

Fisherova lineární diskriminace – příklad

Bylo provedeno měření objemu hipokampu a amygdaly (v cm^3) u 3 pacientů s Alzheimerovou chorobou () a 3 kontrolních subjektů () (označení D – diseased, H – healthy). Naměřené hodnoty objemu hipokampu a amygdaly u pacientů (resp.) a kontrol (resp.) byly zaznamenány do matic resp. :

Určete, zda testovací subjekt patří do skupiny pacientů či kontrolních subjektů pomocí Fisherovy lineární diskriminace.

Řešení:

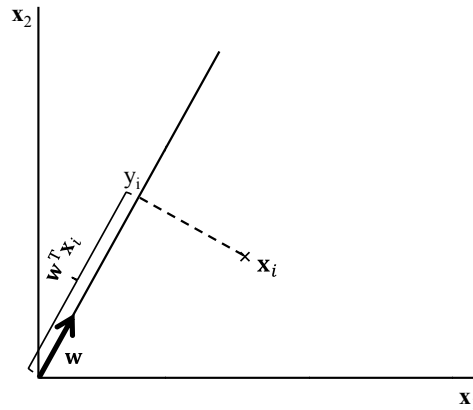
Principem Fisherovy lineární diskriminace je transformace do jednorozměrného (1D) prostoru tak, že chceme maximalizovat vzdálenost skupin (odráží se v čitateli Fisherova diskriminačního kritéria) a minimalizovat variabilitu uvnitř skupin (odráží se ve jmenovateli Fisherova diskriminačního kritéria). Fisherovo diskriminační kritérium je tedy ve tvaru:

$$\frac{(\mathbf{x} - \mathbf{c}_D)^T \mathbf{w}}{\sqrt{\mathbf{w}^T \mathbf{S}_D \mathbf{w}}}$$

kde je projekce centroidu pacientů do 1-D prostoru, je projekce centroidu kontrol, je rozptyl uvnitř třídy pacientů po projekci do 1-D prostoru a je rozptyl uvnitř třídy kontrol. Centroidy jsou vícerozměrné průměry pro třídu pacientů a kontrol:

$$\bar{\mathbf{x}}_D = \frac{1}{n_D} \sum_{i=1}^{n_D} \mathbf{x}_i, \quad \bar{\mathbf{x}}_H = \frac{1}{n_H} \sum_{i=1}^{n_H} \mathbf{x}_i$$

kde je hodnota první proměnné u -tého subjektu a je počet proměnných. Projekce centroidů do 1-D prostoru mohou být vypočítány jako a , kde je váhový vektor udávající směr 1-D prostoru, do něhož promítáme. Obecně může být průmět jakéhokoliv bodu do 1D prostoru vypočítán jako a znázorněn pomocí *Obrázku 1*.



Obrázek 1. Znázornění projekce bodu do 1-D prostoru daného směrovým vektorem w . Bod reprezentuje i -tý subjekt a y_i je jeho projekce. Osy x_1 a x_2 odpovídají dvěma proměnným.

Rozptyl uvnitř třídy pacientů po projekci do 1-D prostoru () lze vypočítat jako čtverec vzdáleností projekcí bodů odpovídajících jednotlivým pacientům od projekce centroidu:

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

kde Σ je kovarianční matice pacientů. Obdobně je možné rozptyl uvnitř třídy kontrol po projekci do 1-D prostoru () vypočítat jako:

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

kde Σ je kovarianční matice kontrol.

Dále si rozepíšeme součet rozptylů uvnitř jednotlivých tříd po transformaci do 1D prostoru, který se vyskytuje ve jmenovateli Fisherova diskriminačního kritéria:

kde S je suma čtverců variability uvnitř skupin a lze ji vypočítat jako: $S = \sum_{i=1}^n (x_i - \bar{x})^2$. V obecném případě, kdy nejsou vyvážené počty subjektů ve skupinách, se počítá vážená suma čtverců variability uvnitř skupin jako $\sum_{i=1}^n w_i (x_i - \bar{x})^2$.

Čítenel Fisherova diskriminačního kritéria si můžeme rozepsat jako:

kde je suma čtverců variability mezi skupinami.

Fisherovo diskriminační kritérium tedy můžeme vyjádřit jako: _____

Chceme maximalizovat _____, proto _____ zderivujeme a položíme výraz roven 0:

$$\frac{\partial}{\partial \mathbf{w}} \left(\frac{\sum_{i=1}^k n_i (\mathbf{x}_i - \mathbf{c}_i)^T \mathbf{w}}{\sum_{i=1}^k n_i} \right) = 0$$

Víme, že _____ má směr _____, protože _____, kde _____ je nějaký skalár. U vektoru _____ nás nezajímá jeho modul (tzn. velikost), jen jeho směr, proto můžeme pominout skalární členy _____ a _____. Dostáváme tedy:

Po odvození vzorečku pro výpočet váhového vektoru _____ do něj můžeme dosadit konkrétní hodnoty centroidů (vícerozměrných průměrů) pro třídu pacientů a kontrol, tzn. _____, _____. Pro výpočet sumy čtverců variability mezi skupinami _____ využijeme výběrové kovarianční matice _____ a _____ (výpočet vícerozměrných průměrů a výběrových kovariančních matic lze nalézt ve Cvičení 1). Suma čtverců variability mezi skupinami bude tedy spočítána jako _____ a její inverze jako _____.

Váhový vektor (diskriminační směr) _____ poté tedy můžeme spočítat následujícím způsobem:

Protože nás nezajímá modul váhového vektoru, ale jen jeho směr, můžeme váhový vektor přeškálovat na: _____. Nyní můžeme vypočítat průměty centroidů do 1D prostoru:

A následně vypočteme průmět hraničního bodu v 1D prostoru: _____

Hraniční bod lze vypočítat i takto: $-\quad -$
 $-\quad -\quad -\quad -\quad -\quad -$

— (protože jsme váhový vektor přeškálovali pomocí vynásobení $-$, musíme vynásobit i $-\quad -$ a pak získáváme -31).

Pokud chceme zařadit nový subjekt $-\quad -$ do jedné z daných tříd, musíme nejprve vypočítat jeho průmět do 1-D prostoru:

Průmět následně srovnáme s hraničním bodem: protože $-\quad -$, subjekt zařadíme do skupiny kontrolních subjektů (kontrolní subjekty leží nalevo od hraničního bodu, protože centroid kontrolních subjektů má menší (=více negativní) hodnotu než hraniční bod).

Po výpočtu váhového vektoru a hraničního bodu můžeme určit obecnou rovnici hranice (normálou hraniční přímky je váhový vektor $-\quad -$):

Pro vykreslení hranice je vhodné vyjádřit hranici ve tvaru:
 Nová osa, do níž se promítá, má směr odpovídající váhovému vektoru $-\quad -$ (je kolmá k hranici) a prochází počátkem $-\quad -$ a lze ji tedy vyjádřit obecnou rovnicí jako:

Pokud nás zajímají souřadnice hraničního bodu $-\quad -$ v původním prostoru, využijeme znalosti, že hraniční bod je průsečík hranice a nové osy:

$$\begin{aligned} &----- \\ &----- \\ &----- \end{aligned}$$

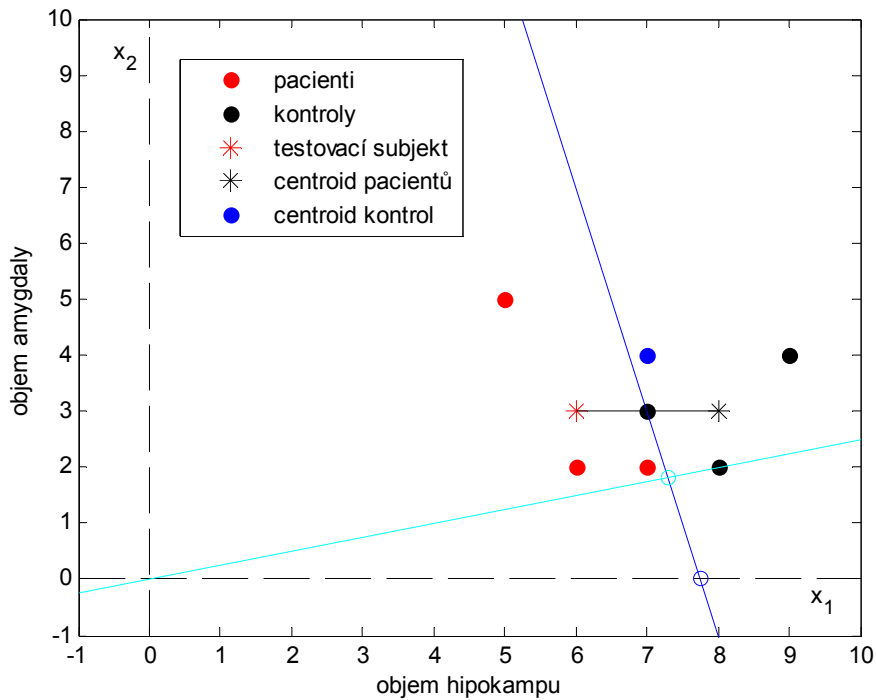
Souřadnici $-\quad -$ pak vypočítáme z druhé rovnice jako: $-\quad -$

Souřadnice hraničního bodu v původním prostoru jsou tedy: $-\quad -$

Ověření, že po projekci hraničního bodu dostanu hodnotu -31 : $-----$
 $-----$

$$-----$$

Klasifikaci pomocí Fisherovy lineární diskriminační analýzy si na závěr znázorníme pomocí *Obrázku 2*.



Obrázek 2. Znázornění klasifikace pomocí Fisherovy lineární diskriminační analýzy. Klasifikační hranice je znázorněna tmavě modře, nová osa, do níž se promítá, světle modře a hraniční bod je vyznačen tmavě modrým prázdným kolečkem. Původní osy x_1 a x_2 odpovídající dvěma proměnným (objemu hipokampu a amygdaly) jsou znázorněny čárkovanými čarami.

Poznámka: Pokud bychom váhový vektor znormovali, hraniční bod by přímo ležel ve vzdálenosti od počátku:

_____ (tzn. hraniční bod leží ve vzdálenosti _____ od počátku v původních souřadnicích)