
Klasyfikatory liniowe

Linear classifiers



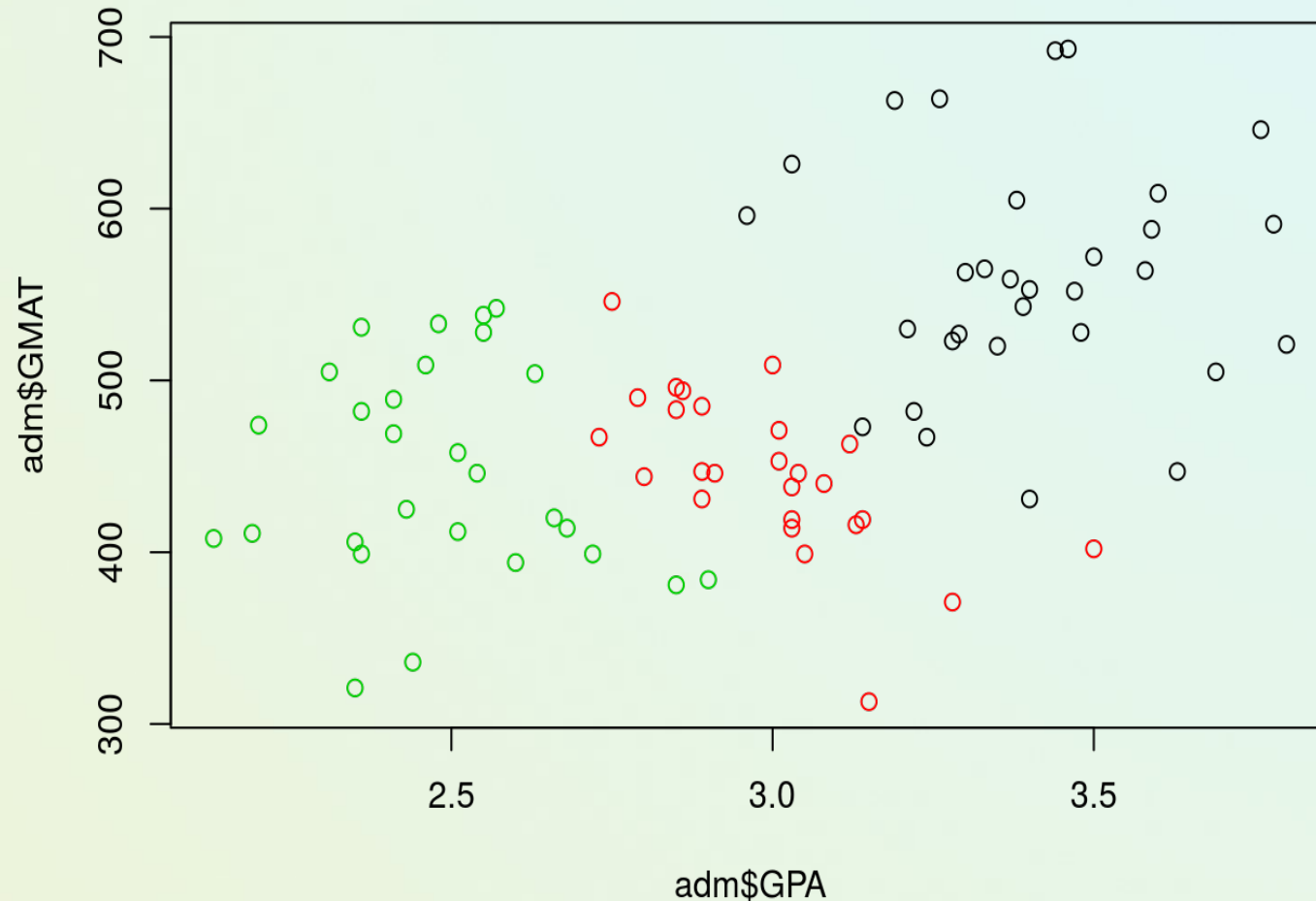
JERZY STEFANOWSKI
Institute of Computing Sciences,
Poznań University of Technology

UM – slajdy wykładu
Wersja specjalna dla 2019

Plan

1. Liniowe klasyfikatory
2. Klasyczne liniowa analiza dyskryminacyjna
3. Sformułowanie probabilistyczne
4. Inne zagadnienia – uczenie perceptronu
5. W stronę SVM

Przykład – dokumentacja R nt. discriminant analysis [za G.Martos]



- Admission data for applicants to graduate schools in business. The objective is to use the GPA and GMAT scores to predict the likelihood of admission (admit, notadmit, and borderline).
- '<http://www.biz.uiowa.edu/faculty/jledolter/DataMining/admission.csv>'

Formalizacja problemu klasyfikacji

- W przestrzeni danych (ang. measurement space) Ω znajdują się wektory danych \mathbf{x} stanowiące próbkę uczącą D , należące do dwóch lub więcej K klas

$$D = \left\{ (\mathbf{x}_i, c_i) \mid \mathbf{x}_i \in R^p, c_i \in \{C_1, \dots, C_k\} \right\}_{i=1}^N$$

- Klasyfikacja jest dokonywana na podstawie funkcji będącej liniową kombinacją p cech i parametrów

$$y = f(\mathbf{x}, \mathbf{w})$$

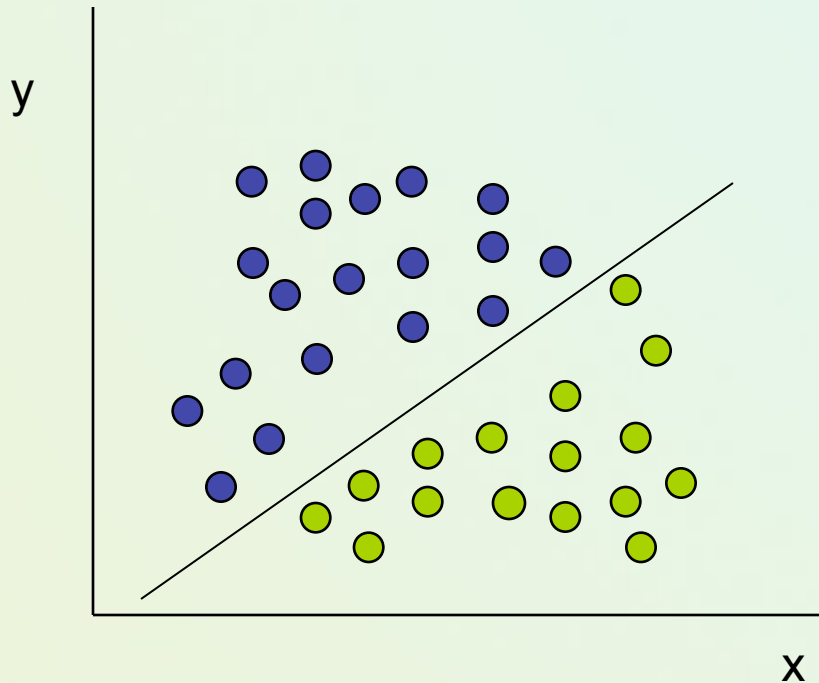
- Dążymy do sytuacji

$$y_i = f(\mathbf{x}_i, \mathbf{w}) = c_i$$

- i/lub minimalizacji błędów klasyfikacji

$$y_i \neq c_i$$

Liniowa funkcja separująca (graniczna)



- Szukamy klasyfikatora pozwalającego na podział całej przestrzeni na obszary odpowiadające klasom (dwie lub więcej) oraz pozwalającego jak najlepiej klasyfikować nowe obiekty x do klas
- Podejście opiera się na znalezieniu tzw. granicy decyzyjnej między klasami
→ $f(\mathbf{x}) = \mathbf{w}^T \cdot \mathbf{x}$

$$y = \begin{cases} f(\mathbf{x}_i) > T & \mathbf{x}_i \in C_1 \\ f(\mathbf{x}_i) < T & \mathbf{x}_i \in C_2 \end{cases}$$

$$w_0 + w^T x = 0$$

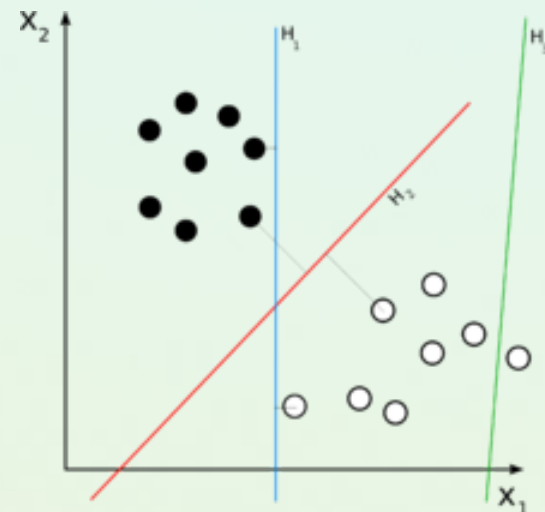
$$w_0 + w^T x > 0$$



$$w_0 + w^T x < 0$$

Różne podejścia do budowy klasyfikatorów liniowych

- Podejścia generatywne (probabilistyczne)
 - Analiza dyskryminacyjna
 - Fisher-owska analiza i tzw. LDA
 - Wersja klasyfikacji Bayesowskiej (dwumianowy rozkład)
- Podejścia wykorzystujące własności zbioru uczącego
 - Perceptron liniowy Rosenblata (iteracyjne poprawki wag)
 - Metoda wektorów nośnych (max. marginesu klasyfikatora)
 - Regresja logistyczna



Co jest celem analizy dyskryminacyjnej

- ▶ • Podejście statystyczne do problemów klasyfikowania obiektów (term. ang. *Discriminant Analysis*)
 - Oryginalnie wprowadzona przez R.A.Fishera (1936) dla funkcji liniowych (2 klasy),
 - Metody probabilistyczne – B.Welch .
- Dostępna w wielu programach, np. SAS, SPSS, R lub Statistica,...
- Liczne zastosowania
- ...

Liniowa analiza dyskryminacyjna

- Problem wprowadzony przez R.A.Fishera w 1936 dla wielowymiarowej przestrzeni atrybutów (zmiennych liczbowych) – dyskryminacja 2 klas
- Fisher oryginalnie zaproponował poszukiwanie kierunku projekcji, na którym można dobrze rozdzielić rzutowane obie klasy
 - Średnie w klasach są dostatecznie oddalone od siebie
 - Obszary rozrzutu (rozproszenia, zmienności) obu klas nie nakładają się zbyt mocno.

Intuicja projekcji w Fisher's Linear Discriminant [EST]

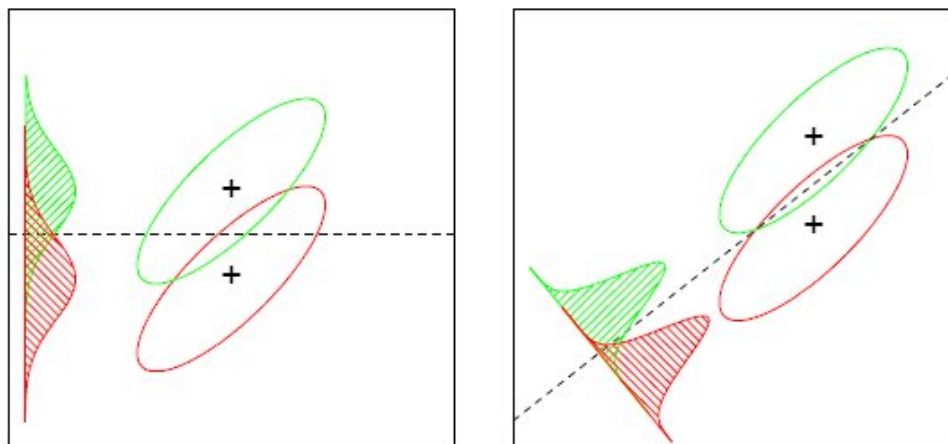


Figure 4.9: *Although the line joining the centroids defines the direction of greatest centroid spread, the projected data overlap because of the covariance (left panel). The discriminant direction minimizes this overlap for Gaussian data (right panel).*

„From training set we want to find out a direction where the separation between the class means is **high** and overlap between the classes is **small**”

Trochę uwag matem. o projekcji

- Dysponujemy przykładami uczącymi opisanymi p -cechami $\mathbf{x}=[x_1, x_2, \dots, x_p]^T$ należącymi do dwóch klas C_1 i C_2 (odpowiednio n_1 i n_2)
- Wektory p -wymiarowe \mathbf{x} są rzutowane na prostą (kierunek związany z parametrami \mathbf{w}). Algebraicznie odpowiada to zastąpieniu ich skalarom $z = \mathbf{w}^T \cdot \mathbf{x}$. Celem jest taki dobór \mathbf{w} aby na podstawie nowej zmiennej z przykłady z obu klas były jak najlepiej rozdzielone.

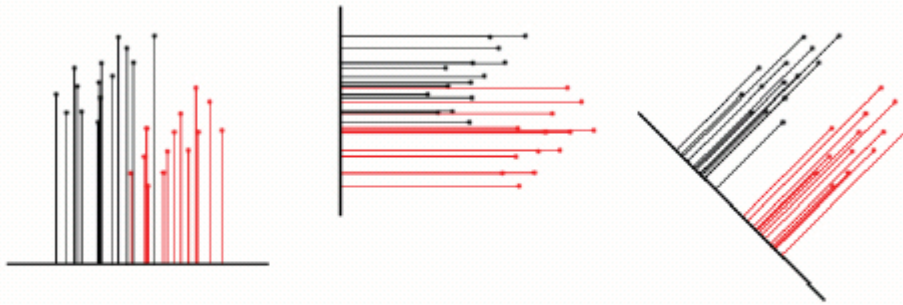
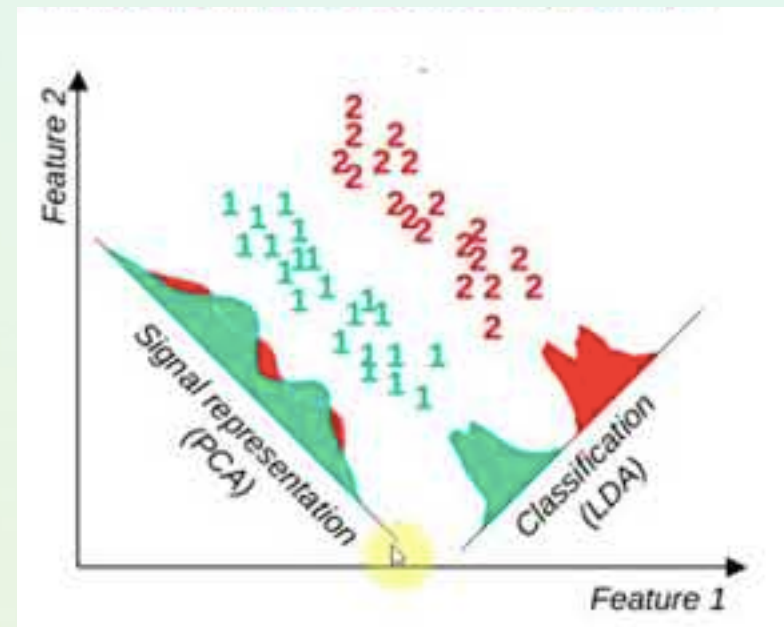


Figure 1: Projection of data from two classes onto various lines.



Założenia co do danych

- Fisher – dość ograniczone założenia: wektor p wartości oczekiwanych $E(\mathbf{x})$ oraz rozproszenie charakteryzowane przez macierz kowariancji $\Sigma = \text{Cov}(\mathbf{x}) = E[(\mathbf{x} - E(\mathbf{x})) \cdot (\mathbf{x} - E(\mathbf{x}))^T]$
- Estymatory

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}) \cdot (x_i - \bar{x})^T$$

- Wariancja po rzutowaniu \mathbf{x} na prosta o wektorze kierunkowym \mathbf{w}

$$\text{Var}(\mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \Sigma \mathbf{w}$$

Sformułowanie problemu Fisher LDA

Cel

- Maksymalizuj odległość rzutowanych średnich klas
- Minimalizuj wariancje wewnątrz klasową
- Odległość między rzutami średnich

$$(\mathbf{w}^T \bar{\mathbf{x}}_1 - \mathbf{w}^T \bar{\mathbf{x}}_2)^2$$

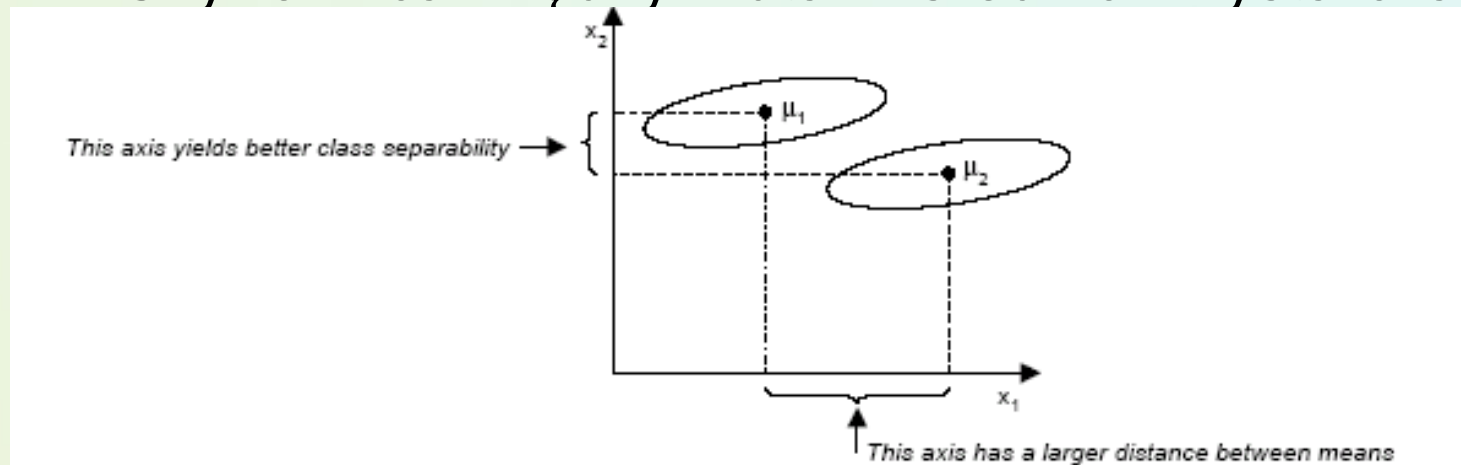
- Fisher założył, że obie klasy mają taką samą macierz kowariancji $S=S_1+S_2$. Dlatego wskaźnik zmienności wewnątrzgrupowej (wspólnej dla obu klas) zdefiniowany jest jako:

$$S_W = \frac{1}{n-2} \sum_{k=1}^2 (n_k - 1) S_k$$

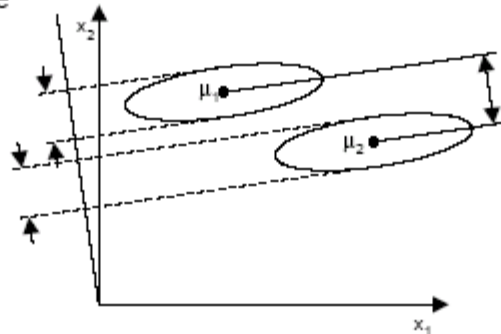
- Pamiętaj, że po rzutowaniu mamy $\mathbf{w}^T S_W \mathbf{w}$

Co optymalizować?

- Czy różnica między rzutami średnich wystarczy?



Therefore, we will be looking for a projection where examples from the same class are projected very close to each other and, at the same time, the projected means are as farther apart as possible



Sformułowanie problemu Fisherian LDA

- W celu maksymalizacji odległości rzutów średnich klas i minimalizacji wariancji wewnątrzklasowej należy poszukiwać wektora \mathbf{w} który maksymalizuje następujące wyrażenie:

$$J(\mathbf{w}) = \frac{(\mathbf{w}^T \bar{\mathbf{x}}_1 - \mathbf{w}^T \bar{\mathbf{x}}_2)^2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

- Po znalezieniu kierunku maksymalizującego $J(\mathbf{w})$ można stosować zasadę klasyfikacji na rzutowanej prostej. Przypisz \mathbf{x} do klasy j dla której

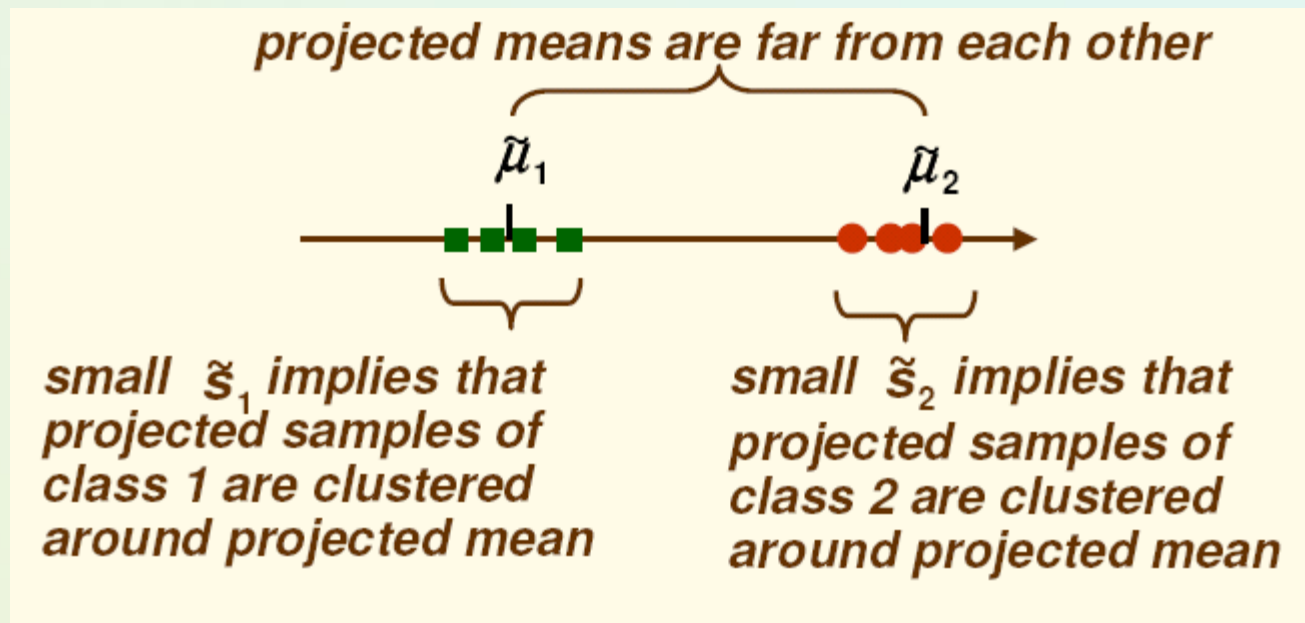
$$\left| \tilde{\mathbf{w}}^T \mathbf{x} - \tilde{\mathbf{w}}^T \bar{\mathbf{x}}_j \right| < \left| \tilde{\mathbf{w}}^T \mathbf{x} - \tilde{\mathbf{w}}^T \bar{\mathbf{x}}_k \right|$$

- Można wykazać, że ten wektor jest proporcjonalny

$$\tilde{\mathbf{w}} \propto S_W^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$$

Uwagi o konstrukcji wskaźnika

- Zwiększanie $J(w)$ ma gwarantować dobrą separację klas i ich rzutów



Hiperpłaszczyzna separująca

- Wyraz wolny to środek odcinka między rzutami średnich

$$m = \frac{1}{2}(\mathbf{w}^T \bar{\mathbf{x}}_1 - \mathbf{w}^T \bar{\mathbf{x}}_2) = \frac{1}{2}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T S_W^{-1}(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2)$$

- Stąd liniowa funkcja dyskryminacyjna Fishera

$$y = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T S_W^{-1}[\mathbf{x} - \frac{1}{2}(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2)]$$

- Więcej informacji, np.
 - J.Koronacki, J.Cwik: Statystyczne systemy uczące się
 - M.Krzyśko et al.: Systemy uczące się

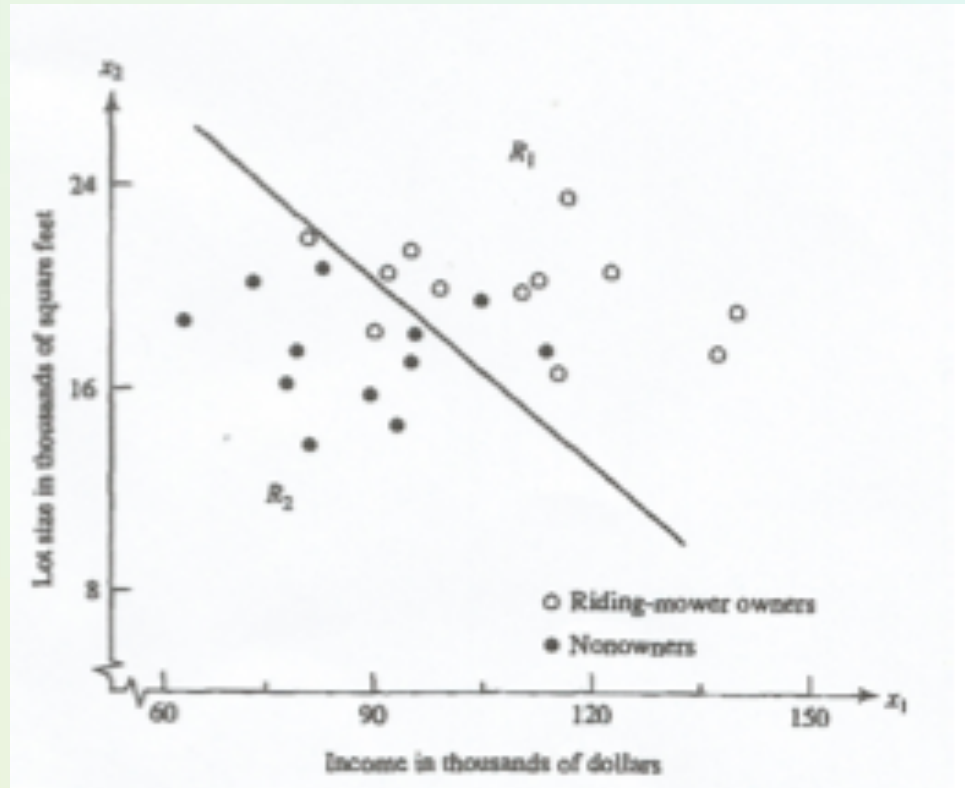
Przykład ilustracyjny [za A.Gołota i in.]

Rozróżnienie rodzin, jako: potencjalnych posiadaczy i nie posiadaczy na podstawie x_1 =przychody i x_2 =powierzchnia działki. [analiza marketingowa / SAS]

Uzyskane wartości przedstawia tabela:

π_1 : posiadacze		π_2 : nie posiadacze	
x_1 (przychody w \$1000s)	x_2 (powierzchnia działki w 1000ft ²)	x_1 (przychody w \$1000s)	x_2 (powierzchnia działki w 1000ft ²)
90.0	18.4	105.0	19.6
115.5	16.8	82.8	20.8
94.8	21.6	94.8	17.2
91.5	20.8	73.2	20.4
117.0	23.6	114.0	17.6
140.1	19.2	79.2	17.6
138.0	17.6	89.4	16.0
112.8	22.4	96.0	18.4
99.0	20.0	77.4	16.4
123.0	20.8	63.0	18.8
81.0	22.0	81.0	14.0
111.0	20.0	93.0	14.8

Przykład posiadaczy kosiarek



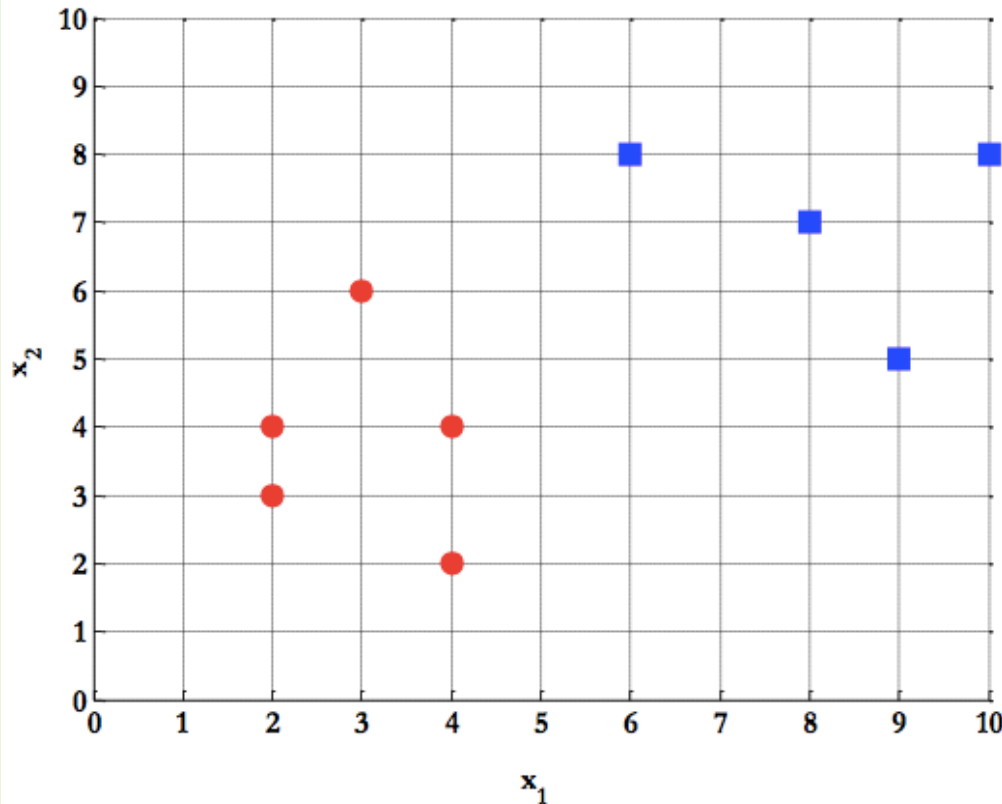
Wiecej w A.Gołota i in. wykład Pol. Gdan. pt. Klasyfikacja i dyskryminacja. Statystyka w SAS

Dla bardziej zainteresowanych detalami

Więcej informacji o rozwiązaniu analitycznym oraz przykład obliczeniowy (krok po kroku) - Sh. Elhabian, A.Farag Tutorial Linear Discriminant Analysis (WWW, patrz załącznik)

– Samples for class ω_1 : $\mathbf{X}_1=(x_1,x_2)=\{(4,2),(2,4),(2,3),(3,6),(4,4)\}$

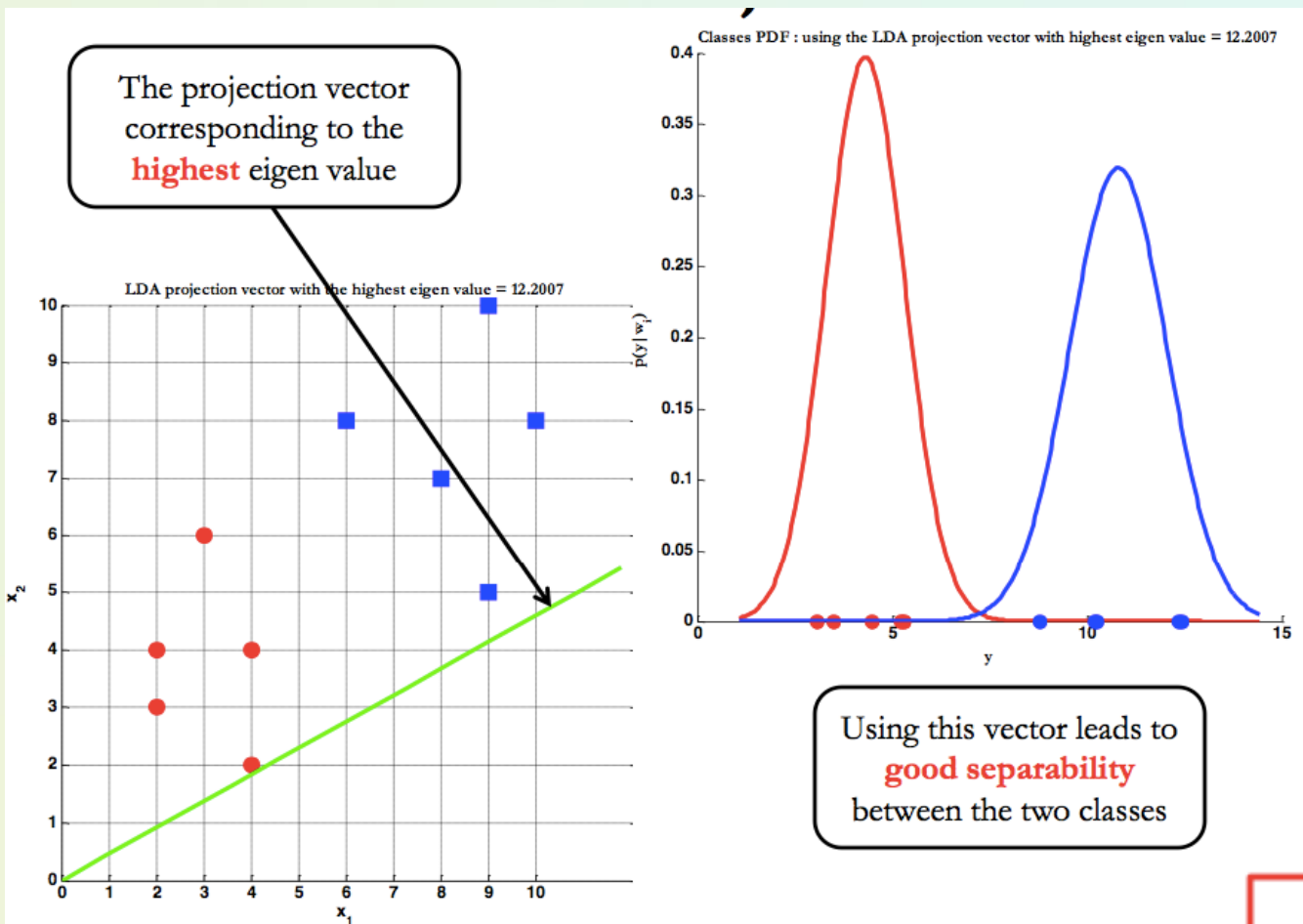
– Sample for class ω_2 : $\mathbf{X}_2=(x_1,x_2)=\{(9,10),(6,8),(9,5),(8,7),(10,8)\}$



```
% samples for class 1
X1 = [4,2;
      2,4;
      2,3;
      3,6;
      4,4];

% samples for class 2
X2 = [9,10;
      6,8;
      9,5;
      8,7;
      10,8];
```

Odnaleziony kierunek projekcji w^*



Więcej w załączonych materiałach

$$w_2 = \begin{pmatrix} 0.9088 \\ 0.4173 \end{pmatrix} = w^*$$

Przypadek wielu klas ($K > 2$)

- Rozwiązanie Fishera uogólniono dla większej liczby K klas (C.Rao 1948)

- Średnia w próbie uczącej $\bar{\mathbf{x}} = \frac{1}{n} \sum_{j=1}^K \bar{\mathbf{x}}_j$

- Macierz zmienności wewnątrzklasowej

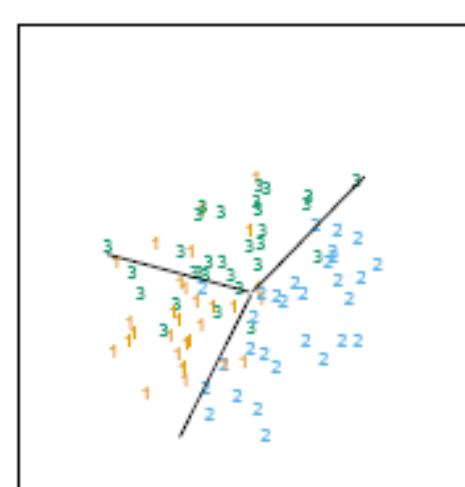
$$S_W = \frac{1}{n - K} \sum_{j=1}^K \sum_{i \in C_j} (x_i - \bar{x}_j)(x_i - \bar{x}_j)^T$$

- Macierz zmienności międzyklasowej

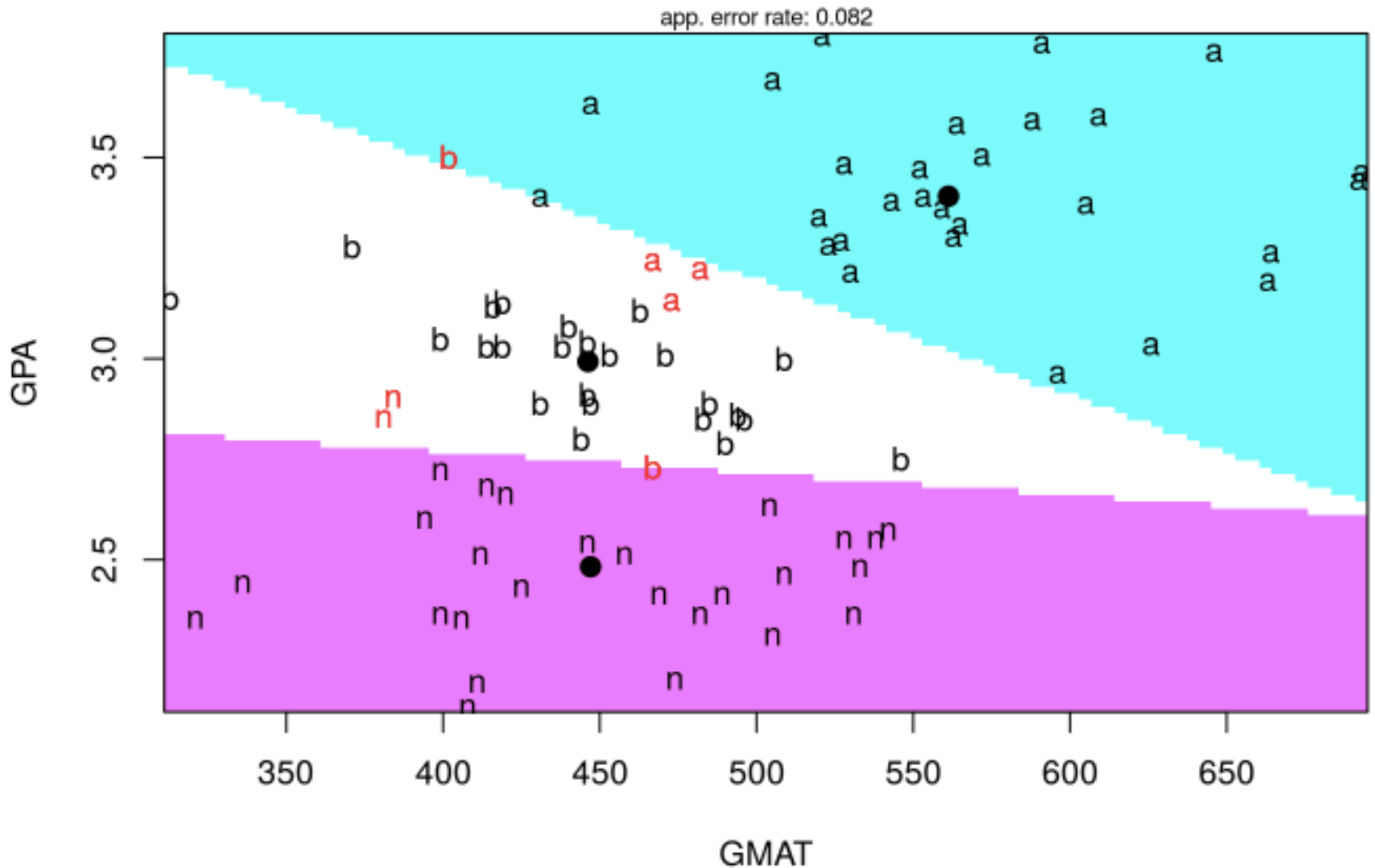
$$S_B = \frac{1}{K - 1} \sum_{j=1}^K n_j (\bar{x}_j - \bar{x})(\bar{x}_j - \bar{x})^T$$

- Poszukuj wektora \mathbf{w} maksymalizującego

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

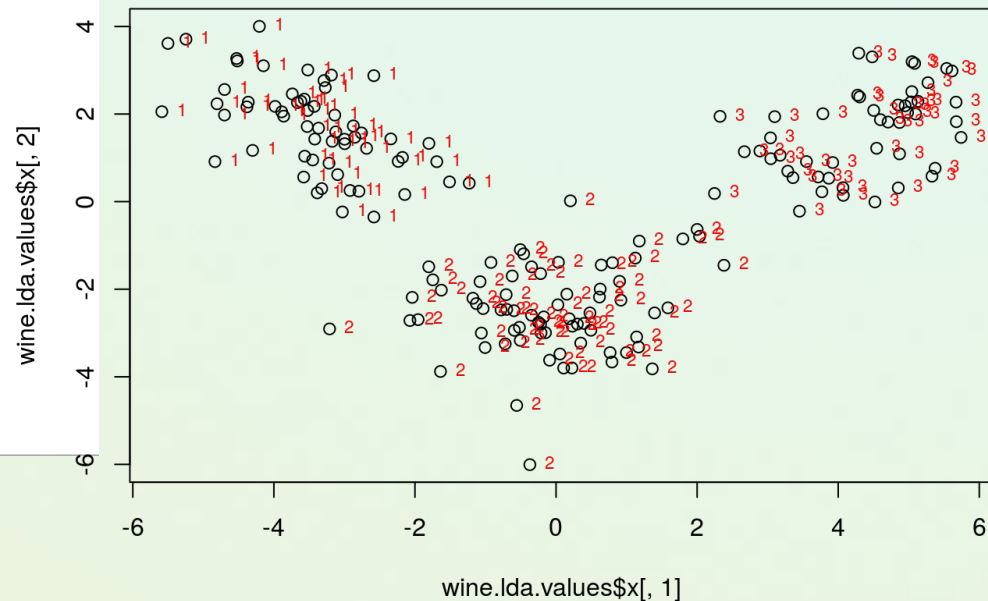
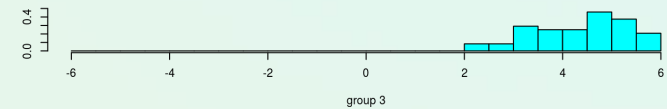
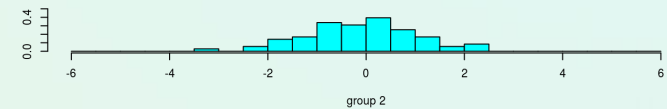
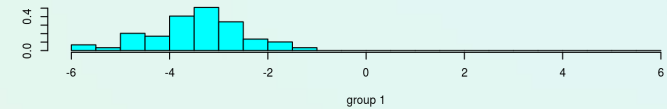


Dane o przyjęciach do uczelni biznesowych

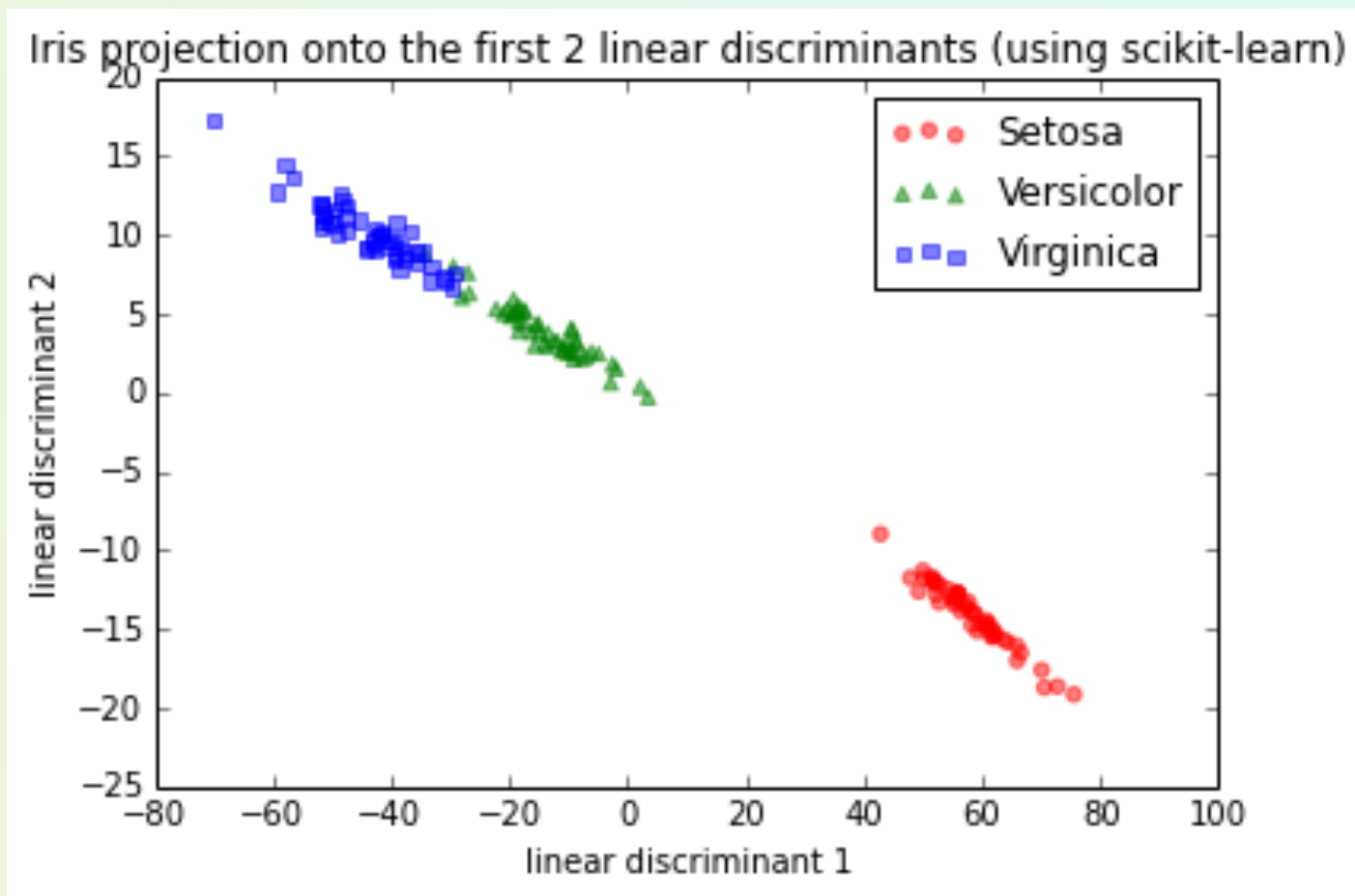


LDA w R (wine data – 3 klasy)

```
## Call:
## lda(Type ~ ., data = wine)
##
## Prior probabilities of groups:
##      1      2      3
## 0.3314607 0.3988764 0.2696629
##
## Group means:
##      Alcohol      Malic      Ash Alcalinity Magnesium Phenols Flavanoids
## 1 13.74475 2.010678 2.455593 17.03729 106.3390 2.840169 2.9823729
## 2 12.27873 1.932676 2.244789 20.23803 94.5493 2.258873 2.0808451
## 3 13.15375 3.333750 2.437083 21.41667 99.3125 1.678750 0.7814583
##      Nonflavanoids Proanthocyanins      Color      Hue Dilution      Proline
## 1 0.290000 1.899322 5.528305 1.0620339 3.157797 1115.7119
## 2 0.363662 1.630282 3.086620 1.0562817 2.785352 519.5070
## 3 0.447500 1.153542 7.396250 0.6827083 1.683542 629.8958
##
## Coefficients of linear discriminants:
##
##              LD1              LD2
## Alcohol      -0.403399781  0.8717930699
## Malic         0.165254596  0.3053797325
## Ash          -0.369075256  2.3458497486
## Alcalinity    0.154797889 -0.1463807654
## Magnesium    -0.002163496 -0.0004627565
## Phenols      0.618052068 -0.0322128171
## Flavanoids   -1.661191235 -0.4919980543
## Nonflavanoids -1.495818440 -1.6309537953
## Proanthocyanins 0.134092628 -0.3070875776
## Color        0.355055710  0.2532306865
## Hue         -0.818036073 -1.5156344987
## Dilution    -1.157559376  0.0511839665
## Proline     -0.002691206  0.0028529846
```



Iris – 3 klasy (rzutowanie projekcji LDA)



Przykład z scikit learn Python

Podjęcia opisowe i probabilistyczne

- ▶ • Stochastyczne / probabilistyczne
 - Zbiór obserwacji jest próbą losową pobraną z k podpopulacji $\pi_1, \pi_2, \dots, \pi_k$; celem jest taki podział aby podpopulacje odpowiadały właściwym k klasom C_1, C_2, \dots, C_k
- Opisowe
 - Nie rozważa się losowości próby, zakłada się że posiadany zbiór zawiera przykłady z k klas C_1, C_2, \dots, C_k ; zadanie polega na poprawnym podziale zbioru na klasy

Sformułowanie probabilistyczne z Tw. Bayesa

- Obiekty $\mathbf{x} \in \mathbb{R}^p$ i wielowymiarowy rozkład prawdopodobieństwa – funkcja gęstości $f(\mathbf{x}|C_i)$
- Każda klasa C_i opisana prawdopodobieństwa apriori p_i
- Bayesowska reguła klasyfikowania

$$P(C | x) = P(x | C)P(C)/P(x)$$

- • Przydziel nowy obiekt \mathbf{x} do tej klasy C_i dla której prawdopodobieństwo a posteriori jest największe:

$$P(C_j | \underline{x}) = p_j \cdot f(\underline{x} | C_j) / \sum_{i=1}^K p_i \cdot f(\underline{x} | C_i)$$

Rozwiązanie probabilistycznej reguły klasyfikacji

- Załóżmy, że rozkłady wektora x w poszczególnych klasach są p -wymiarowymi rozkładami normalnymi:

$$f(\underline{x} | C_i) = (2\pi)^{-0,5p} |\Sigma_i|^{-0,5} \exp\left[-0,5(\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i)\right]$$

- Wykorzystując logarytmiczne przekształcenie twierdzenia Bayesa, obiekt x jest przydzielany do tej klasy C_j dla której funkcja dyskryminująca osiąga maksimum:

$$\delta_j(x) = -0,5(\underline{x} - \underline{\mu}_j)^T \Sigma_j^{-1} (\underline{x} - \underline{\mu}_j) - 0,5 \log |\Sigma_j| + \log p_j$$

- Jest to kwadratowa funkcja dyskryminująca (QDA)

Liniowa funkcja

- Założenie równości macierzy kowariancji Σ

$$\delta_j(x) = \underline{x}^T \Sigma^{-1} \underline{\mu}_j - 0,5 \underline{\mu}_j^T \Sigma^{-1} \underline{\mu}_j + \log p_j$$

- Dla dwóch klas – przekształcenie log-ratio

$$\begin{aligned} \log \frac{\Pr(y = k|\mathbf{x})}{\Pr(y = l|\mathbf{x})} &= \log \frac{f_k(\mathbf{x})}{f_l(\mathbf{x})} + \log \frac{p_k}{p_l} \\ &= \log \frac{p_k}{p_l} - \frac{1}{2} (\underline{\mu}_k + \underline{\mu}_l)^T \Sigma^{-1} (\underline{\mu}_k - \underline{\mu}_l) \\ &\quad + \underline{x}^T \Sigma^{-1} (\underline{\mu}_k - \underline{\mu}_l) \end{aligned}$$

- Więcej w Krzyśko ... lub Hastie et al. Elements of Statistical Learning

Implementacje funkcji w Statistica

The screenshot displays the Statistica software interface with a data table and a menu of statistical functions. The data table has columns for torque, summer_cons, winter_cons, oil_cons, horsepower, D1, and D2. The menu lists various statistical methods such as Basic Statistics/Tables, Multiple Regression, ANOVA, Nonparametrics, Distribution Fitting, Advanced Linear/Nonlinear Models, Multivariate Exploratory Techniques, Industrial Statistics & Six Sigma, Power Analysis, Neural Networks, Data-Mining, QC Data Mining & Root Cause Analysis, Text & Document Mining, Web Crawling, Statistics of Block Data, STATISTICA Visual Basic, and Probability Calculator.

	1	2	3	4	5	6	7	8	9	10	11
	id	Ms									
1	1										
2	2										
3	3										
4	4										
5	5										
6	6										
7	7										
8	8										
9	9										
10	10										
11	11										
12	12										
13	13										
14	14										
15	15										
16	16	65	2,22	67	402	22	23,9	2,3	103	2	3
17	17	90	2,48	51	468	22	26,5	1,2	138	1	1
18	18	90	2,6	15	488	20	23,2	0,1	150	1	1
19	19	76	2,39	65	428	27	33,4	2	116	2	3
20	20	85	2,42	50	454	21,5	26,3	1,3	129	1	2
21	21	85	2,41	58	450	22	25,5	1,5	126	1	2
22	22	88	2,47	48	458	22,4	25,1	1,1	130	1	1
23	23	60	1,93	90	400	24	28,7	4	95	2	3
24	24	64	2,2	71	420	23,1	25,2	2,6	105	2	3
25	25	75	2,39	64	432	22,2	25,1	1,7	114	2	2
26	26	74	2,36	64	420	21,9	25,4	1,9	110	2	2
27	27	68	2,15	70	400	22	26	2,6	100	2	3
28	28	70	2,2	65	412	22,8	25,3	2,1	102	2	3

Workbook2* - Classification Functions; grouping: D1 (autobusyplainpopraw.sta)

Variable	Classification Functions	
	G_1:1 p=,60000	G_2:2 p=,40000
MaxSpeed	-2,49	-2,83
Compr_pressure	823,79	833,98
blackings	-3,55	-3,71
torque	14,19	14,01
summer_cons	22,14	22,68
winter_cons	7,48	7,51
oil_cons	257,56	259,33
horsepower	-15,74	-15,85
Constant	-3569,67	-3482,11

Discriminant Function Analysis Results: autobusyplainpopraw.sta

Number of variables in the model: 8

Wilks' Lambda: ,2423013 approx. F (8,71) = 27,75296 p < ,0000

Quick | Advanced | Classification

Classification functions: Use selection conditions to classify selected cases only

Classification matrix:

Classification of cases:

Squared Mahalanobis distances:

Posterior probabilities:

Save scores:

A priori classification probabilities:

- Proportional to group sizes
- Same for all groups
- User defined

Score to save for each case:

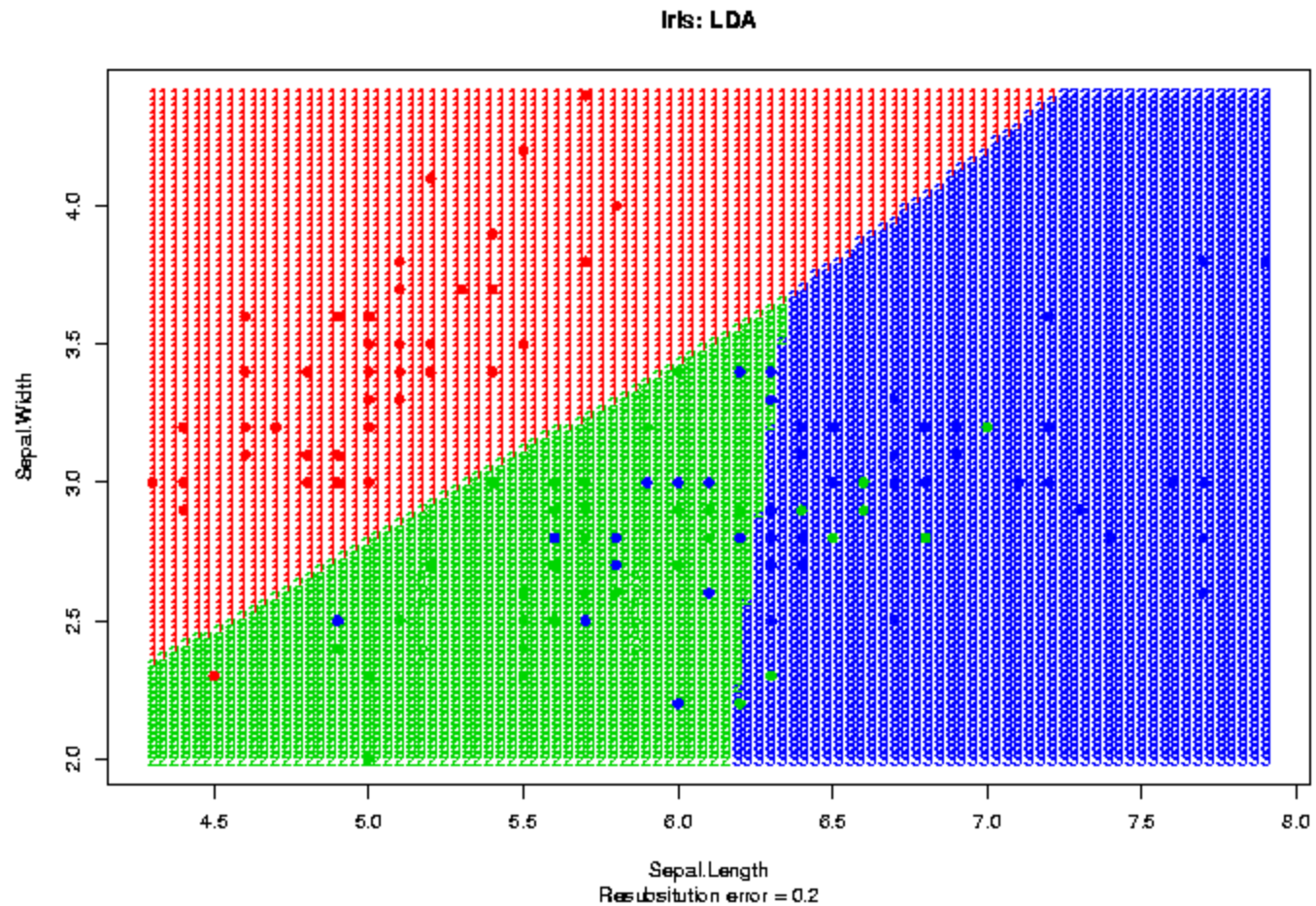
- Save classification for case
- Save distance for case
- Save posterior probability for case

Max. number of cases in a single results spreadsheet: 100000

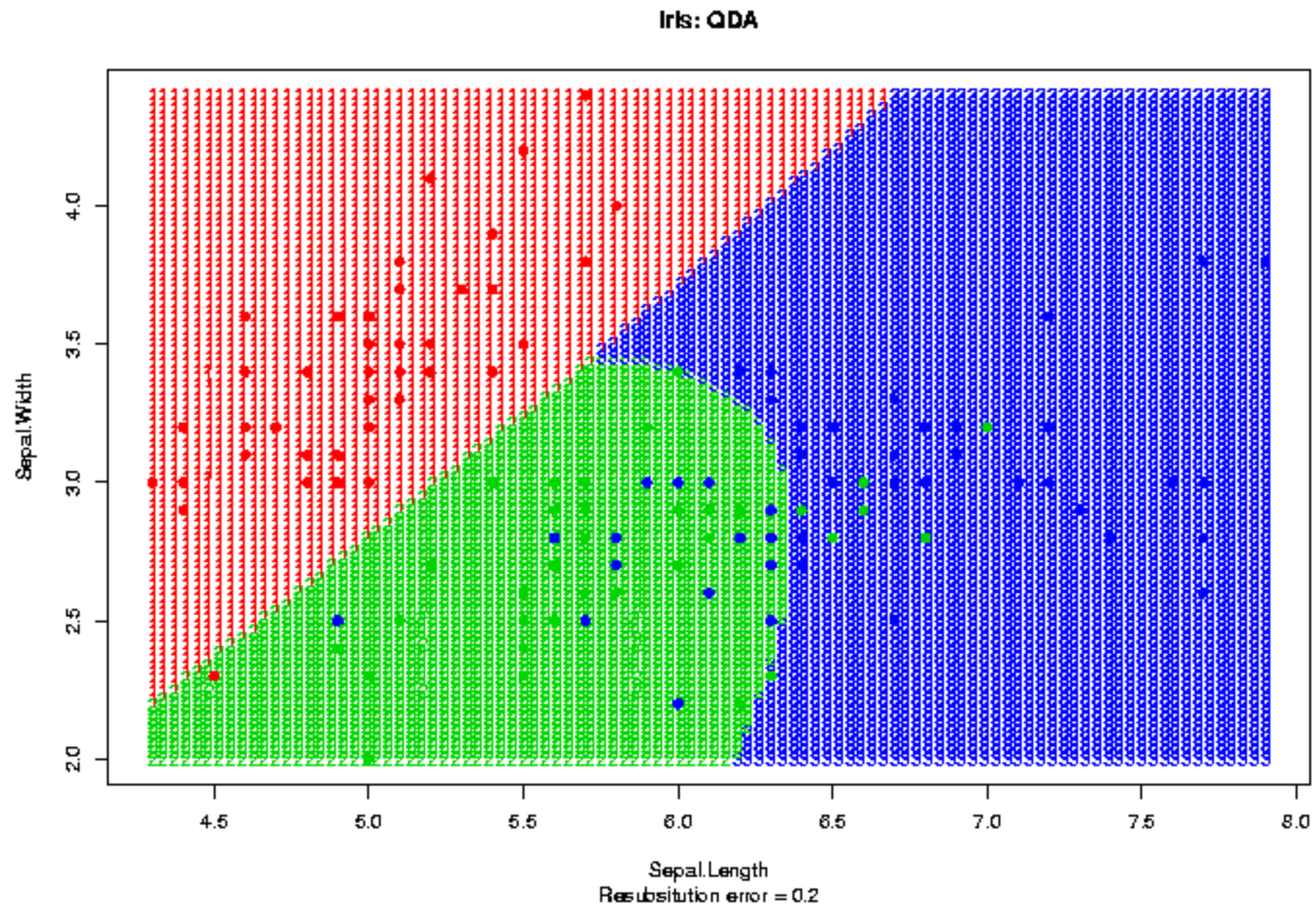
Summary | Cancel | Options

	1	2	3	4	5	6	7	8	9	10	11
	id	MaxSpeed	Compr_pressure	blackings							
25	25	75	2,39	64	432	22,2	25,1	1,1	130	1	1
26	26	74	2,36	64	420	21,9	25,4	4	95	2	3
27	27	68	2,15	70	400	22	26	2,6	105	2	3
28	28	70	2,2	65	412	22,8	25,3	1,7	114	2	2
								1,9	110	2	2
								2,6	100	2	3
								2,1	102	2	3

Example: Linear discriminant analysis



Example: Quadratic discriminant analysis



Porównanie rozwiązań LDA i QDA

- Wybrany zbiór danych (za Hastie et al. Elements of Statistical Learning)

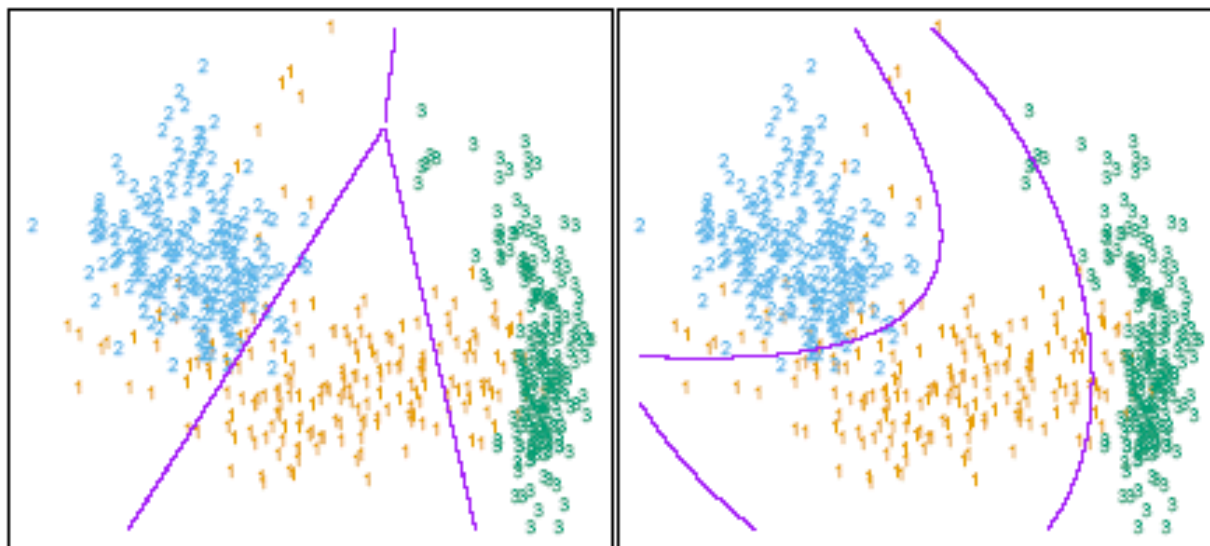


FIGURE 4.1. The left plot shows some data from three classes, with linear decision boundaries found by linear discriminant analysis. The right plot shows quadratic decision boundaries. These were obtained by finding linear boundaries in the five-dimensional space $X_1, X_2, X_1X_2, X_1^2, X_2^2$. Linear inequalities in this space are quadratic inequalities in the original space.

Wymogi stosowania modeli AD

- Zmienne wyrażone na skalach liczbowych
 - Specjalne podejścia dla zmiennych jakościowych (binaryzacja, model lokacyjny,...)
- Zmienne mają wielowymiarowy rozkład normalnych
- Macierze kowariancji dla poszczególnych klas są równe → jeśli nie, to bardziej złożone funkcje kwadratowe dyskryminujące.
- Problem doboru właściwych zmiennych.

Selekcja zmiennych

- W funkcji dyskryminującej uwzględniaj zmienne o dobrych właściwościach dyskryminujących
- Przykład kryterium jakości dyskryminacji:

$$\lambda = \frac{|S_w|}{|S_W + S_B|}$$

gdzie macierz zmienności wewnątrzklasowej

$$S_W = \frac{1}{n - k} \sum_{j=1}^k \sum_{i \in C_j} (x_i - \bar{x}_j)(x_i - \bar{x}_j)^T$$

a macierz zmienności międzyklasowej

$$S_B = \frac{1}{k - 1} \sum_{j=1}^k n_j (\bar{x}_j - \bar{x})(\bar{x}_j - \bar{x})^T$$

Inne zagadnienia

- Pojęcie zmiennych kanonicznych – kierunki które dobrze separują k klasy (także ich wizualizacja)
- Dyskryminacja oparta na regresji liniowej i logistycznej
- Uogólnienie modeli liniowych – elastyczna dyskryminacja (FDA)
- Ad a metoda wektorów nośnych (SVM)
- Powiązanie z metodą PCA
- Odniesienia do Analizy Korespondencji

Typowe obszary zastosowań

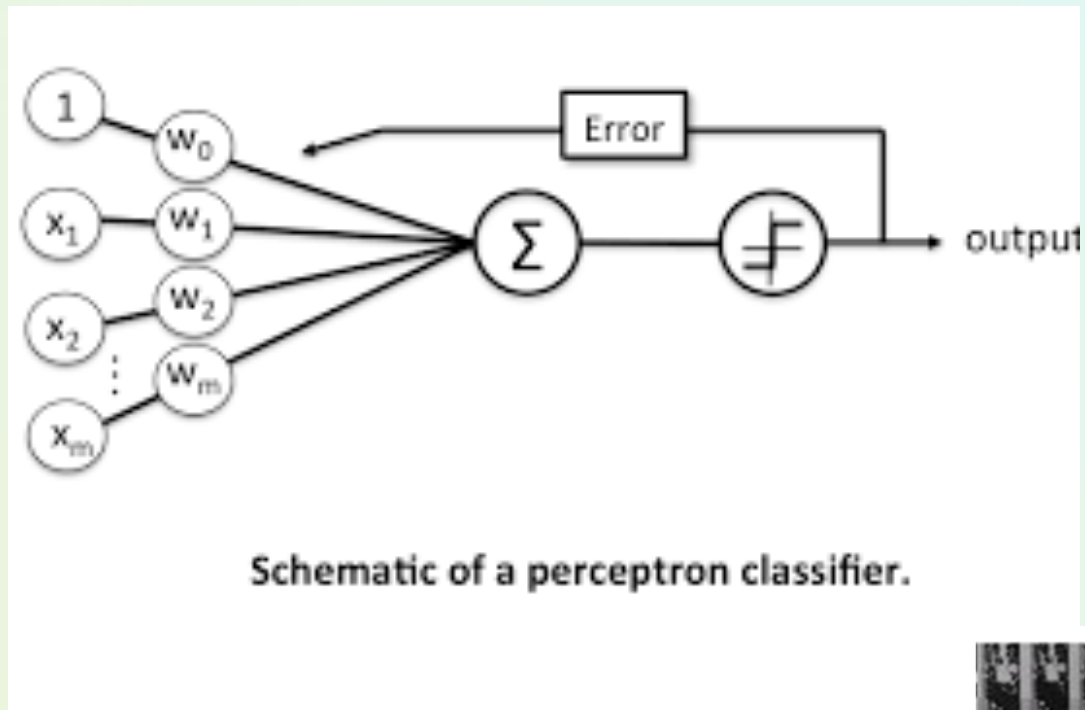
- Analiza danych finansowych (zwłaszcza banki, polityka kredytowa, predykcja bankructw)
- Badania marketingowe
 - Także identyfikacja czynników różnicujących klasy klientów
- Badania danych medycznych, biologicznych lub innych powiązanych nauk
- Rozpoznawania twarzy na obrazach

Więcej

Przeczytaj literaturę

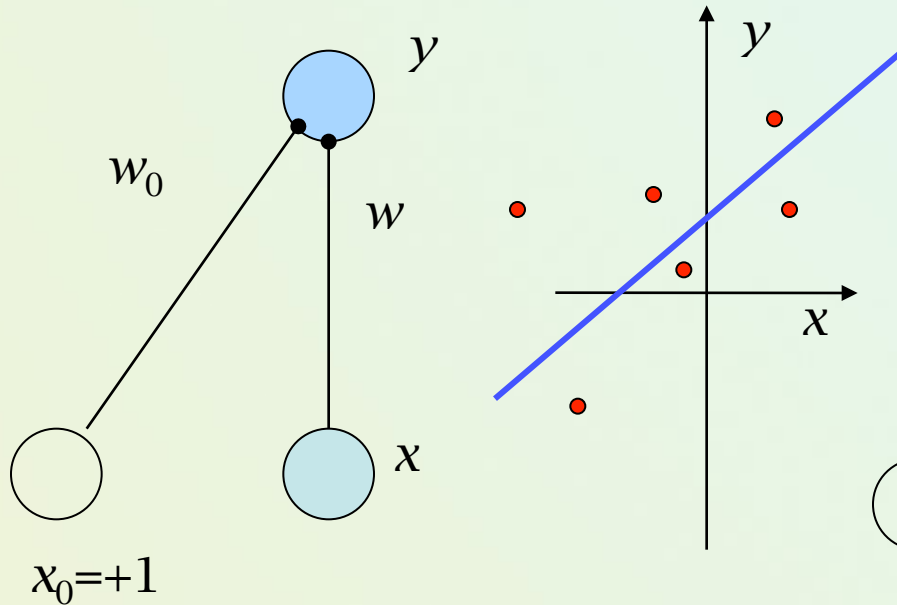
- T.Hastie, R.Tibshirani, J.Friedman: The Elements of Statistical Learning. Springer (zwłaszcza rozdz. 4) → poszukaj wersji elektronicznej pdf
- J.Koronacki, J.Ćwik: Statystyczne systemy uczące się (rozdz. 1 oraz o FDA w rozdz. 6)
- M.Krzyśko, W.Wołyński, T.Górecki, M.Skorzybut: Systemy uczące się. + wcześniejsze prace M.Krzyśko o analizie dyskryminacyjnej
- Angielska Wikipedia „Linear discriminant analysis”
- McLachlan, G. J. (2004). Discriminant Analysis and Statistical Pattern Recognition. Wiley.
- Duda, R. O.; Hart, P. E.; Stork, D. H. (2000). Pattern Classification (2nd ed.). Wiley

Inne podejście do separowalności - ANN

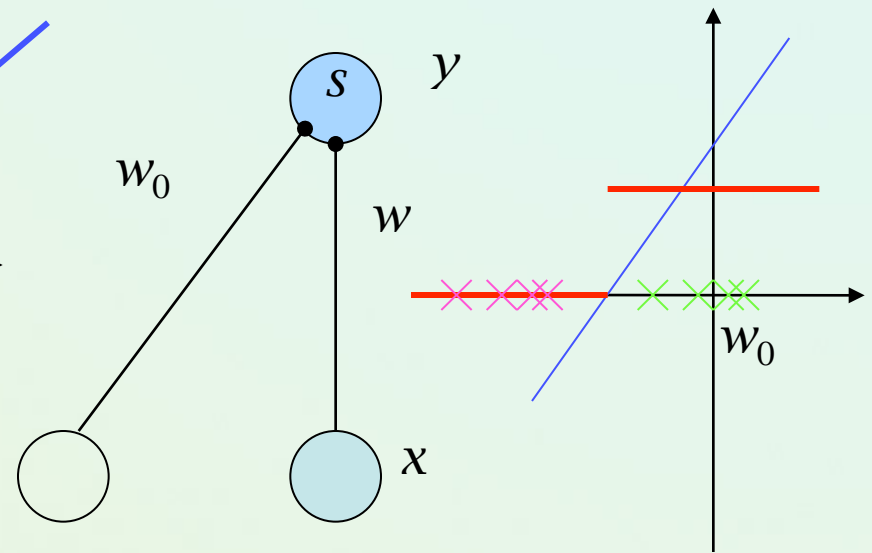


Działanie pojedynczego neuronu (perceptron liniowy)

- Regression: $y = wx + w_0$



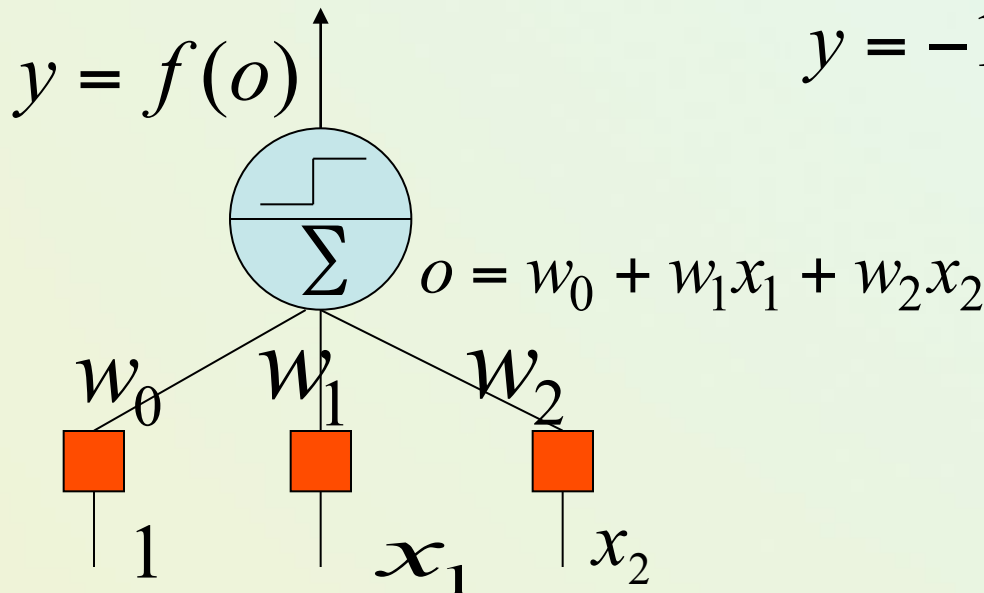
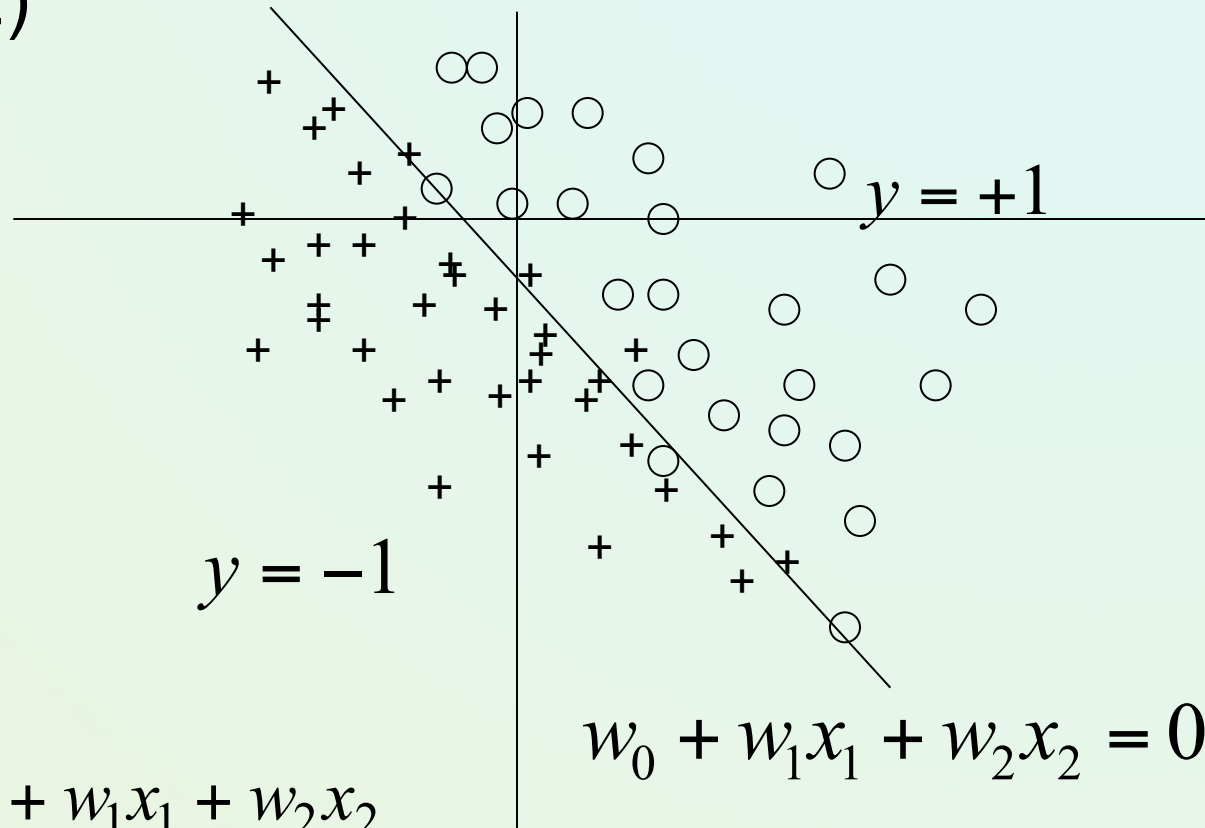
- Classification: $y = 1$ if $(wx + w_0 > 0)$



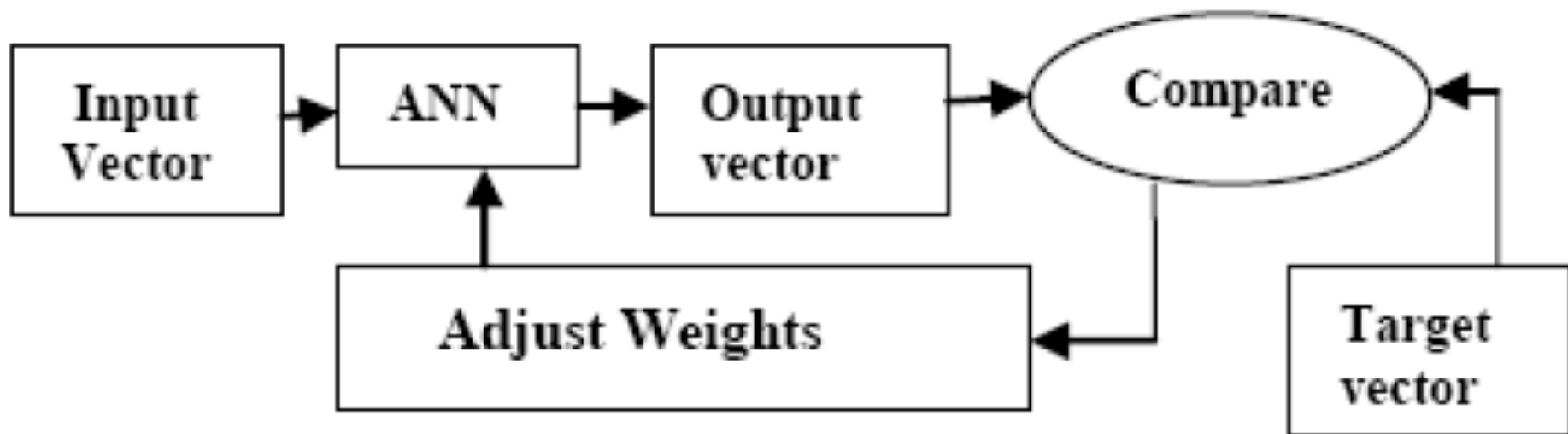
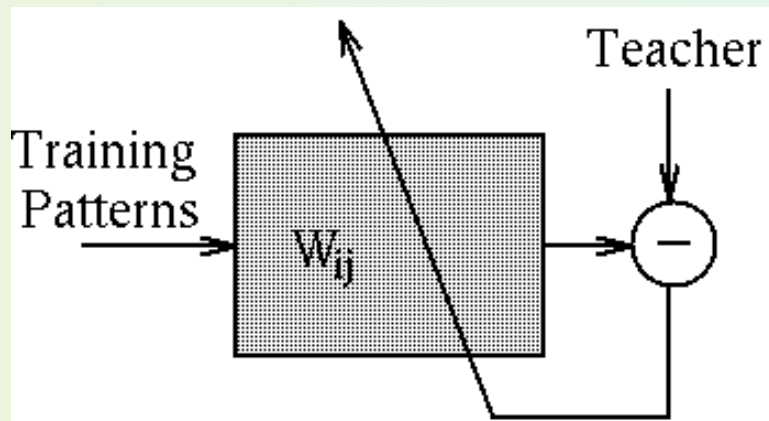
$$y = \text{sigmoid}(o) = \frac{1}{1 + \exp[-\mathbf{w}^T \mathbf{x}]}$$

Perceptron jako liniowy klasyfikator

- Rosenblatt (1962)
- Separowalność liniowa
- Wejścia $x - \mathbb{R}$
- Wyjście :1 or -1

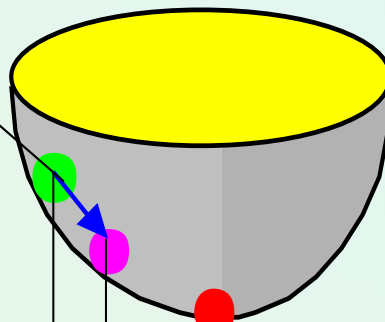


Ogólna idea procesu uczenia neuronów



Błąd

Istota uczenia polega na szukaniu **miejsca**, w którym błąd jest minimalny



minimum funkcji błędów

waga w_2

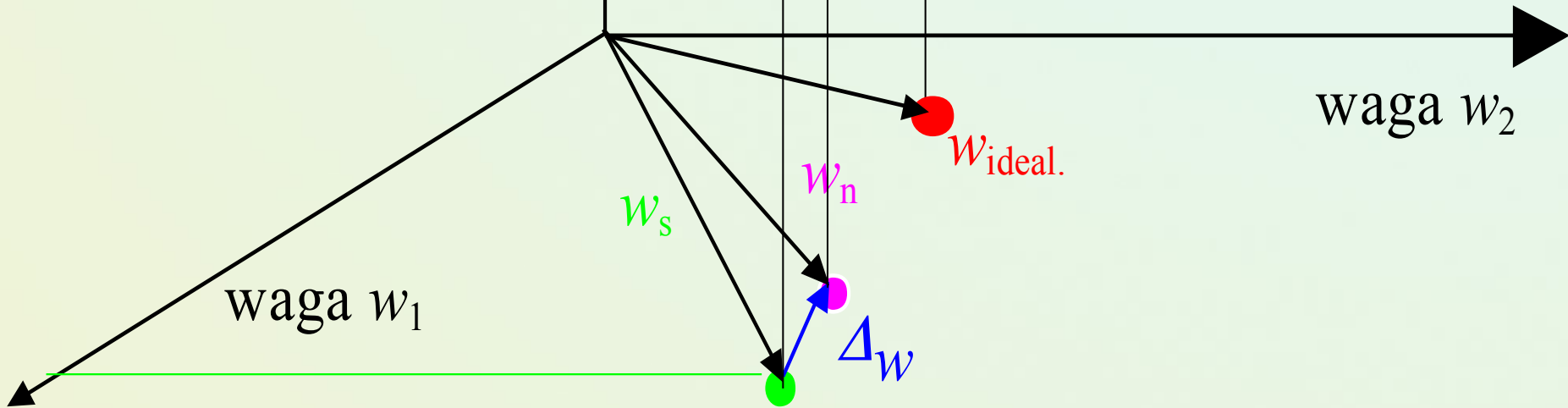
$w_{ideal.}$

w_s

w_n

waga w_1

Δw

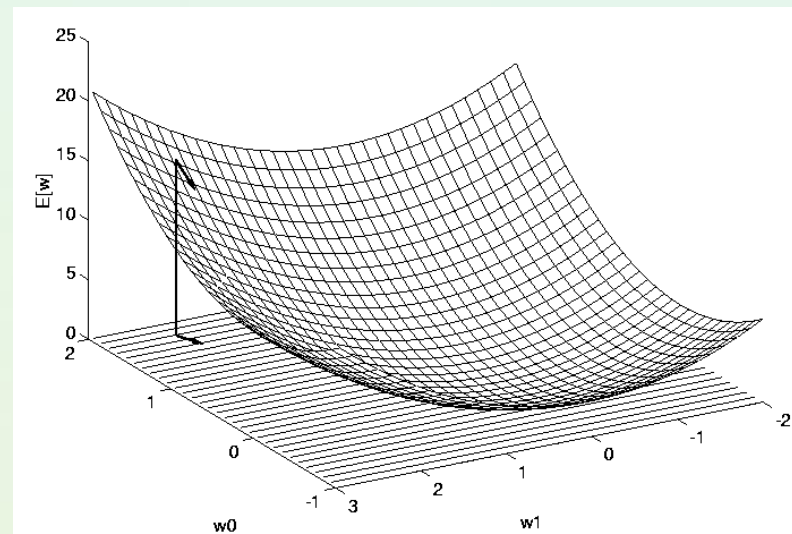


Algorytm spadku gradientu

- Simple Gradient Descent Algorithm
 - Applicable to different type of learning (nie tylko ANN)
- Algorithm *Train-Perceptron* ($D \equiv \{ \langle x, o(x) \equiv d(x) \rangle \}$)
 - Initialize all weights w_i to random values
 - WHILE not all examples correctly predicted DO
 - FOR each training example $x \in D$
 - Compute current output $o(x)$**
 - FOR $i = 1$ to n**
 - $w_i \leftarrow w_i + r(t - o)x_i$ //delta perceptron learning rule**

- Definition: Gradient

$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$



Uwagi nt. uczenia pojedynczego neuronu

- Błąd popełniony przez neuron przy prezentacji j -tego przykładu \mathbf{x}^j

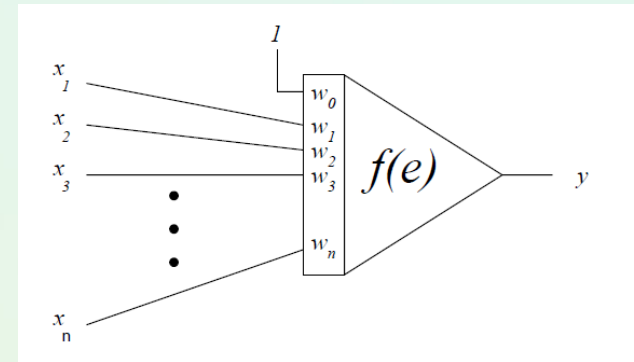
$$\delta^j = z^j - y^j$$

- Agregacja – błąd średniokwadratowy MSE

$$Q = \frac{1}{2} \sum_{j=1}^N (z^j - y^j)^2 = \sum_{j=1}^N Q^j, \quad Q^j = \frac{1}{2} (\delta^j)^2$$

- Korekta wag (Widrow, Hoff 1962):

$$\Delta w_i = \eta \cdot \delta^j \cdot x_i^j$$

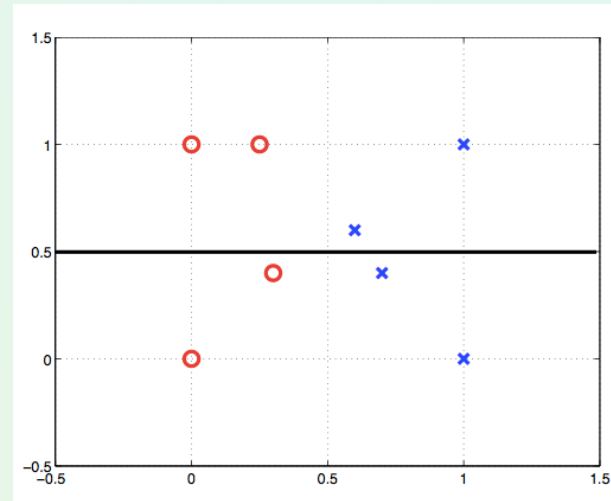
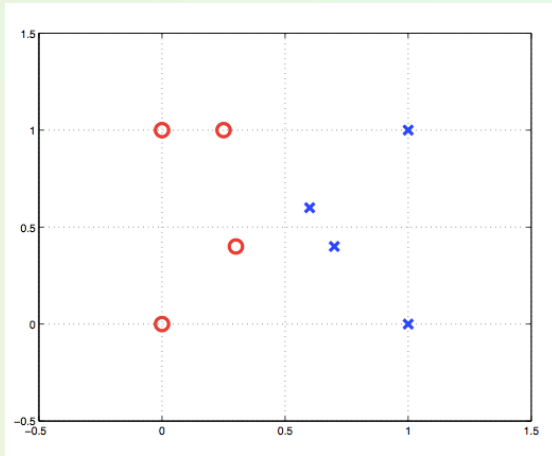


- Reguła delta – neurony nieliniowe

$$\Delta w_i = \eta \cdot \delta^j \cdot (1 - y^j) \cdot y^j \cdot x_i^j$$

Perceptrony uczone są przyrostowo

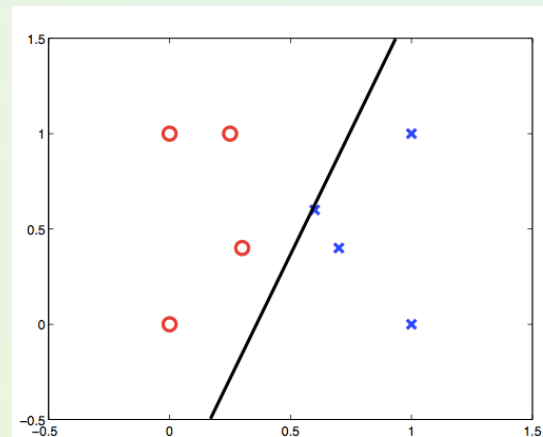
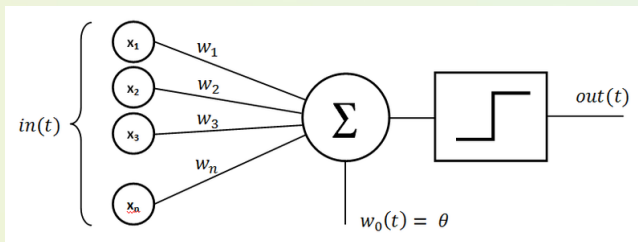
Krokowe znajdowanie odpowiednika płaszczyzny granicznej



początkowy układ wag: 0 ; 1; -0.5

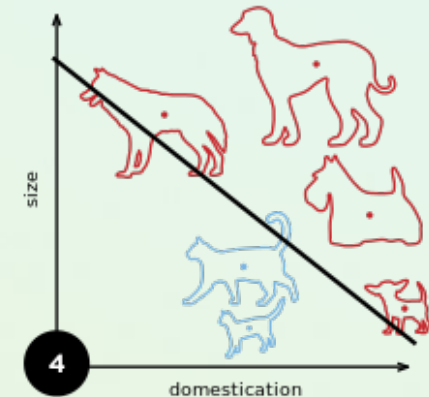
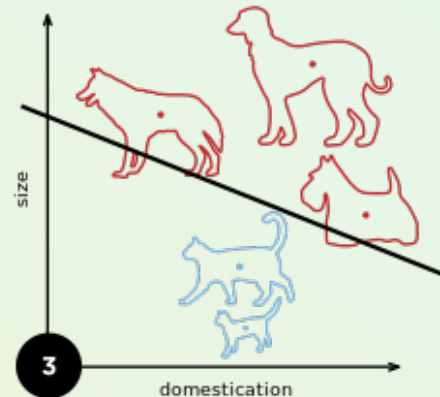
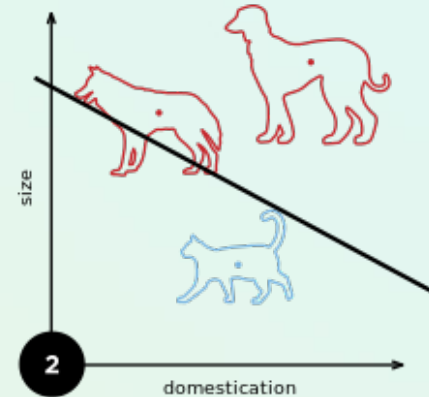
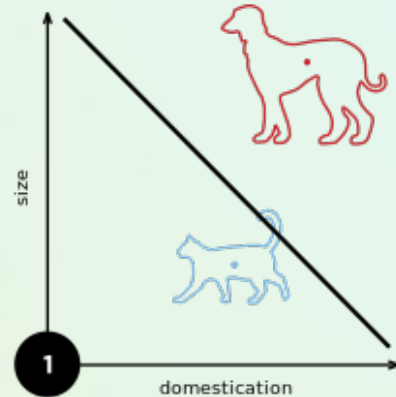
Class 1 has points $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$, $\begin{bmatrix} 0.6 \\ 0.6 \end{bmatrix}$, $\begin{bmatrix} 0.7 \\ 0.4 \end{bmatrix}$

Class 2 has points $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $\begin{bmatrix} 0.25 \\ 1 \end{bmatrix}$, $\begin{bmatrix} 0.3 \\ 0.4 \end{bmatrix}$



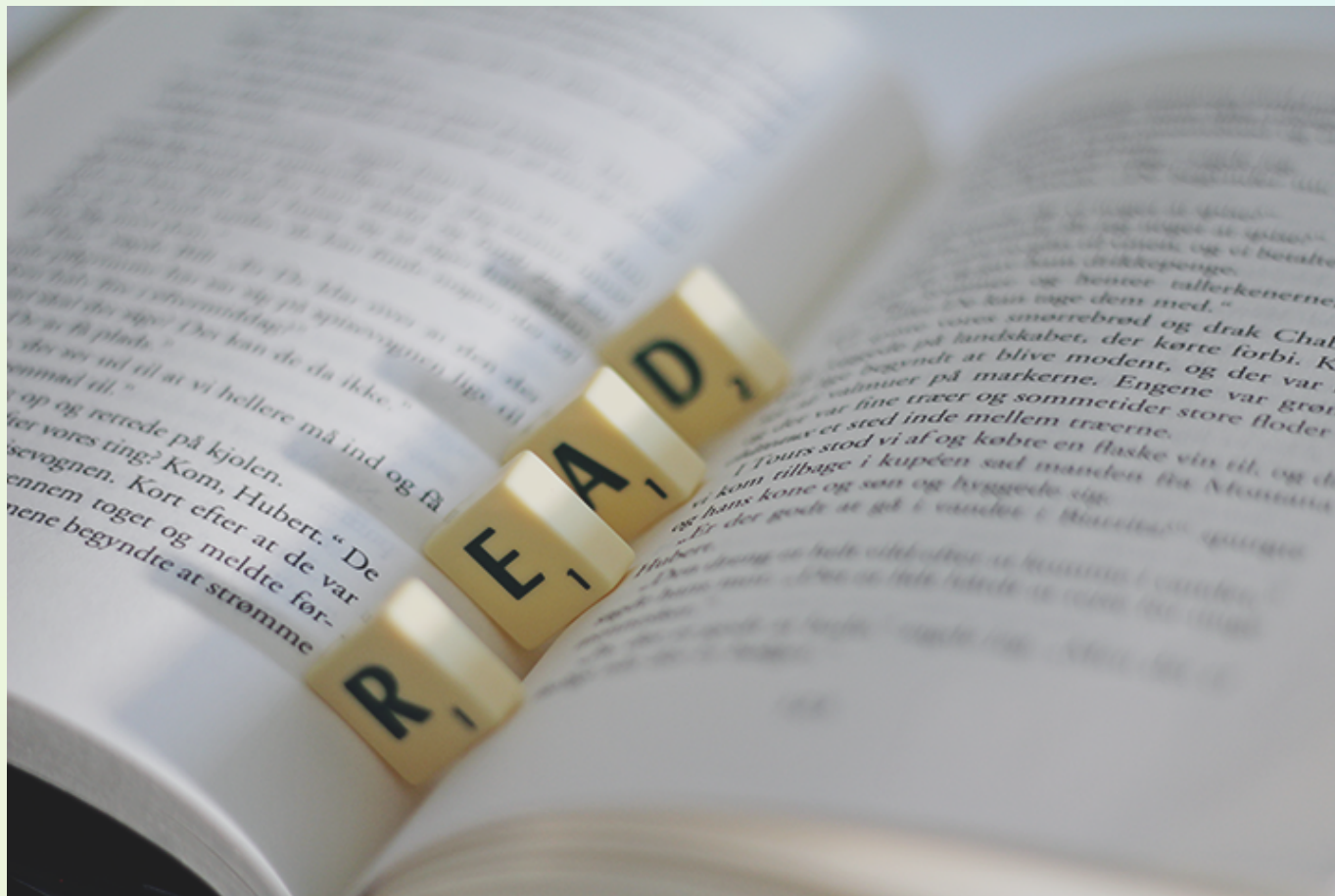
Płaszczyzna po pierwszej epoce przyrostowego pokazania przykładów

Przyrostowy wpływ kolejnych przykładów



- za Wikipedia

Koniec tej części wykładu



Przypomnij sobie wcześniejsze wykłady

Czytaj samodzielnie dodatkowe materiały