



LUND  
UNIVERSITY

## L-4: Linear classification methods

Ted Kronvall, Cristian Sminchisescu



# Admin

## Moodle:

- Register for "Machine Learning 2019" in Moodle by the 14th of April.
- Use enrollment key "clustering2019".
- Sign up for the Q & A sessions for Assignment 1 - the first one is tomorrow!



LUND  
UNIVERSITY

# Small recap from L-3

- In convex optimization, local minima are global. This is not true in general for non-convex problems.
- MSE, LASSO, ridge regression, SVM, logistic regression are convex problems, neural network training is not.
- One may check if a problem is convex using the definition, the second-order derivative, the epigraph, or by identifying it as a composition of convex functions.
- Optimization problems may be either unconstrained or constrained.
  - For an unconstrained problem, the optimal point has its derivative equal to zero.
  - For a constrained problem, the derivative of the Lagrangian function w.r.t. the primal variables is zero at the optimal point.
- The KKT conditions are necessary and sufficient for a point to be optimal in constrained optimization.
- Coordinate descent updates the coordinates in the variable one at a time, without needing to calculate the derivative.
- Gradient and steepest descent use the first-order derivative. Gradient descent is the most common approach in deep learning.

# Today

- ~~Basic concepts; tasks, training, validation, regularization.~~
- ~~Probability theory and optimization~~
- ~~Linear regression~~



- Linear methods for classification



- Clustering methods
- Neural networks
- Convolutional and recurrent NNs



## Generative modeling



- Reinforcement learning
- Natural language processing
- Deep learning for computer vision



LUND  
UNIVERSITY

# The classification model



LUND  
UNIVERSITY

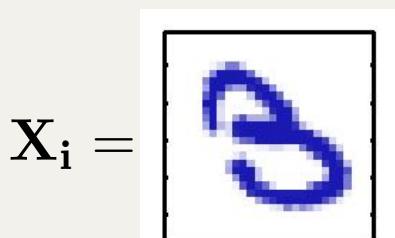
# Classification model

## NOTATIONAL CONVENTIONS

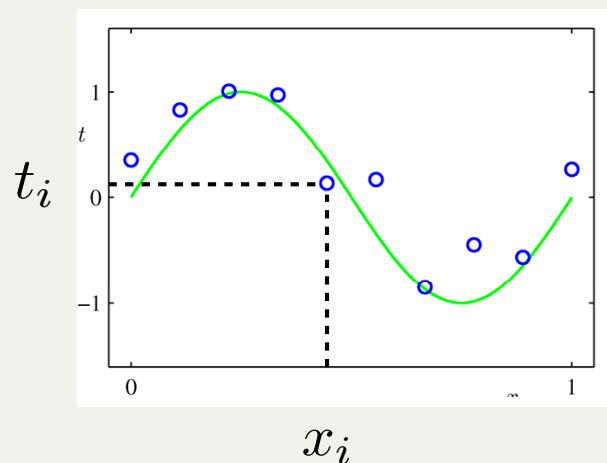
- Recall how notation is used in machine learning problems:

$t_i, \mathbf{t}, \mathbf{T}$	Targets	= Data
$x_i, \mathbf{x}, \mathbf{X}$	Inputs	
$y_i, \mathbf{y}, \mathbf{Y}$	Model output	
$w_i, \mathbf{w}, \mathbf{W}$	trainable variables	
$\lambda, \mu, \rho, \dots$	(capacity) hyperparameters	

no targets = unsupervised data

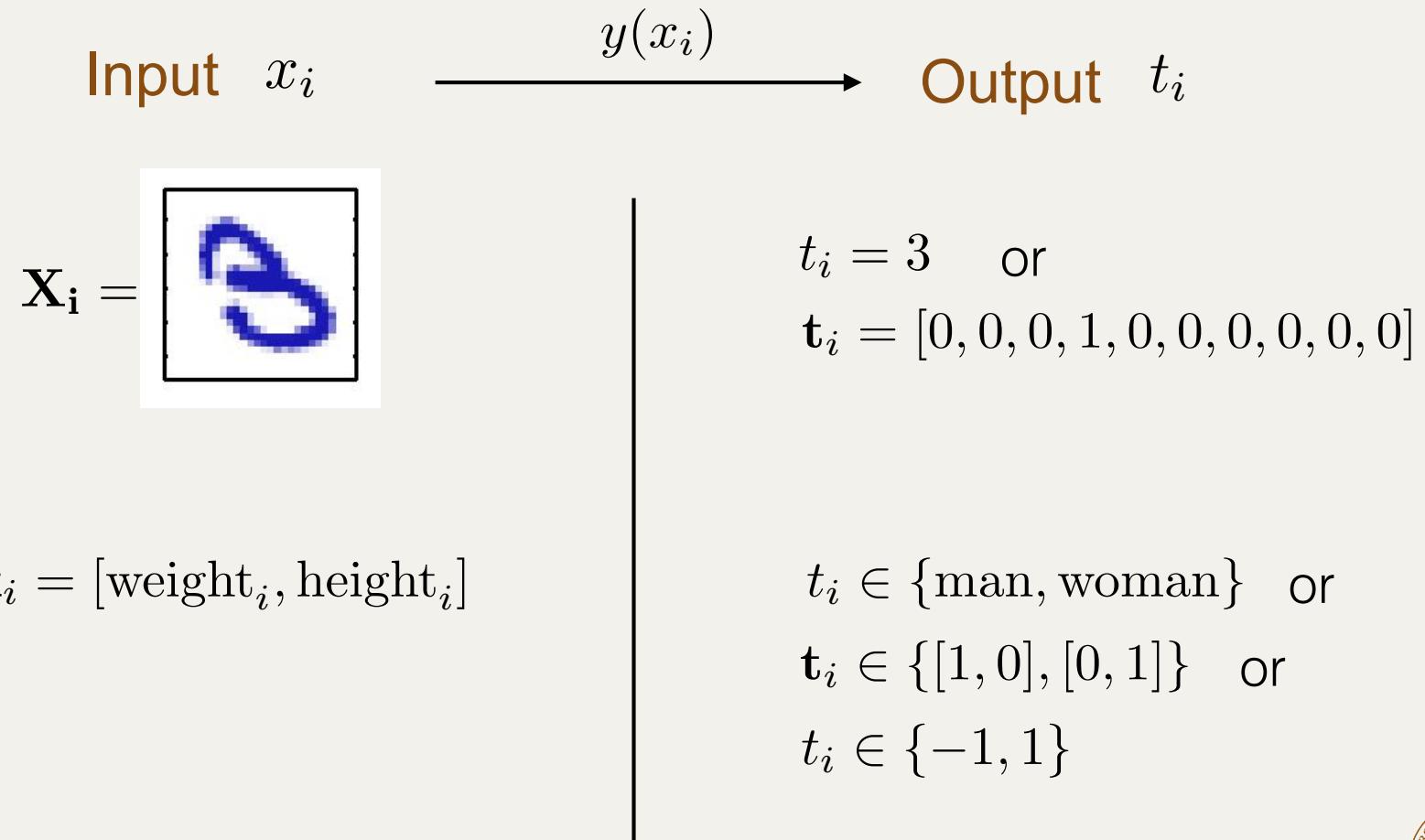


$$t_i = 3$$



LUND  
UNIVERSITY

# Classification model



LUND  
UNIVERSITY

# Classification model

## Linear Classification

- Classification is intrinsically non-linear because of the training constraints that ***place non-identical inputs in the same class***
- Differences in the input vector sometimes causes 0 change in the answer
- Linear classification means that the adaptive part  $\mathbf{w}$  is linear
  - The adaptive part is cascaded with a fixed non-linearity  $f$
  - It may also be preceded by a fixed non-linearity  $\phi$  when nonlinear basis functions are used

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + w_0, \quad c = f(y(\mathbf{x}))$$

fixed non-linear functions

adaptive linear parameters      decision



LUND  
UNIVERSITY

# Classification model

## Approach 1: Discriminant Function

- Use discriminant functions directly, and do not compute probabilities
- Convert the input vector into one or more real values so that a simple process (thresholding, or a majority vote) can be applied to assign the input to the class
- The real values should be chosen to maximize the useable information about the class label present in the real value
- Given discriminant functions  $f_1(\mathbf{x}), \dots, f_K(\mathbf{x})$   
Classify  $\mathbf{x}$  as class  $C_k$ , iff  $f_k(\mathbf{x}) > f_j(\mathbf{x}), \forall j \neq k$



LUND  
UNIVERSITY

# Classification model

## Approach 2: Class-conditional Probabilities

- Infer conditional class probabilities  $p(C_k|\mathbf{x})$
- Use conditional distribution to make optimal decisions, e.g. by minimizing some loss function
- Example, 2 classes

$$p(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \phi), \quad p(C_2|\mathbf{x}) = 1 - p(C_1|\mathbf{x})$$
$$\text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$



LUND  
UNIVERSITY

# Classification model

## Represent Target Values: *Binary* vs. *Multiclass*

- **Two classes ( $N=2$ ):** typically use a single real valued output that has target values of 1 for the *positive class* and 0 (or -1) for the *negative class*  $t_i \in \{-1, 1\}$ 
  - For probabilistic class labels, the target can be the probability of the positive class and the output of the model can be the probability the model assigns to the positive class
- **For the multiclass ( $N>2$ ),** we use a vector of  $N$  target values containing a single 1 for the correct class and zeros elsewhere
  - For probabilistic labels we can then use a vector of class probabilities as the target vector

$$\mathbf{t}_i = [0, 0, 0, 1, 0, 0, 0, 0, 0]$$



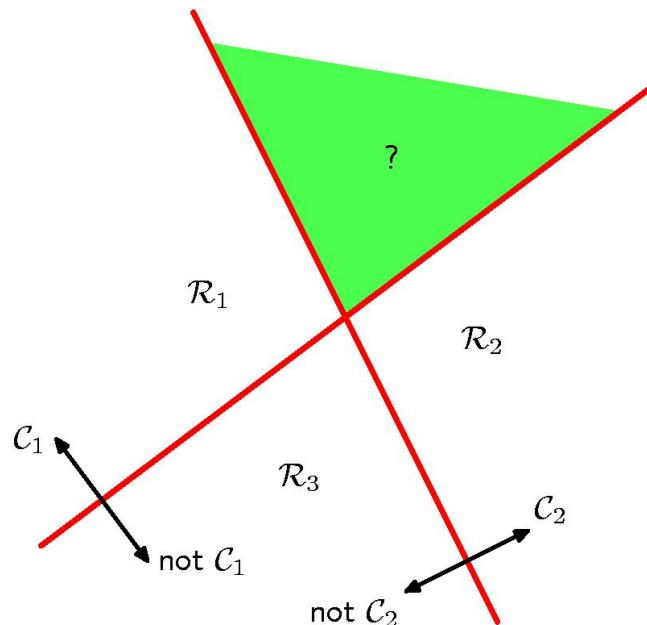
LUND  
UNIVERSITY

## Discriminant Functions for Multiclass

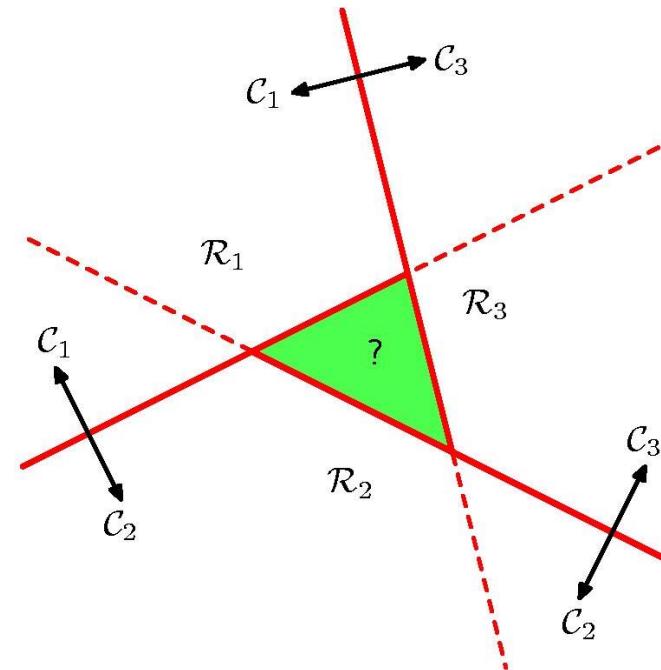
- One possibility is to use  $N$  binary (2-way) discriminants
  - Each function separates one class from the rest
- Another possibility is to use  $\frac{N(N-1)}{2}$  binary (2-way) discriminants
  - Each function discriminates between two specific classes. We have 1 discriminant for each class pair
- Both methods have ambiguities



# Problems with Multi-class Discriminant Functions Constructed from Binary Classifiers



**1 - vs. - all**



**1 vs. 1**

If we base our decision on binary classifiers, we can encounter ambiguities



**LUND**  
UNIVERSITY

# Simple Solution

Use  $N$  discriminant functions,  $y_i, y_j, y_k, \dots$ , and  
take the **max** over their response. Why?

- Consider linear discriminants  $y$

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

- The decision boundary between class  $k$  and  $j$   
is given by the  $D - 1$  hyperplane

$$(\mathbf{w}_k^T - \mathbf{w}_j^T) \mathbf{x} + (w_{k0} - w_{j0}) = 0$$

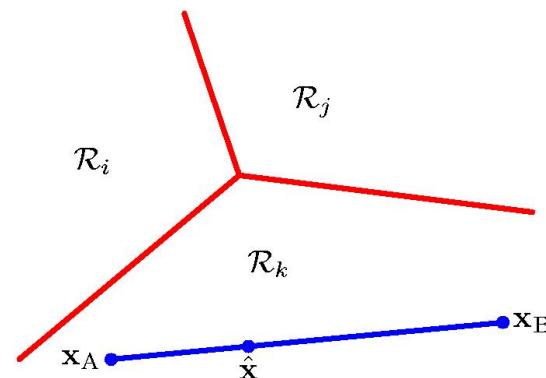
*In this linear case the decision regions are convex*

$$\mathbf{x}_A, \mathbf{x}_B \in R_K, \hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B, 0 \leq \lambda \leq 1$$

$$\text{From the linearity of } y \Rightarrow y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B)$$

$$\begin{aligned} \text{But } y_k(\mathbf{x}_A) &> y_j(\mathbf{x}_A) \text{ and } y_k(\mathbf{x}_B) > y_j(\mathbf{x}_B) \quad \forall j \neq k \Rightarrow y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}}) \\ &\Rightarrow \hat{\mathbf{x}} \text{ also lies inside } R_k \end{aligned}$$

Hence  $R_k$  is convex



LUND  
UNIVERSITY

# Least squares for classification



LUND  
UNIVERSITY

# Classification with LS

## Least Squares Classification

Assume each class has its own linear model:  $y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$

Then we can write:

$y(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}$ , with  $k$ -th column a  $D + 1$ -dim vector  $\tilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T$ ,  $\tilde{\mathbf{x}} = (1, \mathbf{x}^T)^T$   
*decision surfaces are  $D$ -dimensional hyperplanes passing through the origin of the  $D + 1$ -dimensional expanded input space.*

Given  $\{\mathbf{x}_n, \mathbf{t}_n\}, n = 1, \dots, N$ ;  $n$ -th row of  $\mathbf{T}$  is  $\mathbf{t}_n^T$ ;  $\tilde{\mathbf{X}}$ 's  $n$ -th row is  $\tilde{\mathbf{x}}_n^T$

The sum of squares error function for classification is:

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr}\{(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})\}$$
$$\frac{\partial E_D}{\partial \tilde{\mathbf{W}}} = 0 \Rightarrow \tilde{\mathbf{W}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{T} = \tilde{\mathbf{X}}^+ \mathbf{T} \quad \text{X}^+ \text{ is the pseudoinverse of X}$$
$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \mathbf{T}^T (\tilde{\mathbf{X}}^+)^T \tilde{\mathbf{x}} \quad \leftarrow \text{Closed-form solution}$$

**Property:** every vector in the training set and the model prediction for any value of  $\mathbf{x}$ , satisfy some linear constraint:

$$\mathbf{a}^T \mathbf{t}_N + b = 0, \mathbf{a}^T \mathbf{y}(\mathbf{x}) + b = 0, \text{for some constants } \mathbf{a}, b.$$



LUND  
UNIVERSITY

# Classification with LS

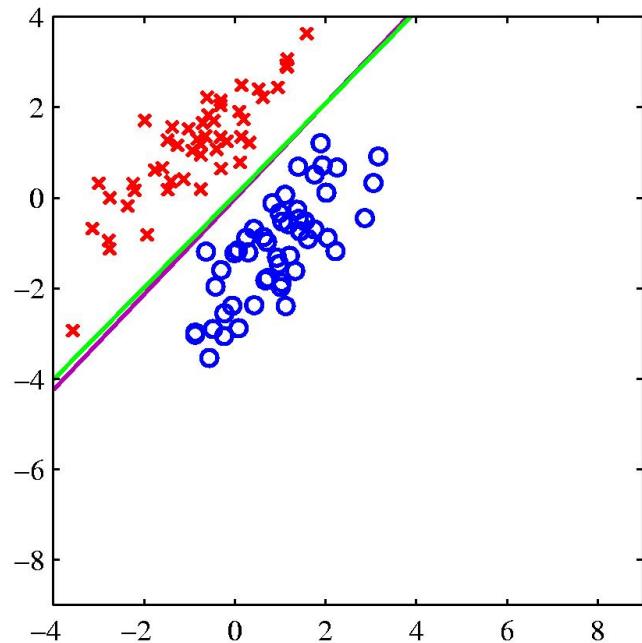
[Matlab example]



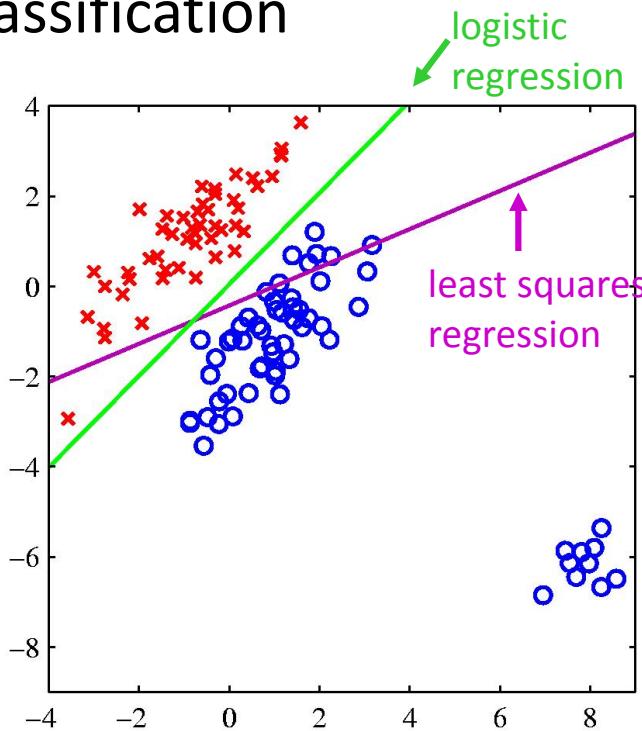
LUND  
UNIVERSITY

# Classification with LS

## Problems with using least squares for classification



Least squares solutions lack robustness to outliers



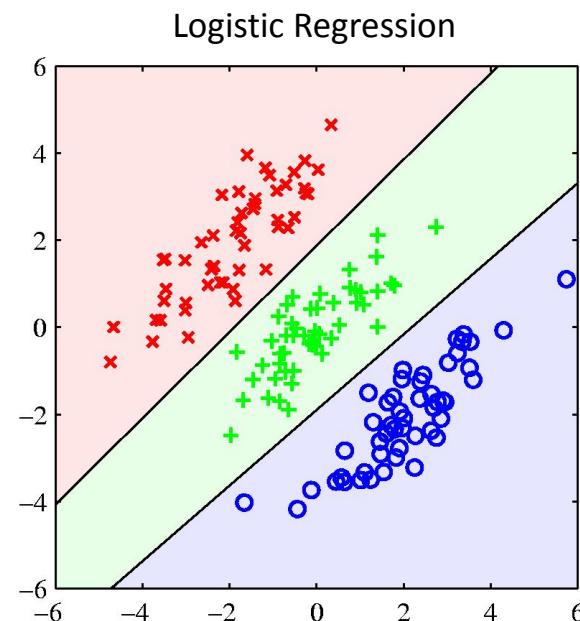
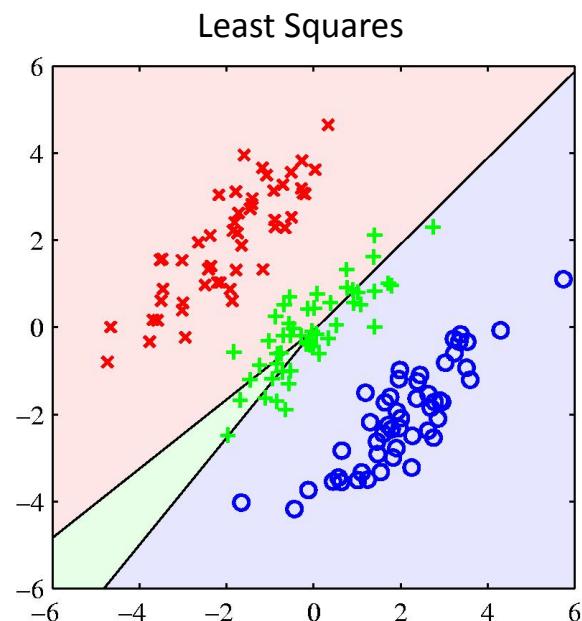
If the right answer is 1 and the model says 1.5, it loses, so it changes the boundary to avoid being “too correct”



LUND  
UNIVERSITY

# Classification with LS

For non-Gaussian targets, least squares regression gives poor decision surfaces



Remember that least squares corresponds to the Maximum Likelihood under a Gaussian conditional distribution  
Clearly the binary target vectors have a distribution that is far from Gaussian



LUND  
UNIVERSITY

# Fisher's linear discriminant (LDA)



LUND  
UNIVERSITY

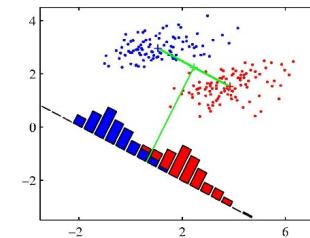
# Linear discriminant analysis

## Fisher's Linear Discriminant

- We can view classification in terms of dimensionality reduction
- A simple linear discriminant function is a projection of the  $D$  – dimensional data  $\mathbf{x}$  down to 1 dimension

Project:  $y = \mathbf{w}^T \mathbf{x}$ ;

Classify: if  $y \geq -w_0$  then  $C_1$  else  $C_2$



- However projection results in loss of information. Classes well separated in the original input space may strongly overlap in 1d
- We will adjust the projection (weight vector  $\mathbf{w}$ ) to achieve the best separation among classes. But what do we mean by *best separation*?



LUND  
UNIVERSITY

# Linear discriminant analysis

## Fisher's View of Class Separation (I)

- The simplest measure of class separation when projected onto  $\mathbf{w}$  is the separation of the projected class means. This suggests choosing  $\mathbf{w}$  so to maximize

$$m_2 - m_1 = \mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1), m_k = \mathbf{w}^T \mathbf{m}_k, \mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}_n, \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}_n$$

- This can be made arbitrarily large by increasing  $||\mathbf{w}||$ . We could handle this by imposing unit norm constraints using Lagrange multipliers. We get

$$\max_{\mathbf{w}} \mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1), \text{ s.th. } ||\mathbf{w}|| = 1$$

$$\Rightarrow \mathbf{w} = -\frac{1}{2\lambda}(\mathbf{m}_2 - \mathbf{m}_1) \propto \mathbf{m}_2 - \mathbf{m}_1$$

- However, if the main direction of variance in each class is not orthogonal to the direction between means, we will not get good separation (see next slide)



LUND  
UNIVERSITY

# Linear discriminant analysis

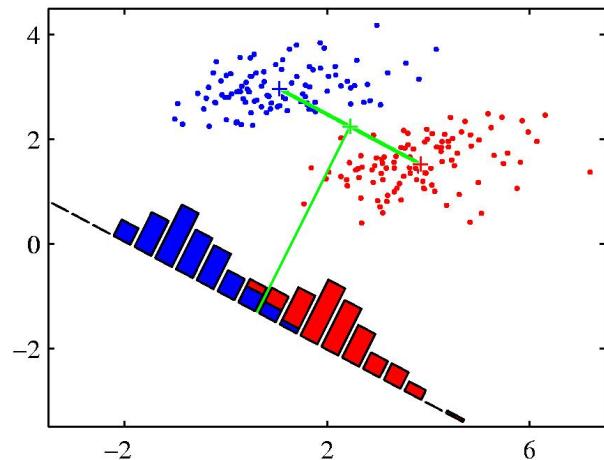
[Matlab example]



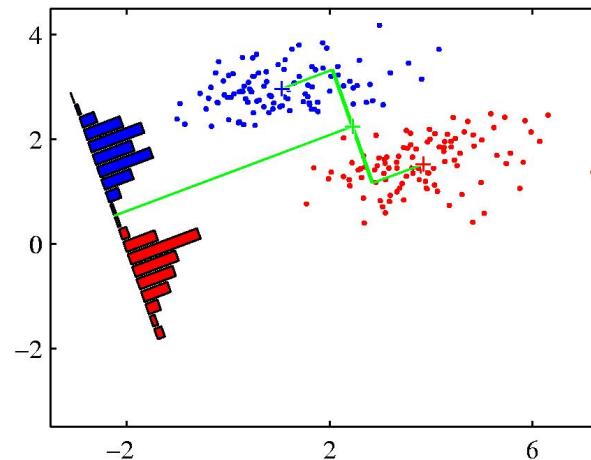
LUND  
UNIVERSITY

# Linear discriminant analysis

## Advantage of using Fisher's Criterion



When projected onto the line joining the class means, the classes are not well separated



Fisher chooses a direction that makes the projected classes much tighter, even though their projected means are less far apart



LUND  
UNIVERSITY

# Linear discriminant analysis

## Fisher's View of Class Separation (II)

- Fisher
  - maximize a function that gives a large separation between the projected class means,
  - while also giving a small variance within each class
  - thereby minimizing class overlap
  - Choose direction maximizing the **ratio** of **between class** variance to **within class** variance
  - This is the direction in which the projected points contain the most information about class membership (under Gaussian assumptions)



LUND  
UNIVERSITY

# Linear discriminant analysis

## Fisher's Linear Discriminant

- We seek a linear transformation that is best for discrimination
- The projection onto the vector separating the class means seems right

$$\mathbf{w} \propto \mathbf{m}_2 - \mathbf{m}_1$$

- But we also want small variance within each class

$$s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2, \quad m_k = \mathbf{w}^T \mathbf{m}_k$$

- Fisher's objective function

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

← Between class  
← Within class



LUND  
UNIVERSITY

# Linear discriminant analysis

## Criterion Derivation

$$\begin{aligned}(m_2 - m_1)^2 &= (\mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1))^2 \\&= \mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T\mathbf{w} \\&= \mathbf{w}^T\mathbf{S}_B\mathbf{w}.\end{aligned}$$

$$\begin{aligned}s_1^2 + s_2^2 &= \sum_{n \in \mathcal{C}_1} (y_n - m_1)^2 + \sum_{k \in \mathcal{C}_2} (y_k - m_2)^2 \\&= \sum_{n \in \mathcal{C}_1} (\mathbf{w}^T(\mathbf{x}_n - \mathbf{m}_1))^2 + \sum_{k \in \mathcal{C}_2} (\mathbf{w}^T(\mathbf{x}_k - \mathbf{m}_2))^2 \\&= \sum_{n \in \mathcal{C}_1} \mathbf{w}^T(\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T\mathbf{w} \\&\quad + \sum_{k \in \mathcal{C}_2} \mathbf{w}^T(\mathbf{x}_k - \mathbf{m}_2)(\mathbf{x}_k - \mathbf{m}_2)^T\mathbf{w} \\&= \mathbf{w}^T\mathbf{S}_W\mathbf{w}.\end{aligned}$$



LUND  
UNIVERSITY

# Linear discriminant analysis

## Fisher's Linear Discriminant Derivations

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

where

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

$$\mathbf{S}_W = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T$$



LUND  
UNIVERSITY

# Linear discriminant analysis

## Fischer's Linear Discriminant Computation

However, the objective  $J(\mathbf{w})$  is invariant to rescaling  $\mathbf{w} \rightarrow \alpha\mathbf{w}$ . We can choose the denominator to be unity. We can then minimize

$$\min_{\mathbf{w}} -\frac{1}{2}\mathbf{w}^T \mathbf{S}_B \mathbf{w}$$

$$\mathbf{w}^T \mathbf{S}_W \mathbf{w} = 1$$

This corresponds to the (primal) Lagrangian

$$L_P = -\frac{1}{2}\mathbf{w}^T \mathbf{S}_B \mathbf{w} + \frac{1}{2}\lambda(\mathbf{w}^T \mathbf{S}_W \mathbf{w} - 1)$$

From the *KKT conditions*

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w} \Rightarrow \mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}$$

Generalized eigenvalue problem, as  $\mathbf{S}_W^{-1} \mathbf{S}_B$  not symmetric



LUND  
UNIVERSITY

# Linear discriminant analysis

[Matlab example]



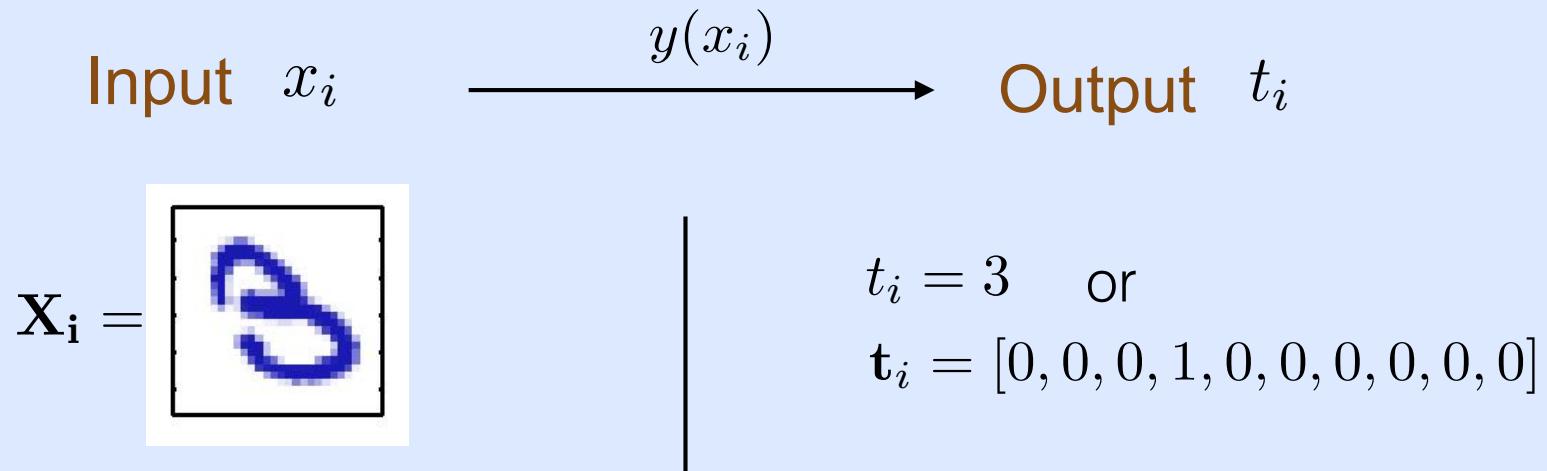
LUND  
UNIVERSITY

# Logistic regression



LUND  
UNIVERSITY

# Logistic regression



- Let's model the outputs as the probabilities for each class!
- One class can be removed, as probabilities should sum to one.



LUND  
UNIVERSITY

# Logistic regression

## Probabilistic Generative Models

- Use a class prior and a separate generative model of the input vectors for each class, and compute which model makes a test input vector most probable
- The posterior probability of class 1 is given by:

$$p(C_1 | \mathbf{x}) = \frac{p(C_1)p(\mathbf{x} | C_1)}{p(C_1)p(\mathbf{x} | C_1) + p(C_2)p(\mathbf{x} | C_2)} = \boxed{\frac{1}{1+e^{-a}}}$$

where  $a = \ln \frac{p(C_1)p(\mathbf{x} | C_1)}{p(C_2)p(\mathbf{x} | C_2)} = \boxed{\ln \frac{p(C_1 | \mathbf{x})}{1-p(C_1 | \mathbf{x})}}$



z is called the logit and is given by the log odds



LUND  
UNIVERSITY

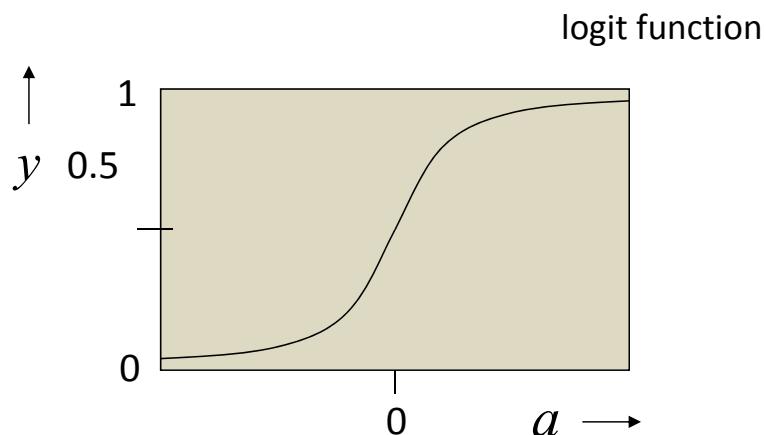
# Logistic regression

## The Logistic Sigmoid (due to S-shape)

- This is also called a squashing function because it maps the entire real axis into a finite interval
- For classification, the output  $a$  is a smooth function of the inputs and the weights  $\mathbf{w}$   $\longrightarrow a = \mathbf{w}^T \mathbf{x} + w_0$
- Properties

$$\sigma(-a) = 1 - \sigma(a), \quad a = \ln\left(\frac{\sigma}{1-\sigma}\right)$$

$$y = \sigma(a) = \frac{1}{1+e^{-a}}$$



$$\frac{\partial a}{\partial w_i} = x_i \quad \frac{\partial a}{\partial x_i} = w_i$$

$$\frac{dy}{da} = y(1-y)$$



LUND  
UNIVERSITY

# Logistic regression

## Multiclass Model (Softmax)

$$\begin{aligned} p(C_k | \mathbf{x}) &= \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|C_k)p(C_k)}{\sum_j p(\mathbf{x}|C_j)p(C_j)} \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \end{aligned}$$

where  $a_k = \ln p(\mathbf{x}|C_k)p(C_k)$

- This is known as the *normalized exponential*
- Can be viewed as a multiclass generalization of the logistic sigmoid
- It is also called a *softmax function* (it is a smoothed version of 'max')

if  $a_k \gg a_j \forall j \neq k$ , then  $p(C_k|x) \simeq 1$  and  $p(C_j|x) \simeq 0$



LUND  
UNIVERSITY

# Logistic regression

## Probabilistic Discriminative Models *Logistic Regression*

- In our discussion of generative approaches, we saw that under general assumptions, the class posterior for  $C_1$  can be written as a logistic sigmoid acting on a linear function of the feature vector  $\phi$
- In logistic regression, we use the functional form of the generalized linear model explicitly

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi), \text{ where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$



LUND  
UNIVERSITY

# Logistic regression

## Maximum Likelihood for *Logistic Regression*

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi); \sigma(a) = \frac{1}{1 + \exp(-a)}, \frac{d\sigma}{da} = \sigma(1 - \sigma)$$

For dataset  $\{\phi_n \equiv \phi(\mathbf{x}_n), t_n\}$ , with  $t_n \in \{0,1\}, n = 1, \dots, N$ ,

$$\mathbf{t} = (t_1, t_2, \dots, t_N)^T, \quad y_n = p(C_1|\phi_n) = \sigma(a_n), a_n = \mathbf{w}^T \phi_n$$

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

Cross-entropy error

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

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

Similar form as the gradient  
of the sum-of-squares  
regression model



LUND  
UNIVERSITY

# Logistic regression

[Matlab example]



LUND  
UNIVERSITY

## **Self study:**

- **Goodfellow et al 2016 (Deep Learning):**
  - -
- **Bishop 2006 (Pattern Recognition and Machine Learning)**
  - Chapter 4: 4.1 (except 4.1.7), 4.3 up to 4.3.2, and 4.3.4.
- **Hastie et al. 2019 (Elements of Statistical Learning)**
  - Chapter 4: 4.1-4.4



**LUND**  
UNIVERSITY

# Small recap from L-3

[Blackboard examples]



LUND  
UNIVERSITY



LUND  
UNIVERSITY

## L-4: Linear classification methods

Ted Kronvall, Cristian Sminchisescu

