

سوال ۴ تعاریف:

diminution: وسعت قابل اندازه گیری از یک نوع خاص را بعد می گویند؛ مانند طول موج و دمای

outlier: یک متغیر از dataset یا outlier که "داده بیرون" می گویند که مقدارش خیلی
بزرگتر متغیرها متفاوت باشد.

independent var: متغیر مستقل متغیری است که وابستگی ای به دیگر متغیرها ندارد و عنوان

مثال در رابطه زیر $y = 2x$ یک متغیر وابسته است.

Sampling stratified: برای انتخاب نمونه ای که فاصله های گره های مختلف باشد، از این نمونه گیری استفاده می کنیم

دین پردازش و در داده کاوی به دستکاری اجزای داده ها پس از استفاده به منظور افزایش عملکرد دین پردازش می گویند.

سوال ۵

PCA یکی از تکنیک‌های کاهش ابعاد خطی است که مجموعه‌ای از متغیرهای همبسته (p) را به تعداد k (که $k < p$) تا از متغیرهای غیرهمبسته

(که همان principal component هستند) تبدیل می‌کند و می‌تواند در عین حال تاجایی که ممکن است تغییرات را در مجموعه داده اولیه حفظ می‌کند. در ML این الگوریتم راغبندان یک روش unsupervised برای کاهش بعد داده می‌کنند.

8 برای مثال PCA، از کاهش ابعاد Iris flowers استفاده می‌کنیم:

9
10 که این مثال را با زمان پایتون نویسیم در screen shot آن را در صفحه بعدی قرار
11 خواهیم داد و خروجی آن را نیز قرار خواهیم داد.

12 ابتدا داده‌ها را به دسته‌های data و targets تقسیم می‌کنیم؛ سپس با استفاده از تابع
13 PCA (که آن را اول import کردیم) عملیات کاهش ابعاد را انجام می‌دهیم.

14 ملاحظه می‌کنیم که این تابع دریافت می‌کند به عنوان ورودی تعداد ویژگی‌ها و تعداد
15 کاهش ویژگی‌ها.

16
17 حال همانطور که در کد مشخص است، اگر خروجی بگیریم، بعدها از ۴ به ۲ کاهش داده می‌شوند.

18
19 خروجی نیز توسط کتابخانه‌ی math plot نشان داده می‌شود.

```
pca-example.py > ...  
1  from sklearn import datasets  
2  from sklearn.decomposition import PCA  
3  from matplotlib import pyplot as plt  
4  
5  iris = datasets.load_iris()  
6  x = iris.data  
7  y = iris.target  
8  
9  Dimr = PCA(n_components = 2)  
10 Dimr.fit(X)  
11 trans = Dimr.transform(X)  
12  
13 plt.scatter(trans[:,0],trans[:,1],c= y)  
14 plt.show()  
15 |
```

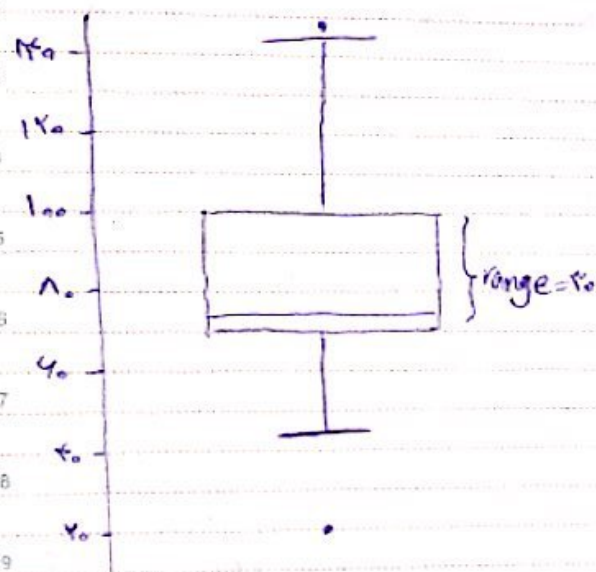


سوال ۶) اینکه همبستگی بین دو متغیر منفی باشد به این معناست که بین متغیرهای رانجی خطی وجود ندارد اما این به این معناست که لزوماً متغیرها مستقل باشند بلکه ممکن است وابسته باشند ولی رانجی بین خطی نداشته باشند.

سوال ۷

نام روش	فرمول اصلی	غالب برای بعد بالا	حساس به مقادیر بزرگ	بازه خروجی
Euclidean Distance	$\sqrt{\sum_{i=1}^n (A_i - B_i)^2}$	X	✓	$(-\infty, +\infty)$
Cosine Similarity	$\frac{A \cdot B}{\ A\ \times \ B\ }$	✓	X	$[-1, 1]$
Hamming Distance	$\sum_{i=1}^n \text{idk}$	✓	X	$(0, +\infty)$
Manhattan Distance	$\sum_{i=1}^n A_i - B_i $	✓	✓	$(0, +\infty)$

12
13
14
15
16
17
18
19



$Y_0 = \min$, \bar{Y}_{\min}

$144 = \max$, \bar{Y}_{\max}

$Vf = \text{median}$, \bar{Y}_{median}

سوال ٨

11 سوال 4 (a) در یک مدل regression خطای که دو متغیر مستقل بسیار با یکدیگر ارتباط داشته باشند

12 multicolinearity اتفاق می افتد؟ در واقع بتوانیم در یک مدل regression متغیرات مستقل را از

13 متغیرات مستقل دیگر ~~تجزیه و تحلیل~~ جدا کنیم؟

14

2 March 2013 = ۱۳۹۲ ربيع الثانی

$$Z = \frac{x - \mu}{\sigma}$$

mean جمعیت μ
sd جمعیت σ

(b) فرمول استفاده از روش z-score به صورت مقابل است:

میک z-score موقعیت یک امتیاز را بر حسب حاصله آن از میانگین

زمانی که در واحدهای انحراف استاندارد اندازه گیری می شود،

توصیف می کند به این صورت که اگر مقدار بالاتر از میانگین باشد، امتیاز مثبت و اگر کمتر از میانگین باشد منفی است.

روش IQR نیز، ۵۰٪ مقادیر متوسط را هنگامی که از کمترین به بالاترین مرتب می شود توصیف می کند.

برای یافتن IQR، ابتدا میانی نهی را پس از نهی بالایی داده ها را برای نهی و

$$IQR = Q_3 - Q_1$$

چارک بالا - چارک پایین

(c) روش اول استفاده از mean/median مقادیر

این روش به این صورت کاری کند که مقدار mean/median برای مقادیر non-missing

محاسبه می کند و با هر مقدار missing values جایگزین می شود به صورت مستقل و جدا از هم در هر ستون

مضایب: استفاده از آن آسان است و بسیار سریع

➔ برای داده های عددی کوچک، بسیار خوب کار می کند

مضایب: - این روش دقیق نیست (و فقط برای داده های عددی مضایب است)

- نتایج فضایی برای داده های encode شده به دست می آورد

- همبستگی بین ویژگی ها را در نظر نمی گیرد و تنها بر سطح ستون کاری کند

8

روش دوم) استفاده از مقادیر پرتکرار یا zero constant :

9

این نیز روش دیگری برای راهپوهای آماری برای استنباط مقادیر miss شده و این روش

برای ویژگی‌های طبقه بندی نیز کاربرد دارد (یعنی هم داده های عددی یا string)

و به این صورت کاری کند که داده های miss شده را با داده های پرتکرار در همان column

جایگزین می کند. مزایا: + با categorical features خوب کاری کند.

معایب: - این روش نیز ارتباط بین ویژگی ها را در نظر نمی گیرد

ممکن است اشکال bias را در data ها ایجاد کند

(یعنی برخی element ها وزن بسیار زیادی در dataset با

داشته باشند)

روش سوم) استفاده از K-NN :

از این الگوریتم برای طبقه بندی های ساده استفاده می شود. این الگوریتم از شباهت ویژگی ها برای

پیش بینی دسته های جدید در هر نقطه استفاده می کند. این روش برای پیش بینی مقادیر miss

شده می تواند بسیار کاربرد داشته باشد.

مزایا: + با توجه به نوع dataset می تواند به سرعت از دو روش قبلی عمل کند.

معایب: - این روش با استفاده از ذخیره کل training data ها، هزینه های

زیادی را به پای میزدارد.

- این روش به نقاط پرت بسیار حساس است.