

Subject :

Year :

Month :

Date :

مقام خدا

نمبرن: ۳۴ چلا Palaeohi در علم دادن اہیت دارد!

درس: مباحث ویزہ

استاد: احمدزادہ

اعضای گروہ: میر صالحی - سمانہ باری

نمبرن: ۱۴۰۲

رشتہ: کاپیوٹر

۱- چرا Data cleaning در علم داده اهمیت دارد؟

Data cleaning یا تسطیح داده، عملی است که در آن داده‌های نادرست یا ناقص را حذف می‌کنیم تا بتوانیم از داده‌ها برای تحلیل و مدل‌سازی استفاده کنیم.

زیادتی داده‌ها افزایش دقت مدل‌های یادگیری ماشین را به همراه می‌آورد. اما داده‌های نادرست یا ناقص می‌توانند باعث

ایجاد نتایج نادرست یا بی‌معنی شوند. بنابراین، تسطیح داده‌ها یک مرحله ضروری در فرآیند یادگیری ماشین است.

داده‌های اضافی را در مراحل بعدی فرآیند تحلیل و مدل‌سازی ایجاد نکنیم. با تسطیح داده‌ها، می‌توانیم از داده‌های تمیز و دقیق برای

در مراحل اولیه این فرآیند استفاده کنیم. این کار باعث می‌شود که نتایج مدل‌سازی ما دقیق‌تر و قابل اعتمادتر باشد.

۲- Missing values چگونه مدیریت می‌شوند؟

مدیریت مقادیر گمشده (Missing values) یکی از مراحل مهم در تسطیح داده‌ها است. این مقادیر می‌توانند به دلیل

اگرچه داده‌های گمشده برای مدیریت این مقادیر وجود دارد که بسته به نوع داده و نیاز از تحلیل

می‌توان آن‌ها استفاده کرد. در داده‌ها به حذف این مقادیر می‌توانیم از روش‌های مختلفی استفاده کنیم.

مقادیر گمشده است. اگر می‌خواهیم حذف داده‌ها را از مجموعه داده‌ها حذف کنیم، باید مطمئن شویم که حذف این مقادیر

بسته به نوع داده‌ها می‌تواند به نتایج مدل‌سازی ما تأثیر بگذارد. بنابراین، باید با دقت و احتیاط در حذف این مقادیر عمل کنیم.

یکی از روش‌های رایج برای مدیریت مقادیر گمشده، حذف داده‌ها است. اما این کار می‌تواند به نتایج مدل‌سازی ما تأثیر بگذارد.

و حذف می‌کند.



۳- **Outliers** چیست و چگونه می توانید آن ها را تشخیص دهید؟
outlier با دادن های بیرون مقدار بیرون هستند که به طور قابل توجهی با سایر داده ها در یک مجموعه متفاوت هستند این مقدار بیرون می توانند ناشی از خطاهای اندازه گیری ورودی یا بی اعتباری **outliers** در تحلیل داده ها بسیار مهم است زیرا می توانند تاثیر زیادی بر نتایج مدل ها و تحلیل ها داشته باشند.

۴- **Data Transformation** چه کاربردی دارد؟
Data Transformation یکی از مراحل کلیدی در فرآیندهای داده پردازشی و تحلیل

داده است این فرآیند به دلایل مختلفی کاربرد دارد که برخی از آن ها اشاره می کنیم به عنوان

۱۵- **مقیاسیت** داده های خام معمولاً شامل فاصله ها و مقیاس های متفاوتی هستند با استفاده

از تکنیک های تغییر دهنده می توان داده ها را تصحیح یا یک سازی و یکپارچه کرد تا کیفیت

آن ها بهبود یابد.

۱۶- **Label Encoding** یا **one-hot Encoding** چه تفاوتی دارند؟

تفاوت های **Label Encoding** و **one-hot Encoding** در بازنمایی متغیرهای کیفی (کاتگوری) است

Label Encoding به فرمت عددی در نظر گرفته می شود و در تکنیک های مدل

در این زمینه **Label Encoding** هستند در ادامه به تفاوت ها و ویژگی های هر یک

Senobar

۱- Model Building Feature چیست دارد؟

Feature selection به انتخاب ویژگی‌ها از مراحلی که می‌تواند در فرآیند ساخت مدل‌های یادگیری

ماشین است و اهمیت آن به دلایل مختلفی بر می‌گردد. - بهبود دقت مدل - انتخاب ویژگی

مهم و مفید می‌تواند دقت مدل را افزایش دهد و ویژگی‌های غیر ضروری یا رانده‌ها را حذف

کند کند. - دقت بیشتر بینش بهتری می‌دهد.

۲- duplicated data چگونه در پایگاه داده‌ها حذف می‌شود؟

حذف duplicated data (داده‌های تکراری) از پایگاه داده‌ها یکی از وظایف مهم در مدیریت

داده‌ها و پردازش داده‌ها است. دلیل این امر آن است که داده‌های تکراری می‌تواند باعث

بازگشت داده‌ها و داده‌های تکراری می‌تواند به وجود بیاید. چندین روش برای حذف داده‌های

متعدد برای شناسایی و حذف داده‌های تکراری در پایگاه داده‌ها استفاده می‌شود.

۳- Irrelevant Feature چیست؟ چه مشکلاتی را در بینش می‌ایجاد می‌کند؟

Machine Learning

وجود Irrelevant data (داده‌های بی‌ربط) در داده‌های آموزشی می‌تواند باعث کاهش

دقت مدل شود. - حذف ویژگی‌های بی‌ربط می‌تواند به بهبود عملکرد مدل کمک کند.

از مشکلات اینجاست که حذف ویژگی‌های بی‌ربط می‌تواند به کاهش دقت مدل

۱- در Data Imputation برای پر کردن values missing کاربرد دارد.
Data Imputation با پر کردن مقادیر گمشده (values = missing) به عنوان یک

تکنیک در علم داده و یادگیری ماشین کاربرد دارد دلایل آن به شرح زیر است. افزایش

دقت مدل - مقادیر گمشده می تواند باعث کاهش کیفیت و دقت مدل ماشین می شود با

پر کردن این مقادیر می توان عملکرد مدل بهبود بخشید و همچنین حاصل کرد که تمام داده

ها در فرآیند آموزش و ارزیابی در نظر گرفته می شود.

۲- چگونه می توانید Normality را در داده های عددی بررسی کنید؟

بررسی نرمال بودن (Normality) داده های عددی یکی از مراحل مهم تحلیل داده ها

است به ویژه زمانی که قصد داریم از آزمون های آماری مبتنی بر فرض نرمال بودن

استفاده کنیم برای بررسی نرمال بودن داده های عددی می توانیم از روش ها و آزمون

های زیر استفاده کنیم - توزیع همبسته گرام - یک همبسته گرام از داده ها رسم کنید و به

شکل توزیع آن توجه کنید اگر داده ها به طور تقریبی به شکل زنگوله ای Normal Distribution

قرار بگیرند احتمالاً نرمال هستند. نتیجه گیری ترکیب چندین روش می تواند به شما دید بهتری

زیر بارها نرمال بودن داده های تان به عدد با این حال همیشه باید همیشه حجم نمونه و نوع داده ها را در نظر
نزد توجه داشته باشید و بر اساس آن نتیجه گیری کنید.

۱۴۰۲/۱۱/۱۰ outliers چیست و چگونه می توان آن ها را تشخیص داد؟

outliers یا داده های بیرون خط، داده هایی هستند که به طور قابل توجهی با سایر داده ها در یک

مجموعه متفاوت هستند. این ها می توانند ناشی از خطاهای اندازه گیری و ورودی اشتباهی

outliers در تحلیل داده ها بسیار مهم است زیرا می تواند تفاوتی غیر قابل توجهی در نتایج مدل

ها و مقادیر آنها داشته باشد.

۱۰ Data Transformation یا کاربرد دارد؟

Data Transformation یکی از مراحل کلیدی در فرآیند آماده سازی و تحلیل

داده است. این فرآیند به دلایل مختلفی کاربرد دارد که به برخی از آن ها اشاره می کنیم.

۱- بهبود کیفیت داده های خام معمولاً شامل خطاها، نویز و اطلاعات ناقص

هستند. با استفاده از تکنیک های تغییر داده می توان داده ها را تصحیح یا کم سازی و یکپارچه کردن کیفیت آن ها بهبود یابد.

۲- Encoding (Label Encoding) یا one-hot Encoding چه تفاوتی دارند؟

تفاوت های Label one-hot در یادگیری ماشین برای تبدیل ویژگی های

کلامی (Categorical Features) به فرمت عددی در نظر گرفته می شود و در تکنیک پردازش

در این زمینه Label Encoding هستند. در داده به تفاوت ها و ویژگی های هر یک

می پردازیم.