

## Welcome to Applied Data Science with Venture Applications

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

### **Course History:**

- Advanced project course for data science systems
- Concepts: theory, tools, project.
- IEOR 135 – Undergrad, IEOR 290 – Grad
- Currently impacted, but we are working on scale and offering every semester

### **Prerequisite:** Students should have:

- a working knowledge of Python
- completed a fundamental probability / statistics course
- 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

#### **Lead Faculty:**

Ikhlaq Sidhu  
Arash Nourian

#### **GSIs and Notebook Areas:**

Lillian Dong – Fundamental Data Science Tools  
Ishaan Malhi – Systems  
Zhi Li – Natural Language Processing and Special Topics  
Deirdre Quillen – Deep Learning

#### **Project Lead and Mentor :**

Nattaphol Vimolchalao



## 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
- Answer Questions

#### **Key Dates:**

Class Starts today

Final projects: During end of class during reading week

Top project will showcase at Collider Cup event



## Course Philosophy

Data-X



Make the Tools



Use State-of-the-Art  
Open Source Tools



Architect  
the System

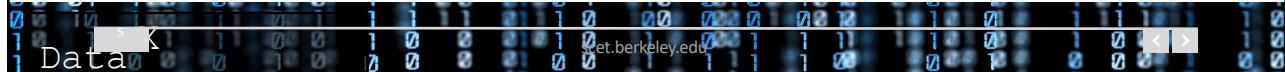


Sell, market, and  
pitch the product

Most CS / Math

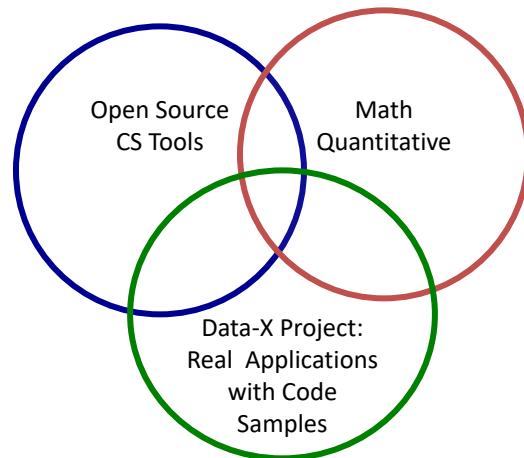
Data-X

Business Topics



## What is in this course

- Course Materials
- Applied Project
- Holistic Perspective:  
Industry, Social  
Applications,  
Customer Driven



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 examples: Recommender systems, Blockchain, Spark etc.



### Course Overview

**Insightful Story** → **Solution**

**Lectures: Quant Models and CS Code Examples**

Brainstorm Challenge and Validate (4)

Propose Low Tech Solution (1)

Execute \* Iterate  
BMoE Reflections Agile Sprint (8)

Demo or Die (1)

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: [www.data-x.blog/projects](http://www.data-x.blog/projects)

## Project Types

**Business or Consumer Use Case**

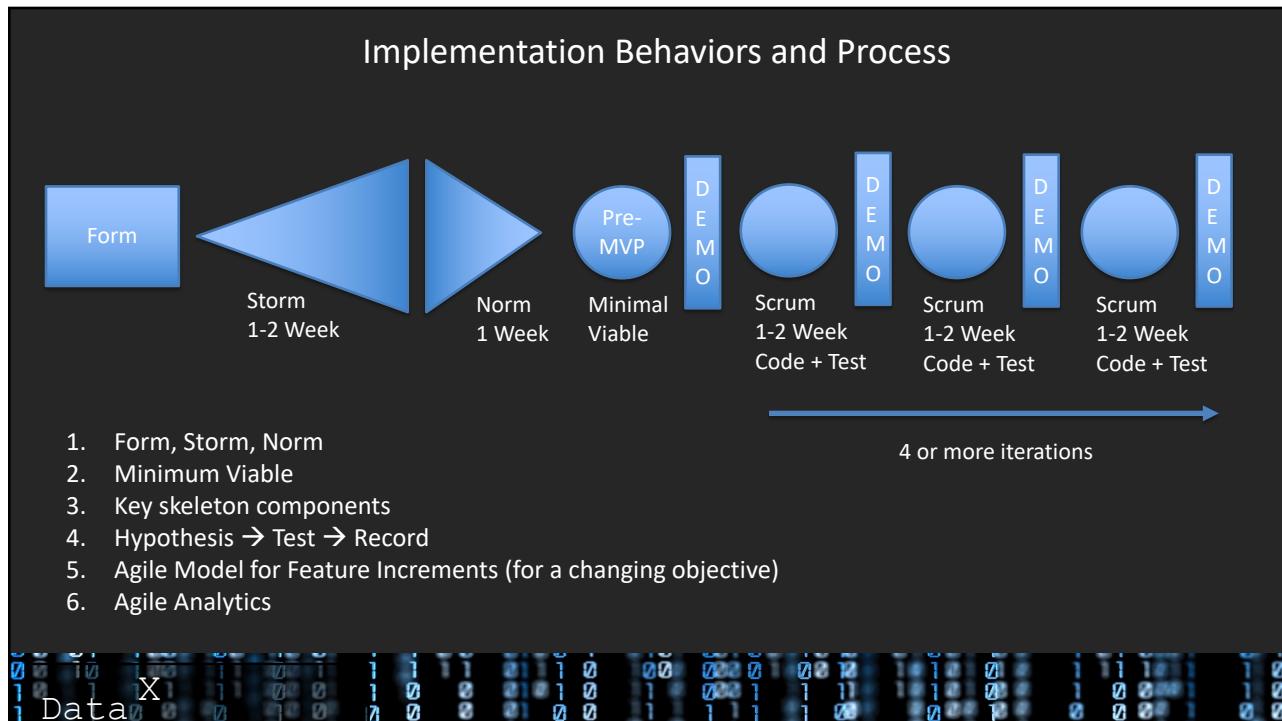
(or improve part of a data pipeline or work towards a research result)

**Social Impact**

**Its Just Cool**

**Project Categories:**

- Completely by student team
- Industry proposed projects
- Data-X library

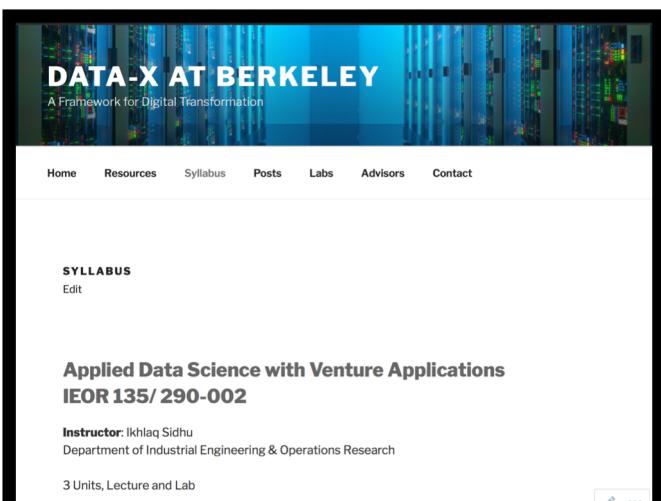


### Many Course Resources Are Already Available at [data-x.blog](http://data-x.blog)

For students and mentors

**See [data-x.blog](http://data-x.blog)**

- Lectures and Slides
- Code Samples
- Articles and Readings
- Projects
- Mentors



**DATA-X AT BERKELEY**  
A Framework for Digital Transformation

Home Resources Syllabus Posts Labs Advisors Contact

**SYLLABUS**  
Edit

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

Instructor: Ikhlilq Sidhu  
Department of Industrial Engineering & Operations Research

3 Units, Lecture and Lab

**Data X**

## Project Ideation

- 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

The screenshot shows a navigation bar with links: Home, Resources, Syllabus, Posts, Labs, Advisors, Contact. Below the navigation bar, there is a section titled "Project Concept Links:" which lists various projects such as "New Venture Success" and "Concept: Blockchain based social currency to regulate social platform such as Twitter". There is also a section titled "Extended Mentor Network:" featuring a single entry for Amir Najan.



## Homework For Week 1

Project Guidelines: <https://data-x.blog/project-guideline/>

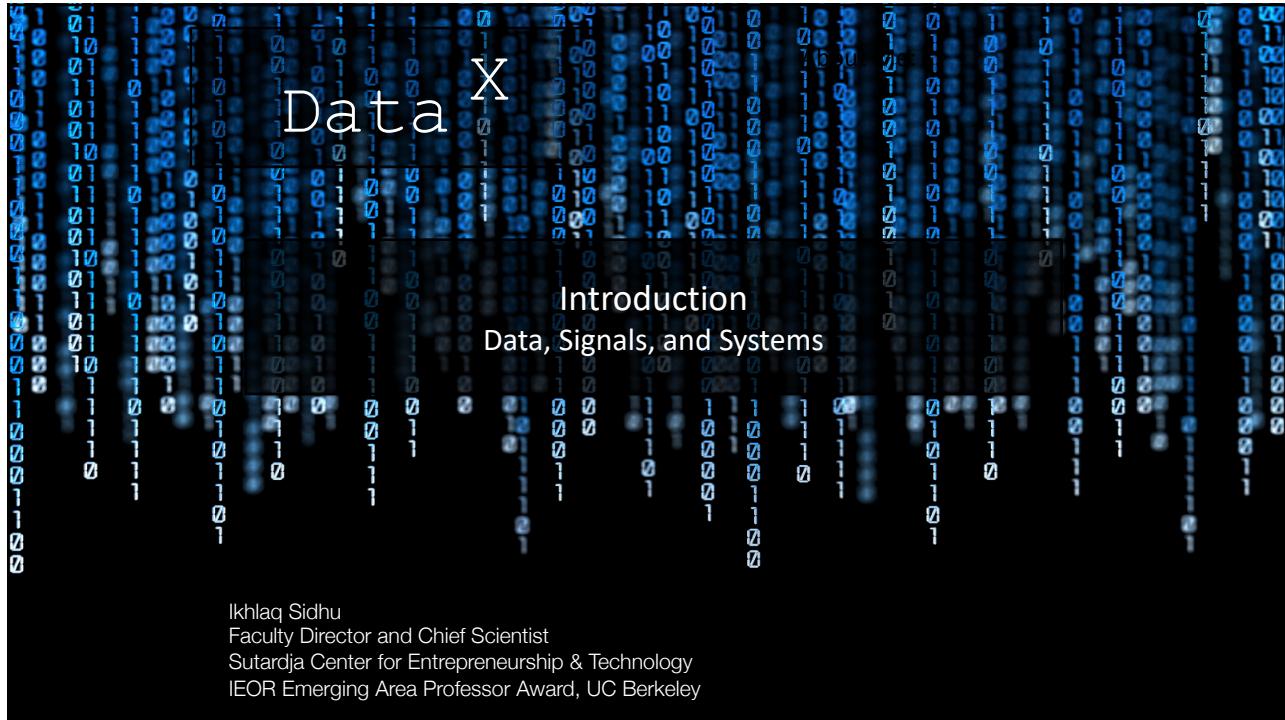
### By Next Week

- Project Assignment: write down 3 ideas a project, 1-3 sentences each. Due within one. Submit with a google form [here](#) and [shared spreadsheet](#) to that everyone can see the concepts submitted. (5 points)
  - A systems or application you will build
  - **Communicate:** WHO the project is for, WHAT will it do, WHY this is needed/valuable.
- Reading: Chapter 1 and Caviar case. Write  $\frac{1}{2}$  page to 1 page reflection about which concepts from the chapter and case might be most relevant for your project. You may also critique concepts. Due in one week. (5 points)

### Code and Theory HW:

- To be assigned by Arash Nourian





## An Overview of Data and AI Applications



## Basic Concept of Working with Data



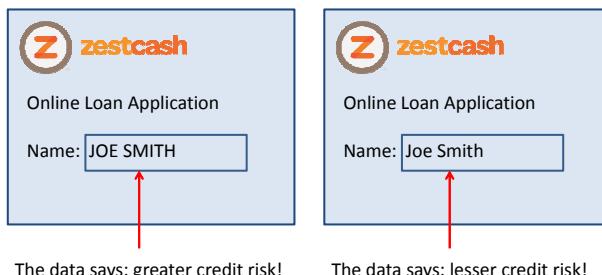
- Data Wrangling
- In Production

Data X1

### Example: Data and Information is a competitive advantage

#### Real-life Example: ZestCash

- “All data is credit data”



Reference: Shomit Ghose

Data X1



**Harrah's Casino: Knowing your customer**

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

**PLAY & WIN ▶**

Reference: Supercrunchers

Why: More Simply

Customer Insight/ Engagement	Operations: Reliable & Predictable	Security & Fraud
 TARGET	 	
		Financial Firms
		Network Security

**Data X**

## Implementation: SW Tools / Stack

0 0 1 0 1 1 0 1 1 1 0 0 0 0 0 1 1 0 1 0 1 1 1 1 0 1 1 1 1 1 1 1  
1 Data X 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
0 0

## 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>

0 0 1 0 1 1 0 1 1 1 0 0 0 0 0 1 1 0 1 0 1 1 1 1 0 1 1 1 1 1 1 1  
1 Data X 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
0 0

## Where Does Data Come From?



## Where Does Data Come From?

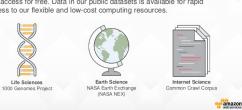
### Real-life Example: ZestCash

- "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



## Web Scraping

# Web Scraping



**Extract data from any website**

```

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'})

```

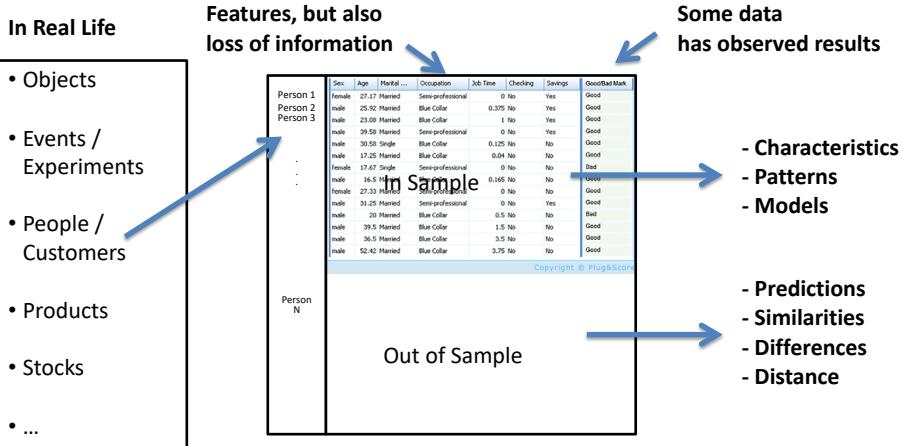
<https://github.com/ikhlaqsidhu/data-x>

[https://github.com/ikhlaqsidhu/data-x/tree/master/03-tools-webscraping-crawling\\_api\\_af0](https://github.com/ikhlaqsidhu/data-x/tree/master/03-tools-webscraping-crawling_api_af0)

## Formatting Data

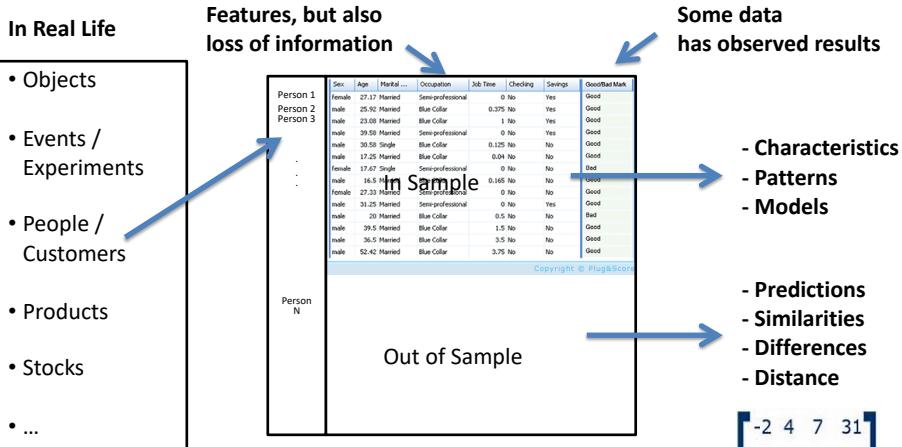


# An ML High Level Framework



0 0 1 1 0 0 1 1 1 0 0 0 0 1 1 0 0 0 0 1 1 1 1 1 1 1 1 1 1  
1 1 Data X 1 1 0 0 1 0 1 0 0 0 1 0 1 0 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1

# An ML High Level Framework



CS: Table

Math: Matrix  $X$ , with  $N$  rows – each person  
 $m$  columns, each feature (age, salary, ...)

$$X = \begin{bmatrix} -2 & 4 & 7 & 31 \\ 6 & 9 & 12 & 6 \\ 12 & 11 & 0 & 1 \\ 9 & 10 & 2 & 3 \end{bmatrix}$$

0 0 1 1 0 0 1 1 1 0 0 0 0 1 1 0 0 0 0 1 1 1 1 1 1 1 1 1 1  
1 1 Data X 1 1 0 0 1 0 1 0 0 0 1 0 1 0 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1

## 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
..	..	..	..

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

Data  $X$

## 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
..	..	..	..

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

Data  $X$

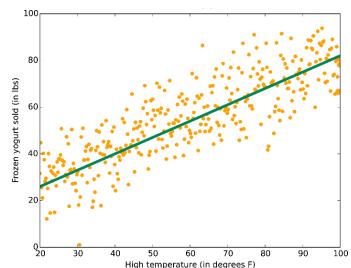
```
#Setting up for Supervised learning
# First clean: use mapping + buckets

# 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.
```

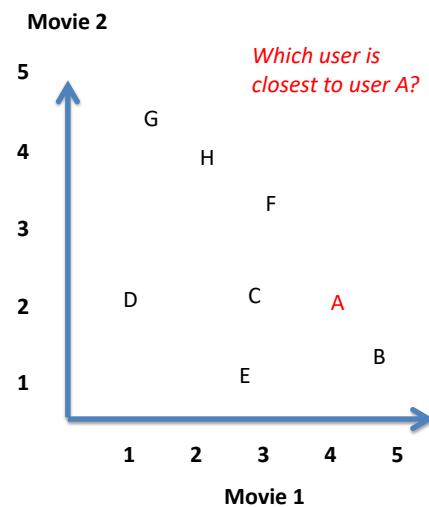
Illustration Source: <https://docs.microsoft.com/en-us/azure/machine-learning/machine-learning-algorithm-choice>

**Data**  $X_1$

## A Fundamental Idea: From Table to N- Dimensional Space

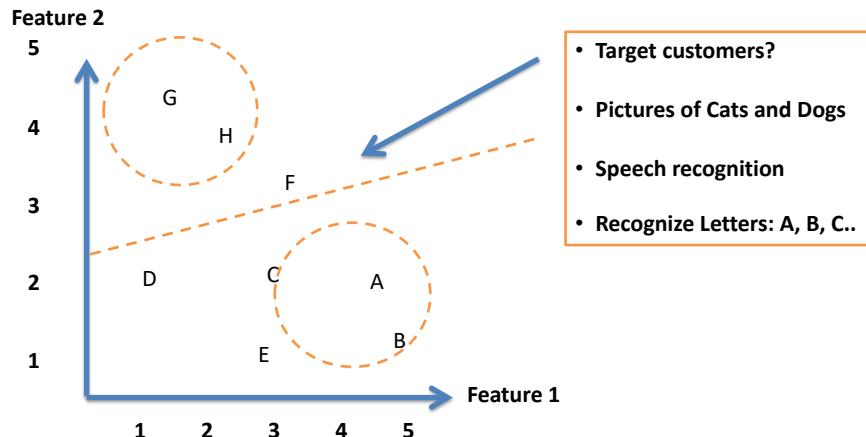
Element	F1	F2	F3
A	4	2	2
B	4.5	1.5	3
C	3	3	5
D	1	2	2
E	3	1.5	5
F	3.5	3.5	1
..	..	..	..

**X =**



**Data**  $X$

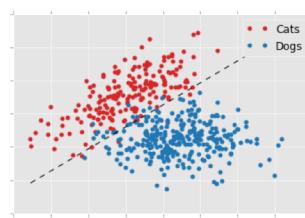
## Clustering to Classification



**Data  $X_1$**

## Traditionally 2 Tasks: Classification & Predictive Scoring

Extracted Data  
often in  
Table  
Format



Classification:  
Cats and Dogs, Speech Recognition  
Movie Recommendation



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

**Data  $X_2$**

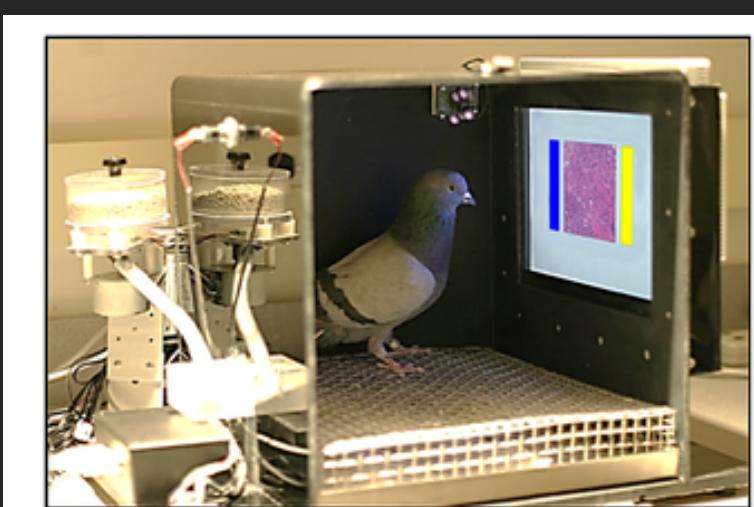
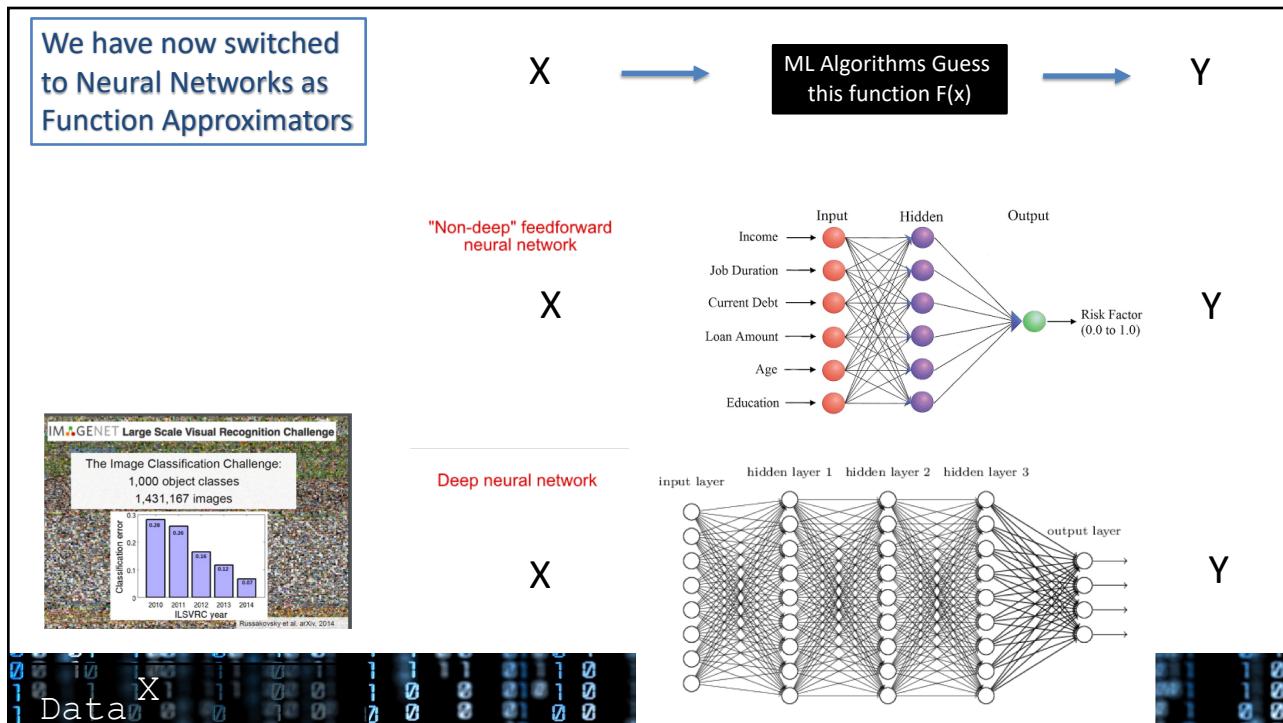


Fig 1. The pigeons' training environment.



## Current Hot Topics in AI/Data



Reinforcement Learning

### Unsupervised Image to Image

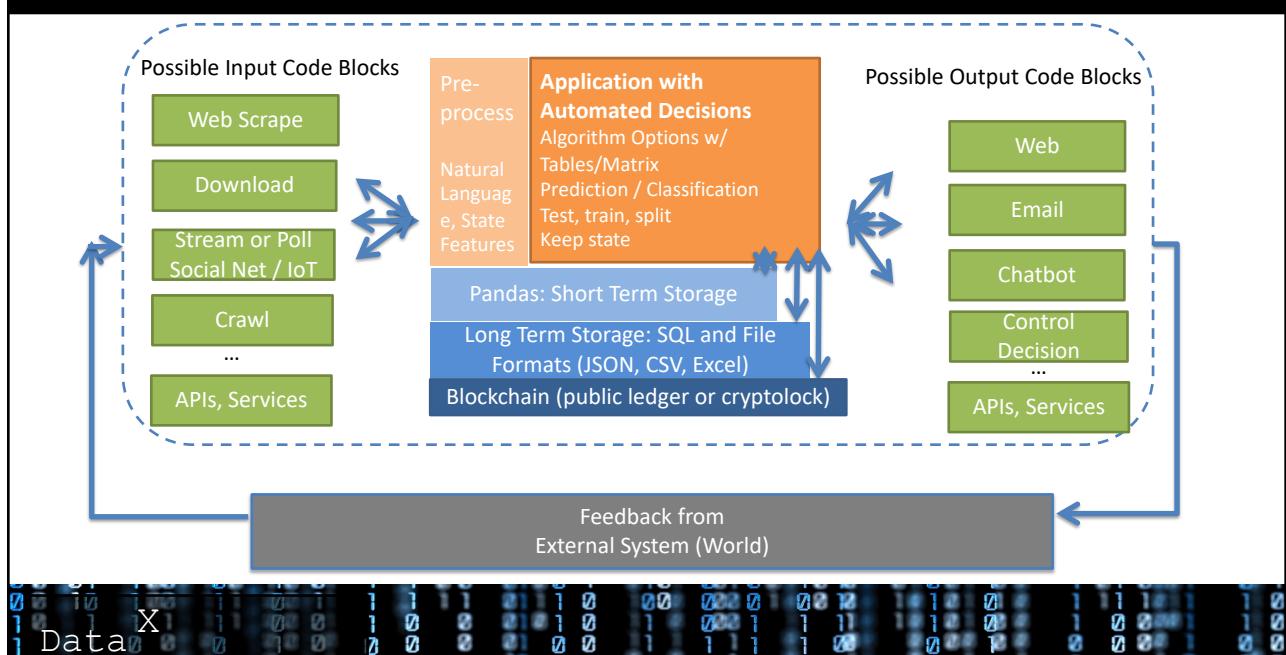


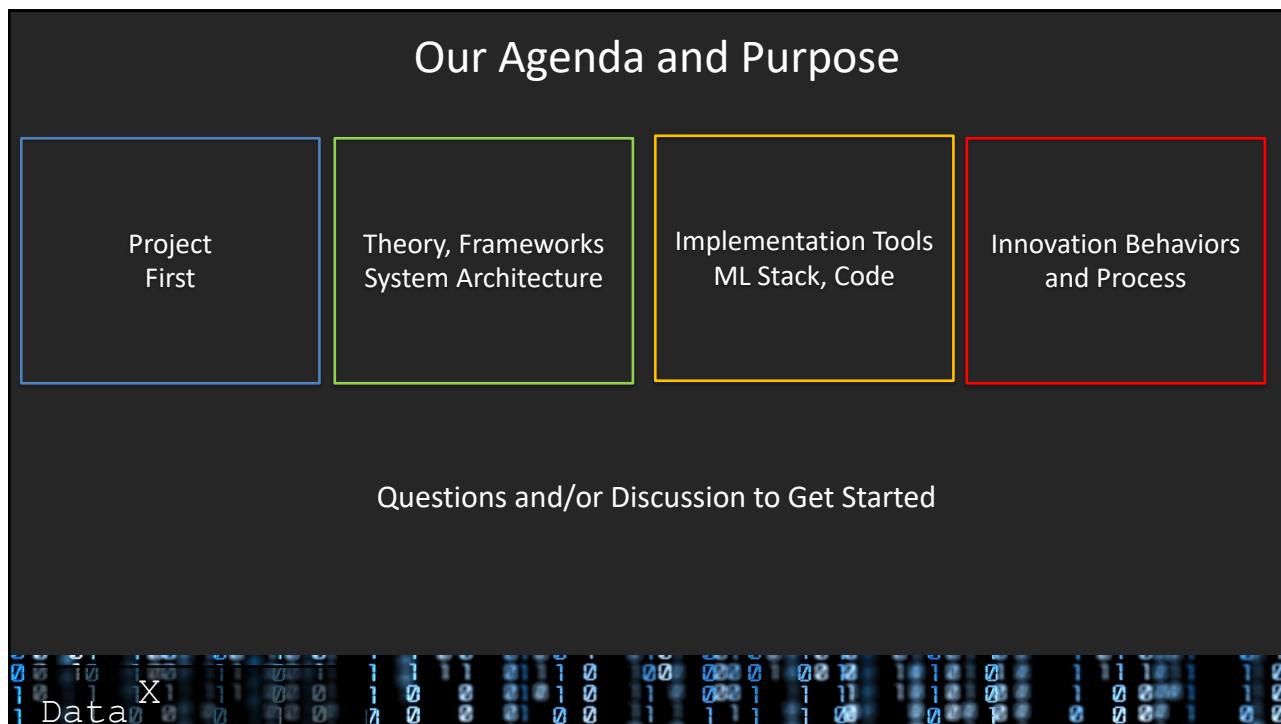
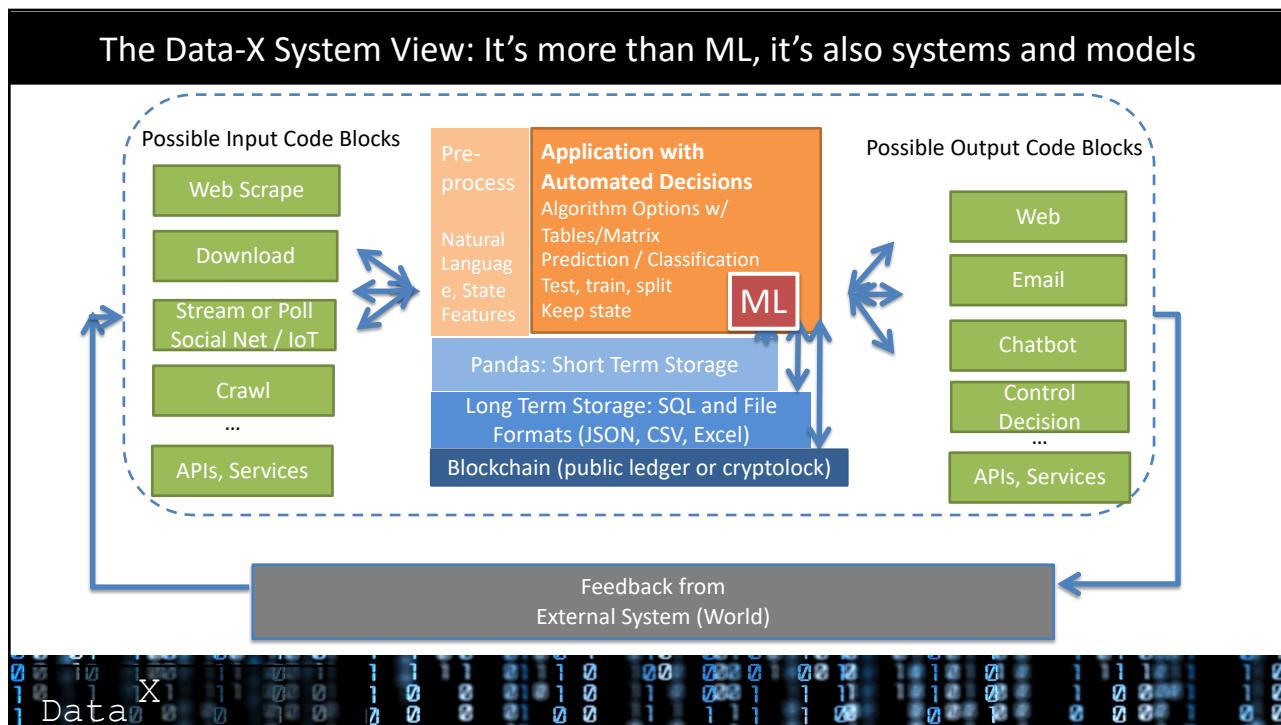
[CycleGAN: Zhu, Park, Isola & Efros, 2017] Peter Abbeel – UC Berkeley | Gradscope | ConvNetAI

Generative Networks and Deep Fakes



## The Data-X System View





End of Section

0 0 1 0 1 0 0 1 0 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0  
1 1 0 1 X1 1 0 0 1 0 1 0 0 0 1 0 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0  
1 Data 0 0 1 0 0 0 0 1 0 0 1 1 1 1 1 0