



# Machine Learning

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir



<https://www.aparat.com/mehran.safayani>



[https://github.com/safayani/machine\\_learning\\_course](https://github.com/safayani/machine_learning_course)



Department of Electrical and computer engineering, Isfahan university of technology, Isfahan, Iran

# Classification

Dr Mehran Safayani

# Classification

همانند مسئله رگرسیون است با این تفاوت که  $y$  مقدار گسسته می گیرد.

Binary classification (دسته بندی دودویی):

$y \in \{c_1, c_2\}$  ,  $c_i$  : برچسب کلاس

$y \in \{-1, +1\}$  ,  $y \in \{0, +1\}$

هیچ ترتیبی بین دو کلاس  
وجود ندارد

Multi-class classification (دسته بندی چند کلاسه):

$y \in \{c_0, c_2, \dots, c_{k-1}\}$

$y \in \{0, 1, 2, \dots, k-1\}$

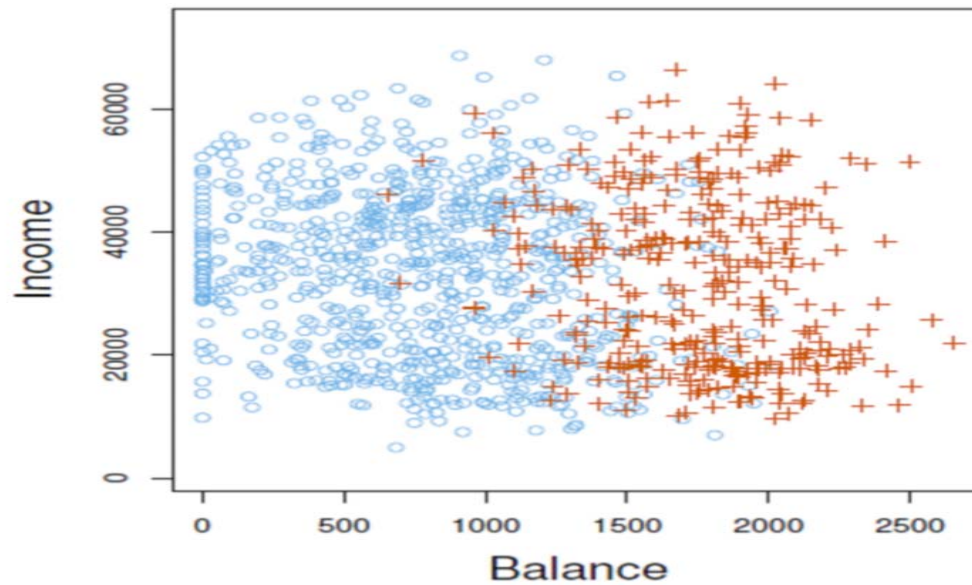
## مثال

می خواهیم بدانیم که یک تراکنش کارت اعتباری جعلی است یا اصلی. با استفاده از اطلاعاتی نظیر تراکنش های گذشته و ....

مثالی دیگر:

نارنجی : بد حساب ( مبلغ ماهیانه را به موقع پرداخت نمی کند)

آبی: خوش حساب



## دسته بند

یک دسته بند فضای ورودی را به ناحیه هایی متعلق به هر کلاس تقسیم میکند. مرز این ناحیه ها را مرز تصمیم (Decision boundry) میگویند.

دسته بند میتواند خطی یا غیر خطی باشد.

دسته بند خطی از ترکیب خطی ویژگی ها استفاده میکند.

$$y = f(\sum \omega_j x_j)$$

$f$  تابعی است که ترکیب خطی را به ویژگی ها تبدیل میکند.

$$f(x) = \begin{cases} 1 & \text{if } \omega^T x > \theta \\ 0 & \text{otherwise} \end{cases}$$

$$w_1 x_1 + w_2 x_2 + b = 0$$

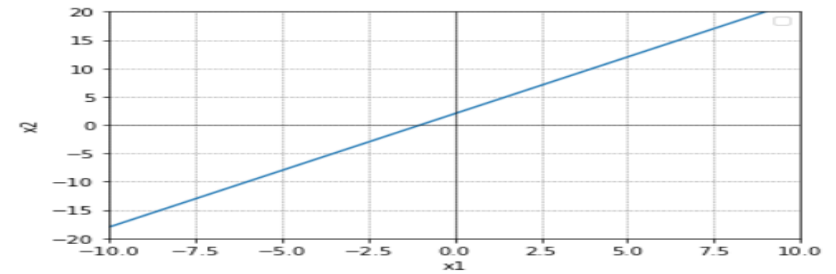
$$w_1 x_1 + w_2 x_2 + b = \theta$$

$$b - \theta = 0$$

$\longleftrightarrow$   
 $b'$

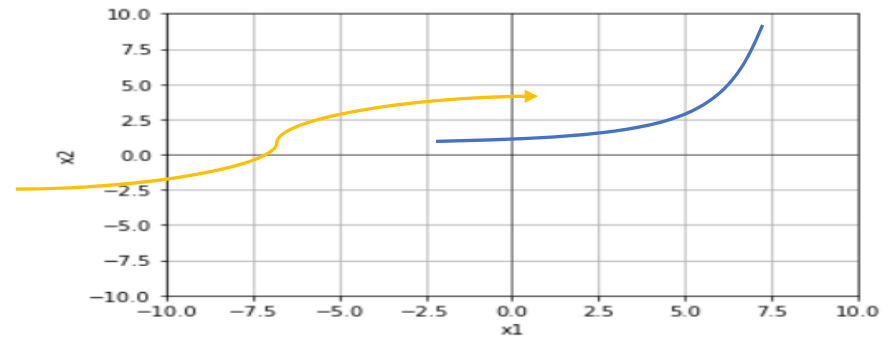
$$w_1 x_1 + w_2 x_2 + b = 0$$

$$x_2 = -\frac{w_1}{w_2} x_1 - \frac{b}{w_2}$$



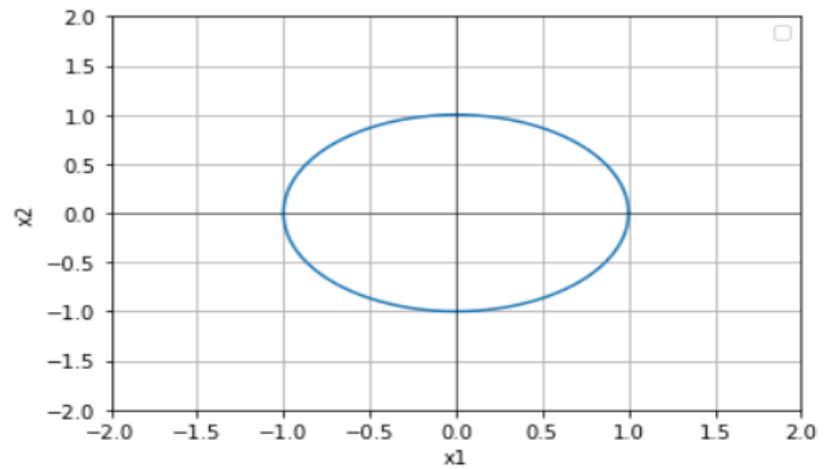
$$w_1 x_1 + w_2 x_2 + w_3 x_1^2 = 0$$

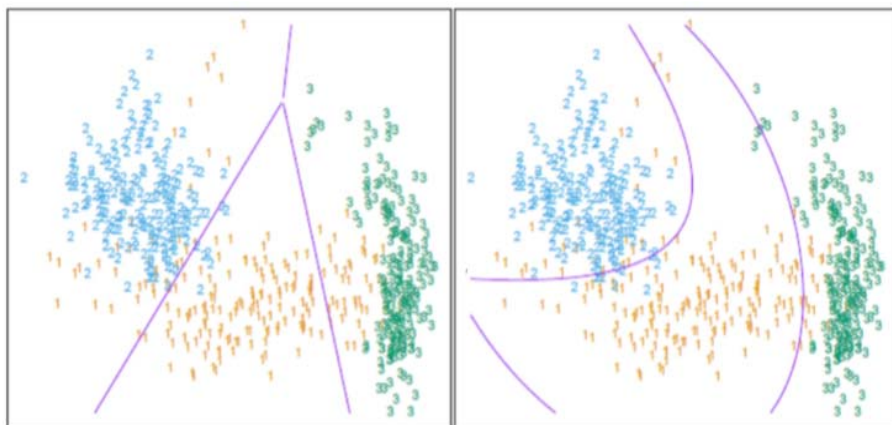
مرز تصمیم



$$w_1 x_1^2 + w_2 x_2^2 + b = 0, \quad w_1 = 1, \quad w_2 = 1, \quad b = -1$$

$$x_1^2 + x_2^2 = 1$$





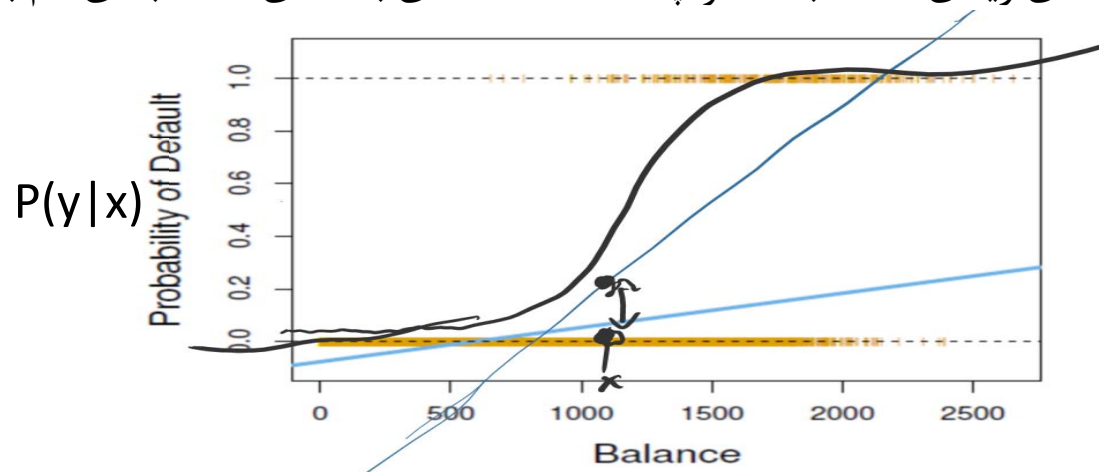
**FIGURE 4.1.** The left plot shows some data from three classes, with linear decision boundaries found by linear discriminant analysis. The right plot shows quadratic decision boundaries. These were obtained by finding linear boundaries in the five-dimensional space  $X_1, X_2, X_1X_2, X_1^2, X_2^2$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

نکته: مرز تصمیم در کلاس بند خطی در فضای ویژگی اولیه (نه فضای ویژگی توسعه یافته) لزوماً خطی نیست.



# دسته بندی به صورت حالت خاص یک مسئله رگرسیون

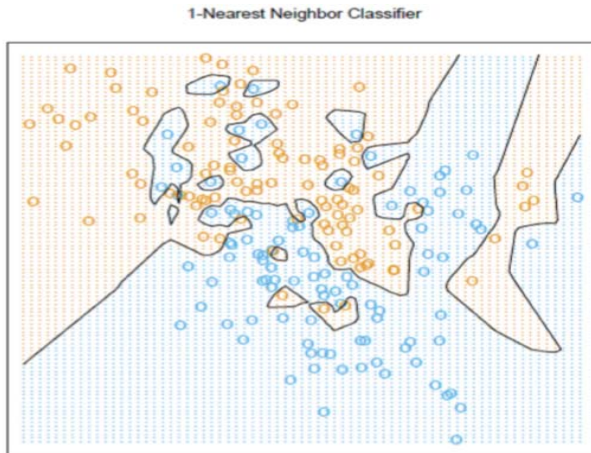
- $y=0, y=1$  را میتوان احتمال تعلق به هر یک از کلاس ها در نظر گرفت.
- خروجی مدل میتواند اعداد منفی یا بزرگتر از یک باشد.
- تعداد کمی نقطه میتواند جهت خط را به میزان زیادی تغییر دهد مانند اضافه کردن چند نقطه با بالانس خیلی زیاد
- SE یک تابع هزینه مناسب برای هدف دسته بندی نیست
- در مسئله دسته بندی بایستی تعداد نقاط با دسته بندی اشتباه حداقل شود.
- SE فاصله برچسب ها تا پیش بینی را حداقل میکند.
- اگر MSE خیلی کم باشد میتوان تضمین کرد که خطای دسته بندی هم کم است ولی برعکسش صحیح نیست یعنی یک تابع رگرسیون ممکن است خطای زیادی داشته باشد هر چند تعداد داده های با خطای دسته بندی کم باشد.



# Nearest Neighbor

## دسته بند نزدیکترین همسایگی

- نمونه های نزدیک برچسب یکسانی دارند.
- نزدیکی را با فاصله اقلیدوسی میتوان اندازه گرفت.
- با داشتن یک مجموعه آموزشی  $S_{train}$  و یک داده  $x$  به دنبال  $x^*$  که نزدیکترین نقطه به  $x$  است میگردیم و برچسب را برابر  $y^*(x^*)$  (برچسب) با قرار میدهیم.



**FIGURE 2.3.** The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1), and then predicted by 1-nearest-neighbor classification.

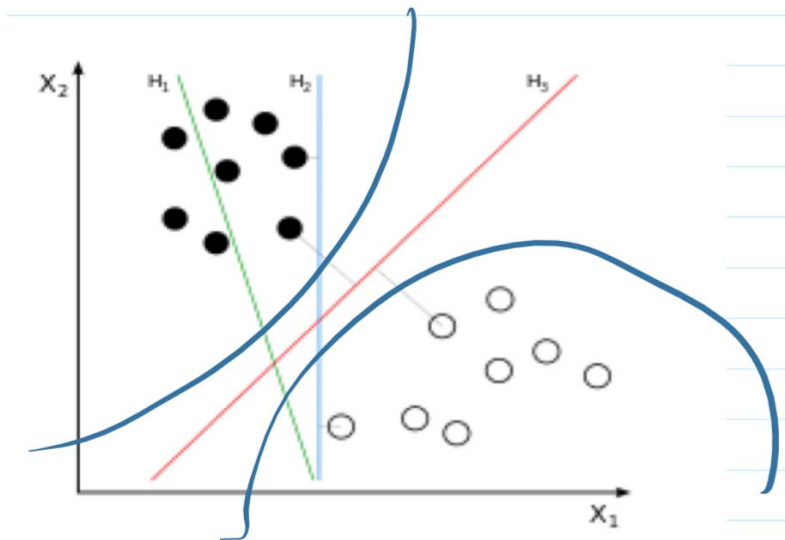
# مرزهای تصمیم خطی

## Linear decision boundaries

- فرض کنید که مرز تصمیم خطی است. **hyperlinear**

- فرض کنید که مسئله تفکیک پذیر خطی است.

- ایده ماشین بردار پشتیبان (**Support vector machine(svm)**)



# مرز تصمیم غیر خطی

- استفاده از کرنل

- ویژگی های quadratic

## دسته بند بهینه بیز

فرض کنید  $X$  داده و  $y$  برچسب است. اگر  $P(y|x)$  را به صورت دقیق بدانیم چگونه دسته بندی میکنیم؟؟

$$\hat{Y} = \underset{Y}{\operatorname{argmax}} P(Y = y \mid X = x) \quad (\text{دسته بند بهینه بیز})$$

Bayes rule:

The diagram shows the Bayes' rule formula with four labels and arrows pointing to the corresponding parts of the equation:

- likelihood** points to  $P(X \mid Y)$  in the numerator.
- prior** points to  $P(Y)$  in the numerator.
- posterior** points to  $P(Y \mid X)$  on the left side of the equation.
- normalizer** points to  $P(X)$  in the denominator.

$$\text{posterior} \rightarrow P(Y \mid X) = \frac{P(X \mid Y)P(Y)}{P(X)}$$

چرا یادگیری دسته بند بهینه مشکل است؟؟

Sky	Temp	Humid	Wind	Water	Forecast	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cold	Change	Yes

فرض کنید  $k$  کلاس داریم. فرض کنید  $x$  شامل  $d$  ویژگی باینری است.

prior:  $P(Y)$  , likelihood:  $P(X | Y)$

توزیع بر روی مقادیر ویژگی ها:

$$P(X = x | Y = y)$$

$$P(S = s, T = w, H = n, W = S, wa = w, F = s | E = y)$$

تعداد پارامترهای مدل:

$$k(2^d - 1)$$

تعداد زیاد پارامتر که نیاز به تعداد زیاد داده برای تخمین دارد

## Conditional independence (استقلال شرطی)

$$X \perp Y \Rightarrow p(x, y) = p(x)p(y)$$

$$\forall i, j, k \quad p(X = i \mid Y = j, Z = k) = p(X = i \mid Z = k)$$

X به شرط دانستن Z از Y مستقل است.

$$P(\text{thunder} \mid \text{Rain, Lightening}) = P(\text{thunder} \mid \text{Lightening})$$

$$T \perp R \quad \times, \quad T \perp R \mid L \quad \checkmark$$

$$P(T, R \mid L) = P(T \mid R, L)P(R \mid L) = \frac{P(T, R, L)}{P(L)} = \frac{P(T, R, L)}{P(R, L)} \frac{P(R, L)}{P(L)} = P(T \mid L)P(R \mid L)$$

$$P(X, Y \mid Z) = P(X \mid Z)P(Y \mid Z)$$

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A | B) \propto P(B | A)P(A)$$

چرا مخرج حذف شده است؟؟

$$\text{IF } P(A=0) = P(A=1) \longrightarrow p(A|B) \propto p(B|A)$$

احتمال lightning چقدر است؟؟

$$P(L | T, R) \propto P(T, R | L) * P(L) \longleftarrow \text{یکنواخت}$$

$$P(T, R | L = 0): \quad P(0,0 | L=0), P(0,1 | L=0), P(1,0 | L=0), P(1,1 | L=0)$$

$$P(T, R | L = 1): \quad P(0,0 | L=1), P(0,1 | L=1), P(1,0 | L=1), P(1,1 | L=1)$$



**Estimate:**  $P(T, R | L)$   $2(2^2 - 1) = 6$  params

با فرض استقلال شرطی:

$$P(T, R | L) = P(T | L) P(R | L) \quad 2(2 - 1) + 2(2 - 1) = 4 \text{ params}$$

# Naive Bayes assumption

ویژگی ها به شرط دانستن کلاس مستقل هستند:

$$P(x_1, x_2 | Y) = p(x_1 | x_2, y)P(x_2 | Y) = P(x_1 | Y)p(x_2 | Y)$$

به طور کلی:

$$P(x_1, \dots, x_d | Y) = \prod_{j=1}^d P(x_j | Y)$$

حال چند پارامتر خواهیم داشت؟

اگر فرض استقلال شرطی داشته باشیم:

$$K(2^d - 1)$$

K: تعداد کلاس

2: تعداد حالت

d: تعداد ویژگی ها

With Naive bayes:

$$K(d)$$

کاهش قابل ملاحظه پارامترها. شاید حتی زیادی کم شد.

## تعریف دسته بند Naive bayes

Prior:  $P(y)$  , ویژگی مستقل به شرط کلاس  $d$  ,  $P(x_j | y)$  : برای هر  $x_j$

قانون تصمیم گیری:

$$\hat{y} = f_{NB}(x) = \underset{y}{\operatorname{argmax}} P(y) \prod_j P(x_j | y)$$

$$\text{prior: } P(Y = y) = \frac{\text{count}(Y = y)}{N}$$

$$\text{likelihood: } P(X_j = x_j | Y = y) = \frac{\text{count}(X = x_j, Y = y)/N}{\text{count}(Y = y)/N}$$

$$P(X | Y) = \frac{P(x, y)}{P(y)}$$

یک مشکل ممکن است رخ دهد:

$$P(x_1 | Y=b) = 0$$

آنگاه:

$$P(Y = b | X) = P(Y = b | X_1)P(Y = b | X_2) \cdots P(Y | X_d) = 0$$

راه حل: Smoothing

$$P(Y = b) = \frac{\text{count}(Y = b) + \lambda}{N + k\lambda}$$

K: تعداد کلاس  
 $\lambda: 1$

$$P(x_1 = a | y = b) = \frac{\text{count}(x_1 = a, y = b) + \lambda}{\text{count}(y = b) + A\lambda}$$

A: تعداد حالات ویژگی  $x_1$   
 $\lambda: 1$

## مثال

Text classification / classify email {spam , not spam} / classify news article topic

هر مقاله حداقل از 1000 کلمه تشکیل شده و هر کلمه میتواند 10000 حالت داشته باشد.  
 $P(x|y)$  بسیار بزرگ است.

تعداد حالات:

$$\propto 2(10000)^{1000} = 2(10)^{4000}$$

With NB:

$$\propto 2 * 10000 * 1000 = 2(10)^7$$

- With NB:

- $f_{NB}(x) = \underset{y}{\operatorname{argmax}} P(Y = y) \prod_{j=1}^{\text{lengthdoc}} P(x_j | Y = y)$

- $P(x_j | Y = y) = \frac{\frac{\text{count}(x_j=x, Y=y)}{N}}{P(Y=y) = \frac{\text{count}(Y=y)}{N}}$

تعداد کل منابع

کل منابع در پایگاه داده

$$\frac{\text{count}(X = x_j, Y = y)}{\text{count}(Y = y)}$$

تعداد منابع در تایپیک Y

# Bag of words(BOW)

$x_1$	$x_2$	...	$x_{1000}$
-------	-------	-----	------------

## داده های پیوسته

- گسسته سازی
- استفاده از توزیع های پیوسته

$$P(x_1, x_2, \dots, x_n \mid c_k) \propto P(c_k) \cdot \prod_{i=1}^n P(x_i \mid c_k)$$

$$P(x_i \mid c_k) = N(\mu_i, \sigma_i)$$

تخمین به وسیله روش ML