# CS 4803 / 7643: Deep Learning

Topics:
- Linear Classifiers
- Loss Functions

Dhruv Batra

Georgia Tech

# Administrativia

- ## Notes and readings on class webpage
  - [https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/](https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/)

- ## Issues from PS0 submission
  - Instructions not followed = not graded

1. We will be using Gradescope to collect your assignments. Please read the following instructions for submitting to Gradescope carefully! Failure to follow these instructions may result in parts of your assignment not being graded. We will not entertain regrading requests for failure to follow instructions.

   - For Section 1: Multiple Choice Questions, it is mandatory to use the LaTeX template provided on the class webpage (`https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/assets/ps0.zip`). For every question, there is only one correct answer. To mark the correct answer, change `\choice` to `\CorrectChoice`
   - For Section 2: Proofs, each problem/sub-problem is in its own page. This section has 5 total problems/sub-problems, so you should have 5 pages corresponding to this section. Your answer to each sub-problem should fit in its corresponding page.
   - For Section 2, LaTeX'd solutions are strongly encouraged (solution template available at `https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/assets/ps0.zip`), but scanned handwritten copies are acceptable. If you scan handwritten copies, please make sure to append them to the pdf generated by LaTeX for Section 1.

# What is the collaboration policy?

- Collaboration

3. We generally encourage you to collaborate with other students. You may talk to a friend, discuss the questions and potential directions for solving them. However, you need to write your own solutions and code separately, and *not* as a group activity. Please list the students you collaborated with.
**Exception: PS0 is meant to serve as a background preparation test. You must NOT collaborate on PS0.**

- Zero tolerance on plagiarism

## Academic Misconduct Process

### ⓘ PROCESS

- How Does the Academic Misconduct Process Begin?
- Who Can Hear My Case?
- What Can I Do To Prepare?
- Office of Student Integrity Meeting Process
- Possible Outcome of the Process
- Faculty Notifications

Any person may file a complaint against a student for violation of the Student Code of Conduct. The complaint should be sent to OSI using the incident referral form. An OSI staff member may contact you during the investigation of the case for more information and to keep updated on the status of the process. Alternatively, the instructor of record for the course may hold a Faculty Conference (refer to Faculty Conference page for more information). The complaint should be submitted as soon as possible after the event takes place or when it is reasonably discovered, no later than thirty (30) business days following the discovery of the incident. In extraordinary circumstances, OSI may waive this timeline.

Students who wish to report an alleged violation of the Student Code of Conduct should notify their instructor. Students may also speak to a member of the Honor Advisory Council or direct questions to staff members in the Office of Student Integrity.

# Recap from last time

# **Image Classification**: A core task in Computer Vision

(assume given set of discrete labels)
{dog, cat, truck, plane, ...}

⟶ cat

# An image classifier

```python
def classify_image(image):
    # Some magic here?
    return class_label
```
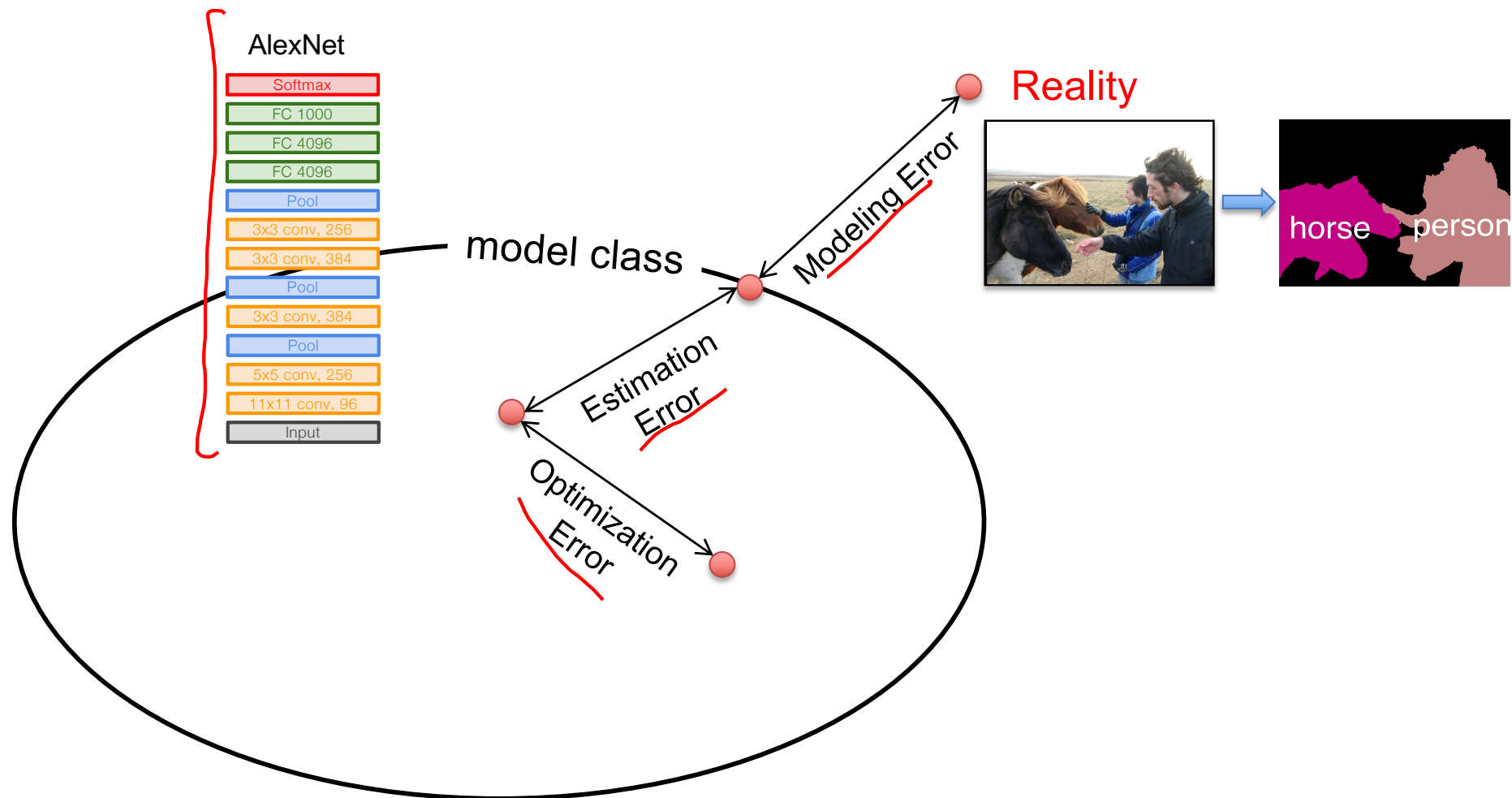
Unlike e.g. sorting a list of numbers,

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

# Supervised Learning

- Input: x                                    (images, text, emails…)
- Output: y                                   (spam or non-spam…)

- (Unknown) Target Function
  - $f: X \rightarrow Y$                       (the "true" mapping / reality)

- Data
  - $\{ (x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N) \}$

- Model / Hypothesis Class
  - $H = \{h: X \rightarrow Y\}$
  - e.g. $y = h(x) = \text{sign}(w^T x)$

- Loss Function
  - How good is a model wrt my data D?

- Learning = Search in hypothesis space
  - Find best h in model class.

# Error Decomposition



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

# First classifier: **Nearest Neighbor**

```python
def train(images, labels):
    # Machine learning!
    return model
```

Memorize all data and labels

```python
def predict(model, test_images):
    # Use model to predict labels
    return test_labels
```

Predict the label of the most similar training image

# Nearest Neighbours

# Instance/Memory-based Learning

Four things make a memory based learner:

- *A distance metric* $d(\vec{x}_i, \vec{x}_j)$

- *How many nearby neighbors to look at?*

- *A weighting function (optional)*

- *How to fit with the local points?*

# Hyperparameters

| Your Dataset |
|:---:|

**Idea #4**: **Cross-Validation**: Split data into **folds**,
try each fold as validation and average the results

| fold 1 | fold 2 | fold 3 | fold 4 | fold 5 | test |
|:---:|:---:|:---:|:---:|:---:|:---:|
| fold 1 | fold 2 | fold 3 | fold 4 | fold 5 | test |
| fold 1 | fold 2 | fold 3 | fold 4 | fold 5 | test |

Useful for small datasets, but not used too frequently in deep learning

# Problems with Instance-Based Learning

- 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 neighbour search

- Doesn't work well when large number of irrelevant features
  - Distances overwhelmed by noisy features

- Curse of Dimensionality
  - Distances become meaningless in high dimensions
  - (See proof in next lecture)

# Plan for Today

- Linear Classifiers
  - Linear scoring functions

- Loss Functions
  - Multi-class hinge loss
  - Softmax cross-entropy loss

# Linear Classification

Neural Network
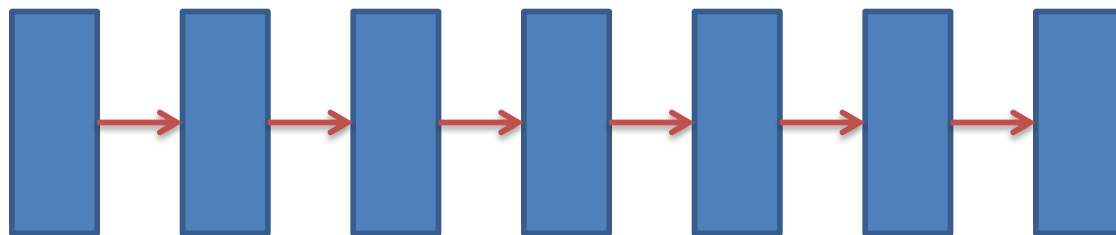
Linear classifiers



This image is CC0 1.0 public domain

# Visual Question Answering



Image Embedding (VGGNet)

Convolution Layer + Non-Linearity · Pooling Layer · Convolution Layer + Non-Linearity · Pooling Layer · Fully-Connected MLP

4096-dim

$\vec{i} \in \mathbb{R}^{d_1}$

Neural Network
Softmax
over top K answers

$h_1^{(2)}$, $h_2^{(2)}$, $h_3^{(2)}$, +1

Input (Features II) · Softmax classifier

$P(y = 0 \mid x)$
$P(y = 1 \mid x)$
$P(y = 2 \mid x)$

Question Embedding (LSTM)

"How many horses are in this image?"

$\vec{q} \in \mathbb{R}^{d_2}$
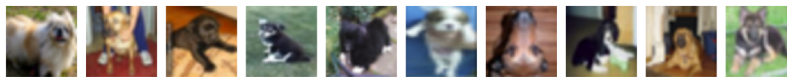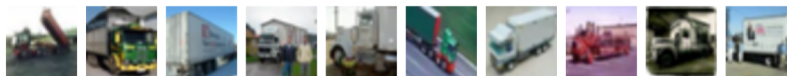
# Recall CIFAR10

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

**50,000** training images
each image is **32x32x3**

**10,000** test images.

# Parametric Approach

Image



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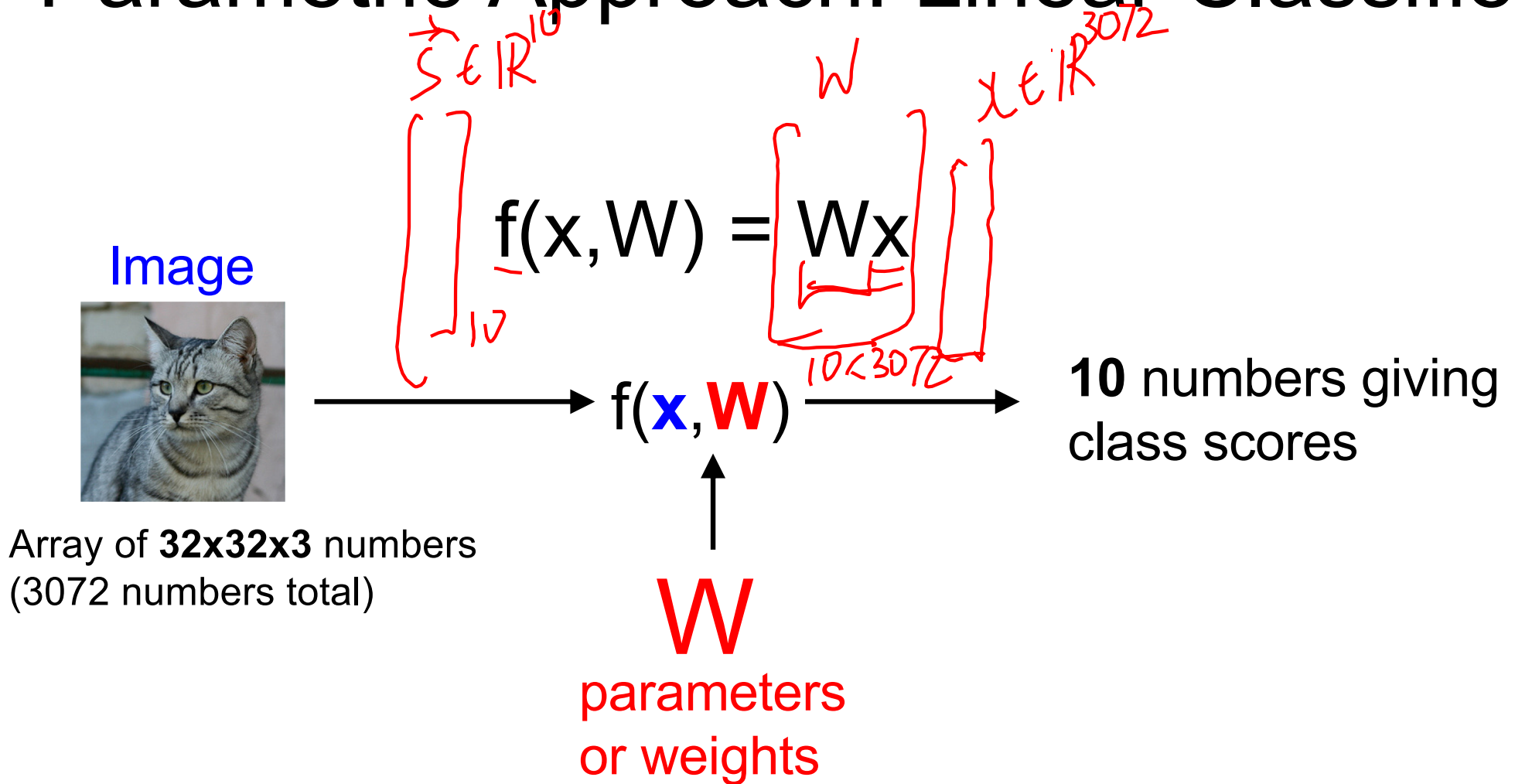(3072 numbers total)

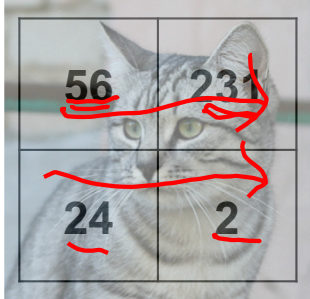$x \in \mathbb{R}^d$

f(**x**, **W**)

↑

**W**
parameters
or weights

**10** numbers giving
class scores

# Parametric Approach: Linear Classifier

**Image**



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

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

$\vec{s} \in \mathbb{R}^{10}$

$W$

$x \in \mathbb{R}^{3072}$

$10$

$10 \times 3072$

f(**x**,**W**) $\longrightarrow$ **10** numbers giving class scores

↑

**W**
parameters
or weights

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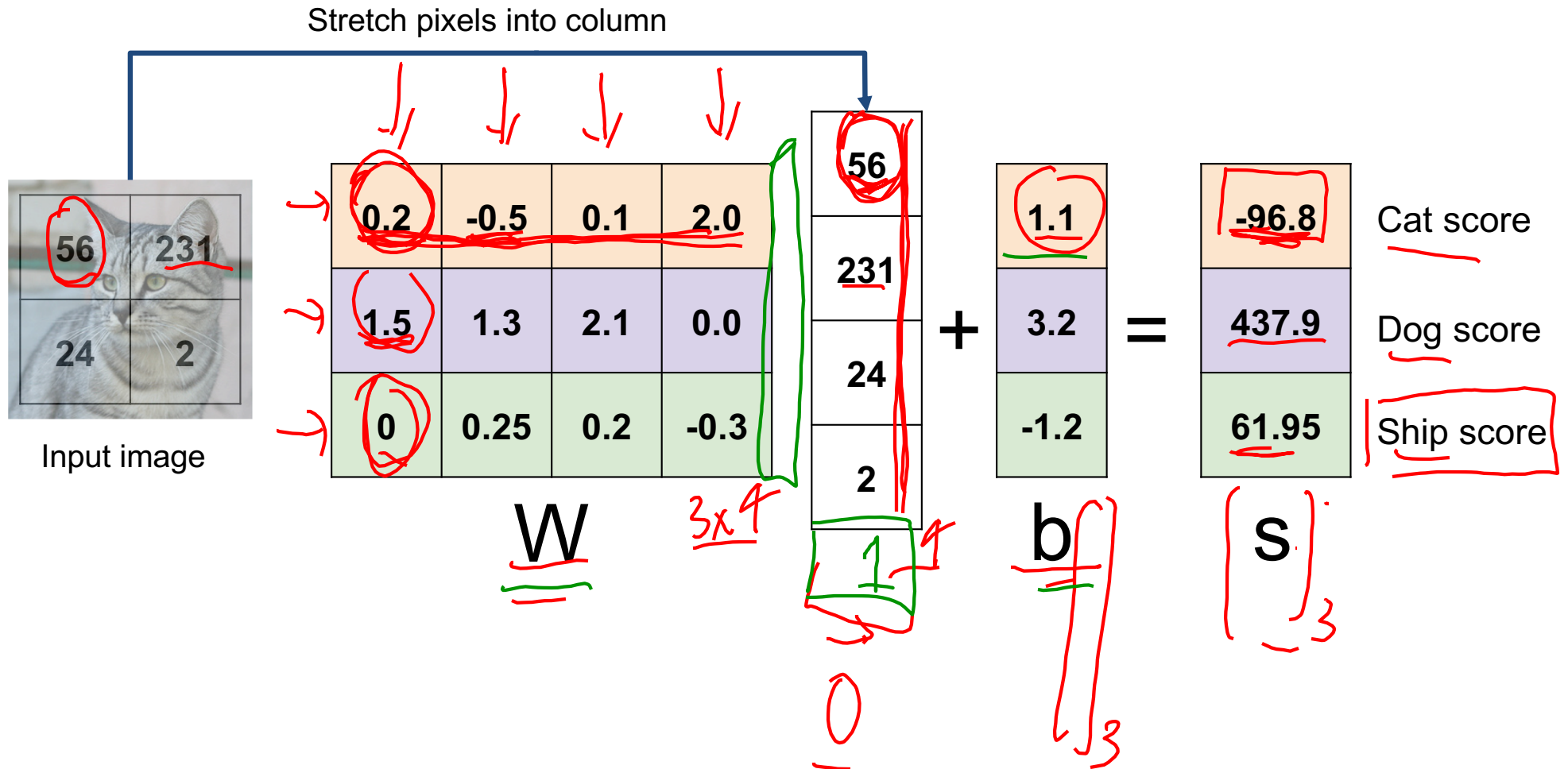
Stretch pixels into column



Input image

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

# 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        3x4

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

1

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

b        3

+

=

| -96.8 |
|------|
| 437.9 |
| 61.95 |

s        3

Cat score

Dog score

Ship score

56   231

24   2

| 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 |

$$x_i$$

$$+$$

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

$$b$$

$$\longleftrightarrow$$

| 0.2 | -0.5 | 0.1 | 2.0 | 1.1 |
|-----|------|-----|-----|-----|
| 1.5 | 1.3 | 2.1 | 0.0 | 3.2 |
| 0 | 0.25 | 0.2 | -0.3 | -1.2 |

$$W \qquad b$$

new, single W

| 56 |
|----|
| 231 |
| 24 |
| 2 |
| 1 |

$$x_i$$

Image Credit: Andrej Karpathy, CS231n

# Parametric Approach: Linear Classifier

**3072x1**

$$f(x,W) = Wx + b$$  **10x1**

**10x1**    **10x3072**

Image



f(**x**,**W**) → **10** numbers giving class scores

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

**W**
parameters
or weights

# Error Decomposition



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

# Error Decomposition



Reality

Modeling Error

Multi-class Logistic Regression

Softmax

FC HxWx3

Input

Estimation Error

Optimization Error = 0

model class

horse    person

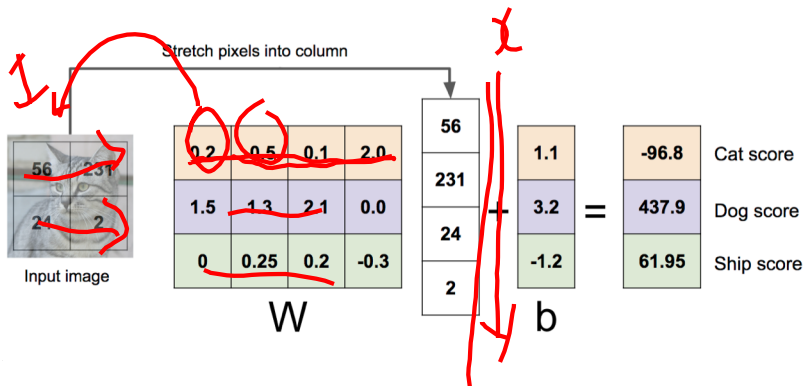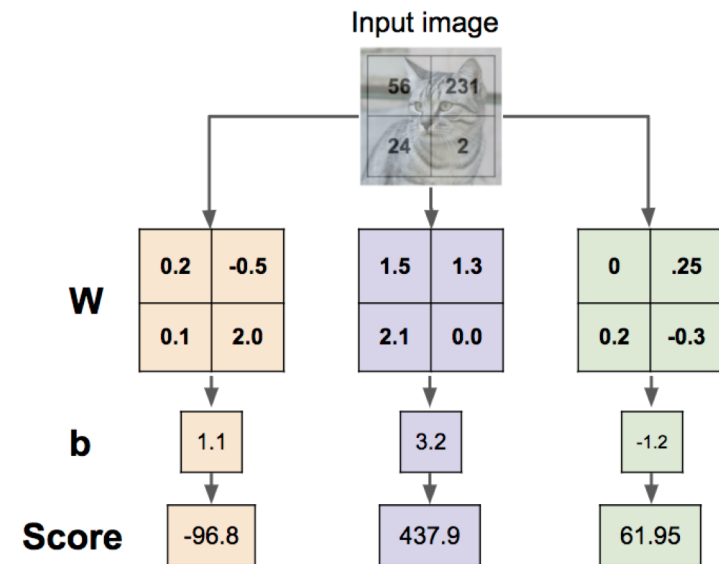# Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Algebraic Viewpoint

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

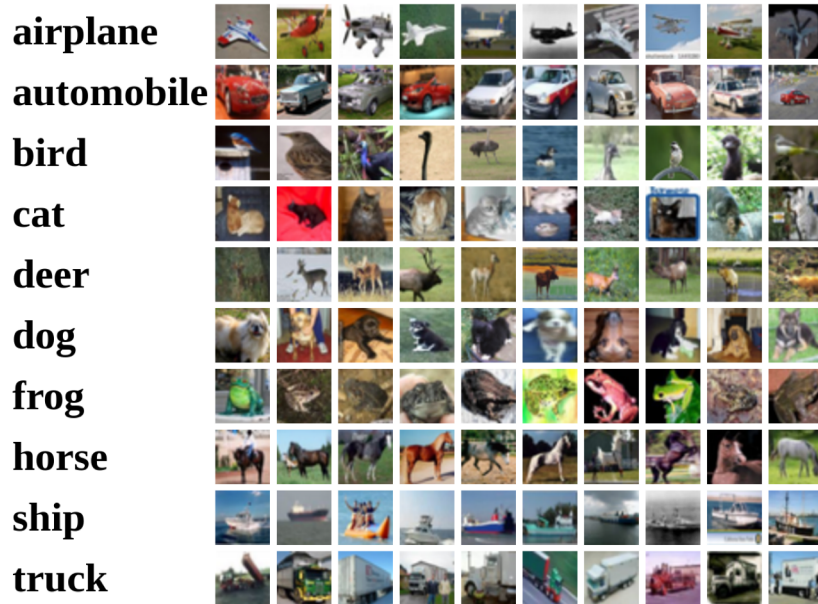# Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

## Algebraic Viewpoint

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



Input image

W

b

Score

# Interpreting a Linear Classifier

# Interpreting a Linear Classifier: Visual Viewpoint



airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

Input image

| 56 | 231 |
| 24 | 2 |

W

| 0.2 | -0.5 |
| 0.1 | 2.0 |

| 1.5 | 1.3 |
| 2.1 | 0.0 |

| 0 | .25 |
| 0.2 | -0.3 |

b

| 1.1 |

| 3.2 |

| -1.2 |

Score

| -96.8 |

| 437.9 |

| 61.95 |

plane    car    bird    cat    deer    dog    frog    horse    ship    truck

# Interpreting a Linear Classifier: Geometric Viewpoint



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

airplane classifier

car classifier

deer classifier

$w^T x = 0$

$W_{car}^T x + b_{car}$

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

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n
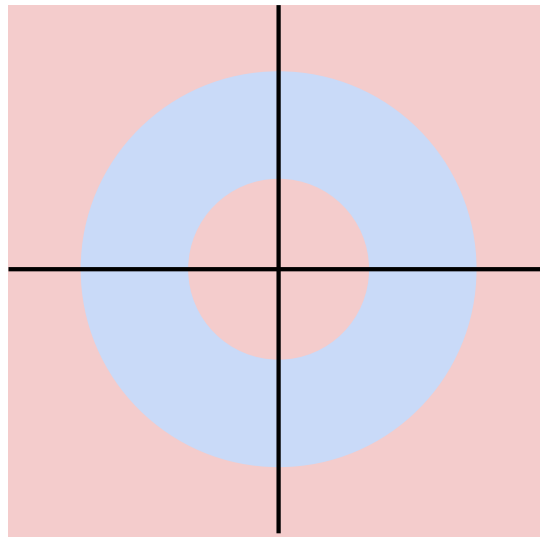
# Hard cases for a linear classifier

**Class 1**:
First and third quadrants

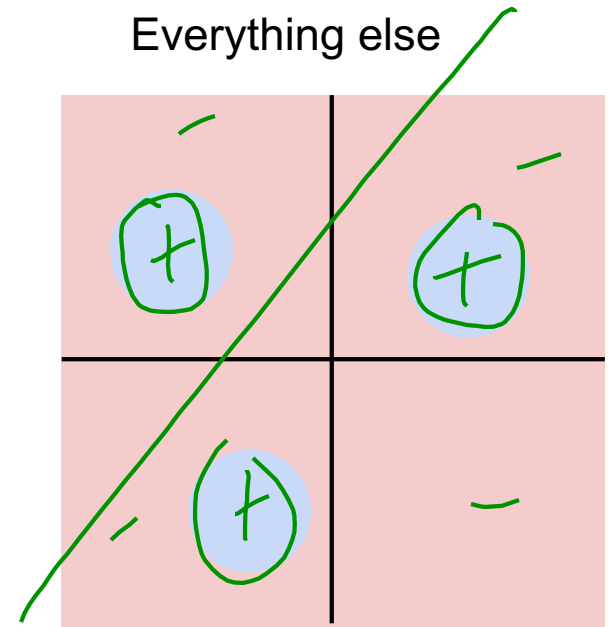**Class 2**:
Second and fourth quadrants

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

**Class 2**:
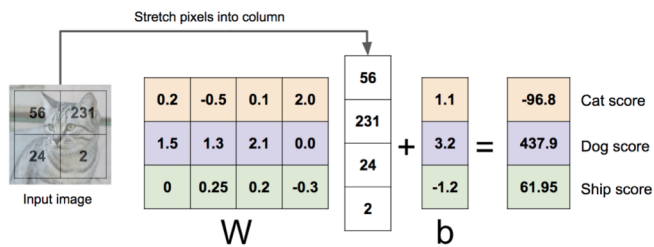Everything else

**Class 1**:
Three modes

**Class 2**:
Everything else

# Linear Classifier: Three Viewpoints

## Algebraic Viewpoint
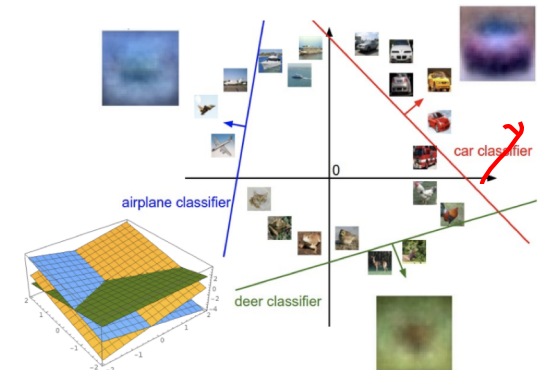
$$f(x,W) = Wx$$



## Visual Viewpoint

One template per class



## Geometric Viewpoint

Hyperplanes cutting up space

# **So far**: Defined a (linear) score function

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

Example class scores for 3 images for some W:

How can we tell whether this W is good or bad?

|  | | | |
|---|---|---|---|
| airplane | -3.45 | -0.51 | 3.42 |
| automobile | -8.87 | **6.04** | 4.64 |
| bird | 0.09 | 5.31 | 2.65 |
| cat | **2.9** | -4.22 | 5.1 |
| deer | 4.48 | -4.19 | 2.64 |
| dog | 8.02 | 3.58 | 5.55 |
| frog | 3.78 | 4.49 | **-4.34** |
| horse | 1.06 | -4.37 | -1.5 |
| ship | -0.36 | -2.09 | -4.79 |
| truck | -0.72 | -2.93 | 6.14 |

# **So far**: Defined a (linear) <u>score function</u>

| | | | |
|---|---|---|---|
| airplane | -3.45 | -0.51 | 3.42 |
| automobile | -8.87 | **6.04** | 4.64 |
| bird | 0.09 | 5.31 | 2.65 |
| cat | **2.9** | -4.22 | 5.1 |
| deer | 4.48 | -4.19 | 2.64 |
| dog | 8.02 | 3.58 | 5.55 |
| frog | 3.78 | 4.49 | **-4.34** |
| horse | 1.06 | -4.37 | -1.5 |
| ship | -0.36 | -2.09 | -4.79 |
| truck | -0.72 | -2.93 | 6.14 |

TODO:

1. Define a **loss function** that quantifies our unhappiness with the scores across the training data.

2. Come up with a way of efficiently finding the parameters that minimize the loss function. **(optimization)**

# Supervised Learning

- Input: x                                    (images, text, emails…)
- Output: y                                   (spam or non-spam…)

- (Unknown) Target Function
  - f: X → Y                                  (the "true" mapping / reality)

- Data
  - $(x_1, y_1)$, $(x_2, y_2)$, …, $(x_N, y_N)$

- Model / Hypothesis Class
  - {h: X → Y}
  - e.g. $y = h(x) = \text{sign}(w^T x)$

- Loss Function
  - How good is a model wrt my data D?

- Learning = Search in hypothesis space
  - Find best h in model class.

# Loss Functions

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



|       | cat   | car   | frog  |
|-------|-------|-------|-------|
| cat   | **3.2**   | 1.3   | 2.2   |
| car   | 5.1   | **4.9**   | 2.5   |
| frog  | -1.7  | 2.0   | **-3.1**  |

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



| | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

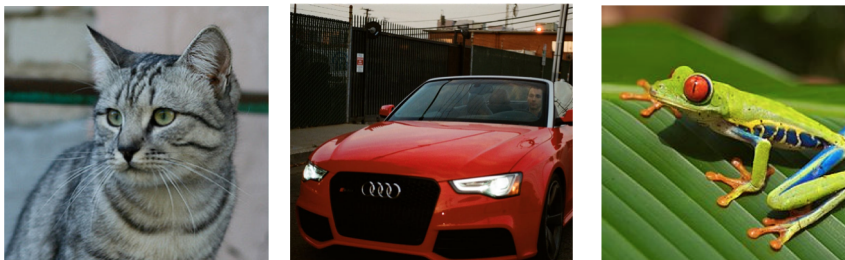A **loss function** tells how good our current classifier is

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^{N}$$

Where $x_i$ is image and $y_i$ is (integer) label

Loss over the dataset is a sum of loss over examples:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



|       |         |         |         |
|-------|---------|---------|---------|
| cat   | **3.2** | 1.3     | 2.2     |
| car   | 5.1     | **4.9** | 2.5     |
| frog  | -1.7    | 2.0     | **-3.1**|

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$
$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat

car

frog

$$
\begin{bmatrix} 3.2 \\ 5.1 \\ -1.7 \end{bmatrix}
$$

**Multiclass SVM loss:**

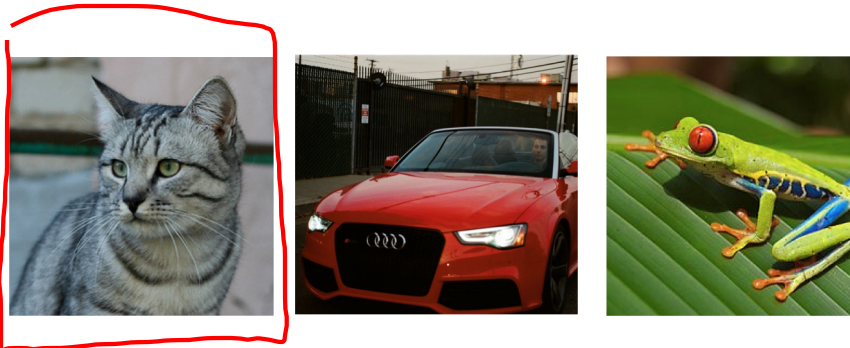Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

$$
L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}
$$

$$y = \max(0, x)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

| | |
|---|---|
| cat | **3.2** |
| car | 5.1 |
| frog | -1.7 |

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



**Multiclass SVM loss:**



"Hinge loss"

|      |      |      |      |
|------|------|------|------|
| cat  | **3.2** | 1.3 | 2.2 |
| car  | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

**Multiclass SVM loss:**



"Hinge loss"

|        |       |       |       |
|--------|-------|-------|-------|
| cat    | **3.2** | 1.3   | 2.2   |
| car    | 5.1   | **4.9** | 2.5   |
| frog   | -1.7  | 2.0   | **-3.1** |

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$
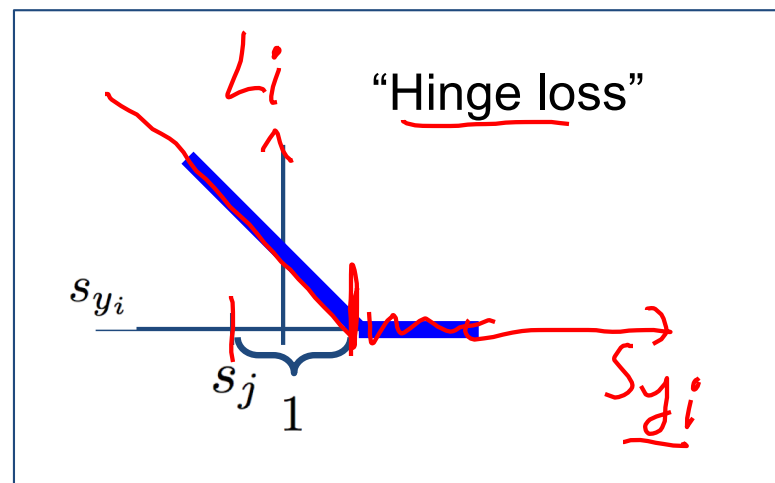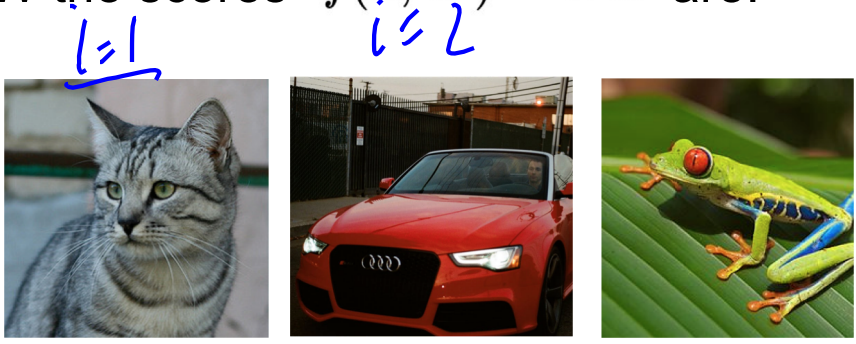
$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

delta

score

scores for other classes

score for correct class

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$



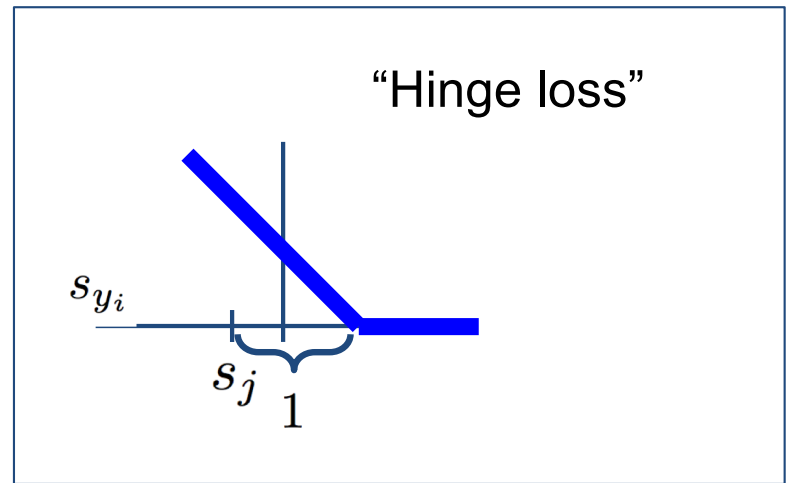|  | | | |
|------|------|------|------|
| cat  | **3.2** | 1.3 | 2.2 |
| car  | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.
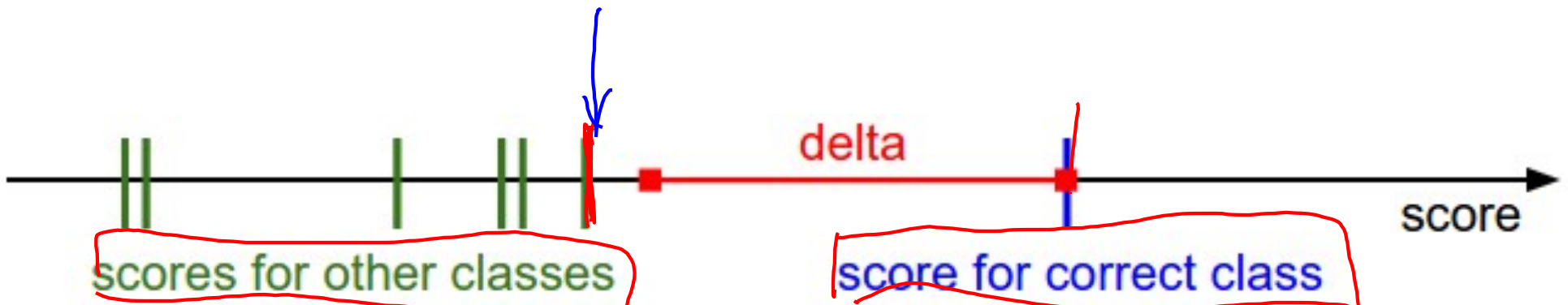With some W the scores $f(x, W) = Wx$ are:



|  |  |  |  |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 |  |  |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 5.1 - 3.2 + 1)
  +max(0, -1.7 - 3.2 + 1)
= max(0, 2.9) + max(0, -3.9)
= 2.9 + 0
= 2.9

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

*in 2*



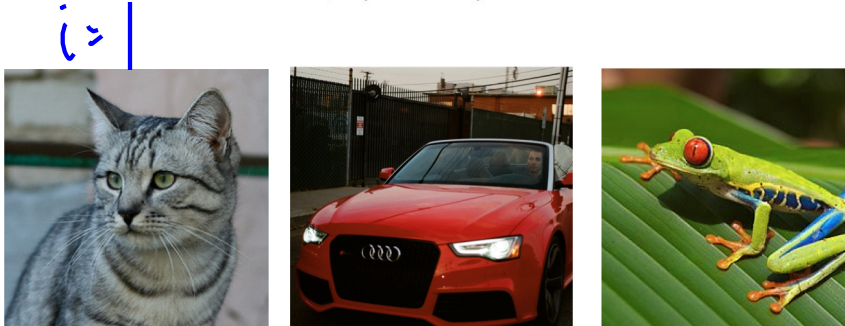| | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 1.3 - 4.9 + 1)
    +max(0, 2.0 - 4.9 + 1)
= max(0, -2.6) + max(0, -1.9)
= 0 + 0
= 0

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

|       |      |      |      |
|-------|------|------|------|
| cat   | **3.2** | 1.3 | 2.2 |
| car   | 5.1 | **4.9** | 2.5 |
| frog  | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 12.9 |

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 2.2 - (-3.1) + 1)
   +max(0, 2.5 - (-3.1) + 1)
= max(0, 6.3) + max(0, 6.6)
= 6.3 + 6.6
= 12.9

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



|        | cat  | car  | frog |
|--------|------|------|------|
| cat    | **3.2**  | 1.3  | 2.2  |
| car    | 5.1  | **4.9**  | 2.5  |
| frog   | -1.7 | 2.0  | **-3.1** |
| Losses: | 2.9  | 0    | 12.9 |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Loss over full dataset is average:

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i$$

L = (2.9 + 0 + 12.9)/3

= **5.27**

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



|      |      |      |      |
|------|------|------|------|
| cat  | **3.2** | 1.3  | 2.2  |
| car  | 5.1  | **4.9** ±ξ | 2.5  |
| frog | -1.7 | 2.0  | **-3.1** |
| Losses: | 2.9 | 0 | 12.9 |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
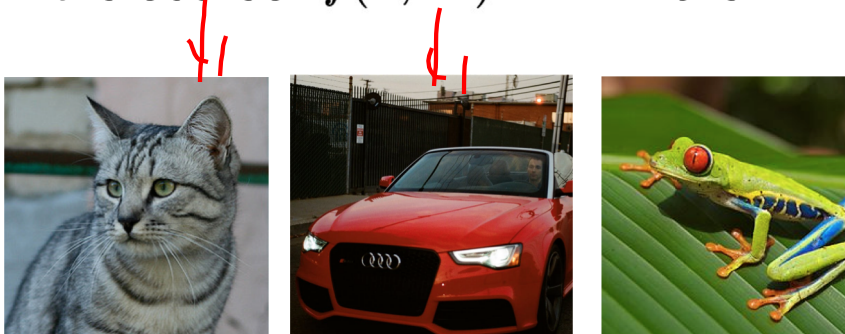where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

$L2$ the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q: What happens to loss if car image scores change a bit?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



| | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 12.9 |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

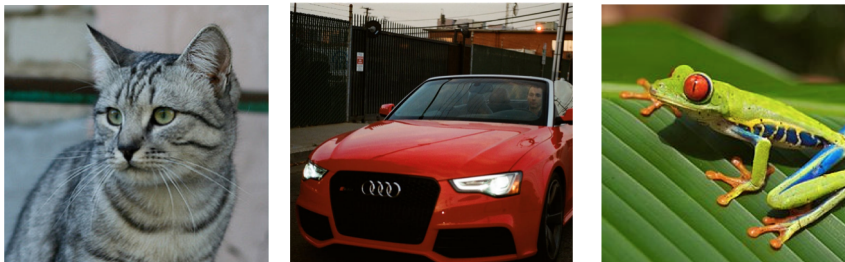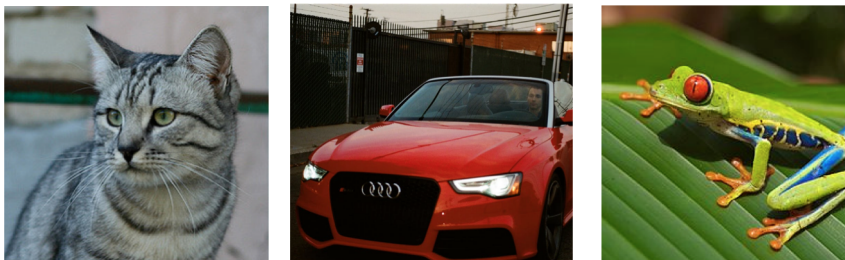$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q2: what is the min/max possible loss?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$



|  | | | |
|---|---|---|---|
| cat | **3.2** 0 | 1.3 | 2.2 |
| car | 5.1 0 | **4.9** | 2.5 |
| frog | -1.7 0 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 12.9 |

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q3: At initialization W
is small so all s ≈ 0.
What is the loss?

#classes - 1

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



|  | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 12.9 |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q4: What if the sum
was over all classes?
(including j = y_i)

$$L = \frac{1}{N} \sum_i L_i$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$



| | | | |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 12.9 |

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q5: What if we used mean instead of sum?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,
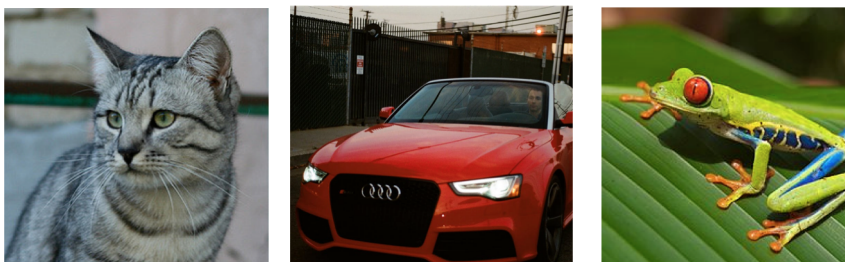
and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

| | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 12.9 |

Q6: What if we used

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)^2$$

$$2(S_j - S_{y_i}) < 0$$

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

$$L = \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_i} \max\left(0, \overset{2}{\boxed{f(x_i; W)_j}} - \overset{2}{f(x_i; W)_{y_i}} + 1\right)$$

E.g. Suppose that we found a W such that L = 0.

Q7: Is this W unique? $2W$

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

$$L = \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

E.g. Suppose that we found a W such that L = 0.
Q7: Is this W unique?

No! 2W is also has L = 0!

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

**Before:**
= max(0, 1.3 - 4.9 + 1)
   +max(0, 2.0 - 4.9 + 1)
= max(0, -2.6) + max(0, -1.9)
= 0 + 0
= 0

**With W twice as large:**
= max(0, 2.6 - 9.8 + 1)
   +max(0, 4.0 - 9.8 + 1)
= max(0, -6.2) + max(0, -4.8)
= 0 + 0
= 0

| | | | |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | |

# Multiclass SVM Loss: Example code

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

```python
def L_i_vectorized(x, y, W):
    scores = W.dot(x)
    margins = np.maximum(0, scores - scores[y] + 1)
    margins[y] = 0
    loss_i = np.sum(margins)
    return loss_i
```

# Softmax Classifier (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

cat **3.2**

car 5.1

frog -1.7

# **Softmax Classifier** (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax function

cat     **3.2**

car     5.1

frog    -1.7

# Softmax Classifier (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ Softmax Function

Probabilities
must be >= 0

| | | exp | |
|-----|------|-----|-------|
| cat | **3.2** | → | **24.5** |
| car | 5.1 | | 164.0 |
| frog | -1.7 | | 0.18 |

unnormalized
probabilities

# **Softmax Classifier** (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

Probabilities must be >= 0

Probabilities must sum to 1

|       | score | exp   | normalize |     |
|-------|-------|-------|-----------|-----|
| cat   | **3.2**   | **24.5**  | **0.13**      | 1   |
| car   | 5.1   | 164.0 | 0.87      | 0   |
| frog  | -1.7  | 0.18  | 0.00      | 0   |

unnormalized probabilities

probabilities

# Softmax Classifier (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ Softmax Function

Probabilities must be >= 0

Probabilities must sum to 1

|       |      |     |        |
|-------|------|-----|--------|
| cat   | **3.2** | **24.5** | **0.13** |
| car   | 5.1  | 164.0 | 0.87 |
| frog  | -1.7 | 0.18  | 0.00 |

exp → normalize →

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

# Softmax Classifier (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

cat **3.2**

car 5.1

frog -1.7

$$L_i = -\log P(Y = y_i | X = x_i)$$

in summary:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

Maximize log-prob of the correct class =
Maximize the log likelihood =
Minimize the negative log likelihood

# Softmax Classifier (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**
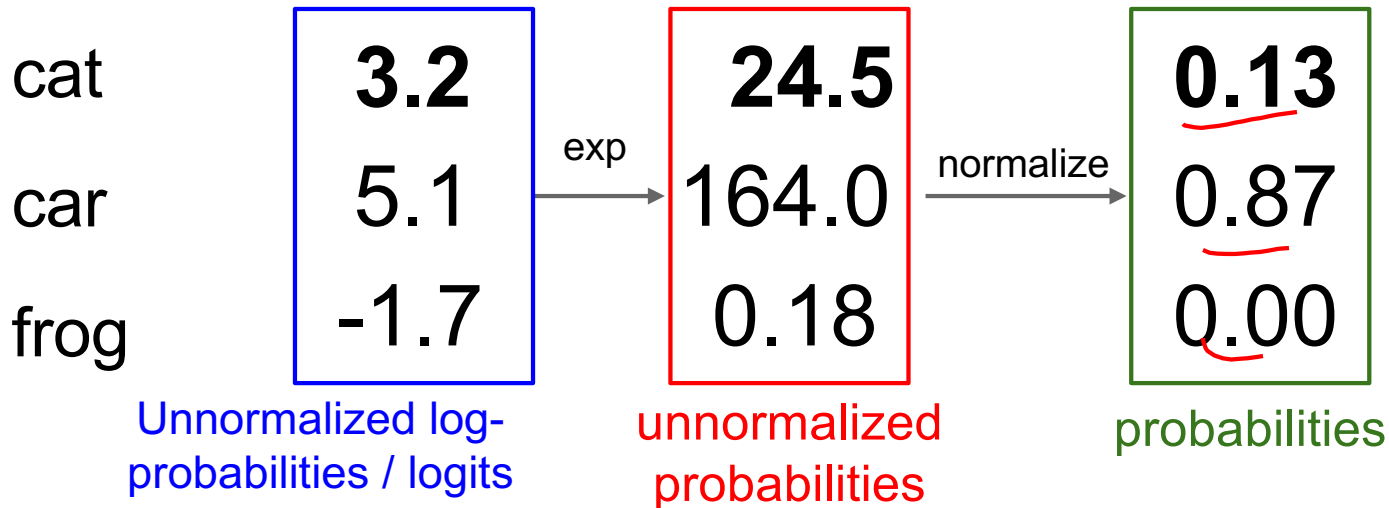
$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

Probabilities must be >= 0

Probabilities must sum to 1

|  | | | |
|---|---|---|---|
| cat | **3.2** | **24.5** | **0.13** |
| car | 5.1 | 164.0 | 0.87 |
| frog | -1.7 | 0.18 | 0.00 |

exp  →  normalize  →

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

# Softmax Classifier (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ Softmax Function

$$L_i = -\log P(Y = y_i | X = x_i)$$

Probabilities must be >= 0

Probabilities must sum to 1

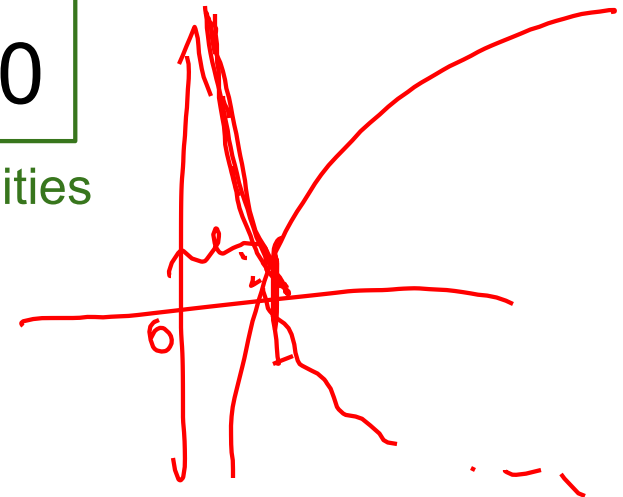| | cat | car | frog |
|---|---|---|---|
| logits | 3.2 | 5.1 | -1.7 |
| exp | 24.5 | 164.0 | 0.18 |
| normalize | 0.13 | 0.87 | 0.00 |

→ L_i = -log(0.13) = **2.04**

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

# Log-Likelihood / KL-Divergence / Cross-Entropy

$$p^* \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \overset{GT}{\underset{y_i}{}}$$

$$\hat{p} = \begin{bmatrix} P(Y=1 \mid \vec{x}_i, \omega) \\ \vdots \\ P(Y=k \mid \vec{x}_i, \omega) \end{bmatrix}$$

$$-1\log \hat{p}(y_i)$$

$$nary\, KL\left(p^* \| \hat{p}\right) = \sum_y p^*(y) \log \frac{p^*(y)}{\hat{p}(y)}$$

$$= \underbrace{\sum p^*(y) \log p^*(y)}_{-H(p^*)} - \underbrace{\sum_y p^*(y) \log \hat{p}(y)}_{+ \underset{min}{} H(p^*, \hat{p})}$$

$$-H(p^*) \rightarrow 0$$

# Log-Likelihood / KL-Divergence / Cross-Entropy

# Log-Likelihood / KL-Divergence / Cross-Entropy

# Softmax Classifier (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ Softmax Function

Maximize probability of correct class

$$L_i = -\log P(Y = y_i | X = x_i)$$

Putting it all together:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

cat **3.2**

car 5.1

frog -1.7

# **Softmax Classifier** (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ Softmax Function

Maximize probability of correct class

$$L_i = -\log P(Y = y_i | X = x_i)$$

Putting it all together:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

cat **3.2**

car 5.1

frog -1.7

Q: What is the min/max possible loss L_i?

# **Softmax Classifier** (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

Maximize probability of correct class

$$L_i = -\log P(Y = y_i | X = x_i)$$

Putting it all together:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

cat **3.2**

car 5.1

frog -1.7

Q: What is the min/max possible loss L_i?
A: min 0, max infinity

# Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ Softmax Function

Maximize probability of correct class

$$L_i = -\log P(Y = y_i | X = x_i)$$

Putting it all together:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

cat    **3.2**

car    5.1

frog    -1.7

Q2: At initialization all s will be approximately equal; what is the loss?

# **Softmax Classifier** (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**

$$\boxed{s = f(x_i; W)}$$

$$\boxed{P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}}$$ Softmax Function

Maximize probability of correct class

$$L_i = -\log P(Y = y_i | X = x_i)$$

Putting it all together:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

cat **3.2**

car 5.1

frog -1.7

Q2: At initialization all s will be approximately equal; what is the loss?
A: log(C), eg log(10) ≈ 2.3

# Softmax vs. SVM

## matrix multiply + bias offset

| 0.01 | -0.05 | 0.1 | 0.05 |
|------|-------|-----|------|
| 0.7 | 0.2 | 0.05 | 0.16 |
| 0.0 | -0.45 | -0.2 | 0.03 |

$$W$$

| -15 |
|-----|
| 22 |
| -44 |
| 56 |

$$x_i$$

$+$

| 0.0 |
|-----|
| 0.2 |
| -0.3 |

$$b$$

$$y_i \quad \boxed{2}$$

### hinge loss (SVM)

| -2.85 |
|-------|
| 0.86 |
| 0.28 |

max(0, -2.85 - 0.28 + 1) +
max(0, 0.86 - 0.28 + 1)
=
**1.58**

### cross-entropy loss (Softmax)

| -2.85 |
|-------|
| 0.86 |
| 0.28 |

*exp* →

| 0.058 |
|-------|
| 2.36 |
| 1.32 |

*normalize* →
(to sum to one)

| 0.016 |
|-------|
| 0.631 |
| 0.353 |

- log(0.353)
=
**0.452**

# Softmax vs. SVM

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \qquad L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$
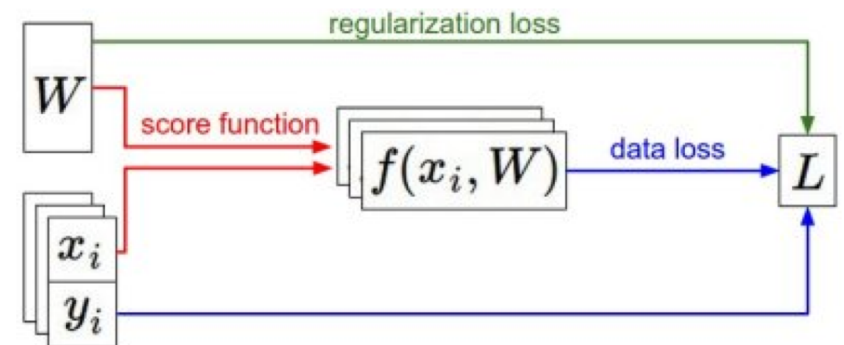
# Recap

- We have some dataset of (x,y)
- We have a **score function:** $\quad s = f(x; W) \overset{\text{e.g.}}{=} Wx$
- We have a **loss function:**

Softmax
$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

SVM
$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$L = \frac{1}{N}\sum_{i=1}^{N} L_i + R(W) \quad \text{Full loss}$$

# Recap

## How do we find the best W?

- We have some dataset of (x,y)
- We have a **score function:** $s = f(x; W) \overset{\text{e.g.}}{=} Wx$
- We have a **loss function**:

Softmax
$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

SVM
$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Full loss
$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + R(W)$$