



BIG DATANIGHTS - TIME SERIES

2.11- 7.12

Wednesdays
18:30 - 20:30



MAYA BERCOVITCH
Director of Data Science
Anaplan



LIRON MORGENSTERN
Data Scientist
Anaplan



HILA PAZ HERSZFANG
Data Science Team Leader
Zscaler



ALICE FRIDBERG
Data Science Team Leader
SKAI



NAOMI FRIDMAN
Deep Learning Consultant
Freelance



ZIV FREUND
Head of AI Research Group
Elbit Systems



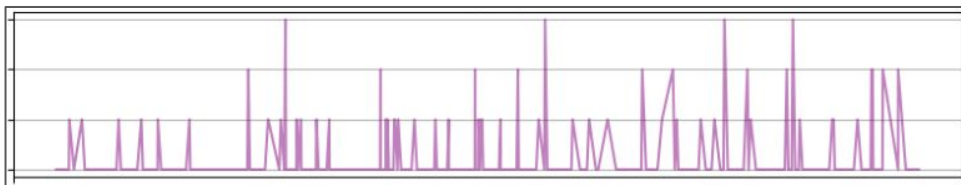
RAZ TAMIR
Data Science Team Leader
Owlytics Healthcare



Amnon Wahle
Program Leader
Data Scientist
Wiliot

SUIT Method

SAMPLE

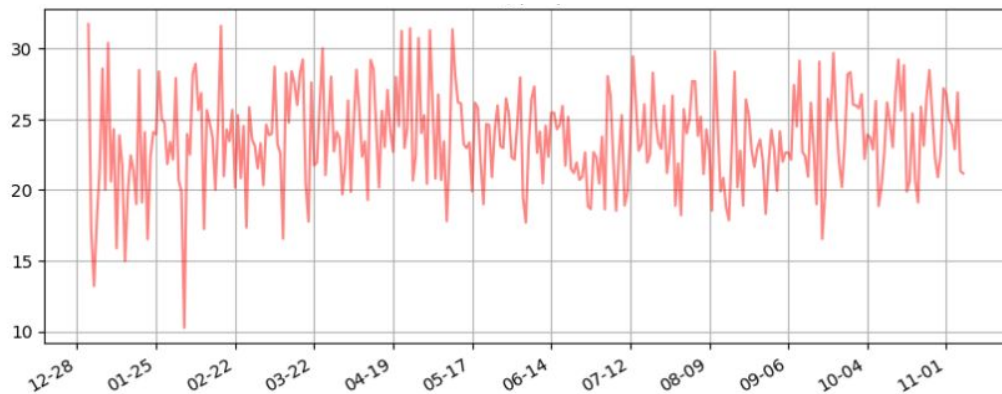


UNIVARIATE



ISOLATE

TRANSFORM



OUR DATA



Search accounts and videos



+ Upload

Log in



pazpazthecoder

Hila Paz Herszfang

Follow

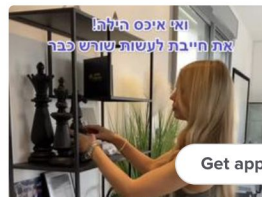
247 Following 14.1K Followers 257.8K Likes

!לא להפריע

אני כותבת קוד 🇮🇱

Videos

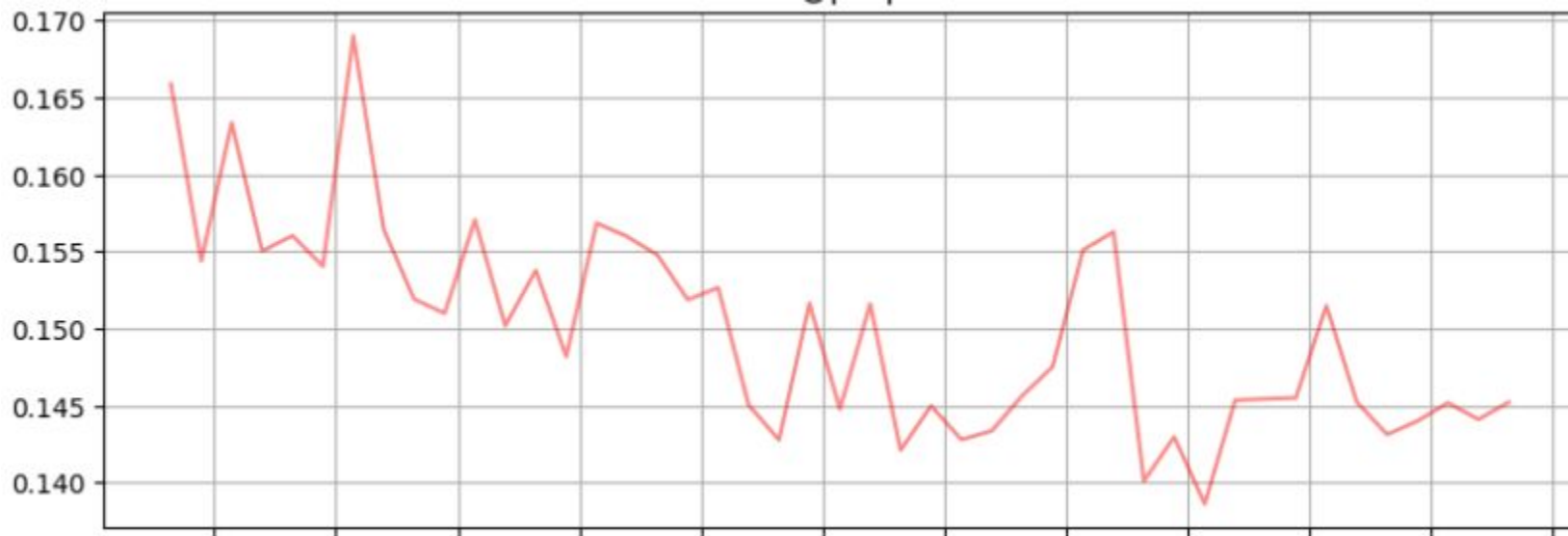
🔒 Liked



Get app

OUR GOAL

Understand why the like ratio is decreasing



DATA TYPES

<code>datetime</code>	<code>followers</code>
<code>2022-01-01 00:00:00</code>	<code>112.0</code>
<code>2022-01-01 00:30:00</code>	<code>NaN</code>
<code>2022-01-01 01:00:00</code>	<code>102.0</code>
<code>2022-01-01 01:30:00</code>	<code>106.0</code>
<code>2022-01-01 02:00:00</code>	<code>105.0</code>
<code>2022-01-01 02:30:00</code>	<code>105.0</code>
<code>2022-01-01 03:00:00</code>	<code>111.0</code>
<code>2022-01-01 03:30:00</code>	<code>108.0</code>
<code>2022-01-01 04:00:00</code>	<code>108.0</code>
<code>2022-01-01 04:30:00</code>	<code>101.0</code>

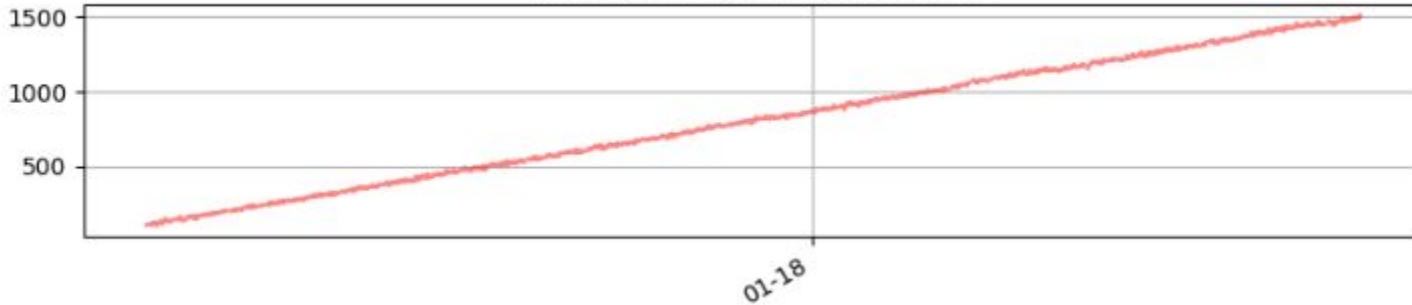
<code>date</code>	<code>views</code>	<code>likes</code>
<code>2022-01-01</code>	<code>1431</code>	<code>157</code>
<code>2022-01-02</code>	<code>1241</code>	<code>124</code>
<code>2022-01-03</code>	<code>1547</code>	<code>185</code>
<code>2022-01-04</code>	<code>1435</code>	<code>200</code>
<code>2022-01-05</code>	<code>1407</code>	<code>196</code>
<code>2022-01-06</code>	<code>1304</code>	<code>104</code>
<code>2022-01-07</code>	<code>1052</code>	<code>115</code>
<code>2022-01-08</code>	<code>1382</code>	<code>193</code>
<code>2022-01-09</code>	<code>1559</code>	<code>265</code>
<code>2022-01-10</code>	<code>1615</code>	<code>209</code>

DATA TYPES

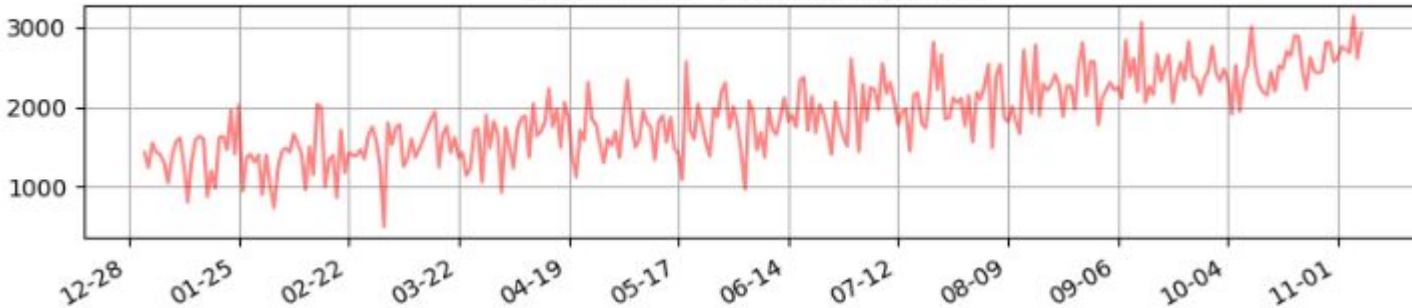
datetime	views	likes	city	age_group
2022-01-01 00:00:00	1	0	rishon	18-21
2022-01-01 00:00:00	1	0	kadima	13-15
2022-01-01 00:00:00	1	0	ramat-gan	18-21
2022-01-01 00:00:00	1	0	tel-aviv	Prefer not to Say
2022-01-01 00:00:27	1	0	tel-aviv	21-25
2022-01-01 00:08:19	1	0	tel-aviv	21-25
2022-01-01 00:10:12	1	0	tel-aviv	18-21
2022-01-01 00:10:22	1	0	tel-aviv	18-21
2022-01-01 00:12:20	1	0	tel-aviv	30+
2022-01-01 00:21:02	1	0	ramat-gan	21-25

DATA TYPES

Tiktok Data - @pazpazTheCoder

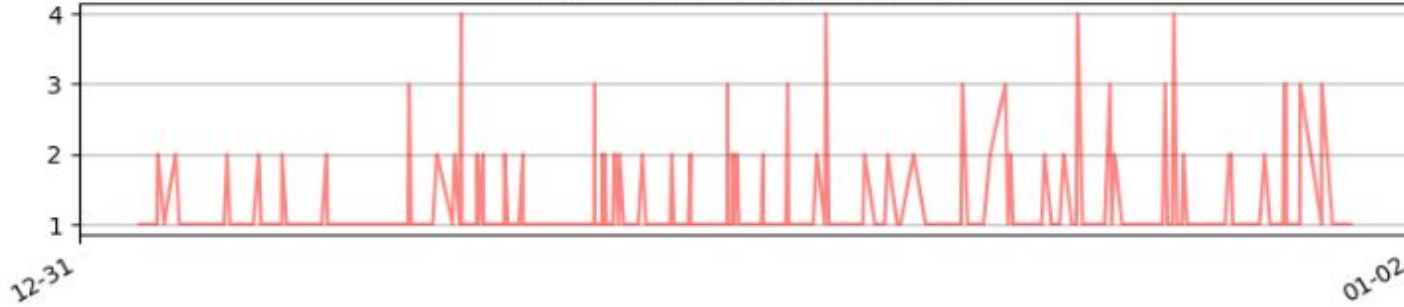


Tiktok Data - @pazpazTheCoder

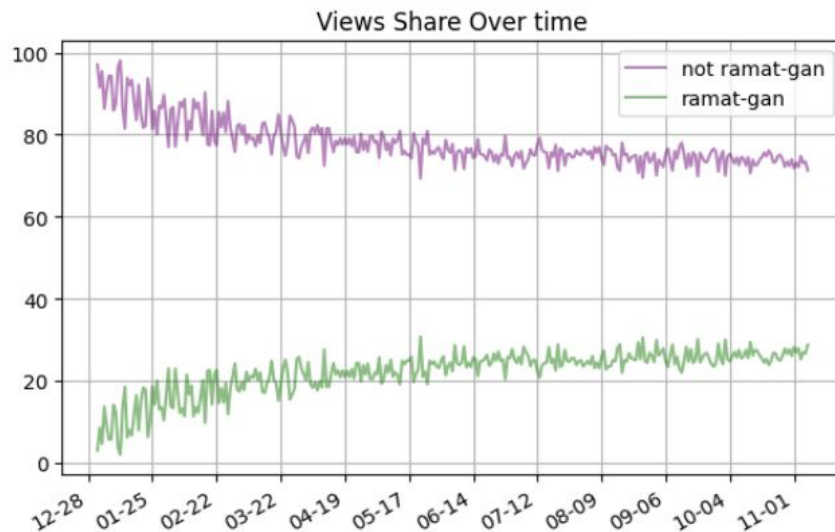
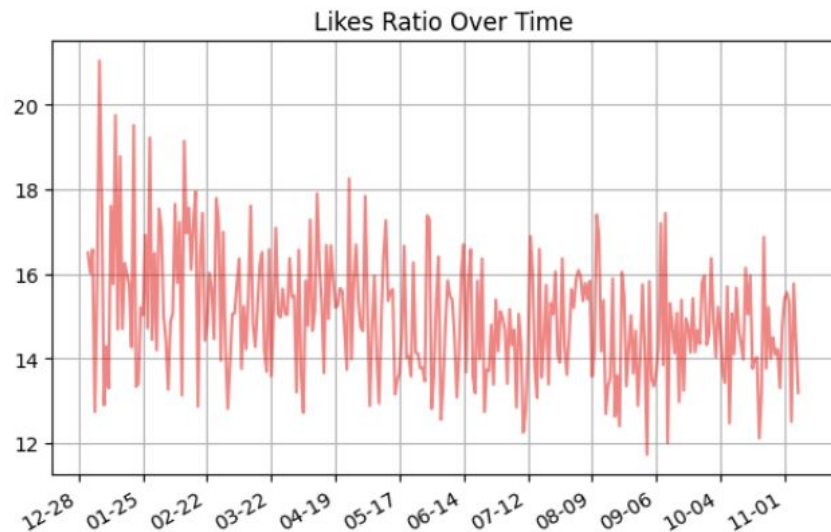


DATA TYPES

Tiktok Data - @pazpazTheCoder

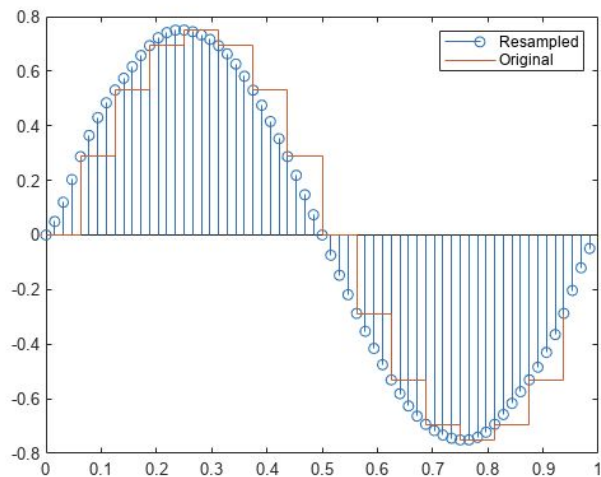


APPLY SUIT



SAMPLING

2.1 DownSampling Cumulative Data



UNIVARIATE

```
1 def get_tlv_likes(df: pd.DataFrame) -> int:
2     tlv = df[df['city'] == 'tel-aviv']
3     likes = sum(tlv['likes'])
4     return likes
5
6 events.resample('1D').apply(get_tlv_likes)
```

▼

datetime	
2022-01-01	48
2022-01-02	55
2022-01-03	94
2022-01-04	44
2022-01-05	69

ISOLATION

```
1 groupers = [pd.Grouper(freq='1D'),'city']
2
3 by_date_and_city = events.groupby(groupers).likes.sum()
4 by_date = events.groupby(pd.Grouper(freq='1D')).likes.sum()
5
6 (by_date_and_city / by_date).unstack()
```

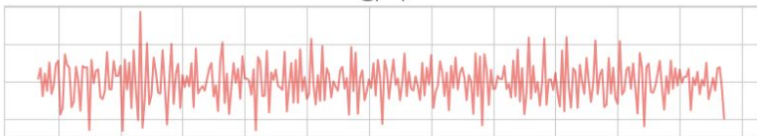
▼

city datetime	kadima	ramat-gan	rishon	tel-aviv
2022-01-01	0.061538	0.000000	0.200000	0.738462
2022-01-02	0.120482	0.012048	0.204819	0.662651
2022-01-03	0.025210	0.008403	0.176471	0.789916

311 rows x 4 columns [Open in new tab](#)

TRANSFORMATION

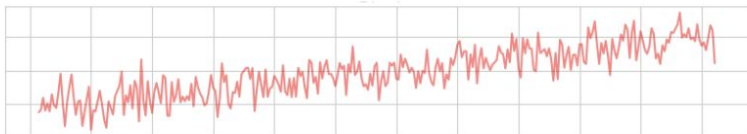
Stationary Data



The unconditional joint probability distribution does not change when shifted in time

E.g. mean and variance are constant

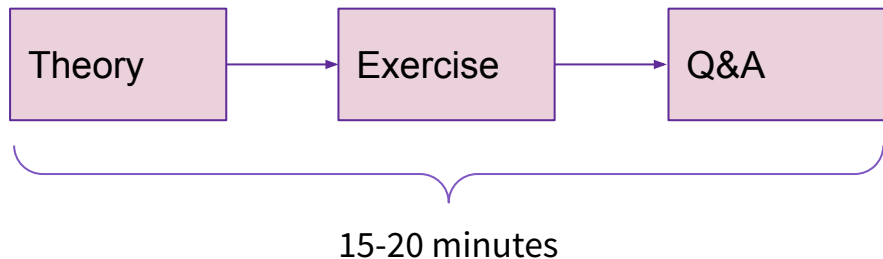
Non-Stationary Data



META STRUCTURE

warmup() *# We are here*

for i **in** range[1,4]:



closure()

CHEAT SHEET

TIME SERIES PREPROCESSING WITH THE SUIT METHOD - CHEAT SHEET

@pazpazthecoder, github.com/suit

November 2022

SAMPLE

Just-right frequency

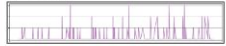
1. Delta



2. Cumulative



3. Events



PHIDS

Plot, Head, Info, Describe,
Set Index

As Freq

Only for cumulative data
`df[col].asfreq(freq='1D')`

Resample

Events -> Sampling

`df.resample('1T').agg(...)`

aligner

reducer

UNIVARIATE

Vanilla Feature Engineering

```
df.resample('1D').aggregate(...)  
df.resample('1D').apply(lambda x: ...)  
  
def get_lambda(x):  
    return ...  
data.resample('1D').apply(get_lambda)
```

Datetime Features

```
resampler.index.isocalendar()  
resampler.index.month
```

ISOLATE

Combine resampling with other features!

```
groupers = [pd.Grouper(freq='1D'), col]  
by_date_and_col = df.groupby(groupers)...
```

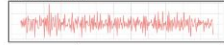
aligner reducer

We can manipulate different samplers together!

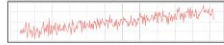
```
by_date = df.groupby(pd.Grouper(freq='1D'))...  
(by_date_and_col / by_date).unstack()
```

TRANSFORM

Augmented Dickey-Fuller test
stationary



non-stationary



```
from statsmodels.tsa.  
stattools import adfuller
```

```
adf = adfuller(data['x'])  
if adf[1] > 0.05:  
    print("data is not  
         stationary")
```

Transformation Techniques


"Features of Features"

lags / diffs
second derivative
combined transformation



GIT REPO

<https://github.com/hip023/suit>

A terminal window with a light pink title bar. The title bar contains three colored window control buttons (red, yellow, green) on the left, a folder icon followed by the text "DataspellProjects — hila@Hila — -zsh — 64x16" in the center, and a plus sign button on the right. Below the title bar, the current directory is shown as "..spellProjects". The main area of the terminal displays a command prompt with a green arrow icon, followed by the text "DataspellProjects git clone git@github.com:hip023/suit.git" in a monospaced font. A grey cursor block is positioned at the end of the command.

```
DataspellProjects — hila@Hila — -zsh — 64x16  
..spellProjects  
→ DataspellProjects git clone git@github.com:hip023/suit.git
```


GIT REPO

[hip023/suit](#) Private

[Code](#) [Issues](#) [Pull requests](#) [Actions](#) [Projects](#) [Security](#) [Insights](#) [Settings](#)

main 1 branch 0 tags

[Go to file](#) [Add file](#) [Code](#)

hilahersz feat: add exercises 79eeafe 13 minutes ago 🕒 8 commits

assets	feat: add assets	21 minutes ago
data	feat: lecture code	19 minutes ago
exercises	feat: add exercises	13 minutes ago
lecture_code	feat: lecture code	19 minutes ago
lecture_pdf	feat: lectures_pdf	14 minutes ago
solutions	feat: add solutions	17 minutes ago
utilities	feat: add utility code	18 minutes ago
.gitignore	feat: hello world	28 minutes ago
Pipfile	feat: hello world	28 minutes ago
Pipfile.lock	feat: hello world	28 minutes ago
README.md	fix: readme	15 minutes ago

README.md

KAGGLE DATA

The screenshot shows the Kaggle Data page for user Hila Paz. The browser address bar displays 'kaggle.com'. The navigation bar includes links for Home, Competitions, Datasets (6), Code (2), Discussion, and Followers (1), along with an 'Edit Public Profile' button. The 'Your Work' section is active, showing three datasets: 'tiktok_followers', 'tiktok_daily_views', and 'tiktok-events'. Each dataset entry includes a thumbnail, the dataset name, the user 'Hila Paz', the update time 'Updated 6 hours ago', and the file details 'Usability 1.2 · 1 File (CSV)'. The file sizes are 62 kB, 2 kB, and 3 MB respectively. Each entry also features a share icon and a counter showing '0' shares.

Browser address bar: kaggle.com

Search bar: Search

Navigation: Home Competitions **Datasets (6)** Code (2) Discussion Followers (1) ... [Edit Public Profile](#)

Section: Your Work Shared With You Bookmarks Updated ▾

tiktok_followers
Hila Paz · Updated 6 hours ago
Usability 1.2 · 1 File (CSV) · 62 kB

tiktok_daily_views
Hila Paz · Updated 6 hours ago
Usability 1.2 · 1 File (CSV) · 2 kB

tiktok-events
Hila Paz · Updated 6 hours ago
Usability 1.2 · 1 File (CSV) · 3 MB

KAGGLE DATA

Make sure that the name of the downloaded file is kaggle.json (if it's not, rename it :))

5. Upload the kaggle.json file to the working directory of the exercise folder - and that's it! you're good to go.



udacity_learn_view
exercises
kaggle.json

```
1 # Run the following cell to collect the first exercise data from kaggle
2 ! mkdir ~/.kaggle
3 ! mv kaggle.json ~/.kaggle/
4 ! chmod 600 ~/.kaggle/kaggle.json
```

```
1 # Run the following commands to import the datasets from Kaggle
```

KAGGLE DATA

Make sure that the name of the downloaded file is kaggle.json (if it's not rename it :))

5. Upload the kaggle.json file to the working directory of the exercise folder - and that's it! you're good to go.



udacity_learn_view
exercises
kaggle.json

```
1 # Run the following cell to collect the first exercise data from kaggle
2 ! mkdir ~/.kaggle
3 ! mv kaggle.json ~/.kaggle/
4 ! chmod 600 ~/.kaggle/kaggle.json
```

```
1 # Run the following commands to import the datasets from Kaggle
```

YOUTUBE CODE

VIDEO ON

MIC Closed

WARM UP

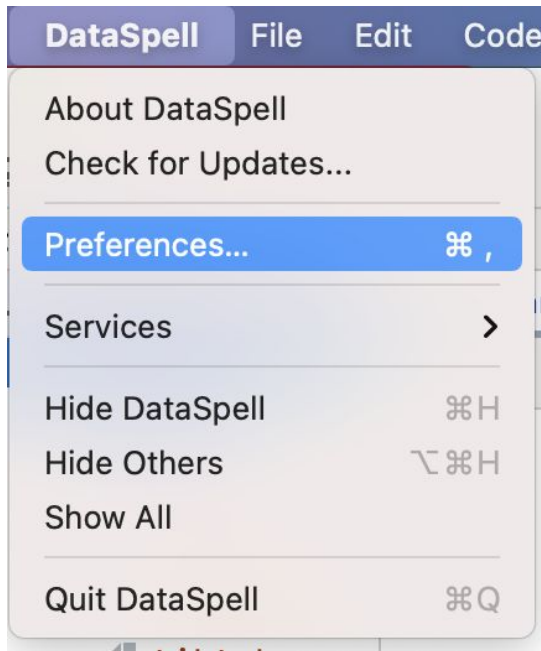


A terminal window titled "DataspellProjects — hila@Hila — -zsh — 64x16". The window shows the current directory as "..spellProjects". A green arrow points to the command "DataspellProjects git clone git@github.com:hip023/suit.git".

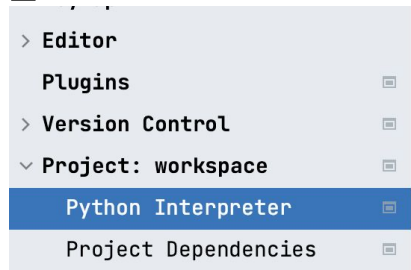
```
DataspellProjects — hila@Hila — -zsh — 64x16  
..spellProjects  
→ DataspellProjects git clone git@github.com:hip023/suit.git
```

WARM UP

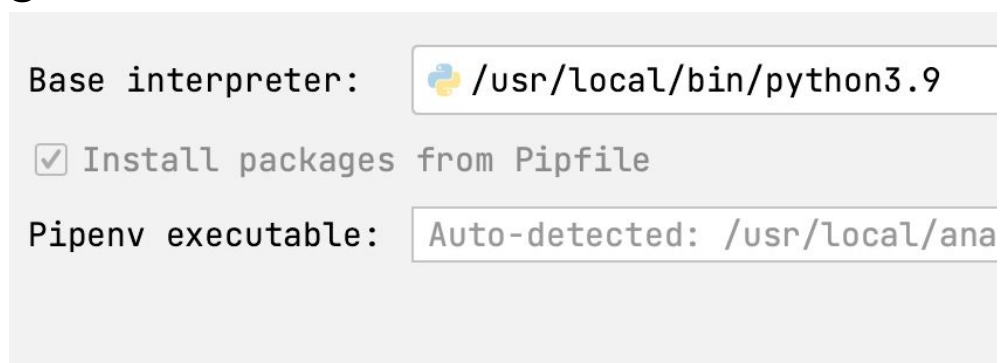
1



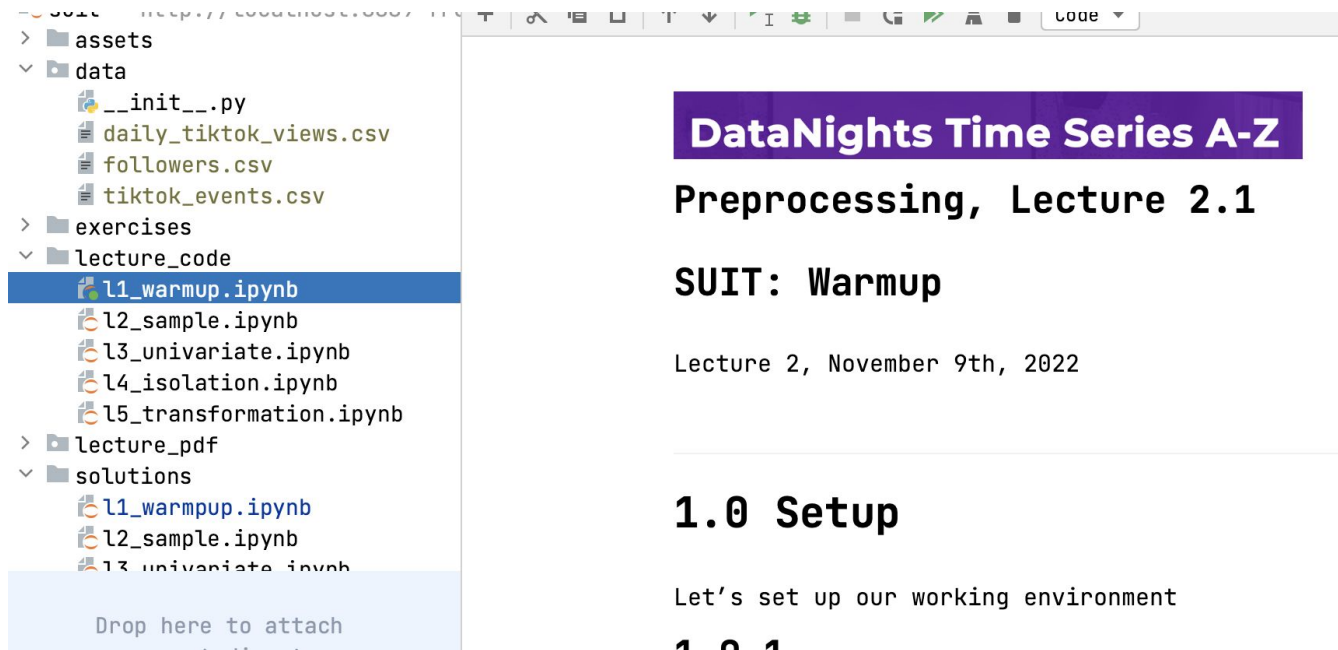
2



3



WARM UP



The screenshot shows a JupyterLab environment. On the left is a file explorer with the following structure:

- assets
- data
 - __init__.py
 - daily_tiktok_views.csv
 - followers.csv
 - tiktok_events.csv
- exercises
- lecture_code
 - l1_warmup.ipynb** (selected)
 - l2_sample.ipynb
 - l3_univariate.ipynb
 - l4_isolation.ipynb
 - l5_transformation.ipynb
- lecture_pdf
- solutions
 - l1_warmup.ipynb
 - l2_sample.ipynb
 - l3_univariate.ipynb

At the bottom of the file explorer is a light blue box with the text "Drop here to attach".

The right pane shows the content of the selected file, **l1_warmup.ipynb**. It has a purple header bar with the text "DataNights Time Series A-Z". Below this is the title "Preprocessing, Lecture 2.1" and the subtitle "SUIT: Warmup". The text "Lecture 2, November 9th, 2022" is displayed below the subtitle. A horizontal line separates the header from the main content area, which begins with the section "1.0 Setup". The first line of the setup section is "Let's set up our working environment".

WARM UP

`l1_warmup.ipynb`

1. Move `kaggle.json`
2. Download 3 datasets from kaggle

EXERCISE

`exercises/l1_warmup.ipynb`

DataNights Time Series A-Z

Preprocessing, Exercise 1 (With Hila)

Warmup

Lecture 2, November 9th, 2022

1.0 Setup

Make sure that you import datasets as explained in `lecture_code/l1_warmup`.

1.1 Let's import some packages

```
1 import pandas as pd
```

EXERCISE

exercises/l1_warmup.ipynb

Time for PHIDS!

1. plot

```
dataframe.column.plot()
```

2. head

```
dataframe.head()
```

3. info

```
dataframe.info()
```

4. describe

```
dataframe.describe()
```

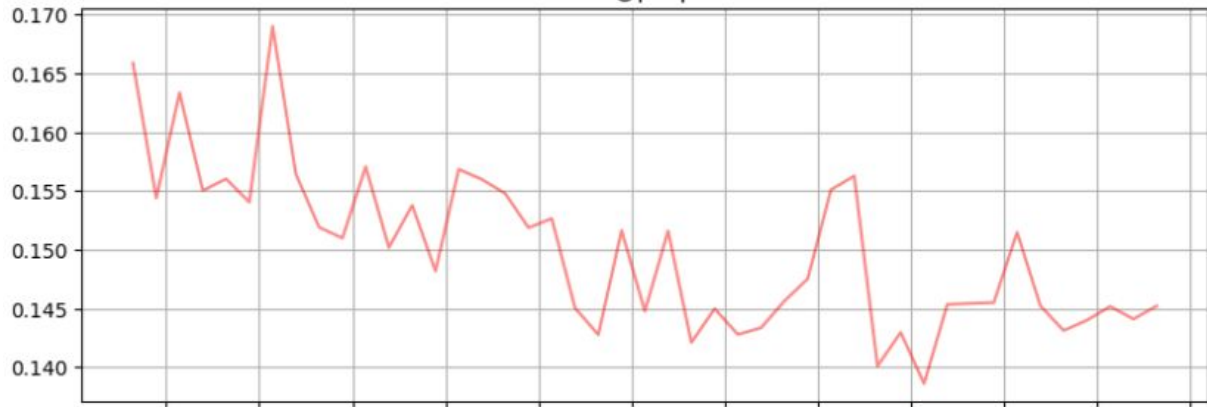
5. set index

```
dataframe.set_index("column")
```

L1 SUMMARY

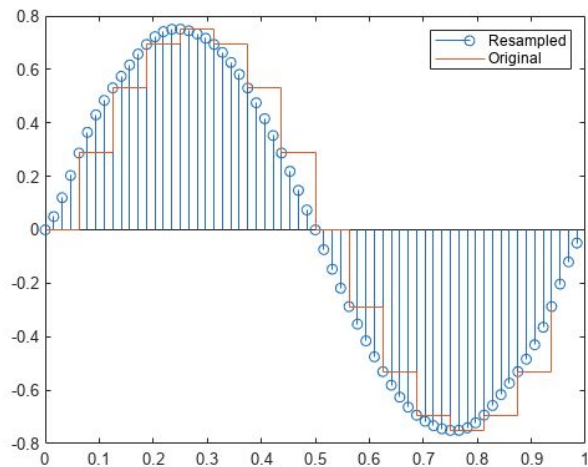
SUIT method for Time Series Preprocessing

All material is available at github.com/hip023/suit



SAMPLE

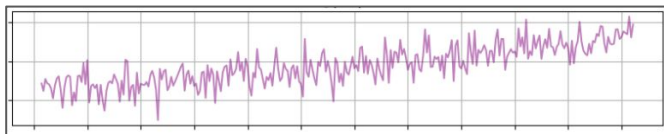
2.1 DownSampling Cumulative Data



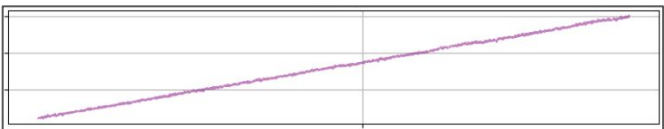
SAMPLE

Just-right frequency

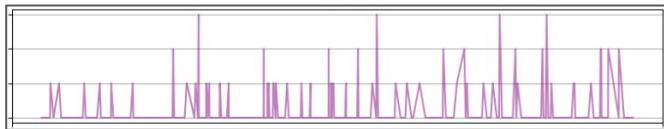
1. Delta



2. Cumulative



3. Events



PHIDS

Plot, Head, Info, Describe,
Set Index

As Freq

Only for cumulative data

```
df[col].asfreq(freq='1D')
```

Resample

Events -> Sampling

```
df.resample('1T').agg(...)
```

aligner reducer

TIME SERIES PREPROCESSING WITH THE **SUIT** METHOD – CHEAT SHEET

[@pazpazthecoder](#), github.com/hip023/suit

November 2022

TIME SERIES PREPROCESSING WITH THE SUIT METHOD – CHEAT SHEET

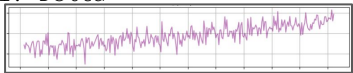
@pazpazthecoder, github.com/hip023/suit

November 2022

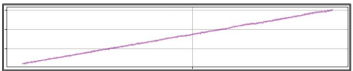
SAMPLE

Just-right frequency

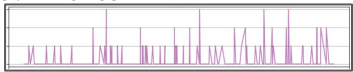
1. Delta



2. Cumulative



3. Events



PHIDS

Plot, Head, Info, Describe,
Set Index

As Freq

Only for cumulative data

```
df[col].asfreq(freq='1D')
```

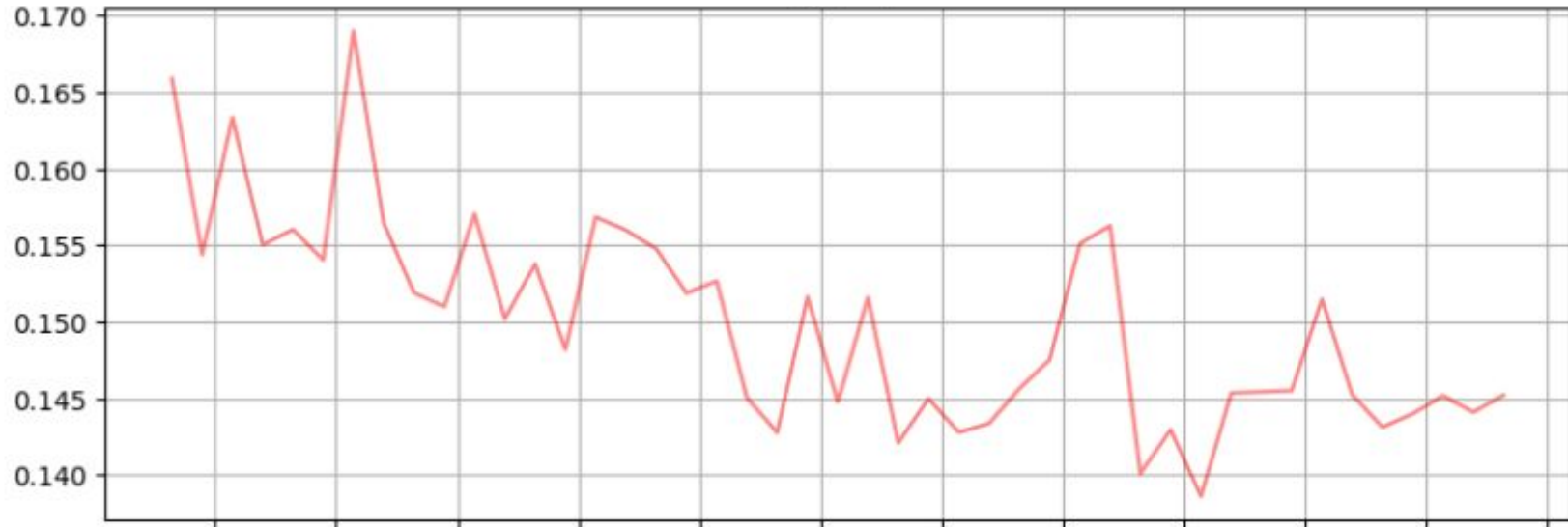
Resample

Events -> Sampling

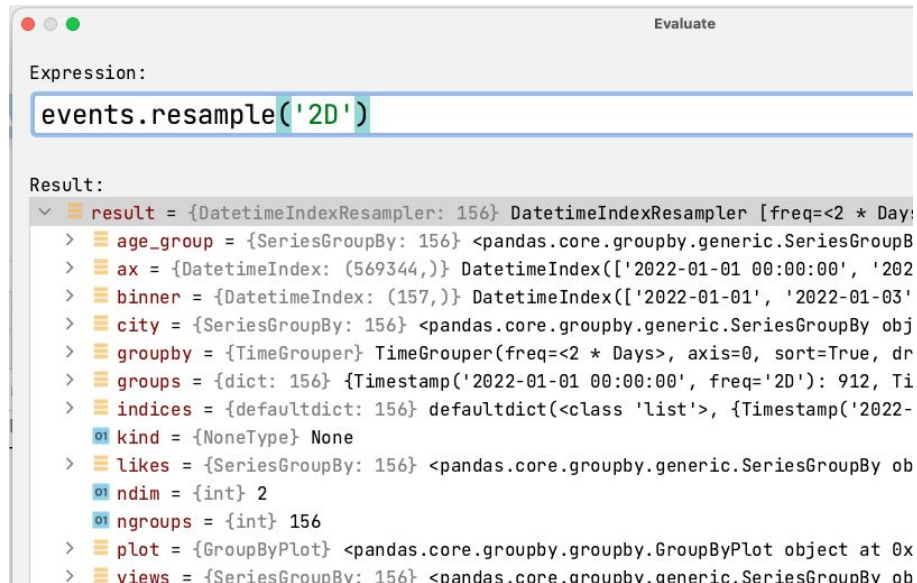
```
df.resample('1T').agg(...)
```

aligner reducer

UNIVARIATE



RESAMPLER



The screenshot shows a Jupyter Notebook window with a title bar containing three colored circles (red, yellow, green) and the word "Evaluate". The "Expression:" input field contains the code `events.resample('2D')`. The "Result:" output field displays a detailed representation of a `DatetimeIndexResampler` object. The object's attributes include `age_group`, `ax`, `binner`, `city`, `groupby`, `groups`, `indices`, `kind`, `likes`, `ndim`, `ngroups`, `plot`, and `views`. The `ngroups` attribute is highlighted with a blue selection box.

```
Expression:
events.resample('2D')

Result:
result = {DatetimeIndexResampler: 156} DatetimeIndexResampler [freq=<2 * Day:
> age_group = {SeriesGroupBy: 156} <pandas.core.groupby.generic.SeriesGroupB
> ax = {DatetimeIndex: (569344,)} DatetimeIndex(['2022-01-01 00:00:00', '202
> binner = {DatetimeIndex: (157,)} DatetimeIndex(['2022-01-01', '2022-01-03'
> city = {SeriesGroupBy: 156} <pandas.core.groupby.generic.SeriesGroupBy obj
> groupby = {TimeGrouper} TimeGrouper(freq=<2 * Days>, axis=0, sort=True, dr
> groups = {dict: 156} {Timestamp('2022-01-01 00:00:00', freq='2D'): 912, Ti
> indices = {defaultdict: 156} defaultdict(<class 'list'>, {Timestamp('2022-
01 kind = {NoneType} None
> likes = {SeriesGroupBy: 156} <pandas.core.groupby.generic.SeriesGroupBy ob
01 ndim = {int} 2
01 ngroups = {int} 156
> plot = {GroupByPlot} <pandas.core.groupby.groupby.GroupByPlot object at 0x
> views = {SeriesGroupBy: 156} <pandas.core.groupby.generic.SeriesGroupBy ob
```

RESAMPLER

Expression:

```
list(events.resample('2D'))
```

Result:

```
1 result = {list: 156} [(Timestamp('2022-01-01 00:00:00', freq='2D'),
2
3   000 = {tuple: 2} (Timestamp('2022-01-01 00:00:00', freq='2D'),
4
5   > 0 = {Timestamp} 2022-01-01 00:00:00
6
7   > 1 = {DataFrame: (912, 4)} views likes city age_grou
8
9   01 __len__ = {int} 2
10
11  > 001 = {tuple: 2} (Timestamp('2022-01-03 00:00:00', freq='2D'),
12
13  > 002 = {tuple: 2} (Timestamp('2022-01-05 00:00:00', freq='2D'),
14
15  > 003 = {tuple: 2} (Timestamp('2022-01-07 00:00:00', freq='2D'),
16
17  > 004 = {tuple: 2} (Timestamp('2022-01-09 00:00:00', freq='2D'),
18
19  > 005 = {tuple: 2} (Timestamp('2022-01-11 00:00:00', freq='2D'),
```

RESAMPLING AGAIN

```
In 24 1 def get_tlv_likes(df: pd.DataFrame) -> int:
      2     tel_aviv = df[df['city'] == 'tel-aviv']
      3     likes = tel_aviv['likes'].sum()
      4     return likes
      5
      6 events.resample('1D').apply(get_tlv_likes)
```

```
Out 24  datetime
        2022-01-01    48
        2022-01-02    55
        2022-01-03    94
        2022-01-04    44
        2022-01-05    40
```

RESAMPLING AGAIN

```
In 24 1 def get_tlv_likes(df: pd.DataFrame) -> int:
      2     tel_aviv = df[df['city'] == 'tel-aviv']
      3     likes = tel_aviv['likes'].sum()
      4     return likes
      5
      6 events.resample('1D').apply(get_tlv_likes)
```

```
Out 24  datetime
        2022-01-01    48
        2022-01-02    55
        2022-01-03    94
        2022-01-04    44
        2022-01-05    40
```

UNIVARIATE

Vanilla Feature Engineering

```
df.resample('1D').aggregate(...)
```

```
df.resample('1D').apply(lambda x: ...)
```

```
def get_lambda(x):  
    return ...
```

```
data.resample('1D').apply(get_lambda)
```

Datetime Features

```
resampler.index.isocalendar()
```

```
resampler.index.month
```

TIME SERIES PREPROCESSING WITH THE SUIT METHOD – CHEAT SHEET

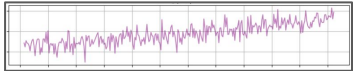
@pazpazthecoder, github.com/hip023/suit

November 2022

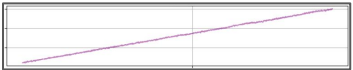
SAMPLE

Just-right frequency

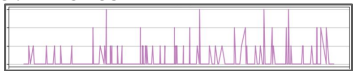
1. Delta



2. Cumulative



3. Events



PHIDS

Plot, Head, Info, Describe,
Set Index

As Freq

Only for cumulative data

```
df[col].asfreq(freq='1D')
```

Resample

Events -> Sampling

```
df.resample('1T').agg(...)
```

aligner reducer

TIME SERIES PREPROCESSING WITH THE SUIT METHOD – CHEAT SHEET

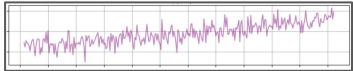
@pazpazthecoder, github.com/hip023/suit

November 2022

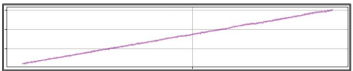
SAMPLE

Just-right frequency

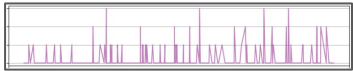
1. Delta



2. Cumulative



3. Events



PHIDS

Plot, Head, Info, Describe,
Set Index

As Freq

Only for cumulative data

```
df[col].asfreq(freq='1D')
```

Resample

Events -> Sampling

```
df.resample('1T').agg(...)
```

aligner reducer

UNIVARIATE

Vanilla Feature Engineering

```
df.resample('1D').aggregate(...)
```

```
df.resample('1D').apply(lambda x: ...)
```

```
def get_lambda(x):  
    return ...
```

```
data.resample('1D').apply(get_lambda)
```

Datetime Features

```
resampler.index.isocalendar()
```

```
resampler.index.month
```


ISOLATION

city datetime	kadima	ramat-gan	rishon	tel-aviv
2022-01-01	0.061538	0.000000	0.200000	0.738462
2022-01-02	0.120482	0.012048	0.204819	0.662651
2022-01-03	0.025210	0.008403	0.176471	0.789916
2022-01-04	0.044118	0.014706	0.294118	0.647059
2022-01-05	0.031579	0.000000	0.242105	0.726316
2022-01-06	0.047619	0.009524	0.323810	0.619048
2022-01-07	0.061947	0.000000	0.061947	0.876106
2022-01-08	0.027027	0.027027	0.324324	0.621622

ISOLATION

Evaluate

Expression:

```
list(events.groupby(groupers))
```

Result:

```
> result = {list: 1244} [((Timestamp('2022-01-01 00:00:00', freq='D'), 'kadima'),
0000 = {tuple: 2} ((Timestamp('2022-01-01 00:00:00', freq='D'), 'kadima'),
> 0 = {tuple: 2} (Timestamp('2022-01-01 00:00:00', freq='D'), 'kadima')
> 1 = {DataFrame: (30, 4)} views likes city age_group [datetime
01 __len__ = {int} 2
```

ISOLATE

Combine resampling with other features!

```
groupers = [pd.Grouper(freq='1D'), col]  
by_date_and_col = df.groupby(groupers)...
```

aligner reducer

We can manipulate different samplers together!

```
by_date = df.groupby(pd.Grouper(freq='1D'))...  
(by_date_and_col / by_date).unstack()
```

TIME SERIES PREPROCESSING WITH THE SUIT METHOD – CHEAT SHEET

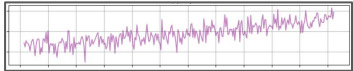
@pazpazthecoder, github.com/hip023/suit

November 2022

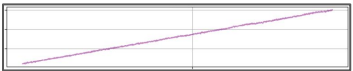
SAMPLE

Just-right frequency

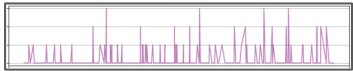
1. Delta



2. Cumulative



3. Events



PHIDS

Plot, Head, Info, Describe,
Set Index

As Freq

Only for cumulative data

```
df[col].asfreq(freq='1D')
```

Resample

Events -> Sampling

```
df.resample('1T').agg(...)
```

aligner reducer

UNIVARIATE

Vanilla Feature Engineering

```
df.resample('1D').aggregate(...)
```

```
df.resample('1D').apply(lambda x: ...)
```

```
def get_lambda(x):  
    return ...
```

```
data.resample('1D').apply(get_lambda)
```

Datetime Features

```
resampler.index.isocalendar()
```

```
resampler.index.month
```

TIME SERIES PREPROCESSING WITH THE SUIT METHOD – CHEAT SHEET

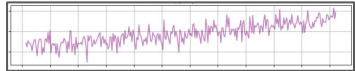
@pazpazthecoder, github.com/hip023/suit

November 2022

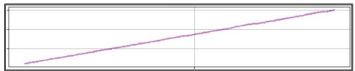
SAMPLE

Just-right frequency

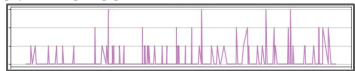
1. Delta



2. Cumulative



3. Events



PHIDS

Plot, Head, Info, Describe,
Set Index

As Freq

Only for cumulative data

```
df[col].asfreq(freq='1D')
```

Resample

Events -> Sampling

```
df.resample('1T').agg(...)
```

aligner reducer

UNIVARIATE

Vanilla Feature Engineering

```
df.resample('1D').aggregate(...)
```

```
df.resample('1D').apply(lambda x: ...)
```

```
def get_lambda(x):  
    return ...
```

```
data.resample('1D').apply(get_lambda)
```

Datetime Features

```
resampler.index.isocalendar()
```

```
resampler.index.month
```

ISOLATE

Combine resampling with other features!

```
groupers = [pd.Grouper(freq='1D'), col]
```

```
by_date_and_col = df.groupby(groupers)...
```

aligner reducer

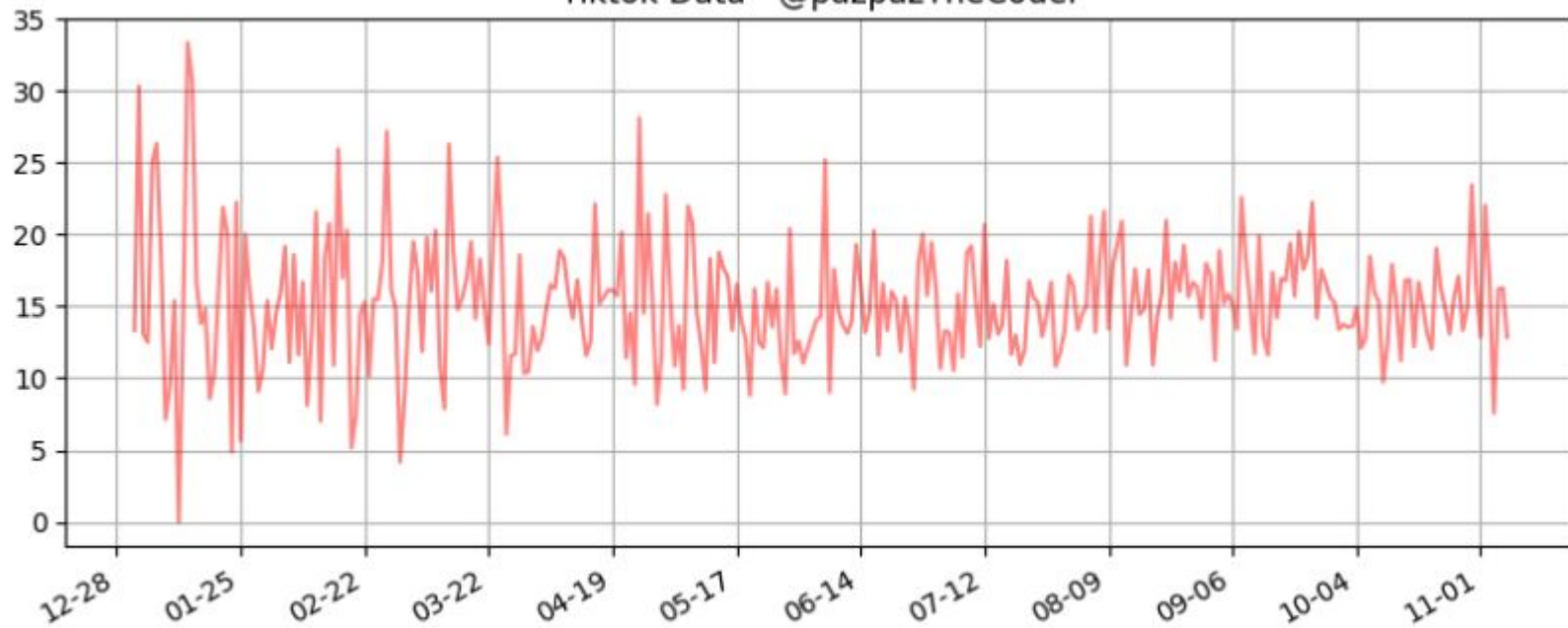
We can manipulate different samplers together!

```
by_date = df.groupby(pd.Grouper(freq='1D'))...
```

```
(by_date_and_col / by_date).unstack()
```

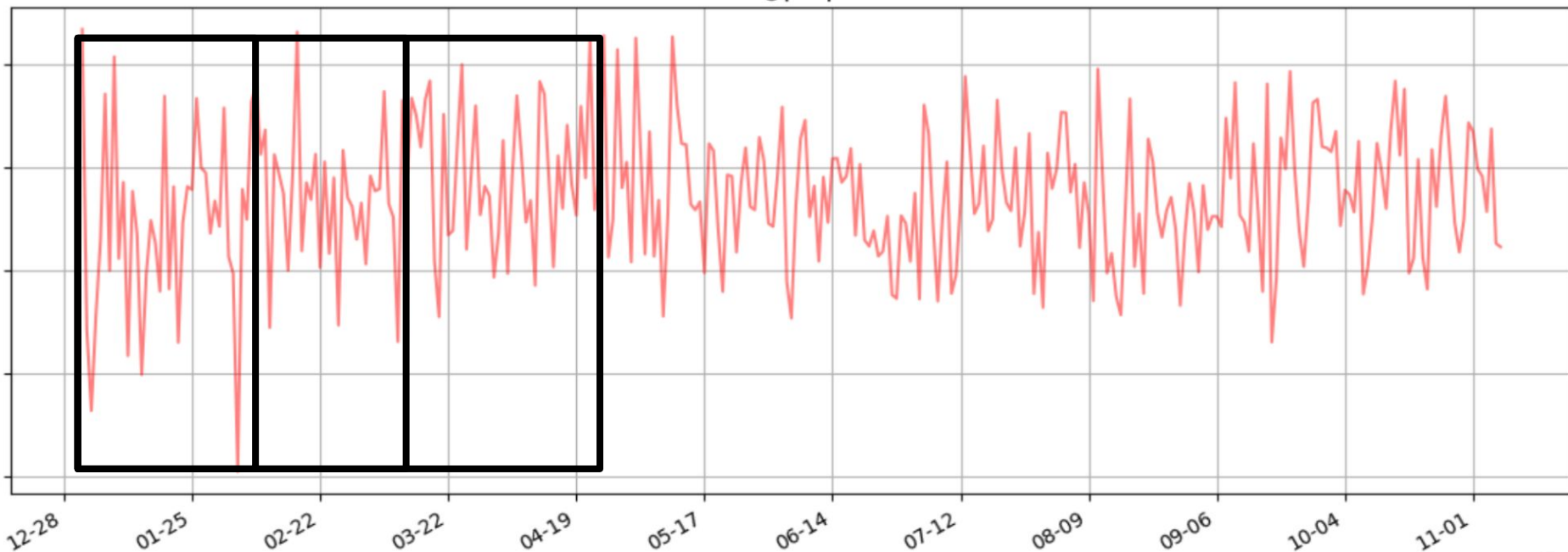
TRANSFORMATION

Tiktok Data - @pazpazTheCoder

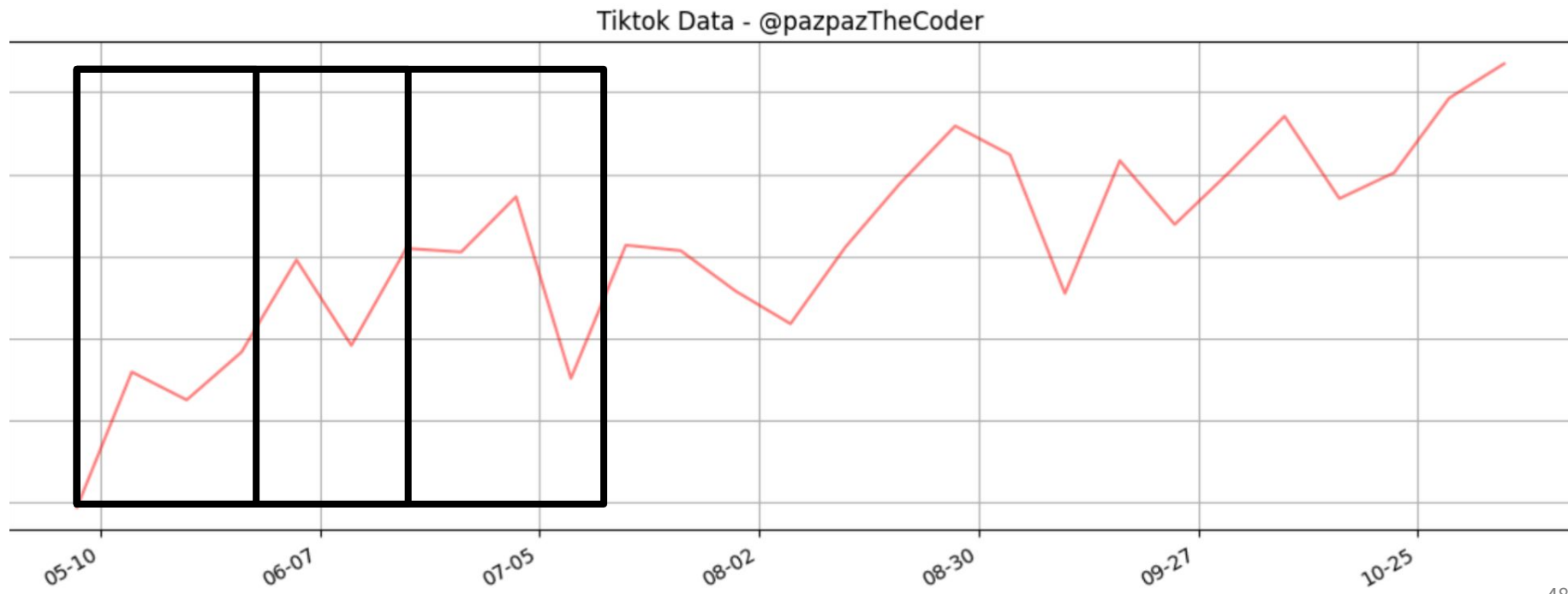


STATIONARY DATA

Tiktok Data - @pazpazTheCoder



NON STATIONARY



STATIONARY TEST

The distribution of

$$(y_t, y_{t+1}, \dots, y_{t+s})$$

For every s , is independent of t .

STATIONARY TEST

Mean is constant

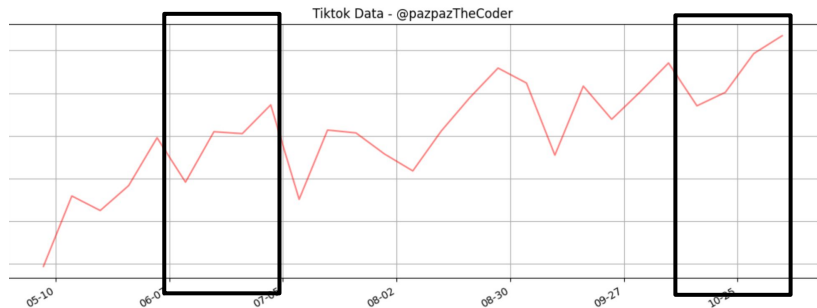
Std is constant

No seasonality

$$y_t = \mu + u_t$$

STATIONARY TEST

Global- local test



- $\mu \setminus \sigma$ change
- seasonality

$$(y_{t_1}, \dots, y_{t_1+s}) \quad (y_{t_2}, \dots, y_{t_2+s})$$

STATIONARY TEST

Production ready?

ADF TEST

Augmented **D**ickey-**F**uller Test



DF TEST

Dickey-Fuller Test

$$y_t = \mu + u_t$$

$$y_t = \mu + \rho y_{t-1} + u_t$$

Null Hypothesis

$$\rho = 1$$

Alternative Hypothesis

$$\rho < 1$$

ADF TEST

Augmented **D**ickey-**F**uller Test

$$y_t = \mu + u_t$$

$$y_t = \mu + \sum_{i=1} \rho_i \cdot y_{t-i} + u_t$$

Null Hypothesis

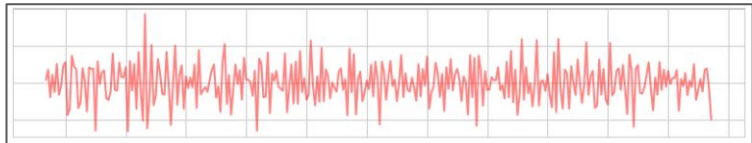
$$\rho = 1$$

Alternative Hypothesis

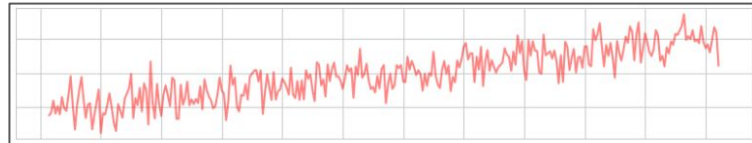
$$\rho < 1$$

TRANSFORM

Augmented Dickey-Fuller test stationary



non-stationary



```
from statsmodels.tsa.  
stattools import adfuller
```

```
adf = adfuller(data['x'])  
if adf[0] > 0.05:  
    print("data is not  
          stationary")
```

Transformation Techniques

“Features of Features”

lags / diffs
second derivative
combined transformation

TIME SERIES PREPROCESSING WITH THE SUIT METHOD – CHEAT SHEET

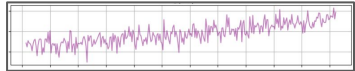
@pazpazthecoder, github.com/hip023/suit

November 2022

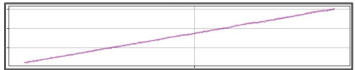
SAMPLE

Just-right frequency

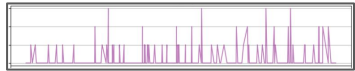
1. Delta



2. Cumulative



3. Events



PHIDS

Plot, Head, Info, Describe,
Set Index

As Freq

Only for cumulative data

```
df[col].asfreq(freq='1D')
```

Resample

Events -> Sampling

```
df.resample('1T').agg(...)
```

aligner reducer

UNIVARIATE

Vanilla Feature Engineering

```
df.resample('1D').aggregate(...)
```

```
df.resample('1D').apply(lambda x: ...)
```

```
def get_lambda(x):  
    return ...
```

```
data.resample('1D').apply(get_lambda)
```

Datetime Features

```
resampler.index.isocalendar()
```

```
resampler.index.month
```

ISOLATE

Combine resampling with other features!

```
groupers = [pd.Grouper(freq='1D'), col]
```

```
by_date_and_col = df.groupby(groupers)...
```

aligner reducer

We can manipulate different samplers together!

```
by_date = df.groupby(pd.Grouper(freq='1D'))...
```

```
(by_date_and_col / by_date).unstack()
```

TIME SERIES PREPROCESSING WITH THE SUIT METHOD – CHEAT SHEET

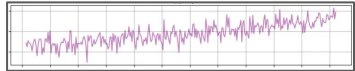
@pazpazthecoder, github.com/hip023/suit

November 2022

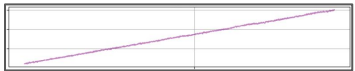
SAMPLE

Just-right frequency

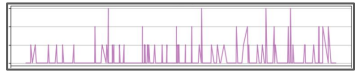
1. Delta



2. Cumulative



3. Events



PHIDS

Plot, Head, Info, Describe,
Set Index

As Freq

Only for cumulative data

```
df[col].asfreq(freq='1D')
```

Resample

Events -> Sampling

```
df.resample('1T').agg(...)
```

aligner

reducer

UNIVARIATE

Vanilla Feature Engineering

```
df.resample('1D').aggregate(...)
```

```
df.resample('1D').apply(lambda x: ...)
```

```
def get_lambda(x):  
    return ...
```

```
data.resample('1D').apply(get_lambda)
```

Datetime Features

```
resampler.index.isocalendar()
```

```
resampler.index.month
```

ISOLATE

Combine resampling with other features!

```
groupers = [pd.Grouper(freq='1D'), col]
```

```
by_date_and_col = df.groupby(groupers)...
```

aligner

reducer

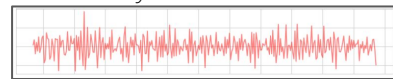
We can manipulate different samplers together!

```
by_date = df.groupby(pd.Grouper(freq='1D'))...
```

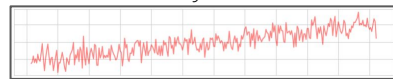
```
(by_date_and_col / by_date).unstack()
```

TRANSFORM

Augmented Dickey-Fuller test
stationary



non-stationary



```
from statsmodels.tsa.  
stattools import adfuller
```

```
adf = adfuller(data['x'])  
if adf[0] > 0.05:  
    print("data is not  
        stationary")
```

Transformation Techniques

“Features of Features”

lags / diffs

second derivative

combined transformation

...