

Project: Data Science :: Persistency of a Drug

Week 13: Deliverables

Presented by: Sammy Suliman, Olivia Foster, and Tahsin Azad

Table of Contents:

Team Member Detail.....	2
Problem Statement.....	3
Project Deadlines.....	3
Data Intake Report.....	4
Data Cleaning.....	5
Data Preprocessing.....	6
Exploratory Analysis.....	6
Model Training.....	7
Flask Application.....	7
Links.....	14

Team Member Details:

Group Name: Call it Version 1.0

Member 1:

- **Name:** Sammy Suliman
- **Email:** sammysuliman@gmail.com
- **Country:** USA
- **College/Company:** UC Santa Barbara
- **Specialization:** Data Science

Member 2:

- **Name:** Olivia Foster
- **Email:** livia.n.foster@gmail.com
- **Country:** USA
- **College/Company:** NT Logistics
- **Specialization:** Data Science

Member 3:

- **Name:** Tahsin Azad
- **Email:** tahsinazad31@gmail.com
- **Country:** USA
- **College/Company:** University of California, Santa Barbara
- **Specialization:** Data Science

Problem Statement:

One of the challenges faced by Pharmaceutical companies is the persistence of a drug (that is, the extent to which a patient will act in accordance with the prescribed time interval, and dose of a medication) as the physician prescribed it. In this problem, we will automate the process of classifying factors that determine the persistence of a drug through Machine Learning and Python.

Drug persistence is a task of classifying different disorders and a patient's medical history to determine the dose and length of dose. In order to train our model, we will need to classify risk factors, medical histories, and disorders. To do this, we will be using a dataset based on over 3000 patients' records.

Project Deadlines:

Week:	Due Date:	Plan:
Week 7	1/19/2023	Problem Statement, Data Collection, Data Report
Week 8	1/26/2023	Data Preprocessing
Week 9	2/2/2023	Feature Extraction
Week 10	2/9/2023	Building Model
Week 11	2/16/2023	Model Result Evaluation
Week 12	2/23/2023	Flask Development and Web Application
Week 13	2/28/2023	Final Submission (Report + Code + Final Submission)

Data Intake Report:

Name: Final Project -- Data Science:: Healthcare - Persistency of a drug:: Group Project

Report date: 1/17/2023

Internship Batch: LISUM16

Version:<1.0>

Data intake by: Olivia Foster

Data intake reviewer: Tahsin Azad

Data storage location:

https://github.com/LiviaNFoster/DataGlacierFinalProject/blob/main/Healthcare_dataset.xlsx

Tabular data details:

Total number of observations	3424
Total number of files	1
Total number of features	69
Base format of the file	xlsx
Size of the data	922 KB

Proposed Approach:

- Unique Row identified using with the “Ptid” column, a unique ID for each patient
- No duplicates of Unique ID confirmed
- Assumptions
 - No false positives or errors during any of the testing
 - No dishonest and unbiased data recorded

Data Cleaning:

We began this process last week and continued to move forward with it.

Null Values:

There were no null values.

Yes versus No:

We had over half of the data be defined as either Y or N, “Yes” or “No” respectively. We replaced Y and N with 1 and 0 respectively.

Renaming Columns:

Because many of the columns had long and difficult names or were not great descriptors, we went ahead and renamed many of them to make referencing easier and to add clarity.

Case sensitive:

Made all words lowercase to make referencing easier without having to worry about case sensitivity.

Acknowledgement of outliers:

In exploring our data, we discovered that because our data is mostly binary (such as above in the ‘Yes versus No’ section) that there were no significant outliers at this time.

‘Age_Bucket’ addressed:

In the original data, age was not presented in the same manner as in the description. Instead of producing quantifiable numbers for analysis, the values were grouped together as objects. To rectify this, each ‘bucket’ was replaced by a normal random output of values that fit the parameters of each ‘bucket.’

‘Unknown’s:

Later in the exploration of our data, we discovered Nulls that were disguised as ‘Unknown’s. To address this, we used hotkey encoding to deal with the unknown values. Depending on the contents of the column however, we would elect to drop the whole column if we did not find the contents of it relevant.

Data Preprocessing:

Encoding Categorical Data:

The data had certain categorical values. These values were changed into integer format, which will later be used for the model. Target Encoding, Label Encoding, One-Hot Encoding and Smoothing were observed.

Key Errors fixed:

Small key errors fixed, allowing for clean, correct encoding.

Exploratory Analysis:

Correlation:

Before building our model, we prepared several correlation maps to investigate what would have the greatest influence on our model for persistence. We found several higher level influencing factors but that reach at most a correlation 0.35.

Distribution charts:

A distribution chart for the age variable was prepared, displaying the mean and median age for the patients. This helps to visualize the data.

Box plot:

A box plot was created to look at any outliers. Since nearly all of the data is binary, only the boxplot for the number of “Risk_Counts” was observed. The boxplot proved that the “Risk_Counts” column had outliers.

Pie Chart:

Pie charts of Race, Ethnicity, Region and Gender were made to show the percentages and relationships between the different variables.

Bar plots:

Since much of our data was binary, we presented a majority of it as bar plots with percentages. This also was another great way to confirm that all the data was used and accounted for.

Model Training:

Flask using Logistic Regression:

Using Logistic Regression, the model was deployed on Flask and pushed to a server such as pythonanywhere.com.

Ensemble Models:

Different types of ensemble models were trained and observed. These include Decision Tree Classifier, Random Forest Classifier, Ada Boosting, Bootstrap Aggregation, and Voting Tree Classifier

Stacking Model:

Logistic Regression, Decision Tree Classifier, and Random Forest Classifier were “stacked” to yield a better result.

Flask Application:

Below is a visual guide of how we implemented our flask application. We used pythonanywhere.com as our server.

Pythonanywhere dashboard starting point:

The screenshot shows the Pythonanywhere dashboard interface. At the top, there's a header with the Pythonanywhere logo, 'by ANACONDA', and navigation links for Dashboard, Consoles, Files, Web, Tasks, and Databases. A welcome message 'Welcome, [liviafoster](#)' is on the right, along with an 'Upgrade Account' button. Below the header, there are four main sections: 'Recent Consoles', 'Recent Files', 'Recent Notebooks', and 'All Web apps'. The 'Recent Consoles' section shows a single entry: 'Bash console 27543864'. The 'Recent Files' section shows entries like '/home/liviafoster/mysite/FinalProject.py' and '/var/www/liviafoster_pythonanywhere_com_wsgi...'. The 'Recent Notebooks' section has a note: 'Your account does not support Jupyter Notebooks. [Upgrade your account](#) to get access!'. The 'All Web apps' section shows a single entry: 'liviafoster.pythonanywhere.com'. There are also buttons for 'View all' in the consoles section and 'Open Web tab' in the web apps section.

PythonAnywhere server set up:

 pythonanywhere
by ANACONDA

liviafoster.pythonanywhere.com Configuration for [liviafoster.pythonanywhere.com](#)

[Add a new web app](#)

Reload: [Reload liviafoster.pythonanywhere.com](#)

Best before date:

We're happy to host your free website – and keep it free – for as long as you want to keep it running, but you'll need to log in at least once every three months and click the "Run until 3 months from today" button below. We'll send you an email a week before the site is disabled so that you don't forget to do that. [See here for more details.](#)

This site will be disabled on **Friday 26 May 2023**

[Run until 3 months from today](#)

Paying users' sites stay up forever without any need to log in to keep them running.

Traffic:

How busy is your site?

This month (previous month)	29	(0)
Today (yesterday)	18	(6)
Hour (previous hour)	0	(0)

Want some more data? [Paying accounts get pretty charts :\)](#)

Files homepage:

 pythonanywhere
by ANACONDA
[/home/](#)  [liviafoster](#)

[Dashboard](#) [Consoles](#) **Files** [Web](#) [Tasks](#) [Databases](#)

[Open Bash console here](#) **23% full** – 116.8 MB of your 512.0 MB quota [More Info](#)

Directories

Enter new directory name [New directory](#)

[.cache/](#) 
[.local/](#) 
[.virtualenvs/](#) 
[mysite/](#) 

Files

Enter new file name, eg hello.py [New file](#)

 .bashrc	   2023-02-26 03:02 560 bytes
 .gitconfig	   2023-02-26 03:02 266 bytes
 .profile	   2023-02-26 03:02 79 bytes
 .pythonstartup.py	   2023-02-26 03:02 77 bytes
 .vimrc	   2023-02-26 03:02 4.6 KB
 README.txt	   2023-02-26 03:02 232 bytes

[Upload a file](#)
100MiB maximum size

'mysite' file:

 pythonanywhere
by ANACONDA
[/home/liviafoster/](#)  [mysite](#)

[Dashboard](#) [Consoles](#) **Files** [Web](#) [Tasks](#) [Databases](#)

[Open Bash console here](#) **23% full** – 116.8 MB of your 512.0 MB quota [More Info](#)

Directories

Enter new directory name [New directory](#)

[_pycache_/](#) 
[statics/](#) 
[templates/](#) 

Files

Enter new file name, eg hello.py [New file](#)

 FinalProject.py	   2023-02-26 03:34 57.2 KB
 Healthcare_dataset.csv	   2023-02-26 03:33 892.0 KB
 app_final.py	   2023-02-26 02:57 955 bytes

[Upload a file](#)
100MiB maximum size

'templates' file:

The screenshot shows the PythonAnywhere Files interface. At the top, there are navigation links: Dashboard, Consoles, **Files**, Web, Tasks, Databases. Below that, a status bar indicates "Open Bash console here" and "23% full - 116.8 MB of your 512.0 MB quota" with a "More Info" link. The main area is divided into "Directories" and "Files". In the "Directories" section, there is a text input field "Enter new directory name" and a "New directory" button. In the "Files" section, there is a text input field "Enter new file name, eg hello.py" and a "New file" button. A file named "project.html" is listed with details: download icon, edit icon, delete icon, timestamp "2023-02-24 23:30", size "21.9 KB". Below the file list is a "Upload a file" button and a note "100MB maximum size".

'statics/' file:

The screenshot shows the PythonAnywhere Files interface. At the top, there are navigation links: Dashboard, Consoles, **Files**, Web, Tasks, Databases. Below that, a status bar indicates "Open Bash console here" and "23% full - 116.8 MB of your 512.0 MB quota" with a "More Info" link. The main area is divided into "Directories" and "Files". In the "Directories" section, there is a text input field "Enter new directory name" and a "New directory" button. In the "Files" section, there is a text input field "Enter new file name, eg hello.py" and a "New file" button. A note "No files here" is displayed. Below the note is a "Upload a file" button and a note "100MB maximum size".

'css/' file:

The screenshot shows the PythonAnywhere Files interface. At the top, there are navigation links: Dashboard, Consoles, **Files**, Web, Tasks, Databases. Below that, a status bar indicates "Open Bash console here" and "23% full - 116.8 MB of your 512.0 MB quota" with a "More Info" link. The main area is divided into "Directories" and "Files". In the "Directories" section, there is a text input field "Enter new directory name" and a "New directory" button. In the "Files" section, there is a text input field "Enter new file name, eg hello.py" and a "New file" button. A file named "final_project.css" is listed with details: download icon, edit icon, delete icon, timestamp "2023-02-24 23:30", size "5.4 KB". Below the file list is a "Upload a file" button and a note "100MB maximum size".

Flask application when implemented:

Persistence of a Drug

General Patient Information:

1. Patient's Age:

2. Patient's Sex: Female

3. Patient's Race: Caucasian

4. Patient's Region: Midwest

5. Is the attending physician a specialist? No

6. Did the patient adhere to prescription instructions? No

7. Is the patient in the IDN system? No

Nontuberculous Mycobacteria (NTM) Section:

Nontuberculous Mycobacteria (NTM) Section:

8. Were glucocorticoids recorded before the patient was diagnosed with NTM? No

9. Were glucocorticoids recorded for the patient's prescription? No

10. How often did a patient receive a DEXA scan under prescription?

11. Did the patient have any fragility fractures before their first continuous therapy? No

12. Did the patient have any fragility fractures during their first continuous therapy? No

13. Did the patient have any risks before the treatments began? No

14. What was the patient's T-score before the beginning of prescription start date? No

15. What was the change in T-score during treatment? No

16. Were there any injectable drugs taken during treatment? No

Patient's Pre-existing conditions and/or Comorbidities:

Patient's Pre-existing conditions and/or Comorbidities:

17. Does the patient have indicators for Malignant Neoplasms? No
18. Does the patient have indicators for immunization? No
19. Did the patient have any complaints during their GENERAL exam without a suspected or reported diagnosis? No
20. Does the patient have any indications of a vitamin D deficiency? No
21. Does the patient have a history of any other joint disorders that are not elsewhere classified? No
22. Did the patient have any complaints during their SPECIALITY exam without a suspected or reported diagnosis? No
23. Is the patient currently under a regimen of a long-term drug therapy? No
24. Does the patient present with symptoms of dorsalgia? No
25. Does the patient have any personal history of other diseases and/or conditions? No
26. Does the patient have any history of other bone density and/or structure disorders? No

27. Does the patient have lipoprotein metabolism disorder and/or any other lipidemia disorders? No
28. Does the patient have osteoporosis without a current pathological fracture? No
29. Does the patient have a personal history of malignant neoplasm? No
30. Does the patient have a history of Gastro-esophageal reflux disease? No

Concomitant Treatments:

All drugs are recorded that occurred within the prior 365 days from first prescription date.

31. Prior to current treatment, were cholesterol and triglyceride regulating preparations made? No
32. Prior to current treatment, were narcotics prescribed? No
33. Prior to current treatment, were systemic corticosteroids prescribed? No
34. Prior to current treatment, were anti-depressants and mood stabilisers prescribed? No

Concomitant Treatments:

All drugs are recorded that occurred within the prior 365 days from first prescription date.

31. Prior to current treatment, were cholesterol and triglyceride regulating preparations made? No
32. Prior to current treatment, were narcotics prescribed? No
33. Prior to current treatment, were systemic corticosteroids prescribed? No
34. Prior to current treatment, were anti-depressants and mood stabilisers prescribed? No
35. Prior to current treatment, were fluoroquinolones prescribed? No
36. Prior to current treatment, were cephalosporins prescribed? No
37. Prior to current treatment, were macrolides and/or similar prescribed? No
38. Prior to current treatment, were any broad spectrum penicillins prescribed? No
39. Prior to current treatment, were any general anaesthetics prescribed? No
40. Prior to current treatment, did the patient receive any viral vaccines? No

Health Risks:

41. Is the patient at risk for type 1 insulin dependent diabetes? No
42. Is the patient at risk for osteogenesis imperfecta? No
43. Is the patient at risk for rheumatoid arthritis? No
44. Is the patient at risk for untreated chronic hyperthyroidism? No
45. Is the patient at risk for untreated chronic hypogonadism? No
46. Is the patient at risk for untreated early menopause? No
47. Does the patient have a family history of someone fracturing their hip? No
48. Does the patient smoke tobacco products? No
49. Does the patient have a history of chronic malnutrition or malabsorption? No
50. Does the patient have a history of chronic liver disease? No
51. Does the patient have a family history of osteoporosis? No
52. Does the patient have a history of low calcium intake? No
53. Does the patient have a history of vitamin D insufficiency? No
54. Does the patient have a history of poor health or frailty? No
55. Does the patient have a history of excessive thinness? No
56. Did the patient have a hysterectomy or oophorectomy performed? No

chronic hypogonadism? No

46. Is the patient at risk for untreated early menopause? No

47. Does the patient have a family history of someone fracturing their hip? No

48. Does the patient smoke tobacco products? No

49. Does the patient have a history of chronic malnutrition or malabsorption? No

50. Does the patient have a history of chronic liver disease? No

51. Does the patient have a family history of osteoporosis? No

52. Does the patient have a history of low calcium intake? No

53. Does the patient have a history of vitamin D insufficiency? No

54. Does the patient have a history of poor health or frailty? No

55. Does the patient have a history of excessive thinness? No

56. Did the patient have a hysterectomy oophorectomy performed? No

57. Does the patient have a history of estrogen deficiency? No

58. Does the patient have a history of immobilization? No

59. Does the patient have a history of recurring falls? No

Predict Persistence

Predicted output:

Note: it is in black script at the very bottom.

58. Does the patient have a history of immobilization? No
59. Does the patient have a history of recurring falls? No

Predict Persistence

Yes, the drug is persistent

Links:

General Repository:

- <https://github.com/LiviaNFoster/DataGlacierFinalProject.git>

Final Jupyter Notebook:

Please note that this file includes all of our work, where we each worked on it one at a time.

- <https://github.com/LiviaNFoster/DataGlacierFinalProject/blob/main/FinalProject.ipynb>

Presentation:

- <https://github.com/LiviaNFoster/DataGlacierFinalProject/blob/main/BusinessPresentation.pdf>

Flask Application:

- liviafoster.pythonanywhere.com