

انواع داده } Quantitative (کمی) : قابل شمارش و اندازه گیری  
به صورت عدد بیان می شود

Qualitative (کیفی) : Categorical در دسته های  
تخلف قناری گیرند

معمولا فا صله و ترتیب ندارند  
اگر ترتیب داشته باشند (Ordinal)  
اگر ترتیب نداشته باشند (صوری-اسمی)  
(Nominal)

آماره: برای خلاصه سازی داده ها و یا بدست آوردن یک انگلی خاص از روابط آماری

استفاده می شود. > و نوع شاخص آماری وجود دارد } مرکزی  
پراکنش

$$\mu_A = \frac{1}{n} \sum_{i=1}^n x_i$$

حسابی } میانگین

$$\mu_G = \left( \prod_{i=1}^n x_i \right)^{\frac{1}{n}}$$

هندسی

$$\mu_H = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$$

هارمونیک

$$\mu_H \leq \mu_G \leq \mu_A$$

$$\mu_W = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

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

Date:



8

میانه را بدون صورت کردن بدست آورید.

Max - Min

دامنه تغییرات

پراکندگی

واریانس و انحراف معیار

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_A)^2$$

$$s = \sqrt{s^2} \quad \text{standard deviation}$$

کواریانس: ارتباط میان دو ویژگی را بیان می کند.

$$\text{COV}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{Ax}) (y_i - \mu_{By})$$

$$\text{COV}(x, y) = E(x \cdot y) - \bar{x} \bar{y}$$

$$\text{CORR}(x, y) = \frac{\text{COV}(x, y)}{\sqrt{\text{COV}(x, x) + \text{COV}(y, y)}}$$

همبستگی

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

$$TP + TN + FP + FN$$

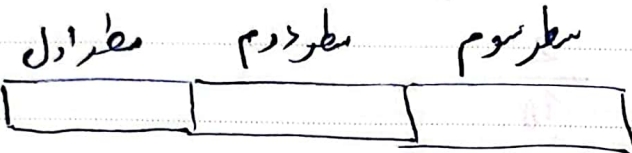
$Covmat(i, j) =$

ماتریس کوواریانس

$Cov(Feature_i, Feature_j)$

Classification

① مستقیم سطرها را کنار هم بچینیم تا آرایه شود بدین به



هم بچینیم تا آرایه شود بدین به

Nearest Neighborhood

② کوواریانس سطرهای تصویر با خودش را محاسبه

کنیم سپس Nearest-Neighbor می دهیم

③ Color histogram  $\Rightarrow$  Nearest-Neighbor

N سطر - ۲۵۵ ویژگی

test hold out - Cross-validation -

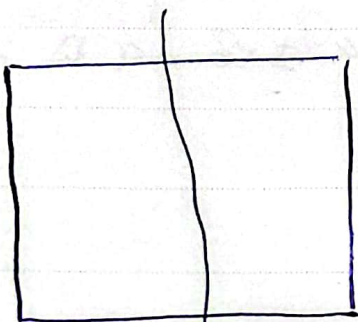
train

K-fold

leave one out - hold out

برای داده های کم هر بار یکی را برای test جدا می کنند. perform metrics for Evaluating Classifier  $TP, FN, FP, TN$

$$Accuracy = \frac{\text{true classified}}{\text{num all data}} = \frac{TP + TN}{P + N}$$



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

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

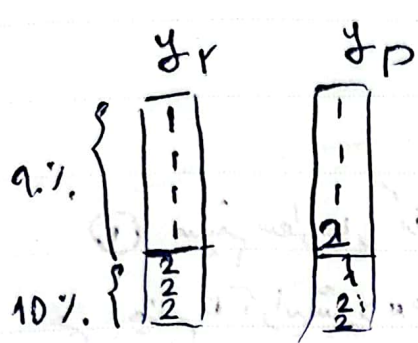


Date:

Sub:

F1-score

برای مسائل چندکلاس به هم را + در نظر بگیرند بقیه داده ها همه را - در نظر

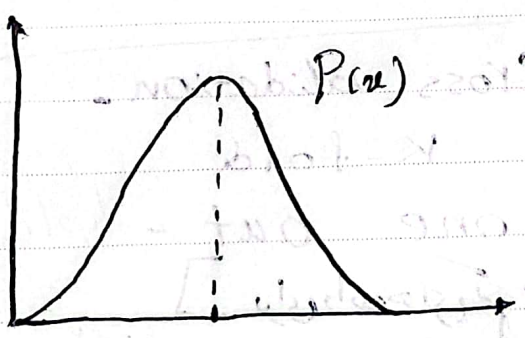


precision =  $\frac{2}{3}$

recall =  $\frac{2}{10}$

میانگین

F1 score



$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\mu \pm \sigma$  68%

$\mu \pm 2\sigma$  95%

Date:

Sub:

$$\left. \begin{aligned} \mu &= 0 \\ \sigma &= 1 \end{aligned} \right\} \text{توزیع استاندارد}$$

توزیع نرمال ← استاندارد داده‌ها منهای میانگین تقسیم بر اختلاف معیار

کیم انحراف معیار نیز برابر می‌شود.

$$y = \alpha + \beta x$$

رگرسیون خطی

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

میانگین  $\bar{y}$  های  $\bar{x}$  داریم  
(مثلاً نقطه)

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

## principal Component Analysis

PCA: یک روش کاهش ابعاد داده است.

مثلاً  $X$  داریم که شامل  $n$  داده هست که هر کدام  $D$  بعد (مثلاً ۱۰۰ بعد)

$$X_{n \times D} \times W_{D \times d} = X_{new} \ n \times d \text{ دارند}$$



Date:

Sub:

الگو رسم:

واده حارا ارسال می کنیم چگونه؟

Linear Scaling:

$$\textcircled{1} \quad \frac{x_i - \min(X)}{\max(X) - \min(X)} \times 2 - 1$$

$$\textcircled{2} \quad \begin{matrix} \neq \\ \text{normalized} \end{matrix} \quad x_i = \frac{x_i - \mu}{\sigma}$$

ماتریس کوواریانس و ویژگی حارا حساب می کنیم

بر دار ویژه و مقدار ویژه C را بدست می آوریم

$$x C \times E_{\text{vector}} = \lambda \alpha E_{\text{vector}}$$

بردار ویژه هارا بر اساس مقدار ویژه  $\lambda$  مرتب می کنیم

$$\lambda_1 > \lambda_2 > \lambda_3$$

$$x \quad \left| \lambda I - C \right| = 0$$

←  
دترمینان

چند جلدی درجه D - D ریشه دارد

هر مقدار ویژه  $\lambda_1$  را در رابطه  $x$  قرار می دهیم  $D \times x$  بردار ویژه  $E_{\text{vector}}$  بدست می آید

بردارهای ویژه بدست می آیند (متعامد هستند)

Date:

Sub:

$$W_{d \times D}$$

برای هر  $\lambda$  یک Eigenvalue داریم که  $D$  استند هستند

$d$  تا از آنها را انتخاب می‌کنیم و  $W$  را قرار می‌دهیم در ستون اول  $\lambda$  با بیشترین مقدار را قرار می‌دهیم.

$d$  مناسب برای هر مسئله چقدر باشد خوب است؟

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^D \lambda_i} \geq 0.95$$

کوچکترین  $k$  که در معادله صدق کند

توزیع در ابعاد جدید 95 درصد الگوریتمی که حفظ

PCA یک روش کاهش بعد خطی است.

روش‌های دسته بندی

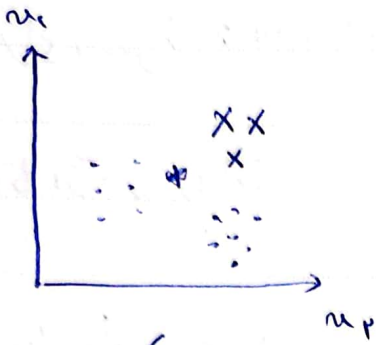
① K-Nearest-Neighbor -  $\epsilon$ -nearest-neighbor

② support vector machine (SVM)

Date:

Sub:

روش های دسته بندی



• کلاس ۱

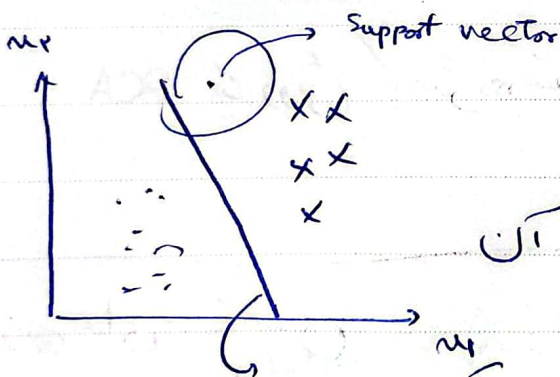
X کلاس ۲

\* ست

برای دسته بندی → کلاس بندی های → دسته بندی

در فاصله  $\epsilon$  همایی را اندازه گیری می کنند

\* برای دسته بندی های Binary ← دسته های زوج انتخاب نمی شود چرا که طبق است  
۵۰ ۵۰



: Support vector machine

در SVM فضا را بهترین سطح margin آن دسته بندی!

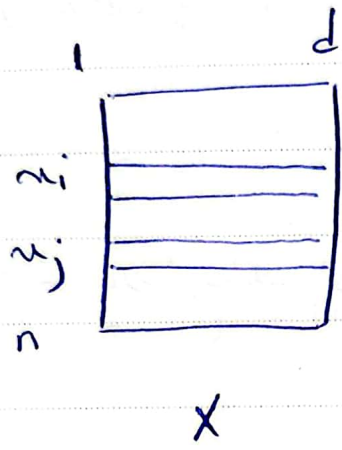
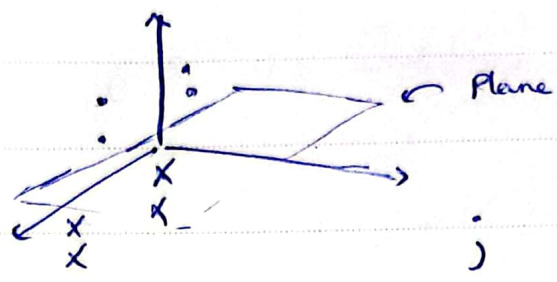
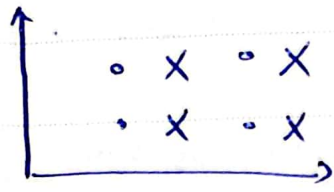
Sum  
فاصله بین کلاس  
و علامه فضا را می بیند

نقطه جدی آن Hyper plane نام دارد

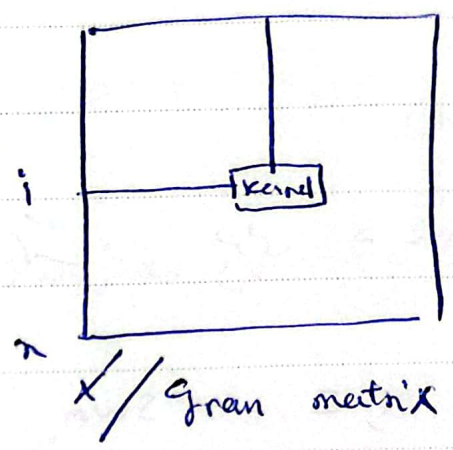
با استفاده از دسته های جدی آن را می بیند و اگر بزرگی این داده ها در بزرگی margin  
همه بزرگ و فاصله دومی می کند



برای همسازان آن برداری آن kernel میزنه در  $\phi$   
 آن با همسازهای وضع



فضای صغیره



$$\text{kernel}(\vec{u}_i, \vec{u}_j) \Rightarrow X'$$

شبهت ۲ بردار را در  
 کلاستر

مکن است همواره فضا را درسته برپا کنه یک لب  $\epsilon$  margin  
 های وضع در آن فضا درسته استوار نزلدر

شبهت ۲ بردار صفت ؟ ( kernel شبهت ۲ بردار )

- kernel
- ① linear  $\rightarrow \vec{a} \cdot \vec{b} = |a| \cdot |b| \cos \phi$
  - ② Poly  $\rightarrow (ax^2 + b + c)^d \rightarrow$  عموماً ۲ عموماً
  - ③ RBF  $\rightarrow \exp\left(\frac{\|a-b\|^2}{2\sigma^2}\right)$
- شبهت بیار  $\rightarrow \sigma$   
 کلاسه شبهت  $\rightarrow a$   
 حافظه  $\rightarrow \sigma^2$

عموماً جواب های دیگری هم هست

عموماً

Weighted k nn

Σ

برای دسته بندی مطالب بیشتر از ۲ مطالب ، دسته ها بیشتر

! one vs one

! one vs other

میان دسته ها ، دسته ها  $m$  مطالب داشته باشند ، برای  $m=5$

در زیر جدول

1, 2 → SVM

1, 3 → SVM

2, 3 → SVM

1, 4 → SVM

2, 4 → SVM

3, 4 → SVM

1, 5 → SVM

2, 5 → SVM

3, 5 → SVM

4, 5 → SVM

۱. SVM به ازای  $m$  مطالب

$m(m-1) \rightarrow$  Sum  
دارم

1 → +

2, 3, 4, 5 → -

روش دوم : تب SVM برای آن سافه آن را

Sum ای خارج

2 → +

1, 3, 4, 5 → -

در بعضی ترتیب پس شروع

3 → +

1, 2, 4, 5 → -

Sum  $m$  دارم



روش دسته بندی بیز Bayesian

$$P(C = C_i) \times \prod_{j=1}^n P(A_j = a_j | C = C_i)$$

به مقدار ویژگی‌ها (موتور)

مقدار احتمال وابسته مقابل را محاسبه می‌کنیم.

این رابطه را برای داده test روی هر کلاس  $C_i$  بدست می‌آید، در نهایت

برای کلاس  $C_i$  که بیشترین مقدار مقدار را دارد، داده test به آن

می‌شود

$A_1, A_2, \dots, A_n$  ویژگی‌ها

$\langle a_1, a_2, \dots, a_n \rangle$  مقادیر ویژگی‌های یک داده نمونه

شماره ردیف ویژگی‌ها

ID	A	B	C
----	---	---	---

برچسب کلاس

مثال:

1	a	d	y
2	a	e	y
3	b	f	y
4	c	e	y
5	b	f	y
6	b	f	n
7	b	e	n
8	c	d	n
9	n	f	n

test  $\langle a, f \rangle$

①

$$P(C = y) \times \prod_{j=1}^2 P(A_j = a_j | C = y) =$$

②

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



$$\textcircled{1} P(C=y) P(A=a | C=y) \times P(B=f | C=y) =$$

$$\frac{1}{2} \times \frac{2}{5} \times \frac{2}{5} = \frac{2}{25} = \frac{8}{100}$$

$$\textcircled{2} P(C=n) P(A=a | C=n) \times P(B=f | C=n) =$$

$$\frac{1}{2} \times \frac{1}{5} \times \frac{2}{5} = \frac{1}{25} = \frac{4}{100}$$

 $\frac{8}{100}$ 
 $\frac{4}{100}$ 

$$\frac{P(A|C)}{P(C=n)} = \frac{\frac{2}{5}}{\frac{1}{2}} = \frac{2}{5} \times \frac{2}{1} = \frac{4}{5}$$

چون خروجی واحدی داریم احتمال دقیق در نیاید  
 مستقلاً و ویرجی یکسان بر حسب متنایست  
 ویرجی ها کافی نیست

اثبات رابطه نیز در کلاس گفته شد

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

$$P(C=C_i | A_1=a_1, \dots, A_n=a_n)$$

$$P(A_1=a_1, \dots, A_n=a_n | C=C_i) P(C=C_i)$$

$$P(A_1=a_1, \dots, A_n=a_n)$$

برای یک کلاس  $C_i$  این مقادیر یکسان هستند آنرا از مخرج حذف کنیم

Date:

Sub:

$$\Rightarrow P(A_1 = a_1, \dots, A_n = a_n \mid C = C_i)$$

افراز کریم

$$P(A_1 = C_1 \mid A_2 = a_2, \dots, A_n = a_n, C = C_i) \times$$

$$P(A_2 = a_2, \dots, A_n = a_n \mid C_i)$$

مجموع  $A_n$  است

$$P(A_1 = a_1 \mid A_2 = a_2, \dots, A_n = a_n, C = C_i) =$$

$$P(A_1 = a_1 \mid C = C_i)$$

$$P(A_1 = a_1) \times \prod_{j=1}^n P(A_j = a_j \mid C = C_i) \Rightarrow$$

$$P(C = C_i) \times \prod_{j=1}^n P_{A_j}$$

glam matrix - sum -

normalization

روبرتی سازی

# درخت تقسیم Decision Tree

برای ساختن درخت چگونه عمل کنیم؟ فرض می‌کنیم داده‌ها Category مقدار داده

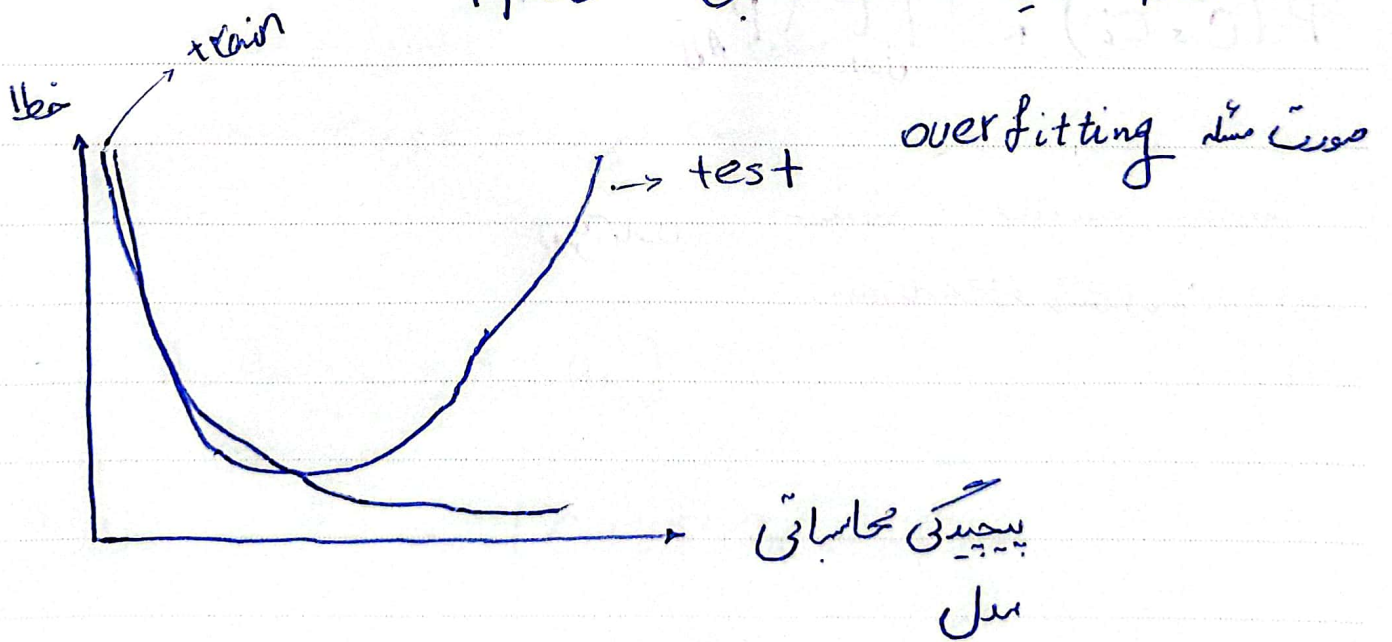
- شده اند یا کسبه هستند
- ① در هر مرحله چه ویژگی feature انتخاب شود.
  - ② محق درخت چقدر باید باشد؟

④ ۱- زمانی که داده‌های برگ در داخل  $k_{min}$  بر حسب یکسان داشته باشند

۲- برای محق یک حد بالا تعیین کنیم ۳- حد بالا برای تعداد گره‌ها تعیین کنیم

۴- تعداد داده‌های موجود در گره خیلی کم شود (۲ یا ۳ داده) به گسترش این گره‌ها

۵- وقت بررسی Validation متوقف می‌کنیم (این مقدار را از قبل مشخص می‌کنیم)





Date:

Sub:

referral  $\Rightarrow$  get user-id input referrals with  
~~show user-id~~ with depth show them  
+ num

add depth  $\Rightarrow$  admin panel

watch users

e.g.

$$T = RC$$

$$R = \frac{C}{T}$$

$$R = \frac{T}{C}$$

Date:

Sub:

برای رفع overfitting خود داده های train + را به دو قسمت تقسیم

می کنیم }  $train \ 70\%$   
 $validation \ 10\%$   
 $test \rightarrow 20\%$

هر بار قبل از افزایش پیچیدگی مدل  
 دقت مدل را بررسی validation  
 امتحان می کنیم اگر کاهش یافت آموزش  
 و افزایش پیچیدگی مناسبی مدل را متوقف می کنیم

برای مثال شرطی که برای عمق درخت می زنیم می تواند از overfitting جلوگیری کند.

x ————— x

### ① 1. information gain

$$IG(H) = Entropy(D) - Entropy_A(D)$$

$\downarrow$                        $\downarrow$   
 ویژگی H                      A

$$Entropy(D) = - \sum_{i=1}^C P_i \times \log_2 P_i$$

$$Entropy_A(D) = \sum_{j=1}^V \frac{|D_j|}{|D|} \times Entropy(D_j)$$

تعداد سطح های ویژگی A  $\rightarrow V$       تعداد کلاس  $\leftarrow C$

$P_i$   $\leftarrow$  احتمال نمونه ای از داده متعلق از داده خاصه ویژگی  $D_j$   
 متعلق به کلاس  $i$  باشد      A آنها مقدار  $V$  داشته باشد.  
 تعداد، اندازه  $|D|$



معیارهای انتخاب ویژگی :

1 - Information gain

2 - Gini index

3 - Gini ratio

4 - Likelihood ratio

معیار گره ریشه درخت تقسیم

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

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$

ویرجی (سئون) A

C تعداد لایس  $p_i$  احتمال وجود داده در لایس

در  $D_1$  و  $D_2$  مجدداً  $Gini(D_1)$  را محاسبه می کنیم  
 به دو مجموعه  $D_1$  و  $D_2$  تقسیم می شود  $A \rightarrow$   
 برای همین تعداد

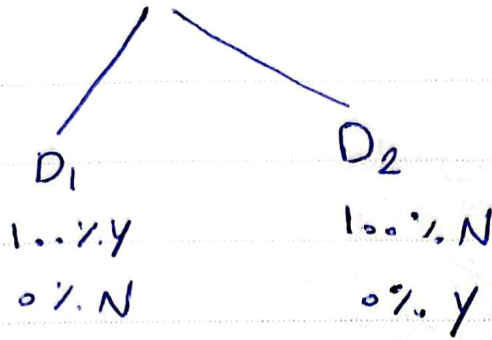
$$Gini(A) = Gini(D) - Gini_A(D)$$

ویرجی که این مقدار را  $\max$  کند  
 انتخاب می کنیم



Date:

Sub:

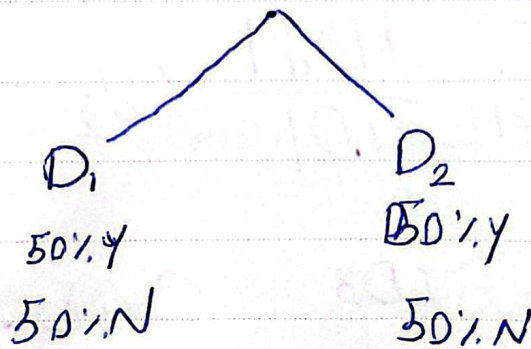


$$\alpha \text{Gini}(D_1) + (1-\alpha) \text{Gini}(D_2) = 50$$

$$\downarrow \qquad \qquad \qquad \downarrow$$

$$\frac{1 - (1^2 + 0^2)}{0} \qquad \qquad \qquad \frac{1 - (0^2 + 1^2)}{0}$$

اگر داده ها کاملاً از هم تفکیک شوند  $\text{Gini} = 0$  می شود.



$$\alpha \text{Gini}(D_1) + (1-\alpha) \text{Gini}(D_2) = 5$$

$$\downarrow \qquad \qquad \qquad \downarrow$$

$$\frac{1 - ((\frac{1}{2})^2 + (\frac{1}{2})^2)}{0.5} \qquad \qquad \qquad 0.5$$

عدد های  $\frac{1}{2}$

$$\frac{|D_1|}{|D_2|} \alpha \frac{1}{2} + \frac{|D_2|}{|D_1|} (1-\alpha) \frac{1}{2} = 5$$

$$\frac{1}{2} \left( \frac{|D_1| + |D_2|}{|D_1|} \right) = \frac{1}{2}$$

$D_1$  و  $D_2$  بر حسب های و اثری هستند مثلاً یکی از بر حسب های سن A را  $D_1$  و بقیه بر حسب هارا  $D_2$  را در نظری بگیریم این کار را برای تمامی

حالتها انجام می دهیم

Date:

Sub:

gain ratio  $\rightarrow$  info gain  
را انتخاب می کند

زیرا  $I_G$  در بعضی موارد اگر تعداد سطح های ویژگی زیاد باشد آنرا انتخاب می کند

$$\text{gain Ratio}(D) = \frac{\text{Info Gain}(A)}{\text{split Info}(A)}$$

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

likelihood ratio ID3 C4.5  $\rightarrow$  برای مقایسه

CART



# Random Forest جنگل تصادفی

overfit بر اساس وزن - معمولاً

Boosting

Bagging

stacking

روش های مختلف ترکیب دسته بندیها  
 }  
 Combination of  
 Classifiers

Bagging : دسته بندی داده ها به تعداد کسبه مجزا تقسیم می کنیم - هر داده های Train

Classifier ها داده نشود - (حذف ویژگی در Bag - حذف داده در Bag)

Random Forest نوعی Bagging است - برای ترکیب از میانگین استفاده می کنیم

stacking : نوع دسته بندی را عوض می کند مثلاً DT با عمق های مختلف یا

SVM با Kernel های مختلف استفاده می کنیم - pattern label

دسته بندیها را بدست می آوریم

	predict(D <sub>i</sub> )		
Classifier 1	1	0	1
Classifier 2	0	1	1
	1	1	1
Classifier 3	1	1	0
⋮			
10.0	0	0	1

Date: Classifier 1

خروجی

Sub:

Classifier 2

labels

1	0	1	1	0	1
0	1	1	1	0	0
1	1	1	0	1	1

X'

دوباره X' را به عنوان ویژگی‌های بدست آمده به یک

Classifier دیگر می‌دهیم تا آموزش ببیند pattern داده‌ها را به

صفر یا یک نگاه است کند. به این Classifier اصطلاحاً "meta-learner"

می‌گوئیم. می‌توانیم کنار ستون‌های خروجی n دسته‌بندی D ویژگی اولیه X را نیز تکرار

دهیم و بعد meta learner را آموزش دهیم.

شبکه عصبی :



چرا در صورت مسئله‌های چند کلاسه استفاده از precision, recall و F1

بتر از accuracy است و زیرا هر بار precision و recall هر کلاس را

در نظر بگیریم (هر بار کلاس هدف مثلا 1 یا 2 یا 3 و ... ) را به عنوان

Positive در نظر گرفته و سایر کلاس‌ها را - در نظر قرار می‌گیریم

استفاده Confusion matrix برای درک بهتر label هایی که به اشتباه تشخیص داده می‌شوند.

صورت مسئله Clustering به دو مدل است 1- تعداد cluster از

قبل مشخص باشد 2- روش Clustering خودش به تعداد مناسب Cluster ها

پی ببرد.

داده‌ها بدون برچسب هستند می‌خواهیم آنها را در چند طبقه تقسیم کنیم



Date:

K-means

Sub:

K-medoid

تعداد خوشه ها از قبل داده شده اند

روش های خوشه بندی

سلسله مراتبی

Top-down

bottom-up

DB scan

الگوریتم تعداد خوشه های مناسب را باید

الگوریتم

۱- به تعداد K مرکز خوشه تقاضی درباره ویژگی های داده ها تو لیدی کنیم (تقاضی)

۲- داده را بر اساس نزدیکی به این K مرکز خوشه برچسب گذاری می کنیم

۳- مراکز خوشه جدید را بر اساس میانگین برچسب های محاسبه محاسبه می کنیم

چون ممکن است نتیجه نسبت به K تعداد خوشه تقاضی متفاوت باشد یعنی ارتقا

مرحله ۱ را با شرط K تا از بدترین نقاط از یکدیگر اجزای کنیم

Top down: همه داده ها را یک خوشه در نظر می گیریم سپس split می کنیم

Bottom up: هر داده را یک خوشه می گیریم و سپس دو خوشه با کمترین فاصله را با هم

merge می کند



Date:

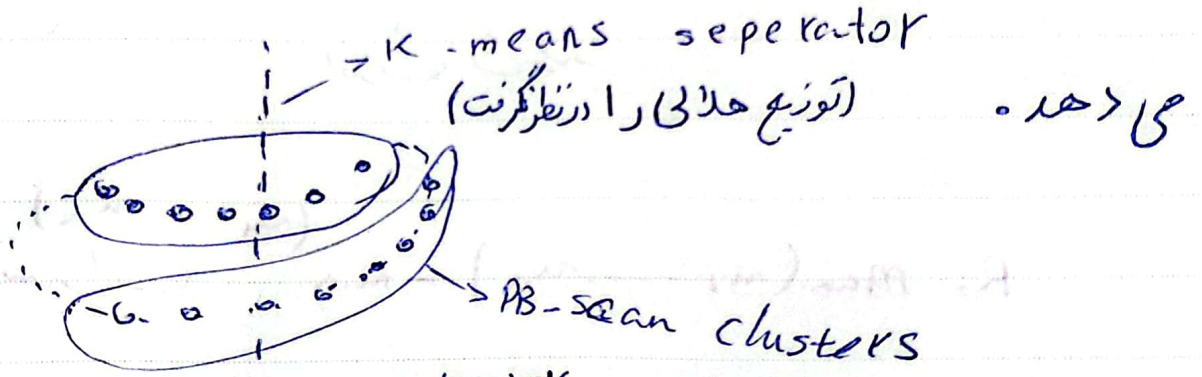
Sub:

top-down: Complete linkage  $\max d_B(a,b) \quad a \in A, b \in B$

bottom-up: single linkage  $\min d_B(a,b) \quad a \in A, b \in B$

dendrogram tree  $\rightarrow$  bottom-up درختی که  
 $\rightarrow$  top-down در این دوروس می سازیم و  
K را با استفاده از همین تعیین می کنیم

این روش ها توزیع داده ها را در نظر نمی گیرند اما روش DBSCAN این کار را انجام می دهد.



ورود n: حداقل تعداد داده در cluster

1- انتخاب راه ها در فاصله  $\epsilon$  از هم هم خوشه هستند

2- از داده هایی که خوشه بندی نشده اند یکی را انتخاب می کنیم

از DBSCAN برای پیدا کردن توزیع نیز استفاده می شود. مثلا مقدار این دو از

n کم است آنرا حذف می کنیم