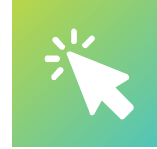


دوره آموزشے داده کاوی کاربردے

معمدامین ملاحسینے اردکانے | mollahoseini.ir



خوش آمدید

محمد امین ملاحسینی اردکانی | mollahoseini.ir

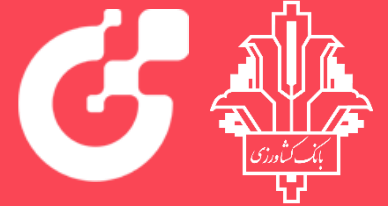
قُلْ هَلْ يَسْتَوِي الَّذِينَ يَعْلَمُونَ وَالَّذِينَ لَا يَعْلَمُونَ
بگو: آیا کسانی که می دانند با کسانی که نمی دانند برابرند؟
سوره زمر آیه ۹

محمد امین ملاحسینی اردکانی



هم بنیان گذار گروه تحلیل اطلاعات مالے بینا

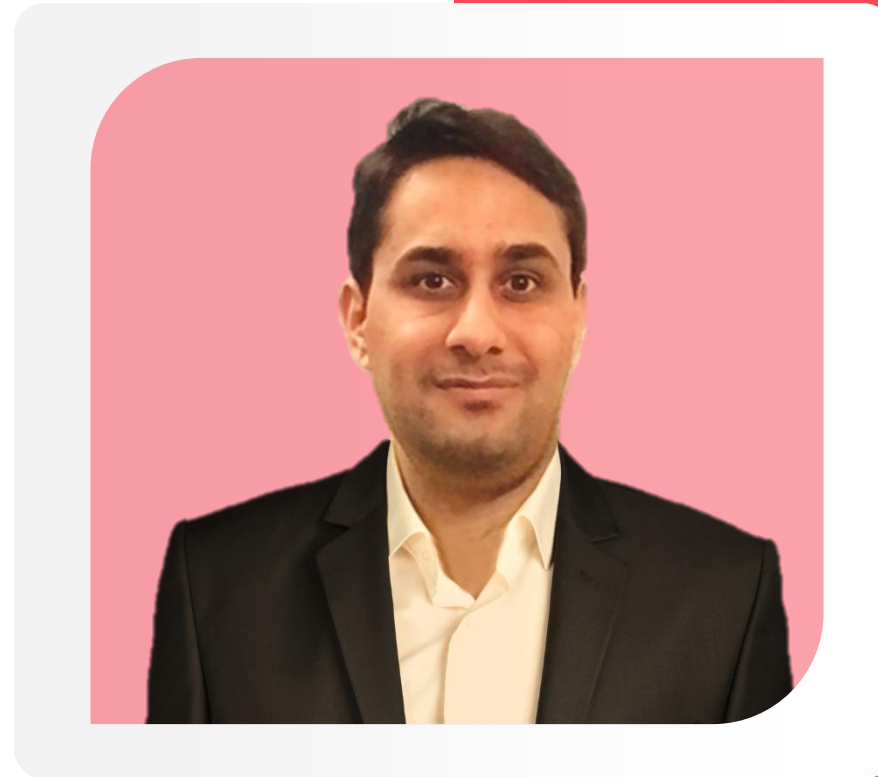
گسترش فناوری های نوین
شرکت سهامی خاص
Hi-Tech Solutions Co.



کارشناس کلان داده در شرکت گسترش
فناوری های نوین بانک کشاورزی



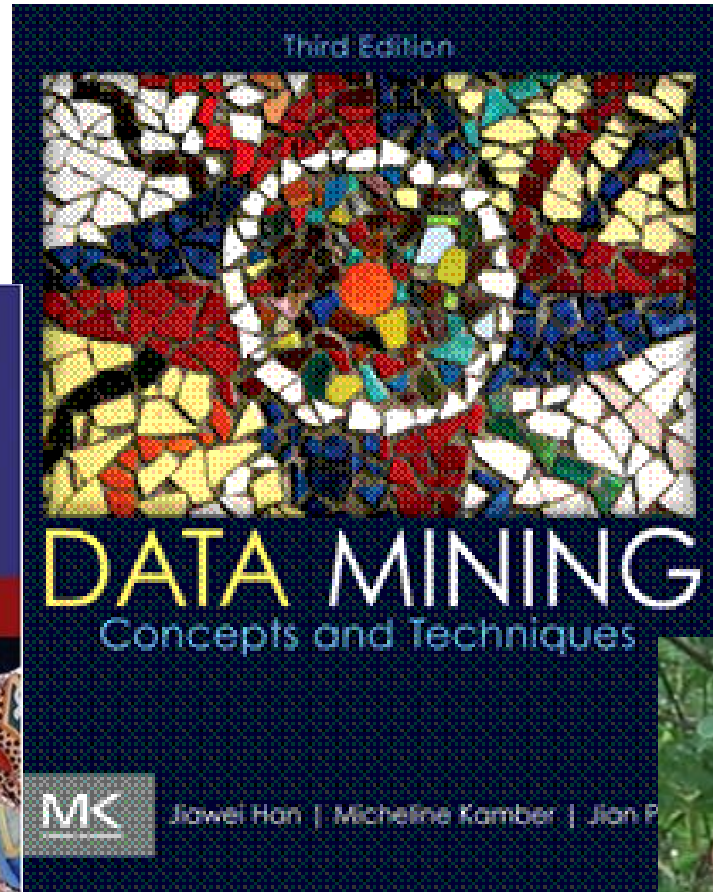
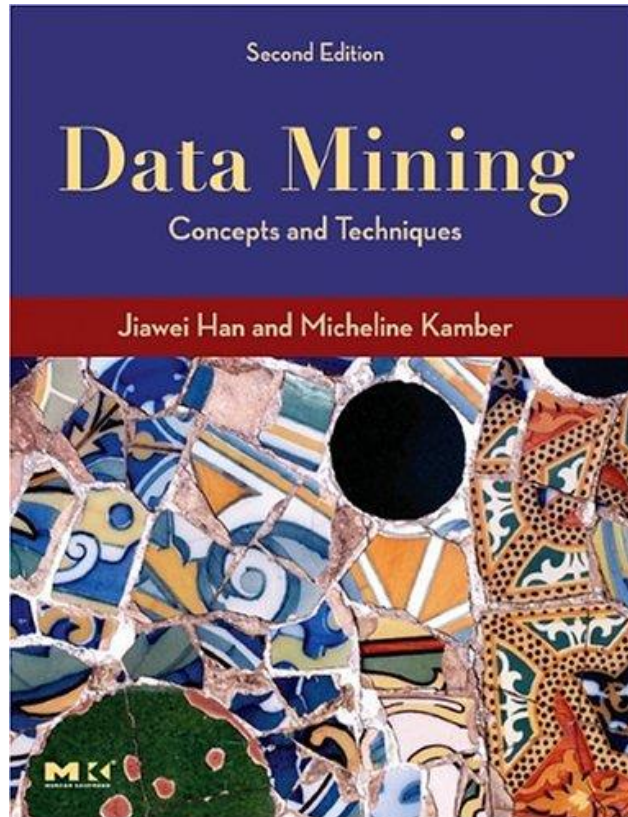
هم بنیان گذار گروه آموزش کانور



www.mollahoseini.ir

m.amin@mollahoseini.ir

منابع دوره



سرفصل‌ها

دوره آموزشی داده کاوی کاربردی



اجرای پروژه داده کاوی

انتخاب مسئله . جمع آوری داده ها . آماده سازی داده ها . انتخاب تکنیک داده کاوی . ارزیابی مدل



تکنیک‌های داده کاوی

کشف الگو . طبقه بندی . خوشه بندی . پیش بینی



آماده سازی و پیش پردازش

پاک سازی داده ها . جمعیت داده ها . تبدیل داده ها . کاهش داده ها . گسسته سازی داده ها



مفاهیم و کاربردها

تعریف داده کاوی . کاربردهای داده کاوی در کسب و کار . مزایای داده کاوی . فرآیند داده کاوی . عوامل موثر . انواع داده ها . تسک‌های اساسی . چالش‌ها



مفاهیم و کاربردها

- تعریف داده کاوی
- کاربردهای داده کاوی در کسب و کار
- مزایای داده کاوی
- فرآیند داده کاوی
- عوامل موثر در داده کاوی
- انواع داده‌ها
- تسک‌های اساسی در داده کاوی
- چالش‌ها

چرا داده کاوی؟

داده کاوی چیست؟

کاربردهای داده کاوی در کسب و کار

هتل و هتل داری

بررسی رفتار مشتریان
شناسایی فرصت‌های بازاریابی
بهبود تجربه مشتری
کاهش هزینه‌های عملیاتی

ت. صنایع غذایی

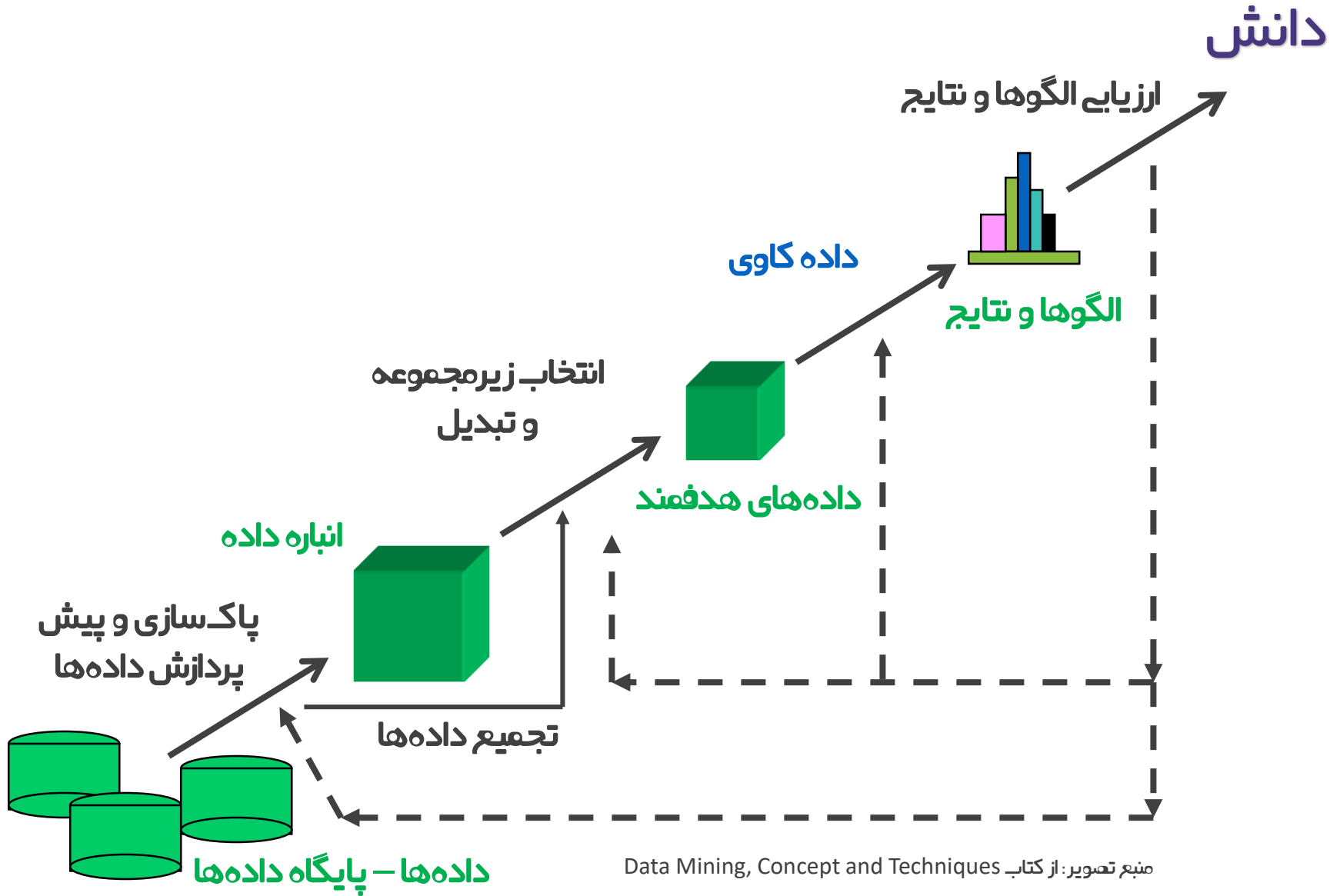
پیش‌بینی تقاضا
شناسایی فرصت‌های بازاریابی
بهبود کیفیت محصولات غذایی
کاهش هزینه‌های تولید

منابع انسانی

پیش‌بینی عملکرد کارکنان
شناسایی استعدادها
جذب و استخدام کارکنان
ارزیابی عملکرد کارکنان

بازارهای مالی

پیش‌بینی قیمت
شناسایی فرصت‌های سرمایه‌گذاری
کشف تقلب
مدیریت ریسک



مزایای داده کاوی

- بهبود تصمیم‌گیری
- افزایش بهره‌وری
- کاهش هزینه‌ها
- افزایش رضایت

فرآیند داده کاوی

عوامل موثر در داده کاوی

چه داده‌ای؟ تکنولوژی؟
نیازمندی؟ کاربرد؟

انواع داده‌ها

داده‌های رابطه‌ای
داده‌های تراکنش
داده‌های جریان
داده‌های سری زمانه
داده‌های زمانه

داده‌های ترتیبی
داده‌های گراف
داده‌های جغرافیایی
داده‌های زمانه مکانی

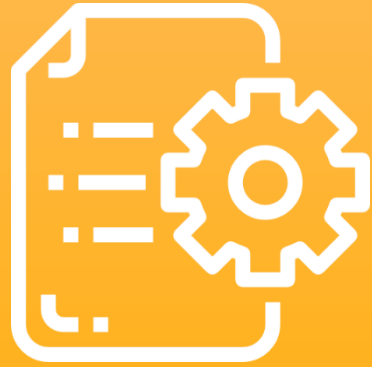
تسک‌های اساسی در داده کاوی

- ✓ توصیف داده‌ها
- ✓ پیدا کردن پترن‌ها، همبستگی‌ها و وابستگی‌ها
- ✓ رده‌بندی
- ✓ خوشه‌بندی
- ✓ داده‌های پرت

چالش‌ها

- حجم بالا (توزیع و موازی سازی – مقیاس پذیری – افزایش)
- واقعی هستند ○ ارزیابی نتایج
- دانش خارجی ○ حریم خصوصی

- ✓ تعریف داده کاوی
- ✓ کاربردهای داده کاوی در کسب و کار
- ✓ مزایای داده کاوی
- ✓ فرآیند داده کاوی
- ✓ عوامل موثر در داده کاوی
- ✓ انواع داده‌ها
- ✓ تسک‌های اساسی در داده کاوی
- ✓ چالش‌ها



آماده سازی و پیش پردازش

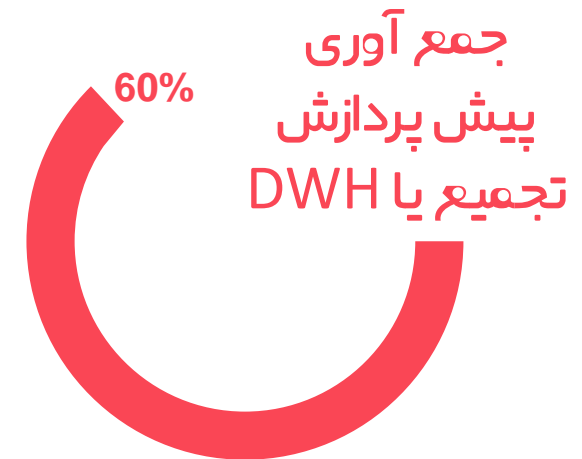
- پاک سازی داده ها
- جمع و یکپارچه سازی داده ها
- تبدیل داده ها
- کاهش داده ها
- گستره سازی داده ها

چرا نیاز به پیش پردازش داریم؟



منبع تصویر: ابزارهای هوش مصنوعی

کامل نباشند
نویزی باشند
متناقض باشند
تکراری باشند

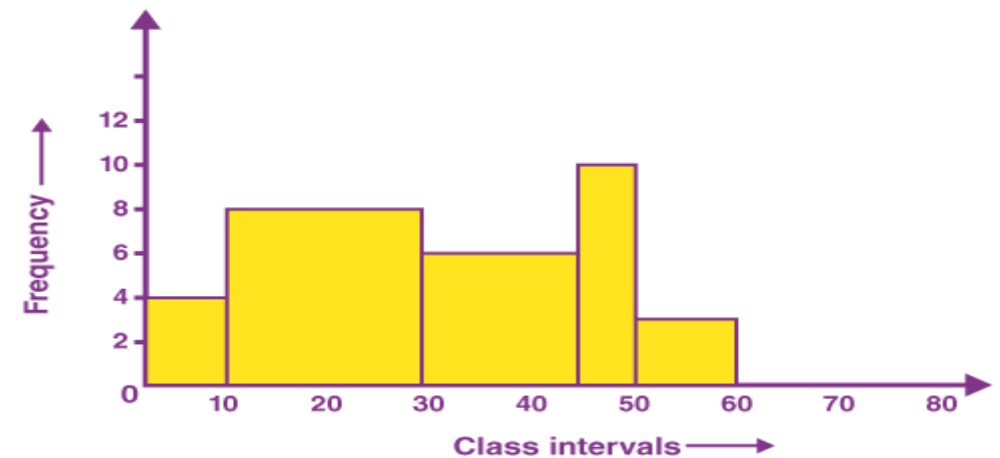
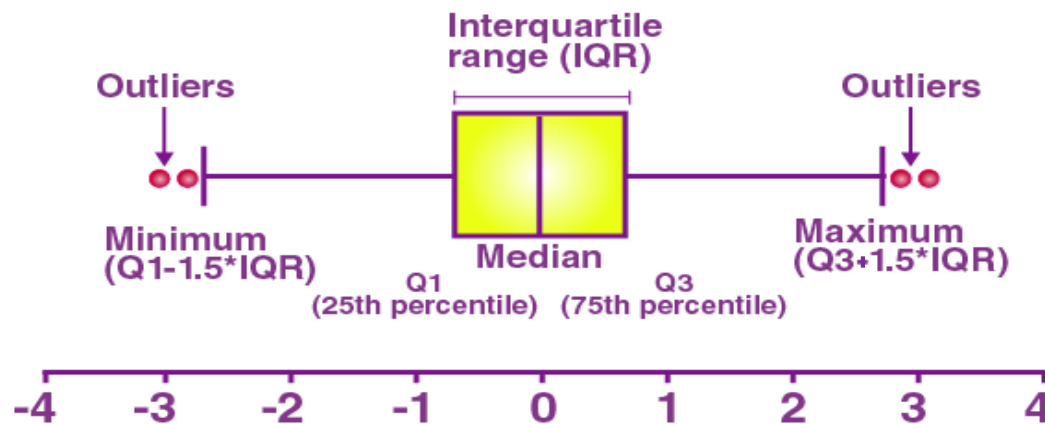


عمل‌های اصلی در پیش پردازش

- ۱- بررسی و توصیف
- ۲- پاک‌سازی
- ۳- تجمیع و یکپارچه‌سازی
- ۴- تبدیل
- ۵- کاهش
- ۶- گسسته سازی

بررسی و توصیف داده‌ها

- ✓ کمینه، بیشینه، میانگین، مد، میانه و ...
- ✓ نمودار جعبه یا Box plot
- ✓ نمودارهای مختلف مانند هیستوگرام (Histogram)



منبع تصاویر: <https://byjus.com>

پاک سازی داده ها

با داده های ناکامل چیکار کنیم؟

با داده های نویزی چطور؟

Binning – Equal width

Binning – Equal depth

Regression

Clustering



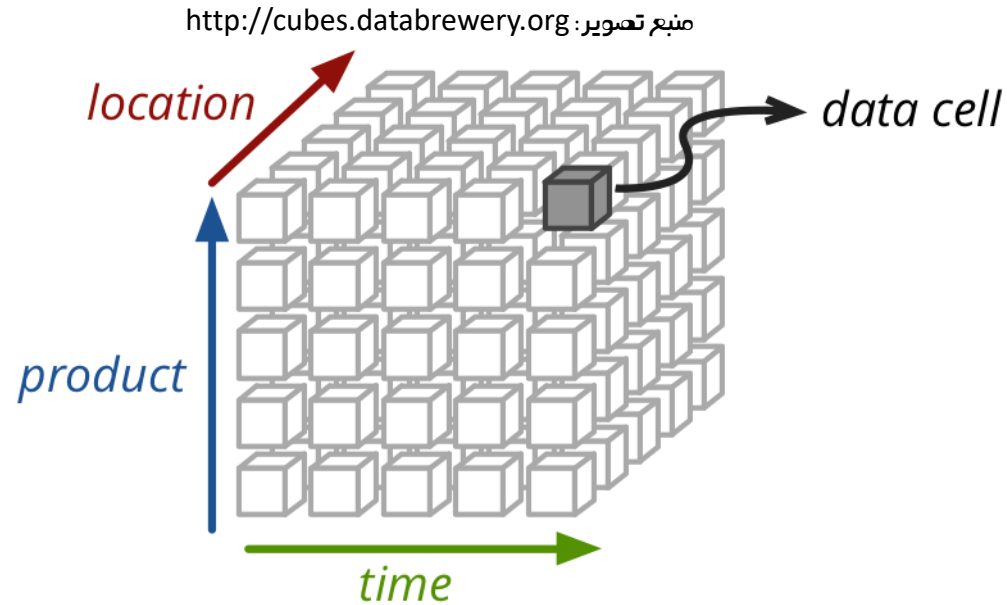
منبع تصویر: <https://medium.com>

تجمیع داده‌ها

- یکپارچه‌سازی اسکیما
- مقادیر مختلف از یک موجودیت
- تناقض‌ها
- افزودن داده‌ها
- همبستگی و وابستگی (Pearson، Chi-squared و غیره)

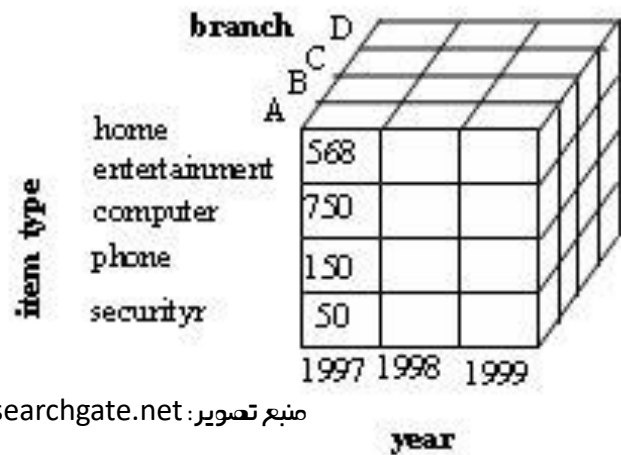
تبدیل داده‌ها

- تجمیع
- ایجاد سلسله مراتب
- نرمال کردن
- تبدیل ویژگی



کاهش داده‌ها

- معکب داده‌ای
- انتخاب زیر مجموعه از ویژگی‌ها
- فشرده سازی
- نمونه برداری



منبع تصویر: <http://researchgate.net>

گسسته سازی داده‌ها

- به کمک خوشه بندی
- با استفاده از بی نظمی

✓ پاک‌سازی داده‌ها

✓ جمع‌بندی و یکپارچه سازی داده‌ها

✓ تبدیل داده‌ها

✓ کاهش داده‌ها

✓ گسترده سازی داده‌ها



تکنیک‌های داده‌کاوی

- کشف الگو
- رده بندی و پیش بینی
- خوشه بندی

کشف الگو یعنی؟

- کدام محصولات باهم خریداری شده اند؟
- ترتیب به چه صورتی بوده تا به آیتم بعدی تخفیف ارائه نماییم؟

itemset و k-itemset

مجموعه آیتم های مکرر و قوانین انجمنی

- مایع ظرفشویی، دستکش
- لپ تاپ، نرم افزار و ویندوز، پرینتر

هدف: یافتن قوانین $A \rightarrow B$

- شرط min support
- شرط min confidence

support: هم شامل A و هم شامل B

confidence: اگر شامل B باشد حتما A را هم داشته باشد

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

شماره تراکنش	آیتم‌های خریداری شده
۱۰	نان، مغزیجات، پنیر
۲۰	نان، قهوه، پنیر
۳۰	نان، پنیر، تخم مرغ
۴۰	مغزیجات، تخم مرغ، شیر
۵۰	مغزیجات، قهوه، پنیر، تخم مرغ، شیر

min support = min confidence = 50%

مثال

itemset های مکرر:

نان: ۳ مغزیجات: ۳

پنیر: ۴ تخم مرغ: ۳

{نان، پنیر}: ۳

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

$A \rightarrow B$

✓ Support = 60% , confidence = 100%

✓ Support = 60% , confidence = 75%

قوانین استخراج شده:

(۱) نان \rightarrow پنیر

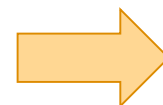
(۲) نان \rightarrow پنیر

فرآیند کشف قوانین و الگوها

یافتن
itemset های
مکرر



ایجاد
قوانین انجمنی



بررسی
شروط

الگوریتم Apriori

➤ اگر itemset ای مکرر باشد، باید تمام زیر مجموعه هایش هم مکرر باشد

- ۱- ابتدا 1-itemset ها رو محاسبه می‌کنیم
- ۲- itemset های کاندید با طول $k+1$ را از روی itemset های با طول k را می‌سازیم
- ۳- برای itemset های کاندید مکرر بودن را بررسی می‌کنیم
- ۴- به صورت تکراری مراحل ۲ و ۳ را تا زمان پایان itemset های مکرر انجام می‌دهیم

مجموعه داده‌ها

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

1st scan

Itemset	sup
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

L_1

Itemset	sup
{A}	2
{B}	3
{C}	3
{E}	3

مثال

min support = 2

L_2

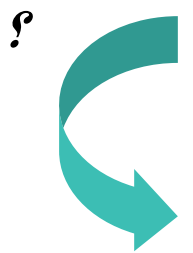
Itemset	sup
{A, C}	2
{B, C}	2
{B, E}	3
{C, E}	2

C_2

Itemset
{A, B}
{A, C}
{A, E}
{B, C}
{B, E}
{C, E}

2nd scan

Itemset	sup
{A, B}	1
{A, C}	2
{A, E}	1
{B, C}	2
{B, E}	3
{C, E}	2



C_3

Itemset
{B, C, E}

3rd scan

Itemset	sup
{B, C, E}	2

L_3

مشکلات الگوریتم Apriori

- حجم بالا و تعداد اسکن پایگاه داده
- تعداد زیاد کاندید
- محاسبه support

برخی از راهکارهای ارائه شده

- ✓ Partition: رفع مشکل حجم بالا و تعداد اسکن پایگاه داده
- ✓ DIC: رفع مشکل اسکن زیاد پایگاه داده
- ✓ DHP: رفع مشکل تعداد زیاد کاندید و تلاش به محاسبه آسان support

روش Partition

۱. پارتیشن بندی کل داده‌ها
۲. محاسبه itemset های مکرر هر پارتیشن
۳. یافتن itemset هایی که حداقل در یکی از پارتیشن ها مکرر باشد
۴. جمع‌بندی itemset ها و اسکن پایگاه داده

مجموعه داده‌ها

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

(۱)

P1) A, C, D
B, C, E

P2) A, B, C, E
B, E

(۲)

Itemset	sup	partition
{C}	2	P1
{B}	2	P2
{E}	2	P2
{B, E}	2	P2

مثال

min support = 75%

(۳)

{B, E}

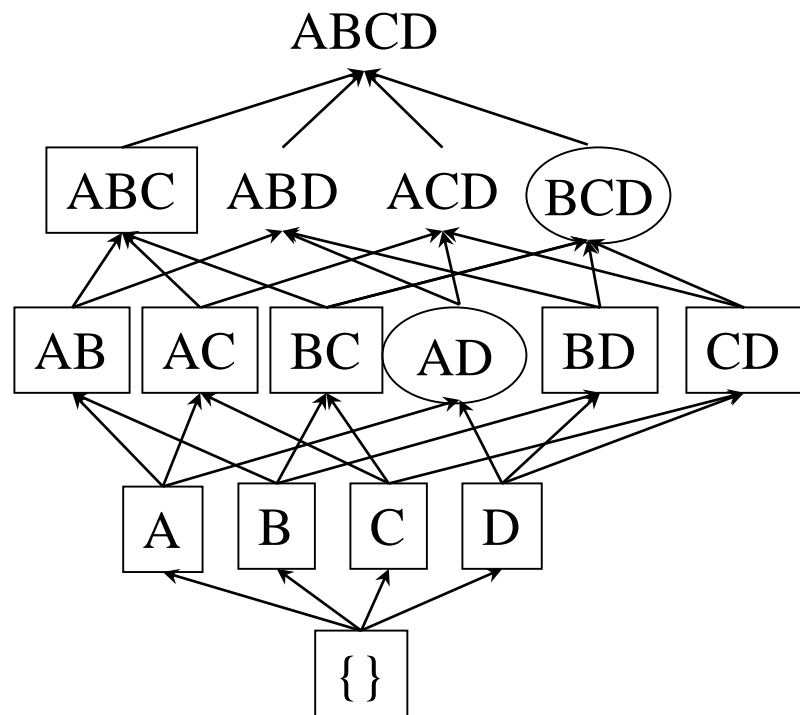
(۴)



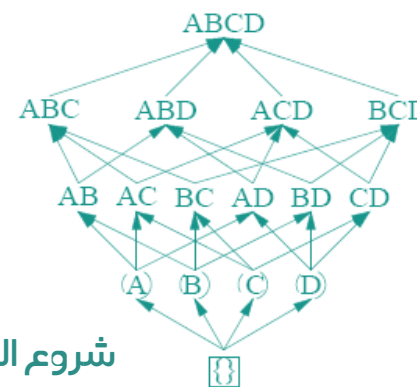
روش DIC

ایده: ساخت شبکه ای از itemset ها

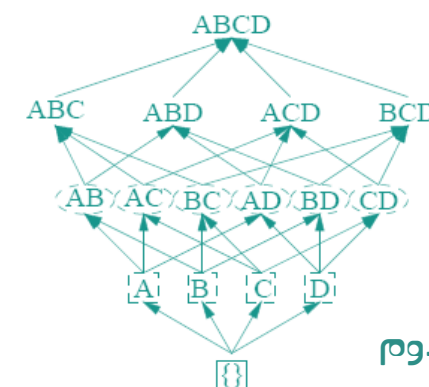
شبکه ای از itemset ها



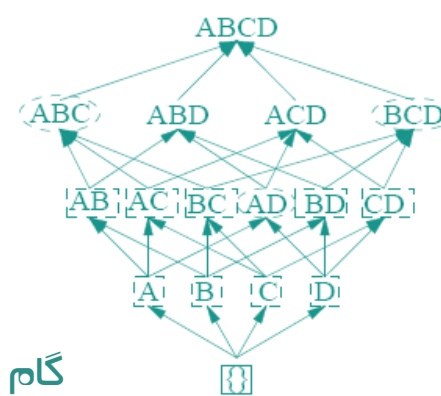
شروع الگوریتم



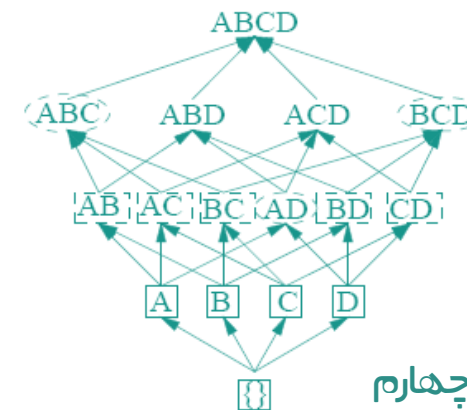
گام دوم



گام سوم



گام چهارم

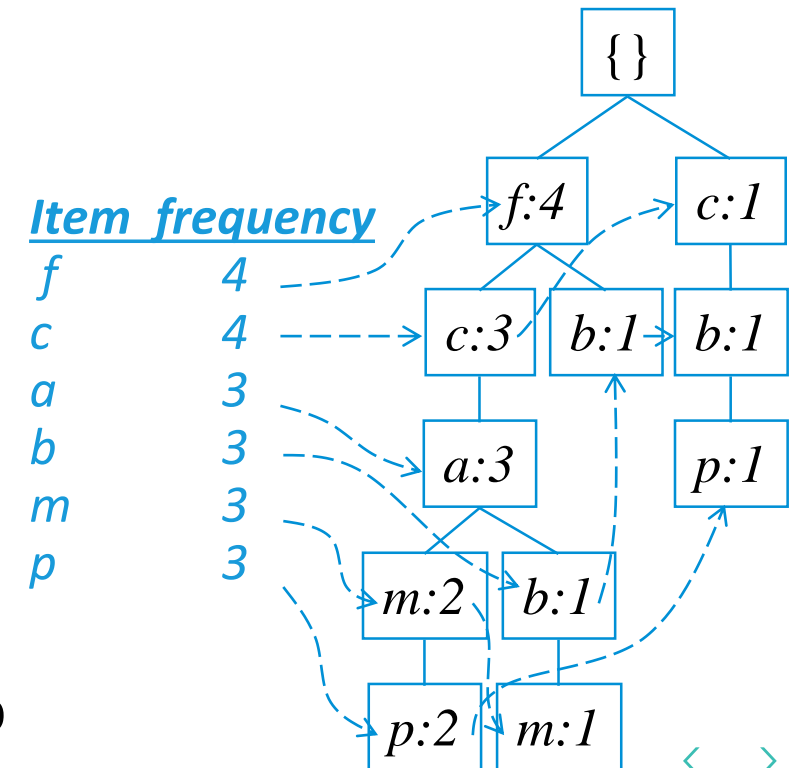


کاوش الگوهای مکرر بدون ایجاد کاندید

۱. پیدا کردن آیتم‌های مجاز
۲. مرتب سازی آیتم‌ها (ایجاد F-list)
۳. ایجاد ساختار درختی (FP-tree) و درج به صورت ترتیبی

min support = 3

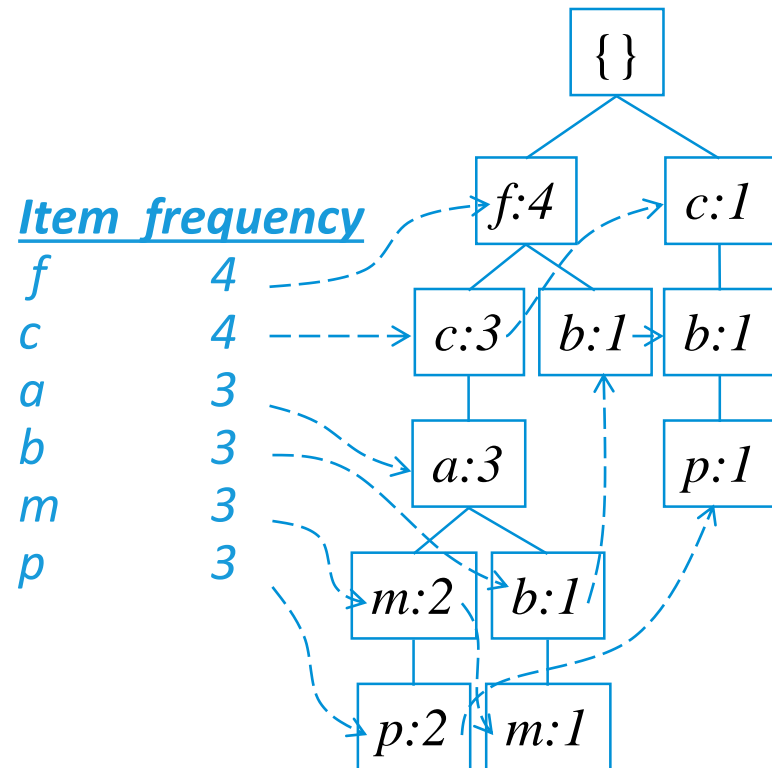
<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o, w}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}



کاوش الگوهای مکرر بدون ایجاد کاندید

پارتیشن بندی الگوها

- الگوهای شامل p
- الگوهای شامل m که شامل p نباشد
- الگوهای شامل b که شامل p و m نباشد
-
- الگوی فقط شامل f



پارتیشن بندی الگوها

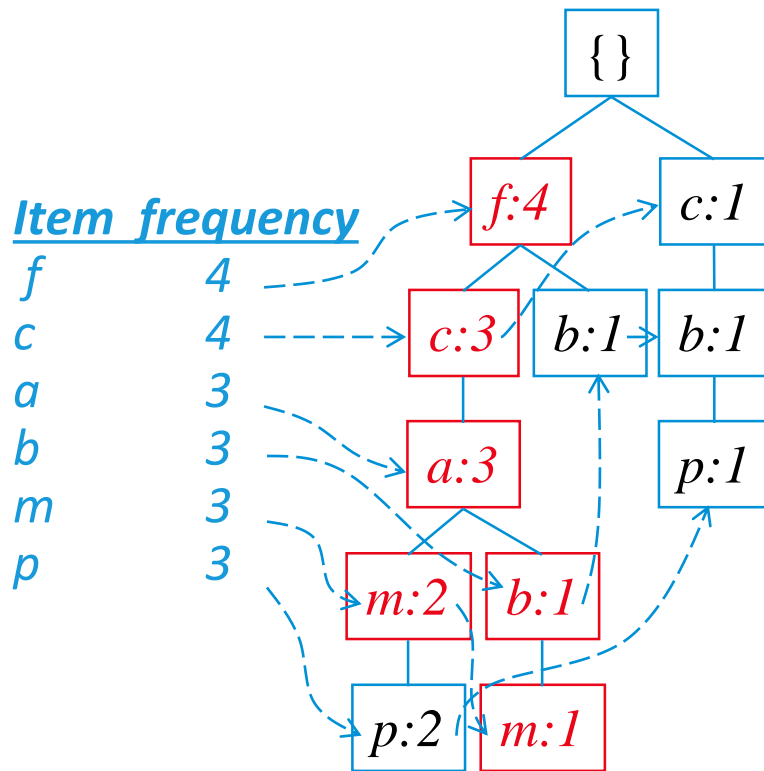
<u>item</u>	<u>cond. pattern base</u>
c	f:3
a	fc:3
b	fca:1, f:1, c:1
m	fca:2, fcab:1
p	fcam:2, cb:1

TID	Items bought	(ordered) frequent items
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o, w}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

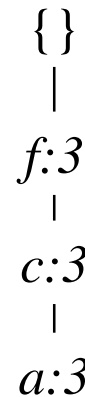
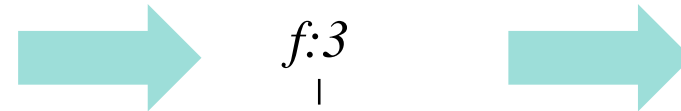
کاوش الگوهای مکرر بدون ایجاد کاندید

پارتیشن بندی الگوها

- الگوهای شامل p
- الگوهای شامل m که شامل p نباشد
- الگوهای شامل b که شامل p و m نباشد
-
- الگوی فقط شامل f



پارتیشن الگو مبتنی بر وجود m
fca:2, fcab:1



m-conditional FP-tree

همه الگوهای قابل ایجاد با m

- m,
- fm, cm, am,
- fcm, fam, cam,
- fcam

پارتیشن الگو مبتنی بر وجود cm $f:3$

{ }

|

 $f:3$ cm -conditional FP-treeالگو: fc

{ }

|

$f:3$

|

$c:3$

|

$a:3$

m -conditional FP-tree

پارتیشن الگو مبتنی بر وجود am $fc:3$

{ }

|

 $f:3$

|

 $c:3$ am -conditional FP-treeپارتیشن الگو مبتنی بر وجود cam $f:3$

{ }

|

 $f:3$ cam -conditional FP-treeالگو: $fcam$

چرا FP-Growth موفق است؟

- ✓ استراتژی D&C
- ✓ نیاز به ایجاد کاندید و تست کاندیدها ندارد
- ✓ پایگاه داده را فشرده نگه می‌دارد
- ✓ نیاز به اسکن چندباره پایگاه داده ندارد

آیا معیارهای min support و min confidence کافی است؟

	دستگاه بازی	نخریدن دستگاه بازی	مجموع (سطری)
دستگاه ویدیو	۴۰۰۰	۳۵۰۰	۷۵۰۰
نخریدن دستگاه ویدیو	۲۰۰۰	۵۰۰	۲۵۰۰
مجموع (ستونی)	۶۰۰۰	۴۰۰۰	۱۰۰۰۰

min support = 40% , min confidence = 60%

قوانین جهت بررسی:

- (۱) خرید دستگاه بازی → خرید دستگاه ویدیو
 (۲) خرید دستگاه ویدیو → خرید دستگاه بازی

قانون یک: support = $4000/10000 = 40\%$, confidence = $4000/7500 = 53.3\%$ ❌

قانون دو: support = $4000/10000 = 40\%$, confidence = $4000/6000 = 66.7\%$ ✅

۷۵ درصد به هر حال دستگاه ویدیو رو میخرن و وابسته به خرید دستگاه بازی نیست.

معیار lift: بررسی میزان همبستگی آیتم‌ها

$$A \rightarrow B \quad \text{lift}(A, B) = \frac{P(A \cap B)}{P(A) \times P(B)}$$

اگر $\text{lift} = 1$ مستقل
اگر $\text{lift} > 1$ همبستگی مثبت
اگر $\text{lift} < 1$ همبستگی منفی

	دستگاه بازی	نخریدن دستگاه بازی	مجموع (سطری)
دستگاه ویدیو	۴۰۰۰	۳۵۰۰	۷۵۰۰
نخریدن دستگاه ویدیو	۲۰۰۰	۵۰۰	۲۵۰۰
مجموع (ستونی)	۶۰۰۰	۴۰۰۰	۱۰۰۰۰

قوانین جهت بررسی:

قانون دو: $\text{lift} = \frac{4000/10000}{6000/10000 \times 7500/10000} = 0.89$ ❌

قانون سه: $\text{lift} = \frac{2000/10000}{6000/10000 \times 2500/10000} = 1.33$ ✅

۲) خرید دستگاه ویدیو \rightarrow خرید دستگاه بازی
۳) نخریدن دستگاه ویدیو \rightarrow خرید دستگاه بازی

symbol	measure	range	formula
ϕ	ϕ -coefficient	-1 ... 1	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
Q	Yule's Q	-1 ... 1	$\frac{P(A,B)P(\bar{A},\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A},\bar{B}) + P(A,\bar{B})P(\bar{A},B)}$
Y	Yule's Y	-1 ... 1	$\frac{\sqrt{P(A,B)P(\bar{A},\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A},\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}}$
k	Cohen's	-1 ... 1	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
PS	Piatetsky-Shapiro's	-0.25 ... 0.25	$P(A, B) - P(A)P(B)$
F	Certainty factor	-1 ... 1	$\max\left(\frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)}\right)$
AV	added value	-0.5 ... 1	$\max(P(B A) - P(B), P(A B) - P(A))$
K	Klogsen's Q	-0.33 ... 0.38	$\sqrt{P(A, B) \max(P(B A) - P(B), P(A B) - P(A))}$
g	Goodman-kruskal's	0 ... 1	$\frac{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
M	Mutual Information	0 ... 1	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\sum_i P(A_i) \log P(A_i) + \sum_j P(B_j) \log P(B_j) - \sum_i \sum_j P(A_i, B_j) \log P(A_i, B_j)}$
J	J-Measure	0 ... 1	$\max(P(A, B) \log\left(\frac{P(B A)}{P(B)}\right) + P(\bar{A}\bar{B}) \log\left(\frac{P(\bar{B} \bar{A})}{P(\bar{B})}\right))$
G	Gini index	0 ... 1	$P(A, B) \log\left(\frac{P(A B)}{P(A)}\right) + P(\bar{A}\bar{B}) \log\left(\frac{P(\bar{A} \bar{B})}{P(\bar{A})}\right)$
s	support	0 ... 1	$\max(P(A)[P(B A)]^2 + P(\bar{B} A)^2 + P(\bar{A})[P(B \bar{A})]^2 + P(\bar{B} \bar{A})^2 - P(B)^2 - P(\bar{B})^2,$
c	confidence	0 ... 1	$P(B)[P(A B)]^2 + P(\bar{A} B)^2 + P(\bar{B})[P(A \bar{B})]^2 + P(\bar{A} \bar{B})^2 - P(A)^2 - P(\bar{A})^2)$
L	Laplace	0 ... 1	$P(A, B)$
IS	Cosine	0 ... 1	$\max(P(B A), P(A B))$
γ	coherence(Jaccard)	0 ... 1	$\max\left(\frac{NP(A,B)+1}{NP(A)+2}, \frac{NP(A,B)+1}{NP(B)+2}\right)$
α	all_confidence	0 ... 1	$\frac{P(A,B)}{\sqrt{P(A)P(B)}}$
o	odds ratio	0 ... ∞	$\frac{P(A)P(\bar{B})}{P(\bar{A})P(B)}$
V	Conviction	0.5 ... ∞	$\frac{\max(P(A), P(B))}{P(A) + P(B) - P(A, B)}$
λ	lift	0 ... ∞	$\frac{P(A, B)}{P(A)P(B)}$
S	Collective strength	0 ... ∞	$\max\left(\frac{P(A)P(\bar{B})}{P(\bar{A}\bar{B})}, \frac{P(B)P(\bar{A})}{P(\bar{B}\bar{A})}\right)$
χ^2	χ^2	0 ... ∞	$\frac{P(A, B) + P(\bar{A}\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A, B) - P(\bar{A}\bar{B})}$
			$\sum_i \frac{(P(A_i) - E_i)^2}{E_i}$

سایر معیارها

اعمال محدودیت ها

الگوهای زیاد و نامتمرکز

اعمال محدودیت‌ها و درخواست

های کاربر مثلا مجموع قیمت

کمتر از مبلغی باشد

محدودیت‌ها و شرایط کاربر باید

بررسی شوند

item	price
A	1
B	2
C	3
D	4
E	5

مجموعه داده‌ها

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

1st scan

Itemset	sup
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

L_1

Itemset	sup
{A}	2
{B}	3
{C}	3
{E}	3

مثال

min support = 2
sum(price) < 5

L_2

Itemset	sup
{A, C}	2
{B, C}	2
{B, E}	3
{C, E}	2

2nd scan

Itemset	sup
{A, B}	1
{A, C}	2
{A, E}	1
{B, C}	2
{B, E}	3
{C, E}	2

C_2

Itemset
{A, B}
{A, C}
{A, E}
{B, C}
{B, E}
{C, E}



تمرین

مجموعه داده‌ها

شناسه	اقلام تراکنش
T1	سوسیسی، نان ساندویچی، سس گوجه
T2	سوسیسی، نان ساندویچی
T3	سوسیسی، نوشابه، چیپس
T4	چیپس، نوشابه
T5	چیپس، سس گوجه
T6	سوسیسی، نوشابه، چیپس

مجموعه آیتم‌های مکرر را بیابید.

چند قانون یا الگو متناسب با شرایط زیر پیشنهاد دهید.

min support = 33.33% , min confidence = 60%

رده بندی یا Classification

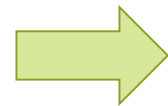
❖ هدف چیست؟

❖ بررسی مفاهیم مانند تاپل، کلاس یا برچسب و ...

❖ کاربردها

فرآیند کلی

آماده سازی
داده‌ها و پیش
پردازش



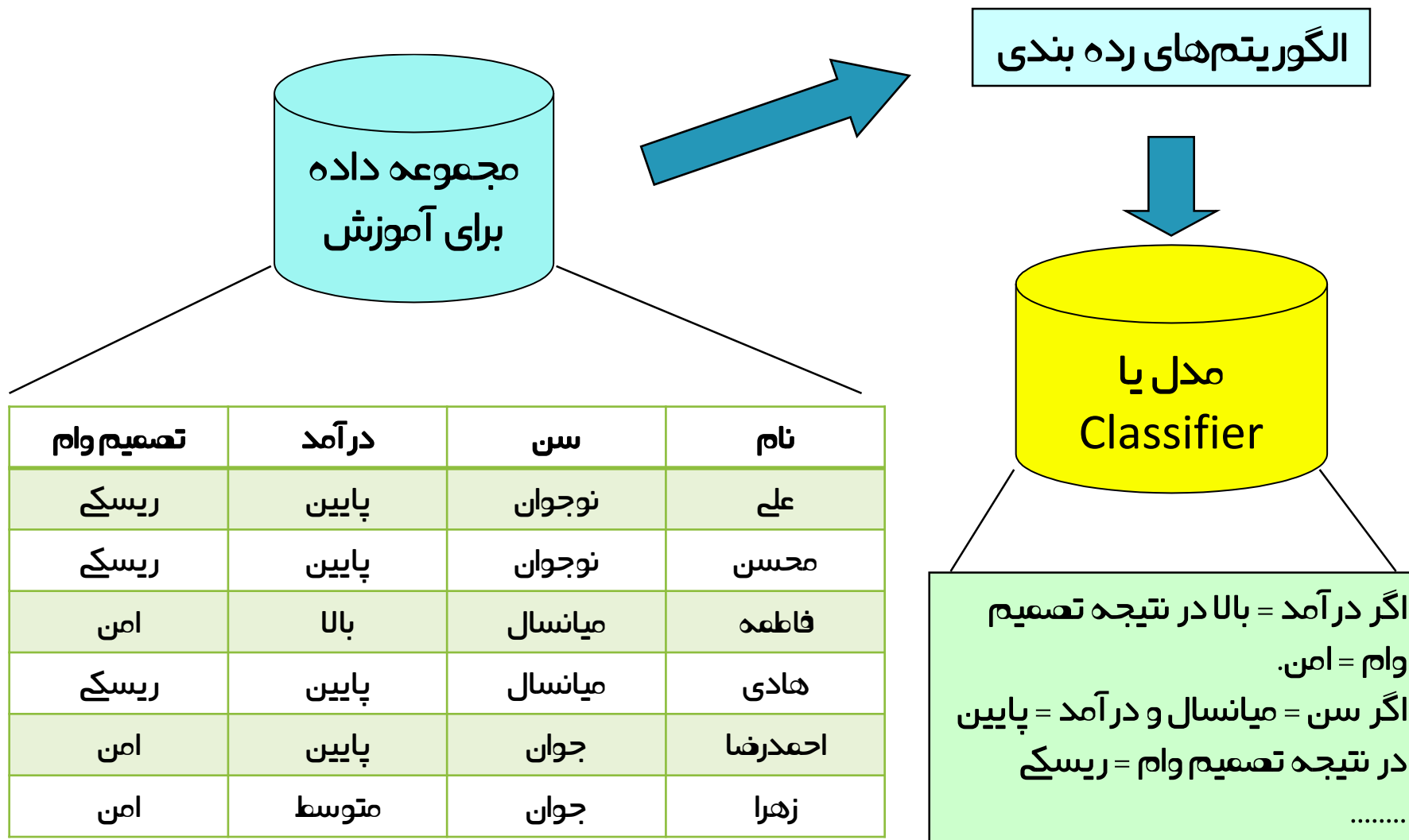
آموزش مدل یا
یادگیری مدل



ارزیابی مدل

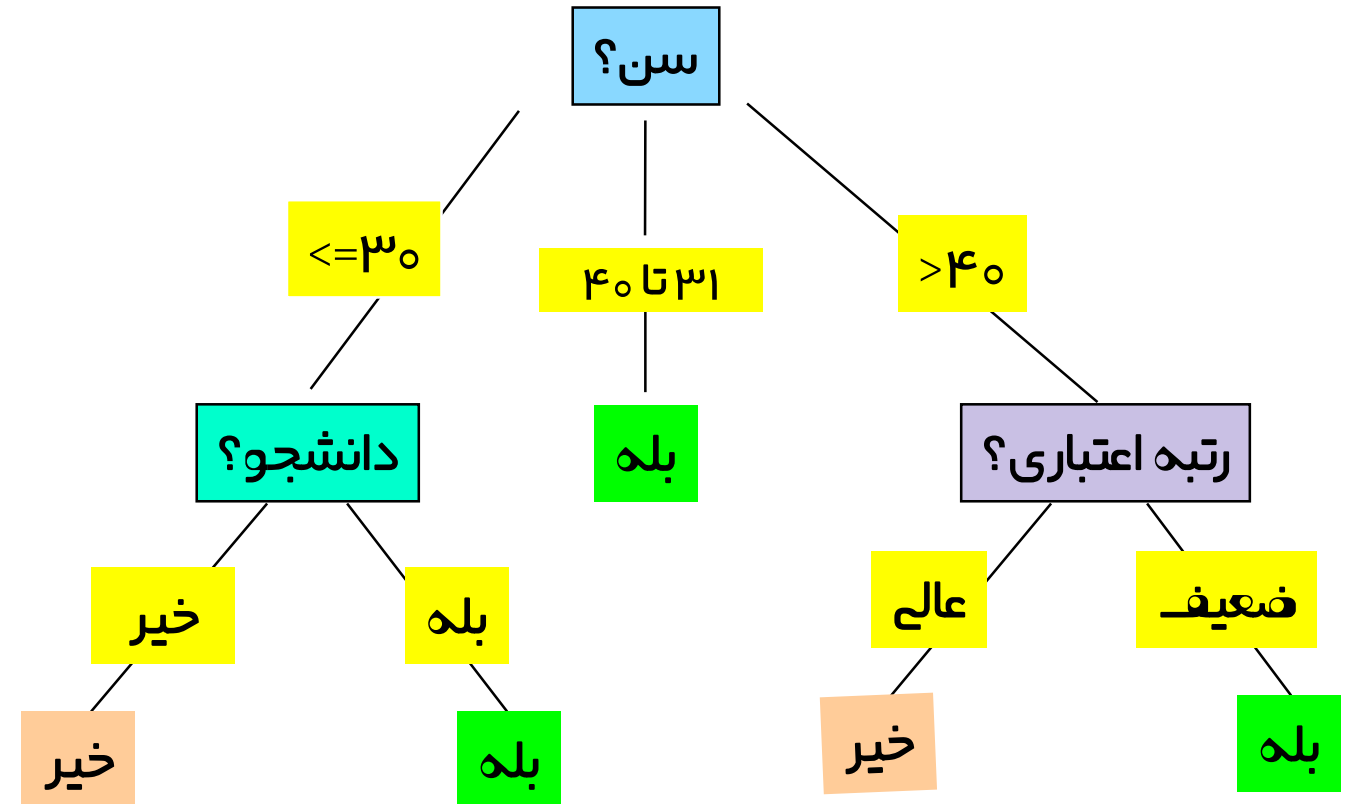
مثال

وام بدهیم یا خیر؟



سن	درآمد	دانشجو	رتبه اعتباری	خرید کامپیوتر
≤ 30	بالا	خیر	ضعیف	خیر
≤ 30	بالا	خیر	عالی	خیر
31 تا 40	بالا	خیر	ضعیف	بله
> 40	متوسط	خیر	ضعیف	بله
> 40	پایین	بله	ضعیف	بله
> 40	پایین	بله	عالی	خیر
31 تا 40	پایین	بله	عالی	بله
≤ 30	متوسط	خیر	ضعیف	خیر
≤ 30	پایین	بله	ضعیف	بله
> 40	متوسط	بله	ضعیف	بله
≤ 30	متوسط	بله	عالی	بله
31 تا 40	متوسط	خیر	عالی	بله
31 تا 40	بالا	بله	ضعیف	بله
> 40	متوسط	خیر	عالی	خیر

درخت تصمیم - ID3



درخت تصمیم – ID3

○ چرا ابتدا سن انتخاب شد؟

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

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times I(D_j)$$

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

خرید کامپیوتر	رتبه اعتباری	دانشجو	درآمد	سن
خیر	ضعیف	خیر	بالا	≤ 30
خیر	عالی	خیر	بالا	≤ 30
بله	ضعیف	خیر	بالا	۳۱ تا ۴۰
بله	ضعیف	خیر	متوسط	> 40
بله	ضعیف	بله	پایین	> 40
خیر	عالی	بله	پایین	> 40
بله	عالی	بله	پایین	۳۱ تا ۴۰
خیر	ضعیف	خیر	متوسط	≤ 30
بله	ضعیف	بله	پایین	≤ 30
بله	ضعیف	بله	متوسط	> 40
بله	عالی	بله	متوسط	≤ 30
بله	عالی	خیر	متوسط	۳۱ تا ۴۰
بله	ضعیف	بله	بالا	۳۱ تا ۴۰
خیر	عالی	خیر	متوسط	> 40

درخت تصمیم - ID3 چرا ابتدا سن انتخاب شد؟

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940 \quad (1)$$

age	p_i	n_i	$I(p_i, n_i)$
≤ 30	2	3	0.971
31...40	4	0	0
> 40	3	2	0.971

$$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694 \quad (2)$$

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246 \quad \checkmark$$

$$Gain(income) = 0.029 \quad (3)$$

$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048 \quad (4)$$

مشکل الگوریتم ID3

- در IG، تلاش در کاهش آنتروپی و در نتیجه کاهش عدم خلوص
- بر اساس کد ملی مشتریان چطور؟
- ✓ الگوریتم C4.5: نرمال سازی – جریمه فیلد به ازای تنوع در مقادیر
- ✓ الگوریتم CART: استفاده از معیار Gini و ایجاد درخت باینری

Overfitting و هرس کردن

رده بندی بیزین

❖ مبتنی بر احتمالات و فرمول بیز

❖ احتمال تعلق یک نمونه به کلاس – < کلاس با بیشترین احتمال

❖ مثال نمونه X: سن ۳۰ و درآمد ۴۰ میلیون تومان

اگر بدانیم فرضیه برقرار است،
احتمال مشاهده نمونه X

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

احتمال خرید یا نخریدن کامپیوتر

احتمال مشاهده نمونه

احتمال هر یک از فرضیات به شرط X

مثلا: احتمال خرید کامپیوتر به شرط X

احتمال نخریدن کامپیوتر به شرط X

$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i)P(C_i)}{P(\mathbf{X})}$$

$$\rightarrow P(C_i | \mathbf{X}) = P(\mathbf{X} | C_i)P(C_i)$$

$$P(\mathbf{X} | C_i) = \prod_{k=1}^n P(x_k | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \times \dots \times P(x_n | C_i)$$

خرید کامپیوتر	رتبه اعتباری	دانشجو	درآمد	سن
خیر	ضعیف	خیر	بالا	≤ 30
خیر	عالی	خیر	بالا	≤ 30
بله	ضعیف	خیر	بالا	31 تا 40
بله	ضعیف	خیر	متوسط	> 40
بله	ضعیف	بله	پایین	> 40
خیر	عالی	بله	پایین	> 40
بله	عالی	بله	پایین	31 تا 40
خیر	ضعیف	خیر	متوسط	≤ 30
بله	ضعیف	بله	پایین	≤ 30
بله	ضعیف	بله	متوسط	> 40
بله	عالی	بله	متوسط	≤ 30
بله	عالی	خیر	متوسط	31 تا 40
بله	ضعیف	بله	بالا	31 تا 40
خیر	عالی	خیر	متوسط	> 40

مثال از رده بند Naïve Bayes

سوال: پیش بین کنید که اگر یک مشتری دانشجو با سن کمتر از 30، درآمد متوسط و دارای رتبه اعتباری ضعیف باشد، کامپیوتر می‌خرد؟

گام اول محاسبه $P(C_i)$ ها:

$$P(\text{خریدن کامپیوتر}) = 9/14 = 0.643$$

$$P(\text{نخریدن کامپیوتر}) = 5/14 = 0.357$$

گام دوم محاسبه $P(X|C_i)$ برای هر کلاس:

$$P(\text{خریدن کامپیوتر} | \text{سن کمتر از } 30) = 2/9 = 0.222$$

$$P(\text{نخریدن کامپیوتر} | \text{سن کمتر از } 30) = 3/5 = 0.6$$

$$P(\text{خریدن کامپیوتر} | \text{درآمد متوسط}) = 4/9 = 0.444$$

$$P(\text{نخریدن کامپیوتر} | \text{درآمد متوسط}) = 2/5 = 0.4$$

$$P(\text{خریدن کامپیوتر} | \text{دانشجو بودن}) = 6/9 = 0.667$$

$$P(\text{نخریدن کامپیوتر} | \text{دانشجو بودن}) = 1/5 = 0.2$$

$$P(\text{خریدن کامپیوتر} | \text{رتبه اعتباری ضعیف}) = 6/9 = 0.667$$

$$P(\text{نخریدن کامپیوتر} | \text{رتبه اعتباری ضعیف}) = 2/5 = 0.4$$

خرید کامپیوتر	رتبه اعتباری	دانشجو	درآمد	سن
خیر	ضعیف	خیر	بالا	≤ 30
خیر	عالی	خیر	بالا	≤ 30
بله	ضعیف	خیر	بالا	۳۱ تا ۴۰
بله	ضعیف	خیر	متوسط	> 40
بله	ضعیف	بله	پایین	> 40
خیر	عالی	بله	پایین	> 40
بله	عالی	بله	پایین	۳۱ تا ۴۰
خیر	ضعیف	خیر	متوسط	≤ 30
بله	ضعیف	بله	پایین	≤ 30
بله	ضعیف	بله	متوسط	> 40
بله	عالی	بله	متوسط	≤ 30
بله	عالی	خیر	متوسط	۳۱ تا ۴۰
بله	ضعیف	بله	بالا	۳۱ تا ۴۰
خیر	عالی	خیر	متوسط	> 40

مثال از رده بند Naïve Bayes

$$P(\text{خریدن کامپیوتر}) = 9/14 = 0.643$$

$$P(\text{نخریدن کامپیوتر}) = 5/14 = 0.357$$

$$P(\text{سن کمتر از } 30 \mid \text{خریدن کامپیوتر}) = 2/9 = 0.222$$

$$P(\text{سن کمتر از } 30 \mid \text{نخریدن کامپیوتر}) = 3/5 = 0.6$$

$$P(\text{درآمد متوسط} \mid \text{خریدن کامپیوتر}) = 4/9 = 0.444$$

$$P(\text{درآمد متوسط} \mid \text{نخریدن کامپیوتر}) = 2/5 = 0.4$$

$$P(\text{دانشجو بودن} \mid \text{خریدن کامپیوتر}) = 6/9 = 0.667$$

$$P(\text{دانشجو بودن} \mid \text{نخریدن کامپیوتر}) = 1/5 = 0.2$$

$$P(\text{رتبه اعتباری ضعیف} \mid \text{خریدن کامپیوتر}) = 6/9 = 0.667$$

$$P(\text{رتبه اعتباری ضعیف} \mid \text{نخریدن کامپیوتر}) = 2/5 = 0.4$$

گام سوم) محاسبه احتمال نهایی:

$$P(X \mid C_i) : P(X \mid \text{خریدن کامپیوتر}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$$

$$P(X \mid \text{نخریدن کامپیوتر}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$$

$$P(X \mid C_i) * P(C_i) : P(X \mid \text{خریدن کامپیوتر}) * P(\text{خریدن کامپیوتر}) = 0.028$$

$$P(X \mid \text{نخریدن کامپیوتر}) * P(\text{نخریدن کامپیوتر}) = 0.007$$

بنابراین نظر ماشین با این مجموعه داده، خریدن کامپیوتر است.

معایب رده بند بیزین

❖ فرض مستقل بودن: مثال بیمار

❖ انواع داده‌ها به ویژه داده‌های متنی

رده بندی با قوانین

❖ قوانین را چگونه به دست آوریم؟

❖ برای تداخل‌ها چه اقدامی می‌توان انجام داد؟

✓ اعمال محدودیت‌های بیشتر

✓ مرتب سازی کلاس‌ها بر اساس فرکانس

✓ مرتب سازی قوانین با نظر کارشناس یا معیارهای آماری

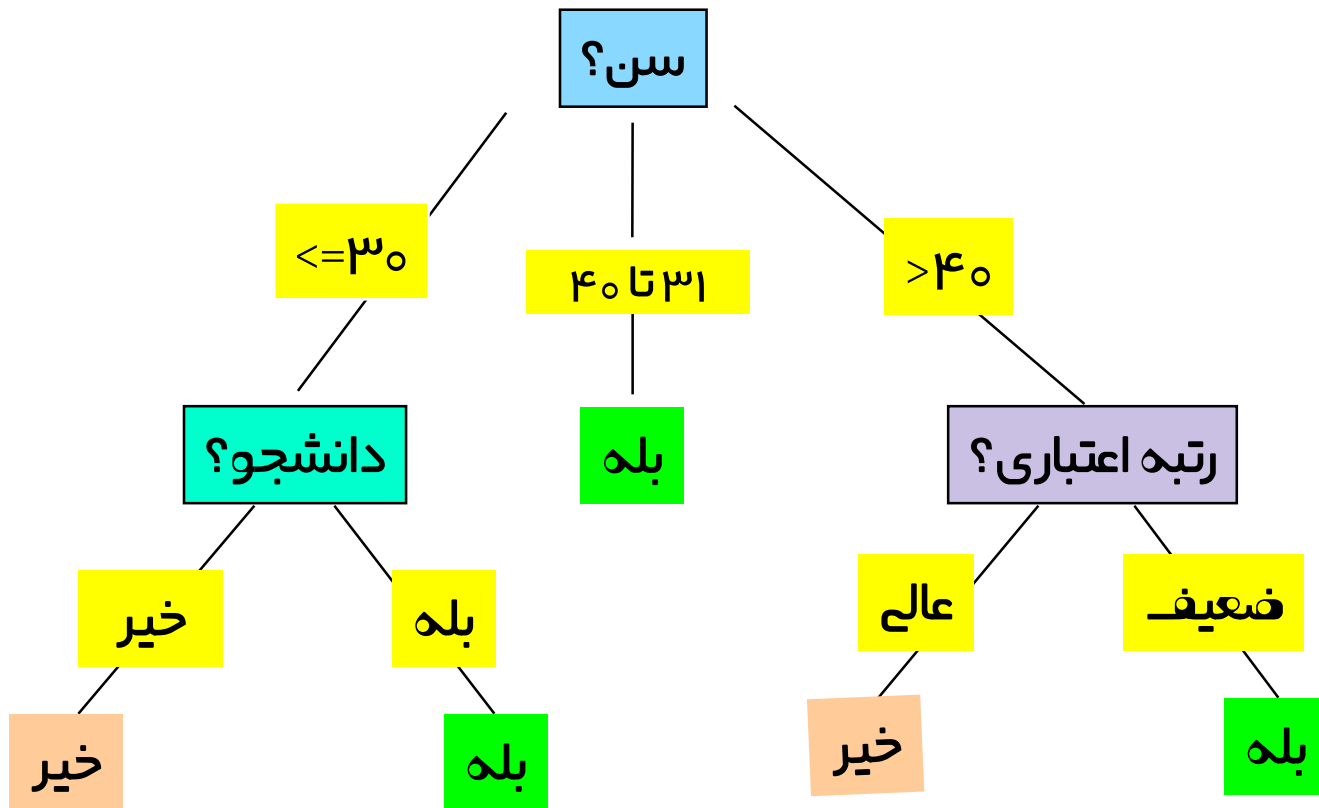
❖ عدم تطابق قوانین با نمونه جدید

استخراج قوانین از درخت تصمیم

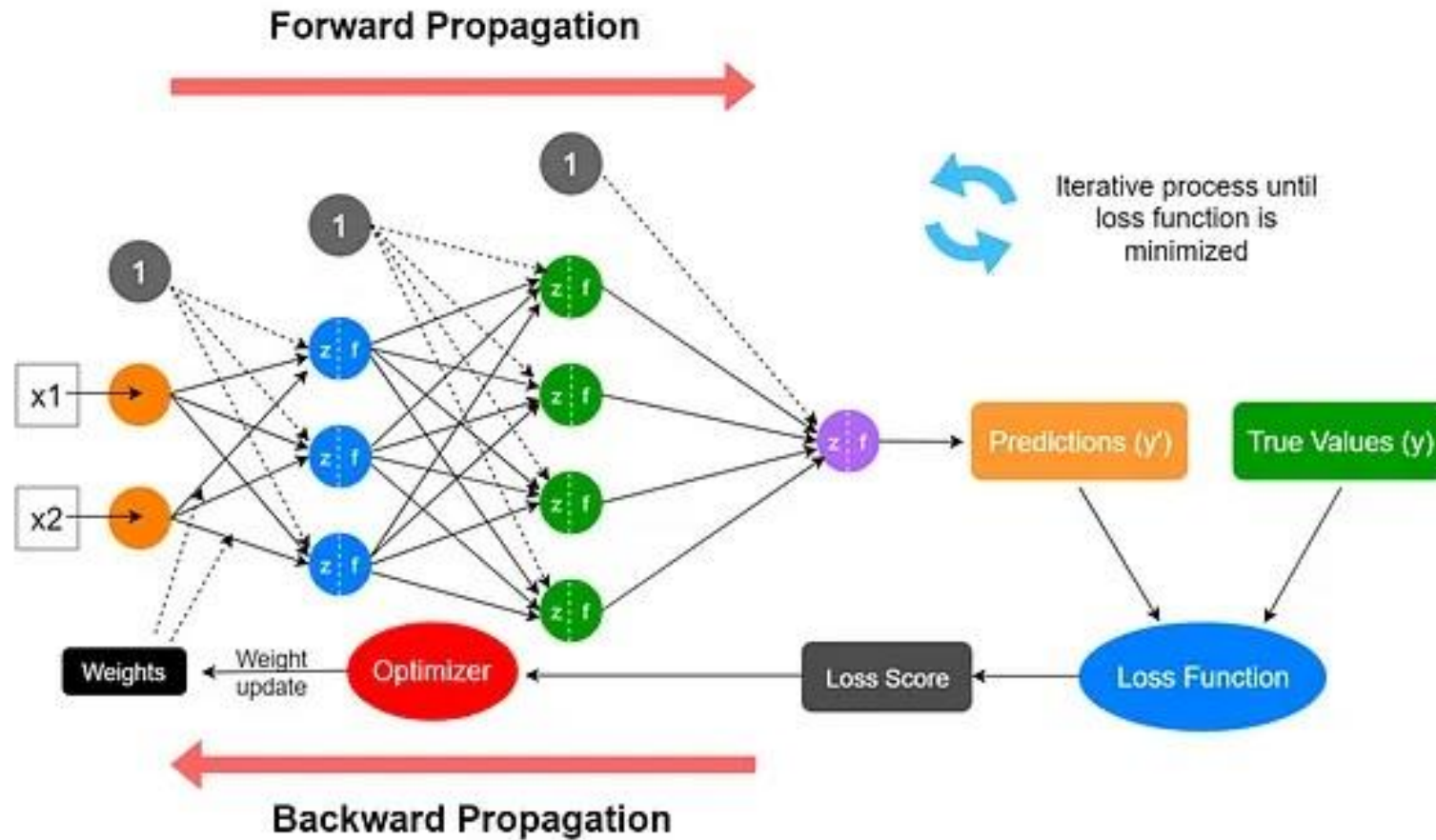
نمونه‌هایی از قوانین:

اگر سن کمتر از ۳۰ و دانشجو نبود، در نتیجه کامپیوتر نهی خرد
 اگر سن کمتر از ۳۰ و دانشجو بود، در نتیجه کامپیوتر می‌خرد
 اگر سن بین ۳۱ تا ۴۰ بود، در نتیجه کامپیوتر می‌خرد

استخراج قوانین از قوانین انجمنی



رده بندی با شبکه‌های عصبی



- زمان آموزش بالا
- طراحی مدل
- عدم تفسیر مدل ایجاد شده

- مقاوم در برابر نویز
- استفاده در اکثر کاربردها

منبع تصویر: <https://miro.medium.com>

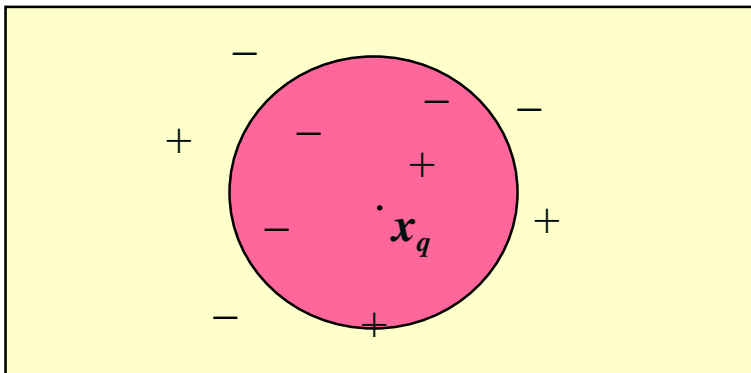
رده بندی مبتنی بر نزدیک‌ترین همسایه‌ها یا k -Nearest Neighbor

➤ یک روش lazy است

➤ نزدیک‌تر بودن - وزن بیشتر

➤ انتخاب k مناسب

➤ مقاوم در برابر نویز



رده بندی رگرسیون

○ هدف و کاربرد

○ ارتباط بین متغیر و پاسخ

ماتریس در هم ریختگی (confusion)

مقدار پیش بین شده/مقدار واقعی	(P) C_1	(N) $\sim C_1$
(P) C_1	True Positives (TP)	False Negatives (FN)
(N) $\sim C_1$	False Positives (FP)	True Negatives (TN)

TP : کلاس مثبت و به درستی پیش بین شده
 FP : کلاس منفی و به اشتباه پیش بین شده
 TN : کلاس منفی و به درستی پیش بین شده
 FN : کلاس مثبت و به اشتباه پیش بین شده

مثال: آیا فرد مبتلا به کرونا است؟

	بله	خیر
بله	TP(71)	FN(4)
خیر	FP(10)	TN(15)

ارزیابی

توانایی مدل در پیش بین درست

❖ دقت

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

خطاهای مدل در پیش بین

❖ نرخ خطا

$$error\ rate = \frac{FP + FN}{TP + TN + FP + FN}$$

○ اگر داده‌ها بالانس نباشند؟

ماتریس در هم ریختگی (confusion)

مقدار پیش بین شده/مقدار واقعی	(P) C_1	(N) $\sim C_1$
(P) C_1	True Positives (TP)	False Negatives (FN)
(N) $\sim C_1$	False Positives (FP)	True Negatives (TN)

TP : کلاس مثبت و به درستی پیش بین شده
 FP : کلاس منفی و به اشتباه پیش بین شده
 TN : کلاس منفی و به درستی پیش بین شده
 FN : کلاس مثبت و به اشتباه پیش بین شده

مثال: آیا فرد مبتلا به کرونا است؟

	بله	خیر
بله	TP(71)	FN(4)
خیر	FP(10)	TN(15)

ارزیابی

موارد مثبتی که مدل درست پیش بین کرده است

❖ حساسیت

$$sensitivity = \frac{TP}{P}$$

موارد منفی که مدل درست پیش بین کرده است

❖ تشخیص پذیری (خاصیت)

$$specificity = \frac{TN}{N}$$

ماتریس در هم ریختگی (confusion)

مقدار پیش بین شده/مقدار واقعی	(P) C_1	(N) $\sim C_1$
(P) C_1	True Positives (TP)	False Negatives (FN)
(N) $\sim C_1$	False Positives (FP)	True Negatives (TN)

TP : کلاس مثبت و به درستی پیش بین شده
 FP : کلاس منفی و به اشتباه پیش بین شده
 TN : کلاس منفی و به درستی پیش بین شده
 FN : کلاس مثبت و به اشتباه پیش بین شده

مثال: آیا فرد مبتلا به کرونا است؟

	بله	خیر
بله	TP(71)	FN(4)
خیر	FP(10)	TN(15)

ارزیابی

درست بودن نتیجه، وقتی مثبت پیش بین می شود

صحت ✦

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

پوشش مدل در پیش بین

پوشش ✦

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

نحوه ارزیابی مدل

❖ روش Holdout

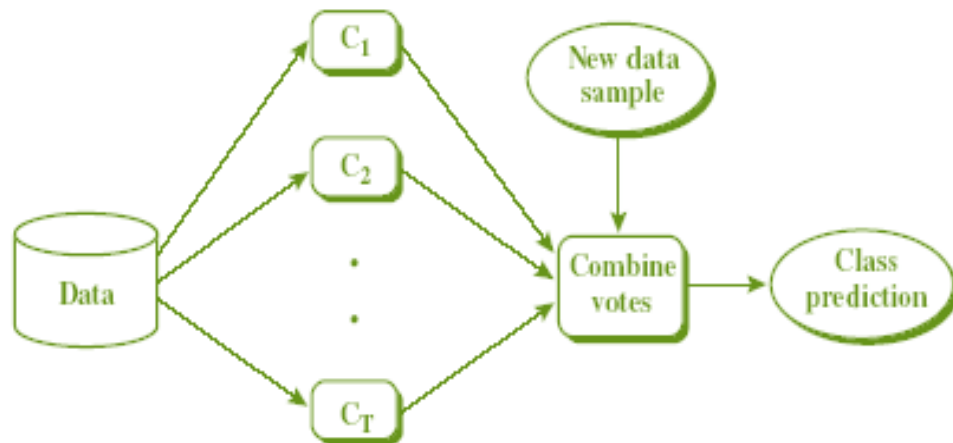
- به نسبت مشخص
- به صورت تصادفی بدون برگرداندن

❖ روش Cross-validation

- K پارتیشن
- برچسب اختصاصی در یک پارتیشن

❖ روش Bootstrap



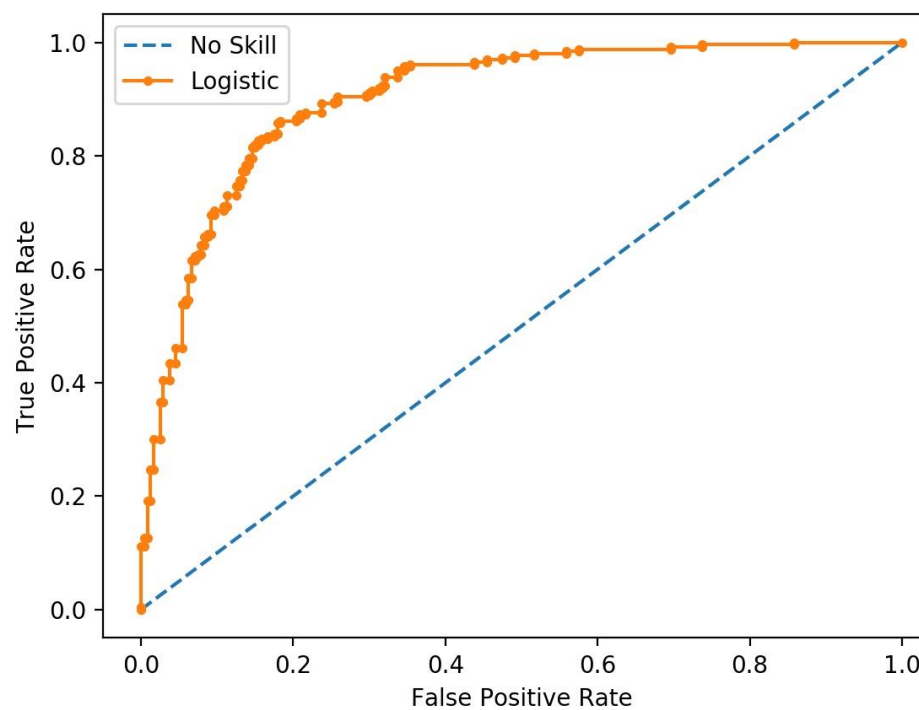


روش گروهی یا Ensemble

❖ نظر چند رده بند

❖ Bagging

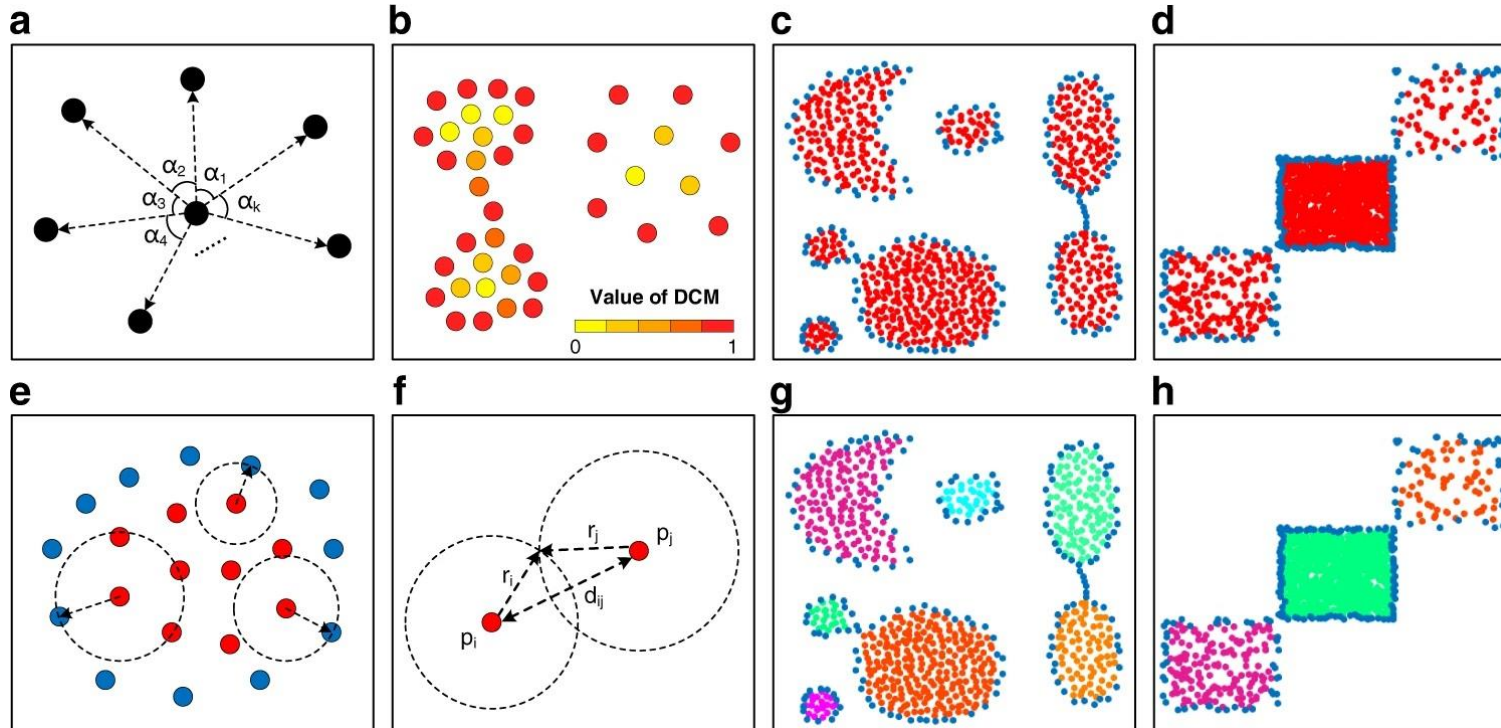
❖ Boosting



نمودار ROC

❖ سطح زیر نمودار

خوشه بندی



❖ هدف

❖ به عنوان تسک (بدون برچسب)

❖ به عنوان پیش پردازش

❖ کاربرد

❖ شباهت = فاصله

❖ خوشه بندی خوب

منبع تصویر: مقاله

Clustering by measuring local direction centrality for data with heterogeneous density and weak connectivity

		Object j		
		1	0	sum
Object i	1	a	b	$a+b$
	0	c	d	$c+d$
sum		$a+c$	$b+d$	p

شباهت یا Similarity

❖ بر اساس نوع فیلد
 باینری
 اسفے
 عددی
 ترتیبی

Object ID	Test-1 (categorical)
1	Code-A
2	Code-B
3	Code-C
4	Code-A

اسفے

Object ID	Test-1 (ordinal)
1	Excellent
2	Fair
3	Good
4	Excellent

ترتیبی

شباهت خوشه‌ها

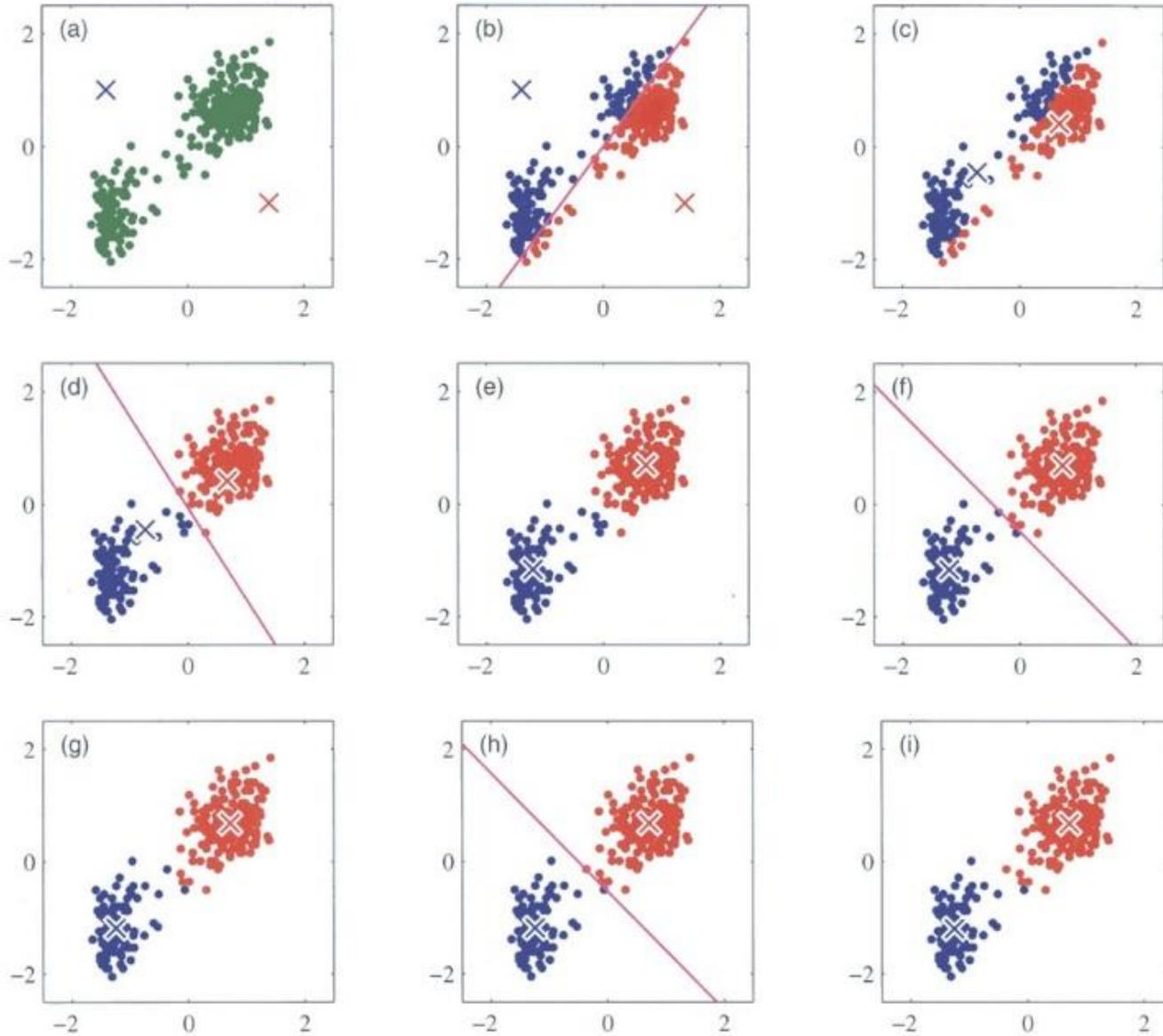
Single Link ❖ کمترین فاصله بین عناصر دو خوشه

Complete Link ❖ بیشترین فاصله بین عناصر دو خوشه

Average ❖ متوسط فاصله بین عناصر دو خوشه

Centroid ❖ فاصله مرکزهای خوشه

Medoid ❖ فاصله عناصر مرکزی خوشه‌ها



منبع تصویر: dendroid.sk

خوشه بندی – پارتیشنینگ

k-means ❖

k-medoids ❖

CLARA ❖

الگوریتم k-means

۱. انتخاب k مرکز و خوشه بندی
۲. پیدا کردن مرکز جدید و بررسی جابه جایی
آبجکت‌ها
۳. تکرار مرحله دو، تا زمانی که آبجکت‌ها
جابه جا نشود

معایب الگوریتم k-means

- اگر خصیصه عددی نبود
- مشخص نمودن k یا تعداد خوشه‌ها
- داده‌های نویز و پرت

سطح

l = 0

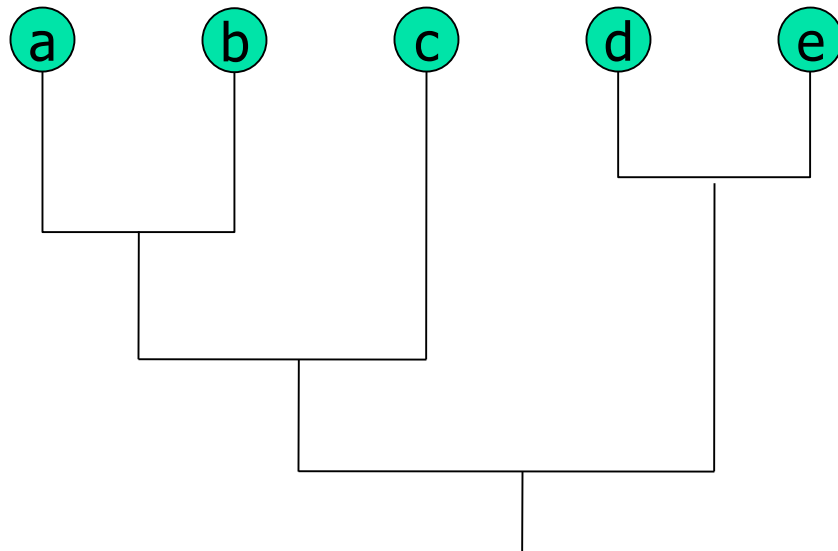
l = 1

l = 2

l = 3

l = 4

Dendrogram



شباهت

1.0

0.8

0.6

0.4

0.2

0.0

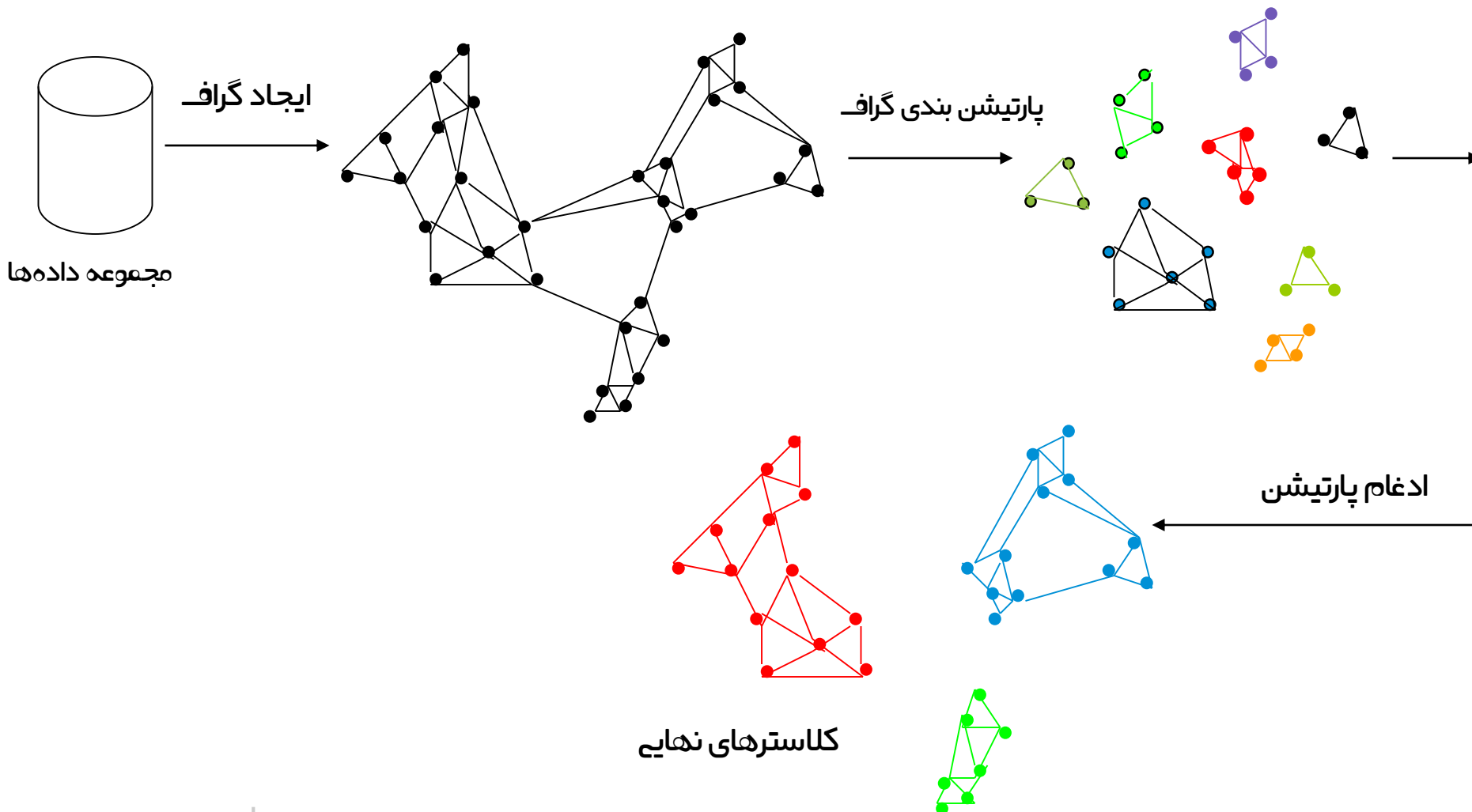
خوشه بندی – سلسله مراتبی

- عدم نیاز به مشخص کردن k
- نیاز به شرط خاتمه
- پایین به بالا – بالا به پایین

CHAMELEON – ROCK – BIRCH – AGNES – DIANA

الگوریتم CHAMELEON

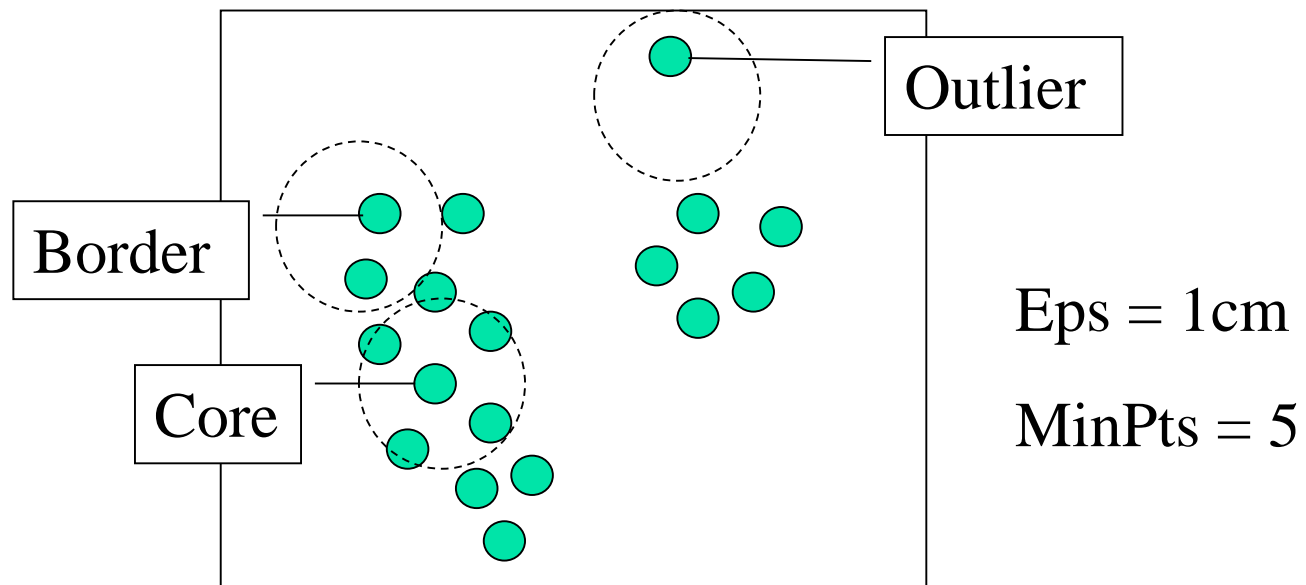
سلسله مراتبی ❖



خوشه بندی - مبتنی بر چگالی

❖ استفاده از چگالی نقاط و همسایگی

❖ حساس به پارامترها



DBSCAN

OPTICS

DENCLUE

CLIQUE

خوشه بندی – متناسب با محدودیت ها

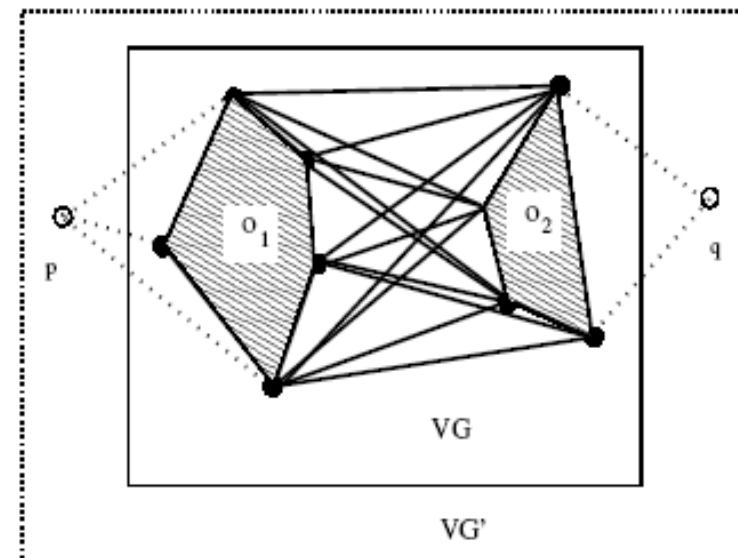
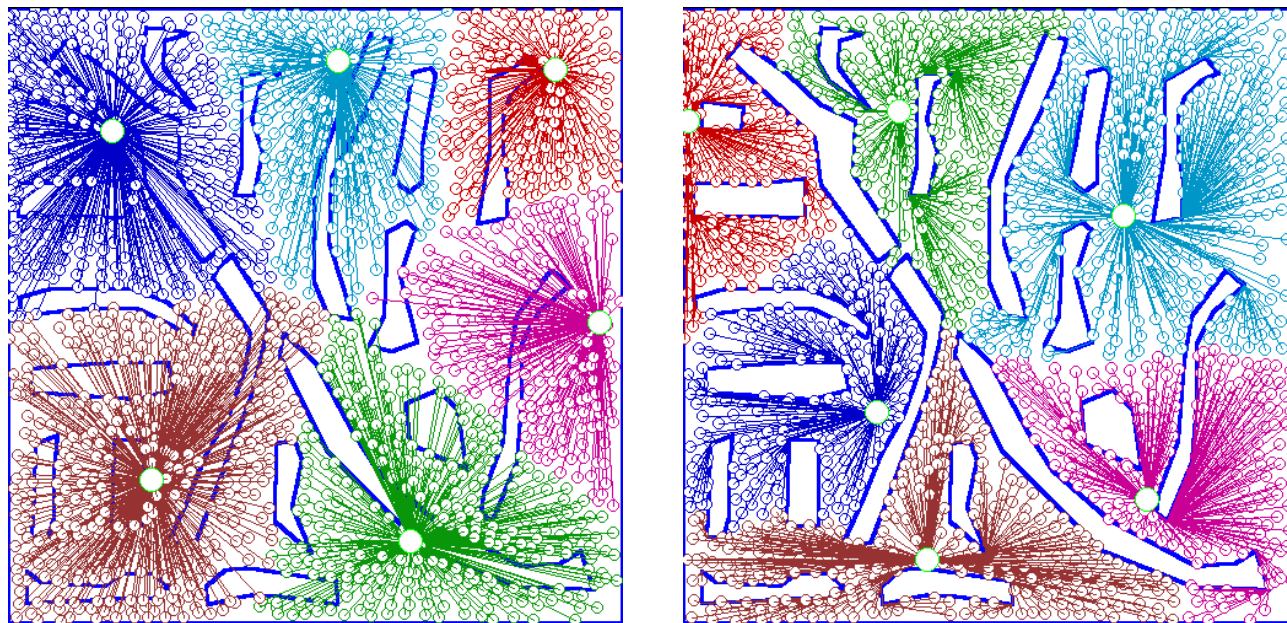
COP-k-means

CVQE

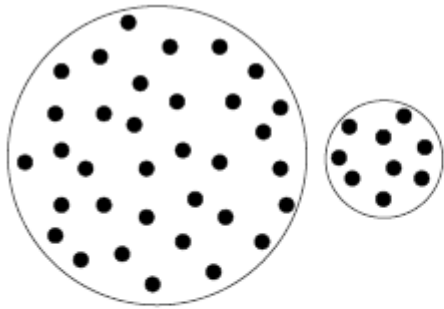
❖ بر اساس دانش کاربر

❖ مثال: کمینه کردن فاصله شهروندان از دستگاه های ATM

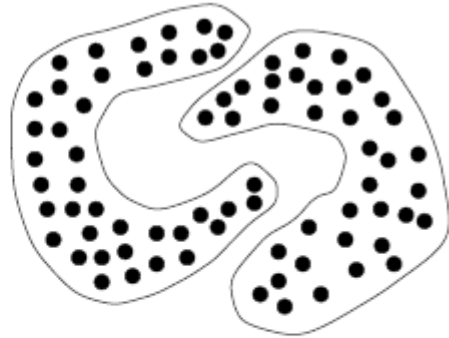
❖ سخت و نرم



روش خوشه بندی برتر؟



a) Clusters of widely different sizes

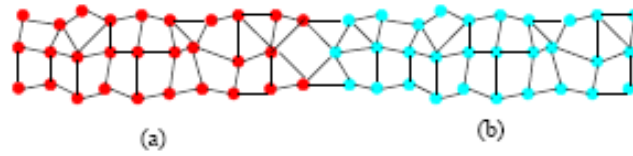


b) Clusters with convex shapes

(۱)

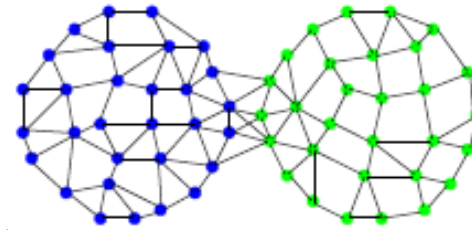
Data sets on which centroid and medoid approaches fail.

(۲)



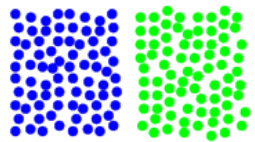
(a)

(b)



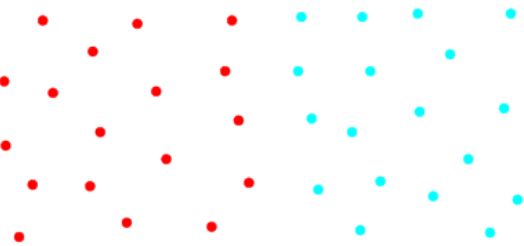
(c)

(d)



(a)

(b)



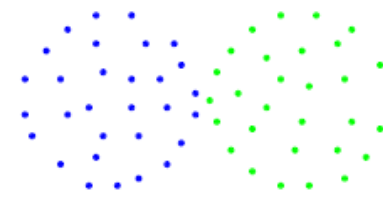
(c)

(d)



(a)

(b)



(c)

(d)

(۳)

(۴)

ارزیابی خوشه‌ها

❖ داده‌ها خوشه پذیر هستند؟

❖ مجموعه استاندارد

❖ تست Hopkins

w : فاصله آجکت با همسایگان نزدیک

q : فاصله نزدیکترین آجکت‌ها با نقاط انتخاب شده

$$H = \frac{\sum_{i=1}^n q_i}{\sum_{i=1}^n q_i + \sum_{i=1}^n w_i}$$

❖ اگر H کوچک باشد و نزدیک به 0.5 : یکنواخت

❖ اگر H نزدیک به یک باشد: ساختار خوشه بندی

ارزیابی خوشه‌ها

❖ تعیین تعداد خوشه‌ها

بر اساس دانش قبلی

بر اساس تجربه

روش Elbow

❖ اندازه گیری کیفیت خوشه‌ها

✓ همگن بودن

✓ کامل بودن

✓ Rag bag

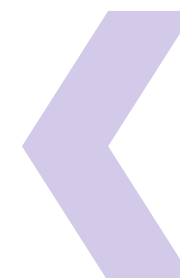
✓ خوشه‌های کوچک

BCubed precision

BCubed recall

روش‌های بررسی نزدیکه

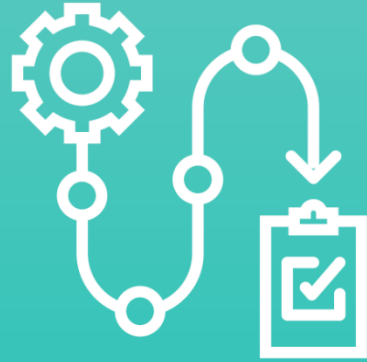
آیتم‌ها داخل خوشه‌ها



✓ کشف الگو

✓ رده بندی و پیش بینی

✓ خوشه بندی



اجرای پروژه داده گاوی

- معرفی چند نرم افزار کاربردی
- فرآیند CRIPS
- پروژه و ارزیابی
- پروژه عملی

معرفی نرم افزارهای کاربردی



نرم افزار Weka



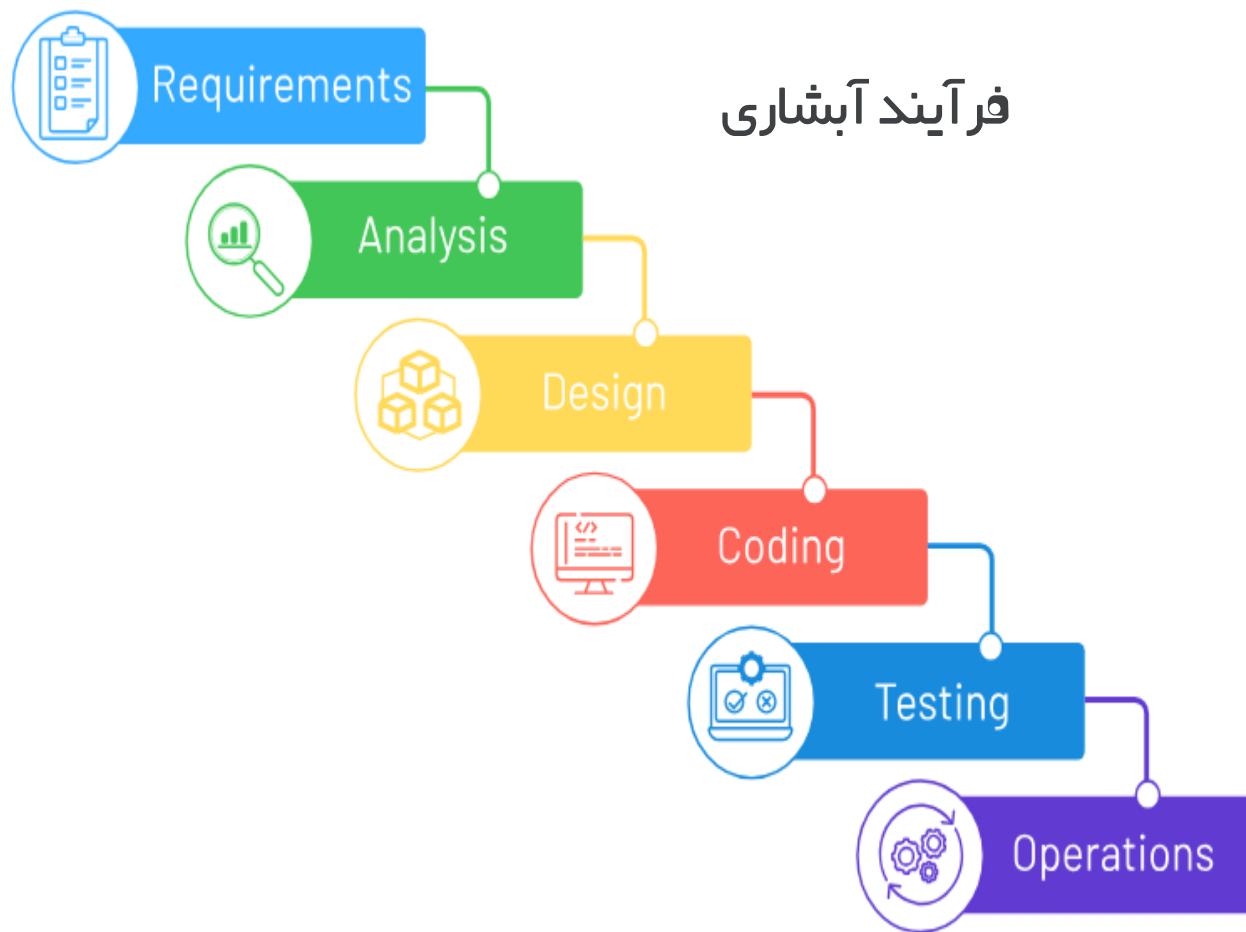
نرم افزار RapidMiner



نرم افزار Orange



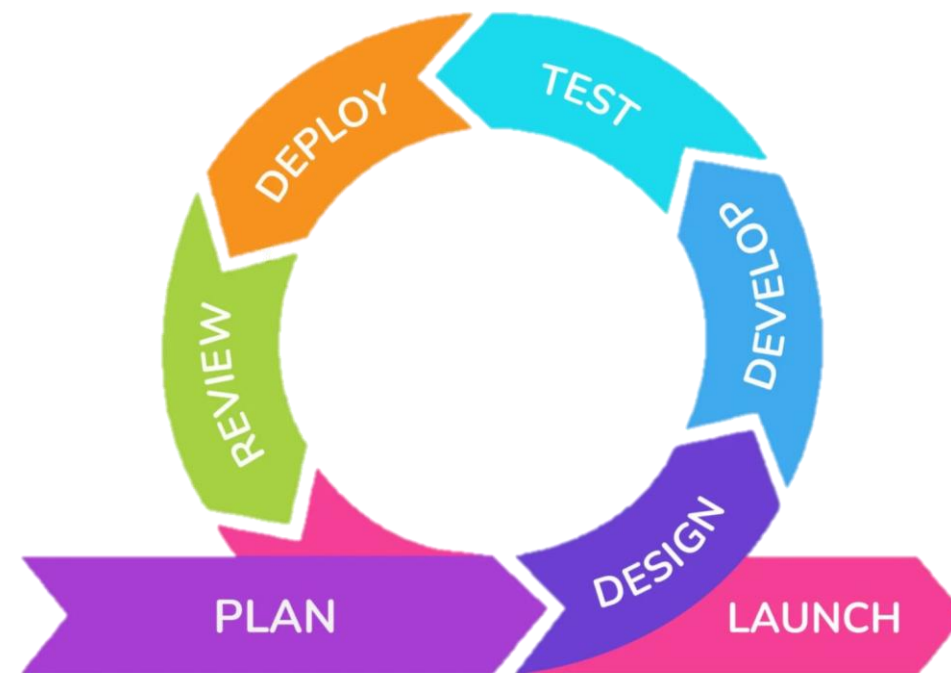
نرم افزار KNIME



منبع تصویر: <https://www.actitime.com>

فرآیند کریسپ CRIPS

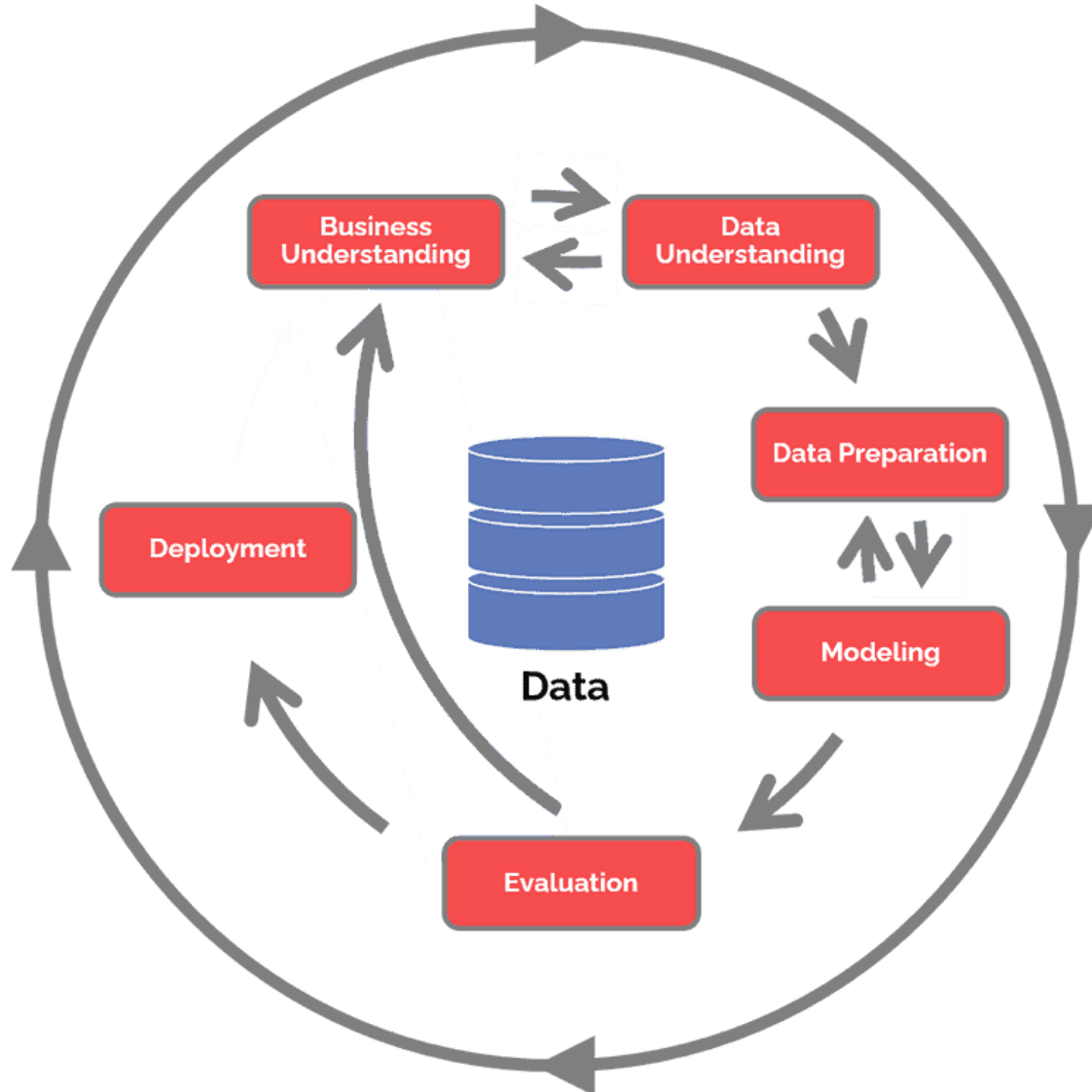
فرآیند چابک



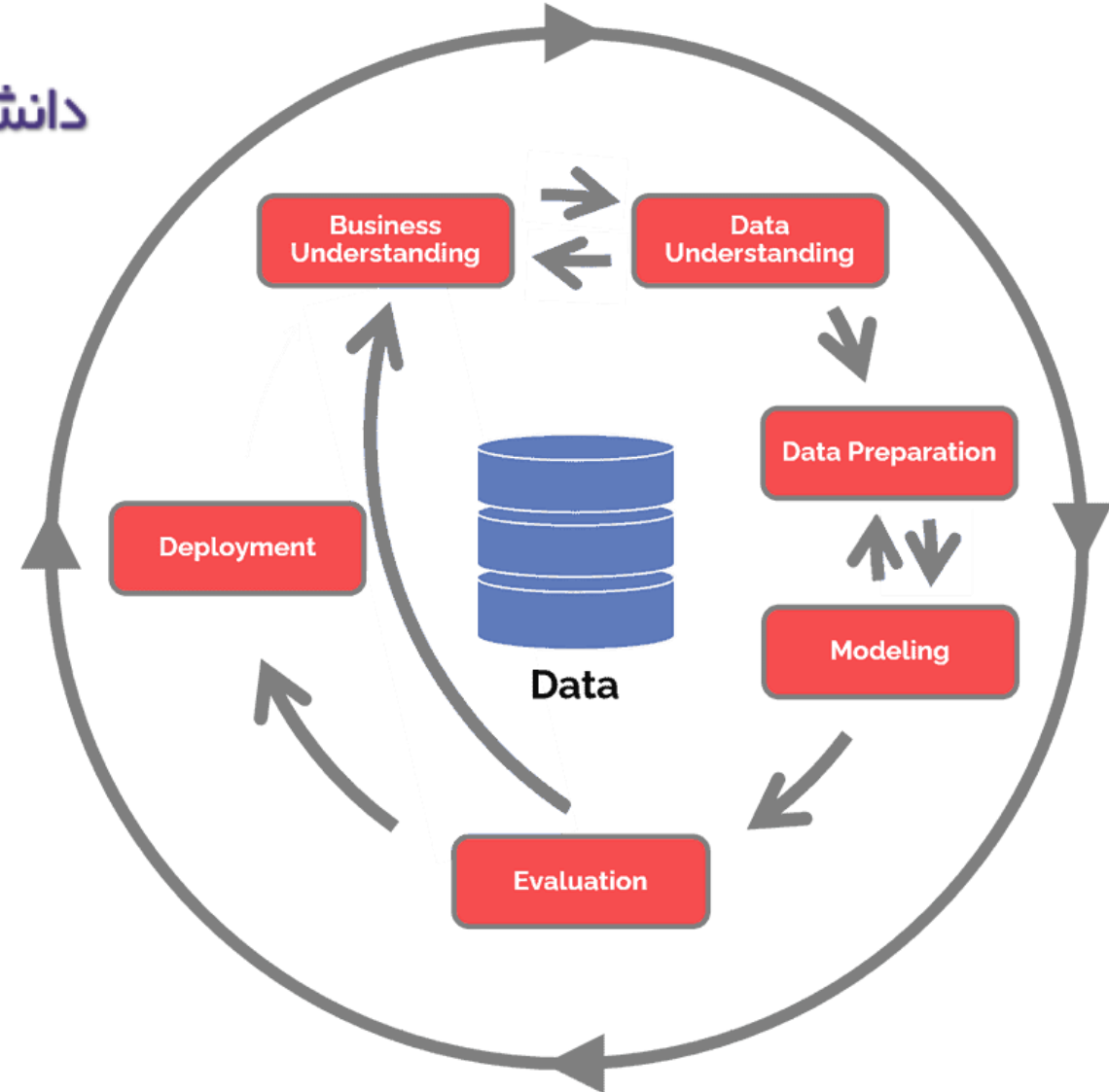
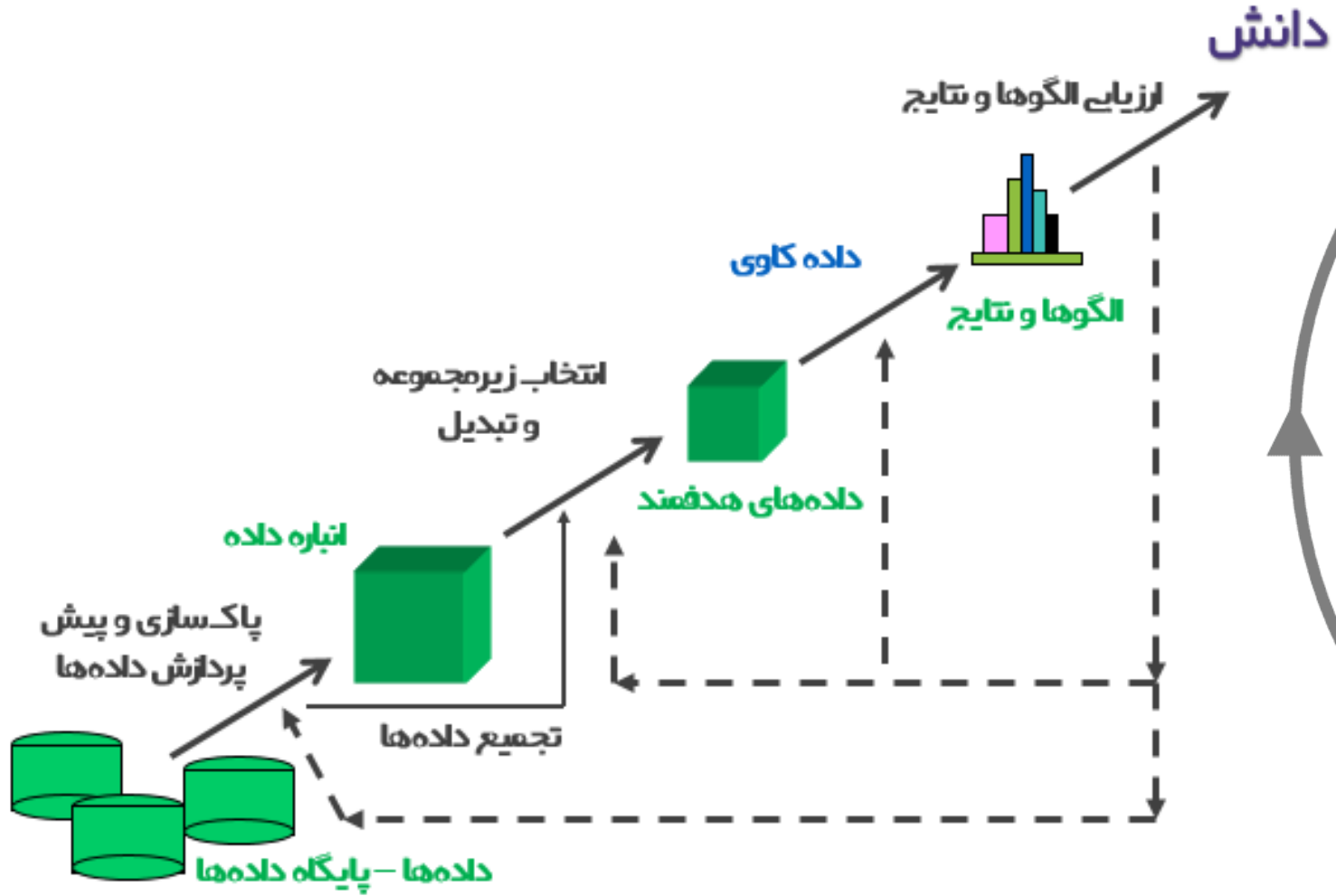
منبع تصویر: <https://www.krasamo.com>

فرآیند کریسپ CRIPS

۱. فهم کسب و کار (Business Understanding)
۲. فهم داده (Data Understanding)
۳. آماده سازی داده (Data Preparation)
۴. مدل سازی (Modeling)
۵. ارزیابی (Evaluation)
۶. پیاده سازی و انتشار (Deployment)



منبع تصویر: <https://www.datascience-pm.com>



پروژه و ارزیابی

۱. دوبخش (حضور و مشارکت (۲۰ درصد) + **انجام پروژه (۸۰ درصد)**)
۲. بر اساس کریسپ، مرحله یک و دو را مستند نمایید و شرح دهید.
۳. داده‌های کثیف را در داده‌های در دسترس شناسایی و راهکاری برای تمیز شدن آن‌ها ارائه نمایید.
۴. پیش پردازش‌های لازم، جمع‌ها و تبدیل‌های مورد نیاز را همراه با علت شرح دهید.
۵. ده کاربرد از داده‌های در دسترس به همراه نحوه‌ی اجرا یا ایجاد مدل را شرح دهید.
۶. **امتیازی:** تولید مدل و ارزیابی مدل متناسب به نیازمندی

مهلت ارسال: ۳۰ دی ماه ۱۴۰۲ **ارسال به ایمیل:** m.amin@mollahoseini.ir

پروژه عملی

❖ **موضوع:** بررسی و پردازش داده‌ها با استفاده از نرم افزار رپیدماینر

دانلود و نصب نرم افزار رپیدماینر

<https://download.ir/rapidminer-studio-developer/>

بررسی کامپوننت‌ها در نرم افزار

پایان
مشکرہ

