

Subject

پیشرفت

نام درس:

پیشرفت

درس مقدماتی:

Data Preprocessing

نام درس: مقدماتی

واژه‌ها: دروس مقدماتی ۱۴۵۳

نام مدرس: محسن استاد احمدزاده

Subject: تمرینات مباحث ویژه

۸. چرا Dafa Cleaning در علم داده اهمیت دارد؟

پاک سازی داده ها مرحله ای اساسی در علم داده است زیرا داده های Dafa Cleaning نادرست یا ناقص می توانند منجر به نتایج نادرست و تصمیم گیری ها غلط شوند. پاک سازی داده ها شامل شناسایی و اصلاح یا حذف داده های نادرست و ناقص و تکرارها یا غیر مرتبطها است.

با وجود داده های کمیز و بیضه و مدل ها یادگیری ماشین می توانند عملکرد بهتری داشته باشند و پیش بینی های دقیق تر ارائه دهند.

۹. Missing values چگونه مدیریت می شود؟

Missing values یا مقادیر گمشده به داده هایی اشاره دارد که در یک مجموعه داده وجود ندارد. مدیریت این مقادیر می تواند شامل چندین روش باشد:

حذف رکوردها: رکورد هایی که شامل مقادیر گمشده هستند می توانند حذف شوند.

بجایگزینی: مقادیر گمشده می توانند با میانگین، میان یا مد داده های موجود جایگزین شوند.

مدل سازی: استفاده از الگوریتم های خاص برای پیش بینی مقادیر گمشده بر اساس سایر ویژگی ها.

۱۰. Outliers چیست و چگونه می توان آن ها را تشخیص داد؟

Outliers یا نقاط دور افتاده داده هایی هستند که به طور قابل توجهی از سایر نقاط داده متفاوت هستند. این نقاط می توانند ناشی از خطا در جمع آوری داده ها موارد واقعی باشند. تشخیص Outliers می تواند با استفاده از روش های زیر انجام

شود.



نمودار جعبه‌ای (Box Plot): برای شناسایی نقاط دور افتاده، روش  $Z$ -Score مناسب است.  $Z$ -Score برای شناسایی داده‌هایی که بیش از ۳ انحراف استاندارد از میانگین فاصله دارند.

روش IQR: مناسب فاصله بین کوارتایل‌های اول و سوم و شناسایی نقاطی که خارج از محدوده مشخص شده قرار دارند.

### D. Data Transformation

یا تبدیل داده‌ها به فرآیند مشخصی شکل داده‌ها به منظور بهبود کیفیت  $Data Transformation$  و قابلیت استفاده آن‌ها اشاره دارد. این فرآیند شامل استاندارد سازی، نرمال سازی و دسته بندی به مدل‌ها کمک می‌کند تا بهتر یاد بگیرند و عملکرد  $Data Transformation$  داده‌ها شود. بهتر دانسته باشند، همچنین می‌تواند سرعت پردازش را نیز افزایش دهد.

E. <sup>needing</sup> Label Encoding, Encoding Techniques (One-Hot): به تفاوتی دارند. در این روش هر دسته به یک عدد صحیح اختصاص داده می‌شود. این روش برای داده‌های ترتیبی مناسب است.

One-Hot Encoding: در این روش هر دسته به یک بردار یابینری تبدیل می‌شود که تنها یک عنصر آن ۱ و بقیه ۰ هستند. این روش برای داده‌های غیر ترتیبی مناسب است و مانع از ایجاد ترتیب نادرست بین دسته‌ها می‌شود.

F. Feature Selection: در Model-building اهمیت دارد.

**Feature Selection**: یا انتخاب ویژگی ها به فرایند انتخاب زیر مجموعه ای از ویژگی ها بر مبنای با هدف برای استفاده در مدل ساز اشاره دارد.  
این عمل به کاهش ابعاد داده ها افزایش دقت مدل و کاهش زمان آموزش کمک میکند.  
همچنین با کاهش نویز و جلوگیری از **overfitting** به بهبود کارایی مدل کمک میکند.  
**6. Duplicate Data**: چگونه در پایگاه داده ها حذف می شود؟

**Duplicate Data**: داده های تکراری می توانند با استفاده از روش ها زیر در پایگاه داده ها حذف شوند:  
**3.4** استفاده از توابع گروهی: مانند **GROUP BY** در **SQL** برای شناسایی و حذف رکورد های تکراری.

استفاده از روش ها با برنامه نویسی: مانند استفاده از پایتون برای شناسایی و حذف داده های تکراری **DataFrame**.  
**4. Irrelevant Data**: مشکلاتی را در پیش جینی های **Machine Learning** ایجاد میکند؟

یا داده های غیر مرتبط می توانند باعث کاهش دقت مدل ها یا جلوگیری از شناسایی **Irrelevant Data** شوند و زمان آموزش را افزایش دهند **overfitting** شوند. این داده ها می توانند موجب همپنین به مدل را دچار سردرگمی میکنند و ممکن است به نتایج نادرست منجر شوند.  
**1. Data Imputation** چرا برای چر کردن **Missing Values** کاربرد دارد؟  
یا چر کردن مقادیر گم شده تکنیکی است که برای بهبود کیفیت داده ها و **Data Imputation** جلوگیری از مشکلات ناشی از مقادیر گم شده استفاده می شود. این کار باعث



عقد اطلاعات و افزایش دقت مدل های یادگیر ما شین ما شود.

ل. پیگوند ما توانید  $Normal$  را در داده های عددی بر ما کردی

بر ما بر  $Normal$  یا نرمال بودن داده های عددی ما توان از روش های زیر استفاده کرد:

شمار هیس توگرام: بر ما مشاهده توزیع داده ها.

شمار  $Q-Q$ : بر ما مقایسه توزیع داده ها با توزیع نرمال.

آزمون های آماری: مانند آزمون  $Shapiro-wilk$  یا  $Kolmogorov-Smirnov$  بر ما نرمال بودن داده ها.

در نهایت  $Duclaux$  و مدیر است داده ها به عنوان مراحل حیاتی در علم داده یاد گیر ما شین به شمار ما آیند و تاثیر قابل توجهی بر عملکرد مدل ما دارند.