

Homework 3: Data Preparation

CPE232 Data Models

WISIT SUWANNAO 67070501042

Project setup

```
In [1]: # !pip install matplotlib
```

```
In [2]: import pandas as pd

df = pd.read_csv('bike_sharing_demand.csv')
```


```
In [3]: df.head()
```

	season	year	month	hour	holiday	weekday	workingday	weather	temp	feel_temp	humidity	windspeed	count
0	spring	0	1	0	False	6	False	clear	9.84	14.395	0.81	0.0	16
1	spring	0	1	1	False	6	False	clear	9.02	13.635	NaN	0.0	40
2	spring	0	1	2	False	6	False	clear	9.02	13.635	0.80	0.0	32
3	spring	0	1	3	False	6	False	clear	9.84	14.395	0.75	0.0	13
4	spring	0	1	4	False	6	False	clear	9.84	14.395	0.75	0.0	1



```
In [4]: url = "https://knutt.me/"
```

The Secret URL Challenge!

Welcome, brave explorer! Your mission, should you choose to accept it, is to uncover a hidden phrase scattered across the questions below. Each question holds a vital clue—a word or phrase—that will bring you closer to unlocking the **Secret URL!**

Once you have gathered all the hidden words, combine them **in order** and attach them to this URL:
 [https://knutt.me/\[your_combined_phrase\]](https://knutt.me/[your_combined_phrase])

For example, if you discover the words `["quest", "begin"]`, your final URL will be:
 <https://knutt.me/questbegin> 

Are you ready to solve the mystery and reveal the secret link? Let the adventure begin!  

```
In [5]: df.describe()

Out[5]:
```

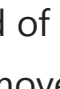
	year	month	hour	weekday	temp	feel_temp	humidity	windspeed	count
count	200.0	200.0	200.000000	200.000000	200.000000	200.000000	170.000000	200.000000	200.000000
mean	0.0	1.0	11.455000	3.160000	9.389000	11.689600	0.559059	13.745452	53.950000
std	0.0	0.0	6.832377	2.235933	3.713618	4.580663	0.176368	8.637962	48.931472
min	0.0	1.0	0.000000	0.000000	3.280000	3.030000	0.280000	0.000000	1.000000
25%	0.0	1.0	6.000000	1.000000	6.560000	9.090000	0.422500	7.001500	12.000000
50%	0.0	1.0	11.000000	3.000000	8.200000	10.985000	0.510000	12.980000	47.000000
75%	0.0	1.0	17.000000	5.000000	10.660000	13.635000	0.690000	19.250775	76.000000
max	0.0	1.0	23.000000	6.000000	18.860000	22.725000	1.000000	36.997400	219.000000











Clue 1: A Note from the Keeper of the Winds

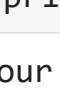


"Traveler, the first clue hides in the mist! To uncover it, follow these steps carefully!"

- Find the moment when the wind was strongest during misty weather.
- Look at that row and gather the numbers hidden in the hour and count columns.
- Add 67 to each number and turn them into letters, but divide count by 3.
- Arrange them in the order given by hour and count to reveal the hidden phrase!

"Solve this mystery, and you will take the first step toward unlocking the secret URL!"  

Monkey Mode Activated! 

- Ooo ooo! Find rows where weather is 'mist'!  
- Pick the row with the BIGGEST windspeed! 
- Grab hour and count columns and divide count by 3!  
- Add 67 to each number!   
- Turn those numbers into LETTERS!  

 Ooo OOO! Secret phrase unlocked!  

```
In [6]: # Find the moment when the wind was strongest during misty weather.
max_wind_speed_in_misty_weather = df[df["weather"].str.contains("mist", case=False, na=False)][["windspeed"]].max()
target_row = df[(df["weather"].str.contains("mist", case=False, na=False) & (df["windspeed"] == max_wind_speed_in_misty_weather))]

# get the hour and count of the target row
hour, count = target_row["hour"].values[0] + 67, target_row["count"].values[0]//3 + 67




# Just change the hour and count to the corresponding ascii character
result = str(chr(int(hour))) + str(chr(int(count)))

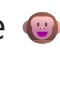
# concatenate the result to the url
url = url + result
print("your current url is: ", url)






your current url is: https://knutt.me/NA
```

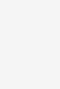

Clue 2: The Hidden Words in the Weather

The next piece of the puzzle lies in the unique weathers that were observed! To find the clue:

- Look at all the different weather conditions recorded in the dataset.
- Take the last two word of each unique weather type you find.
- The combination of these words will lead you to the next step in your adventure!
-  Unravel this mystery, and you'll be one step closer to the secret URL!  

Monkey Mode 

- Ooo ooo! Find all the different weather types! 
- Get the LAST TWO word of each one!  
- Combine the words to move closer to the secret!  

 Monkey magic will lead you to the next clue! 

```
In [7]: # get the unique values of the target column
unique_values = df["weather"].unique()

# get the last two characters of each unique value
last_two_character = [x[-2:] for x in unique_values]

# join all the last two characters
result = "".join(last_two_character)

# concatenate the result to the url
url = url + result
print("your current url is: ", url)

your current url is: https://knutt.me/MaTyin
```

Clue 3: The missing Humidity

Someone tried to hide a secret message in the humidity levels! you need to see this!

```
In [8]: df["humidity"].plot()

Out[8]:
```





```
In [9]: df["humidity"].mean()

Out[9]: np.float64(0.5598588235294117)
```

Missing value in the `humidity` column make their average weird.
Find the missing numbers and combine them to reveal the next part of the secret URL!

Monkey Mode 

- Ooo ooo! Find the missing numbers in the humidity column!  
- Combine the missing numbers to reveal the next part of the secret URL!  

 This is too easy for us. You too you also can do it!  

```
In [10]: # get the number of missing values in humidity column
missing_values = df["humidity"].isnull().sum()

# concatenate the missing values to the url
url = url + str(missing_values)

print("your current url is: ", url)

your current url is: https://knutt.me/MaTyin30
```

Clue 4: Make the Hum(dity)an back!

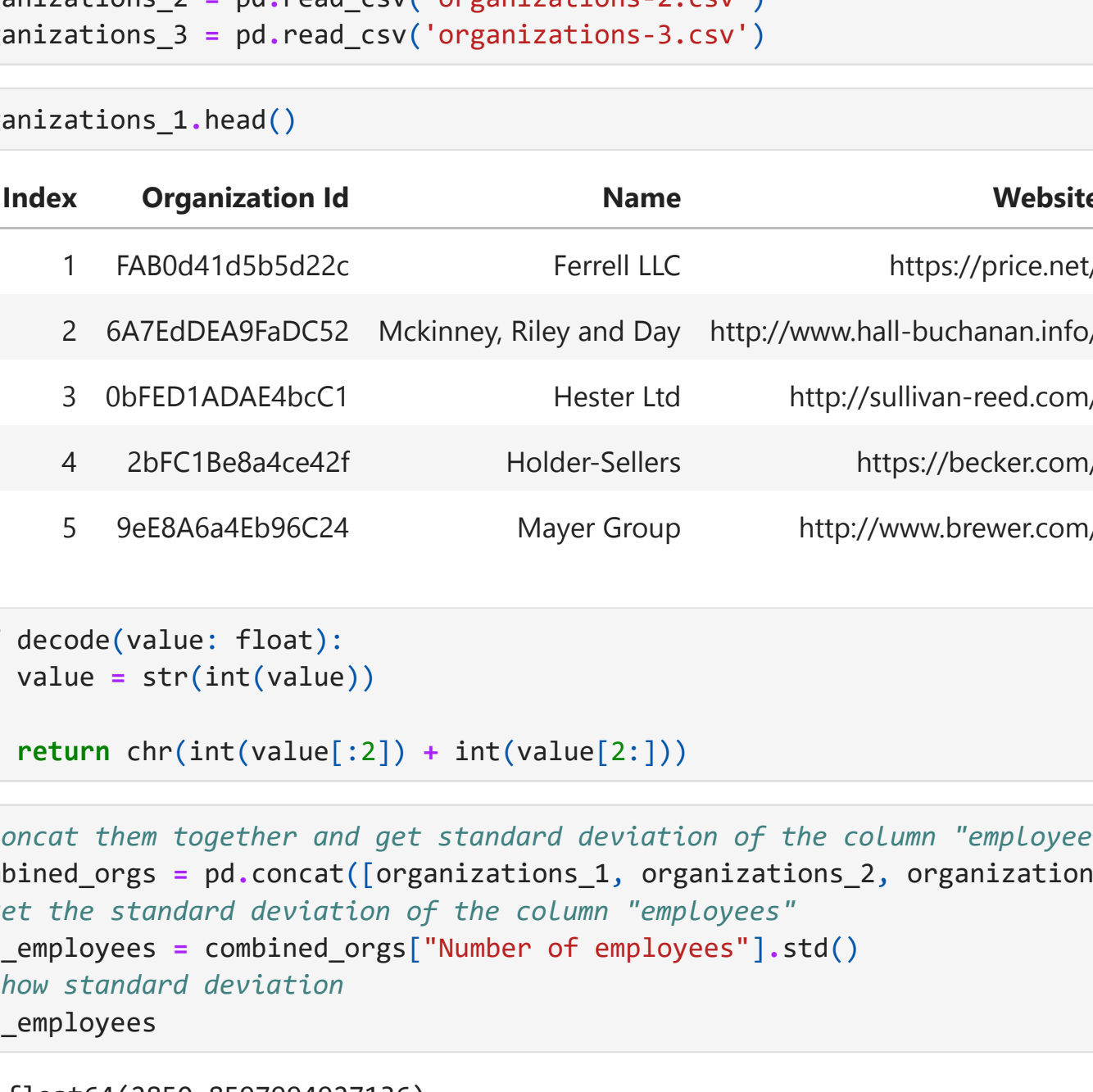
Yes! we got a number of missing humidity from the previous clue. Now, we need to make it back to the original data. This is too hard? *Don't worry about it* you can do it without my help.

```
In [11]: # do it by yourself
# create function that interpolate the missing values in humidity column
for i in range(len(df)):

    # check if the value is missing
    if pd.isnull(df.loc[i, "humidity"]): # checking by using pd.isnull() function
        # if the value is missing, interpolate it with the average of the previous and next value
        if i > 0 and i < len(df) - 1:
            df.loc[i, "humidity"] = (df.loc[i - 1, "humidity"] + df.loc[i + 1, "humidity"]) / 2

In [12]: df["humidity"].plot()

Out[12]:
```



now, find the average of the humidity column and add it to the missing value. Then, you will find the next part of the secret URL!

```
In [13]: average_humidity = df["humidity"].mean()
average_humidity

Out[13]: np.float64(0.5575249999999999)

oh, I forgot to tell you. We only use first 2 decimal places of the average value.

In [14]: # get first 2 decimal of the average humidity
result = str(int(average_humidity*100))


# concatenate the result to the url
url = url + result

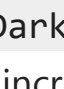


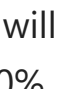

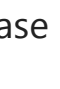
print("your current url is: ", url)




your current url is: https://knutt.me/MaTyin3055L
```

Clue 5: The Secret Message from the different weathers

We almost there! Find an average of each weather type in the dataset. Then use the ascii number of the sum between `clear` weather and difference of `misty` and `rain` weather to reveal the next part of the secret URL!

Monkey Mode 

- Find the average of each weather type!  
- Use the ASCII number of the sum between `clear` weather and difference of `misty` and `rain` weather!  
- Combine the numbers to reveal the next part of the secret URL!  

 You're almost there! Keep going!  

```
In [15]: # use groupby to get the average count of each weather
average_count = df.groupby("weather")["count"].mean()

# get the average count of clear, misty, and rain weather
clear_avg = average_count["clear"]
misty_avg = average_count["misty"]
rain_avg = average_count["rain"]

# get the groupby_character follow by instructions
groupby_character = chr(int(clear_avg - (misty_avg - rain_avg)))

# concatenate the groupby_character to the url
url = url + groupby_character

print("your current url is: ", url)

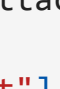
your current url is: https://knutt.me/MaTyin3055L

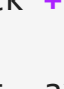
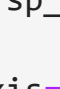
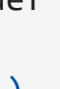
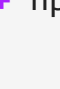
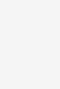
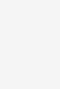
In [16]: print("your final url is: ", url)




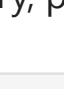
your final url is: https://knutt.me/MaTyin3055L
```

Clue 6: Fusion!

You've made it this far! Now, you just need to combine the dataframe and get the standard deviation of `Number of employees` column, then put it in `decode` tools to reveal the final part of the secret URL!

Monkey Mode 

- Combine the dataframe and get the standard deviation of `Number of employees` column!  
- Use the standard deviation as a phrase to unlock the final part of the secret URL!  
- Put the phrase in the `decode` tools to reveal the final part of the secret URL!  

 Don't be afraid! We will stay with you!   

```
In [17]: organizations_1 = pd.read_csv('organizations-1.csv')
organizations_2 = pd.read_csv('organizations-2.csv')
organizations_3 = pd.read_csv('organizations-3.csv')

In [18]: organizations_1.head()

Out[18]:
```

	Index	Organization Id	Name	Website	Country	Description	Founded	Industry	Number of employees
0	1	FAB0d41d5b5d2zc	Ferrell LLC	https://price.net/	Papua New Guinea	Horizontal empowering knowledgebase	1990	Plastics	3498
1	2	6A7EdEA9AdCS2	McKinney, Riley and Day	http://www.hall-buchanan.info/	Finland	User-centric system-worthy leverage	2015	Glass / Ceramics / Concrete	4952
2	3	0bFED1ADA84bc1C	Hester Ltd	http://sullivan-reed.com/	China	Switchable scalable moratorium	1971	Public Safety	5287
3	4	2bKc1B8bA8c42f	Holder-Sellers	https://becker.com/	Turkmenistan	De-engineered systemic artificial intelligence	2004	Automotive	921
4	5	9eEBA6A6B95C24	Mayer Group	http://www.brewer.com/	Mauritius	Synchronized needs-based challenge	1991	Transportation	7870

```
In [19]: def decode(value: float):
value = str(int(value))

return chr(int(value[:2]) + int(value[2:]))

In [20]: # concat them together and get standard deviation of the column "employees"
combined_orgs = pd.concat([organizations_1, organizations_2, organizations_3], ignore_index=True)
# get the standard deviation of the column "employees"
std_employees = combined_orgs["Number of employees"].std()
# show standard deviation
std_employees

Out[20]: np.float64(2850.8597994927136)

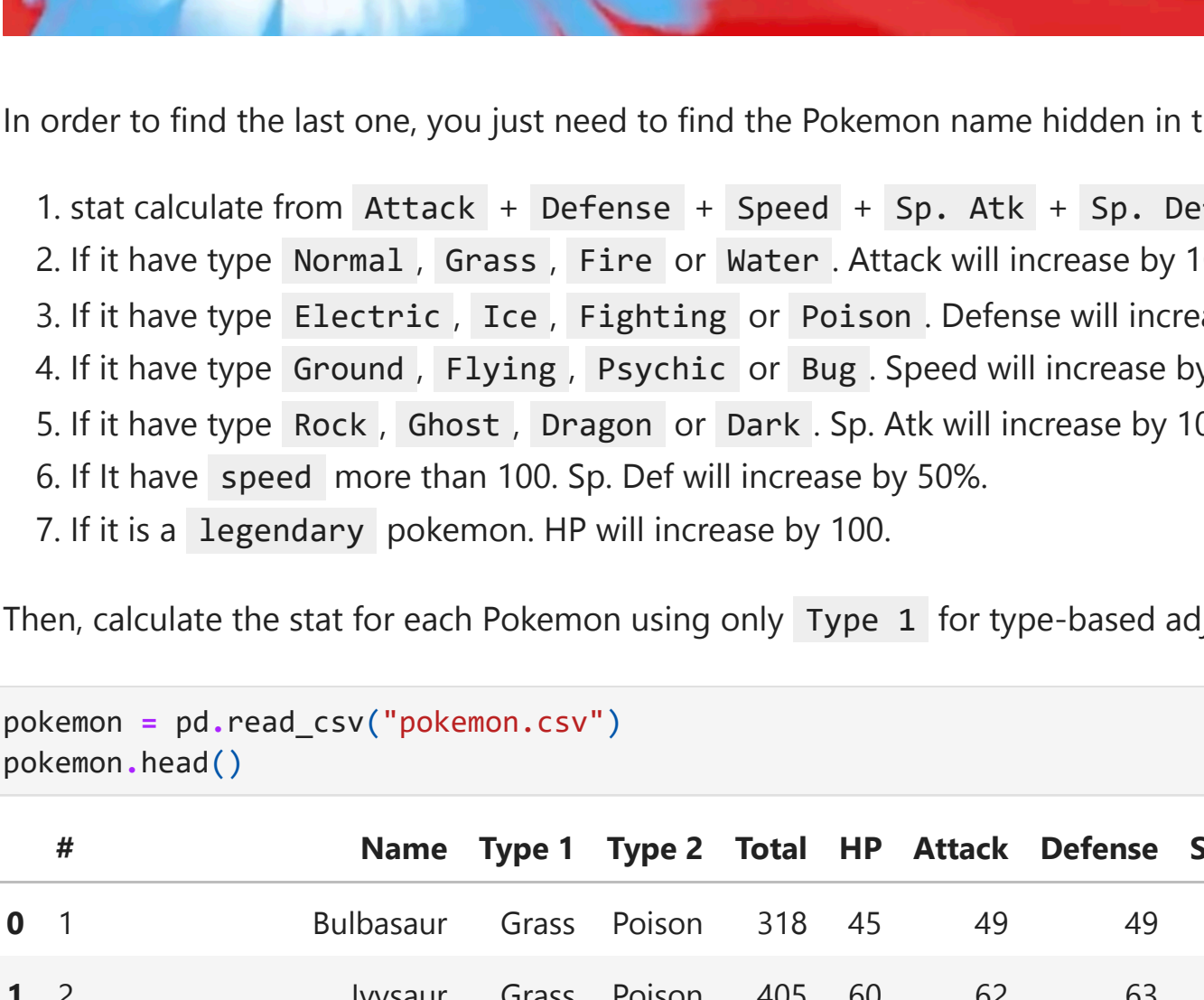
In [21]: url = url + decode(std_employees) # your variable that contains the standard deviation

print("your current url is: ", url)

your current url is: https://knutt.me/MaTyin3055LM45Spikachu
```

Final Clue: WHO'S THAT POKEMON?

Who's that Pokémon



In order to find the last one, you just need to find the Pokemon name hidden in the URL by adding a new column called `stat` that follows the conditions below:

- stat calculate from `Attack` + `Defense` + `Speed` + `Sp. Atk` + `Sp. Def` + `HP`
- If it have type `Normal`, `Grass`, `Fire` or `Water`: Attack will increase by 10%.
- If it have type `Electric`, `Ice`, `Fighting` or `Poison`: Defense will increase by 10%.
- If it have type `Ground`, `Flying`, `Psychic` or `Bug`: Speed will increase by 10%.
- If it have type `Rock`, `Ghost`, `Dragon` or `Dark`: Sp. Atk will increase by 10%.
- If it have `Speed` more than 100, Sp. Def will increase by 50%.
- If it is a `legendary` pokemon, HP will increase by 100.

Then, calculate the stat for each Pokemon using only `Type 1` for type-based adjustments. After that, find the average of the `stat` column across all Pokemon, and append it to the previous URL as an integer.

```
In [22]: pokemon = pd.read_csv('pokemon.csv')
pokemon.head()

Out[22]:
```

	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0	1	Bulbasaur	Grass	Poison	318	45	49	49	65	65	45	1	False
1	2	Ivysaur	Grass	Poison	405	60	62	63	80	80	60	1	False
2	3	Venusaur	Grass	Poison	525	80	82	83	100	100	80	1	False
3	3	VenusaurMega Venusaur	Grass	Poison	625	80	100	123	122	120	80	1	False
4	4	Charmander	Fire	NaN	309	39	52	43	60	50	65	1	False

```
In [23]: # Complete the Final Clue
pokemon["stat"] = pokemon["Attack"] + pokemon["Defense"] + pokemon["Speed"] + pokemon["Sp. Atk"] + pokemon["Sp. Def"] + pokemon["HP"]

def adjust_stat(row):
    attack = row["Attack"]
    defense = row["Defense"]
    speed = row["Speed"]
    sp_atk = row["Sp. Atk"]
    sp_def = row["Sp. Def"]
    hp = row["HP"]
    t1 = row["Type 1"]
    t_group_1 = ['Normal', 'Grass', 'Fire', 'Water']
    t_group_2 = ['Electric', 'Ice', 'Fighting', 'Poison']
    t_group_3 = ['Ground', 'Flying', 'Psychic', 'Bug']
    t_group_4 = ['Rock', 'Ghost', 'Dragon', 'Dark']

    if t1 in t_group_1:
        attack *= 1.10
    if t1 in t_group_2:
        defense *= 1.10
    if t1 in t_group_3:
        speed *= 1.10
    if t1 in t_group_4:
        sp_atk *= 1.10
    if row["Speed"] > 100:
        sp_def *= 1.50
    if row["Legendary"]:
        hp += 100

    return attack + defense + speed + sp_atk + sp_def + hp

pokemon["stat"] = pokemon.apply(adjust_stat, axis=1)

In [24]: pokemon["stat"].mean()

Out[24]: np.float64(455.58074999999997)

In [25]: url = url + str(int(pokemon["stat"].mean()))

Whoa! Now you're about to find out what it is!

In [26]: print("who's that Pokemon?: ", url)

Who's that Pokemon?: https://knutt.me/MaTyin3055LM455

After you get the name from the greatest discovery, put it here to unlock the final URL:



In [27]: # Pokemon name
# Please replace the '...' with the name found in the URL printed above!
pokemon_name = "Pikachu".lower()

url = url + pokemon_name


In [28]: print("Your Final URL is: {}".format(url))

Your Final URL is: https://knutt.me/MaTyin3055LM455Pikachu

Final Mission (Optional)

Access the secret URL and complete your quest!  

Question: What is the final secret URL?
Ans: https://knutt.me/NWartyn3055LM455Pikachu

Enjoy the adventure! 
```