



## دسته بندی بر مبنای نزدیکترین همسایه ها

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

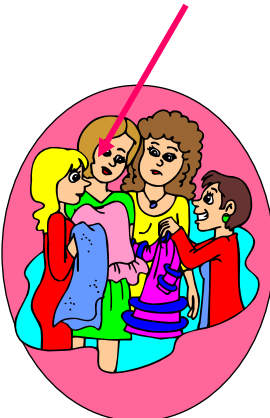
## رئوس مطالب

- مفاهيم اوليه
- انواع classifier ها
- روال كلي KNN
- مشكلات و راه حلهاي ارائه شده

سميه عليزاده هيات علمي دانشكده صنايع  
دانشگاه خواجه نصير طوسي

## مفاهيم

Tell me who your friends are and I'll tell you who you are!



سميه عليزاده هيات علمي دانشكده صنايع  
دانشگاه خواجه نصير طوسي

## انواع دسته بندیها

### Eager •

- ساخت مدل از نمونه های آموزشی قبل از classification
- درخت تصمیم گیری ، نمونه ای از eager classifier ها.

### Lazy

- ذخیره سازی صرف نمونه های آموزشی ، بدون ساخت مدل از آنها

#### Instance based learner •

- تعویق یادگیری تا زمان classification
- KNN نمونه ای از lazy classifier ها.

- تفاوت دو روش در زمان train و test

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## روال کلی KNN

1. تعیین پارامتر  $K$  (تعداد نزدیکترین همسایه ها)
2. محاسبه فاصله نمونه ورودی با تمام نمونه های آموزشی
3. مرتب کردن نمونه های آموزشی بر اساس فاصله و انتخاب  $K$  همسایه نزدیک
4. کلاسی که اکثریت را در همسایه های نزدیک دارد بعنوان تخمینی برای کلاس نمونه ورودی بکار ببر (combine)

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## تابع فاصله

- برای محاسبه فاصله میتوان از تابع فاصله اقلیدسی استفاده کرد.
- فاصله اقلیدسی بین دو تاپل  $X_1$  و  $X_2$  از رابطه زیر بدست می آید.

$$X_2 = (x_{21}; x_{22}; \dots; x_{2n}) \text{ و } X_1 = (x_{11}; x_{12}; \dots; x_{1n})$$

$$dist(X_1, X_2) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}$$

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مثال

- در یک بررسی ، پرسشنامه ای برای دسته بندی کاغذ ها به دو دسته خوب و بد ، بر اساس دو ویژگی مقاومت در برابر اسید و دوام انجام شد. 4 نمونه **training** در زیر دیده میشود:

X1 = مقاومت در برابر اسید = (seconds)	X2 = دوام (kg/square meter)	Y = Classification
7	7	Bad
7	4	Bad
3	4	Good
1	4	Good

کارخانه ، کاغذ جدیدی تولید میکند که تست آزمایشگاه  $x_1=3$  و  $x_2=7$  را برای آن تعیین کرده است. میخواهیم بدون تحقیق پرهزینه ، دسته بندی این کاغذ را بدانیم.

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مثال (ادامه...)

1. تعیین پارامتر  $K$

فرض میکنیم  $K = 3$

2. محاسبه فاصله نمونه ورودی با تمام نمونه های آموزشی

با در نظر گرفتن  $(3, 7)$  بعنوان ورودی فاصله آنرا با تمام نمونه های train محاسبه میکنیم.

فاصله اقلیدسی با نمونه $(3, 7)$	دوام $X2 =$	مقاومت در برابر اسید $X1 =$
	(kg/square meter)	(seconds)
$\sqrt{(7-3)^2 + (7-7)^2} = 4$	7	7
$\sqrt{(7-3)^2 + (4-7)^2} = 5$	4	7
$\sqrt{(3-3)^2 + (4-7)^2} = 3$	4	3
$\sqrt{(1-3)^2 + (4-7)^2} = \sqrt{13}$	4	1

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مثال (ادامه...)

3. مرتب کردن نمونه های آموزشی بر اساس فاصله و انتخاب  $K$  همسایه نزدیک

رتبه فاصله (فاصله اقلیدسی)	فاصله اقلیدسی با نمونه $(3, 7)$	دوام $X2 =$	مقاومت در برابر اسید $X1 =$	جزو ۳ همسایه نزدیک هست؟
۳	$\sqrt{(7-3)^2 + (7-7)^2} = 4$	7	7	بله
۴	$\sqrt{(7-3)^2 + (4-7)^2} = 5$	4	7	خیر
۱	$\sqrt{(3-3)^2 + (4-7)^2} = 3$	4	3	بله
۲	$\sqrt{(1-3)^2 + (4-7)^2} = \sqrt{13}$	4	1	بله

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مثال (ادامه...)

4. کلاسی که اکثریت را در همسایه های نزدیک دارد بعنوان تخمینی برای کلاس نمونه ورودی بکار ببر.

کلاس نزدیکترین همسایه	جزو ۳ همسایه نزدیک هست ؟	رتبه (فاصله اقیلمسی)	فاصله اقیلمسی یا نمونه (3, 7)	دوام = X2 (kg/square meter)	مقاومت در برابر اسید = X1 (seconds)
Bad	بله	۳	$\sqrt{(7-3)^2 + (7-7)^2} = 4$	7	7
-	خیر	۴	$\sqrt{(7-3)^2 + (4-7)^2} = 5$	4	7
Good	بله	۱	$\sqrt{(3-3)^2 + (4-7)^2} = 3$	4	3
Good	بله	۲	$\sqrt{(1-3)^2 + (4-7)^2} = \sqrt{13}$	4	1

چون دو Good و یک Bad داریم و  $1 < 2$  نتیجه میگیریم که کاغذی که در آزمایشگاه اعداد  $x_1=3$  و  $x_2=7$  را بدست آورده است در دسته خوب قرار میگیرد.

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## بررسی دقیق تر KNN

- مسائل مربوط به تابع فاصله
- انتخاب تابع combine
- انتخاب مقدار K

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مسائل مربوط به تابع فاصله

1. ویژگیهای غیر عددی
2. تفاوت در مقیاس اندازه گیری
3. عدم وجود مقدار برای ویژگی
4. انتخاب تابع فاصله
5. ویژگیهای نامرتب

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مسائل مربوط به تابع فاصله

- مقایسه مقادیر ویژگیهای غیر عددی

- ساده ترین روش مقایسه این است که اگر مقدار ویژگی در دو نمونه برابر است تفاوت 0 و در غیر اینصورت تفاوت 1 در نظر گرفته شود ؟
- وجود روشهای پیچیده تر

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مسائل مربوط به تابع فاصله (ادامه...)

- تفاوت در مقیاس اندازه گیری ویژگیها

– محو اثر بعضی ویژگیها

- ابتدا همه مقادیر نرمال شده و سپس مقایسه صورت میگیرد
- نرمالسازی ویژگی  $A$  با مقدار  $U$  به مقدار  $U'$  در بازه  $[0,1]$

$$v' = \frac{v - \min_A}{\max_A - \min_A}$$

$\min_A$  و  $\max_A$  مینیمم و ماکزیمم روی مجموعه  $\text{train}$  محاسبه میشود.

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مسائل مربوط به تابع فاصله (ادامه...)

- ویژگی در یک (یا دو) نمونه مقدار ندارد

– حداکثر مقدار ممکن بعنوان تفاوت نظر گرفته میشود.

تفاوت	ویژگی
۱	غیر عددی
۱	عددی
مقدار بزرگتر $U$ و $1-U$ در نظر گرفته میشود.	هیچکدام مقدار ندارد در یکی مقدار $U$ دارد و دیگری مقدار ندارد

– مثال : ویژگی  $A$  در یک نمونه مقدار  $0.4$  و دیگری مقدار ندارد  
 $\max(0.4, 1 - 0.4) = 0.6$

$0.6$  بعنوان تفاوت در نظر گرفته میشود

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی



## مسائل مربوط به تابع فاصله (ادامه...)

- انتخاب تابع فاصله

فاصله بین دو نمونه  $\mathbf{X1}=(x_{11}; x_{12}; \dots ; x_{1n})$  و  $\mathbf{X2}=(x_{21}; x_{22}; \dots ; x_{2n})$  میتواند از روابط زیر محاسبه شود :

1. Manhattan

$$dist(\mathbf{x}_1, \mathbf{x}_2) = \sum_{i=1}^n |x_{1i} - x_{2i}| \quad \text{2. تابع اقلیدسی}$$

$$dist(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}$$

- سخت گیری
- صرف نظر از ریشه
- استفاده از توانهای دیگر
- تابع مرسوم

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مسائل مربوط به تابع فاصله (ادامه...)

- یکسان گرفتن اهمیت ویژگیها در تابع فاصله

– وجود ویژگیهای نامرتب

- لزوم در نظر گرفتن وزن برای ویژگیها

$$Euclidean(X1, X2) = \sqrt{\sum_{i=1}^n w_i (x_{1i} - x_{2i})^2}$$

- روشهای تعیین وزن

– نظر خبره

– الگوریتمی

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## مسائل مربوط به تابع فاصله (ادامه...)

- روش های تعیین وزن الگوریتمی
  - Cross-validation
    - مقدار تصادفی اولیه برای وزن ویژگیها ، انجام دسته بندی و تغییر برای مینیمم کردن نرخ خطا
    - الگوریتم ژنتیک
  - روش Aha
    - با یک مقدار اولیه شروع میکند و بعد از هر بار دسته بندی ، وزن ها را عوض میکند.
  - شرط درستی روش
    - با اضافه کردن ویژگیهای نامرتبط ، کیفیت دسته بندی افت نکند

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## روشهای ارزیابی دسته بندی ها

- (Accuracy) صحت پیش بینی
- زمان لازم برای ساخت مدل و زمان لازم برای استفاده از مدل
- پایداری (Robustness): توانایی برخورد با داده های غیر معمول (noise) و یا مقادیر جا افتاده (missing values)
- قابلیت تفسیر (Interpretability): قابل فهم بودن مدل و فراهم سازی بینش لازم توسط آن
- جمع و جور بودن (Compactness) مدل : اندازه درخت یا تعداد قواعد و ...

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## شاخص دقت در دسته بندی

Classes	<i>buys_computer = yes</i>	<i>buys_computer = no</i>	Total	Recognition (%)
<i>buys_computer = yes</i>	6,954	46	7,000	99.34
<i>buys_computer = no</i>	412	2,588	3,000	86.27
Total	7,366	2,634	10,000	95.52

	$C_1$	$C_2$
$C_1$	true positives	false negatives
$C_2$	false positives	true negatives

حساسیت  $sensitivity = \frac{t\_pos}{pos}$

شفافیت  $specificity = \frac{t\_neg}{neg}$

دقت  $precision = \frac{t\_pos}{(t\_pos + f\_pos)}$

صحت

$$accuracy = sensitivity \frac{pos}{(pos + neg)} + specificity \frac{neg}{(pos + neg)}$$

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی