

INTRODUCTION TO DATA SCIENCE (LABS)

13/10/2020

WFAiS UJ, Informatyka Stosowana
I stopień studiów

<http://th-www.if.uj.edu.pl/~erichter/dydaktyka/Dydaktyka2020/DataScience-2020/index.html>

Ready, but I can still have few ideas to make course more interesting.

2

Introduction to Data Science (Wprowadzenie do analityki danych)

Wydział Fizyki, Astronomii i Informatyki Stosowanej,
Uniwersytet Jagielloński w Krakowie

Rok akademicki 2020/2021

Office hours: Tuesday 13:00 - 14:00; office G-0-10;
COVID-19: Please use MITKama system or email if applicable.

Lectures & Assignments:

Week	Lecture slides	Lab slides	Python scripts with assignments	Datasets	Tutorials
13-10-2020	Introduction Data exploration	Introduction lab	assignment-0-mathon assignment-0-mumpy assignment-0-matplotlib assignment-0-mandas assignment-1	the_browse_data_explain info on kaggle.com	How to Start on abstract numpy.pdf on abstract matplotlib.pdf on abstract pandas.pdf on abstract pandas.pdf
20-10-2020	Regression-Primer Regression-Advanced-I		assignment-2		sklearn scikit-learn-machine-learning-book python-cho-ai-journey python-cho-ai-journey
27-10-2020	Regression-Advanced-II		assignment-3		numpy tutorial numpy tutorial
3-11-2020	Regression-Advanced-III Classification-Primer				
10-11-2020	Classification-Advanced-I		assignment-4	dataamazon_baby.csv.zip	sklearn LogisticRegression
17-11-2020	Classification-Advanced-II				
24-11-2020	Classification-Advanced-III				
1-12-2020	Clustering-Primer		assignment-5	dataamazon_nikki.csv.zip	
8-12-2020	Clustering&Retrieval-Advanced-I				
15-12-2020	Clustering&Retrieval-Advanced-II Clustering&Retrieval-Advanced-III Recommending-System-Primer		its time to start developing your personal Data Science project: Select one from the list below or create your own: FBI_ThreatIntelligence FBI_FraudInvestigation FBI_Terrorism FBI_Visualization FBI_CyberSecurity FBI_CyberSecurity FBI_CyberSecurity FBI_CyberSecurity	reading club data science data-science-at-stanford-ed data-science data-science data-science data-science	
13-01-2021	Modeling_simulation_tests Causal methods		assignment-6		
18-01-2021	Covariational Inference		assignment-7		Sebastian Raschka A Fan at Stanford Uni.
26-01-2021	Multivariate Analysis and Artificial Neural Network				

Lectures are based on the materials from Coursera:

Dr. Mina Cetingöz-Ründel - "Data Analysis and Statistical Inference"
Dr. Guestrin and Dr. El Fero - "Machine Learning Specialization"
Foundations link
Regression link
Classification link
Clustering and Retrieval link

Related interactive material from Coursera:

Dr. Peng, J. Lesk and B. Caffo - "Exploratory Data Analysis"
J. Leskovec, A. Rajaraman and J. Ullman - "Mining Massive Datasets"
B. Caffo, H. D. Peng and J. Lesk - "Regression Models"

Data Science applications in physics:

B. Nachman - "Advanced Machine Learning for Classification, Regression, and Generation in Jet Physics"

4. Skop

ML applications in Q&A:

ML techniques in HEP: Workshop, Berkeley Laboratory, 11 - 13 December 2018

<https://indico.fnal.gov/event/4282/>

Collection of datasets:

<https://openml.org/dataset>

<https://www.kaggle.com/datasets>

<https://www.kaggle.com/datasets>

<https://www.kaggle.com/datasets>

Useful links:

<https://turi.com/download/install-graphlab-creates-svc-coursera.html>

<https://turi.com/download/academic.html>

<https://github.com/turi-ods/Frame>

Clustering:

https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVecorizer.html

https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html

Boosting:

https://turi.com/learn/userguide/userguide-learnings/boosted_trees_classification.html

<https://homes.cs.washington.edu/~tqchen/pdf/BoostedTrees.pdf>

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingClassifier.html>

[1] <https://class.coursera.org/statistics>

[2] <http://www.openintro.org/stat/ta-tbook.php>

[3] <https://class.coursera.org/e-data-008>

[4] <https://class.coursera.org/mmde>

[5] <https://ocw.mit.edu/>

[6] <http://www.cmu.edu/~swm/tutorials.html>

Additional materials

<http://www.stat.cmu.edu/~cshalizi/ADAFaEPoV/ADAFaEPoV.pdf>

<http://statweb.stanford.edu/~tibs/ElemStatLearn/>

http://en.wikipedia.org/wiki/Determining_the_number_of_clusters_in_a_data_set

<http://www.youtube.com/watch?v=wQhVWUcXM0A>

Link to lectures given in 2017

<http://th-www.if.uj.edu.pl/~erichter/dydaktyka/Dydaktyka2017/AiSAD-2017/index.html>

Link to lectures given in 2014

<http://th-www.if.uj.edu.pl/~erichter/dydaktyka/Dydaktyka2014/AiSAD-2014/index.html>

Ostatnia modyfikacja: 7 October 2020

Wydział Fizyki: Wb

Witaj!

LABS

3

- **Normal times: run in person**
 - ▣ **Time for you to write your code and (for me) to discuss with each student her/his progress with assignments.**
- **COVID-19 times: run on-line**
 - ▣ **We will go through content of assignments, present (mostly you) results of analyses and observations, have short oral presentations (please be active)**
 - ▣ **Occasion to share with everybody problems, exchange snippets of the code or interesting observations.**

Assignments, Project, Short Presentations

4

This is not a course of programming, but you will be expected to write programs.

- ▣ **Baseline is python + anaconda libraries.**
- ▣ **You can use also R or other Data Science specific programming language/library**

I will not be teaching you programming or helping to debug your code, you are on your own ...

For labs you will be graded with:

- ▣ **completed assignments: 7 + optional - max 72 scores**
- ▣ **personalised project - max 25 scores**
- ▣ **short topical presentations - max 25 scores**

Graded will be not (necessarily) quality of the code, but maturity of how you analyse and interpret the data.

To pass the course you need to collect at least **65 scores.**

Assignments, Project, Short Presentations

5

PEGAZ system:

This system we will use to collect your assignments/projects/short presentations

- ▣ I will be sending you back comments
- ▣ **You will see your grades there**

Please don't use email to send me your scripts!

MSTeams:

This system we will use for on-line classes

- ▣ you can use it also for communication among yourself, eg. setting up chats/meetings within a team
- ▣ for communication with me preferably use emails

PEGAZ system:

use to submit your assignments/projects/presentations.

6

The screenshot displays the PEGAZ system interface with a list of course activities. The interface is organized into sections, each with a plus sign icon and a title. Each section contains a list of activities, each with a plus sign icon, a document icon, and a title with a deadline and a pencil icon for editing. To the right of each activity list, there are 'Modyfikuj' (Modify) buttons with a dropdown arrow and a user icon. At the bottom of each section, there is a blue button with a plus sign and the text 'Dodaj aktywność lub zasób' (Add activity or resource).

- Ogłoszenia** (Announcements)
- Forum towarzyskie** (Social forum)
- Data exploration**
 - Assignment-1: due 18.10.2020, final deadline 25.10.2020
 - Assignment-1-opt: due 25.10.2020, final deadline 1.11.2020
 - Assignment-0-opt: to complete during the course
- Machine Learning**
 - Assignment-2: due 1.11.2020, final deadline 8.11.2020
 - Assignment-3: due 15.11.2020, final deadline 22.11.2020
 - Assignment-4: due 29.11.2020, final deadline 6.12.2020
 - Assignment-5: due 13.12.2020, final deadline 20.12.2020
- Personalised project and presentations**
 - Project: due 10.01.2021, final deadline 17.01.2021
 - Short topical presentations

PEGAZ system:

use to submit your assignments/projects/presentations.

7

The screenshot displays a user interface for the PEGAZ system. At the top, there is a header with a plus icon and the text "Problem modeling, MC methods, Statistical Inference". To the right of this header is a "Modyfikuj" button with a dropdown arrow. Below the header, there are two entries, each with a plus icon, a document icon, and a link: "Assignment-6: due 17.01.2021, final deadline 24.01.2021" and "Assignment-7: due 24.01.2021, final deadline 31.01.2021". To the right of each entry is a "Modyfikuj" button with a dropdown arrow and a user icon. At the bottom center, there is a plus icon and the text "Dodaj temat". On the right side, there is a blue button with a plus icon and the text "Dodaj aktywność lub zasób".