

Assignment on Dataset using Numpy and Panda

Name :- Soham Kand

PRN :- 202401070215

Roll No :- 63

Division :- ET2

**Subject :- EDS Assignment On Dataset using
Numpy and Panda**

Dataset :- Twitter US Airline Sentiment.

**Formulate 20 problem statements for a given
dataset using Numpy and Pandas and Apply**

Numpy and pandas methods to find the solution for the formulated problem statements. Datasheet is Twitter US Airline Sentiment.

First, a brief about the Twitter US Airline Sentiment dataset:

This dataset contains tweets about airlines in the United States labeled as *positive*, *negative*, and *neutral*. Main columns include:

- `tweet_id`
- `airline_sentiment`
(positive/negative/neutral)
- `airline`
- `airline_sentiment_confidence`
- `negativereason` (if sentiment is negative)
- `negativereason_confidence`
- `text`
- `tweet_created`
- `user_timezone` etc.

❖ Here's the list:

N
o
.

Solution using Pandas/Numpy

**Find total
number**

1 of tweets `len(df)` or `df.shape[0]`
in the
dataset.

**Count the
number
of**

2 positive, `df['airline_sentiment'].value_`
negative, `counts()`
and
neutral
tweets.

**Find the
airline**

3 with the
most
tweets.

`df['airline'].value_counts().i`
`dxmax()`

**Find the
mean**

4 sentiment
confidenc
e score

`df['airline_sentiment_confiden`
`ce'].mean()`

N
o
.

Solution using Pandas/Numpy

for all tweets.

Find the airline with the 5 highest average sentiment confidence.

```
df.groupby('airline')['airline  
_sentiment_confidence'].mean()  
.idxmax()
```

Find the tweet with the 6 highest sentiment confidence.

```
df.loc[df['airline_sentiment_c  
onfidence'].idxmax()]
```

List top 5 reasons 7 for negative tweets.

```
df['negativeresason'].value_  
counts().head(5)
```

Calculate the 8 df['airline_sentiment'].value_ counts(normalize=True) * 100

N
o
.

Solution using Pandas/Numpy

percentag
e of
tweets for
each
sentiment

.

Find how
many
tweets

9 have df['negativereson'].isna().su
missing m()
negativ
ereason
values.

Find the
most

1 common df['user_timezone'].value_coun
0 user ts().idxmax()
timezone
among all
tweets.

1 Find df['airline_sentiment_confiden
1 standard ce'].std()
1 deviation

N
o
.

of
sentiment
confidenc
e.

Identify
airlines
1 which
2 have only
positive
tweets.

For each
airline,
count
1 how
3 many
negative
tweets
exist.

Calculate
the total
1 number
4 of tweets
created
per day.

Solution using Pandas/Numpy

```
df.groupby('airline')[['airline  
_sentiment']].nunique() (then  
filter where only 'positive')
```

```
df[df['airline_sentiment']=='n  
egative'].groupby('airline').s  
ize()
```

```
df['tweet_created'].dt.date.va  
lue_counts()
```

N

o

Solution using Pandas/Numpy

.

Find the
tweet(s)

1 posted at `df.loc[df['tweet_created'].idxmin()]`
earliest time.

Find
average
negativer

1 reason_code == 'negative'
2 confidence == 'confidence'
3 where reason_code == 'confidence'.mean()
sentiment is negative.

Find the
number
of unique tweet_ids
users who
retweeted.

1 which reason_code == 'negative'
8 reason_code == 'negative'.value_counts().idxmax()

N
o
.

Solution using Pandas/Numpy

most
common
for
'United'
airline.

Check if
there are

any
1 duplicate df['text'].duplicated().sum()
9 tweets
based on
text.

Find
which
airline
has the
2 highest (group airline & calculate ratio
0 proportion positive/total)
n of
positive
tweets.

Here's an example of code to perform this:

```
python
CopyEdit
import pandas as pd
import numpy as np

# Load the dataset
df = pd.read_csv('Tweets.csv')      #
Adjust if your file name differs

# 1. Total tweets
total_tweets = len(df)

# 2. Sentiment count
sentiment_counts =
df['airline_sentiment'].value_counts()

# 3. Airline with most tweets
most_tweeted_airline =
df['airline'].value_counts().idxmax()

# 4. Mean sentiment confidence
mean_confidence =
df['airline_sentiment_confidence'].mea
n()

# 5. Airline with highest avg
sentiment confidence
```

```
best_airline_confidence =
df.groupby('airline')['airline_sentiment_confidence'].mean().idxmax()

# 6. Tweet with highest sentiment confidence
highest_confidence_tweet =
df.loc[df['airline_sentiment_confidence'].idxmax()]

# 7. Top 5 negative reasons
top_negative_reasons =
df['negativereson'].value_counts().head(5)

# 8. Sentiment percentage
sentiment_percentage =
df['airline_sentiment'].value_counts(normalize=True) * 100

# 9. Missing negativereson
missing_negativereson =
df['negativereson'].isna().sum()

# 10. Most common timezone
common_timezone =
df['user_timezone'].value_counts().idxmax()
```

```
# 11. Std deviation of sentiment
confidence
std_confidence =
df['airline_sentiment_confidence'].std
()

# 12. Airlines with only positive
tweets
positive_airlines =
df[df['airline_sentiment']=='positive'
]['airline'].value_counts()

# 13. Negative tweets per airline
negative_tweets_airline =
df[df['airline_sentiment']=='negative'
].groupby('airline').size()

# 14. Tweets per day
df['tweet_created'] =
pd.to_datetime(df['tweet_created'])
tweets_per_day =
df['tweet_created'].dt.date.value_coun
ts()

# 15. Earliest tweet
earliest_tweet =
df.loc[df['tweet_created'].idxmin()]
```

```
# 16. Average negative reason confidence
negativereson_confidence
avg_neg_reason_conf =
df[df['airline_sentiment']=='negative']
]['negativereson_confidence'].mean()

# 17. Unique users
unique_users =
df['tweet_id'].nunique()

# 18. Most common negative reason for United
common_neg_reason_united =
df[(df['airline']=='United') &
(df['airline_sentiment']=='negative')]
['negativereson'].value_counts().idxmax()

# 19. Duplicate tweets based on text
duplicate_tweets =
df['text'].duplicated().sum()

# 20. Airline with highest positive ratio
positive_ratio =
(df[df['airline_sentiment']=='positive']
```

```
'].groupby('airline').size() /  
df.groupby('airline').size()).idxmax()
```