

CA3

رضا چهرقانی

۸۱۰۱۴۰۱

1. اتوانکودر

1.1. یکی از کاربردهای اتوانکودر کاهش نویز در تصاویر می باشد. از آنجایی که داده ها اغلب دارای نویز هستند ما از نوعی از اتوانکودرها برای کاهش نویز استفاده می کنیم. اتوانکودر برای حذف نویزها ابتدا تصویر ورودی را encode می کند و هنگامی که می خواهد آن را decode کند سعی می کند تصویر را به گونه ای بازسازی کند که نویز نداشته باشد.

یکی دیگر از کاربردهای اتوانکودرها کاهش ابعاد داده ی ورودی است. خیلی از زمان ها برای سادگی کار بعد داده های ورودی را کاهش می دهیم و فقط یک سری از ویژگی های را که می خواهیم نگه می داریم. همان طور که در صورت پروژۀ گفته شده encoder داده های ورودی را به فضای پنهان تبدیل می کند که این فضا دارای بعد کمتر می باشد و ویژگی های مهم و معنادار از آن استخراج شده است. برای این کار پس از آموزش اتوانکودر قسمت decoder را حذف می کنید و خروجی ما در واقع همان خروجی encoder می شود.

1.2. اتوانکودر سعی در کاهش تفاوت بین تصویر ورودی و تصویر بازسازی شده می کند. اما به دلیل آنکه برای ساخت تصویر بین پیکسل ها میانگین می گیرد که موجب تار شدن تصویر می شود. همچنین میزان تار بودن تصویر خیلی بستگی به اندازه ی فضای پنهان دارد. وقتی فضای پنهان کوچک باشد موجب از دست رفتن جزئیات و افزایش تاری تصویر بازسازی شده می شود. زیرا مدل باید تمام اطلاعات لازم را در فضای اندکی جای دهد. در فضای پنهان بزرگ پتانسیل ذخیره ی جزئیات بیشتر وجود دارد که موجب کاهش تاری می شود اما ممکن است به داده و محاسبات بیشتری نیازمند باشد.

1.3. در آزمون فرض به روش p-value اگر مقدار محاسبه شده بسیار کوچک باشد آزمون رد می شود. در اینجا هم این مقدار تقریباً برابر صفر می باشد. پس داده های MSE از توزیع نرمال پیروی نمی کنند.

2. معادله رگرسیون

2.1. نقطه پرت نقطه داده ای است که پاسخ y از روند کلی بقیه داده ها پیروی نمی کند و تفاوت قابل توجهی با سایر نقاط داده دارد. اگر یک نقطه داده دارای مقادیر x پیش بینی کننده افراطی باشد نقطه اهرمی است. به نقطه ای با هر دو ویژگی نقطه تاثیر گذار می گویند. این نقاط می توانند منجر به پیش بینی های نادرست شود. نقاط پرت می توانند باعث شوند که خط رگرسیون به سمت آنها منحرف شود. نقاط اهرمی می تواند باعث شود که خط رگرسیون به سمت آنها شیفیت پیدا کند. نقاط تاثیر گذار می توانند تاثیر قابل توجهی بر روی خط رگرسیون داشته

باشند و باعث شوند که در جهتی حرکت کند که نمایشگر کل داده ها نیست. اگر این نقاط مشکل ساز تشخیص داده شوند حذف آن ها مهم است. البته ممکن است همه ی این نقاط مشکل ساز نباشند و داده ای معتبری باشند که باید در تحلیل استفاده شوند.

2.2. ضریب تعیین نسبتی از واریانس برحسب متغیر وابسته است که از متغیرهای مستقل قابل پیش بینی باشد (در بخش بعد از چنین فرمولی برای محاسبه ی آن استفاده می کنیم). هنگامی که برای بررسی نیکوئی برارزش رگرسیون خطی استفاده می شود مقدار ضریب تعیین بین 0 و 1 قرار می گیرد. هنگامی که برابر صفر است یعنی مدل همان میانگین \bar{y} ها را پیش بینی می کند ولی وقتی یک می شود یعنی همه ی پیش بینی ها درست است. در رگرسیون خطی مقدار ضریب تعیین برابر با مجذور ضریب همبستگی بین داده های ورودی و مقادیر پیش بینی شده است. در کل ضریب تعیین ابزاری برای ارزیابی درستی پیش بینی یک مدل است.

2.3. در رگرسیون بر پایه هشت داده اصلی ضریب تعیین به عدد یک نزدیک است و این یعنی نشانه ی درستی بالای پیش بینی است. همانطور هم که از نمودار پیداست نقاط نزدیک خط رگرسیون می باشند. در رگرسیون دومی با اضافه کردن نقطه دور افتاده ضریب تعیین به شدت کاهش می یابد. دلیل آن از نمودار پیداست. زیرا نقطه ای که اضافه کرده ایم به شدت از خط رگرسیون دور است. در سومی هم به همین شکل است زیرا نقطه اهرمی از خط رگرسیون فاصله دارد. اما در چهارمی وقتی که نقطه ای دور افتاده-اهرمی را اضافه می کنیم این نقطه بسیار به خط رگرسیون نزدیک است و مدل پیش بینی کننده ی ما را منحرف نمی کند و ضریب تعیین در این حالت نزدیک به یک خواهد بود. پس در نتیجه دو نقطه ای دور افتاده و اهرمی مشکل ساز هستند و برای افزایش ضریب تعیین و در نتیجه دقت پیش بینی باید آن ها را حذف کرد اما نقطه ای دور افتاده-اهرمی مشکل ساز نیست و حتی می توان به عنوان داده ی معتبر نیز در نظر گرفت.

2.4. چند نوع مدل رگرسیونی داریم که هر کدام به نحوی تاثیر این نقاط دور افتاده یا اهرمی را کاهش می دهد. یکی از آن ها رگرسیون قوی است. که روش هایی برای محدود کردن تأثیر نقض مفروضات صحیح توسط فرآیند تولید داده های اصولی بر تخمین های رگرسیون دارد. برای مثال تابع Huber loss جایگزینی قوی برای روش Least Squares است که هنگام محاسبه سهم داده های پرت را کاهش می دهد و در نتیجه تأثیر آن ها را بر تخمین رگرسیونی کاهش می دهد. یکی دیگر از انواع رگرسیون، رگرسیون وزن دار است که به هر داده وزن مشخصی در محاسبات می دهد. در این نوع رگرسیون اگر به داده ها وزن مناسبی داده شود می توان خطا را به حداقل رساند و به حداکثر دقت رسید. یکی دیگر از کارهایی که می توان کرد تشخیص داده های مشکل ساز و حذف آن ها است. همچنین اگر متغیر وابسته خیلی منحرف شده بود با در نظر گرفتن این تغییر می توان تأثیر نقاط پرت را کاهش داد.

3. قضیه حد مرکزی و نمونه گیری

3.1. به نظرم می توانیم از داده های سایر بازیکنان میانگین بگیریم و حاصل را برای بازیکنانی که اطلاعات آنها را نداریم قرار دهیم.

3.2. مقدار min برابر با کوچکترین داده و max برابر با بزرگترین داده است. حال اگر اعداد را از کوچک به بزرگ

مرتب کنید عدد وسط برابر با میانه Q2 می شود (اگر تعداد داده ها زوج بود میانه برابر میانگین دو عدد وسط می شود). حال اگر اعداد بین min و میانه را در نظر بگیریم میانه ی این اعداد برابر است با Q1 و به طور مشابه میانه ی اعداد بین میانه و max برابر Q3 می باشد. البته تابع sns.boxplot که از آن در رسم نمودار جعبه ای استفاده کردیم بر اساس IQR تعدادی از داده ها را به عنوان داده ی پرت در نظر می گیرد.

3.3. ب. به ازای هر نقطه ی (x, y) در نمودار Q-Q مقادیر x برابر توزیع اول (معمولاً توزیع نظری) و مقادیر y

برابر توزیع دوم می شود. البته اندازه ی نمونه دو توزیع باید برابر باشد و باید آن ها مرتب شده باشند.

ج. همان طور که از نمودار پیداست نقاط نمودار Q-Q با تقریب خوبی روی خط $y=x$ قرار می گیرند. به همین دلیل می توان نتیجه گرفت که توزیع وزن بازیکنان تقریباً دارای توزیع نرمال است.

د. همان طور که از مقادیر پیداست مقدار p-value تقریباً برابر نیم است که مقدار زیادی محسوب می شود پس آزمون فرض پذیرفته می شود. پس طبق این می توان نتیجه گرفت وزن بازیکنان از توزیع نرمال پیروی می کند.

ه. با توجه به نمودار Q-Q با افزایش n توزیع وزن بازیکنان انطباق بیشتری با توزیع نرمال پیدا می کند. یعنی توزیع بازیکنان به سمت توزیع نرمال حرکت می کند. پس می توان نتیجه گرفت توزیع وزن بازیکنان، نرمال است. همچنین می توان نتیجه گرفت وقتی هر چه اندازه ی نمونه بزرگتر باشد نتیجه ی ما دقیق تر خواهد بود و می توانیم پارامترهای توزیع آن ها را به صورت دقیق تری بدست بیاوریم.

3.4. همان طور که از مقادیر p-value پیداست با افزایش n مقدار آن به سمت صفر میل می کند. زیرا با افزایش

نمونه ها معلوم می شود که آن ها متعلق به توزیع نرمال نیستند (زیرا متعلق به توزیع پواسون هستند). در $n=5$ مقدار p-value بیشتر از 0.05 است (مقدار significance level برابر 0.05 است). پس نمی توانیم فرض صفر را رد کنیم. اما در $n=50, 5000$ مقدار p-value کمتر از 0.05 است پس می توانیم آن ها را رد کنیم و بگوییم توزیع آن ها نرمال نیست. قضیه حد مرکزی می گوید وقتی تعداد زیادی متغیر تصادفی با هم جمع می شوند توزیع متغیر تصادفی حاصل نرمال خواهد بود. اما ما در اینجا هیچ جمعیتی انجام نداده ایم.
