



” ”

# داده کاوی

“ “



دکتر حسام عمرانپور



مراجع: Data mining concepts and techniques, Jiawei Han, Michelin

Cumber and Jian Pei

داده کاوی مهدی اسماعیلی

تنظیم: فاطمہ حیدری

حدیثہ پورعلی

# فصل 3

دسته بندی

# دسته بندی (classification)

ورودی: داده به همراه برچسب

#	طول	ارتفاع	نوع
1	7	4	اتوبوس
2	6,5	4,5	اتوبوس
3	7,5	4,5	اتوبوس
4	9	4,5	اتوبوس
5	3	1,5	پراید
6	2,5	1,7	پراید
7	2	1,6	پراید

	خانه دارد	تعداد فرزندان	انگوشمال دارد	حقوق دریافتی	وام راپس داده است؟
#1	1	2	1	800	بله
#2	0	1	0	750	بله
#3	0	2	1	700	بله
#4	1	0	1	650	خیر
	⋮	⋮	⋮	⋮	⋮

# دسته بندی (classification)

دسته بندی دو فاز دارد:

1. آموزش (train):

یک مجموعه داده به همراه برچسب آن ها به مدل داده می شود (یا اینکه داده ها ذخیره می شوند) و براساس این داده ها پارامتر های مدل یادگرفته می شوند.

2. آزمایش (test):

یک داده بدون برچسب داده می شود و باید مشخص شود در چه کلاسی قرار می گیرد.

# روش های دسته بندی

## 1) روش درخت تصمیم decision tree

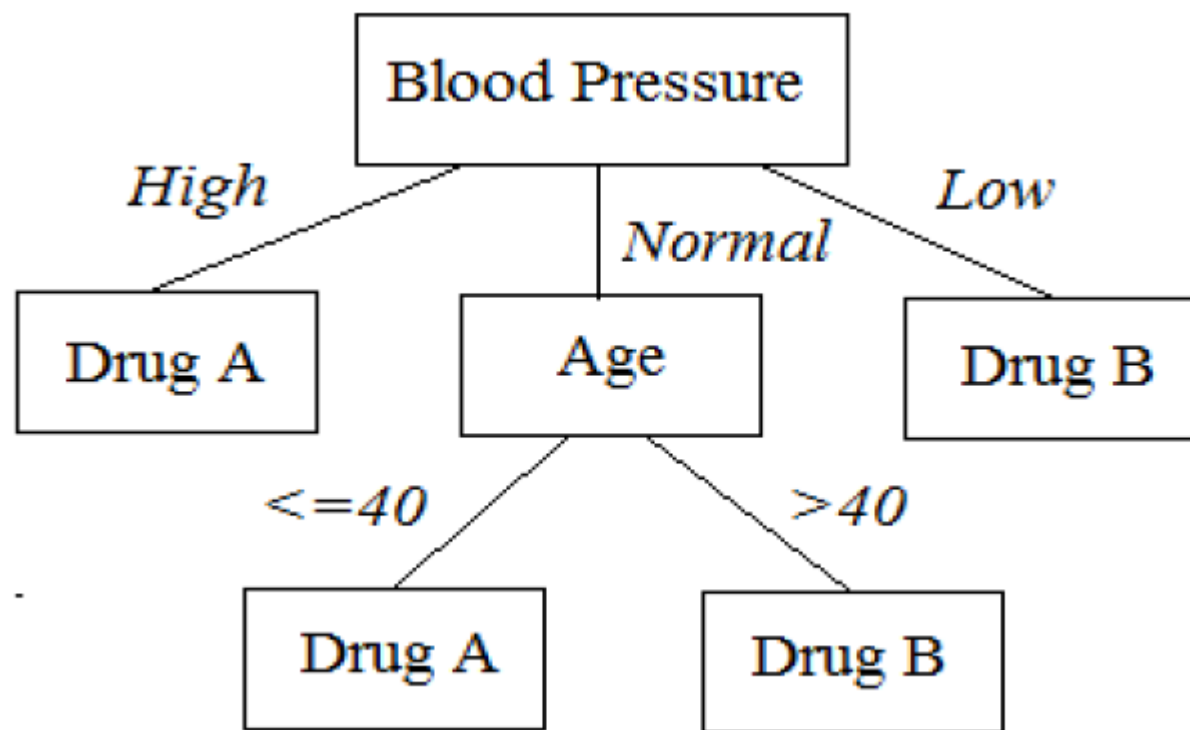
یک درخت است که :

1. روی هر گره یک ویژگی ذخیره می شود (داده های category)

2. روی هر گره یک ویژگی به همراه یک (یا چند) آستانه ذخیره می شود.

# درخت تصمیم

ID	Sex	Age	Blood P.	Drug
1	Male	20	Normal	A
2	Female	73	Normal	B
3	Male	37	High	A
4	Male	33	Low	B
5	Female	48	High	A
6	Male	29	Normal	A
7	Female	52	Normal	B
8	Male	42	Low	B
9	Male	61	Normal	B
10	Female	30	Normal	A
11	Female	26	Low	B
12	Male	54	High	A





# درخت تصمیم

- سعی براین است که درخت تصمیم کمترین پیچیدگی را داشته باشد
- در داده کاوی در فاز train با افزایش پیچیدگی مدل، از مکانی به بعد خطای تست افزایش می یابد. در واقع اگر مدل زیاد آموزش ببیند و یا زیاد پیچیده شود، قابلیت تعمیم خود را در داده های test از دست می دهد. به این عمل اصطلاحا بیش برآزش داده ها (overfitting) می گویند.



# درخت تصمیم

پیچیدگی محاسباتی در مدل درخت تصمیم می تواند عمق درخت یا تعداد نود ها (ویژگی ها) باشد.

# درخت تصمیم

■ دو مسئله:

■ 1. شرط پایان الگوریتم:

➤ همه ی نمونه ها از مجموعه آموزش در یک کلاس باشند.

➤ به حداکثر عمق درخت برسیم.

➤ تعداد نمونه های موجود در گره از حداقل تعدادی که به عنوان آستانه مشخص می شود کمتر باشد.

➤ مقادیر محاسبه شده برای انتخاب ویژگی برای هیچ ویژگی بیشتر از حد آستانه نیست.

■ 2. مناسب ترین ویژگی در هر مرحله:

# معیارهای انتخاب ویژگی در درخت تصمیم:

## 1. Information Gain:

- $Information\ Gain_A(D) = Entropy(D) - Entropy_A(D)$

- $Entropy(D) = - \sum_{i=1}^c p_i * \log_2 p_i$

- $Entropy_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} * Entropy(D_j)$

- $c$  تعداد کلاس ها :

- $p_i$  احتمال وجود داده در کلاس :

- $A_v$  تعداد اعضا (سطح های) ویژگی :

- $D_j$  می باشد  $v_j$  در آن ها برابر با  $A$  که ویژگی  $D$  قسمتی از داده های اولیه :

# معیارهای انتخاب ویژگی در درخت تصمیم:

مثال از information gain:

ID	Age	Income	Job	Computer
1	Old	Medium	Student	No
2	Middle	High	Teacher	No
3	Old	Low	Teacher	No
4	Young	Medium	Teacher	Yes
5	Young	Low	Teacher	Yes
6	Old	Medium	Student	Yes
7	Middle	Medium	Student	Yes
8	Young	High	Teacher	No
9	Old	High	Student	No
10	Middle	High	Student	No

# معیارهای انتخاب ویژگی در درخت تصمیم:

ادامه مثال از information gain

از آنجا که از 10 نمونه ی موجود در داده ها 4 نمونه دارای برچسب Yes و 6 نمونه ی دیگر

دارای برچسب No هستند پس داریم:

$$Entropy(D) = -\frac{4}{10} \text{Log}_2\left(\frac{4}{10}\right) - \frac{6}{10} \text{Log}_2\left(\frac{6}{10}\right) = 0.970$$

سطوح هر یک از ویژگی سن، درآمد و شغل عبارت است از:

$$\text{Domain}(\text{Age}) = \{\text{Old}, \text{Middle}, \text{Young}\}$$

$$\text{Domain}(\text{Income}) = \{\text{High}, \text{Medium}, \text{Low}\}$$

$$\text{Domain}(\text{Job}) = \{\text{Teacher}, \text{Student}\}$$

## معیار های انتخاب ویژگی در درخت تصمیم:

$$\begin{aligned} Entropy_{Age}(D) &= \frac{4}{10} \times \left( -\frac{3}{4} \text{Log}_2\left(\frac{3}{4}\right) - \frac{1}{4} \text{Log}_2\left(\frac{1}{4}\right) \right) + \\ &\quad \frac{3}{10} \times \left( -\frac{2}{3} \text{Log}_2\left(\frac{2}{3}\right) - \frac{1}{3} \text{Log}_2\left(\frac{1}{3}\right) \right) + \\ &\quad \frac{3}{10} \times \left( -\frac{1}{3} \text{Log}_2\left(\frac{1}{3}\right) - \frac{2}{3} \text{Log}_2\left(\frac{2}{3}\right) \right) = 0.875 \end{aligned}$$

$$\begin{aligned} Entropy_{Income}(D) &= \frac{4}{10} \times \left( -\frac{4}{4} \text{Log}_2\left(\frac{4}{4}\right) - \frac{0}{4} \text{Log}_2\left(\frac{0}{4}\right) \right) + \\ &\quad \frac{4}{10} \times \left( -\frac{1}{4} \text{Log}_2\left(\frac{1}{4}\right) - \frac{3}{4} \text{Log}_2\left(\frac{3}{4}\right) \right) + \\ &\quad \frac{2}{10} \times \left( -\frac{1}{2} \text{Log}_2\left(\frac{1}{2}\right) - \frac{1}{2} \text{Log}_2\left(\frac{1}{2}\right) \right) = 0.524 \end{aligned}$$

$$\begin{aligned} Entropy_{Job}(D) &= \frac{5}{10} \times \left( -\frac{3}{5} \text{Log}_2\left(\frac{3}{5}\right) - \frac{2}{5} \text{Log}_2\left(\frac{2}{5}\right) \right) + \\ &\quad \frac{5}{10} \times \left( -\frac{3}{5} \text{Log}_2\left(\frac{3}{5}\right) - \frac{2}{5} \text{Log}_2\left(\frac{2}{5}\right) \right) = 0.970 \end{aligned}$$

ادامه مثال از information gain:

سپس باید برای این 3 ویژگی مقدار آنروپی را محاسبه کنیم:

# معیار های انتخاب ویژگی در درخت تصمیم:

ادامه مثال از information gain: /

پس از آن مقدار Information Gain برای همه ی ویژگی ها محاسبه می شود: /

$$InformationGain(Age)=0.970-0.875=0.095$$

$$InformationGain(Income)=0.970-0.524=0.446$$

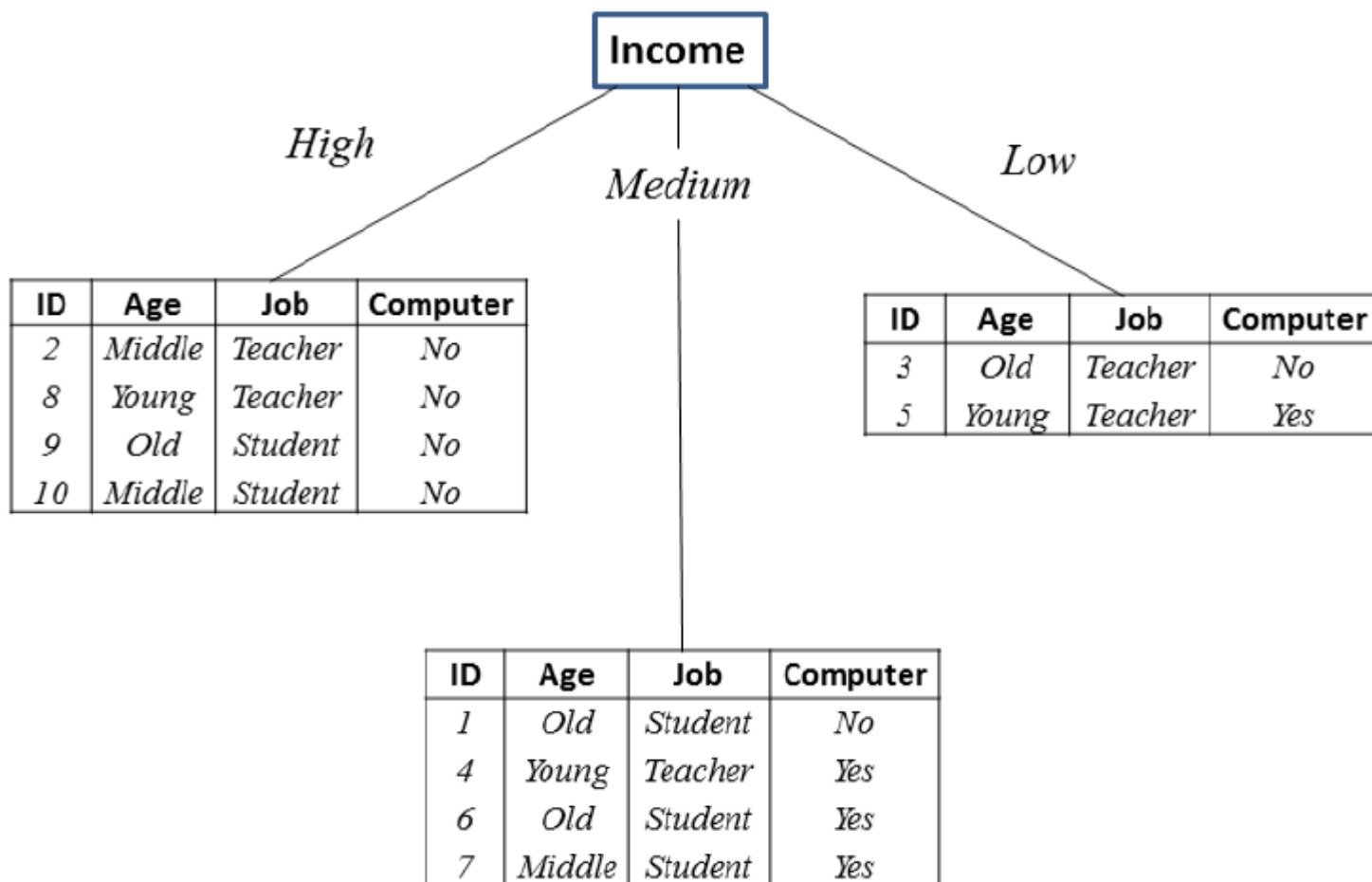
$$InformationGain(Job)=0.970-0.970=0$$

ویژگی که دارای بیشترین Information Gain است انتخاب میشود. /

پس ویژگی در آمد که دارای بیشترین مقدار است برای ریشه درخت انتخاب میشود(شکل زیر): /



# معیارهای انتخاب ویژگی در درخت تصمیم:



# معیار های انتخاب ویژگی در درخت تصمیم:

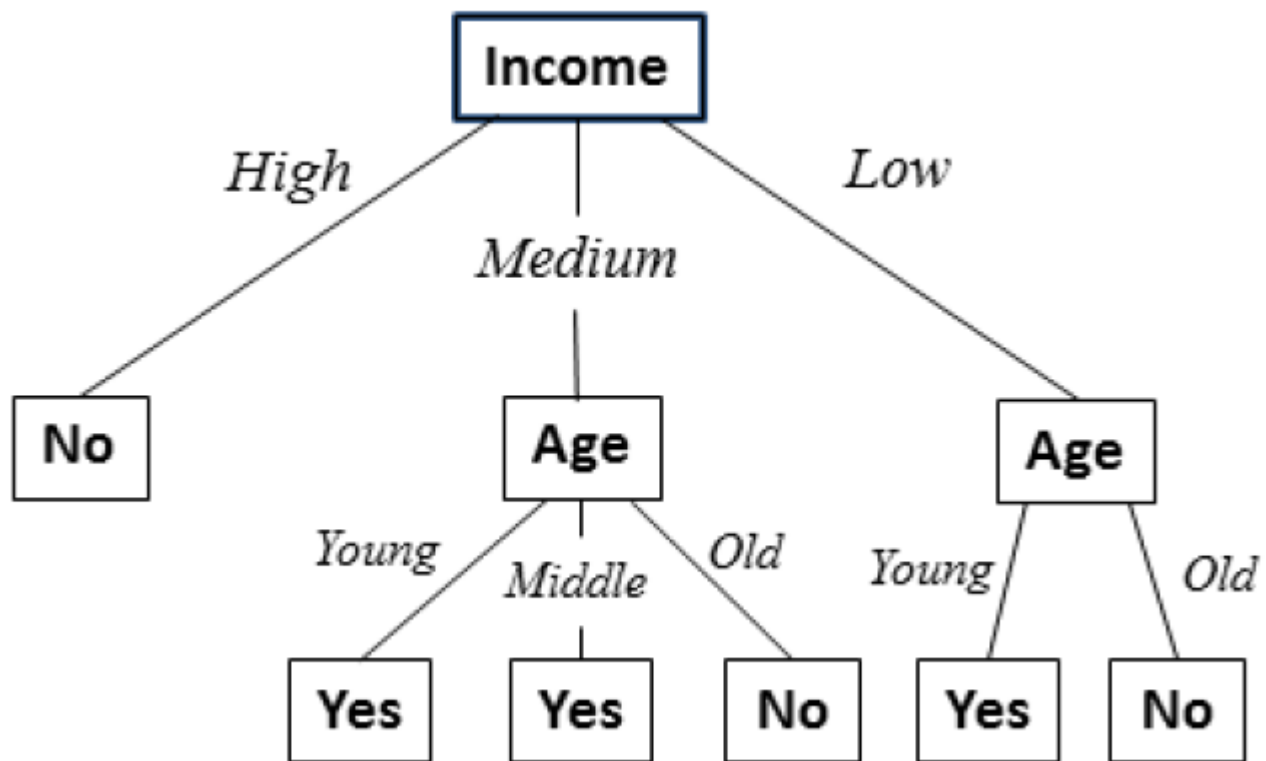
ادامه مثال از information gain:

درواقع داده های آموزشی به 3 زیرمجموعه افراز می شوند. برای هر یک از این زیرمجموعه ها همانند قبل بهترین ویژگی جهت انشعاب محاسبه می شود. این بار ویژگی درآمد در محاسبات شرکت نمی کند و از میان دیگر ویژگی ها یکی انتخاب می شود. این کار تا هنگامی که یکی از شروط توقف الگوریتم محقق شود، ادامه می یابد.

درخت تصمیم نهایی به شکل زیر است:

# معیارهای انتخاب ویژگی در درخت تصمیم:

ادامه مثال از information gain:



# معیارهای انتخاب ویژگی در درخت تصمیم:

## 2. Gini Index:

$$Gini(D) = 1 - \sum_{i=1}^c p_i^2$$

$$Gini_A(D) = \frac{|D_1|}{|D|} * Gini(D_1) + \frac{|D_2|}{|D|} * Gini(D_2) + \dots = \sum_{j=1}^v \frac{|D_j|}{|D|} * Gini(D_j)$$

ویژگی با Gini Index کوچکتر انتخاب می شود و یا MAX تغییر درجه ناخالص انتخاب می شود.

$$Gini(D) - Gini_A(D) =$$

معمولا  $v=2$  در نظر گرفته می شود.

# معیارهای انتخاب ویژگی در درخت تصمیم:

3. Gain Ratio: است Information Gain نرمال شده ی

$$\text{Gain Ratio}_A(D) = \frac{\text{Information Gain}_A(D)}{\text{splite info}(A)}$$

$$\text{splite info}(A) = - \sum_{j=1}^v \frac{|D_j|}{|D|} * \log_2 \frac{|D_j|}{|D|}$$

# معیارهای انتخاب ویژگی در درخت تصمیم:

## 4. Likelihood Ratio:

- $$Likelihood_A(D) = 2 * \ln 2 * |D| * Information\ Gain(A)$$

## 5. DKM: برای مسئله ی دو کلاسه است:

- $$DKM_A(D) = 2 * \sqrt{\frac{|D_1|}{|D|} * \frac{|D_2|}{|D|}}$$

# الگوریتم های درخت تصمیم

## 1. ID3

- از معیار Information Gain برای انتخاب ویژگی استفاده می کند.
- شرط توقف در هر برگ:
  - نمونه های باقی مانده در برگ متعلق به یک کلاس باشند.
  - برگ برای تمام ویژگی ها کوچکتر مساوی صفر Information Gain باشد.

# الگوریتم های درخت تصمیم

## 3. C4.5

از معیار Gain Ratio برای انتخاب ویژگی استفاده می کند

توقف هنگامی رخ می دهد که:

1. تعداد نمونه ها کمتر از مقدار مشخص شده ای باشد.

2. متعلق به یک کلاس باشند.



# ادامه روش های دسته بندی

## 2) روش K: K Nearest Neighbor (KNN) نزدیکترین همسایه

1. نزدیکترین داده به داده تست را از مجموعه آموزش پیدا میکند K

2. سپس رای گیری می کند (براساس برچسب آن ها)

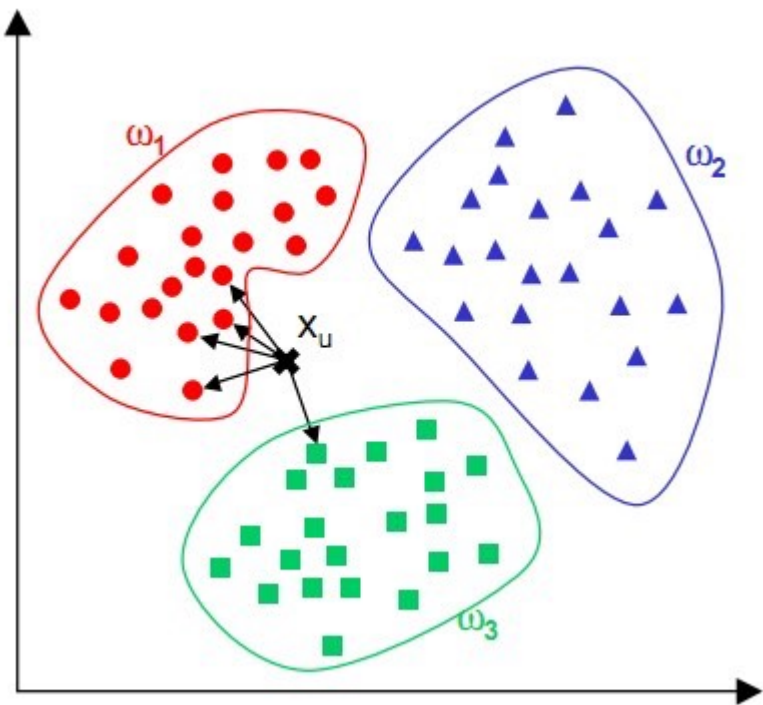
3. در نهایت برچسب خروجی Test مشخص می شود.

معیار فاصله ی اقلیدسی است:

$$dis(x, \hat{x}) = \sqrt{\sum_{i=1}^d (x_i - \hat{x}_i)^2}$$

# KNN

- چرا به جای  $k$  از یک نزدیکترین همسایه استفاده نمیکنیم؟
- اگر تعداد رای ها مساوی شود چی می کنند؟



# ادامه روش های دسته بندی

## 3) روش $\epsilon$ : ENN همسایگی

یک عدد از قبل مشخص شده است. برای داده های تست، در شعاع  $\epsilon$  تعداد برچسب ها به رای گیری گذاشته می شوند و برچسب داده داده های تست مشخص می شود.

**سوال:** اگر در شعاع مشخص شده در داده های train هیچ داده ای موجود نباشد یا تعداد

رای ها مساوی باشد چه اتفاقی می افتد؟

**پاسخ:** شعاع را کم کم بزرگ می کنیم

# ادامه روش های دسته بندی

## 4 روش Naive Bayes (بیز):

- در این روش با استفاده از احتمالات، احتمال رخداد نمونه را در همه ی کلاس ها بدست می آوریم سپس آن کلاسی که بیشترین احتمال رخداد را برای این نمونه داشت به عنوان کلاس مورد نظر برای این نمونه ی تست در نظر گرفته می شود.
- فرض براین است که تاثیر مقدار ویژگی روی برچسب کلاس، مستقل از مقادیر دیگر ویژگی ها است (استقلال شرطی ویژگی ها Conditional Independence)

# Naive Bayes

- $d = \langle a_1, \dots, a_n \rangle$

- $$p(c = c_i | d) = p(c = c_i | A_1 = a_1, A_2 = a_2, \dots, A_n = a_n) = \frac{p(A_1 = a_1, \dots, A_n = a_n | c = c_i) * p(c = c_i)}{p(A_1 = a_1, \dots, A_n = a_n)} = \frac{p(A_1 = a_1, \dots, A_n = a_n | c = c_i) * p(c = c_i)}{\sum_{k=1}^m p(A_1 = a_1, \dots, A_n = a_n | c = c_k) * p(c = c_k)}$$

- با توجه به اینکه مخرج کسر بالا براي همه ي کلاس ها روي اين نمونه یک عدد مي شود مي توان آن را در نظر نگرفت.

# Naive Bayes

- $p(c = c_i | d) = p(c = c_i) * p(A_1 = a_1 | A_2 = a_2, \dots, A_n = a_n, c = c_i) * p(A_2 = a_2, \dots, A_n = a_n | c = c_i)$
- $p(c = c_i | d) = p(c = c_i) * p(A_1 = a_1 | A_2 = a_2, \dots, A_n = a_n, c = c_i) * p(A_2 = a_2 | A_3 = a_3, \dots, A_n = a_n, c = c_i) * p(A_3 = a_3, \dots, A_n = a_n | c = c_i) =$   
 $p(c = c_i) * p(A_1 = a_1 | c = c_i) * p(A_2 = a_2 | c = c_i) * \dots * p(A_n = a_n | c = c_i)$
- $p(c = c_i | d) = p(c = c_i) * \prod_{j=1}^n p(A_j = a_j | c = c_i)$

# Naive Bayes

زمانی که ممکن است  $p(A_j = a_j | c = c_i) = 0$  شود می توان به جای رابطه زیر :

$$p(A_j = a_j | c = c_i) = \frac{n_{ji}}{n_i}$$

از رابطه زیر استفاده کرد: ( $\lambda$  می تواند 1 یا  $\frac{1}{2}$  باشد)

$$p(A_j = a_j | c = c_i) = \frac{n_{ji} + \lambda}{n_i + \lambda * d_j}$$

$d_j$  تعداد اعضا دامنه ویژگی  $A_j$  :

$n_{ji}$  : تعداد نمونه هایی که در آن  $A_j$  برابر  $a_j$  است در کلاس  $c_i$  هستند.

$n_i$  : تعداد داده ها با برچسب  $c_i$

# Naive Bayes

ID	A	B	C
1	a	d	y
2	a	e	y
3	b	f	y
4	c	e	y
5	b	f	y
6	b	f	n
7	b	e	n
8	c	d	n
9	c	f	n
10	a	d	n

مثال: برچسب نمونه تست  $\text{data}=\langle a, f \rangle$

$$P(C=y) = 5/10 = 1/2$$

$$P(C=n) = 5/10 = 1/2$$

$$P(A=a|C=y) = 2/5$$

$$P(A=b|C=y) = 2/5$$

$$P(A=c|C=y) = 2/5$$

$$P(A=a|C=n) = 1/5$$

$$P(A=b|C=n) = 2/5$$

$$P(A=c|C=n) = 2/5$$

$$P(B=d|C=y) = 1/5$$

$$P(B=e|C=y) = 2/5$$

$$P(B=f|C=y) = 2/5$$

$$P(B=d|C=n) = 2/5$$

$$P(B=e|C=n) = 1/5$$

$$P(B=f|C=n) = 2/5$$

$$P(C=y) \times \prod_{j=1}^2 P(A_j = a_j | C=y) = \frac{1}{2} \times \frac{2}{5} \times \frac{2}{5} = \frac{4}{50} = 0.08$$

$$P(C=n) \times \prod_{j=1}^2 P(A_j = a_j | C=n) = \frac{1}{2} \times \frac{1}{5} \times \frac{2}{5} = \frac{2}{50} = 0.04$$



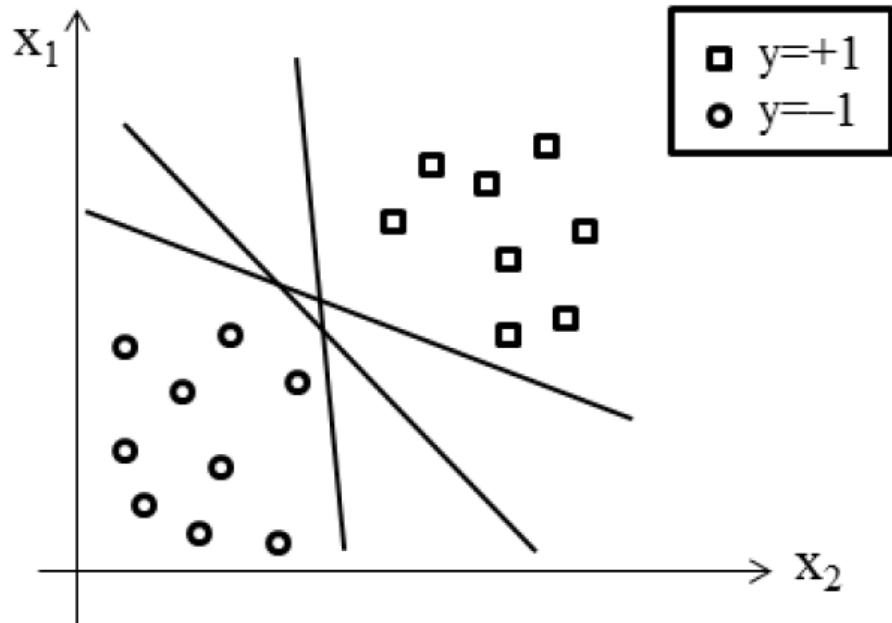
# ادامه روش های دسته بندی

## 5) روش دسته بندی ماشين بردار پشتيبان (SVM):

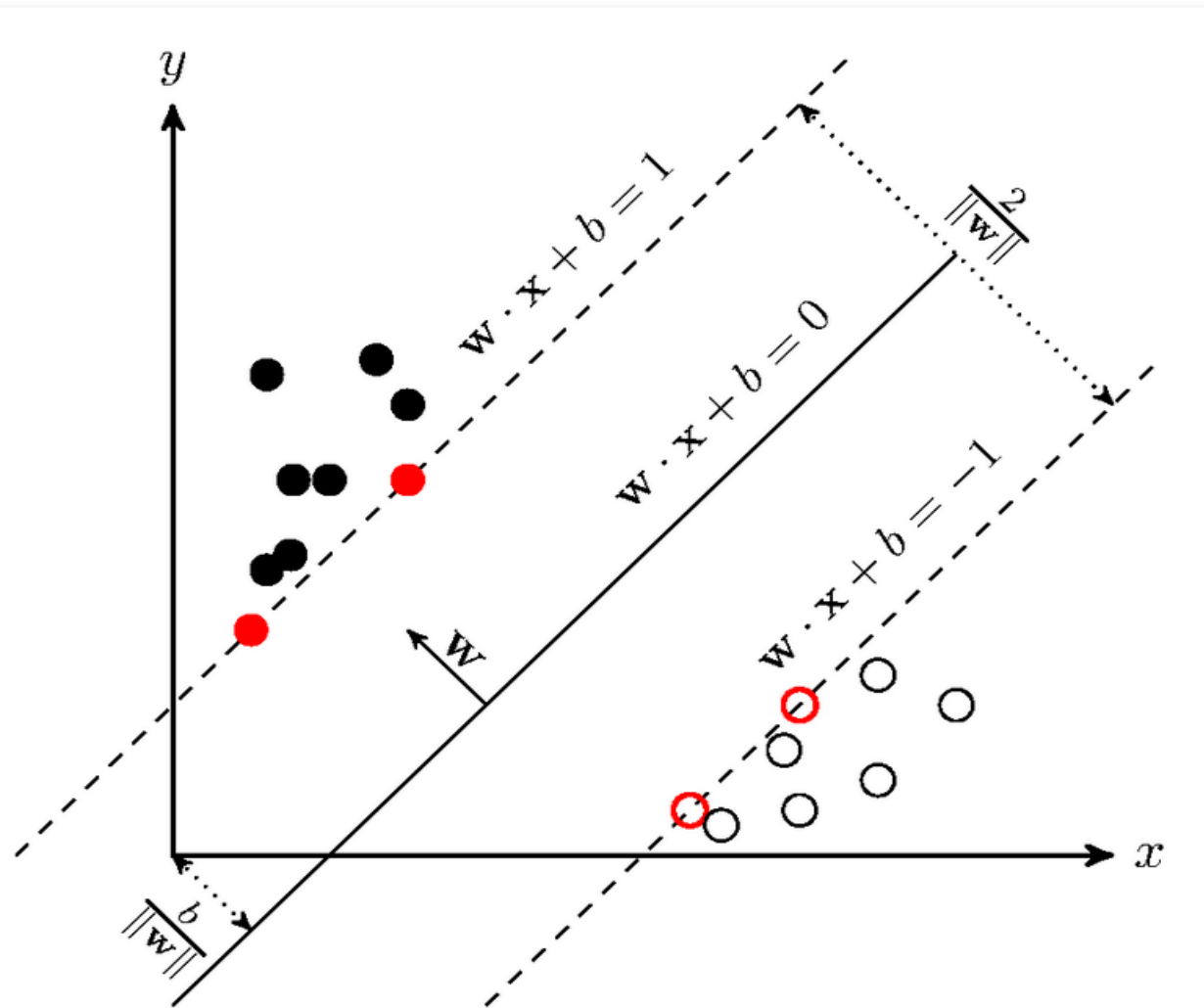
داده ها را با يك خط به دو قسمت تقسيم كن (دنيا ساده است).

اين روش براي صورت مسئله هاي دو کلاسه تعريف شده است ولي براي صورت مسئله هاي چند

کلاسه گسترش داده مي شود.



# SVM (Support Vector Machine)



# SVM

- $\|w\| = \sqrt{w_1^2 + w_2^2 + \dots + w_d^2}$
- اگر  $d$  ،  $x$  بعدي باشد معادله خط جداکننده به معادله ابرصفحه (hyper plane) تبدیل می شود.
- $WX + b = 0 \rightarrow w_1x_1 + w_2x_2 + \dots + w_dx_d + b = 0$
- خط جداساز توسط داده های آموزشی پیدا می شود سپس در مرحله ی تست:  
 $wx + b \geq 0 \rightarrow y = 1$   
 $wx + b \leq 0 \rightarrow y = -1$

# SVM

$$\text{برچسب } X = \text{sign}\left(\sum_{i=1}^s y_i \alpha_i X_i X^T + b\right)$$

$y_i$  برچسب بردار پشتیبان :

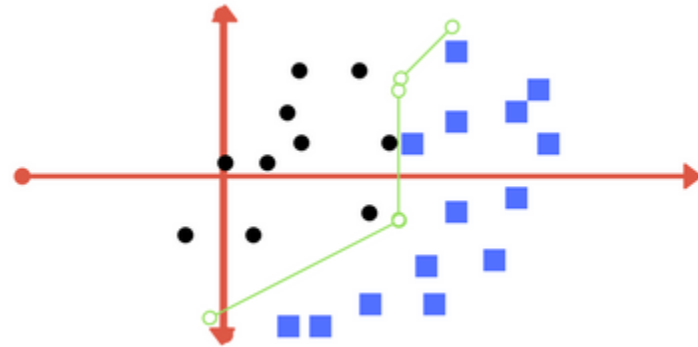
$X_i$  بردار پشتیبان :

$\alpha_i$  مقدار مناسب آن پیدا train عدد ثابتی نیست که باید در مرحله ی

شود.

# SVM

در بعضي مواقع داده ها به صورت خطي از هم جدا پذير نيستند در اين مواقع با اعمال توابع کرنل روي داده ها آن ها را از فضاي  $d$  بعدي به فضاي  $n$  بعدي مي بريم که تفکیک پذيري با خط ممکن شود.



# SVM

کرنل هایی که بیشتر مورد استفاده قرار می گیرند:

$$\text{polynomial} \rightarrow k(\vec{x}_i, \vec{x}_j) = (\vec{x}_i * \vec{x}_j + 1)^d$$

$$\text{Gaussian - Radial Basis Function (RBF)} \rightarrow k(\vec{x}_i, \vec{x}_j) = \exp\left(-\frac{\|\vec{x}_i - \vec{x}_j\|^2}{2\delta^2}\right)$$

$$\text{hyperbolic tangant} \rightarrow k(\vec{x}_i, \vec{x}_j) = \tanh(\alpha \vec{x}_i \cdot \vec{x}_j + c)$$

$$\text{linear} \rightarrow k(\vec{x}_i, \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j$$

- $Gram\ Matrix(i, j) = kernel(\vec{x}_i, \vec{x}_j)$
- در ستون زام از  $Gram\ Matrix$  بيان مي كنيم كه هر کدام از داده ها چقدر به داده ي زام شباهت دارند (مقايسه مي شوند).

# SVM

زمانی که یک دسته بند برای صورت مسئله ی دو کلاسه طراحی شده است چگونه آن را به مسئله ی چند کلاسه گسترش دهیم؟

1. One to one :  $1 و 1$  ، ... ،  $3 و 1$  ،  $2 و k$  ،  $2 و 2$  ، ... ،  $4 و 2$  ،  $3 و k$  ، ..... ،  $k و k-1$

2. One to other: 1 2,3 مقابل  $1, \dots, k$

2 در مقابل  $1, 3, \dots, k$

و...

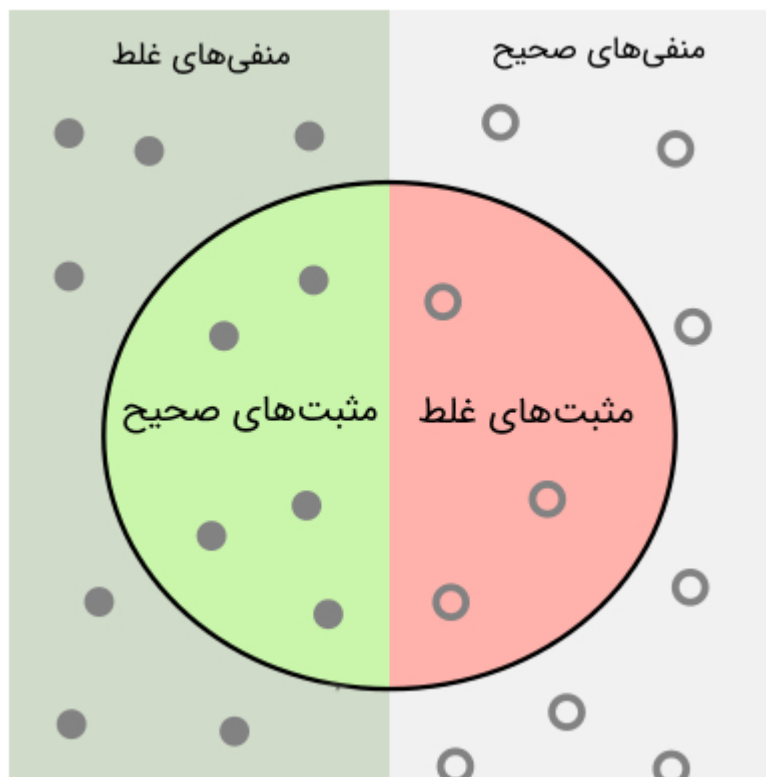


# معیار ارزیابی روش های دسته بندی

معیار دقت در روش های دسته بندی برای مرحله ی test :

در این معیار ها برچسب های پیش بینی شده برای داده ها با برچسب های واقعی

مقایسه می شوند.



# معیار ارزیابی روش های دسته بندی

- $precision = \frac{TP}{TP + FP}$

- $recall = \frac{TP}{FN + TP}$

- $F - Measure = \frac{\alpha * precision * recall}{recall + precision} \quad \alpha=2$

- $ACC = \frac{TP + TN}{TP + TN + FP + FN}$

# خوشه بندي clustering

- هدف: کنار هم قرار دادن داده های مشابه به صورتی که داده های که هم خوشه نیستند در خوشه های متفاوت قرار بگیرند.
- در واقع الگوریتم هایی که برای خوشه بندي ارائه می شوند دو هدف زیر را دنبال می کنند:
  1. فاصله ی درون کلاسی کم
  2. فاصله ی برون کلاسی زیاد
- بعضی مواقع تعداد خوشه ها (برای خوشه بندي) از قبل مشخص است در بعضی مواقع مشخص نیست

# clustering

تعداد خوشه ها مشخص باشد 2 روش وجود دارد:

## 1. k-means:

مراحل:

- 1) برچسب گذاری تصادفی داده ها (اینکه متعلق به کدام خوشه هستند)
- 2) محاسبه میانگین و قرار دادن آن به عنوان مرکز خوشه
- 3) برچسب گذاری جدید بر اساس فاصله تا نزدیکترین مرکز خوشه
- 4) تکرار 2 و 3 تا شرط پایان

## 2. k-medoids