

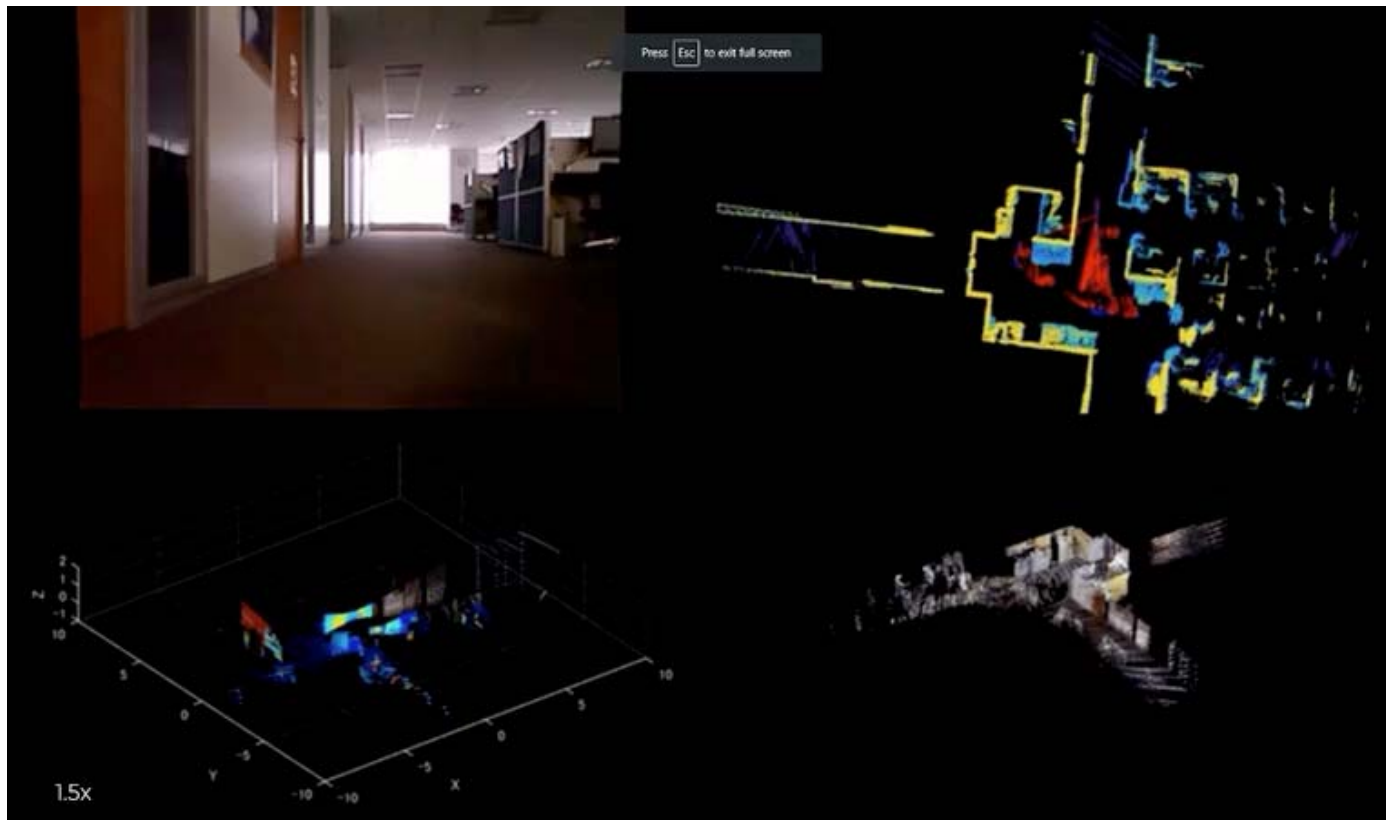
CMPE 185 Autonomous Mobile Robots

Perception: Learning Based Object Classification and Detection

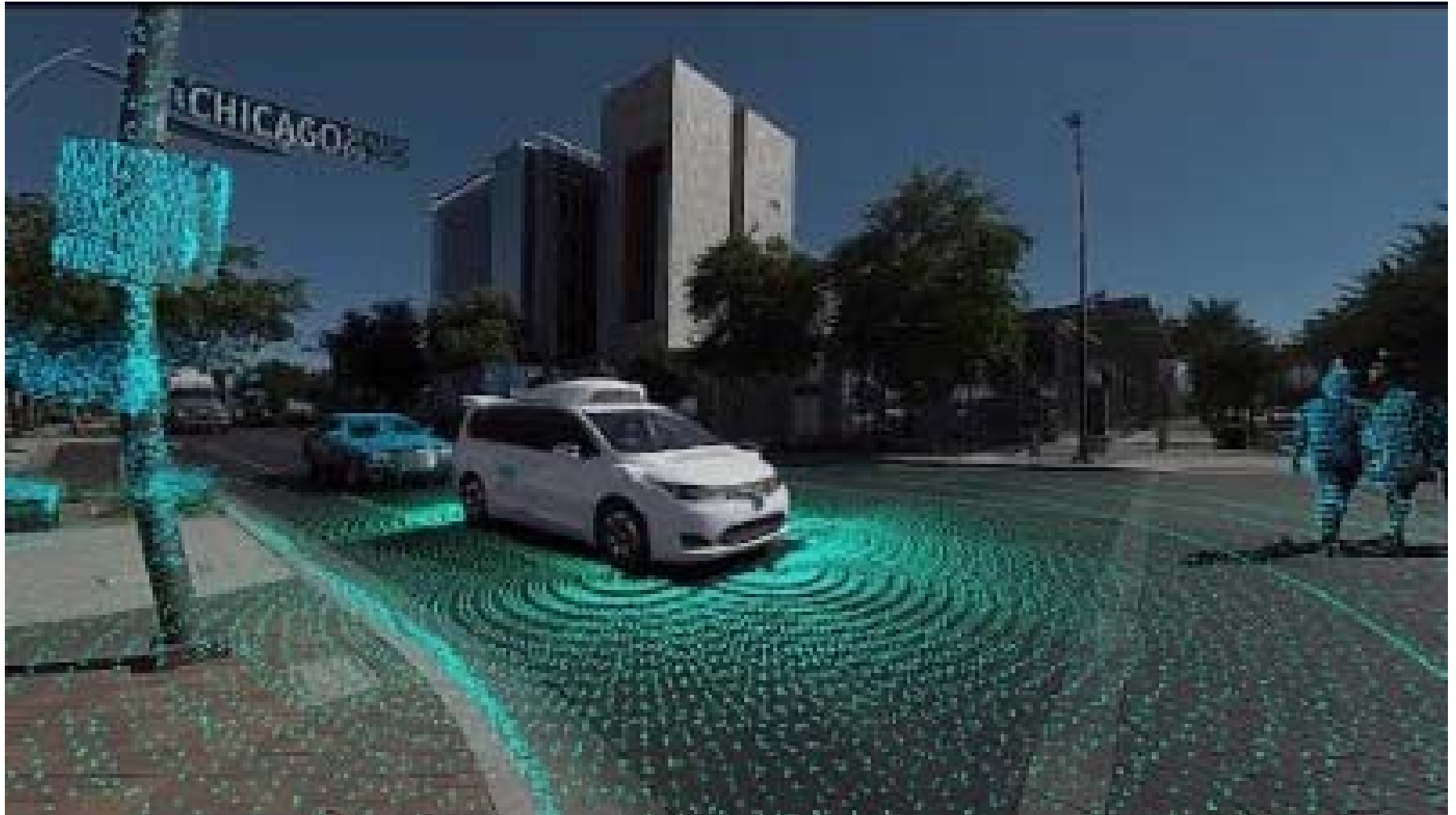
Dr. Wencen Wu
Computer Engineering Department
San Jose State University

Computer Vision

- **Enable robot vision to build environment maps and localize your mobile robot**



Example: Waymo Experience – Sensor Fusion

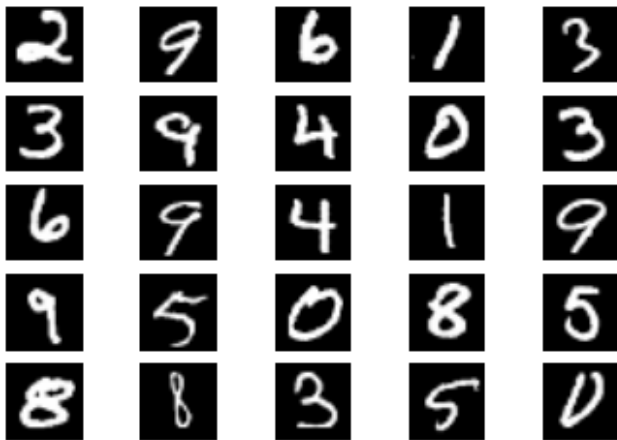


Example: Autonomous Parking



Image Classification

- K classes
- Task: assign correct class label to the whole image



Digit classification (MNIST)



Object recognition (Caltech-101)

Classification v.s. Detection

Classification



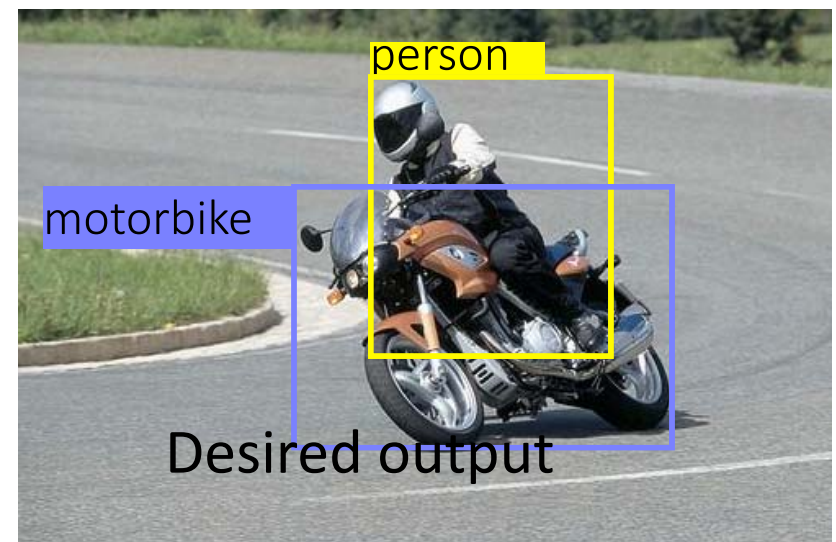
Detection



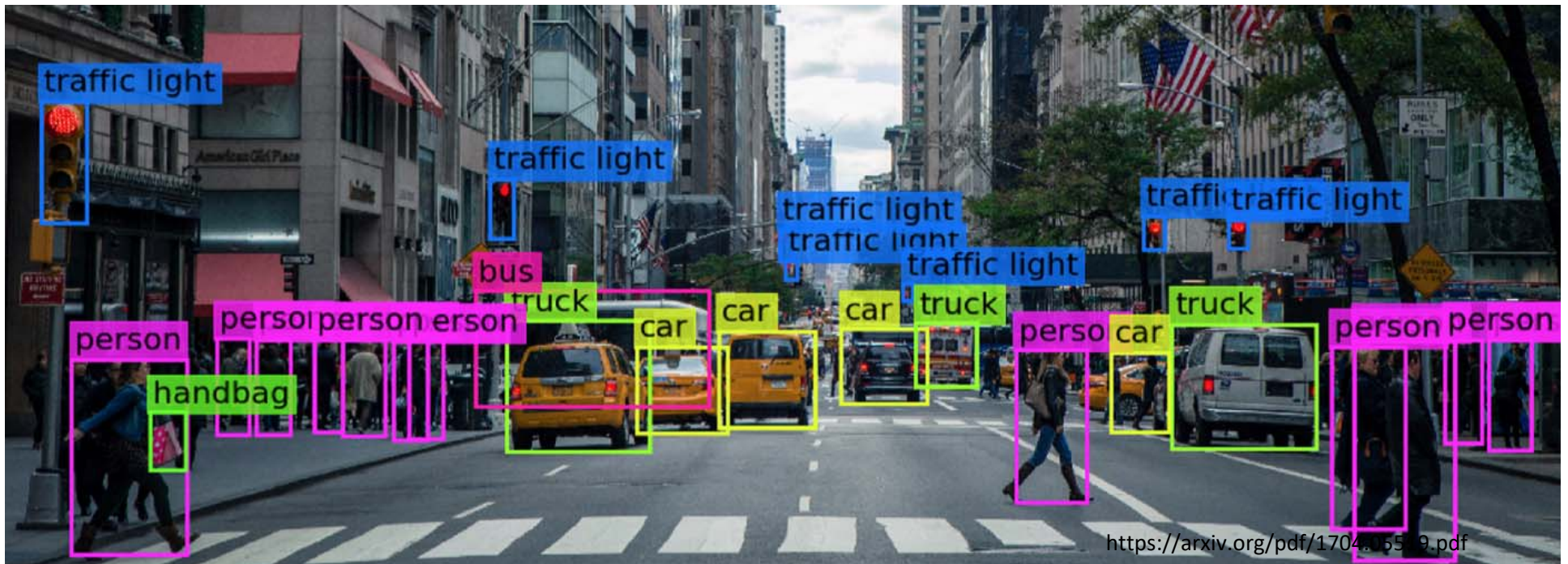
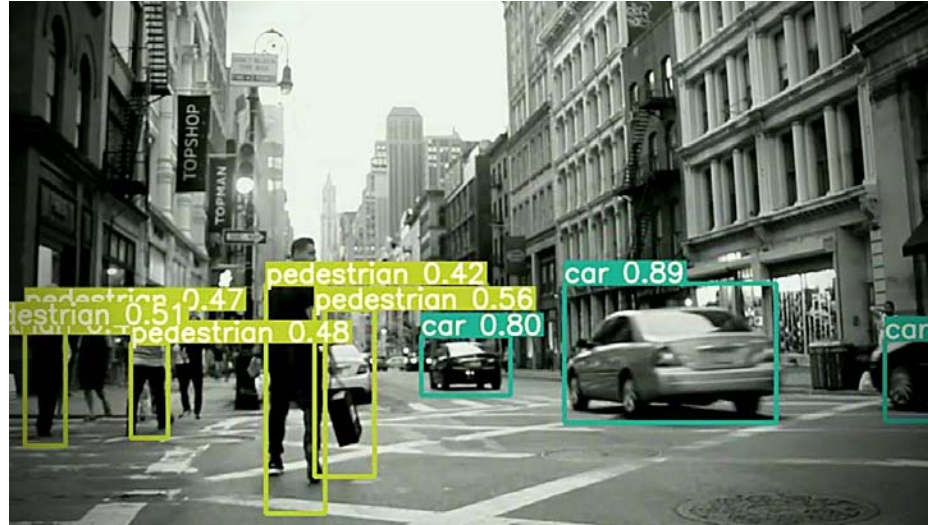
Problem Formulation

- **When performing object detection, we wish to obtain:**
 - A **list of bounding boxes**, or the (x, y) -coordinates for each object in an image
 - The **class label** associated with each bounding box
 - The **probability/confidence score** associated with each bounding box and class

{ airplane, bird, motorbike, person, sofa }



Object Detection in Autonomous Driving



Evaluating a Detector



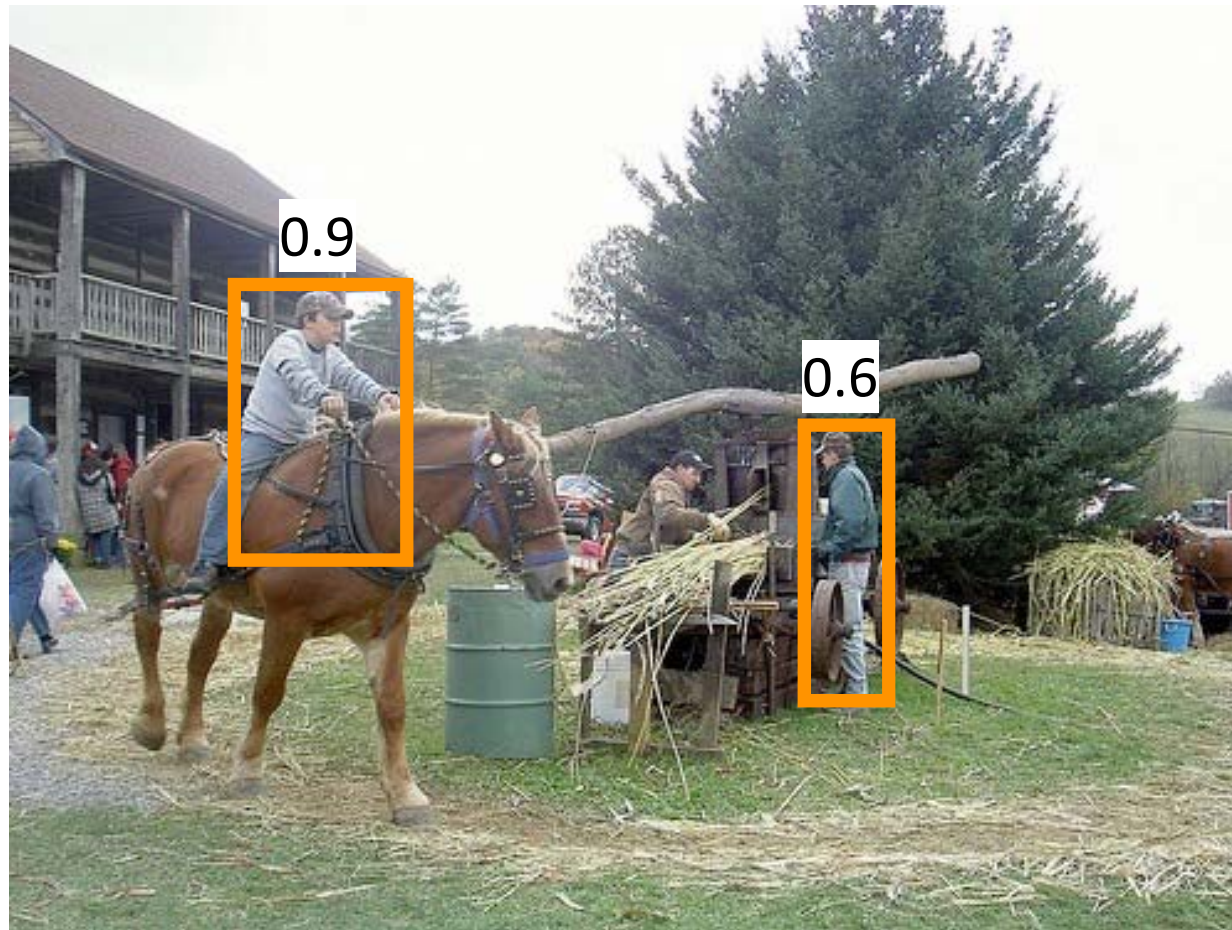
Test image (previously unseen)

First Detection



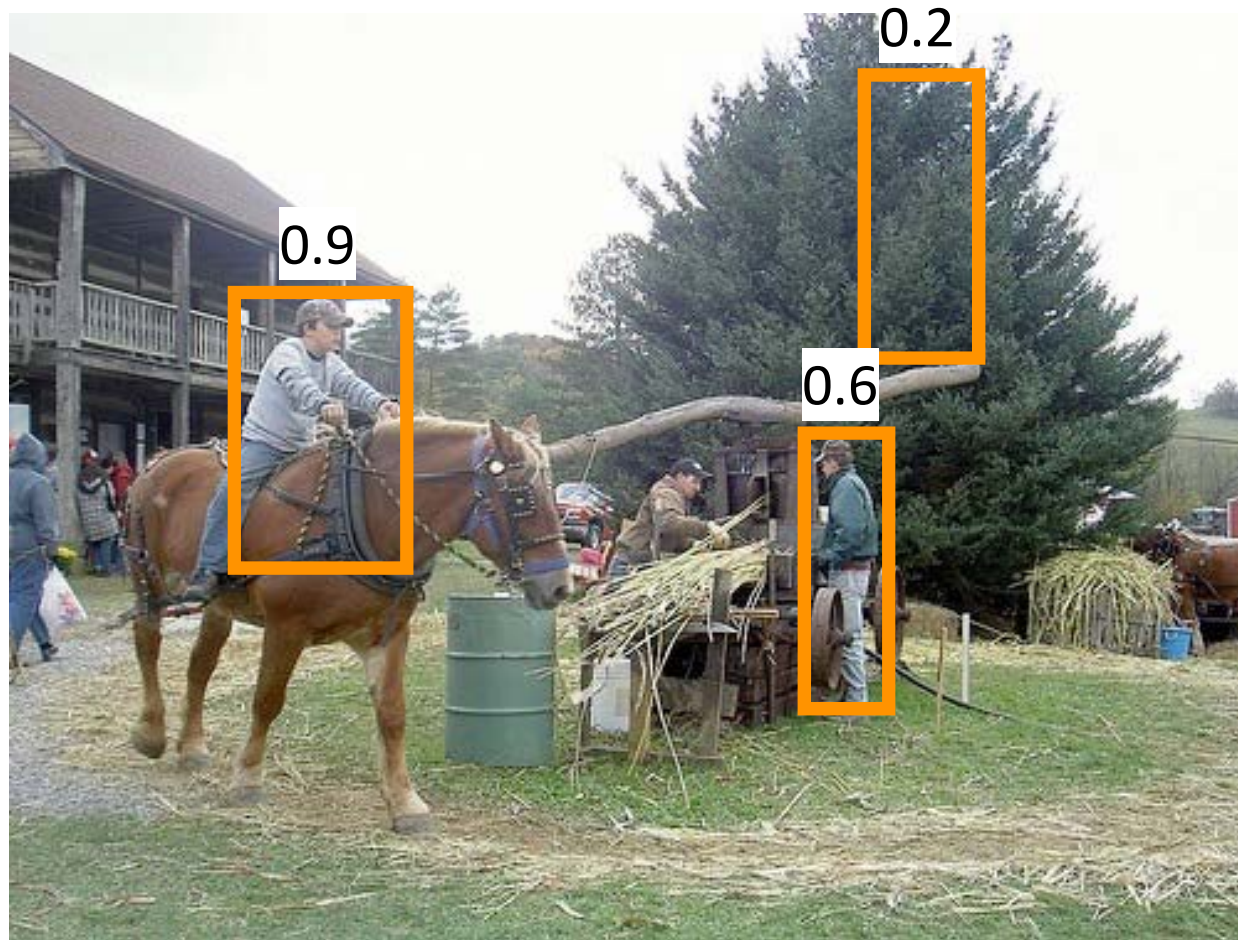
 'person' detector predictions

Second Detection



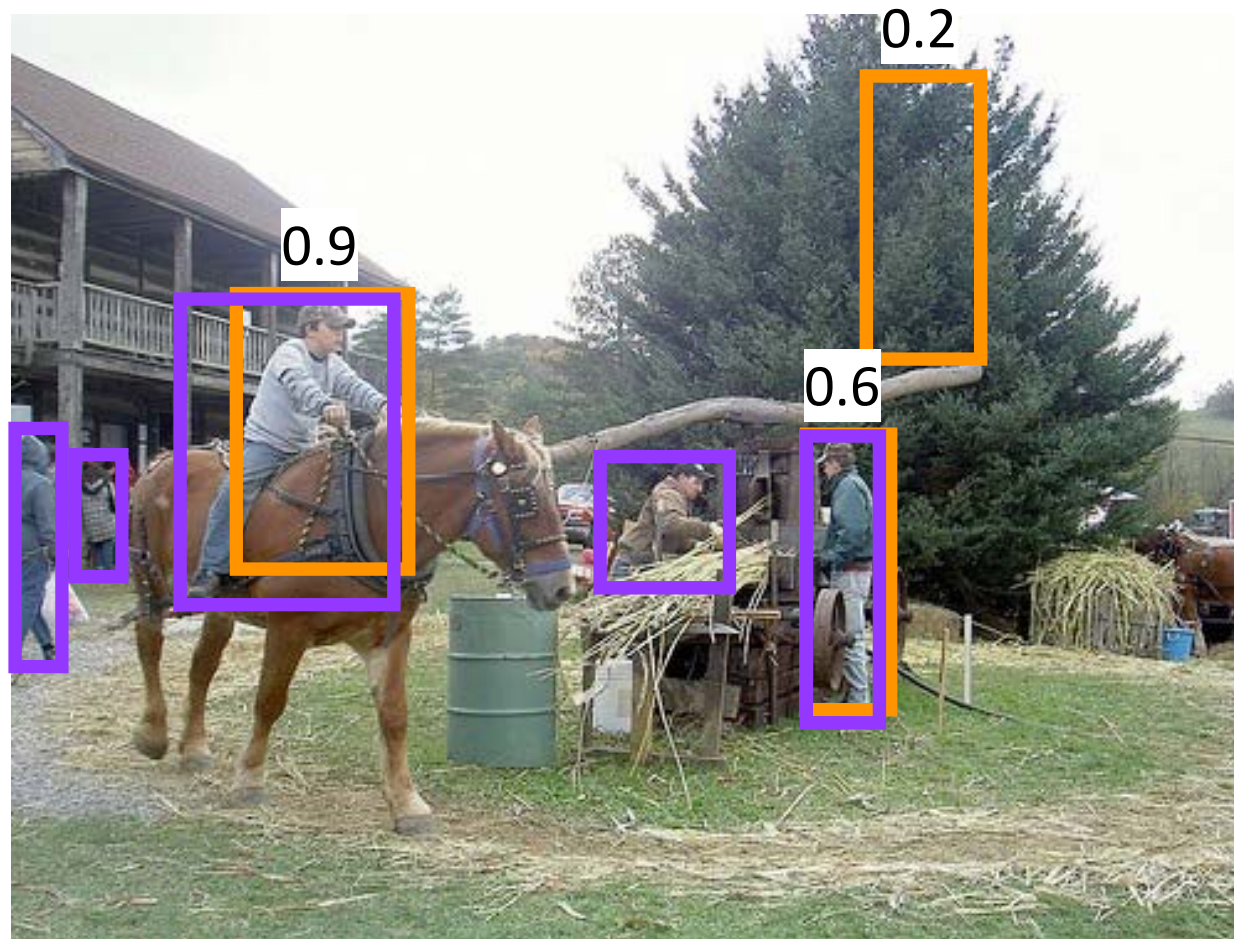
 'person' detector predictions

Third Detection



 'person' detector predictions

Compare to Ground Truth








'person' detector predictions

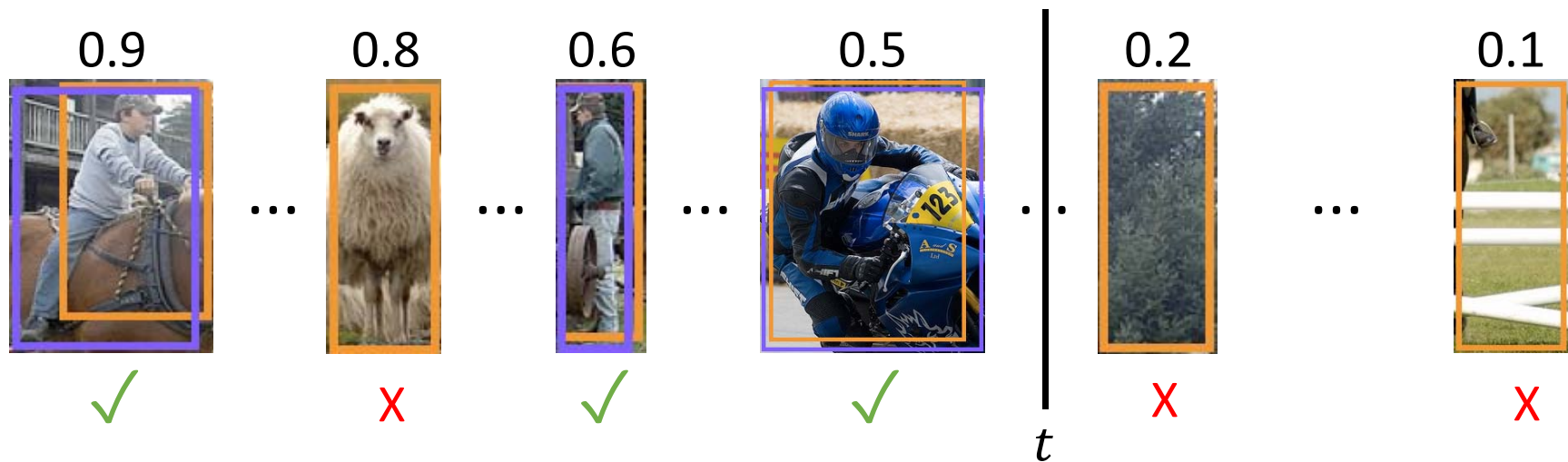


ground truth 'person' boxes

Sort by Confidence

0.9	...	0.8	...	0.6	...	0.5	...	0.2	...	0.1
										
✓		✗		✓		✓		✗		✗
true								false		
positive								positive		
(high overlap)								(no overlap, low overlap, or duplicate)		

Evaluation Metric



$$precision@t = \frac{\#true\ positives@t}{\#true\ positives@t + \#false\ positives@t}$$

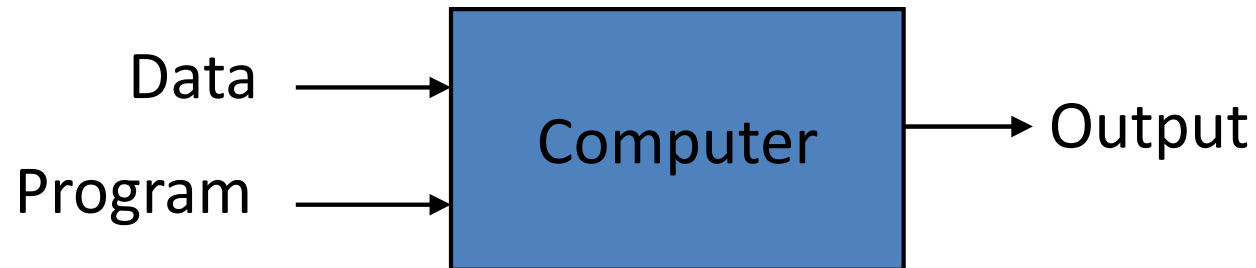
$$\frac{\checkmark}{\checkmark + \times}$$

$$recall@t = \frac{\#true\ positives@t}{\#ground\ truth\ objects}$$

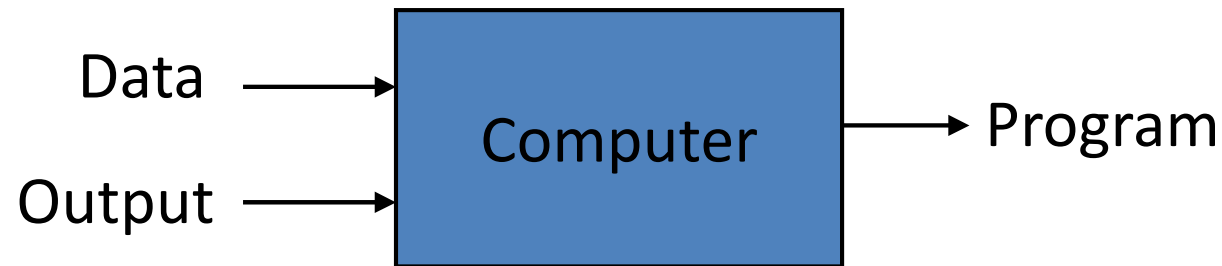
Machine Learning Based Object Detection

Traditional Programming v.s. Machine Learning

- Traditional Programming



- Machine Learning



Why Machine Learning is Hard?

Is this a tree?



A brown trunk moving upwards
and branching with leaves...?

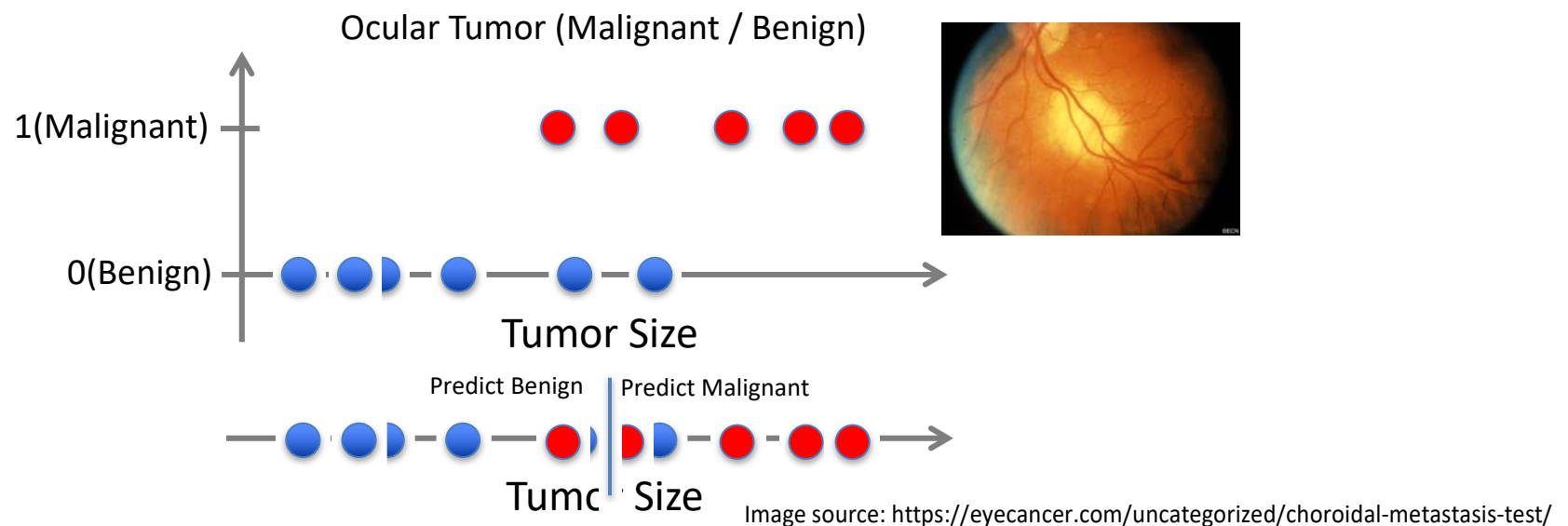


Types of Learning

- **Supervised (inductive) learning**
 - Given: training data + desired outputs (labels)
- **Unsupervised learning**
 - Given: training data (without desired outputs)
- **Semi-supervised learning**
 - Given: training data + a few desired outputs
- **Reinforcement learning**
 - Rewards from sequence of actions

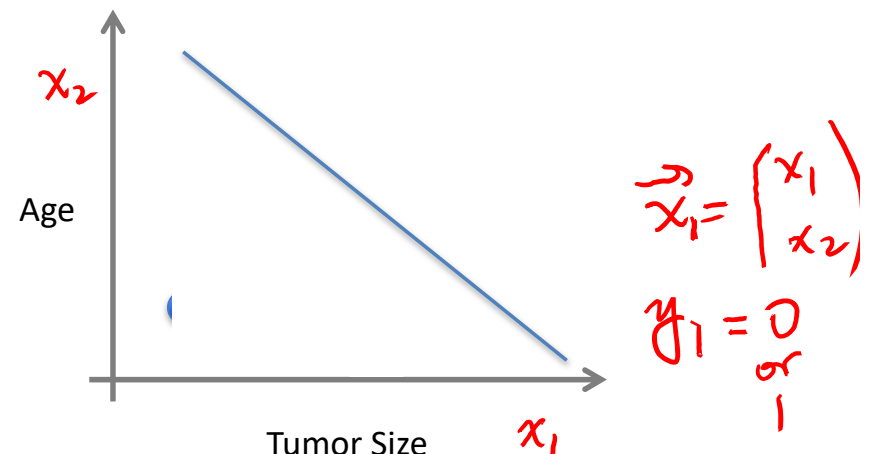
Supervised Learning: Classification

- Given $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Learn a function $f(x)$ to predict y given x
 - y is categorical == classification



Supervised Learning: Classification

- x can be multi-dimensional
 - each dimension corresponds to an attribute:
 - clump thickness
 - color
 - distance from optic nerve
 - ...

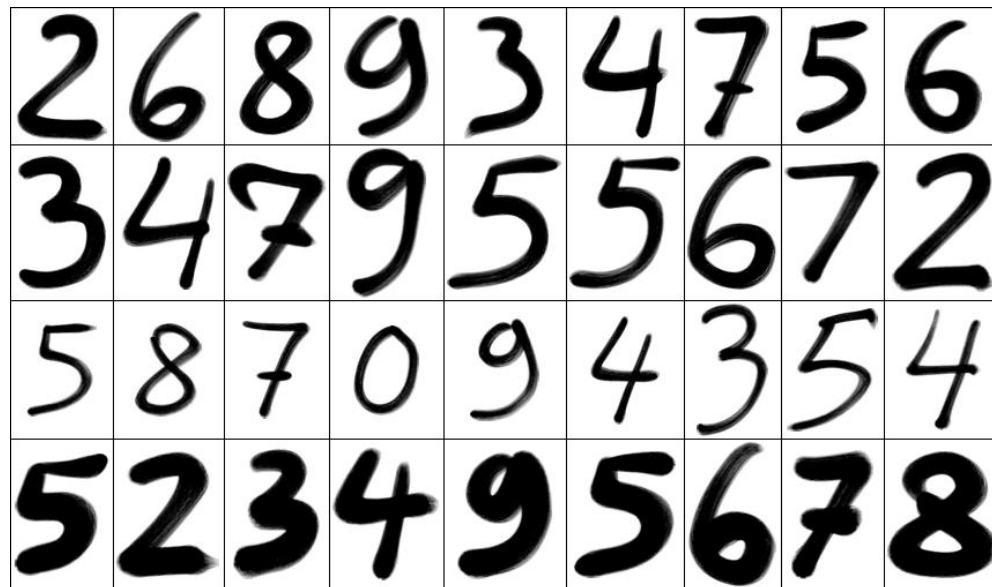


- Cell type is the most telling feature, but it's risky to do a biopsy of the eye
 - ML can help determine *when* a feature is needed



Supervised Learning: Multiclass Classification

- Multiclass classification problems
 - Written digits $\rightarrow 0, 1, \dots, 9$
 - Pictures \rightarrow apple, orange, strawberry
 - Emails \rightarrow spam, primary, social, promotion



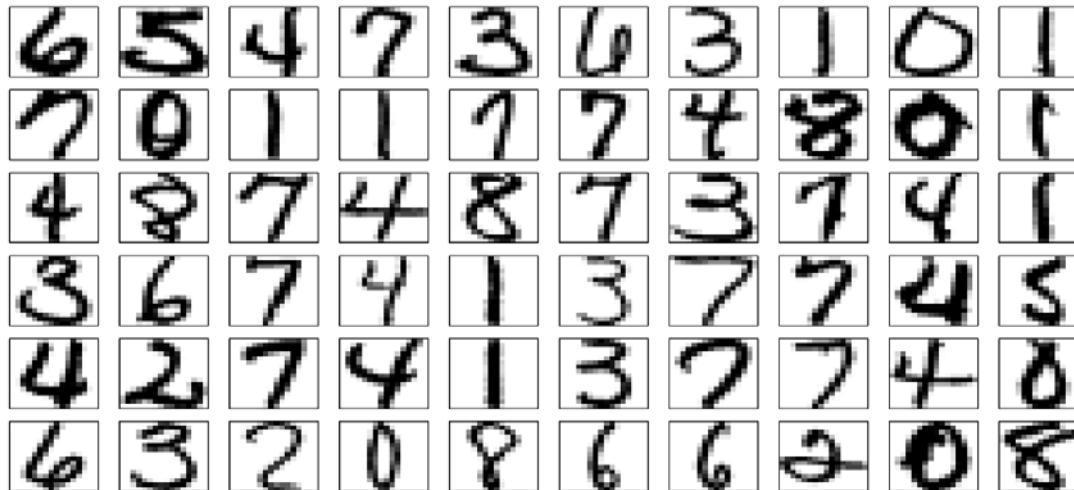
Example Data from ImageNet



Machine Learning Example

- Digit recognition

Each digit is a 16×16 image.

[illegible]

$$\mathbf{x} = (1, x_1, \dots, x_{256}) \leftarrow \text{input}$$

$$y = b$$

The Key Players

- Pictures
 - input $\mathbf{x} \in \mathbb{R}^d = \mathbf{X}$
- Classes: cat, dog, desk, etc.... ,
 - output $y \in \{1, 2, \dots\} = \mathbf{Y}$
- True relationship between \mathbf{x} and y
 - target function $f: \mathbf{X} \rightarrow \mathbf{Y}$
- Data
 - data set $\mathbf{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
 - $y_n = f(\mathbf{x}_n)$

We learn the function f from the data \mathbf{D}

$$\underline{f(\mathbf{x})} \rightarrow y \quad \checkmark$$

- \mathbf{X} , \mathbf{Y} , and \mathbf{D} are given by the learning problem; the target f is fixed but **unknown**

Learning

- Start with a set of candidate hypotheses H which you think are likely to represent f

$$H = \{h_1, h_2, \dots\}$$

H is called the **hypothesis set** or **model**

- Select a hypothesis h from H . The way we do this is called a **learning algorithm**
- Use h for new input. We hope $h \approx f$
- Again, X , Y , and D are given by the learning problem; the target f is fixed but **unknown**

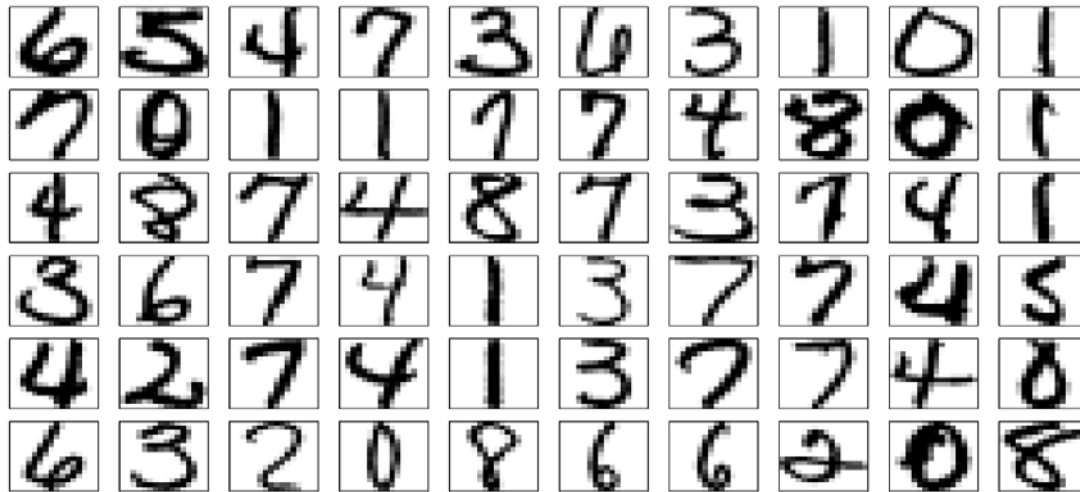
We choose H and the learning algorithm

This is a very general setup (e.g. choose H to be all possible hypotheses)

Revisit: Digit Recognition Problem

- Digit recognition

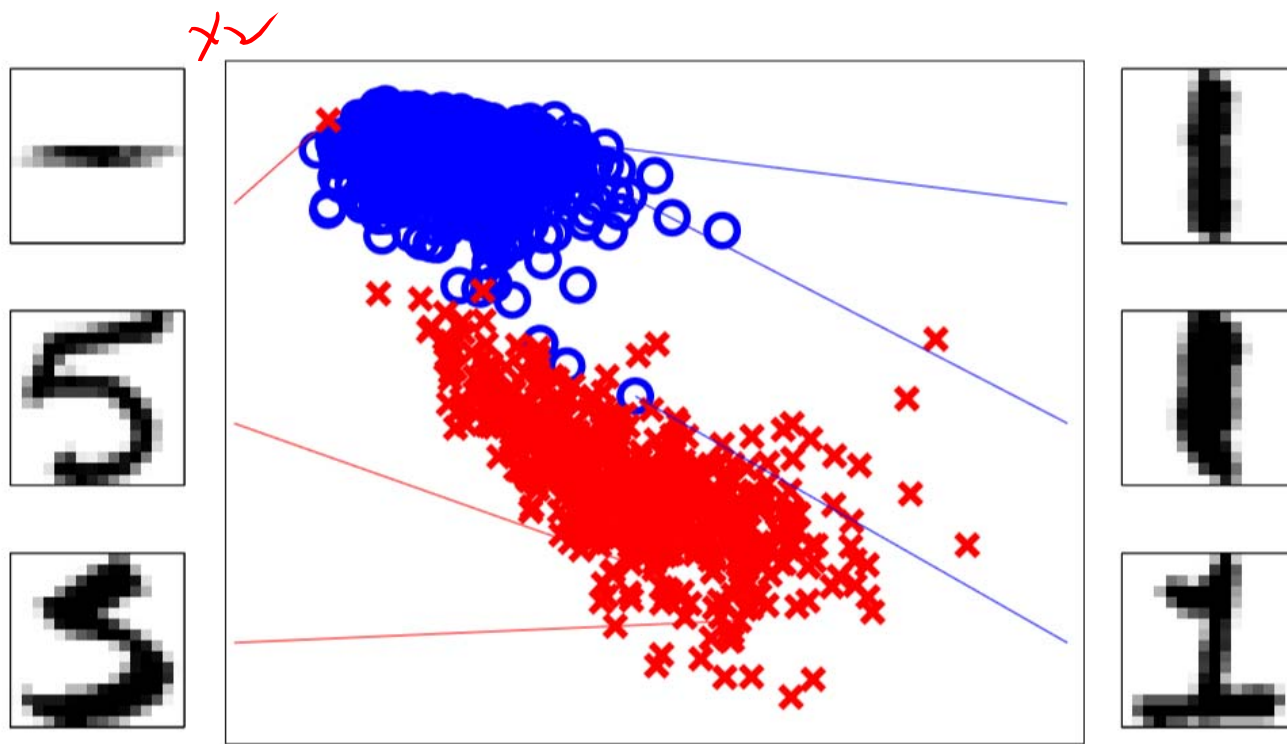
Each digit is a 16×16 image.

[illegible]
$$y = b$$
$$\mathbf{x} = (1, x_1, \dots, x_{256}) \leftarrow \text{input}$$
$$\mathbf{w} = (w_0, w_1, \dots, w_{256}) \leftarrow \text{linear model}$$

Hypothesis: $h = g(\mathbf{w}^T \mathbf{x})$
e.g.,: $h = \text{sign}(\mathbf{w}^T \mathbf{x})$

Machine Learning Example

- Feature: an important property of the input that you think is useful for classification, e.g.,



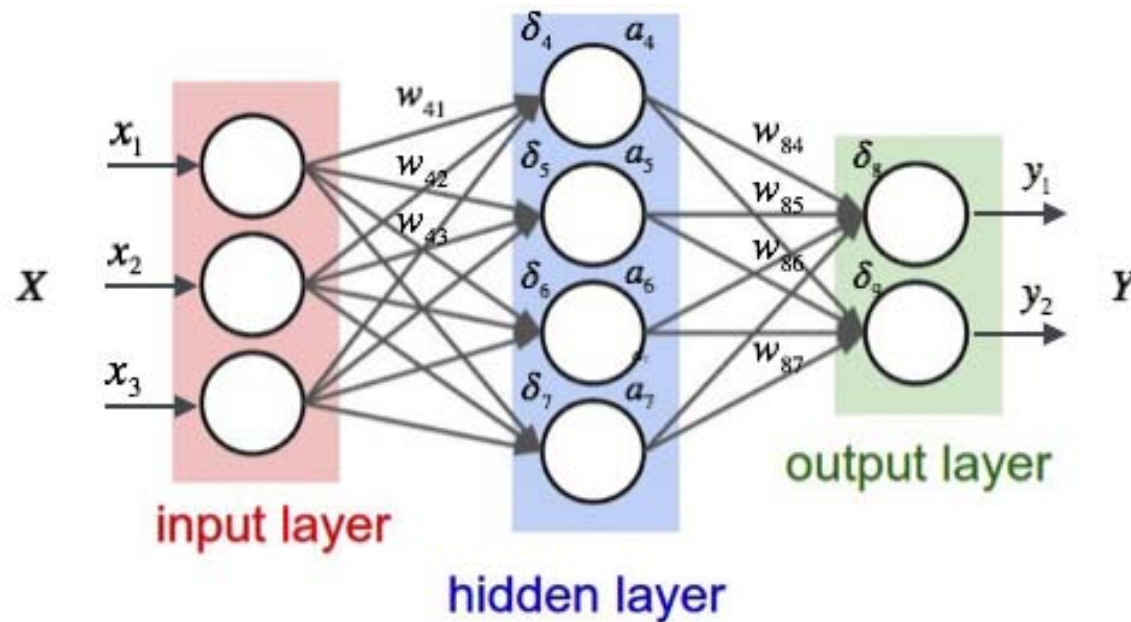
$$\mathbf{x} = (x_0, x_1, x_2)$$

x_1 : intensity

x_2 : symmetry

Neural Network

- How does a neural network work?



Input

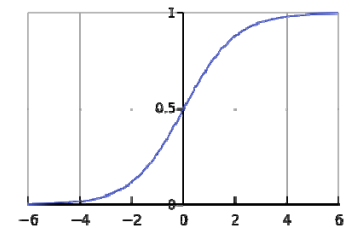
$$\mathbf{x} = (x_1, x_2, \dots)^T$$

Activation function

$$a_j = g(\mathbf{w}_j^T \mathbf{x})$$

Example:

$$g(t) = \frac{1}{1 + e^{-t}}$$



Goal: learn \mathbf{w} !

Neural Network for Object Classification



Pedestrian



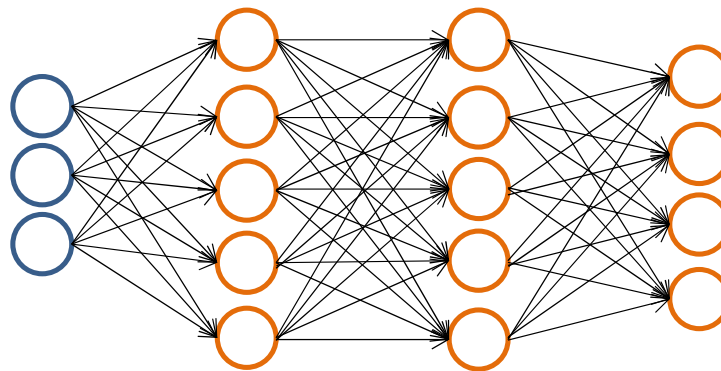
Car



Motorcycle



Truck



$$h_{\Theta}(\mathbf{x}) \in \mathbb{R}^K$$

- We want

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

when pedestrian

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

when car

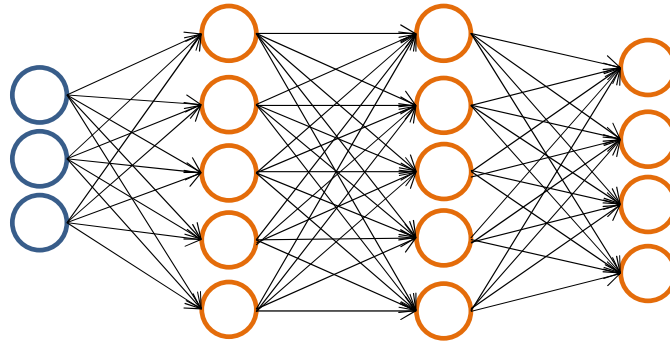
$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

when motorcycle

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

when truck

Neural Network for Object Classification



$$h_{\Theta}(\mathbf{x}) \in \mathbb{R}^K$$

- We want

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

when pedestrian

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

when car

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

when motorcycle

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

when truck

- Given $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- Must convert labels to 1-of-K representation
 - e.g., $y_i = [0 \ 0 \ 1 \ 0]^T$ when motorcycle, $y_i = [0 \ 1 \ 0 \ 0]^T$ when car

- Thank You!