

EN.601.482/682 Deep Learning

# Basics Part I:

## Image Features, Regression, and Classification

Mathias Unberath, PhD

Assistant Professor

Dept of Computer Science

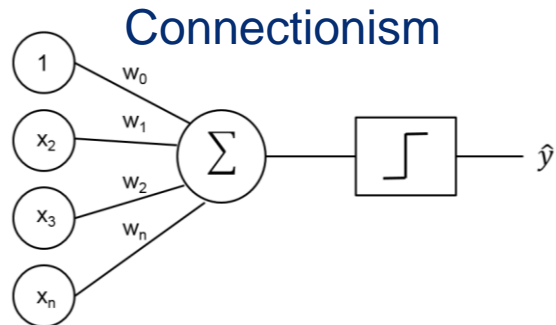
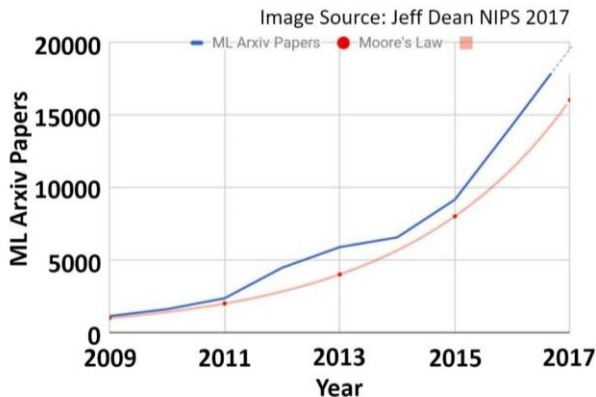
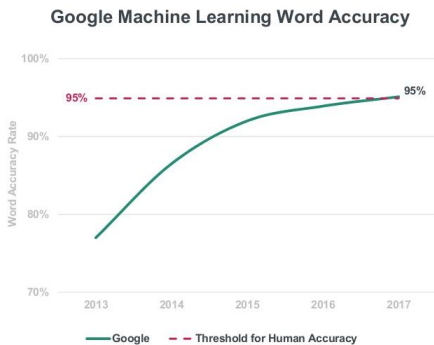
Johns Hopkins University

# Reminder

- Sign up on **Piazza** (access code dlF23)
- Read the **syllabus**
- **Homework assignment 1 will be released today** (due next Wednesday)

# Reminder

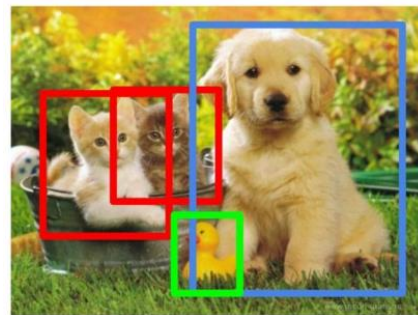
Neural networks are taking / have taken over!



Language processing



Object Detection



CAT, DOG, DUCK

# Reminder

- What is the difference between unsupervised and supervised learning?
- What is the difference between classification and regression?



# Today's Lecture

## Linear Classification

- Problem Statement
- Image Features
- The SVM Loss

## Logistic Regression

- The Softmax Function
- A Little Bit on Maximum Likelihood

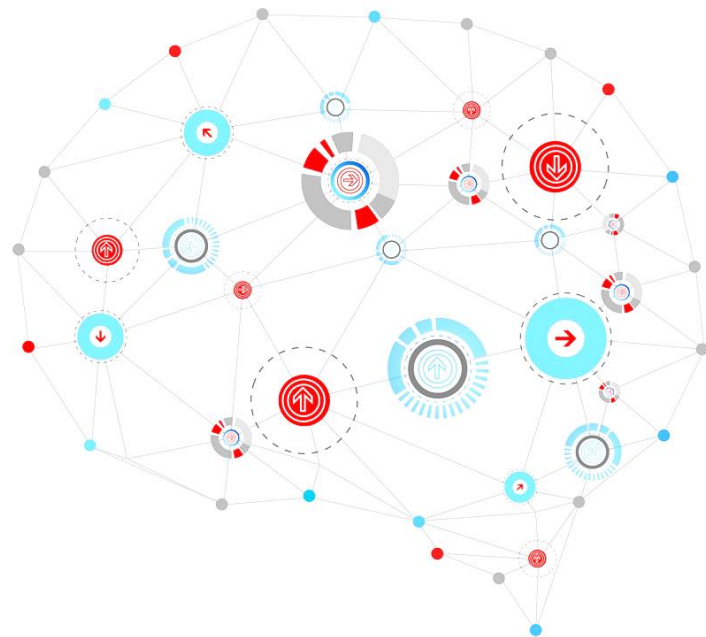


Image Features, Regression, and Classification

# Problem Statement



# From Classification to Instance Segmentation

**Classification**



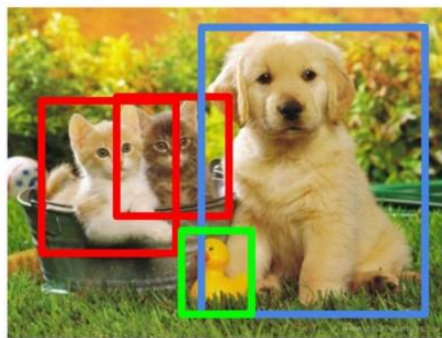
CAT

**Classification  
+ Localization**



CAT

**Object Detection**



CAT, DOG, DUCK

**Instance  
Segmentation**



CAT, DOG, DUCK

---

Single object

---

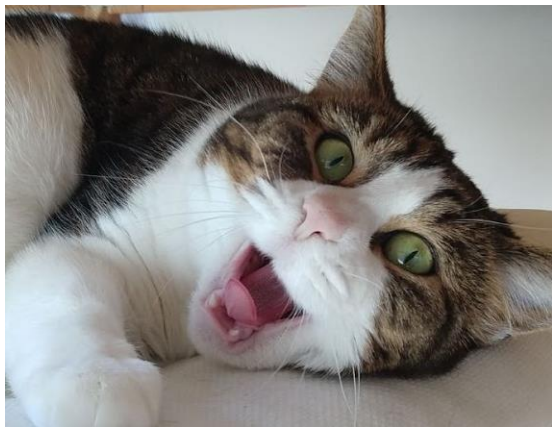
Multiple objects

[Ouaknine, A. \(2018\) Review of DL for Object Detection Blog Post](#)

# Image Classification

Assume set of discrete labels is known, e.g. {dog, cat, car, plane, ...}

Classification describes a mapping of image onto label

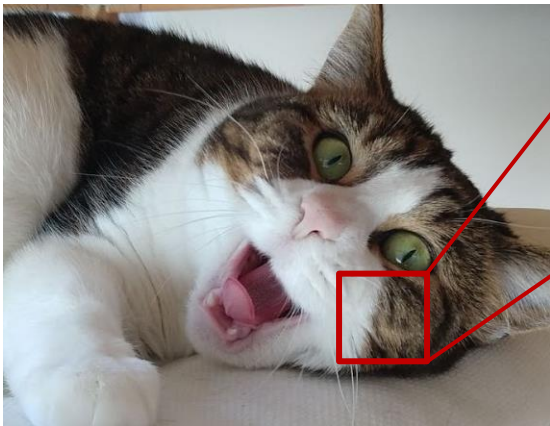


cat



# Image Classification

Why is this challenging?  
**Semantic gap**



What the human sees

1185	113	180	122	104	99	186	89	96	103	112	113	104	97	93	871
91	90	102	104	104	79	90	103	99	105	123	120	100	94	891	
76	85	90	105	120	105	87	98	99	115	112	100	105	99	891	
88	81	81	83	100	131	127	100	90	98	102	99	98	93	103	881
106	91	81	84	89	91	88	86	101	107	108	88	76	84	86	861
114	100	85	55	65	60	64	54	64	87	112	120	90	74	84	811
133	137	147	102	65	61	60	65	52	54	74	84	102	93	85	821
110	137	144	148	109	95	86	78	62	65	63	53	60	73	86	1011
1125	133	140	137	119	121	117	84	65	70	88	85	54	64	72	801
1127	125	131	147	133	137	136	111	96	88	75	61	64	72	841	
1115	114	100	123	130	140	113	110	113	100	100	92	74	65	72	791
89	93	90	97	100	147	133	138	113	114	133	109	100	90	77	801
89	77	80	81	77	79	102	129	117	113	117	129	125	150	125	871
82	86	82	88	78	71	68	101	120	136	138	101	107	114	131	1101
83	65	75	88	89	75	62	81	120	130	125	105	81	90	130	1101
87	65	71	87	106	96	68	45	76	130	126	107	92	84	105	1101
110	97	82	86	117	123	116	68	41	51	65	83	89	95	102	1071
1104	146	113	88	82	130	124	104	76	40	45	88	80	101	102	1031
1157	170	157	128	93	86	114	132	112	97	69	55	70	82	98	1041
1130	120	134	103	139	100	100	110	121	134	134	87	65	53	89	1011
1129	132	96	137	150	144	120	135	104	107	102	93	87	81	72	791
1129	107	90	88	83	117	135	149	122	109	104	75	80	107	112	901
1122	121	102	88	82	88	94	127	108	148	103	102	80	78	82	1071
1122	104	148	103	75	66	78	83	93	102	110	120	102	61	68	1011

What the computer sees

An image is essentially a grid of numbers (often, but not always) between [0,255].

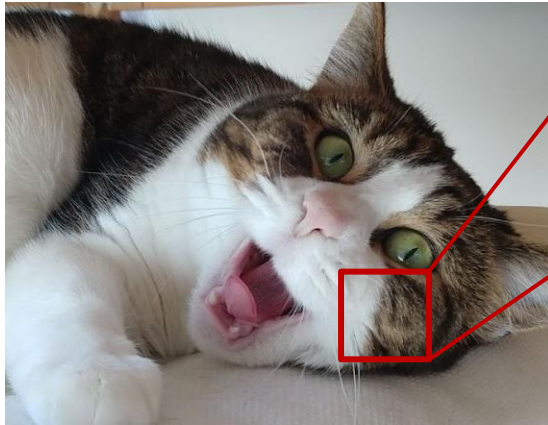
Very common:

CV: 640 x 480 x 3 (RGB channels, 8 bit)

X-ray: 1240 x 960 x 1 (Grayscale, 14 bit)

# Image Classification

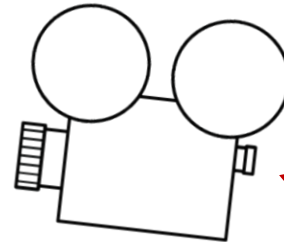
Why is this challenging?  
**Viewpoint**



What the human sees

1185	112	180	122	104	99	106	89	96	103	112	118	104	97	93	871
91	90	102	104	104	79	90	103	99	105	123	120	110	105	94	851
76	85	90	105	120	105	87	98	95	99	115	112	100	105	99	851
88	81	81	83	100	131	127	100	90	98	102	99	98	93	103	841
106	91	81	84	89	91	88	86	101	107	108	88	76	84	86	861
114	100	85	55	65	60	64	54	64	87	112	120	100	74	84	811
133	137	147	102	65	61	60	65	52	54	74	84	102	93	85	821
110	137	144	140	109	95	86	70	62	65	63	53	60	73	86	1011
1125	133	140	137	119	121	117	84	65	70	88	85	54	64	72	801
1127	125	131	147	133	137	136	111	96	88	75	61	64	72	84	811
1125	114	100	123	130	140	113	110	113	100	100	92	74	65	72	791
89	93	90	97	100	147	133	138	113	114	133	109	100	90	77	801
89	77	80	81	77	79	102	129	117	119	117	129	129	150	125	871
82	86	82	88	70	71	68	101	120	136	138	101	107	114	131	1101
83	65	76	88	80	75	62	81	120	130	126	105	81	90	130	1101
87	65	71	87	100	96	60	45	76	130	126	107	82	84	105	1101
110	97	82	86	117	123	116	68	41	51	65	83	80	95	102	1071
104	146	113	88	82	130	124	104	76	40	45	88	80	101	102	1031
1157	170	157	128	93	86	114	132	112	97	69	55	70	82	98	941
1100	120	134	101	139	100	100	138	121	134	134	87	65	53	89	901
1129	132	96	137	150	144	120	135	104	107	102	93	87	81	72	791
1129	107	90	88	83	117	153	149	122	109	104	75	80	107	132	901
1122	121	102	88	82	88	94	127	100	148	103	102	80	78	82	1071
1122	104	148	103	75	66	70	83	93	102	110	120	102	61	80	1011

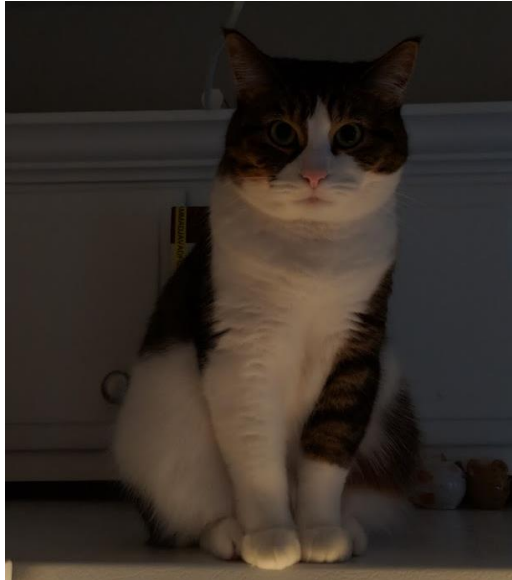
What the computer sees



When camera moves,  
all pixel values change!

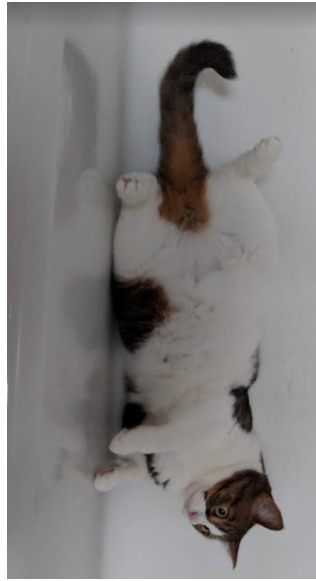
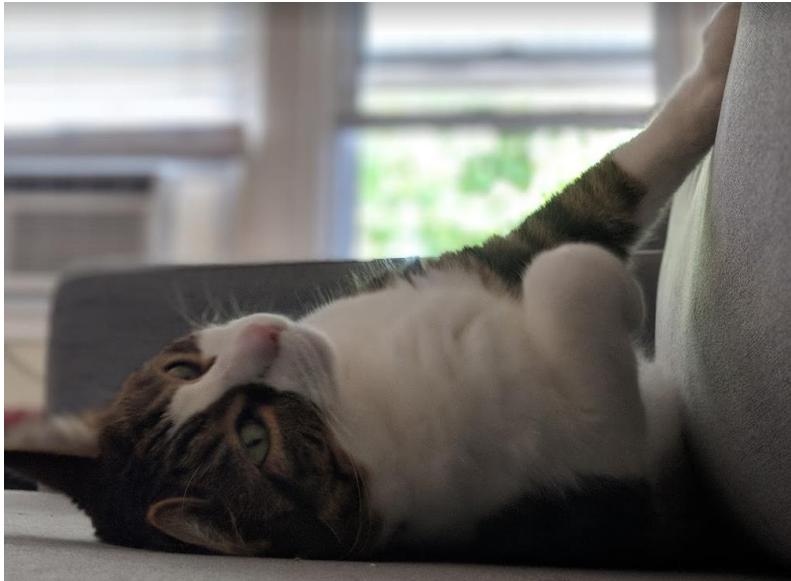
# Image Classification

Why is this challenging?  
**Lighting conditions**



# Image Classification

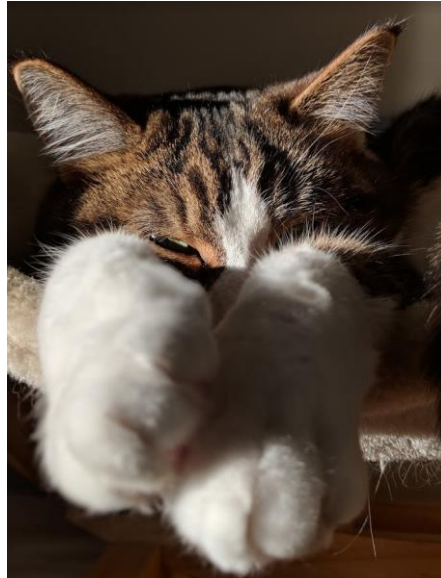
Why is this challenging?  
**Deformation**



# Image Classification

Why is this challenging?

**Occlusion**





# Image Classification

Why is this challenging?  
**Background clutter**



# Image Classification

Why is this challenging?  
**Intraclass variation**



# An Example Classifier

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```

There is no obvious way to hard-code such classification algorithm.



# The Machine Learning Approach

1. Curate a (large) dataset of images with corresponding annotation (i.e. labels)
2. Train a classifier using machine learning
3. Evaluate the classifier on unseen images

Example training set



# The Machine Learning Approach

1. Curate a (large) dataset of images with corresponding annotation (i.e. labels)
2. Train a classifier using machine learning
3. Evaluate the classifier on unseen images

Example dataset: **CIFAR10**

10 classes

50,000 training images

10,000 testing images

Each image: 32 x 32 x 3



[Krizhevsky, A. \(2009\) Learning Multiple Layers of Features from Tiny Images](#)

# Image Classification

If there are multiple classes, e.g. {cat, car, frog}  
→ Response for every label

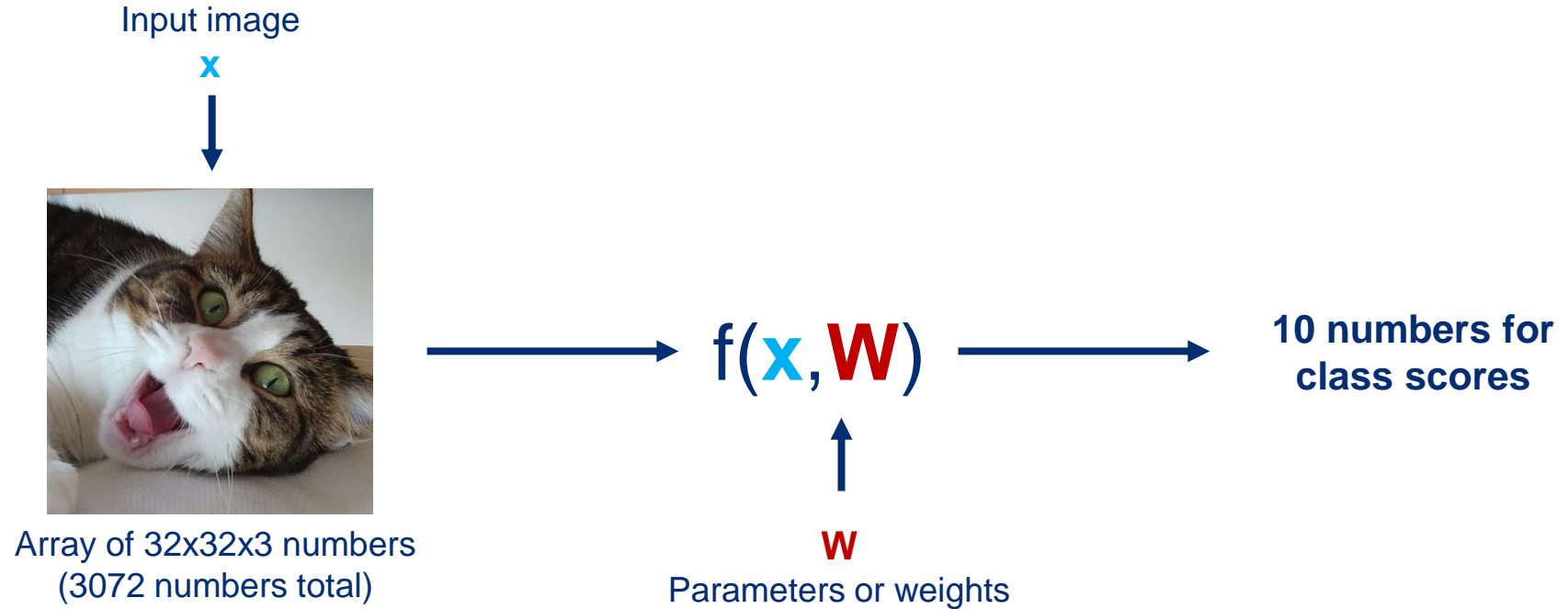


5.1	cat
2.3	car
-0.7	frog

1
0
0

One-hot encoding

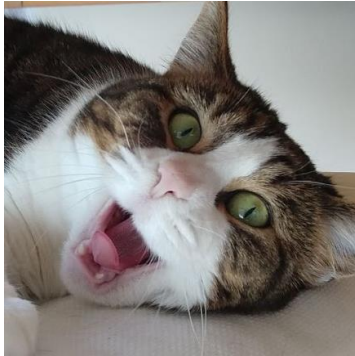
# Parametric Approach



# Parametric Approach

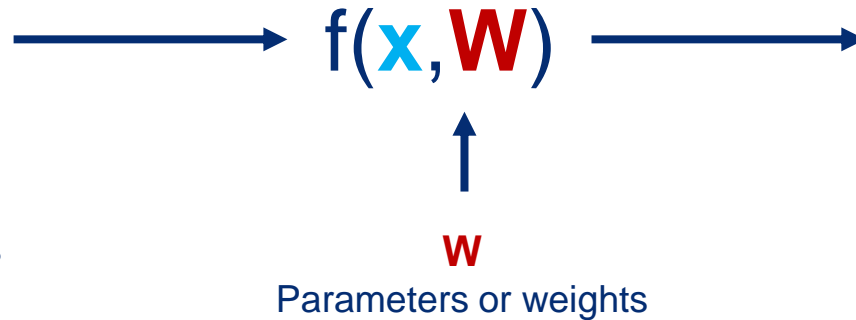
Input image

$x$



Array of 32x32x3 numbers  
(3072 numbers total)

Goal in machine learning:  
**Find this  $f(x, W)$ !**

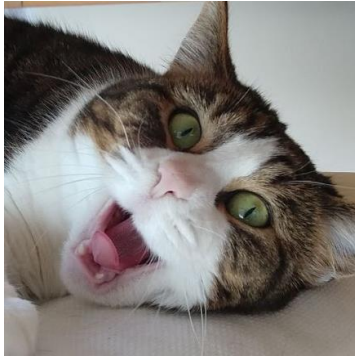


10 numbers for  
class scores

# Linear Classifier

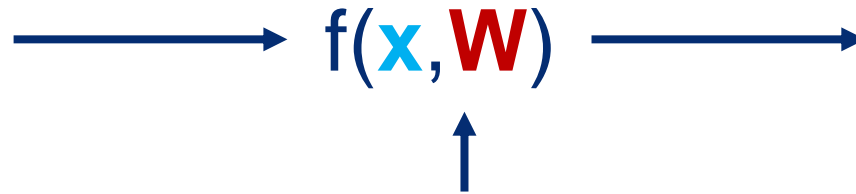
Input image

$x$



Array of 32x32x3 numbers  
(3072 numbers total)

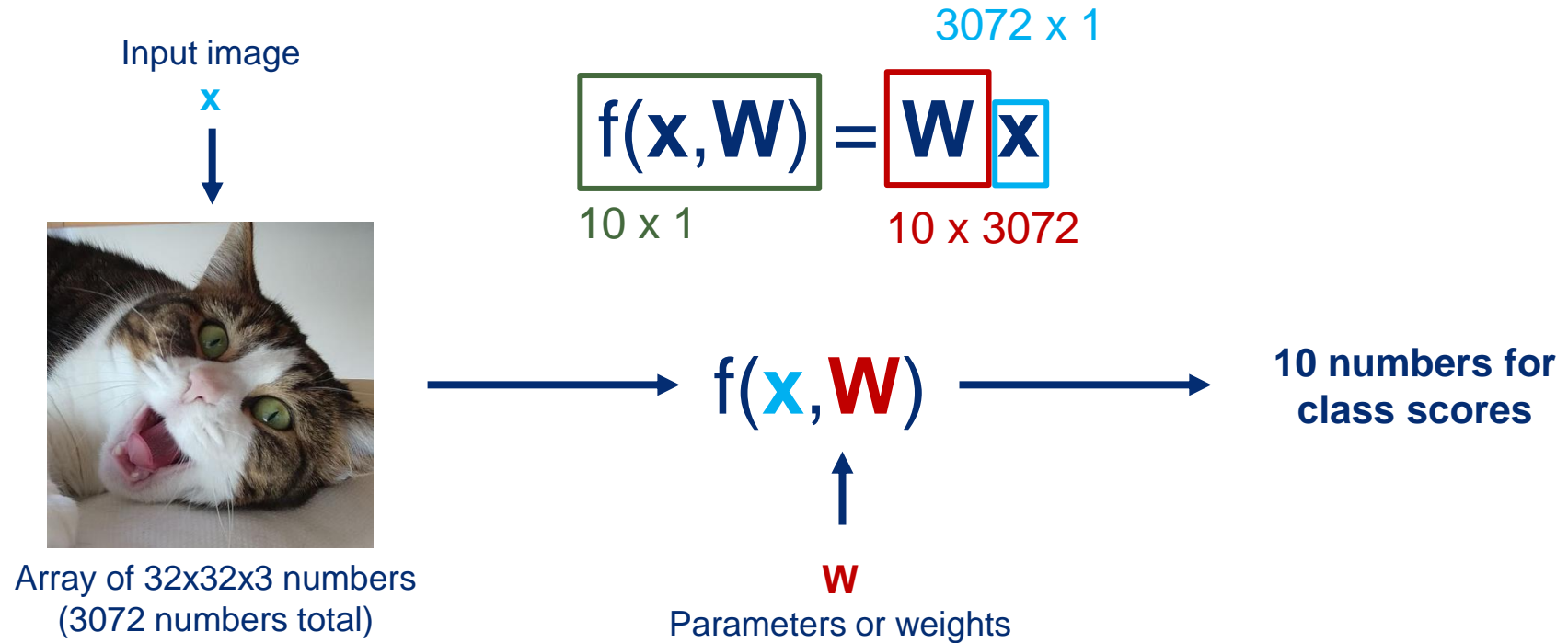
$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x}$$



$\mathbf{W}$   
Parameters or weights

10 numbers for  
class scores

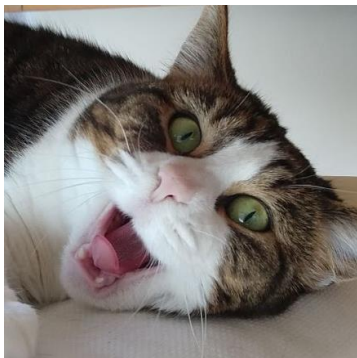
# Linear Classifier



# Linear Classifier

Input image

$x$

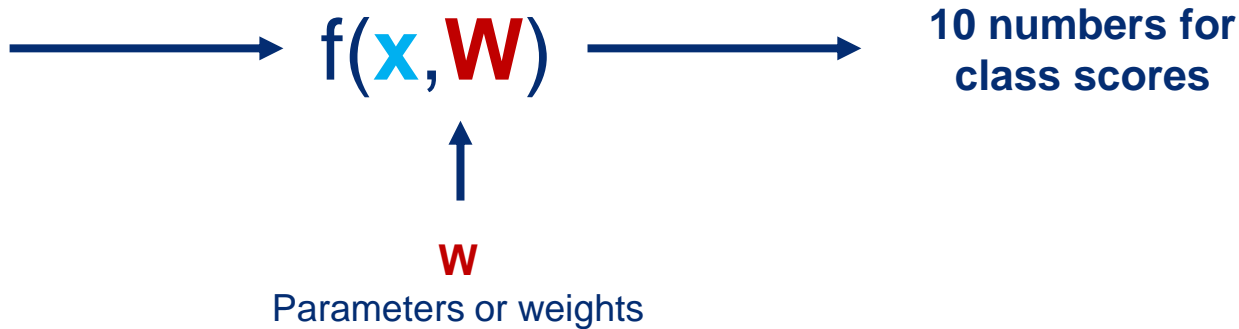


Array of 32x32x3 numbers  
(3072 numbers total)

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

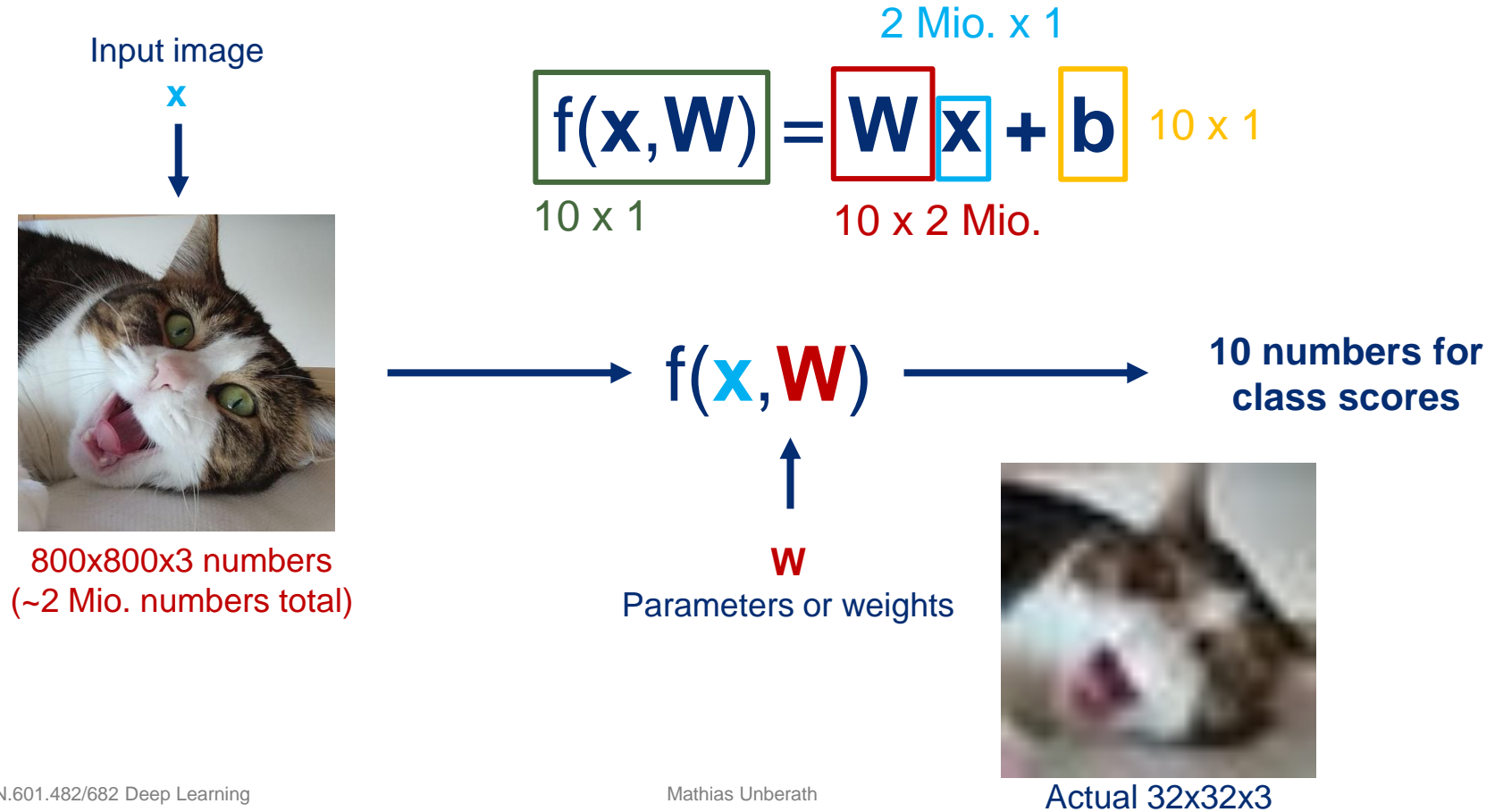
$10 \times 1$        $3072 \times 1$        $10 \times 3072$        $10 \times 1$

Bias term: constant!  
→ Input independent preference

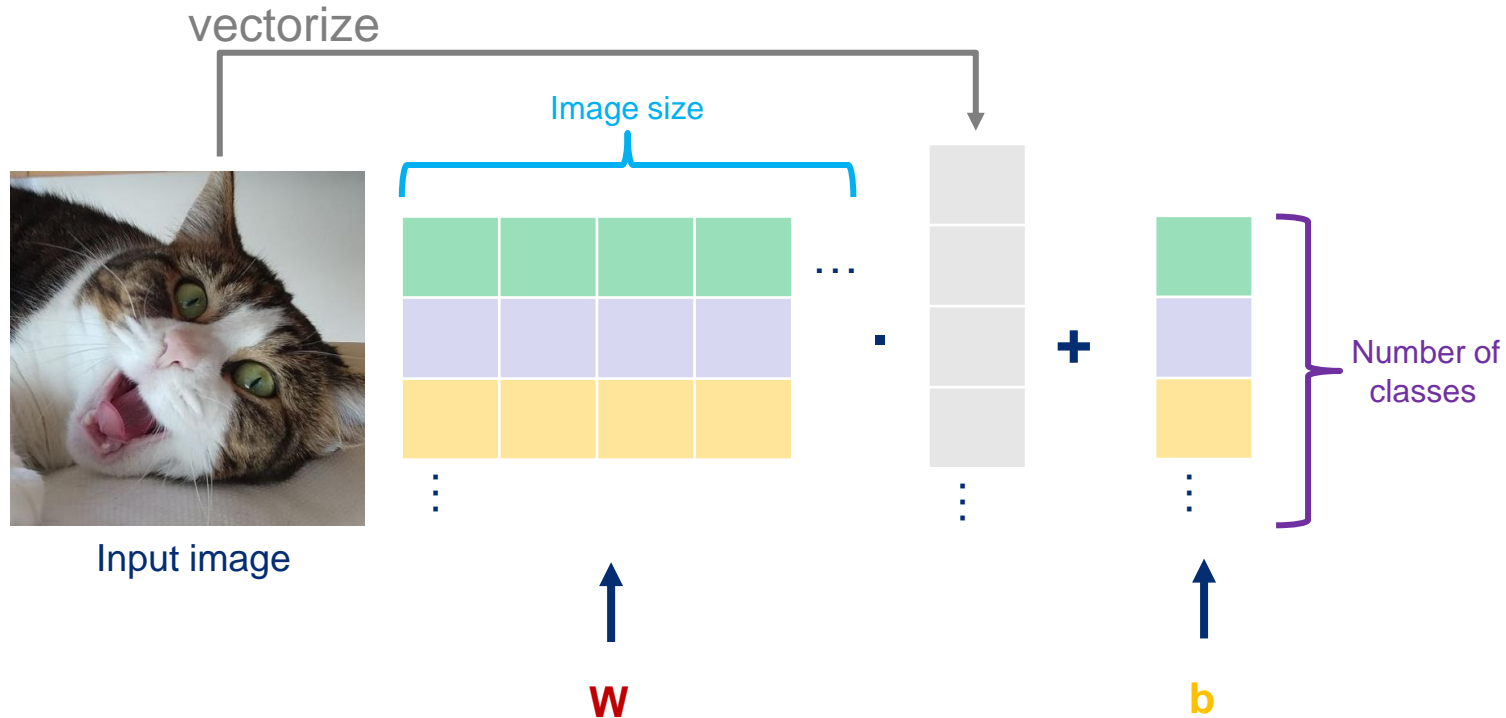




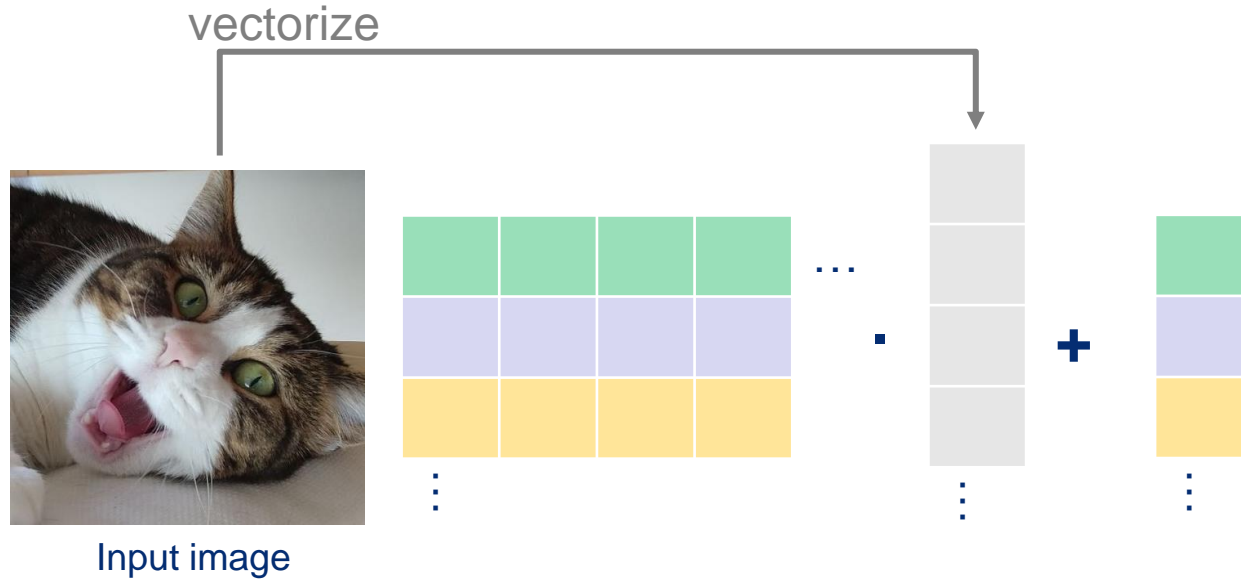
# Linear Classifier



# Interpretation of Linear Classification in this Formalism



# Interpretation of Linear Classification in this Formalism



- Every row  ... of  $W$  acts like a “template” for the corresponding class

# Interpretation of Linear Classification in this Formalism

**Visual  
Interpretation**



Input image

vectorize

airplane

automobile

bird

cat

deer

dog

frog

horse

ship

truck



plane

car

bird

cat

deer

dog

frog

horse

ship

truck



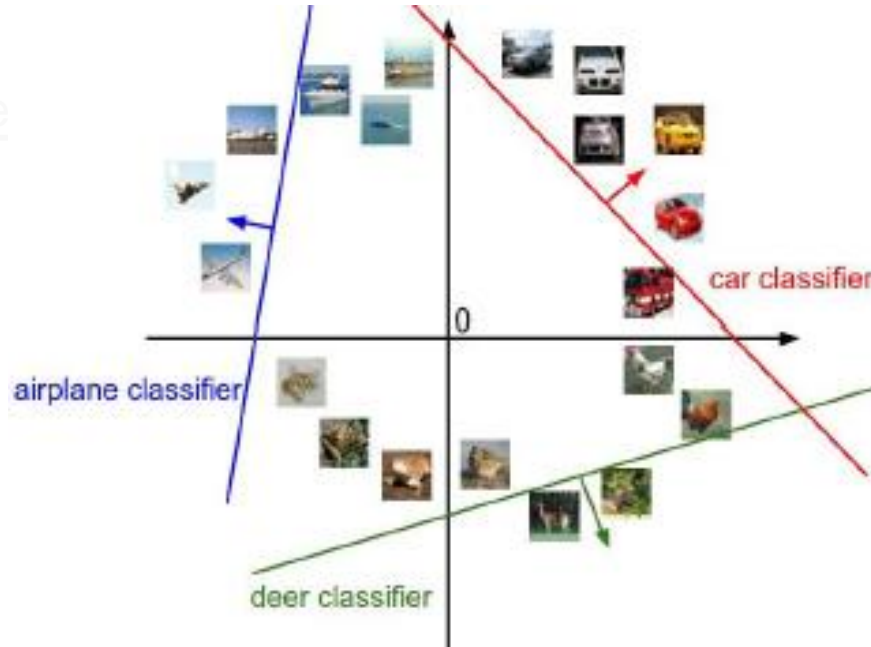
# Interpretation of Linear Classification in this Formalism

## Geometric Interpretation



Input image

vectorize



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



# An Obvious Challenge

Template assumption on  
pixel level is strong!

We need some form of abstraction  
→ Enter: **Features**!

Seagull



# Where are we now? Linear classifier

- Algebraic formalism:  $f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} + \mathbf{b}$
- Visual interpretation:  
Rows of  $\mathbf{W}$  form templates for classes
- Geometric interpretation:  
Instances  $\mathbf{x}$  are points in high-dimensional space  
Rows of  $\mathbf{W}$  define hyperplanes (linear decision boundaries)
- Challenge:  
Images exhibit lots of variation! Template assumption difficult to justify.  
→ Higher-level representation of images

Image Features, Regression, and Classification

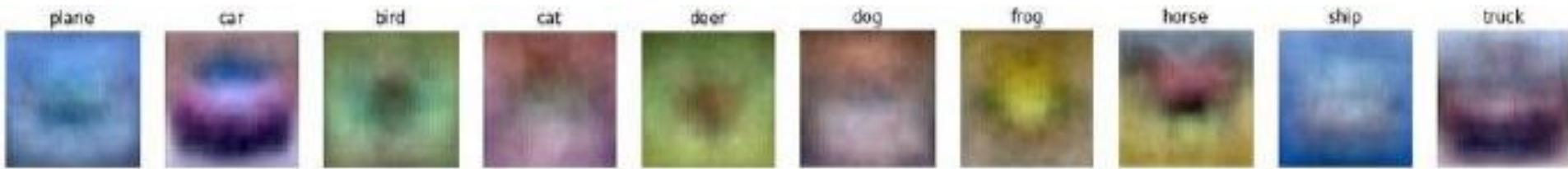
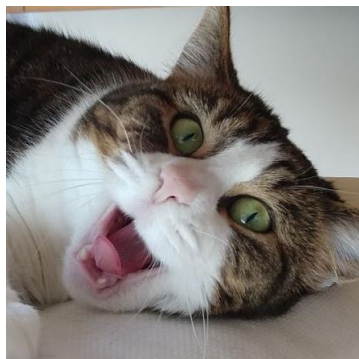
# Image Features





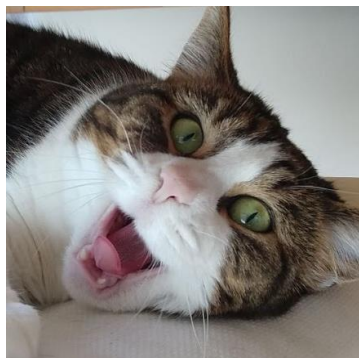
# Image Features

- What are the image features in this example? Pixel intensities!



# Parametric Approach

- Higher-level representation via feature extraction



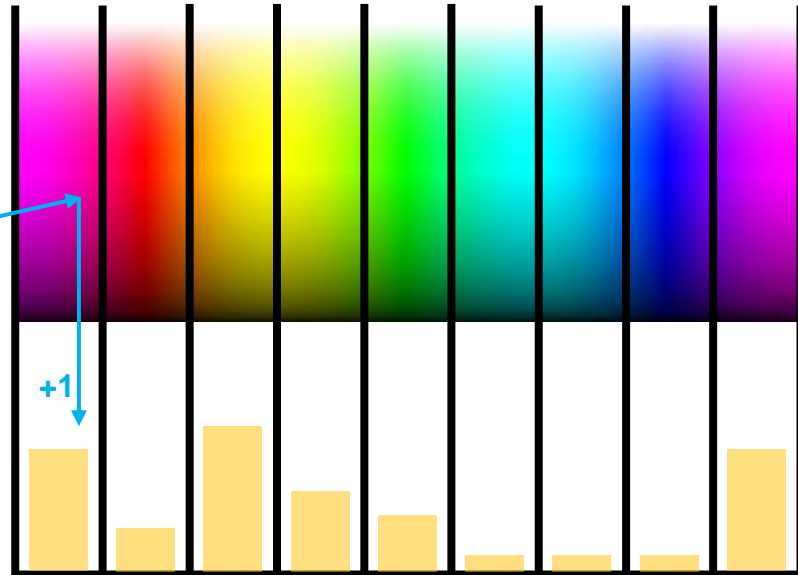
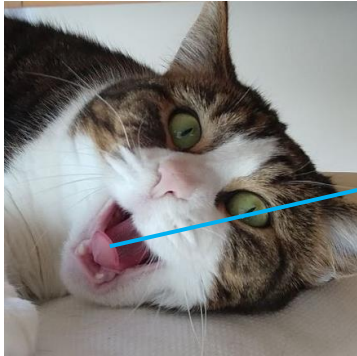
Feature  
extractor



class scores

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

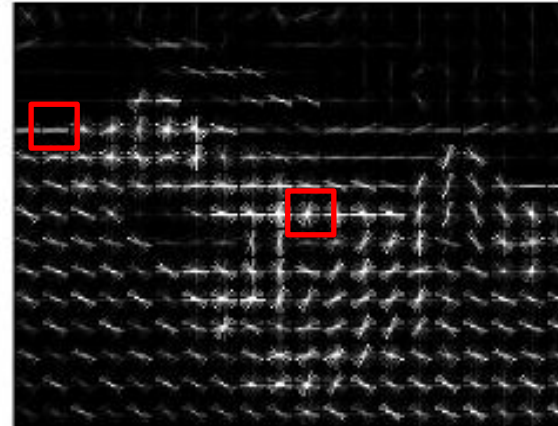
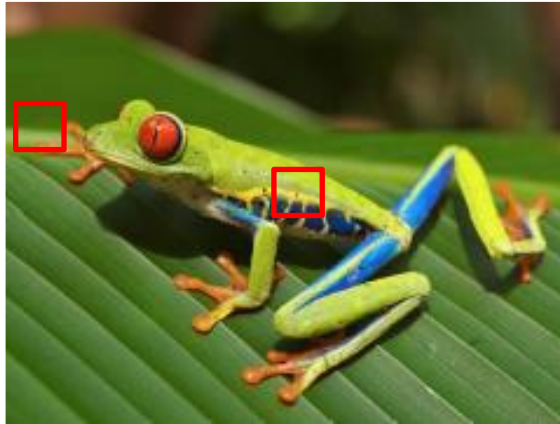
# Feature Extractor: Color Histogram



Color histogram is a *global* descriptor  
Feature vector has  $\#(\text{hist bins})$  elements

# Feature Extractor: Histogram of Oriented Gradients

- Divide into small regions
- Quantize gradient direction into bins



Example numbers:

Input image 320x240; Patch size: 8x8; 9 gradient bins.

→ 40x30 patches, each patch is 9-bin histogram: Feature vector has  $40 \times 30 \times 9 = 10,800$  elements!

[Dalal, N., & Triggs, B. \(2005\) Histograms of oriented gradients for human detection. CVPR](#)

Images from Stanford cs231n.

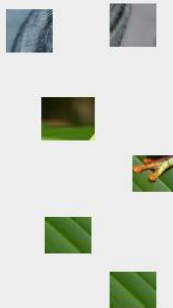


# Feature Extractor: Bag of Visual Words

## Step 1: Build codebook



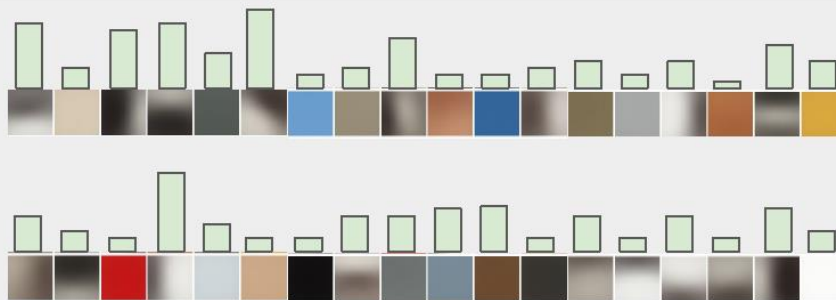
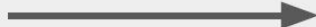
Extract random patches



Cluster patches to form "codebook" of "visual words"



## Step 2: Encode images



# Some Observations

- For all feature extractors, we have some design choice (hyperparameters)
  - E.g. for color histogram: Number of bins
- Feature engineering is not trivial
  - Requires strong domain knowledge
- Previously:



Guyon, I., & Elisseeff, A. (2003) An introduction to variable and feature selection. *Journal of machine learning research*, 3(Mar), 1157-1182.

# Some Observations

- For all feature extractors, we have some design choice (hyperparameters)
  - E.g. for color histogram: Number of bins
- Feature engineering is not trivial
  - Requires strong domain knowledge
- Teaser:

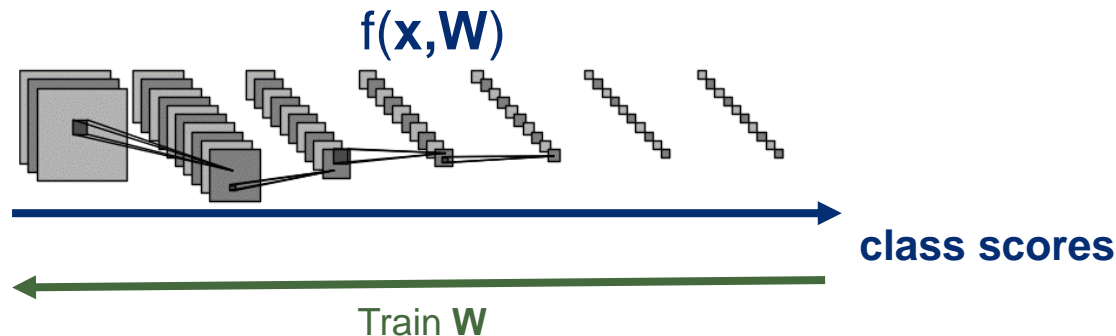
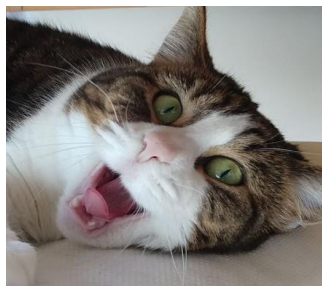


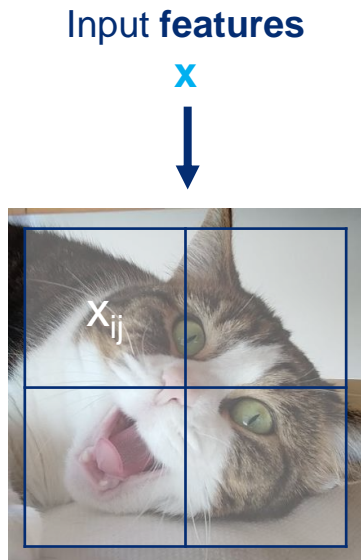
Image Features, Regression, and Classification

# The SVM Loss

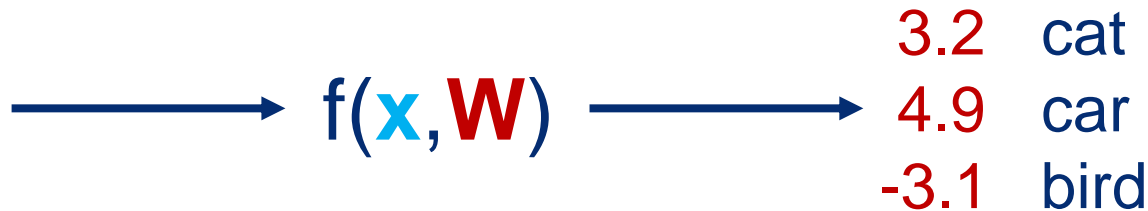




# Unhappy?



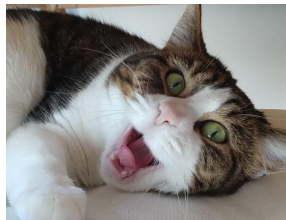
$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} + \mathbf{b}$$



- Given some  $\mathbf{W}$ , how to express our (un)happiness with the current prediction?  
→ Loss functions
- How to minimize our unhappiness? Update  $\mathbf{W}$  → Optimization

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat

**3.2**

1.3

2.2

Car

5.1

**4.9**

2.5

Bird

-1.7

2.0

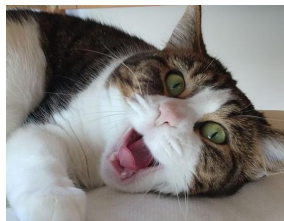
**-3.1**

These images are CC0 1.0 public domain.



# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat

**3.2**

1.3

2.2

Car

5.1

**4.9**

2.5

Bird

-1.7

2.0

**-3.1**

Loss function:

Quantifies classifier performance

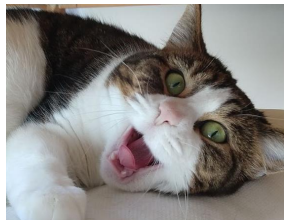
Given a dataset  $\{(x_i, y_i)\}_{i=1}^N$  where  
 $x_i$  the image (or its feature representation)  
 $y_i$  the corresponding label (integer)

Loss over dataset = sum over all examples

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

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat

**3.2**

1.3

2.2

Car

5.1

**4.9**

2.5

Bird

-1.7

2.0

**-3.1**

Given a dataset  $\{(x_i, y_i)\}_{i=1}^N$  where  
 $x_i$  the image (or its feature representation)  
 $y_i$  the corresponding label (integer)

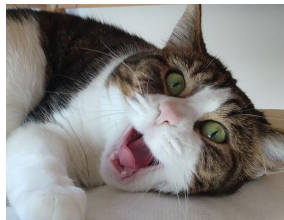
Using:  $s = f(x_i, W)$

The Support Vector Machine (SVM) loss:

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

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat

**3.2**

1.3

2.2

Car

5.1

**4.9**

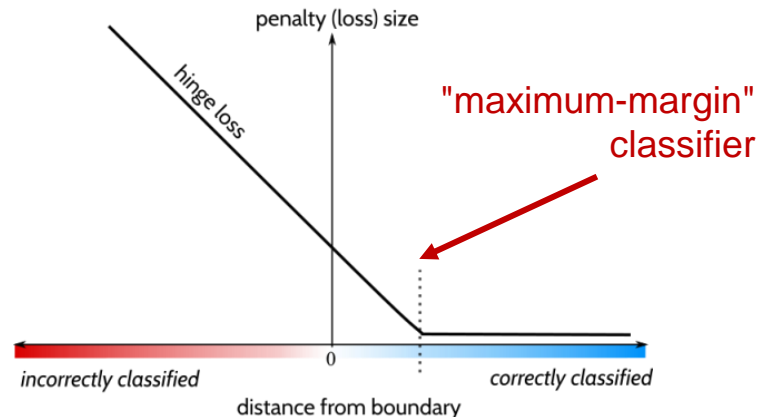
2.5

Bird

-1.7

2.0

**-3.1**

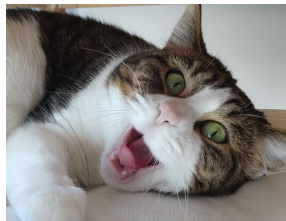


The Support Vector Machine (SVM) loss:

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

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat	<b>3.2</b>	1.3	2.2
Car	5.1	<b>4.9</b>	2.5
Bird	-1.7	2.0	<b>-3.1</b>
Loss	<b>2.9</b>		

The Support Vector Machine (SVM) loss:

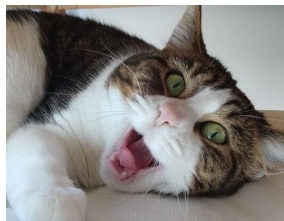
$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

**Example: Cat**

$$\begin{aligned} L_i &= \max(0, 5.1 - 3.2 + 1) \\ &\quad + \max(0, -1.7 - 3.2 + 1) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 \\ &= 2.9 \end{aligned}$$

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat	3.2	1.3	2.2
Car	5.1	4.9	2.5
Bird	-1.7	2.0	-3.1
Loss	2.9	0	

The Support Vector Machine (SVM) loss:

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

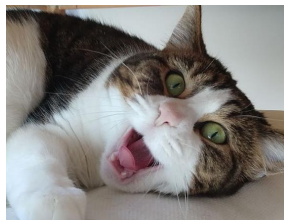
**Example: Car**

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



# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat	3.2	1.3	2.2
Car	5.1	4.9	2.5
Bird	-1.7	2.0	-3.1
Loss	2.9	0	12.9

The Support Vector Machine (SVM) loss:

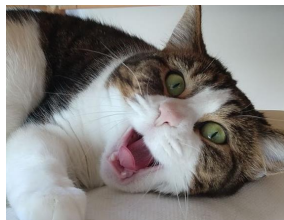
$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

**Example: Bird**

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

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat	<b>3.2</b>	1.3	2.2
Car	5.1	<b>4.9</b>	2.5
Bird	-1.7	2.0	<b>-3.1</b>
Loss	<b>2.9</b>	<b>0</b>	<b>12.9</b>

The Support Vector Machine (SVM) loss:

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

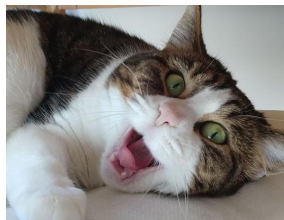
Overall loss:

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

$$L = 1/3 * (2.9 + 0 + 12.9) = \mathbf{5.27}$$

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat	<b>3.2</b>	1.3	2.2
Car	5.1	<b>4.9</b>	2.5
Bird	-1.7	2.0	<b>-3.1</b>
Loss	<b>2.9</b>	<b>0</b>	<b>12.9</b>

The Support Vector Machine (SVM) loss:

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

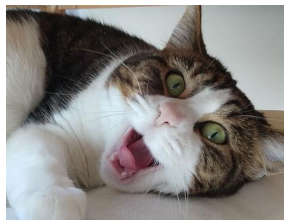
Overall loss:

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

Q1: What happens to  $L$  if car score changes a bit?

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat	<b>3.2</b>	1.3	2.2
Car	5.1	<b>4.9</b>	2.5
Bird	-1.7	2.0	<b>-3.1</b>
Loss	<b>2.9</b>	<b>0</b>	<b>12.9</b>

The Support Vector Machine (SVM) loss:

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

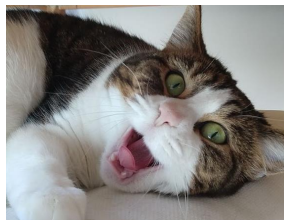
Overall loss:

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

Q2: What is the min/max for  $L$ ?

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat	<b>3.2</b>	1.3	2.2
Car	5.1	<b>4.9</b>	2.5
Bird	-1.7	2.0	<b>-3.1</b>
<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>

The Support Vector Machine (SVM) loss:

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

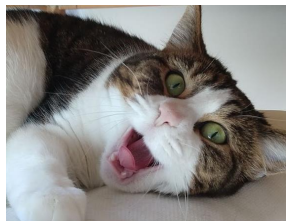
Overall loss:

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

Q3: If  $W$  is small and all  $s \approx 0$ , what is  $L$ ?

# The SVM Loss

- 3 training examples and 3 classes: cat, car, bird
- $W$  has been determined, the scores are:



Cat	<b>3.2</b>	1.3	2.2
Car	5.1	<b>4.9</b>	2.5
Bird	-1.7	2.0	<b>-3.1</b>
Loss	<b>2.9</b>	<b>0</b>	<b>12.9</b>

The Support Vector Machine (SVM) loss:

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

True label for cat (one-hot encoding).

<b>1.00</b>
0.00
0.00

Q4: Why not use the L2 norm?

Image Features, Regression, and Classification

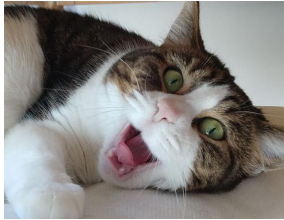
# Logistic Regression





# Multinomial Logistic Regression

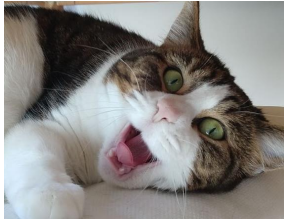
- Predictions as per linear regression are unbounded
- For classification, we would like to interpret scores as **probabilities**



Cat	<b>3.2</b>
Car	5.1
Bird	-1.7

# Multinomial Logistic Regression

- Again:  $s = f(x_i, W)$
- Then:  $P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$  Softmax function

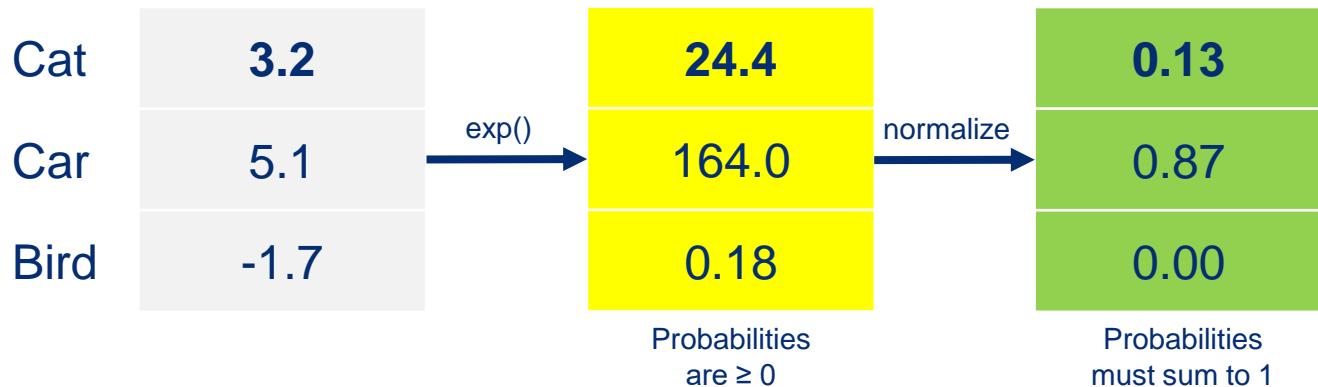
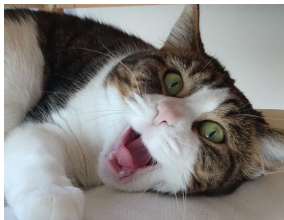


Cat	3.2	exp() →	24.4
Car	5.1		164.0
Bird	-1.7		0.18

Probabilities  
are  $\geq 0$

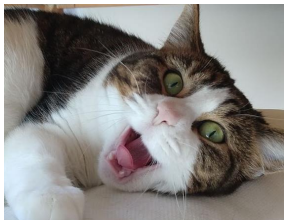
# Multinomial Logistic Regression

- Again:  $s = f(x_i, W)$
- Then:  $P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$  Softmax function

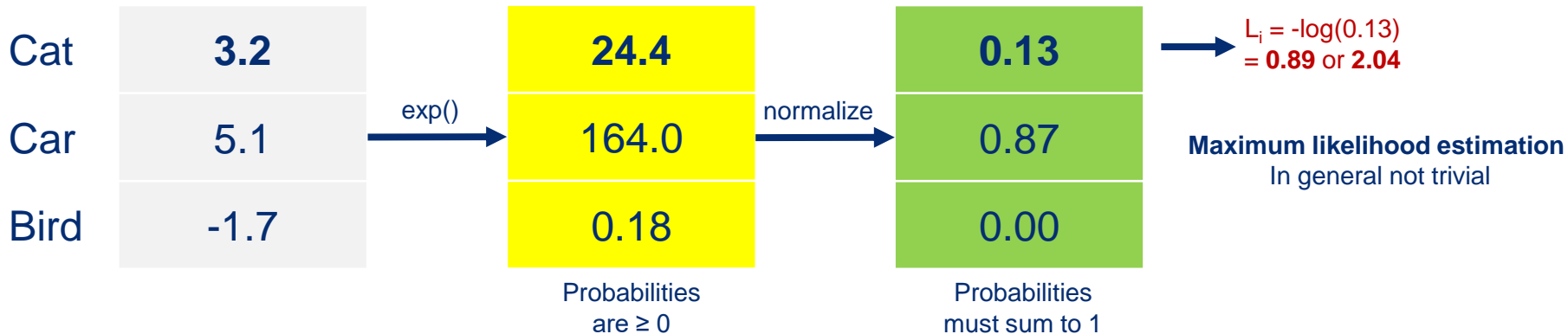


# Multinomial Logistic Regression

- Again:  $s = f(x_i, W)$
- Then:  $P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$  Softmax function

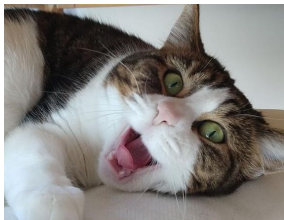


Loss:  $L_i = -\log P(Y = y_i|X = x_i) \rightarrow$  Log-likelihood of true class

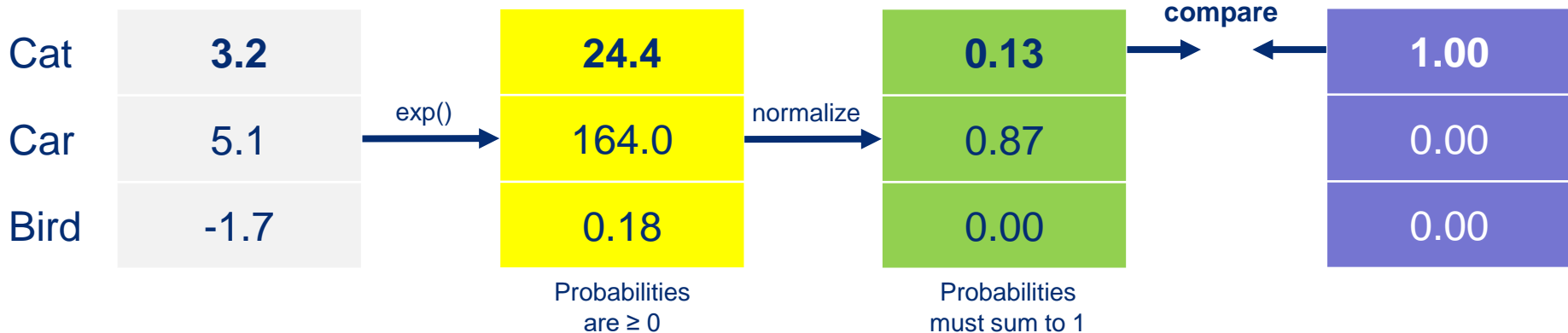


# Multinomial Logistic Regression

- Again:  $s = f(x_i, W)$
- Then:  $P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$  **Softmax function**

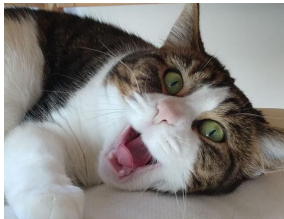


Loss:  $D_{KL}(Q, P) = \sum_k Q(k) \log \frac{Q(k)}{P(k)}$  → Kullback-Leibler Divergence



# Multinomial Logistic Regression

- Again:  $s = f(x_i, W)$
- Then:  $P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$  Softmax function



Maximize probability of correct class:

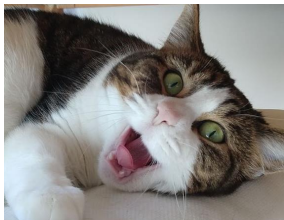
$$L_i = -\log P(Y = y_i|X = x_i) = -\log \frac{e^{s_{y_i}}}{\sum_j e^{s_j}}$$

Cat	3.2
Car	5.1
Bird	-1.7

**Q1: What is min/max of this loss?**

# Multinomial Logistic Regression

- Again:  $s = f(x_i, W)$
- Then:  $P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$  Softmax function



Maximize probability of correct class:

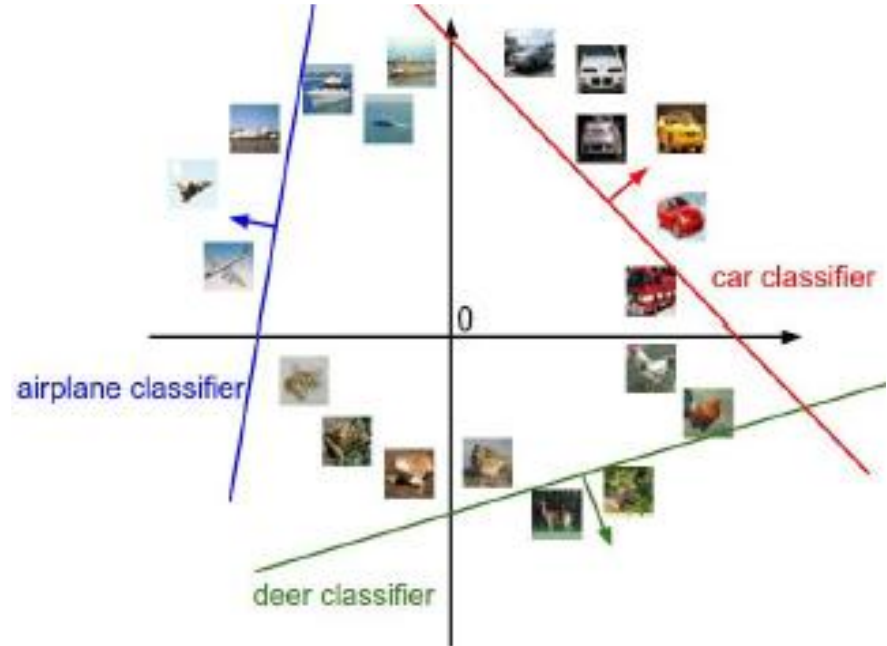
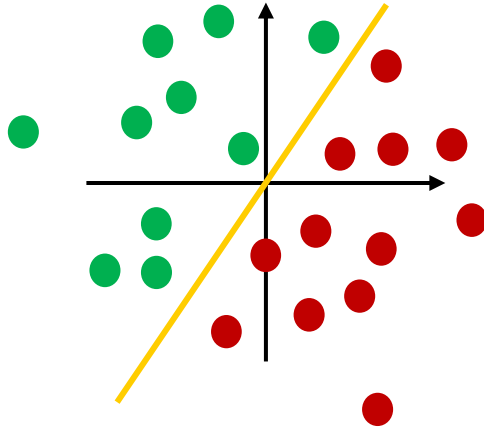
$$L_i = -\log P(Y = y_i|X = x_i) = -\log \frac{e^{s_{y_i}}}{\sum_j e^{s_j}}$$

Cat	3.2
Car	5.1
Bird	-1.7

**Q2: At initialization, all scores are similar and close to zero?**

# Decision Boundaries of Linear Regression

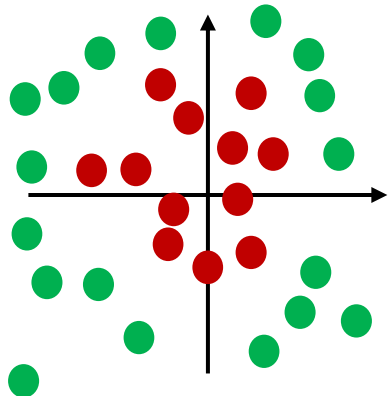
- Scores:  $s = f(x_i, W) = Wx_i + b$
- SVM loss: “Maximum margin” at decision boundary



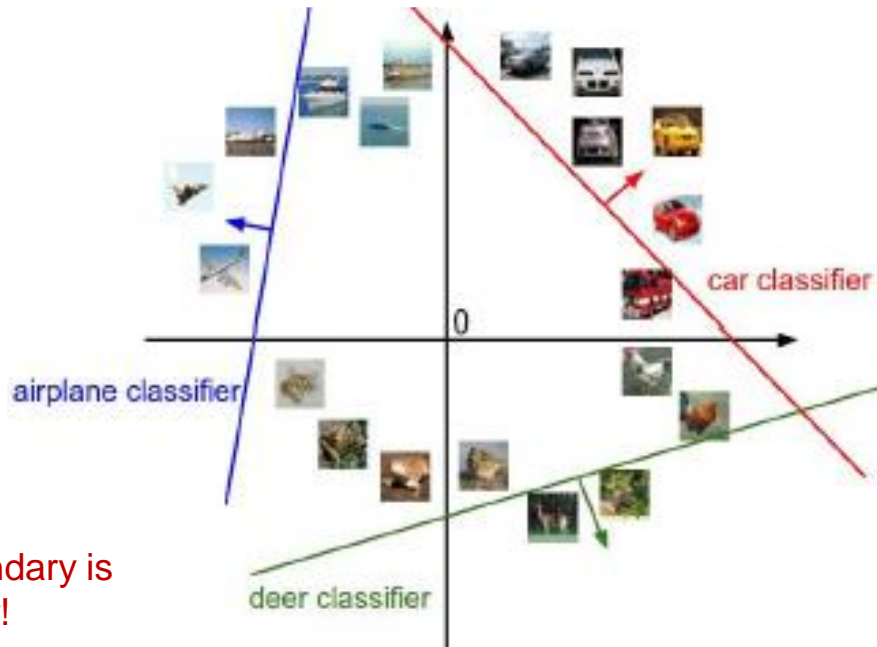


# Decision Boundaries of Linear Regression

- Scores:  $s = f(x_i, W) = Wx_i + b$
- SVM loss: “Maximum margin” at decision boundary



Required decision boundary is clearly non-linear!



**Q: Can I solve this with Logistic Regression?**

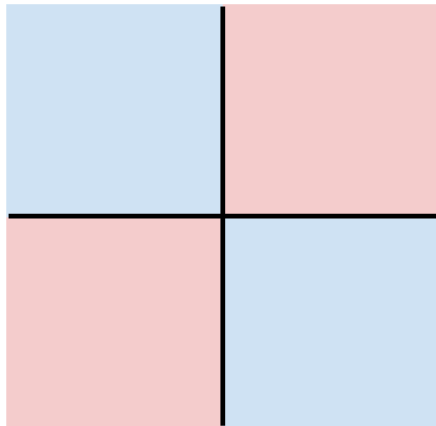
# Hard Cases for Linear Classification

**Class 1:**

First and third quadrants

**Class 2:**

Second and fourth quadrants

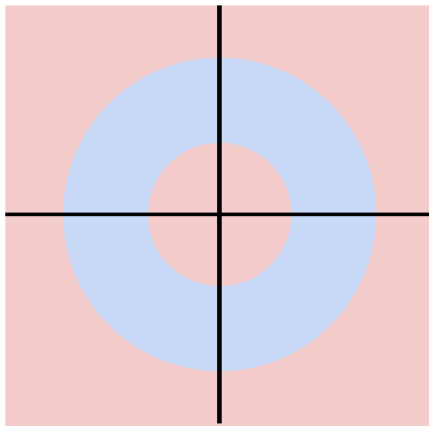


**Class 1:**

$1 \leq \text{L2 norm} \leq 2$

**Class 2:**

Everything else

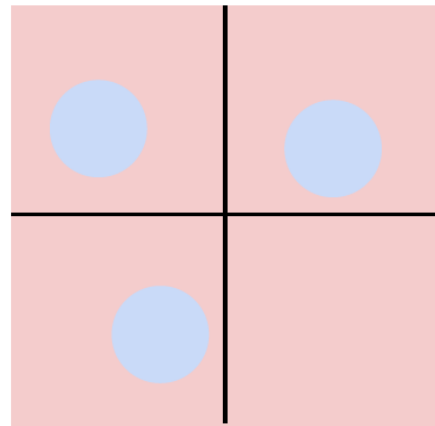


**Class 1:**

Three modes

**Class 2:**

Everything else



- Clearly, the two classes are not linearly separable!

# Teaser

## Neural Network

Linear  
classifiers



**Next lecture:**  
- **Regularization**  
- **Optimization**

Image Features, Regression, and Classification

**Questions?**

