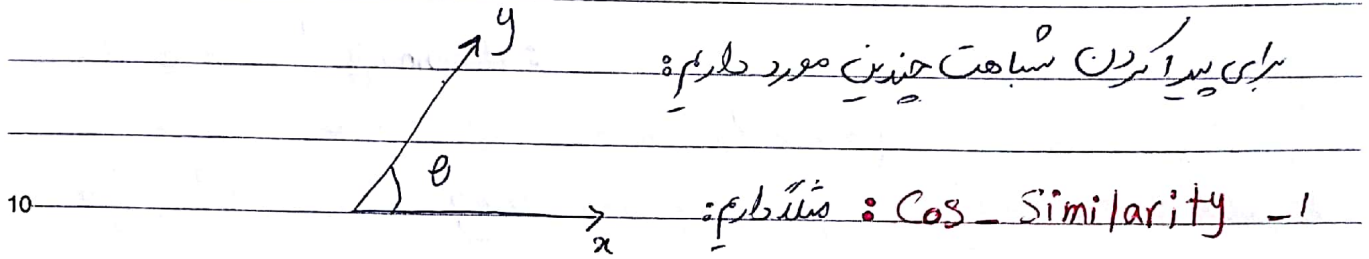


مثلاً $\frac{1}{3}$ متن شده apple. حال یک فرمول پیدا کرده شباهت بین دو vector را بنویسیم.

apple در متن 1، 10 بار تکرار شده و در متن دوم 30.

	Apple	Micro soft	Obama	Electro
5 D1	10/30	20/30	0	0
D2	30/90	60/90	0	0

برای پیدا کردن شباهت چندین مورد داریم:



سپید رنگند → عمود بر هم → $\theta = 90$
 $\cos \theta = 0$

سپید رنگند → بر هم منطبق اند → $\cos \theta = 1$
 $\theta = 0$
تصویر 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.

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

2- distance: مقداری که هر چقدر object ها متفاوت اند هر چه به صفر نزدیک باشد شباهت بیشتر.

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

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

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

د) Hamming distance:

تعداد بیت‌هایی که با هم تفاوت دارند (در دو کد)
10 چگونه شباهت بین String ها را پیدا کنیم؟ در آنتی DNA کاربرد دارد

Edit distance: تعداد Insert و Delete های که با هم متفاوت هستند باید بشود هم بارند.

فصل 4 کتاب پیدا کردن قوانین انجمنی و الگوهای پرتکرار

مفاهیم: الگوی پرتکرار: مثلاً در یک فروشگاه کالاهای است که بیشتر از همه خریداری شده‌اند.

و یا الگوهای پرتکرار را باید چگونه با هم

20 یکی از سوالات: بدانیم چه کالاهای، چه کالاهای می‌خریم. یعنی توالی در خرید کالا مثلاً می‌دانیم هر کس کامپیوتر بخرد، بعد از آن پرینتر می‌خرد. پس می‌توانیم یک تخفیف برای پرینتر می‌داریم که صرفاً خریداری می‌شود.

۱- توالی در خرید آنتیم‌ها

2- چه آنتیم‌هایی با هم خریداری شده‌اند. (مثلاً سرود ما با هم خریداری شده‌اند پس اینها

3- DNA تکرار یک چیز بیرونی چیز دیگر کنار هم می‌زنیم

مثلاً یک بار در وی DNA های مختلف جدا می‌کنیم و آزمایش می‌کنیم
با مثلاً چندین فروشنده چنانچه در زمینه فروش داشته

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

جای زیادی کاربرد دارد مثل سیستم فروش و بازاریابی، علم پزشکی،
پهنای صفحات وب ← در هر دو این ها نیاز داریم اشخاص را یکبار با هم ببینیم

10 سوال: اگر کسی کامپیوتر خرید بعد از خرید آن به چه احتمالی، آگهی و فروش می‌خورد؟

قانون: باید یک قانون پیدا کنیم که این اتفاق افتاده

15 $Computers \Rightarrow Software$ [support = 2% , confidence = 60%]

$A \Rightarrow B$

باید روابط با هم

20 فرض کنیم فروشنده n آگهی به شکل زیر دارد که می‌خواهم یک کالا را A بخردم.
بعد از خرید A ، صفت B را می‌خردم که اشتراک این دو رو می‌خواهم
همه وقت هست است و هیچ نباید اشتراک داشته باشند

T: مجموعه آگهی ها I: item

$A \Rightarrow B$

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

$A \cap B = \emptyset$

25

Support: (درجه پشتیبان) فرض کنید قانون در مجموعه تراکنش S, D است support است. در این صورت S درصدی از تراکنش D است که حاوی A و B است و با احتمال $P(A \cup B)$ باشد. هر چند $support(A \Rightarrow B) \Rightarrow A \cup B$

5 یعنی جایی که روگالا با هم خرید کرده اند.

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

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

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

این support که می بینیم min حالت است.

15 این قانون که $min = 1/2$ یعنی جایی جاها A و B با هم اند پس جایی قانون ما مثل گرفته ایم. آن وقتی می بینیم 98% support باشد. پس داریم قانون ما ضعیف است می بینیم. یعنی جایی که می بینیم هر کجا فاصله بین 98% باشد این رو با هم خریداری کرده باشند.

اگر قانون هیچ چیزی پیدا نکند پس این قانون بردن خورد.

20 یا مثلاً قانون ما این بگیریم کسانی که امروز یک کالا خریده اند. این احتمال ضعیف زیاد است. ← از آن چیزی خواهم بدست آورم که صحتش بردن خورد!

پس زیاد رتبه خارج کردن و کم رتبه خارج کردن بردن خورد. به چیز متوسط من خواهم خریدار به نظر آید که چیزی بداند و نه زیاد آزار می دهم.

25 ItemSet مجموعه اکتهم ها

Item

K-Itemset کمر Itemset کا عضویت

فراوانی (تکرار) : نفع شراکت Item در مجموع رہتا ہے

5 support → Absolute Support تعداد رکھتا ہے کہ A, B ہر دو کا مجموعہ دے دینا

→ Relative Support احتمال A اور B ہر دو کا مجموعہ دے دینا

Relative confidence: فرق ہر دو کا مجموعہ دے دینا

$$\text{Relative} = \frac{\text{تعداد رکھتا ہے A, B}}{\text{تعداد کل رکھتا ہے}}$$

سوال: فرض کریں کہ فروشنہ خریدتے تھے

TID	Items	Notes
10	Coke, nuts, diaper	min-support = 2 قرار دے رہے ہیں
15 20	Coke, coffee, diaper	یہ ہر دو چیزیں 2 یا 2 آگے ہونے
30	Coke, diaper, egg	یا تو ہونے
40	Nut, Egg, milk	ان (3) چیزوں (Itemset) 1 کا ہے
50	Nut, coffee, diaper, egg, milk	یہاں کہہ رہے ہیں

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

{egg} {milk}

یہاں ہر دو چیزیں خریدتے تھے

min support = 4 ہے

یہاں ہر دو چیزیں خریدتے تھے

2 - Itemset = {Coke, diaper}

{nut, egg} {egg, milk}

3-Itemset = {Nut, egg, milk}

4-Itemset ^{چیزی پیدا نمی شود} ^{چون تراش می خورند و پیدا نمی شود} ^{اما اگر min support = 1 باشد این 4 تا هم پیدا می شود}

5 ^{این min support را جوری در نظر می گیرند که این مجموعه ای که بدست می آید متناسب باشد}
support confidence

Coke ⇒ diaper ($\frac{3}{5} = 40\%$ و $\frac{3}{3} = 100\%$)

diaper ⇒ Coke ($\frac{3}{5} = 40\%$ و $\frac{3}{4} = 75\%$) ^{اینجا diaper خریده و بعد Coke سفرد}

تعیین ضرایب Support و Confidence

10 ^{این محضه باره کاری خیلی نزدیک است و این روش ها کم میارند}

15 اگر یک مجموعه پرستار باشد ^{آیا زیر مجموعه 3 آن پرستار است؟ بله} فرض کنید یک مجموعه 100 تا می پیدا کردیم که پرستار است.

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

$$\binom{100}{1} + \binom{100}{2} + \binom{100}{3} + \dots + \binom{100}{100} = 10^{30}$$

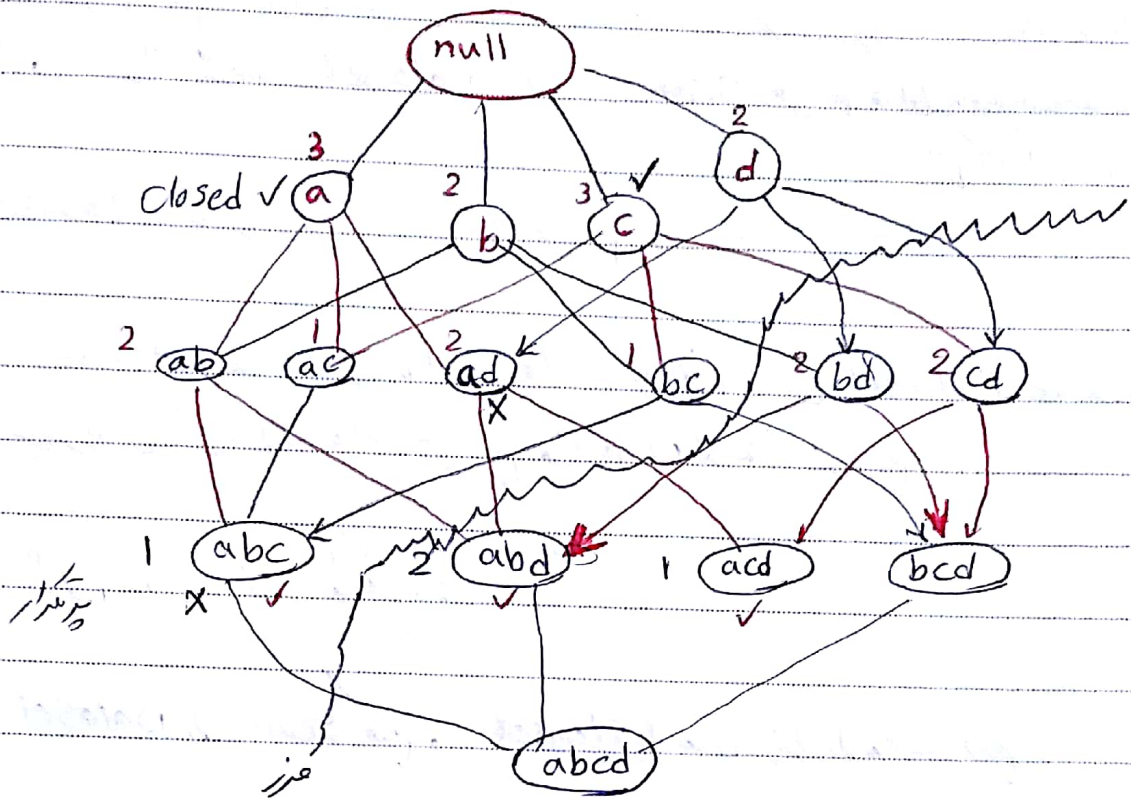
18 ^{این 10³⁰ مجموعه پرستار پیدا کردیم} اما باینی نسبت همه این ها را از هزینه کنیم ^{فقط 400 مجموعه 100 تا می را از هزینه می کنیم}

20 **Closed Frequent Itemset**: ^{بسته} یک Itemset مانند α که Closed است ^{به گونه ای که}

25 هیچ Super Itemset ای مانند α وجود نداشته باشد ^{نسبت به support آن} ^{آنگاه نسبت بالاتر از خودش - فراتر از خودش} ^{صاف α باشد}

مثلاً ارسال قبل از مجموعه 3-Itemset ^{چیزی بیشتر پیدا نکردیم}

Soroush ^{این مجموعه Closed است}



در این گراف یک فرزند بین مرتکبها و یک مرتکبها می‌کنند
 اما این فرزند ممکن است دقیق نباشد
 اگرچه هر چه که روی فرزند هستند support ها بیان لا بررسی می‌کنیم
 اگر مجموعه زیر مجموعه‌ها در سمت چپ فرزند باشد که آن ما کسب فرزند بودن
 دنیا است می‌تواند مثل abc و ad

۳ تا الگوریتم برای پیدا کردن مجموعه‌های مرتکبها:

Apriori - FP Growth - ECLAT

Apriori بر اساس دانش قبلی ما یک الگوی مرتکبها استخراج می‌کنند
 اگر مجموعه‌ها مرتکبها باشد که زیر مجموعه‌ها آن مرتکبها است
 این الگوریتم از این قانون استفاده می‌کنند:

اگر یک K -Itemset مرتکبها داشته باشیم از روی آن می‌توانیم یک $(K+1)$ -Itemset
 مرتکبها پیدا کنیم.

مراحل الگوریتم:
 ۱- در ابتدا مجموع Dataset و Scan می کنیم: اول باید دسته‌بندی را مشاهده کنیم
 سپس کجا این اتفاق افتاده

۲- Itemset ۱- های پرتکرار را استخراج می کنیم: ۱

۳- Itemset ۲- ها را با استفاده از join داخلی Itemset ۱- ها بدست می آوریم
 بعد پرتکرارهای آن را استخراج می کنیم و اسم آن را C2 می نامیم.

چنان مثال فروشنده را در نظر بگیریم:

۱- Dataset و Scan کنیم و Itemset ۱- های پرتکرار را بدست آوریم:

$$I_1 = \{ \{cok\} \{nut\} \{diaper\} \{coffee\} \{egg\} \{milk\} \}$$

$$I_2 = \{ \{cok, nut\} \{cok, diaper\} \{cok, coffee\} \{cok, egg\} \{cok, milk\} \}$$

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

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

برای پیدا کردن Itemset ۳- ها I_2 را با خودش join می کنیم (عندها از جدول قبلی)

پرتکرارهای I_2 را با خودش join می کنیم \rightarrow اگر ۴ تا یا بیشتر نویسیم فقط سه تا می ماند پرتکرار

$$I_3 = \{ cok, nut, diaper \}$$

این مثال کامل نیست

این روند را تا جایی ادامه می دهیم که دیگر مجموع پرتکرار پیدا نکنیم و به Closed می رسد.

مسئله: مجموعه‌های زیر را با بهنجاری و بریدن، با استفاده از الگوریتم.

TID	Items	Support
10	A, C, D	2
20	B, C, E	3
30	A, B, C, E	3
40	B, E	1

minsupport = 2

$L_1 = \{ \{A\}, \{B\}, \{C\}, \{D\}, \{E\} \}$

هری می‌کنیم Pruning چون minsupport = 2 است

$$L_2 = \{ \{A, B\}, \{A, C\}, \{A, E\}, \{B, C\}, \{B, E\}, \{C, E\} \}$$

Supports: 1, 2, 1, 2, 3, 2

$$L_3 = \{ \{A, B, C\}, \{A, C, E\}, \{B, C, E\} \}$$

Supports: 1, 1, 2

مجموعه $\{B, C, E\}$ بهنجاری است ← زیرمجموعه‌های آن هم بهنجاری است.
 مجموعه D حذف شد و در مراحل بعدی D را نادیده بگیریم.
 در این مرحله مجموعه $\{B, C, E\}$ نیز بهنجاری است.

مسئله: * درصفتی که کاربرد دارد.

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

minsupport = 2

عیب این روش این است که هر بار باید کل Dataset را Scan کنیم
 یا مثلاً یک سری مجموعه‌های اضافه داریم که باید آن‌ها را بررسی کنیم یا در حین
 تولید که به مرور بخوریم و هم تولید کنند
 در دیتا سیت‌های بزرگ خیلی بدتر شود ← یعنی زمان بر - زیاد اینها مثلاً دیتا سیت کوچک بود

اصلی‌ترین چالش‌های محاسبه:

- ۱- مرور چندباره پایگاه داده ترانسج
 - ۲- تعداد زیاد کاندیداها تولید شده
 - ۳- محاسبه Support برای کاندیداها زمان‌گیر است.
- به دنبال راه حل هستیم که این چالش‌ها را برطرف کنیم.
 راه حل ۱: ← ۱- کاهش مرور (اسکن) پایگاه داده
 ۲- کاندیداها را مفید تولید شوند
 ۳- محاسبه سریع Support ها

روش‌های بهبود Apriori:

- ۱- درم سازی
- ۲- کاهش ترانسج ها
- ۳- پارتیشن بندی
- ۴- نمونه گیری
- ۵- شمارش پویا

درم سازی: روش یک مثال یاد می‌کنیم.

- ایده اصلی: افزایان با این Itemset F-ها تولید می‌شوند 2-Itemset ها هم تولید می‌شوند
- 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 آن‌ها پرترنگار و frequent ساخته می‌شود

TID Items
 I_1, I_2, I_3 سوال: این سوال در صفحه قبل نوشته شده.

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

با توجه به این تقسیم بر 7 می تواند اعداد 0 تا 6 باشد. پس این حالت ها هم می تواند باشد.

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

bucket Address	0	1	2	3	4	5	6	Frequent
bucket Count	2	2	4	2	1	4	4	
bucket Content	{I1, I4}	{I1, I5}	{I2, I3}	{I2, I4}	{I2, I5}	{I1, I2}	{I1, I3}	
محتویات bucket	{I3, I5}	{I1, I5}	{I2, I3}	{I2, I4}		{I1, I2}	{I1, I3}	
			{I2, I3}			{I1, I2}	{I1, I3}	
			{I2, I3}			{I1, I2}	{I1, I3}	

اگر $\text{minSupport} = 3$ ← مجموع هر یک bucket آن ها بیشتر مساوی 3 باشد Frequent است.

تابع hash یک سری تراول دارد. روی چیزهایی است که Frequent خوانند. این روشی خطا دارد. هر ترانسیم تابع hash را مجوری اینجا می کشیم که تا حدی این تراول بر طرف شود.

تراول ← مثلاً ۳ تا مجموع داریم اما این ۴ تا مجموع هم میشه مثلاً $\{I_1, I_2\}$ و $\{I_1, I_5\}$ و $\{I_2, I_3\}$ تراول اینها فرودین Frequent را چندین باریم.

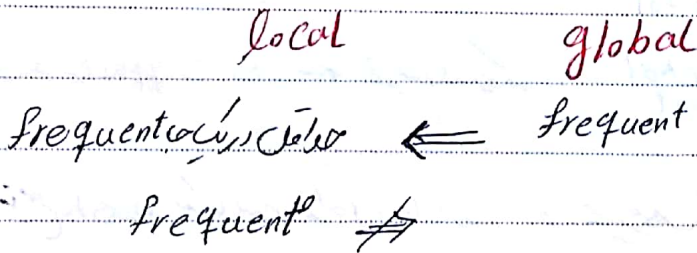
خاص ترانسی: اگر ترانسی در مجموع Itemset k الگوی بیشتر ندارد پس در مجموع k+1 Itemset هم عموماً الگوی بیشتر ندارد بنابراین آن ترانسی را حذف می کنیم.

مثلاً $\{I_1, I_2\}$ بیشتر است ← مربوط به کدام ترانسی است؟ تراول 200
 پس ترانسی 200 را حذف می کنیم و فقط عددی که آن را بیشتر می کشیم
 اما ترانسی 400 را حذف نمی کنیم! چون I_1 را هم دارد.
 ترانسی های را حذف می کنیم که فقط I_2, I_4 را دارند.

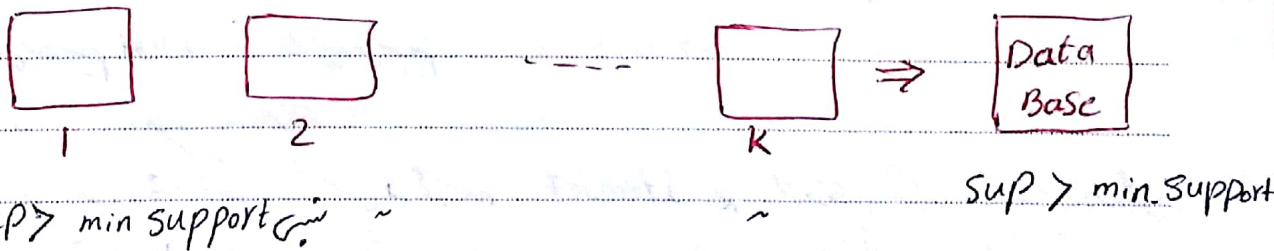
پارتیشن بندی: فرض کنید دیتابیس داریم که 10G است اما فضای که داریم 1G است. اگر قرار باشد هر دفعه 1G از دیتابیس را بیاریم و تو حافظه و Itemset-1 پیدا کنیم. این شکل باید 10 بار از دیتابیس اطلاعات را بیاریم توی حافظه اصلی. مثلاً فرض کنید $min\ support = 20\%$ باشد.

دیتابیس ها حافظه داریم که هر مجموعی یا ایتیم در هر قسمت 2 بار تکرار شده باشد در حافظه در 20 قسمت 20 بار تکرار می شود و به تکرار است اما ممکن است توزیع داده ها نواق باشد مثلاً یک جا 10 بار تکرار شده و یک جا 1 بار. پس می توان نتیجه گرفت اگر در یک قسمت به تکرار باشد حتی در حافظه هم تکرار است.

local ← در یک پارتیشن global ← در کل
 اگر یک Item در کل به تکرار باشد در یک قسمت یا local هم تکرار می شود.
 این قسمت مثل قسمت های قبل نیست که با Data کار داشته باشد بلکه با انتقال از دیتابیس حافظه کار دارد.



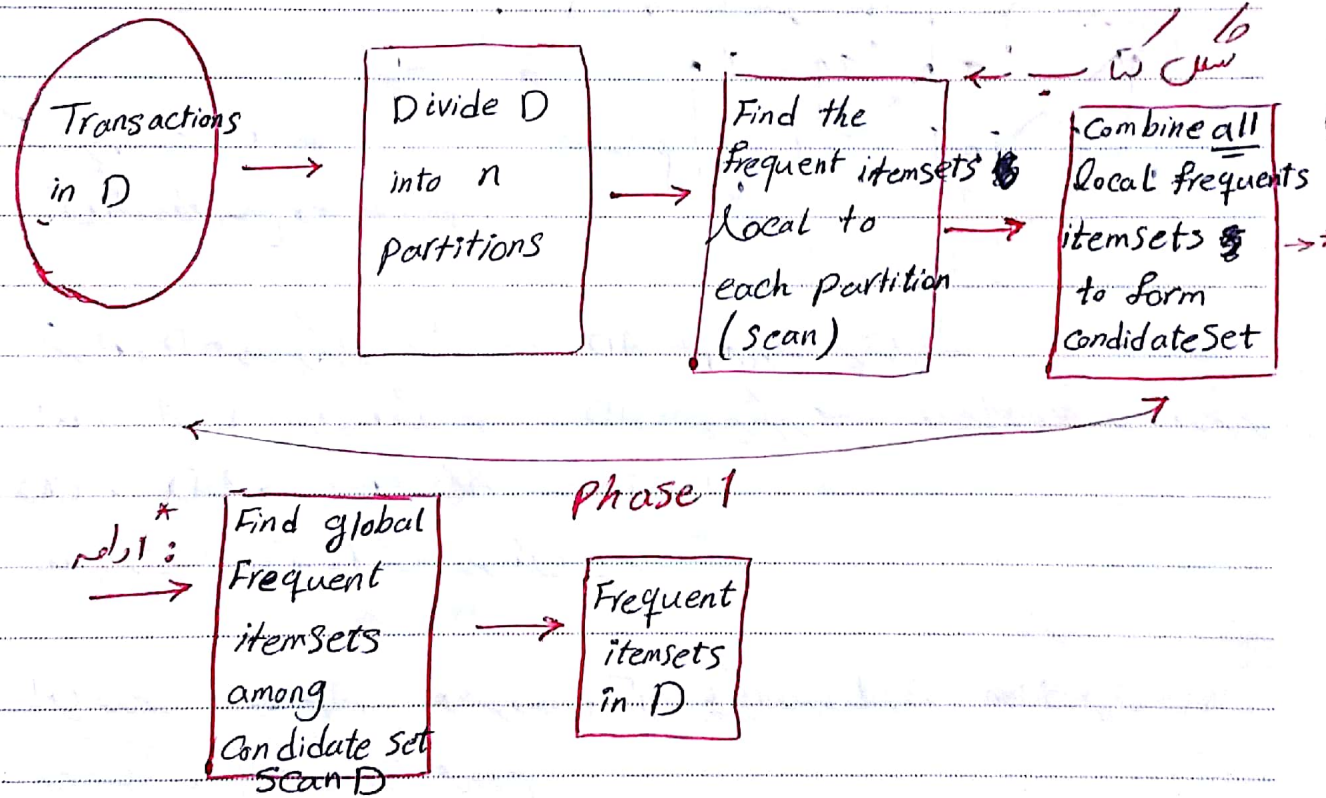
پارتیشن بندی در دوفاز انجام می شود:
 در فاز اول: داده ها به چند بخش تقسیم می شوند (به طوری که هر بخش بتواند در حافظه اصلی قرار گیرد)



تعریف $min\ support$ نسبی: باید به بسیاری نسبی = D باشد $\Leftarrow min$ تکرار برای هر پارتیشن برابر است با $min\ support \times$ تعداد تکرارهای همان یک بخش

مثلاً $min\ support = 20\%$ است و دیتابیس 10 قسمت تقسیم می شود به طوری که در هر قسمت 2 تکرار باشد.

جدید تعریف کنیم ← 1. انکوائی پر تکرار ہر عنصر (انکوائی میں) ممکن است نسبت ہر کل
 پائیاہ دارہ پر تکرار بائیاہ بائیاہ
 2. ہر Itemset کہ در پائیاہ پر تکرار بائیاہ باید حداقل در کسی از پارتیشن های پائیاہ دارہ
 پر تکرار بائیاہ
 3. نتیجہ ہر itemset پر تکرار در یک عنصر ہر تکرار در مجموعہ کل پائیاہ دارہ بائیاہ
 (یعنی عنصر کا تکرار بائیاہ و ایک پر تکرار بائیاہ در فائز دوم مستحق ہر تکرار)
 فائز دوم ← ہر ہر پائیاہ دارہ
 برائے ہر itemset کا تکرار باید Support واقف ہر تکرار
 در واقع این روش ہم دارہ پائیاہ دارہ ہر تکرار ہر تکرار



Support کل ہر itemset ہر
 min support کل ہر تکرار

نمونه برداری :

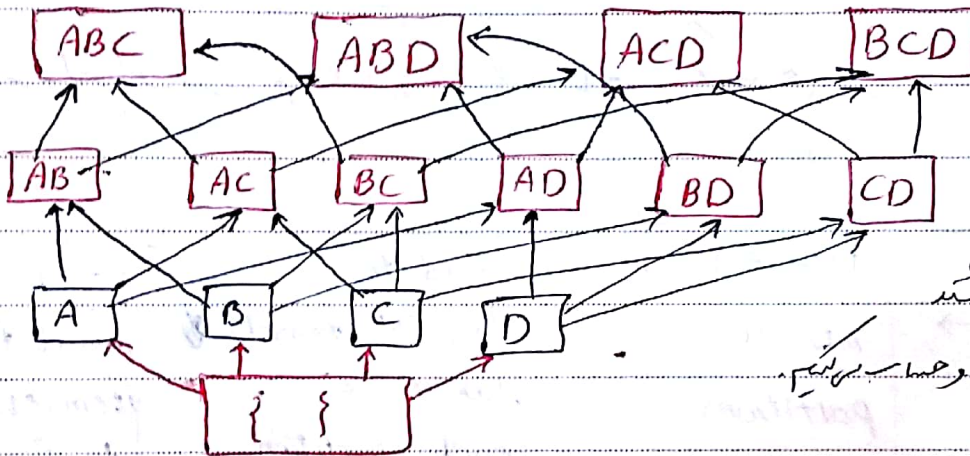
مشکل : در حالتی نمونه برداری می‌کنیم که نمونه‌ها در فاصله‌ها هم نزدیک باشند اما اگر نمونه‌ها از هم فاصله داشته باشند این کار قابل قبول نیست.

سپاریس پویا : حرف این کار برای خاصیت تعداد است.

مثال : اطلاق کار null است. همچنین داریم بعد 1-Itemset کارها می‌نویسیم.

ABCD

مرحله بعد 2-Itemset کارها می‌نویسیم.



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

اگر A و D تکرار باشند $\leftarrow AD$ هم می‌تواند تکرار باشد.

اما اگر A و D تکرار نباشد $\leftarrow AD$ 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 مجموعه‌ها زیر آن تکرار باشند.

این روش Apriori و محسودهای آن پایه‌داره بازید Scan می‌کنند

پس سرانجام روش‌های این موردیم