

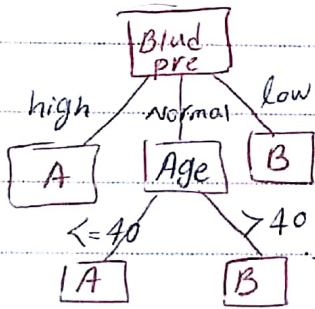
برای Age می توانیم Info حساب کنیم چون دسته های آن زیاد است  
 روش های دسته سازی را می بینیم روی این ایمان می کنیم  
 مثلاً دو دسته کوچکتر از 40 و بزرگتر از 40

روش های دسته بندی: } Classification  
 1- روش تقسیم  
 2- روش پیرین  
 3- rate-based

ID	Sex	Age	Blood-Pre	Drug	شماره
1	Male	20	Normal	A	
2	Female	73	Normal	B	
3	M	37	High	A	
4	M	33	low	B	
5	F	48	H	A	
6	M	29	N	A	
7	F	52	L	B	
8	M	42	N	B	
9	M	61	N	B	
10	F	30	L	A	
11	F	26	H	B	
12	M	54		A	
13	M	56	Normal	B	

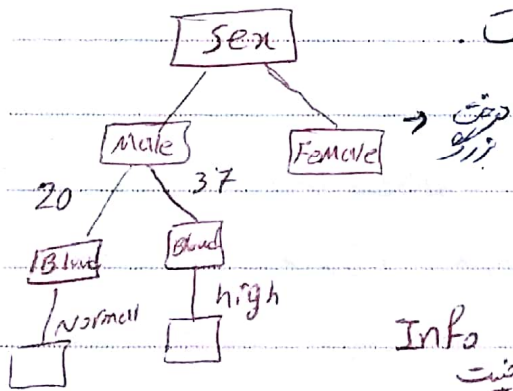
مثلاً چون دو دسته داریم بر این اساس هر کدوم را دروی A بخورد یا دروی B هر خواصم عددی در دست کنیم که بگوید بیمار دروی A بخورد یا دروی B و برای بیمارانی که عددشان کم  
 برای سن پیر و دسته بندی استفاده کنیم.

دنبال یک روش مناسب هستیم که تمام این موارد را در آن بگنجاند.  
 درخت تصمیم را می بینیم مثلاً عمق درخت کم باشد  
 ساختار زیاد باشد ← باید صورتش را داشته باشد  
 در این درخت جنبه را در نظر نمی گیریم چون تا می تونه درصدا کم ندارد



حالا یک درخت تقسیم دیگری

این که چه ویژگی در هر سطح انتخاب شود خیلی مهم است.  
 ما تا معیار را تقسیم (مقدار تقسیم)



$$Info_{(D)} = - \left[ \frac{6}{13} \log_2 \frac{6}{13} + \frac{7}{13} \log_2 \frac{7}{13} \right]$$

→ A                      → B

$$Info = \frac{8}{13} \times \left[ \frac{4}{8} \log_2 \frac{4}{8} + \frac{4}{8} \log_2 \frac{4}{8} \right]$$

نوعاد مردها                      مردانی که دروی A                      مردهایی که دروی B

$$+ \frac{5}{13} \times \left[ \frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} \right]$$

نوعاد زنان

ID	Age	Job	House	Credit	دانشجویان
1	old	F	T	Excellent	yes
2	old	F	T	Good	yes
3	Middle	F	F	Fair	No
4	Middle	T	T	Good	Yes
5	Young	F	F	Fair	No
6	old	F	F	Fair	No
7	young	F	T	Excellent	Yes
8	Young	T	F	Good	Yes
9	Middle	T	T	Fair	Yes
10		F	F	Good	Yes
					No

$$Gini(D) = 1 - \left( \left( \frac{6}{10} \right)^2 + \left( \frac{4}{10} \right)^2 \right) = 0.48$$

$$Gini(Job) = \frac{7}{10} \times \left[ 1 - \left( \frac{4}{7} \right)^2 - \left( \frac{3}{7} \right)^2 \right] + \frac{3}{10} \times \left[ 1 - \left( \frac{2}{3} \right)^2 + \left( \frac{1}{3} \right)^2 \right]$$

= 0.342

$$Gini(House) = \frac{5}{10} \times \left[ 1 - \left( \frac{4}{5} \right)^2 - \left( \frac{1}{5} \right)^2 \right] + \frac{5}{10} \times \left[ 1 - \left( \frac{0}{5} \right)^2 + \left( \frac{5}{5} \right)^2 \right]$$

= 0.12

برای Credit ۳ دسته داریم نمی توانیم از Gini استفاده کنیم - فقط  
دو دسته را می توانیم حساب کنیم پس هم منظور اما اینطور نیست پس می تواند حساب کنیم.

House ← اولین شرط و جهت

حالا برای اینکه بتوانیم از Gini استفاده کنیم باید هر کدام را به دو دسته تقسیم کنیم

مثلاً سن ← old, middle, young - old, middle, young یا

برای هر دو دسته را حساب کنیم و Gini حساب کنیم ببینیم بهترین دسته بندی کدام است برای

credit هم منظور

اول کاربِ Dataset داریم باید عملیات زیر را انجام بدهیم:  
 Pre Processing بر اساس شاخص های فریزی برای اینکه اکتاف انجام شود  
 روش اول ۴ مرحله زیر می باشد {از داده های ناقص داریم  
 ۲ داده های پرت و نویز

برای داده های ناقص می توانیم آن ها را حذف کنیم  
 وقتی داده های ما زیادند بهتر حذف کردن تا همگی در کل دیده شوند آن را حذف می کنیم  
 یک راه دیگر این است که یک میانگین در نظر بگیریم و به جای داده های ناقص بذاریم  
 اینکه کدام راه را انتخاب کنیم بستگی به مقدار داده های ما دارد  
 می توانیم میانگین یک دسته یا کلاس را قرار دهیم مثلاً از روی بقیه می بینیم که  
 می بینیم که داده های ما در این کلاس قرار دارد  
 اینکه میانگین را قرار دهیم یا چیزی دیگر یک میز تحریر است

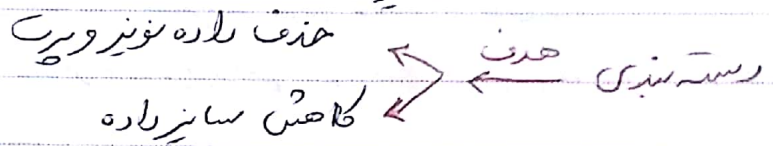
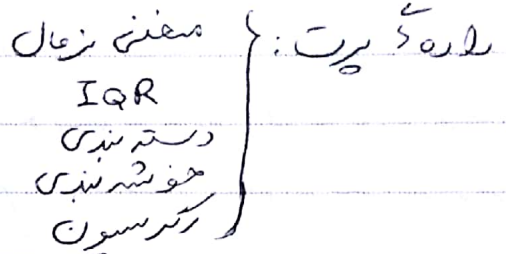
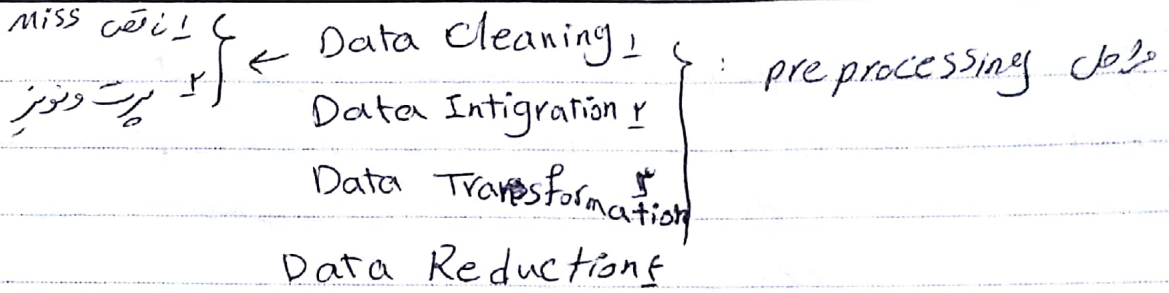
داده های پرت را حذف می کنیم که این هم تحریر است به خودمان بستگی دارد

روش های تشخیص داده های پرت:   
 ۱ توزیع نرمال   
 ۲ IQR چارت   
 ۳ روش تشخیص مبتنی بر فاصله از میانگین و انحراف معیار   
 در کلاس بندی   
 در کلاس خوب است.

۴ یک داده داریم و می خواهیم ببینیم پرت است یا نه؟  
 فاصله این دیتا با همدین نمونه را حساب می کنیم. (با P نمونه)  
 اگر فاصله آن با همدین داده ها بیشتر از یک مقدار threshold(d) باشد می بینیم  
 که پرت است. d و P به صورت تحریر تعیین می شوند

۳۵ داده

انفر نمونه   
 ۱۰   
 ۱۵   
 ۲۰   
 ۱۳   
 ۷   
 > threshold(d)

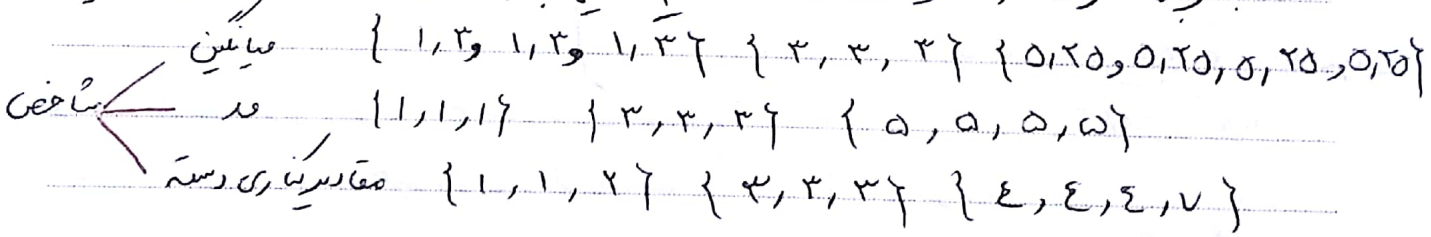


مثال: { 3, 2, 1, 5, 4, 3, 1, 7, 5, 3 }

اول داده ها را مرتب کنیم: { 1, 1, 2, 3, 3, 4, 5, 5, 7 }

بعد از آن دسته ها را انتخاب کنیم: { 1, 1, 2 } { 3, 3, 3 } { 4, 5, 5, 7 }

بعد برای هر دسته یک شاخص می گذاریم. بهترین شاخص برای هر دسته شاخص لایه گذاری است.



حجم این است که چندتا دسته درست کنیم و طول هر دسته چقدر باشد؟

$K = 1 + 3.3 \log n$

$K \leftarrow$  تعداد دسته ها

$n \leftarrow$  تعداد کل نمونه ها

به جای قدر مطلق می بینیم است

میانگین با میانگین دسته - مقادیر دسته

Min  $\rightarrow$   $E_{min} = 1$

حجم و حجم گفته شود

سوال: { 4, 8, 1, 2, 9, 2, 2, 8, 1, 5 }

دسته بندی به شکل

الف - { 1, 1, 2 } { 2, 2, 5 } { 4, 8, 8, 9 }

{ 1, 1, 1 } { 2, 2, 2 } { 8, 8, 8, 8 }

عدد { 1, 2, 8 }

$$\text{Error} = |1-1| + |1-1| + |2-1| + |2-2| + |5-2| + |4-8| + |8-8| + |8-8| + |9-8| = 7$$

دسته بندی به شکل

ب = { 1, 1, 2, 2, 2 } { 5, 4 } { 8, 8, 9 }

عدد 2 5 8

$$\text{Error} = |1-2| + |1-2| + |2-2| + |2-2| + |2-2| + |5-5| + |6-5| + |8-8| + |8-9| = 4$$

Error کمتر شد پس دسته بندی به شکل ب بهتر است.

خوشه بندی: داده ها می توانند به هم را در یک دسته بگذاریم و صفی خوشه یا table یا کلاس بنامیم.

کاربرد: مثلاً مشتری های شهر به هم را در یک دسته قرار می دهیم / مشتری جدید را در یک از این دسته ها قرار می دهیم (سبابت بین آن ها با هم است)

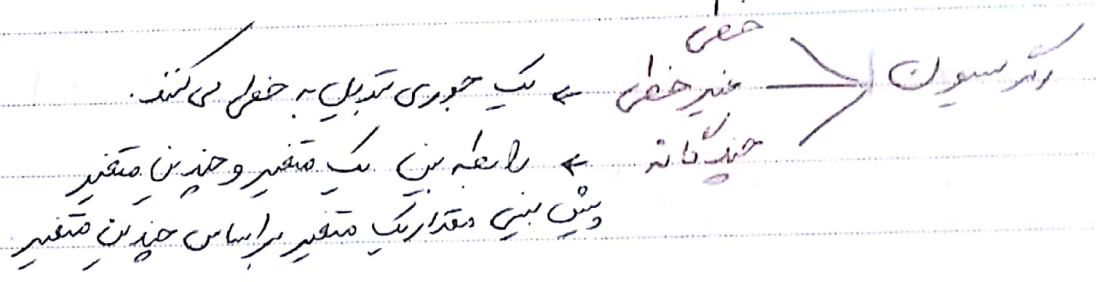
عین بینی کنیم مشتری جدید چه چیزهایی می خرد / کاربرد دیگر: پردازش تصویر، تشخیص الگو، تشخیص مشتری ها

اشیا شهر به هم را در یک دسته قرار می دهیم / مثلاً میوه های شهر به هم را در یک دسته قرار می دهیم

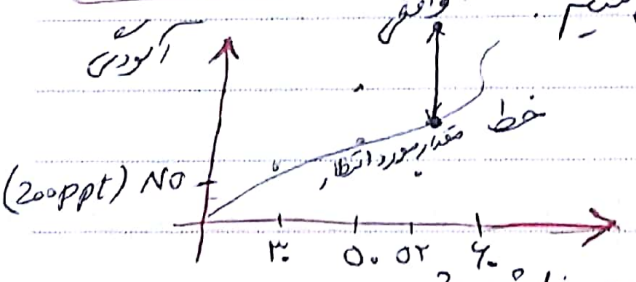
اگر داده ها مرتبط بود و به هم گام از دسته ها شهر نبود پس می گوییم داده مرتبط است / پس برای تشخیص داده های مرتبط کاربرد دارد

رگرسیون: برای اینکه مقادیر متغیر را بر اساس متغیر دیگر پیش بینی کنیم مثلاً رابطه بین دما و آلودگی

اگر رگرسیون برای دو متغیری استفاده می شود این دو متغیر با هم ارتباط دارند یا نه یعنی ارتباط مثبت یا ارتباط منفی اما ارتباط صفر نباشد (هم یعنی مستقل) پس حتی پس این دو متغیر جنسیت همبستگی دارد. (Correlation - Covariance)



دما | آلودگی



در خواص رابطه بین آلودگی و دما بررسی کنیم.  $y = \alpha x + \beta$  باید خط پیدا کنیم.

بعضی از داده ها روی خط نمی افتند

که به فاصله این داده ها از خط Error گوئیم. خطی که انتخاب می کنیم Error نداشته باشد.

رگرسیون چندگانه  $\rightarrow y = \alpha_1 x_1 + \alpha_2 x_2 + \beta$

$$\alpha = \bar{y} + \beta \bar{x}$$
$$\beta = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

## Inteigration یکپارچه سازی

چندین دیتا بیس داریم که هم خواصم آن ها را یکپارچه کنیم مثلاً ویترین قد در یک دیتا بیس بر حسب Foot و در یک دیتا بیس بر حسب سنتیمتر است. یا سازگاری یا داده های تکراری داشته باشیم. پس در یکپارچه سازی هدف حذف افزونگی و ناسازگاری است که تبدیل می شود به اینبار داده. حالا فرض کنیم یک اینبار داده داریم. مثلاً داده های مابین ۵۰ تا ۵۵ - هست اما بعضی صدها داده کاوی نیاز دارند که عدد مابین [۱۰] باشد که نه ما لانز

## Transformation تفسیر شکل داده : نه ما لانز کردن یا استاندارد سازی

روش های استاندارد سازی :

① حرکت نقطه اعشار مثلاً 3.76 - 800 = data

$$\frac{\text{اعداد}}{\text{توان } 10^3} \rightarrow \frac{\text{همه اعداد}}{10^3} \quad \text{مثلاً } 800 \approx 1000 \rightarrow 10^3 = \text{قدر مطلق بزرگترین data}$$

② موضوع اعداد در بازه [Min, MAX] قرار می دهند

$$\frac{\text{عدد} - \text{Min}}{\text{Max} - \text{Min}} (\text{NewMax} - \text{NewMin}) + \text{NewMin}$$

③ میانگین - عدد

$$\frac{\text{میانگین} - \text{عدد}}{\text{انحراف معیار}}$$

معموداً قبل از اعمال شبکه عصبی نه ما لانز سازی انجام می دهیم.

## data Reduction کاهش داده ها

حالا در خواصم یک سری داده ها حذف کنیم. یا مثلاً داده های ما ویترین های خیلی زیادی دارند مثلاً دیتا بیس از دانشجویان داریم اگر سن - شماره دانشجویی - اسم - محل و ... وقتی در خواصم مطلع می شود دانشجویان را حساب کنیم و ویترین های اکس و اسم ... تا شمریدارند پس آن ها را حذف می کنیم. حالا مثال دیگر در خواصم تا شمریدارند پس در خواصم حساب کنیم - مقهورای - مقهورای روشن - مقهورای تیره مقهورای خیلی تیره ... می توانیم روسته در نظر بگیریم که مقهورای روشن و مقهورای تیره

H4MKELASI

- ۱- کاهش تعداد درگیری ها
  - ۲- کاهش تعداد نمونه ها
  - ۳- کاهش مقادیر یک درگیری
- ۱- انتخاب تصادفی که می تواند با جایگزینی یا بدون جایگزینی باشد  
 ۲- تقسیم به شدت از ۱۰ تا ۱۰۰۰ اولیگی را انتخاب کن از ۰ تا ۱۰۰۰  
 ۳- بر اساس کلاس های مختلف انتخاب کنیم  
 شدت فر حواصم از دخترا و پسر ها به یک مقدار انتخاب کنیم  
 که جنسیت در بررسی های ما اثر می گذارد

صل تبدیل یک طبقه به ۳ دسته

کاهش مقادیر یک درگیری فرض کنید سن از ۱۳۰ - ۵ سال باشد  
 در خواصم دسته بندی کنیم دریا روشن داریم ← کای دو  $7^2$   
 گسته سازی کنیم ← روشن که گسته سازی اند ← آنتروپی

ID	ویژگی سن	class
1	1	A
2	3	B
3	7	A
4	8	A
5	9	A
6	11	B
7	23	B
8	37	A
9	39	B
10	45	A
11	46	A
12	50	A

این دیتا بیست حالت خروجی کلاس می تواند  
 دخترا یا پسر باشد  
 خنود جان شروع به دسته بندی می کنیم  
 تعداد این دسته ها زیاد است  
 بعضی از دسته ها را می توان ادغام کرد  
 در خواصم بینم چه جور می توان این دسته ها را  
 با هم ادغام کرد

$[0, 2)$   $[2, 5)$   $[5, 7.5)$   $[7.5, 8.5)$   $[8.5, 10)$

$[10, 17)$   $[17, 30)$ , ...

اسم کلاس ها را می نویسیم

	A	B	
$[7.5, 8.5)$	$A_{11} = 1$	$A_{12} = 0$	$\rightarrow R_1 = 1$
$[8.5, 10)$	$A_{21} = 1$	$A_{22} = 0$	$\rightarrow R_2 = 1$
	$C_1 = 2$	$C_2 = 0$	$N = 2$

اسم دسته ها را می نویسیم  
 نگاه می کنیم  
 از این کلاس در این  
 دسته جدا می شوند  
 بعد این ها را سطری و ستونی با هم جمع می کنیم



### فصول گامی دو:

$$\chi^2 = \sum_{i=1}^2 \sum_{j=1}^k \frac{(A_{ij} - E_{ij})^2}{E_{ij}}$$

$\downarrow$  دریا دسته داریم  
 $\downarrow$  تعداد نمونه؟ دریا اصل دسته اول - امین کلاس  
 $\downarrow$  فراوانی مورد انتظار

هر خواصیم بسنیم دریا دسته را  
بصورت ارقام کنیم یا نه؟

N: تعداد کل نمونه ها

$$E_{ij} = R_i \times C_j / N$$

دریا دسته اول → دریا دسته اول → دریا دسته اول  
دریا دسته اول = تعداد کلاس = دریا دسته اول

یک جدول گامی دو داریم که بر اساس دریا دسته اول مقدار گامی دو دریا دسته اول نوشته شده

$$\chi^2 = (1-1)^2/1 + (0-0.1)^2/0.1 + (1-1)^2/1 + (0-0.1)^2/0.1$$

$$= 0.2 < 2.706 \rightarrow \text{مقدار } \chi^2 \text{ از اساس جدول}$$

چون مقدار  $\chi^2$  از مقدار  $\chi^2$  که باید باشد کمتر است پس این دو دسته را ارقام می کنیم. اینقدر این کار را تکرار می کنیم تا به این نتیجه برسیم هیچ دو دسته ای با هم ارقام نمی شوند

(7.5, 1.0)

$$E_{11} = 1 \quad E_{12} = 0 \Rightarrow 0.1 \leftarrow \text{برای این مثال صفر باشد}$$

$$E_{21} = 1 \quad E_{22} = 0 \Rightarrow 0.1$$

کسته سازی متنی بر اساس روشی: مجموعه D را به دو مجموعه  $D_1$  و  $D_2$  می تقسیم می کنیم  
حالا می خواصیم بسنیم این تقسیم بندی خوب است یا نه؟



$$Info_A(D) = \frac{|D_1|}{|D|} \times Entropy(D_1) + \frac{|D_2|}{|D|} \times Entropy(D_2)$$



اصول کلاس K-اسم در D

مشدد: بزرگساز آتا 59 است

$$Entropy = - \sum_{i=1}^K P_i \log_2 P_i$$

فان این را به در دست تقدیم کنیم

از آتا 30 و از آتا 59

30 نقطه انفعال هر خواصیم بینیم آیا 30 نقطه خوبی است؟

حرف به حرف ستون ها (حرف بعضی از دستگیری ها) خوب است و بر کاربرد  
 کلاس حج نمونه صلح به به کاربرد اما کاربردش تر از هم حرف دستگیری است  
 اما روش دستگیری مهم نیست به کاربردش نیست (کلاس مقادیر دستگیری)

حرف بعضی از دستگیری ها به واضح است که بعضی از دستگیری ها در تقسیم خوبی یا پس می آید  
 عا اثری ندارند.

دستگیری X	دستگیری Y	نوع کلاس
3	7	A
2	9	B
6	6	A
5	5	A
8	7	B
4	9	A

هر خواصیم بر اساس دستگیری X در

دستگیری Y نوع کلاس را مشخص کنیم

هر خواصیم بینیم برای تقسیم کلاس

آیا ضروری دستگیری را نیاز داریم یا طبق

دستگیری می توانیم نوع کلاس را مشخص کنیم

$$فردول = \frac{|میانگین B - میانگین A|}{\sqrt{\frac{واریانس A}{N_1} + \frac{واریانس B}{N_2}}}$$

برای ردیف X این فردول را جدا

و برای ردیف Y این فردول را جدا

حساب می کنیم

$N_1$  و  $N_2$  تعداد هر کدام

اصول کلاس A  
 $X_A = \{3, 6, 5, 4\}$      $Mean(X_A) = 4.5$      $Var(X_A) = 1.25$

اصول کلاس B  
 $X_B = \{2, 8\}$      $Mean(X_B) = 5$      $Var(X_B) = 9$

$$فردول = \frac{4.5 - 5}{\sqrt{\frac{1.25}{1} + \frac{9}{2}}} = 0.2279 < 0.10$$

فاصله میهن :  $d(o_1, o_2) = |4-7| + |5-8| + |4-3| + |8-1| = 4$

$d(o_2, o_3) = 12$

$d(o_1, o_3) = 8$

$$\sqrt[h]{\sum_{i=1}^k (x_{ik} - x_{jk})^h}$$

عیون کوشش : برابر  $h$  مختلف

اما  $h$  را بیشتر از 2 نماند از آن معمولاً

Mutual Neighbor Distance (دارد 2 عدد)

فرض کنید چند نمونه داریم که دارای ویژگی  $x$  و  $y$  است

	X	Y
A	4	2
B	6	4
C	8	8
D	3	2
E	3	1
F	2	2

فرض کنید می خواهیم عدم تمایز  $A$  و  $C$  را بدست آوریم.

$MND = NN(A, C) + NN(C, A) = 2 + 2 = 4$

$NN(A, C) = 2$  اول شباهت  $A$  با  $C$  را

$d(A, C) = 7.21$  هر چه عدد بعدی شماره  $A$  در رتبه

$d(B, C) = 4.47$  شباهت  $C$  با  $B$  شباهت دارد. اولین

$d(D, C) = 7.81$  شباهت  $A$  را  $B$  دارد و بعد  $A$ .

$d(E, C) = 8.60$

$d(F, C) = 8.48$

$NN(C, A) = 5$

حالا شباهت  $A$  با  $C$  را بدست می آوریم.

$d(A, B) = 2.83$

اول عدد  $A$  می رویم یعنی  $A$  را

$d(A, C) = 7.21$

حالا  $A$  و  $C$  رتبه 5 ام را دارد

$d(A, D) = 1.00$

$d(A, E) = 1.41$

$d(A, F) = 2.00$

بسنی - (روغالی)

تساہ برار لاسہ باسنی:

	قد	وزن	Lips	Hair	Hand	Gender
O <sub>1</sub>	155	65	thin	Curly	R	F
O <sub>2</sub>	172	75	thick	straight	R	M
O <sub>3</sub>	162	68	thin	curly	R	F
O <sub>4</sub>	168	62	thick	curly	R	F
O <sub>5</sub>	175	80	thick	curly	L	M

روغالی، ستون کتا باسنی تہند آن در باسنی ہر کتہم

$\Rightarrow 1 > 70$  (توزن)  
 $\Rightarrow 0 < 70$   
 165  $\Rightarrow$  1  
 165  $\Rightarrow$  0  
 حالات خواصم تساہ نمونہ کتا سو کتا سہرا تہم

	O <sub>3</sub>		
O <sub>4</sub>	1	a=2	b=2
		c=1	d=1
	3	3	6

$$P = \frac{a+d}{a+b+c+d}$$

تاکتہ مصافحہ ہر تعداد =  $\frac{2(a+d)}{b+c+2(a+d)}$  حالت در کتہاں

تاکتہ ہر حالت آ غنہم سان =  $\frac{a+d}{2(b+c)+a+d}$

Similarity و Dissimilarity

$$Sim = \frac{1}{1+Dis}$$

$$D = \sqrt{2(1-Sim)}$$

$$S = 1 - Dis$$

عداز سہرا کتہاں تساہ سو کتا سہرا تہم

این مقدار کمترین است که آن را باید مقدار ترنسولو مقایسه کنیم که این مقدار است نه  
 را خودمان مشخص می کنیم که اینجا تقسیم ۱۵ باشد  
 حالا که این مقدار کمتر از ۱۵ باشد آن در این حذف می کنیم. یعنی مقدار  $\alpha$  در تعیین  
 نوع کلاس نقش ندارد و با در نظر گرفتن  $\alpha$  نوع کلاس تعیین می شود.

$$Y_A = \{7, 6, 5, 9\} \Rightarrow \text{Mean}(Y_A) = 6.75 \quad \text{var}(Y_B) = 2.20$$

$$Y_B = \{9, 7\} \Rightarrow \text{Mean}(Y_B) = 8 \quad \text{var}(Y_B) = 1$$

$$\frac{6.75 - 8}{\sqrt{\frac{2.20}{4} + \frac{1}{2}}} = 1.2198 > 0.5$$

$$\sqrt{\frac{2.20}{4} + \frac{1}{2}} = 1.0227$$

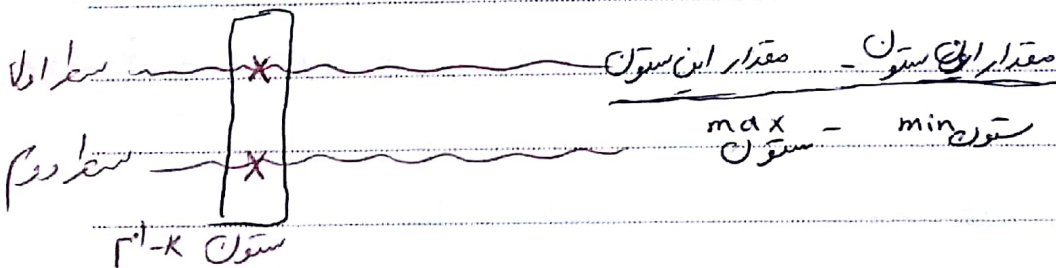
حذف بعضی از ستون ها ← بررسی تمام ویژگی ها بصورت مجزا  
 روش آنتروپی ← روش آنتروپی

روش آنتروپی: می خواهیم بین دو یا سطر از دیتا ست مقایسه کنیم تفاوت دارند.

$$S_{ij} = \alpha \sum_k D_{ik} \quad \alpha = 0.5$$

$\alpha$  بصورت تجربه تعیین می شود اما ۰.۵ بهترین حالت است.  
 برابر داده های

$$D_{ij} = \sqrt{\sum_{k=1}^n (x_{ik} - x_{jk})^2 / (\max_k - \min_k)^2}$$



$$S_{ij} = \left( \sum_{k=1}^n (|x_{ik} - x_{jk}|) / n \right)$$

برابر داده های هستند

$$Entropy = \sum_{i=1}^{N-1} \sum_{j=i+1}^N [s_{ij} \times \log s_{ij} + (1 - s_{ij}) \times \log (1 - s_{ij})]$$

با استفاده از روش AI حوش مصنوعی این تابع آنتروپی را تابع هدف قرار میدهیم هر چه آنتروپی کمتر باشد بهتر است.

**Clustering:** میخواهیم داده‌ها را در کلاس‌ها دسته‌بندی کنیم. هدف ما این است که هر نمونه را به یکی از کلاس‌ها اختصاص دهیم. ما به هر نمونه یک label می‌دهیم. هر چه label کمتر باشد، بهتر است.

0	d(2,1)	
1	d(3,1)	d(3,2)
...	...	...
n		

فرض کنید نمونه‌ها 1 تا n نام دارند. ما به هر نمونه یک label می‌دهیم. هر چه label کمتر باشد، بهتر است. ما به هر نمونه یک label می‌دهیم. هر چه label کمتر باشد، بهتر است.

دسته	Att1	Att2	Att3	Att4
O1	4	5	4	8
O2	7	8	3	1
O3	3	4	5	3

مقدار ویژگی‌ها را از نوع عددی روش اقلیدسی می‌توانیم در فضاها محاسبه کنیم.

$$\sum_{i=1}^k \sqrt{(x_{ik} - x_{jk})^2}$$

$$d(O_1, O_2) = \sqrt{(4-7)^2 + (5-8)^2 + (4-3)^2 + (8-1)^2} = 8.25$$

$$d(O_2, O_3) = 6.32$$

هر چه این عدد کمتر باشد، یعنی فاصله کمتر است. یعنی دسته‌ها به هم نزدیکتر است.

$$d(O_1, O_3) = 5.29$$

موضوعی نمونه‌ها را در K دسته‌بندی می‌کنیم و n را رکورد داریم. از این n دسته‌ها را به صورت تصادفی انتخاب می‌کنیم و مرکزهای دسته‌ها را می‌دانیم.

O<sub>2</sub> را می‌خواهیم در یک دسته‌بندی کنیم. شباهت O<sub>2</sub> را با مرکزهای حساب می‌کنیم و به مرکزهای شبیه‌تر بود در همان دسته قرار می‌دهیم. تا آخر نمونه‌ها همه را می‌بندیم و دسته‌بندی می‌کنیم. وقتی همه نمونه‌ها دسته‌بندی شدند، دفعه‌های مشابهت هر دسته را می‌دانیم و می‌توانیم آن دسته‌ها را به روش K-Means هم نامیم.

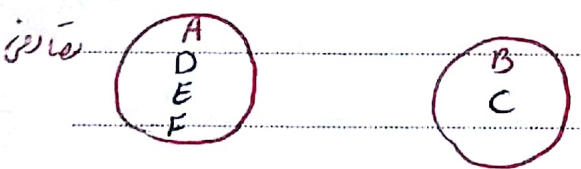
موضوعی دسته‌ها را به هم می‌بندیم و در آن K کردن. و دسته‌ها را عوض می‌کنیم دوباره همان کار را از اول انجام می‌دهیم تا هیچ زمانی این کار اراقد می‌دهیم؟ نمونه‌ها در دسته‌ها جا نمی‌شوند.

	X	Y
A	4	2
B	6	4
C	8	8
D	3	2
E	3	1
F	2	2

سوال: k=2

اولی‌ترین شباهت را می‌کنیم:

	A	B	C	D	E	F
A	0					
B	2.80	0				
C	7	4.47	0			
D	1	3.60	7.80	0		
E	1.40	4.24	8.66	100	0	
F	2	4.47	8.48	1	1.40	0



$$\text{mean}_{1x} = \frac{4+3+3+2}{4} = 3$$

$$\text{mean}_{1y} = \frac{2+2+1+2}{4} = 1.75$$

$$\text{mean}_{2x} = \frac{6+8}{2} = 7$$

می‌بینیم دسته اول را می‌بینیم روی x

می‌بینیم دسته اول را می‌بینیم روی y

$$\text{mean}_{2y} = \frac{4+8}{2} = 6$$

x	y
(3, 1.75)	(7, 6)
A	B
D	C
E	

دسترده جدید

فاصله اقلیم:

$$d(A, \text{Mean1}) = 1.03$$

$$d(A, \text{mean2}) = 5$$

$$d(B, \text{Mean1}) = 3.75$$

$$d(B, \text{Mean2}) = 2.23$$

بر این افراز ک می بیند یعنی به ک تا دسته تقسیم کردیم چون میانگین تحت تأثیر داده های بیرون است از میانگین و در استفاده می کنیم

**K-medoids** اول K تا دسته درست کن K تا نمونه را به صورت تصادفی در این دسته قرار بدهد و n-k تا دسته بقیه این K دسته تقسیم می کند تفاوت K-Means در انتخاب نمونه برای انتخاب نمونه جدید یک نمونه از یکی از دسته ها به صورت تصادفی انتخاب می کند و به عنوان نمونه جدید هر یک از دسته دیگر قرار می دهد

برای هر دسته یک نمونه جدید و نمونه قدیم داریم. برای اینکه بینم نمونه جدید بهتر است یا نه یک مدل داریم به نام مدل خطی

$$\text{مدل خطی} = \sum_{i=1}^K \sum_{P \in C_i} (P - O_i)$$

نمونه هر دسته



روانگه: برای تعیین نمونه جدید  $O_{new}$  به جای نمونه کنونی  $O_{old}$  (مردسته کنونی است) عملیات زیر باید به هر نمونه غیر نمونه مثل P انجام شود اگر P متعلق به نمونه ای باشد که  $O_{old}$  نمونه آن خوشه باشد، چنانچه با جایگزینی  $O_{new}$  به جای  $O_{old}$  فاصله P به نمونه جدید نزدیکتر شود پس P را در خوشه به نام  $O_{new}$  قرار می دهیم در غیر این صورت P را در خوشه  $i$  قرار می دهیم به نحوی که فاصله  $O_{new}$  به خوشه  $i$  با P حداقل باشد و  $old \neq i$  است



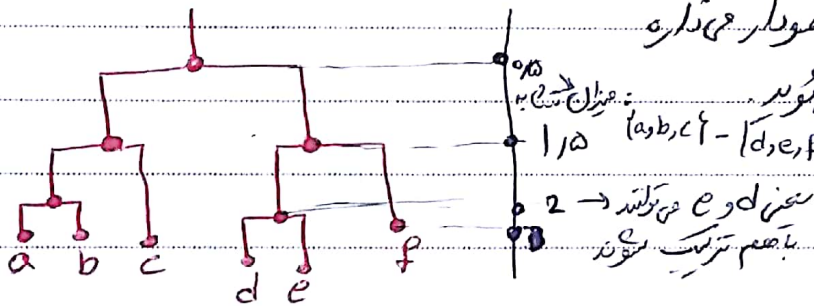
گروه P متعلق به خوشه ای باشد که  $O_i$  مانده کند باید چنانچه با جانشینی  $O_{new}$  جای  
 $O_{old}$  همپایان فاصله P با مانده خوشه خود یعنی  $O_i$  حداقل است هیچ تغییری  
 صورت نمیگیرد در غیر این صورت P در خوشه ای قرار میگیرد که مانده آن  $O_{new}$  است

$n$  = تعداد نمونه  $k$  = تعداد دسته ها  $t$  = تعداد تکرار الگوریتم

$$O(nkt) \approx O(n)$$

سلسله مراتبی: در رویه داریم از بالا به پایین یک دسته ضمنی برقرار داریم و بعد از آن  
 آن دسته جدید در هر گروهیم

یا از پایین به بالا اول تعداد زیادی دسته داریم بعد از آن که می توان ادغام کرد را ادغام می کنیم



یک محور عمودی کنار این نمودار می تازد  
 که هر عدد نشان دهنده دسته را می نویسد  
 $\{a,b,c\} - \{d,e,f\}$   
 یعنی d و e می تازند  $\rightarrow 2$   
 با هم ترکیب شوند

با دو مفهوم تقسیم کردن و ادغام کردن رویه بروسیم

معمولاً از روشی پایین به بالا استفاده می کنند علاوه بر این به سبب دو طرفه  
 تنها به سبب دو طرفه را بررسی می کنیم

	x	y
$O_1$	1	2
$O_2$	5	2
$O_3$	6	1
$O_4$	1	1
$O_5$	4	1

مثال:

فهرست شماره:	$O_1$	$O_2$	$O_3$	$O_4$	$O_5$
$O_1$	0				
$O_2$	4	0			
$O_3$	5.1	1.4	0		
$O_4$	1.0	4.1	5.0	0	
$O_5$	3.1	1.4	2.0	3.0	0

$\{O_1\}$   $\{O_2\}$   $\{O_3\}$   $\{O_4\}$   $\{O_5\}$

انتخاب دو نمونه با بیشترین شباهت  
 $O_4$  و  $O_5$  زیرا از همه با هم شبیه اند

Subject:

Year. Month. Date. ( )

	$\{0_1, 0_4\}$	$0_2$	$0_3$	$0_5$
$\{0_1, 0_4\}$	0			
$0_2$	4	0		
$0_3$	5	1.4	0	
$0_5$	3	1.4	2	0

نوبت به نوبت تمام نودها را به دست می آوریم

بهترین نتایج را می گیریم؟

چون می توانیم نتایج  $0_2$  و  $0_3$  و  $0_5$  را با هم بگیریم

هر بار نودها را با هم می گیریم

$$d(C_i, C_j) = \text{Min} \{ d(0_1, 0_2) \}$$

$$= d(C_k, C_i \cup C_j) = \text{Min} \{ d(C_k, C_i) \text{ و } d(C_k, C_j) \}$$

$$d(0_2, \{0_1, 0_4\}) = \text{Min} \{ d(0_1, 0_2), d(0_4, 0_2) \} = \text{Min} \{ 4, 4 \} = 4$$

	$\{0_1, 0_4\}$	$\{0_2, 0_3\}$	$\{0_5\}$
$\{0_1, 0_4\}$	0		
$\{0_2, 0_3\}$	4	0	
$\{0_5\}$	3	1.4	0

یعنی  $0_2$  و  $0_3$  و  $0_5$  را با هم می گیریم

بهترین نتایج

	$\{0_1, 0_4\}$	$\{0_2, 0_3, 0_5\}$
$\{0_1, 0_4\}$	0	
$\{0_2, 0_3, 0_5\}$	3	0

