

Subject :

Year :

Month :

Date :

مقام خدا

نمبرن: ۳۴ پلا Palaeohi در علم داده اهمیت دارد!

درس: مباحث ویژه

استاد: احمدزاده

اعضای گروه: میر صالحی - سمانه باری

نمبرن: ۱۴۰۲

رشته: کامپیوتر

۱- چرا Data cleaning در علم داده اهمیت دارد؟

Data cleaning یا تسطیح داده‌ها یکی از مراحل کلیدی در علم داده است که دلیل زیر اهمیت

زیادتی دارد، افزایش دقت مدل‌های یادست، ناقص یا سازگاری می‌توانند باعث

ایجاد نتایج بیهوده می‌باشد. ۲- کاهش هزینه‌های داده‌ها نامناسب می‌توانند هزینه

های اضافی را در مراحل بعدی غیرآیند تحلیل و مدل‌سازی ایجاد کنند یا با تسطیح داده‌ها

در مراحل اولیه این هزینه‌ها کاهش می‌یابد.

۲- Missing values چگونه مدیریت می‌شوند؟

مدیریت مفادیرگشده (Missing values) یکی از مراحل مهم در با تسطیح داده‌ها است

روش‌های مختلفی برای مدیریت این مفادیر وجود دارد که بسته به نوع داده و نیاز تحلیل

می‌توان آن‌ها استفاده کرد. در داده به حذف روش‌های مختلف

مفادیرگشده اشاره می‌شود. حذف داده‌ها اگر در محدوده‌ای از داده‌ها گشده

باشد، می‌توان ردیف‌های کامل مجموعه‌های داده را حذف کرد. حذف ستون‌ها اگر

یک ستون دارای مقدارگشده زیاد باشد و اطلاعات کمی ارائه دهد می‌توان آن ستون

را حذف کرد.



۳- outliers چیست و چگونه می توانید آن ها را تشخیص دهید؟
 outlier با دادن نهای 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 و ورودی است یا outliers در تحلیل داده ها بسیار مهم است زیرا می تواند تأثیر زیادی بر نتایج مدل ها و تحلیل ها داشته باشد.

۴- Detect and Format Data vs 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

۱- کیفیت داده های خام معمولاً شامل خطاها و ناهمگنی و اطلاعات ناقص هستند و با استفاده

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

آن ها بهبود یابد.

۵- one-hot Encoding (Label Encoding) چیست و تفاوتی دارند؟

تفاوت های Label Encoding در یادگیری ماشین برای تبدیل ویژگی های کلاسیک

(Categorical Feature) به فرمت عددی در نظر گرفته می شود و تکنیک استاندارد

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

Senobar

۱- Model Building Feat.

Feature selection و ویژگی‌ها را از مراحل کلیه در فرآیند ساخت مدل‌های یادگیری

ماشین است و اهمیت آن بدلیل معضلی بر می‌گردد. - ایجاد دقت مدل با انتخاب ویژگی

جمع و مفید می‌تواند دقت مدل را افزایش داد و ویژگی‌های غیر ضروری را حذف می‌تواند.

باعث کاهش کیفیت پیش‌بینی‌ها می‌شود.

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 داده‌های بزرگ است. بدلیل بار خوانی و پردازش داده‌ها در سیستم‌ها

معمول است داده‌ها داده‌های تکراری می‌توانند به وجود می‌آیند و باید حذف شوند و تکرارهای

متمم برای شناسایی و حذف داده‌های تکراری در پایگاه داده‌ها شروع می‌شود.

۳- Irrelevant Data - چه مشکلاتی را در پیش‌بینی‌های ابعادی می‌کند؟

Machine Learning

وجود Irrelevant data (داده‌های بی‌ربط) در داده‌های آموزشی مدل‌های یادگیری

می‌تواند مشکلات و چالش‌های زیادی را در پیش‌بینی‌های ابعادی ایجاد کند و در نتیجه

از مشکلات اصلی ناشی از وجود داده‌های بی‌ربط اشاره می‌شود.

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 های مبتنی بر یاد

پر کردن این مقادیر می توانیم با کمک مدل بهبود بخشید و امکان حاصل کرد که تمام داده

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

۲- چگونه می توانیم Normality را در داده های عددی بررسی کنیم؟

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 Distrb

توزیع بگیرند احتمالاً نرمال هستند. نتیجه گیری ترکیب چندین روش می تواند به شما دید بیشتری

زیر بارها نرمال بودن داده های تان به عدد با این حال همیشه باید همیشه حجم نمونه و نوع داده ها و

نوع توزیع داشته باشید و بر اساس آن نتیجه گیری کنید.

۲- 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، ورودی اشتباهی

outliers در تحلیل داده ها بسیار مهم است زیرا می توانند نتایج زیادی بر نتایج مدل

ها و تحلیل ها داشته باشند.

۳- Data Types Formation یا کاربرد دارد؟

Data Types Formation یکی از مراحل کلیدی در فرآیند معای داده پردازشی 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، همبستگی و اطلاعات ناقص.

۲- دسته بندی استفاده از تکنیک های تغییر داده می توان داده ها را تصحیح یا

سازی و یکپارچه کردن کیفیت آن ها بهبود یابد.

۳- one-hot Encoding (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 ادامه به تفاوت ها و ویژگی های هر یک

می پردازیم.

