



Interactive Weather Data Visualizations with Plotly

[LATEST](#)

Exploring the NOAA Global Surface Summary of the Day

[EDITOR'S PICKS](#)

Will Norris

[DEEP DIVES](#)

Mar 26, 2021 10 min read

[CONTRIBUTE](#)

Making Sense of Big Data

[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

Photo by [Michael D](#) on [Unsplash](#)

The Global Surface Summary of the Day (GSOD) is a dataset from over 9,000 weather stations dating back to 1901. It is created and maintained by the National Oceanic and Atmospheric Administration (NOAA).

Administration (NOAA) and generates a daily summary of hourly surface measurements for 18 climate variables like mean dew point, mean wind speed, and min/max temperature.

This is a great data set to use for practicing visualizing spatial data since it has many stations across the globe with a relatively large time series to play around with. It also is updated daily which makes it a great candidate for practicing with data visualization.

While Python might not be the first language or library that comes to mind for visualizing spatial data, there are several libraries available that can help you get started. Python has become a very well rounded interactive plotting library over the last few years, and its geographic plotting capabilities have continued to improve. Because of its excellent documentation and popularity I will focus only on Plotly in this post, but please note that there are other libraries available.

This data set is somewhat large at 4.2 GB so loading it into memory will take a little effort; it is not a massive data set, but it is larger than the average "practice" data set. We will start by walking through how to access the data and explore visualizing the processed data on interactively. Make a cup of coffee, get comfortable, and let's dive into this larger weather data set and visualize our findings.

File Structure

There are multiple ways to obtain this data set. One way is to download version of it is available from [Kaggle](#), or you can get it updated daily like the [NOAA repositories](#) and the Kaggle dataset.

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

difficult to format in Python. This is how the NOAA repository is structured:

1929.tar.gz/

1930.tar.gz/

 03005099999.csv

 03026099999.csv

 ...

 ...

LATEST

EDITOR'S PICKS

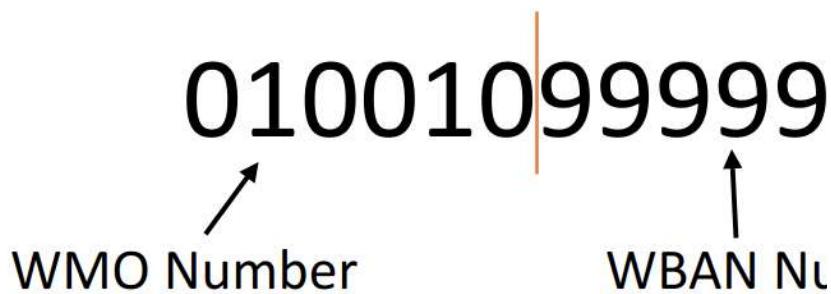
DEEP DIVES

CONTRIBUTE

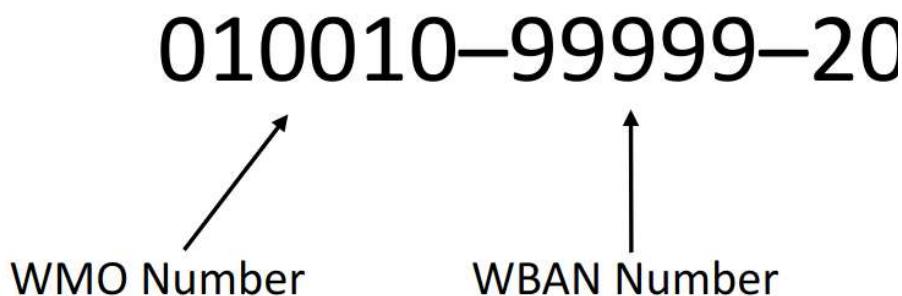
NEWSLETTER

[Sign in](#)

[Contributor Portal](#)



NOAA CSV File Structure (Created By Author)



Kaggle Gzip File Structure (Created By Author)

WMO Number

The World Meteorological Organization (WMO) is an agency of the United Nations responsible for weather, climate, and water resources. The first two digits are the country code, the next three digits are set by the National Meteorological Service (NMS). If the last number is a zero that means the WMO number is, or was in the past, official.

WBAN Number

WBAN stands for Weather-Bureau-Army-Navy. It is an identifier for digital data storage that started in 1962. It was created by the US Weather Bureau, Canada's Transportation Department, and the three major military branches to share data with each other. It has since expanded to include some German and Korean stations.

Note: The important piece to remember is that WBAN numbers create a unique station ID that is our primary key for referencing each weather station once loaded into our database.

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

Downloading Our Data

The best way to get this data is directly from NCDC. They provide a tarball for each year that contains all of the data for every station active during that year. I have included a Python module for downloading and processing all of the data. You can find it in the documentation section at the end.

```
1 import requests
2 import re
3 import os
4
5 def get_data(directory="noaa_gsod"):
6     """
7         Download all data from GSOD that is currently in the Bulk Download
8     """
```

```

9     Parameters
10    -----
11    directory: str
12        The directory you want to download the csv files to. Will create it if it
13        doesn't exist
14    """
15    if not os.path.exists(directory):
16        os.makedirs(directory)
17    base_url = "https://www.ncei.noaa.gov/data/global-summary-of-the-day/archive/"
18    rt = requests.get(base_url).text
19    years = re.findall("(?<!\\>)\\d{4}", rt)                                LATEST
20    for year in years:
21        url = base_url+str(year)+'.tar.gz'                                     EDITOR'S PICKS
22        r = requests.get(url, stream=True)
23        if r.status_code == 200:
24            with open('noaa_data/'+str(year)+'.tar.gz', 'wb') as
25                f.write(r.raw.read())

```

get_data.py hosted with ❤ by GitHub

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

To download the NOAA data I decided to make a function using requests. It starts on the home dataset and uses regex to find all of the years available on the base URL page. With each year in a list it is then the final URLs and download each tarball. Remember the code is [available on Github](#) and can be imported in Python.

Data Processing

The most important question to ask yourself when datasets is what data you want to have access to into a dataframe. There were over 11,000 stations which means ~10 years of daily data would be quite a lot of work with. Instead you may want to do some aggregation. In this article I decided to aggregate the data by month and get daily data for one day of each year. Aggregating the

will greatly reduce the size of our final dataframe and still let us find trends over a decently long time series.

It may seem weird to keep one specific day un-aggregated, and it is! The real reason we want to save a single day's data is so that we can load in the most recent day to a live dashboard with Dash in a future post. For now it will still be useful for creating similar plots to what we will explore with Dash later.

[LATEST](#)
[EDITOR'S PICKS](#)
[DEEP DIVES](#)
[CONTRIBUTE](#)
[NEWSLETTER](#)
[Sign in](#)
[Contributor Portal](#)

Step 0: Imports

Since processing this data is a little more involved than the previous practice datasets, I have included a GitHub repository that contains a module that can be imported locally into your Python project.

```
import gsodpy.gsodProcess as gsod
import gsodpy.gsodDownloader as gsodDown
import datetime
```

In order to import modules locally like this, the file must be in your project directory.

Step 1: Define Time Series

The first thing we need to do is decide how many years and which day we will store un-aggregated. We will use the following code to create our desired file list.

```
num_years = 10
target_day = datetime.datetime(2020,3,20)
years, files = get_years_files(num_years)

1 def get_years_files(num_years):
2     """
3         Create list of files and list of most recent num_years
```

```

4
5     Parameters
6     -----
7     num_years: int
8         number of years in time series
9     """
10    files = os.listdir("noaa_data")
11    files.sort()
12    files = files[-num_years:]
13    years = [int(re.findall('\d+', file)[0]) for file in files]
14
15    return years, files

```

get_years.py hosted with ❤ by GitHub

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

Once we have defined our time series, we are ready to start cracking open tarballs and processing .csv files! These are the functions I wrote to unpack our data.

```

1  def process_all_years(yearfiles, target_day):
2      """
3          Open each zip file in tar, check if station was operating in
4
5          Parameters
6          -----
7          yearfiles: str
8              List of zipped tarballs to unpack
9          target_day: datetime.datetime
10             Target day to store unaggregated data
11
12
13    df = pd.DataFrame([])
14    df_day = pd.DataFrame([])
15    for yearfile in yearfiles[:-1]: # skip most recent year for now
16        print("Processing file: {}".format(yearfile))
17        tar = tarfile.open("noaa_data/" + yearfile, "r")
18        print("Number of stations in file: {}".format(len(tar.getnames())))
19        for member in tar:
20            de = pd.read_csv(io.BytesIO(tar.extractfile(member).read()))
21            de = process_df(de)
22            df_day = df_day.append(de[(de['MONTH'] == target_day.month) & (de['YEAR'] == target_day.year)])
23            de = de.groupby(['STATION', 'WMO', 'WBAN', 'YEAR', 'MONTH'])
24            df = df.append(de)
25    tar.close()

```

```

25     df = add_meta(df)
26     df_day = add_meta(df_day)
27     df_day.reset_index(inplace = True, drop=True)
28     return df, df_day
29
30
31 def process_df(df):
32     """
33     Clean and process the raw weather station dataframes
34
35     Parameters
36     -----
37     df: pd.DataFrame
38         Raw dataframe from station csv file to clean and process
39     """
40     df['WMO'] = df['STATION'].apply(lambda x: str(x)[0:6])
41     df['WBAN'] = df['STATION'].apply(lambda x: str(x)[6:11])
42     df['MAX'] = df['MAX'].apply(lambda x: str(x).strip('*'))
43     df['MIN'] = df['MIN'].apply(lambda x: str(x).strip('*'))
44     df[['STATION', 'TEMP', 'DEWP', 'WDSP', 'MAX', 'MIN']] = df[['STATION']]
45     df['DATE'] = pd.to_datetime(df['DATE'], format='%Y%m%d', errors='coerce')
46     df['YEAR'] = pd.DatetimeIndex(df['DATE']).year
47     df['MONTH'] = pd.DatetimeIndex(df['DATE']).month
48     df['DAY'] = pd.DatetimeIndex(df['DATE']).day
49     df.drop(['TEMP_ATTRIBUTES', 'DEWP_ATTRIBUTES', 'MAX_ATTRIBUTE',
50              'PRCP_ATTRIBUTES', 'SLP_ATTRIBUTES', 'MIN_ATTRIBUTES',
51              'SLP', 'STP', 'STP_ATTRIBUTES'], axis=1, inplace=True)
52     return df
53
54
55 def add_meta(df):
56     """
57     Add metadata column to dataframe for plotly text boxes
58
59     Parameters
60     -----
61     df: pd.DataFrame
62         The final weather station dataframe to add metadata to
63     """
64     df['META'] = df['NAME']
65     df['ELEV_LABEL'] = df['ELEVATION'].apply(lambda x: 'Elevation: ' + str(x))
66     df['META'] = df[['META', 'ELEV_LABEL']].apply(lambda x: x.str.cat(sep=', '))
67     df['addmeta'] = df['TEMP'].apply(lambda x: "Temp: {} C".format(x))
68     df['META'] = df[['META', 'addmeta']].apply(lambda x: x.str.cat(sep=', '))
69     df = df.drop(['NAME', 'ELEV_LABEL', 'addmeta'], axis=1)
70     return df

```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

To save you some time, here is the basic psuedo-code for this workflow:

- Loop through each tarball (years)
 - Loop through each csv file stored in tarball (stations)
 - Open .csv file for current station/year with LATEST
 - Call process_df() to clean the dataframe
 - Append all data for this station on target_col
 - Aggregate data by month
 - Append monthly aggregate to df
 - Call add_meta() to create metadata column on df and
 - Return completed df and df_day

[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

Overall the process is not complex, and I'm sure more concise given more time. Remember that the number of stations in each year it takes quite a bit of time to process this data.

We can go ahead and call our processing function on the dataframes:

```
df, df_day = gsod.process_all_years(files, target_day)
```

Note: If you notice that your machine is struggling to process all the data try starting with a short time series. If you want a longer time series you can also select a subset of stations in the tarball to reduce the number of stations you need to process.

Step 3: Save Processed Data

I highly recommend pickling your finished dataframes so that you can load them in much faster later.

```
import pandas as pd

# Save DataFrames to pickle in current directory
df.to_pickle("df_monthly.pkl")
df_day.to_pickle("df_daily.pkl")
```

LATEST

```
# Read them back in later:
```

EDITOR'S PICKS

```
df = pd.read_pickle("df_monthly.pkl")
df_day = pd.read_pickle("df_daily.pkl")
```

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

Visualizations

Now that we have some meaningful data loaded get to the fun part and start making visualizatio

If you use the free version of Plotly you are limit 512kb. This isn't a problem if you are making yo However, for me to share interactive plots here Plotly's Chart Studio to upload them to their clc am limited to roughly 300 weather stations that randomly with this function:

```
1 def sample_stations(df_month, num_samples=300):
2     stationID = df_month["STATION"].unique().tolist()
3     station_subset = random.sample(stationID, num_samples)
4     return df_month[df_month['STATION'].isin(station_subset)]
```

sample_stations.py hosted with ❤ by GitHub

Basic Maps

To start off, lets take a look at one month in our time series to see where some of our stations are on the map. This is about as basic of a map we can create with Plotly and this data; considering how little configuration it took, I think it looks pretty good.

Plotly's main data structure is the `Figure`, which is an instance of the `plotly.graph_objects.Figure` class. The figure is rendered via Plotly.js using JavaScript, which is why we are all LATEST looking maps in Python.

EDITOR'S PICKS

```

1  month = 2
2  year = 2016
3  df_month = sample_stations(df, 300)
4  tempsdf = df_month[(df_month['MONTH']==month) & (df_month['YEAR'] ]
5
6  fig = go.Figure(data=go.Scattergeo(
7      text = tempsdf['META'],
8      lat = tempsdf['LATITUDE'],
9      lon = tempsdf['LONGITUDE'],
10     mode = 'markers',
11     marker_color = tempsdf['TEMP'],
12   ))
13
14 fig.update_layout(
15     title = 'Temperature on {},{}'.format(month,year),
16     geo_scope='world',
17     geo = dict(
18         showcountries = False,
19         showcoastlines=False,
20         resolution = 110)
21   )
22 fig.show()

```

basicmap.py hosted with ❤ by GitHub

DEEP DIVES

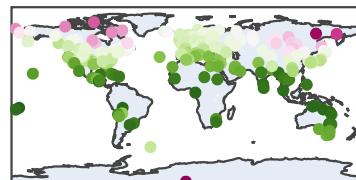
CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

Temperature on 2, 2016

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

I love plotting spatial data with Plotly because of the responsive interface that can be embedded almost anywhere. Quality base maps have always been difficult for Python packages, but Plotly's are great.

While this plot may be fine, there are absolutely changes I would make to it. The coastlines are currently too thin. The color scale could use adjusting. Overall it could be improved so that our data points stand out more. Luckily all of these changes can be accomplished with Plotly!

Styled Plots: Extreme Temperatures on Given Day

One reason I chose to keep a single day of the month instead of taking monthly aggregates was so that we could quickly find the warmest and coldest places on a given day.

First we need to find our extreme temperatures

```
1 def create_extremes_df(df, day, num_extremes):  
2     """
```

```

3     Create extremes dataframe to plot coldest and warmest places on a given day
4     """
5
6     extremes = pd.DataFrame([])
7
8     extremes = extremes.append(df[df['DATE']==day].sort_values(by="TEMP", ascending=False).head(1))
9     extremes = extremes.append(df[df['DATE']==day].sort_values(by="TEMP", ascending=False).tail(1))
10
11    return extremes

```

extremes.py hosted with ❤ by GitHub

[view raw](#)

This plot uses the same structure as our first or more time customizing the marker and figure pr notice that I am able to pass a dictionary to cer marker to customize them further. In most cases are objects that are being initialized by the infor dictionary. For example, the `marker` parameter in `plotly.graph_objects.[scattergeo.Marker]` (https://plotly.github.io/plotly.py-docs/generated/plotly.graph_objects.scattergeo.html#plotly.graph_objects.scattergeo.Marker) object. While this increases comp from the `scattergeo.Marker` documentation how gives us over the markers on our plot.

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

```

1  scl = [0,"rgb(150,0,90)"],[0.125,"rgb(0, 0, 200)"],[0.25,"rgb(0,
2) [0.375,"rgb(0, 152, 255)"],[0.5,"rgb(44, 200, 150)"],\
3 [0.75,"rgb(255, 234, 0)"],[0.875,"rgb(255, 111, 0)"],[1,"rgb(255,
4
5 fig = go.Figure(data=go.Scattergeo(
6     lat = extremes['LAT'],
7     lon = extremes['LON'],
8     text = extremes['META'],
9     marker = dict(
10         color = (extremes['TEMP']-32)*5/9,
11         colorscale = scl,
12         cmin = -50,
13         cmax = 50,
14         size = 5,
15         colorbar = dict(
16             thickness = 10,
17             titleside = "right",
18             outlinecolor = "rgba(68, 68, 68, 0)",
19             tickvals = [-50,-30,-15,0,15,30,50],

```

```

20         ticks = "outside",
21         ticklen = 3,
22         ticksuffix = " C",
23         showticksuffix = "all"
24     )
25   )
26 ))
27
28 fig.update_layout(
29     geo = dict(
30         scope = 'world',                                     LATEST
31         showland = True,
32         landcolor = "rgb(212, 212, 212)",                EDITOR'S PICKS
33         showlakes = True,
34         lakecolor = "rgb(255, 255, 255)",                DEEP DIVES
35         showsubunits = False,
36         showcountries = False,                            CONTRIBUTE
37         showcoastlines=False,
38         resolution = 110
39     ),
40     title='{} Coldest and Hotest Temperatures on {}'.format(num_ε)
41
42 fig.show()

```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

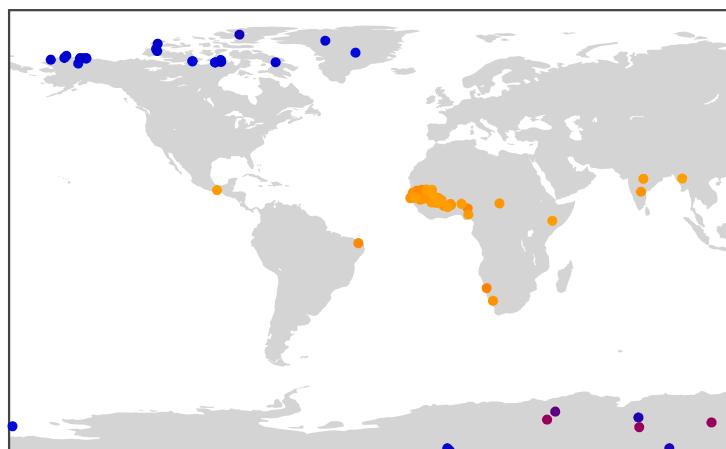
NEWSLETTER

Sign in

Contributor Portal

scattergeo.py hosted with ❤ by GitHub

40 Coldest and Hostest Temperatures on (3,20,202)



Mapbox Integration

I have always been a sucker for Mapbox. It has great base maps and is super responsive. Plotly has limited Mapbox integration that can enhance your geographical plots. Especially when zooming in you will notice a pleasing amount of detail in the base map compared to the base Plotly version. Mapbox has many effects on the way Plotly maps are rendered, so I will be covering [Mapbox Integration in Plotly](#) in an entire post.

[EDITOR'S PICKS](#)

```

1 mapbox_access_token = open(".mapbox_token").read()
2
3 fig = go.Figure(go.Scattermapbox(
4     lat=extremes['LATITUDE'],
5     lon=extremes['LONGITUDE'],
6     mode='markers',
7     marker = dict(
8         color = (extremes['TEMP']-32)*5/9,
9         colorscale = scl,
10        cmin = -50,
11        cmax = 50,
12        size = 5,
13        colorbar = dict(
14            thickness = 10,
15            titleside = "right",
16            outlinecolor = "rgba(68, 68, 68, 0)",
17            tickvals = [-50,-30,-15,0,15,30,50],
18            ticks = "outside",
19            ticklen = 3,
20            ticksuffix = " °C",
21            showticksuffix = "all"
22        )
23    ),
24    text=extremes['META'],
25 ))
26
27 fig.update_layout(
28     autosize=True,
29     hovermode='closest',
30     mapbox=dict(
31         accesstoken=mapbox_access_token,
32         bearing=0,
33         pitch=0,
```

[DEEP DIVES](#)

[CONTRIBUTE](#)

[NEWSLETTER](#)

[Sign in](#)

[Contributor Portal](#)

```

34         zoom=0
35     ),
36 )
37
38 fig.show()

```

scattermap.py hosted with ❤ by GitHub

[view raw](#)

40 Coldest and Hottest Temperatures on (3,20,202

LATEST



EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

[Sign in](#)[Contributor Portal](#)

Adding Sliders: Average Temperatures by Month

One feature that I really like about Plotly is how adding a slider to view how data changes over time lets take a look at the average monthly tempera

You will notice that I am using a slightly different this plot. This is technically the "old" method of But for more complex plots like this it is nice to the data separately.

```

1  data = [ dict(
2      visible = False,
3      name = '',
4      type = 'scattergeo',
5      text = df_month[(df_month['MONTH'] == month) & (df_month['YEAR']

```

```

6      lat = df_month[(df_month['MONTH']==month) & (df_month['YEAR']==year)][['LATITUDE']],
7      lon = df_month[(df_month['MONTH']==month) & (df_month['YEAR']==year)][['LONGITUDE']],
8      marker = dict(
9          color = (df_month[(df_month['MONTH']==month) & (df_month['YEAR']==year)][['TEMP']] - 32) * 5 / 9,
10         colorscale = scl,
11         cmin = -50,
12         cmax = 50,
13         opacity = 0.5,
14         size = 5,
15         colorbar = dict(
16             thickness = 10,                                              LATEST
17             titleside = "right",
18             outlinecolor = "rgba(68, 68, 68, 0)",
19             tickvals = [-50, -30, -15, 0, 15, 30, 50],
20             ticks = "outside",
21             ticklen = 3,
22             ticksuffix = " °C",
23             showticksuffix = "all"
24         )
25     )
26 ) for year in [years[0]] for month in range(1,13)]                                              EDITOR'S PICKS
27 data[-1]['visible'] = True
28
29 steps = []
30 for i in range(len(data)):
31     step = dict(
32         method = 'restyle',
33         args = ['visible', [False] * len(data)],
34         label = [str(month)+". "+str(year) for year in [years[-1]]]
35     )
36     step['args'][1][i] = True
37     steps.append(step)
38
39 sliders = [dict(
40     active = len(steps)-1,
41     currentvalue = {"prefix": "Month and year: "},
42     steps = steps
43 )]
44
45 layout = dict(
46     sliders=sliders,
47     geo = dict(
48         scope = 'world',
49         showland = True,
50         landcolor = "rgb(212, 212, 212)",
51         showlakes = True,
52         lakecolor = "rgb(255, 255, 255)",
```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

```

53     showsubunits = False,
54     showcountries = False,
55     showcoastlines=False,
56     resolution = 110
57   ),
58 )
59 fig = dict( data=data, layout=layout )

```

avg_month_temp.py hosted with ❤ by GitHub

[view raw](#)

LATEST

EDITOR'S PICKS

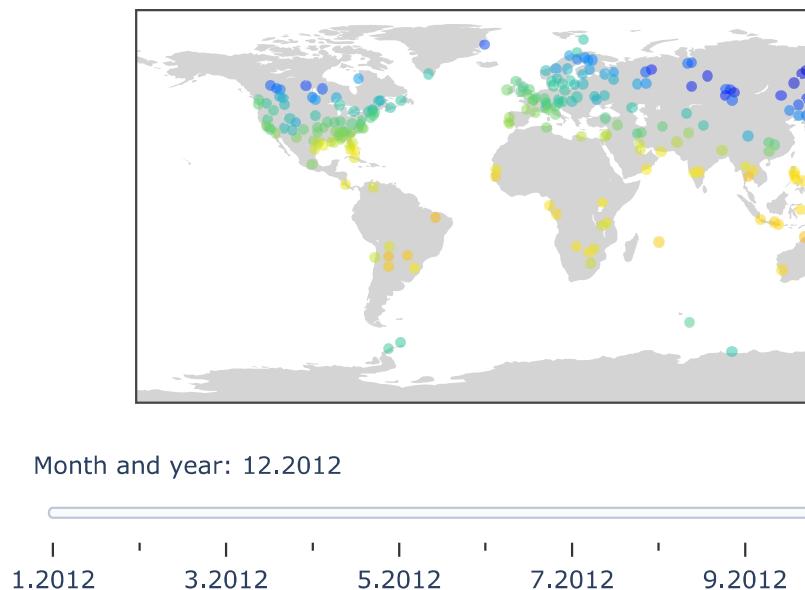
DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal



Creating this plot looks more complex than the the bulk of the code is simply restructuring data groupings for Plotly to easily send to its JavaScript have ever worked with json data, you will probal of dictionaries can easily be converted to json.

Wrapping Up

Spatial data is never the easiest to work with in multidimensional, often quite large, and normal are not well suited to plotting geographical coor

removes a lot of the headache by providing easy to use functionality with quality base maps that I generally had to use javascript to access a couple years ago. When it comes to spatial plotting in Python, Plotly is generally my first choice.

There is so much more to learn with plotting spatial data with Plotly so don't be afraid to sift through the documentation and examples online. The same developers also make Dash, w LATEST for building live dashboards.

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

Documentation

- [Github Repo for Processing GSOD Data](#)
- [GSOD Data Repository \(NOAA\)](#)
- [Kaggle GSOD Data](#)
- [Plotly](#)
- [Mapbox](#)

• • •

WRITTEN BY

Will Norris

See all from Will Norris

Topics:

Editors Pick

Geospatial

Making Sense Of Big D

Python

Share this article:

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal



ARTIFICIAL INTELLIGENCE

What Do Large Language Models “Understand”?

A deep dive on the meaning of understanding and how it applies to LLMs

Tarik Dzekman

August 21, 2024 31 min read

DATA SCIENCE

Must-Know Bivariate Normal Explained

Derivation and intuition behind this powerful concept

Luigi Battistoni

August 14, 2024



DATA SCIENCE



DEEP LEARNING

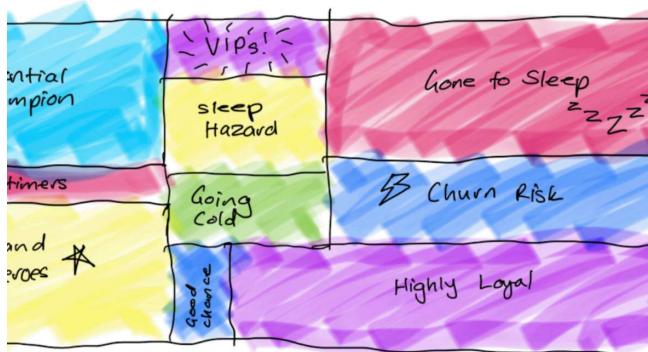
Deep Dive in by Hand



Optimizing Marketing Campaigns with Budgeted Multi-Armed Bandits

With demos, our new solution, and a video

Vadim Arzamasov



ANALYTICS

Methods for Modelling Customer Lifetime Value: The Good Stuff and the Gotchas

Part three of a comprehensive, practical guide to CLV techniques and real-world use-cases

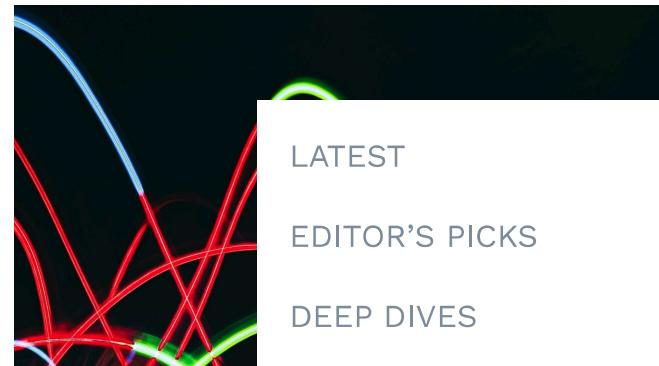
Katherine Munro

November 17, 2023 12 min read

Explore the wisdom of LSTM leading into xLSTMs - a probable competition to the present-day LLMs

Srijanie Dey, PhD

July 9, 2024 13 min read



DATA SCIENCE

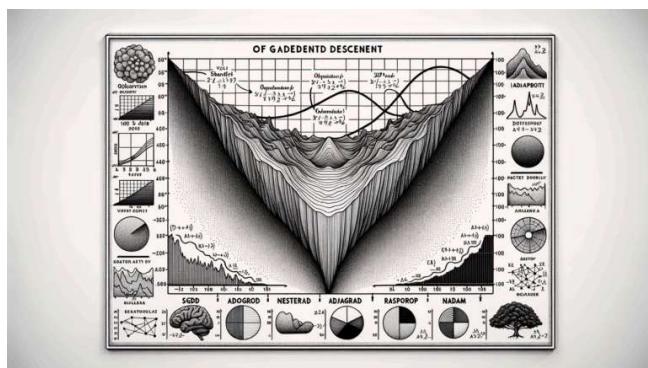
Squashing the Gradient into Penalized Regression for

How to build a penalized regression model
Álvaro Méndez Císcar
August 16, 2024

NEWSLETTER

Sign in

Contributor Portal



DATA SCIENCE

The Math Behind Keras 3 Optimizers: Deep Understanding and Application

This is a bit different from what the books say.

Peng Qian

August 17, 2024 9 min read



LATEST

Your home for data science and AI. The world's leading publication for data analytics, data engineering, machine learning, and artificial intelligence professionals.

[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[ABOUT · ADVERTISE · PRIVACY POLICY · TERMS](#)[Contributor Portal](#)[COOKIES SETTINGS](#)