

data projects →

Data is better with a little magic: New updates to Hex's AI tools for data analysis →



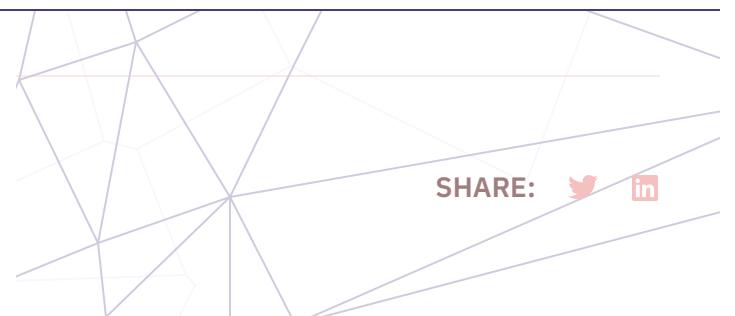
G

HEX

R Click below to copy this free template
to your workspace.

Copy this template

SHARE:



HEX | Dumpling Forecast

ALL SOURCES

Demo Snowflake (2 rows)

Our public demo data stack, including dbt integrations, a dbt project, and more!

Hex Content + 0. Query

Search DBs, schemas, tables, or

ANALYTICS > PROD

ORDER_ITEM_ID NUMBER

Unique identifier for an order item, orders can have one or more items

ORDER_STATUS VARCHAR

One of completed or cancelled

CUSTOMER_ID NUMBER

Unique identifier for a Dumpling Shack customer

ORDERED_AT TIMESTAMPTZ

Date and time the order was placed

ORDER_DATE DATE

Date of the order

MENU_ITEM VARCHAR

String describing the menu item

MENU_ITEM_ID NUMBER

Unique identifier for a menu item

CATEGORY VARCHAR

Menu item category, one of Quick Bites, Veggies, Dumplings, Noodles, or Sweets

PRICE DOUBLE

Cost of a menu item in USD

IS_SPICY BOOLEAN

True if the menu item is considered spicy

IS_ALLERGEN_FREE BOOLEAN

True if the menu item does not contain gluten, shellfish, fish, peanuts, eggs, soy,...

Add cell

Dataframe ✓ 609 rows 0 seconds 49.23 KB Cached October 13th 2022, 5:06 pm

	MONTH	CATEGORY	IS_SPICY	ORDER_TOTAL
0	2019-12-01 00:00:00+00:00	Quick Bites	False	554.5
1	2017-08-01 00:00:00+00:00	Veggies	False	352.0
2	2019-05-01 00:00:00+00:00	Dumplings	True	171.0
3	2018-09-01 00:00:00+00:00	Dumplings	False	3793.0
4	2017-01-01 00:00:00+00:00	Sweets	False	909.0
5	2016-08-01 00:00:00+00:00	Dumplings	False	4072.5
6	2016-12-01 00:00:00+00:00	Sweets	False	1157.0
7	2016-10-01 00:00:00+00:00	Dumplings	False	3948.5
8	2017-01-01 00:00:00+00:00	Veggies	False	348.0
9	2019-10-01 00:00:00+00:00	Dumplings	False	4108.5

dumpling_orderx 1

SQL 1

SOURCE Demo Snowflake

View template

```
1 select month, category, sum(order_total) from dumpling_orders group by
```

Dataframe

ADD CELL | Python | SQL | Text | Transform | Display | Input parameters

Izzy Miller Hex Content





Click below to copy this free template
to your workspace.



Copy this template



vs.

customers, understand their
es, leading to improved satisfaction



- ~~Finance~~. Track market sentiment towards specific stocks or industries, informing investment decisions and risk management.

- **Politics:** Analyze public opinion on political candidates, policies, and events, providing valuable insights for campaigns and governance.
- **Healthcare:** Understand patient sentiment towards treatments, doctors, and healthcare experiences, improving patient care and communication.

How do we unlock this power of sentiment analysis? Python is to go-to tool for building a own sentiment analysis engine. Its ecosystem of libraries and frameworks like NLTK, TextBlob, and spaCy make it ideal for natural language processing tasks.

Here, we want to show you how to use Python and Hex to perform sentiment analysis. We'll cover everything from the foundational concepts to the practical



Click below to copy this free template
to your workspace.



Copy this template



cages like NLTK, TF-IDF, SpaCy, or
use statements in SQL directly

p review data. The data is stored in
the Snowflake warehouse that can be accessed within the Hex environment.
We'll focus on two main aspects of the Yelp review data:

- 1. Narrative Structure:** Analyzing the narrative structure of a review, and its correlation with sentiment. We hypothesize that reviews often tell a story about the reviewer's experience and that this narrative structure might be associated with the sentiment of the review.
- 2. Cultural Trends:** Analyzing all reviews over time to uncover cultural trends. We hypothesize that certain types of food or dining experiences might be associated with more positive reviews during their peak in popularity (ex: avocado toast).

Let's start with downloading the necessary dependencies for sentiment analysis.



Click below to copy this free template
to your workspace.

need to install some additional d
downloaded with the help of the fo



Copy this template



5m

Next, we need to load all the dependencies that are going to be used in this article.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.feature_extraction.text import TfidfVectorizer
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize
import nltk
import re
```



HEX

Click below to copy this free template to your workspace.



Copy this template



```
cuda/cudart_stub.cc:28] Could not find cu
/platfrom/cpu_feature_guard.cc:182] This T
X512F FMA, in other operations, rebuild Te
iler/tf2tensorrt/utils/py_utils.cc:38] TF-
e/hexuser/nltk_data...
ome/hexuser/nltk_data...
eptron_tagger to
eptron_tagger.zip.
```

```
[nltk_data] Downloading package stopwords to
[nltk_data]     /home/hexuser/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
```

As you can see in the above code, a lot of text preprocessing, modeling, and visualization libraries are imported. Then the additional data such as `punkt`, `wordnet`, `averaged_perceptron_tagger`, and `stopwords` is downloaded from the NLTK library. Finally, the tokenizer, tagger, parser, NER, and word vectors are loaded from the spacy library which will be used further for sentiment analysis.

Data Loading and Preprocessing

Now, we need to load the dataset from the warehouse and process it. We need to perform the text cleaning (removing punctuation, converting to lowercase, etc.) and calculate the sentiment scores for each review as part of the data

The screenshot shows the Hex platform interface. At the top, there are several orange icons: a square with a diagonal line, a triangle, a hexagon, a double arrow, a circle with a dot, and a circle with a minus sign. Below these is a large red 'HEX' logo. On the right side, there are more icons: a circle with a dot, a circle with a minus sign, a circle with a plus sign, and a circle with a question mark. A vertical bar labeled 'G' is also visible.

Below the icons, there is a section with a purple border containing text and two circular buttons. The text reads: "Click below to copy this free template to your workspace." Below this text are two purple circular buttons with white icons: one with a left arrow and one with a right arrow. Between them is the text "Copy this template".

To the right of this section is a dark grey rectangular area containing a table with five columns of review data. The columns are labeled: "0:00", "Definitely come for Happy hour! Prices are amazing, sak", "0:00", "Nobuo shows his unique talents with everything on the", "0:00", "The oldish man who owns the store is as sweet as can b", "0:00", "Wonderful Vietnamese sandwich shoppe. Their baguett", and "0:00", "They have a limited time thing going on right now with r". At the bottom right of this area, there are two small icons: a clipboard and a downward arrow, followed by the text "2,29,907 rows".

As you can see in the above image, the dataset consists of columns like `review_id`, `stars`, `date`, `text`, and so on.

Sentiment Analysis

Once the dataset is loaded, we need to start parsing and classifying the reviews in this dataset. For this, we need to use the TextBlob library of Python that returns two metrics for every text that is passed to it as input. The first one is Polarity, which indicates the positivity/negativity in the sentiment of the text. The second one is subjectivity which refers to objective info/facts versus personal opinions or emotions.

You simply need to load the textblob library and pass your text as input to get the polarity and subjectivity as follows:



HEX

Click below to copy this free template
to your workspace.



Copy this template



				give bad review...				
IESLBzqUCLdSzSqm0eCSxQ	4	2012-06-14	love the gyro plate. Rice is so good and I als...	review	6oRAC4uyJCsjl1X0WZpVSA	0.566667	0.733333	
G-WvGaISbqqaMHlNnByodA	5	2010-05-27	Rosie, Dakota, and I LOVE Chaparral Dog Park!...	review	_1QQZuf4zZ0yFCvXc0o6Vg	0.608646	0.700000	

As you can see in the above code, the `get_sentiment()` method returns the polarity and subjectivity scores, Then we store these scores for all the reviews in the `polarity` and `subjectivity` columns of our original dataframe.

Narrative Structure



Click below to copy this free template
to your workspace.



Copy this template



sentiment arc of a review



sentences

return polarities

[polarity for s in sentences]

```
# Calculate the sentiment arcs for the first 100 reviews
arcs = reviews['text'].head(100).map(get_sentiment_arc)

# Plot the sentiment arcs
plt.figure(figsize=(10, 6))
for arc in arcs:
    plt.plot(np.linspace(0, 1, len(arc)), arc, alpha=0.1, color='blue')
plt.xlabel('Position in review')
plt.ylabel('Sentiment polarity')
plt.title('Sentiment Arcs of Yelp Reviews')
plt.show()
```



Click below to copy this free template to your workspace.



Copy this template



using the `get_sentiment_arc()` function on the first 100 reviews. Then we created a

between the sentiment arc and the overall sentiment, this can be done with the help of the following code:

```
# Define a function to calculate the correlation between the sentiment arc and the overall sentiment
def get_correlation(arc, overall_sentiment):
    return np.corrcoef(arc, np.linspace(0, 1, len(arc)) * overall_sentiment)[0, 1]

# Calculate the correlations for the first 100 reviews
correlations = [get_correlation(arc, sentiment) for arc, sentiment in zip(arcs, reviews['polarity'].head(100))]
```

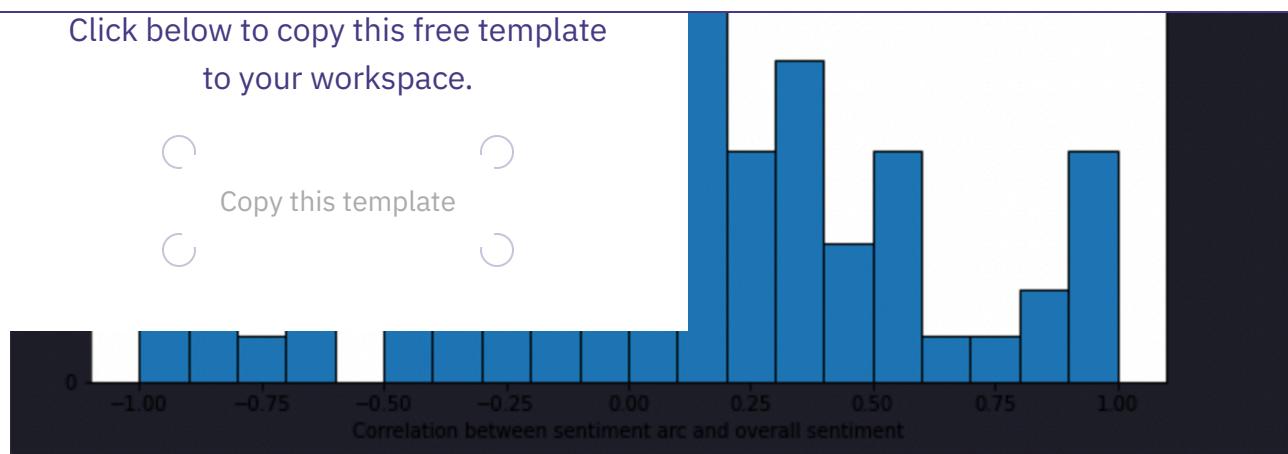
In the above code, we are simply iterating over the sentiment arc and polarity score (sentiment) of each review to calculate a correlation coefficient value. We are storing the correlation into a list named correlations. Next, you can plot a simple bar graph to visualize these correlations as follows:

**HEX**

Click below to copy this free template
to your workspace.



Copy this template



Upon inspection, the sentiment arcs of reviews are diverse and vary widely. There's a weak positive correlation between the sentiment arc and an overall positive sentiment— meaning that for overall positive reviews, the positivity of the review increases from the beginning to the end of the review and vice versa for negative ones.

However, it's a pretty weak effect, and the wide distribution means there are probably other factors at play in determining the overall sentiment of a review.

Cultural Trends



G



Click below to copy this free template

to your workspace.



ch term

Copy this template



| .str.contains(term, case=False)

```
# display the first few rows of the dataframe  
reviews.head()
```

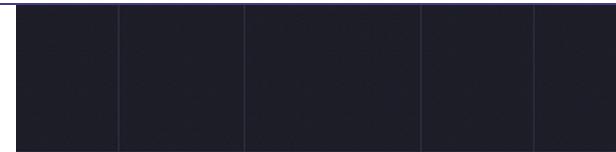


HEX

Click below to copy this free template
to your workspace.



Copy this template



te a new column in the dataset that
ns appeared. Then we will calculate
an aggregated mean of the polarity score for the review containing these terms.
Finally, we will create a line chart for visualizing the average sentiment of Yelp
reviews mentioning cultural terms.

```
# Convert the date column to datetime
reviews['date'] = pd.to_datetime(reviews['date'])

# Extract the year from the date
reviews['year'] = reviews['date'].dt.year

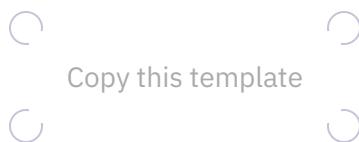
# Calculate the average sentiment for each term for each year
average_sentiments = {}
for term in terms:
    average_sentiments[term] = reviews[reviews[term]].groupby('year')[['polarity']].mean()

# Plot the average sentiment over time
plt.figure(figsize=(10, 6))
```



HEX

Click below to copy this free template
to your workspace.



2005 2006 2007 2008 2009 2010 2011 2012 2013

As you can see, the sentiment of the reviews containing the term `avocado` is increasing with time while with the `kale` term, it is decreasing.

you can also visualize the sentiment polarity for each term with the help of the following code:

```
# Calculate the sentiment polarity for each term
sentiment_polarities = {}
for term in terms:
    sentiment_polarities[term] = reviews[reviews[term]]['polarity']

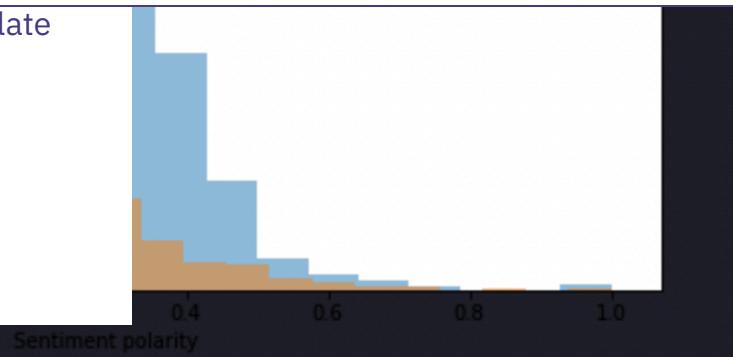
plt.figure(figsize=(10, 6))
for term, polarity in sentiment_polarities.items():
    plt.hist(polarity, bins=20, alpha=0.5, label=term)
```

**HEX**

Click below to copy this free template
to your workspace.



Copy this template



Review Ratings

Finally, we'll quickly look at the reviews-of-reviews given by other readers on whether a review was useful, funny, or cool. Do these correlate to sentiment—or to one another? To do so, you can simply use the `corr()` method to calculate the correlation among features like `polarity`, `useful`, `funny`, and `cool`.

```
reviews[['funny', 'useful', 'cool']] = pd.DataFrame(reviews['votes'].tolist())
```

```
# Calculate the correlation between the sentiment polarity and the number of votes
correlations = reviews[['polarity', 'useful', 'funny', 'cool']].corr()
```



Click below to copy this free template
to your workspace.

ation between `polarity` and `stars`

 Copy this template

sentiment polarity and the numbe

 `stars']]`.corr()

```
# Display the correlation
correlation
```

	polarity	stars
polarity	1.000000	0.496123
stars	0.496123	1.000000
 2 rows		

As you can observe, these two features are correlated which makes sense as sentiment reviews directly reflect the ratings that users provide.

This is it! now you have your own Yelp data sentiment analysis system.

Python Packages for Sentiment Analysis



Click below to copy this free template
to your workspace.



Copy this template



basic preprocessing steps where it
cleaning, tokenization, and

ential features from the text and
converting it into a format suitable for analysis. Techniques like bag-of-
words or TF-IDF (Term Frequency-Inverse Document Frequency) can be
effortlessly implemented using NLTK's functionalities.

- Pre-trained models and classifiers in NLTK categorize text into positive, negative, or neutral sentiments. Custom classifiers can be trained as well.
- NLTK includes lexicons, such as the WordNet sentiment lexicon, which assigns sentiment scores to words. By utilizing these lexicons, NLTK enables a lexicon-based approach to sentiment analysis.
- It seamlessly integrates with machine learning models (e.g., Naive Bayes), allowing users to train models on labeled datasets for improved sentiment predictions.



Click below to copy this free template
to your workspace.



Copy this template



ke breaking down sentences into

features from the text, providing

- ~~its parsing capability helps understand~~ how words in a sentence relate.
- Spacy identifies and categorizes entities (like names or locations), enhancing analysis with specific elements.

TextBlob

TextBlob is a simplified and user-friendly natural language processing library in Python. It incorporates pre-trained models and lexicons, making sentiment analysis accessible without the need for extensive configuration. A standard approach to perform the sentiment analysis with TextBlob is as follows:

- Begin by creating a TextBlob object with the text data you want to analyze. TextBlob automatically processes the text and prepares it for sentiment analysis.

Click below to copy this free template
to your workspace.



Copy this template



ended for text found on social
that it is sensitive to the strength and
it. This tool is a useful resource for
platforms since it excels at handling
' encountered there. The steps of

using Vader for sentiment analysis are listed below.

- Begin by installing the VADER sentiment analysis library in your Python environment. You can install it using a package manager like [pip](#).
- Import the VADER module and initialize the sentiment analyzer.
- Utilize the VADER sentiment analyzer to score the sentiment of the text. VADER provides a compound score that represents the overall sentiment intensity.
- Interpret the compound score to understand the sentiment. Positive scores indicate positive sentiment, negative scores denote negativity, and scores around zero suggest neutrality.
- Optionally, fine-tune the analysis by adjusting parameters or thresholds based on specific requirements or domain-specific characteristics.



Click below to copy this free template
to your workspace.



Copy this template



ures suitable for machine learning.

Scikit-Learn, such as Naive Bayes or
entiment analysis.

Load dataset using Scikit-Learn's
straightforward API, providing text features and sentiment labels.

- Assess the model's performance using Scikit-Learn's functions for metrics like accuracy, precision, recall, and F1 score.
- Once trained, use the model to predict sentiments for new text data, enabling efficient sentiment analysis.

Gensim

Gensim is used for document similarity analysis and topic modeling. Although Gensim's main focus is on unsupervised learning tasks such as topic modeling, it can also be used for sentiment analysis tasks because it has text processing and vector space modeling functionalities. Gensim is especially helpful in situations where sentiment analysis goals are aligned with a topic extraction



Click below to copy this free template
to your workspace.



Copy this template



specific techniques or models.

learning library [TensorFlow](#). Tasks

involving sentiment analysis that require the capability of deep learning architectures are a good fit for TensorFlow. You can build a simple deep-learning neural network for sentiment analysis as follows:

- Prepare a labeled dataset with text samples and corresponding sentiment labels (positive, negative, or neutral).
- Use TensorFlow's text preprocessing tools to convert the text into numerical representations. Common techniques include [tokenization](#) and [padding](#).
- Design a neural network architecture suitable for sentiment analysis. This typically involves [embedding layers](#), [recurrent](#) or [convolutional layers](#), and output layers.
- Compile the model by specifying the [optimizer](#), [loss function](#), and [metrics](#). This step configures the training process.



HEX

Click below to copy this free template
to your workspace.



Copy this template



and corresponding sentiment labels.
onvert text into numerical



analysis, tailoring it to project

requirements.

- Set up the sentiment analysis model using PyTorch's neural network modules.
- Train the model on the dataset using PyTorch's automatic differentiation for efficient backpropagation.
- Evaluate model performance on a validation set using metrics like accuracy, precision, recall, and F1 score.

Sentiment analysis is one of those techniques it is great to have in your toolbox as a data analyst. It allows you to understand how your business, product, or team are perceived by users and public and then act on that information. Given the wealth of tools available to you to build not just basic good/bad sentiment analysis but sophisticated sentiment engines, there is no excuse not to have this technique running on your data all the time.