# Data: Things to Consider

## CSCI-P556 Applied Machine Learning
## Lecture 3

**D.S. Williamson**

# Agenda and Learning Outcomes

## Today's Topics

- **Topics**:

  - Data and Why is it needed?

  - What challenges do data present?

  - Where to get data? Gather or Collect?

- **Learning Outcomes**: At the end of today's class, students should:

  - Describe the bad data problem

  - Understand how to define needed data for an experiment

# TopHat Questions

# What is Data?

**Anything can be considered as data**

- Anything that can be observed and digitized, can function as data for Machine Learning

- Depending on the learning algorithm (supervised or unsupervised), additional meta data may be needed (e.g. attributes/features/inputs and labels)
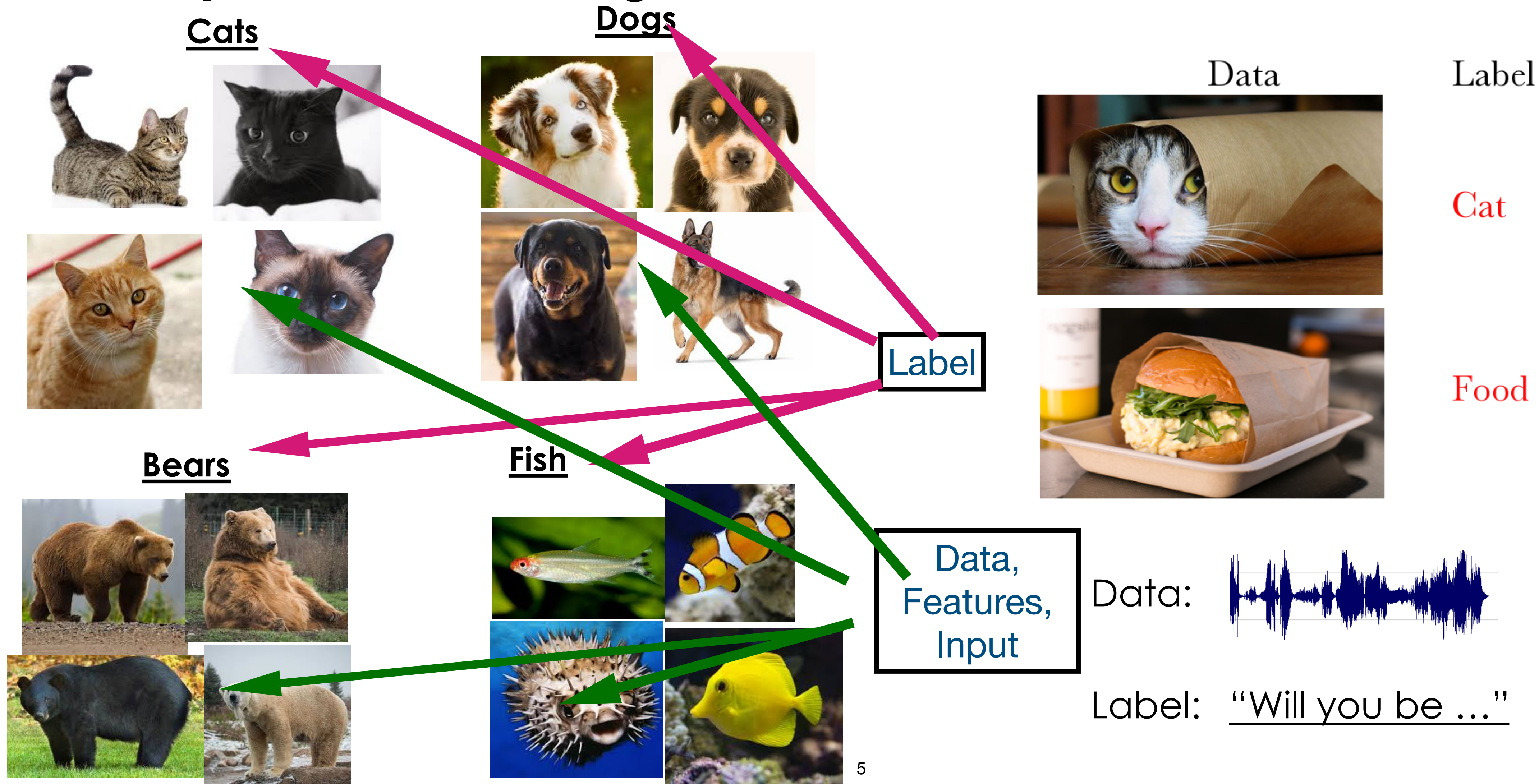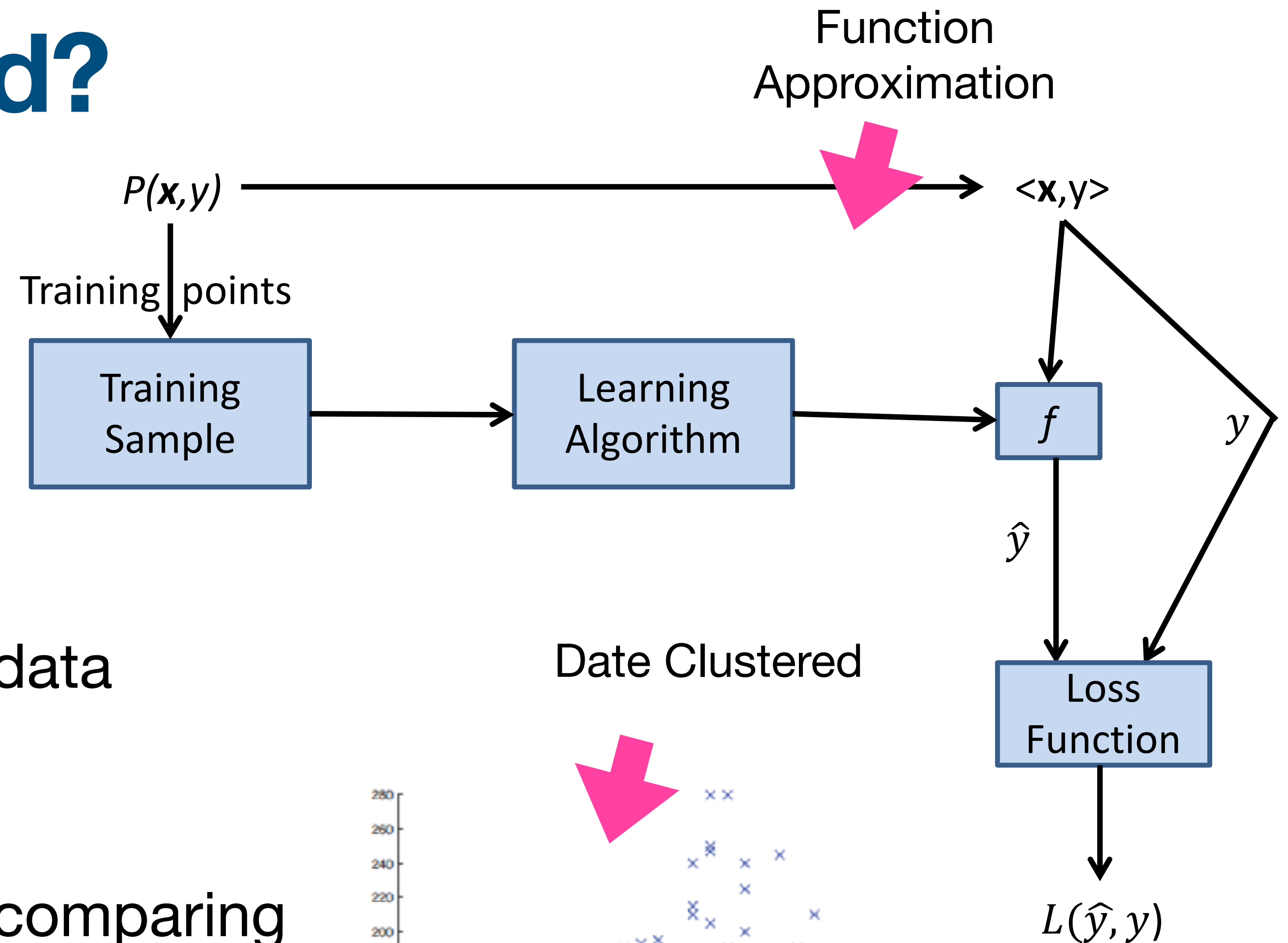
Health Records

Images
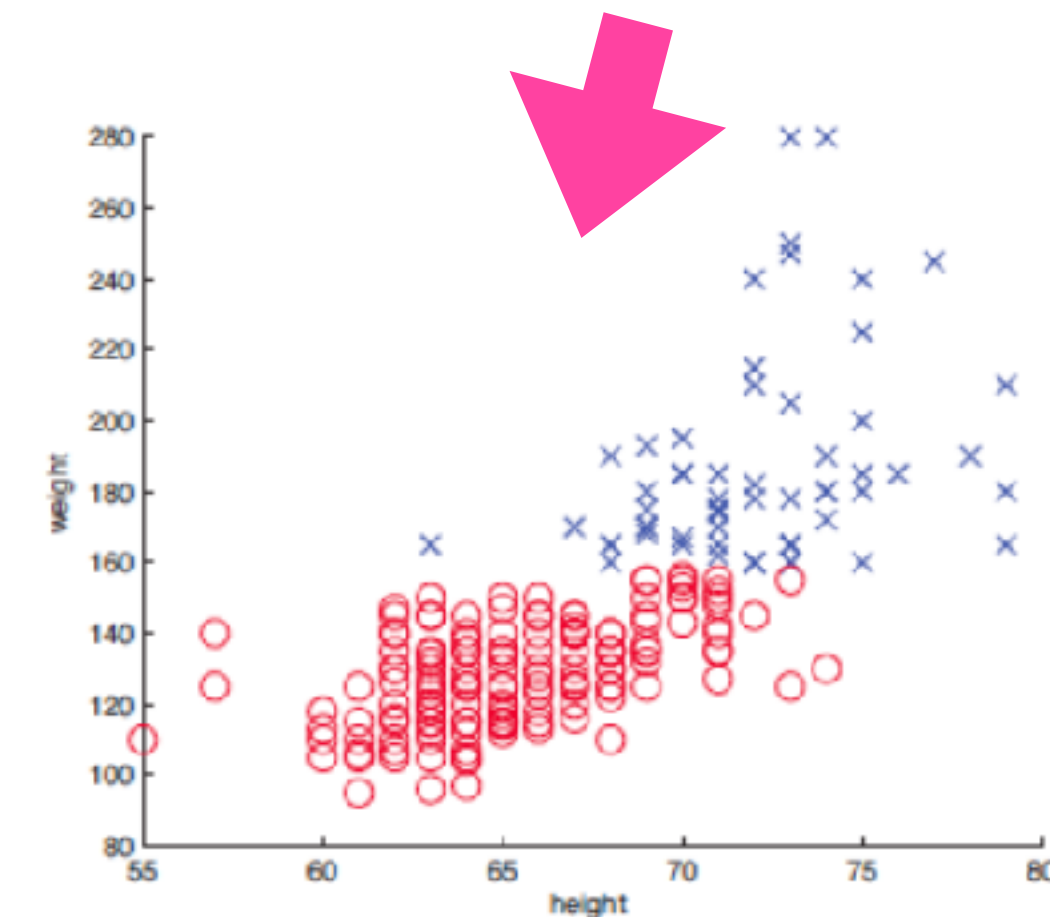
Tweets

# Recall: Attributes and Labels of Data

## For Supervised Learning

**Cats**

**Dogs**

**Bears**

**Fish**

Food

Label

Data, Features, Input

Data:

Label: "Will you be ..."

# Why is Data Needed?

## No data no learning!

Function Approximation

$P(\mathbf{x}, y)$ → $\langle \mathbf{x}, y \rangle$

Training points

| Training Sample | Learning Algorithm | $f$ | $y$ |

$\hat{y}$

- Machine Learning "learns" from data

  - Either a function for prediction

  - Or patterns for clustering and comparing

Date Clustered

Loss Function

$L(\hat{y}, y)$

# Main Challenges of Machine Learning
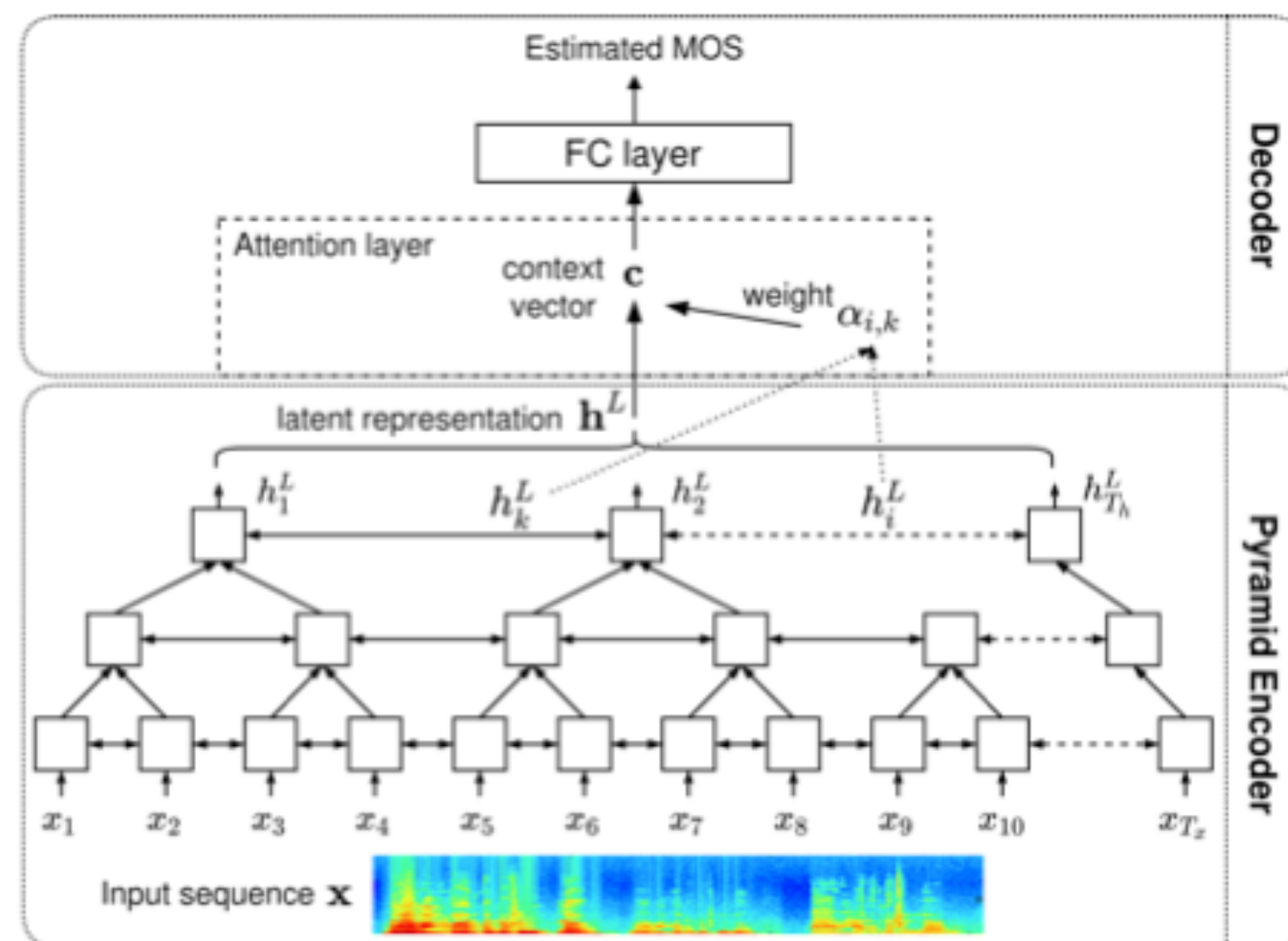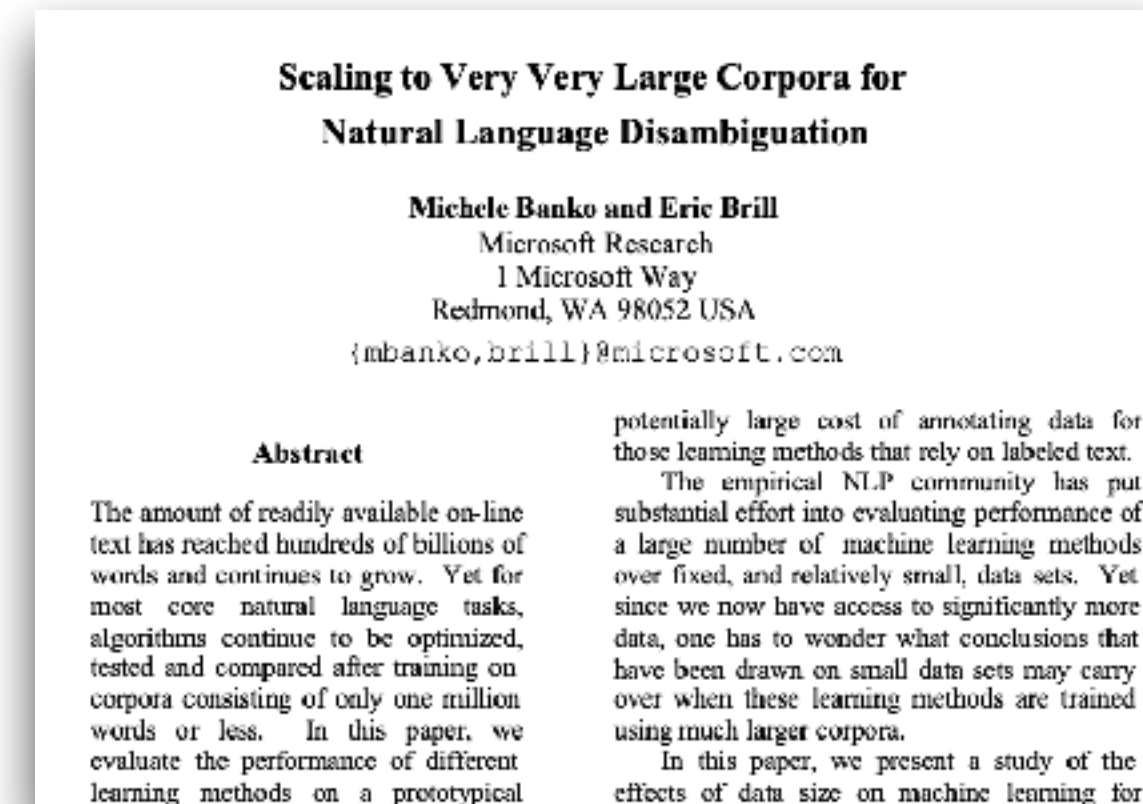## Nothing is perfect!

- Our main task, as machine learning scientists, is to select a learning algorithm and train it on some data

- Three potential challenges:

  - **<u>Bad algorithm</u>**: The algorithm is not good enough for this particular problem

  - **<u>Bad data</u>**: The data is not good enough

  - **<u>Bad evaluation</u>**: Metrics aren't indicative of future usage

- Today, we'll talk about the ***<u>Bad data problem</u>***

# The Unreasonable Effectiveness of Data
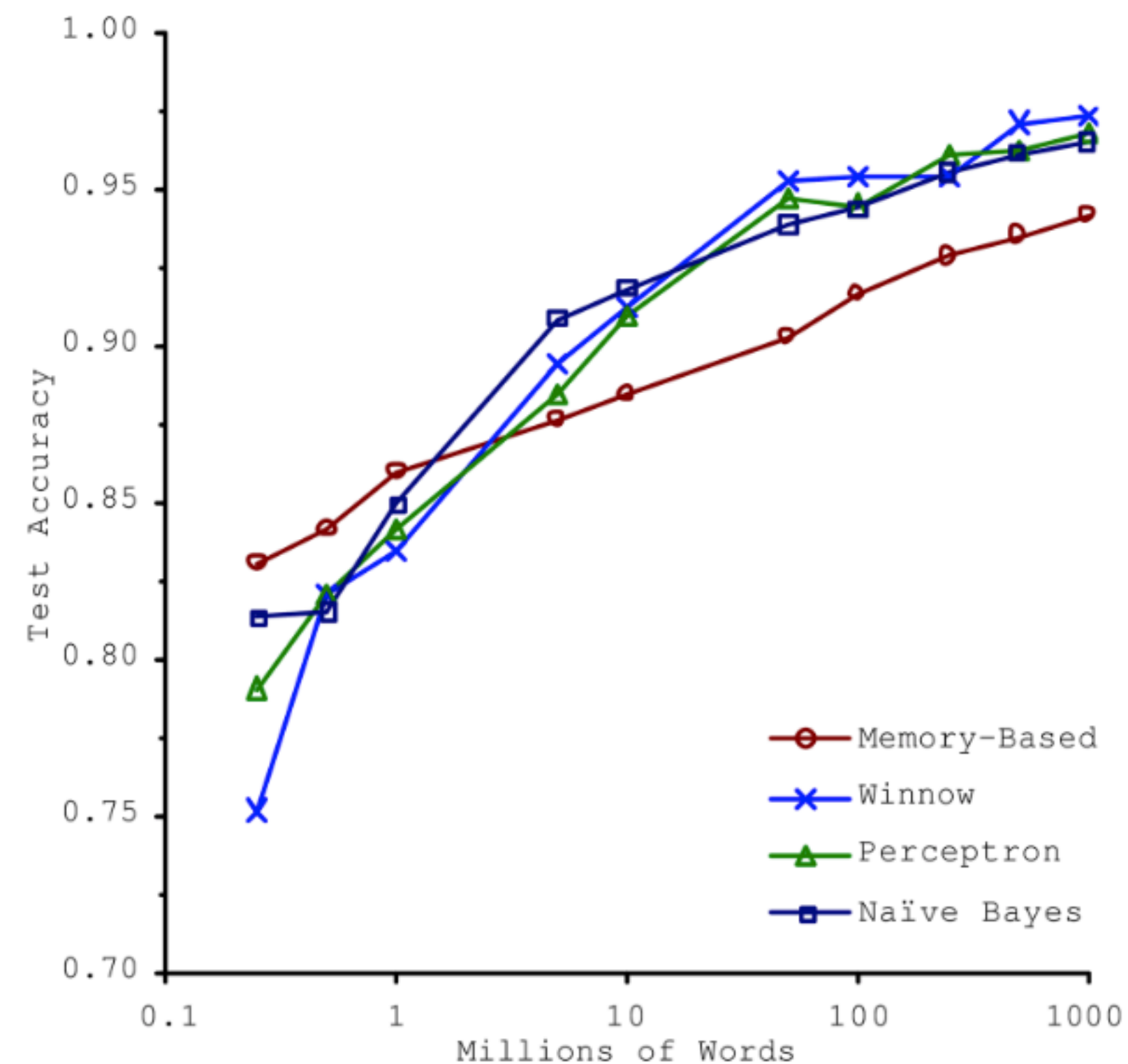
## Focus at School vs. Work

- In school, algorithms are king!

- In work, data often matters most!



- Microsoft researchers in 2001 showed that different algorithms performed almost identically all on a complex NLP problem once enough data was given.

- "These results suggest that we may want to reconsider the trade-off between spending time and money on algorithm development versus spending it on corpus development."

- This idea was further popularized by Google researchers.

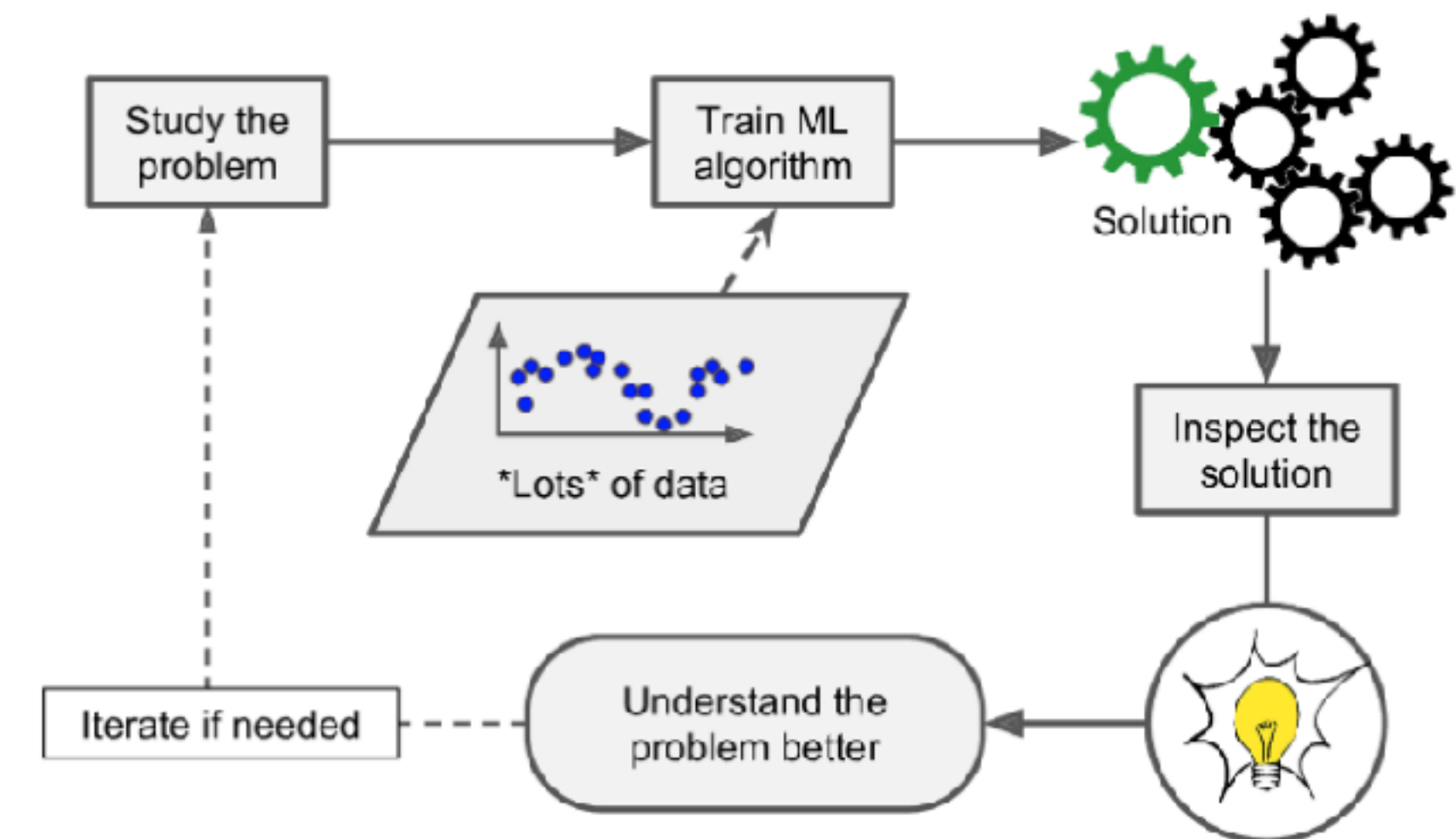# The Importance of Data Vs. Algorithm
## Data Matters



- The performance of different learning algorithms as data varies

- Higher test accuracy indicates better performance

- Notice the main trends,

  - All algorithms perform **"bad" with less data**.

  - All algorithms perform **"well" with large quantities of data.**

- Data is too often overlooked by computer and data scientists!

# Main Challenges of Machine Learning

## #1 Insufficient Quantity of Training Data



- Babies and Toddlers can learn from few examples

  - Instantly learn how to eat, who are their parents.

  - Easily learn objects (e.g. spoon, bottle, etc.)

- **Machines are not as good!**

  - A **lot** of data is needed for algorithms to work
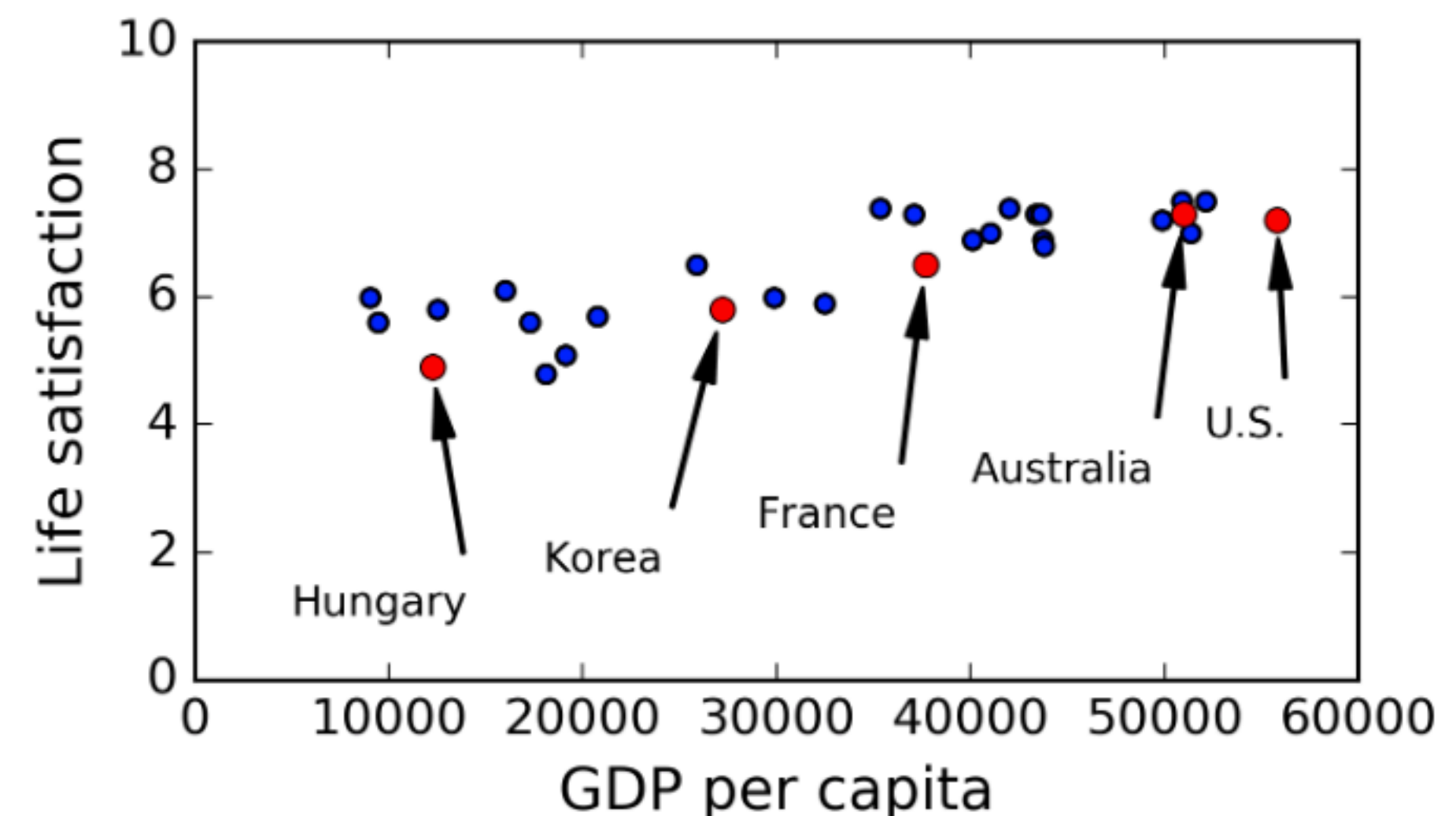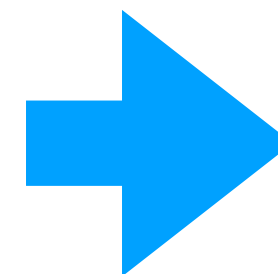
# Main Challenges of Machine Learning

## #2 Nonrepresentative Training Data

- Data you **_have_** must be representative of data you **_expect_** to see.

- Suppose you want to know if money makes people happy. You have data that shows the GDP per capita for multiple countries and user rated Life Satisfaction

The current data implies "happiness" increases linearly with money!

Table 1-1. Does money make people happier?

| Country | GDP per capita (USD) | Life satisfaction |
|---------|---------------------|-------------------|
| Hungary | 12,240 | 4.9 |
| Korea | 27,195 | 5.8 |
| France | 37,675 | 6.5 |
| Australia | 50,962 | 7.3 |
| United States | 55,805 | 7.2 |



11
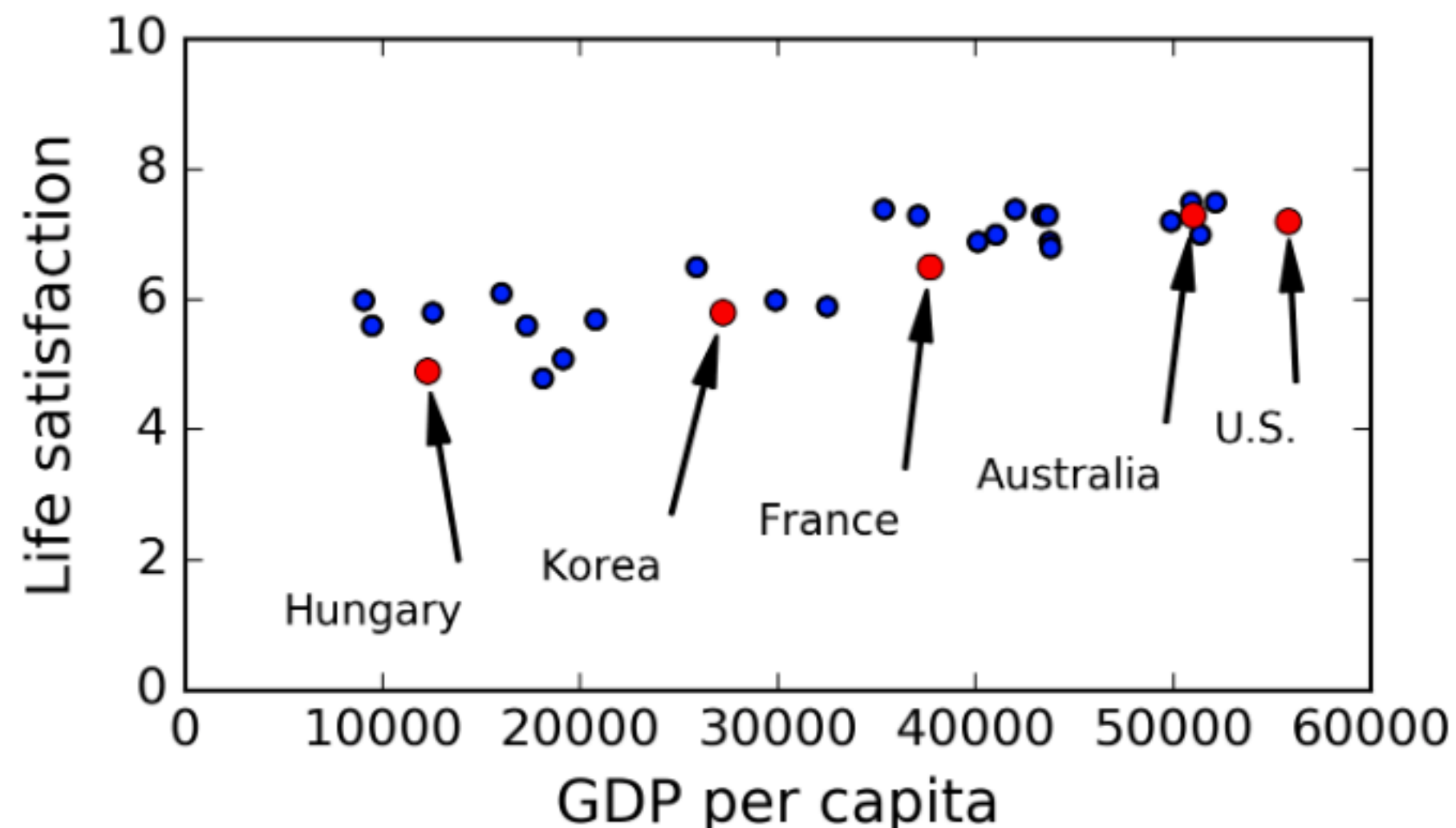
# Main Challenges of Machine Learning

## #2 Nonrepresentative Training Data

- Data you **_have_** must be representative of data you **_expect_** to see.

- Suppose you want to know if money makes people happy. You have data that shows the GDP per capita for multiple countries and user rated Life Satisfaction
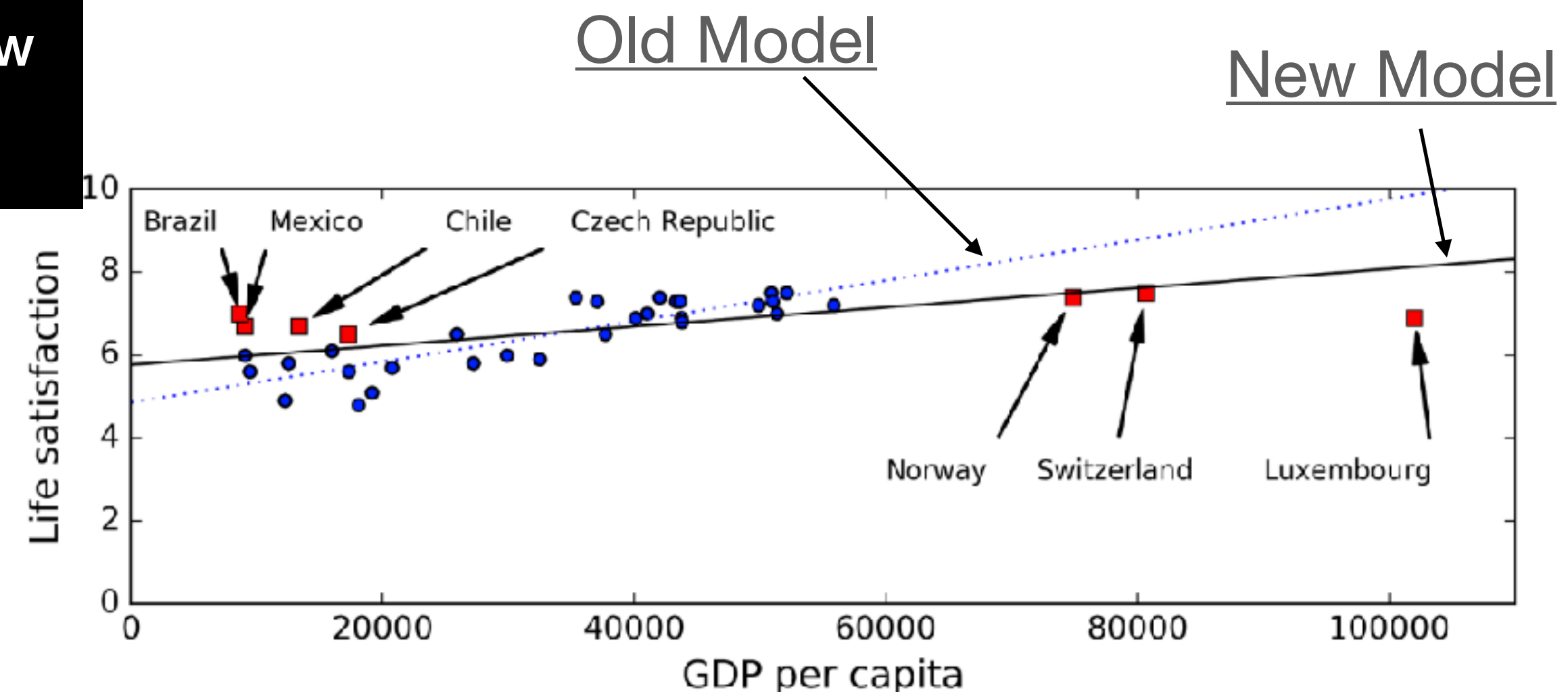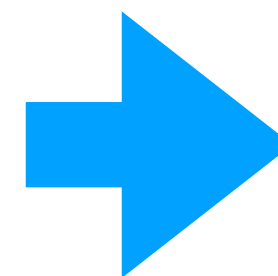
New data doesn't follow the old trends!

Very rich countries are not happier than moderately rich countries (may be unhappy)

Some poor countries seem happier than rich

Consider New Data

Old Model

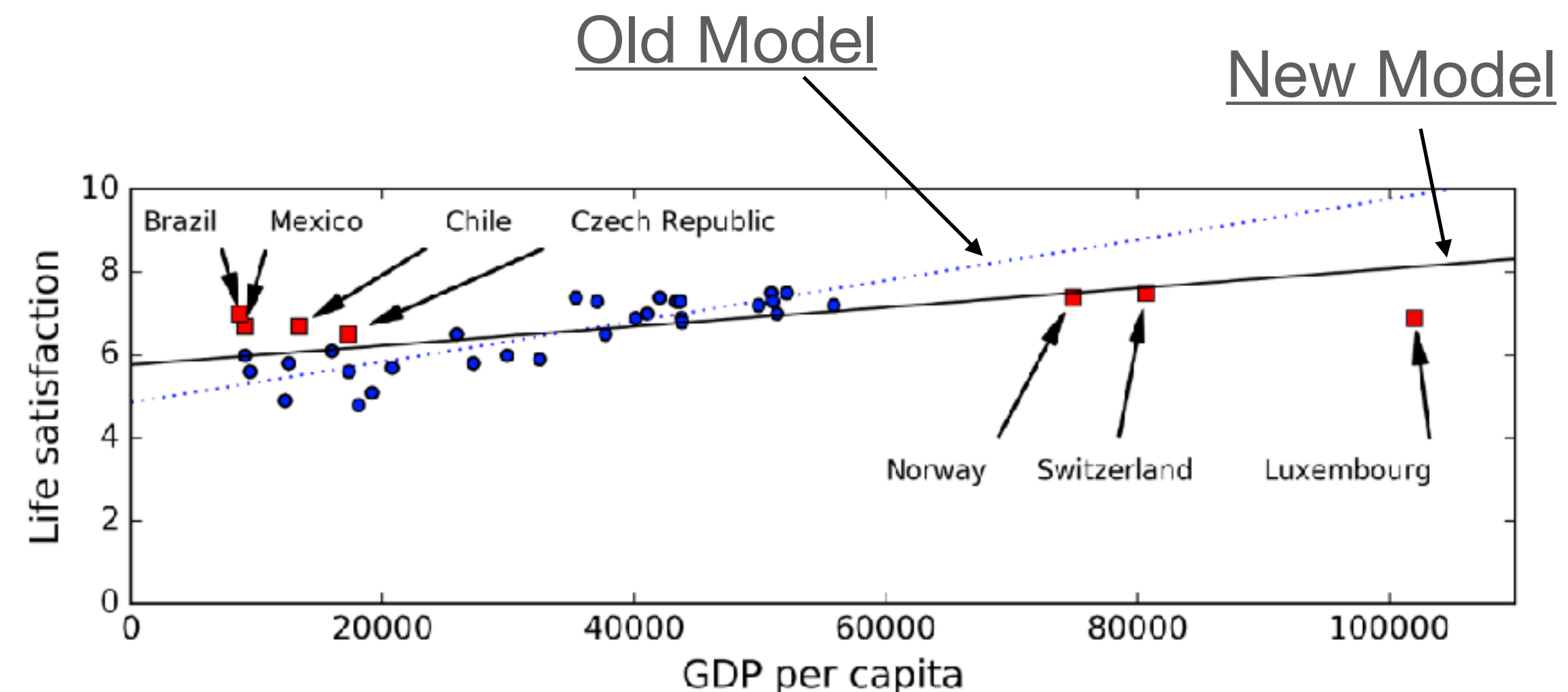New Model

# Main Challenges of Machine Learning

## #2 Nonrepresentative Training Data

- Data you **_have_** must be representative of data you **_expect_** to see.

- Suppose you want to know if money makes people happy. You have data that shows the GDP per capita for multiple countries and user rated Life Satisfaction

- Your training data set **_MUST_** be representative of the cases you want to **_generalize_** to!

  - In this case, what type of data do you need?

New data doesn't follow the old trends!

Very rich countries are not happier than moderately rich countries (may be unhappy)

Some poor countries seem happier than rich

Old Model

New Model



13

# Main Challenges of Machine Learning

## #2 Nonrepresentative Training Data - Sampling Bias

### A Famous Example of Sampling Bias

Perhaps the most famous example of sampling bias happened during the US presidential election in 1936, which pitted Landon against Roosevelt: the *Literary Digest* conducted a very large poll, sending mail to about 10 million people. It got 2.4 million answers, and predicted with high confidence that Landon would get 57% of the votes.

Instead, Roosevelt won with 62% of the votes. The flaw was in the *Literary Digest's* sampling method:

- First, to obtain the addresses to send the polls to, the *Literary Digest* used telephone directories, lists of magazine subscribers, club membership lists, and the like. All of these lists tend to favor wealthier people, who are more likely to vote Republican (hence Landon).

- Second, less than 25% of the people who received the poll answered. Again, this introduces a sampling bias, by ruling out people who don't care much about politics, people who don't like the *Literary Digest*, and other key groups. This is a special type of sampling bias called *nonresponse bias*.

Here is another example: say you want to build a system to recognize funk music videos. One way to build your training set is to search "funk music" on YouTube and use the resulting videos. But this assumes that YouTube's search engine returns a set of videos that are representative of all the funk music videos on YouTube. In reality, the search results are likely to be biased toward popular artists (and if you live in Brazil you will get a lot of "funk carioca" videos, which sound nothing like James Brown). On the other hand, how else can you get a large training set?
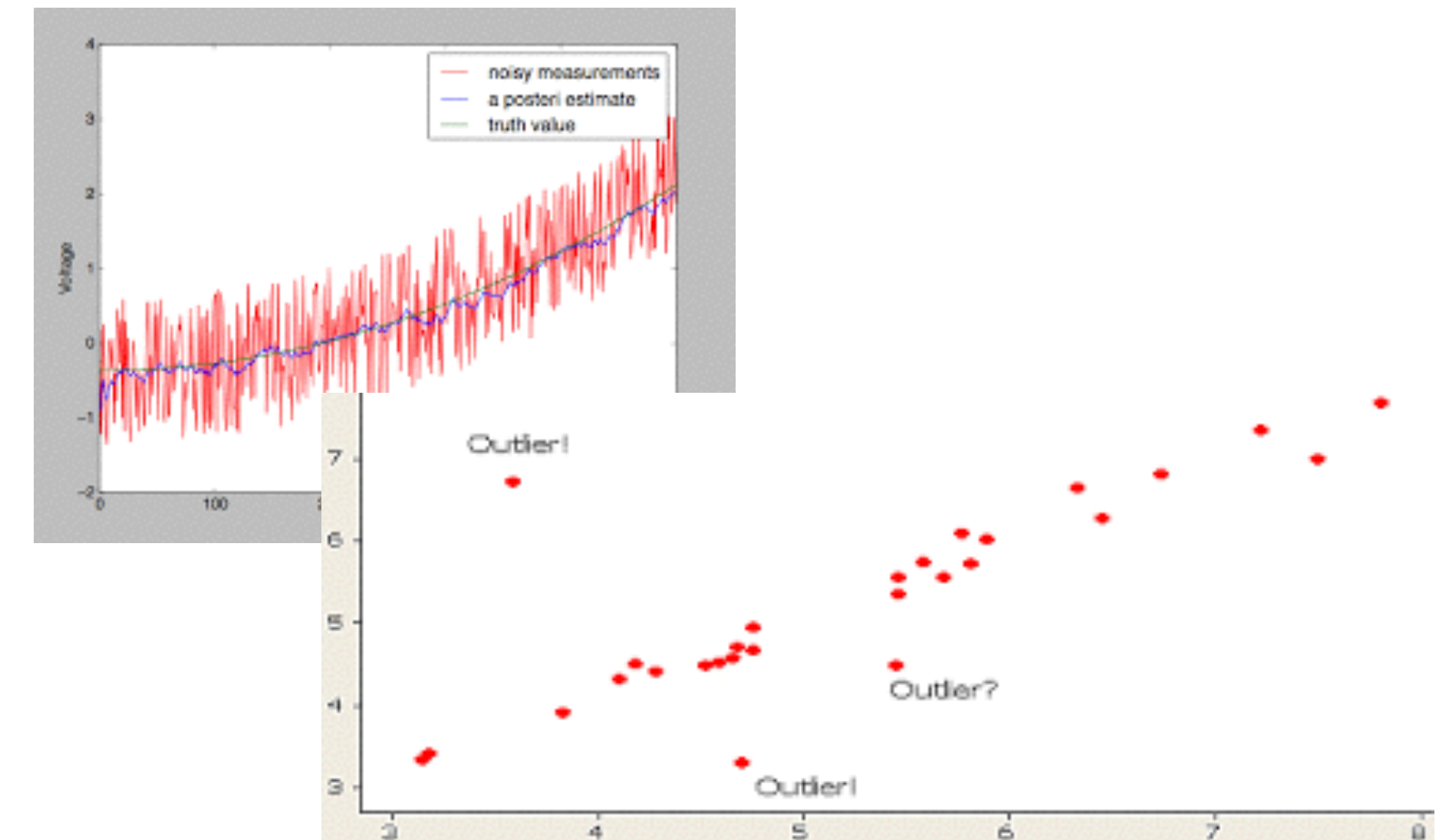
- Can you think of or imagine any other instances of sampling bias?

14

# Main Challenges of Machine Learning

## #3 Poor Quality Data and #4 Irrelevant Features
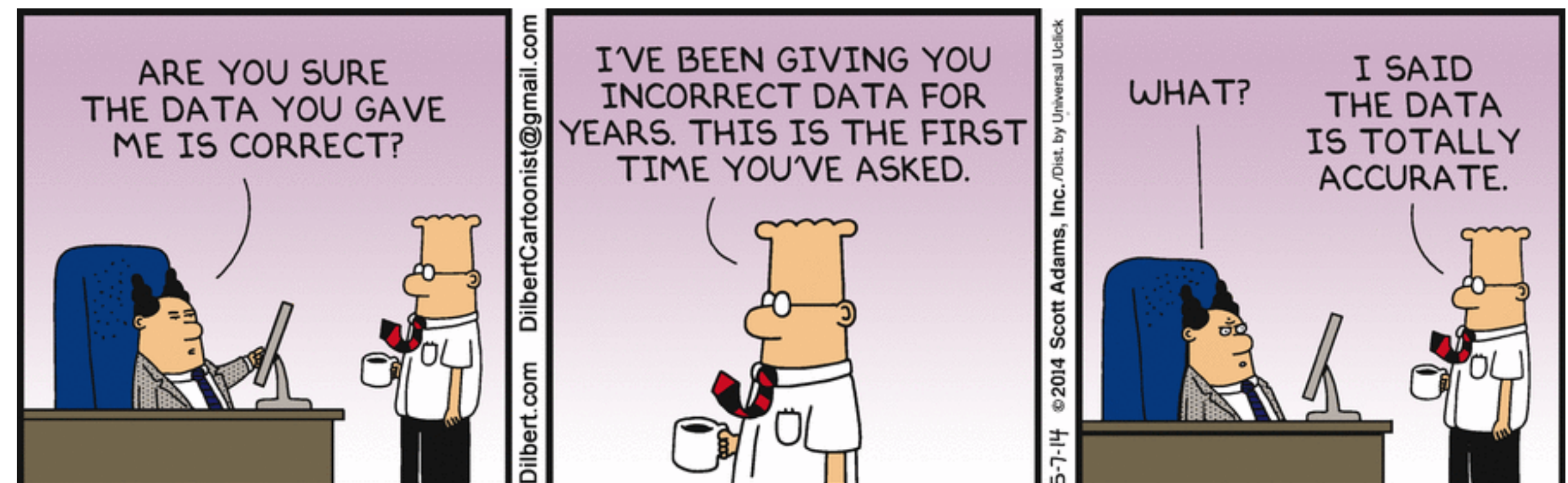
- **Poor Quality Data**

  - Errors, outliers and noise (poor sensors) make machine learning more difficult

  - Often need to clean data in this scenario

    - Discard or remove outliers

    - Ignore missing features/attributes, a particular instance, or fill in values

    - Find/collect "replacement" data

- **Irrelevant Features**

  - Garbage In, Garbage Out

  - Feature engineering may be needed:

    - **Selection**: select most useful features

    - **Extraction**: combine existing feaures

    - **Creation**: create new features from new data

# Data Selection for Experimental Design
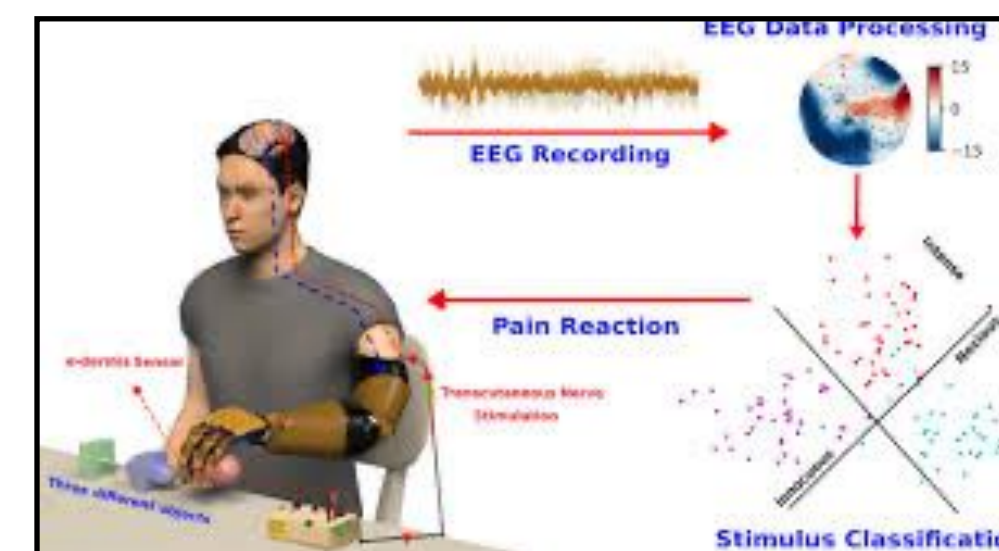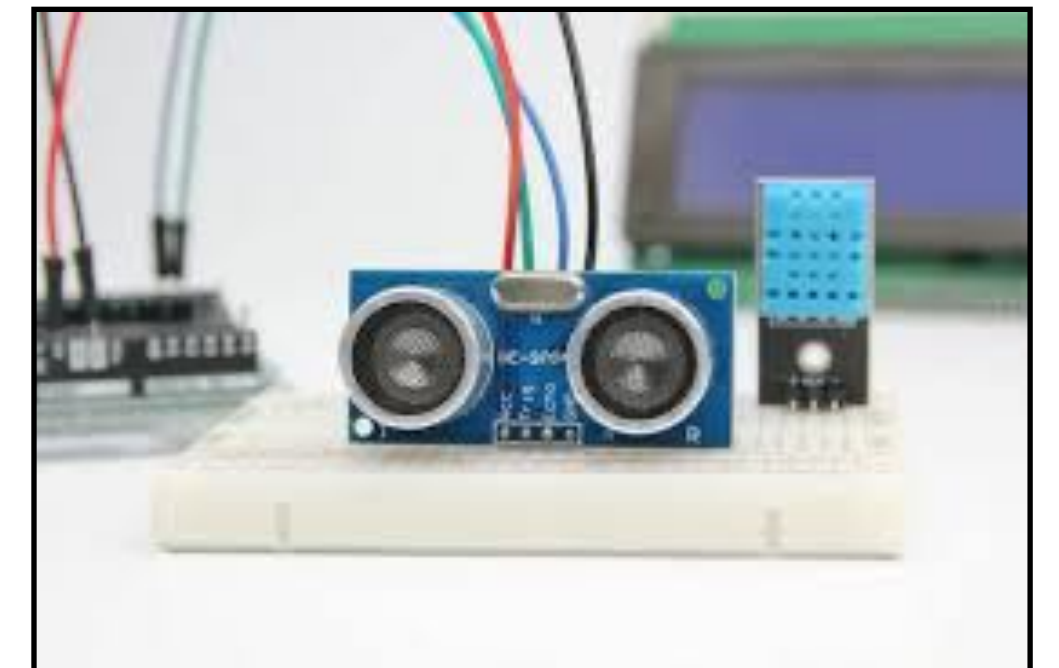## Group practice - Zoom Break out rooms

- **Scenario**: You are a machine learning scientists and you want to develop an approach that can classify whether a given fruit on a farm is ripe or not. You want your approach to work for as many fruits and farms as possible.

- **Task**: You need to gather/collect data to train your system. In a group, think of the type of data that you need, including how/where/when it is collected or gathered from. Be sure that you consider the bad data problem.

  - In each group, one person take notes to summarize your approach

  - We'll come back and discuss in 5-10 minutes

- **Assumptions**:

  - You have the perfect classifier to use and metrics for evaluation

  - A drone can be used to fly over the farm(s) to see the fruits

  - Data will be annotated accurately once gathered

# How do we get Data?
## Collect it Ourselves

- Sensors for data collection are everywhere!

- Carefully designed experiments must be conducted.

- Must consider data challenges

  - **Inaccurate, noisy and/or missing data**: Over collect when possible

  - **Data imbalance and bias**: must ensure "enough" data from ALL classes, scenarios, and environments are collected

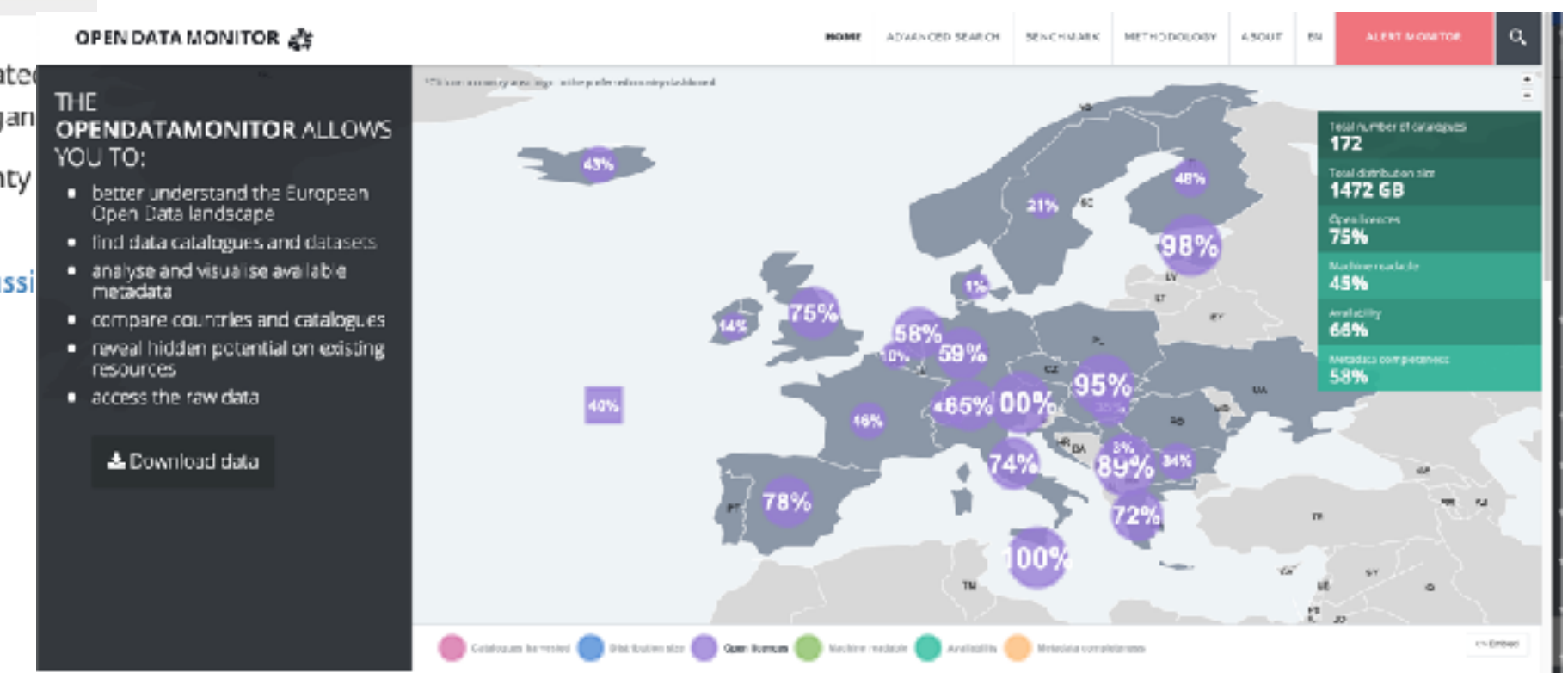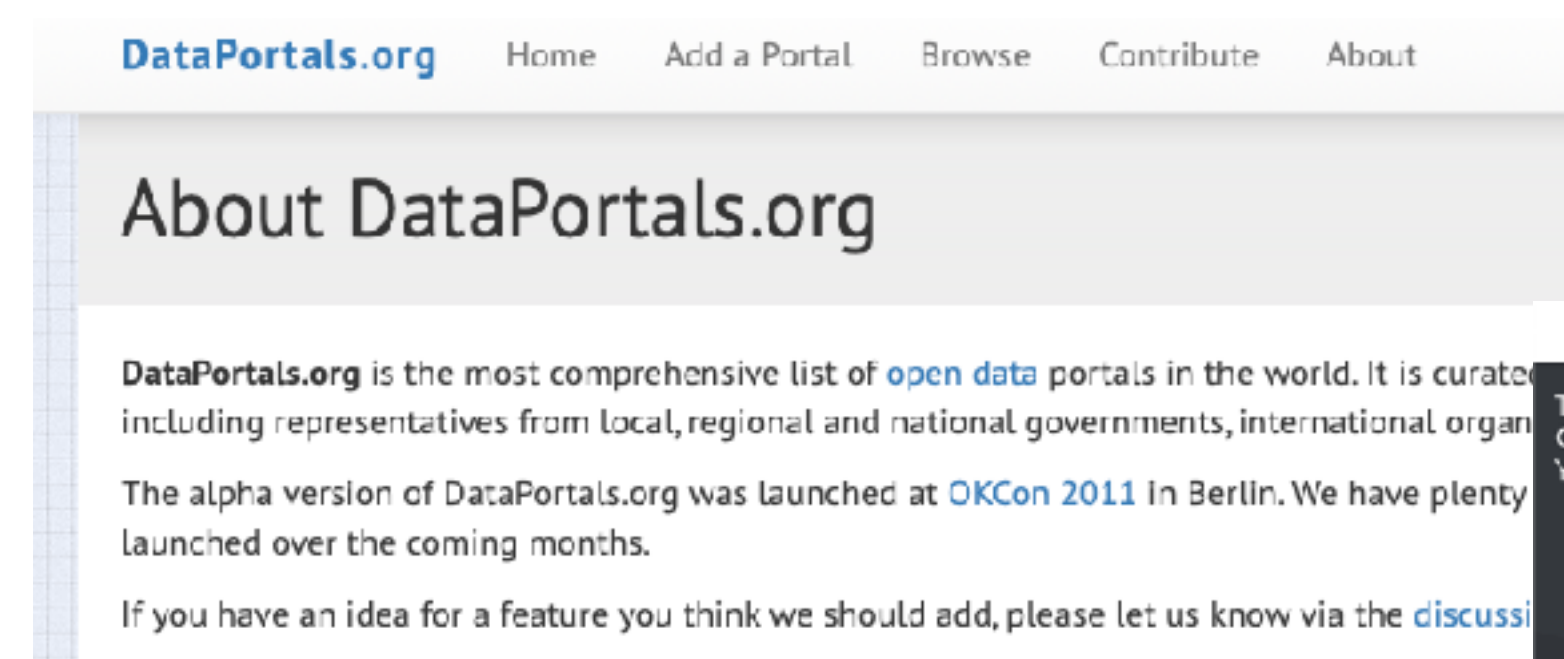# How do we get Data?
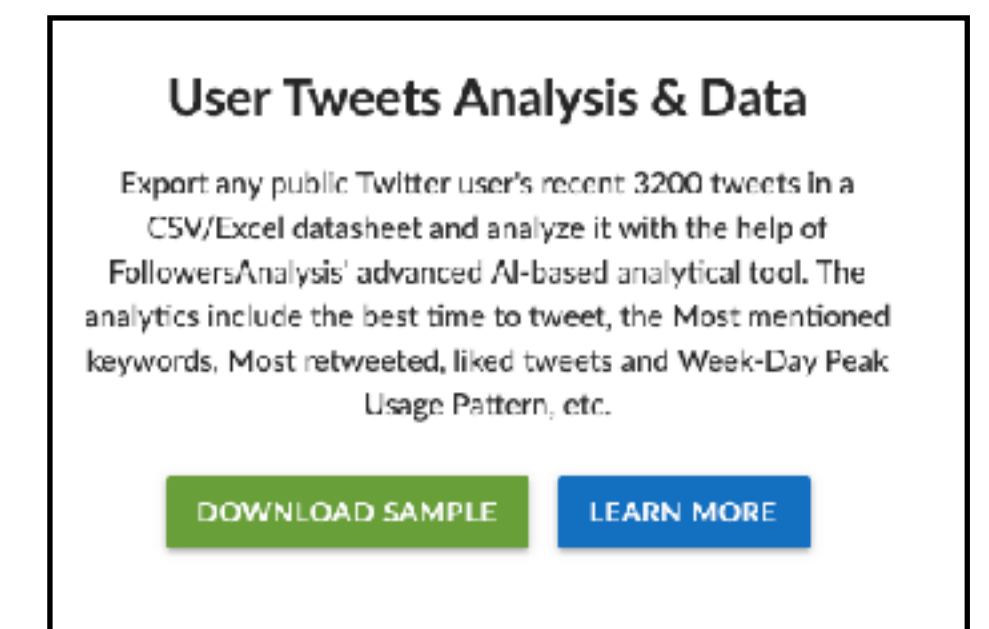## Collect it Ourselves



- Now that you have data, it must be annotated!

- This is the hard, costly and time-consuming part!

- For example, objects in pictures must be denoted. Transcriptions of audio must be provided

- You may have to do this yourself! Luckily online crowdsourcing has helped with this.

  - This is perfect, right?

# How do we get Data?

## Online Sources

- Thankfully, there are many downloadable datasets that we can start with!
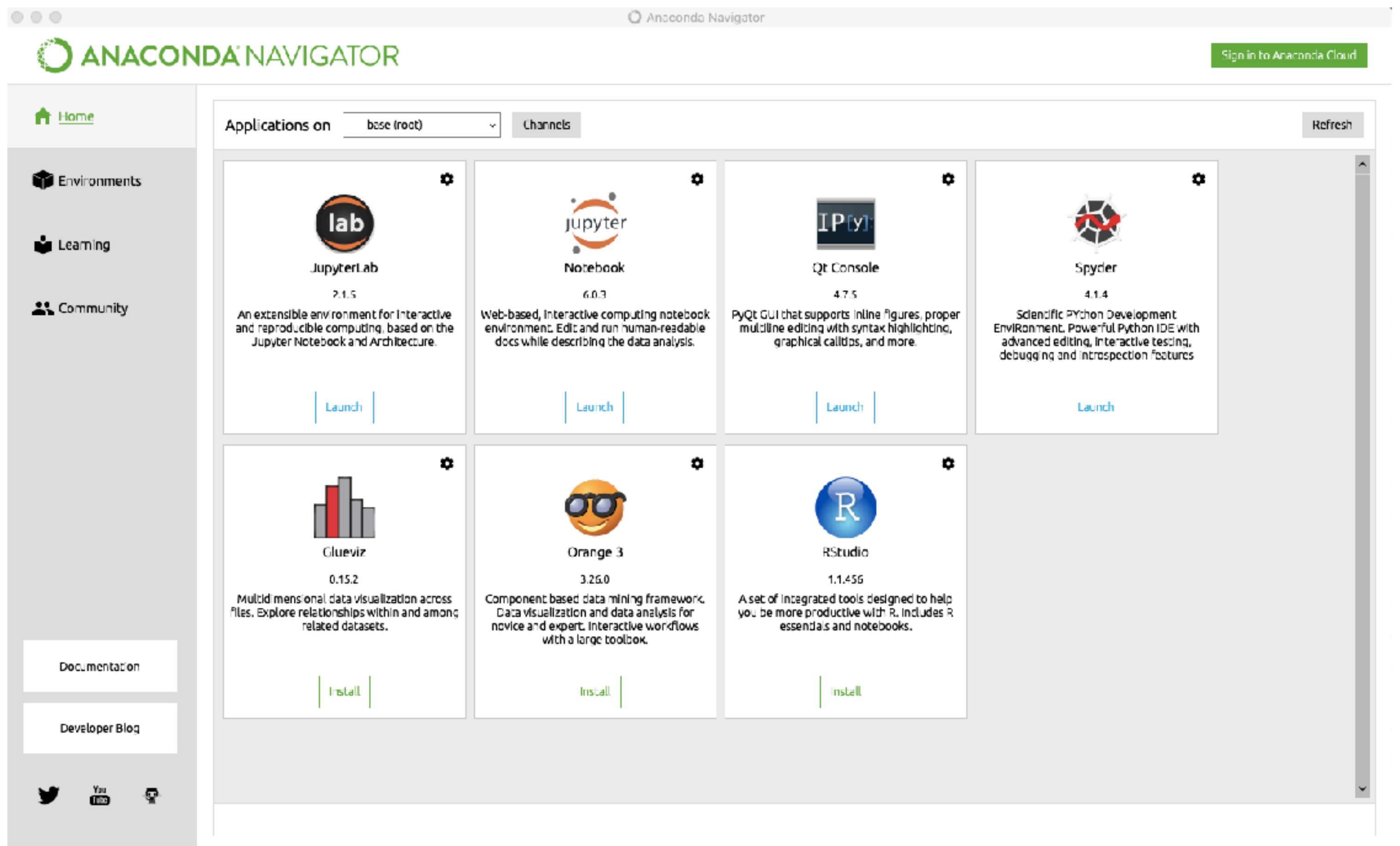
# Data Pre-processing

## Now that we have data, what's next? An Example Case

- Suppose you are a Data Scientist at a Housing Corporation. Your boss wants you to build a prediction model of median housing prices in California using their census data

- **Data has info about**: population, median income, median housing prices, … for each block group or district in California.

- **How should this problem be framed**?

  - Supervised Learning, Unsupervised learning, Reinforcement Learning? Why?

  - Classification, Regression, Other? Why?

  - Batch vs. Online?

# Python Environment Setup

## Option 1. Anaconda Environment

- Use one of the following links for instructions for working with Anaconda

  - https://docs.anaconda.com/anaconda/

  - https://docs.anaconda.com/anaconda/navigator/tutorials/manage-environments/

- Easy and intuitive programming environment

# Python Environment Setup
## Option 2. Isolated/Local Environment

- Follow specific instructions for your system (Windows, Linux, Apple) from the site: https://packaging.python.org/guides/installing-using-pip-and-virtual-environments/

- Basic steps (Python 3.3 or later on Mac):

  - Install pip:
    ```
    python3 -m pip install --user --upgrade pip
    ```

  - Check the installation:
    ```
    python3 -m pip --version
    pip 9.0.1 from $HOME/.local/lib/python3.6/site-packages (python 3.6)
    ```

  - Create Virtual Environment:
    ```
    python3 -m venv env
    ```

  - Activate Virtual Environment:
    ```
    source env/bin/activate
    ```
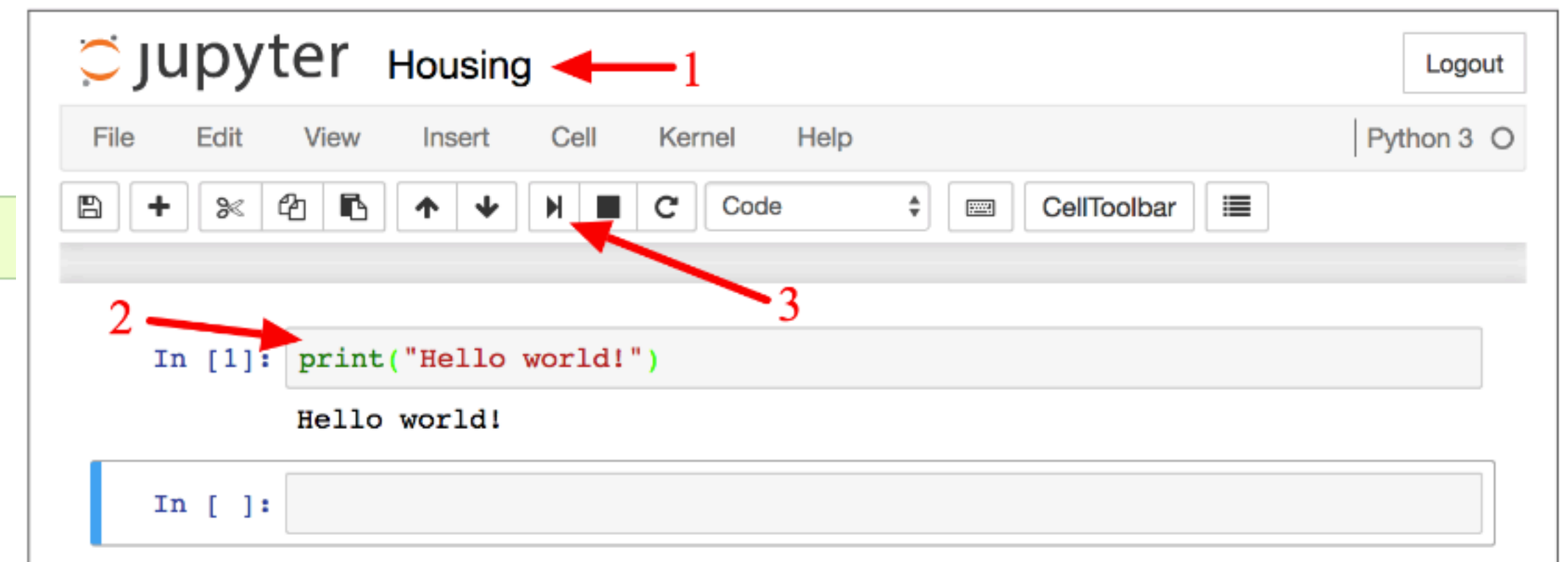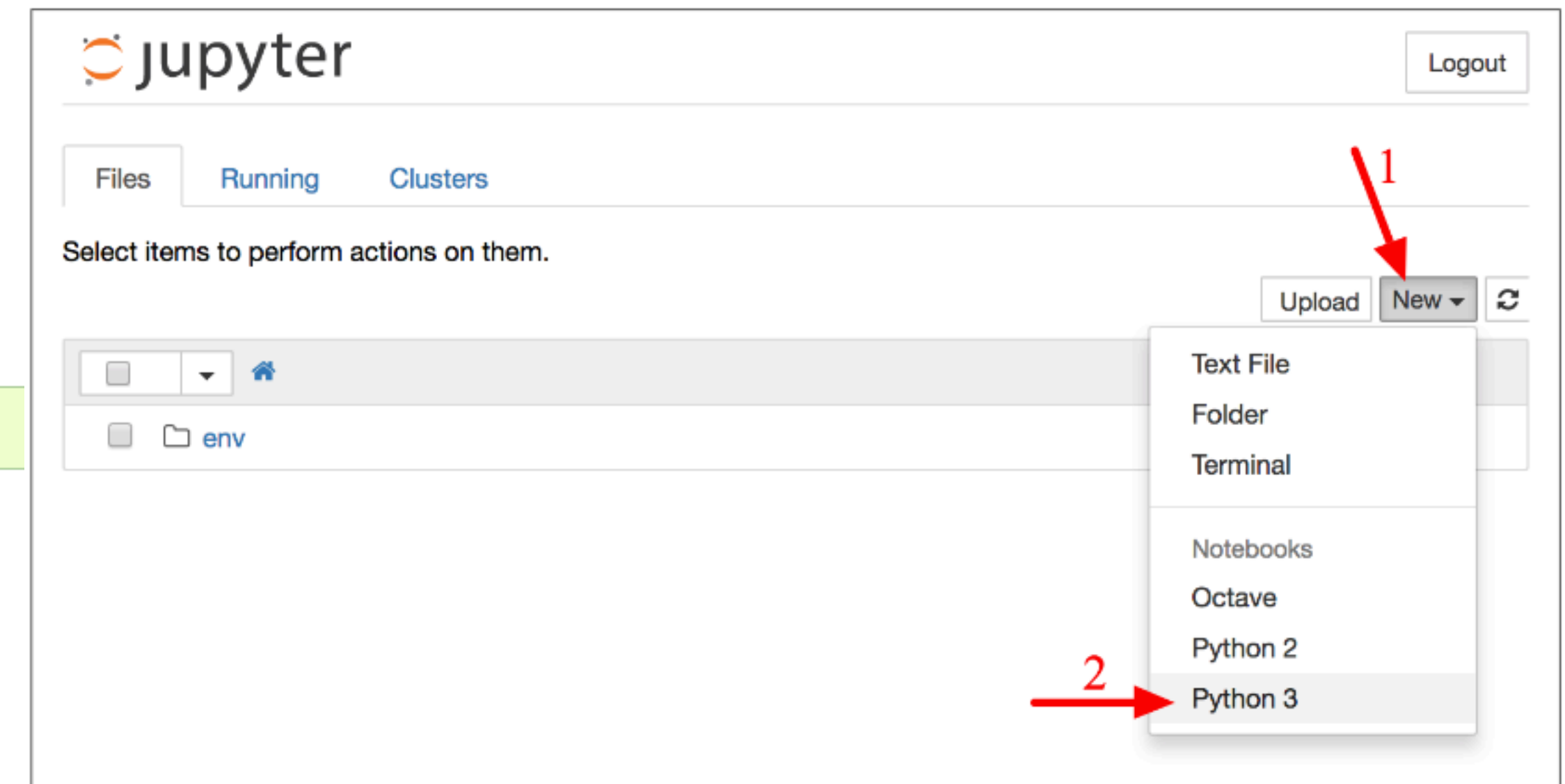
# Python Environment Setup
## Option 2. Isolated/Local Environment

- Now you an install any/all required modules

```
python3 –m pip install -U jupyter matplotlib dumpy pandas spicy scikit-learn
```

- You can now work in jupyter notebook (or commandline)

```
jupyter notebook
```

# Next Class

**Be sure your Python environment is setup appropriately**
  **- Post issues to Piazza**