

نام و نام خانوادگی: حسنا اویار حسینی	تمرین اول داده کاوی
شماره دانشجویی: ۹۸۲۳۰۱۰	تاریخ: فروردین - ۱۴۰۲

سوال (۱)

الف) نويز را می توان به عنوان نمونه هایی با برچسب نادرست (نويز کلاس) یا نمونه های که در مقادیر ویژگی آن ها خطایی وجود دارد (نويز صفت) تعریف کرد. اما داده پرت مفهوم گسترده تری است که نه تنها شامل داده هایی می شود که دچار خطا شده اند بلکه شامل داده های معتبر اما ناسازگاری که ممکن است از تغییرات طبیعی در جامعه یا فرآیند ناشی شود هم می شود در واقع آنها اغلب خیلی کوچک یا خیلی بزرگ هستند و عموماً از پیروی از یک الگو که در بقیه داده ها برقرار است سرکشی می کنند.. شاید بتوان گفت داده پرت خارج از محدوده داده ای است که ما انتظار داریم اما نويز داده ای ناخواسته و اشتباه است و باید حذف شود.

ب) داده پرت می تواند اطلاعات ارزشمندی را در یک آزمایش در اختیار ما بگذارد. آنها چیزی منحصر به فرد در مورد یک موقعیت به ما می گویند. اگر ما بفهمیم که چرا یک امر دور افتاده رخ می دهد، به ما کمک می کند تا یک مشکل را به روشی بهتر حل کنیم. موارد پرت اغلب می توانند شروع چیز جدیدی را نشان دهند. برای مثال داده های پرت باعث می شوند تا به سوالات زیر فکر کنیم. آیا هنگام اندازه گیری این مشاهدات، اتفاق غیرعادی مانند قطع برق، شرایط آزمایشی غیرعادی یا هر چیز خارج از حد معمول رخ داده است؟ آیا خطاهای اندازه گیری یا ورود داده رخ داده است؟ و در صورت وجود ناهنجاری متوجه آن شویم و اگر آسیب زا است آن را برطرف کنیم. برای مثال در شناسایی تقلب با حالت خیلی خاص (مانند استعداد برتر) داده پرت مفید است.

ج) بله. یک نويز میتواند یک داده پرت باشد. برای مثال اگر در هنگام اندازه گیری مقادیر دستگاه مورد استفاده برای اندازه گیری دچار خطا شود و یک داده با مقدار خیلی کوچک یا خیلی بزرگ تولید کند یک نويز داریم که داده پرت نیز می باشد.

سوال (۲)

- انبار داده یا **Data Warehouse** مجموعه بزرگی از داده های تجاری است که اطلاعات بسیار ساختار یافته را از منابع مختلف برای ایجاد نوعی دیدگاه تجاری ذخیره میکند. انبار داده ها برای اتصال، گزارش دهی، بررسی و تحلیل داده های تجاری از منابع مختلف مورد استفاده قرار می گیرد و هسته اصلی سیستم هوش تجاری به شمار می رود. می توانید انبار داده ها را یک بانک اطلاعاتی در نظر بگیرید که داده های گذشته و فعلی را در یک مکان واحد جمع آوری می کند. حجم زیادی از داده های موجود در انبارهای داده از منابع مختلفی جمع آوری می شوند؛ برنامه های کاربردی داخلی مانند بازاریابی، فروش و امور مالی

نمونه‌هایی از این منابع هستند. در واقع میتوان گفت انبارهای داده، مخازن مرکزی داده‌های یکپارچه از یک یا چند منبع پراکنده هستند که برای پشتیبانی از طیف گسترده ای از تصمیم سازی ها در یک سازمان خاص طراحی شده است.

- تفاوت های کلیدی بین پایگاه داده و انبار داده عبارتند از:

تفاوت ها	Database	Data warehouse
تعریف	یک پایگاه داده داده های فعلی مورد نیاز برای یک برنامه را ذخیره می کند.	یک انبار داده، داده های جاری و قدیمی را از یک یا چند سیستم در یک طرح واره از پیش تعریف شده و ثابت ذخیره می کند، که به تحلیلگران تجاری و دانشمندان داده اجازه می دهد تا به راحتی داده ها را تجزیه و تحلیل کنند.
بار کاری	عملیاتی	تحلیلی
انعطاف پذیری طرحواره	طرحواره سفت یا انعطاف پذیر بسته به نوع پایگاه داده	تعریف طرحواره از پیش تعریف شده و ثابت
به روز بودن داده ها	Real time	ممکن است بر اساس فراوانی فرآیندهای ETL به روز نباشد
کاربران	توسعه دهنده ها	تحلیلگران کسب و کار و دانشمندان داده
کاربرد	پایگاه های داده برای تراکنش های کوچک و اتمی بیشترین کاربرد را دارند. انجام تغییرات و تعداد درخواست های بالا	انبارهای داده برای پرسش های تجاری بزرگ تر که به سطح بالاتری از تجزیه و تحلیل داده ها نیاز دارند، مناسب تر هستند. خواندن داده ها و درخواست های کمتر

داده ها	پایگاه‌های داده معمولاً فقط حاوی به‌روزترین اطلاعات هستند، که جستجوهای تاریخی را غیرممکن می‌سازد	انبارهای داده از ابتدا برای اهداف گزارش دهی و تجزیه و تحلیل با استفاده از داده های تجاری تاریخی که مرتبط است طراحی شده اند و داده های قدیمی را نیز نگهداری میکنند.
منبع داده ها	داده ها معمولاً از یک منبع گردآوری شده اند.	منبع داده ها معمولاً چندین مورد بوده است.

• شباهت ها:

- هر دو وظیفه نگهداری و ذخیره داده های ساختار یافته یا نیمه ساختار یافته را دارند.
- هر دو را می توان پرس و جو کرد و با تراکنش ها به روز کرد.
- هر دوی آنها حاوی داده هایی درباره یک یا چند نهاد مانند مشتریان و محصولات هستند.
- به طور کلی، لایه پایین انبار داده یک سیستم پایگاه داده رابطه ای است. پایگاه های داده نیز سیستم پایگاه داده رابطه ای هستند. سیستم های DB رابطه ای از سطرها و ستون ها و مقدار زیادی داده تشکیل شده است.
- هر دو از دسترسی چند کاربره پشتیبانی می کنند. یک نمونه واحد از پایگاه داده و انبار داده می تواند توسط بسیاری از کاربران در یک زمان قابل دسترسی باشد.

سوال (۳)

صدک یا percentile مقداری است در یک توزیع نرمال که درصد مشخصی از مشاهدات زیر آن است. در اصل روش توزیع نرمال و صدک یک روش آماری است که برای تشخیص داده‌های پرت در یک مجموعه داده استفاده می‌شود. در این روش، ابتدا توزیع نرمال (یا گاوسی) برای داده‌های مجموعه تعریف می‌شود. توزیع نرمال یک توزیع پیوسته است که به شکل منحنی بلند و باریکی می‌باشد.

سپس با استفاده از توزیع نرمال، با استفاده از صدک ها می‌توان داده‌های پرت را تشخیص داد. صدک ها یک مقداری است که n درصد از داده‌های مجموعه کوچکتر یا مساوی آن هستند. به عبارت دیگر، اگر $n=95$ ، صدک

۹۵ برابر با مقداری است که ۹۵ درصد از داده‌های مجموعه کوچکتر یا مساوی آن هستند. برای تشخیص داده‌های پرت، داده‌هایی که خارج از محدوده صدک مشخصی هستند، به عنوان داده‌های پرت شناخته می‌شوند. به عنوان مثال، نقاط داده ای که از ۹۹ percentile درصد فاصله دارند و کمتر از ۱ percentile هستند، نقطه پرت در نظر گرفته می‌شوند.

و یا در حالتی دیگر روش ابتدا داه ها را مرتب میکنیم

سپس Q_1, Q_2, Q_3 را در بین داده های مرتب شده پیدا میکنیم (که به ترتیب برابرند با داده ای که یک چهارم داده ها از آن کوچکترند - داده ای که نیمی از داده ها از آن کوچکترند - داده ای که سه چهارم داده ها از آن کوچکترند می باشند)

در مرحله بعد $IRQ = Q_3 - Q_1$ را محاسبه میکنیم.

در مقدار $Q_1 - 1.5 \times IRQ, Q_3 + 1.5 \times IRQ$ را به عنوان ابتدا و انتهای بازه محاسبه میکنیم.

داده هایی که خارج از بازه مرحله قبل می باشند یعنی کوچکتر از $Q_1 - 1.5 \times IRQ$ و بزرگتر از $Q_3 + 1.5 \times IRQ$ را به عنوان داده پرت معرفی میکنیم.

سوال (۴)

الف) پاکسازی داده ها (Data cleaning)، شامل شناسایی و رفع خطاهای احتمالی داده‌ها برای بهبود کیفیت آنهاست. در این فرآیند، داده‌های «کثیف» را شناسایی، بررسی، تجزیه و تحلیل، اصلاح یا حذف می‌کنیم تا مجموعه داده‌های خود را پاکسازی کنیم. داده‌های کثیف به معنی ناهماهنگی‌ها و خطاها هستند که می‌توانند از هر بخش فرآیند تحقیق، مانند طراحی ضعیف، اندازه گیری غلط، ورود داده‌های ناقص و... به دست آیند.

این فرایند میتواند شامل مراحل زیر باشد:

- مشاهدات تکراری یا غیرمرتبط را حذف کنید

- رفع خطاهای ساختاری

- داده های دور افتاده نامطلوب را فیلتر کنید

- رسیدگی به داده های گمشده

- اعتبارسنجی و اطمینان از کیفیت مجموعه داده

ب) اهمیت تجسم داده ها در تجزیه و تحلیل داده های پیچیده، شناسایی الگوها و استخراج بینش های ارزشمند است. ساده سازی اطلاعات پیچیده و ارائه آنها به صورت بصری، تصمیم گیرندگان را قادر می سازد تا تصمیمات

آگاهانه و موثر را سریع و دقیق اتخاذ کنند. اما یک از چالش های آن تعداد زیاد داده، محدوده گسترده برای هر ویژگی از داده و ابعاد بالای داده می باشد.

ج) زیرا با انجام پاکسازی داده دقت و صحت و مورد اعتماد بودن تجسم داده ها افزایش می یابد. برای مثال اگر پاکسازی داده انجام ندهیم و داده های پرت در میان داده های باقی بماند محدوده مقدار مربوط به یک ویژگی از داده ها به علت وجود چند داده پرت ممکن است گسترش بیابد و دقت در نمایش کاهش بیابد و اکثر داده ها که در بازه کوچک تری هستند به صورت فشرده تر و با دقت کمتری نمایش داده شوند. از طرفی اگر در پیش پردازش missing value ها را برطرف و یا حذف نکنیم در هنگام نمایش ممکن است به علت نداشتن مقادیر مورد نیاز به مشکل بخوریم.

سوال (۵)

(الف)

- Cosine similarity:

$$\frac{[(-3 * 9) + (5 * 20) + (8 * 16) + (-2 * 8) + (1 * 2) + (2 * 10) + (3 * -6) + (-5 * -15) + (10 * 25) + (-1 * -2)]}{\sqrt{[(-3)^2 + 5^2 + 8^2 + (-2)^2 + 1^2 + 2^2 + 3^2 + (-5)^2 + 10^2 + (-1)^2] * [(9^2 + 20^2 + 16^2 + 8^2 + 2^2 + 10^2 + (-6)^2 + (-15)^2 + 25^2 + (-2)^2]}} = \frac{516}{659.082695874} = 0.782906308$$

- Correlation:

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

$$\bar{x} = \frac{-3 + 5 + 8 + -2 + 1 + 2 + 3 + -5 + 10 + -1}{10} = 1.8$$

$$\bar{y} = \frac{9 + 20 + 16 + 8 + 2 + 10 + -6 + -15 + 25 + 2}{10} = 7.1$$

$$\begin{aligned} \sum (x_i - \bar{x})^2 &= (-3 - 1.8)^2 + (5 - 1.8)^2 + (8 - 1.8)^2 + (-2 - 1.8)^2 + (1 - 1.8)^2 \\ &+ (2 - 1.8)^2 + (3 - 1.8)^2 + (-5 - 1.8)^2 + (10 - 1.8)^2 + (-1 - 1.8)^2 \\ &= 209.6 \end{aligned}$$

$$\begin{aligned} \Sigma(y_i - \bar{y})^2 &= (9 - 7.1)^2 + (20 - 7.1)^2 + (16 - 7.1)^2 + (8 - 7.1)^2 + (2 - 7.1)^2 \\ &+ (10 - 7.1)^2 + (-6 - 7.1)^2 + (-15 - 7.1)^2 + (25 - 7.1)^2 \\ &+ (2 - 7.1)^2 = 1290.9 \end{aligned}$$

$$\begin{aligned} \Sigma(x_i - \bar{x})(y_i - \bar{y}) &= (-3 - 1.8) * (9 - 7.1) + (5 - 1.8) * (20 - 7.1) + (8 - 1.8) \\ &* (16 - 7.1) + (-2 - 1.8) * (8 - 7.1) + (1 - 1.8) * (2 - 7.1) \\ &+ (2 - 1.8) * (10 - 7.1) + (3 - 1.8) * (-6 - 7.1) + (-5 - 1.8) \\ &* (-15 - 7.1) + (10 - 1.8) * (25 - 7.1) + (-1 - 1.8) * (2 - 7.1) \\ &= 389.2 \end{aligned}$$

$$r = \frac{\Sigma(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{(\Sigma(x_i - \bar{x})^2)(\Sigma(y_i - \bar{y})^2)}} = \frac{384.2}{\sqrt{(20.6 * 129.9)}} = .7386$$

- **Mutual Information:**

Joint probability این دو بردار در واقع فراوانی هر ترکیبی از مقادیر را نشان می دهد:

Joint probability =

[illegible]

$$P(a) = [\cdot.5, \cdot.1, \cdot.5, \cdot.1, \cdot.5, \cdot.5, \cdot.1, \cdot.1, \cdot.5, \cdot.1]$$

$$P(b) = [\cdot.5, \cdot.1, \cdot.5, \cdot.1, \cdot.5, \cdot.5, \cdot.1, \cdot.1, \cdot.5, \cdot.1]$$

$$MI(a, b) = \sum_i \sum_j p(a_i, b_j) \log \left[\frac{p(a_i, b_j)}{(p(a_i)p(b_j))} \right] = 1.7467$$

(ب)

- **Cosine similarity**: بر اساس شباهت کسینوس، با توجه به اینکه شباهت متوسط تا زیاد بین G^1 و G^2 وجود دارد. نشان می‌دهد که این دو ژنوم مستقل از یکدیگر نیستند، زیرا تا حدی شباهت دارند.
- **Correlation**: با توجه به اینکه بازه **correlation** بین ۱- و ۱ می‌باشد و **correlation** محاسبه شده برای این دو ژن به نسب زیاد است میتوان گفت این دو ژن با یکدیگر ارتباط دارند و تا حدی وابسته اند. البته به طور کلی، در حالی که همبستگی می‌تواند اطلاعات مفیدی در مورد رابطه بین دو متغیر ارائه دهد، برای ایجاد یک رابطه علی کافی نیست. تجزیه و تحلیل و تفسیر بیشتر برای تعیین ماهیت و قدرت رابطه بین دو متغیر ضروری است.
- **Mutual Information**: مقدار ۱.۷۴۶۷ نسبتاً بالا است که نشان می‌دهد وابستگی متوسط تا قوی بین دو بردار وجود دارد. با این حال، توجه به این نکته مهم است که این فقط یک تخمین بر اساس داده‌های موجود است و قدرت وابستگی می‌تواند با نمونه‌های مختلف یا در زمینه‌های مختلف تغییر کند. بنابراین، تحلیل و تفسیر اضافی برای درک کامل ماهیت و اهمیت رابطه بین دو بردار ضروری است.

(ج) اطلاعات متقابل، شباهت کسینوس و همبستگی می‌تواند نتایج متفاوتی در مورد وابستگی بین دو بردار بدهد. زیرا:

- اطلاعات متقابل میزان اطلاعاتی را که یک بردار در مورد بردار دیگر ارائه می‌دهد اندازه‌گیری می‌کند. این می‌تواند هر نوع وابستگی، از جمله روابط غیر خطی را جذب کند. بر خلاف همبستگی، رابطه خطی بین دو بردار را در نظر نمی‌گیرد. بنابراین، اگر رابطه بین دو بردار غیر خطی باشد، می‌تواند نتایج متفاوتی نسبت به همبستگی بدهد.
- تشابه کسینوس، کسینوس زاویه بین دو بردار را اندازه‌گیری می‌کند. از ۱- تا ۱ متغیر است، که ۱ نشان می‌دهد که بردارها در یک جهت هستند، ۱- نشان می‌دهد که آنها در جهت مخالف هستند و ۰ نشان می‌دهد که آنها متعامد هستند. شباهت کسینوس اندازه بردارها را در نظر نمی‌گیرد، فقط جهت آنها را در نظر می‌گیرد. بنابراین، اگر بزرگی بردارها مهم باشد، می‌تواند نتایج متفاوتی نسبت به اطلاعات متقابل و همبستگی بدهد.

- همبستگی رابطه خطی بین دو بردار را اندازه گیری می کند. از ۱- تا ۱ متغیر است که ۱ نشان دهنده رابطه خطی مثبت کامل، ۰ نشان دهنده عدم وجود رابطه خطی و ۱- نشان دهنده رابطه خطی منفی کامل است. همبستگی فرض می کند که رابطه بین دو بردار خطی است. بنابراین، اگر رابطه بین دو بردار غیر خطی باشد، می تواند نتایج متفاوتی نسبت به اطلاعات متقابل بدهد.

به طور خلاصه، اطلاعات متقابل، شباهت کسینوس و همبستگی می توانند نتایج متفاوتی در مورد وابستگی بین دو بردار به دست دهند، زیرا آنها جنبه های مختلف رابطه بین بردارها را اندازه گیری می کنند. انتخاب این که از کدام معیار استفاده شود به ماهیت داده ها و سؤال خاص تحقیق بستگی دارد. البته که در این مثال خاص نتایج یکسان شده اند.

سوال ۶)

- **Data aggregation** در داده کاوی فرآیند یافتن، جمع آوری و ارائه داده ها در قالب خلاصه شده برای انجام تجزیه و تحلیل آماری طرح های تجاری یا تجزیه و تحلیل الگوهای انسانی است. هنگامی که داده های متعدد از مجموعه داده های مختلف جمع آوری می شود، جمع آوری داده های دقیق برای ارائه نتایج قابل توجه بسیار مهم است. تجمیع داده ها می تواند به تصمیم گیری محتاطانه در بازاریابی، مالی، قیمت گذاری محصول و غیره کمک کند. گروه های داده انباشته با استفاده از خلاصه های آماری جایگزین می شوند. داده های انبوه موجود در انبار داده می تواند به حل مسائل منطقی کمک کند که به نوبه خود می تواند فشار زمانی را در حل پرس و جو از مجموعه داده ها کاهش دهد.
- نمونه گیری داده ها یک تکنیک تجزیه و تحلیل آماری است که برای انتخاب، دستکاری و تجزیه و تحلیل زیرمجموعه ای نماینده از نقاط داده برای شناسایی الگوها و روندها در مجموعه داده های بزرگتر مورد بررسی استفاده می شود.

سوال ۷)

الف)

- انتخاب ویژگی: انتخاب ویژگی فرآیند کاهش تعداد متغیرهای ورودی هنگام توسعه یک مدل پیش بینی است. کاهش تعداد متغیرهای ورودی برای کاهش هزینه محاسباتی مدل سازی و در برخی موارد برای بهبود عملکرد مدل مطلوب است.

- استخراج ویژگی: استخراج ویژگی فرآیند استخراج خودکار ویژگی های معنی دار از داده های خام با استفاده از الگوریتم های یادگیری ماشین است. این روش اغلب در هنگام برخورد با داده های بدون ساختار، مانند تصاویر یا متن، که در آن داده های خام به شکلی نیستند که در حالت یادگیری ماشینی به راحتی قابل استفاده باشد، استفاده می شود. استخراج ویژگی را می توان با استفاده از تکنیک هایی مانند تجزیه و تحلیل مؤلفه اصلی (PCA)، تجزیه و تحلیل تشخیص خطی (LDA) و t-SNE انجام داد.

- مهندسی ویژگی: مهندسی ویژگی یک تکنیک یادگیری ماشینی است که از داده ها برای ایجاد متغیرهای جدیدی که در مجموعه آموزشی نیستند استفاده می کند. می تواند ویژگی های جدیدی را هم برای یادگیری تحت نظارت و هم برای یادگیری بدون نظارت، با هدف ساده سازی و سرعت بخشیدن به تبدیل داده ها و در عین حال افزایش دقت مدل ایجاد کند. به عنوان مثال، ایجاد یک ویژگی جدید که نشان دهنده نسبت دو ویژگی دیگر است، یا رمزگذاری متغیرهای طبقه بندی شده با استفاده از one hot encoding.

به طور خلاصه، مهندسی ویژگی شامل ایجاد ویژگی های جدید از داده های موجود، انتخاب ویژگی شامل انتخاب مرتبط ترین ویژگی ها برای یک مدل، و استخراج ویژگی شامل استخراج خودکار ویژگی های معنی دار از داده های خام است.

(ب)

- الگوریتم های کاهش بعد، یعنی روش هایی که برای کاهش تعداد ویژگی ها یا متغیرهای موجود در داده ها به کار می روند، به دو دسته الگوریتم های کاهش بعد خطی و غیر خطی تقسیم می شوند. این دو دسته الگوریتم به شکل زیر تفاوت دارند:

الگوریتم های کاهش بعد خطی: این الگوریتم ها با استفاده از تبدیل خطی، مثل تجزیه ماتریسی (PCA)، عاملیت ماتریسی (FA) و تجزیه ماتریس مثبت-نیمه معین (NMF)، ویژگی های موجود در داده ها را کاهش می دهند. در این روش ها، تغییر در مقدار هر ویژگی با ترکیب خطی ویژگی های دیگر بیان می شود.

الگوریتم های کاهش بعد غیر خطی: این الگوریتم ها با استفاده از تبدیل غیر خطی، مثل کرنل PCA و تجزیه ماتریس مثبت-نیمه معین با کرنل (Kernel NMF)، ویژگی های موجود در داده ها را کاهش می دهند. در این روش ها، تبدیل غیرخطی برای تبدیل داده ها به فضای بیشتری استفاده می شود تا بتوان از مزایای غیر خطی آن ها بهره گرفت.

در کل، الگوریتم های کاهش بعد خطی و غیر خطی هر دو برای کاهش تعداد ویژگی های موجود در داده ها به کار می روند، اما الگوریتم های کاهش بعد خطی برای داده هایی که دارای ویژگی های خطی هستند مفید هستند، در حالی که الگوریتم های کاهش بعد غیر خطی برای داده هایی که دارای ویژگی های غیر خطی هستند مفید هستند.

• PCA:

۱. استانداردسازی داده ها
۲. به دست آوردن بردارهای ویژه و مقدارهای ویژه از ماتریکس کواریانس (Covariance matrix) یا ماتریس همبستگی (Correlation Matrix)، یا انجام «تجزیه مقدارهای منفرد (Singular Vector Decomposition)
۳. مرتب سازی مقدارهای ویژه به ترتیب نزولی و انتخاب k بردار ویژه ای که متناظر با K بزرگ ترین مقدار ویژه هستند K . تعداد ابعاد زیرفضای ویژگی جدید است. ($k \leq d$)
۴. ساخت ماتریکس تصویر W از K بردار ویژه انتخاب شده
۵. تبدیل مجموعه داده اصلی X به وسیله W ، برای به دست آوردن زیرفضای K بعدی Y

• t-SNE:

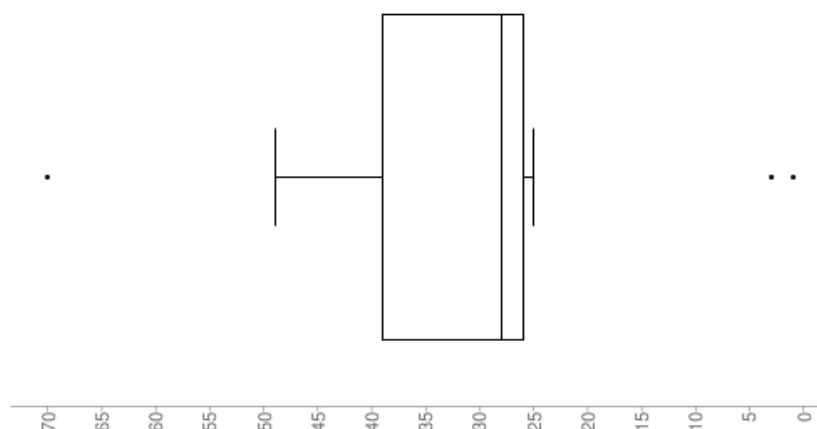
الگوریتم t-SNE شامل دو مرحله به شرح زیر است: مرحله اول بدین صورت است که، یک توزیع احتمال بر روی نقاط در ابعاد بالاتر ایجاد می کند به طوری که به اشیاء مشابه احتمال بالاتر و اشیاء غیر مشابه احتمال کمتر اختصاص داده می شود.

مرحله دوم بدین صورت است که، همان توزیع احتمال را در ابعاد پایین تر به طور مکرر تکرار می کند تا زمانی که واگرایی کولبک-لیبلر به حداقل برسد. به بیانی دیگر، واگرایی کولبک-لیبلر معیاری برای اندازه گیری تفاوت بین توزیع های احتمال مرحله اول و دوم است.

در واقع برای انجام این کار، t-SNE از یک توزیع احتمالاتی برای مدل کردن روابط بین داده ها استفاده می کند. در این روش، ابتدا یک توزیع احتمالاتی چگالی برای داده ها در فضای اصلی محاسبه می شود، سپس توزیع احتمالاتی دیگری برای داده ها در فضای کاهش یافته (معمولاً دو یا سه بعدی) محاسبه می شود. سپس، با کمک گرادینان نزولی، پارامترهای این توزیع احتمالاتی به گونه ای تنظیم می شوند که توزیع احتمالاتی داده ها در فضای کاهش یافته، بهترین توزیعی باشد که داده ها را به خوبی در آن نمایش دهد.

سوال ۸)

Median: ۲۸
 Minimum: ۱
 Maximum: ۷۰
 First quartile: ۲۶
 Third quartile: ۳۹
 Interquartile Range: ۱۳
 Outliers: ۱, ۳, ۷۰



سوال ۹)

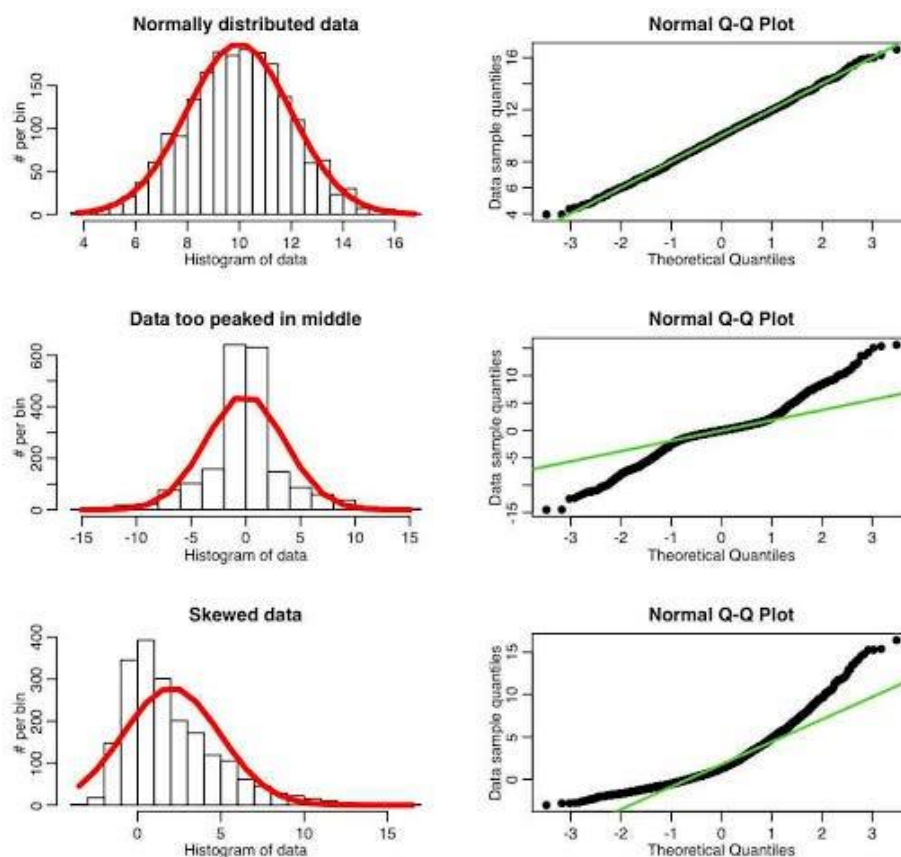
الف) Q-Q plot یا quantile-quantile plot یک نمودار است که برای بررسی تطابق توزیع داده های یک مجموعه با توزیع دیگری (مثلاً توزیع نرمال) استفاده می شود. در این روش، داده های مورد بررسی از کمترین تا بیشترین مقدار مرتب می شوند و برای هر یک از مقادیر، مقدار متناظر آن در توزیع دیگری که می خواهیم با آن مقایسه کنیم محاسبه می شود. سپس دو دنباله از این مقادیر در نمودار قرار داده می شوند: دنباله اول مقادیر مرتب شده از داده های واقعی است و دنباله دوم مقادیر متناظر در توزیع مورد نظر (مثلاً توزیع نرمال) است. در صورتی که داده ها با توزیع دیگری کاملاً تطابق داشته باشند، دو دنباله در یک خط قرار می گیرند.

ب) همانطور که در شکل زیر دیده می شود ۳ حالت را بررسی می کنیم:

۱. در این نمودار، نقاط بر روی خط $x=y$ قرار دارند که نماینده توزیع نرمال است. به عبارت دیگر، داده ها با توزیع نرمال کاملاً تطابق دارند.

۲. زمانی که هر دو انتهای نمودار Q-Q انحراف از خط مستقیم دارد و مرکز آن از یک خط مستقیم پیروی می کند، توزیع داده ها به نحوی است که در دو انتها تعداد داده ها کم و اکثر داده ها در میانه قرار گرفته اند.

۳. زمانی که، نقاط به طور کلی در یک خط قرار دارند، اما در برخی نقاط از خط کمی فاصله دارند. این نشان دهنده تطابق ناقص داده ها با توزیع نرمال است. اگر انتهای پایین نمودار Q-Q از خط مستقیم منحرف شود اما انتهای بالایی منحرف نشود، به وضوح می توانیم بگوییم که توزیع دنباله بلندتری در سمت چپ خود دارد یا به سادگی به سمت چپ انحراف (یا منحنی منفی) است و برعکس.



از روی q-q plot میتوان:

- ارزیابی کنید که آیا نمونه ای از داده ها از توزیع نظری خاصی پیروی می کند یا خیر.
- انحرافات از یک توزیع نظری را شناسایی کنید، که ممکن است نشان دهنده نیاز به روش های آماری جایگزین باشد و درباره skewness داده ها نظر داد.

- تجسم میزان تناسب داده ها با توزیع نظری، که می تواند به محققان کمک کند فرآیند تولید داده های اساسی را درک کنند و تصمیمات آگاهانه ای در مورد تجزیه و تحلیل های آماری بگیرند.

سوال ۱۰)

الف) بازه اعداد: ۰ تا ۱

$$\text{min} - \text{max}: \frac{\text{value} - \text{min}}{\text{max} - \text{min}}$$

ب) بازه اعداد: محدوده مقادیر پس از نرمال سازی امتیاز Z به توزیع داده های اصلی بستگی دارد. اگر داده ها از توزیع نرمال پیروی کنند، محدوده معمولاً بین 3σ تا -3σ می باشد.

$$z - \text{score}: \frac{\text{value} - \mu}{\sigma}$$

ج)

Decimal scaling normalization یکی از روش های نرمال سازی داده ها است که در آن اعشار عدد را از محل آن به سمت چپ حرکت داده و تا جایی که لازم است، با صفر پر شده و عدد نهایی به دست می آید. برای اعمال این روش، از فرمول زیر استفاده می شود:

$$X_{\text{new}} = \frac{X}{10^j}$$

در این فرمول، X نشان دهنده عدد اولیه است و j تعداد ارقام اعشاری در X است. با استفاده از این فرمول، داده ها به صورتی نرمال شده و مقادیر آنها در محدوده [۰,۱] قرار می گیرند. این روش معمولاً برای داده هایی که با مقادیر بزرگ و کوچک زیادی سر و کار دارند و یا واحدهای مختلفی دارند، استفاده می شود.

بازه اعداد: محدوده نرمال سازی مقیاس دهی به تعداد ارقام اعشاری مورد استفاده در فرآیند نرمال سازی بستگی دارد. اگر از اعشار j استفاده شود، در صورت مثبت یا منفی بودن مقادیر داده اصلی، محدوده مقادیر نرمال شده بین -۱ و ۱ خواهد بود. با این حال، اگر مقادیر داده های اصلی همه مثبت باشند، محدوده مقادیر نرمال شده بین ۰ و ۱ خواهد بود. شایان ذکر است که محدوده مقادیر نرمال شده ممکن است دقیقاً بین -۱ و ۱ یا ۰ و ۱ نباشد، اما بسته به تعداد ارقام اعشاری استفاده شده در فرآیند عادی سازی، تقریباً در این محدوده خواهد بود.

سوال (۱۱)

(الف)

(۱۱ الف)

هدف ما پیدا کردن β به گونه‌ای است که S (نقطه به در زیر تعریف شده) به کمین برسد.

$$S = \|y - X\beta\|^2$$

برای این کار یک روش استفاده از مشتق است. به این معنی که $\frac{\partial S}{\partial \beta}$ را محاسبه و برابر با صفر قرار دهیم.

S را می‌توان به صورت زیر نوشت:

$$S = (y - X\beta)^T (y - X\beta)$$

این تعریف به این معنی است:

$$u = f(\beta) = y - X\beta, \quad g(u) = u^T u$$

خواهیم داشت:

$$\frac{\partial S}{\partial \beta} = \underbrace{\frac{\partial g(u)}{\partial u}}_{2u^T} \cdot \underbrace{\frac{\partial f(\beta)}{\partial \beta}}_{-X} = (2u^T)(-X) = 0$$

$$-2(y - X\beta)^T X = 0$$

$$-2X^T y + 2X^T X \beta = 0$$

$$2X^T y = 2X^T X \beta \Rightarrow \boxed{\beta = (X^T X)^{-1} X^T y}$$

Normal Equation ←

(ب)

$$X^T X B = X^T y$$

$$X = \begin{bmatrix} 1 & 2 \\ 1 & 4 \\ 1 & 1 \\ 1 & 3 \\ 1 & 2 \\ 1 & 6 \end{bmatrix}, y = \begin{bmatrix} 5 \\ 6 \\ 3 \\ 6 \\ 3 \\ 1.0 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 2 & 4 & 1 & 3 & 2 & 6 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 1 & 4 \\ 1 & 1 \\ 1 & 3 \\ 1 & 2 \\ 1 & 6 \end{bmatrix} = \begin{bmatrix} 6 & 18 \\ 18 & 70 \end{bmatrix}$$

$$X^T y = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 2 & 4 & 1 & 3 & 2 & 6 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \\ 3 \\ 6 \\ 3 \\ 1.0 \end{bmatrix} = \begin{bmatrix} 33 \\ 121 \end{bmatrix}$$

$$B = \frac{1}{-44} \begin{bmatrix} 70 & -18 \\ -18 & 6 \end{bmatrix} \begin{bmatrix} 33 \\ 121 \end{bmatrix} = \begin{bmatrix} 1.375 \\ 1.375 \end{bmatrix}$$

$$e = y - XB = \begin{bmatrix} .875 \\ -.875 \\ .25 \\ .5 \\ -1.125 \\ .375 \end{bmatrix}$$

(ج)

$$X^T X B = X^T y$$

$$X = \begin{bmatrix} 1 & 2 & 4 \\ 1 & 4 & 16 \\ 1 & 1 & 1 \\ 1 & 3 & 9 \\ 1 & 2 & 4 \\ 1 & 6 & 36 \end{bmatrix}, y = \begin{bmatrix} 5 \\ 6 \\ 3 \\ 6 \\ 3 \\ 1.0 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 6 & 18 & 70 \\ 18 & 70 & 324 \\ 70 & 324 & 1506 \end{bmatrix}$$

$$X^T y = \begin{bmatrix} 33 \\ 121 \\ 545 \end{bmatrix}$$

$$B = \begin{bmatrix} 2.22584 \\ 0.75056 \\ 0.08764 \end{bmatrix}$$

$$e = y - XB = \begin{bmatrix} 0.92248 \\ -0.63032 \\ -0.6404 \\ 0.73372 \\ -1.07752 \\ 0.11576 \end{bmatrix}$$

سوالات بخش عملی (در نوت بوک هم آمده):

(Q۵)

داده‌افزایی (افزایش داده)، در داده‌کاوی، به سازوکارهایی برای افزایش شمارِ داده‌ها گفته می‌شود. داده‌های تازه با ساختن رونوشت‌ها و نمونه‌هایی از داده‌های کنونی (دردسترس)، یا ساخت داده‌هایی با الگوگرفتن از داده‌های کنونی (کمی ناهمسان)، ساخته می‌شوند. شیوه‌ی داده‌افزایی را می‌توان یک رگولارایزر (همترازکننده) دانست، که راه‌کاری است برای چالشِ بیش‌برازش در زمان آموزش یک الگوی (مدل) یادگیری ماشین.

-افزایش داده‌های صوتی

توزیع نویز: برای بهبود عملکرد مدل، نویز گاوسی یا تصادفی را به مجموعه داده صوتی اضافه کنید.

جابجایی: انتقال صدا به چپ (سریع به جلو) یا راست با ثانیتهای تصادفی.

تغییر سرعت: سری‌های زمانی را با نرخ ثابتی افزایش می‌دهد.

تغییر زیر و بم: به طور تصادفی زیر و بم صدا را تغییر دهید.

-افزایش داده‌های متنی

تغییر کلمه یا جمله: تغییر تصادفی موقعیت یک کلمه یا جمله.

جایگزینی کلمه: جایگزین کلمات با مترادف.

دستکاری درخت نحو: جمله را با استفاده از همان کلمه بازنویسی کنید.

درج کلمه تصادفی: کلمات را به صورت تصادفی درج می‌کند

حذف تصادفی کلمه: کلمات را به صورت تصادفی حذف می‌کند.

-تقویت تصویری

تغییرات هندسی: به طور تصادفی ورق زدن، برش، چرخش، کشش و بزرگنمایی تصاویر. باید مراقب اعمال چندین تغییر شکل روی تصاویر مشابه باشید، زیرا این کار می تواند عملکرد مدل را کاهش دهد.

تغییر فضای رنگ: به طور تصادفی کانال های رنگ، کنتراست و روشنایی RGB را تغییر دهید
فیلترهای هسته: به طور تصادفی وضوح یا تاری تصویر را تغییر می دهند
پاک کردن تصادفی: بخشی از تصویر اولیه را حذف کنید.
اختلاط تصاویر: ترکیب و ترکیب چندین تصویر.

تقویت داده ها فقط در مجموعه آموزشی انجام می شود زیرا به تعمیم و استحکام مدل کمک می کند. بنابراین هیچ فایده ای برای افزایش مجموعه تست وجود ندارد.

(Q۶)

- آپ سمپلسنگ رویه ای است که در آن نقاط داده تولید شده مصنوعی (مرتبط با کلاس اقلیت) به مجموعه داده تزریق می شود. پس از این فرآیند، تعداد هر دو برچسب تقریباً یکسان است. این روش یکسان سازی از تمایل مدل به سمت طبقه اکثریت جلوگیری می کند. علاوه بر این، تعامل (خط مرزی) بین کلاس های هدف تغییر می کند. و همچنین، مکانیسم **upsampling** به دلیل اطلاعات اضافی، سوگیری را به سیستم وارد می کند.
- پایین نمونه سازی مکانیزمی است که تعداد نمونه های آموزشی را که در طبقه اکثریت قرار می گیرند کاهش می دهد. زیرا به افزایش تعداد دسته های هدف کمک می کند. با حذف داده های جمع آوری شده، اطلاعات ارزشمند زیادی را از دست می دهیم.

(Q۷)

• **SMOTETomek**

این روش ترکیبی از توانایی **SMOTE** برای تولید داده های مصنوعی برای کلاس اقلیت و توانایی **Tomek Links** برای حذف داده هایی است که به عنوان پیوندهای **Tomek** این روش روشی برای انجام **down sampling, up sampling** می باشد. از کلاس اکثریت (یعنی نمونه هایی از داده های کلاس اکثریت نزدیک ترین به داده های کلاس اقلیت هستند). فرآیند پیوندهای **SMOTE-Tomek** به شرح زیر است.

- ۱- شروع **SMOTE** داده های تصادفی را از کلاس اقلیت انتخاب کنید.
- ۲- فاصله بین داده های تصادفی و **k** نزدیکترین همسایه آن را محاسبه کنید.

۳- اختلاف را با یک عدد تصادفی بین ۰ و ۱ ضرب کنید سپس نتیجه را به عنوان نمونه مصنوعی به کلاس اقلیت اضافه کنید.

۳- مرحله شماره ۲-۳ را تکرار کنید تا نسبت مورد نظر طبقه اقلیت برآورده شود). پایان(SMOTE

(۴- شروع پیوندهای Tomek) داده های تصادفی را از کلاس اکثریت انتخاب کنید.

۵- اگر نزدیکترین همسایه داده های تصادفی، داده های کلاس اقلیت است (یعنی ایجاد پیوند Tomek)، سپس پیوند Tomek را حذف کنید.

• Smoteenn

این روش ترکیبی از توانایی SMOTE برای تولید نمونه های مصنوعی برای کلاس اقلیت و توانایی ENN برای حذف برخی مشاهدات از هر دو کلاس است که دارای کلاس متفاوتی بین کلاس مشاهده و کلاس اکثریت-K نزدیک ترین همسایه آن هستند. فرآیند SMOTE-ENN را می توان به صورت زیر توضیح داد.

(۱- شروع (SMOTE) داده های تصادفی را از کلاس اقلیت انتخاب کنید.

۲- فاصله بین داده های تصادفی و k نزدیکترین همسایه آن را محاسبه کنید.

۳- اختلاف را با یک عدد تصادفی بین ۰ و ۱ ضرب کنید سپس نتیجه را به عنوان نمونه مصنوعی به کلاس اقلیت اضافه کنید.

۴- مرحله شماره ۲-۳ را تکرار کنید تا نسبت مورد نظر طبقه اقلیت برآورده شود). پایان(SMOTE

(۵- شروع K (ENN) را به عنوان تعداد نزدیکترین همسایگان تعیین کنید. اگر تعیین نشد، $3K$.

۶-K- نزدیک ترین همسایه مشاهده را از بین مشاهدات دیگر در مجموعه داده پیدا کنید، سپس کلاس اکثریت را از K-nearest همسایه برگردانید.

۷- اگر کلاس مشاهده و کلاس اکثریت از-K نزدیک ترین همسایه مشاهده متفاوت باشد، مشاهده و-K نزدیک ترین همسایه آن از مجموعه داده حذف می شود.

۸- مراحل ۲ و ۳ را تکرار کنید تا نسبت مورد نظر هر کلاس برآورده شود. (پایان ENN)

- در روش ترکیبی افزایش اندازه و کاهش اندازه معمولاً به این صورت انجام می شود که ابتدا سیگنال ورودی با نرخ نمونه برداری بالا به نرخ نمونه برداری کمتری کاهش داده می شود، سپس پردازشی روی سیگنال کاهش یافته اعمال می شود و در نهایت سیگنال پردازش شده با نرخ نمونه برداری اولیه دوباره بزرگ شده و بازیابی می شود.

- دو روش SMOTEENN و SMOTETomek هر دو از ترکیب دو روش oversampling و undersampling برای مقابله با مشکل ایجاد تعادل در داده‌های دوتایی استفاده می‌کنند. در هر دو روش، ابتدا روش undersampling برای حذف نمونه‌های اکثریتی (مثلا نمونه‌های اقلیتی را حفظ می‌کنند و نمونه‌های اکثریتی را حذف می‌کنند) اعمال می‌شود و سپس با استفاده از روش oversampling (مثلا SMOTE)، نمونه‌های اقلیتی جدیدی ایجاد می‌شوند.