

Subject:

Year. Month. Date. ()

اول کار یک Dataset داریم باید عملیات زیر را انجام بدهیم:

Pre Processing بر اساس مشخصات فریزر پراکنده که خاص انجام می شود

روش اول ۳ 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 هم در کل دنیا جا ندارد آن را حذف می کنیم

یک راه دیگر این است که یک میانگین در نظر بگیریم و به جایی داده های ناقص نداریم

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 مقیاس دیگر که

می بینیم که داده های ما در این کلاس قرار دارد

اینکه میانگین را قرار دهیم یا چیزی دیگر یک چیز تجربه است

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

روش ۱ تشخیص داده های پرت: ۱ توزیع نرمال

۲ IQR چارک

۳ روش تشخیص مبتنی بر فاصله برای جایی که کلاستر بندی در نظر می گیریم خوب است

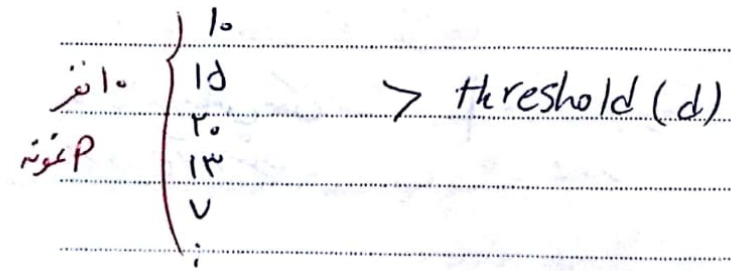
۳ یک داده داریم و می خواهیم ببینیم پرت است یا نه؟

فاصله این دنیا با همسایه نمونه را حساب می کنیم (با P نمونه)

اگر فاصله آن با همسایه داده ها بیشتر از یک مقدار $threshold(d)$ باشد می گوئیم

یک داده پرت است. d و P به صورت تجربه تعیین می شوند

۳۵ داده



Subject:

Year. Month. Date. ()

مرحله آماده سازی داده :

Data Cleaning 1 ← { 1 ناقص miss
 2 پرست و تویز }
 Data Intigration 2
 Data Transformation 3
 Data Reduction 4

داده های پرست :

IQR
 دسته بندی
 خوشه بندی
 اکتگرسیون

حذف داده نویز و پرست

کاهش سایز داده

حذف دسته بندی

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

اول به درجه ها 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 از آن هر دسته یک شاخص می گذاریم یعنی به جای هر دسته شاخص می گذاریم

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

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

$K \leftarrow$ تعداد دسته ها

$n \leftarrow$ تعداد کل نمونه ها

$$K = 1 + 3.3 \log n$$

به جای قدر مطلق ممکن است
 مجبوریم جمع گفته شود
 میانگین یا میانگین دسته - مقادیر دسته
 Min \rightarrow $E_{W, V}$

Subject:

Year. Month. Date. ()

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

مسئله:

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 بین دسته بندی به شکل ب بهتر است.

خوشه بندی: داده ها مساب به هم را در یک دسته بگذاریم و صغیر خروجی یا label یا کلاس بنامیم.

کاربرد: مثلا مشتری های شبیه به هم را در یک دسته قرار میدهیم (شباهت بین آن ها باید داریم)

در این بین گفتیم مشتری جدید چه چیزهایی می خرد.

کاربرد دیگر: پرزاش تصویر، تشخیص الگو، تشخیص مشتری ها

اسیاست شبیه به هم را در یک دسته قرار میدهیم.

مثلا میوه های شبیه به سیب را در یک دسته قرار میدهیم

اگر داده جاریت بود و به جمع کلام از دسته ها شبیه نبودیم در هر کلاس داده برت است

بی برای تشخیص داده های جاریت کاربرد دارد

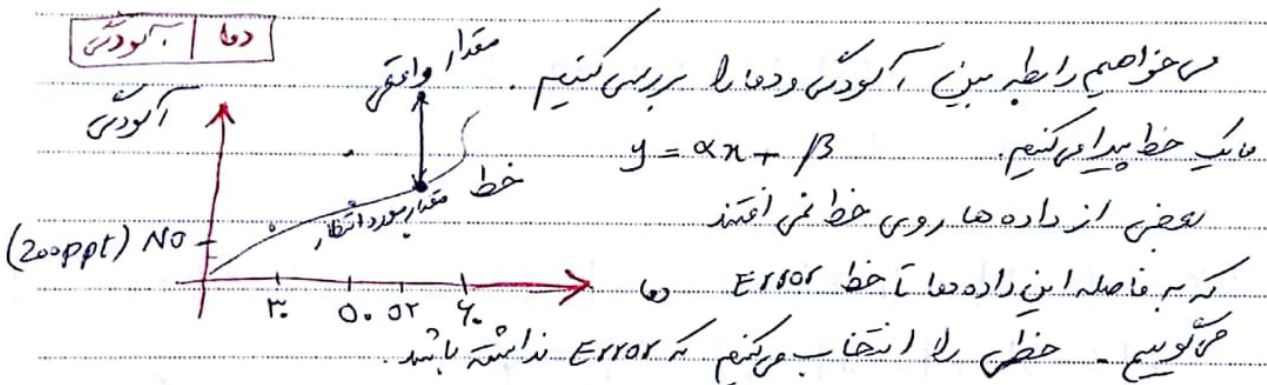
Subject:

Year. Month. Date. ()

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

رگرسیون برای دو متغیری استفاده می شود که این دو متغیر با هم ارتباط داشته باشند یعنی ارتباط مثبت یا ارتباط منفی اما ارتباط صفر نباشد (همبستگی مستقل) بین آنها $(Correlation - Covariance)$ این دو متغیر صریح همبستگی دارند.

رگرسیون \rightarrow خط خطی \leftarrow یک جدولی تبدیل به خطی می کنند
مهندسی \leftarrow رابطه بین یک متغیر و چندین متغیر
وین بینی مقدار یک متغیر بر اساس چندین متغیر



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

$$\alpha = \bar{y} + \beta \bar{x}$$

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

Subject:

Year. Month. Date. ()

Inteageration

چندین دیتا سٹس دیکھ کر ہم جواہر آں حالہ دیکھا کہ کئی صفحہ ویئر کے دریا دیتا سٹس پر حسب Foot و دیگر دیتا سٹس پر حسب صورت آہٹ ہے نا سازگاری یا دادہ ہاں شکراری دانستہ باقی ہیں دریا کہی سازگی حذف افزوی و نا سازگاری است سے تبدیل ہوتی ہے ایسا روادہ۔

حالا فرض کنیم یک ایسا روادہ داریم۔ مثلاً دادہ ہاں مابین ۵۰-۵۵-۵۵-۵۵-۵۵ اما بعض صفحہ ہاں دادہ گاوی نیاز دارند کہ علا مابین [۱۰] باقی ہے نفع لانی

Transformation

رویش صفحہ استاندر سازی:

① حرکت نقطہ اعشار مثلاً $data = 800 \pm 3.76$

اعداد توان \rightarrow عدد اعداد $\rightarrow 10^3$
 10^3
 قدر مطلق بزرگترین data

②
$$\frac{NewMax - NewMin}{Max - Min} (عدد - Min) + NewMin$$

من خواص اعداد در بازه [Min, MAX] قرار گیرند.
 NEW MIN NEW MAX

③
$$\frac{میانگین - عدد}{انحراف معیار}$$

معمولاً قبل از اعمال سلبی نفع سازگی انجام میدهم.

data Reduction

دادہ کی ماخیز زیادند۔ یک دیتا سٹس خفی بزرگ

حالا من خواص یک سری دادہ حالہ حذف کنیم۔
 یا مثلاً دادہ ہاں ما ویئر ہاں خفی زیادند مثلاً دیتا سٹس از دانشجویان دیکھ کر
 آکرں۔ شماره دانشجوی۔ اسم محلہ و۔ وقتی من خواص سطح علم دانشجویان را حساب کنیم
 ویئر ہاں آکرں واسم۔ تاثر پذیرند پس آں حالہ حذف ہم کنیم۔
 حالا مثال دیگر من خواص تاثر پذیر را بررس کنیم۔ رتبہ علمی۔ مقوہ ای۔ مقوہ ای تیرہ
 مقوہ ای خفی تیرہ۔ من توانیم رودستہ در نظر بگیریم۔ مقوہ ای روشن و مقوہ ای تیرہ

Subject:

Year. Month. Date. ()

- ۱- کاهش تعداد ویژگی‌ها
 ۲- کاهش تعداد نمونه‌ها
 ۳- کاهش مقادیر یک ویژگی
- ۱- انتخاب تصادفی که می‌تواند با جایگزینی یا بدون جایگزینی باشد
 ۲- نظم و شش از ۱۰ تا ۱۰۰ اول یکی را انتخاب کن از ۱ تا ۱۰۰ تکمیل
 ۳- بر اساس کلاس‌های مختلف انتخاب کنیم.
 مثلاً می‌خواهیم از دخترها و پسرها به یک تعداد انتخاب کنیم
 که نسبت در پسران و دختران معادل باشد.

مثال تبدیل یک سیستم به ۳ دسته
 کاهش مقادیر یک ویژگی: فرض کنید سن از ۱۳-۵ سال باشد
 می‌خواهیم ۳ دسته بندی کنیم. در آن روش داریم 7^2 ← کلاس دو 7^2
 کسبه سازی کنیم ← روش ۳ کسبه سازی اند ← آنتروپی

ID	ویژگی سن	class نوع سبک
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
R12	50	A

این دسته‌ها است نوع کلاس می‌تواند
 دختر یا پسر باشد
 خودتان شروع به دسته بندی می‌کنیم.
 تعداد این دسته‌ها زیاد است.
 بعضی از دسته‌ها را می‌توان ادغام کرد.
 می‌خواهیم ببینیم چه جور می‌توان این دسته‌ها را
 با هم ادغام کرد.

- $[0, 2)$ $[2, 5)$ $[5, 7.5)$ $[7.5, 8.5)$ $[8.5, 10)$
 $[10, 17)$ $[17, 30)$, ...

	A	B	اسم کلاس‌ها در نویسی
$[7.5, 8.5)$	$A_{11} = 1$	$A_{12} = 0$	اسم دسته‌ها در نویسی $\rightarrow R_1 = 1$ جمع
$[8.5, 10)$	$A_{21} = 1$	$A_{22} = 0$	نقشه من کنیم $\rightarrow R_2 = 1$ جمع از این کلاس در این دسته جدا هستند
	$C_1 = 2$ جمع	$C_2 = 0$	$N = 2$ بعد این معادله‌ها را می‌توانیم با هم جمع می‌کنیم.

Subject:

Year. Month. Date. ()

فردی کاری رویه

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

\downarrow درجه آزادی داریم
 \downarrow مقدار نمونه؟ درجه آزادی دسته اول - امین کلاس
 \downarrow فراوانی مورد انتظار
 هر خواصی بینیم دوریا دسته را با هم ارقام کنیم یا نه؟

N : مقدار کل نمونه ها

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

\rightarrow ۱ = درجه آزادی \rightarrow در این مثال \rightarrow ۱ - تعداد کلاس = درجه آزادی

یک جدول کلین دو طرفیم که بر اساس درجه آزادی مقدار کلین بود در این نوشته شده

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

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

چون مقدار χ^2 کمتر از مقدار χ^2 که باید باشد شد پس این دو دسته را ارقام نمیکنیم. اینقدر این کار باید کرد که مختلف نگاریم تا بر این نتیجه برسیم هیچ دو دسته ای با هم ارقام نمیکنند

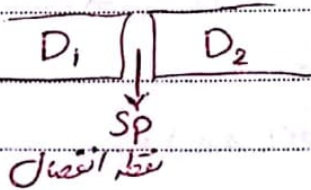
(7.5, 1.4)

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

صفر نباید باشد

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

بسته سازی صحنی بر آنتروپی: مجموع D را به دو مجموع D_1 و D_2 میکنیم
 حاصلی خواصیم بینیم این تقسیم بندی خوب است یا نه؟



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

Subject:

Year. Month. Date. ()

اصول کلاس - K ام در D

عدد "بزرگترین" از آن 59 است

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

حالتی را به درستی تقسیم کنیم

از آن 30 و از آن 59

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

حرف ϵ حذف ستون ϵ (حذف بعضی از ویژگی‌ها) خوب است و بر کاربرد
کامپیوتر حج نمونه صمیم بد به کار صبره اما کاربردش تراکم حذف ویژگی‌ها است
اما روش دسته بندی مهم نیست ϵ کاربردش نیست (کامپیوتر مقادیر ϵ ویژگی‌ها)

حذف بعضی از ویژگی‌ها ϵ و امانع است که بعضی از ویژگی‌ها در تعیین خروجی یا پس‌زمینه ϵ
علا اثری ندارند

ویژگی X	ویژگی 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 ویژگی X و
ویژگی Y نوع کلاس را مشخص کنیم
هر خواصیم بینیم برای تعیین کلاس
آیا هر دو ویژگی را نیاز داریم یا یک
ویژگی می‌توانیم نوع کلاس را مشخص کنیم

$$| \text{میانگین B} - \text{میانگین A} |$$

برای ردیف X این فرمول را جدا

$$\sqrt{\frac{\text{واریانس A}}{N_1} + \frac{\text{واریانس B}}{N_2}}$$

و برای ردیف Y این فرمول را جدا
حساب می‌کنیم

N_1 و N_2 تعداد هر کدام

خارجی هر کلاس X
حسند A $X_A = \{3, 6, 5, 4\}$ $Mean(X_A) = 4.5$ $Var(X_A) = 1.25$

خارجی هر کلاس X
کلاس B حسند $X_B = \{2, 9\}$ $Mean(X_B) = 5$ $Var(X_B) = 9$

$$\text{فرمول} = \frac{4.5 - 5}{\sqrt{\frac{1.25}{4} + \frac{9}{2}}} = 0.2279 < 0.15$$

Subject:

Year. Month. Date. ()

این مقدار که بدست آمد آن را با یک مقدار نرمال شده مقایسه می کنیم که این مقدار استاندارد را افزودن منطقی می کنیم که اینجا گفتیم ۱۵ باشد. در حالی که این مقدار کمتر از ۱۵ باشد آن درجه را حذف می کنیم. یعنی مقدار x در تعیین نوع کلاس نقش ندارد و با درجه y نوع کلاس تعیین می شود.

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

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

$$\frac{6.75 - 8}{\sqrt{\frac{2.20}{4} + \frac{1}{2}}} = 1.2198 > 0.8$$

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

حذف بعضی از ستون ها ← در صورتی که درجه ها به صورت مجزا
روشن است و درجه

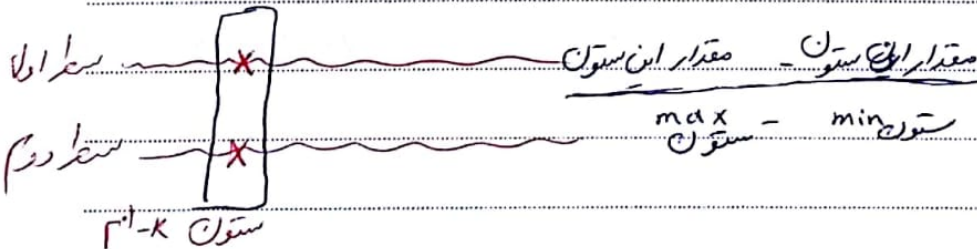
روشن است و درجه به این دو به سطر از دیانیت میزنیم تا در آنجا تفاوت دارند

$$S_{ij} = e^{\alpha D_{ij}} \quad \alpha = 0.15$$

α به صورت تجربی تعیین می شود اما ۰.۱۵ بهترین حالت است.

بلبر داده عددی

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



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

بلبر داده عددی
نیستند

Subject:

Year. Month. Date. ()

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

با استفاده از روش های AI حوسه مصنوعی این تابع آنتروپی را تابع هدف قرار می دهیم هر چه آنتروپی بیشتر باشد بهتر است

Subject:

Year: Month: Day: ()

در داده کاوی با داده و Data سروکار داریم، روی آن کار انجام می دهیم
ویک نتیجه (Result) می گیریم - page: ()

Data set مجموعه داده

چندین مدل برابر مجموعه داده داریم:

۱- رکورد Record: روی ما را با صورت رکورد ذخیره می کنیم
رکورد را می توان به ۳ نوع تقسیم کرد:

الف. حالتی: ۳ صورت فانتزی از ربات وجود دارد.

ب. اسناد متنی: به عنوان مثال یک متن داریم و یک سری کلمات که تعداد کلمات آن ها زیاد است را استخراج می کنیم و می توانیم هر کلمه چند بار در این متن ها تکرار شده است

کلمه متن	تعداد	مسابقات
سند ۱	۵	۱
سند ۲	۵	۲
سند ۳	۱	۵

ج. تراکنش Transaction

در آن یک سری ID, Item وجود دارد.

T ID	item
۱	سند ۱
۲	سند ۲
۳	سند ۳

۲- گراف و شبکه: گراف مثل وب است، می توان هر صفحه از وب را یک Node یا نو در نظر گرفت و با این آن ها را با هم متصل کرد

کاربرد گراف و شبکه را می توان در وب، شبکه های اجتماعی، ساختار مولکولی دید.

مثلاً DNA

مثلاً در ساختار مولکولی از روی Data یک سری الگوریتم می آورند
بعد می بینند که آن را فرد بعدی باید داشته باشد

Soroush

انجمن علمی علوم کامپیوتر
دانشگاه کاشان

t.me/KUCSSA

۳- ترتیبی یا ordered :

جایی که ترتیب داده ها مهم است مثل فیلم و ویدیو
بعضی مواقع زمان هم مهم است که به آن ها (نسبت به زمان) (Temporal) می گویند

۴- فضایی یا Spatial : که مجموعه نوسه ها است که به بدنه دارند

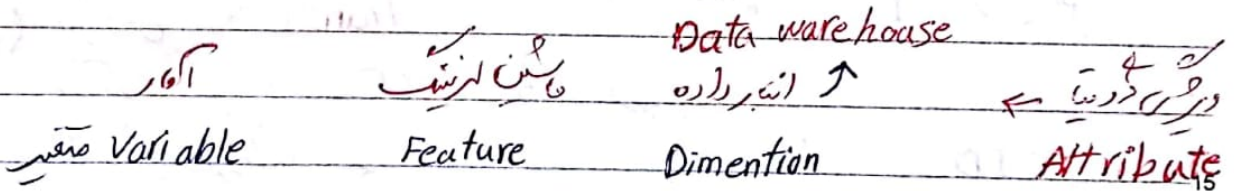
کاربرد آن در نقشه و تصاویر ماهواره ای است

Data و ویژگی ها آن :

به ویژگی Attribute می گویند

10

۱- ابعاد داده : ویژگی که داده را به یک ابعاد داده می شناسیم.
در هر کدام از موقعیت های این ویژگی را با ابعاد زیر بیان می شناسیم:



۲- پراکنگی داده یا Sparsity : که در آن بسیاری از داده ها

۳- رفع ابهامات Data یا Resolution :

داده ها مهم اند چون به مقیاس وابسته اند برای رفع ابهام باید مقیاس (scale) مشخص کنیم

۴- توزیع داده یا distribution : توضیح ندارد

۲۵- اگر یک Table داشته باشیم ، سطرها مجموع دیتا و ستون ها ویژگی ها هستند

در هر Dataset به هر رکورد یک شی داده می گویند (data object)

entity یا موجودیت هم می گویند که به آن

مثلاً در یک بیمارستان ← بیمار، دکتر، پرستار و ... موجودیت اند.
در سطح دانشگاه ← دروس، اساتد، دانشجو و ... موجودیت اند.

اینکه در یک موجودیت ← instance - tuple - object - example
Sample و

مثلاً ← می‌تواند یک Sample از یک Database

گفتیم که ستون یک Table و یک Attribute هستند که می‌توانند
مقدارهای زیر را داشته باشند:

10 - 1- Nominal اسمی: که به آن Name of thing و Category هم می‌گویند.
مثل رنگ صوف و وضعیت تاهل و کدستی، شماره دانشجویی از این نوع اند.
میانم و طبع، تفریق و میانگین برای این نوع ویژگی‌ها معنی ندارد.
اما حد برای آن‌ها قابل تعیین است.

15 - 2- Binary بوردویی: که مقدار صفر و یک مقدار یک می‌گیرد. و دو نوع اند:
1- متقارن Symmetric: اهمیت صفر و یک یکسان است
2- نامتقارن Asymmetric: برای فاکتوریت‌ها اولویت دارند و هم است که
کدام یک باشد و کدام صفر. مثلاً اهمیت یک بهتر از اهمیت صفر باشد.

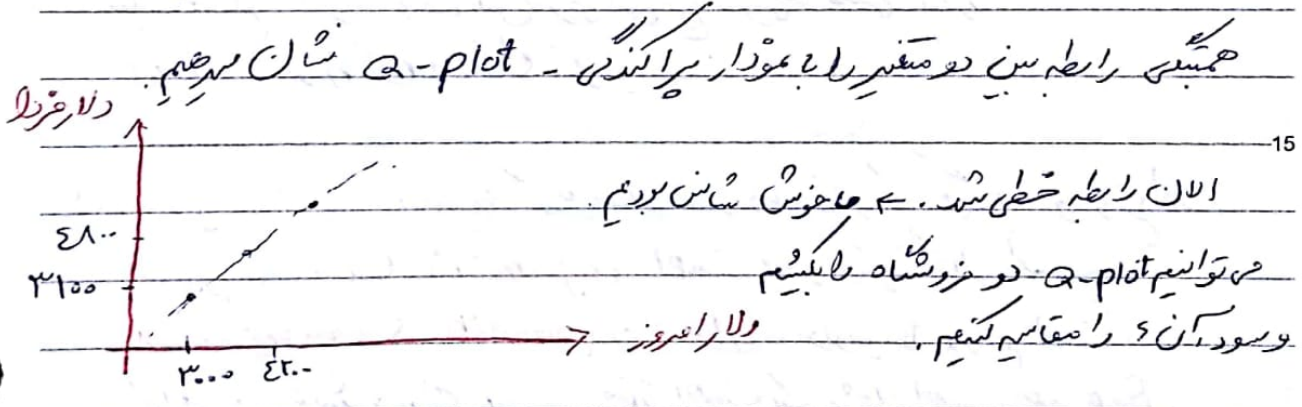
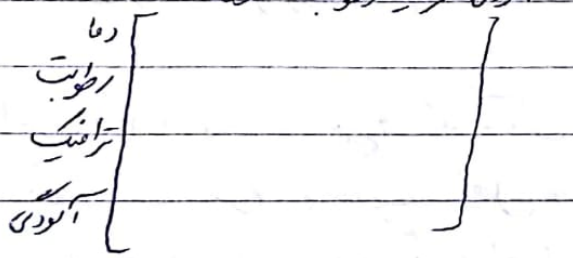
20 - 3- ترتیبی Ordinal:

25 - 4- Quantitative عددی:
1- interval فاصله: فاصله بین نمرات یکسان باشد
در اینجا مفهوم بینا صفر نداریم (صفر ذاتی نداریم)
می‌توانیم بگوییم این چند برابر آن است.

۲- Ratio (نسبت): صفر معنی ندارد (صفر دانه داریم)
 نسبت ها را می توانیم نسبت به هم بگیریم. معنی می توان مقایسه کرد.

مثال: COV (کواریانس) دو جفت عدد
 آلودگی | ترانزیستور | رطوبت هوا | دمای امروز

۵
 هر جوازم بینیم برابر یک نمونه
 رابطه بین دما و رطوبت یا دما و آلودگی را پیدا کنیم. (مثلاً بر سر یک نمونه ما)
 هر جوازم بینیم کدام یک از این جوازم را رابطه دارند؟ پس شکل خاتریس کواریانس می رسم
 آلودگی - ترانزیستور - رطوبت - دما



۲۰
 رابطه آلودگی نیز برابر تعیین همبستگی بین دو متغیر داریم. عددی که مثبت تر، آلودگی زیاد است

$$\text{Correlation} = \frac{\text{Cov}(x, y)}{\sqrt{\text{Cov}(x, x) \text{Cov}(y, y)}}$$

۰ ← همبستگی ندارند. مستقل اند.
 ۱ ← همبستگی مثبت ← (رابطه مستقیم)
 -۱ ← همبستگی منفی (برعکس)

ما با پایا داره دهنی بزرگ کار می کنیم.

بسیار روشن حسابی که می توان این اطلاعات را بدست آورد؟ **نه**

زیرا روشن حسابی که می توان برای داره های محدود است.

تحلیل رتبه ها رویش حسابی که می

مرکزیت رتبه ها اول باید بینیم این رتبه مرکزیت دارد.

ساخته که مرکزیت آگاهی
ساخته که مرکزیت آگاهی

باید بینیم رتبه را که می دارد باشد. ساخته که آگاهی

ساخته حسابی مرکزیت داره ها میانگین

Mean = $\bar{x} = \frac{\sum x_i}{n}$ اگر رتبه ها وزن داشته باشد وزن \times مقدار مجموع وزن ها

عیب میانگین در داره های پرت روی آن تأثیر می زند

یک داره می پرت میانگین را خارج جامه کند پس میانگین معیار خوبی برای مرکزیت

داره ها نیست

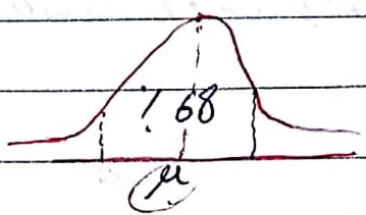
پس میانگین را حرس کنیم

Trimmed Mean: رتبه ها را مرتب می کنیم و از هر دو اول و آخر عدد از کف

داره ها حذف می کنیم به این شکل داره های پرت حذف می شوند

توزیع نرمال: بیشتر افراد جامعه اصراف میانگین قرار می گیرند.

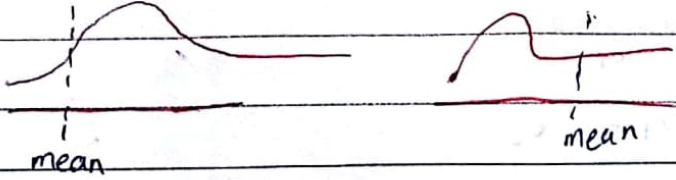
در حالت نرمال میانگین معیار خوبی است



اما ممکن است نمودار چوکی داشته باشد

یعنی یک سمت منحنی کشیده باشد

در این دو حالت میانگین معیار خوبی نیست



Median: وسطی است

داده ها را مرتب می کنیم، داده ای که در وسط قرار می گیرد همان است.

افراد در این سن

0-20

20-40

40-60

60-80

80-100

فرض کنید یک سری داده را داریم:

همانند در دسترس قرار می گیرد

$$Median = l_1 + \frac{N/2 - \sum_{j=1}^{k-1} f_{reqj}}{f_{reqk}}$$

کلاس پایین دسته ای
همان در آن قرار دارد

f_{reqm}

فراوانی پایه

فراوانی پایه های کوچکتر از پایه میانه (دسته های قبلی)

10

همچون داده ها دارند، آن ها را دسته بندی می کنیم

Mode: مودو

uni mode
no mode

داده ای که تکرار آن از بقیه بیشتر باشد

uni mode ← از هر داده یکی داریم * اگر دسته ای که بیشترین تکرار را دارد دسته مد است

Bi Mode ← از هر داده ۲ تا داریم یعنی ۲ داده که بیشترین تکرار را دارد فراوانی دسته قبل از مد

No Mode ← داده های ما عدد ندارد. *فردی که در آن دسته است*

* اگر صورت دسته بندی بود عدد صحیحی مشخص کنیم: $M = l + \frac{D_{before} \times w}{D_{after}}$

mean - mode = 3 x (mean - median)

Quartile: هر کدومی داده ها چگونه است

داده ها را مرتب کرده و آن ها را به ۴ دسته تقسیم می کنیم

۲۵٪ داده ها Q₁

۵۰٪ داده ها Q₂

۷۵٪ داده ها Q₃

۱۰۰٪ داده ها Q₄

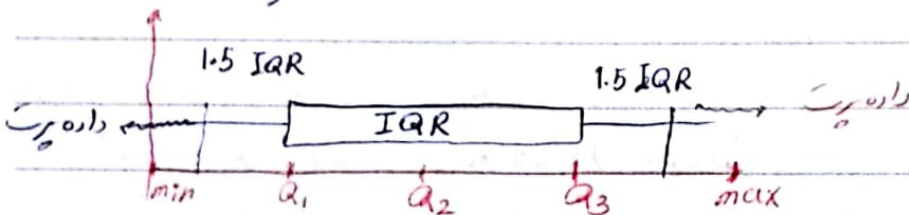
25

مهمترین توزیع در آمار است و استفاده می کنند، توزیع نرمال است

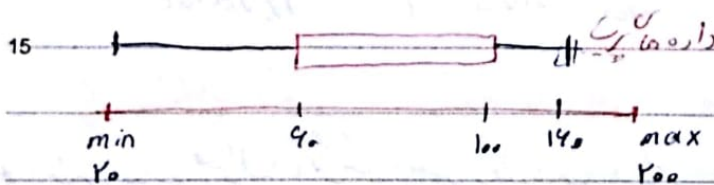
Inter Quartile Range = IQR = Q3 - Q1

Five Num: Q1 - Q2 - Q3 - Min - Max ← آماره های پنج گانه

Box plot ← عرض کند داده ها را از min تا max را رسم



استفاده از Boxplot چگونه داده ها را نمایش می دهد؟
1.5 برابر IQR را حساب کنید.
بعد از Q3 به اندازه 1.5 برابر IQR را قبول کنید و بقیه را داده های پرت بنامید.
قبل از Q1



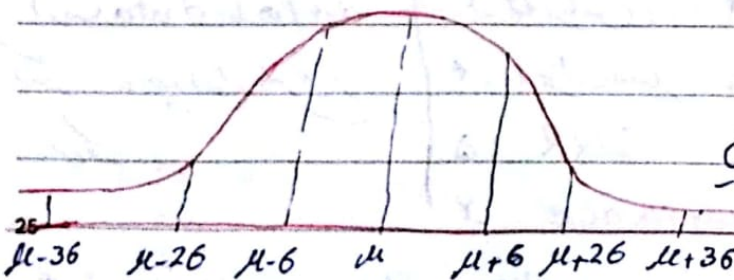
مثال:

IQR = 100 - 40 = 60 1.5 IQR = 90

قبل از Q1 ← 60 - 60 = 0 اما داده ها نباید از 0 کمتر پس از این طرف
داده های پرت نداریم. و بعد از Q3 به اندازه 1.5 IQR علاوه بر 140 از 160
به بعد داده های پرت می باشد.

توزیع نرمال:

بنابر این داده ها در سه خارج از این
بازه اند داده پرت اند.



Soroush

انجمن علمی علوم کامپیوتر
دانشگاه کاشان

95% داده ها → 97% داده ها

99% داده ها

اگر دو متغیر مستقل باشند:

$$\text{Cov} = 0$$

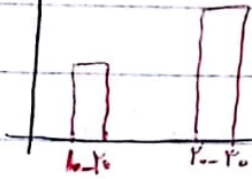
Cov > 0

Page: ()

Subject:

Year: Month: Day: ()

y



معمولاً: Histogram

این گنبد که داده‌ها نام صورت

بازه، بازه در نظر می‌گیرد، گنبد سازی می‌کند n

5

می‌توانیم بفهمیم که نمودارهایی که ضلعی بلند یا حتی کوتاه‌اند، پرت‌اند.

تعداد فروش قیمت

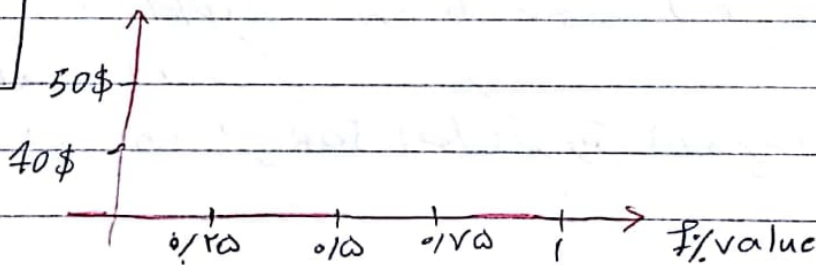
40\$	300
50\$	120

نمودار: Q-Plot : Quantile Plot

محور x آن ها F/Value می‌باشد.

کی نمودار داریم که از هر کالایی چه قدر فروخته‌ایم؟

10



15

؟ مثال بزنیم

* برابر بررسی میزان تا مرتبه متغیر از داده‌ها و برابر بررسی همبستگی دو متغیر از داده‌ها این استعاره می‌کنیم

1. دامنه تغییرات: Max - Min

$$2. \text{انحراف جابجایی داده: } \frac{1}{n} \sum_{i=1}^n |x_i - \mu|$$

قد مطلق در درجه صحت این سریم سرانج معیار دیگر است و از این

$$3. \text{واریانس: } \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$4. \text{انحراف معیار: } \sigma$$

5. IQR

20 شاخص دیگر از این است

از این data تا data دیگر با چه انداز نام تغییر می‌کند.

25

$$6. \text{Covariance ضریب همبستگی: } \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

* مثلاً از روی دمای امروز می‌خواهیم دمای فردا را پیش‌بینی کنیم. دهم به از واریانس استفاده می‌کنیم.

Soroush اما دمای داریم هر طوبی بعد از روی این که می‌خواهیم دمای فردا را پیش‌بینی کنیم به توانیم

شد رابطه رطوبت امروز با دمای فردا چیست؟ اثر

Quantile - Quantile plot Q-Q plot نمودار؟

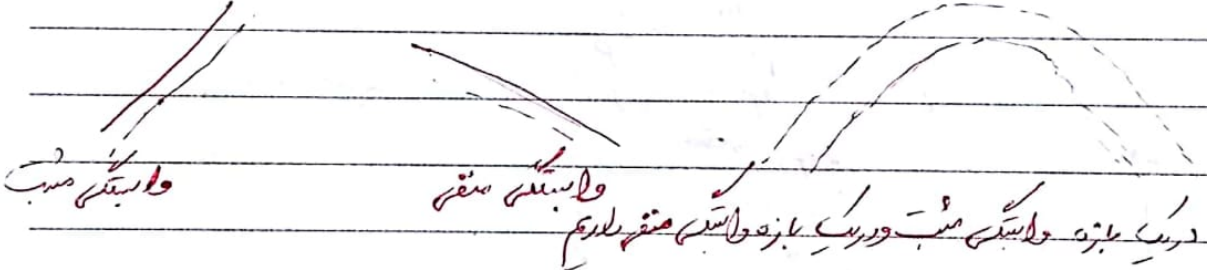
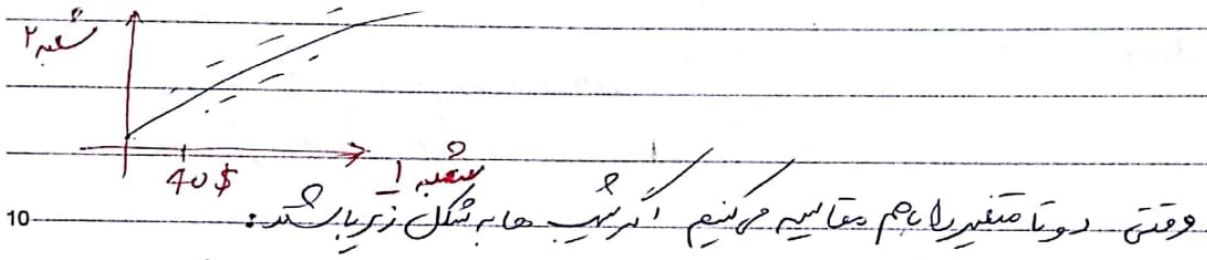
اقلیم که در فروتنگاه 1 به فروتنی وقت نام میبرد

و برابر هر کدام یک Plot رسم می کنند

در واقع نسبت به هم می گویند مثلا تعداد فروتنی 80 دلار در شعبه 1 معادل

صورت شود با تعداد فروتنی 90\$ در شعبه 2

این محبت صفا با بررسون است در بررسون هر دو هم ارتباط بین دو تابع غیر وابسته



similarity ساهت dissimilarity عدم ساهت مفهوم؟

ما باید بیوانیم در سنجش این دو هم به یکی از تکنیک های اصلی این درس داده های مشابه را در یک کلاس با دسته قرار دهیم چون داده ها جنس زیادند فرض کنید یک پایگاه داده جنس بزرگ از کل دانشگاه ها کشور داریم

و هر یک مورد شامل یک سری صفات است. حال ما خواهم رکوردی مشابه را در یک کلاس میذاریم به بحث میپردازیم هر شود ساهت از کدام نظر؟

فرض می کنیم یک ماتریس داده داریم

n رکورد - m ویژگی

از روی این ماتریس عدم ساهت را

تکلیف می دهیم

x_{11}	x_{12}	x_{1m}
x_{21}	x_{22}	x_{2m}
\vdots	\vdots		\vdots
x_{n1}	x_{n2}	x_{nm}

ماتریس عدم شباهت:

$$\begin{bmatrix} 0 & & & \\ d_{2,1} & 0 & & \\ d_{3,1} & d_{3,2} & 0 & \\ d_{4,1} & d_{4,2} & d_{4,3} & 0 \end{bmatrix}$$

تعداد کل ویژگی‌ها که نام تطابق دارند = $\frac{P-m}{P}$

تعداد کل ویژگی‌ها = $\frac{P}{P}$

نسبت = $d_{ij} = \frac{P-m}{P}$

فرض کنید یک پایش داریم که فقط دو رنگ دارد اما رنگ مو هم باشد
 اگر دو نفر رنگ مو هم داشته باشند $d_{1,2} = \frac{1-1}{1} = 0$

اگر دو نفر باشند که یکی رنگ مو هم داشته باشد و دیگری رنگ مو هم نداشته باشد $d_{1,2} = \frac{1-0}{1} = 1$

شباهت $S = 1 - \frac{P-m}{P} = \frac{m}{P}$

مثال: فرض کنید سری داده به شکل زیر داریم:

ID	test 1	test 2	test 3
1	A	Excellent	45
2	B	Fair	22
3	C	Good	64
4	A	Excellent	28

فرض کنید ماتریس عدم شباهت برای داده های زیر:

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

test 4

Brown

Blod

Brown

Brown

$d_{1,2} = \frac{2-0}{2} = 1$

- صفا $1,1 \rightarrow 9$
- $0,0 \rightarrow T$
- $1,0 \rightarrow 2$
- $0,1 \rightarrow 5$

فرض کنید سری داده با نری داریم:

علم سلامت بین داده های باینری متعارف:

$$dis(i, j) = \frac{r+s}{q+r+s+t}$$

باینری متعارف = اصیت اینکه جواب منفی شود یا مثبت شود، یکسان است.

5 علم سلامت بین داده های باینری نامتعارف:

$$dis(i, j) = \frac{r+s}{q+r+s}$$

تعداد داده های نامتعارف:

$$Sim-Jaccard = \frac{q}{q+r+s}$$

10 فرض کنید دنبال افرادی که HIV دارند میگردیم تعداد آن ها چندیم است و اما پیدا کردن آن ها برابر با چندیم است.

مثلاً هر دو اصیم رابطه بین HIV در بابت دارویی داریم پس جابری هم است که هر دو اینها باشند. مثلاً رابطه بین جنسیت و رخسار مثلاً زنان سفیدی پس جابری است که مردان سفیدی باشند اصیت ندارد. سرفه تب

15 مثال:

	Gender	Fear	Cough	test 1	جواب آزمون
Jack	M	Y	N	P	مثبت - منفی
Mary	F	Y	N	P	داده از جنبی باینری اند.
Jim	M	Y	P	PV	

20

$$dis(Jack, Mary) = \frac{0+1}{2+0+1} = 0.33$$

1 ← P
0 ← N

25

اگر داده های نامتعارف اند. (جنسیت را در نظر نمیگیریم) متعارف است

test 2	test 3	test 4
N	N 0	N
N	P 1	N
N	N	N

r = 0
s = 1
q = 2

Jack, Mary با هم اختلاف زیادی ندارند نسبت تقریباً زیاد است.

dis (Jack, Jim) =

نرم افزار WIKI و rapid miner

اسلایدها : Data Mining - lecture 5

در این جلسه مرخصی Similarity و dissimilarity را برای هر رکورد بیایم

5 مخاطب حلیه دلال مرخصی مفهوم مقدار رو object مرخواند سیم هم باشه

اگر خریدار یک کالا خرید میکنیم کتاب کلاسی دیگر را امکان دارد بخرد

خریدارهایی که مشابه خرید میکنند در یک گروه دسته بندی میکنیم و تبلیغات را به همه افراد

آن دسته میفرستیم

بگویم نه متن ها رو ب را دسته بندی میکنیم. دسته های وب سیم به هم را پیدا میکنیم

10 دنبال تراکس های متفاوت هستیم که اصلاً آن ها مشکل دارند

برای همین معنی Similarity و dissimilarity را می بینیم

Similarity: مقدار عددی، این مقدار رو object ما سیم به هم اند

روسی را مقایسه میکنیم و یک مقدار می دهیم. هر چه نزدیکتر باشد، شباهت بیشتر است.

15 اگر این عدد نزدیک صفر باشد شباهت کم است

مقدار می تواند 1 تا 100 هم باشد

یعنی دو object سیم هم اند. $S(p, q) = S(q, p)$ اگر

Jac - similarity = $\frac{\text{اشتراک دو مجموعه}}{\text{اجتماع دو مجموعه}}$

20

بدست آوردن شباهت بین Vector ها:

فرض کنید 2 تا document داریم. مرخصی Similarity بین این دو را

اندازه بگیریم. تعداد شباهت هایی که بین این دو وجود دارد را بدست آورده

و تقسیم بر تعداد کل کلمات میکنیم

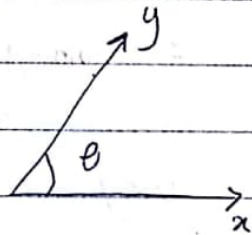
25

مثلاً $\frac{1}{3}$ متن شده apple. حل یک فرمول پیدا کرده شباهت بین دو vector را بفهمیم.

apple در متن ۱، ما با یک کلمه شده در متن ۳۰

	Apple	MicroSoft	Obama	Electro
D1	10/30	20/30	0	0
D2	30/90	60/90	0	0

برای پیدا کردن شباهت چندین مورد داریم:



۱- Cos-Similarity: مثلاً داریم:

سینه زینت → عمود بر هم → $\theta = 90$
 $\cos \theta = 0$

سینه اند → بر هم منطبق اند → $\cos \theta = 1$
 $\theta = 0$

نقشه هندسی اندازه شباهت: $\cos(x, y)$ هر شود شباهت بیشتر

برای مقایسه Document ها استفاده از Cos رایج تر است. به شرطی که بتوان از Cos استفاده کرد که از قبل داده ها را زوال سازی کنیم.

$$\cos(d_1, d_2) = \frac{d_1 \cdot d_2}{\|d_1\| \|d_2\|}$$

۲- distance: مقداری که هر فهمیم دو object چقدر متفاوت اند. هر چه به صفر نزدیک باشد شباهت بیشتر.

برای پیدا کردن فاصله چندین رویش داریم:
الف) اقلیدسی:
$$\text{فاصله در اقلیدسی} = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots}$$

ب) منصف:
$$\text{فاصله منصف} = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots$$

ج) مین فونسی:
$$\text{فاصله در مین فونسی} = \left(|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots \right)^{\frac{1}{h}}$$

د) Hamming distance:

تعداد بیت‌هایی که با هم تفاوت دارند (در دو کد)
10) خطری شباهت بین string ها را پیدا کنیم؟ در آنتن DNA کاربرد دارد

Edit distance: تعداد Insert و Delete و Replace که باید
شده هم باشند.

فصل 4 کتاب - پیدا کردن قوانین انجمنی و الگوهای تکرار

مفاهیم: الگوی تکرار: مثلا در یک فروشگاه کالای است که بیشتر از همه خریداری شده اند.

چرا الگوهای تکرار را باید یاد کنیم؟

20) یکی از سوالات: بعد از چه کالایی، چه کالایی می‌خریم. یعنی توانی در خرید کالا مثلا می‌دانیم هر کس کامپیوتر بخرد، بعد از آن پرینتر می‌خرد. پس ما می‌توانیم تخفیف برای پرینتر می‌داریم که صرفا خریداری شوند.

1- توانی در خرید کنیم

2- چه کنیم تا با هم خریداری شده اند. (مثلا سرویس ما بهم خرید شده اند پس اینها)

3- DNA تکرار چند پروی چند تکرار (تکرار می‌کنیم)

مثلاً یک بار روی DNA های مختلف چه تأثیری نداشته
با مثلاً چه میزان فروشنده چنانچه در میزان فروش داشته

۴- طبقه بندی اسناد و مثلاً این خبرند اینها باید کنار هم باشد
مثلاً هر کسی اخبار ایران را میخواند اخبار افغانستان هم میخواند پس
این دو باید کنار هم باشند

جایی زیادی کاربرد دارد مثل سیستم فروش و بازاریابی، علم ترسیمی
و این صفحات وب به دردی اینها نیاز داریم اشخاص دیگر را با هم لینک

10

سوال: اگر کسی کامپیوتر خرید بعد از خرید آن به چه اصطلاحی گفته میشود؟

قانون: باید یک قانون را پیدا کنیم که این اتفاق افتاده

15

Computers \Rightarrow Software [support = 27, confidence = 60%]

$$A \Rightarrow B$$

باید روابط با هم

20

فرض کنیم فروشنده n سیستم به شکل زیر دارد که هر کدام یک کالا A میخریم.
بعد از خرید A، صفت B را میخریم که اشتراک این دو در آن است.
هر وقت سمت راست وجه نباید اشتراک داشته باشند

T: مجموعه اکتیوها

I: item

25

$$A \Rightarrow B$$
$$A \cap B = \emptyset$$

$$T = \{I_1, I_2, \dots, I_n\}$$

Support: (درجه پشتیبان) فرض کنید قانون در مجموعه ترانسکشن S و D باشد. S در این صورت S در D درصدی از ترانسکشن D است که حاوی A و B است و با احتمال $P(A \cup B)$ پیدایش می‌دهند

$$\text{support}(A \Rightarrow B) \Rightarrow A \cup B$$

یعنی جایی که دو کالا با هم خریداری شده‌اند.

Confidence: اگر طرف اول true باشد (یعنی A رخ دهد) چه قدر احتمال دارد که طرف دوم نتیجه درست باشد.

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

یعنی اگر A رخ دارد چه قدر احتمال دارد B رخ دهد؟

این support که می‌بینیم min حالت است.

پس در آن $\text{min} = 4\%$ یعنی خیلی جاها A و B با هم اند پس خیلی قانون باطل گرفته ایم. اما وقتی می‌بینیم 98% یا پورت باشد پس داریم قانون با ضمیمه جهت می‌بینیم یعنی خیلی کم می‌بینیم هرگز فاصله 98% با این دو با هم خریداری شده باشند.

اگر قانون هیچ چیزی پیدا کنند پس این قانون بدر دهنی خورد.

باید مثلاً قانون را این بگیریم کسانی که امروز یک کالا خریده‌اند. این احتمال خیلی زیاد است. از آن چیزی خواص بدست آوریم که هیچی بدر دهنی خورد!

پس زیاده‌ریا خارج کردن و کم‌ریا خارج کردن بدر دهنی خوردند. به چیز متوسط می‌خواهیم نه زیاد بقیه‌مانند که چیزی پیدا کنند و نه زیاده‌ریا که از آن بی‌بهره است.

Itemset مجموعه اکتسم‌ها

تعداد ترانسکشن‌ها = N و صورت ترانسکشن = A اینم

$$(2^k - 1) \times N \times A \rightarrow$$

تعداد اکتسم‌هایی که می‌توانیم (استفاده کنیم)

Item

K-Itemset اگر Itemset ها K عضو داشته باشند

فراوانی (تکرار) : نرخ تکرار یک Item در مجموعه دیتا است

5 support → Absolute Support تعداد رخ داده های که A و B همزمان با هم رخ میدهند

→ Relative Support احتمال اینکه A و B با هم رخ میدهند

Relative Confidence: فنون دارد چون Confidence مستقیماً رخ داده و -

Relative احتمال = $\frac{\text{تعداد رخ داده A و B}}{\text{تعداد کل رخ داده ها}}$

مسئله: فرض کنید یک فروشنده چندین تراسینگ دارد

TID	Items	Notes
10	Coke, nuts, diaper	2 = min-support قرار می دهیم
15 20	Coke, coffee, diaper	پس باید بیشتر از 2 بار آمده باشد
30	Coke, diaper, egg	تا ما بنویسیم
40	Nut, Egg, milk	این مجموعه Itemset - 1 است
50	Nut, coffee, diaper, egg, milk	باید کنیم که برشمارند

1 - Itemset = {Coke} {nuts} {diaper} {coffee}

{egg} {milk}

پس این روش خنجر کزانی دارد

لذا این support = 4 باشد

که هیچ چیزی بر این مورد حتی برای 1-Itemset پس فیلتر می کنیم

25 2 - Itemset = {Coke, diaper}

{nut, egg} {egg, milk}

3-Itemset = {nut, egg, milk}

4-Itemset و همچنین پیدا کنیم
پس قانون پشتیبانی هر آیتم خرید باید پیدا کنیم
اما اگر $min\ support = 1$ باشد این آیتم هم پیدا می شود.

5
پس $min\ support$ را جوری در نظر بگیرند که این مجموعه ای به دست می آید
متناسب باشد. $support\ confidence$

$Coke \Rightarrow diaper$ ($\frac{3}{5} = 60\%$, $\frac{3}{3} = 100\%$)

اینجا $diaper \Rightarrow Coke$ خرید و بعد $(\frac{3}{5} = 60\%$, $\frac{3}{4} = 75\%$)
Coke سفرد

تین ضرایب $Support$ و $Confidence$.

او محطه باره کاری چند بزرگ است و این روش ها کم میارند

15 اگر یک مجموعه پرستار باشد، آیا زیر مجموعه آن هم پرستار است؟ بله
فرض کنید یک مجموعه 100 آیتمی پیدا کردیم که پرستار است.

تعداد زیرمجموعه های این مجموعه:
زیرمجموعه 2 عضوی
زیرمجموعه 1 عضوی

$$\binom{100}{1} + \binom{100}{2} + \binom{100}{3} + \dots + \binom{100}{100} = 2^{100}$$

پس 2^{100} مجموعه پرستار پیدا کردیم.

20 اما نیازی نیست همه این ها را ذخیره کنیم فقط 400 مجموعه 100 آیتمی را ذخیره می کنیم.

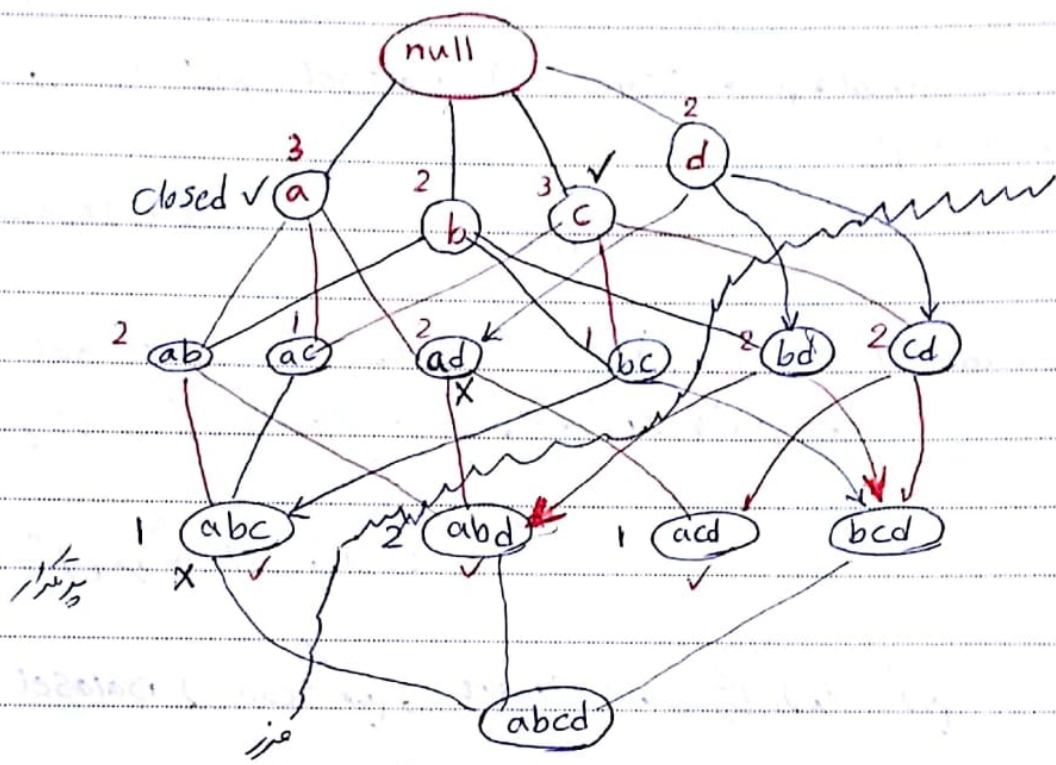
Closed Frequent Itemset: یک Itemset مانند X Closed است

هیچ Super Itemset ای مانند X وجود نداشته باشد که $support$ آن

بیشتر است بالاتر از خود X - فراتر از خود X - مناسب X باشد.

مثلاً در مثال قبل از مجموعه 3-Itemset ای چیزی بیشتر پیدا نکردیم

Soroush پس این مجموعه Closed است.



در این گراف یک فرزند بین مرتکبها و یک تکرارها می کشد. کشیدن کن ضریب واحد نیست
 اما این فرزند ممکن است دقیق نباشد.
 آنهمهایی که روی فرز هستند support های آن را بررسی می کنیم.
 اگر آنهمهایی که زیر مجموعه های در سمت پایین فرز باشد، آن ها کسبای فری گویند
 دنیا سمت فر گویند مثل abc و ad

۳ تا الگوریتم برای پیدا کردن مجموعه های مرتکبها:

Apriori - FP Growth - ECLAT

Apriori بر اساس دانستن متدی و یک الگوی مرتکبها استخراج می کنند.
 اگر مجموعه ای مرتکبها باشد زیر مجموعه های آن مرتکبها است.
 این الگوریتم از این قانون استفاده می کنند.
 اگر یک K-Itemset داشته باشیم از روی آن می توانیم یک (k+1) Itemset
 مرتکبها پیدا کنیم.

Subject:

Year. Month. Date. ()

مراحل الگوریتم:

۱- در ابتدا مجموع Dataset و Scan می کنیم: اول باید دسته بندی را نگاه کنیم پسیم که این اتفاق افتاده

۲- Itemset 1- های پر تکرار را استخراج می کنیم: L_1

۳- Itemset 2- ها را با استفاده از join داخلی Itemset 1- ها بدست می آوریم بعد پر تکرار های آن را استخراج می کنیم و اسم آن را C_2 می نامیم.

همان مثال فردی که در نظر می گیریم:

۱- Scan Dataset و Itemset 1- های پر تکرار را بدست می آوریم:

$$L_1 = \{ \{cok\} \{nut\} \{diaper\} \{coffee\} \{egg\} \{milk\} \}$$

$$L_2 = \{ \{cok, nut\} \{cok, diaper\} \{cok, coffee\} \{cok, egg\} \{cok, milk\} \}$$

$$\{nut, diaper\} \{nut, coffee\} \{nut, egg\} \{nut, milk\} \}$$

minsupport = 2 \Rightarrow هر دسته ای که تعداد تکرار آن از minsupport کمتر باشد حذف می کنیم.

برای بدست آوردن Itemset 3- ها L_2 را با خودش join می کنیم (عصبه از وسطه قطع)

پر تکرار های L_2 را با خودش join می کنیم \rightarrow اگر ۴ تایی بود یعنی نویسیم فقط سه تایی ها را می نویسیم

$$L_3 = \{ cok, nut, diaper \}$$

این مثال کامل نیست

این روند را تا جایی ادامه می دهیم که دیگر مجموع پر تکرار پیدا نکنیم و به نویسیم Closed می شود.

Subject:

Year. Month. Date. ()

سوال: مجموعہ های پر تکرار را بدست آورید، با استفاده از الگوریتم.

TID	Items
10	A, C, D
20	B, C, E
30	A, B, C, E
40	B, E

minsupport = 2

$\mathcal{L}_1 = \{ \{A\}, \{B\}, \{C\}, \{D\}, \{E\} \}$
 هر یک در یک سبک
 Pruning چون minsupport = 2 است

$$\mathcal{L}_2 = \{ \{A, B\}, \{A, C\}, \{A, E\}, \{B, C\}, \{B, E\}, \{C, E\} \}$$

1 2 1 2 3 2

$$\mathcal{L}_3 = \{ \{A, B, C\}, \{A, C, E\}, \{B, C, E\} \}$$

1 1 2

این مجموعہ {B, C, E} پر تکرار است ← زیر مجموعہ های آن هم پر تکرار است
 مجموعہ D حذف شد و در مراحل بعدی D را ندیدیم.
 این اگر تکرار مجموعہ پر تکرار نبود در مراحل بعدی نیز پر تکرار نیست.

سوال: * در صفحه بعد کار بردارو.

TID	Items
T100	I ₁ , I ₂ , I ₅
T200	I ₂ , I ₄
T300	I ₂ , I ₃
T400	I ₁ , I ₂ , I ₄
T500	I ₁ , I ₃
T600	I ₂ , I ₃
T700	I ₁ , I ₃
T800	I ₁ , I ₂ , I ₃ , I ₅
T900	I ₁ , I ₂ , I ₃

minsupport = 2

Subject:

Year. Month. Date. ()

عیب این روش این است که هر بار باید کل Dataset را Scan کنیم
 یا مثلاً یک سری مجموعه‌ها را اضافه داریم که باید آن‌ها را بررسی کند پس بارهایی
 تولید کند که در درخت ذخیره باشند و عمل تولید کند
 در نهایت پس‌ها بزرگ خنجر بزرگ شود یعنی زمان بر - زیاد اینجا مثلاً در نهایت کوچک بود

اصول ترین چالش‌های معسایره :

- ۱- مرور چندباره پایگاه داده تراکنش
 - ۲- تعداد زیاد کاندیداهای تولید شده
 - ۳- معیار Support برای کاندیداهای زمان تراکنش
- به دنبال راه حل هستیم که این چالش‌ها را برطرف کنیم
 راه حل که : ← ۱- کاهش مرور (اسکن) پایگاه داده
 ۲- کاندیداهای مفید تولید شوند
 ۳- معیار سریع Support ها

روش‌های عبور Apriori : ۱- درم سازی ۲- کاهش تراکنش ها ۳- پارسی کردن
 ۴- نمونه گیری ۵- شمارش پویا

درم سازی : روش یک مثال یاد می‌کنیم

- ایده اصلی : هر زمان با این Itemset ها تولید شوند ۲-Itemset ها هم تولید شوند
- ۱- تابع درم سازی تعریف می‌کنیم و عنصر 2-Itemset را به مقدار تابع hash
 - ۲- هر مجموعه یک bucket گفته می‌شود
 - ۳- اگر شمارش اعضای bucket از min-support بیشتر باشد
 - ۴- عنوان اکتوی برقرار یا frequent ساخته می‌شود

Subject:

Year. Month. Date. ()

TID

Items

مثال: این مثال در صفحه قبل نوشته شده.

I_1, I_2, I_3

$$\text{hash} = (\text{order of } x) \times 10 + (\text{order of } y) \% 7$$

بافته فایده تقسیم بر ۷ می تواند اعداد ۰ تا ۶ باشد. پس این حالت ها را می توانیم داشته باشیم

$$\text{مثال } (I_1, I_4) \Rightarrow 1 \times 10 + 4 = 14 \% 7 = 0$$

bucket Address	0	1	2	3	4	5	6	Frequent
bucket Count	2	2	4	2	1	4	4	
bucket Content	{ I_1, I_4 }, { I_3, I_5 }	{ I_1, I_5 }	{ I_2, I_3 }, { I_2, I_3 }, { I_2, I_3 }	{ I_2, I_4 }	{ I_2, I_5 }	{ I_1, I_2 }, { I_1, I_2 }, { I_1, I_2 }	{ I_1, I_3 }, { I_1, I_3 }, { I_1, I_3 }	✓

این bucket Count آن ها بیشتر از ۳ باشد Frequent است. $\text{minSupport} = 3$ ← مجموع های که

تابع hash یک سری تداخل دارد. دوتا چیز ممکن است یک Frequent بخورند. این روش خطا دارد. می توانیم تابع hash را جور دیگری انتخاب کنیم که تداخل بر طرف شود.

تداخل ← مثلا ۳ مجموعه داریم اما این ۳ مجموعه ممکن نیستند. دوتا { I_1, I_2 } و دوتا { I_1, I_2, I_3 } با هم اشتراک روینیم Frequent را چندین باریم.

کلاس تراکنش: اگر تراکنش در مجموعه Itemset k است و اگر یک آیتم تکرار ندارد پس در مجموعه Itemset k+1 هم هست. اگر آیتم تکرار ندارد بنابراین آن تراکنش را حذف می کنیم.

مثلا { I_2, I_4 } تکرار نیست. مربوط به کدام تراکنش است؟ تراکنش 200. پس تراکنش 200 را حذف می کنیم و رفت بعدی آن را بررسی می کنیم. اما تراکنش 400 را حذف نمی کنیم! چون I_1 را هم دارد. تراکنش های را حذف می کنیم فقط I_2, I_4 را دارند.

Subject:

Year. Month. Date. ()

پارتیشن بندی: فرض کنید دیتابیس داریم که 10G است اما فضای که داریم 1G است. اگر قرار باشد هر دفعه 1G از دیتابیس را بیاریم و نگاه کنیم، Itemset-1 پیدا کنیم. این شکل باید 10 بار از دیتابیس اطلاعات را بیاریم و در حافظه اصلی مثلاً فرض کنید $minSupport = 6$ 20 بار.

دیتابیس ها ه ا قیمت داریم اگر هر چیزی با آنتیم در هر قسمت 2 بار تکرار شده باشد در حافظه در ه قسمت 20 بار تکرار می شود و به تکرار است اما ممکن است توزیع داده ها زغال نباشد مثلاً یک جا 10 بار تکرار شده و یک جا 1 بار پس می توان نتیجه گرفت اگر در یک قسمت تکرار باشد حتی در حافظه هم تکرار است

Local ← در یک پارتیشن Global ← در کل

اگر یک Item در کل تکرار باشد در یک قسمت با Local هم تکرار می شود این قسمت مثل قسمت های قبلی نیست که با Data کار داشته باشد بلکه با انتقال از دیتابیس به حافظه کار دارد.

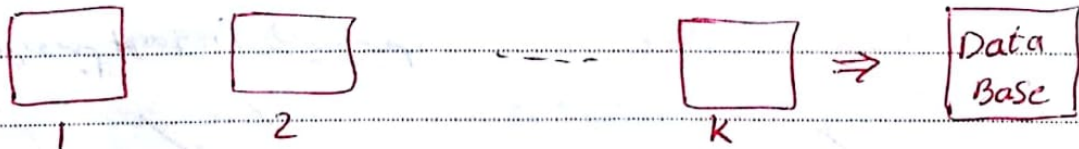
Local Global

Frequent ← همان در یک قسمت

Frequent ↗

پارتیشن بندی در دو فاز انجام می شود:

دو فاز اول: داده ها به چند بخش تقسیم می شوند (به طوری که هر بخش بتواند در حافظه اصلی قرار گیرد)



نیست $Sup > min support$

$Sup > min support$

تعریف $min support$ نیست: داریم $Sup > min support$ = (باید) min تکرار برای هر پارتیشن برابر است با $min support \times$ تعداد تکرار آن در هر یک بخش

مثلاً $min support = 20$ است و دیتابیس 10 قسمت تقسیم می شود به طور متوسط در هر قسمت 20 تکرار باید

HAMKELASI

انجمن علمی علوم کامپیوتر

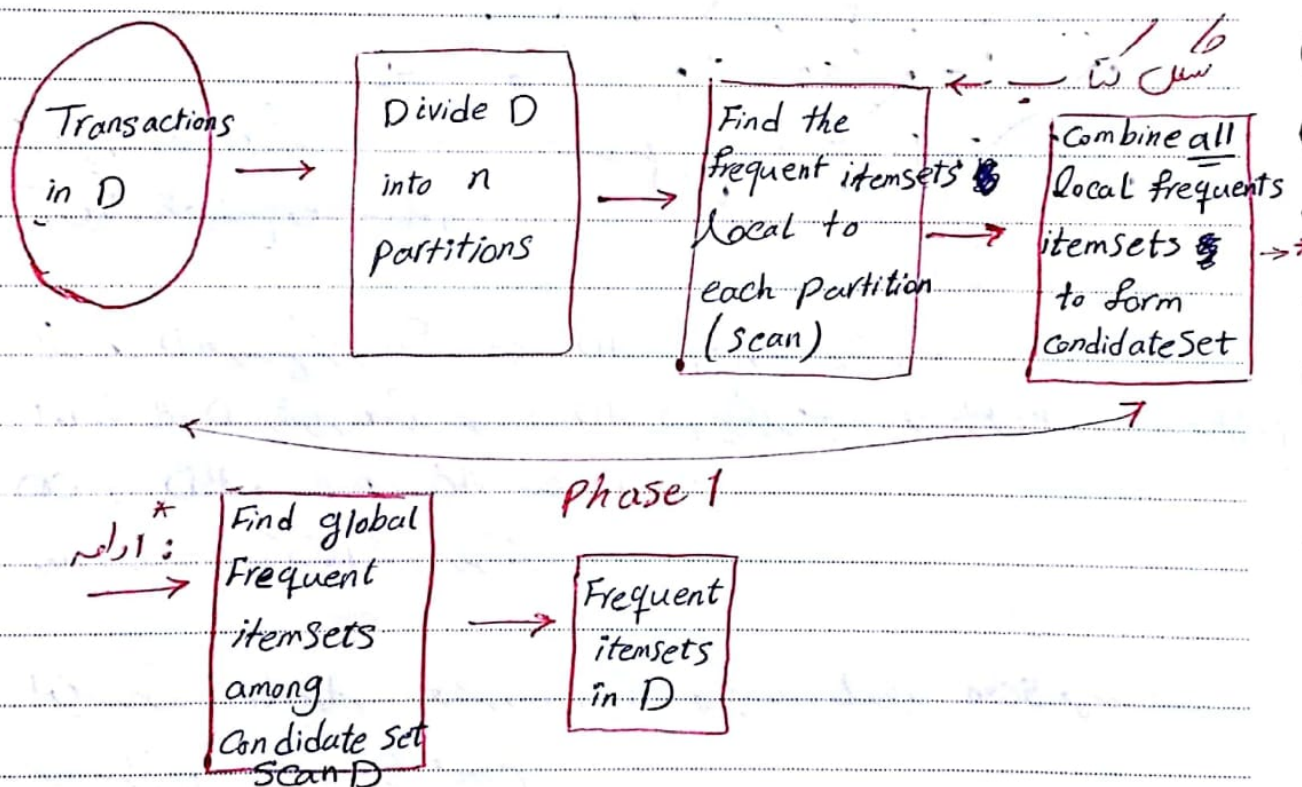
دانشگاه کاشان

نیست $min support$

t.me/KUCSSA

- ۱. الگوهای پر تکرار هر بخش (الگوهای مهم) ممکن است نسبت به کل پایگاه داده پر تکرار باشند یا نباشند
 - ۲. هر Itemset که در پایگاه پر تکرار باشد باید حداقل در یکی از پارتیشن‌های پایگاه داده پر تکرار باشد.
 - ۳. نتیجه هر itemset پر تکرار در یک بخش می‌تواند در مجموع کل پایگاه داده باشد. (یعنی هنوز کاندید است و اینکه پر تکرار باشد یا نه در فاز دوم مشخص می‌شود)
- فاز دوم ← مرور روی پایگاه داده
- برای هر itemset کاندید باید Support واقعی را محاسبه کنیم

در واقع این روش هم دوره پایگاه داده را مرور می‌کند چندین بار



د. itemset کل support
 و min support کل می‌تواند کنیم

Subject:

Year. Month. Date. ()

نمونه برداری :

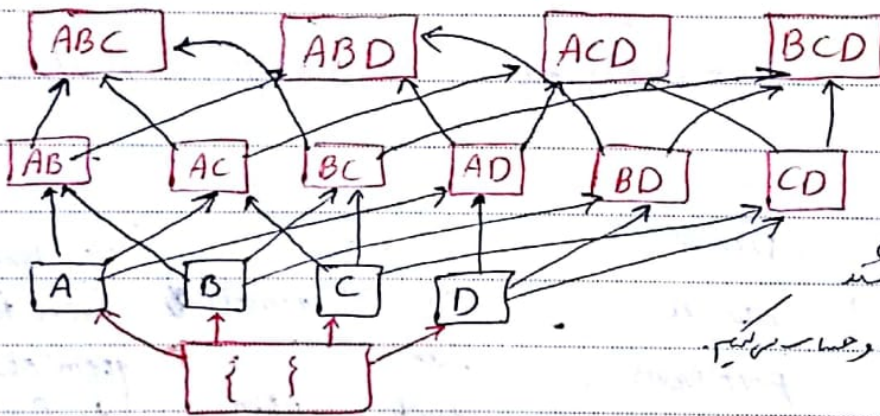
مسئله : در حالتی نمونه برداری کنیم که نمونه‌های حاصلی به هم نزدیک باشند اما اگر نمونه‌ها از هم فاصله داشته باشند این کار قابل قبول نیست.

سپارشی پویا : هدف این کار برای کاهش تعداد اسکن‌هاست.

مسئله : اول کار null است، هیچی نداریم. بعد 1-Itemset کارها می‌نویسیم.

ABCD

بعد 2-Itemset کارها می‌نویسیم.



اگر زیری‌ها تکرار باشند اسکن‌ها در نظر نمی‌گیریم و حساب نمی‌کنیم.

اگر A و D هر دو تکرار باشند $AD \Leftarrow AD$ هم می‌تواند تکرار باشد.
 اما اگر A و D تکرار نباشد AD هم تکرار نیست و در مراحل بعد آن را در نظر نمی‌گیریم.
 $BCD \Leftarrow AC$ و BC و BD و CD
 زمانی تکرار است که تمام مجموعه زیر آن تکرار باشد.

این روش Apriori و وجودهای آن پایه داده‌های زیاد Scan می‌کنند
 پس سرانجام روش‌های دور می‌رویم.

Subject:

Year. Month. Date. ()

روش FP-Growth روشی برای ساخت الگوهای پر تکرار از طریق استخراج الگوها

Frequent pattern Growth Approach

این الگوریتم به صورت افزایشی رشد میکند Growth

هدف آن یافتن Itemset های پر تکرار بدون تولید مجموعه های کاندید

اینجا الگوها آن به دلیل معایب روش Apriori بوده است

معایب ← از حقیقی بودن به سطح؛ تراکشن های یکبار بر سرش می کشند و پس از آن مجموعه های کاندید

تولید می کنند و این از سمت مجموعه های کاندید پر تکرار بودن را مشخص می کند

۲- تعداد کاندیدها زیاد بود

FP-Growth ← ۱- حقیقی بودن عمق

۲- مجموعه های کاندید را سریع تولید نمی کند

فازهای روش FP-Growth ← ۱- درخت FP را تولید کنند

۲- Conditional fp-tree را مشخص کنند

نحوه ساخت درخت FP ←

۱- عناصر Itemset 1 ها را پیدا کرده و Support آن ها را محاسبه کنند

۲- بر اساس مقدار Support مرتب کنند

1-Itemset Support

سوال * در چند صفحه قبل یاد گرفتیم

اول Itemset 1 ها را تعیین کنند

I₁

6

مرحله ①

I₂

7

I₃

6

I₄

2

I₅

2

مرحله ②
مرتب کردن

I₂

7

I₁

6

I₃

6

I₄

2

I₅

2

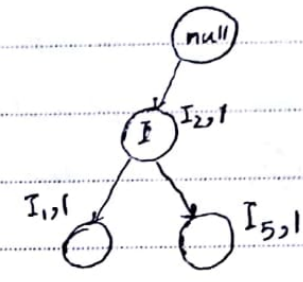
Subject:

Year. Month. Date. ()

ساخت درخت تراش حاصل مرتب کردن

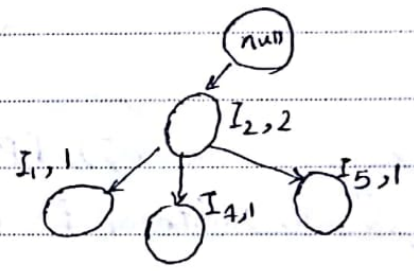
TID	Items
100	$I_1, I_2, I_5 \rightarrow I_2, I_1, I_5$
200	$I_2, I_4 \rightarrow I_2, I_4$
300	$I_2, I_3 \rightarrow I_2, I_3$
400	$I_1, I_2, I_4 \rightarrow I_2, I_1, I_4$
500	$I_1, I_3 \rightarrow I_1, I_3$
600	$I_2, I_3 \rightarrow I_2, I_3$
700	$I_1, I_3 \rightarrow I_1, I_3$
800	$I_1, I_2, I_3, I_5 \rightarrow I_2, I_1, I_3, I_5$
900	$I_1, I_2, I_3 \rightarrow I_2, I_1, I_3$

ساخت درخت

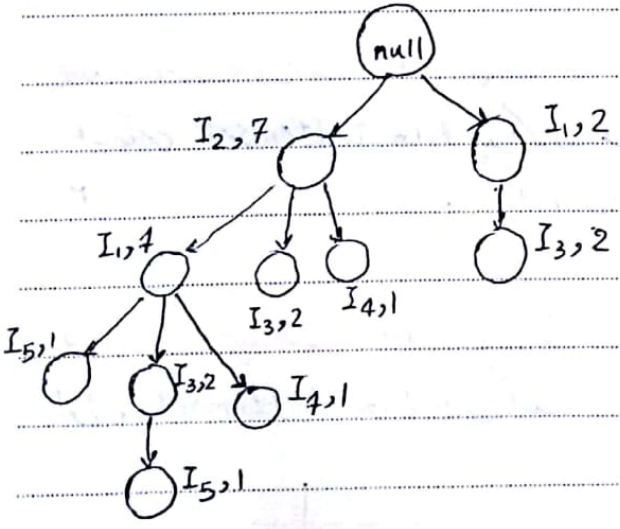


دوره اول

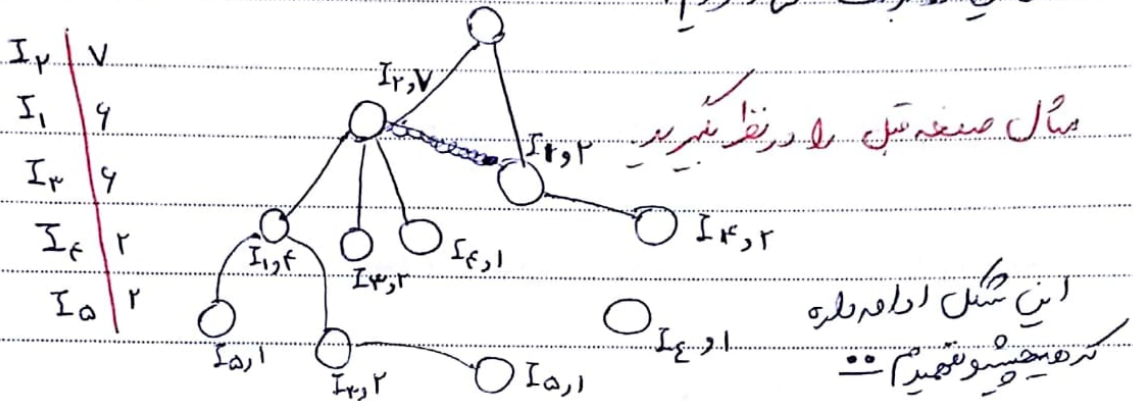
دوره دوم



همین ترتیب بقیه تراش حاصل به درخت اضافه می کنیم



نمونه فصل الگوریتم Apriori این است که سطح به سطح کار می کند
 الگوریتم FP-Growth بر طبق سطح کار کردن، عمق کار می کند و از روی درخت
 قوانین را به دست می آوردیم



از ریزش های درخت شروع به حذف می کنیم
 (I₂, I₁, I₃, I₅) (I₂, I₁, I₅)

I₅ را حذف کرده و به عنوان پیوند آنتو قرار می دهیم بنابراین دو مسیر پیوسته زیر حاصل می شود

به هم های آنتو های شرطی Conditional pattern Base
 1) (I₂, I₁)
 2) (I₂, I₁, I₃)
 الگوریتم Eclat

کاوش آنتو به تکرار بر اساس جدا داده بصورت عمودی
 در الگوریتم Eclat قالب داده ها را عوض می کنیم (افقی به عمودی)
 در حالت عمودی برای هر تراشه مجموع کالاها را بصورت سفارشی در مقابل آن می نویسیم
 به این چنین اطلاعات قالب افقی یا horizontal می نویسیم
 در مقابل آن قالب عمودی یا vertical به صورت زیر است
 تمام آنتو ها را استیکن و پیوسته تراشه تراشه قرار دارد؟

I₁ { T100, T400, T500, T700, T800, T900 }

I₂ { T100, T200, T300, T800, T900 }

I₃ { T300, T500, T700, T800, T900 }

Subject:

Year. Month. Date. ()

I_f { T200, T400 }

I_d { T100, T800 }

: 2- ItemSet

{ I_1, I_2 } { T100, T400, T800, T900 }

{ I_1, I_3 } { T500, T700, T800, T900 }

{ I_1, I_4 } { T400 } حذف می‌کنیم

1- ItemSet برشماره راحت تر می‌باشد.

{ I_1, I_5 } { T100, T800 }

{ I_2, I_3 } { T300, T400, T800, T900 }

{ I_2, I_4 } { T200, T400 }

{ I_2, I_5 } { T100, T800 }

مجموعه‌های مختلف را بدست می‌آوریم

میزان دفعه‌ها که در کاشن می‌دهد

{ I_3, I_4 } { T800 } حذف می‌کنیم

چهار تا که نبودن را می‌نویسیم

: 3- Itemset

{ I_1, I_2, I_3 } { T100, T900 }

{ I_1, I_2, I_4 } { T100, T800 }

آیا تمام الگوها برشماره مورد علاقه باشد؟

Subject:

Year. Month. Date. ()

۱۰,۰۰۰ تراشع خرید
 ۴۰۰۰ خرید بازی کامپوٹری
 ۴۰۰۰ خرید بازی کامپوٹری و فلم ویدیوی
 ۷۵۰۰ خرید فلم طای ویدیوی

min support = 30 Confidence = 40

روپا راندر صورت منفی مع دلیل شده اند
 (فلم ویدیو x) خرید → (بازی کامپوٹری و تراشع x) خرید

Support = 40 Confidence = 74

قانون Strong (قوی) باید افرینیم و به مدیر فروشگاه رسیدیم
 قانون گراه کننده است زیرا اگر یک نفر ویدیوی بخرد ۷۵٪ احتمال دارد

ویدیوی → کامپوٹری 74

احتمال خرید هم شده

بسی حرفه‌ای قوی منبر می‌شود به این که یک قانون جذاب پیدا کردیم
 برای پیدا کردن قوانین جذاب باستند علاوه بر Support و Confidence معیار دیگر
 نسبت را هم در نظر گرفت که مسائل همبستگی بین دو صفت را باید
 از چند روش می‌توانیم همبستگی را بدست آوریم:

۱ Lift با بازی ۲ توزیع (کای دو) ۳ all Confidence
 ۴ ماکسیم - اطمینان ۵ کوئر منسکی ۶ کسینوسی

مثال: ۱۰۰۰ دانشجو ۴۰۰ دانشجو (گوزر است) S
 ۷۰۰ دانشجو (دوچرخه سواری) B ۴۲۰ دانشجو (استا و دوچرخه سواری)
 SNB

$$P(SNB) = \frac{420}{1000} = 0.42$$

$$P(S) = \frac{400}{1000} = 0.4$$

$$P(S) \times P(B) = 0.4 \times 0.7 = 0.28$$

$$P(B) = 0.7$$

Subject:

Year. Month. Date. ()

$P(S \cap B) = P(S) \times P(B)$ مستقل

$P(S \cap B) > P(S) \times P(B)$ ارتباط مثبت

$P(S \cap B) < P(S) \times P(B)$ ارتباط منفی

$lift(A, B) = \frac{P(A \cup B)}{P(A) \times P(B)}$

$lift = 1 \Rightarrow$ استقلال مستقل

$lift < 1 \Rightarrow$ اعداد A بصورت منفی با رخداد B در ارتباط است

$lift > 1 \Rightarrow$ A و B بصورت مثبت با یکدیگر هم سبب هستند

پیدا کردن عمل صافی که طوره دلایله پیدا کند و درونی آن و عملیات انجام دهم:

1- دسته بندی دینا Classification

2- پیش بینی predict

مثلاً در پزشکی ما دسته داریم A, B, C یا در طبقه بندی کلمات می باشد

که وقتی فرد دیگری مبتلا شد بدانیم در کدام دسته قرار گرفته است

مثلاً در بانک از روی سیاه یا سپید بدانیم که آیا باید به آن فرد وام بدهیم یا نه؟ risk/safe

تیب درمانی و سلام و دینار تجارت مثال دری هم هستند.

مشتری جدید یا این ویژگی ها (خصوصیات) حقیقتاً احتمال دارد که کامپیوتر بپذیرد یا نه؟ yes/no

1- Classification

روا فاز داریم: 1- فاز یادگیری / آموزش / training / learning

2- فاز طبقه بندی

Subject:

Year. Month. Date. ()

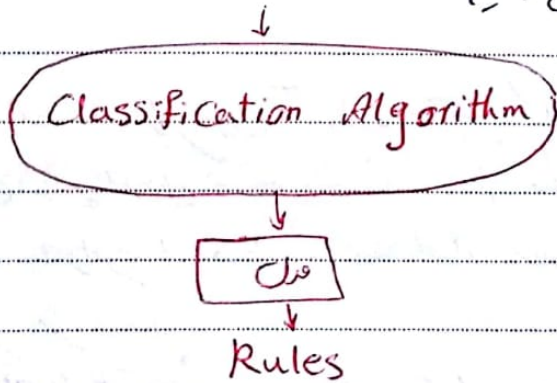
باید روی درست کنیم که دسته یا طبقه بندی کند classifier

باید به فاز اول برسیم که کار طبقه بندی را انجام دهد بعد در فاز دوم سنت می کنیم که درست کار می کند یا نه؟
و در کدام دسته قرار می گیرد؟

فاز اول: Training Data Set

name	Age	income	loan decision
Sandy	youth	low	risky
Bill	"	low	risky
Cavolina	middle Aged	high	safe
Rick	"	low	risky
Susan	se	low	safe
Clarie	"	med	safe
Joe	middle	high	safe
Ali	youth	med	risky
Reza	"	low	safe

مجموعه داده را به یک الگوریتم می دهیم، این الگوریتم یک مدل در دسترس آورد در این مدل یک سری قوانین استخراج می کنیم



قوانین: $\text{if age} = \text{youth then loan_decision} = \text{risky}$
 $\text{if incom} = \text{high} \sim \sim \sim = \text{Safe}$
 $\text{if age} = \text{middle-aged} \text{ ~~then~~ and incom} = \text{low}$
 the loan decision = risky

این قوانین نباید در مجموع رعایت می‌شود مثال نقض داشته باشند
 این رضایت قانون را نقض می‌کند. یعنی قانون درست نیست در مورد اون
 پس این قانون ها باید در حدی اصلاحی درست بودن را دارند \leftarrow **دست**

دو نوع یادگیری داریم: ۱- با نظارت *Supervised learning*

۲- بی نظارت *non supervised*

۱- اگر فیلد که نظارت می‌کند در مجموع باشد یادگیری با نظارت می‌شود
 دنیا را به چندین دسته تبدیل می‌کنیم.

۲- فیلد که نظارت می‌کند نداریم باید *base* ویژگی‌های مشترک را تعیین کنیم
 بر اساس ویژگی‌های مشترک دسته بندی می‌کنیم هر کدام که جدید می‌کند شنیده به کدام دسته است
 در کدام دسته قرار می‌گیرد. بحث خوشه بندی

یادگیری با نظارت: داده‌های آموزشی که شامل مساعدات، اندازه گیری و ... همراه با
 برچسب هستند که کلاس آن مساعدات را نشان می‌دهد داده‌های جدید بر اساس
 مجموعه آموزشی دسته بندی می‌شود.

یادگیری بدون نظارت: برابر تاپل که داده برچسب کلاس نداشته‌اند اغلب
 تعداد و نوع کلاس نیز از قبیل مشخص نیست؛ طور مثال وضعیت وام گرفته
 مشخص نباشد بنابراین از لحاظ شباهت افراد می‌توان ریسک دانش وام را
 بررسی کرد.

هر سطر *tuple* نمونه می‌شود *instance / sample*

ممکن است؛ جایی قوانین یک درجهت با هم برخورد

Subject:

Year. Month. Date. ()

اینکه صدقه به ازای هر تایل یک قانون استخراج کنیم اصلاً خوب نیست

فاز اول: ارتباط بین داده و کلاس را کما حق اولیات باید سری قانون درخت تصمیم و باید فرمولی را بتوان نشان داد.

این قاعده ها چه منظوری مورد استفاده قرار می گیرند؟

- ۱- قاعده ها قوانین پیدا شده و برای دسته بندی داده در حیدر می تواند کار را کند
- ۲- دید واضح تری نسبت به داده در موجود می دهند
- ۳- نمایش داده را فشرده سازی می کنند

بعد از ایجاد Classifier بسته یک مجموعه داده جدید بنام test انتخاب کرده و جهت test و ارزیابی به مدل بدهیم

accuracy دقت: از معاینه نتایج مدل و نتایج واقعی مجموعه داده test بدست می آید.
تایم داده ها را به train و test

درصدی از مجموعه test که classifier درست طبقه بندی می کند
نابراین زمانی که صفت Classifier تأیید می شود می تواند برای تشخیص به چسب داده های جدید بدون label استفاده نمود.

چگونه مدل را بسازیم؟

Decision Tree

درخت تصمیم گیری

ساختار فلوچارتی دارد هرگز به جز بزرگ ها کاریک شرط را انجام می دهد

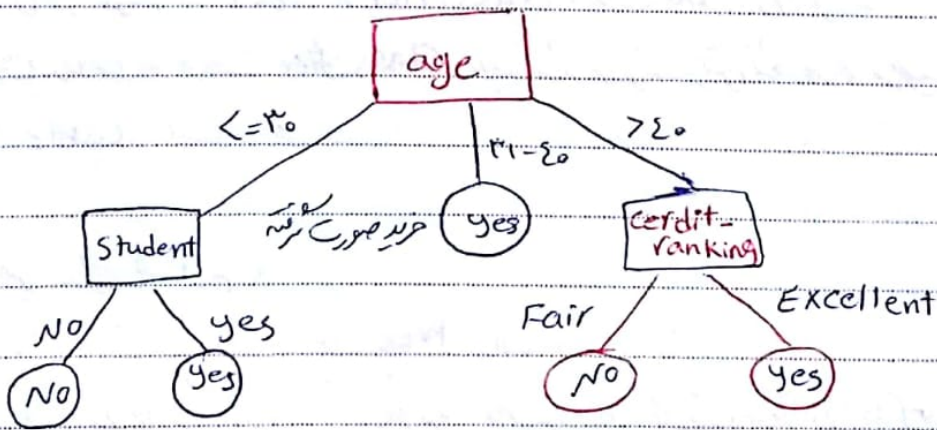
به عبارتی به روی یک از ویژگی ها test را انجام می دهد

هر ساند یا این پس از شرط نتیجه شرط می باشد و هر بزرگ به چسب یا label کلاس است.

Subject:

Year. Month. Date. ()

age	income	Student	credit-ranking	buy-com
≤ 30	hi	no	Fair	no
≤ 30	hi	no	Excellent	no
31-50	hi	no	F	yes
> 50	med	no	F	yes
> 50	low	yes	F	yes
> 50	low	yes	Ex	No
31-50	low	yes	Ex	yes
≤ 30	med	no	F	No
≤ 30	low	yes	F	yes
> 50	med	yes	F	yes
≤ 30	med	yes	Ex	No
31-50	med	no	Ex	yes
31-50	hi	yes	F	yes
> 50	med	no	Ex	No



دہری ایجنٹ بننے پر فیصلہ کرنے کا سوچنا ہے۔ اگر عمر 30 سے کم ہے تو اس کے لیے طلبہ کی حیثیت سے فیصلہ کیا جائے گا۔ اگر عمر 31 سے 50 کے درمیان ہے تو اس کے لیے طلبہ کی حیثیت سے فیصلہ کیا جائے گا۔ اگر عمر 50 سے زیادہ ہے تو اس کے لیے طلبہ کی حیثیت سے فیصلہ کیا جائے گا۔

اگر یک جابری پیدا کنیم به همه افراد جز یک دسته ، کلاس با بکند در آن جا توقف داریم

ضرایب و ویژگی‌های درخت: ۱- نیاز به دانستن و تخصص ندارد ۲- هر تواننده که با تعداد زیاد
پوشش دهد ۳- درک آن آسان است ۴- آموزش دسته بندی ساده است

۵- وقت خوب ۶- کاربردش
مسئله ۱: تعیین ویژگی‌ها بر اساس نیاز درخت

۲- زمانی که درخت ساخته شده ، بعضی از شاخه‌ها طوری تولید می‌شوند ، طوری که با بکند

۳- نمونه این شاخه‌ها را استن (دستم و حفری کنیم)

الگوریتم درخت تقسیم: یک روش خردمند بدون بازگشت است ، عقب به صورت بالا به پایین
و تقسیم و حل است.

نحوه ساخت: در آغاز همه ریشه و نمونه که الگوریتم را در ریشه قرار می‌دهد و در مرحله ۲

یک ویژگی را انتخاب می‌کنیم و بر اساس آن سمت را انجام می‌دهیم اگر ویژگی چند مقدار
داشته باشد ، چند شاخه داریم.

مرحله ۳: داده‌ها بر اساس ویژگی داده شده پارتیشن بندی می‌شوند (کل داده که الگوریتم
بر شاخه انتقال می‌دهد)

مرحله ۲ و ۳ به صورت بازگشت انجام می‌دهیم تا به شرط توقف برسیم

شرایط توقف: ۱- نمونه‌ها برای یک گروه داده شده متعلق به یک کلاس هستند

۲- هیچ ویژگی‌ای برای پارتیشن بندی بقیه باقی مانده باشد.

۳- هیچ نمونه‌ی دیگری نداشته باشیم

همچنین مسئله در درخت تقسیم این است که ویژگی‌ها را چگونه انتخاب کنیم

در هر سطح چه ویژگی‌ها را انتخاب کنیم

برای این که کدام ویژگی را انتخاب کنیم ۳ تا الگوریتم داریم.

برای انتخاب بهترین ویژگی در هر سطح: ID3 - C45 - CART

Subject:

Year. Month. Date. ()

- 1 فرض کنید بهترین ویژگی تقسیم برای تریه بعدی A باشد
- 2 A را به عنوان ویژگی تقسیم برای تریه قرار ده
- 3 برای هر مقدار از A یک فرزند جدید ایجاد کن
- 4 مرتب سازی نمونه های آموزشی برای هر تریه با توجه به مقدار ویژگی ساخته
- 5 اگر همی نمونه های آموزشی کلمات طبقه بندی کنند (حال مقدار ویژگی حرف) متوقف شو، در غیر این صورت برای تریه های جدید تکرار کن

در هر الگوریتم یک معیار داریم که معض می کند در هر سطح کدام ویژگی انتخاب شود

مثلاً در ID3 معیار Information gain

Gain Ratio C4.5

Gini CART

Age: *حالی که خواصم برابر دنیا می بینم در این جهت تقسیم را بازنویس*

$$\text{Information gain}_A = \text{Info}(D) - \text{Info}_A(D)$$

$\text{Info}(D)$ ریس نمونه A

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2 p_i$$

Info Info *بله* *بله* *چندتا* *کلاس* *داریم* *الان* *دو تا* *کلاس* *داریم* *No, yes*

$$p(\text{yes}) = \frac{9}{14} \quad p(\text{No}) = \frac{5}{14}$$

احتمال هر کدام چقدر هست؟

$$\text{Info}(D) = - \frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$\text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Info}(D_j) \quad A = \text{Age}$$

چند دسته Age داریم؟ ۳ تا ۳ بار فصول بازنویس

برای هر کدام yes و No ما حساب می کنیم

$$\text{Info}_A(D) = \frac{5}{14} \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{9}{14} \left(-\frac{4}{9} \log_2 \frac{4}{9} - \frac{1}{9} \log_2 \frac{1}{9} \right)$$

$$+ \frac{5}{14} \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right)$$

H4MKELASI

انجمن علمی علوم کامپیوتر
دانشگاه کاشان

t.me/KUCSSA

Subject:

Year. Month. Date. ()

$$Info(A) = Info(D) - Info_A(D) = 0.264$$

$$Info_{(income)} = 0.029 \quad Info(Stu) = 0.151 \quad Info(Cr) = 0.048$$

Info(Age) از حد بیشتر است پس اولین ویژگی در درخت تقسیم Age است.

مقدار اطلاعاتی که نیاز هست برای اینکه تقسیم بندی صورت بگیرد
عبارت این روش: اکثر این دنیا بین ID راسته باشد و اول کار روی ID ساختن
مکنند و مقدار ساختن خیلی زیادتر شود.

اگر مقدار یک فیلد پیوسته باشد مثلاً Age پیوسته باشد و هر کس سن واقعی اش داشته باشد.
در این صورت باید Split Point پیدا کنیم **نقطه تقسیم**

معایب Information gain برابر مقدار پیوسته:

ماست نقطه Split point برای ویژگی پیوسته A تعیین کنیم بنابراین مقدار A را به صورت
صغوری مرتب کرده به طور معمول هر توان نقطه وسط بین هر جهت از مقدار پیوسته
عنوان نقطه تقسیم در نظر بگیریم.

$$\left(\frac{a_i + a_{i+1}}{2} \right) \text{ و این نقطه ای با کمترین expected Infor. requirement}$$

پیدا کردن این نقاط کار راحتی نیست. مثلاً هر بار باید با هم بین دو تا عدد یکی را به عنوان نقطه تقسیم
در نظر بگیریم. بهترین کار این است که هر زمان نسبت به داده ها ساخت داشته باشیم.

Split Info

$$A = income \Rightarrow -\frac{4}{14} \log \frac{4}{14} - \frac{6}{14} \log \frac{6}{14} - \frac{4}{14} \log \frac{4}{14} = 1.557$$

$$\frac{0.029}{1.557} = \text{Gain Ratio}$$

Subject:

Year. Month. Date. ()

هر چه Gain Ratio بیشتر باشد، این ویژگی را انتخاب می‌کنیم.

Gini : هر چه درخت را به دو قسمت تقسیم می‌کنند و برای هر کدام Gini را بدست می‌آورند

اگر مجموعه A، V مقدار داشته باشد 2^V مقدار می‌توان از آن داشت:

$$\text{income} = \{ \text{low}, \text{med}, \text{high} \}$$

$$\{ \text{low}, \text{med} \} \quad \{ \text{high} \} \rightarrow \text{این مجموعه را به دو قسمت تقسیم می‌کنیم}$$

$$\{ \text{low}, \text{high} \} \quad \{ \text{med} \} \rightarrow \text{gini} = 0.458$$

$$\{ \text{med}, \text{high} \} \quad \{ \text{low} \} \rightarrow \text{gini} = 0.450$$

با مجموعه‌ها باید برابر شوند
 $\left. \begin{array}{l} \{ \text{low} \} \\ \{ \text{med} \} \\ \{ \text{high} \} \end{array} \right\}$

چون می‌خواهیم به ۲ دسته تقسیم شود

$$\text{gini}(D) = 1 - \sum_{j=1}^n P_j^2$$

$$\text{gini}_A(D) = \frac{|D_1|}{|D|} \text{gini}(D_1) + \frac{|D_2|}{|D|} \text{gini}(D_2)$$

$$\text{gini}(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

$$\text{gini}_A(D) = \frac{10}{14} \left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right) = 0.443$$

$$0.459 - 0.126 = 0.333$$

هر چه gini عدد باشد آن تقسیم نیمی بهتر است.

حجرتين دستيابي سن: $gini = 0.357$ → {youth, Senior & middle}

St → 0.367 C-R → 0.429

Age - gini ازجه كندارت بين Age برابر دستيابي استقاره كنيم عيب اين روش محاسبات زياد است.

اذا باوجود همين اينها درخت تقسيم گيري روش خوبي است.

مشكلات درخت تقسيم: زياد شدن بغير سادگي؛ كه بايد همين كنيم.

مكن است برابر مدل داده آموزشي درخت تقسيم درست كارند اما اين محاسبه داده تست درست كارند يعني درخت overfitting باشد كه ناسازگار است.

1- تعداد ساختن درخت تقسيم زيادتر باشد 2- درخت تقسيم برابر محاسبه داده تست باخوبي عمل نمي كند

راه حل: هرس كردن pruning

1- prepruning در هنگام ساخت درخت اگر نياز به هرس كردن كنيم

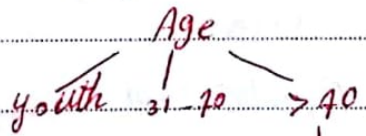
2- post pruning بعد از اينكه درخت ساخته شد آن را هرس مي كند كه اين روش خوبي است.

در صورت تقسيم بگيريم كه اگر تقسيم مناسب نباشد اراده ندهيم.

در واقع در صورت كه نتيجتي داده نشود، تقسيم بگيريم ساختن درخت نخواهد

داشت آن ساخته هرس شده و دست تقسيم انجام نخواهد شد. با سيزه آن گره را

تبديل به برگ نموده و كلاس آن ساخته را بر اساس بيشترين تعداد اعضاي آن محاسبه كرد.



Credit -



استفاده تقسيم بهتر مي شود.

اين را سبب برگ مي كنيم

دسته به اينكه تعداد Yes ها و No ها كمتر باشد بگيريم گره پاهان مي زنيم

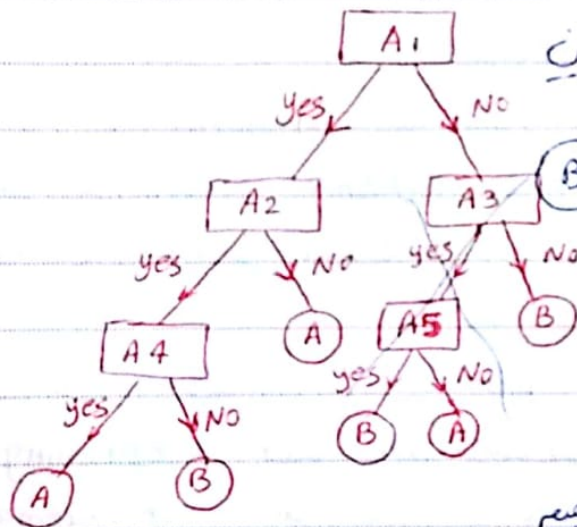
Subject:

Year. Month. Date. ()

Credit_rang = Fair NO
income = high yes

حالا یک مثال تعریف می‌کنیم مثلاً Fair high NO
تقسیم بندی به نفعان نیست

معیار تقسیم بندی: انتخاب یک مقدار استاندارد بر مبنای information gain



هرس کردن برای این درخت شروع می‌شود و بر اساس معیار هرس شدن یا نبودن یک ساختار تعیین می‌شود.

یک مجموعه داده است به آن می‌گویند
بنیم روی چند داده از داده‌ها خط می‌کشیم که
داده خط

مثلاً میزان خط از ۱۰ درصد بیشتر است یا مثلاً ۹۰ درصد خط دارد این ساختار
هرس می‌کنیم

به این مجموعه داده که بر مبنای تقسیم داده خط استفاده می‌شود مجموعه pruning set می‌گویند

مشکل درخت تقسیم: Replication! زیر درخت تکاری

یک زیر درختی توی درخت عیناً داده می‌گردد
repetition و تکرار تکاری همین بار روی یک درخت شرط نداریم
درختی‌ها تکاری در عمق‌های مختلف
بر مبنای این مشکلات از هرس کردن میزان استفاده کرد

HAKKELASI

انجمن علمی علوم کامپیوتر
دانشگاه کاشان

t.me/KUCSSA

Subject:

Year. Month. Date. ()

برای حل این مشکلات از multi-attribute و multi-variant استفاده می‌کنیم
۲ صورت ترکیبی استفاده می‌کنند. شرط‌ها را روی هر دو قرار می‌دهند

مشکل بعدی: اگر داده‌ها زیاد باشند در حافظه اصلی جا نمی‌گنجد
اصطلاحاً مقیاس پذیر نیست (Scalable)

راه حل ←

۱- **Rain Forest** یک سری نسبت‌ها در دسترس

Attribute value class table

این به جای اینکه کل اطلاعات را بیاریم، روی حافظه اصلی یک سری نسبت‌ها می‌سازیم و آن‌ها را بیاریم

ArcL - list on Age

BOAT - ۲

Age	Buy - comp	
	yes	No
<30	2	3
30 < <40	4	0
>40	3	2

از معیار آسانی

Boot strapping

سرعت در دسترس

روشن دست‌نزدی در آیه

روشن‌ترین: بین آن‌ها قانون‌ها احتمال و قانون‌ها پیوسته است
مجموعه داده قبلی را داریم. اگر یک تایل داشته باشیم، بتوانیم احتمال تعلق
آن‌ها به یک کلاس خاص را تعیین کنیم.

ساده‌ترین **naïve Bayesian classifier** ← Independent

فرض اولیه: ویژگی‌ها از یکدیگر مستقل هستند و هیچ تایل‌هایی بر روی یکدیگر ندارند
رغم‌ترین نسبت این روش‌ها ← سرعت بالا (چون باید فقط متغیرها را بشماریم)
و معمولاً وقت بالایی ندارد

اگر X وجود دارد یا رخ داده احتمال پیدا کنیم کامپیوتر بخرد

$P(H|X)$ **evident** X : تایل

posterior probability H : فرضیه (مثلاً کامپیوتر بخرد)

H4MKELASI

Subject:

Year. Month. Date. ()

likelihood

$$P(H|X) = \frac{P(X|H) P(H)}{P(X)}$$

prior probability ← P(X) , P(H)

مثال: دکترا مراد که منتزیت در ۵۰ درصد مواقع باعث خوشگرددن می شود
 H: فرضیه: خوشگرددن
 X: جوکر ناسی از منتزیت.

$$P(H|X) = \frac{5\% \times \frac{1}{50000}}{\frac{1}{10}}$$

تایل X یک ویژگی نیست
 عندئذ X باید ←

Age < 30

incom = low , Fair , ...

حالات احتمال ایند کامپیوتری فرزند است؟

$$P(C|A_1 A_2 \dots A_n) = \frac{P(A_1 A_2 \dots A_n | C) P(C)}{P(A_1 A_2 \dots A_n)}$$

هر یک از حالات در P(C) ضرب می شود

$$= \frac{[P(A_1|C) P(A_2|C) \dots P(A_n|C)] P(C)}{P(A_1 A_2 \dots A_n)}$$

max کنیم صورت و در نهایت در نظر می گیریم ←

روتا کلاس در نظر می گیریم ←

C₁ = yes C₂ = No

X = (age <= 30 , income = med , student = yes , credit = Fair)

هر یک از حالات را حساب می کنیم

$$P(C_1) = \frac{9}{14} \quad P(C_2) = \frac{5}{14}$$

$$P(\text{age} <= 30 | C_1) = \frac{2}{9} = 0.222$$

$$P(\text{income} = \text{med} | C_1) = \frac{4}{9} = 0.444$$

$$P(\text{Student} \neq \text{yes} | C_1) = \frac{6}{9} = 0.667$$

$0.047 \times P(C_1)$ ضرب این

$$P(\text{cr} = F | C_1) = \frac{6}{9} = 0.667$$

تک تک این در هر صفا، احتمال را برابر

کلاس C₁ حساب کردم. حالا هر ریم سرای کلاس C₁

$$P(\text{age} \leq 30 | C_2) = \frac{3}{5} = 0.6$$

$$P(\text{income} = \text{med} | C_2) = \frac{2}{5} = 0.4$$

حالا بعد این احتمال ها را در هم ضرب کن

$$P(\text{stu} = \text{yes} | C_2) = \frac{1}{5} = 0.2$$

و در P(C) ضرب کن 0.019 ضرب

$$P(\text{cr} = F | C_2) = \frac{2}{5} = 0.4$$

این $\times P(C_2) = 0.007$

از C₁ بهتر است پس این تا این عضو کلاس C₁ یعنی yes است

عیب این روش: اگر یکی از احتمال ها صفر باشد وقتی در بقیه ضرب شود، حاصل ضرب صفر می شود. و کل احتمال را صفر می کند و باعث می شود تا این عضو یکی دیگر از کلاس ها شود. چون فرض بر این است که ویژگی ها مستقل هستند بنابراین تا این متعلق به آن کلاس نیست.

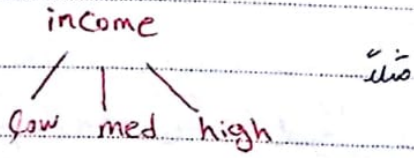
البته حل این مشکل **Dataset** ضعیف تر است مثلاً همان احتمال

Laplace Correction در این روش برای صفر شدن احتمال، یک تا این به صورت

صفر وارد می شود. چون مجموع داده بسیار بزرگ است اضافه کردن یک تا این تأثیری نخواهد داشت.

$$\text{original: } P(A|C) = \frac{N_{ic}}{N_c}$$

$$\text{Laplace correction} = \frac{N_{ic} + 1}{N_c + 1}$$



مثلاً $\text{income} = \text{low}$ کامیوتر خریده و احتمال آن صفر است

بنابراین برای اینکه عدالت رعایت شود به هر کلاس یک تا این اضافه می کنیم پس صورت 1+ می شود

و مخرج 3+ می شود $\frac{+1}{+3}$ برای وقتی که مجموع دیتا کم باشد

و نیاز به وقت بیشتری دارم

Subject:

Year. Month. Date. ()

مقایسه نسبت به درخت تصمیم راحت تره (موضوع نرمه) که باید سه سرعت هم زیاد شده

۱- پیاده سازی آسان

۲- سرعت بالا

۳- نسبت به داده نوزی رپت حساس است

۴- داده Miss: (بعضی از اطلاعات قابل از دست رفتن نیستند)

مغایب: ۱- فرض مستقل بودن باید کاهش دقت خواهد شد

۲- در عمل میان متغیرها وابستگی وجود دارد

Bayesian Belief Network

روش دست بندی: ۱- درخت تصمیم ۲- روش بیزین ۳- روش مبتنی بر قاعده

Rule-Based: قوانین را استخراج میکنیم

if (age = youth, 88 student = yes) then buy-comp = yes

حالا می توانیم برای مثالی که جدید می آید بگویم جزو کدام کلاس است

(age = youth) 88 (student = yes) => buy-comp = yes

استفاده از قوانین - تولید قوانین (تایم) مقدم (شرط)

Condition

پوشش Cover: اگر شرطی که باید تایید در یک قانون صدق کند در یک Rule

تایید یا Cover می کند Rule cover Tuple

روش ارزیابی: این می توانیم مقدار تایید را زیاد پیدا کنیم که این قانون آن در پوشش دهد

$$Coverage = \frac{n_{cover}}{D}$$

به شرطی که درست باشد Cover کند

(پوشش)

$$\frac{n_{correct}}{n_{cover}}$$

(دقت)

شرطی که ارزیابی یک قانون

حالاتی که برای خرید X می آید. میخواهم مشخص کنم که کدام کلاس است؟
اگر شرط در آن در جدول قانون صدق نکند؟

$$X = (age < 30, income = med, stu = yes, cre = f) = ?$$

$$if (age = youth \ \& \ cr = F)$$

IF (age=youth && Studen=yes) then buy-comp=yes

مخصوصاً برای این قوانین مشخص کنیم که X صفتی به کدام دسته است؟

۱- اگر R₁ تنها قانون برگزیده شده باشد، قانون با بالاترین رتبه برای X دسته آن را تعیین می کند

۲- اگر بیش از یک قانون trigger شود نیاز به برطرف کردن تعارض داریم چون هر قانون ممکن است یک دسته مختلف برای X مشخص کند

برای رفع اختلاف در استراتژی وجود دارد

۱- size ordering: بالاترین اولویت را به قانونی می دهد که رتبه بالاتری دارد. شرط را برگزیده کرده باشد یعنی بیشترین تعداد شرط ها را داشته باشد

از بالا به پایین قانون ها را مرتب می کنیم (اول قانونی که قدا ۴ شرط بعدی ۳ و ...)

۲- Rule ordering: قانونی قابل قبول است که کلاس آن شیخ تر است

prevalent

Class based ordering

اگر برای X در هیچ قانونی برگزیده نشود در این حالت باید یک رفتار را تعیین کنیم به احتمال آن ضمیمه کنیم

این نحوه استفاده از قوانین بود

درخت تقسیم

نحوه تولید قوانین: با استفاده از قوانین را استخراج کنیم

حالا چرا در روش rule-base از درخت تقسیم استفاده می کنیم؟

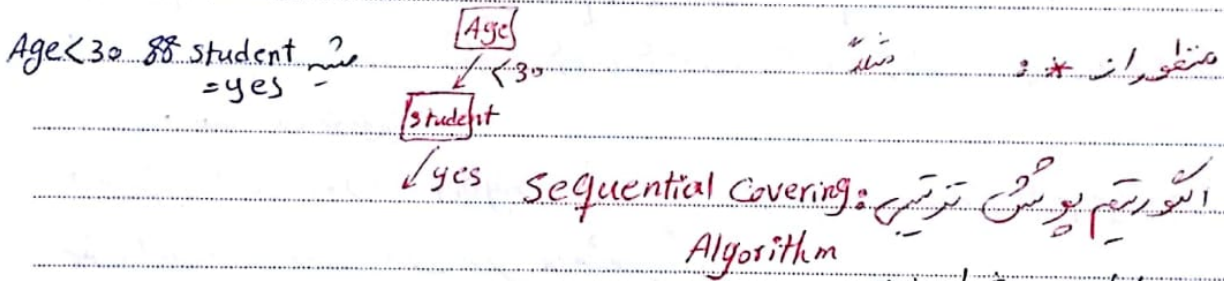
چون ممکن است درخت تقسیم پیچیده شود. ما آن را تبدیل به قانون می کنیم

Subject:

Year. Month. Date. ()

در دخت تقسیم: برای هر مقدار از Age یک قانون ایجاد می شود
* هر زوج مقدار خصوصیت در مقدار Age یک ترکیب مختلف می سازد
بزرگ، بر حسب کلاس است.

عمیق ترین مرتبه: قانون ها در Age ظاهر می شوند *Mutual exclusive*
یعنی یا این قانون یا آن - دو قانون *trigger* نمی شوند برای یک تابع
جامع است و کامل یعنی هر چه خواصیم می توانیم در این قانون جا بدهیم
مثال *sizeorder* و ... می داریم یا قانون ریاضیات.



مرحله اول هیچ شرطی نمی داریم
همه اطلاعات را می بینیم
هر مرحله یک شرط می داریم.
برای مثال وام.

IF — THEN
loan_decision = accept

if income = high then
loan = accept

if income = med then
loan = accept

if income = high & cr = exc
then loan = accept

در سطح 2 کدام قانون را قبول کنیم؟
قانونی که مقدار تابعی بیشتر را با دقت خوب پوشش دهد
از بالا به پایین پوشش کاهش می یابد
اما دقت افزایش می یابد
تا کجا شرط می داریم؟
ما اینکه از تابعی هایی که داریم هیچ معنایی نباشند.
بین دقت و پوشش به یک تعادلی برسند
در هر مرحله یک سری شرط به آن اضافه می کنیم
قوانین را به صورت مستقیم از داده های آموزشی استخراج می کنیم

Subject:

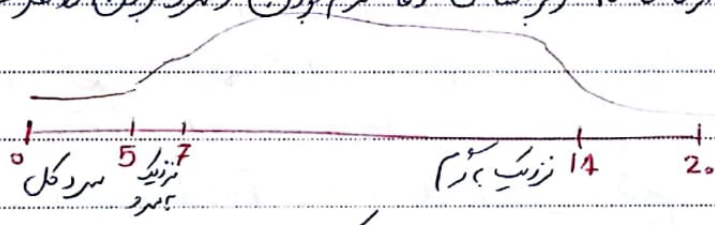
Year. Month. Date. ()

کلمه: التورتم

- ۱- از مجربه داده آموزش مرحله به مرحله قانون ایجادند
- ۲- هر قانونی که موضوعه همسود تالی علیه که توسط آن قانون موضوعه فریبند حذف فریبند
- این کارها تا زمان شرط یا بیان ادامه دارند.
- شرط جاری یا بیان: تا جاری که هیچ مثال آموزش دیگری باقی مانده باشد
- تا جاری لغت قانون ایجاد شده کمتر از ۵۰ صفتن نباشد.

پوس و رفت

منطقه غازی: مثلاً می گویند هوا سرد یا سرد که این رسم یا سرد بودن منتهی است یک معیاس دما می داریم از ۰ تا ۲۰ در اساس دما گرم بودن و سرد بودن را تعریف می کنیم.



سب منطقه فقط صفر و یک نیست و می توان درجه بندی کرد.

مثلاً: در یک دستورالعمل می خواهند یک فرآیند بدهند. تحت تأثیر دو عامل سردی و لغت غذا بر اساس ترکیب این دو مقدار انجام لامنه می کنیم. این ها قوانین غازی می گویند.

Service	Food	Tip
poor	delicious ✓	
avg		
good		

کتاب صعبه امزجاری داریم بر این غازی

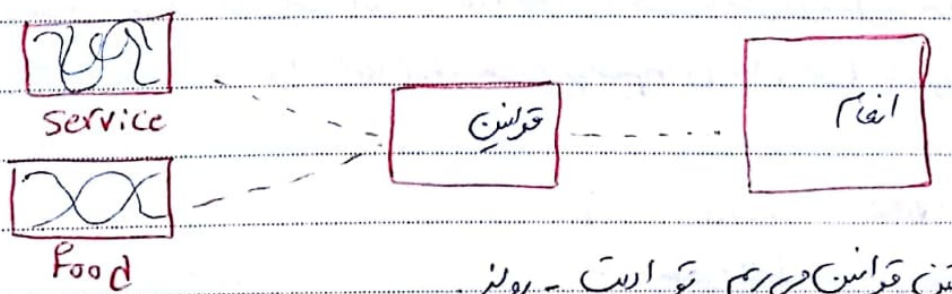
مثلاً: اینجا می توانیم بر سر از دستورالعملمان لایزیم. هر زینم fuzzy تو صفر کاندو ویندو اینجا وارد غازی قول می کنیم

نمبره ۲-۱۵ که همان فازی تول است. یک زیرمجموعه داریم که می توانیم این پورت را کپی کنیم.

الان می خواهیم مسئله انفاک را یاد بگیریم از سر سر است متغیر اضافه نمی کنیم
۲ ورودی و یک خروجی داریم. متغیر اول سرویس. متغیر دوم غذا. خروجی T
در هر کدام باید ریج شریف کنیم و شکل ها را مشخص کنیم.

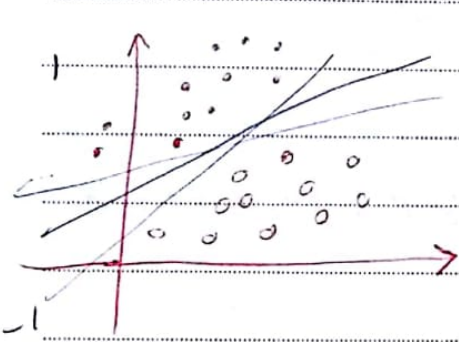
این متغیرها یک سری توابع دارند می ایم بازه ها را تعیین می کنیم سرویس از ۰ تا ۵ Pool of
از ۰ تا ۵ هتویط است.

می توانیم این ریج ها را چاپ کنیم (از ورودی شکل نویی ترا افترار)
اگر این ها را اشتباه کردیم می توانیم روی آن کلیک راست کنیم و Remove
می توانیم مشکل های دیگری هم اضافه کنیم سر سر است از زیر در. که این شکل ها
فرمت های متفاوتی دارند. مثلاً حلالی - ذوزنقه ای
این ریج ها را برای Food و service مشخص می کنیم.
برای انفاک هم باید ریج را مشخص کنیم.
قوانین که فازی روی آن ها کار می کند بیان می کنیم



برای نتن قوانین می رسم تو ایت - روز
ارتباط بین دو عامل را سمت چپ مشخص می کنیم. مثلاً ارتباط or باشد یا And
اینکه قوانین به چه شکلی باشد بستگی به کاربرد.
قبل از مشخص کردن قوانین باید یک سری نیت به یاد داشته باشیم
می توانیم نمودار این قوانین را بنویسیم.
این plat به نشان می دهد که خروجی های ما روی چه حالتی هستند.

مسئله Free paper - لایه های که از دو دسته جدا می شوند با لایه های جدا می کنند
و اینجا سبب می کنیم به این صورت معادله را می بینیم ..



SVM: یعنی از بهترین حالت دسته بندی است

فرض کنید فقط دو دسته داریم n_1 و n_2 داریم
ساده ترین راه دسته بندی این داده ها این است که یک
خط بکشیم تا آنکه بهترین دسته بندی کارشان این است.

روتا کلاس اول کلاس اول و کلاس اول
روتا کلاس اولی جدا کنیم که توی دو تا کلاس بقیه

تعداد می بخشد خط بین این دو کلاس می توانیم داشته باشیم که دسته بندی را درست انجام می دهد
در تمام اینورتم های قبلی محتم نبود که به کدام خط می رسم و فقط رسیدن به یک خط محتم بود.

اما در SVM رسیدن به یک خط محتم است؟ اما خط محتم کدام است؟
یک خط وسط - که حداکثر فاصله را از هر دسته داشته باشد.
فاصله خط از اولین ریاضی این کلاس بیشتر باشد.

در هر کلاس داده هایی که به هم خیلی نزدیک اند برای ما مهم هستند. چون این؟ یعنی
می کنند که خط کجا باشد؟

نیاز است که صفاً تا دین روی این خط اثر ندارند (مداخله) هر دین مان یک کلاس
اما در نقاط بهترین هم پیدا کنیم، مثلاً تا بی نهایت، اشکال ندارد.

هم ریاضی حق ندارد در margin قرار بگیرد که این سخت تر است hard margin
margin یک داشته است.



ما الان صفاً با ساده سازی کردیم و گفتیم تکلیف پذیر است با این خط
اما در واقع اینگونه نیست

soft margin - یک درصد خط توی ریاضی که می کشیم داریم. یعنی یک سری داده هایی
margin هستند

معادله خط $w_0 + w^t x = y^t$ (برای اینکه توی بودا سستی داریم)

این همان label های ما هستند.



Subject:

Year. Month. Date. ()

حالا دوتا فصول ارائه میدیم برای مثالون ۱: $if\ w_n + w_0 > +1$ class +1

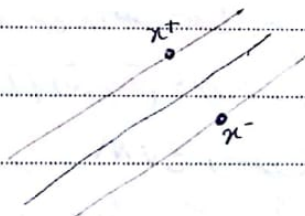
$if\ w_n + w_0 < -1$ class -1

هر چند به margin بزرگترین برسیم. محبت است به دست تر کلاس بندی می کنیم.

این دوتا معادله را تبدیل به یک معادله می کنیم، طرفین آن را در یک y ضرب می کنیم
برای راحتی مساوی هم بر می داریم.

$$\left. \begin{array}{l} w_n + w_0 > 1 \\ w_n + w_0 < -1 \end{array} \right\} \Rightarrow y(w_n + w_0) > 1$$

فرض کنید دوتا نقطه داریم که روی هم اند (مساره ترین حالت) چگونه خط بین آن دو پیدا کنیم که margin آن ماکزیم باشد؟



این مسر خط موازی اند هم به برابر زغال می آید است.
 x^+ و x^- داره هستند.

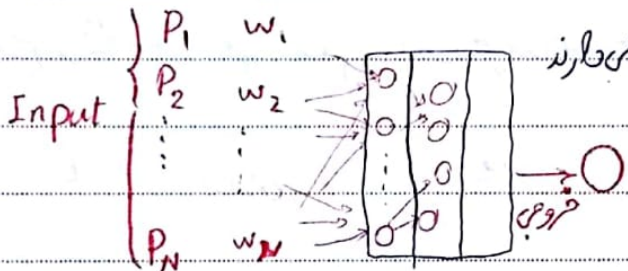
Subject:

Year. Month. Date. ()

شبکه عصبی: فرض کنید میزان آلودگی هوا - همه عوامل را برای یک سال متوالی داریم در کنارش دمای هوا ، تعداد ماشین ها و چند تا پارکتر دیگر هم داریم به میزان: به اسم آن پارکترها حرکت می کنند چند تا ورودی داریم و یک خروجی

یک مدل استخراج کنیم که بتوانیم میزان آلودگی فردا ، یک هفته بعد و ... را کنیم مثل هواشناسی ریاضی یک سال قبل ، دو سال قبل را به عنوان آموزش می دهیم

یک سری ورودی داریم



هر کدام از ورودی ها یک آلودگی در خروجی دارند که این را با وزن شان می دهیم

شبکه عصبی مجموعه ی نورون ها است

چند تا لایه داریم و در هر لایه چندین نورون داریم و در نهایت آن را به خروجی می رسانیم

نورون ها P_i و w_i ها را می گیرند $P_i \cdot w_i + b_i$ ورودی نورون

هر دفعه که ما آموزش می دهیم کاری می کنند که w_i ها و b_i ها را یاد بگیرند

اینقدر به آن آموزش می دهیم که این w و b ها به گونه ای در دست سیستم کنیم

مثلاً بر اساس تعداد پنجره ها و تعداد افراد درون اتاق دمای اتاق را پیش می گویند

سفروا هم دماشبه ۳۰ درجه اول به پنجره باز می کنیم و بیخ نور می بینیم دماشبه ۳۰ درجه

بعد دوباره پنجره را می بندیم اینقدر با تعداد پنجره ها و تعداد افراد بازی می کنیم

تا به خروجی مد نظر برسیم پس تعداد w_i ها و b_i ها تغییر می کنند تا به خروجی

مطلوب برسیم یعنی هر آموزش می دهیم

مدل از روسی دیبای سال ۲۰۱۷ ساخته شده (پس دیبای ۲۰۱۷ و ۲۰۱۸ دیبای ترین)

دیبای ۲۰۱۸ لایحه جدید دیبای ۲۰۱۸ متوجه دیبای است

هر کدام از این نورون ها یک تابع دارند

تو متلب از داخل این لینک ها پیدا کنیم یا تو گاند می بینیم nn_start

برای حل مسائل چندین روش داریم که این را این نوشته

Subject:

Year. Month. Date. ()

ما فعلاً input-output را می بینیم. fitting tool. مشکل ما اینست که ما می خواهیم به این داده ها Dataset را در Dataset خودمون یک سری Dataset باره

یک سری نمونه باره را به عنوان آموزش می داریم. یک سری برای اعتبار سنجی که مدل ما را آزمون خروجی با مقدار مورد نظر می نبود. ما را تغییر میده تا به میزان دست بریم. یک سری برای تست که به مدل می دم برای تست کردن خروجی های واقعی را داریم و نگاه می کنیم ببینیم که مدل چقدر درصد خطا میده

MSE خطای متوسط مربع

رگرسیون: ارتباط بین خروجی واقعی و خروجی بدست آمده

Simple script که در این دستورات را می بینیم

ارائه SVM: سه تا خط موازی داریم. خط وسطی که خط جدا ساز هر کدام از این ها با هم فاصله دارند که این فاصله P می باشد حالا یا کل فاصله و یا فاصله بین خط وسطی و اول یا خط وسطی و P هر کدوم

(x, y) تا یابی ما می باشد. y کلاس آن

$+1$ -1

برای اینکه ببینیم x در کدام دسته است

این $+1$ و -1 که P است آن را با label استباه نمی کنیم

$$\left. \begin{aligned} w^t x + b &\geq +1 \\ w^t x + b &\leq -1 \end{aligned} \right\} \text{در هر کدام قرار می دهیم}$$

حاجت ما اینست که داده ها در بین این دو بازه قرار نگیرند

که این خطی که اتفاق می افتد

سافت مارجین: امثال نذاره که داده ها بین این دو بازه قرار بگیرند

Subject:

Year. Month. Date. ()

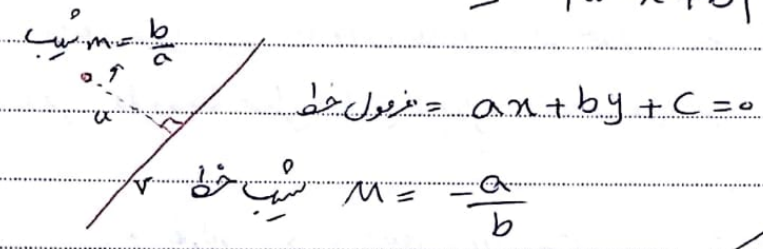
حرفین محدودیت حاصله تی ضرب هر کنیم
 فرد و محدودیت بتبدیل به یکی می شوند.

$$y_i (w^t x + b) \geq +1$$

مهر خواستیم فاصله نقاط تا خط حد اکثر شود :

$$\frac{|w^t x + b|}{\|w\|} \geq \rho$$

$$\Rightarrow |w^t x + b| \geq \rho \|w\|$$



مطلوب نیستیم خودتان حد کنید
 فاصله نقطه و خط = $\frac{|ax + by + c|}{\sqrt{a^2 + b^2}}$

ρ دارد حاصله بود برابر است فاصله بیشتر شود ρ را فاکتور هم می کنیم
 فرض می کنیم

$$\rho \|w\| = 1 \Rightarrow \rho = \frac{1}{\|w\|}$$

یعنی $\|w\|$ باید \min شود

معادله بحیثی سازی $\rightarrow \min \frac{1}{2} \|w\|^2$

s.t $y_i (w^t x + b) \geq 1$
 برابر این از فرمول لاگرانژ استفاده کنیم

وقتی مواضع $\|w\|$ را \min کنیم هر توانیم هر ضریب از آن را \min کنیم

این را به مسئله اضافه می کنیم که جهت نسبت به خطها حاصل می باشد
 خطای هر کدام از آن ها را با ϵ در نظر می گیرند - عنوان خطها را با α
 خطوط اطراف در نظر می گیرند

مجموع خطها را \min می سازیم
 سبب مواضع خطها \min سازی شود برابر \min است نرم ρ را \max کنیم
 $\min \|w\|$ کنیم

Subject:

Year. Month. Date. ()

C سے تیزان خطا کر۔ در نظر میں لیں۔ اگر جو اہم نکتہ لکھ کر لکھیں سے فارغ ہوں
 C کا صفحہ در نظر میں لیں۔ ہر جگہ جو اہم نکتہ لکھ کر لکھیں سے فارغ ہوں
 در نظر میں لیں۔ تاہم صورت تحریر بہ نسبت فراموش

دارہ کے مرکز سے باہر رہیں تاہم ان خط منحنی ہوگا۔
 بہر حال خط سے تو ان دارہ کا اہم جہاں سے خطی جہاں پڑیں
 اہم بعض دارہ کا خطی جہاں پڑیں۔ باقی جہاں تو ان جہاں جہاں



روشن کرتے: تعداد اعداد انتہا باہر ہیں
 کہ خطی جہاں پڑیں ہوں۔
 تو ان باہر صفحہ آگن کا جہاں
 دو ورتوں (x و x) دارہم برابر ہیں این دو ورتوں کو لہذا ہم سازیم

$$\phi : x \rightarrow \phi(x)$$

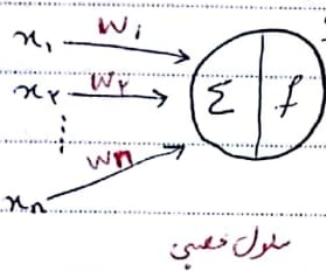
ان ورتوں کو لہذا ہم سازیم
 ان خطی جہاں پڑیں ہوں
 ان سارہ ترین حالت ممکن ہو با اضافہ کریں کہ بعد جواب رسیدیم
 اہم بعض وقتا باید خطی بعد اضافہ لکھیں

Subject:

Year. Month. Date. ()

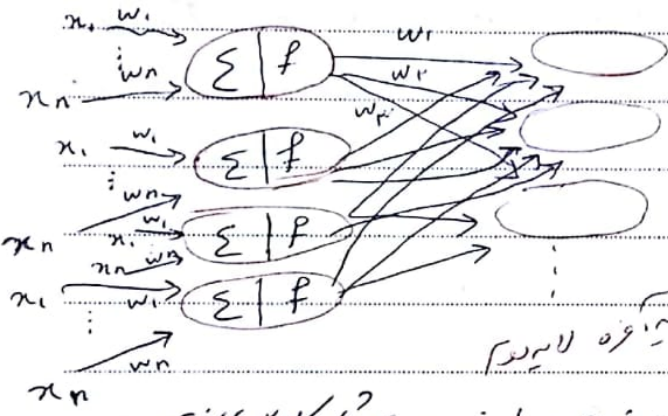
$$w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

شبکه عصبی



مجموعه از سلول‌های عصبی که هر سلول شبیه عصبی
 بر اساس ورودی و محاسبات شبکه عصبی را لایه به لایه در نظر
 می‌گیریم. هدف از این شبکه عصبی یادگیری است.
 این که در هر لایه چندتا سلول داشته باشیم مهم ترین است.

سلول عصبی



w_i ها کلا مهم متفاوتند.

$$f(x) = \frac{1}{1 + e^{-x}}$$

شبکه این لایه به لایه داریم

لایه آخر مثل ۴ تا کلاس داریم
 ۴ تا خروجی در نظر می‌گیریم. هر کدام یک خروجی دارند. یک کلاس بندی
 خروجی بر اساس تعداد کلاس‌های که داریم.

تعداد خروجی و ورودی بر اساس تعداد ورودی تعیین می‌شود.
 بزرگ‌ترین بین این به تعداد خروجی‌های که می‌خواهیم در لایه آخر خروجی‌ها داریم.
 فرض کنید x_1, \dots, x_n همگی صفر باشند در این صورت چه بدام از w ها و b ها
 تأثیر ندارد و خروجی نیز صفر می‌شود.

ورودی بین صفر و یک است. اگر بین صفر و یک نباشد زغالانز می‌کنیم.
 اینجا هم یک بایاس اضافه می‌کنیم تا w که دارد که آن هم تنظیم می‌شود.
 b بایاس مقدار یک ضرب w است. بایاس یک مقدار بین ۰ و ۱ است.

داره های آموزش لایه این شبکه می‌دهیم. مقدار اول این w ها را از رندوم در نظر
 می‌گیریم. خروجی یک عدد بین ۰ و ۱ است. (لایه ۱ و ۲ خنثی = ۰)
 در این ترین یک خروجی می‌دهد.

یک خطی داریم. جواب اصلی باید یک باشد اما خروجی ما لایه است.

Subject:

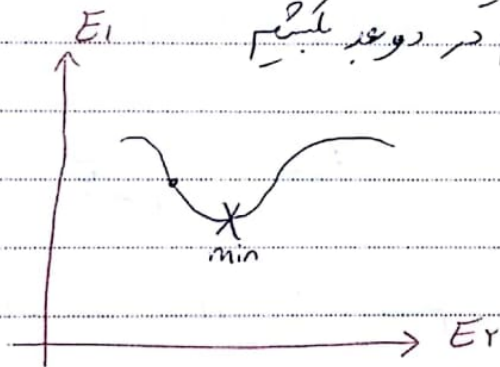
Year. Month. Date. ()

عروضه باید ۵ و ۱ باشد اما نسبت
 این خطه داریم. هر چه قدر خطه کمتر باشد
 یعنی w ها را درست تنظیم کنیم
 این برابر یک داده است. مثلاً حرارتاً داده ردی شبکه را تنظیم می کنیم.

مثلاً $\sum (1 - 0.17) = 0.15 = 0.15$

بر اساس این ها یک نمودار خطه تشکیل می دهد.

مغزین می کنیم رویا و سرتی داریم که بتوانیم در دو بعد بکشیم



E_1 : میزان خطه w_1
 E_2 : میزان خطه w_2

در مرحله ترین خطه لا بدست می آوریم

می خواهیم خطه min شود. هر چه اساس w_1 و w_2 را تغییر دهیم که این به نسبت
 کمینه شدن حرکت کند؟ مشتق می گیریم. باید این نقطه (min) منبسط و مطلق در نظر

بر اساس این ها مشخص می کنیم مقدار w ها چگونه تغییر کند
 این که می بینیم خطه کم شود، چه قدر به min نزدیک شویم به فاصله دارد
 در بعدی زیاد، تعیین این کار سخت می باشد.

یک تعداد تکرار مشخص می کنیم. مثلاً در این تعداد تکرار خطه min کنیم
 یا اینکه مشخص می کنیم مجموع خطه باید از این عدد کمتر باشد.

رنگ شبکه آموزش دیده شبکه مغزین w در تعیین شدند و بعد شبکه را Test می کنیم.
 آموزش در یک مرحله نیست. اینکه چند بار در یک شبکه نشان داده شود. تعیین می کنیم
 در مرحله test جواب مناسبی یا نمی شود. عروضه را ارائه می دهد و با جواب خاصی مقایسه
 می کند تا ببیند چه قدر خطه دارد.

وقتی داده ها را به عنوان تست می داریم. رنگی داده جدید به یاد ماشین می بیند

Subject:

Year. Month. Date. ()

: Apriori / الی

1 { I₁, I₂, I₃, I₄, I₅, I₆ }

2 { I₂, I₃, I₄, I₅, I₆, I₇ }

3 { I₁, I₇, I₈, I₉ }

4 { I₁, I₄, I₆, I₉, I₁₀ }

5 { I₂, I₄, I₅, I₁₀, I₁₁ }

min support = 60%

Support			
{ I ₁ }	3 *	{ I ₁ , I ₂ }	1
{ I ₂ }	3 *	{ I ₁ , I ₄ }	3 **
{ I ₃ }	2	{ I ₁ , I ₅ }	2
{ I ₄ }	5 *	{ I ₁ , I ₆ }	2
{ I ₅ }	4 *	{ I ₂ , I ₇ }	3 **
{ I ₆ }	3 *	{ I ₂ , I ₅ }	3 **
{ I ₇ }	1	{ I ₂ , I ₆ }	2
{ I ₈ }	1	{ I ₄ , I ₅ }	4 **
{ I ₉ }	1	{ I ₄ , I ₆ }	3 **
{ I ₁₀ }	2	{ I ₅ , I ₆ }	2
{ I ₁₁ }	1		

1, 2, 4 0

1, 4, 5 0

1, 4, 6 0

2, 4, 5 3 **

2, 4, 6 0

4, 5, 6 0

انجا closed سہ
درستکار دست نذاریم