

Lecture 3: Linear Classifiers

Reminder: Assignment 1

- <http://web.eecs.umich.edu/~justincj/teaching/eecs498/assignment1.html>
- Due **Sunday September 15, 11:59pm EST**
- We have written a **homework validation script** to check the format of your .zip file before you submit to Canvas:
- <https://github.com/deeplearning-class/tools#homework-validation>
- This script ensures that your .zip and .ipynb files are properly structured; they do not check correctness
- **It is your responsibility** to make sure your submitted .zip file passes the validation script

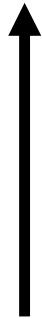
Last time: Image Classification

Input: image



Output: Assign image to one
of a fixed set of categories

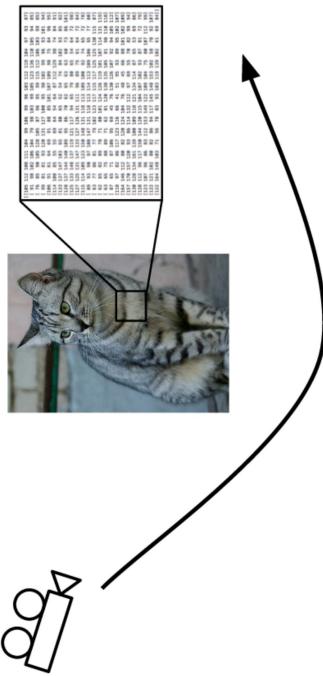
cat bird deer
 dog truck



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

Last Time: Challenges of Recognition

Viewpoint



Illumination



Deformation



Occlusion



This image is [CC0 1.0](#) public domain

This image by [Umberto salvagnini](#) is licensed under [CC BY 2.0](#)

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

Clutter



This image is [CC0 1.0](#) public domain

Intraclass Variation

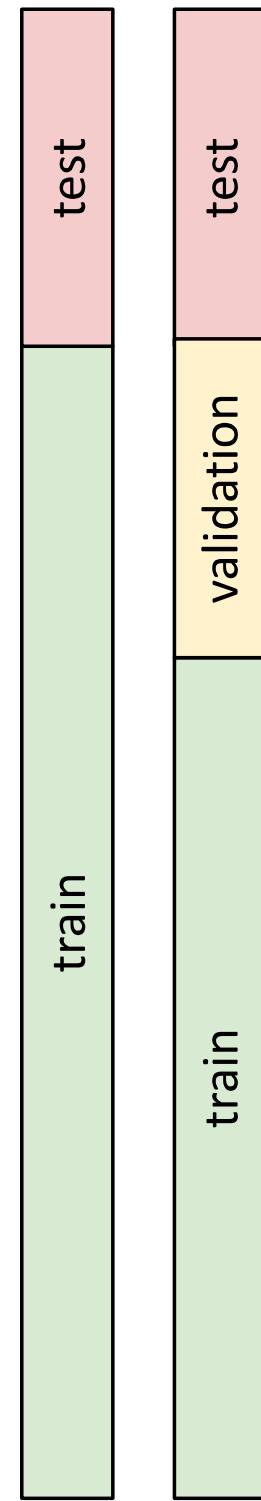
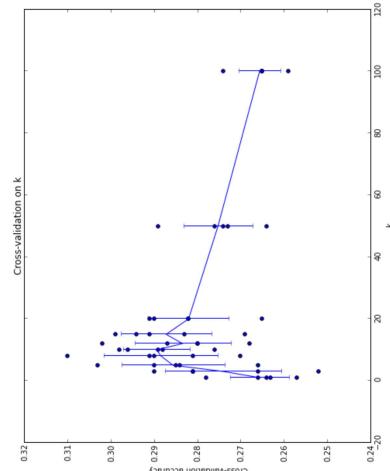
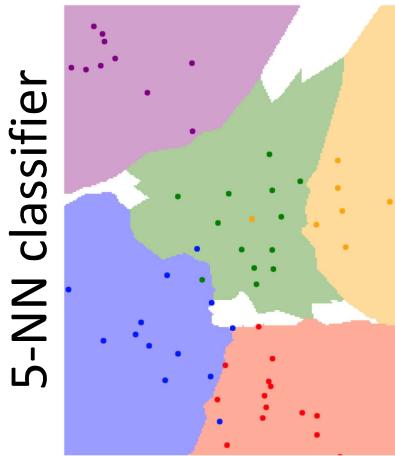


This image is [CC0 1.0](#) public domain

Last time: Data-Drive Approach, kNN



1-NN classifier



Today: Linear Classifiers

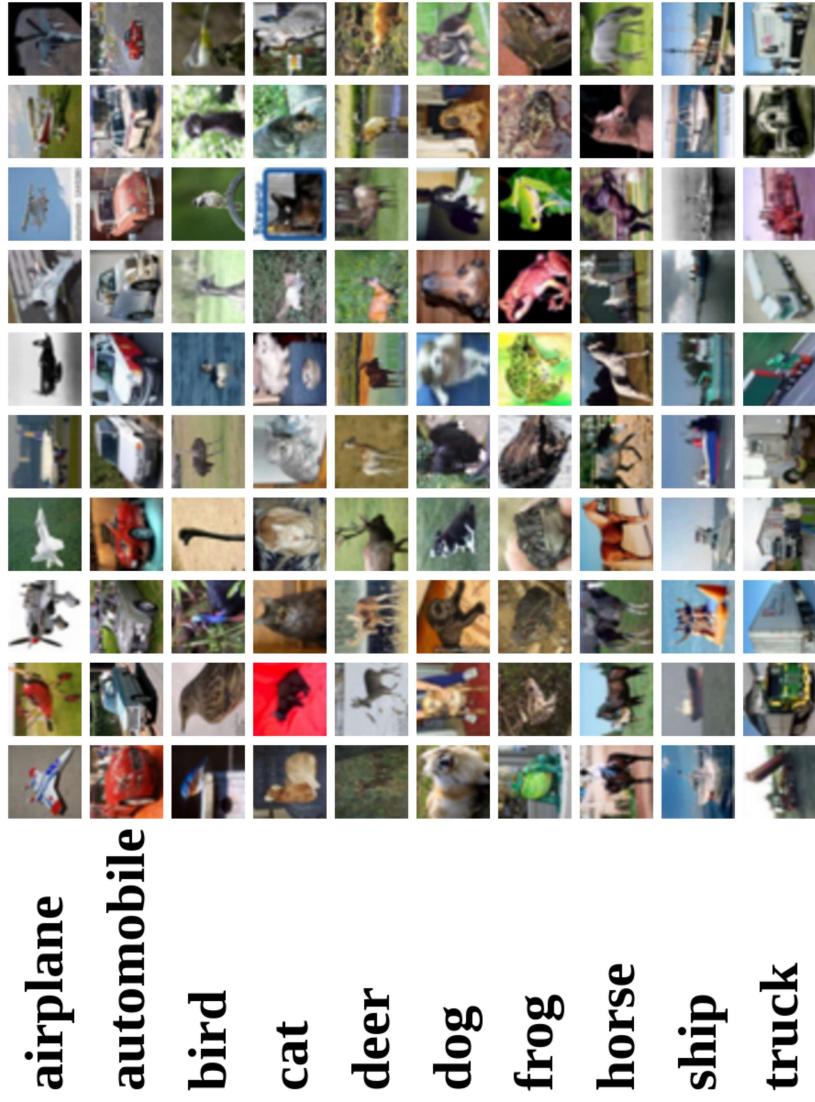
Neural Network



Linear
classifiers

This image is CC0 1.0 public domain

Recall CIFAR10



50,000 training images
each image is 32x32x3

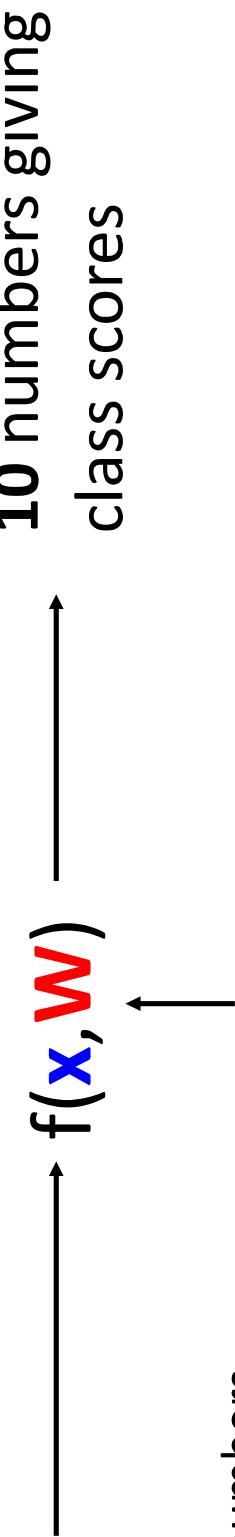
10,000 test images.

Parametric Approach

Image



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



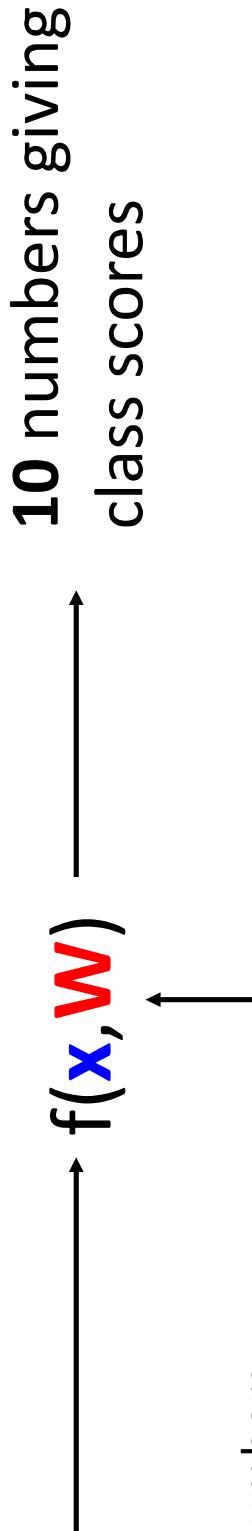
W
parameters
or weights

Parametric Approach: Linear Classifier

Image



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



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

W

parameters
or weights

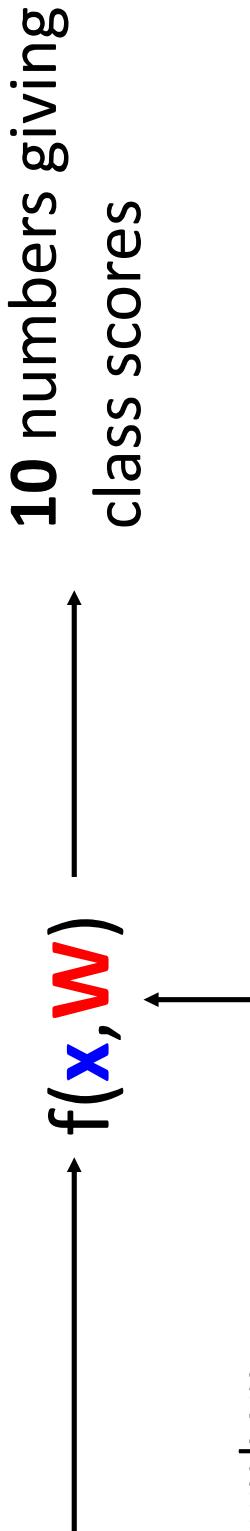
Parametric Approach: Linear Classifier **(3072,)**

Image



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

(10,) (10, 3072)



Array of 32x32x3 numbers
(3072 numbers total)

W

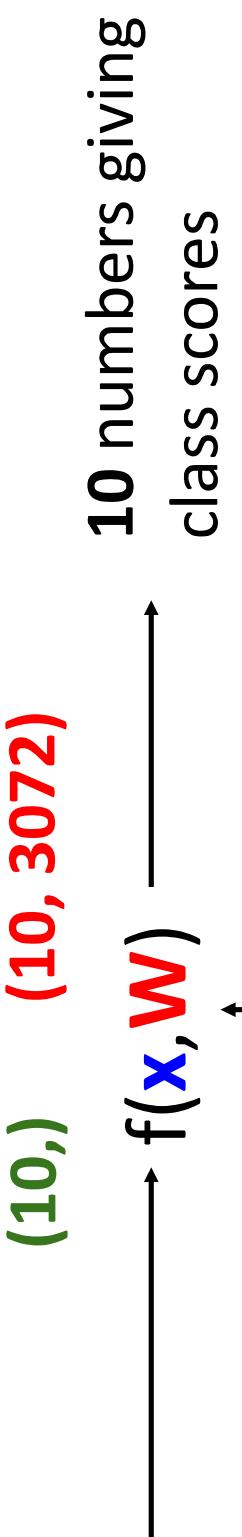
parameters
or weights

Parametric Approach: Linear Classifier

(3072,)

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

Image



Array of 32x32x3 numbers
(3072 numbers total)

W

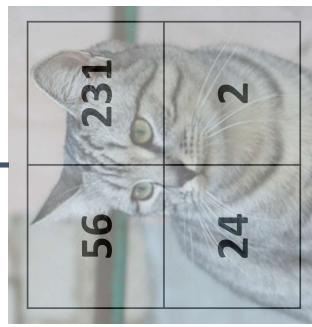
parameters
or weights

Example for 2×2 image, 3 classes (cat/dog/ship)

Stretch pixels into column

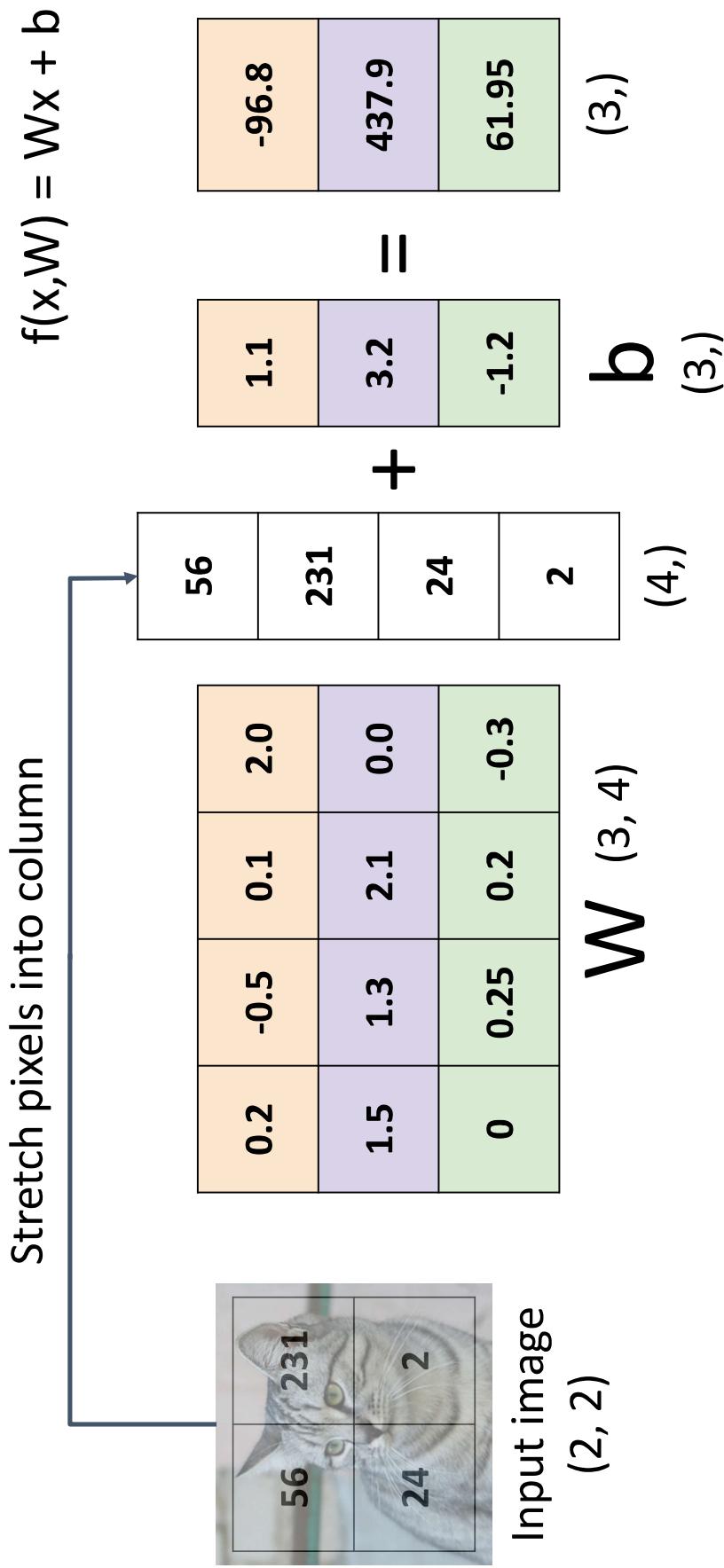
$$f(x, W) = Wx + b$$

56
231
24
2



Input image
(2, 2)
(4,)

Example for 2×2 image, 3 classes (cat/dog/ship)

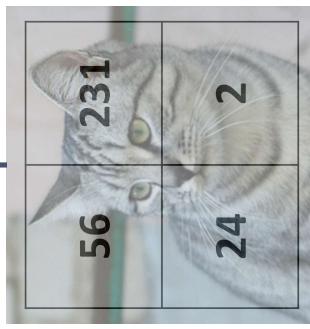


Linear Classifier: Algebraic Viewpoint

Stretch pixels into column

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

Input image $(2, 2)$



$W \quad (3, 4)$

$b \quad (3,)$

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

Diagram illustrating the computation:

The input image is stretched into a column vector:

56
231
24
2

This vector is multiplied by the weight matrix W :

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

The result is a row vector:

1.1
3.2
0

This is then added to the bias b (a column vector of size 3):

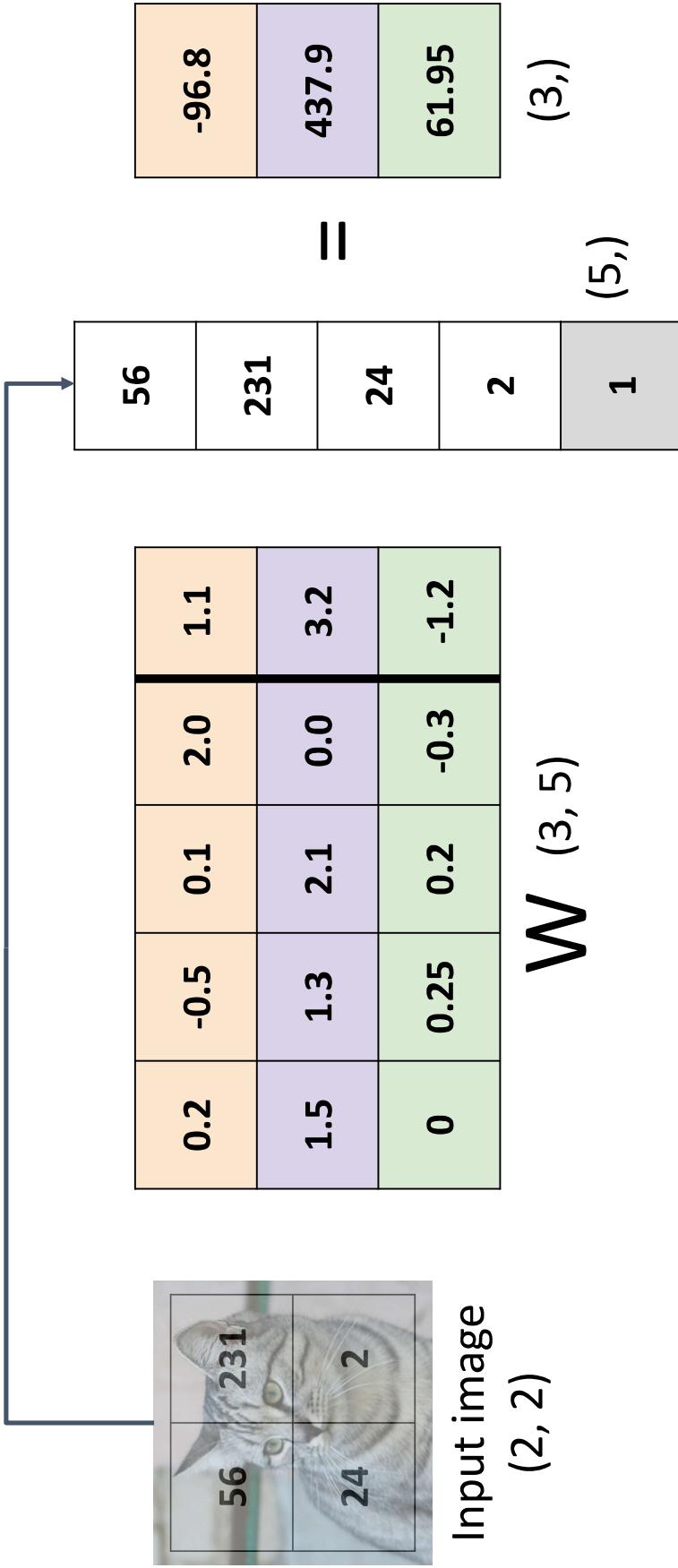
-96.8
437.9
61.95

The final result is:

3.0
3.2
-1.2

Linear Classifier: Bias Trick

Add extra one to data vector;
bias is absorbed into last
column of weight matrix
Stretch pixels into column



Linear Classifier: Predictions are Linear!

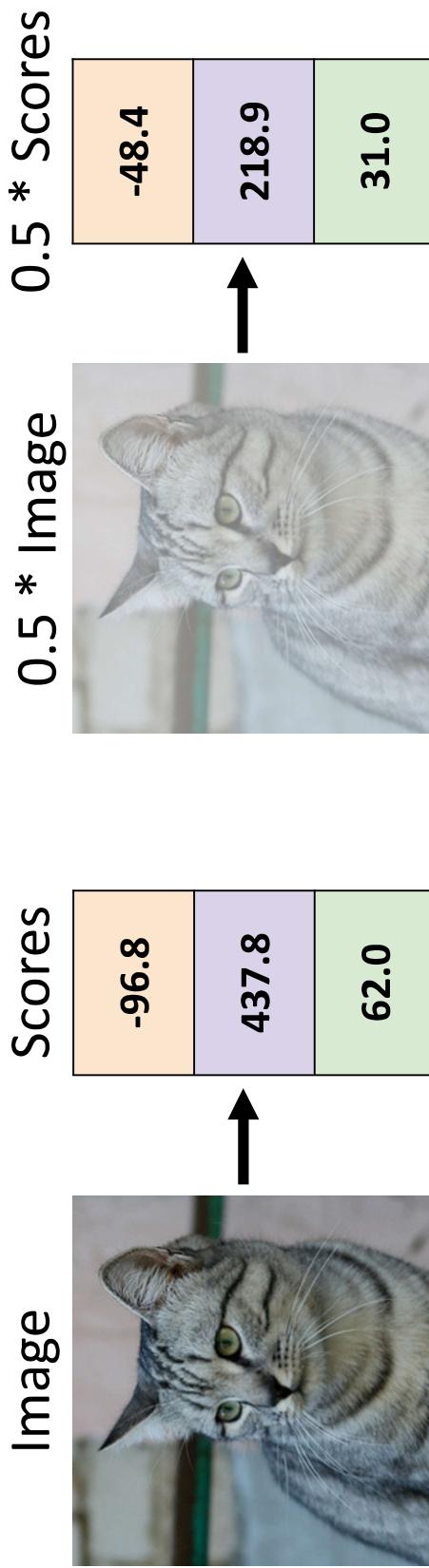
$$f(x, W) = Wx \quad (\text{ignore bias})$$

$$f(cx, W) = W(cx) = c * f(x, W)$$

Linear Classifier: Predictions are Linear!

$$f(x, W) = Wx \quad (\text{ignore bias})$$

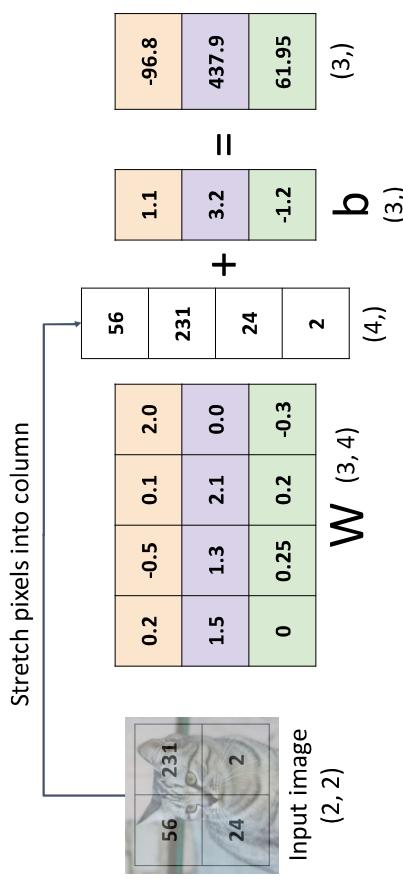
$$f(cx, W) = W(cx) = c * f(x, W)$$



Interpreting a Linear Classifier

Algebraic Viewpoint

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



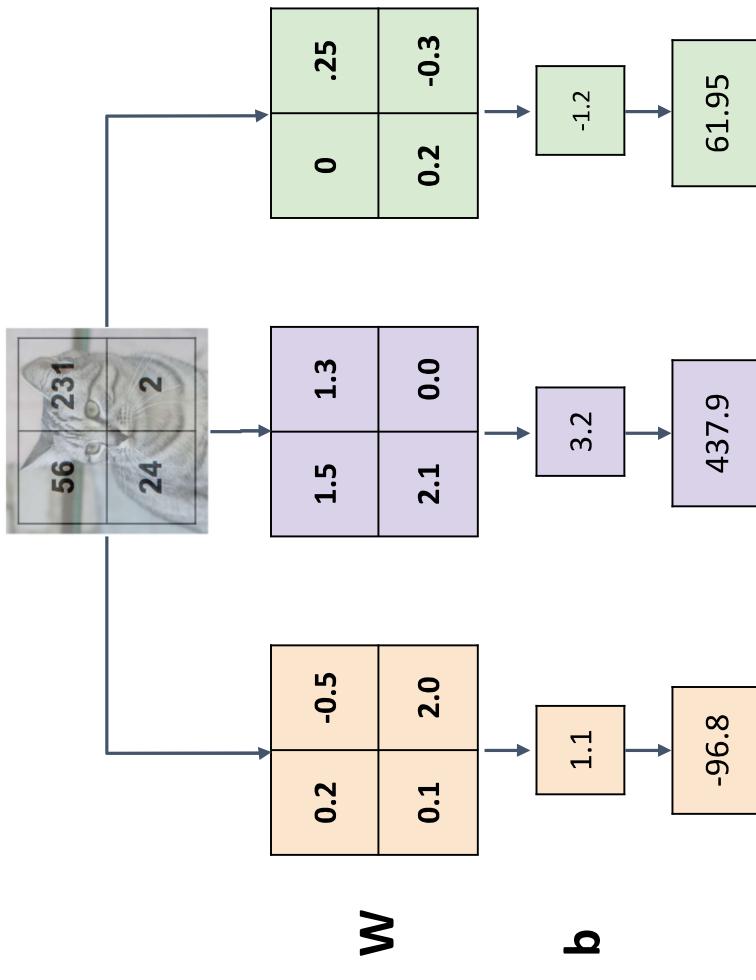
Interpreting a Linear Classifier

Algebraic Viewpoint

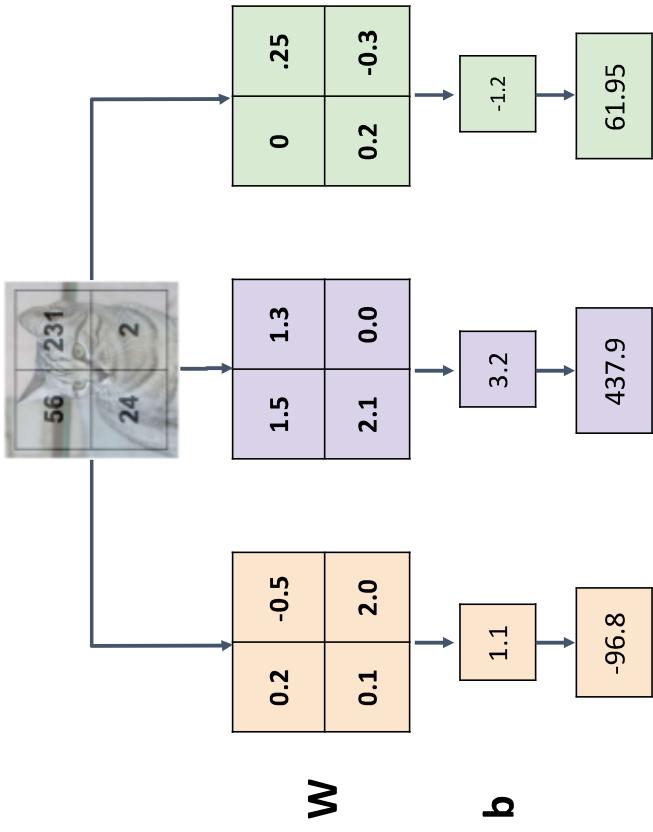
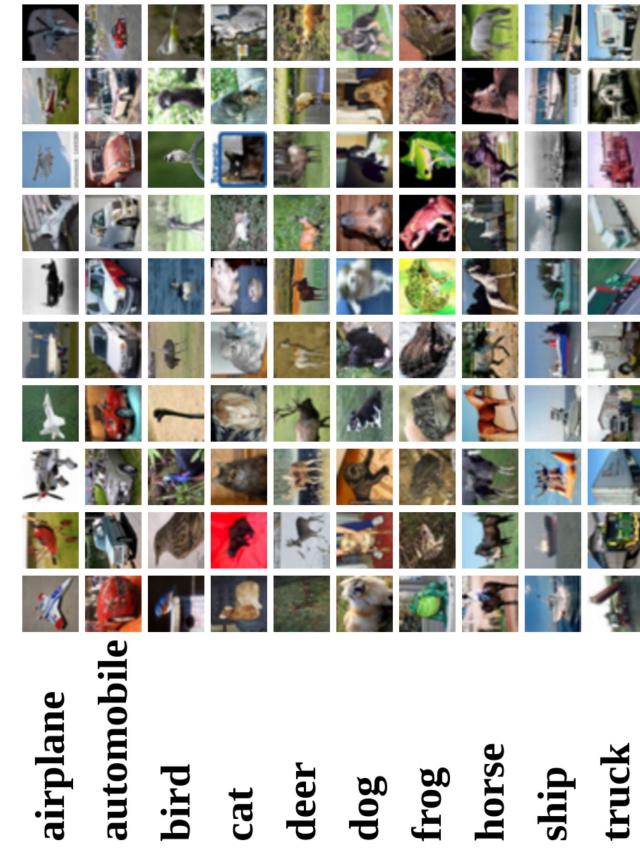
$$f(x, W) = Wx + b$$

Stretch pixels into column

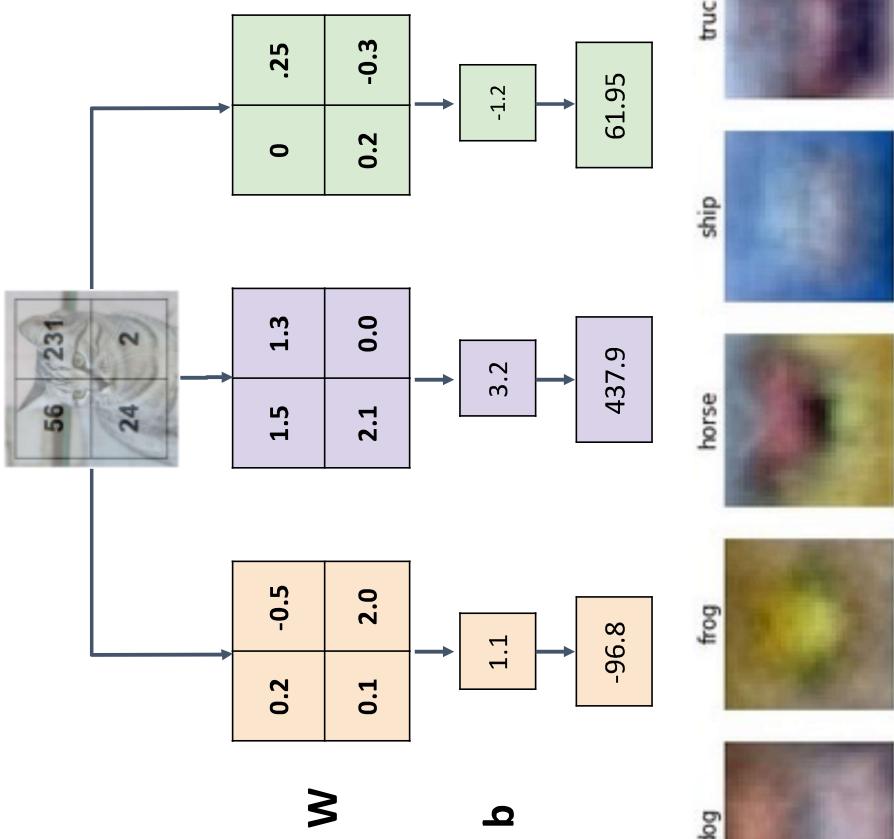
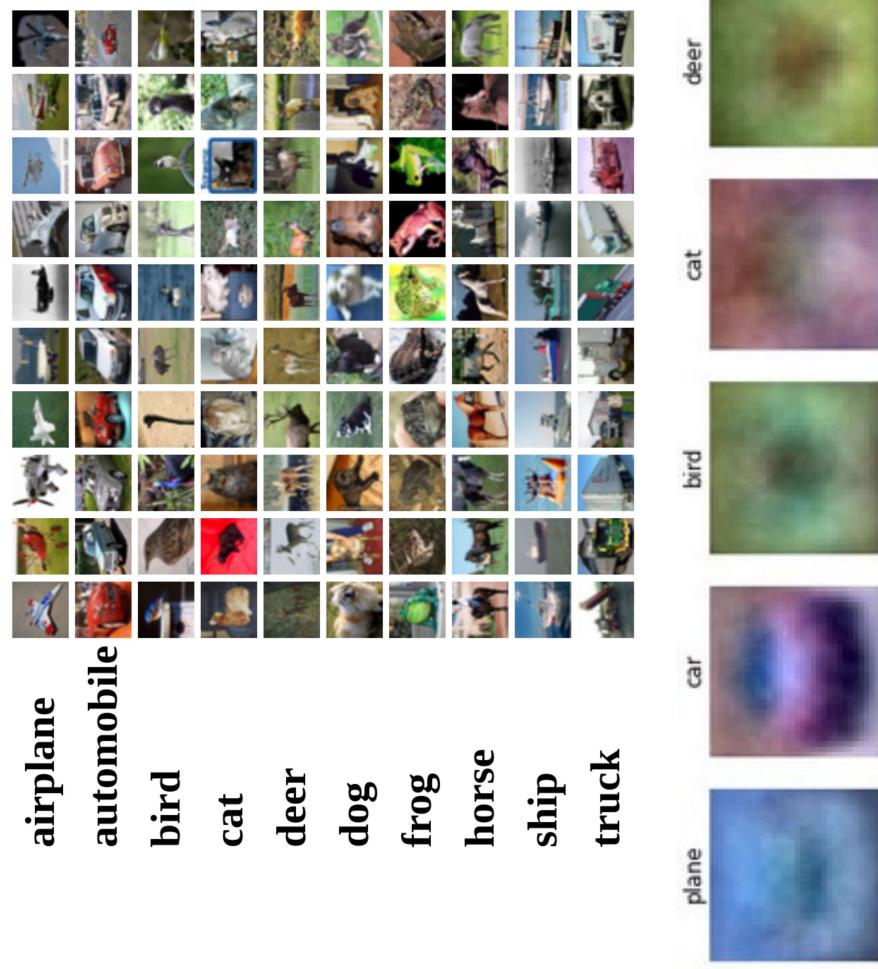
$$\begin{array}{c}
 \text{Input image} \\
 (2, 2)
 \end{array}
 \xrightarrow{\text{Stretch pixels into column}}
 \begin{array}{c}
 W \quad (3, 4) \\
 \begin{matrix}
 56 & 231 & 24 & 2
 \end{matrix}
 \end{array}
 \xrightarrow{\quad}
 \begin{array}{c}
 b \quad (3,) \\
 \begin{matrix}
 1.1 & 2.0 & -0.5 \\
 3.2 & -0.5 & 1.3 \\
 -1.2 & 0.2 & 2.1
 \end{matrix}
 \end{array}
 \xrightarrow{\quad}
 \begin{array}{c}
 -96.8 \\
 437.9 \\
 61.95
 \end{array}
 \xrightarrow{+}
 \begin{array}{c}
 \begin{matrix}
 1.1 & 2.0 & -0.5 \\
 3.2 & -0.5 & 1.3 \\
 -1.2 & 0.2 & 2.1
 \end{matrix} \\
 \begin{matrix}
 1.1 & 2.0 & -0.5 \\
 3.2 & -0.5 & 1.3 \\
 -1.2 & 0.2 & 2.1
 \end{matrix}
 \end{array}
 \xrightarrow{=}
 \begin{array}{c}
 \begin{matrix}
 -96.8 & 437.9 & 61.95
 \end{matrix}
 \end{array}$$



Interpreting an Linear Classifier

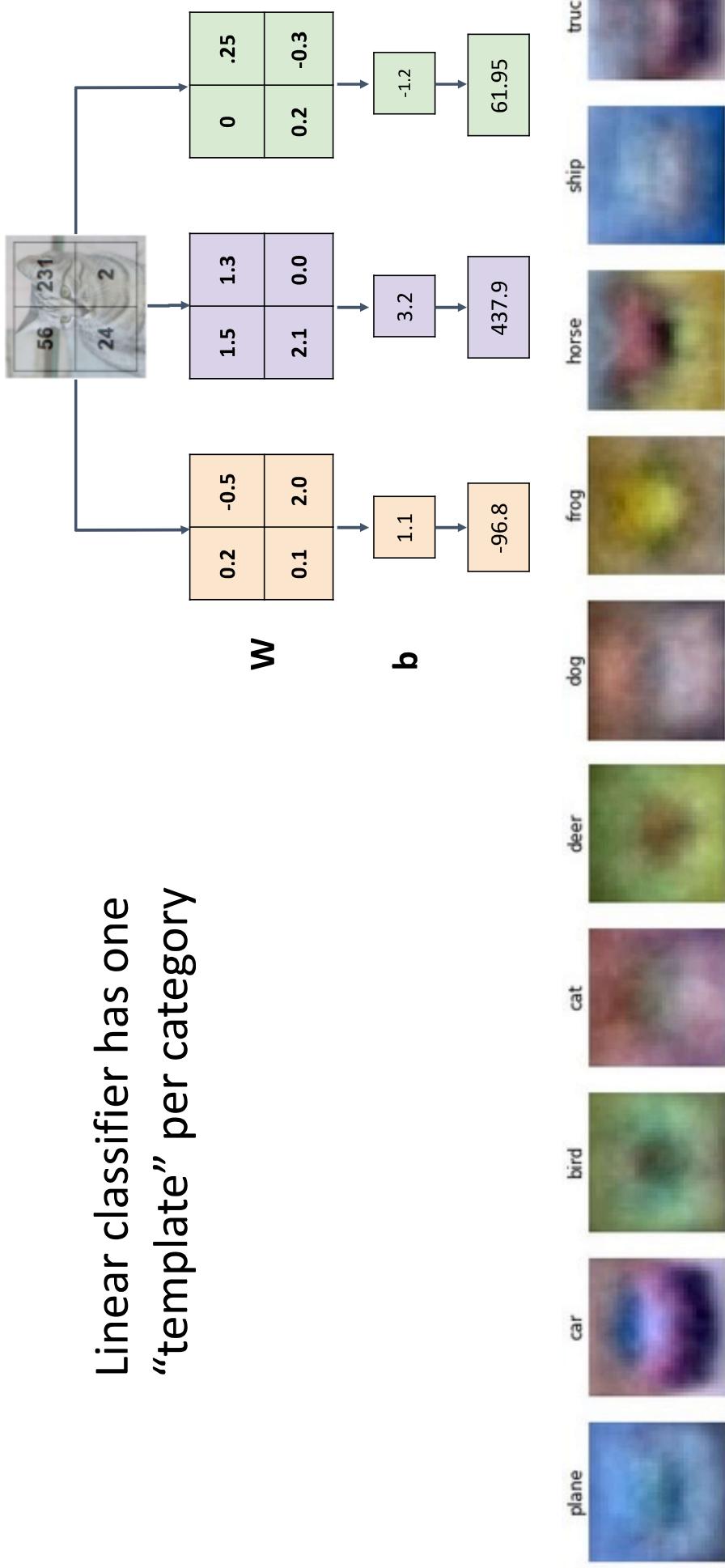


Interpreting an Linear Classifier: Visual Viewpoint



Interpreting an Linear Classifier: Visual Viewpoint

Linear classifier has one
“template” per category

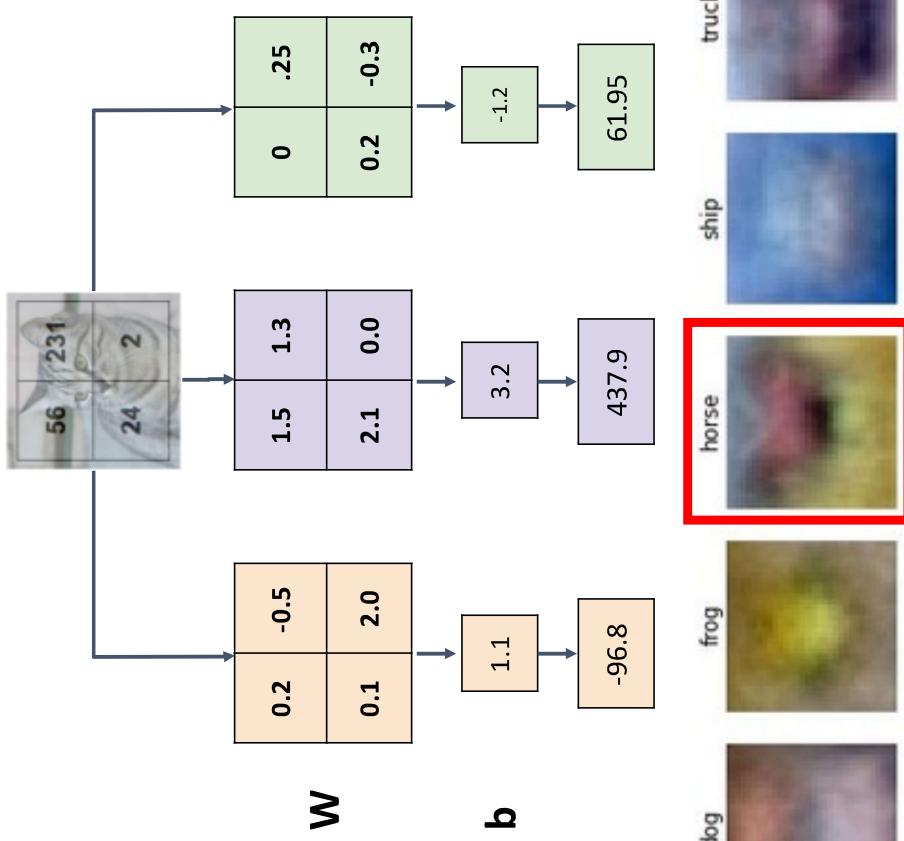


Interpreting an Linear Classifier: Visual Viewpoint

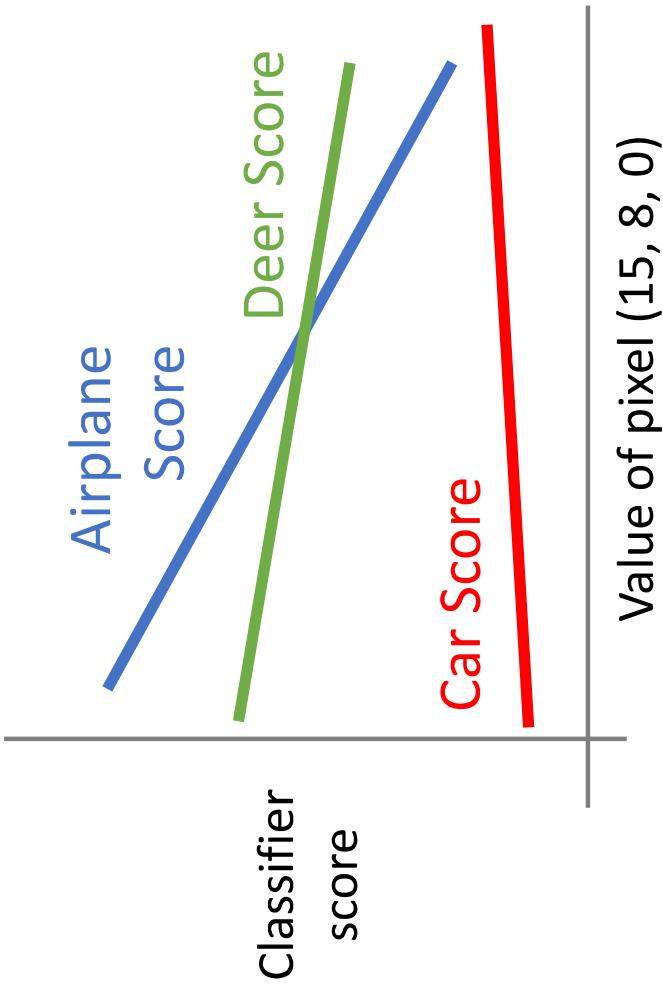
Linear classifier has one
“template” per category

A single template cannot capture
multiple modes of the data

e.g. horse template has 2 heads!



Interpreting a Linear Classifier: Geometric Viewpoint

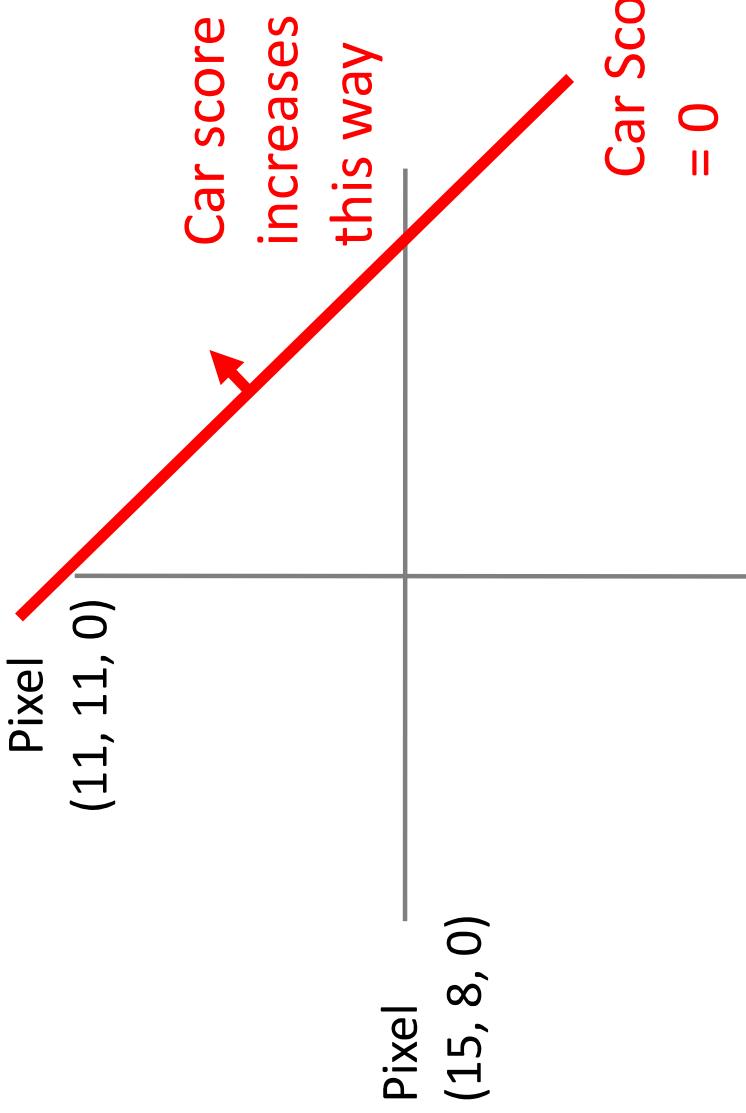


$$f(x, W) = Wx + b$$



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

Interpreting a Linear Classifier: Geometric Viewpoint

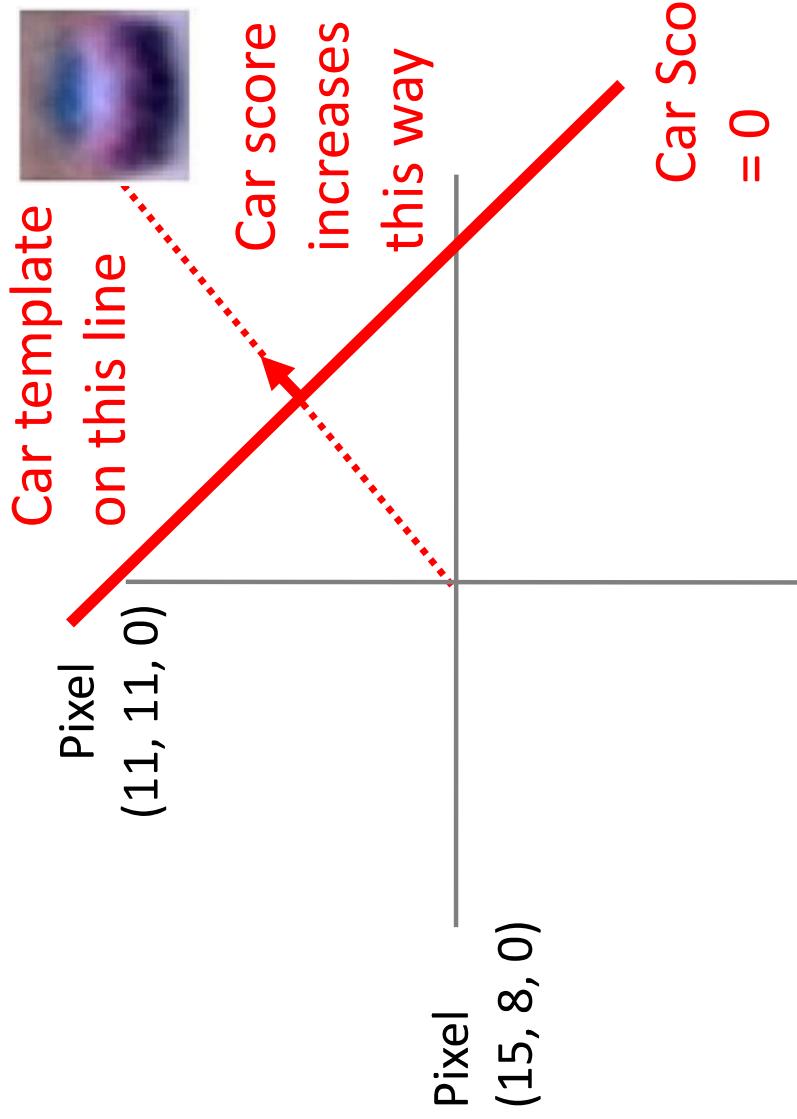


$$f(x, W) = Wx + b$$



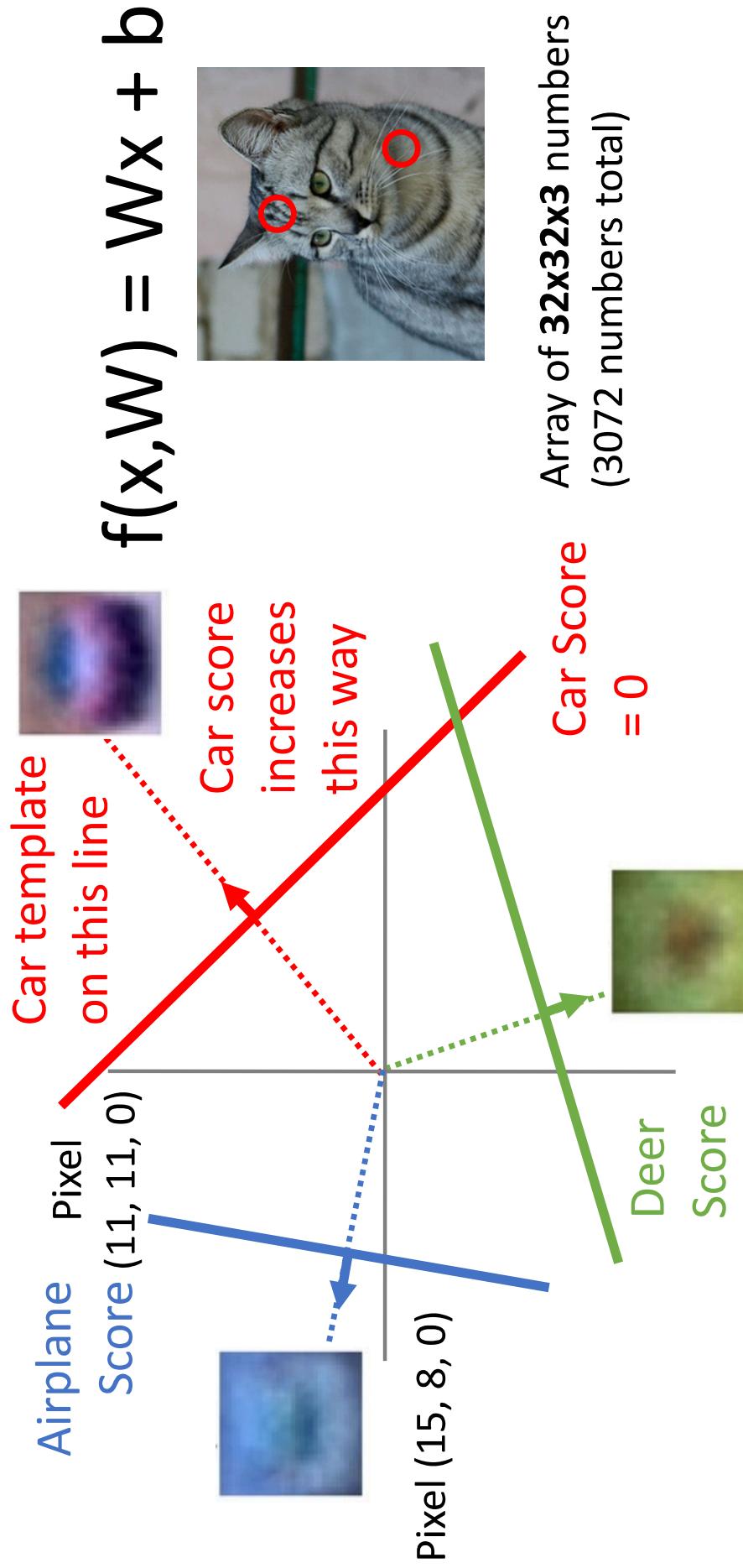
Array of **32x32x3** numbers
(3072 numbers total)

Interpreting a Linear Classifier: Geometric Viewpoint

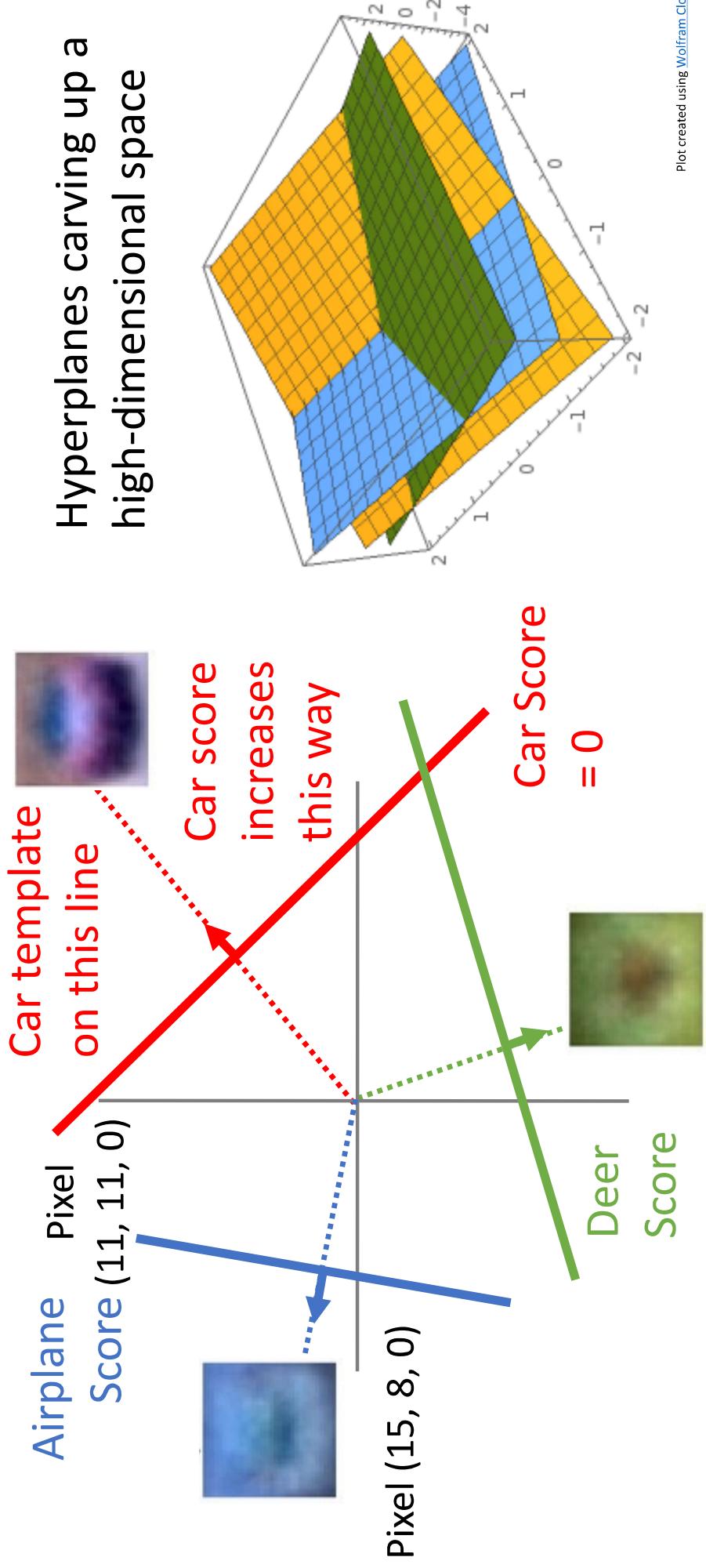


Array of **32x32x3** numbers
(3072 numbers total)

Interpreting a Linear Classifier: Geometric Viewpoint



Interpreting a Linear Classifier: Geometric Viewpoint



Hard Cases for a Linear Classifier

Class 1:

First and third quadrants
 $1 \leq \text{L2 norm} \leq 2$

Class 1:

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

Class 1:

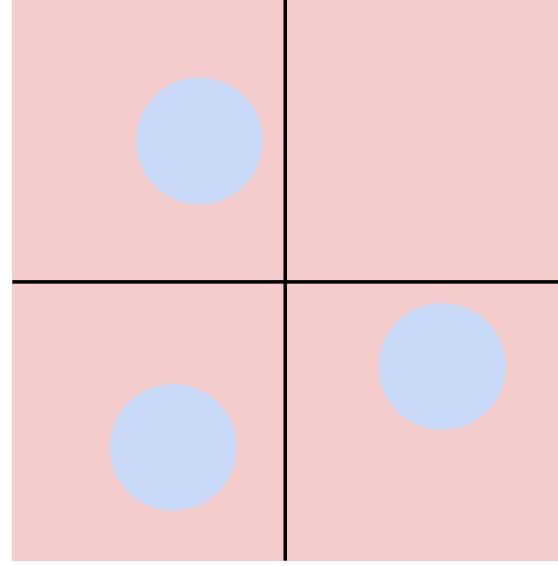
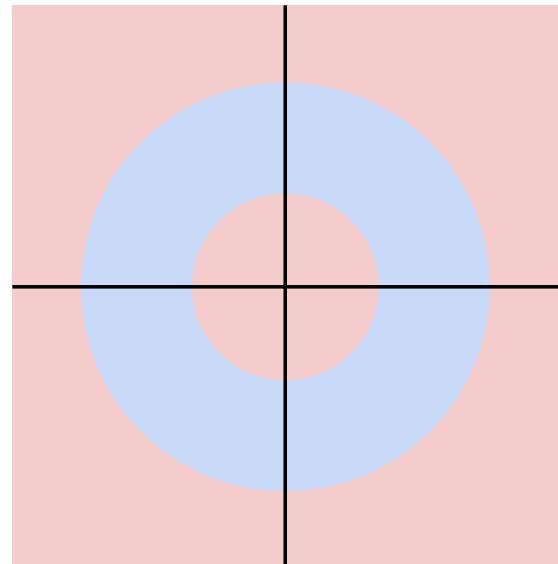
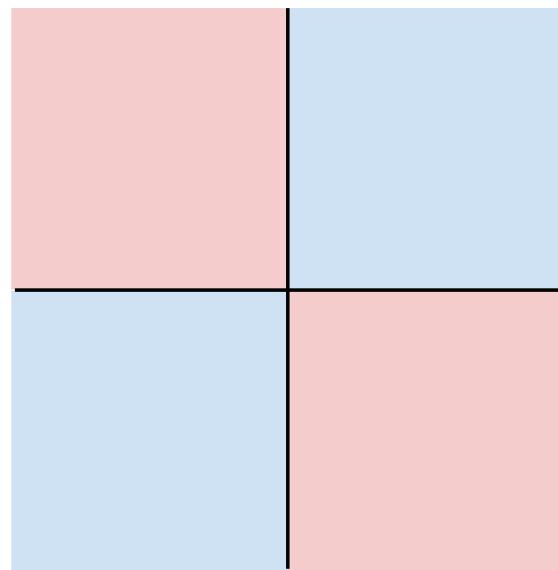
Three modes

Class 2:

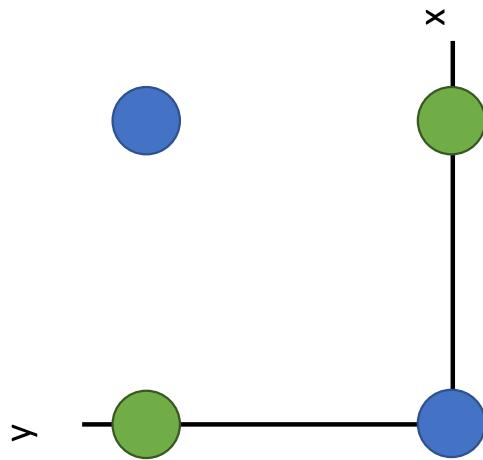
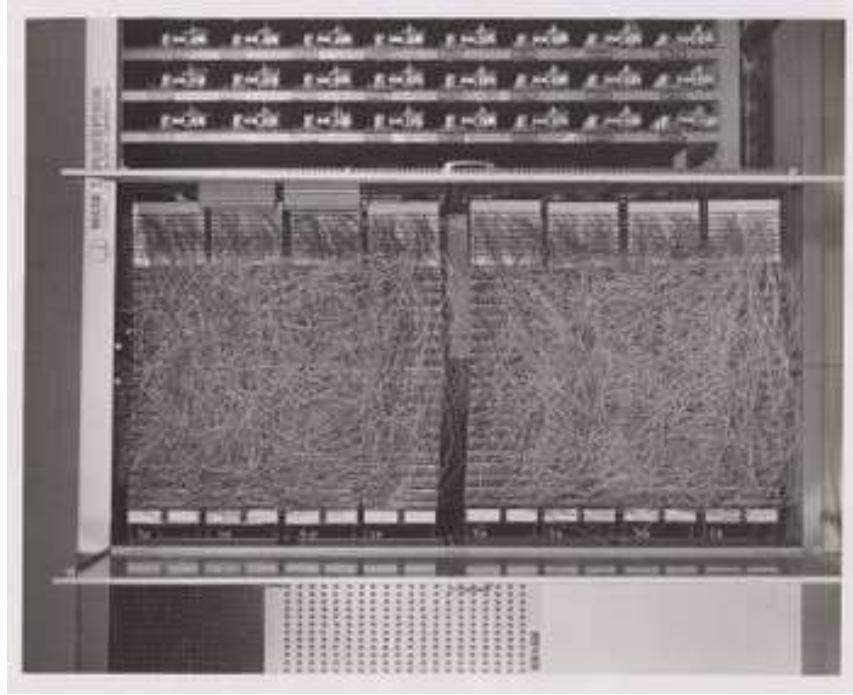
Second and fourth quadrants
Everything else

Class 2:

Everything else



Recall: Perceptron couldn't learn XOR



X	Y	F(x,y)
0	0	0
0	1	1
1	0	1
1	1	0

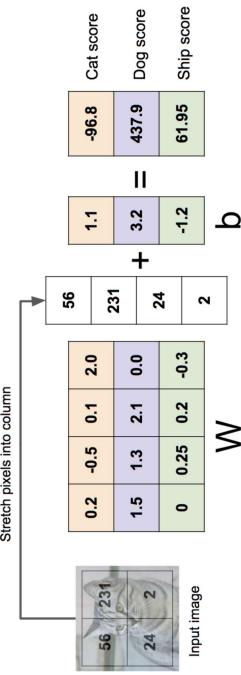
Linear Classifier: Three Viewpoints

Algebraic Viewpoint

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

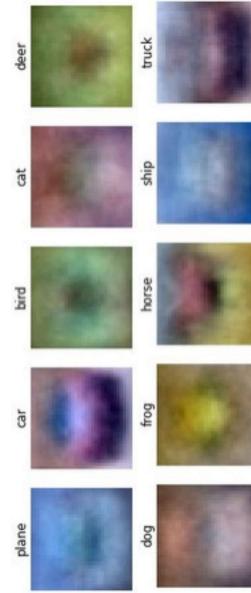
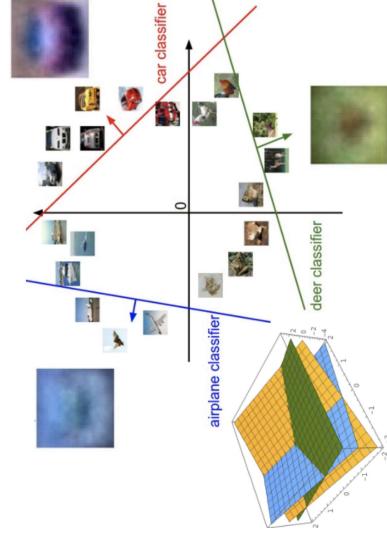
Visual Viewpoint

One template
per class



Geometric Viewpoint

Hyperplanes
cutting up space



So Far: Defined a linear score function

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



Given a W , we can
compute class scores
for an image x .

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

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

Choosing a good W

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

TODO:



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

1. Use a loss function to quantify how good a value of W is
2. Find a W that minimizes the loss function (optimization)

LOSS Function

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

Low loss = good classifier

High loss = bad classifier

(Also called: **objective function**;
cost function)

LOSS Function

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

Low loss = good classifier

High loss = bad classifier

(Also called: **objective function; cost function**)

Negative loss function sometimes called **reward function, profit function, utility function, fitness function**, etc

LOSS Function

Given a dataset of examples

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

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

Low loss = good classifier
High loss = bad classifier

(Also called: **objective function**;
cost function)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc

LOSS Function

Given a dataset of examples

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

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

Low loss = good classifier

High loss = bad classifier

Where x_i is image and
 y_i is (integer) label

Loss for a single example is

$$L_i(f(x_i, W), y_i)$$

(Also called: **objective function**;
cost function)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc

LOSS Function

Given a dataset of examples

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

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

Low loss = good classifier

High loss = bad classifier

Where x_i is image and
 y_i is (integer) label

Loss for a single example is

$$L_i(f(x_i, W), y_i)$$

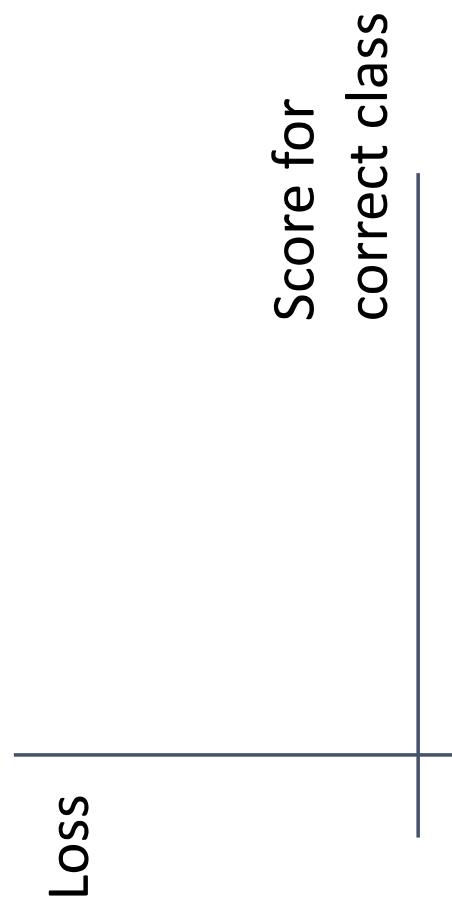
Loss for the dataset is average of per-example losses:

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

Negative loss function sometimes called **reward function, profit function, utility function, fitness function**, etc

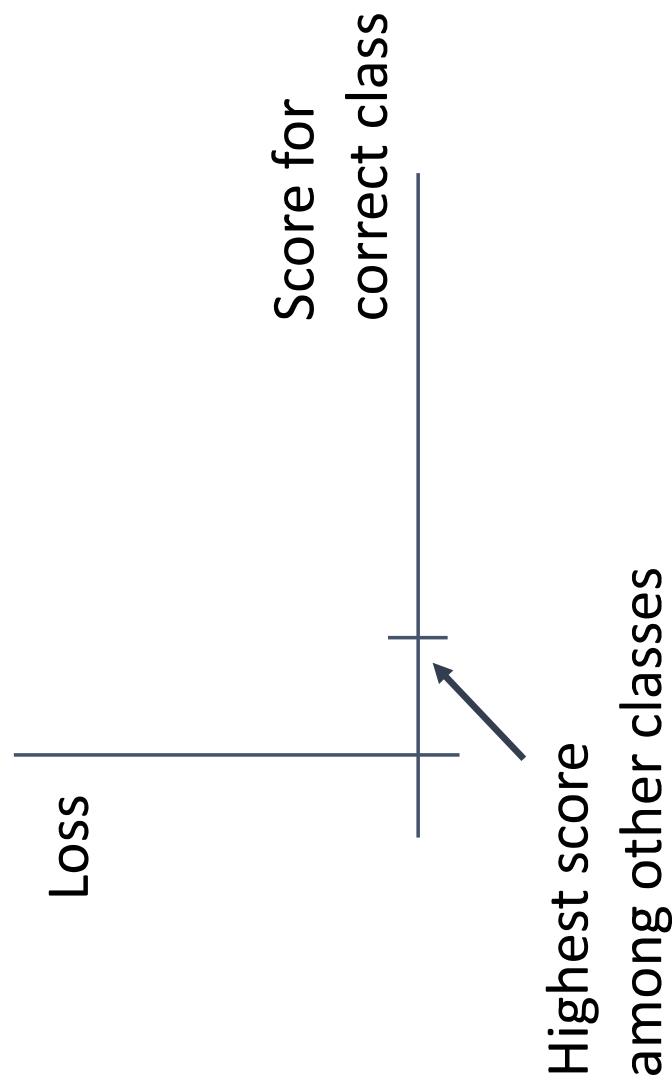
Multiclass SVM Loss

”The score of the correct class should be higher than all the other scores”



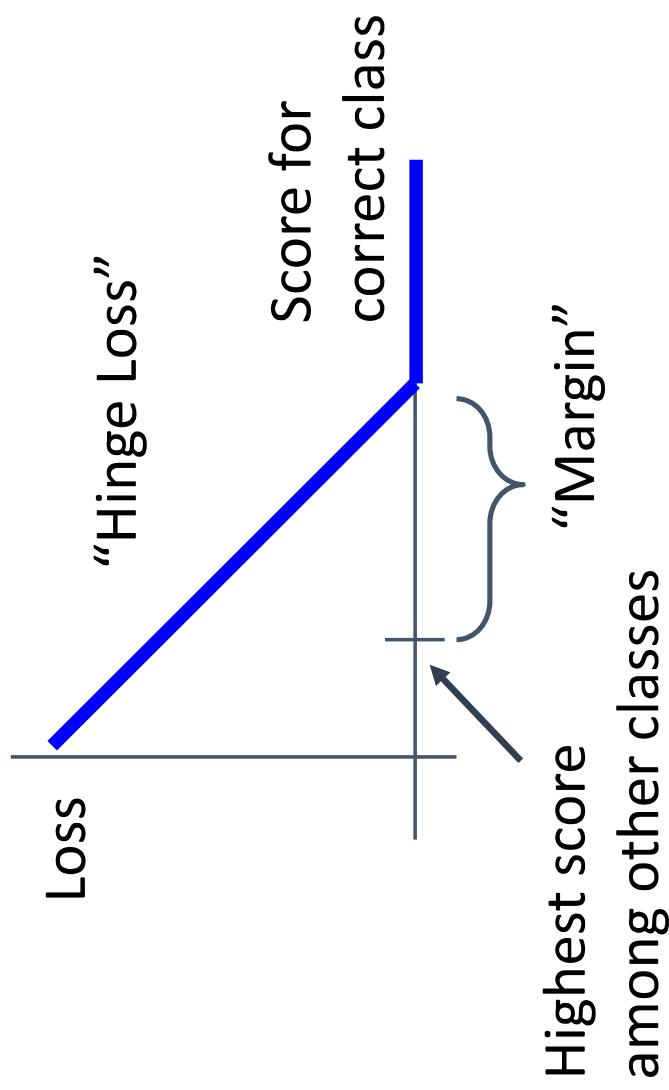
Multiclass SVM Loss

”The score of the correct class should be higher than all the other scores”



Multiclass SVM Loss

”The score of the correct class should be higher than all the other scores”



Multiclass SVM Loss

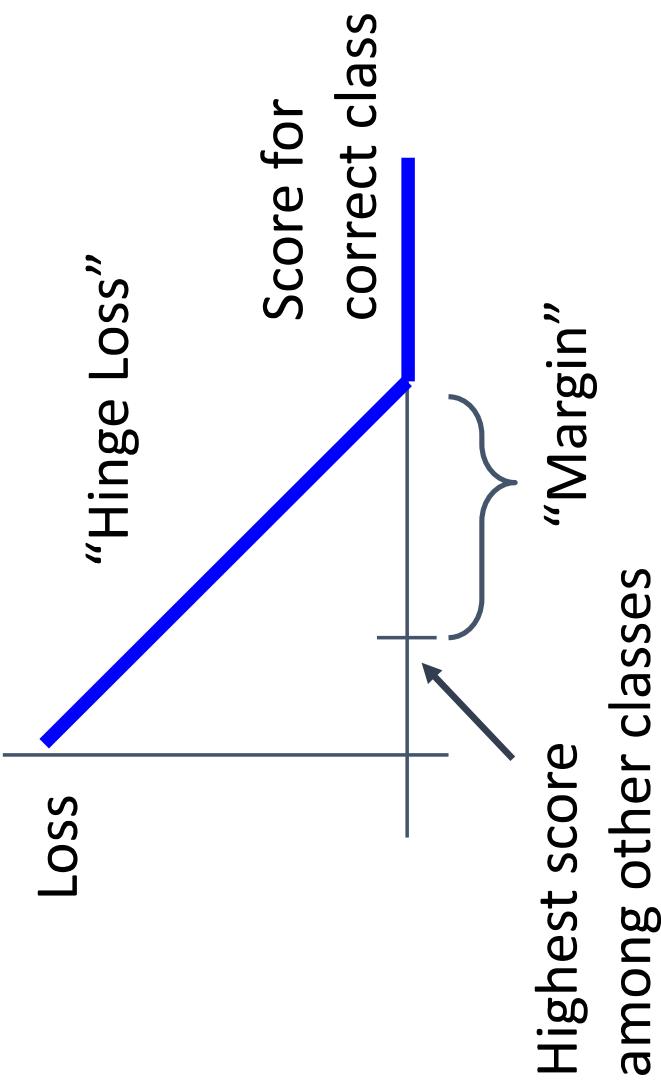
”The score of the correct class should be higher than all the other scores”

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)

Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

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



Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)



Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:
 $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$

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

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)

Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

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



cat	1.3	2.2
car	4.9	2.5
frog	-1.7	-3.1
Loss		2.9

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)



Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

cat	3.2	2.2	$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$
car	5.1	2.5	$= \max(0, 1.3 - 4.9 + 1)$
frog	-1.7	-3.1	$+ \max(0, 2.0 - 4.9 + 1)$
Loss	2.9	2.0	$= \max(0, -2.6) + \max(0, -1.9)$
			$= 0 + 0$
			$= 0$

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)

Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

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



	2.2	2.5	-3.1	12.9
--	-----	-----	------	------

cat	3.2	1.3	4.9	0
car	5.1	4.9	2.0	
frog	-1.7	2.0		
Loss	2.9			

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)



Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

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

Loss over the dataset is:

cat	3.2	1.3	2.2	2.0	-3.1	$L = (2.9 + 0.0 + 12.9) / 3$	12.9
car	5.1	4.9	2.5				
frog	-1.7						
Loss	2.9						

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)

Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

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



cat	3.2	1.3	2.2	2.5	-3.1	12.9
car	5.1	4.9				
frog	-1.7	2.0				
Loss	2.9	0				

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

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)



Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

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

cat	3.2	1.3	2.2	2.5	-3.1	12.9
car	5.1	4.9	2.0			
frog	-1.7					

Q2: What are the min
and max possible loss?

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)

Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

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



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Loss	2.9	0	12.9

Q3: If all the scores were random, what loss would we expect?

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)



Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

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

cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Loss	2.9	0	12.9

Q4: What would happen
if the sum were over all
classes? (including $i = y_i$)

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)



Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

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

cat	3.2	1.3	2.2	2.5	-3.1	12.9
car	5.1	4.9				
frog	-1.7	2.0				
Loss	2.9	0				

Q5: What if the loss used
a mean instead of a sum?

Multiclass SVM Loss

Given an example (x_i, y_i)
 $(x_i$ is image, y_i is label)



Let $s = f(x_i, W)$ be scores

Then the SVM loss has the form:

cat	3.2	1.3	2.2	$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$
car	5.1	4.9	2.5	
frog	-1.7	2.0	-3.1	
Loss	2.9	0	12.9	

Q6: What if we used
this loss instead?

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

Multiclass SVM LOSS

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

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

Q: Suppose we found some W with $L = 0$. Is it unique?

Multiclass SVM LOSS

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

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

Q: Suppose we found some W with $L = 0$. Is it unique?

No! $2W$ is also has $L = 0$!

Multiclass SVM Loss

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

Original W:

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



cat	3.2	2.2	2.5	
car	5.1	4.9	-3.1	
frog	-1.7	2.0	12.9	
Loss	2.9	0		

Using 2W instead:

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

Multiclass SVM Loss

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



How should we choose between
 W and $2W$ if they both perform
the same on the training data?

cat	3.2	2.2	2.5	-3.1	12.9
car	5.1	4.9	2.0	0	
frog	-1.7				
Loss	2.9				

Regularization: Beyond Training Error

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

Data loss: Model predictions
should match training data

Regularization: Beyond Training Error

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too well* on training data

Regularization: Beyond Training Error

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

λ = regularization strength
(hyperparameter)



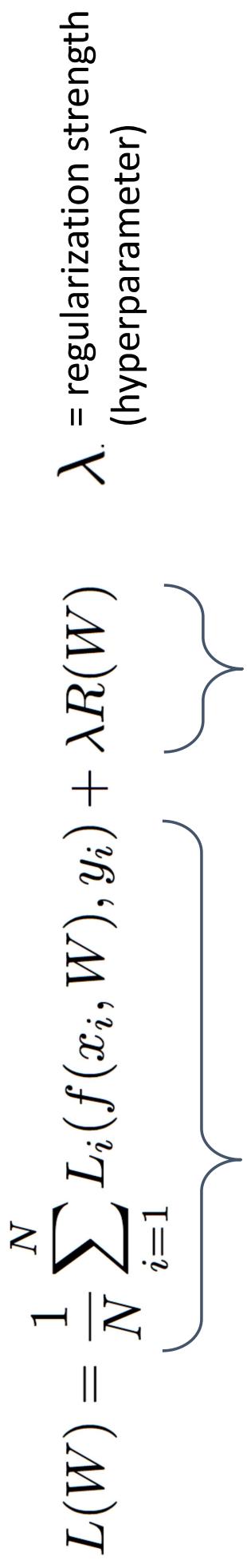
Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too well* on training data

Regularization: Beyond Training Error

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

λ = regularization strength
(hyperparameter)



Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too well* on training data

Simple examples

L2 regularization:

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

L1 regularization:

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

Elastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$ Cutout, Mixup, Stochastic depth, etc...

More complex:

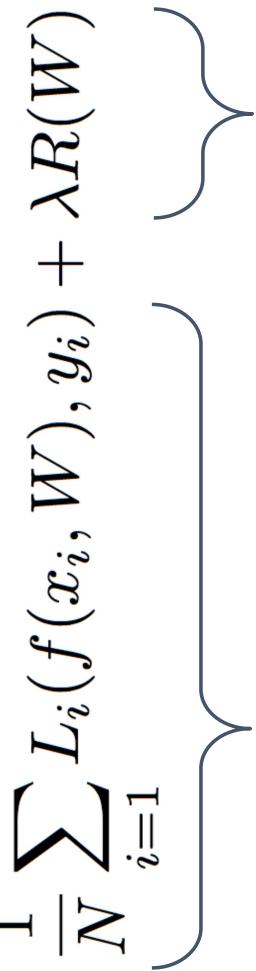
Dropout

Batch normalization

Regularization: Beyond Training Error

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

λ = regularization strength
(hyperparameter)



Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too well* on training data

Purpose of Regularization:

- Express preferences in among models beyond "minimize training error"
- Avoid **overfitting**: Prefer simple models that generalize better
- Improve optimization by adding curvature

Regularization: Expressing Preferences

$$\begin{aligned}x &= [1, 1, 1, 1] \\w_1 &= [1, 0, 0, 0]\end{aligned}$$

L2 Regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

$$w_1^T x = w_2^T x = 1$$

Regularization: Expressing Preferences

$$\begin{aligned}x &= [1, 1, 1, 1] \\w_1 &= [1, 0, 0, 0]\end{aligned}$$

L2 Regularization

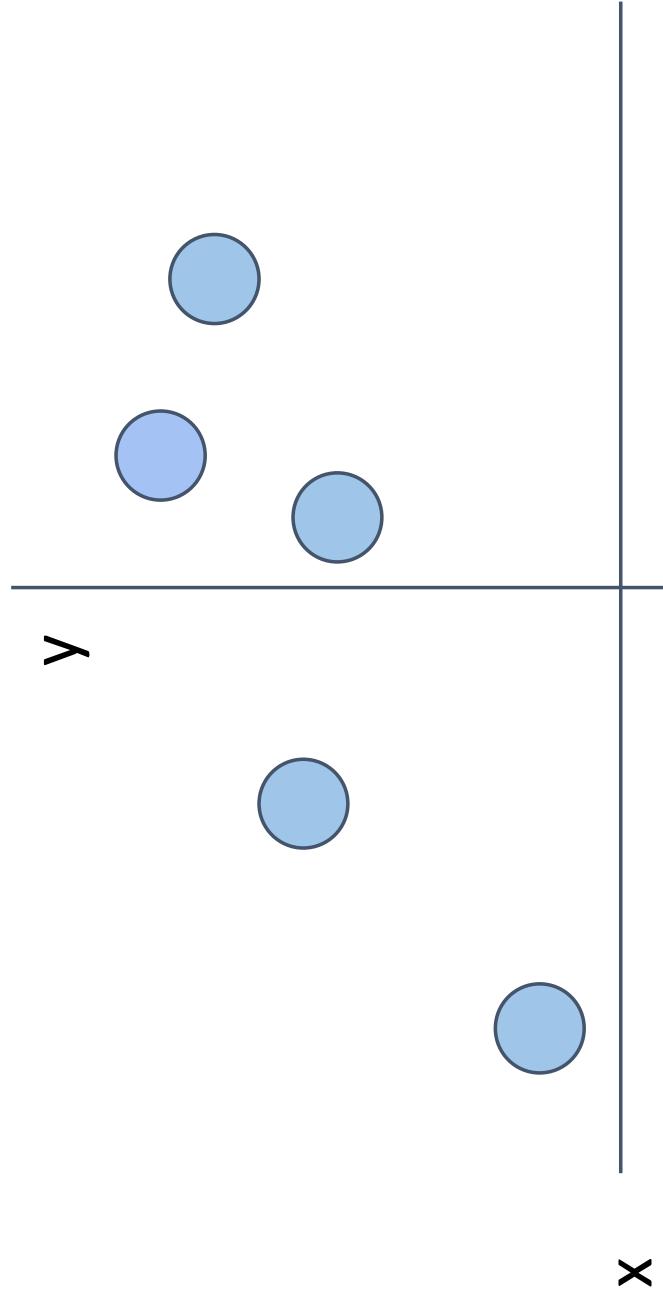
$$R(W) = \sum_k \sum_l W_{k,l}^2$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

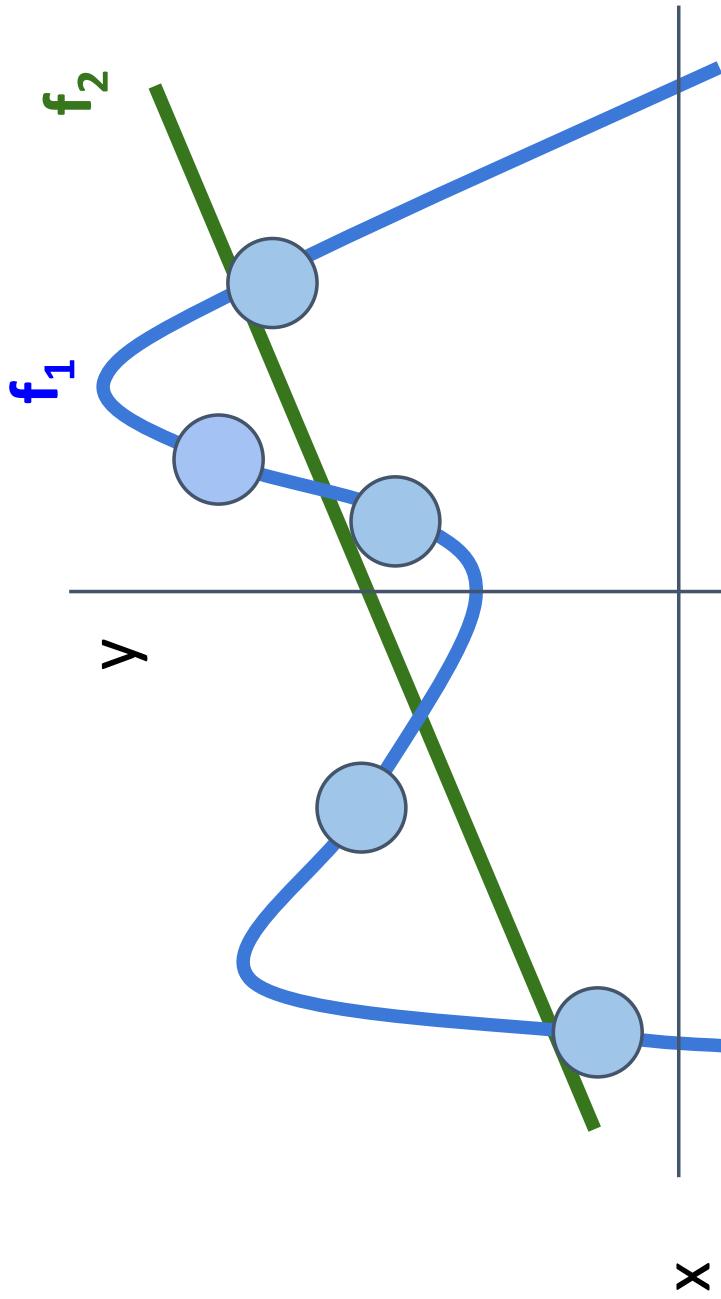
L2 regularization likes to
“spread out” the weights

$$w_1^T x = w_2^T x = 1$$

Regularization: Prefer Simpler Models

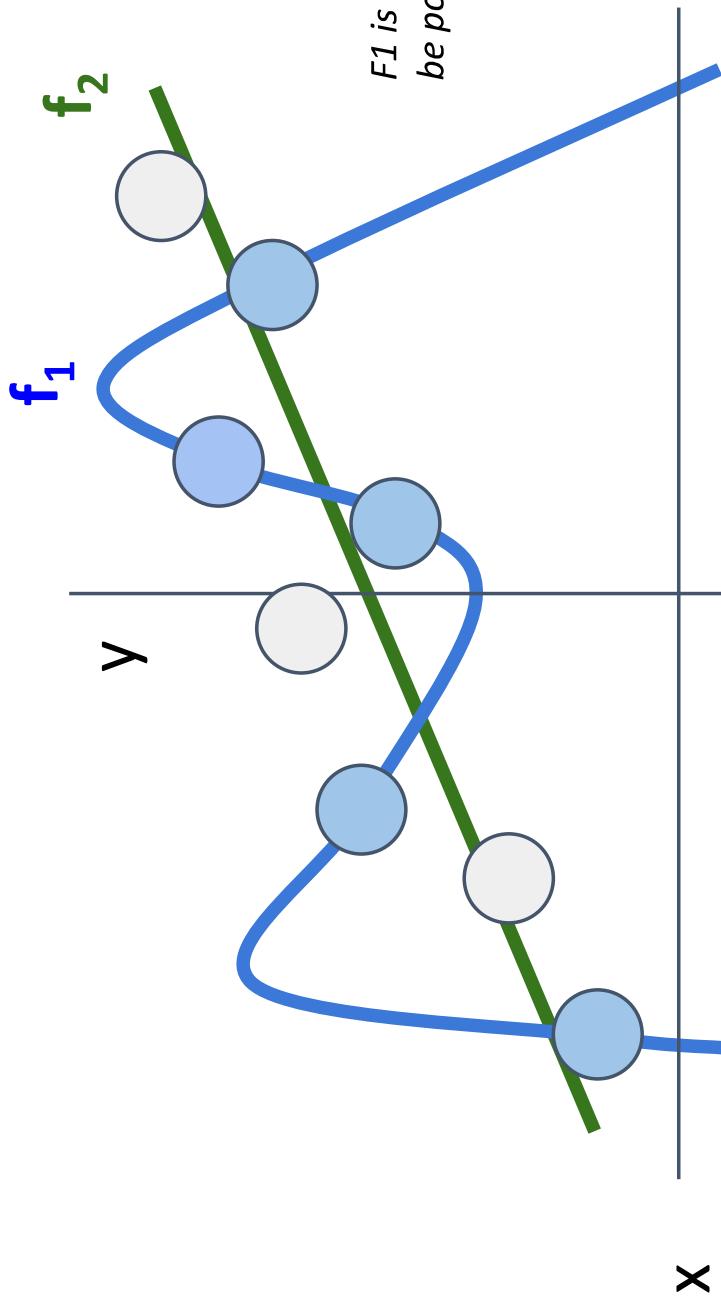


Regularization: Prefer Simpler Models



The model f_1 fits the training data perfectly
The model f_2 has training error, but is simpler

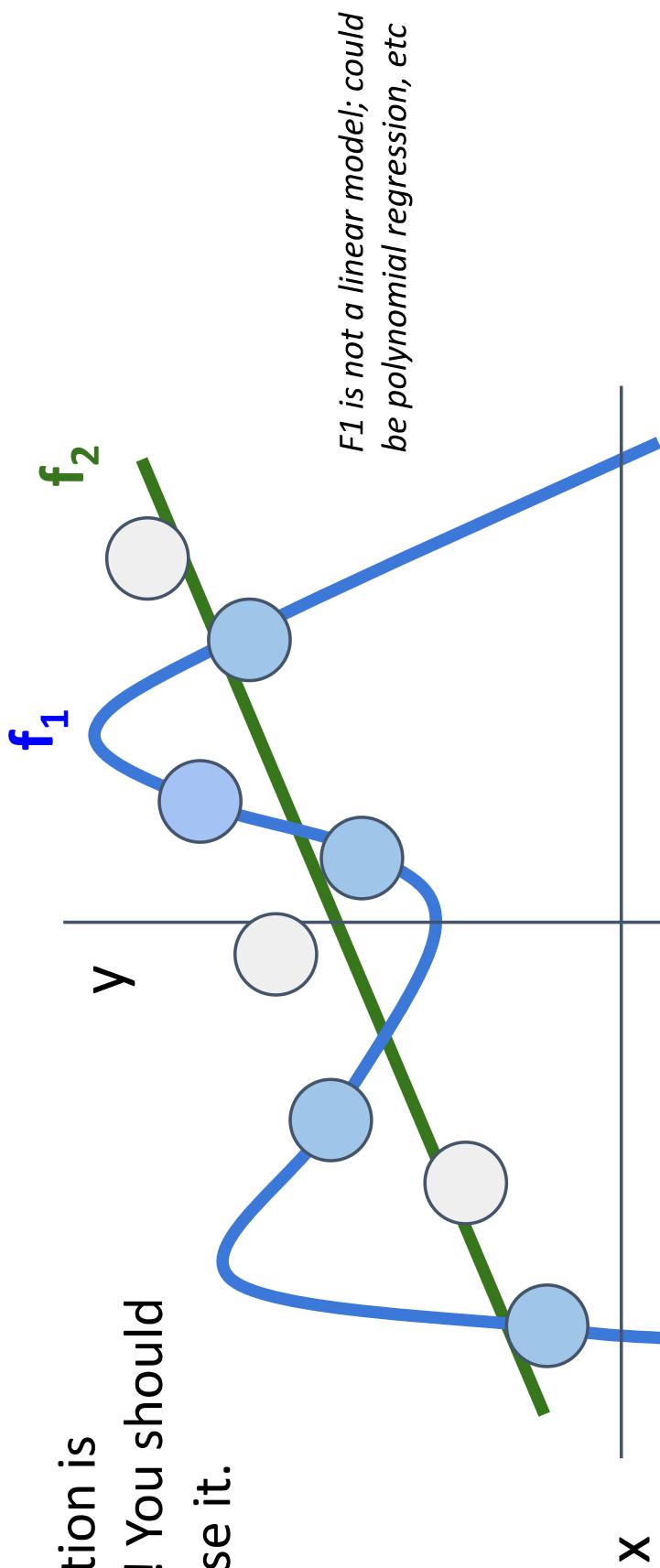
Regularization: Prefer Simpler Models



Regularization pushes against fitting the data
too well so we don't fit noise in the data

Regularization: Prefer Simpler Models

Regularization is important! You should (usually) use it.



Cross-Entropy Loss (Multinomial Logistic Regression)

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



cat **3.2**

car **5.1**

frog **-1.7**

Cross-Entropy Loss (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

Cross-Entropy Loss (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

Unnormalized log-
probabilities / logits

Cross-Entropy Loss (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

$$\boxed{24.5} \quad \boxed{164.0} \quad \boxed{0.18}$$

\exp \rightarrow

$$\boxed{3.2} \quad \boxed{5.1} \quad \boxed{-1.7}$$

Unnormalized log-
probabilities / logits
*unnormalized
probabilities*

Cross-Entropy Loss (Multinomial Logistic Regression)

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



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

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

Probabilities
must be ≥ 0

Probabilities
must sum to 1

0.13
0.87
0.00

24.5
164.0
0.18

3.2
5.1
-1.7

cat
car
frog

normalize
→

unnormalized
probabilities

Unnormalized log-
probabilities / logits

probabilities

Cross-Entropy Loss (Multinomial Logistic Regression)

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



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

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

Probabilities
must be ≥ 0

Probabilities
must sum to 1

Softmax
function

cat
car
frog

$$\begin{array}{r} 3.2 \\ 5.1 \\ -1.7 \end{array}$$

$$\begin{array}{r} 164.0 \\ 0.18 \end{array}$$

$$\begin{array}{r} 24.5 \\ 0.87 \\ 0.00 \end{array}$$

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

$$L_i = -\log(0.13) = 2.04$$

Unnormalized log-
probabilities / logits

unnormalized
probabilities

probabilities

Cross-Entropy Loss (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
car
frog

24.5

5.1

-1.7

normalize

exp

Unnormalized log-
probabilities / logits

0.13

164.0

unnorma-
lized
probabilities

0.87

0.18

0.00

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

$$= 2.04$$

Maximum Likelihood Estimation

Choose weights to maximize the likelihood of the observed data probabilities (See EECS 445 or EECS 545)

Cross-Entropy Loss (Multinomial Logistic Regression)

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

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



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

Probabilities must be ≥ 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

24.5	exp	164.0
0.13	normalize	0.87

0.13	0.87	0.00
------	------	------

1.00	0.00	0.00
------	------	------

Correct
probs

Cross-Entropy Loss (Multinomial Logistic Regression)

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



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

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

Probabilities must be ≥ 0

Probabilities must sum to 1

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

cat
car
frog

$$\begin{array}{r} 3.2 \\ 5.1 \\ -1.7 \end{array}$$

$$\begin{array}{r} 24.5 \\ 164.0 \\ 0.18 \end{array}$$

$$\begin{array}{r} 0.13 \\ 0.87 \\ 0.00 \end{array}$$

$$\begin{array}{r} 1.00 \\ 0.00 \\ 0.00 \end{array}$$

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

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

Kullback–Leibler divergence

Softmax function

Compare

Cross-Entropy Loss (Multinomial Logistic Regression)

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

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



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

Probabilities must be ≥ 0

Probabilities must sum to 1

Softmax function

cat

3.2

5.1

-1.7

normalize

164.0

0.18

0.00

Compare

1.00

0.00

0.00

Cross Entropy

1.00

0.00

0.00

$H(P, Q) =$

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

Correct probs

probabilities

unnormalized probabilities / logits

Unnormalized log-probabilities / logits

Cross-Entropy Loss (Multinomial Logistic Regression)

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

$$\mathbf{s} = f(\mathbf{x}_i; \mathbf{W})$$



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

Softmax function

Maximize probability of correct class
Putting it all together:

$$L_i = -\log P(Y = y_i | \mathbf{X} = \mathbf{x}_i) \quad L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

3.2

car 5.1

frog -1.7

Cross-Entropy Loss (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

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

Putting it all together:

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

car 5.1

frog -1.7

Cross-Entropy Loss (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

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

Putting it all together:

car 5.1

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

A: Min 0, max +infinity

frog -1.7

Cross-Entropy Loss (Multinomial Logistic Regression)

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

$$\mathbf{s} = f(\mathbf{x}_i; \mathbf{W})$$



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

Softmax function

Maximize probability of correct class

$$L_i = -\log P(Y = y_i | \mathbf{X} = \mathbf{x}_i) \quad L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

Putting it all together:

- cat 3.2 car 5.1 Q: If all scores are small random values, what is the loss?
- frog -1.7

Cross-Entropy Loss (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

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

Putting it all together:

car 5.1

Q: If all scores are
small random values,
what is the loss?
A: -log(C)
 $\log(10) \approx 2.3$

frog -1.7

Cross-Entropy vs SVM Loss

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

$y_i = 0$

and

Q: What is cross-entropy loss?
What is SVM loss?

Cross-Entropy vs SVM Loss

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

$y_i = 0$

and

Q: What is cross-entropy loss?
What is SVM loss?

A: Cross-entropy loss > 0
SVM loss = 0

Cross-Entropy vs SVM Loss

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

$y_i = 0$

and

Q: What happens to each loss if I slightly change the scores of the last datapoint?

Cross-Entropy vs SVM Loss

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

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

$y_i = 0$

and

Q: What happens to each loss if I slightly change the scores of the last datapoint?

A: Cross-entropy loss will change;
SVM loss will stay the same

Cross-Entropy vs SVM Loss

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

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

assume scores:

$$\boxed{[10, -2, 3]}$$

$$\boxed{[10, 9, 9]}$$

$$\boxed{[10, -100, -100]}$$

$$\boxed{y_i = 0}$$

and

Q: What happens to each loss if I double the score of the correct class from 10 to 20?

Cross-Entropy vs SVM Loss

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

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

assume scores:

$$\boxed{[10, -2, 3]}$$

$$\boxed{[10, 9, 9]}$$

$$\boxed{[10, -100, -100]}$$

$$\boxed{y_i = 0}$$

and

Q: What happens to each loss if I
double the score of the correct class
from 10 to 20?

A: Cross-entropy loss will decrease,
SVM loss still 0

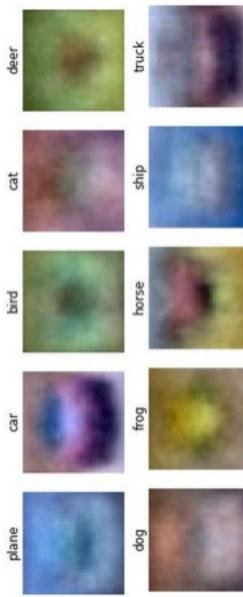
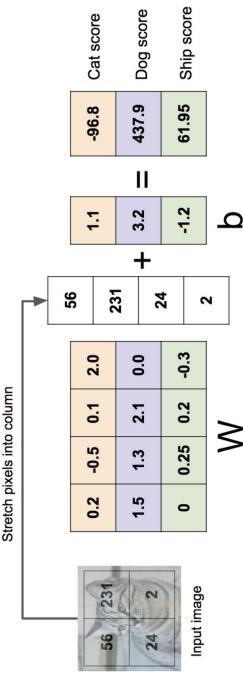
Recap: Three ways to think about linear classifiers

Algebraic Viewpoint

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

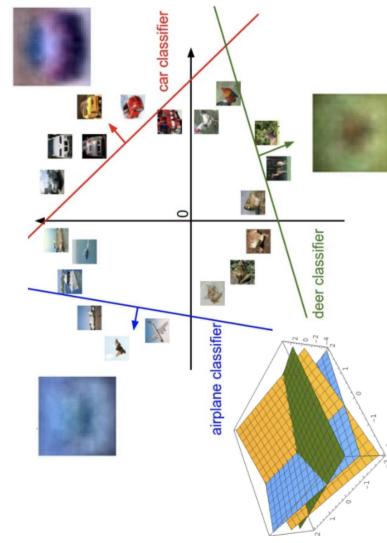
Visual Viewpoint

One template
per class



Geometric Viewpoint

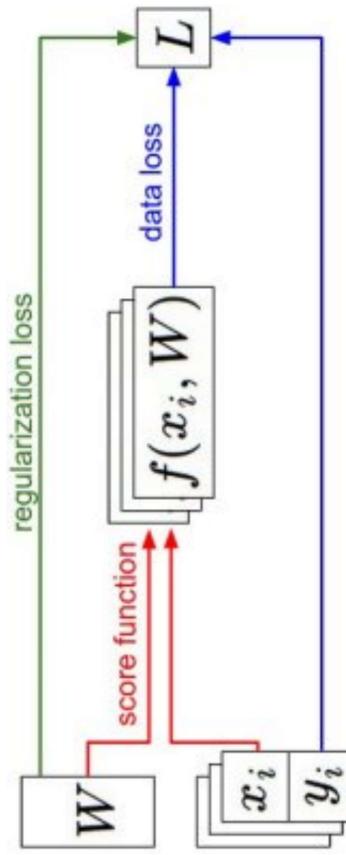
Hyperplanes
cutting up space



Recap: Loss Functions quantify preferences

- We have some dataset of (x, y)
- We have a **score function**:
- We have a **loss function**:

$$\begin{aligned} L_i &= -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \text{ Softmax} \\ L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ L &= \frac{1}{N} \sum_{i=1}^N L_i + R(W) \text{ Full loss} \end{aligned}$$



Recap: Loss Functions quantify preferences

- We have some dataset of (x, y)
 - We have a **score function**:
 - We have a **loss function**:
- Q: How do we find the best W ?**

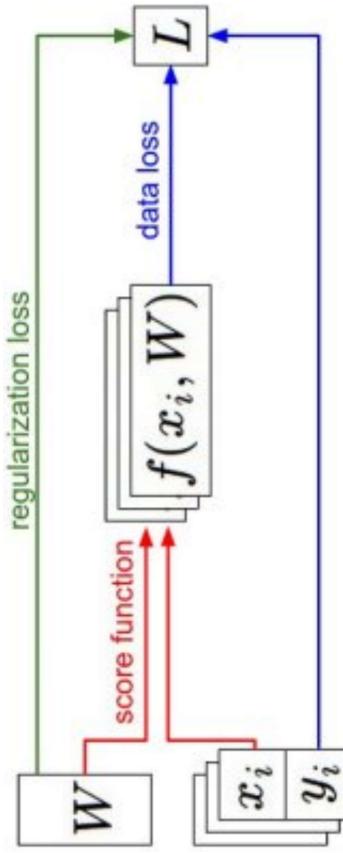
$$s = f(x; W) = Wx$$

Linear classifier

$$L_i = -\log\left(\frac{e^{sy_i}}{\sum_j e^{sj}}\right) \text{ Softmax}$$

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

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



Next time:
Optimization