

DATASET:

business_id											
A	B	C	D	E	F	G	H	I	J	K	L
date	review_id	stars	text	type	cool	useful	funny				
hKry9yApeilPP0Uj1Etmvkg	26-01-2011	FWKX83p0-ka4J5d3c6E5A	5	My wife took me here on my birthday for	review	r1t8r2KDx5vH5nAx9C3q5Q	2	5	0		
RJwVLvYzElq1VAihDhYiow	27-07-2011	iJ233zrXzUx-0XG08NWyA	5	I have no idea why some people give bad	review	0a2kyEL0d3Yb1V6aivbluQ	0	0	0		
oRAC4UyJCSj1X0WZpV5SA	14-06-2012	IESLbZqUdLS5s4m0eCsYQ	4	I love the gyro plate. Rice is so good and I also di	review	0hT2Kt1PhobPvh6cDC8JQg	0	1	0		
1QZQZ4fz2QyCvXco0o6Vg	27-05-2010	G-WGwAlBsqqaMHlNn06oVg	5	Rosie, Dakota, and I LOVE Chaparral D	review	uZ6t9fT0ncRCOG0YfUghhg	1	2	0		
iozyU1r1qPz2Z-1Br0vTw	05-01-2012	1u1Fqzr5QfJG_6EMRCAGo	5	General Manager Scott Petello is a good egg!!!	review	vYmM4Kt3c8Zf0Bg-J5mWkU	0	0	1		
#NAME?	12-12-2007	m2CKSeep1CQRWxUirUxSg	4	Quiescence is, simply put, beautiful. Full	review	oY3N1NtngPbPCTRsMfuZ7g	4	3	1		
pW0Ne_hH8d9KJCjnrw1xA	12-02-2010	rIFQ3XnpP4aFWLk_C5r2a	5	Drop what you're doing and drive here. After l	review	vFwWvHw2lFREZv_dyKz_1g	7	7	4		
ip0NNEAGFGf1AdmrR-n	12-07-2012	1J7G9J9u4YmX7rZs05NfIQ	4	Luckily, I didn't have to travel far to make my	review	1ueiurK57sZeaV_ U15AB13A	1	1	0		
vUueA3XWZD6b3h0Qa0H-g	17-08-2012	XtnfYmmlY71uY5GxiUIA	4	Definitely come for Happy hour! Prices are ama	review	Vh_Dlilzg3G5Qh4qfZ2b6A	0	0	0		
mVHuiYan8e3cOno3PorrnJA	11-08-2010	1JAXA46pU1swYrCdfKtQ	5	Nobuo shows his unique talents with everything	review	sUNKXg8-KfTCMD0V6zRzQg	0	1	0		
uSCv0q_WBqle3mX2IqsOQ	16-06-2010	E11jzP9K5W5K7fUaRWfRw	5	The oldish man who owns the store is as sweet	review	#NAME?	1	3	1		
uA5N4XjdH4j4KtCCQpC_vg	21-10-2011	3PrD0L77rgmE-Urzn0H22w	5	Wonderful Vietnamese sandwich shoppe. Their	review	C1rH3p3dmepNeA7XiouwB6Q	1	1	0		
i53YU1CIDIEFSJCQpk8v1g	11-01-2010	cGnKNX3j9rthE0-TH24-qA	5	They have a limited time thing going on right	review	UPtYsD6fUDU0XqY2Y-6Dcg	1	2	0		
NGNIYMexPxyoWav1APUq7Ja	23-12-2011	FvEEw1_OsrYdwlV5HrliH	4	Good tattoo shop. Clean space, multiple artists	review	Xm8BXE1JHqsc5e5BKf0GfQ	1	2	0		
cSA59H71xIdJA_12mChLA	20-05-2010	pfWvBKYYmAUeXwihDlUcQw	4	I'm 2 weeks new to Phoenix. I looked up Irish	review	JQg-4ge4Bae3lx_szHrR8g	1	1	0		
b9PFC6L6J24PNkLbaFAw	20-03-2011	HxmqdWcerVWO03G56zBr0w	2	Was it worth the 21\$ for a salad and small	review	WvOJ2y7T2V2zeYwHu2uQA	0	2	0		
uipgcPNO91Ko6oLatNTN-g	12-10-2008	HqkD_OUL-FcmA4-F8cQvqaQ	3	We went here on a Saturday afternoon and	review	5BbTlYfYK10MF0vT0tUg	4	4	2		
c510R6e8m9Qydu49U0ATKgc	03-05-2010	J451rZly0WrmW4y4g-Kng	5	okay this is the best place EVER! I grew up shopi	review	u1KwCbPmWkFxEyKZ0YtKq	0	6	0		
uSc5eKR81QqY2T_O0LQ	06-03-2009	vQc2DPnqYKTYCkG5oPcGA	3	I met a friend for lunch yesterday.	review	Uu1Lqf4yKABRMzsd8QzcsA	5	1	0		
Ijz2b5bK9wml0BZWYfuCg	17-11-2011	a0iUCj-z5k_kHQz5_eNqg	4	They've gotten better and better for me in the	review	nDBly08J5UrmrHQ2JCbyiH	1	1	1		
dfN04D3eozpJl0k3s5Zbg	10-08-2008	MuqGtUfSDD1PCZ21V3aQ	3	DVAP....	review	C6i0taaYdLTf5W7dZUyLA	2	4	1		
ckJyKfLMKAsvRjURNOKcg	28-06-2011	lmqVtFH03U213BVKvNjWrXa	5	This place shouldn't even be reviewed - because	review	YN3ZLD0g8kpnf6bChicE2A	1	1	2		
iFA9dqXT5EA_TmRgbo03QQ	03-07-2011	CQ9yC8hgKXv4enAp0kDl0hA	5	first time my friend and I went there... it was d	review	6iG5SR1p23VhVjEBX8JNfA	0	0	0		
J0068bbJfAB6GefBqfB8eQ	19-09-2010	Dx95fUGzn06YQoKijgm-j	1	U can go there n check the car out. If u wanna t	review	RQEDYEDV_HkPqV53hNffw	0	1	0		
hupPwFNfNMjivwWB5druaU	22-05-2011	cFqNkZ2rV2DpYBde_TxmLA	5	I love this place! I have been coming here for	review	13xj6F5YV00rZV5RzSp4w	0	1	0		
v2p2PvSp0dA0h0nJiyMA	26-05-2010	ChBeixKzenFfK0eOoCldbA	4	This place is great. A nice little ol' fashion	review	r1t8r2KDx5vH5nAx9C3q5Q	0	0	0		
uG2B7y8j2pYvYvYvYvYvYv	03-03-2012	1J7AT-ATG-ATG-ATG-ATG-ATG	5	I have been to this place many times and I love	review	6uN1z0y0c0hN0YU1u1u1u1	0	0	0		

EDS Theory assignment1.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

Files

Analyze your files with code written by Gemini Upload

sample_data

- README.md
- anscombe.json
- california_housing_test.csv
- california_housing_train.csv
- mnist_test.csv
- mnist_train_small.csv
- yelp.csv

Disk 70.72 GB available

```
[1] import numpy as np
import pandas as pd

df = pd.read_csv('/content/sample_data/yelp.csv')

#1 What is the distribution of star ratings?
star_distribution = df['stars'].value_counts().sort_index()
print("Star Rating Distribution:")
print(star_distribution)
```

Star Rating Distribution:

stars	
1	749
2	927
3	1461
4	3526
5	3337

Name: count, dtype: int64

+ Code + Text

```
[7] #2 What are the most common words in positive reviews (5-star)?
from collections import Counter
import re

positive_reviews = df[df['stars'] == 5]['text']
words = ' '.join(positive_reviews).lower()
words = re.findall(r'\w+', words)
common_words = Counter(words).most_common(20)
print("\nMost common words in 5-star reviews:")
print(common_words)
```

EDS Theory assignment1.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text

Files

Analyze your files with code written by Gemini Upload

sample_data

- README.md
- anscombe.json
- california_housing_test.csv
- california_housing_train.csv
- mnist_test.csv
- mnist_train_small.csv
- yelp.csv

Disk 70.72 GB available

Most common words in 5-star reviews:

```
[8] #3 What is the average length of reviews by star rating?
df['review_length'] = df['text'].apply(len)
avg_length_by_stars = df.groupby('stars')['review_length'].mean()
print("\nAverage review length by star rating:")
print(avg_length_by_stars)
```

Average review length by star rating:

stars	
1	826.515354
2	842.256742
3	758.498289
4	712.923142
5	624.999101

Name: review_length, dtype: float64

```
#4 Which businesses have the most reviews?
top_businesses = df["business_id"].value_counts().head(10)
print("\nBusinesses with most reviews:")
print(top_businesses)
```

Businesses with most reviews:

business_id	
JokKtdXU7zXHcr20Lrk29A	37
ntN85eu27C04nwyPa8IHtw	37

EDS Theory assignment1.ipynb

File Edit View Insert Runtime Tools Help

Commands

+ Code + Text

Files

Analyze your files with code written by Gemini

Upload

sample_data

README.md

anscombe.json

california_housing_test.csv

california_housing_train.csv

mnist_test.csv

mnist_train_small.csv

yelp.csv

[14] #9 Which users are the most active reviewers?

```

top_reviewers = df['user_id'].value_counts().head(10)
print("\nMost active reviewers:")
print(top_reviewers)

```

Most active reviewers:

user_id	
fczQCSmaWf78toLEmb0Zsw	38
0CMz8Ya03f8xu4KqQgKb9Q	25
90a6z--Curl84aCzZyPsg	22
0mqJhdkEdak_A1FBhFNxqA	18
4ozupHULqcyO42s3zNUzOQ	18
wHg1YkCzdq9wB30TRgxHQ	17
PzSmcfrCjebXSLXR0mngQ	16
0bIXP9quoJegyVzu9ipGgQ	16
uzbtb-u-GVjTa2gtQfry5g	15
joiZw_aUINv8TuGoytrH7g	15

Name: count, dtype: int64

#10 How has the number of reviews changed over time?

```

df['date'] = pd.to_datetime(df['date'])
reviews_over_time = df.set_index('date').resample('M').size()
print("\nReviews over time:")
print(reviews_over_time.head())

```

Reviews over time:

date	
2005-04-30	1
2005-05-31	0
2005-06-30	0

EDS Theory assignment1.ipynb

File Edit View Insert Runtime Tools Help

Commands

+ Code + Text

Files

Analyze your files with code written by Gemini

Upload

sample_data

README.md

anscombe.json

california_housing_test.csv

california_housing_train.csv

mnist_test.csv

mnist_train_small.csv

yelp.csv

Reviews over time:

date	
2005-04-30	1
2005-05-31	0
2005-06-30	0
2005-07-31	2
2005-08-31	0

Freq: ME, dtype: int64

<ipython-input-15-21a79b5be145>:3: FutureWarning: 'M' is deprecated and will be removed in a future version

reviews_over_time = df.set_index('date').resample('M').size()

[16] #11 What is the average star rating over time?

```

avg_stars_over_time = df.set_index('date').resample('M')['stars'].mean()
print("\nAverage star rating over time:")
print(avg_stars_over_time.head())

```

Average star rating over time:

date	
2005-04-30	5.0
2005-05-31	NaN
2005-06-30	NaN
2005-07-31	2.5
2005-08-31	NaN

Freq: ME, Name: stars, dtype: float64

<ipython-input-16-ff5968880cb>:2: FutureWarning: 'M' is deprecated and will be removed in a future version

avg_stars_over_time = df.set_index('date').resample('M')['stars'].mean()

#12 What are the most common bigrams in positive reviews?

```

from nltk import bigrams

```

EDS Theory assignment1.ipynb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text

Files

Analyze your files with code written by Gemini Upload

sample_data

- README.md
- anscombe.json
- california_housing_test.csv
- california_housing_train.csv
- mnist_test.csv
- mnist_train_small.csv
- yelp.csv

Disk 70.72 GB available

[17] #12 What are the most common bigrams in positive reviews?

```
from nltk import bigrams

positive_text = ' '.join(df[df['stars'] == 5]['text'].tolist()).lower()
words = re.findall(r'\w+', positive_text)
common_bigrams = Counter(bigrams(words)).most_common(10)
print("\nMost common bigrams in positive reviews:")
print(common_bigrams)
```

Most common bigrams in positive reviews:

```
[('of', 'the'), 1302], (('in', 'the'), 1133), (('and', 'the'), 1127), (('it', 's'), 1119), (('this', 'plac
```

[13] #13 How does the sentiment of reviews correlate with star ratings?

```
from textblob import TextBlob

df['sentiment'] = df['text'].apply(lambda x: TextBlob(x).sentiment.polarity)
sentiment_by_stars = df.groupby('stars')['sentiment'].mean()
print("\nAverage sentiment by star rating:")
print(sentiment_by_stars)
```

Average sentiment by star rating:

```
stars
1    -0.017906
2     0.096832
3     0.193733
4     0.277061
5     0.333293
Name: sentiment, dtype: float64
```

EDS Theory assignment1.ipynb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text

Files

Analyze your files with code written by Gemini Upload

sample_data

- README.md
- anscombe.json
- california_housing_test.csv
- california_housing_train.csv
- mnist_test.csv
- mnist_train_small.csv
- yelp.csv

Disk 70.72 GB available

[19] #14 What is the distribution of review lengths?

```
print("\nReview length statistics:")
print(df['review_length'].describe())
```

Review length statistics:

```
count    10000.000000
mean       710.738700
std       617.399827
min         1.000000
25%       294.000000
50%       541.500000
75%       930.000000
max      4997.000000
Name: review_length, dtype: float64
```

#15 Which businesses have the highest average star ratings (with at least 5 reviews)?

```
business_stats = df.groupby('business_id').agg(
    avg_stars=('stars', 'mean'),
    review_count=('stars', 'count')
)
top_businesses = business_stats[business_stats['review_count'] >= 5].sort_values('avg_stars', ascending=False)
print("\nBusinesses with highest average ratings (min 5 reviews):")
print(top_businesses)
```

Businesses with highest average ratings (min 5 reviews):

```
business_id  avg_stars  review_count
rIAeltELaGdQKh_LYlLEA  5.000000         5
lVkdZ8StafPlq7fii5row  5.000000         6
Q-Xa9GCEMT65YiBD5TW_hA  5.000000         7
```

EDS Theory assignment1.ipynb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text

Files

Analyze your files with code written by Gemini Upload

sample_data

README.md

anscombe.json

california_housing_test.csv

california_housing_train.csv

mnist_test.csv

mnist_train_small.csv

yelp.csv

Output

21

#16 What is the relationship between 'cool', 'funny', and 'useful' votes?

vote_correlation = df[['cool', 'funny', 'useful']].corr()

print("\nCorrelation between vote types:")

print(vote_correlation)

Correlation between vote types:

cool funny useful

cool 1.000000 0.764342 0.887102

funny 0.764342 1.000000 0.723406

useful 0.887102 0.723406 1.000000

22

#17 How do the longest reviews compare to the shortest in terms of star ratings?

longest_reviews = df.nlargest(10, 'review_length')

shortest_reviews = df.nsmallest(10, 'review_length')

print("\nAverage stars for longest reviews:", longest_reviews['stars'].mean())

print("Average stars for shortest reviews:", shortest_reviews['stars'].mean())

Average stars for longest reviews: 3.5

Average stars for shortest reviews: 4.2

23

#18 What is the distribution of review dates?

EDS Theory assignment1.ipynb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text

Files

Analyze your files with code written by Gemini Upload

sample_data

README.md

anscombe.json

california_housing_test.csv

california_housing_train.csv

mnist_test.csv

mnist_train_small.csv

yelp.csv

Output

23

#18 What is the distribution of review dates?

print("\nReview date range:")

print("Earliest:", df['date'].min())

print("Latest:", df['date'].max())

Review date range:

Earliest: 2005-04-18 00:00:00

Latest: 2013-01-05 00:00:00

24

#19 How many unique businesses and users are in the dataset?

unique_businesses = df['business_id'].nunique()

unique_users = df['user_id'].nunique()

print(f"\nUnique businesses: {unique_businesses}")

print(f"Unique users: {unique_users}")

Unique businesses: 4174

Unique users: 6403

25

#20 What is the average number of words per review by star rating?

df['word_count'] = df['text'].apply(lambda x: len(x.split()))

avg_words_by_stars = df.groupby('stars')['word_count'].mean()

print("\nAverage word count by star rating:")

print(avg_words_by_stars)

Average word count by star rating:

stars

1 153.953271

2 156.435814

EDS Theory assignment1.ipynb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text

Files

Analyze your files with code written by Gemini Upload

<> ↻ 📄 🔍

{x}

sample_data

README.md

anscombe.json

california_housing_test.csv

california_housing_train.csv

mnist_test.csv

mnist_train_small.csv

yelp.csv

Disk 70.72 GB available

[25] #20 What is the average number of words per review by star rating?

df['word_count'] = df['text'].apply(lambda x: len(x.split()))
avg_words_by_stars = df.groupby('stars')['word_count'].mean()
print("\nAverage word count by star rating:")
print(avg_words_by_stars)

Average word count by star rating:
stars
1 153.953271
2 156.435814
3 140.714579
4 131.174135
5 114.463590
Name: word_count, dtype: float64

[] Start coding or generate with AI.