

بہ نام خدا

آمریکہ ملی جہارت میناب

تہریات بخش چہار

نام و نام خانوادگی: فردین صداقت

واحد درسی: مباحث ویرہ

رشتہ: مہندسی کامپیوٹر

مدرس: مہندس محمد احمد زارہ

فروردین ۱۴۰۴

## A - چرا Data cleaning در علم داده اهمیت دارد؟

تجزیه و تحلیل داده‌ها تنها برای تجربه و تحلیل داده‌ها مهم است، همچنین برای مدیریت عمومی کسب و کار یا حاکمیت داده نیز مهم می‌باشد. منابع کلان داده‌های پیرا برده و دائماً در حال تغییر هستند. بنابراین نگهداری منظم پایگاه داده‌ها به شکلی تا همه چیز را اضافه کنند. اگر کیفیت داده‌های شما از سطح پایین تری برخوردار باشد نتیجه هرگونه تجزیه و تحلیل با این داده‌ها نیز ناقص خواهد بود. حتی اگر هر مرحله دیگر از روند تجزیه و تحلیل داده‌ها به درستی دنبال کنید، اگر داده‌های شما نامرزن باشد تفاوتی ایجاد نمی‌کند.

## B - Missing Values چگونه مدیریت می‌شوند؟

حذف اشیاء یا ویژگی‌های داده؟

اگر در یک مجموعه داده تنها چند نمونه دارای مقادیر مفقود شده باشد، حذف آنها ممکن است منجر به اشتباه یا بهر در ارتباط این است که ویژگی‌های دارای مقادیر مفقود شده دست‌خیز حذف کنیم.

برآورد مقادیر مفقود شده؟

یک مجموعه داده وارد نظر بگیریم که دارای نقاط داده بسیار مشابه هستند. در این وضعیت اغلب از مقادیر ویژگی نقاط نزدیک به نقطه با مقدار مفقود شده برای برآورد مقدار از دست رفته استفاده می‌شود.

با استفاده از فرمت‌های دیگر که در هنگام تجزیه و تحلیل؟

اگر یک یا هر دو شیء دارای مقادیر مفقود شده برای برخی از ویژگی‌ها باشند، می‌توان شباهت را محاسبه با استفاده از ویژگی‌های مقادیر مفقود شده ندارند محاسبه کرد.



## C Outlier چیست و چگونه می توانیم آن را تشخیص دهیم؟

برخی اوقات در مجموعه داده ها، داده هایی ثبت شده اند که به طور چشمگیری با سایر

صارد مغایرت هستند. آنها خرد را در یک با چند ویژگی متغیر می کنند. می توانند باعث ایجاد

نامنجاری در نتایج به دست آمده از طریق الگوریتم ها و سیستم های تحلیل شوند.

در مدل های با دقت داده های دور افتاده می توانند مزایای آموزش را از دست بدهند. این

امری نیازمند موجب به طولانی شدن زمان آموزش و یا منجر به ایجاد مدل های با دقت کم شود.

## تشخیص Outliers

- مرتب سازی
- با استفاده از نمودار
- با استفاده از نمره Z

## D Data Transformation چه کاربردی دارد؟

پس از استخراج اطلاعات مورد نیاز از مکان های مورد نظر URL

داده ها در قالب بدون ساختار خواهند بود. سپس می توان آن را به صورت ساختار

یافته مانند CSV، صفحه گسترده یا PDF، برای ارائه یا ذخیره سازی تبدیل کرد.

## E - one hot encoding و Label encoding چه

تفاوتی با هم دارند؟

one hot encoding روشی برای تبدیل متغیر های دسته ای به فرمت

دو بیتی (binary) است. این تکنیک سترن های جدیدی ایجاد می کند که حاوی مقادیر

صفر و یک (برای دسته موجود متغیر اصلی است) هر دسته از سطوح

اصلی به یک صورت سترن جداگانه نمایش داده می شود که مقدار 1 حضور آن در دسته

مقدار 0 نبرد آنرا نشان می دهد. clips™



زمانی که متغیر دسته‌ای داده‌ای ترتیب طبیعی (مانند "پایین"، "متوسط" و "بالا")

بایند بهتر است از Label encoding استفاده کنید. این روش به هر دسته یک

عدد صحیح منحصر بفرد اختصاص می‌دهد.

F - چرا Feature Selection در Model-building اهمیت دارد؟

روش‌های انتخاب ویژگی (Feature selection methods) به منظور مواجهه

با داده‌های ابعاد بالا، به مولفه‌های جدایی ناپذیر از فرایند یادگیری مبدل شده‌اند. یک

انتخاب ویژگی صحیح می‌تواند منجر به بهبود یادگیرنده استقرایی از جهت گوناگون

از جمله سرعت یادگیری، ظرفیت تعمیم و سادگی مدل استخراج شده می‌شود.

در مباحث عملی یک متخصص داده بایستی خود ویژگی‌های مورد نیاز را از میان داده‌گان

استخراج کند. حتی در برخی موارد بایستی به دنبال دیتاست جدید بگردد و داده‌ها را جمع‌آوری کند.

G - Duplicate Data چگونه در پایگاه داده‌ها حذف می‌شود؟

برای حذف داده‌های تکراری از پایگاه داده ابتدا باید مشخص کنیم که آیا تکراری بودن چیست (مثلاً

تکرار در مقرون‌بیل به شماره تلفن) سپس با استفاده از دستورات SQL مانند

ROW-NUMBER() یا Group by داده‌های تکراری شناسایی می‌شوند. بعد از

شناسایی معمولاً فقط یکی از رکورد ها (مثلاً اولین مورد) نگه داشته و بقیه حذف می‌شوند. این

کار ممکن است بایک دستور Delete همراه با شرط انجام شود.

H - Irrelevant Data چه مشکلهایی را در پیش‌بینی‌های

Machine Learning ایجاد می‌کند؟

داده‌های بی‌ربط یا غیر ضروری گاهی هستند که در واقع مورد نیاز نیستند و در چارچوب مسائل

که می‌داریم آنرا حل کنیم، مثلاً سبب نیستند. مثال: اگر ما داده‌های مربوط به سلامت عمومی



مردم را تجزیه تحلیل می کنیم، شاید ذهن دادن بی ربط است. به میزان متغیر دیگر، اگر می خواستیم داده های مربوط به نعل ناز را تجزیه و تحلیل کنیم، اما مجبوریم داده متغیر قدی تراست. باید آن مشاهدات بی ربط را حفظ نکنیم. این امر می تواند تجزیه و تحلیل را، آلوده کند و سرگردمی را از هدف اصلی شما به حداقل برساند.

## 1- چرا Data Imputation برای پر کردن Missing Values

کاربرد دارد؟

1- نامگذاری یا کدگذاری متغیر در زبان برنامه نویسی پایتون:

2- ایجاد اختلال در مجموعه داده

3- تاثیرگذاری بر مدل نهایی

4- تمایل به بازتابی تمامی نمونه های مجموعه داده

## 1- چگونه می توانه Normality را در داده های عددی بررسی کنیم؟

برای بررسی نرمال بودن داده های عددی توان از روش های استفاده کرد:

روش های مثل هیستوگرام و  $Q-Q$  که شکل توزیع داده را نشان می دهند. اگر شبیه

منحنی زنگوله ای باشد یا نقاط  $Q-Q$  روی خط باشد، داده نرمال نیست.

روش های آماری مانند آزمون  $Shapiro-Wilk$ ،  $Kolmogorov-Smirnov$  و

$Anderson-Darling$  که با بررسی  $P$ -Value نتیجه می دهند. اگر  $P > 0.05$

داده نرمال فرض می شود.

ترکیب این روش ها را می توان دقیق تری از نرمال بودن داده ها را می دهد.