

# CS 4803 / 7643: Deep Learning

Topics:

- Linear Classifiers
- Loss Functions

Dhruv Batra  
Georgia Tech

# Administrativa

- Notes and readings on class webpage
  - [https://www.cc.gatech.edu/classes/AY2020/cs7643\\_fall/](https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/)
- HW0 solutions and grades released
- 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 L<sup>A</sup>T<sub>E</sub>X template provided on the class webpage ([https://www.cc.gatech.edu/classes/AY2020/cs7643\\_fall/assets/ps0.zip](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, L<sup>A</sup>T<sub>E</sub>X'd solutions are strongly encouraged (solution template available at [https://www.cc.gatech.edu/classes/AY2020/cs7643\\_fall/assets/ps0.zip](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 L<sup>A</sup>T<sub>E</sub>X for Section 1.

Recap from last time |

# Image Classification: A core task in Computer Vision



This image by [Nikita](#) is licensed under [CC-BY 2.0](#)

→

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



cat

Ⓡ

# An image classifier

```
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)

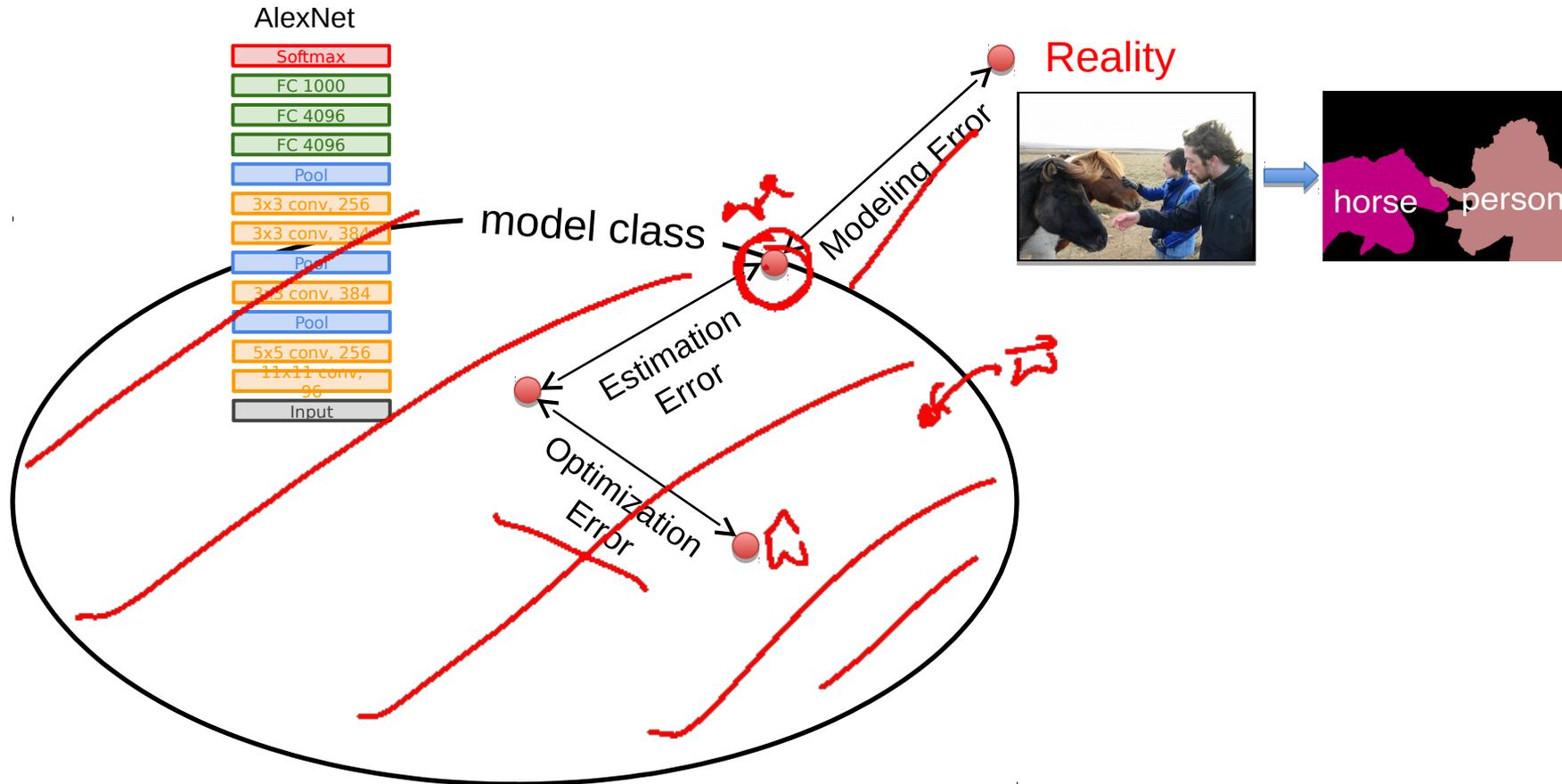
- Data
  - $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

- 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



# First classifier: Nearest Neighbor

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



Memorize all  
data and labels

```
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

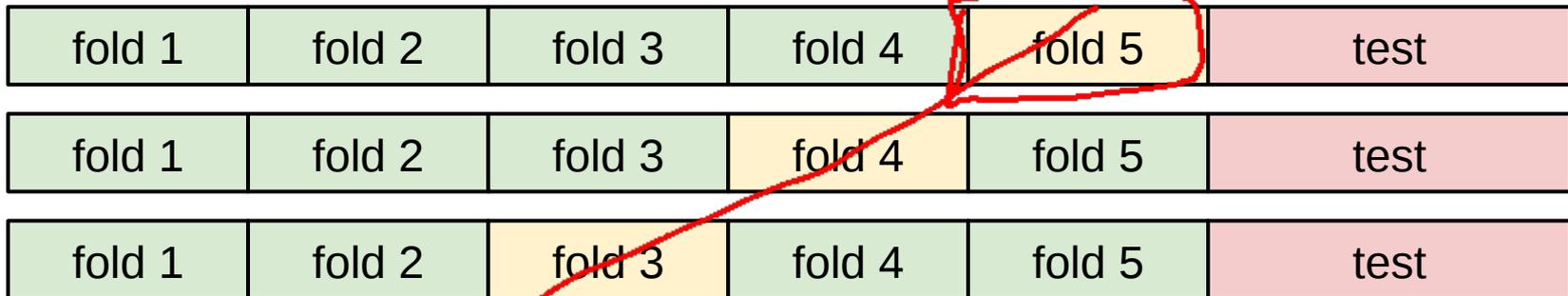
Four things make a memory based learner:

- *A distance metric*
- *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



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

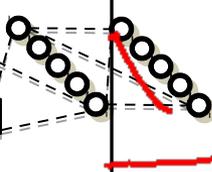
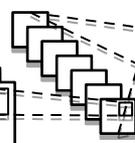
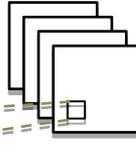
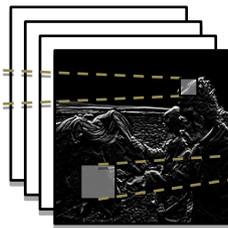
Linear  
classifiers



This image is [CC0.1.0](https://creativecommons.org/licenses/by/4.0/) public domain

# Visual Question Answering

## Image Embedding (VGGNet)



4096-dim

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

I

Convolution Layer + Non-Linearity

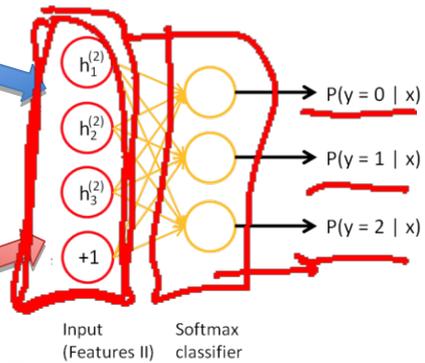
Pooling Layer

Convolution Layer + Non-Linearity

Pooling Layer

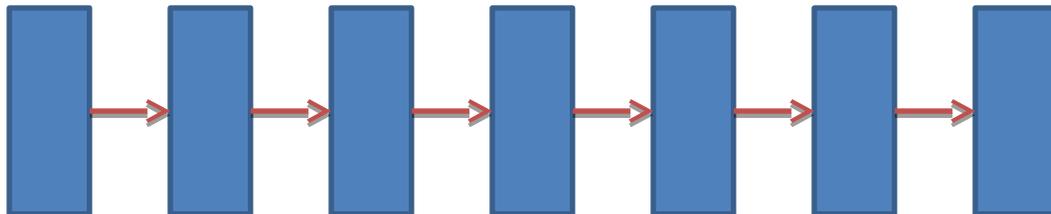
Fully-Connected MLP

Neural Network  
Softmax  
over top K answers



## Question Embedding (LSTM)

"How many horses are in this image?"



$\vec{q} \in \mathbb{R}^{d_2}$

Linear

# 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

$$\mathcal{H} = \{h: x \rightarrow y\}$$

Image



10 numbers giving class scores

$$\vec{s} = \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_{10} \end{bmatrix}$$

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

**W**

parameters  
or weights

# Parametric Approach: Linear Classifier

$$\mathcal{H} = \{h: X \rightarrow Y\}$$

Image



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

$$\mathcal{X} \in \mathbb{R}^{3072}$$

$$\underline{f(x, W)} = \underline{W} \underline{x} + \underline{b}$$

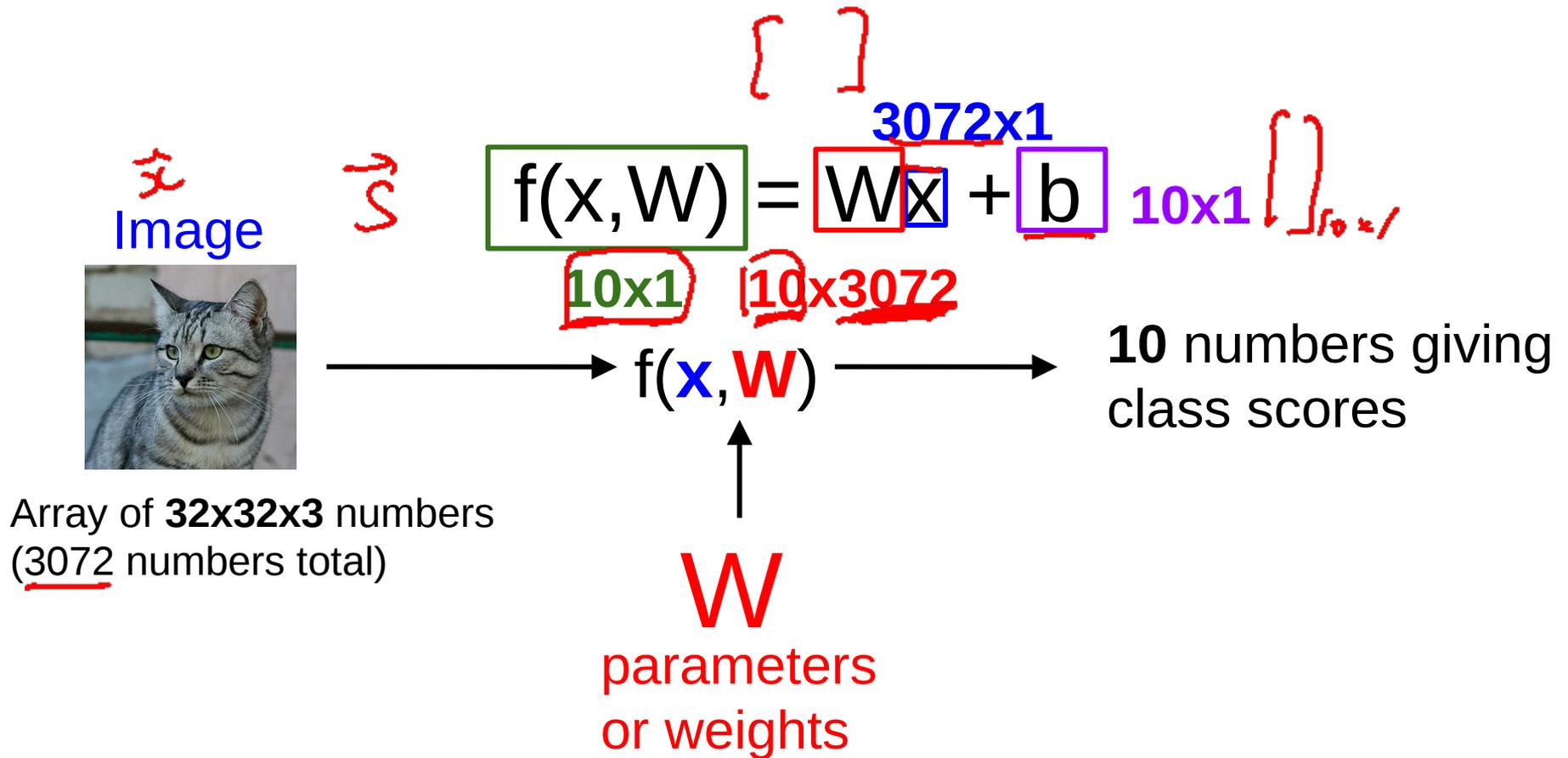


**10** numbers giving  
class scores

**W**

parameters  
or weights

# Parametric Approach: Linear Classifier

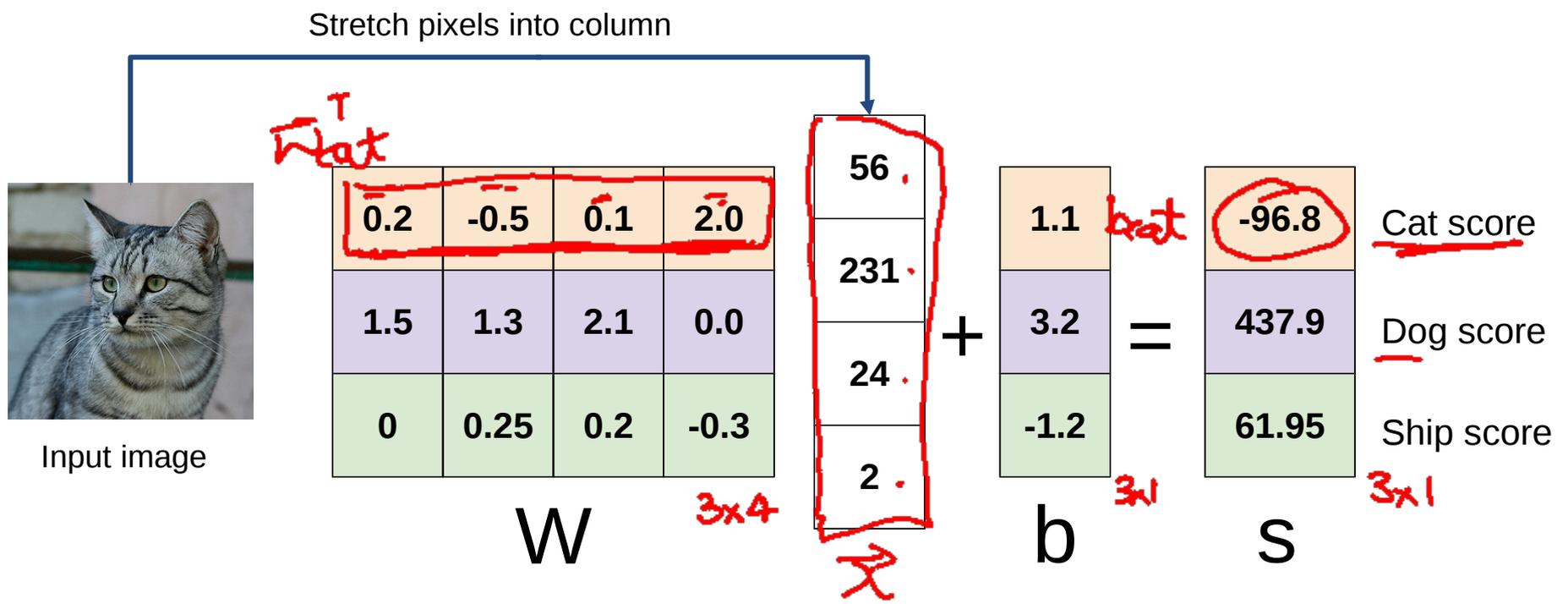


# Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Stretch pixels into column



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



$$\langle W_{cat}, \vec{x} \rangle + b_{cat} = S_{cat}$$

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

$b$

↔

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

$W$        $b$

new, single  $W$

56
231
24
2
1

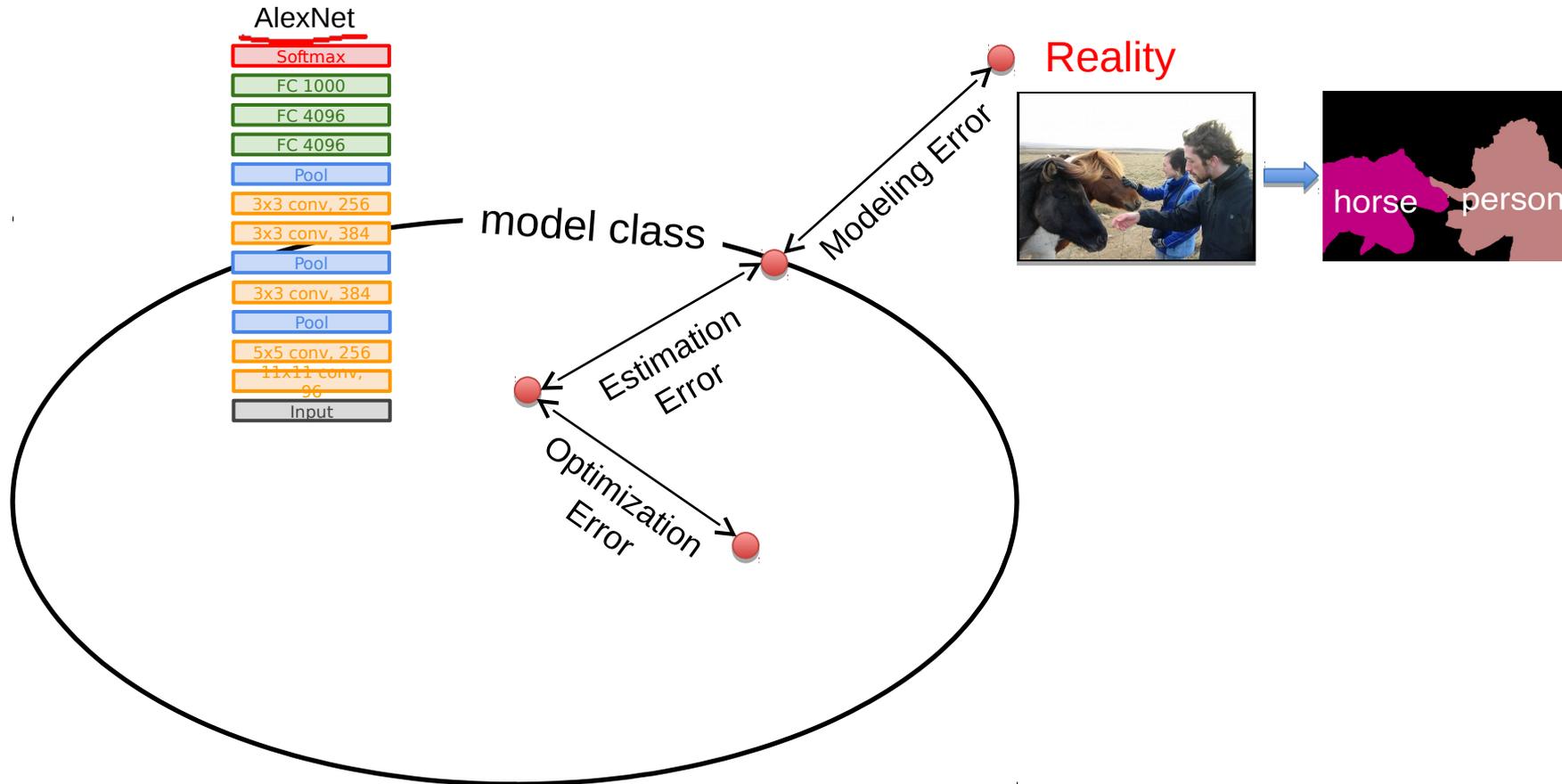
$x_i$

5<sup>th</sup> col

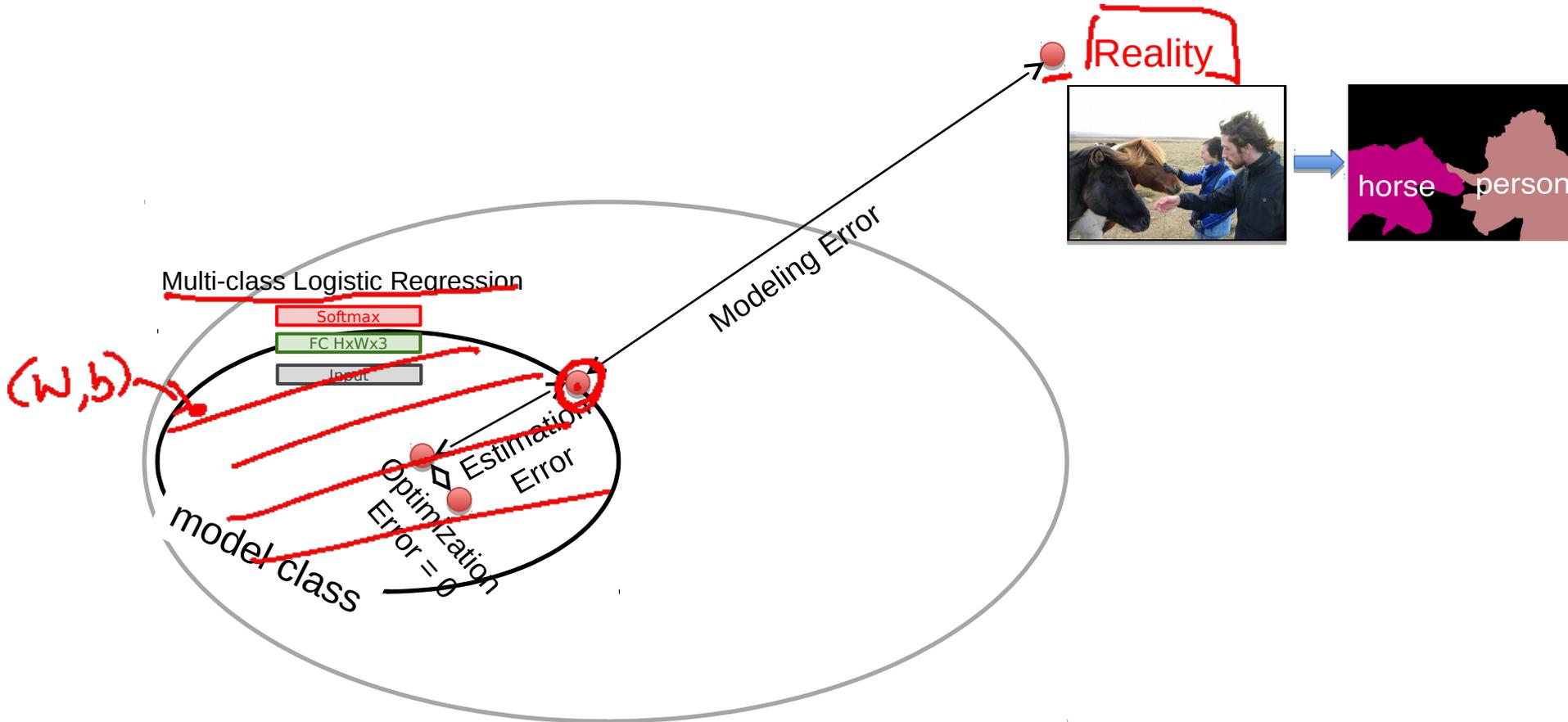


$$\vec{w}^T x + b = \begin{bmatrix} \vec{w}^T & b \end{bmatrix} \begin{bmatrix} \vec{x} \\ 1 \end{bmatrix}$$

# Error Decomposition



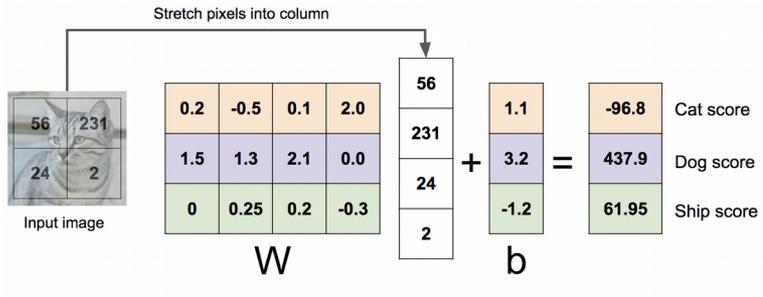
# Error Decomposition



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

Algebraic Viewpoint

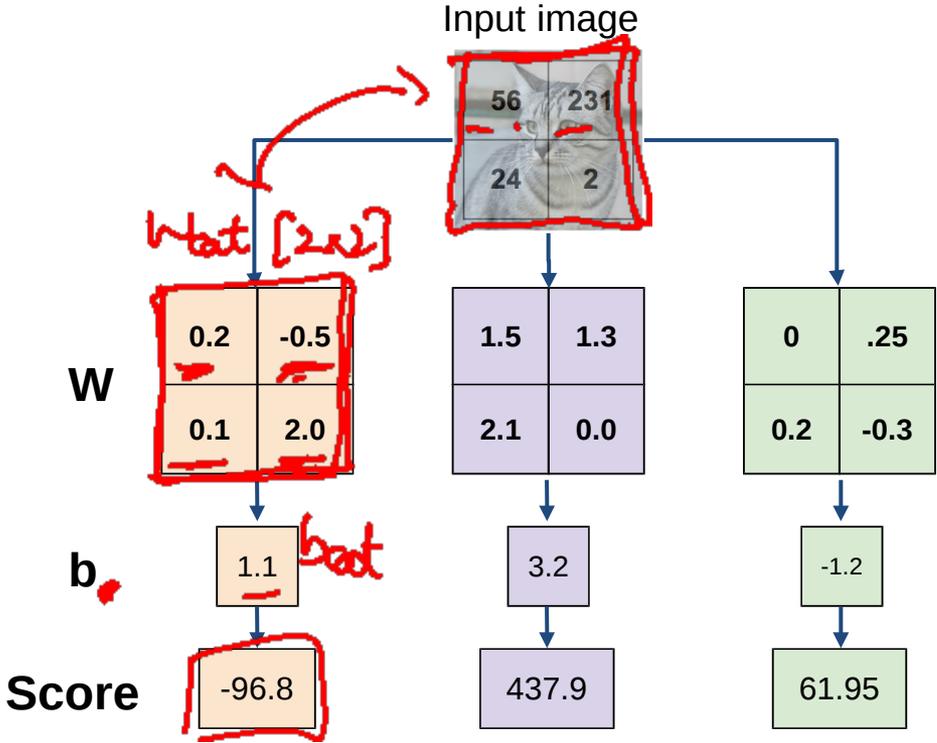
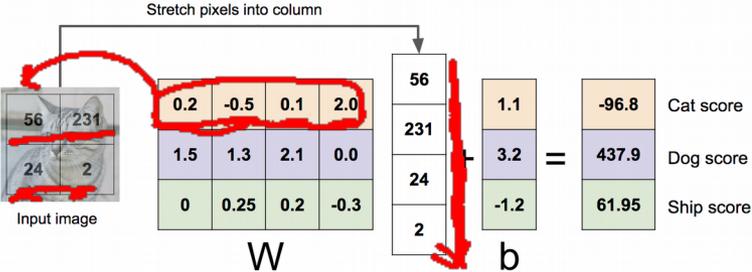
$$f(x, W) = Wx + b$$



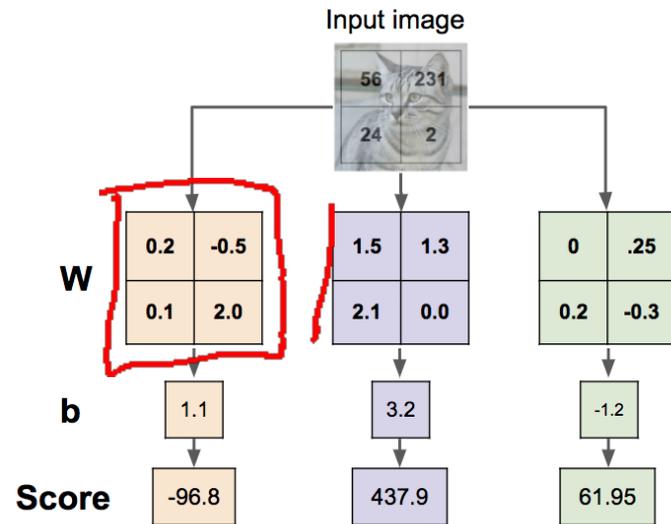
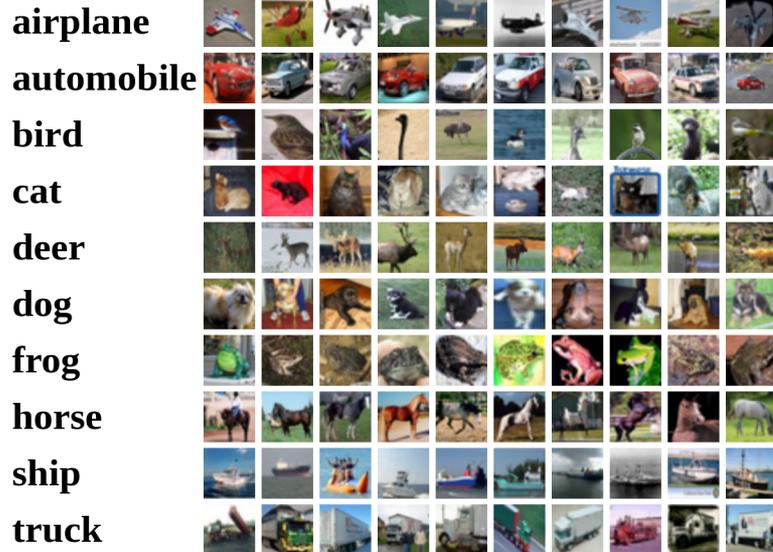
# Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

## Algebraic Viewpoint

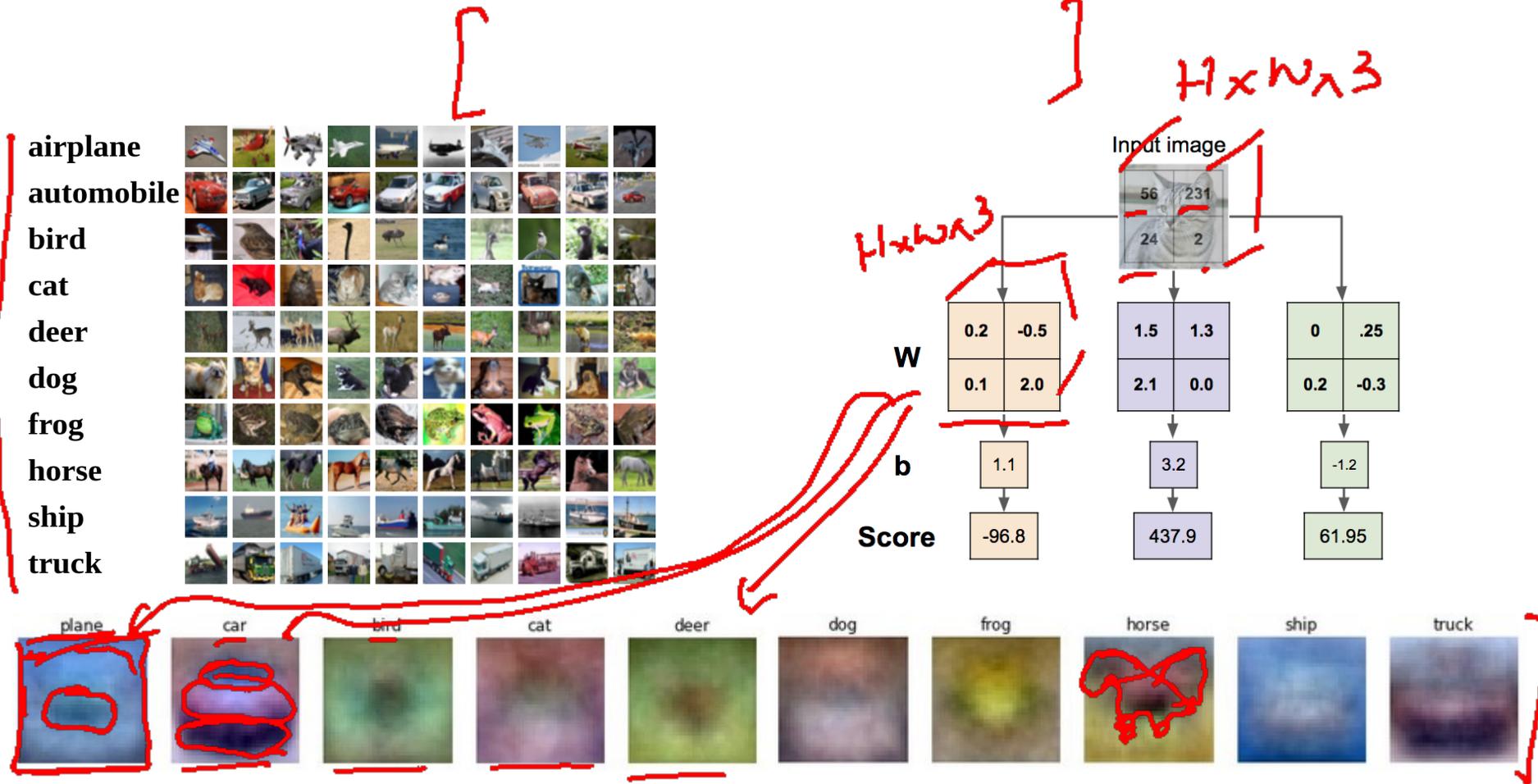
$$f(x,W) = Wx + b$$



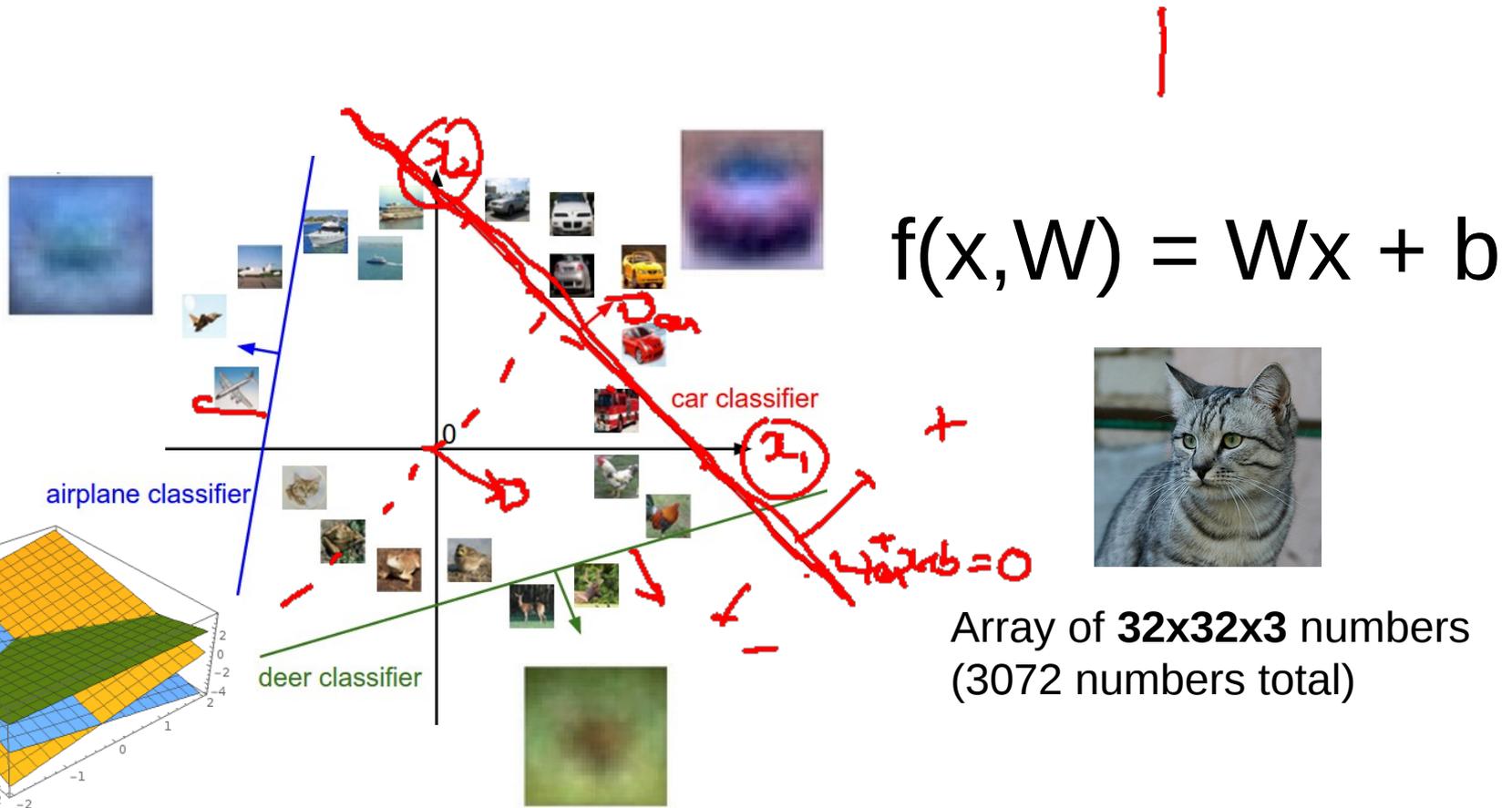
# Interpreting a Linear Classifier



# Interpreting a Linear Classifier: Visual Viewpoint



# Interpreting a Linear Classifier: Geometric Viewpoint



Plot created using [Wolfram Cloud](https://www.wolframcloud.com/)

Cat image by [Nikita](#) is licensed under [CC-BY 2.0](#)

$$\underline{w_{car}^T} \vec{x} + \underline{b}$$

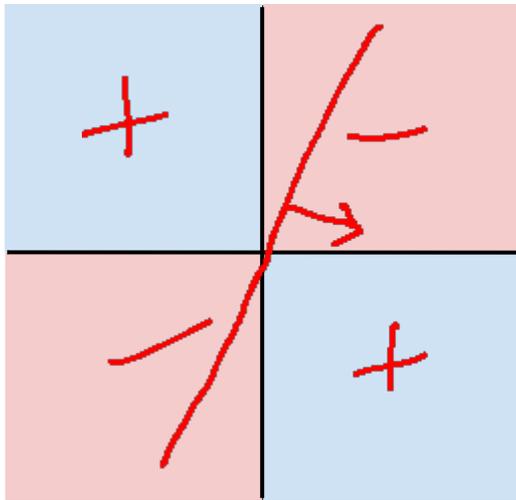
$$\left[ \leftarrow w_{car}^T \rightarrow \right]$$

$w$

# Hard cases for a linear classifier

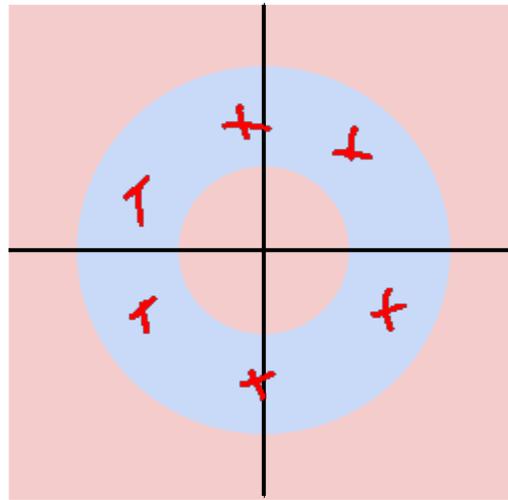
**Class 1:**  
First and third quadrants

**Class 2:**  
Second and fourth quadrants



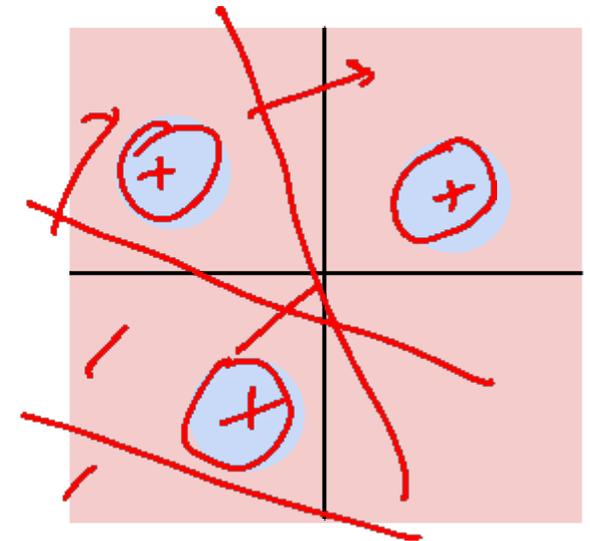
**Class 1:**  
 $1 \leq \text{L2 norm} \leq 2$

**Class 2:**  
Everything else



**Class 1:**  
Three modes

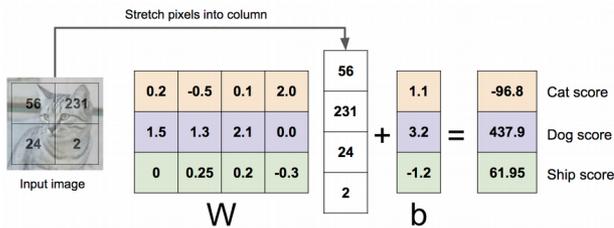
**Class 2:**  
Everything else



# Linear Classifier: Three Viewpoints

## Algebraic Viewpoint

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



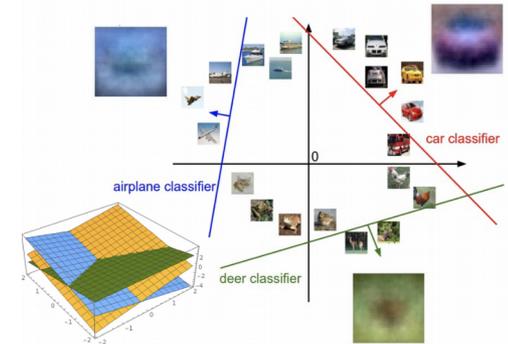
## Visual Viewpoint

One template  
per class



## Geometric Viewpoint

Hyperplanes  
cutting up space



# So far: Defined a (linear) score function

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



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

Example class scores for 3 images for some W:

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

Cat image by Nikita is licensed under CC-BY 2.0. Car image is CC0 1.0 public domain; Frog image is in the public domain

# So far: Defined a (linear) score function



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

TODO:  $L(w, D)$

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)

Cat image by Nikita is licensed under CC-BY 2.0. Car image is CC0 1.0 public domain; Frog image is in the public domain

# 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
  - $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

- Model / Hypothesis Class
  - $\{h: X \rightarrow Y\}$
  - 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.

min Loss

max Reward

# Loss Functions

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



<u>cat</u>	<b>3.2</b>	1.3	2.2
<u>car</u>	5.1	<b>4.9</b>	2.5
<u>frog</u>	-1.7	2.0	<b>-3.1</b>

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



cat	<b>3.2</b>	1.3	2.2
car	<b>5.1</b>	<b>4.9</b>	2.5
frog	<b>-1.7</b>	2.0	<b>-3.1</b>

*(Handwritten red annotations: a large bracket around the first column, a box around 5.1, and a vector symbol  $\vec{s}_i$  under -1.7)*

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:

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

*(Handwritten red annotations: a box around the fraction 1/N, a bracket under the sum term, and a vector symbol  $\vec{s}_i$  under the sum index)*

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

## 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:

$i=1$



### 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:  $s = f(x_i, W)$

<u>cat</u>	$j=1$	<u>3.2</u>
car	$j=2$	<u>5.1</u>
frog	$j=3$	<u>-1.7</u>

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 \\ \bar{s}_j - s_{y_i} + 1 \end{cases}$$

score of correct class

if  $\bar{s}_{y_i} \geq \bar{s}_j + 1$

otherwise

unhappy

$$\max(0, x) = \begin{cases} 0 & x < 0 \\ x & x \geq 0 \end{cases}$$

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**

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

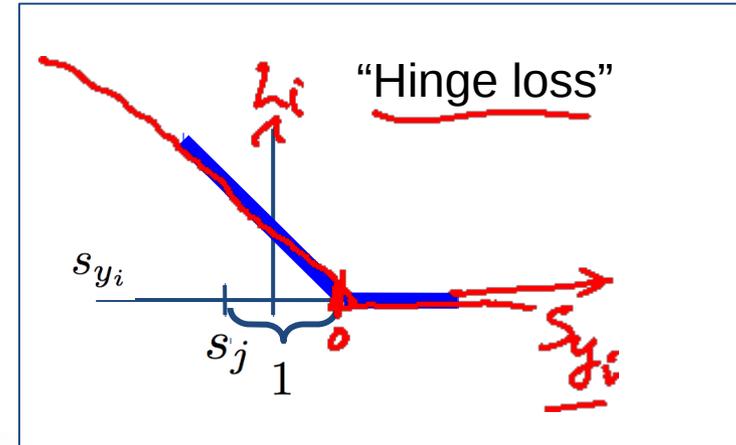
$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

### Multiclass SVM loss:



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

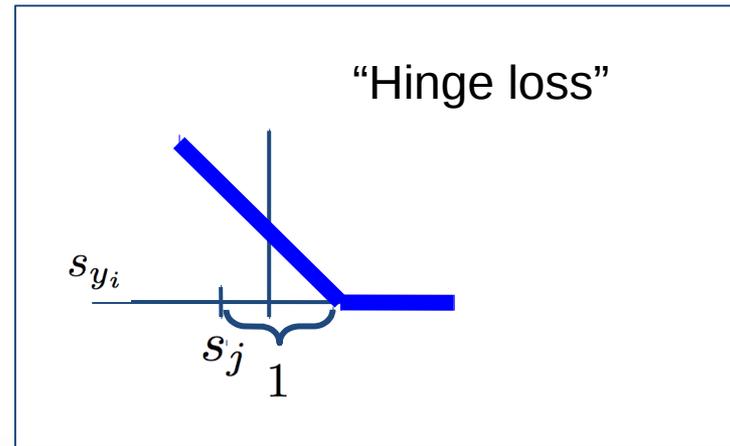
$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

### Multiclass SVM loss:



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

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



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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

## 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:

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

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



cat	<b>3.2</b>	1.3	2.2
car	<b>5.1</b>	<b>4.9</b>	2.5
frog	<b>-1.7</b>	2.0	<b>-3.1</b>
Losses:	<b>2.9</b>		

$L_i$

### 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:

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

$$= \max(0, 5.1 - 3.2 + 1)$$

$$+ \max(0, -1.7 - 3.2 + 1)$$

$$= \max(0, 2.9) + \max(0, -3.9)$$

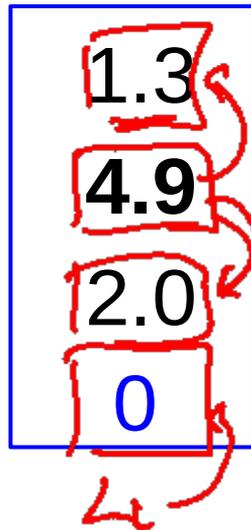
$$= 2.9 + 0$$

$$= 2.9$$

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



cat	<b>3.2</b>	<b>1.3</b>	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	<b>2.0</b>	<b>-3.1</b>
Losses:	2.9	<b>0</b>	



### 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:

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_j - \underline{s_{y_i}} + 1) \\
 &= \max(0, \underline{1.3} - \underline{4.9} + \underline{1}) \\
 &\quad + \max(0, \underline{2.0} - \underline{4.9} + \underline{1}) \\
 &= \max(0, \underline{-2.6}) + \max(0, \underline{-1.9}) \\
 &= \underline{0} + \underline{0} \\
 &= \underline{0}
 \end{aligned}$$

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	2.9	0	<b>12.9</b>

### 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:

$$\begin{aligned}
 L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\
 &= \max(0, 2.2 - (-3.1) + 1) \\
 &\quad + \max(0, 2.5 - (-3.1) + 1) \\
 &= \max(0, 6.3) + \max(0, 6.6) \\
 &= 6.3 + 6.6 \\
 &= 12.9
 \end{aligned}$$

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	<b>2.9</b> $L_1$	<b>0</b> $L_2$	<b>12.9</b> $L_3$

### 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:

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

$$L = \frac{(2.9 + 0 + 12.9)}{3} = 5.27$$

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



cat	<b>3.2</b>	1.3 <sup>no</sup> <sub><math>\epsilon_1</math></sub>	2.2
car	5.1	<b>4.9</b> <sub><math>\epsilon_2</math></sub>	2.5
frog	-1.7	2.0 <sub><math>\epsilon_3</math></sub>	<b>-3.1</b>
Losses:	<b>2.9</b>	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:

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

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	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
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:



cat	<b>3.2</b>	<i>no</i>	1.3	2.2
car	5.1	<i>no</i>	<b>4.9</b>	2.5
frog	-1.7	<i>no</i>	2.0	<b>-3.1</b>
Losses:	<b>2.9</b>		<b>0</b>	<b>12.9</b>

### 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:

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

Q3: At initialization  $W$   
 is small so all  $s \approx 0$ .

What is the loss?

#class - 1

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
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  $j = y_i$ )  $+1$

$$L_i \leftarrow L_i + 1$$

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
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:

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

Q5: What if we used  
 mean instead of  
 sum?

$$f(x, W) = \underline{W}x$$

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

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

Q7: Is this W unique?

$$L(W) = 0$$

$$\Rightarrow L(2W) = 0$$

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

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

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:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	2.9	0	

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

**Before:**

$$\begin{aligned}
 &= \max(0, \underline{1.3} - \underline{4.9} + 1) \\
 &\quad + \max(0, 2.0 - 4.9 + 1) \\
 &= \max(0, -2.6) + \max(0, -1.9) \\
 &= 0 + 0 \\
 &= 0
 \end{aligned}$$

**With  $W$  twice as large:**

$$\begin{aligned}
 &= \max(0, \underline{2.6} - \underline{9.8} + 1) \\
 &\quad + \max(0, 4.0 - 9.8 + 1) \\
 &= \max(0, -6.2) + \max(0, -4.8) \\
 &= 0 + 0 \\
 &= 0
 \end{aligned}$$

# Multiclass SVM Loss: Example code

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

```
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**

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

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

Softmax function

cat	<b>3.2</b>
car	5.1
frog	-1.7

# Softmax Classifier (Multinomial Logistic Regression)



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

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

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

Probabilities  
must be  $\geq 0$

cat

$e^{3.2}$

car

$e^{5.1}$

frog

$e^{-1.7}$

exp

24.5

164.0

0.18

unnormalized  
probabilities

# Softmax Classifier (Multinomial Logistic Regression)



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

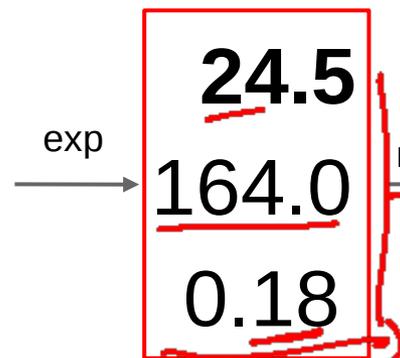
$$s = f(x_i; W)$$

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

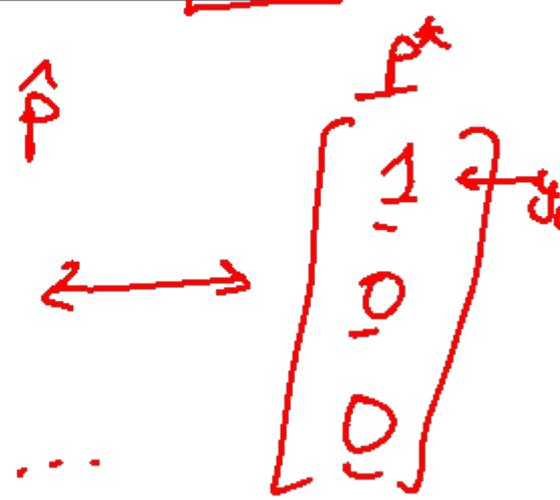
Probabilities must be  $\geq 0$

Probabilities must sum to 1

cat      3.2  
car      5.1  
frog      -1.7



normalize →



# Softmax Classifier (Multinomial Logistic Regression)



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

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

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

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

cat

3.2

car

5.1

frog

-1.7

exp

24.5

normalize

164.0

0.18

0.13

0.87

0.00

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**

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

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

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

cat  
car  
frog

3.2
5.1
-1.7

exp

24.5
164.0
0.18

normalize

0.13
0.87
0.00

Unnormalized log-  
probabilities / logits

unnormalized  
probabilities

probabilities

# Softmax Classifier (Multinomial Logistic Regression)

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

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

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



*Handwritten: cat*

Probabilities must be  $\geq 0$

Probabilities must sum to 1

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

cat

3.2

car

5.1

frog

-1.7

Unnormalized log-probabilities / logits

exp

24.5

164.0

0.18

unnormalized probabilities

normalize

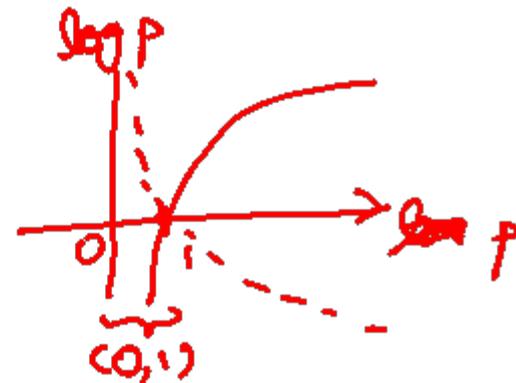
0.13

0.87

0.00

probabilities

→  $L_i = -\log(0.13) = 2.04$



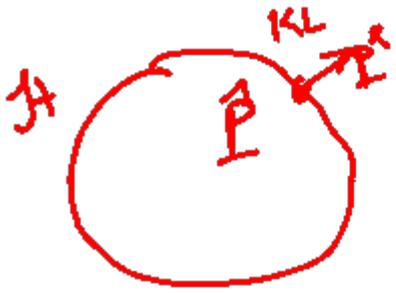
# Log-Likelihood / KL-Divergence / Cross-Entropy

$$p^* = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$y_i$  →

$$\hat{p} = \begin{bmatrix} P_a(Y=1 | x_i, w) \\ P_a(Y=2 | x_i, w) \\ \vdots \\ P_a(Y=k | x_i, w) \end{bmatrix}$$

$$\min_w \text{KL}(p^* || \hat{p}) = \sum_y p^*(y) \log \frac{p^*(y)}{\hat{p}(y)}$$



$$= \boxed{\sum_y p^*(y) \log p^*(y)} - \boxed{\sum_y p^*(y) \log \hat{p}(y)}$$

$$= \boxed{-H(p^*)}$$

$$\min_w \boxed{H(p^*, \hat{p})}$$

↓

$$\min_w \underline{-\log \hat{p}(y_i)}$$

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

# Softmax Classifier (Multinomial Logistic Regression)



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

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

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

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	<b>3.2</b>
car	5.1
frog	-1.7

# Softmax Classifier (Multinomial Logistic Regression)



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

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

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

Maximize probability of correct class

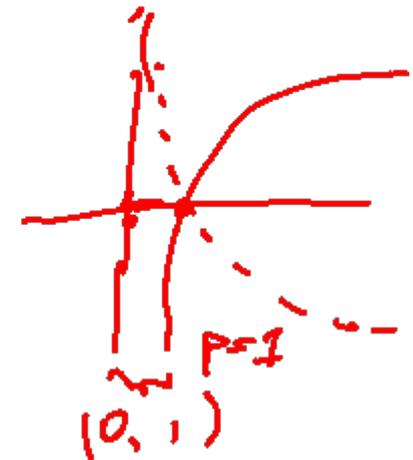
$$L_i = -\log P(Y = y_i | X = x_i)$$

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**

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

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

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	<b>3.2</b>
car	<b>5.1</b>
frog	<b>-1.7</b>

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

A: ~~min~~ 0, max infinity  
*inf*

# Softmax Classifier (Multinomial Logistic Regression)



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

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

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

Maximize probability of correct class

Putting it all together:

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

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

*W=0*

cat

3.2 *20*

car

5.1 *20*

frog

-1.7 *20*

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

*p*

$$\begin{bmatrix} 1/K \\ 1/K \\ 1/K \end{bmatrix}$$

*$-\log(1/K)$   
 $\log(K)$*

# Softmax Classifier (Multinomial Logistic Regression)



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

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

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

Maximize probability of correct class

Putting it all together:

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

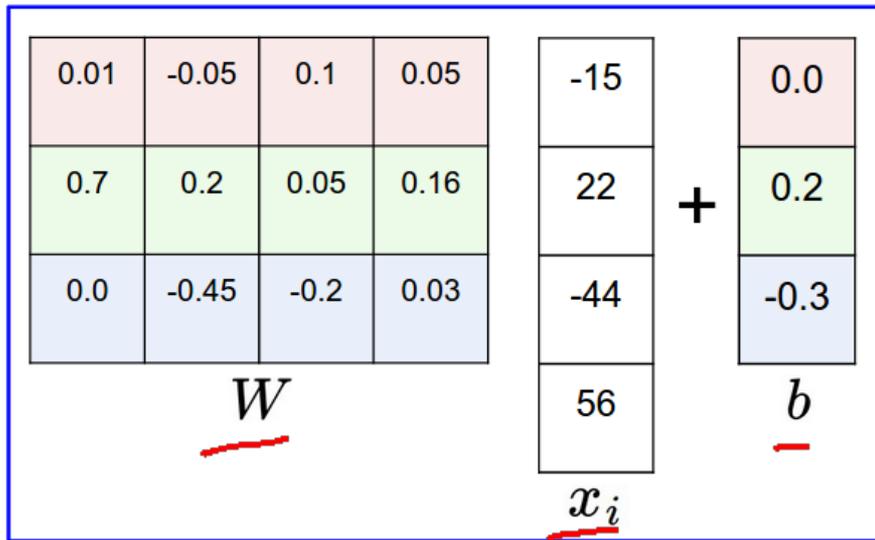
$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

cat	<b>3.2</b>
car	5.1
frog	-1.7

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

# Softmax vs. SVM

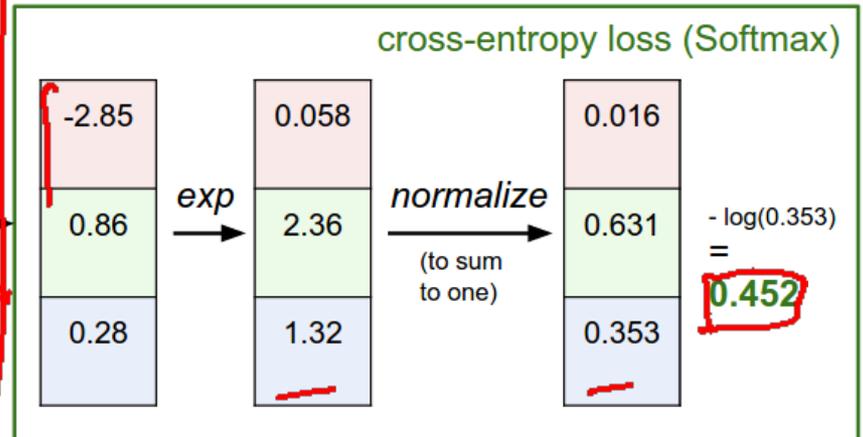
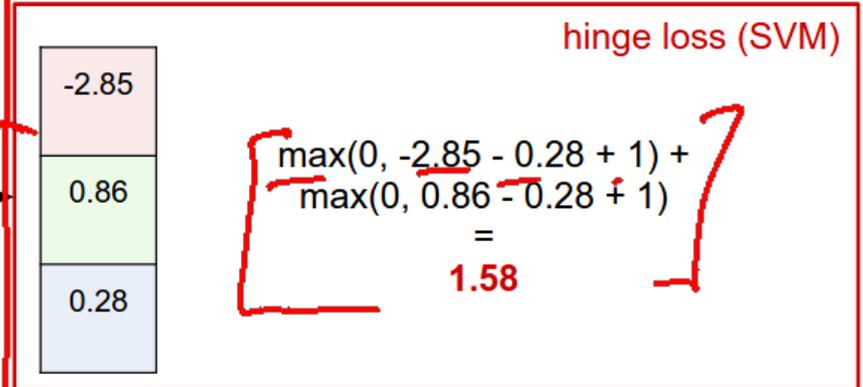
matrix multiply + bias offset



y<sub>i</sub> 2

Model  
class

Model → Loss



## Softmax vs. SVM

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

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

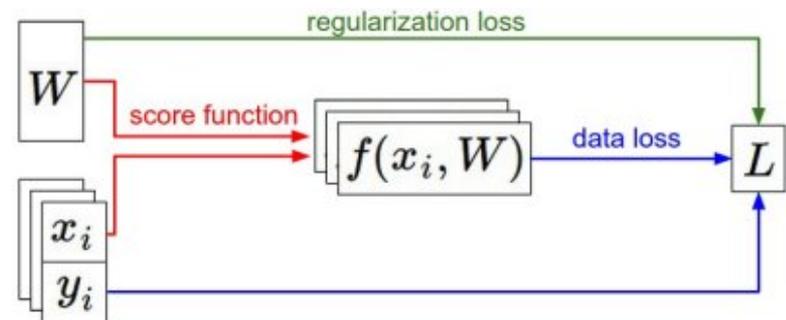
# Recap

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

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

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

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



# Recap

How do we find the best  $W$ ?

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

$W^*$   
min/val

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

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

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

