# Search for a Connection: Energy Demand & Twitter Trending Topics

**EECS 6893**
**Big Data Analytics**
**Final Project Report**

**Group:**
**Kevin Mark Murning (kmm2344)**
**Rifqi Luthfan (rl3154)**
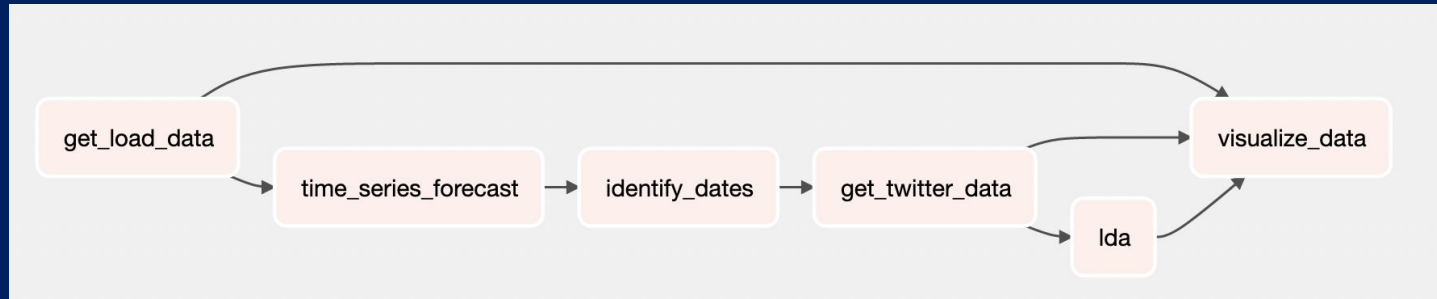**Rohan Raghuraman (rr3417)**

# Contents

1. **Motivation**
2. **System, Experiment Design, Methodology**
   a. **Get Load Data**
   b. **Time Series Forecast**
   c. **Identify Critical Dates**
   d. **Get Twitter Data & do LDA Topic Modelling**
3. **Results and Analysis**
4. **Business Value**
5. **Demo**

Columbia | Engineering
The Fu Foundation School of Engineering and Applied Science
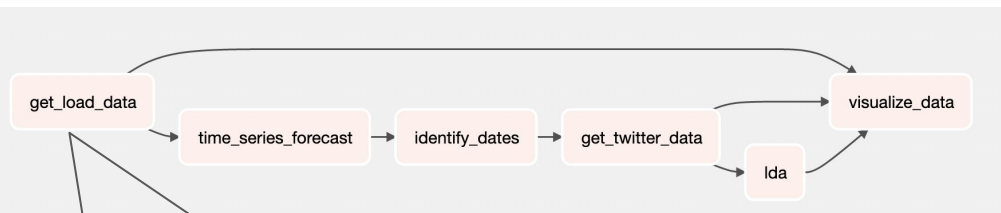
# Motivation

- **Effective energy demand forecasting plays a vital role in power systems**
    - Resource Allocation
    - Economically viable pricing
- **Some trends are easy to infer**
    - An increase in demand in winter months due to use of heating
- **However, outlier events can cause spikes in demand that are hard to predict**
    - How can we predict anomalous changes in demand?
- **Much work has been done on monitoring social media for time series forecasting**
    - Very little has been done in the energy sector
- **Our Question:**
    - *Can Twitter activity formulated into topics be linked in a causal relationship with energy demand spikes?*

# System, Experiment Design and Methodology

# Get Electricity Load Energy Demand Data



1. Web scraping: load data 5 mins granularity from http://dss.nyiso.com/dss_oasis/PublicReports using python requests -> store in Google Cloud Storage as CSV files
2. Create BigQuery database from CSV files in Google Cloud Storage:
   a. base table -> store data as is
   b. Preprocessing data for aggregation table (hourly, daily, weekly, monthly, yearly) -> store data with the needed aggregation so we do not need to always run a heavy query -> used for web data viz
3. **Metrics that we track: size of database**
   a. Hourly aggregation table only take 25 MB of query compared to multiple GBs of raw data
   b. Daily weekly monthly yearly aggregation take at max 1 MB query size

# Time Series Forecast



1. Query from BigQuery using pandas-gbq and convert queried data into Spark RDD
2. Use FBProphet Library for forecasting
   a. Applied forecasting in parallel for different zones in NY State utilizing PySpark's applyInPandas()
3. **Metrics that we track: MAPE, RMSE**
   a. Target MAPE below 10%
   b. RMSE value is used for the next step
4. **Hyperparameters tuned:** changepoint_prior_scale, seasonality_mode, interval_width

# Identify Critical Dates



1. Use the forecast results and evaluation data
   a. Merge forecast and evaluation result, calculate forecast error of each date
   b. Get absolute difference value
2. Analyze trends & seasonality
3. **Create a metric to identify anomalous date based on analysis**
   a. RMSE value is used -> anomalous hour if the difference in actual vs predicted value is more than 5xRMSE value
   b. Anomalous date if 4 or more anomalous hours occurs in a day

Merge forecast and evaluation result, calculate forecast error of each date

```
dates_df = pd.merge(
    forecast_df[(forecast_df['ds']<'2021-11-01')][['ds', 'zone_name', 'y', 'yhat']],
    eval_df[['zone_name', 'mae', 'rmse', 'mape']],
    on=['zone_name'], how='left'
)
dates_df['y_abs_diff'] = np.abs(dates_df['y']-dates_df['yhat'])
dates_df['y_pct_diff'] = dates_df['y_abs_diff']/dates_df['y']

dates_df
```

| | ds | zone_name | y | yhat | mae | rmse | mape | y_abs_diff | y_pct_diff |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 2010-01-01 00:00:00 | CAPITL | 12518.099609 | 12613.607422 | 1338.881470 | 1882.543945 | 0.079469 | 95.507812 | 0.007630 |
| 1 | 2010-01-01 01:00:00 | CAPITL | 13092.299805 | 11989.854492 | 1338.881470 | 1882.543945 | 0.079469 | 1102.445312 | 0.084206 |
| 2 | 2010-01-01 02:00:00 | CAPITL | 12562.200195 | 11588.465820 | 1338.881470 | 1882.543945 | 0.079469 | 973.734375 | 0.077513 |
| 3 | 2010-01-01 03:00:00 | CAPITL | 12101.299805 | 11466.584961 | 1338.881470 | 1882.543945 | 0.079469 | 634.714844 | 0.052450 |
| 4 | 2010-01-01 04:00:00 | CAPITL | 11989.799805 | 11715.201172 | 1338.881470 | 1882.543945 | 0.079469 | 274.598633 | 0.022903 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 1140871 | 2021-10-31 19:00:00 | MILLWD | 3901.063965 | 4167.789062 | 464.994049 | 638.606567 | 0.117714 | 266.725098 | 0.068372 |
| 1140872 | 2021-10-31 20:00:00 | MILLWD | 3805.951416 | 4060.871826 | 464.994049 | 638.606567 | 0.117714 | 254.920410 | 0.066979 |
| 1140873 | 2021-10-31 21:00:00 | MILLWD | 3647.767578 | 3807.325439 | 464.994049 | 638.606567 | 0.117714 | 159.557861 | 0.043741 |
| 1140874 | 2021-10-31 22:00:00 | MILLWD | 3449.541748 | 3461.595947 | 464.994049 | 638.606567 | 0.117714 | 12.054199 | 0.003494 |
| 1140875 | 2021-10-31 23:00:00 | MILLWD | 3095.496094 | 3104.627197 | 464.994049 | 638.606567 | 0.117714 | 9.131104 | 0.002950 |

Trend and seasonality analysis



Based on RMSE

```
filtered_dates = (dates_df[(dates_df['y_abs_diff']>(5*dates_df['rmse']))]
    .groupby(['zone_name', dates_df['ds'].dt.date))
    .agg(anomaly_occurence=('ds', 'nunique'))
    .reset_index())
filtered_dates = filtered_dates[filtered_dates['anomaly_occurence']>3]
filtered_dates#.describe()
```

| | zone_name | ds | anomaly_occurence |
|---|---|---|---|
| 7 | CAPITL | 2011-07-21 | 4 |
| 8 | CAPITL | 2011-07-22 | 5 |
| 14 | CAPITL | 2012-06-20 | 5 |
| 15 | CAPITL | 2012-06-21 | 6 |
| 46 | CAPITL | 2013-07-18 | 5 |
| ... | ... | ... | ... |
| 1954 | WEST | 2013-09-11 | 6 |
| 2028 | WEST | 2016-08-12 | 4 |
| 2071 | WEST | 2019-06-28 | 4 |
| 2073 | WEST | 2019-07-02 | 4 |
| 2120 | WEST | 2021-06-29 | 4 |

# Get Twitter Data & LDA Topic Modelling

1. Get Twitter data on the identified dates
    a. Need full archive search access of Twitter API
    b. Store tweets in Storage (unstructured data)
2. Preprocess Data
    a. Remove: links, username, hashtags, medias
    b. Vectorize: lemmatize, tokenize
3. Do topic modelling with LDA
    a. **Hyperparameters tuned:** number of topics

get_load_data → time_series_forecast → identify_dates → get_twitter_data → visualize_data
get_twitter_data → lda → visualize_data

### Get Tweets

```python
search_url = "https://api.twitter.com/2/tweets/search/all"

# Possible queries:
# 'place:"new york" OR place:"albany" OR place:"niagara" is:verified'

# Use the below links for reference on how to build queries:
# https://developer.twitter.com/en/docs/twitter-api/tweets/search/integrate/build-a-query#availability
# https://developer.twitter.com/en/docs/twitter-api/tweets/search/api-reference/get-tweets-search-all
# https://developer.twitter.com/en/docs/twitter-api/tweets/search/quick-start/full-archive-search

# Optional params: start_time,end_time,since_id,until_id,max_results,next_token,
# expansions,tweet.fields,media.fields,poll.fields,place.fields,user.fields
query_params = {'query': 'place:"manhattan, ny" OR place:"new york, ny" OR place:"brooklyn, ny" OR place:
                'tweet.fields': 'author_id,created_at',
                'max_results': '500',
                'start_time': '2021-02-01T00:00:01Z',
                'end_time': '2021-02-01T23:59:59Z',
                'place.fields': 'country_code',
                'expansions': 'geo.place_id'}


def bearer_oauth(r):
    """
    Method required by bearer token authentication.
    """

    r.headers["Authorization"] = f"Bearer {bearer_token}"
    r.headers["User-Agent"] = "v2FullArchiveSearchPython"
    return r


def connect_to_endpoint(url, params):
    response = requests.request("GET", search_url, auth=bearer_oauth, params=params)
    print(response.status_code)
    if response.status_code != 200:
        raise Exception(response.status_code, response.text)
    return response.json()


json_response = connect_to_endpoint(search_url, query_params)
print(json.dumps(json_response, indent=4, sort_keys=True))
with open('2021-02-01_nyc.json', 'w') as fp:
    json.dump(json_response, fp, indent=4, sort_keys=True)
```

### Store Data on Google Drive

```python
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive

#gauth = GoogleAuth()
drive = GoogleDrive(gauth)

gauth = GoogleAuth()
# Try to load saved client credentials
gauth.LoadCredentialsFile("mycreds.txt")
if gauth.credentials is None:
    # Authenticate if they're not there
    gauth.LocalWebserverAuth()
elif gauth.access_token_expired:
    # Refresh them if expired
    gauth.Refresh()
else:
    # Initialize the saved creds
    gauth.Authorize()
# Save the current credentials to a file
gauth.SaveCredentialsFile("mycreds.txt")

upload_file_list = ['2021-02-01_nyc.json']
for upload_file in upload_file_list:
    gfile = drive.CreateFile({'parents': [{'id': '19t7W5_NMU16A2XYPCbLIiiSr4XH2Zpsj'}]})
    # Read file and set it as the content of this instance.
    gfile.SetContentFile(upload_file)
    gfile.Upload() # Upload the file.
```

1. Preprocess Twitter Data
    a. Remove Links
    b. Remove Usernames
    c. Remove Hashtags
    d. Lemmatization and Tokenization

```python
def run_LDA(dataframe, num_topics):
    text_dict = Dictionary(dataframe.tokens)
    tweets_bow = [text_dict.doc2bow(tweet) for tweet in dataframe['tokens']]
    k = num_topics
    tweets_lda = LdaModel(tweets_bow,
                          num_topics=k,
                          id2word=text_dict,
                          random_state=2,
                          passes=10,
                          iterations=100)

    return tweets_lda, tweets_bow
```

# RESULTS

# Results - Topics that Influenced Electricity Use

## Traffic / Accidents

- Topic of **traffic includes accidents** caused by weather, congestion, etc.
- Traffic topic mostly showed up in **car-heavy city regions of NY state**, like Albany & Buffalo
- NYC had nothing with traffic which was interesting, possibly due to prevalence of public transportation

- Related to energy use because accidents cause increase **strain on emergency resources like police/hospitals** which leads to **increased energy use**
- Traffic congestion causes people to **work later, buildings open longer**

## Road / Regular Construction

- Also, **mostly present in larger cities**
- Could be **related to traffic/accidents** because construction work **congests roads**
- Could also **cause accidents**

- Road and regular construction naturally is **electricity intensive**: electric equipment/machinery
- **Accidents caused by construction** will affect energy use due to above reasons

## Cultural Events

- NYC had **irrelevant topics related to NYC tourism**, like pictures, monuments, etc.
- **Had to filter** the above to show relevant topics
- Most relevant topics included cultural events like **sports games, and Fashion Week (NYFW)**

- Sporting events cause more people to watch games on TV, go to sports bars; **stadiums are massive electric loads**
- NYFW was surprising, however it brings in **added tourism, media, so more electricity use**
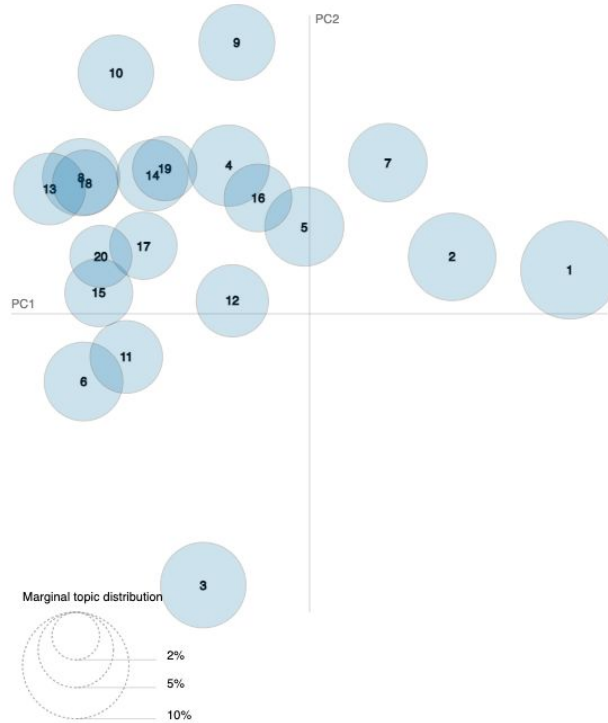
# Interesting Topics from LDA

# Results - Topics that Influenced Electricity Use

## Traffic / Accidents

- Topic of **traffic includes accidents** caused by weather, congestion, etc.
- Traffic topic mostly showed up in **car-heavy city regions of NY state**, like Albany & Buffalo
- NYC had nothing with traffic which was interesting, possibly due to prevalence of public transportation

- Related to energy use because accidents cause increase **strain on emergency resources like police/hospitals** which leads to **increased energy use**
- Traffic congestion causes people to **work later, buildings open longer**

## Road / Regular Construction

- Also, **mostly present in larger cities**
- Could be **related to traffic/accidents** because construction work **congests roads**
- Could also **cause accidents**

- Road and regular construction naturally is **electricity intensive**: electric equipment/machinery
- **Accidents caused by construction** will affect energy use due to above reasons

## Cultural Events

- NYC had **irrelevant topics related to NYC tourism**, like pictures, monuments, etc.
- **Had to filter** the above to show relevant topics
- Most relevant topics included cultural events like **sports games, and Fashion Week (NYFW)**

- Sporting events cause more people to watch games on TV, go to sports bars; **stadiums are massive electric loads**
- NYFW was surprising, however it brings in **added tourism, media, so more electricity use**

# Interesting Topics from LDA

# Results - Topics that Influenced Electricity Use
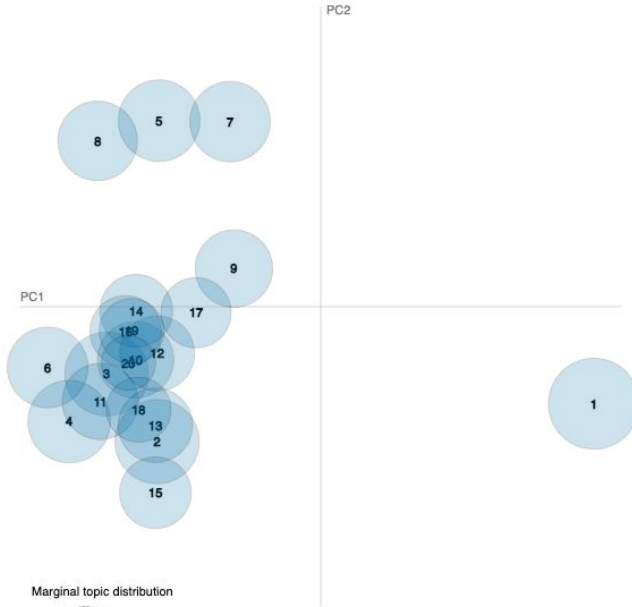
## Traffic / Accidents

- Topic of **traffic includes accidents** caused by weather, congestion, etc.
- Traffic topic mostly showed up in **car-heavy city regions of NY state**, like Albany & Buffalo
- NYC had nothing with traffic which was interesting, possibly due to prevalence of public transportation

- Related to energy use because accidents cause increase **strain on emergency resources like police/hospitals** which leads to **increased energy use**
- Traffic congestion causes people to **work later, buildings open longer**

## Road / Regular Construction

- Also, **mostly present in larger cities**
- Could be **related to traffic/accidents** because construction work **congests roads**
- Could also **cause accidents**

- Road and regular construction naturally is **electricity intensive**: electric equipment/machinery
- **Accidents caused by construction** will affect energy use due to above reasons

## Cultural Events

- NYC had **irrelevant topics related to NYC tourism**, like pictures, monuments, etc.
- **Had to filter** the above to show relevant topics
- Most relevant topics included cultural events like **sports games, and Fashion Week (NYFW)**

- Sporting events cause more people to watch games on TV, go to sports bars; **stadiums are massive electric loads**
- NYFW was surprising, however it brings in **added tourism, media, so more electricity use**

Columbia | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# Interesting Topics from LDA

# Business Value

- **Efficient prediction of anomalous spikes in energy demand is invaluable to utility companies.**
  - Appropriate resource allocation and pricing
  - Fewer losses to be incurred during unexpected spikes in demand.
- **Our analysis has found that utilizing social media to monitor certain topics can be used to predict energy demand anomalies.**
  - Of the topics we found, utility companies only monitor weather.
  - Our analysis has found that traffic incidents and large cultural events are highly correlated with energy demand spikes.
- *We have shown that through the use of big data analytics, a social media sensing pipeline that monitors twitter activity for traffic incidents or large cultural events could be utilized by utility companies to reduce losses incurred by anomalous spikes in demand*

# Demo

1. Filters:
   a. Aggregation: hourly, daily, weekly, monthly, yearly
   b. Start and end date
   c. Area to be observed
2. Visualization
   a. Electricity load energy demand line chart
   b. LDA topics visualization
      i. Bubble chart for intertopic distance
      ii. Bar chart for viewing top terms in the each topic with term frequency sorted by saliency and relevance