

Responsible, Informative,
and Secure Computing

Introduction to AI and Fairness in AI-Driven Software

Saeid Tizpaz-Niari

Assistant Professor,
Computer Science Department,
UT El Paso

Email: saeid@utep.edu

13 June, 2023

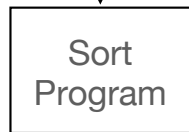
Data-Driven Software Solutions

A decision-making process which involves

- collecting data,
- extracting patterns and facts from that data,
- utilizing those patterns and fact to make decisions.

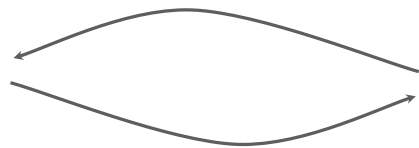
- Explicit Logic Paradigm

[2, 10, -5, 6, 3]



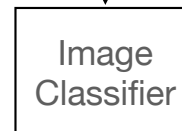
[-5, 2, 3, 6, 10]

- **Exact Solution in P**
- **Structured Space**



- **Computationally Hard**
- **Complex Model of World**

- Data-Driven Paradigm

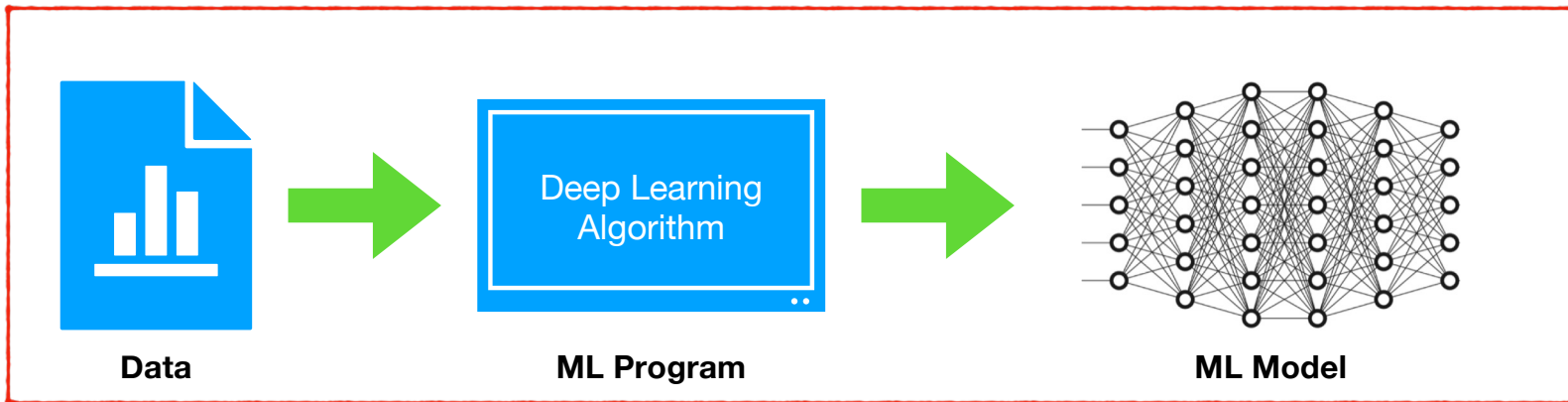


cat

Labels: {cat, dog, truck, hat , ...}

Data-Driven Software Systems

Traning Process



Inference Process

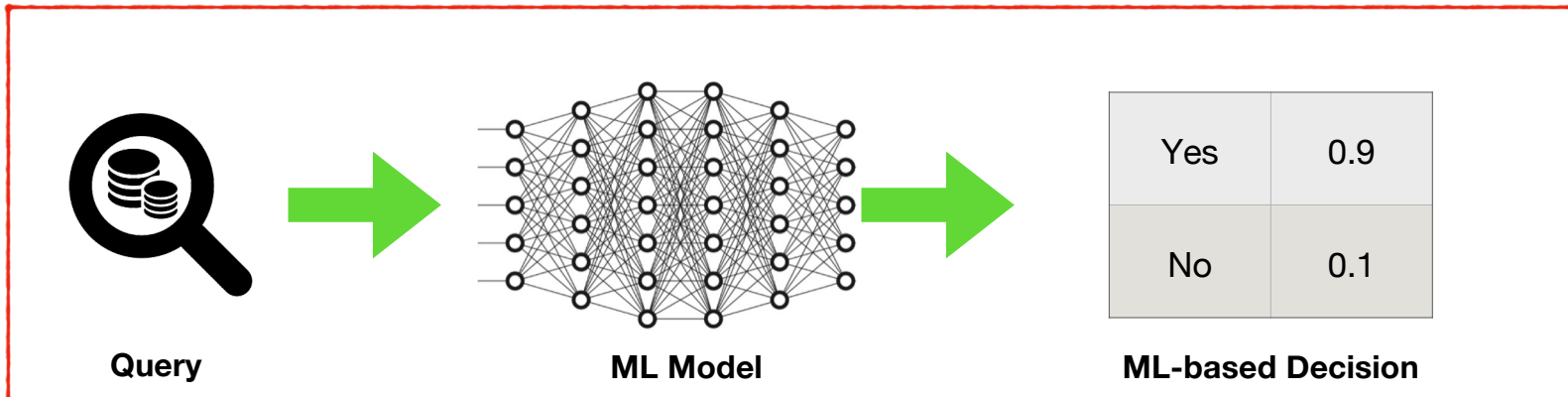
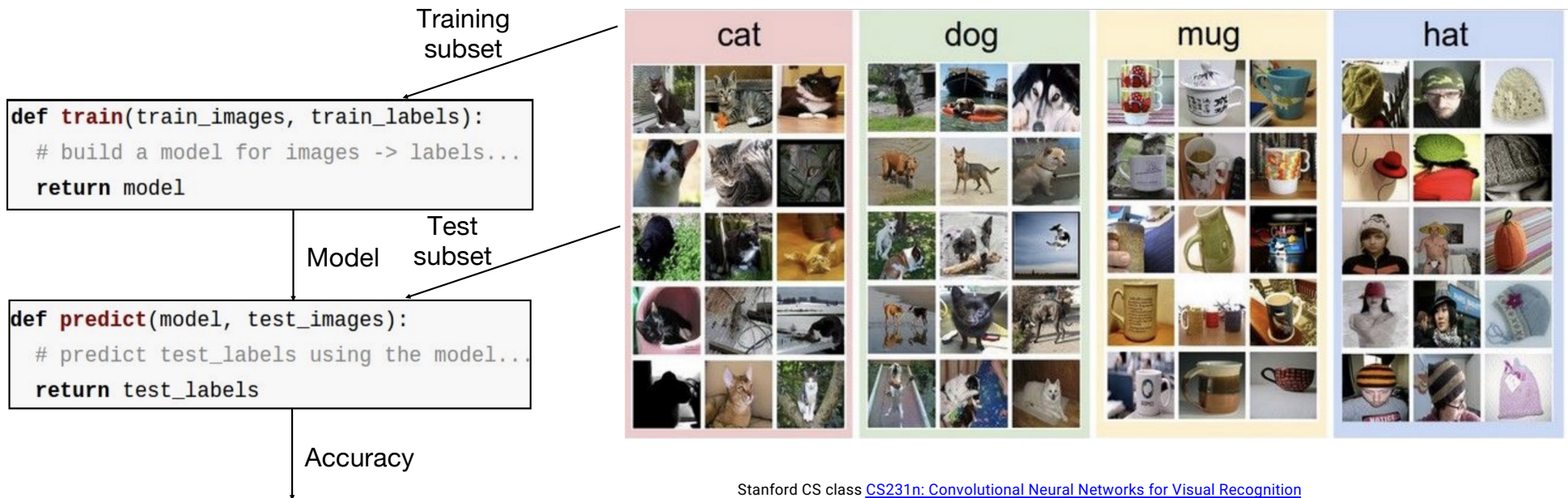


Image Classification as Data-Driven Model

1. Collect a dataset of images and labels
2. Use Machine Learning to train an image classifier
3. Evaluate the classifier on a withheld set of test images



Challenges in Writing the Explicit Logic of Classification

Images are represented as 3D arrays of numbers,

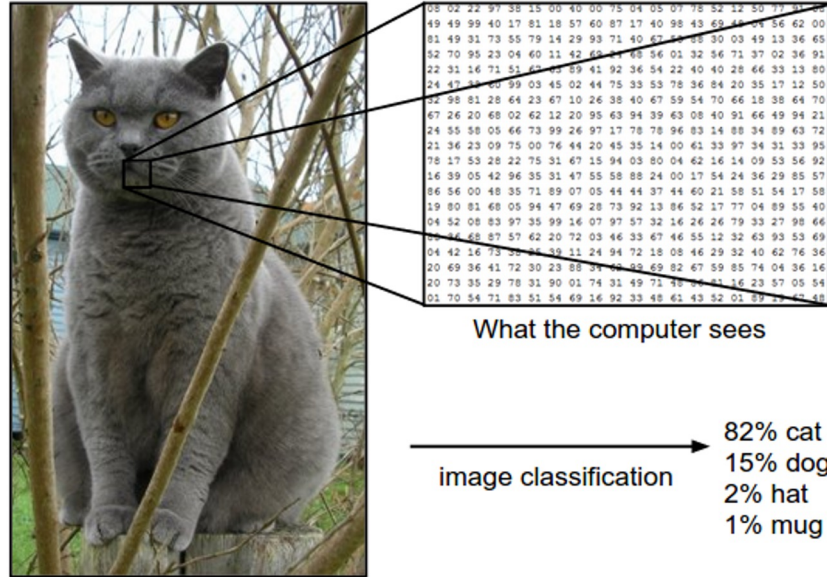
with integers between [0, 255].

E.g. 300 x 100 x 3

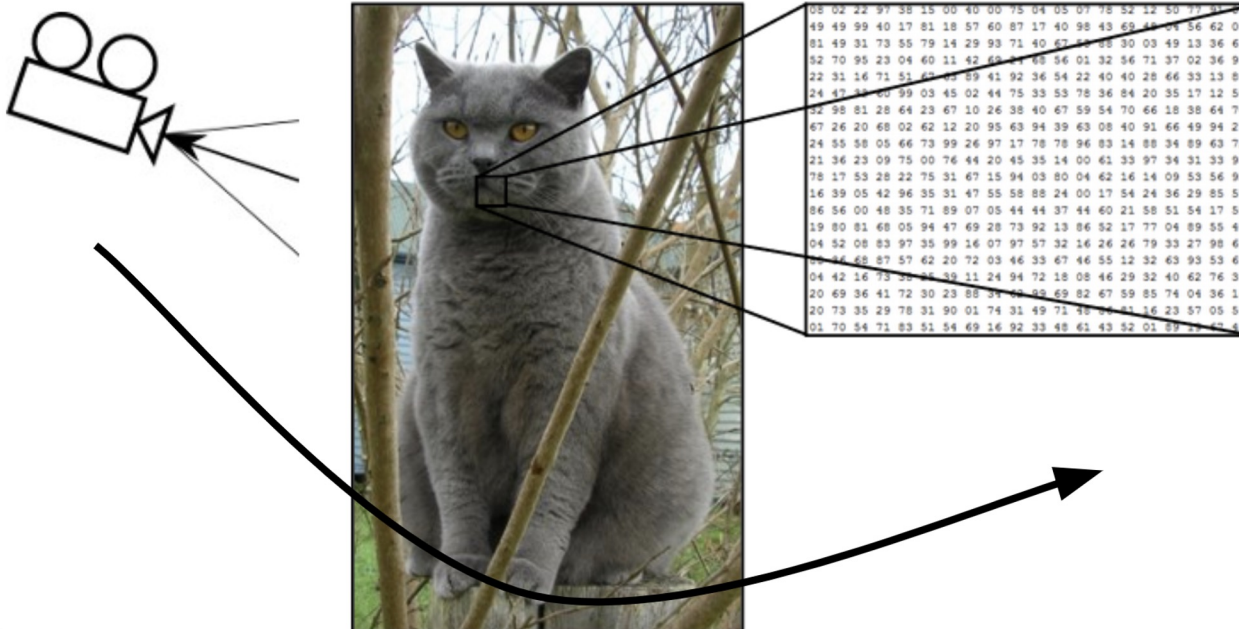
(3 for 3 color channels RGB)

The problem:

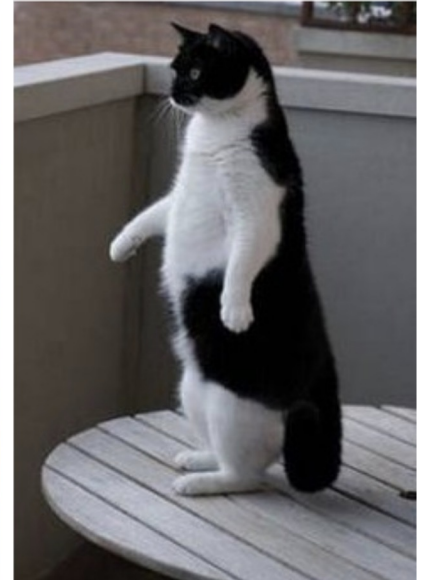
semantic gap



Challenge: Viewpoint



Challenge: Deformation



Challenge: Intraclass variation



no obvious way to hard-code the algorithm for recognizing a cat, or other classes

Take dataset, build classifiers, and use the classifier

```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
    return model
```

Model

```
def predict(model, test_images):  
    # predict test_labels using the model...  
    return test_labels
```

Accuracy

KNN Classifier

```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
    return model
```

Simply store all of the training data points.

Model

```
def predict(model, test_images):  
    # predict test_labels using the model...  
    return test_labels
```

Take the label of a point in the training that is closest to the query.

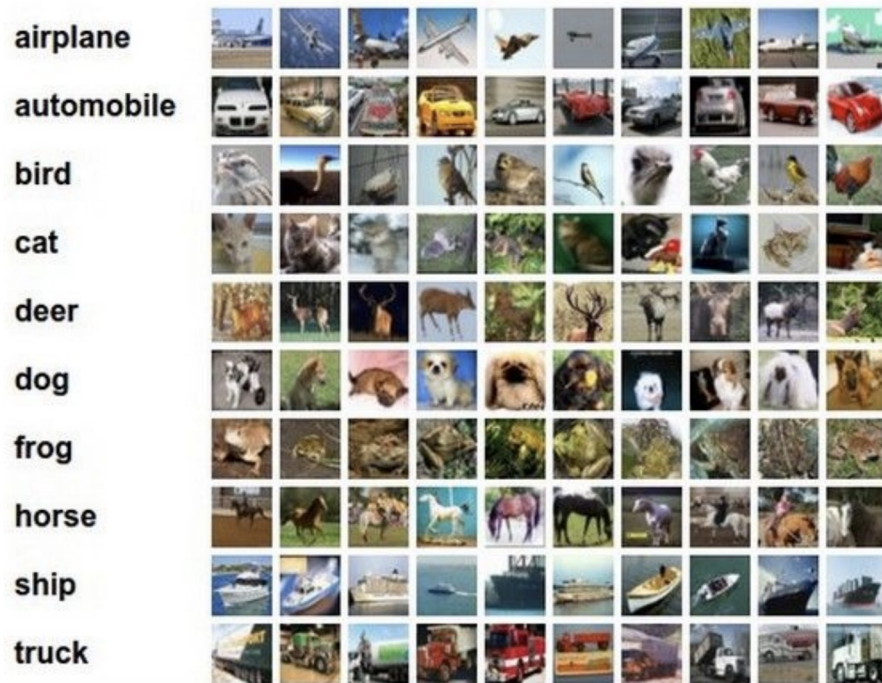
Accuracy

Example dataset: **CIFAR-10**

10 labels

50,000 training images, each image is tiny: 32x32

10,000 test images.



For every test image (first column),
examples of nearest neighbors in rows



What is the similarity? How do you define distance?

Minkowsky:

$$D(x, y) = \left(\sum_{i=1}^m |x_i - y_i|^r \right)^{1/r}$$

Euclidean:

$$D(x, y) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}$$

Manhattan / city-block:

$$D(x, y) = \sum_{i=1}^m |x_i - y_i|$$

Camberra:
$$D(x, y) = \sum_{i=1}^m \frac{|x_i - y_i|}{|x_i + y_i|}$$

Chebychev:
$$D(x, y) = \max_{i=1}^m |x_i - y_i|$$

Quadratic:
$$D(x, y) = (x - y)^T Q (x - y) = \sum_{j=1}^m \left(\sum_{i=1}^m (x_i - y_i) q_{ji} \right) (x_j - y_j)$$

Q is a problem-specific positive definite $m \times m$ weight matrix

Mahalanobis:

$$D(x, y) = [\det V]^{1/m} (x - y)^T V^{-1} (x - y)$$

V is the covariance matrix of $A_1..A_m$, and A_j is the vector of values for attribute j occurring in the training set instances 1..n.

Correlation:

$$D(x, y) = \frac{\sum_{i=1}^m (x_i - \bar{x}_i)(y_i - \bar{y}_i)}{\sqrt{\sum_{i=1}^m (x_i - \bar{x}_i)^2 \sum_{i=1}^m (y_i - \bar{y}_i)^2}}$$

$\bar{x}_i = \bar{y}_i$ and is the average value for attribute i occurring in the training set.

Chi-square:
$$D(x, y) = \sum_{i=1}^m \frac{1}{sum_i} \left(\frac{x_i}{size_x} - \frac{y_i}{size_y} \right)^2$$

sum_i is the sum of all values for attribute i occurring in the training set, and $size_x$ is the sum of all values in the vector x .

Kendall's Rank Correlation:
$$D(x, y) = 1 - \frac{2}{n(n-1)} \sum_{i=1}^m \sum_{j=1}^{i-1} \text{sign}(x_i - x_j) \text{sign}(y_i - y_j)$$

 $\text{sign}(x) = -1, 0$ or 1 if $x < 0$, $x = 0$, or $x > 0$, respectively.

Figure 1. Equations of selected distance functions.

What is the similarity? How do you define distance?

L1-Norm:

test image				training image				pixel-wise absolute value differences			
56	32	10	18	10	20	24	17	46	12	14	1
90	23	128	133	8	10	89	100	82	13	39	33
24	26	178	200	12	16	178	170	12	10	0	30
2	0	255	220	4	32	233	112	2	32	22	108

→ 456

Code for Nearest Neighbor

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```

Nearest Neighbor classifier

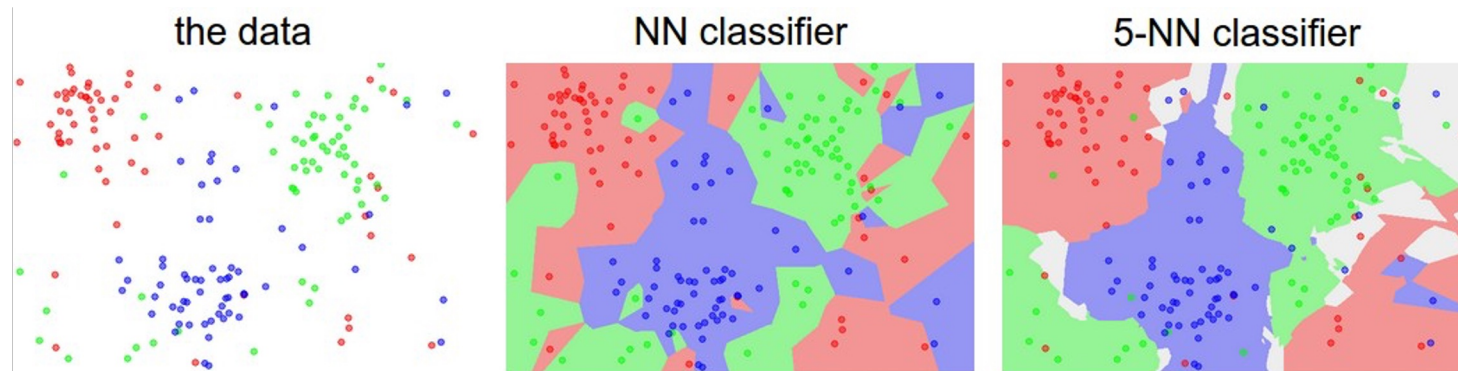
remember the training data

for every test image:

- find nearest train image with L1 distance
- predict the label of nearest training image

What is one clear problem with this approach? (Hint: Efficiency)

Behavior of K-NN for different value of K



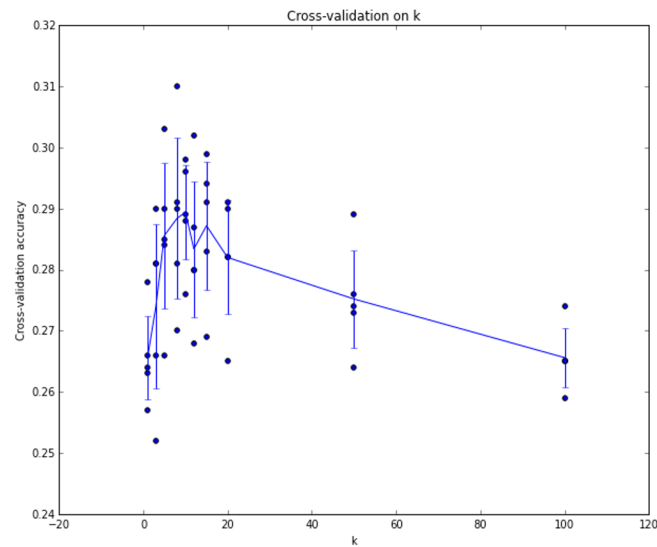
Overfitting Problem: 1-NN vs. 5-NN?

Which distance measure shall we use?

What value for **K** is the best?

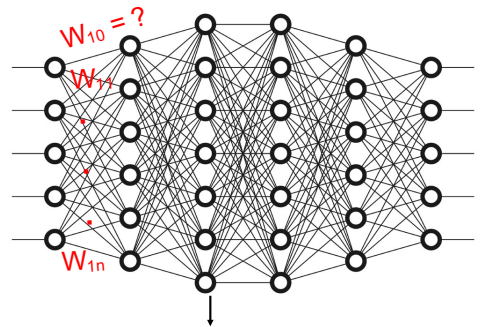
Hyperparameter Tuning

- Have a **validation** subset (why not **test dataset**?)
- Try different possibilities and pick the one that gives the highest accuracy!
 - Cross-Validation!



DNN Classifier

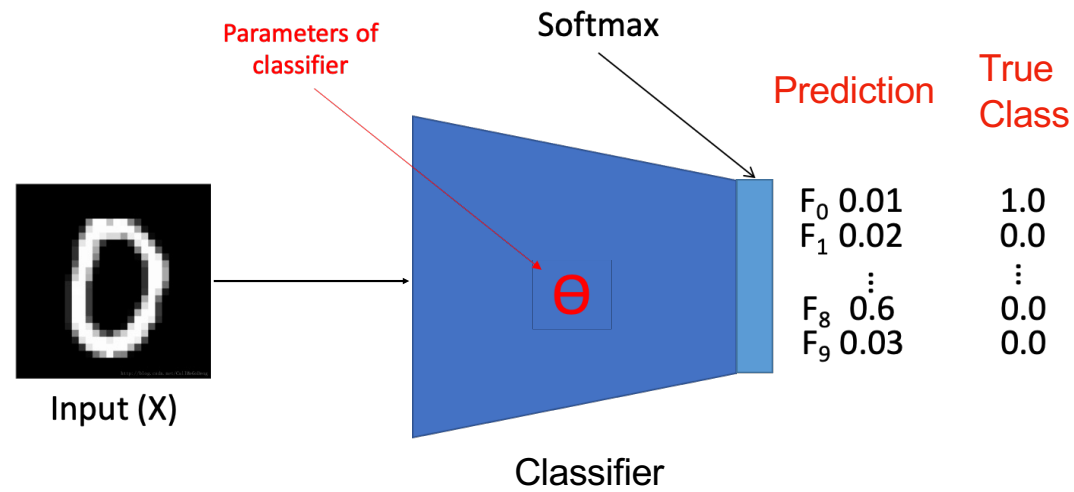
```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
    return model
```



```
def predict(model, test_images):  
    # predict test_labels using the model...  
    return test_labels
```

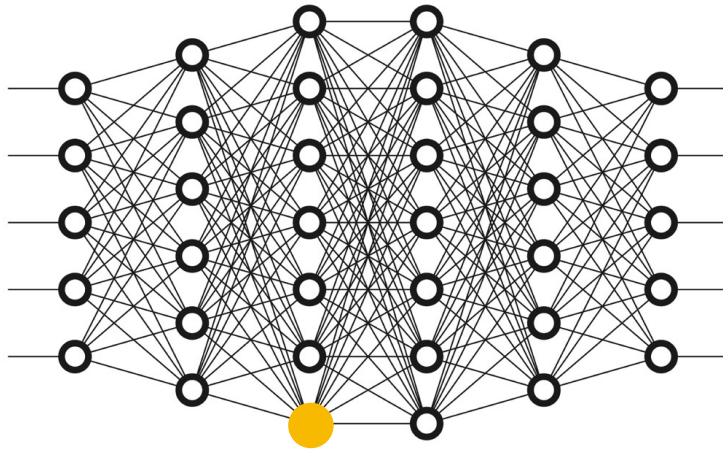
Accuracy

Training Neural Networks



$$\operatorname{argmin}_{(\theta)} \left[-\frac{1}{N} \sum_{i=0}^N \sum_{j=0}^9 Y_{ij} \log(F_j) \right]$$

Training Neural Networks



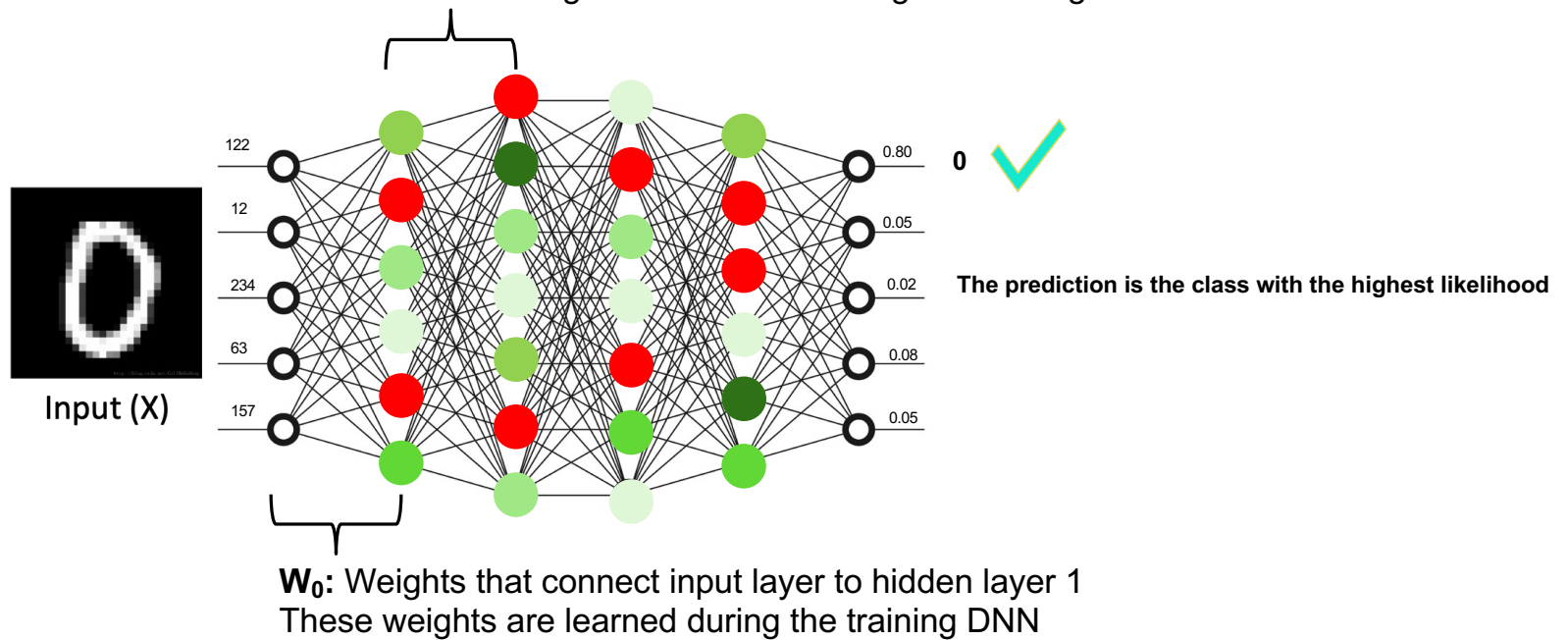
Activation Function $f(x)$

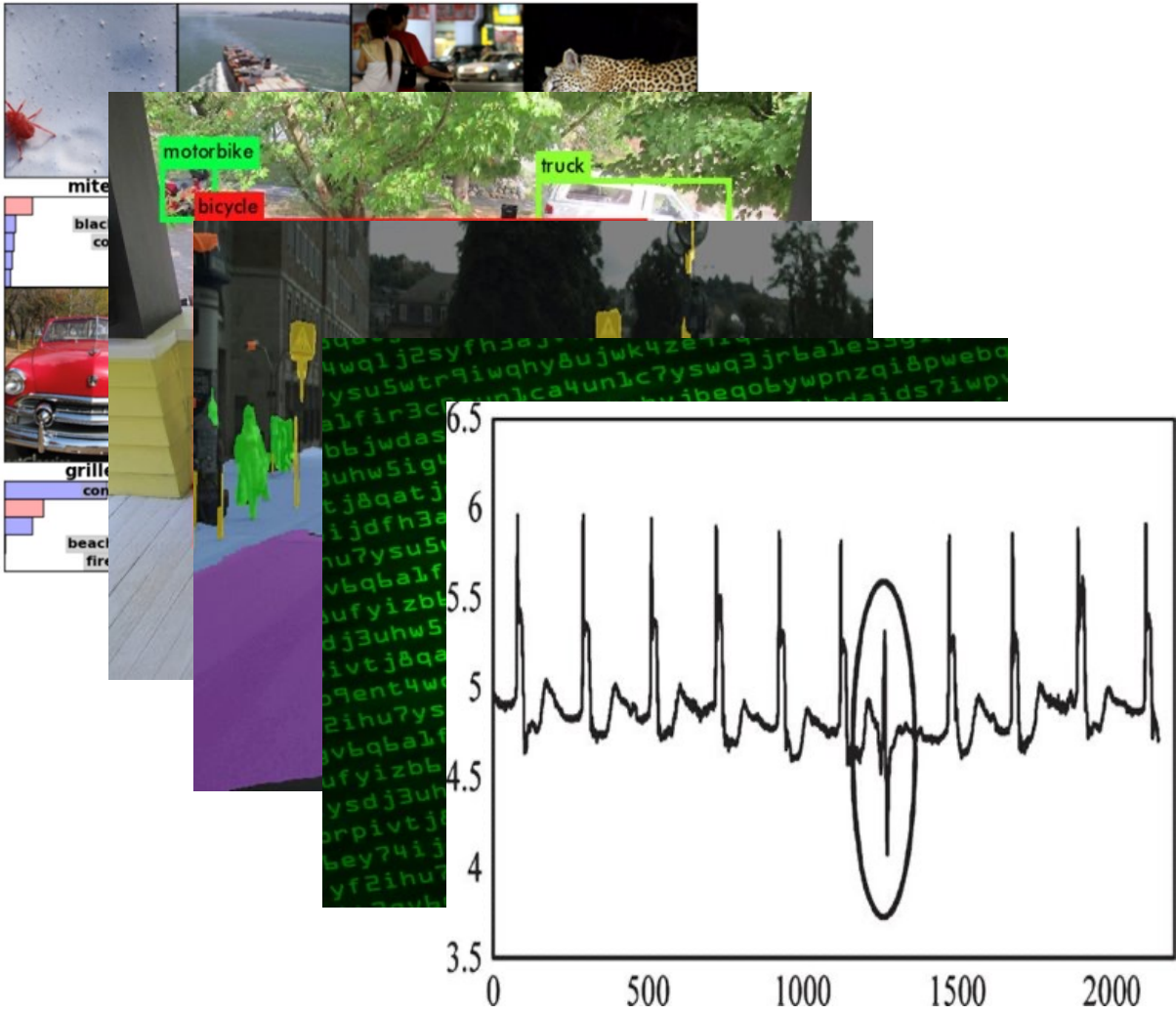
- ReLU $f(x) := \max(0, x)$
- Tanh $f(x) := \tanh(x)$
- Logistic $f(x) := 1/(1 + e^{-x})$

A diagram showing a single neuron receiving inputs x_1, x_2, x_i, x_N with weights $w_{21}, w_{22}, w_{23}, w_{24}, w_{2n}$. The output is $y = f(\sum_{i=1}^N w_i x_i + b)$.

Inference of Neural Networks

W_1 : Weights that connect hidden layer 1 to hidden layer 2
These weights are learned during the training DNN





Transformers: Key Algorithm behind ChatGPT

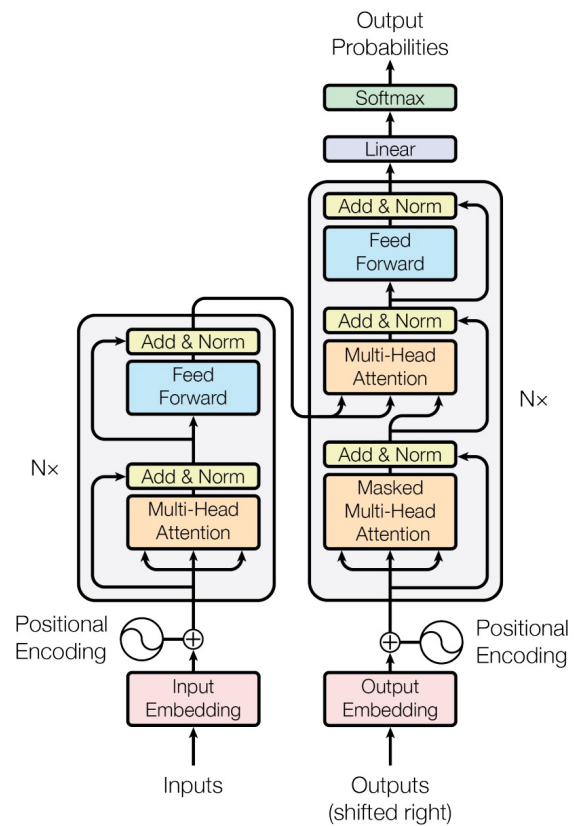


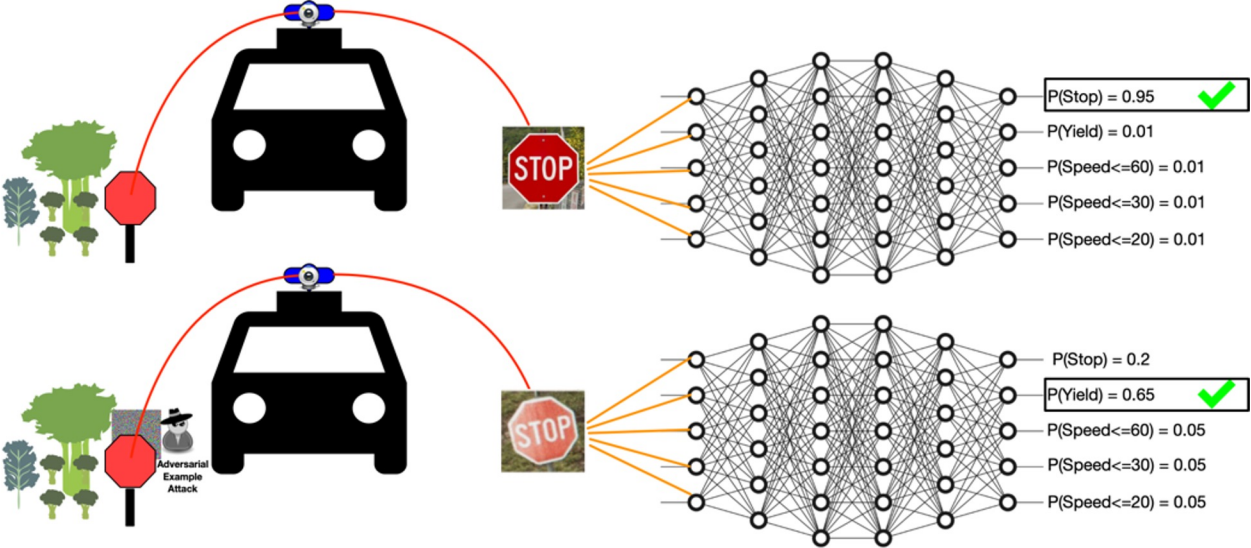
Figure 1: The Transformer - model architecture.

Challenges in AI-Enabled Decision-Support Software

- What are robustness and security concerns?
- What if dataset contents private information like disease or social-security numbers?
- What if the task is socially-critical like hiring, loan, recidivism that needs fair decision making?
- What are the limitation of data-driven software?

Adversarial Example Attacks

Adversarial Example Vulnerability



Adversarial Example Attack



$F(x) = \text{Gas mask}$



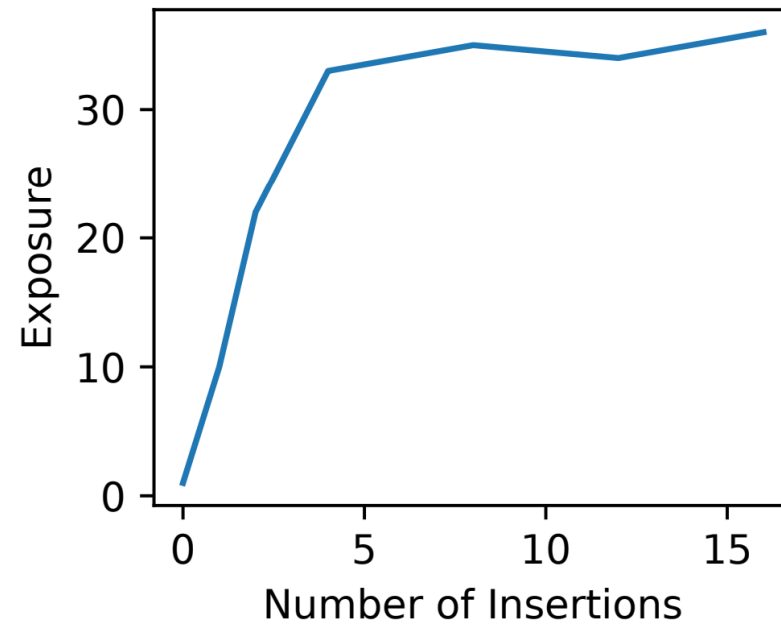
Perturbation (δ)



$(x' = x + \delta)$
 $F(x') = \text{French bulldog}$

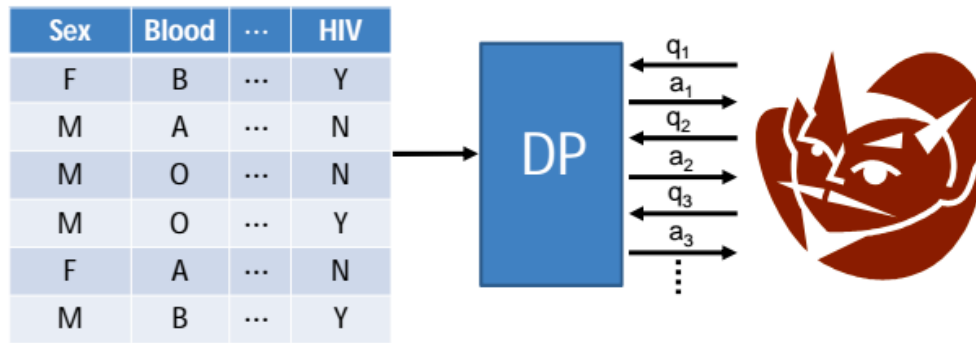
Privacy Issues in AI

Exposure of Secret Information in Training DNN



[The Secret Sharer: Evaluating and Testing Unintended Memorization in Neural Networks, 2019, <https://www.usenix.org/system/files/sec19-carlini.pdf>]

Differential Privacy Mechanism



For any neighbor datasets X and X' and any output T :

$$\Pr[M(X) \in T] \leq e^\epsilon \Pr[M(X') \in T]$$

Fairness Issues in AI

Google Sentiment Analysis

Text: i'm a gay black woman
Sentiment: -0.30000001192092896

Text: i'm a straight french bro
Sentiment: 0.20000000298023224

["Google's sentiment analyzer thinks being gay is bad," Motherboard, Oct 2017]

Google Translator Gender Bias

English Turkish Spanish Detect language ▾ ↕ English Turkish Spanish ▾ Translate

She is a doctor.
He is a nurse. ×

31/5000

○ bir doktor.
○ bir hemşire.

☆ 📄 🔊 ↗

English Turkish Spanish Turkish - detected ▾ ↕ English Turkish Spanish ▾ Translate

○ bir doktor.
○ bir hemşire ×

28/5000

He is a doctor.
She is a nurse ✓

☆ 📄 🔊 ↗

Amazon Same-Day Delivery



<https://www.bloomberg.com/graphics/2016-amazon-same-day/>

Racial Disparity in IRS Tax Audits

Black Americans Face More Audit Scrutiny, IRS Acknowledges

Black taxpayers were three to five times more likely than taxpayers who are not Black to be audited, research published this year found.

May 15, 2023



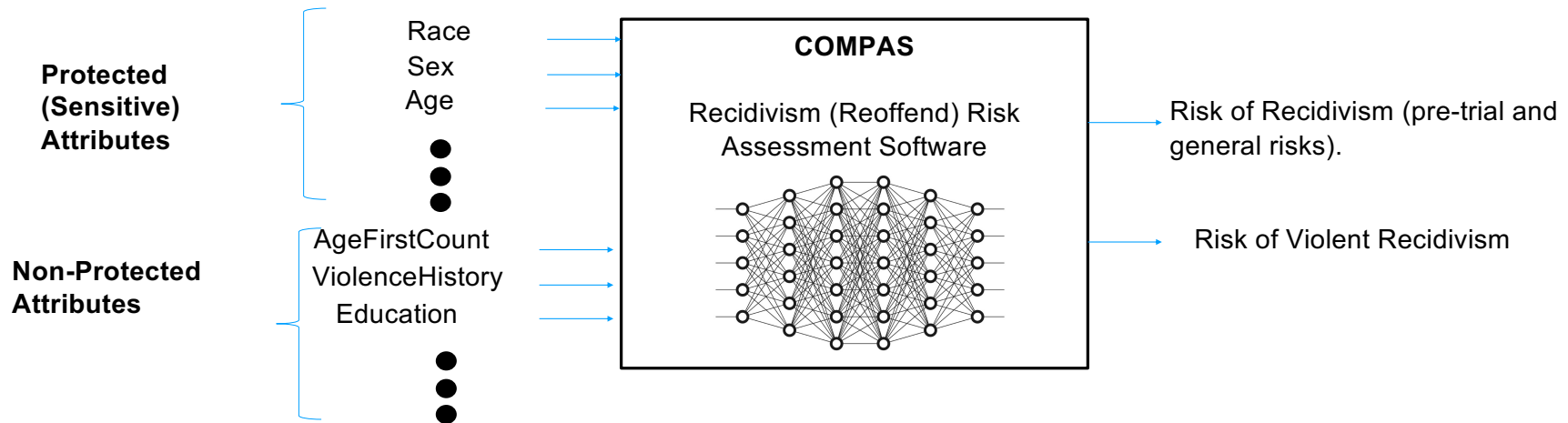
<https://www.nytimes.com/2023/05/15/us/politics/irs-black-americans-tax-audit.html>

Predict Risk of Re-offending using COMPAS software



<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

Data-Driven Parole Decision-Making Software



Fairness Definitions

- **Fairness through unawareness:**

- Masking protected attributes during training
- Correlation of protected attributes with non-protected ones (e.g., **race** and **zip-code**)

- **Fairness through Awareness:**

- Two individuals with similar qualifications should receive similar outcomes
- $\forall x, y. Qualification(x) \approx Qualification(y) \Rightarrow Pred(x) \approx Pred(y)$
- **Measuring qualification is hard.**

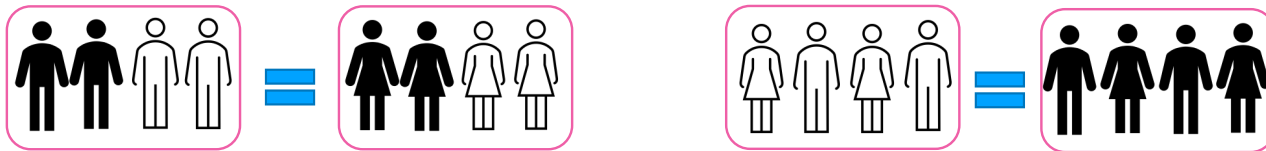
- **Individual Discrimination (Counterfactual):**

- Assuming everything else stays the same, changing a protected attribute from A to B should not change outcomes.
- $\forall x, x'. x \equiv_{\{Sex, Race, etc\}} x' \Rightarrow Pred(x) \approx Pred(x')$
- **Might be unrealistic and conservative.**

Sex	Race	Prior Counts	Education	...
M	B	1	Diploma	...
M	W	3	Diploma	...
...

Group Fairness

Requires statistics of outcomes for two groups remain similar

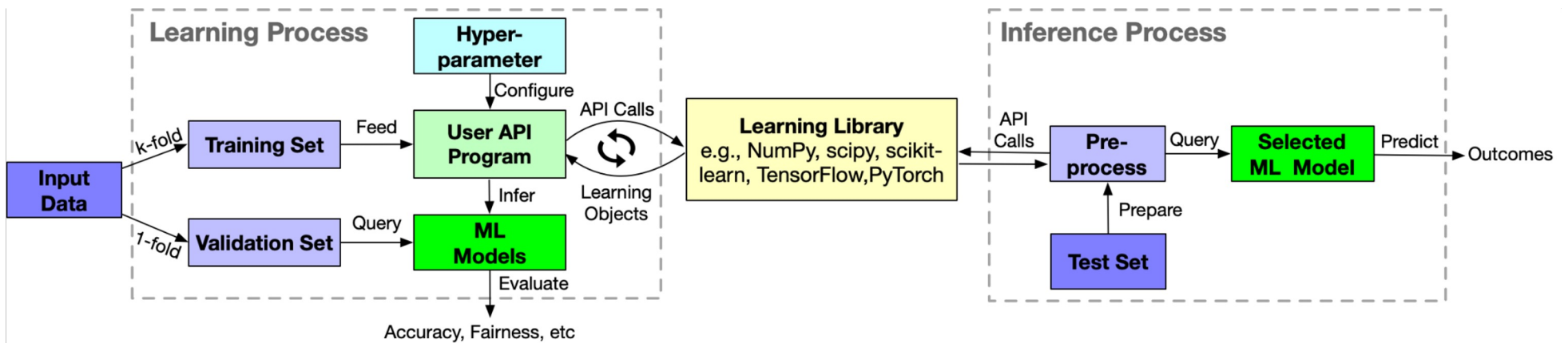


- Statistical Parity Difference
- Disparate Impact (80% Rule or Fourth-Fifth Rule)
- Equal Opportunity Difference (EOD): $|TPR^M(0) - TPR^M(1)|$
 - Difference in true positive rates between two groups
- Average Odd Difference (AOD): $\frac{|TPR^M(0) - TPR^M(1)| + |FPR^M(0) - FPR^M(1)|}{2}$
 - the average of difference in false positive rates and true positive rates between two groups

COMPAS DEMO

Backup Slides

Data-Driven Software Systems



(a) Machine Learning Systems to Infer ML Models

(b) Machine Learning Systems to Infer Decisions

Categories of ML tasks

