

COMS 3251 Spring 2021: Lab 5

YOUR NAME(s): Vicente Farias

YOUR UNI(s): vf2272

COVID-19 Vaccinations

In this lab you will be looking at real-time data tracking the progress of COVID-19 vaccinations in the United States. There are several public data sets; we will be using the one from [Our World in Data](#). Each row of this data set provides vaccination numbers on a given date of a particular state or the US itself; each column indicates metrics such as daily vaccinations or people fully vaccinated.

First import the latest data set by uploading the provided CSV file under Files to the left and then running the code cell below. We will be using [pandas](#), Python Data Analysis Library, to import and read the data.

```
In [1]: import pandas as pd
import numpy as np
import numpy.linalg as nla
import matplotlib.pyplot as plt

ts = pd.read_csv('/content/us_state_vaccinations.csv')
data = ts.values
pd.DataFrame(data)
```

Out[1]:

	0	1	2	3	4	5	6	7	8	9	10	11	12
0	2021-01-12	Alabama	78134	377025	70861	0.15	1.59	7270	1.45	7.69	NaN	NaN	NaN
1	2021-01-13	Alabama	84040	378975	74792	0.19	1.71	9245	1.53	7.73	5906	5906	1205
2	2021-01-14	Alabama	92300	435350	80480	NaN	1.88	NaN	1.64	8.88	8260	7083	1445
3	2021-01-15	Alabama	100567	444650	86956	0.28	2.05	13488	1.77	9.07	8267	7478	1525
4	2021-01-16	Alabama	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	7557	7498	1529
...
4763	2021-03-21	Wyoming	221336	314015	134703	14.77	38.24	85508	23.27	54.26	343	1347	2327
4764	2021-03-22	Wyoming	221452	314015	134800	14.78	38.26	85556	23.29	54.26	116	1279	2210
4765	2021-03-23	Wyoming	235171	324875	142012	15.95	40.63	92303	24.54	56.13	13719	3029	5234
4766	2021-03-24	Wyoming	235424	336475	142147	15.97	40.68	92428	24.56	58.14	253	3058	5284
4767	2021-03-25	Wyoming	242919	350735	146634	16.54	41.97	95709	25.34	60.6	7495	4124	7126

4768 rows × 14 columns

Data Extraction

`data` contains all the information contained in the CSV file. We will look at the portion of the data pertaining to the US as a whole. We will first use `numpy.where` to extract the correct row indices. The example below then further extracts the data of column `8`, which corresponds to the number of people who have received at least one vaccine dose per 100 people. Finally, this information is plotted on a scatter plot.

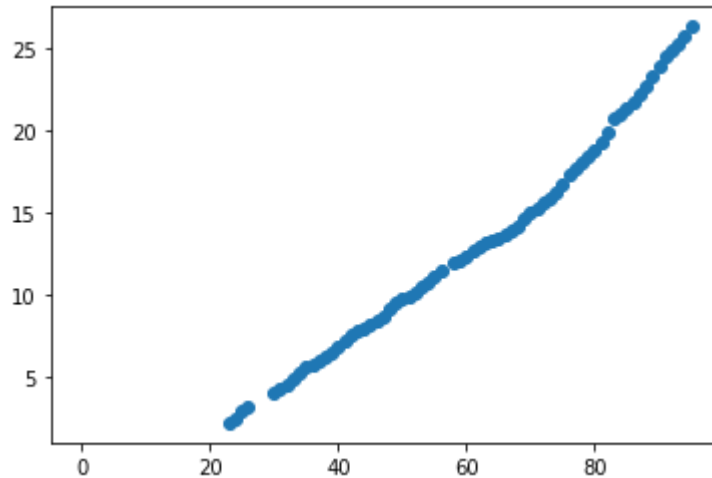
```
In [2]: indices = np.where(data == 'United States')[0]

# people vaccinated per hundred
vax = data[indices,8]
n = vax.size
days = np.arange(n)

print(vax)
plt.scatter(days, vax)
```

```
[nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan
nan nan nan nan nan 2.23 2.41 2.92 3.19 nan nan nan 4.1 4.3 4.53 4.89
5.24 5.57 5.8 5.99 6.23 6.54 6.89 7.25 7.59 7.84 7.96 8.18 8.41 8.71 9.11
9.51 9.74 9.9 10.18 10.46 10.79 11.16 11.53 nan 11.95 12.13 12.36 12.64
12.89 13.14 13.29 13.42 13.63 13.88 14.21 14.59 14.99 15.28 15.59 15.92
16.28 16.73 17.28 17.73 18.07 18.4 18.81 19.3 19.87 20.75 21.02 21.4
21.73 22.19 22.74 23.26 23.91 24.52 24.93 25.28 25.74 26.31]
```

Out[2]: <matplotlib.collections.PathCollection at 0x7f875189da10>



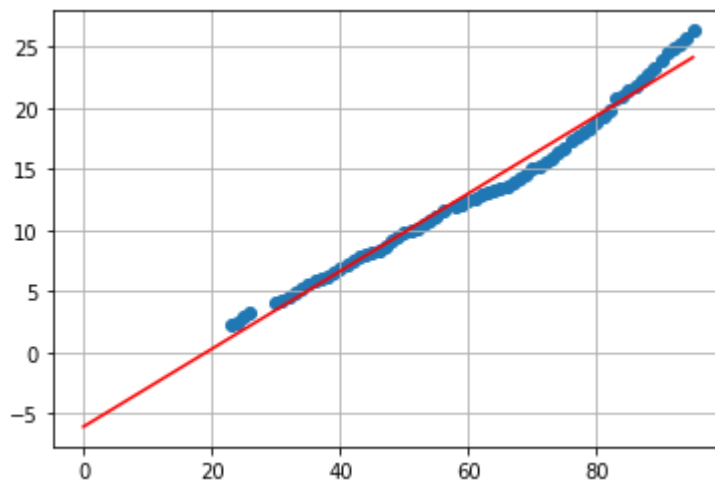
Real-world data is messy. Notice in particular the missing data points, most of which is at the beginning. The value 0 on the horizontal axis is associated with December 20, 2020, while the value 23 is associated with January 12, 2021 and the first available datum.

Linear Regression

We can perform linear regression on this data and fit a trend line, as long as we are careful with the missing data. We first form the design matrix and output vector using the data structures above. We then remove the rows from both structures containing missing data. After solving for the model parameters, we can then compute and plot the fitted line on the data.

```
In [3]: nan_idx = np.where(pd.isnull(vax))
X = np.column_stack((np.ones(n), days))
X = np.delete(X, nan_idx, 0)
b = np.delete(vax, nan_idx)

fit = nla.lstsq(X, b.astype(float), rcond=None)
beta = fit[0]
plt.scatter(days, vax)
plt.plot(days, np.column_stack((np.ones(n), days)) @ beta, 'r')
plt.grid()
```



The fitted line is shown in red above. It should look like a pretty good fit for the most part, with a couple of caveats. The most recent data should appear to be veering away from a linear trend to a steeper rise (a good thing!). On the other end of the domain, you should see that the line extends downward into "negative people vaccinated" at the beginning. Of course this makes no sense, and one takeaway is that this particular data has not always been well described by the same linear trend.

You will be exploring other aspects of this data set in the exercises below.

Exercises

PROBLEM 1 (20 points)

Let's now investigate the number of people who are *fully* vaccinated per hundred. This data is given in column index `5`. If you replicate the process in the example above, you will see that the fit again extends into the negative for the first few weeks.

We can try to fix this problem by making the following modifications. First, set the first value of the data (on day 0) to be 0 so that it's no longer `nan`. Second, a linear fit may not be best. Instead, we can use polynomial or exponential basis functions, the latter of which can be approximated by a sum of increasing powers of the data values.

Use the hypothesis

$$\hat{f}(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4.$$

As a hint, this will require you to include additional columns in the design matrix X . Each additional column should correspond to one of the new basis functions of x . Output a plot with the fit overlaid on top of the scatter plot of points, and then answer the questions below.

```
In [9]: indices = np.where(data == 'United States')[0]
        # people fully vaccinated per hundred
```

```

vax_full = data[indices, 5]
n = vax_full.size
days = np.arange(n)

vax_full[0] = 0

nan_idx = np.where(pd.isnull(vax_full))
X = np.column_stack((np.ones(n), days, days**2, days**3, days**4))
X = np.delete(X, nan_idx, 0)
b = np.delete(vax_full, nan_idx)

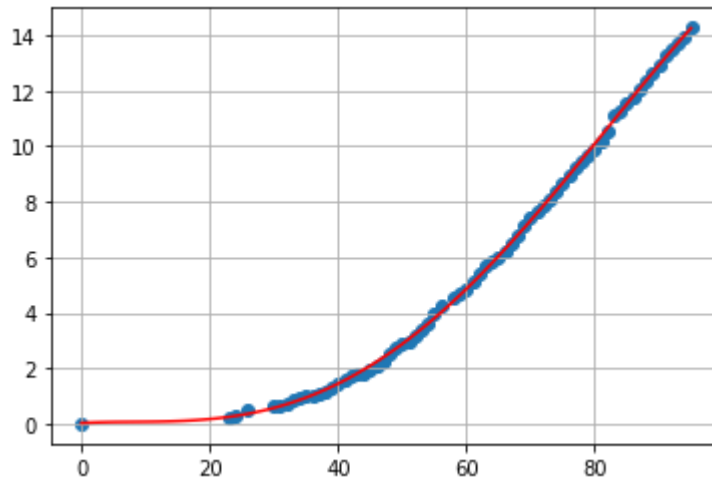
fit = nla.lstsq(X, b.astype(float), rcond=None)
beta = fit[0]
print(beta)
plt.scatter(days, vax_full)
plt.plot(days, np.column_stack((np.ones(n), days, days**2, days**3, days**4)) @ beta,
plt.grid()

```

```

[ 2.49667632e-02  1.20279884e-02 -1.35838967e-03  6.24586014e-05
 -3.46218259e-07]

```



1. Why is this fit able to avoid the problem that the original linear fit encountered for the first 23 days of the domain? Other than this problem, what other metric(s) can you compare to show that this hypothesis may be better than a strictly linear one?
1. Report the β_i coefficients. Looking at their magnitudes, which of the basis functions of the hypothesis carry the most "weight"? If we have to drop one or two of the basis functions, which ones should we drop first?

ENTER YOUR RESPONSES HERE

1. This fit is able to avoid the problem that the linear fit has because there is a data point (0, 0), which causes the quartic function to remain positive. We could calculate the norm of the difference between the predictions and the actual values, which may show that this hypothesis is better than the linear hypothesis.
2. $\beta_0 = 2.49667632 \times 10^{-2}$, $\beta_1 = 1.20279884 \times 10^{-2}$, $\beta_2 = -1.35838967 \times 10^{-3}$, $\beta_3 = 6.24586014 \times 10^{-5}$, $\beta_4 = -3.46218259 \times 10^{-7}$

The lower degree basis functions have the most "weight". If we had to drop terms, we should start with the higher order terms, as they have the least weight.

PROBLEM 2 (20 points)

Regression is typically used to infer correlations and causal relationships between variables, but it can also be used for prediction. Let's now look at total vaccinations (first and second doses combined) per 100 people. This data is given in column index 6.

Replicate the entire procedure of the previous problem on this data, although you will not have to explicitly set the value at 0. Then when overlaying the fitted curve on the scatter plot, extend the curve all the way to day 170 instead of stopping at the last day of the data.

```
In [19]: indices = np.where(data == 'United States')[0]

# total vaccinations
vax_full = data[indices, 6]
n = vax_full.size
days = np.arange(n)

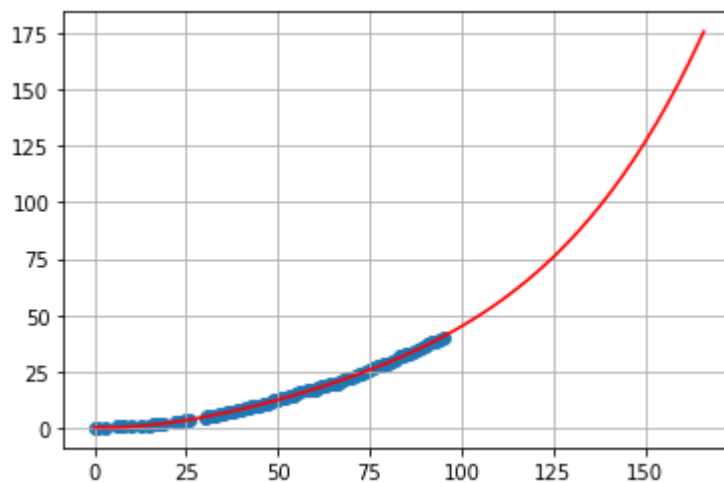
nan_idx = np.where(pd.isnull(vax_full))
X = np.column_stack((np.ones(n), days, days**2, days**3, days**4))
X = np.delete(X, nan_idx, 0)
b = np.delete(vax_full, nan_idx)

fit = nla.lstsq(X, b.astype(float), rcond=None)
beta = fit[0]
print(beta)
plt.scatter(days, vax_full)
days = np.arange(n + 71)
plt.plot(days, np.column_stack((np.ones(n + 71), days, days**2, days**3, days**4)) @ beta)
plt.grid()

days = np.arange(n * 3)
prediction = np.column_stack((np.ones(n * 3), days, days**2, days**3, days**4)) @ beta

print(days[np.where(prediction >= 100)])
print(days[np.where(prediction >= 200)])
```

```
[ 3.86038522e-01 -5.63855544e-02  8.84013059e-03 -7.62711145e-05
 3.81765789e-07]
[139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156
157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174
175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192
193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210
211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228
229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246
247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264
265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282
283 284 285 286 287]
[173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190
191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208
209 210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226
227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244
245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262
263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280
281 282 283 284 285 286 287]
```



1. On approximately what day do we expect to have administered 100 total vaccinations per 100 people in the US? How about 200? Give a reason for why 200 vaccinations may never actually be attained. For reference, all current vaccines require no more than 2 doses, and they have only been approved for adults (which comprise about 74% of the total US population).
2. Briefly describe what you expect the actual vaccination curve to look like over the next 100 days. What should we be mindful of when using regression for prediction, especially for systems with exponential growth?

ENTER YOUR RESPONSES HERE

1. We expect to have administered 100 total vaccinations per 100 people in the US on approximately day 139. We expect to have administered 200 total vaccinations per 100 people in the US on approximately day 173. 200 vaccinations per 100 people implies that every single person in the US has gotten two doses, which is not possible if the vaccine is not approved for all people in the US.

2. We need to be mindful of the fact that there is a limit on the number of people who can get a vaccine and that as more people get vaccinated the rate of vaccination will slow. We expect the actual vaccination curve to look more like a logistic growth curve that levels off as we approach full population vaccination.

PROBLEM 3 (20 points)

Oftentimes, a hypothesis described solely by a few polynomial basis functions is not sufficient. A more flexible hypothesis is the auto-regressive model, which explicitly uses past data to make future predictions:

$$\hat{x}_{t+1} = \beta_1 x_t + \cdots + \beta_M x_{t-M+1}$$

M denotes the memory length of the model; e.g., if $M = 3$ then

$\hat{x}_{t+1} = \beta_1 x_t + \beta_2 x_{t-1} + \beta_3 x_{t-2}$. The data that we are fitting will thus appear in both the output vector as well as the design matrix. The following would be the linear regression problem for $M = 3$:

$$\begin{bmatrix} x_3 & x_2 & x_1 \\ x_4 & x_3 & x_2 \\ \vdots & \vdots & \vdots \\ x_{n-2} & x_{n-3} & x_{n-4} \\ x_{n-1} & x_{n-2} & x_{n-3} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} x_4 \\ x_5 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix}$$

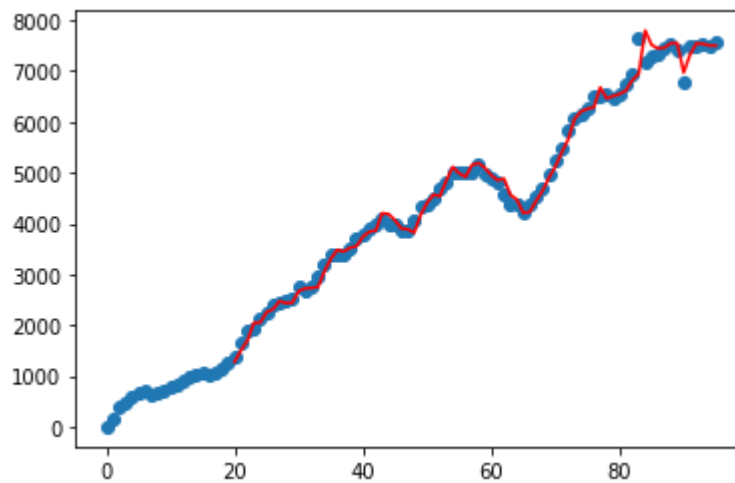
While you can construct these quantities for any M using loops, we are providing you with the `AR` procedure below to form them so that you can focus more on the outputs and analysis. For this last problem, you will look at daily vaccinations per million (column index `12`). The code below will run the regression for a given value of M and plot the output.

```
In [17]: def AR(data, M):
    n = data.size
    idx,_ = np.mgrid[0:n-M, 0:M] + np.arange(0,M)
    X = vax[idx].astype(float)
    b = vax[M:].astype(float)
    return X, b

# daily vaccinations per million
vax = data[indices,12]
vax[0] = 0
n = vax.size
days = np.arange(n)

M = 20
X, b = AR(vax, M)
fit = nla.lstsq(X, b, rcond=None)
beta = fit[0]
plt.scatter(days, vax)
plt.plot(np.arange(M,n), X @ beta, 'r')
```

```
Out[17]: [<matplotlib.lines.Line2D at 0x7f8750e6d150>]
```

You should see that the auto-regressive model generates a very close fit of the data regardless of its trends. One open question is the value of M . The best value is the one that minimizes the root-mean-square (RMS) error, or the norm of the difference between the predicted value and the true value divided by the size of the fitted data. In other words,

$$RMS = \frac{\|\mathbf{b} - X\hat{\beta}\|}{\sqrt{n - M}}$$

where n is the size of the original data set. Note that $n - M$ becomes smaller as M gets larger, since the first $M - 1$ data have no predictions.

Write code below to compute the RMS for values of M ranging from 1 to 50. You will have to use a loop and fit an auto-regressive model for each value. Generate a single plot of RMS versus M as your output.

```
In [21]: # daily vaccinations per million
vax = data[indices,12]
vax[0] = 0
n = vax.size
days = np.arange(n)

def RMS(b, X, beta, n, M):
    return np.linalg.norm(b - X @ beta) / np.sqrt(n - M)

M_arr = np.arange(1, 51)
RMS_arr = []

for M in M_arr:
    X, b = AR(vax, M)
    fit = nla.lstsq(X, b, rcond=None)
    beta = fit[0]
    RMS_arr.append(RMS(b, X, beta, n, M))

plt.scatter(M_arr, RMS_arr)
```

```
Out[21]: <matplotlib.collections.PathCollection at 0x7f8751113c90>
```

