i(\mathbf{x}) = z_j(\mathbf{x})$  DS/DB with similar geom. [CMB]

# Perceptron training (learning $\mathbf{w}$ )

- Stochastic gradient descent:  $\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi_n t_n$ 
  - Only interested in the direction of  $w$ . So  $\eta$  can be set to 1.
  - That is, perceptron function  $y(x, w)$  unchanged if  $w$  scaled by a constant.
- Then, for a single misclassified data point, add  $\phi(x_n)$  to  $\mathbf{w}$  if  $t = +1$ , and subtract it from  $\mathbf{w}$  o.w.
  - This will lead to decrease in error for a single update, because:
$$-\mathbf{w}^{(\tau+1)T} \phi_n t_n = -\mathbf{w}^{(\tau)T} \phi_n t_n - (\phi_n t_n)^T \phi_n t_n < -\mathbf{w}^{(\tau)T} \phi_n t_n$$
  - If data is linearly separable, proof of convergence in finite steps exists.
- Only of historical importance, as it suffers from many issues: training doesn't converge for non-linearly separable dataset, multiple decision surfaces possible for linearly separable datasets, convergence in practice may take very many steps, etc.

# Perceptron training in action



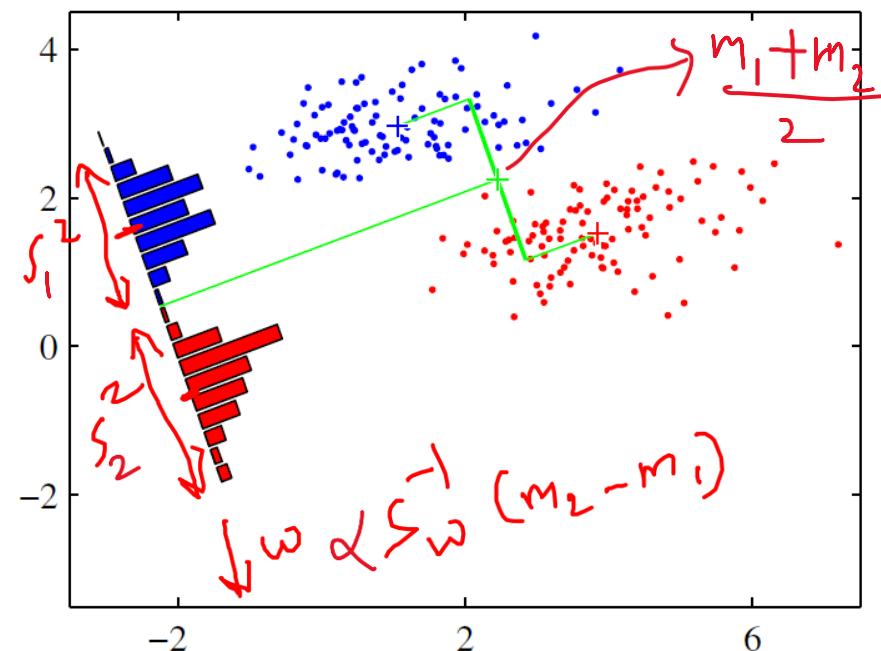
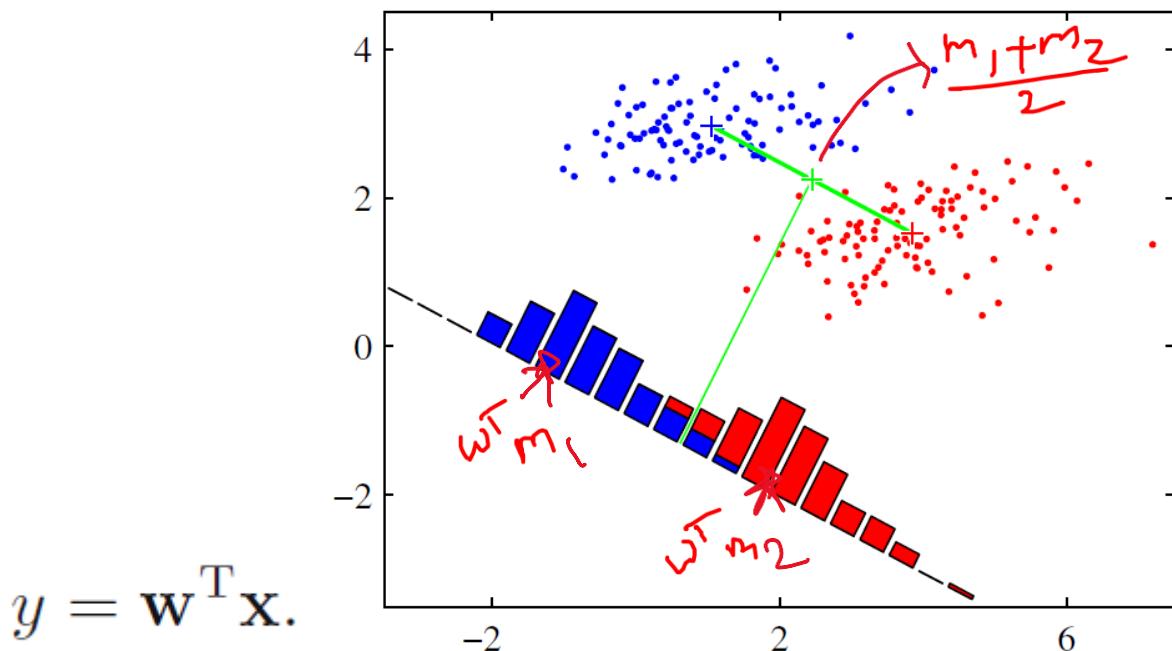
**Figure 4.7** Illustration of the convergence of the perceptron learning algorithm, showing data points from two classes (red and blue) in a two-dimensional feature space  $(\phi_1, \phi_2)$ . The top left plot shows the initial parameter vector  $w$  shown as a black arrow together with the corresponding 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 plot for which all data points are correctly classified.

# Outline for Module M5

- M5. Classification (Linear models)
  - M5.0 Introduction
  - M5.1 Generative model ((Naïve) Bayes classifier)
  - M5.2 Discriminative model (Logistic regression)
  - **M5.3 Discriminant function**
    - Geometry of decision surfaces
    - Perceptron objective function
    - Linear/Fisher's Discriminant Analysis
  - M5.4 Kernel methods
  - M5.5 Concluding thoughts

# Fisher's discriminant analysis

- Reduce points to one dimension using a linear function, and then apply a threshold (step fn.  $f(\cdot)$ ) to classify as in perceptron).
- So same choice of  $f(\cdot)$ ,  $f(\mathbf{w}^T \mathbf{x})$  as in perceptron, but choice of objective function to learn  $\mathbf{w}$  differs (between/within class variance).



[CMB]

# FDA - Choice of obj. fn. – attempt #1

Maximize  $\mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1)$  s.t.  $\|\mathbf{w}\| = 1$

where  $\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n$ ,  $\mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n$ .

yields  $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$

# FDA - Choice of obj. fn. – attempt #2

Maximize  $J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}.$

$\frac{\mathcal{B} \subset \mathcal{W}}{\mathcal{W} \subset \mathcal{W}}$

$$m_k = \mathbf{w}^T \mathbf{m}_k = y(m_k)$$

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2$$

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n. \quad y(x) = \mathbf{w}^T \mathbf{x}$$

[CMB]

# Maximizing $J(\mathbf{w})$

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

where  $\mathbf{S}_B$  is the *between-class* covariance matrix and is given by

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \quad \text{d} \times \text{d} \quad (4.27)$$

and  $\mathbf{S}_W$  is the total *within-class* covariance matrix, given by

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T. \quad (4.28)$$

Differentiating (4.26) with respect to  $\mathbf{w}$ , we find that  $J(\mathbf{w})$  is maximized when

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}. \quad (4.29)$$

# Maxima of $J(w)$ -- identical to generative model (Bayes classifier; cf. next slide)!

From (4.27), we see that  $S_B w$  is always in the direction of  $(m_2 - m_1)$ . Furthermore, we do not care about the magnitude of  $w$ , only its direction, and so we can drop the scalar factors  $(w^T S_B w)$  and  $(w^T S_W w)$ . Multiplying both sides of (4.29) by  $S_W^{-1}$  we then obtain

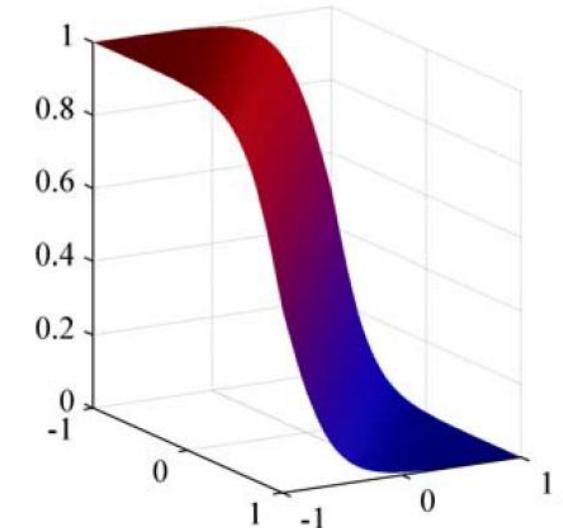
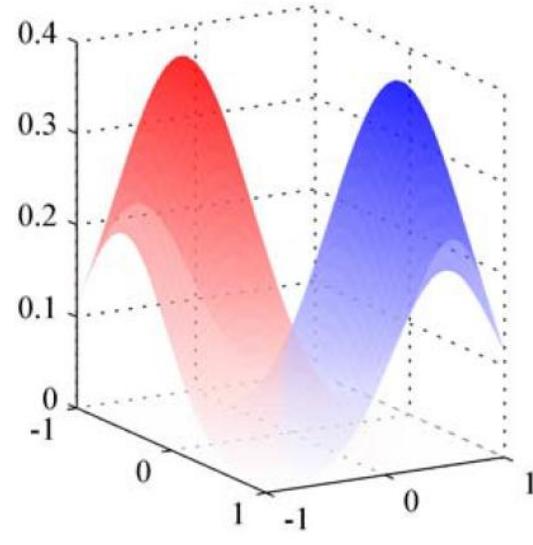
$$w \propto S_W^{-1} (m_2 - m_1). \quad (4.30)$$

- What if data is axis-aligned? (think about identity as well as diagonal  $S_w$ )
- What to do once you've projected data points i.e., what is a good threshold to apply?


$$w^T \underbrace{(m_2 + m_1)}_{2} \quad ( := -w_0 )$$

Recall: K=2 Bayes classifier – Gaussian, general (shared) covariance

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}. \quad p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$



$$\begin{aligned} \mathbf{w} &= \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ w_0 &= -\frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}. \end{aligned}$$

[CMB]

# A comparison: PCA vs FDA

