

# CS 4644 / 7643-A

## DANFEI XU

(SLIDE CREDIT: PROF. ZSOLT KIRA)

Topics:

- Machine learning intro, applications (CV, NLP, etc.)
- Parametric models and their components

- **PS0: This should take less than 2hrs!**
- Please do it now, and give others a chance at waitlist if your background is not sufficient (beef it up and take it next time)
  - Do it even if you're on the waitlist!
- **Piazza: not all enrolled!**
  - Enroll now! <http://piazza.com/gatech/spring2022/cs46447643a/>
  - Make it active!
- **Office hours** start next week

- **Collaboration**
  - Only on HWs and project (not allowed in HW0/PS0).
  - You may discuss the questions
  - Each student writes their own answers
  - Write on your homework anyone with whom you collaborate
  - Each student must write their own code for the programming part
  - Do NOT search for code implementing what we ask; search for concepts
- **Zero tolerance on plagiarism**
  - Neither ethical nor in your best interest
  - Always credit your sources
  - Don't cheat. We will find out.

- **Grace period**
  - 2 days grace period for each assignment (**EXCEPT PS0**)
    - Intended for checking submission NOT to replace due date
    - No need to ask for grace, no penalty for turning it in within grace period
    - Can NOT use for PS0
- **After grace period, you get a 0 (no excuses except medical)**
  - Send all medical requests to dean of students (<https://studentlife.gatech.edu/>)
  - Form: [https://gatech-advocate.symplicity.com/care\\_report/index.php/pid224342](https://gatech-advocate.symplicity.com/care_report/index.php/pid224342)
- **DO NOT SEND US ANY MEDICAL INFORMATION!** We do not need any details, just a confirmation from dean of students

# Learn Numpy!

CS231n Convolutional Neural Networks for Visual Recognition

## Python Numpy Tutorial

This tutorial was contributed by [Justin Johnson](#).

We will use the Python programming language for all assignments in this course. Python is a great general-purpose programming language on its own, but with the help of a few popular libraries (numpy, scipy, matplotlib) it becomes a powerful environment for scientific computing.

We expect that many of you will have some experience with Python and numpy; for the rest of you, this section will serve as a quick crash course both on the Python programming language and on the use of Python for scientific computing.

<http://cs231n.github.io/python-numpy-tutorial/>

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Machine Learning Overview

# What is Machine Learning (ML)?

*“A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .”*

*Tom Mitchell (Machine Learning, 1997)*

# When is Machine Learning useful?

```
algorithm quicksort(A, lo, hi) is
  if lo < hi then
    p := partition(A, lo, hi)
    quicksort(A, lo, p - 1)
    quicksort(A, p + 1, hi)
```

```
algorithm partition(A, lo, hi) is
  pivot := A[hi]
  i := lo
  for j := lo to hi do
    if A[j] < pivot then
      swap A[i] with A[j]
      i := i + 1
  swap A[i] with A[hi]
  return i
```

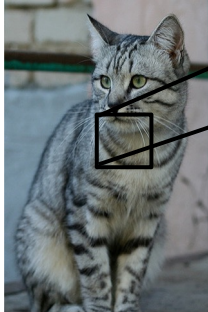


Cat

When it's difficult / infeasible to write a program



# Example: Object Recognition



```

[1095 112 100 111 104 99 996 99 99 100 112 100 104 97 97 871
 1 91 98 102 106 104 79 99 103 99 105 112 116 118 105 94 891
 7 99 89 98 106 120 105 87 96 95 99 115 112 106 103 109 891
 1 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94
 1206 92 62 64 69 92 88 85 101 107 109 98 75 64 91
 1154 100 85 55 55 69 64 54 64 87 112 129 98 74 84 91
 1237 127 144 140 95 81 89 85 52 54 74 84 102 93 85 82
 1208 137 144 140 109 95 86 79 82 85 83 83 88 73 88 101
 1225 127 140 127 110 121 117 94 82 79 80 68 54 79 80
 1227 125 131 147 133 127 126 131 111 96 79 75 61 64 72 84
 115 114 109 123 158 140 132 118 113 100 90 92 74 65 79 81
 1 89 93 98 97 100 147 131 118 113 114 113 109 106 95 77 88
 1 83 77 88 81 77 79 102 123 117 115 117 125 120 115 87
 1 62 65 82 89 78 71 80 101 124 119 118 101 107 114 111 110
 1 63 65 75 88 89 71 62 81 120 118 110 105 81 98 110 110
 1 87 65 71 87 106 95 69 45 76 120 126 107 92 94 105 121
 1138 97 82 86 117 123 116 86 41 51 95 89 95 102 107
 1164 146 112 88 82 120 124 104 76 48 45 66 88 101 102 109
 1157 170 137 128 93 86 114 132 112 97 69 65 78 82 69 94
 1108 120 114 141 139 100 109 118 121 114 87 65 63 68 88
 1208 112 96 117 158 144 120 115 104 107 102 93 87 81 72 79
 1223 107 96 86 83 112 110 120 104 75 88 107 112 99
 1221 131 102 88 82 86 94 117 145 148 153 102 58 78 92 107
 1221 164 148 103 71 56 78 83 93 103 119 139 102 61 69 841]
    
```

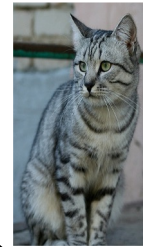
This image by Nikita is licensed under CC-BY 2.0

What the computer sees  
What the computer sees

An image is just a big grid of numbers between [0, 255]:

e.g. 800 x 600 x 3  
(3 channels RGB)

Viewpoint Changes



```

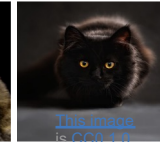
[1095 112 100 111 104 99 996 99 99 100 112 100 104 97 97 871
 1 91 98 102 106 104 79 99 103 99 105 112 116 118 105 94 891
 7 99 89 98 106 120 105 87 96 95 99 115 112 106 103 109 891
 1 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94
 1206 92 62 64 69 92 88 85 101 107 109 98 75 64 91
 1154 100 85 55 55 69 64 54 64 87 112 129 98 74 84 91
 1237 127 144 140 95 81 89 85 52 54 74 84 102 93 85 82
 1208 137 144 140 109 95 86 79 82 85 83 83 88 73 88 101
 1225 127 140 127 110 121 117 94 82 79 80 68 54 79 80
 1227 125 131 147 133 127 126 131 111 96 79 75 61 64 72 84
 115 114 109 123 158 140 132 118 113 100 90 92 74 65 79 81
 1 89 93 98 97 100 147 131 118 113 114 113 109 106 95 77 88
 1 83 77 88 81 77 79 102 123 117 115 117 125 120 115 87
 1 62 65 82 89 78 71 80 101 124 119 118 101 107 114 111 110
 1 63 65 75 88 89 71 62 81 120 118 110 105 81 98 110 110
 1 87 65 71 87 106 95 69 45 76 120 126 107 92 94 105 121
 1138 97 82 86 117 123 116 86 41 51 95 89 95 102 107
 1164 146 112 88 82 120 124 104 76 48 45 66 88 101 102 109
 1157 170 137 128 93 86 114 132 112 97 69 65 78 82 69 94
 1108 120 114 141 139 100 109 118 121 114 87 65 63 68 88
 1208 112 96 117 158 144 120 115 104 107 102 93 87 81 72 79
 1223 107 96 86 83 112 110 120 104 75 88 107 112 99
 1221 131 102 88 82 86 94 117 145 148 153 102 58 78 92 107
 1221 164 148 103 71 56 78 83 93 103 119 139 102 61 69 841]
    
```

All pixels change when the camera moves!

Illumination



This image is CC0 1.0 public



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain

Deformation



This image by Salvagnin is licensed under CC-BY 2.0



This image by Salvagnin is licensed under CC-BY 2.0



This image by sara bear is licensed under CC-BY 2.0

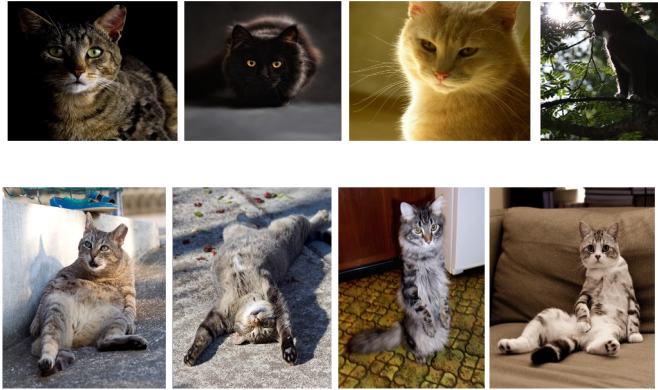


This image by Tom That is licensed under CC-BY 2.0

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

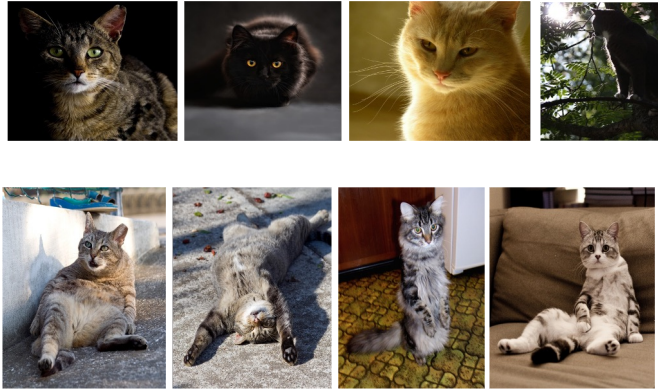
## Why Object Recognition is Hard

# The Power of Machine Learning



Cat

# The Power of Machine Learning



Deep Neural  
Networks



Cat

# The Power of (Deep) Machine Learning

TECHNOLOGY

## A Massive Google Network Learns To Identify — Cats

June 26, 2012 · 3:00 PM ET

Heard on [All Things Considered](#)



4-Minute Listen

+ PLAYLIST

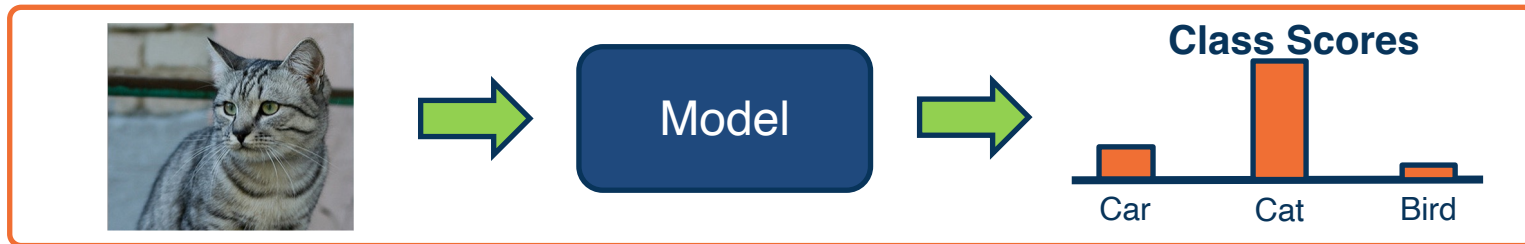


*All Things Considered* host Audie Cornish talks with Andrew Ng, Associate Professor of Computer Science at Stanford University. He led a Google research team in creating a neural network out of 16,000 computer processors to try and mimic the functions of the human brain. Given three days on YouTube, the network taught itself how to identify — cats.

Source: <https://www.npr.org/2012/06/26/155792609/a-massive-google-network-learns-to-identify>

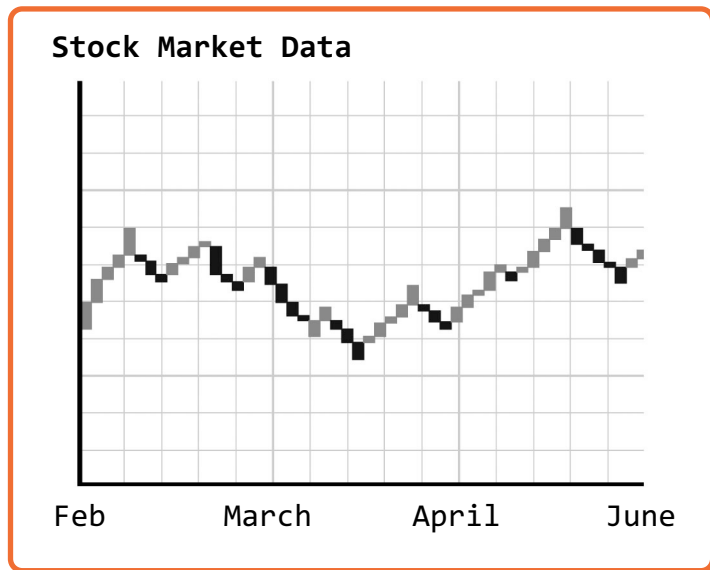
Deep Learning

# Application: Computer Vision

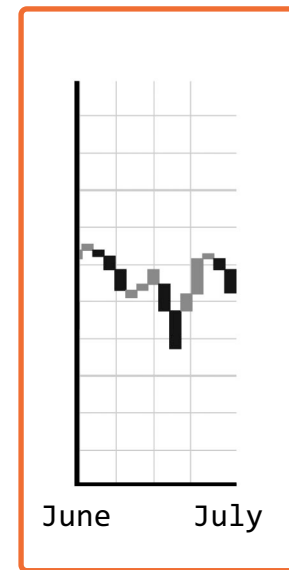
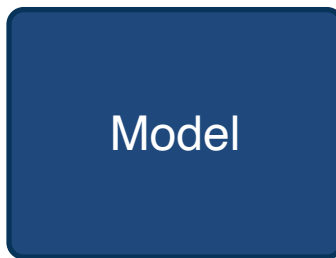


Example: Image Classification

# Application: Time Series Forecasting



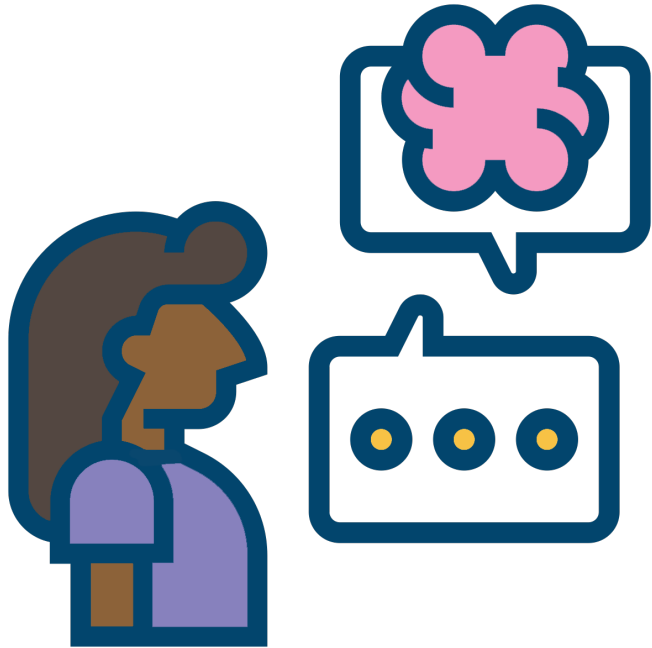
**Input**



**Prediction**

**Example: Time Series Forecasting**

# Application: Natural Language Processing (NLP)



Very large number of NLP sub-tasks:

- ◆ Syntax Parsing
- ◆ Translation
- ◆ Named entity recognition
- ◆ Summarization
- ◆ Similarity / paraphrasing

**Sequence modeling:** Variable length sequential inputs and/or outputs

**Recent progress:** Large-scale Language Models

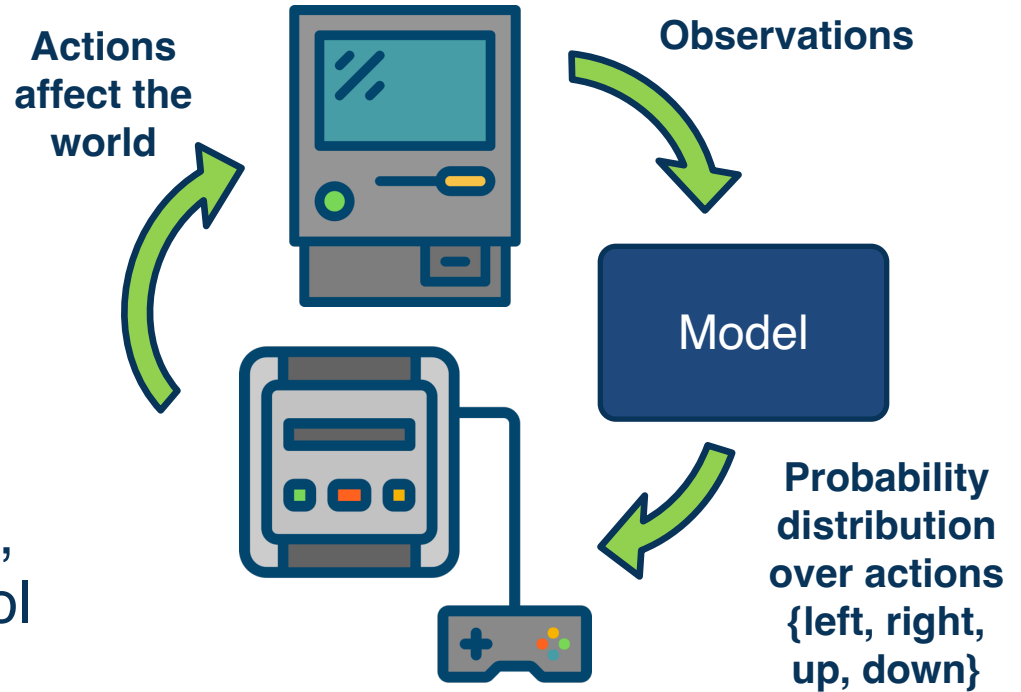
Example: Natural Language Processing (NLP)

# Application: Decision Making

## Example: Video Game

- Sequence of inputs/outputs
- Actions affect the environment

**Examples:** Chess / Go,  
Video Games,  
Recommendation Systems,  
Network Congestion Control  
...





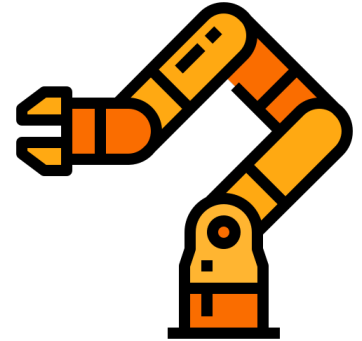
Robotics involves a **combination of AI/ML techniques:**

- ◆ **Sense:** Perception
- ◆ **Plan:** Planning
- ◆ **Act:** Controls

Some things are **learned (perception)**, while others **programmed**

**An evolving landscape**

**Application:**



Rest of the lecture (also next lecture):

- ◆ **Types of Machine Learning Problems**
- ◆ **Parametric Models**
- ◆ **Linear Classifiers**
- ◆ **Gradient Descent**

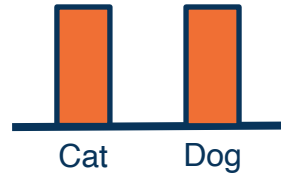
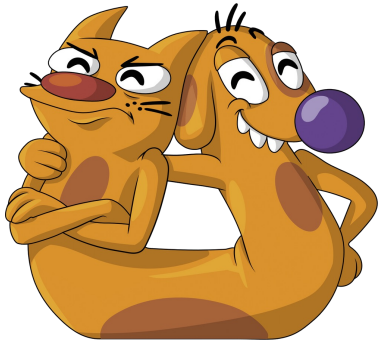
**Supervised  
Learning**

**Unsupervised  
Learning**

**Reinforcement  
Learning**

# Supervised Learning

- Train Input:  $\{X, Y\}$
- Learning output:  $f : X \rightarrow Y$
- Usually  $f$  is a **distribution**, e.g.  $P(y|x)$



<https://en.wikipedia.org/wiki/CatDog>

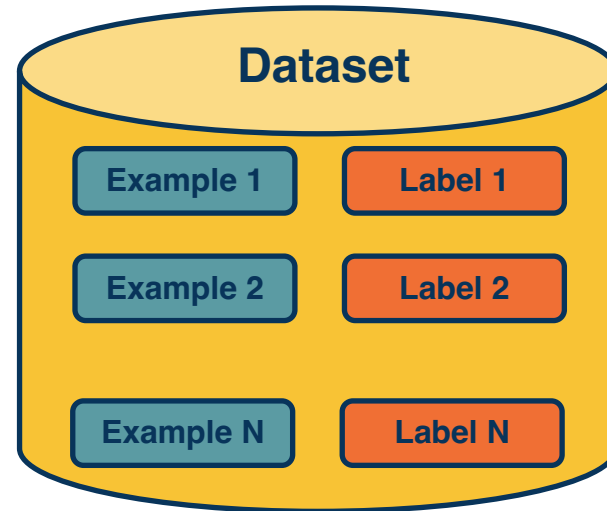
## Dataset

$$X = \{x_1, x_2, \dots, x_N\} \text{ where } x \in \mathbb{R}^d$$

Examples

$$Y = \{y_1, y_2, \dots, y_N\} \text{ where } y \in \mathbb{R}^c$$

Labels



# Supervised Learning

- Train Input:  $\{X, Y\}$
- Learning output:  $f : X \rightarrow Y$ , e.g.  $p(y|x)$

## Terminology:

- Model / Hypothesis Class
  - $H: \{f: X \rightarrow Y\}$
  - Learning is search in hypothesis space

E.g.,  $H = \{f(x) = w^T x \mid w \in \mathbb{R}^d\}$

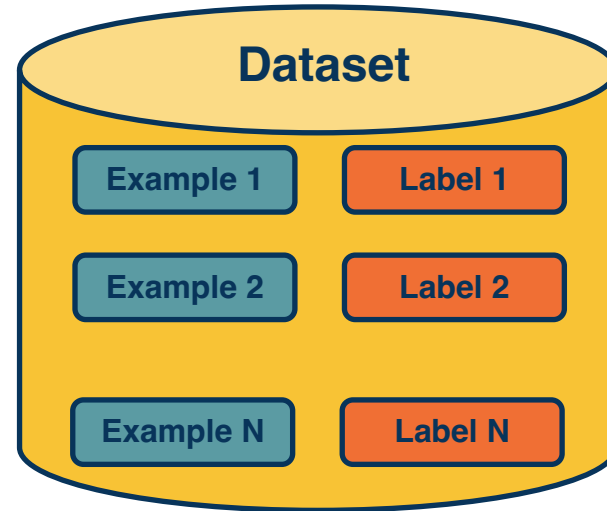
## Dataset

$X = \{x_1, x_2, \dots, x_N\}$  where  $x \in \mathbb{R}^d$

Examples

$Y = \{y_1, y_2, \dots, y_N\}$  where  $y \in \mathbb{R}^c$

Labels



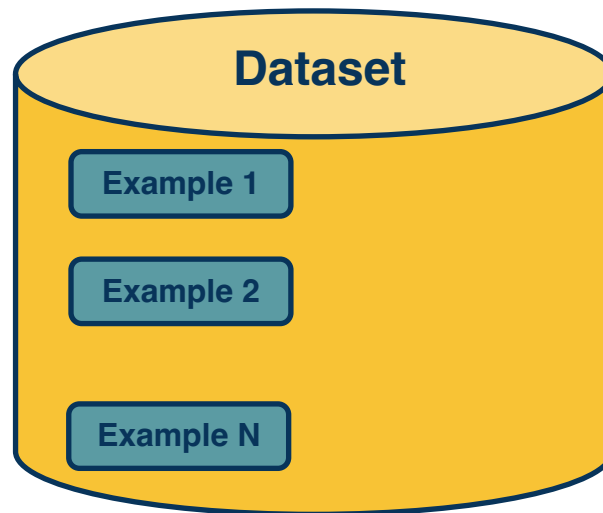
## Unsupervised Learning

- Input:  $\{X\}$
- Learning output:  $p_{data}(x)$
- How likely is  $x$  under  $p_{data}$ ?
- Can we sample from  $p_{data}$ ?
- Example: Clustering, density estimation, generative modeling, ...

## Dataset

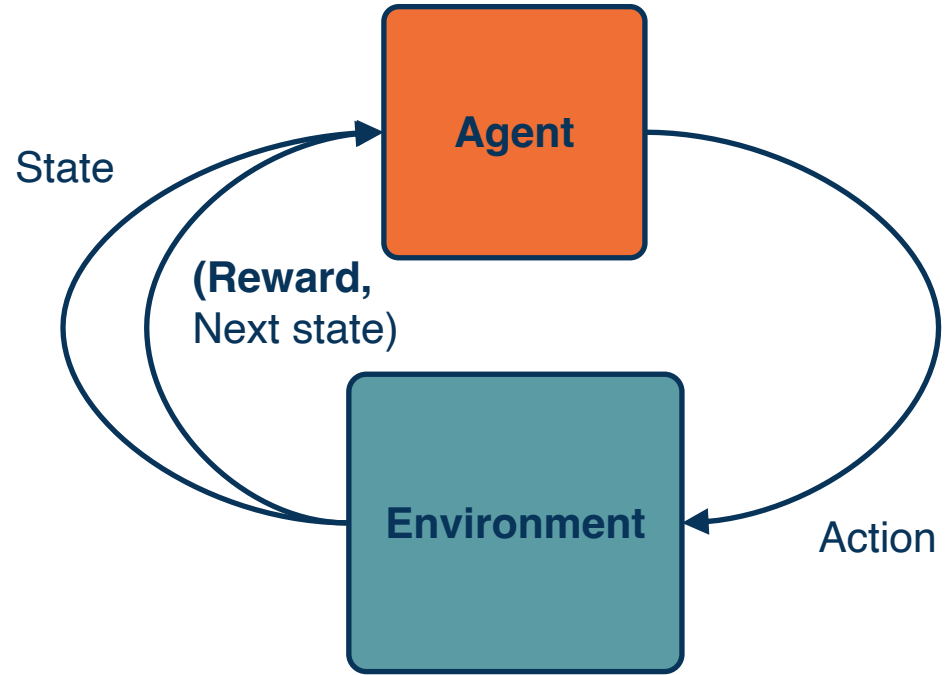
$X = \{x_1, x_2, \dots, x_N\}$  where  $x \in \mathbb{R}^d$

Examples



## Reinforcement Learning

- ◆ Supervision in form of **reward**
- ◆ No supervision on what action to take



*Adapted from: [http://cs231n.stanford.edu/slides/2020/lecture\\_17.pdf](http://cs231n.stanford.edu/slides/2020/lecture_17.pdf)*

## Supervised Learning

- ◆ Train Input:  $\{X, Y\}$
- ◆ Learning output:  
 $f : X \rightarrow Y$ ,  
e.g.  $P(y|x)$

## Unsupervised Learning

- ◆ Input:  $\{X\}$
- ◆ Learning output:  $P(x)$
- ◆ Example: Clustering, density estimation, etc.

## Reinforcement Learning

- ◆ Supervision in form of **reward**
- ◆ No supervision on what action to take

**Very often combined**, sometimes within the same model!



Rest of the lecture (also next lecture):

- ◆ Types of Machine Learning Problems
- ◆ **Parametric Models**
- ◆ Linear Classifiers
- ◆ Gradient Descent

## Non-Parametric Model

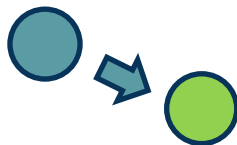
No explicit model for the function,  
**examples:**

- ◆ Nearest neighbor classifier
- ◆ Decision tree

Hypothesis class changes with  
the number of data points

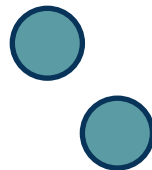
## Non-Parametric – Nearest Neighbor

Example 1, cat



Query

Example 2, dog



Example 4, dog

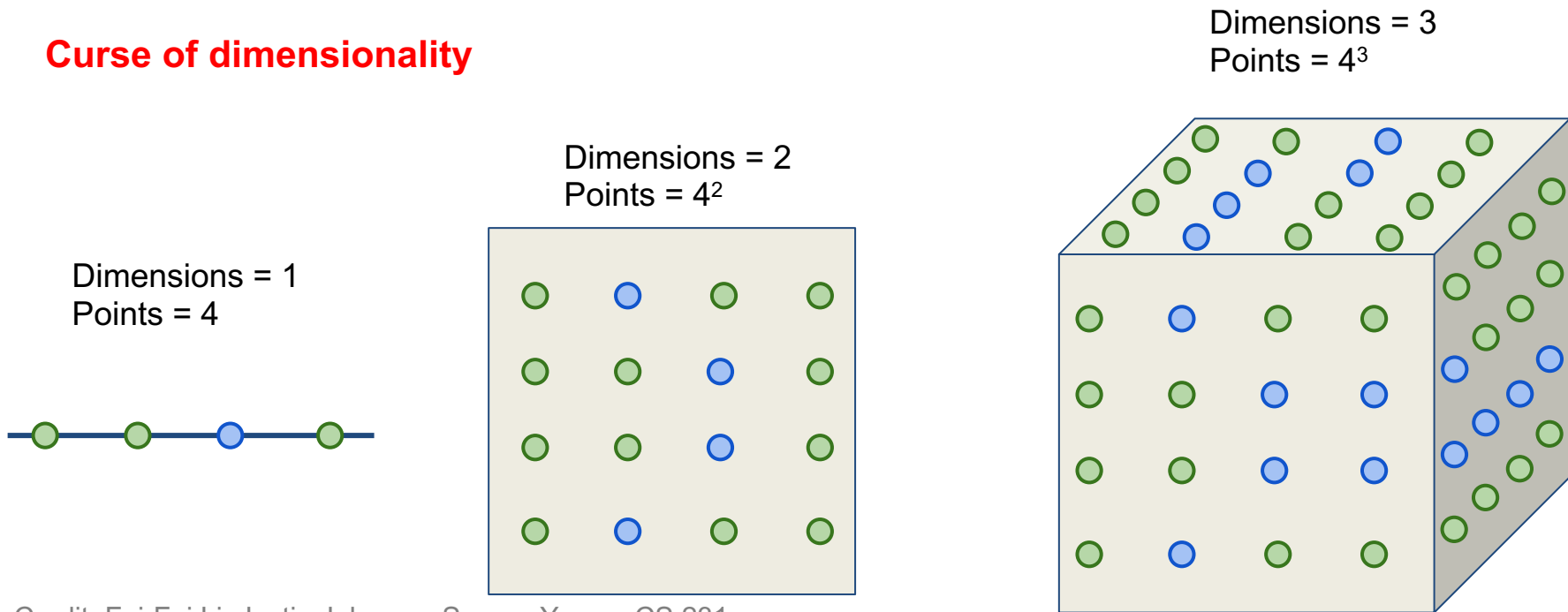


Example 3, car

**Procedure:** Take label of nearest example

k-Nearest Neighbor on high-dimensional data (e.g., images) is *almost never* used.

### Curse of dimensionality



Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

- **Curse of Dimensionality**
  - Data required increases exponentially with the number of dimensions
- **Doesn't work well when large number of irrelevant features**
  - Distances overwhelmed by noisy features
- **Expensive**
  - No Learning: most real work done during testing
  - For every test sample, must search through all dataset – very slow!
  - Must use tricks like approximate nearest neighbor search

## Parametric Model

Explicitly model the function  $f : X \rightarrow Y$  in the form of a parametrized function

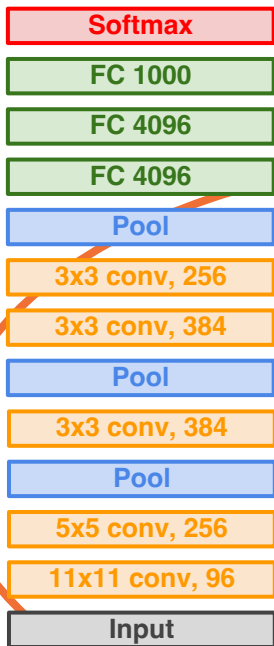
$f(x, W) = y$ , **examples:**

- ◆ Logistic regression/classification
  - ◆ Number of parameters grows **linearly** with the number of dimensions!
- ◆ Neural networks
- ◆ Hypothesis classes doesn't change

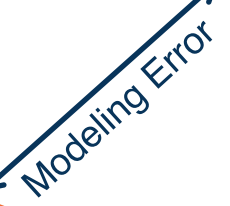
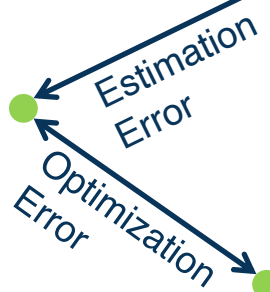
## Parametric – Linear Classifier

$$f(x, W) = Wx + b$$

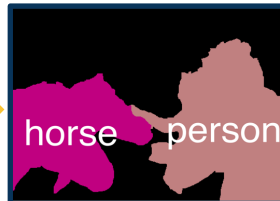
# AlexNet



model class

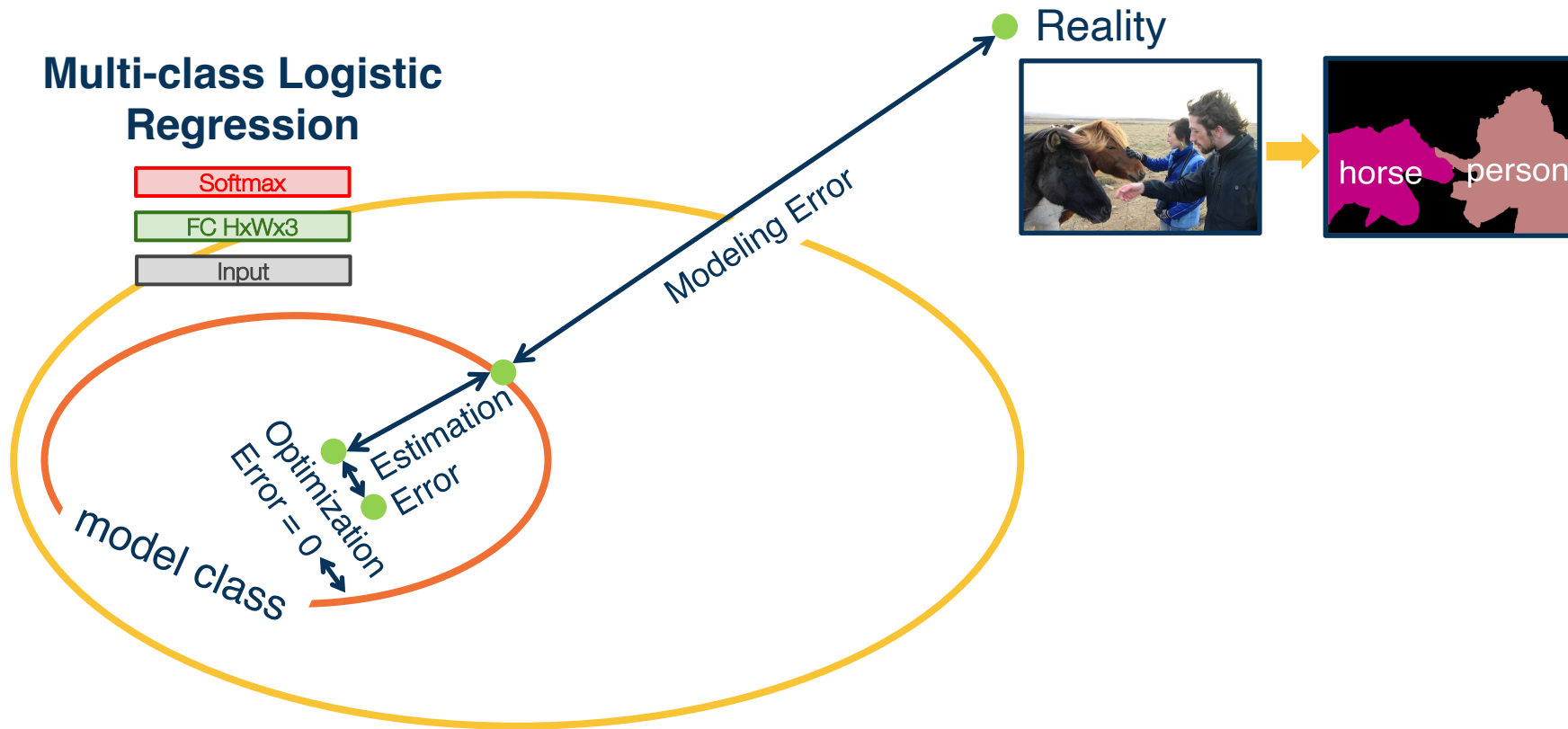
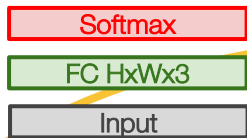


Reality



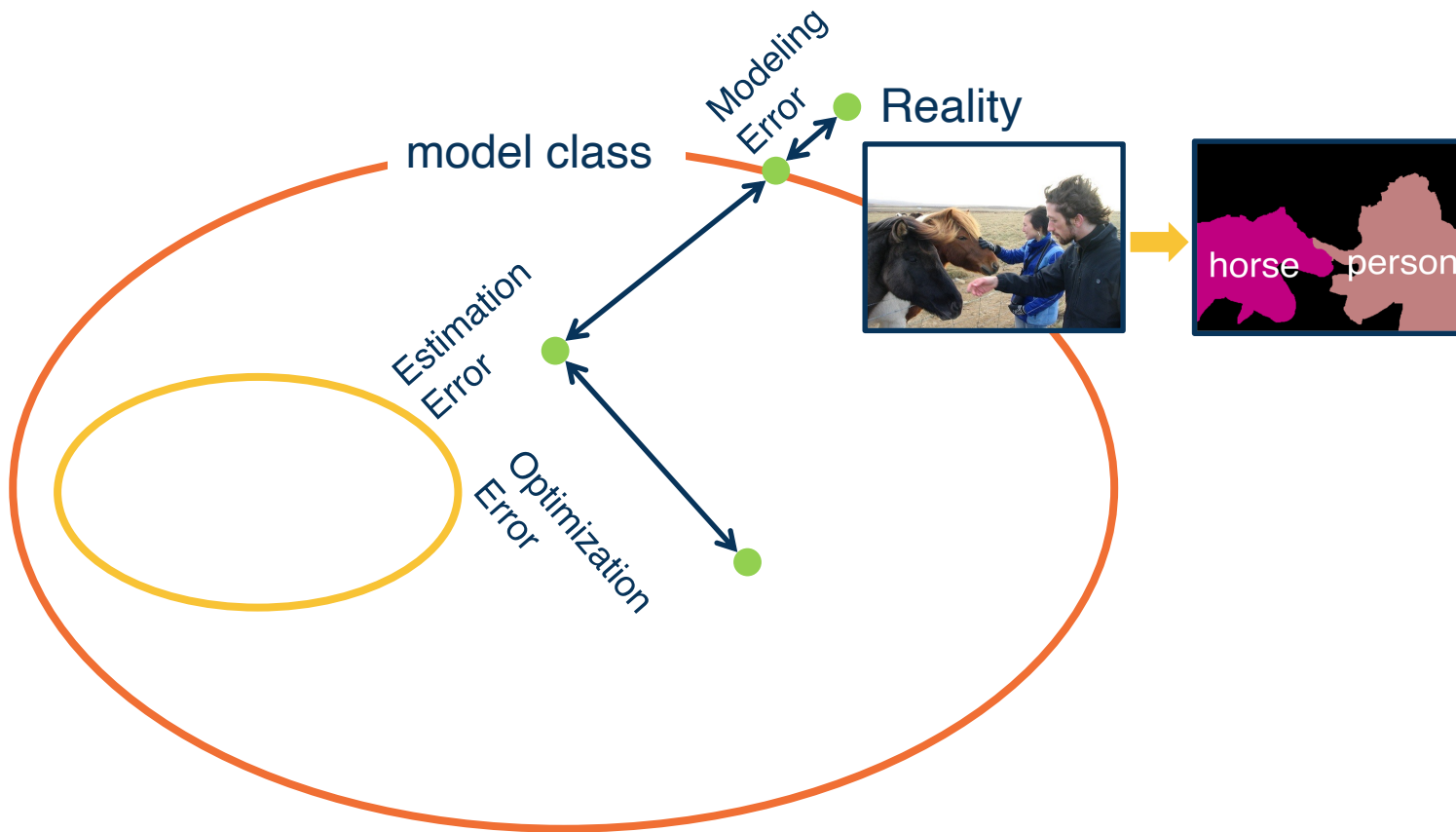
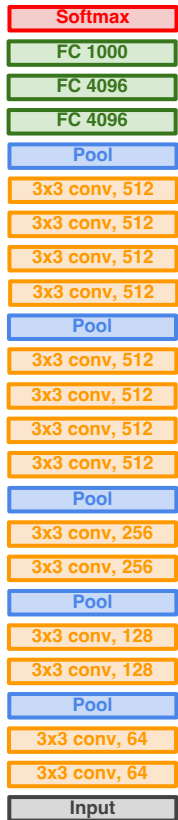
From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Multi-class Logistic Regression



From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# VGG19



From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

## Types of Errors and Generalization



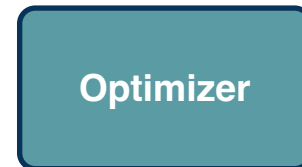
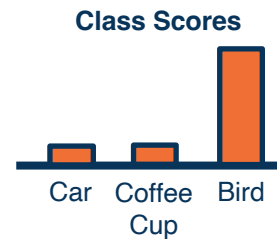
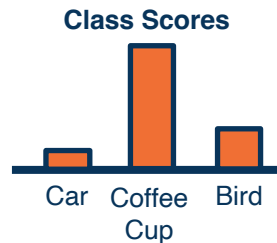
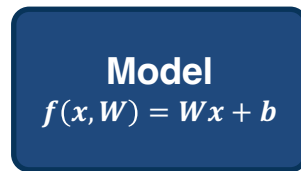
Rest of the lecture (also next lecture):

- ◆ Types of Machine Learning Problems
- ◆ Parametric Models
- ◆ **Linear Classifiers**
- ◆ Gradient Descent

- Input (and representation)
- Functional form of the model
  - Including parameters
- Performance measure to improve
  - Loss or objective function
- Algorithm for finding best parameters
  - Optimization algorithm

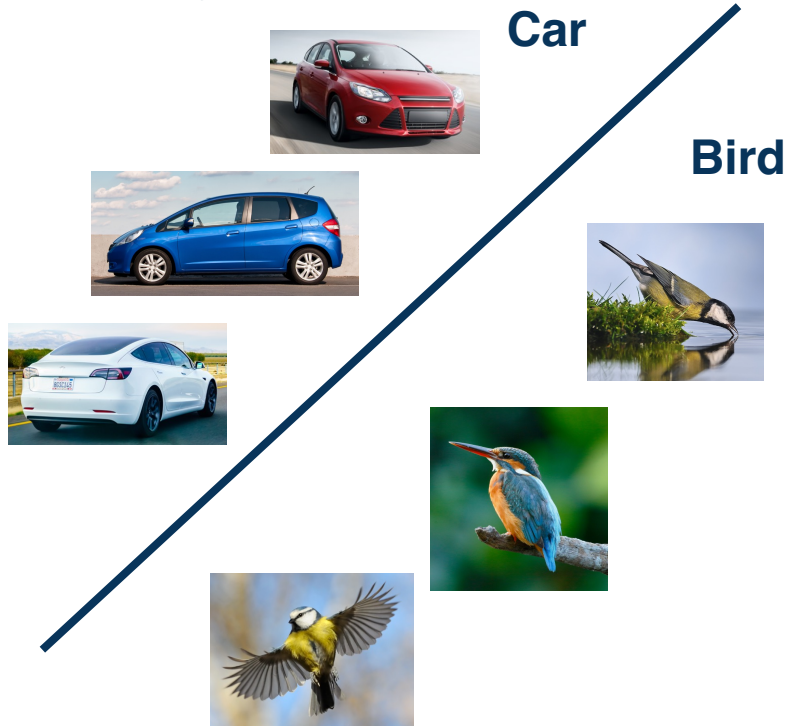


Data: Image



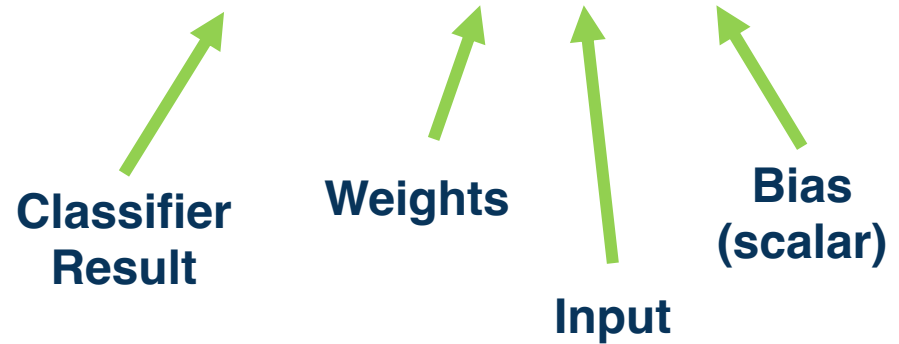
## Components of a Parametric Model

What is the **simplest function** you can think of?



Our model is:

$$f(x, w) = w \cdot x + b$$



(Note if  $w$  and  $x$  are column vectors we often show this as  $w^T x$ )

# Linear Classification and Regression

## Simple linear classifier:

- Calculate score:

$$f(x, w) = w \cdot x + b$$

- Binary classification rule ( $w$  is a vector):

$$y = \begin{cases} 1 & \text{if } f(x, w) \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- For multi-class classifier take class with highest (max) score

$$f(x, W) = Wx + b$$



## Data: Image



**Model**  
 $f(x, W) = Wx + b$



## Class Scores



$$x = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix} \xrightarrow{\text{Flatten}} x = \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{21} \\ x_{22} \\ \vdots \\ x_{n1} \\ \vdots \\ x_{nn} \end{bmatrix}$$

To simplify notation we will refer to inputs as  $x_1 \cdots x_m$  where  $m = n \times n$

## Model

$$f(x, W) = Wx + b$$

Classifier for class 1  $\longrightarrow$

Classifier for class 2  $\longrightarrow$

Classifier for class 3  $\longrightarrow$

$$\begin{bmatrix} W_{11} & W_{12} & \cdots & W_{1m} \\ W_{21} & W_{22} & \cdots & W_{2m} \\ W_{31} & W_{32} & \cdots & W_{3m} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$W$                        $x$                        $b$

(Note that in practice, implementations can use  $xW$  instead, assuming a different shape for  $W$ . That is just a different convention and is equivalent.)

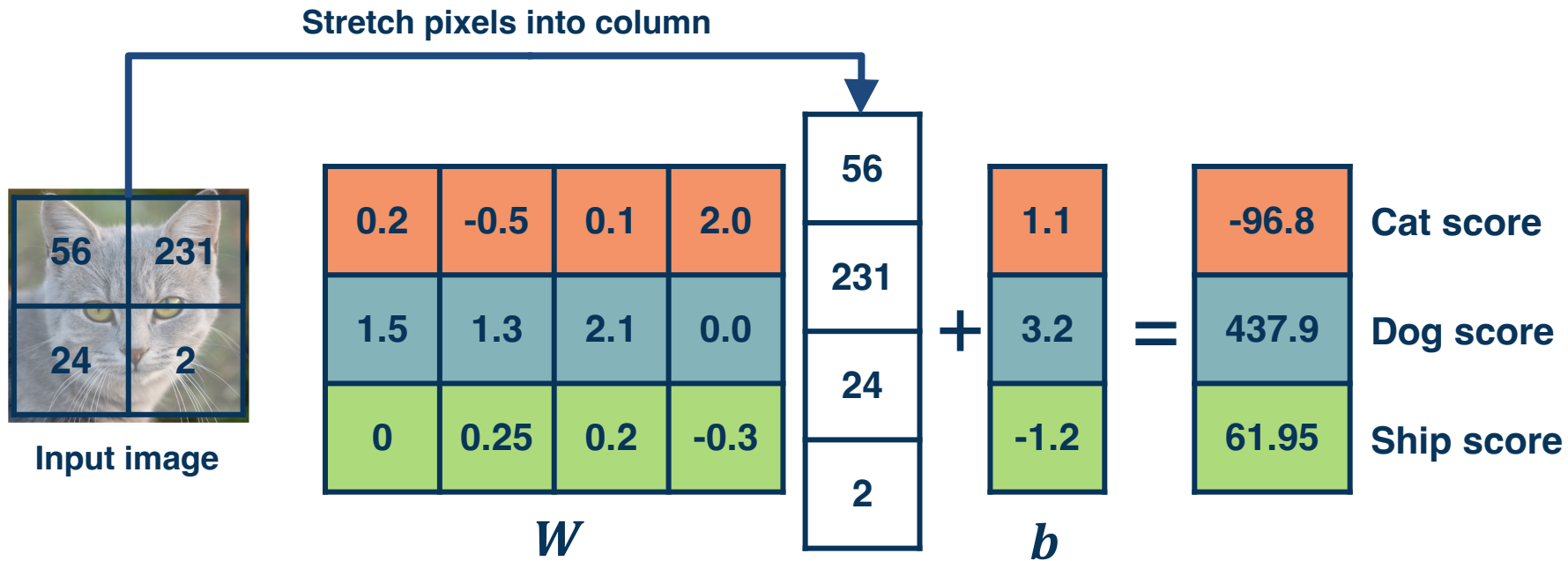
- We can move the bias term into the weight matrix, and a “1” at the end of the input
- Results in **one matrix-vector multiplication!**

$$\text{Model}$$
$$f(x, W) = Wx + b$$

$$\begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b_1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b_2 \\ w_{31} & w_{32} & \cdots & w_{3m} & b_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \\ 1 \end{bmatrix}$$

$W$   $x$

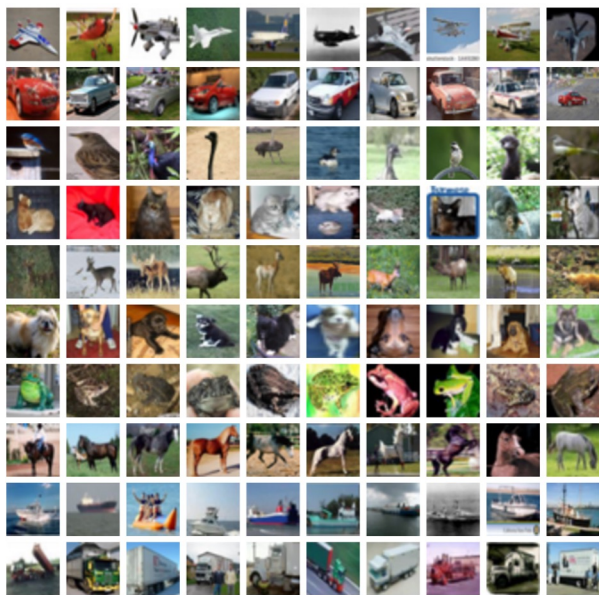
# Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n

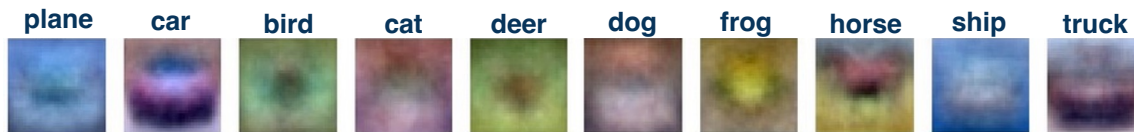


airplane  
automobile  
bird  
cat  
deer  
dog  
frog  
horse  
ship  
truck

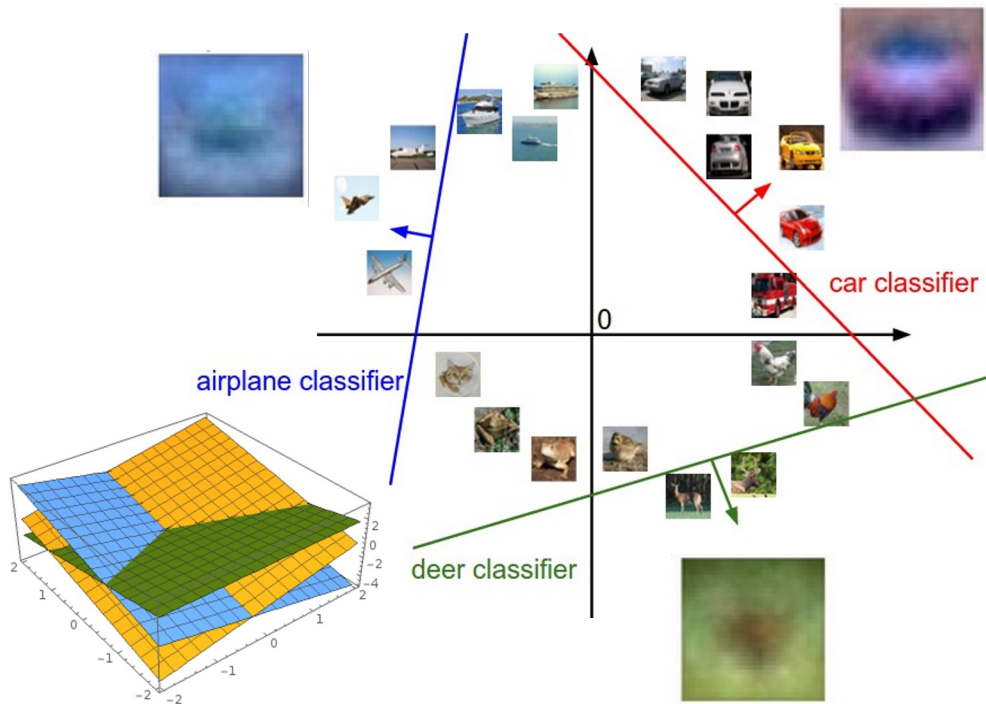


## Visual Viewpoint

We can convert the weight vector back into the shape of the image and visualize



*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*



Plot created using Wolfram Cloud

## Geometric Viewpoint

$$f(x, W) = Wx + b$$

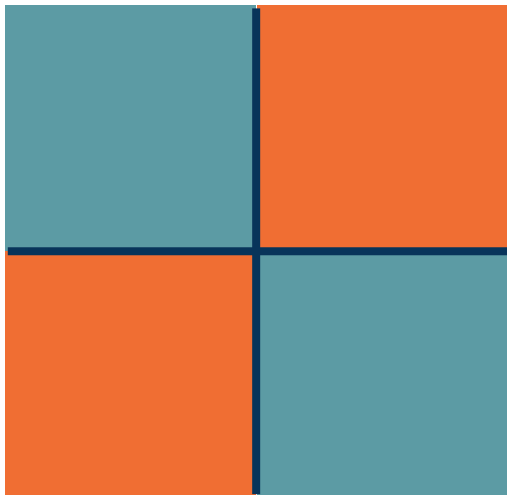


Array of **32x32x3** numbers  
(3072 numbers total)

*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

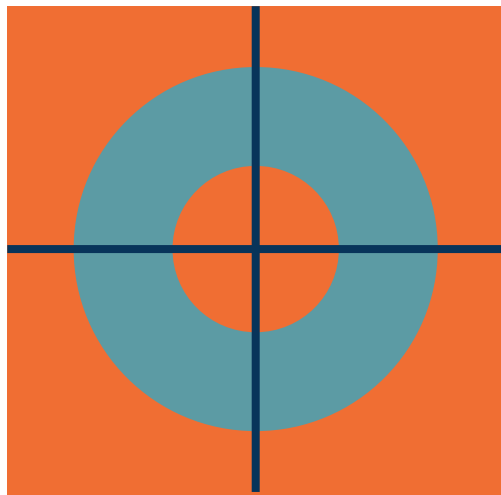
**Class 1:**  
number of pixels  $> 0$  odd

**Class 2:**  
number of pixels  $> 0$  even



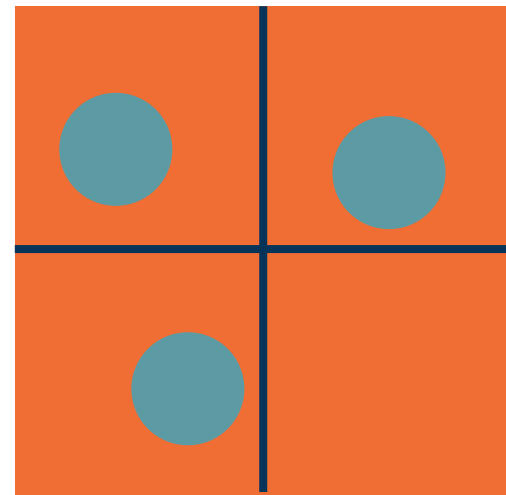
**Class 1:**  
 $1 \leq \text{L2 norm} \leq 2$

**Class 2:**  
Everything else



**Class 1:**  
Three modes

**Class 2:**  
Everything else



*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

# Neural Network

Linear  
classifiers

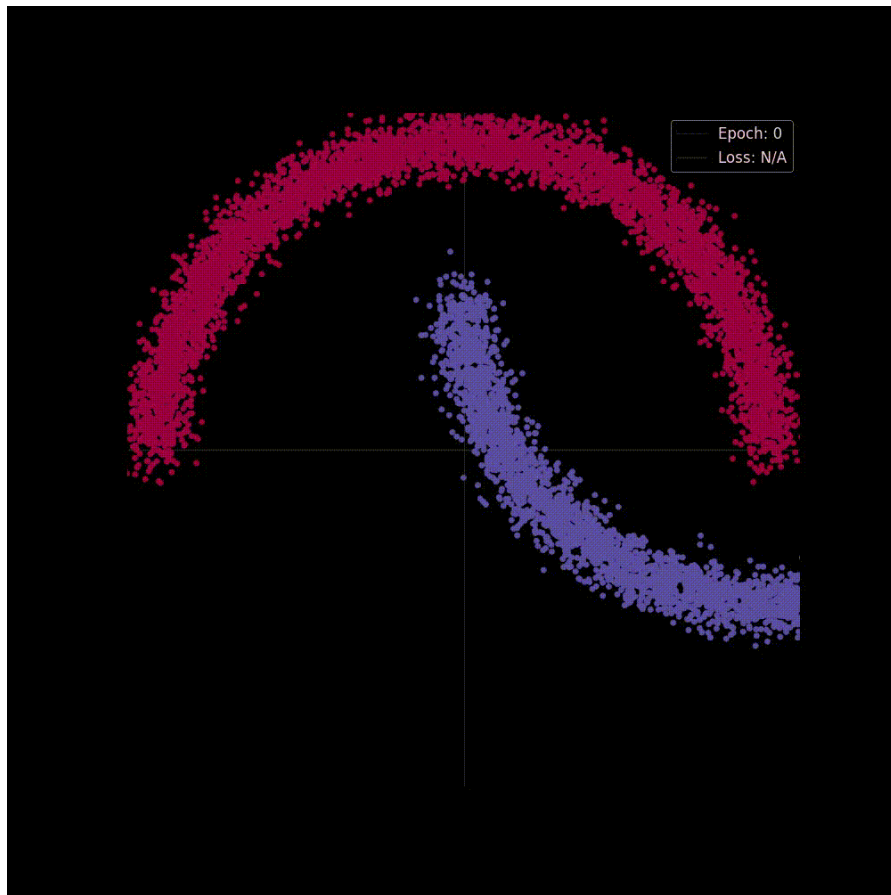
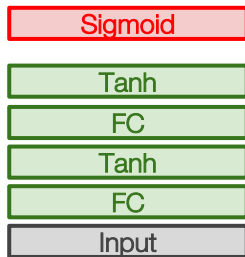


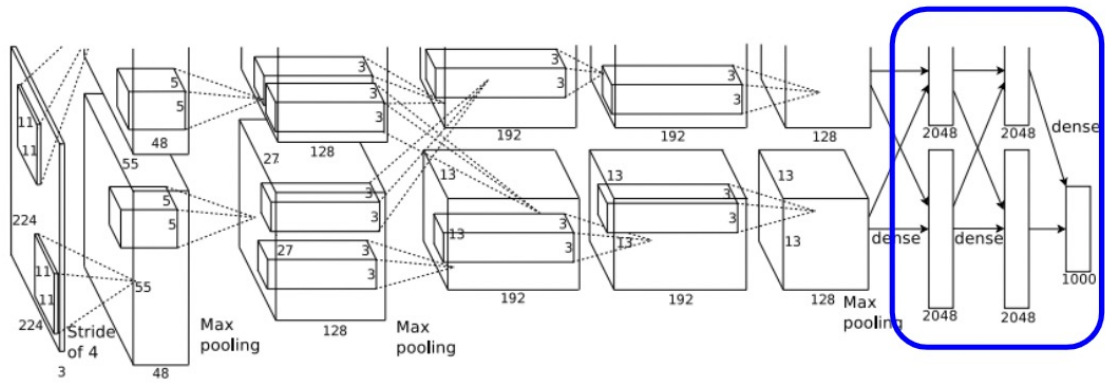
[This image](#) is [CC0 1.0](#) public domain

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# (Deep) Representation Learning for Classification

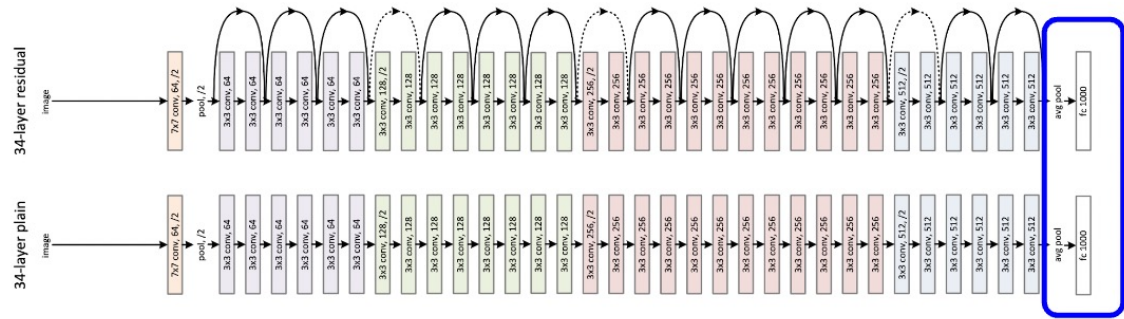
A function that transforms raw data space into a linearly-separable space





[Krizhevsky et al. 2012]

Linear layers

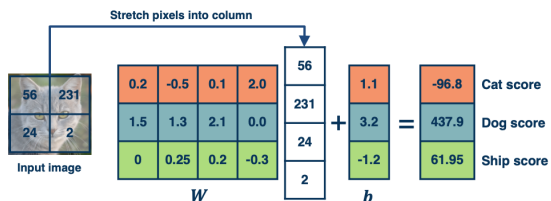


[He et al. 2015]

[This image](#) is [CC0 1.0](#) public domain

## Algebraic Viewpoint

$$f(x, W) = Wx$$



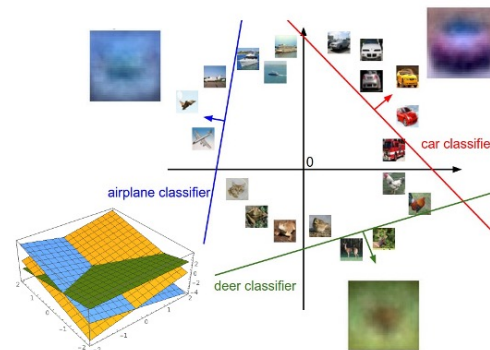
## Visual Viewpoint

One template per class



## Geometric Viewpoint

Hyperplanes cutting up space



Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n

# Next time:

- Input (and representation)
- Functional form of the model
  - Including parameters
- Performance measure to improve**
  - Loss or objective function**
- Algorithm for finding best parameters
  - Optimization algorithm



Data: Image

