

Date: / /

Sat Sun Mon Tue Wed Fri

Subject: -----

بخش ۴

درس : مباحث ویژه

موضوع : Data Preprocessing

نام دانشجو :

فرید تنگلی پور

کبير، حمى زاده

نام استاد: محمد احمدزاده

Date: / /

Sat Sun Mon Tue Thu Wed Fri

Subject: \_\_\_\_\_

A. چرا Data cleaning در علم داده اهمیت دارد؟

Data cleaning یا پاکسازی داده ها اهمیت زیادی در علم داده دارد چون داده ها

نا درست، ناقص یا نامرتب می توانند نتایج تحلیل را به شدت تحت تأثیر قرار بدهند.

وقتی داده ها تمیز نباشند، مدل ها را یادگیر، ماشین ممکن است دقت یا دینی داشته باشند

و تحلیل ما به اشتباه پیفتند. بنابراین، پاکسازی داده ها باعث می شود نتایج دقیق تر و

قابل اعتماد تر بگیریم و به تصمیم ها بهترین برسیم.

B. Missing values چگونه مدیریت می شوند؟

مدیریت مقادیر گم شده در داده های از جمله مهم در پیش بردار می داده ها است.

چند روش برای مدیریت این مقادیر وجود دارد:

1. حذف داده ها: ردیف یا ستون ها را از مقادیر گم شده یا حذف می کنیم

2. جایگزینی: با میانگین، میانگین یا مد میزنیم.

3. پیش بینی: با مدل ها یا دیگر ماشین مقادیر را پیش بینی می کنیم.

Date: / /

Sat Sun Mon Tue Thu Wed Fri

Subject: -----

4. روش‌ها خاص: روش‌های مثل multiple Imputation یا k-Nearest

وجود دارند می‌توانند مقادیر گم شده را به شکل بهتر برگردانند.

5. نگهداری: یک دسته جداگانه به آن مقادیر گم شده تفریق می‌دهند.

c. outliers چیست و چگونه می‌توانید آن را تشخیص دهید؟

outliers یا نقاط پرت به داده‌هایی گفته می‌شود که به طور قابل توجهی از سایر داده‌ها در یک

مجموعه داده فاصله دارند، این نقاط معمولاً می‌توانند نشان‌دهنده خطا در داده‌ها یا

تغییرات طبیعی یا حتی موارد خاص و منحصر به فرد باشند.

برای تشخیص outliers می‌توان از روش‌های زیر استفاده کرد:

1. نمودار جعبه‌ای (Box Plot): نقاط خارج از محدوده جعبه.

2. نمودار پراکنش (Scatter Plot): نقاط دور از سایر نقاط.

3. معیارهای آماری: مانند z-score یا IQR.



Date: / /

Sat. Sun. Mon. Tue. Wed. Fri.

Subject: -----

D. Data Transformation چه کاربردی دارد؟ تبدیل داده‌ها به فرایند تغییر شکل،

ساختار یا فرمت داده‌ها اشاره دارد که به منظور آماده‌سازی آن‌ها برای تجزیه و تحلیل یا

استفاده در سیستم‌ها مختلف انجام می‌شود. این کار چندین کاربردی مهم دارد:

1. بهبود کیفیت داده‌ها: خطاها و نواقص را اصلاح می‌کند.

2. یکپارچه‌سازی داده‌ها: داده‌ها را به فرمت استاندارد تبدیل می‌کند.

3. تحلیل بهتر: سازمان‌های داده‌ها را برای تحلیل ساده‌تر.

4. افزایش کارایی: سرعت پردازش را بالا می‌برد.

5. پشتیبانی از تصمیم‌گیری: اطلاعات بهتر را برای تصمیم‌گیرندگان فراهم می‌کند.

Date: / /

Sat Sun Mon Tue Thu Wed Fri

Subject: -----

(One-Hot Label Encoding), Encoding .E

Encoding Techniques چه تفاوتی دارند؟

تفاوت اصلی بین Label Encoding و one-hot Encoding به نوع نمایش داده‌ها

دستار (categorical Data) مربوط می‌شود:

Label Encoding: هر دسته به یک عدد منحصر به فرد تبدیل می‌شود.

one-hot Encoding: هر دسته به یک بردار باینری تبدیل می‌شود.

F. چرا Feature Selection در building اهمیت دارد؟

Feature Selection یا انتخاب ویژگی‌ها در ساخت مدل خیلی مهم است. به چند دلیل:

1. کاهش ابعاد: مدل را ساده‌تر و سریع‌تر می‌کند.

2. بهبود دقت: نویزی‌ها را حذف می‌کند و دقت مدل را افزایش می‌دهد.

3. جلوگیری از overfitting: کمک می‌کند تا مدل بهتر روی داده‌ها عمل کند.

4. تفسیرپذیری: فهمیدن تاثیر ویژگی‌ها را راحت‌تر می‌کند.

Date: / /

Sat Sun Mon Tue Thu Wed Fri

Subject: -----

داده‌های تکراری (Duplicate Data) چیست؟ چگونه می‌توانیم آن‌ها را حذف کنیم؟

حذف داده‌های تکراری در پایگاه داده‌ها معمولاً با استفاده از چند روش انجام می‌شود.

1. استفاده از دستور `Distinct`: برای انتخاب داده‌های یکتا می‌توان از

دستور `select Distinct` استفاده کرد. این دستور فقط رکوردهای منحصر به فرد را برمی

گرداند.

2. استفاده از `Group By`: می‌توان داده‌ها را بر اساس یک یا چند ستون گروه‌بندی کرده و

رکوردهای تکراری را حذف کرد.

3. استفاده از `Join`: می‌توان با استفاده از یک `Self Join` بین خود جدول و یک

زیرمجموعه رکوردهای تکراری را شناسایی و حذف کرد.

4. `Irrelevant Data`: چه مشکلاتی را در پیش‌بینی‌ها `machine Learning` ایجاد

می‌کند؟ می‌تواند باعث کاهش دقت مدل‌ها `machine Learning` شود. این داده‌ها مربوط

به نویزها و داده‌های بی‌ارتباط (Irrelevant Data) در پیش‌بینی‌ها هستند.



Date: / /

Sat Sun Mon Tue Thu Wed Fri

Subject: -----

مدل‌ها ممکن است نتوانند به درستی یاد بگیرند و عملکردشان ضعیف شود.

۱. چرا Data Imputation برای پر کردن Missing values کاربرد دارد؟

Data Imputation یک تکنیک مهم در علم داده‌ها است که استفاده (missing values)

برای پر کردن مقادیر گم‌شده می‌شود و وقتی که در یک داده، بعضی از مقادیر موجود نیستند

این می‌تواند باعث مشکلاتی در تحلیل مدل‌سازی شود.

دلایل آن برای پر کردن مقادیر:

۱. دقت مدل را افزایش می‌دهد.

۲. اطلاعات مهم را حفظ می‌کند.

۳. سوگیری نتایج را کاهش می‌دهد.

۴. حجم داده‌ها را افزایش می‌دهد.

Date: / /

Sat Sun Mon Tue Wed Fri

Subject: \_\_\_\_\_

ن. چگونه می توانیم Normality را در داده های عددی بررسی کنیم؟

برای بررسی نرمالیتی های داده های عددی ما توان از روش های زیر استفاده کرد:

1. هیستوگرام: بررسی شکل توزیع.

2. نمودار Q-Q: مقایسه نقاط داده با توزیع نرمال.

3. آزمون شاپیرو-ویلک: P-value کمتر از 0.05 نشان دهنده عدم نرمالیتی.

4. آزمون کولموگوروف-اسمیرنوف: بررسی انطباق با توزیع نرمال.

5. آزمون اندرسون-دارلینگ: مشابه آزمون های قبلی.