

به نام خدا

درس: داده کاوی

استاد: دکتر رضا رمضانی

پروژه اول: تحلیل و پیش پردازش

نام و نام خانوادگی: فاطمه مومنی

این پروژه با زبان برنامه نویسی پایتون و با استفاده از کتابخانه های sklearn, pandas, numpy, matplotlib و seaborn پیاده سازی شده است. در ادامه کدهای مربوط به هر بخش، توضیح داده شده اند.

دستورات نشان داده شده در شکل ۱، برای نمایش تمام ستون های مجموعه داده به طور کامل، و دستورات نشان داده شده در شکل ۲، برای خواندن فایل های مجموعه داده، استفاده شده است.

```
questions = pd.read_csv('/content/DM1_dataset/Questions.csv')
answers = pd.read_csv('/content/DM1_dataset/Answers.csv')
```

شکل ۲

```
pd.set_option('display.width', 2000)
pd.set_option('display.max_columns', 61)
pd.set_option('display.max_rows', 200)
```

شکل ۱

اولین دستور شکل ۳، تعداد مقادیر جالفتاده هر ویژگی را می شمارد. تابع isna() این مقادیر را یافته و تابع sum() مجموع تعداد آنها را به دست می دهد. سپس درصد مقادیر جالفتاده را محاسبه کرده و در دیکشنری NaN_percent که کلیدهای آن نام ویژگی ها است، ذخیره می کنیم.

```
#counting number of missing values per column
NaN_num = answers.isna().sum()

#define a dictionary for storing percentage of missing values in each column
NaN_percent = dict()
answers_columns = answers.columns
total_num = len(answers)

#calculating percentage of missing values per column
for num in range(len(NaN_num)):
    NaN_percent[answers_columns[num]] = NaN_num[num] * 100 / total_num
```

شکل ۳

در برخورد با داده های جالفتاده، سه روش به کار گرفته شده است: حذف ویژگی، پر کردن با میانگین مقادیر ویژگی و پر کردن با مد. مطابق شکل ۴، ابتدا دیکشنری حاوی درصد مقادیر جالفتاده را به صورت نزولی مرتب می کنیم. سپس ستون هایی که درصد مقادیر جالفتاده در آنها نسبتا زیاد و بیش از ۴۰٪ است را در لیست max_miss ذخیره می کنیم. در نهایت، ستون های مشخص شده را با استفاده از تابع drop() حذف می کنیم.

```
#sorting list of percentages to find columns with max number of missing values that stored in max_miss list
sort = sorted(NaN_percent.items(), key=lambda x: x[1], reverse=True)
max_miss = list()
for i in sort:
    if i[1] >= 40.0:
        max_miss.append(i[0])
    else:
        break

# deleting columns have many missing values
processed_ans = answers.drop(axis=1, columns=max_miss)
```

شکل ۴

پس از یافتن نوع داده هر ویژگی توسط تابع dtypes، مقادیر جاافتاده ستون‌هایی که از نوع عددی و دسته‌بندی شده هستند را به ترتیب با میانگین و مد سایر مقادیر ستون‌ها، پر می‌کنیم. شکل ۵ کدهای مربوطه را نشان می‌دهد.

```
# finding type of columns
print(processed_ans.dtypes)

# filling missing values with mean for numeric(float64) columns
numeric = ['Age', 'WorkWeekHrs']
for i in numeric:
    processed_ans[i] = processed_ans[i].fillna(processed_ans[i].mean())

# filling missing values with mode for categorical columns
categorical = ['MainBranch', 'Hobbyist', 'Age1stCode', 'CompFreq', 'Country', 'CurrencyDesc', 'CurrencySymbol',
               'EdLevel', 'Employment', 'Ethnicity', 'Gender', 'JobSat', 'JobSeek', 'NEWDevOps', 'NEWDevOpsImpt',
               'NEWEdImpt', 'NEWLearn', 'NEWOFFTopic', 'NEWOnboardGood', 'NEWOtherComms', 'NEWOvertime', 'NEWPurpleLink',
               'OpSys', 'OrgSize', 'PurchaseWhat', 'Sexuality', 'SOAccount', 'SOComm', 'SOPartFreq', 'SOVisitFreq',
               'SurveyEase', 'SurveyLength', 'Trans', 'UndergradMajor', 'WelcomeChange', 'YearsCode', 'YearsCodePro']
for i in categorical:
    processed_ans[i] = processed_ans[i].fillna(processed_ans[i].mode()[0])
```

شکل ۵

همان‌گونه که در شکل ۶ مشاهده می‌شود، پس از یافتن ستون‌های باقی‌مانده، ستون‌های چند مقداری را برحسب کاراکتر جدا کرده و مقدار دارای بیشترین تکرار را برای هر ستون می‌یابیم (به‌دلیل عدم وجود مقادیر جاافتاده در اولین ستون مجموعه داده، یعنی ستون Respondent، حلقه for از عدد یک شروع شده‌است). سپس مقادیر جاافتاده هر ستون را با مقدار به‌دست آمده برای آن ستون، پر می‌کنیم. مقادیر متفاوت ستون LanguageWorkedWith (بدون احتساب تکرار هر مقدار) در لیست LanguageWorkedWith برای استفاده در بخش‌های بعدی، ذخیره شده‌است.

```
# remained columns
remained = list()
for i in answers_columns:
    if i not in (max_miss + numeric + categorical):
        remained.append(i)

multi_value = dict()
for i in range(1, len(remained)):
    multi_value.clear()
    column = processed_ans[remained[i]]
    for row in column:
        values = str(row).split(';')
        for val in values:
            if val != 'nan':
                if val not in multi_value.keys():
                    multi_value[val] = 1
                else:
                    multi_value[val] += 1
    if remained[i] == 'LanguageWorkedWith':
        LanguageWorkedWith = multi_value.keys()
    fill_value = max(multi_value, key=multi_value.get)
    processed_ans[remained[i]] = processed_ans[remained[i]].fillna(fill_value)
```

شکل ۶

برای بررسی وجود یا عدم وجود داده‌های پرت، نمودار جعبه‌ای ویژگی‌ها را رسم می‌کنیم. به این منظور، ابتدا باید مقادیر ستون‌های غیر عددی را به مقادیر عددی تبدیل کنیم که این کار با به‌کارگیری روش Label Encoding و توسط تابع `fit_transform()` انجام شده‌است. شکل ۷ و شکل ۸ به ترتیب دستورات کدگذاری مقادیر ستون‌ها و ترسیم و نمایش نمودار جعبه‌ای را نشان می‌دهند.

```
# show boxplot of all columns
for i in encod_ans.columns:
    plt.boxplot(encod_ans[i])
    plt.title(i)
    plt.show()
```

شکل ۸

```
# use label encoder to encode object data types
encod_ans = processed_ans.copy()
object_type = categorical + remained
label_encoder = preprocessing.LabelEncoder()
for i in object_type:
    encod_ans[i] = label_encoder.fit_transform(encod_ans[i])
```

شکل ۷

با بررسی نمودارهای جعبه‌ای، ویژگی‌های بدون مقادیر جاف‌تاده به‌دست می‌آیند که نام آن‌ها در شکل ۹ مشاهده می‌شود. دستورات نشان داده‌شده در شکل ۱۰ نیز سایر ویژگی‌ها که داده‌های پرت دارند را به‌دست می‌آورد.

```
# features with outlier
outlier_columns = list()
for i in encod_ans.columns:
    if i not in (no_outlier_columns + discrete_columns):
        outlier_columns.append(i)
```

شکل ۱۰

```
# features without outlier
no_outlier_columns = ['Respondent', 'CompFreq', 'Country', 'CurrencyDesc', 'CurrencySymbol', 'JobFactors',
    'JobSat', 'JobSeek', 'LanguageDesireNextYear', 'LanguageWorkedWith', 'MiscTechWorkedWith',
    'NEWCollabToolsWorkedWith', 'NEWDevOps', 'NEWEdImpt', 'NEWJobHuntResearch', 'NEWLearn',
    'NEWOftTopic', 'NEWOnboardGood', 'NEWOtherComms', 'NEWOvertime', 'NEWPurpleLink',
    'NEWStuck', 'OrgSize', 'PlatformDesireNextYear', 'PlatformWorkedWith', 'SOComm',
    'SOVisitFreq', 'SurveyEase', 'YearsCode', 'YearsCodePro']

# discrete columns (no outlier)
discrete_columns = ['Hobbyist', 'Gender', 'PurchaseWhat', 'Sexuality', 'SOAccount', 'SurveyLength', 'Trans',
    'UndergradMajor', 'WelcomeChange']
```

شکل ۹

به‌منظور برخورد صحیح با داده‌های پرت هر ویژگی، توسط دستورات شکل ۱۱ نمودار scatter مقادیر هر یک را رسم می‌کنیم. بررسی نمودارها نشان می‌دهد که مقادیر بیشتر ستون‌ها از الگوی خاصی تبعیت نمی‌کنند و از این‌رو نمی‌توان گفت که داده‌های آن‌ها رفتار غیرطبیعی دارند. تنها داده‌ای در ستون Age نادرست (مقداری بیشتر از ۲۵۰) است که طبق دستورات شکل ۱۲ آن را حذف می‌کنیم. داده‌های پرت ستون‌های دیگر را با استفاده از روش IQR با میانه مقادیر هر ستون جایگزین می‌کنیم.

```
# outlier removal
for row in range(len(encod_ans['Age'])):
    if encod_ans.loc[row, 'Age'] > 250.0:
        remove = encod_ans.loc[row, 'Respondent']
encod_ans = encod_ans.set_index('Respondent').drop(remove)
```

شکل ۱۲

```
# detecting outliers
for i in outlier_columns:
    plt.plot(encod_ans[i], 'o')
    plt.title(i)
    plt.show()
```

شکل ۱۱