

# COMP/EECE 7/8745 Machine Learning

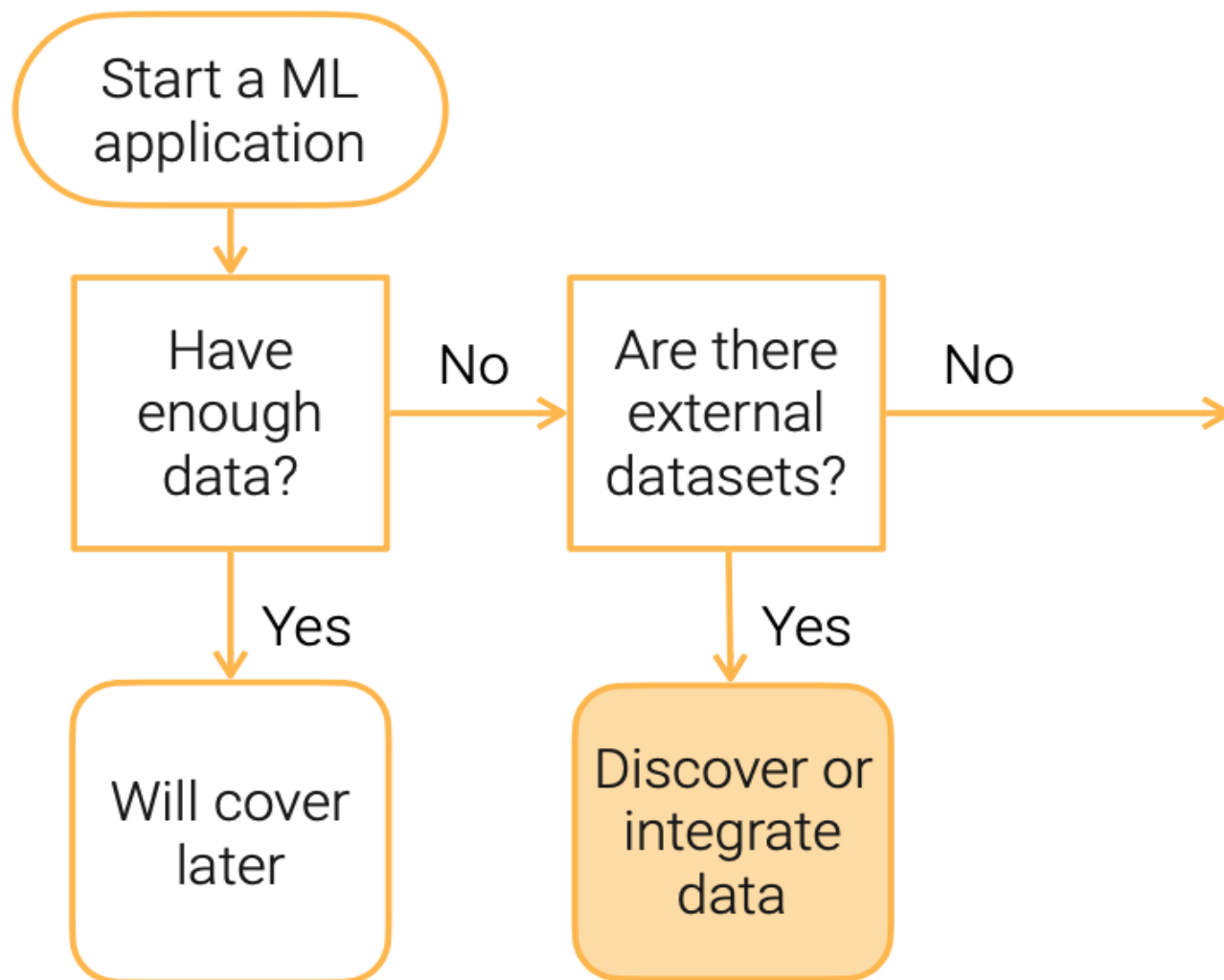
Topics:

## **Data preparations:**

- Data acquisition and labeling
- Web scraping
- Why data preprocessing and clearing?
- Data integration, transformation, and reduction
- Discretization

Md Zahangir Alom  
Department of Computer Science  
University of Memphis, TN

# Flow Chart for Data Acquisition



# Discover What Data is Available

- Identify existing datasets
- Find benchmark datasets to evaluate a new idea
  - E.g. A diverse set of small to medium datasets for a new hyperparameter tuning algorithm
  - E.g. Large scale datasets for a very big deep neural network

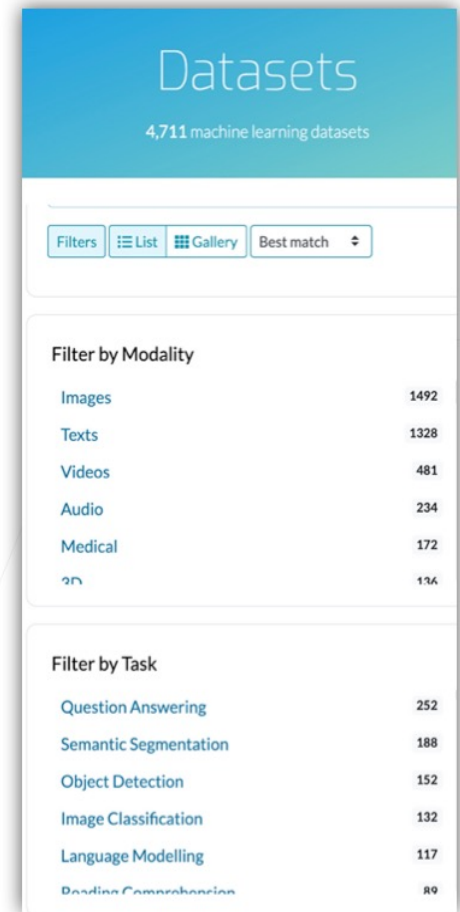
# Popular ML datasets

- MNIST: digits written by employees of the US Census Bureau
- ImageNet: millions of images from image search engines
- AudioSet: YouTube sound clips for sound classification
- LibriSpeech: 1000 hours of English speech from audiobook
- Kinetics: YouTube videos clips for human actions classification
- KITTI: traffic scenarios recorded by cameras and other sensors
- Amazon Review: customer reviews and from Amazon online shopping
- SQuAD: question-answer pairs derived from Wikipedia

More at [https://en.wikipedia.org/wiki/List\\_of\\_datasets\\_for\\_machine-learning\\_research](https://en.wikipedia.org/wiki/List_of_datasets_for_machine-learning_research)

# Where to Find Datasets

- [Paperswithcodes Datasets](#): academic datasets with leaderboard
- [Kaggle Datasets](#): ML datasets uploaded by data scientists
- [Google Dataset search](#): search datasets in the Web
- Various toolkits datasets: [tensorflow](#), [huggingface](#)
- Various conference/company ML competitions
- [Open Data on AWS](#): 100+ large-scale raw data
- Data lakes in your own organization

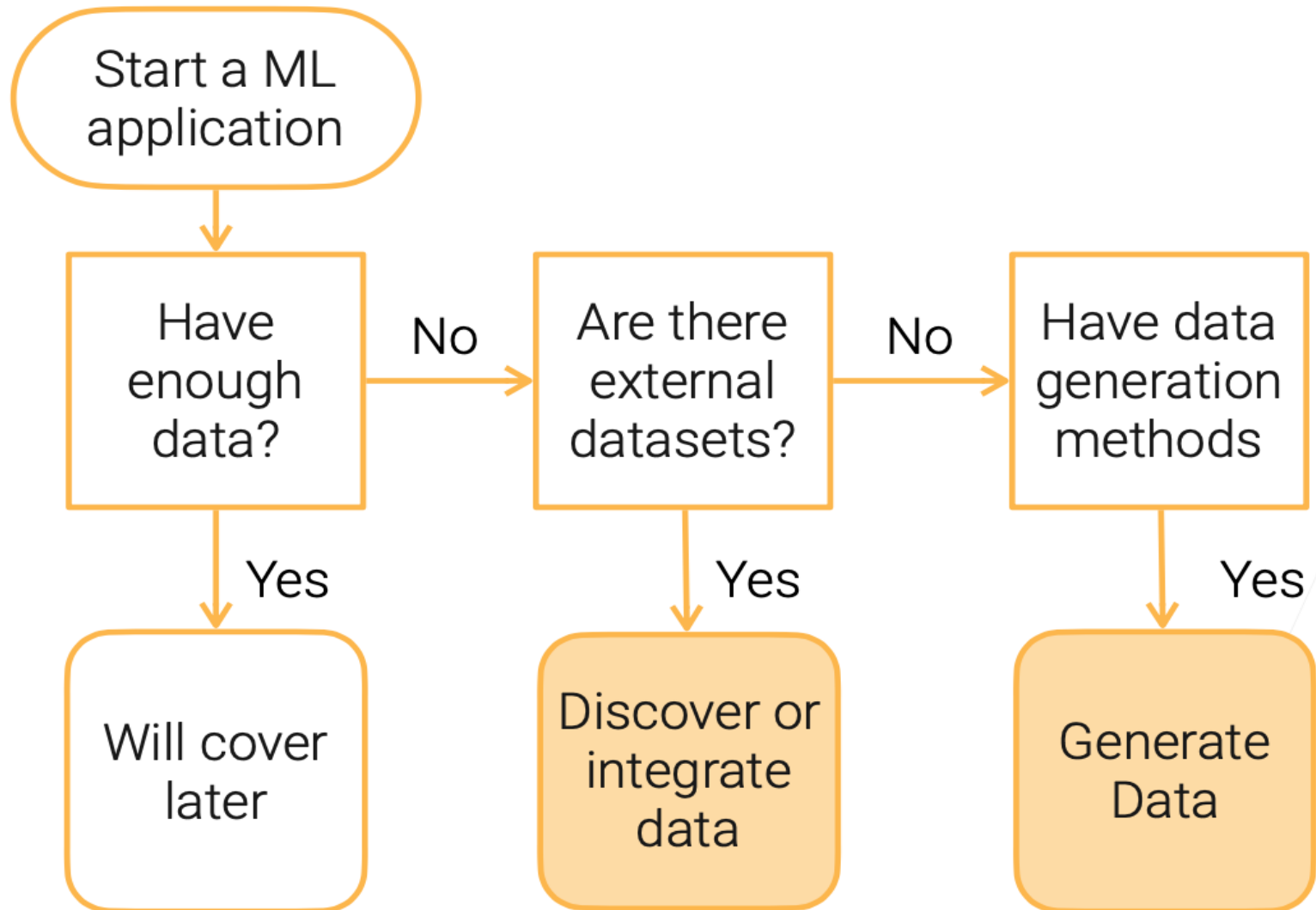


# Where to Find Datasets

	Pros	Cons
Academic datasets	Clean, proper difficulty	Limited choices, too simplified, usually small scale
Competition datasets	Closer to real ML applications	Still simplified, and only available for hot topics
Raw Data	Great flexibility	Needs a lot of effort to process

- You often need to deal with raw data in industrial settings
- Data curation can be a big projection involving multiple teams. Processing pipeline, storage, legal issue, privacy,...

# Flow Chart for Data Acquisition



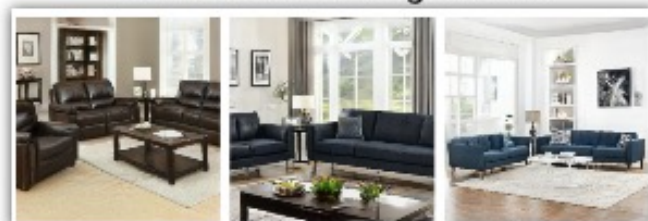
# Generate Synthetic Data

- Use GANs



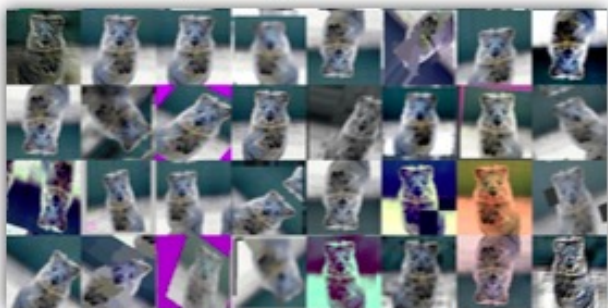
<https://thispersondoesnotexist.com/>

Furnitures in living rooms

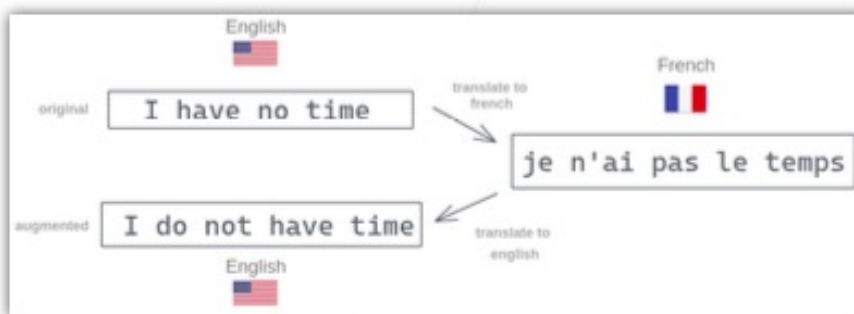


Gadde et al., ICCV'21

- Simulation
- Data augmentations  
Image augmentation



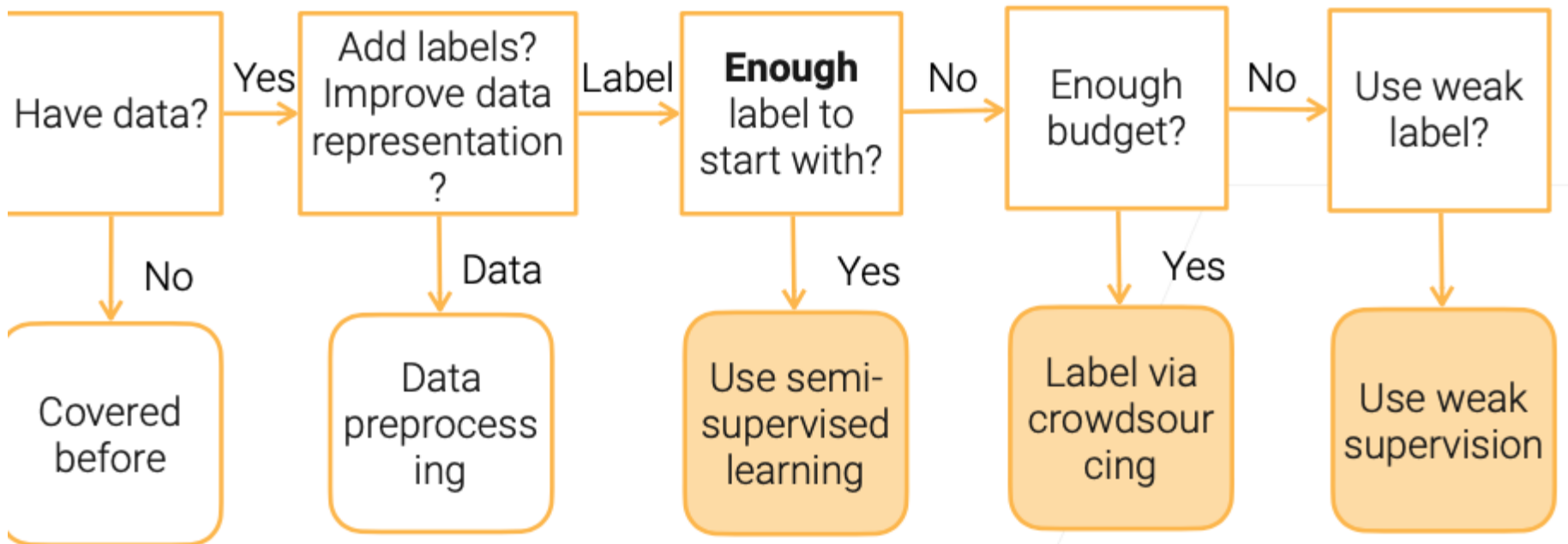
Back Translation



<https://amitniss.com>



# Flow Chart for Data Labelling

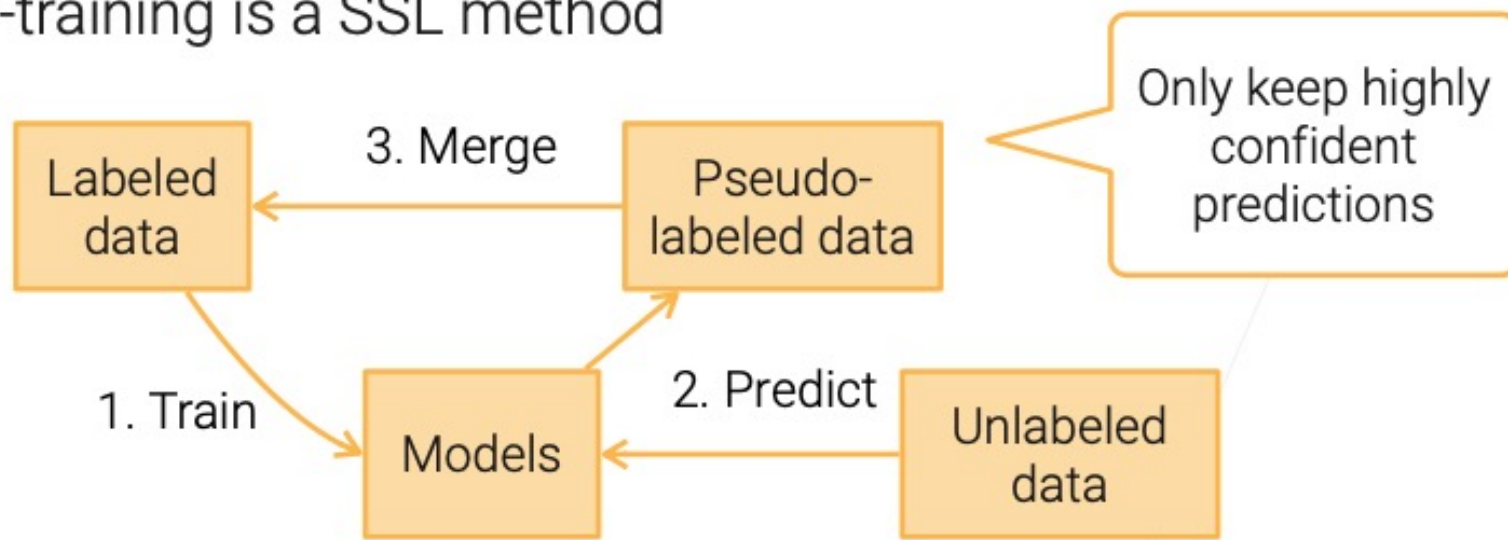


# Semi-Supervised Learning (SSL)

- Focus on the scenario where there is a small amount of labeled data, along with large amount of unlabeled data
- Make assumptions on data distribution to use unlabeled data
  - **Continuity assumption:** examples with similar features are more likely to have the same label
  - **Cluster assumption:** data have inherent cluster structure, examples in the same cluster tend to have the same label
  - **Manifold assumption:** data lie on a manifold of much lower dimension than the input space

# Self-training

- Self-training is a SSL method



- We can use expensive models
  - Deep neural networks, model ensemble/bagging

# Label through Crowdsourcing

- ImageNet labeled millions of images through Amazon Mechanical Turk. It took several years and millions dollars to build
- According to Amazon SageMaker Ground Truth, the estimated price of using Amazon Mechanical Turk:

Image/text classification	\$0.012 per label
Bounding box	\$0.024 per box
Semantic segmentation	\$0.84 per image

# User interaction

- Example of user instruction and labeling task (MIT Place365)

**Start** **Is this a cliff scene?**  
**Definition:** a high, steep or overhanging face of rock.

**Task**  
For each of the **810** images, answer yes or no to the above question. Only answer **Yes** to **real photos**. Always answer **No** to **cartoon, drawing, CG rendering**, or real photos with a **large text overlay** on the photo. Here are some examples:

No Single Object No Text Overlay No Drawing No Screenshot No Graphics No Bad Photo

Not Only Logo No Magazine/Newspaper

No No Yes Yes

Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes

b)

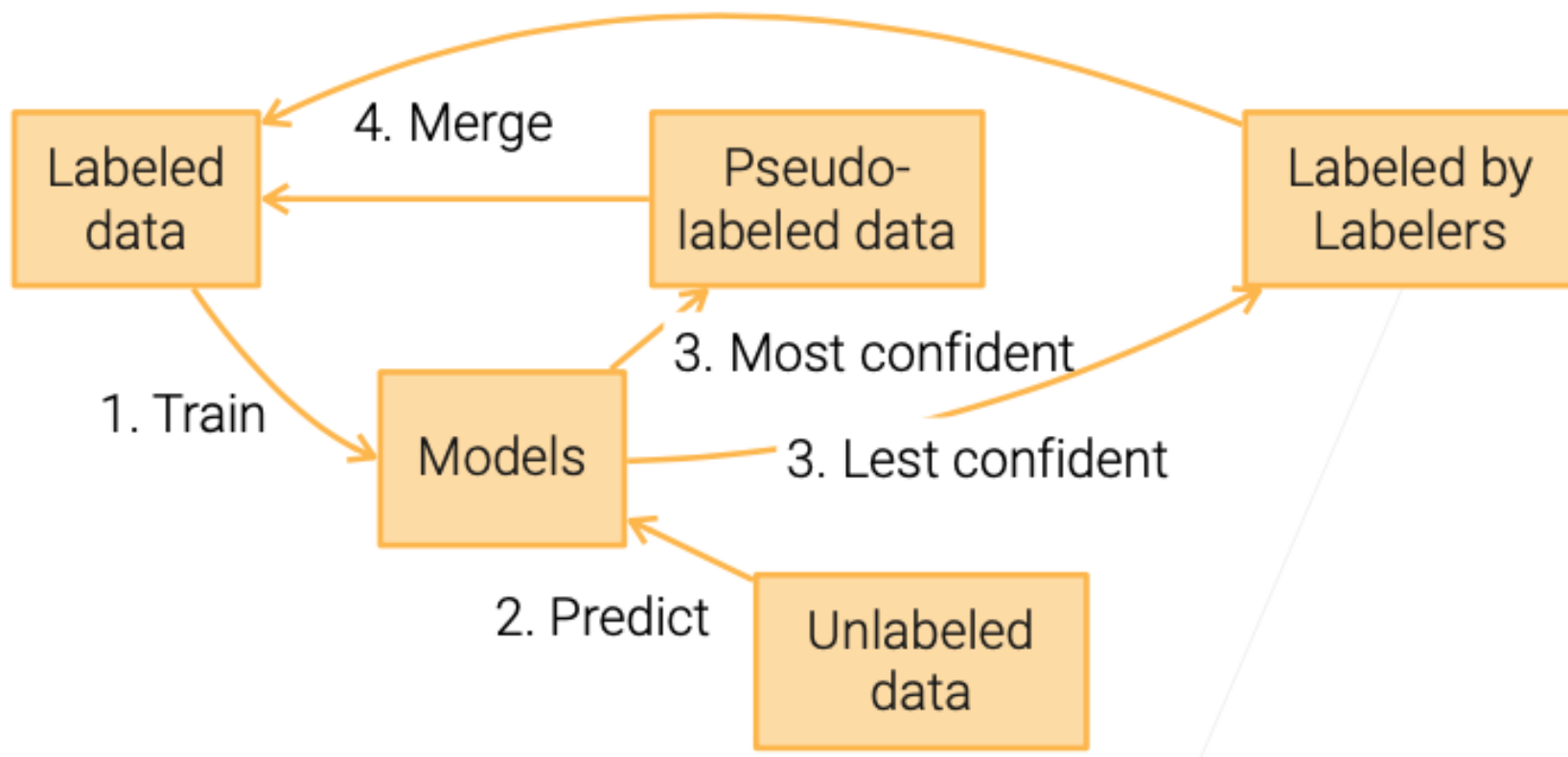
**Instruction** **Is this a cliff scene?** **Submit (799 images left)**  
**Definition:** a high, steep or overhanging face of rock.

**Yes**

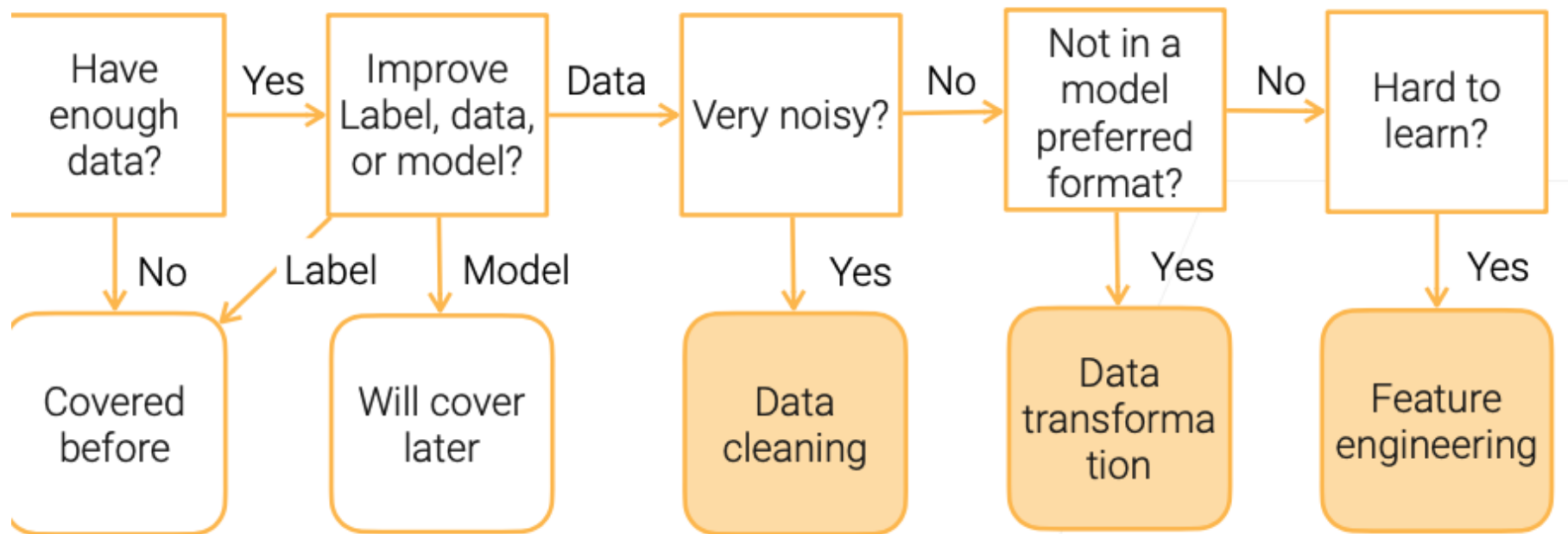
**No** **No**

# Active Learning + Self-training

- These two methods are often used together



# Flow chart for data preprocessing



# Web Scrapping

- The goal is to extract data from website
  - Noisy, weak labels, can be spammy
  - Available at scale
  - E.g. price comparison/tracking website
- Many ML datasets are obtained by web scrapping
  - E.g. ImageNet, Kinetics
- Web crawling VS scrapping
  - Crawling: indexing whole pages on Internet
  - Scrapping: scraping particular data from web pages of a website



Image credit: Aaron Zappia



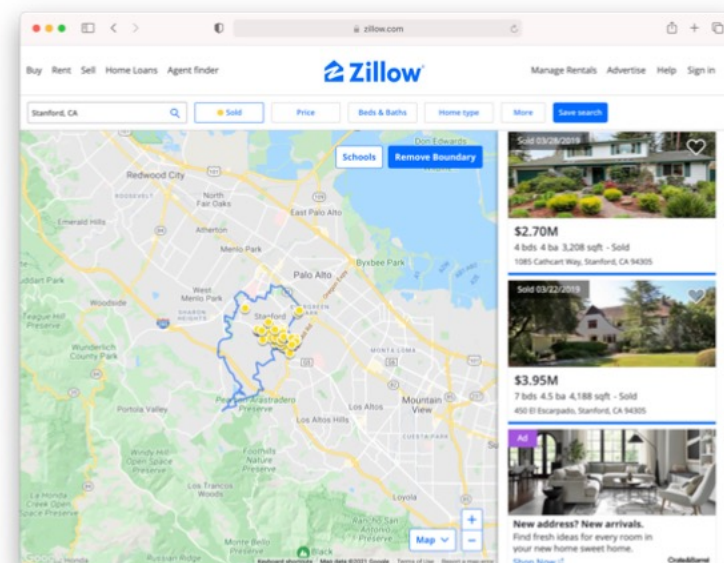
# Tools

- “curl” often doesn’t work
  - Website owners use various ways to stop bots
- Use headless browser: a web browser without a GUI
- You need a lot of new IPs, easy to get through public clouds
  - In all IPv4 IPs, AWS owns 1.75%, Azure 0.55%, GCP 0.25%

```
1 from selenium import webdriver
2
3 chrome_options = webdriver.ChromeOptions()
4 chrome_options.headless = True
5 chrome = webdriver.Chrome(
6     chrome_options=chrome_options)
7
8 page = chrome.get(url)
```

# Case Study

- Query houses sold in near Stanford
  - <https://www.zillow.com/stanford-ca/sold/>
  - <https://www.zillow.com/stanford-ca/sold/2-p/>
  - ...
- You can replace the city and state in the URL for other places

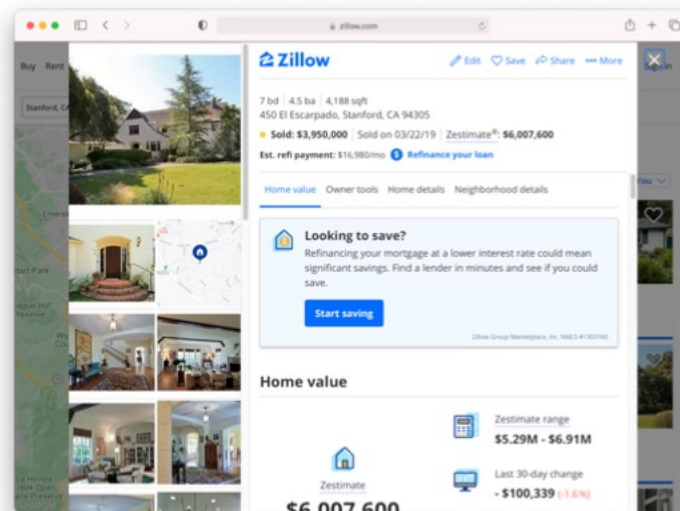


# Craw individual pages

- Get the house IDs from the index pages

```
1 page = BeautifulSoup(open(html_path, 'r'))
2 links = [a['href'] for a in page.find_all(
3     'a', 'list-card-link')]
4 ids = [l.split('/')[2].split('_')[0]
5     for l in links]
```

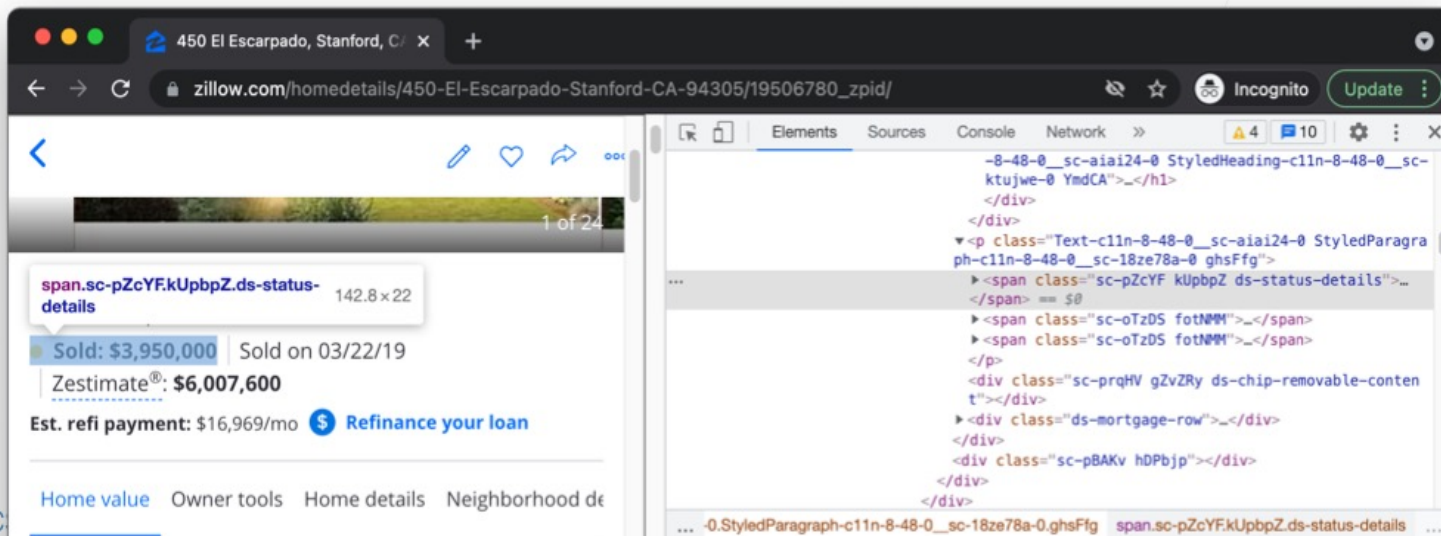
- The house detail page by ID
  - [https://www.zillow.com/homedetails/19506780\\_zpid/](https://www.zillow.com/homedetails/19506780_zpid/)



# Extract data

- Identify the HTML elements through **Inspect**

```
1 sold_items = [a.text for a in page.find(
2     'div', 'ds-home-details-chip').find('p').find_all('span')]
3 for item in sold_items:
4     if 'Sold:' in item:
5         result['Sold Price'] = item.split(' ')[1]
6     if 'Sold on' in item:
7         result['Sold On'] = item.split(' ')[-1]
```



Stanford CA

# Extract data

- Repeat the previous process to extract other field data

**Zillow** Save Share More

4 bd | 4 ba | 2,939 sqft  
44626 Sandia Creek Dr, Temecula, CA 92590  
Sold: \$960,000 | Sold on 02/11/21 | Zestimate®: \$1,119,900  
Est. refi payment: \$4,308/mo Refinance your loan

Home value Owner tools Home details Neighborhood details >

### Overview

This single level home on 9.66 acres land is beautifully maintained and has a lot of features you don't find in a typical tract home! there is a guest home of approx. 566 s.f. (included in the overall square footage) attached to a 4-car garage! great picturesque views of the hills and valleys of the De Luz Custom Home community in the hills to the West of the City of Temecula! Enjoy breathtaking sunsets off a big balcony that features trex decking, main house features 3 bed and 2.5 baths with wood flooring

Listed by:  
George Tektonopoulos DRE# 00999662  
Tekton Real Estate

Source: SDMLS, MLS#: SW21002332 SAN DIEGO |MLS

Zillow checked: September 09, 2021 at 10:32pm  
Data updated: February 12, 2021 at 12:45pm

Bought with: Marcie George  
First Team Real Estate

### Facts and features

**Type:** Detached **Parking:** 4 Garage spaces  
**Year built:** 1986 **Lot:** 9.66 Acres  
**Heating:** Forced Air **Buyer's Agent Fee:** 2.5%  
**Cooling:** Central Forced Air

### Interior details

**Bedrooms and bathrooms** Interior Features  
• Bedrooms: 4 • Interior features: Bedroom Entry

[See more facts and features](#)

### Services availability

Zillow Internet Resource Center  
[See Providers](#) | [Compare Speeds](#) | [Get Deals](#)

### Price and tax history

Price history



```
1 { 'Id': '18173100',
2   'Address': '44626 Sandia Creek Dr,',
3   'Sold Price': '$960,000',
4   'Sold On': '02/11/21',
5   'Summary': 'This single level home on 9.66 acres la
6   'Type': 'SingleFamily',
7   'Year built': '1986',
8   'HOA': '$70 annually',
9   'Lot': '420,790 sqft',
10  'Bedrooms': '4',
11  'Bathrooms': '4',
12  'Full bathrooms': '3',
13  '1/2 bathrooms': '1',
14  'Main level bathrooms': '3',
15  'Association name': 'De Luz Ranchose',
16  'Tax assessed value': '$784,118',
17  'Annual tax amount': '$12,036',
18  'Listed On': '1/5/2021',
19  'Listed Price': '$975,000',
20  'Last Sold On': '3/5/2008',
21  'Last Sold Price': '$840,000',
22  'Elementary School': 'Murrieta Elementary School',
23  'Elementary School Score': '7',
24  'Elementary School Distance': '4.9',
25  'Middle School': 'Thompson Middle School',
26  'Middle School Score': '6',
27  'Middle School Distance': '5.5',
28  'High School': 'Murrieta Valley High School',
29  'High School Score': '8',
```

Stanford CS 329P (2021 Fall) - <https://stanford.edu/cs329p/>

## Cost

- Use AWS EC2 t3.small (2GB memory, 2 vCPUs, \$0.02 per hour)
  - 2GB is necessary as the browser needs a lot memory, CPU and bandwidth are usually not an issue
  - Can use spot instance to reduce the price
- The cost to crawl 1M houses is \$16.6
  - The speed is about 3s per page,
  - 8.3 hours if using 100 instances
  - The extra cost includes storage, restart instances when IP is banned



# Crawl Images

- Get all image URLs

```
1 p = r'https:\\\\photos.zillowstatic.com\\fp\\([\\d\\w\\-\\_]+).jpg'  
2 ids = [a.split('-')[0] for a in re.findall(p, html)]  
3 urls = [f'https://photos.zillowstatic.com/fp/{id}-  
uncropped_scaled_within_1536_1152.jpg' for id in ids]
```

- A house listing has ~20 images
  - The crawling cost is still reasonable: ~\$300
  - Storing these images is expensive: ~\$300 per month
    - You can reduce the image resolutions, or send data back

# Legal Considerations

- Web scraping isn't illegal by itself
- But you should
  - NOT scrape data have sensitive information (E.g. private data involving username/password, personal health/medical information)
  - NOT scrape copyrighted data (E.g. YouTube videos, Flickr photos)
  - Follow the Terms of Service that explicitly prohibits web scraping
- Consult a lawyer if you are doing it for profit



# Data preprocessing

# Why data preprocessing?

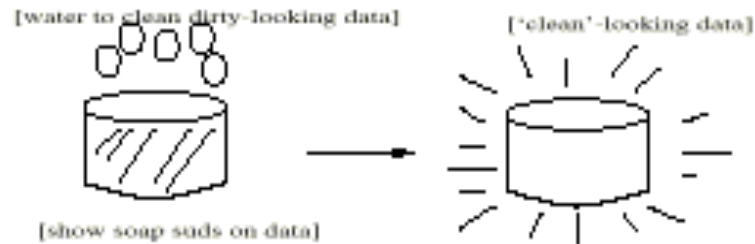
- Data in the real world is dirty
  - **incomplete**: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
  - **noisy**: containing errors or outliers
  - **inconsistent**: containing discrepancies in codes or names
- **No quality data, no quality mining results!**
  - Quality decisions must be based on quality data
- A multi-dimensional measure of data quality:
  - A well-accepted multi-dimensional view:
    - accuracy, completeness, consistency, timeliness, believability, value added, interpretability, accessibility
  - Broad categories:
    - intrinsic, contextual, representational, and accessibility.

# Major tasks in data preprocessing

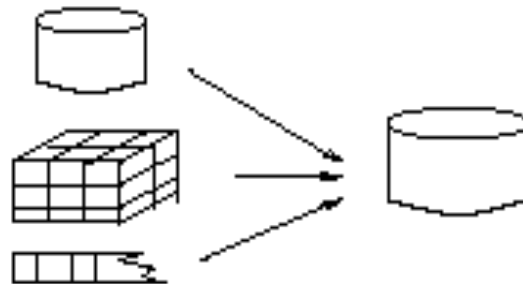
- Data cleaning
  - Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies
- Data integration
  - Integration of multiple databases, data cubes, files, or notes
- Data transformation
  - Normalization (scaling to a specific range)
  - Aggregation
- Data reduction
  - Obtains reduced representation in volume but produces the same or similar analytical results
  - Data discretization: with particular importance, especially for numerical data
  - Data aggregation, dimensionality reduction, data compression, and generalization

# Forms of data preprocessing

## Data Cleaning



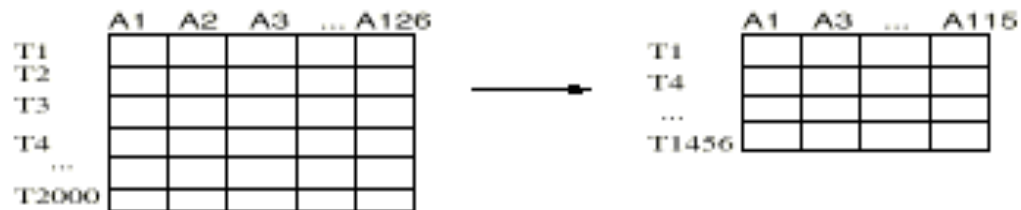
## Data Integration



## Data Transformation

-2, 32, 100, 59, 48 → -0.02, 0.32, 1.00, 0.59, 0.48

## Data Reduction



# Data Cleaning

- Data cleaning tasks
  - Fill in missing values
  - Identify outliers and smooth out noisy data
  - Correct inconsistent data

# Missing Data

- Data is not always available
  - E.g., many tuples have no recorded value for several attributes, such as **customer income in sales data**
- Missing data may be due to
  - equipment malfunction
  - inconsistent with other recorded data and thus deleted
  - data not entered due to misunderstanding
  - certain data may not be considered important at the time of entry
- Missing data may **need to be inferred**

# How to handle missing data?

- **Ignore the tuple:** usually done when class label is missing (assuming the task is classification—not effective in certain cases)
- **Fill in the missing value manually: tedious + infeasible?**
  - Use a **global constant to fill** in the missing value: e.g., “unknown”, a new class?!
  - Use the **attribute mean** to fill in the missing value
  - Use the **attribute mean for all samples of the same class to fill** in the missing value: smarter
  - Use the most **probable value to fill in the missing value**: inference-based such as regression, Bayesian formula, decision tree

# Noisy Data

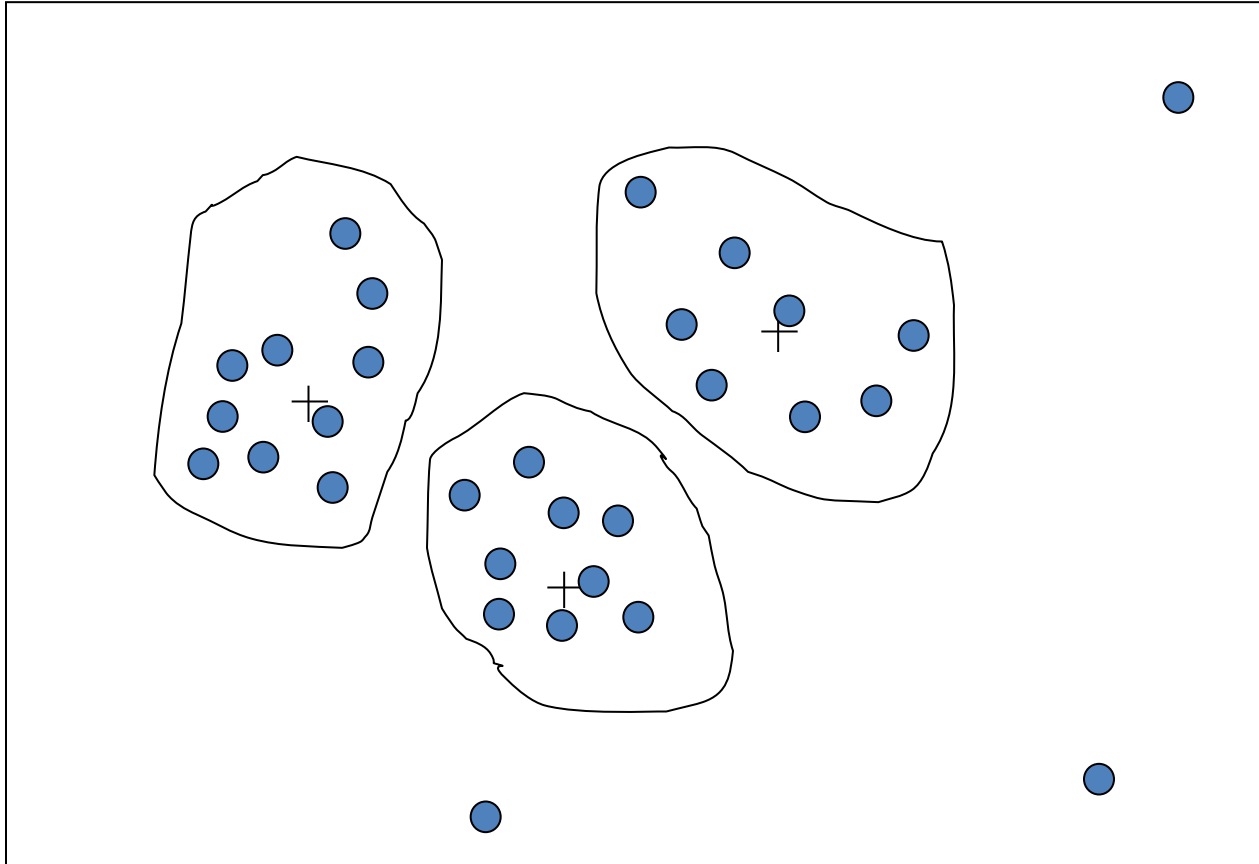
- Q: What is noise?
  - A: Random error in a measured variable.
  - Incorrect attribute values may be due to
    - faulty data collection instruments
    - data entry problems
    - data transmission problems
    - inconsistency in naming convention
  - Other data problems which requires data cleaning
    - duplicate records
    - incomplete data
    - inconsistent data



# How to handle noisy data?

- Binning method:
  - first sort data and partition into (equi-depth) bins
  - then one can smooth by bin means, smooth by bin median, smooth by bin boundaries, etc.
  - used also for discretization
- Clustering
  - detect and remove outliers
- Semi-automated method: combined computer and human inspection
  - detect suspicious values and check manually
- Regression
  - smooth by fitting the data into regression functions (next week)

# Cluster Analysis



# Data integration

- Data integration:
  - combines data from multiple sources into a coherent store
- Schema integration
  - integrate metadata from different sources
  - Entity identification problem: identify real world entities from multiple data sources, e.g.,  $A.cust-id \equiv B.cust-\#$
- Detecting and resolving data value conflicts
  - for the same real world entity, attribute values from different sources are different
  - possible reasons: different representations, different scales, e.g., metric vs. British units, different currency

# Data transformation

- **Normalization:** scaled to fall within a small, specified range
  - min-max normalization
  - z-score normalization
  - normalization by decimal scaling
- Attribute/feature construction
  - New attributes constructed from the given ones

# Data Transformation: Normalization

**Particularly useful for classification (Distance measurements, NN classification, etc)**

- min-max normalization

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new\_max}_A - \text{new\_min}_A) + \text{new\_min}_A$$

- z-score normalization

$$v' = \frac{v - \text{mean}_A}{\text{stand\_dev}_A}$$

- normalization by decimal scaling

$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$

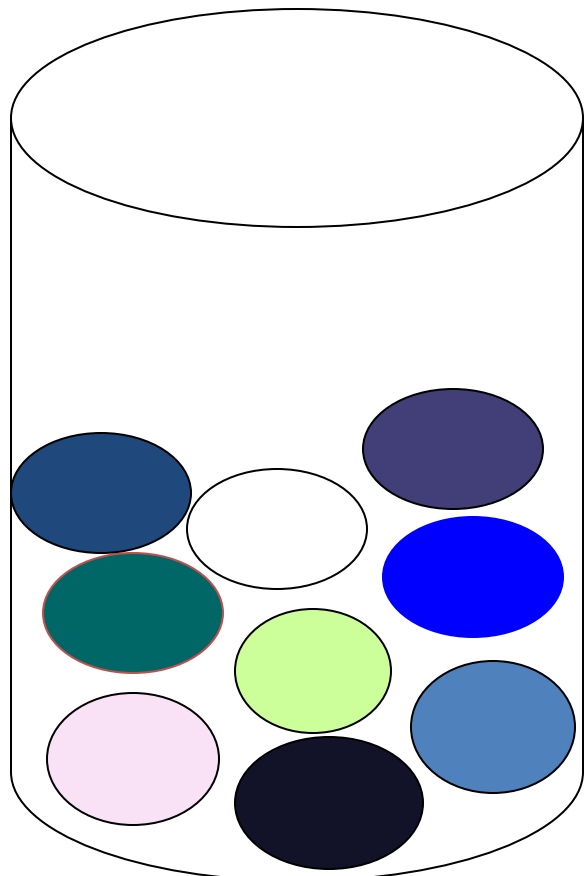
# Data reduction

- Obtains a reduced representation of the data set that is much smaller in volume but yet produces the same (or almost the same) analytical results
- Data reduction strategies
  - Dimensionality reduction
  - Data compression
  - Numerosity reduction
  - Discretization and concept hierarchy generation

# Sampling

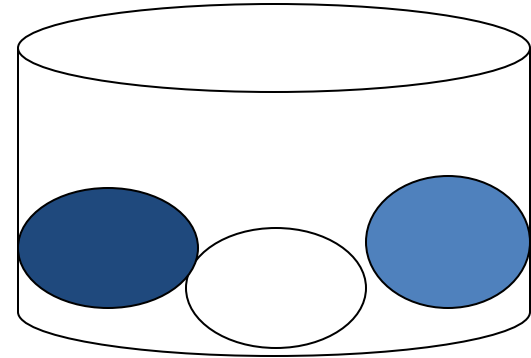
- Allow a mining algorithm to run in complexity that is potentially sub-linear to the size of the data
- Cost of sampling: proportional to the size of the sample, increases linearly with the number of dimensions
- Choose a representative subset of the data
  - Simple random sampling may have very poor performance in the presence of skew
- **Develop adaptive sampling methods**
  - Stratified sampling:
    - Approximate the percentage of each class (or subpopulation of interest) in the overall database
    - Used in conjunction with skewed data
- Sampling may not reduce database I/Os (page at a time).

# Sampling

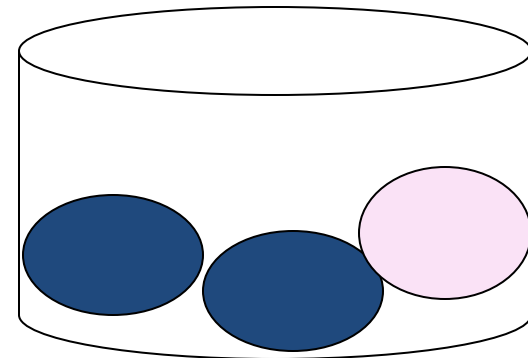


Raw Data

SRSWOR  
(simple random  
sample without  
replacement)



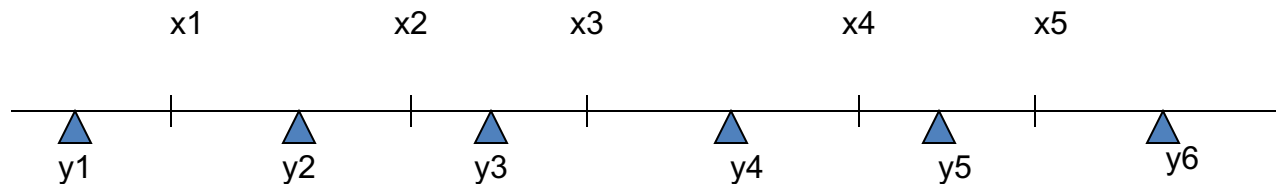
SRSWR





# Discretization/Quantization

- Three types of attributes:
  - Nominal — values from an unordered set
  - Ordinal — values from an ordered set
  - Continuous — real numbers
- Discretization/Quantization:
  - ☞ divide the range of a continuous attribute into intervals



- Some classification algorithms only accept categorical attributes.
- Reduce data size by discretization
- Prepare for further analysis

# Summary

- Data acquisition and preparation
- Data preprocessing includes
  - Data cleaning and data integration
  - Data reduction and feature selection
  - Discretization
- What's next:
  - Model validation and evaluation