

Subject: استاد احمدزاده  
Year: Month: Date:

مباحث و ویژه

بخش چهارم

دو شرفاوری

Data Preprocessing

A Data cleaning در علم داده اهمیت دارد زیرا داده های بیکیفیت بالا امکان ایجاد تحلیل دقیق تر، تصمیم گیری بهتر و کاهش خطاها و ریسک ها را افزایش می دهد و برای بهبود عملکرد سیستم و سایر اهداف مهم می باشد.

B Missing Values چگونه مدیریت می شود؟  
حذف داده ها (حذف سطر یا ستون) و این حذف باید بر مبنای ۲ جایگزینی داده های گم شده (جایگزینی با میانگین یا به تحلیل وابسته) بهترین روش برای مدیریت داده های گم شده مبتنی بر منبع داده خاص است. داده های گم شده برای داده های کم و تعداد کمی مناسب است. حذف داده های پرت و غیره می تواند منجر به تحریف نتایج شود.

C Outliers چیست و چگونه می توان آن ها را تشخیص داد؟  
داده های پرت و داده های استثنای در داده ها می باشد که می تواند به دلیل خطای اندازه گیری یا تغییر در رفتار باشد. این داده ها می توانند به عنوان داده های غیر عادی در نظر گرفته شوند و حذف آن ها می تواند منجر به تحریف نتایج شود. این داده ها با بقیه داده ها متفاوت است و اگر حذف آن ها را در نظر بگیریم، ممکن است منجر به حذف داده های معتبر شود.

داده های پرت و غیره می توانند به دلیل خطای اندازه گیری یا تغییر در رفتار باشد. این داده ها می توانند به عنوان داده های غیر عادی در نظر گرفته شوند و حذف آن ها می تواند منجر به تحریف نتایج شود. این داده ها با بقیه داده ها متفاوت است و اگر حذف آن ها را در نظر بگیریم، ممکن است منجر به حذف داده های معتبر شود.

D Data transformation  
تبدیل و مقایسه تجزیه و تحلیل داده استفاده می شوند این فرآیند به ما امکان می دهد از منابع مختلف یک داده را برای رسم یک تصویر از آن ها در یک نمودار واحد نمایش دهیم. همچنین تصمیم گیری را آسان تر می کند.



## Encoding Techniques (One-Hot Encoding) : E

و Label Encoding با تفاوت دارند؟

### 1. Label Encoding

در این روش هر دسته یک عدد منحصراً اختصاص داده می شود. برای مثال اگر یک ویژگی با مقادیر ۱، قرمز، سبز، آبی داشته باشد، این مقادیر به ترتیب مقادیر ۰، ۱، ۲، ۳ اختصاص می یابند. بنابراین اگر در یک مدل نیاز به این روش می توانیم به الگوریتم ها تبدیل های نامرتبی به عددی را می توانیم. اما مشکل این است که فرض کنید بین مقادیر عددی ارتباطی وجود دارد (مثلاً ۲ بیشتر از ۱ است).

### 2: One-Hot Encoding

در این روش هر دسته یک مقادیر باینری اختصاص داده می شود. برای مثال اگر یک ویژگی ۳ دسته داشته باشد، این مقادیر به صورت ۳ ستون باینری نمایش داده می شود.

■ قرمز:  $[0, 0, 1]$

■ سبز:  $[0, 1, 0]$

■ آبی:  $[1, 0, 0]$

مزاها: این روش هیچ وابستگی از مقدار یا ترتیب دسته ها ندارد.

به علاوه قلا به Label Encoding برای ویژگی های نامرتبی مناسب است.

مناقصات: در حالت One-Hot Encoding برای

ویژگی های غیر ترتیبی بهترین گزینه است.



Subject:

Year.

Month.

Date.

Model building و feature selection  
۱. استوار درجه ۱. اعتبار کمتر ۲. کاهش بیشترین مدل ۳. جلوگیری  
از اُور فیت ۴. افزایش قابلیت تفسیر

Duplicate Data چگونه در پایگاه داده ها حذف می شود  
۱. استفاده از دستور select ۲. استفاده از توابع گروه بندی  
۳. ابزار ها و نرم افزار های مدیریت داده ۴. بررسی و پاک سازی  
به طور کلی حذف داده های تکراری به بهبود کیفیت داده ها و کاهش  
حجم پایگاه داده کمک می کند.

Irrelevant Data چه مشکلاتی را در پیش می آید  
Machine Learning ایجاد می کند ۱. کاهش دقت مدل  
۲. کاهش سرعت سیکل ۳. افزایش زمان آموزش ۴. کاهش تعمیم  
پذیری ۵. خطای بیشتر

Data Imputation برای پر کردن Missing values  
کاربرد دارد ۱. حفظ اطلاعات ۲. بهینه سازی ۳. کاهش  
انحراف ۴. ایجاد مدل در نتیجه کیفیت رگرسیون و سایر مدل ها  
داده های ناقص است به عنوان یک پارامتر و کیفیت داده  
ها کمتر می کند



Subject:

Year.

Month.

Date.

ن. چگونه می توانند  $Normality$  را در داده های عددی بررسی کنند؟  
۱. جستجو برای نمودار چپه ای ۲. نمودار چپه ای ۳. نمودار چپه ای  
۴. آنزومون های آماری: (۱) آنزومون شاپیرو-ویلک (۲) آنزومون  
کولموگوروف-اسمیرنوف این روش ها بیشتر از نمودار چپه ای تحلیل آماری  
مکمل می کنند تا نزدیکی داده های عددی را به خوبی بررسی کنند و نتایج  
مناسب را مورد تحلیل های آماری و مدل سازی قرار دهند.