# CS 4803 / 7643: Deep Learning

Topics:
- Linear Classifiers
- Loss Functions

Dhruv Batra

Georgia Tech

# Administrativia

- ## Notes on class webpage

  - https://www.cc.gatech.edu/classes/AY2019/cs7643_fall/


- ## HW0 Reminder

  - Due: 09/05

# Recap from last time

# **Image Classification**: A core task in Computer Vision
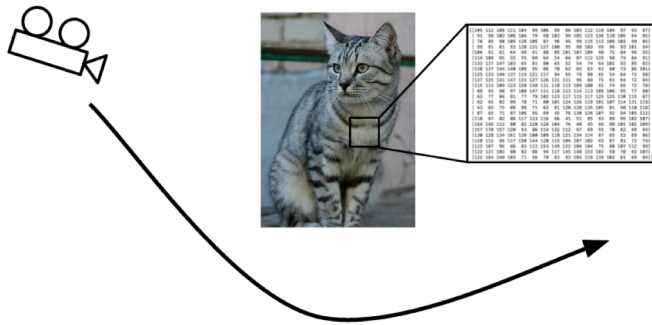
(assume given set of discrete labels)
{dog, cat, truck, plane, ...}

⟶ cat

# Challenges of recognition

Viewpoint



Illumination

Deformation

Occlusion

Clutter

Intraclass Variation

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

# An image classifier

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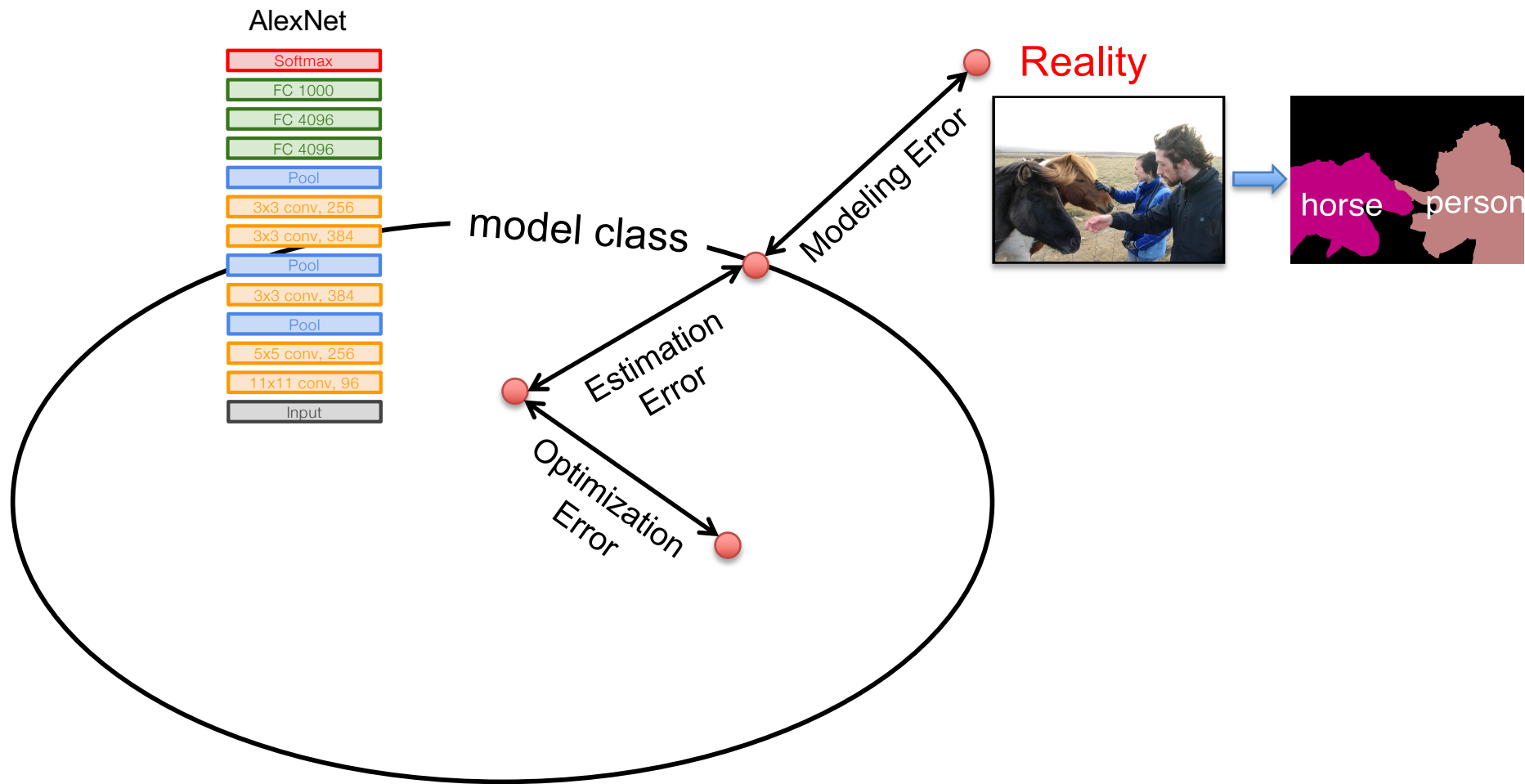
Unlike e.g. sorting a list of numbers,

**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
  - $\{h: X \rightarrow Y\}$
  - 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

# Error Decomposition

- Approximation/Modeling Error
  - You approximated reality with model

- Estimation Error
  - You tried to learn model with finite data

- Optimization Error
  - You were lazy and couldn't/didn't optimize to completion

- Bayes Error
  - Reality just sucks

# 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?*

# Parametric vs Non-Parametric Models

- Does the capacity (size of hypothesis class) grow with size of training data?
  - Yes = Non-Parametric Models
  - No = Parametric Models

# 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)

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

- Very slow at test time
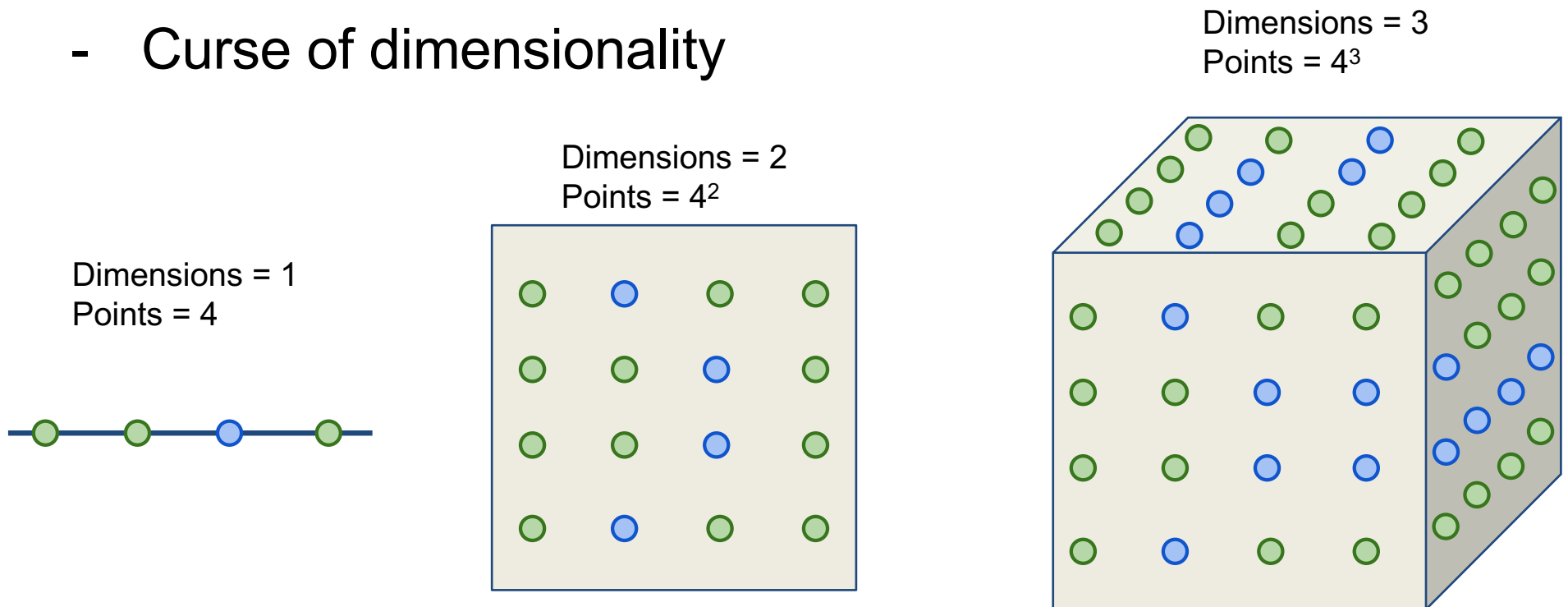- Distance metrics on pixels are not informative

| Original | Boxed | Shifted | Tinted |
| --- | --- | --- | --- |

(all 3 images have same L2 distance to the one on the left)

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

- Curse of dimensionality

Dimensions = 1
Points = 4
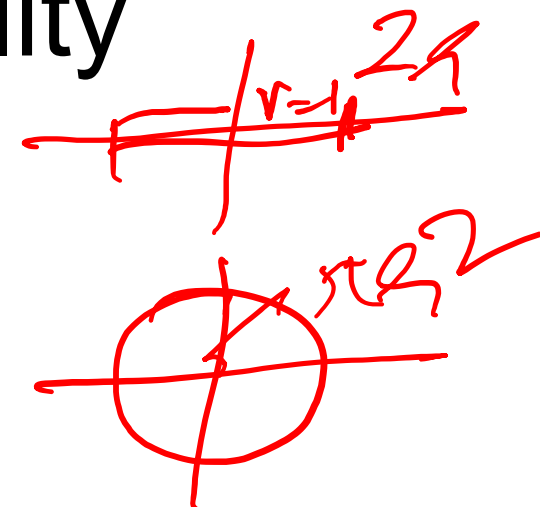
Dimensions = 2
Points = $4^2$

Dimensions = 3
Points = $4^3$

# Curse of Dimensionality

- Consider: Sphere of radius 1 in d-dims

- Consider: an outer ε-shell in this sphere

- What is $\dfrac{\textit{shell volume}}{\textit{sphere volume}}$ ?
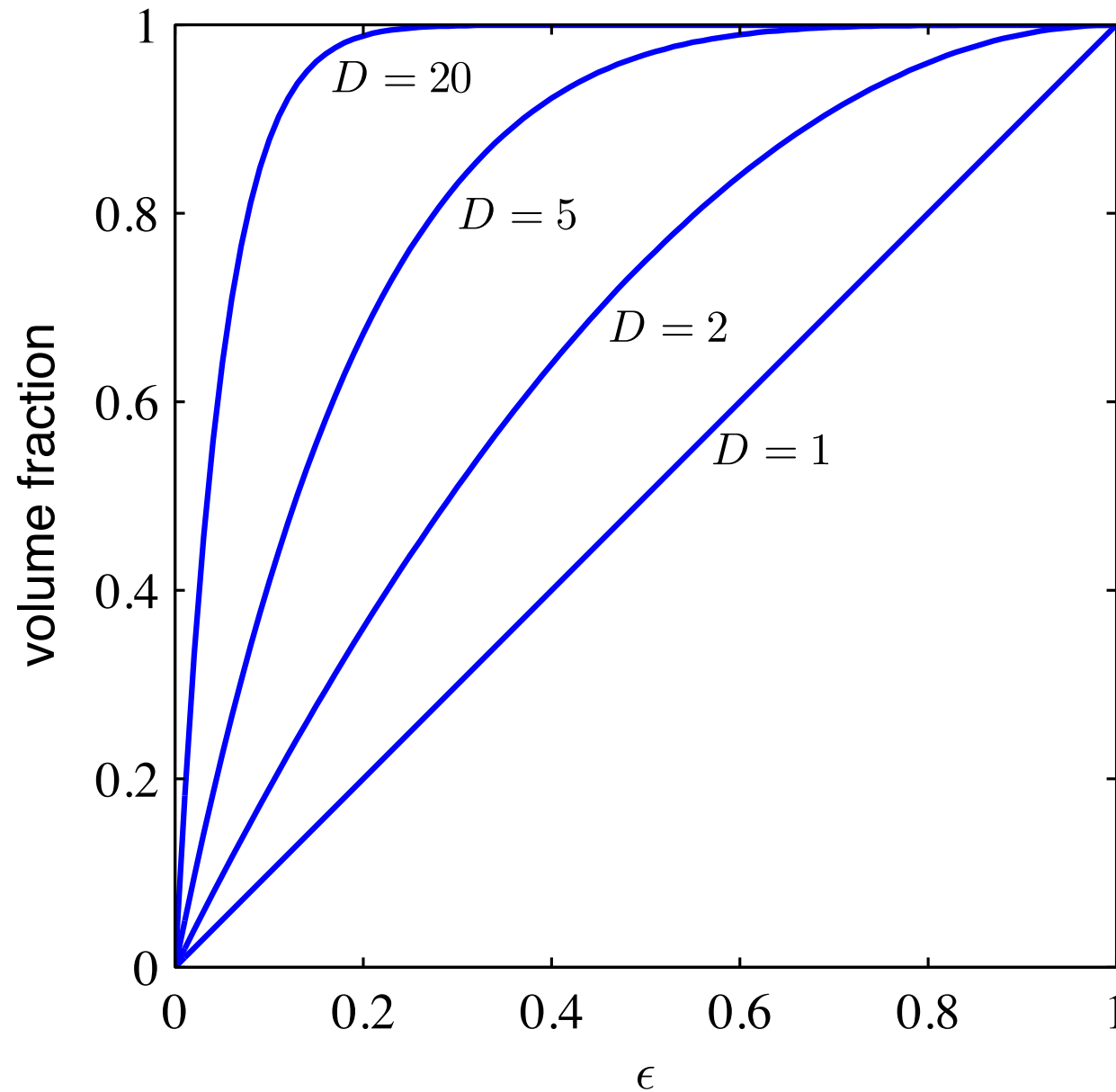
$$\frac{1^d - (1-\varepsilon)^d}{1^d}$$

$$1 - (1-\varepsilon)^d$$

# Curse of Dimensionality

# Plan for Today

- ## Linear Classifiers
  - Linear scoring functions

- ## Loss Functions
  - Multi-class hinge loss
  - Softmax cross-entropy loss
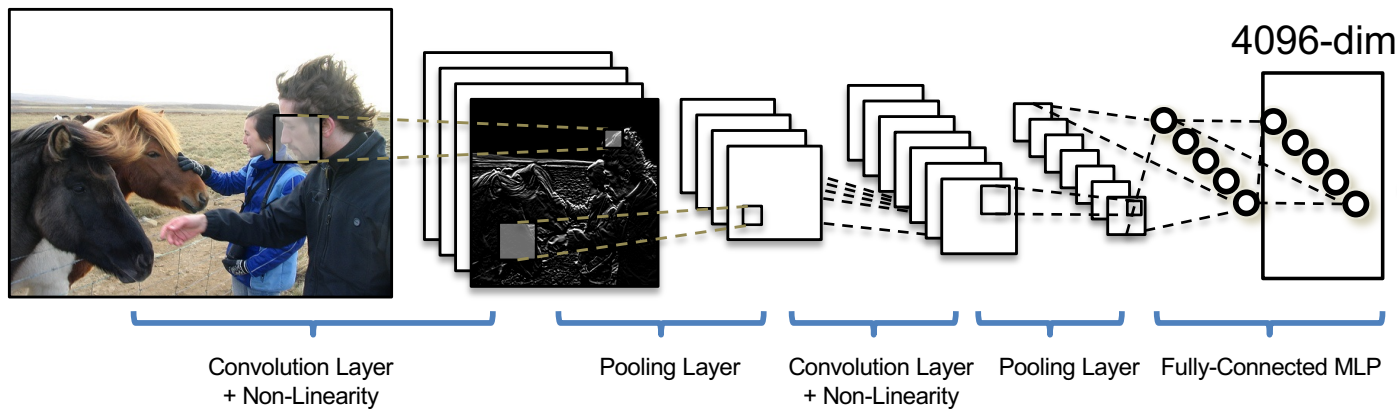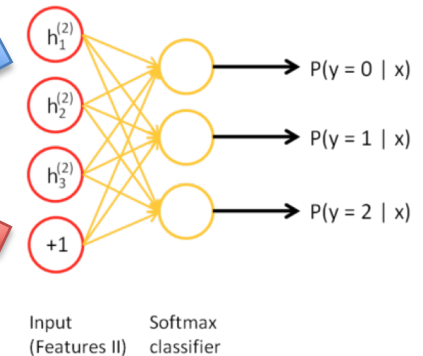
# Linear Classification

Neural Network

Linear classifiers

This image is CC0 1.0 public domain
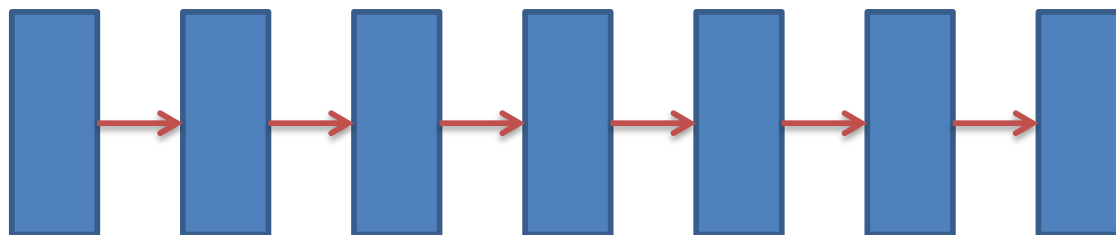
# Visual Question Answering

## Image Embedding (VGGNet)



4096-dim

Convolution Layer + Non-Linearity

Pooling Layer

Convolution Layer + Non-Linearity

Pooling Layer

Fully-Connected MLP

## Neural Network Softmax over top K answers

$P(y = 0 \mid x)$

$P(y = 1 \mid x)$

$P(y = 2 \mid x)$

Input (Features II)

Softmax classifier
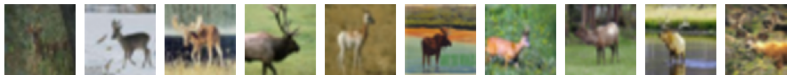
## Question Embedding (LSTM)

"How many horses are in this image?"

# Recall CIFAR10



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

**10,000** test images.

# Parametric Approach

Image



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

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

↑

**W**
parameters
or weights

# Parametric Approach: Linear Classifier

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

Image



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

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

**W**
parameters
or weights

# Parametric Approach: Linear Classifier

**Image**



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

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

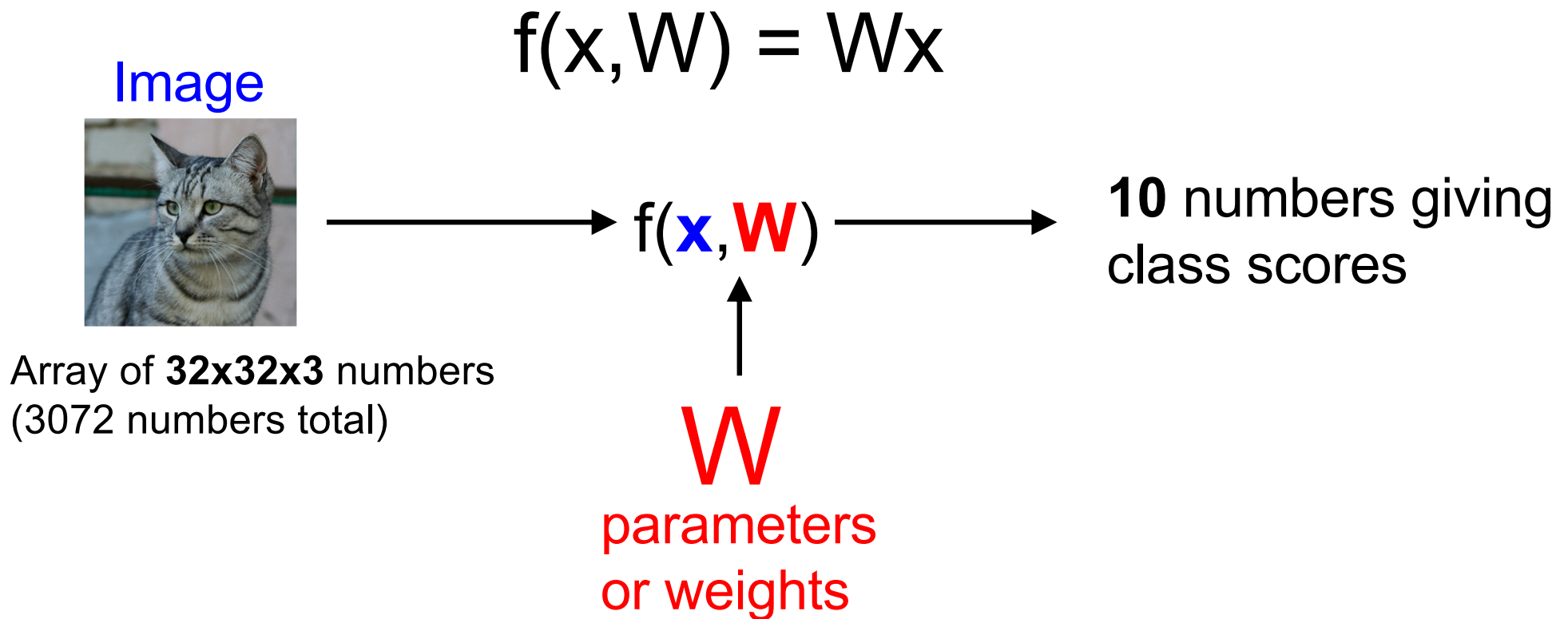**3072x1**

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

$f(\mathbf{x},\mathbf{W})$

**10** numbers giving
class scores

**W**
parameters
or weights

# Parametric Approach: Linear Classifier

**Image**



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

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

3072x1 (for $x$)

10x1 (for $f(x,W)$)  10x3072 (for $W$)

10x1 (for $b$)

$f(\mathbf{x},\mathbf{W})$

**10** numbers giving class scores

**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

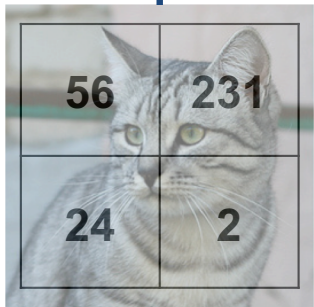# Error Decomposition

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

Stretch pixels into column

56 231
24 2

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 |

$W$

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

$+$

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

$b$

$=$

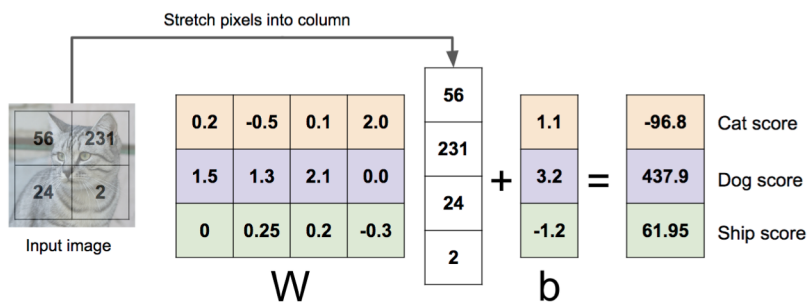| -96.8 | Cat score |
|-------|-----------|
| 437.9 | Dog score |
| 61.95 | Ship score |

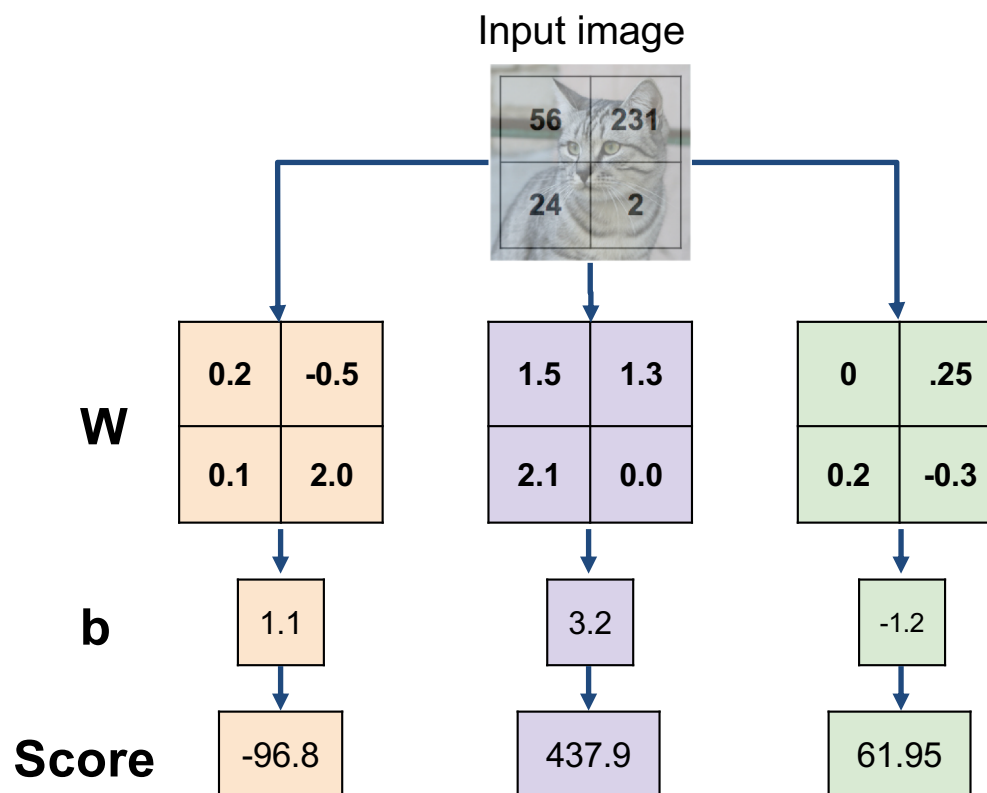# Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Algebraic Viewpoint

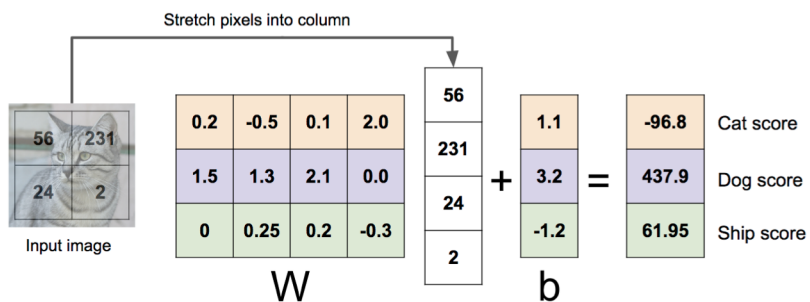$$f(x,W) = Wx$$

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

Algebraic Viewpoint

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



Input image

**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 |

# Interpreting a Linear Classifier

# Interpreting a Linear Classifier:
## Visual Viewpoint

# Interpreting a Linear Classifier: <u>Geometric Viewpoint</u>



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

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

# 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) <u>score function</u>

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



| | | | |
|---|---|---|---|
| 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 |

Example class scores for 3 images for some W:

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

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

# **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 $\rightarrow$ Y                       (the "true" mapping / reality)

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

- Model / Hypothesis Class
  - {h: X $\rightarrow$ Y}
  - e.g. y = h(x) = 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 | **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 | **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:



**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:



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

**Multiclass SVM loss:**

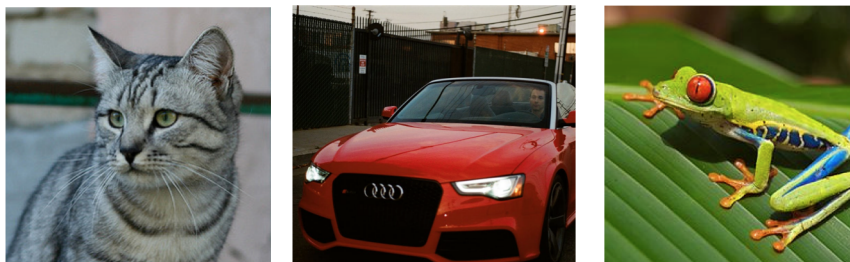Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

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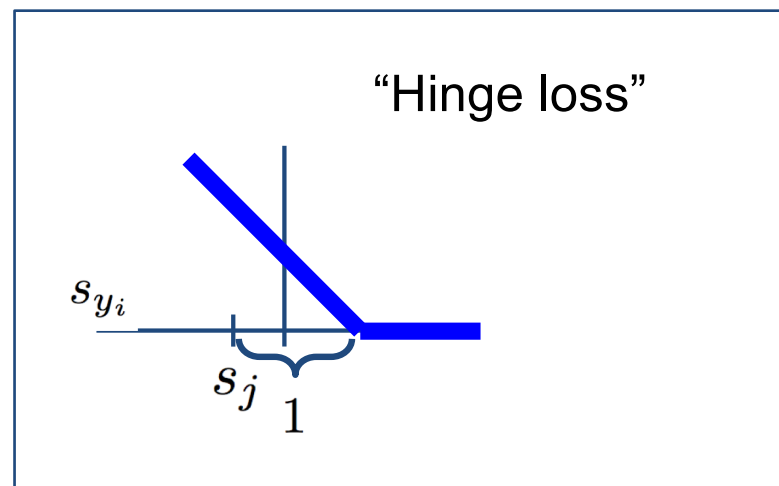
the SVM loss has the form:

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

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



|  | cat | car | frog |
|---|---|---|---|
| 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:

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



|        | 3.2  | 1.3  | 2.2  |
|--------|------|------|------|
| 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    |      |

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:
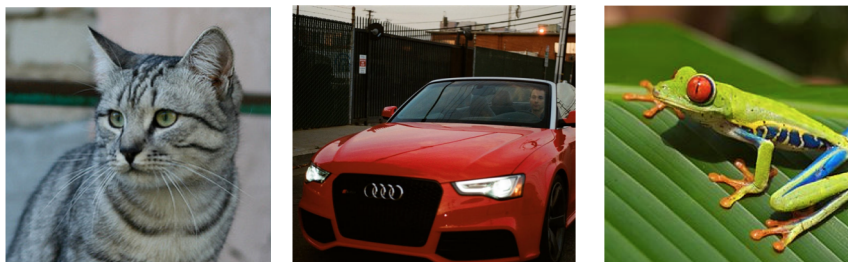
**Multiclass SVM loss:**

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:
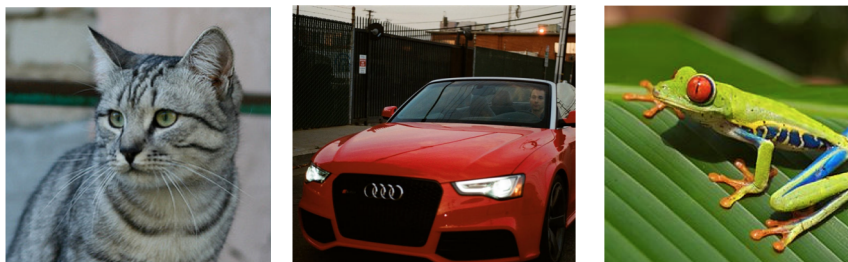
**Multiclass SVM loss:**

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:

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

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:



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

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

the SVM loss has the form:

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

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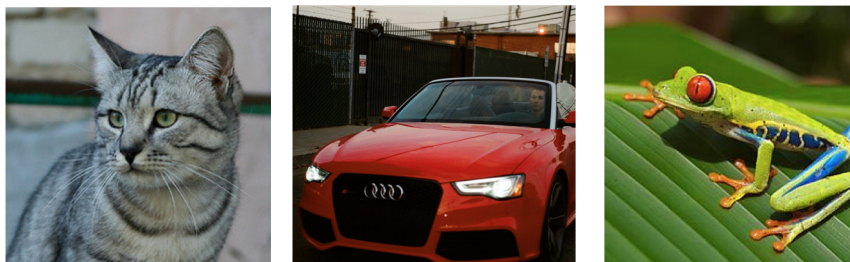
$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q2: what is the
min/max possible
loss?

$0, \infty$

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: $s = f(x_i, W)$
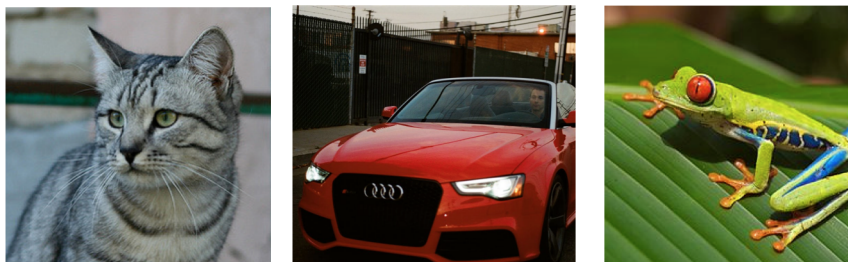
the SVM loss has the form:

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

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

(#classes − 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 |

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    | 0      | 12.9   |

**Multiclass SVM loss:**

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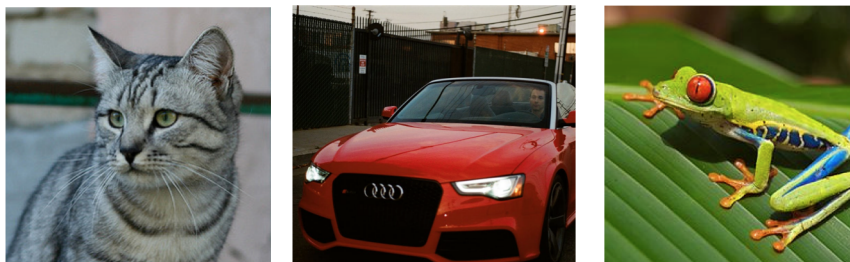where $y_i$ is the (integer) label,

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)

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 |

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
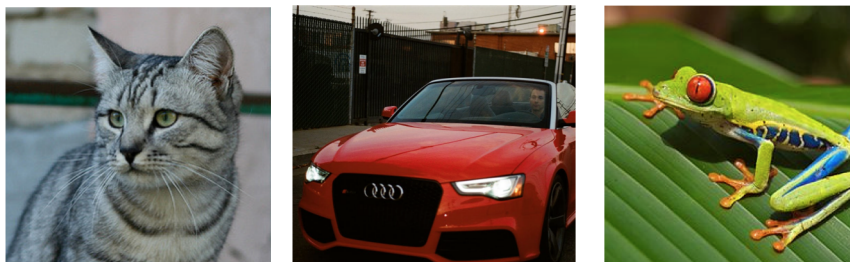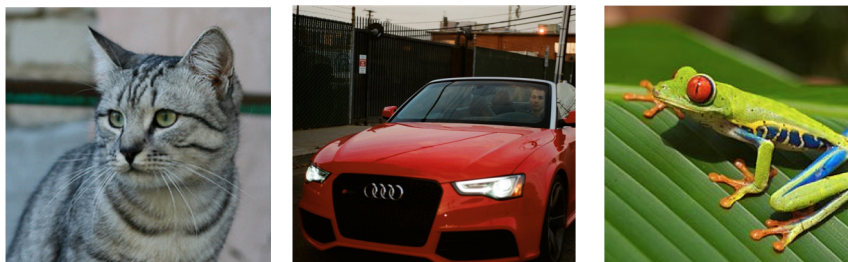scores vector: $s = f(x_i, W)$



|       |        |       |       |
|-------|--------|-------|-------|
| 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:

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

Q6: What if we used

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

# 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
```

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

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

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

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

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

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

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



|  | 🐱 | 🚗 | 🐸 |
|---|---|---|---|
| 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 | |

**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

# **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

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

exp

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

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

exp →

|        |
|--------|
| **24.5** |
| 164.0  |
| 0.18   |

unnormalized probabilities

normalize →

|        |
|--------|
| **0.13** |
| 0.87   |
| 0.00   |

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

| | Unnormalized log-probabilities / logits | exp | unnormalized probabilities | normalize | probabilities |
|---|---|---|---|---|---|
| cat | **3.2** | | **24.5** | | **0.13** |
| car | 5.1 | | 164.0 | | 0.87 |
| frog | -1.7 | | 0.18 | | 0.00 |

# 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

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

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

$L_i = -\log(0.13)$
$= 0.89$

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

$-\log(0)$

$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

$-\log(1)$
$= 0$

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

# Softmax Classifier (Multinomial Logistic Regression)

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

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

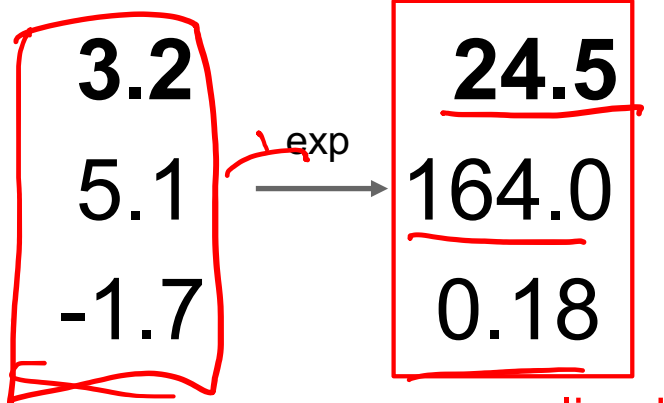$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$ Softmax Function

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

Probabilities must be >= 0

Probabilities must sum to 1

| | cat | car | frog |
|---|---|---|---|
| | 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**

**Maximum Likelihood Estimation**
Choose probabilities to maximize the likelihood of the observed data

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

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

$$D = \{(x_i, y_i)\} \qquad \underline{IID} \sim P^x$$

$$\boxed{\hat{W}_{MLE}} = \max \; \underline{P(\underline{D} \mid w)}$$

$$\equiv \max \; \log P(\underline{D} \mid w)$$

$$= \max \; \sum_i \log P(y_i \mid x_i, \underline{w})$$

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

$$y_i$$

$$
\begin{bmatrix} 0 \\ 0 \\ 0 \\ \boxed{1} \\ 0 \\ 0 \end{bmatrix}
$$

$$P^{gt}$$

$$
\begin{bmatrix} P(y=1|x_i,w) \\ \vdots \\ \vdots \\ P(y=k|x_i,w) \end{bmatrix}
$$

$$\hat{P}$$

$$H(P^{gt}, \hat{P})$$

$$= -\sum_y P^{gt}(y) \, \log \hat{P}(y)$$

$$\approx -\log \hat{P}(y_i | x_i, w)$$

$$D_{KL}(P^{gt}, \hat{P}) = \sum_y P^{gt}(y) \, \log \frac{P^{gt}(y)}{\hat{P}(y)}$$

$$= \sum_y P^{gt}(y) \log P^{gt}(y) - \sum_y P^{gt}(y) \log \hat{P}(y)$$

$$= -H(P^{gt}) + H(P^{gt}, \hat{P})$$

# 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

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

<span style="color:red">Probabilities must be >= 0</span>

<span style="color:green">Probabilities must sum to 1</span>

|      | logits |       | unnorm prob |       | probabilities |       | Correct probs |
|------|--------|-------|-------------|-------|---------------|-------|---------------|
| cat  | **3.2**  | exp → | **24.5**    | normalize → | **0.13**  | compare | **1.00**      |
| car  | 5.1    |       | 164.0       |       | 0.87          |       | 0.00          |
| frog | -1.7   |       | 0.18        |       | 0.00          |       | 0.00          |

<span style="color:blue">Unnormalized log-probabilities / logits</span>

<span style="color:red">unnormalized probabilities</span>

<span style="color:green">probabilities</span>

<span style="color:purple">Correct probs</span>

# 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

|       | Unnormalized log-probabilities / logits | | unnormalized probabilities | | probabilities | | Correct probs |
|-------|------|---|------|---|------|---|------|
| cat   | **3.2**  | exp | **24.5** | normalize | **0.13** | compare | **1.00** |
| car   | 5.1  |   | 164.0 |   | 0.87 |   | 0.00 |
| frog  | -1.7 |   | 0.18 |   | 0.00 |   | 0.00 |

*Kullback–Leibler divergence*

$$D_{KL}(P \| Q) = \sum_y P(y) \log \frac{P(y)}{Q(y)}$$

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

# Softmax Classifier (Multinomial Logistic Regression)
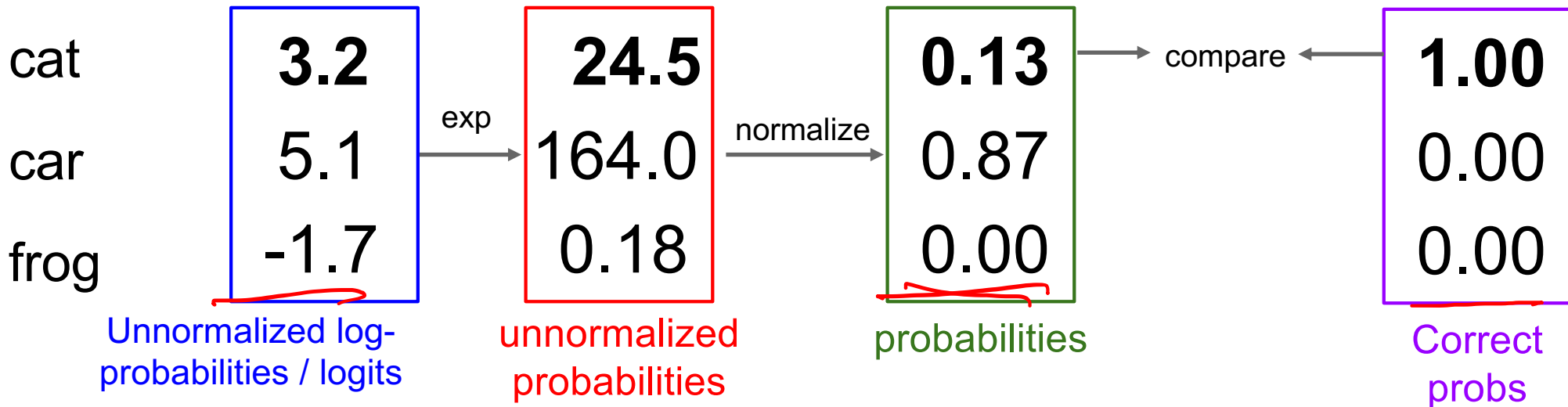
Want to interpret raw classifier scores as **probabilities**

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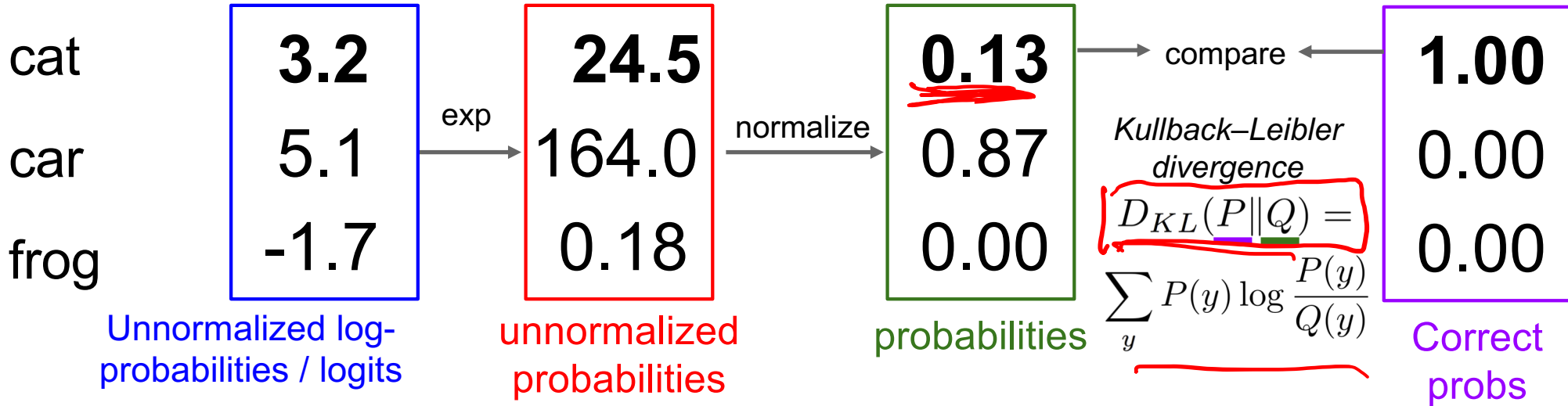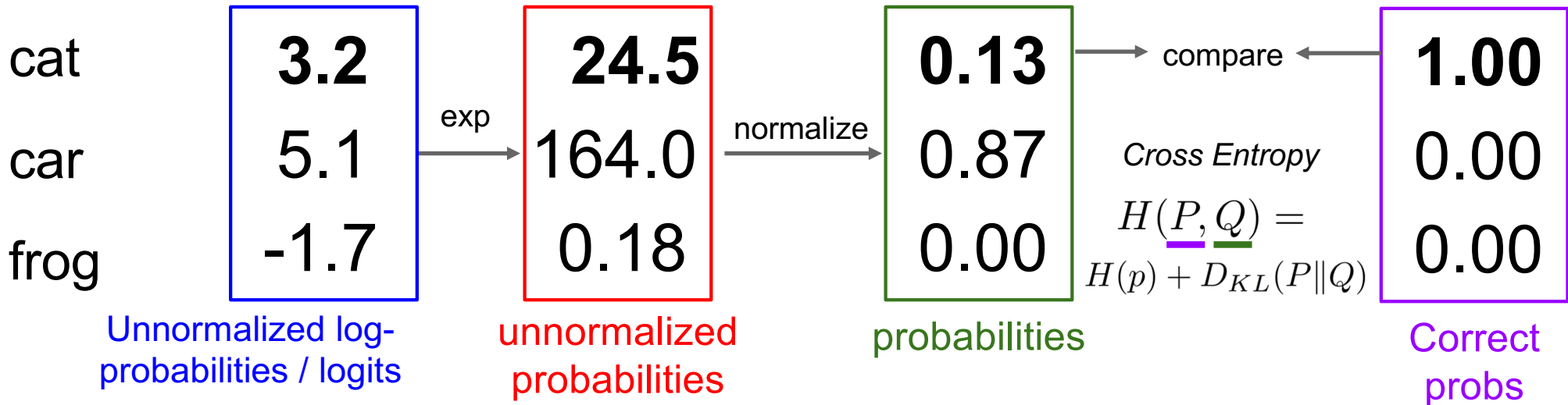
$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

Probabilities must be >= 0

Probabilities must sum to 1

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

|  | logits | unnormalized probs | probabilities | correct probs |
|------|------|-------|------|------|
| cat | **3.2** | **24.5** | **0.13** | **1.00** |
| car | 5.1 | 164.0 | 0.87 | 0.00 |
| frog | -1.7 | 0.18 | 0.00 | 0.00 |

exp → normalize → compare ←

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

Correct probs

Cross Entropy

$$H(P, Q) = H(p) + D_{KL}(P\|Q)$$

# 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?

$-\log(0)$

# **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?

$$-\log\left(\frac{1}{3}\right) = \log(3)$$

# 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

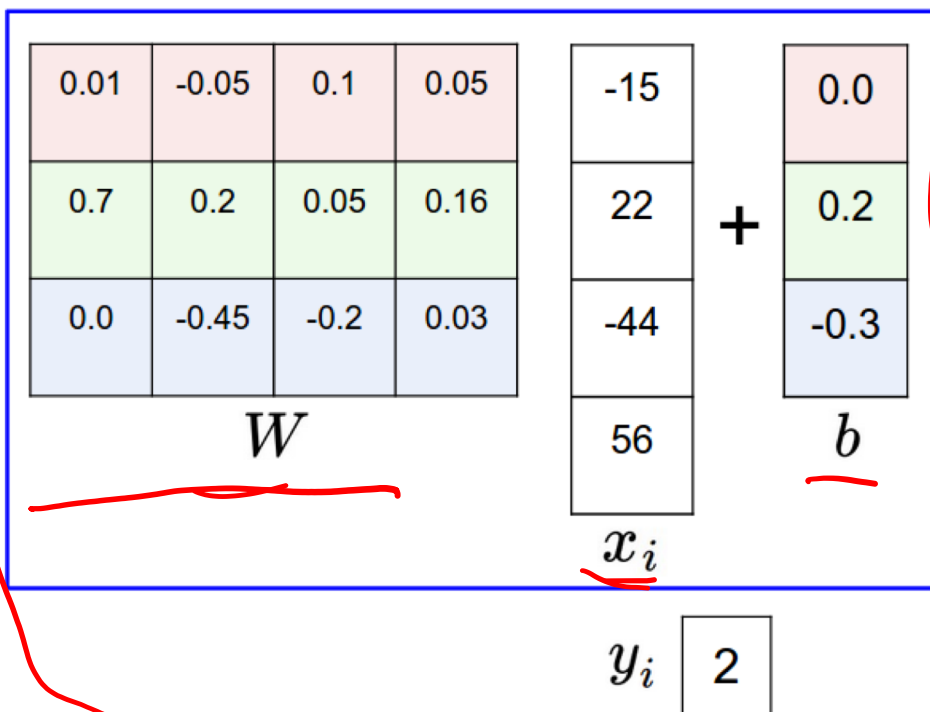$$L_i = -\log P(Y = y_i | X = x_i)$$

Putting it all together:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$
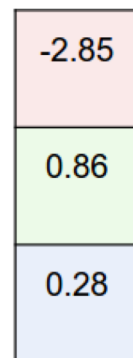
cat 3.2

car 5.1

frog -1.7

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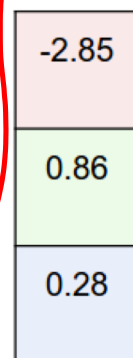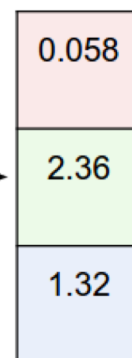
# Softmax vs. SVM

**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 |

$$b$$

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

**Model**

**hinge loss (SVM)**

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

$$\max(0, -2.85 - 0.28 + 1) + \max(0, 0.86 - 0.28 + 1)$$
$$=$$
$$\mathbf{1.58}$$

**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 |

$$-\log(0.353)$$
$$=$$
$$\mathbf{0.452}$$

**Objective**

# Softmax vs. SVM

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

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

# Softmax vs. SVM

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

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

---

assume scores:
[10, -2, 3]
[10, 9, 9]
[10, -100, -100]
and $y_i = 0$

Q: Suppose I take a datapoint and I jiggle a bit (changing its score slightly). What happens to the loss in both cases?
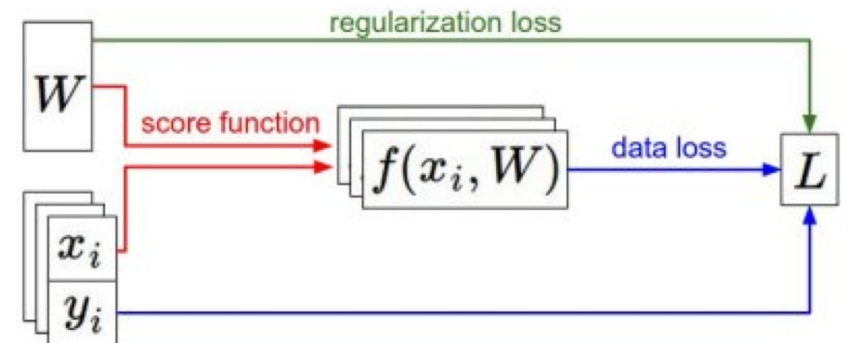
# Recap

- 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

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

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

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

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