Topics:

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

# CS 4644 / 7643-A
# ZSOLT KIRA

**Machine Learning Applications**

- **What's up with the capacity/waitlist?**

- **PS0 due Sunday night!**
  - Please do it!
  - We have fixed some gradescope autograder issues (sorry!)

- **Piazza**: not all enrolled!
  - Enroll now! https://piazza.com/gatech/spring2023/cs46447643/home (Code: DLSPR23 or through canvas)
  - Note: Do NOT post anything containing solutions publicly!
  - Make it active!

- **Office hours** start next week

**Administrivia**

Georgia
Tech

- **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
  - Do NOT search for code implementing what we ask; search for concepts
  - **Each student must write their own code/proofs**

- **Zero tolerance on plagiarism**
  - Neither ethical nor in your best interest
  - Always credit your sources
  - Don't cheat. We will find out.

**Collaboration Policy**

Georgia Tech

- **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

- **DO NOT SEND US ANY MEDICAL INFORMATION!** We do not need any details, just a confirmation from dean of students

**Grace Period**

Georgia Tech

## 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/

## Python + Numpy Tutorial

Georgia Tech

# 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."*
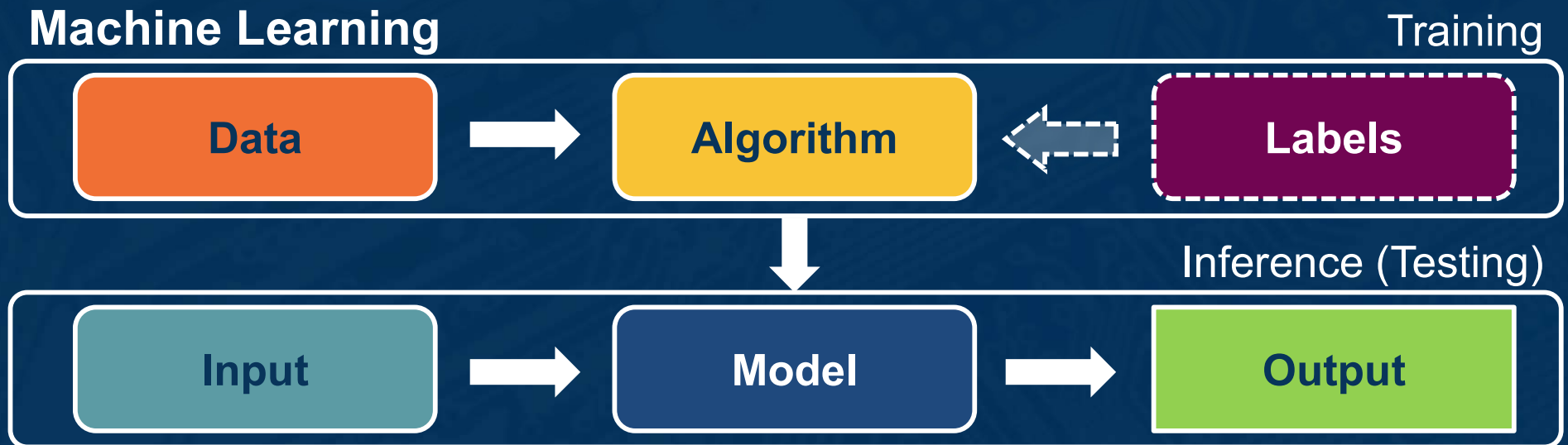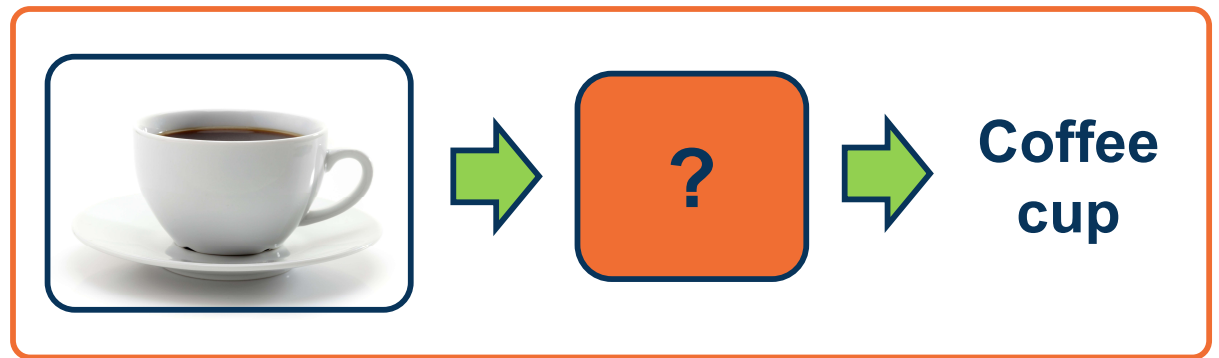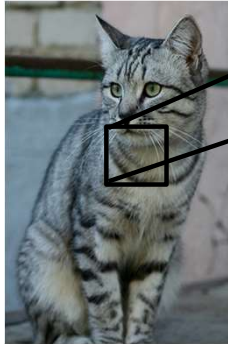
*Tom Mitchell (Machine Learning, 1997)*

Georgia
Tech

Machine learning thrives when it is **difficult to design an algorithm** to perform the task

Applications:

```
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
```



?

**Coffee cup**

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**

All pixels change when the camera moves!

**Illumination**

This image is CC0 1.0 public domain

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**

by Umberto Salvagnin is licensed under CC-BY 2.0

This image Umberto Salvagnin is licensed under CC-BY 2.0

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

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

# Why Image Classification is Hard

Georgia Tech

# The Power of Deep Learning



Adapted from figure by Andrej Karpathy

# Application: Computer Vision



**Class Scores**

Car    Coffee Cup    Bird

**Class Scores**

Normal    Benign    Malignant

**3D Reconstructions**

Kanazawa *et al.*, 2018

Georgia Tech

# Application: Time-Series Forecasting

Given a series of measurements, **output prediction for next time period**



**Input**

Model

**Prediction**

**Application: Natural Language Process (NLP)**

**Very large number of NLP sub-tasks:**

- Syntax Parsing

- Translation

- Named entity  recognition

- Summarization

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

**Recent progress**: Large-scale language  models

**Example: Natural Language Processing (NLP)**

Georgia Tech

# Decision-making tasks

- Sequence of inputs/outputs

- Actions affect the environment

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

**Application:**

Actions affect the world

Observations

Model

Probability distribution over actions {left, right, up, down}

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

- **Sense:** Perception
- **Plan:** Planning
- **Act:** Controls/Decision-Making

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

- Evolving landscape

**Example: Robotics**

Georgia Tech

| Supervised Learning | Unsupervised Learning | Reinforcement Learning |

**Types of Machine Learning**

Georgia Tech

# Supervised Learning

- Train Input: $\{X, Y\}$

- Learning output: $f : X \rightarrow Y$, e.g. a **distribution** $P(y|x)$



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

## Dataset

$$X = \{x_1, x_2, ..., x_N\} \; where \; x \in \mathbb{R}^d \quad \boxed{\text{Examples}}$$

$$Y = \{y_1, y_2, ..., y_N\} \; where \; y \in \mathbb{R}^c \quad \boxed{\text{Labels}}$$

### Dataset

| Example 1 | Label 1 |
|-----------|---------|
| Example 2 | Label 2 |
| Example N | Label N |

# Supervised Learning

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

**Terminology:**

- Model / Hypothesis Class
  - $H: \{h: X \rightarrow Y\}$
    - Learning is search in hypothesis space
- Note **inputs $x_i$ and $y_i$** are each represented as **vectors**

## Dataset

$$X = \{x_1, x_2, \ldots, x_N\} \; where \; x \in \mathbb{R}^d$$ **Examples**

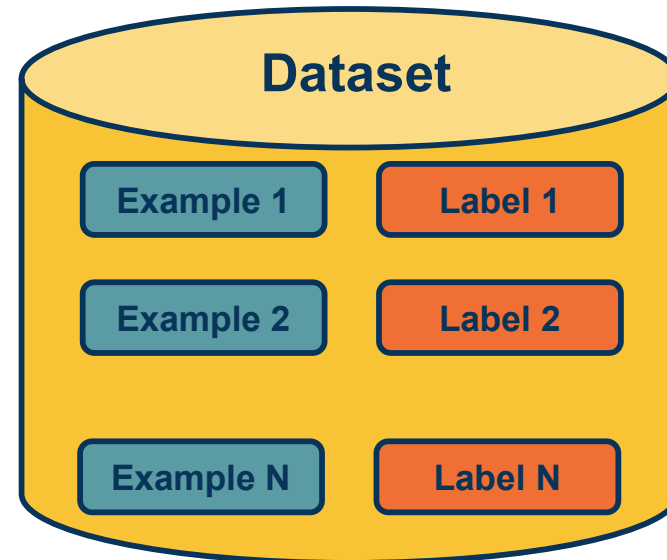$$Y = \{y_1, y_2, \ldots, y_N\} \; where \; y \in \mathbb{R}^c$$ **Labels**

### Dataset

| Example 1 | Label 1 |
| Example 2 | Label 2 |
| Example N | Label N |

Georgia Tech

**Dataset**

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

Examples

## 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, etc.

**Dataset**

Example 1

Example 2

Example N

# Reinforcement Learning

- Supervision in form of **reward**

- No supervision on what action to take

State

**Agent**

Reward
Next state

**Environment**

Action

*Adapted from: http://cs231n.stanford.edu/slides/2020/lecture_17.pdf*

Georgia
Tech

## 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!

**Types of Machine Learning**

Georgia Tech

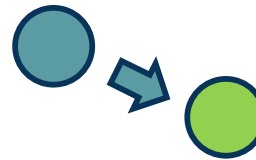## Non-Parametric Model

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

- ⬡ **Nearest neighbor classifier**

- ⬡ Decision tree

Capacity (size of hypothesis class) grow with size of training data!

## Non-Parametric – Nearest Neighbor

Example 1, cat          Example 2, dog

Query

Example 4, dog

Example 3, car

**Procedure:** Take label of nearest example

**Supervised Learning**

Georgia Tech

# k-Nearest Neighbor on images almost **never used.**

- Curse of dimensionality

  - Lots of weird behavior in high-dimensional spaces, e.g. orthogonality of random vectors, percentage of points around shell, etc.

Dimensions = 3
Points = $4^3$

Dimensions = 2
Points = $4^2$

Dimensions = 1
Points = 4

25

**Curse of Dimensionality**

Georgia Tech

- **Curse of Dimensionality**

  - Distances become meaningless in high 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

**Problems with Instance-Based Learning**

Georgia Tech

## 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

- Neural networks

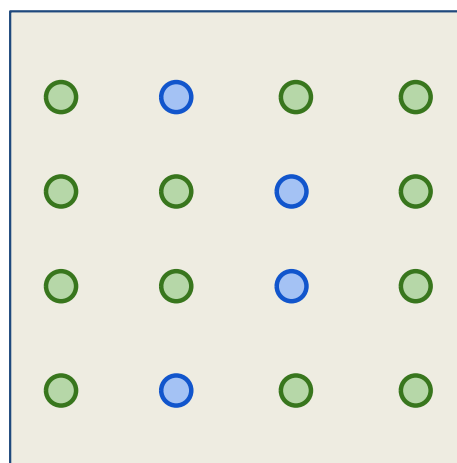Capacity (size of hypothesis class) **does not** grow with size of training data!

Learning is **search**

**Parametric – Linear Classifier**

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

# A Learning Problem



| Example | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---------|-------|-------|-------|-------|-----|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

**No Assumptions means no learning**

**Learning from a Broader Perspective**

Georgia Tech

Training Stage:

      Training Data $\{ (x_i, y_i) \} \rightarrow h$              (Learning)

Testing Stage

      Test Data $x \rightarrow h(x)$      (Apply function, Evaluate error)

**Procedural View of ML**

Georgia Tech

Probabilities to rescue:

X and Y are *random variables*

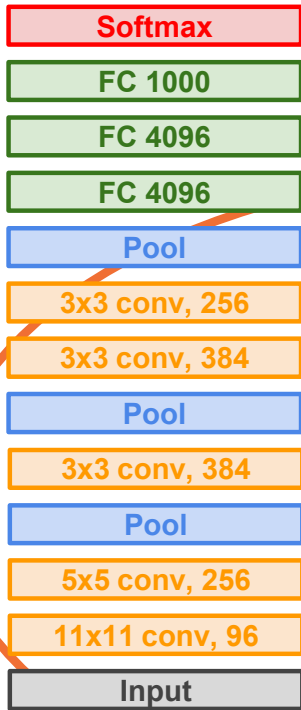$D = (x_1, y_1), (x_2, y_2), ..., (x_N, y_N) \sim P(X, Y)$

IID: Independent Identically Distributed

Both training & testing data sampled IID from $P(X,Y)$

Learn on training set

Have some hope of *generalizing* to test set

**AlexNet**

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 256
3x3 conv, 384
Pool
3x3 conv, 384
Pool
5x5 conv, 256
11x11 conv, 96
Input

model class

Reality

Modeling Error

Estimation Error

Optimization Error

horse    person

**Generalization**

Georgia Tech

**Multi-class Logistic Regression**

Softmax
FC HxWx3
Input

Reality

horse    person

Modeling Error
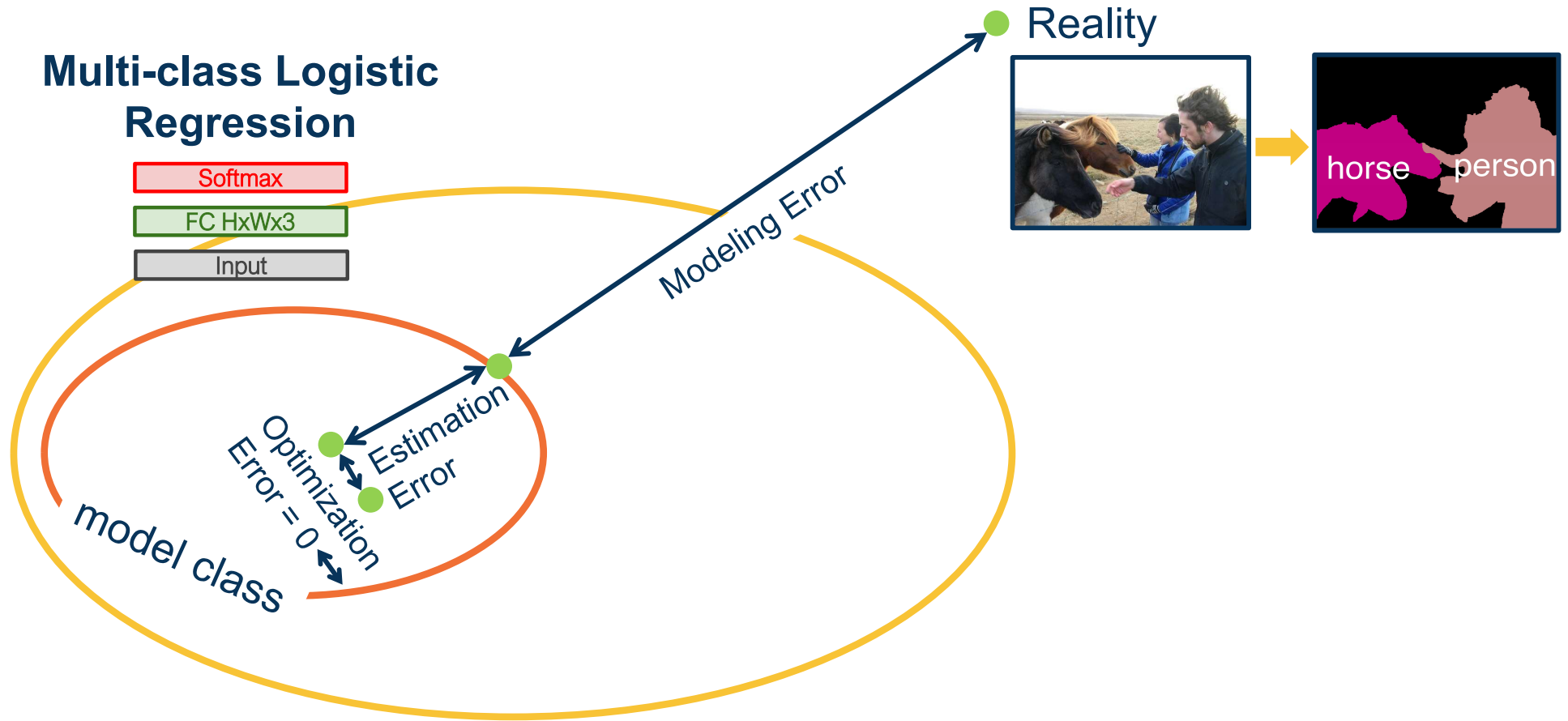
Estimation Error

Optimization Error = 0
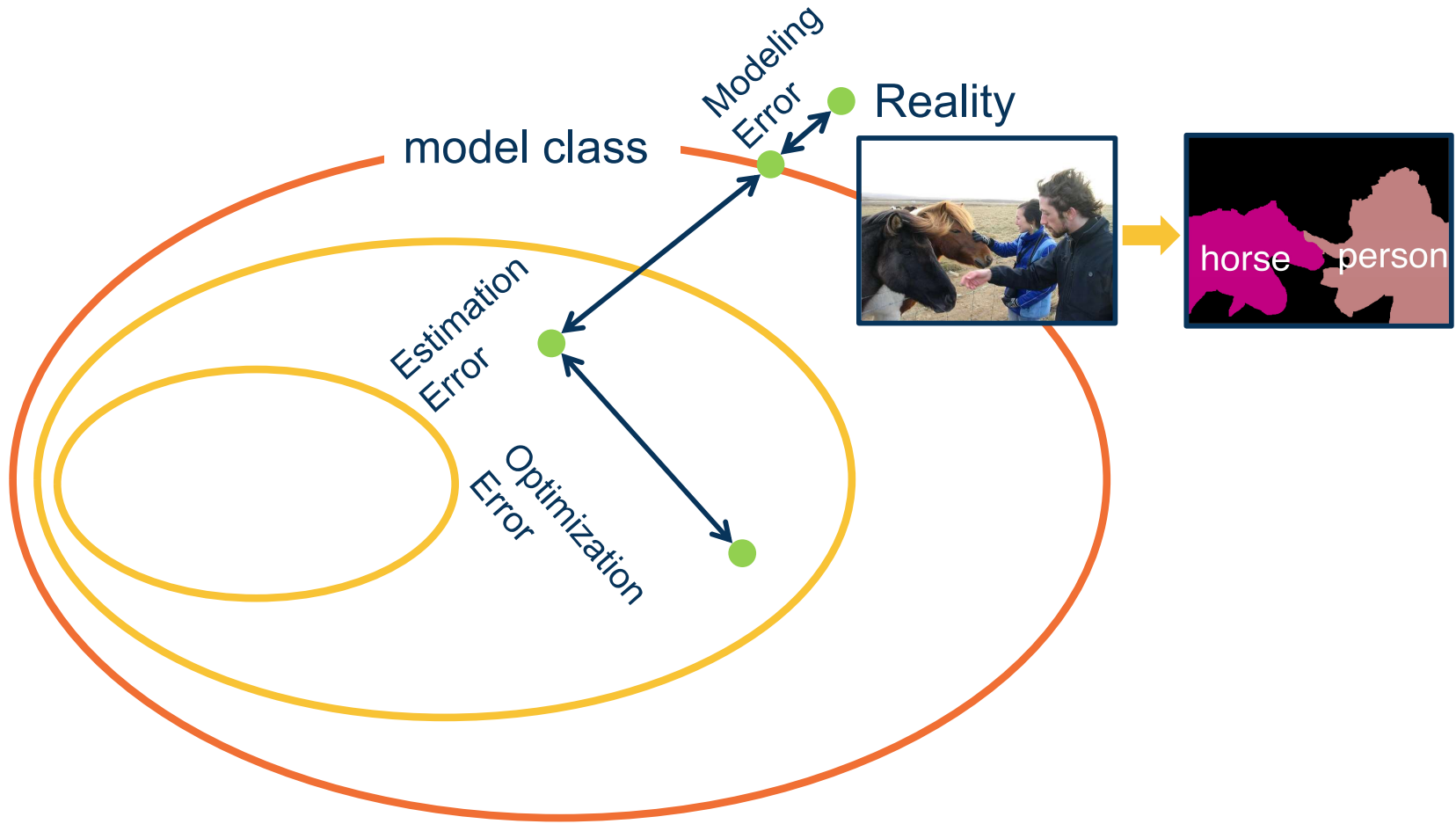
model class

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

**Generalization**

Georgia Tech

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

20 years of research in Learning Theory oversimplified:
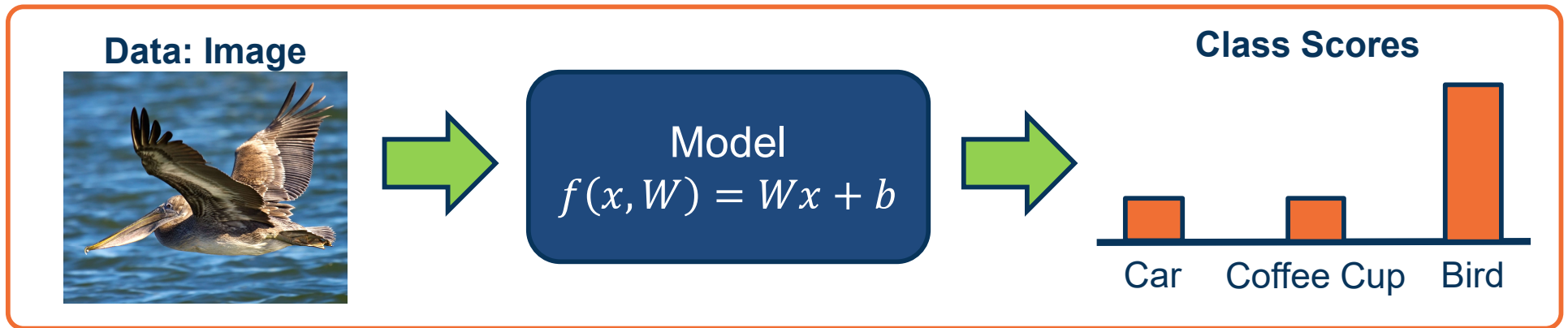
If you have:
  Enough training data D
  and H is not too complex
  then *probably* we can generalize to unseen test data

  **Caveats:** A number of recent empirical results question our intuitions built from this clean separation.

**Zhang et al., Understanding deep learning requires rethinking generalization**

**Guarantees**

Georgia Tech

**Data: Image**

**Class Scores**

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

Car    Coffee Cup    Bird

**Input** $\{X, Y\}$ **where:**

- ◆ $X$ is an image

- ◆ $Y$ is a **ground truth label** annotated by an expert (human)

- ◆ $f(x, W) = Wx + b$ is our model, chosen to be a linear function in this case

- ◆ $W$ and $b$ are the parameters (**weights**) of our model that must be learned

**Example: Image Classification**

Georgia Tech

Input image is **high-dimensional**

⬡ For example $n$=512 so 512x512 image = **262,144** pixels

⬡ Learning a classifier with high-dimensional inputs is hard

Before deep learning, it was typical to perform **feature engineering**

⬡ Hand-design algorithms for converting raw input into a lower-dimensional set of features

**Input Image**



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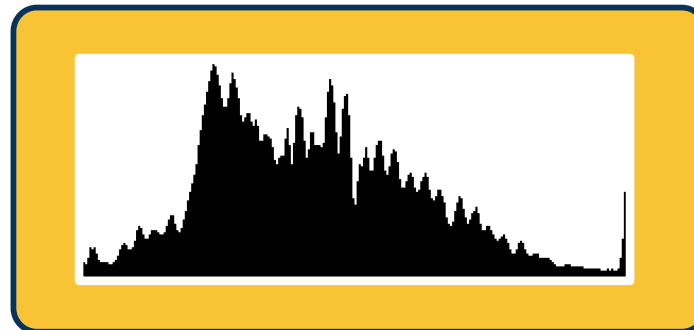
Georgia Tech

# Example: Color histogram

◆ Vector of numbers representing number of pixels fitting within each bin

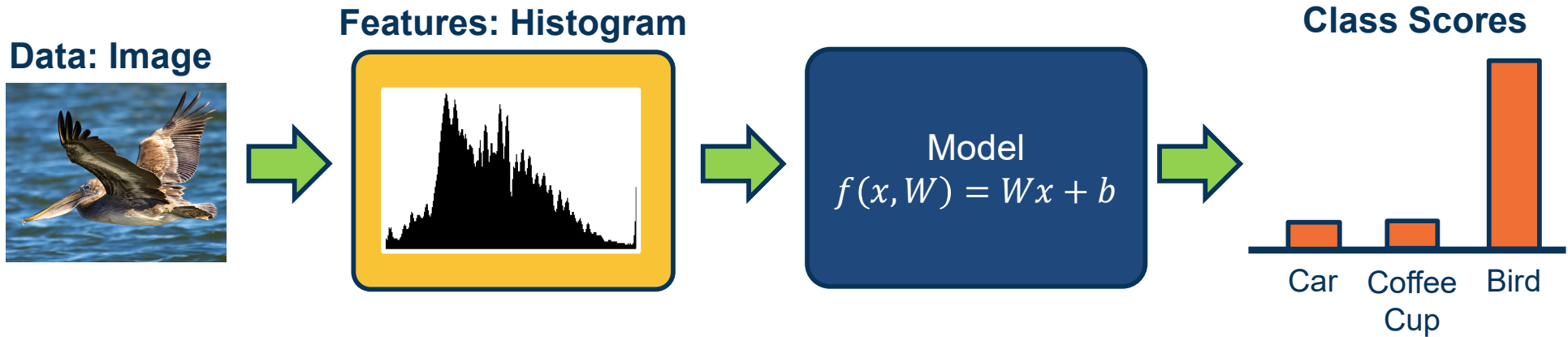◆ We will later see that learning the feature representation itself is much more effective

**Data: Image**          **Features: Histogram**

**Data: Image**

**Features: Histogram**

**Class Scores**

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

Car    Coffee    Bird
       Cup
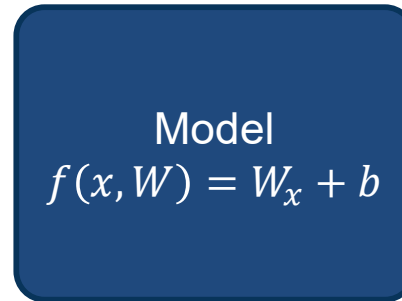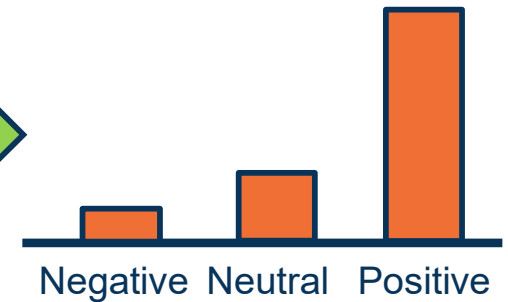
**Input $\{X, Y\}$ where:**

- $X$ is an **image histogram**

- $Y$ is a **ground truth label represented a probability distribution**

- $f(x, W) = Wx + b$ is our model, chosen to be a linear function in this case

- $W$ and $b$ are the **weights** of our model that must be learned

**Example: Image Classification**

Georgia Tech

## Data: Text



TV Fan 88 ✔ @tv_fan1988 · 7h

This television is amazing! The picture quality is extremely high.

💬 646    🔁 1,2k    ♡ 16,4k    ⬆

## Model
$$f(x, W) = W_x + b$$

## Class Scores



Negative  Neutral  Positive

---

**Input $\{X, Y\}$ where:**

- $X$ is a sentence
- $Y$ is a **ground truth label** annotated by an expert (human)
- $f(x, W) = Wx + b$ is our model, chosen to be a linear function in this case
- $W$ and $b$ are the **weights** of our model that must be learned

## Word Histogram

| Word | Count |
|------|-------|
| this | 1 |
| that | 0 |
| is | 2 |
| ... | |
| extremely | 1 |
| hello | 0 |
| onomatopoeia | 0 |
| ... | |

**Example: Image Classification**

Georgia Tech

- 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

Class Scores

Car   Coffee   Bird
      Cup

Loss Function

Optimizer

Class Scores

Car   Coffee   Bird
      Cup

Data: Image

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

Neural Network

Linear classifiers

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

**Deep Learning as Legos**

Georgia Tech

Simple Function

# Linear Classification and Regression

**Simple linear classifier:**

◆ Calculate score:
$$f(x, w) = w \cdot x + b$$

◆ Binary classification rule ($w$ is a vector):
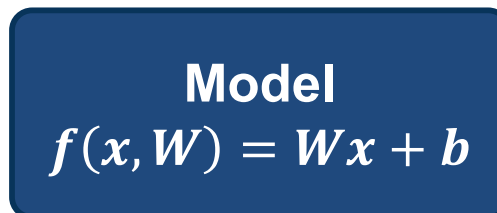
$$y = \begin{cases} 1 & \text{if } f(x, w) >= 0 \\ 0 & \text{otherwise} \end{cases}$$

◆ For multi-class classifier take class with highest (max) score
$$f(x, W) = Wx + b$$

Georgia Tech

**Data: Image**



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

**Class Scores**

Car    Coffee    Bird
          Cup

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

**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 $\longrightarrow$

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

$$W \qquad\qquad x \qquad b$$

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

**Weights**

Georgia Tech

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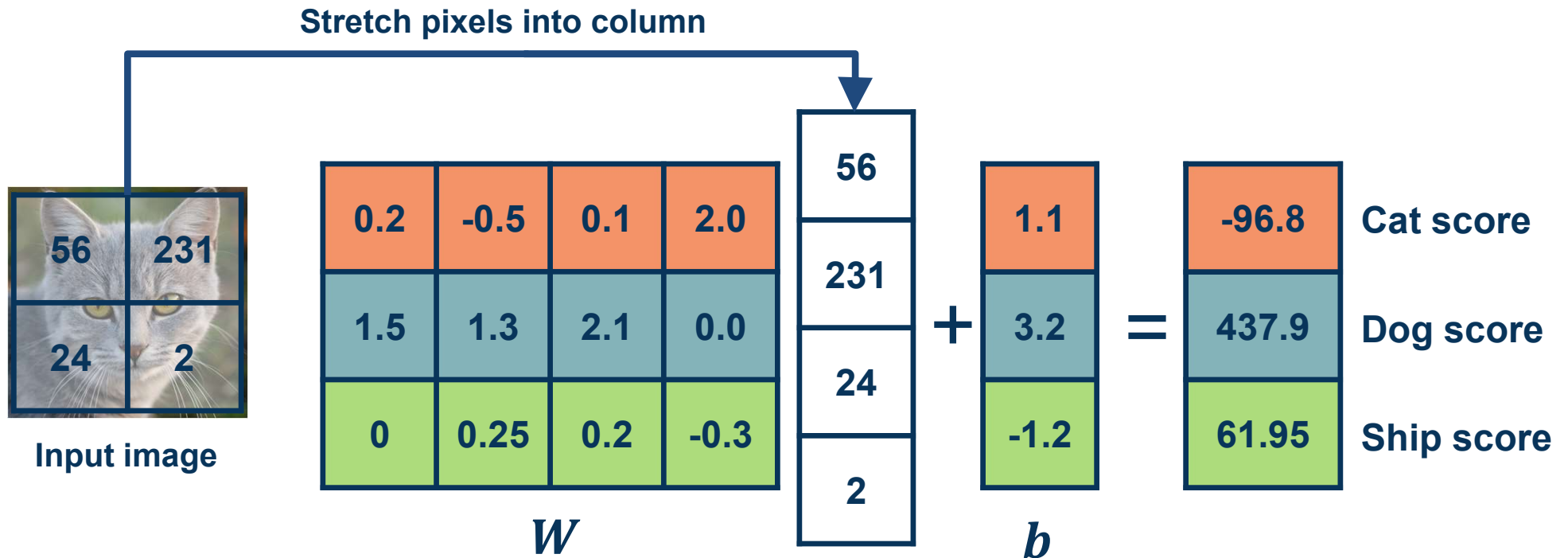
**Model**
$$f(x, W) = Wx + b$$

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

$$W \qquad\qquad x$$

Georgia Tech

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

Stretch pixels into column



Input image

| | | | | |
|---|---|---|---|---|
| 0.2 | -0.5 | 0.1 | 2.0 |
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

**W**

| 56 |
|---|
| 231 |
| 24 |
| 2 |

+

| 1.1 |
|---|
| 3.2 |
| -1.2 |

***b***

=

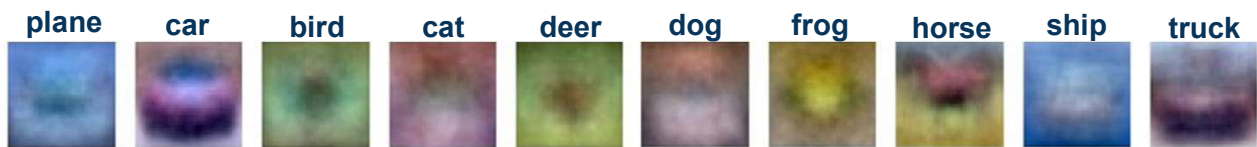| -96.8 | Cat score |
|---|---|
| 437.9 | Dog score |
| 61.95 | Ship score |

**Example**

Georgia Tech

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

| plane | car | bird | cat | deer | dog | frog | horse | ship | truck |

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

## Interpreting a Linear Classifier

Georgia Tech

# Geometric Viewpoint

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

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

*Plot created using Wolfram Cloud*

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

**Interpreting a Linear Classifier**

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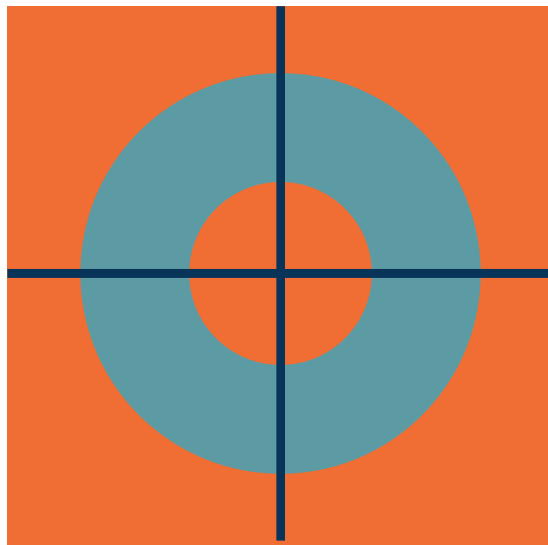number of pixels > 0 odd

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

**Class 1:**
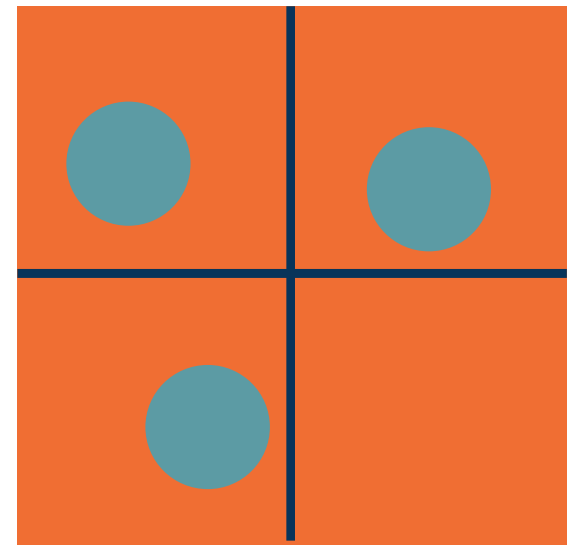1 < = L2 norm < = 2

**Class 2:**
Everything else
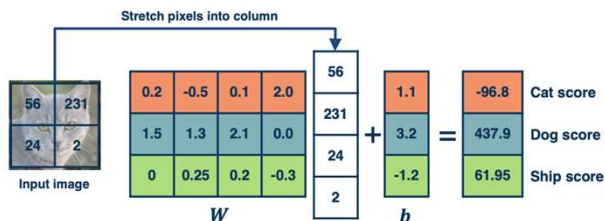
**Class 1:**
Three modes

**Class 2:**
Everything else



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

**Hard Cases for a Linear Classifier**

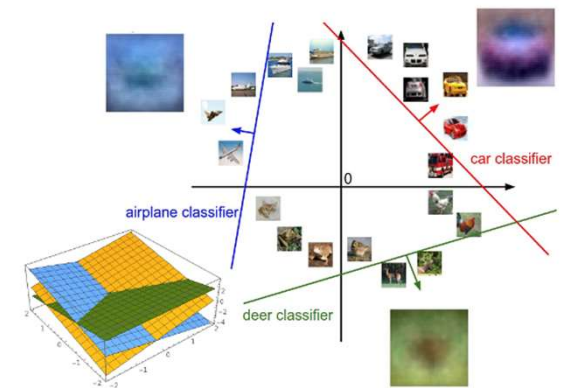Georgia Tech

**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*

**Linear Classifier: Three Viewpoints**

Georgia Tech

- We will learn complex, parameterized functions
    - Start w/ simple building blocks such as linear classifiers

- Key is to learn parameters, but learning is hard
    - Sources of generalization error
    - Add bias/assumptions via architecture, loss, optimizer

- Components of parametric classifiers:
    - Input/Output, Model (function), Loss function, Optimizer
    - Example: Image/Label, Linear Classifier, Hinge Loss, ?

Georgia Tech