

Data X

Introduction Data, Signals, and Systems

Ikhlaq Sidhu

Chief Scientist & Founding Director, Sutardja Center for Entrepreneurship & Technology
IEOR Emerging Area Professor Award, UC Berkeley

Welcome to Applied Data Science with Venture Applications

Sutardja Center for Entrepreneurship & Technology
Industrial Engineering End Operations Research Department

Fall 2018 course information:

Undergrad: [INDENG 135](#) - LEC 001, 28108

Grad Section: [INDENG 290](#) – Lec 2, 17093

Location: TuTh 12:30 pm – 2:00 pm | [Evans 10](#)

Prerequisite: Students should have a working knowledge of Python, have completed a fundamental probability / statistics course, as well as have a basic understanding Linear Algebra.



Welcome to Applied Data Science with Venture Applications

Sutardja Center for Entrepreneurship & Technology
Industrial Engineering End Operations Research Department

- Teaching Team
 - Ikhlaq Sidhu, SCET Faculty Director / Chief Scientist
 - Alexander Fred-Ojala, afo@berkeley.edu (instructor)
 - Sumayah Rahman, rahmans@berkeley.edu (GSI)
 - Tanya Piplani, tanyapiplani@berkeley.edu (GSI)
- Data-X Lab Students also contribute and coordinate: Vanessa Salas, and others..



Course Philosophy

Data-X



Make the Tools

Use State-of-the-Art
Open Source Tools

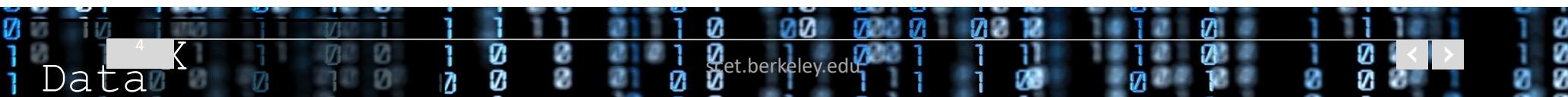
Architect
the System

Sell, market, and
pitch the product

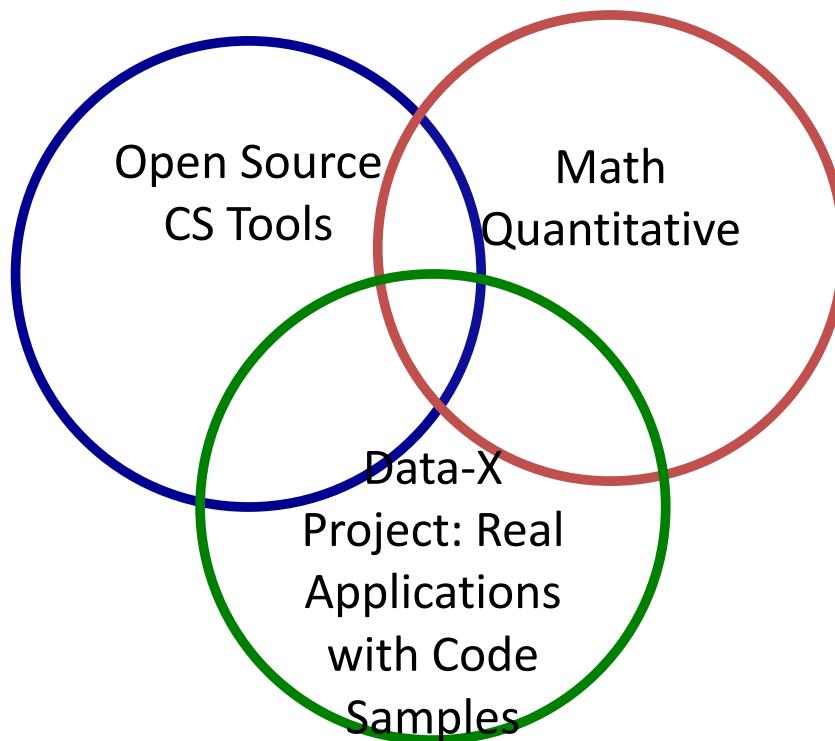
Most CS / Math

Data-X

Business Topics



What is in this course



Holistic Perspective: Industry, Social Applications, Customer Driven

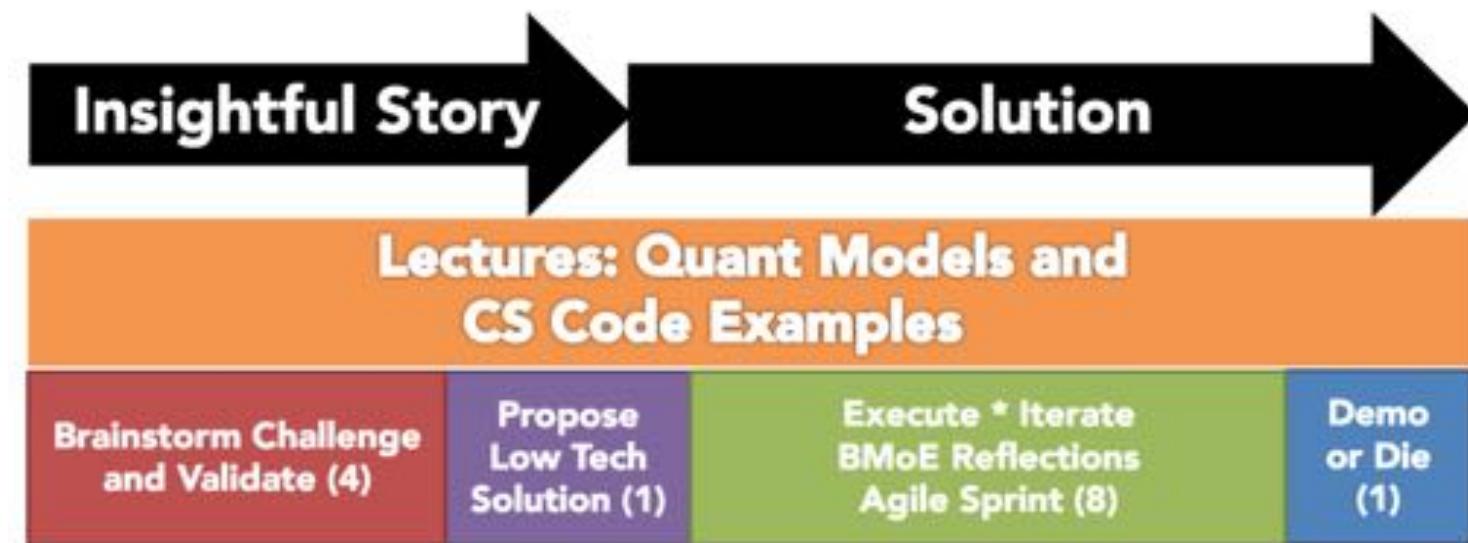


What is taught in the class?

- The ML stack most commonly used in creating ML/AI/Data applications
- Application and systems viewpoint of data and ML
- Implementation, architecture, and relevant processes to build real systems
- Connection with relevant mathematical, statistical foundations (optimization, entropy, correlation, LTI, prediction, classification)
- Practical insight into advanced techniques and tools: (eg. CNNs, NLP, scraping, recurrent networks, etc.)
- System modeling for data applications
- Application talks: Recommender systems, Blockchain, Spark etc.



Course Overview



Open-ended, real-world project: Typically 5 students, with available advisor network

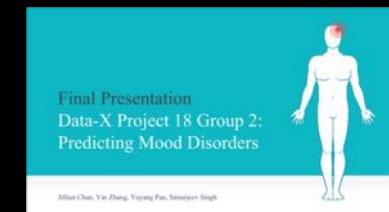
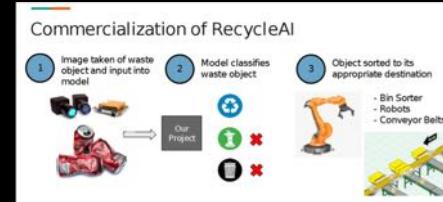


Data-X Project Examples

- Detection of fake news
- Prediction of long-term energy prices
- Automatic recycling through image recognition
- AI for crime detection, traffic guidance, medical diagnostics, etc.
- A version of Zillow that is recalculated with the effects of AirBnB income
- Signal processing and pattern analysis to improve earthquake warning systems
- Early Autism Detection
- Secure Health Records stored on a Blockchain

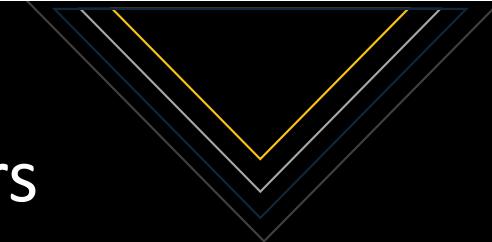
find many, many more at:

Data-X www.data-x.blog/projects



Data-X By the numbers

- 350 alumni students
- 50% avg enrollment increase / semester
- 80+ great projects completed
- 8+ published research papers
- 100+ industry experts in network
- 20+ students got employment as data scientists only because they took Data-X



Many Opportunities even After the course

Amazing testimonials:

I think this class is so awesome because it teaches the tools and concepts that are most commonly used in workplace teams that are involved with data science/machine learning.

Teaching
Data-Lab Projects

Develop New Data-X Materials

Data

Agenda on Day 1

Today:

- Course Introduction High Level Overview of Data, and Data-X Project (1:20 min)

By this week:

- Get your Notebook/development environment working
- Develop initial project ideas
- HW assigned by email/bcourse

Key Dates:

Class Starts today

Final projects: During end of class, possibly reading week

Top 2 projects will showcase at Collider Cup event



Most Resources Are Available at data-x.blog

1. Go to Data-X.blog
 - Syllabus
 - Instructions for SW Install
 - Link to GitHub with Cookbook Code Samples and Slides
2. Download Instructions to Install Python 3.x Anaconda Environment. For now you only need Anaconda, don't worry about other packages that are not already included.
3. Be able to create your own Jupyter notebook
4. Self-Review Python references as needed. See Ref CS01 and as needed BIDS Python Bootcamp.



SYLLABUS
[Edit](#)

Applied Data Science with Venture Applications
IEOR 135/290-002

Instructor: Ikhlaq Sidhu
Department of Industrial Engineering & Operations Research

3 Units, Lecture and Lab



Project Types

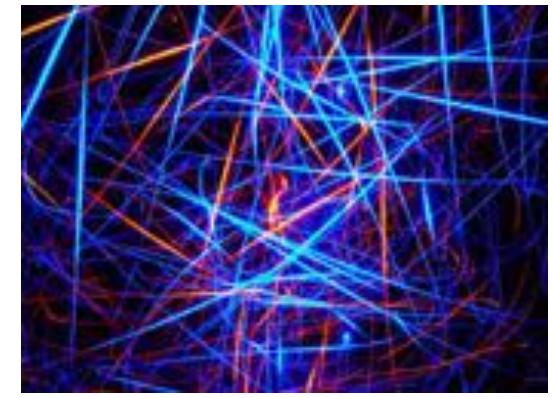


Business or Consumer
Use Case



Social Impact

(or improve part of a data pipeline
or work towards a research result)



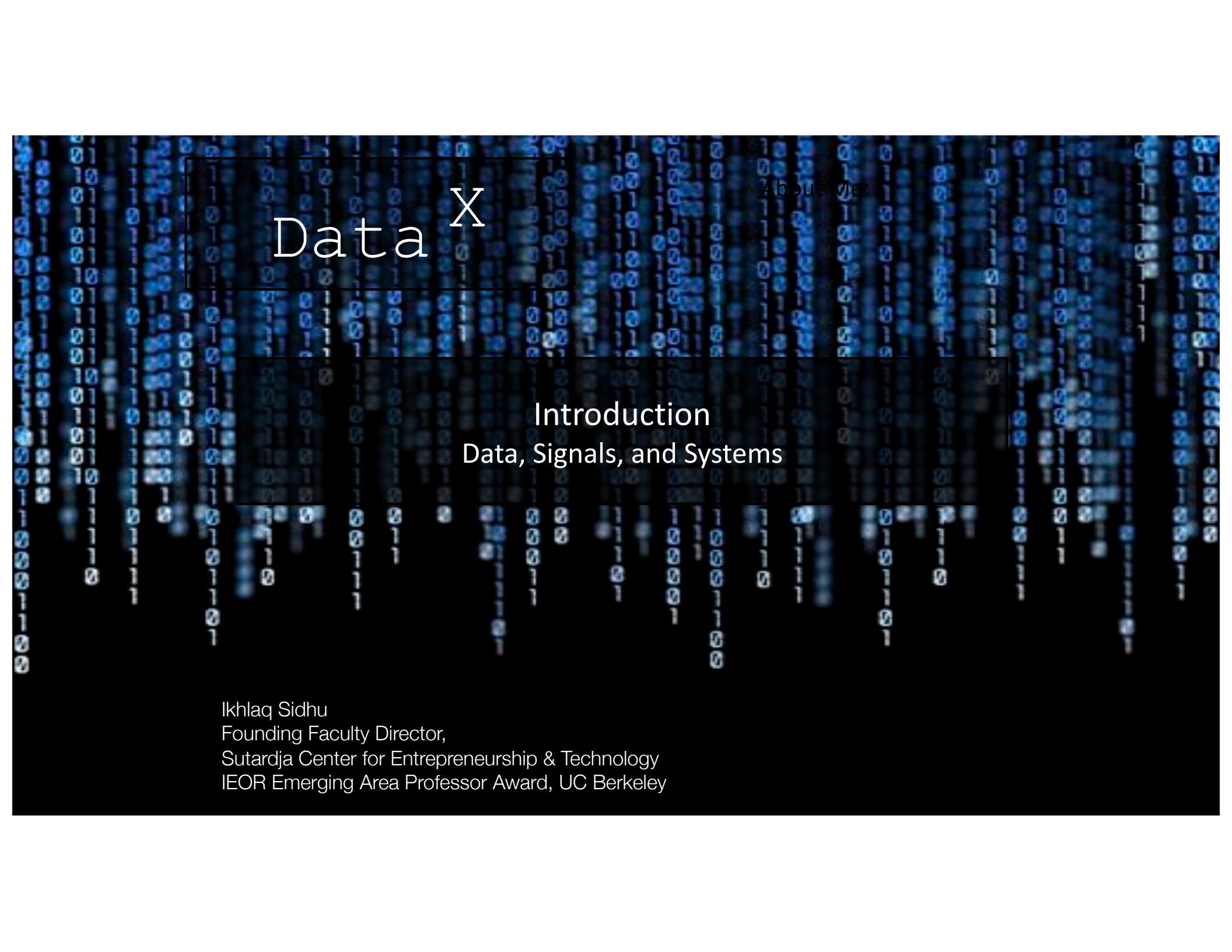
Its Just Cool



Homework For Week 1

- HW Part 1: For Your Project – By Next week
- Come up with 3 ideas for class projects in 1-3 sentences.
- A systems or application you will build
- **Communicate:** WHO the project is for, WHAT will it do, WHY this is needed/valuable.
- ---
- Homework Part II
- Python-based review notebook (Breakout and Homework, BKHW). To sent by email.





Data X

Introduction Data, Signals, and Systems

Ikhlaq Sidhu
Founding Faculty Director,
Sutardja Center for Entrepreneurship & Technology
IEOR Emerging Area Professor Award, UC Berkeley

An Overview of Data and AI Applications

Data X

Basic Concept of Working with Data



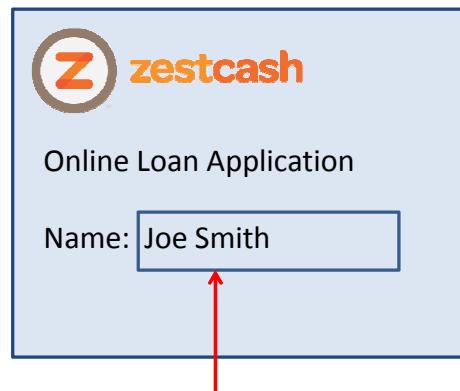
- Data Wrangling
- In Production



Example: Data and Information is a competitive advantage

Real-life Example: ZestCash

- “All data is credit data”



The data says: greater credit risk!

The data says: lesser credit risk!



- Service provider of Gambling and Casinos
- Entry Card
- Pain points
- Intervention



Reference: Supercrunchers

Why: More Simply

Customer
Insight/
Engagement

Operations:
Reliable &
Predictable

Security &
Fraud



Compliance **360°**

Financial Firms

Network Security

Implementation: SW Tools / Stack

Data X

The Most Common Open Source Tools: AI/ML Stack

Start with Python as an interface
Jupyter Notebooks for prototyping

- Python: The interface
- NumPy, SciPy: Working with Arrays
- Pandas: Working in Tables, SQL to Pandas
- Sklearn: ML
- Matplotlib: Visualizing Data
- TensorFlow, Keras: Neural Networks
- SQL to Pandas
- NLP / NLTK: Natural Language
- Spark: For large data sets (GB, TB+)



<https://www.youtube.com/watch?v=Q0jGAZAdZqM>

<https://conda.io/docs/user-guide/install/download.html>



Where Does Data Come From?

Data X

Where Does Data Come From?

Real-life Example: ZestCash

- All data is credit data"



Public datasets on AWS

To enable more innovation, AWS hosts a selection of datasets that anyone can access for free. Data in our public datasets is available for rapid access to our flexible and low-cost computing resources.



Web Scraping



Extract data from any website

Your Own Web Site

Public Data Sets
Stock market, etc.

IOT/Sensors

Other Web Sites

Data X

Web Scraping



Extract data from any website

Web Scraping

```
1 from bs4 import BeautifulSoup
2 import requests
3 page_link = "https://www.website_to_crawl.com"
4 # fetch the content from url
5 page_response = requests.get(page_link, timeout=5)
6 # parse html
7 page_content = BeautifulSoup(page_response.content, "html.parser")
8
9 # extract all html elements where price is stored
10 prices = page_content.find_all(class_="main_price")
11 # prices has a form:
12 # <div class="main_price">Price: $66.68</div>
13 # <div class="main_price">Price: $56.68</div>
14
15 # you can also access the main_price class by specifying the tag of the class
16 prices = page_content.find_all('div', attrs={'class': 'main_price'})
```

<https://github.com/ikhlaqsidhu/data-x>

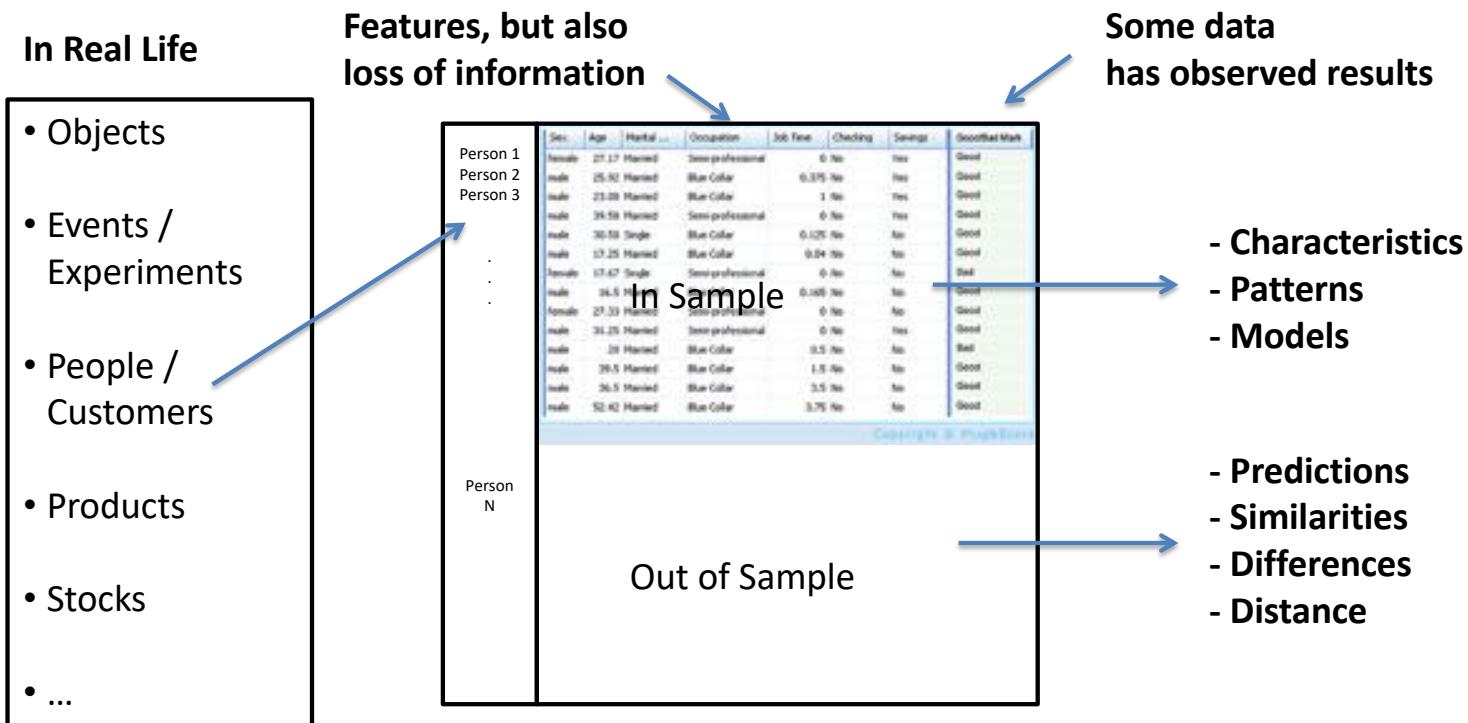
https://github.com/ikhlaqsidhu/data-x/tree/master/03-tools-webscraping-crawling_api_af0

Data X

Formatting Data

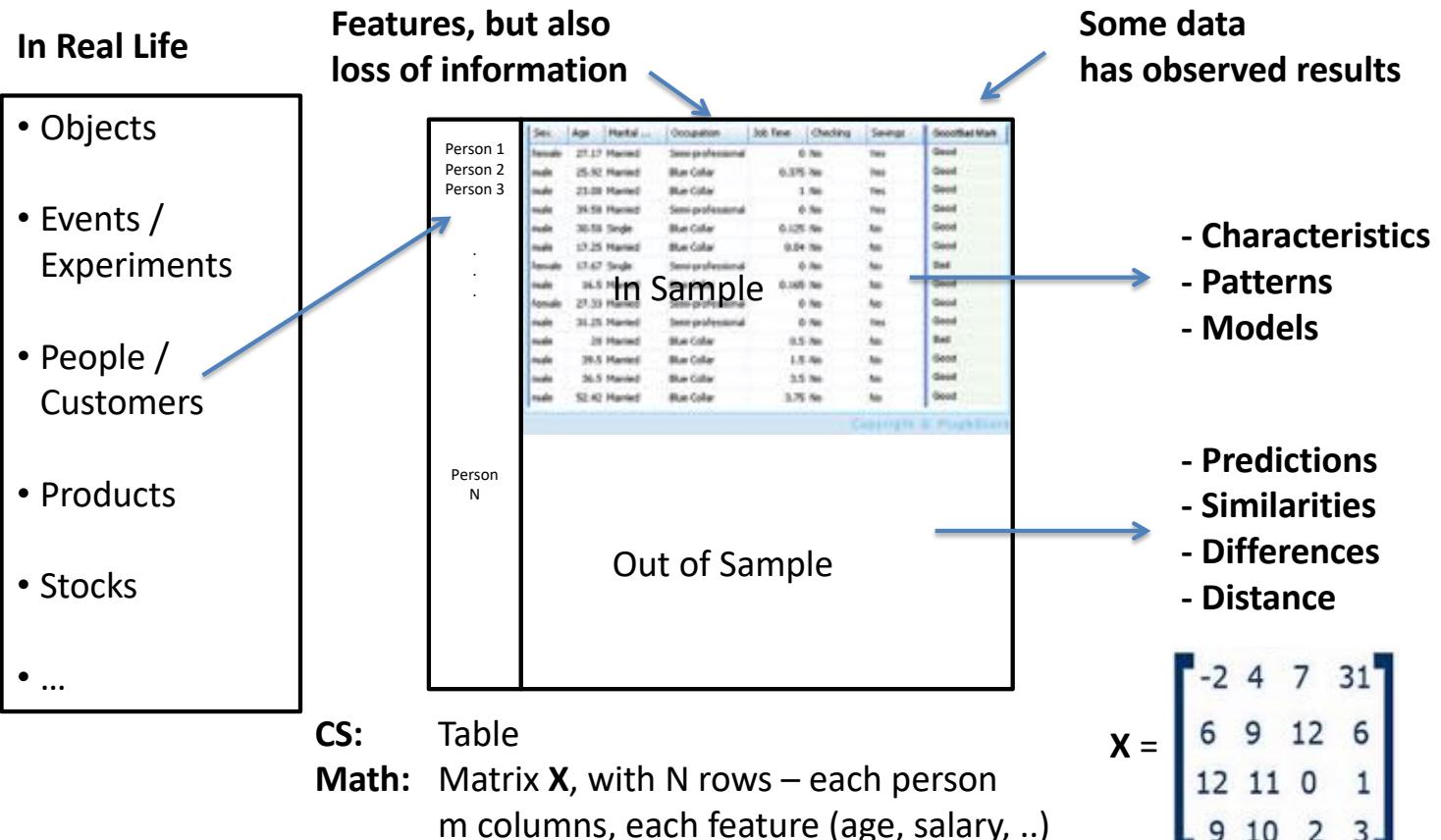
Data X

An ML High Level Framework



Data X

An ML High Level Framework



A Fundamental Idea: From Table to Score

$X =$

Cust	F1	F2	F3
A	4	2	2
B	4.5	1.5	3
C	3	3	5
D	1	2	2
E	3	1.5	5
F	3.5	3.5	1
..

$F(X)$

Cust	Credit Score
A	552
B	381
C	760
D	330
E	452
F	678
..	..



A Fundamental Idea: From Table to Score

X =

Cust	F1	F2	F3
A	4	2	2
B	4.5	1.5	3
C	3	3	5
D	1	2	2
E	3	1.5	5
F	3.5	3.5	1
..

$F(X)$

Cust	Credit Score
A	552
B	381
C	760
D	330
E	452
F	678
..	..

```
#Setting up for Supervised learning  
# First clean: use mapping +  
buckets
```

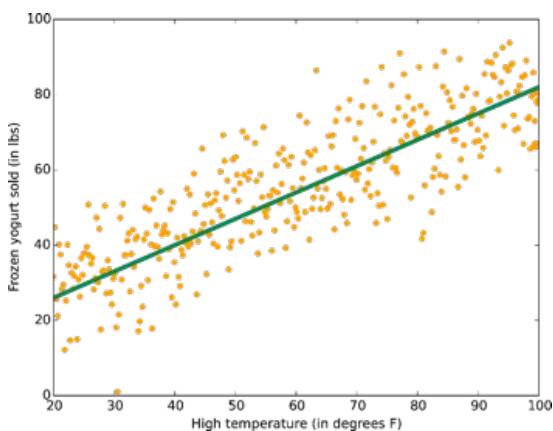
```
# X = matrix of data – e.g 1000 rows  
# Y = In sample responses
```

```
# Typically we want to split in to  
training data and test data
```

```
X_train = X[0:500]  
Y_train = Y[0:500]  
X_test = X[501:1000]  
Y_test = Y[501:1000]
```



Linear Regression Illustration



```
#Setting Linear Regression in sklearn  
from sklearn import linear_model  
  
model= linear_model.LinearRegression()  
model.fit(X_train, Y_train)  
  
Y_pred_train = model.predict(X_train)  
Y_pred_test = model.predict(X_test)  
  
# Compare Y_pred_test with Y_test for  
error.
```

Illustration Source: <https://docs.microsoft.com/en-us/azure/machine-learning/machine-learning-algorithm-choice>

X

Data

A Fundamental Idea: From Table to N- Dimensional Space

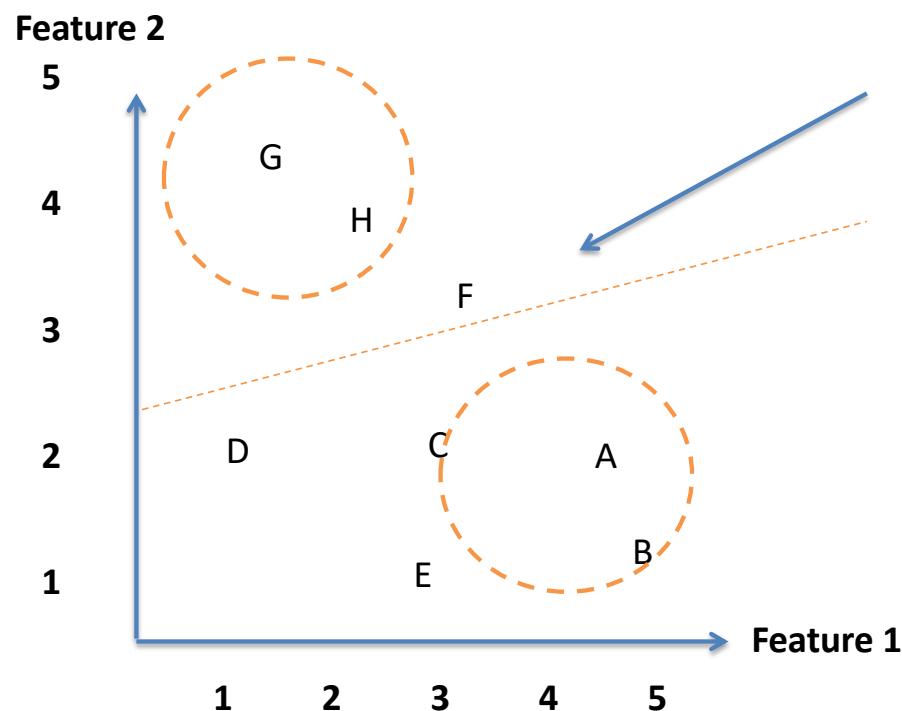
$X =$

Element	F1	F2	F3
A	4	2	2
B	4.5	1.5	3
C	3	3	5
D	1	2	2
E	3	1.5	5
F	3.5	3.5	1
..



Data X

Clustering to Classification

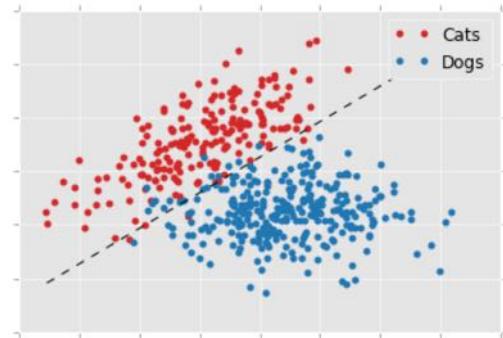


- Target customers?
 - Pictures of Cats and Dogs
 - Speech recognition
 - Recognize Letters: A, B, C..

Data X

Traditionally 2 Tasks: Classification & Predictive Scoring

Extracted Data
often in
Table
Format



Classification:
Cats and Dogs, Speech Recognition
Movie Recommendation



Scoring:

Credit Score: 830
Heat Index: 75

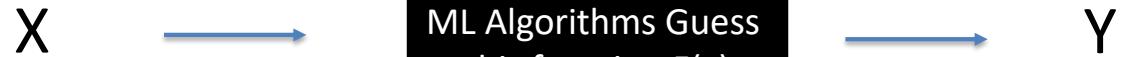
Movie Rating: PG-13
Movie Length: 120 minutes



The most famous
application has been
recommendation:
“which other user is
most like you”

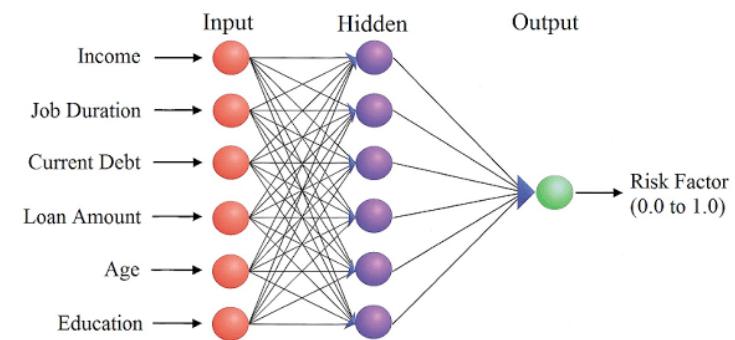
Data X

We have now switched
to Neural Networks as
Function Approximators



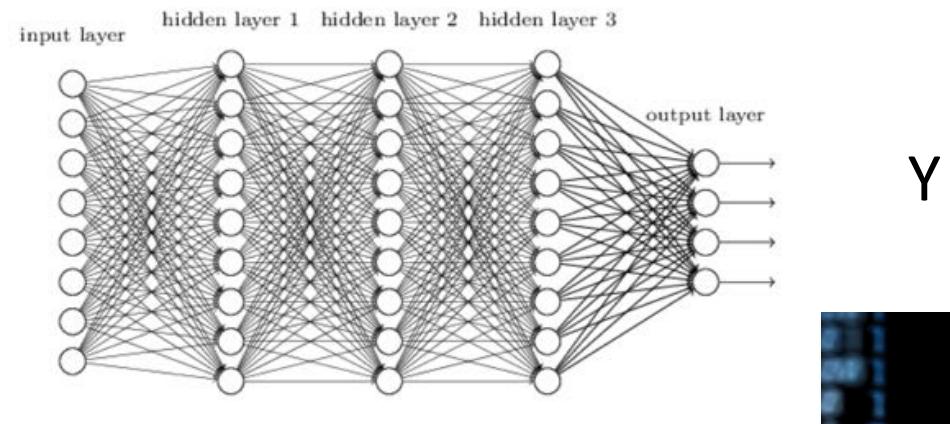
"Non-deep" feedforward
neural network

X



Deep neural network

X

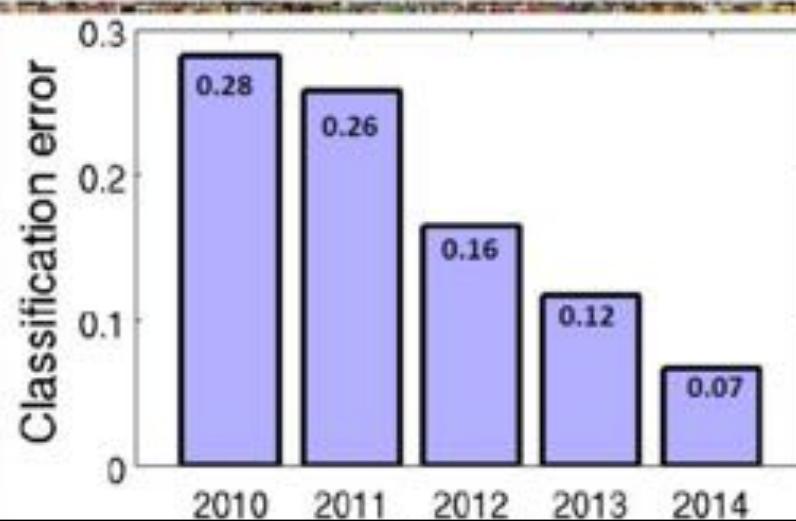


Data X

Y

IMAGENET Large Scale Visual Recognition Challenge

The Image Classification Challenge:
1,000 object classes
1,431,167 images



Neural net results are closest human results

Russakovsky et al. arXiv, 2014

Project Types

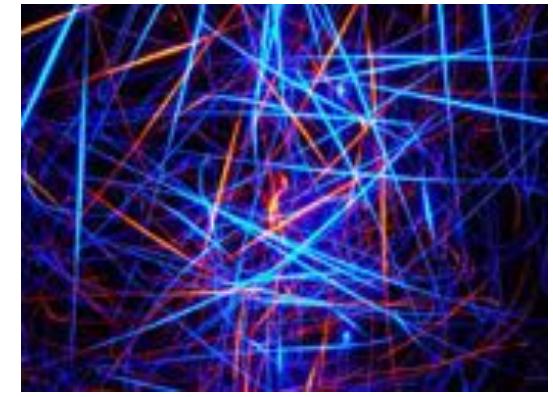


Business or Consumer
Use Case



Social Impact

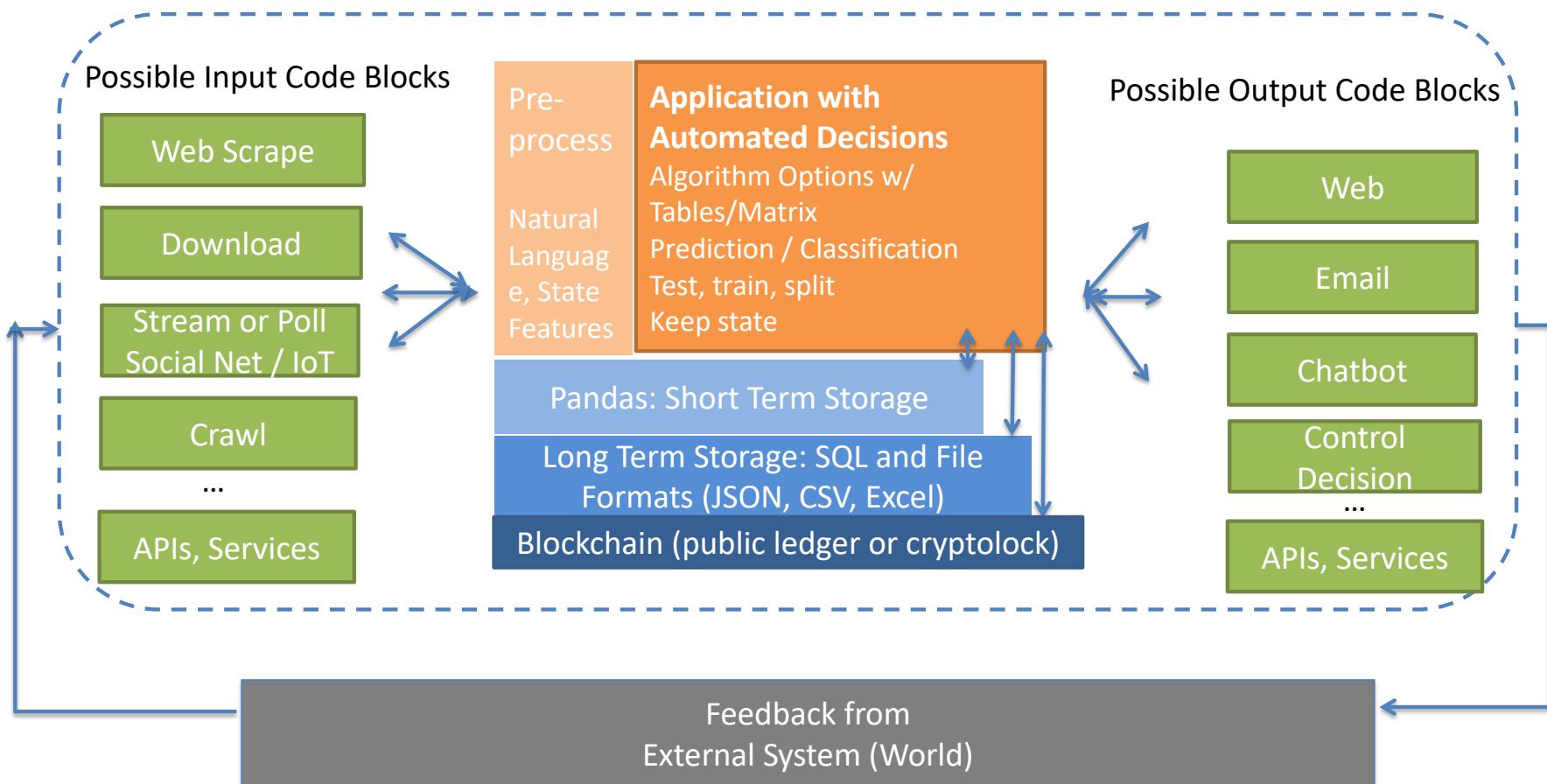
(or improve part of a data pipeline
or work towards a research result)



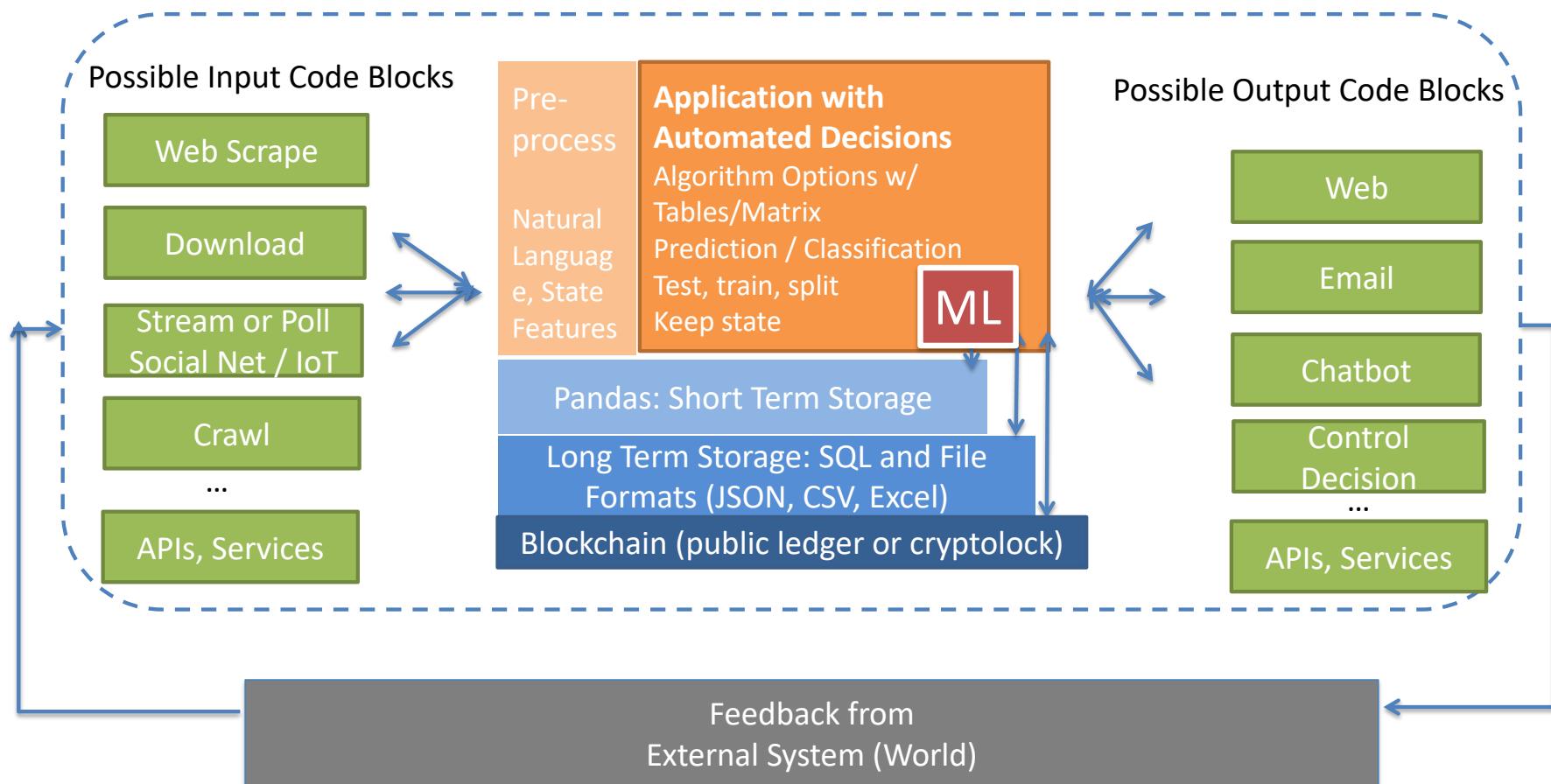
Its Just Cool



The Data-X System View



The Data-X System View: It's more than ML, it's also systems and models



Project Ideation

- Past Projects Concepts:
 - See the Advisor's Tab of data-x.blog
- Past Projects:
 - See the archive on the Posts page and on the Labs page of Data-x.blog
- Combine ideas or extend previous work
- You can also choose to build part of a system,
 - ie, just the part that automatically collects data by web scraping, or
 - just the part that makes a decision based on data already available

Home Resources Syllabus Posts Labs Advisors Contact

* DISCUSSION ADVISING: DISCUSSION AT BERKELEY

Project Concept Links:

- New Venture Success ([link](#))
- Concept: Blockchain based social currency to regulate social platform such as Twitter ([link](#))
- Concept: Personal Genome Hacking ([link](#))
- Concept: Holy Grail of Venture Capital ([link](#))
- AI Music Software development student cooperation opportunity ([link](#))
- Concept: Predicting future outcomes based on historical records ([link](#))
- Concept: US Power Plant project ([link](#))
- Concept: Multi-disciplinary data analysis of common psychological conditions ([link](#))
- Concept: Visualizing investment opportunities in touristic regions (Open Data for Greece 1.0) ([link](#))
- Concept: The University Bot ([link](#))
- Concept: Materials Recycling using Machine Learning
- Concept: Insights from Personal Photos
- Fuzzy Joins – A Modeling Discussion for Probabilistic Joins in Data Tables
- Concept: Faculty Research Matching with NLP and ML
- Concept: Inferred Information via Probabilistic Joins

Extended Mentor Network:

- Amir Najian, Geospatial Data Scientist at RMS – Geospatial Machine Learning and uncertainty modeling in Geocoder systems



Reminder: Homework For Week 1

- HW Part 1: For Your Project – By Next week
- Come up with 3 ideas for class projects in 1-3 sentences.
- A systems or application you will build
- **Communicate:** WHO the project is for, WHAT will it do, WHY this is needed/valuable.
- ---
- Homework Part II
- Python-based review notebook (Breakout and Homework, BKHW). To sent by email.



End of Section

Data^X