

CS163: Deep Learning for Computer Vision

Lecture 3: Linear Neural Networks

Last time: Image Classification

Input: image



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



chair
bed
sofa
table
cabinet

This image by Nikita is
licensed under CC-BY 2.0

Last Time: Challenges of Visual Recognition

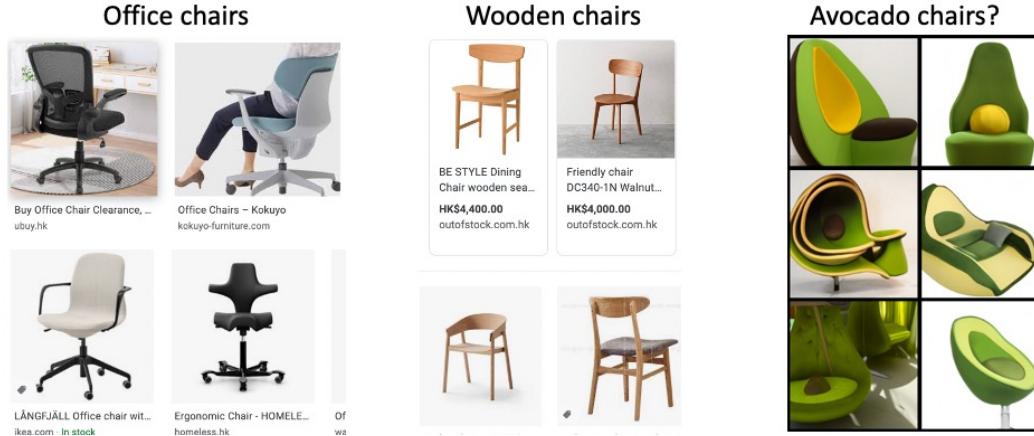
Viewpoint



Domain changes



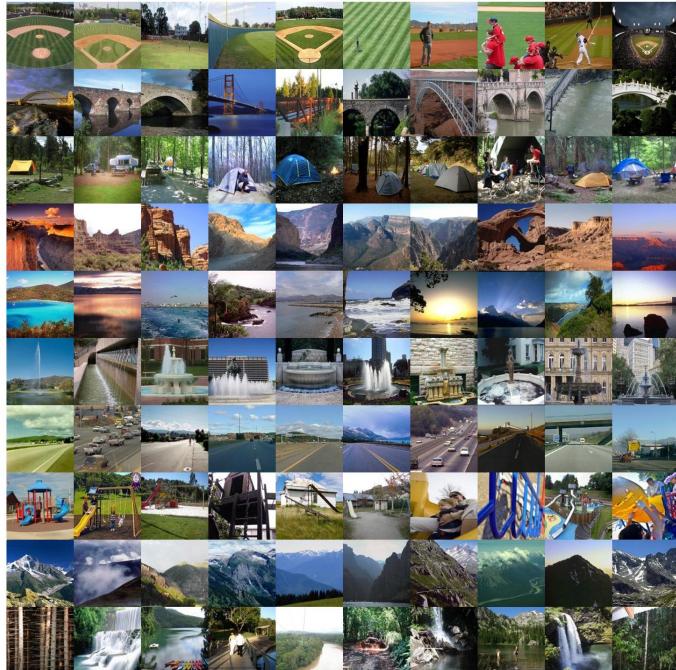
Intraclass variation and subcategories



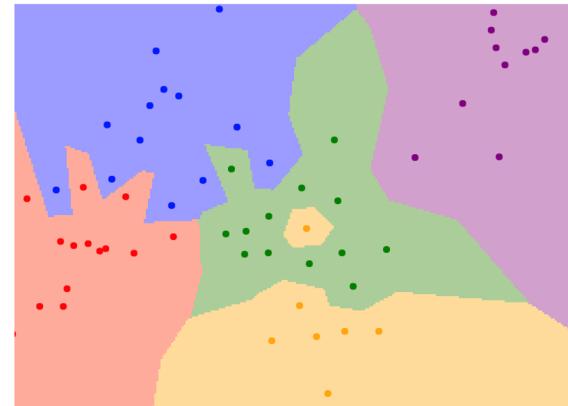
Occlusion & Clutter



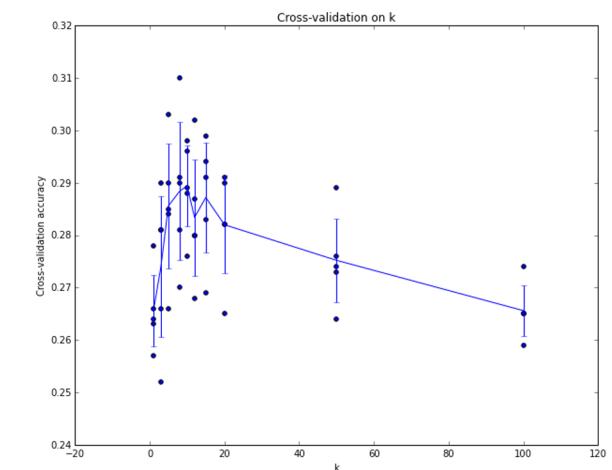
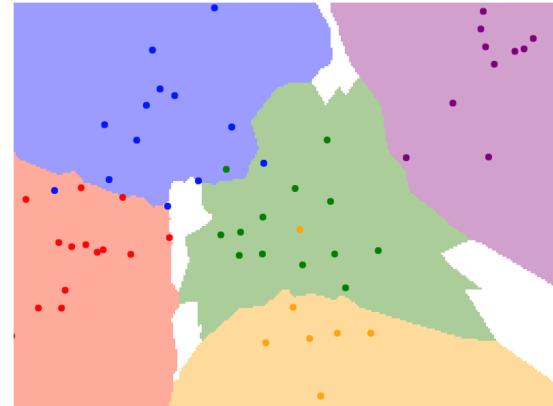
Last time: Data-Drive Approach, kNN



1-NN classifier



5-NN classifier



Today: Linear Neural
Networks/Layers/Classifiers

Linear Layer is the building block of deep neural network



TinyPlaces: 20 scenes categories from Places

10 outdoor scenes



baseball_field

bridge

campsite

canyon

coast

fountain

highway

playground

mountain

rainforest

10 indoor scenes



bathroom

bedroom

bookstore

classroom

dining_room

food_court

kitchen

lobby

living_room

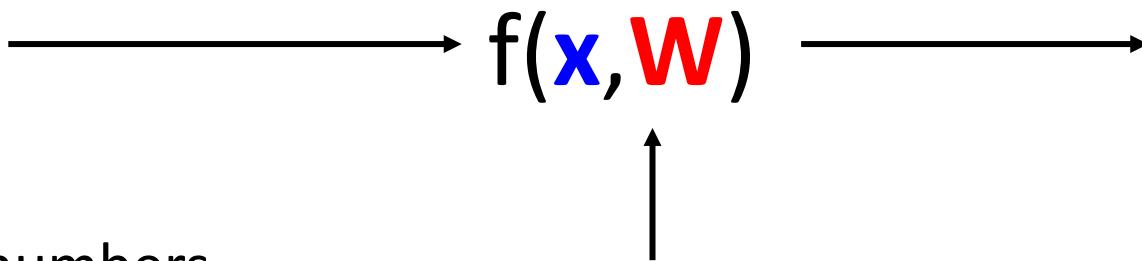
office

20,000 training images
each image is **32x32x3**

2,000 test images.

Parametric Approach: Indoor/Outdoor Classification

Image



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

1 number giving
outdoor-ness
score

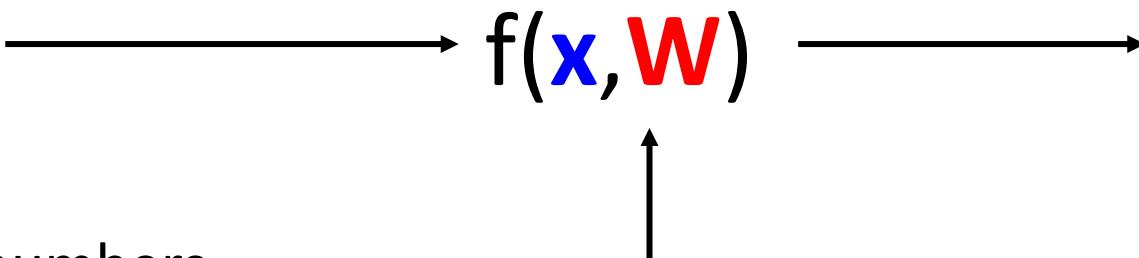
W
parameters
or weights

Parametric Approach: Scene Classification

Image



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



20 numbers giving
class scores

bathroom
bedroom
bookstore
classroom
dining_room
food_court
kitchen
lobby
living_room
Office

....

Parametric Approach: Linear Classifier

Image



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

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

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



W
parameters
or weights

20 numbers giving
class scores

Parametric Approach: Linear Classifier $(3072,)$

Image



Array of $32 \times 32 \times 3$ numbers
(3072 numbers total)

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

(20,) (20, 3072)

20 numbers giving
class scores

W
parameters
or weights

Parametric Approach: Linear Classifier

(3072,)

Image



$$f(x, W) = Wx + b \quad (20,)$$

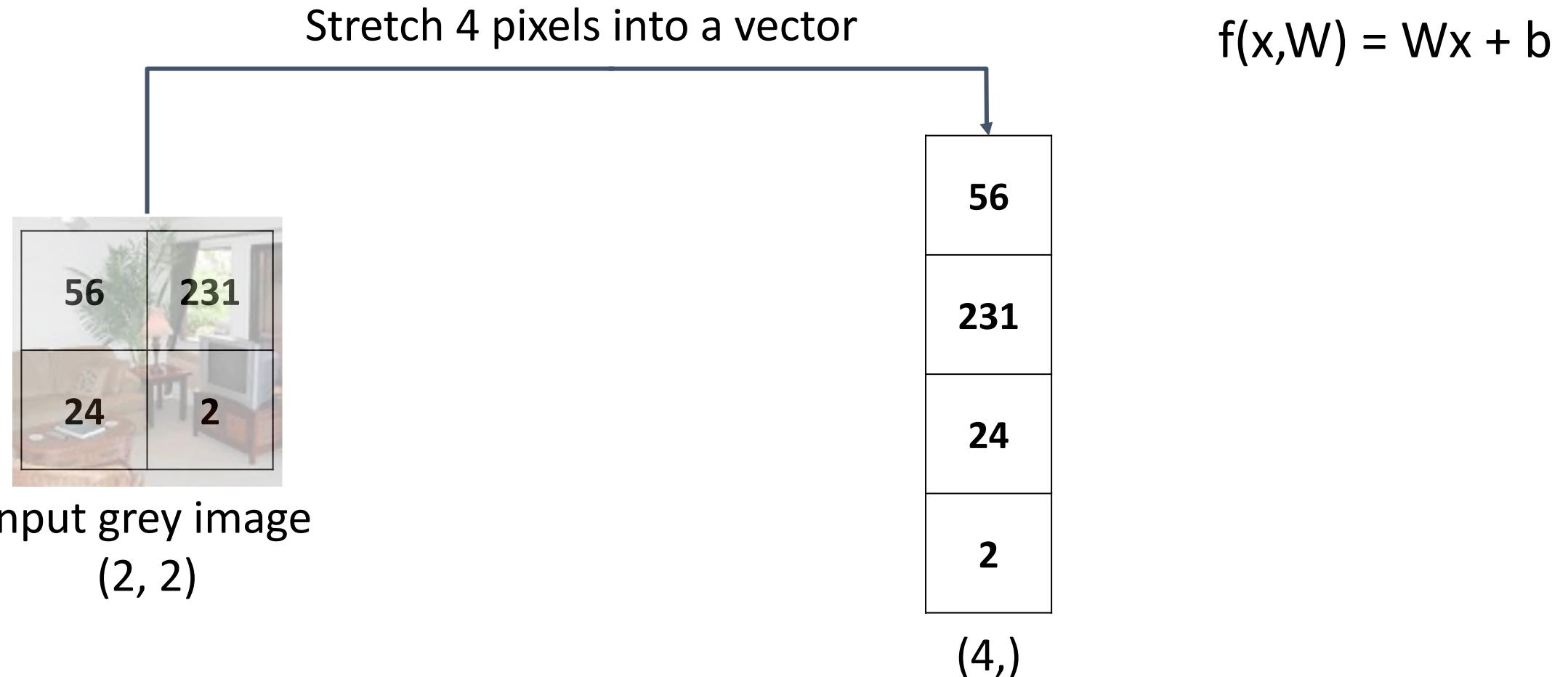
(20,) (20, 3072)

20 numbers giving
class scores

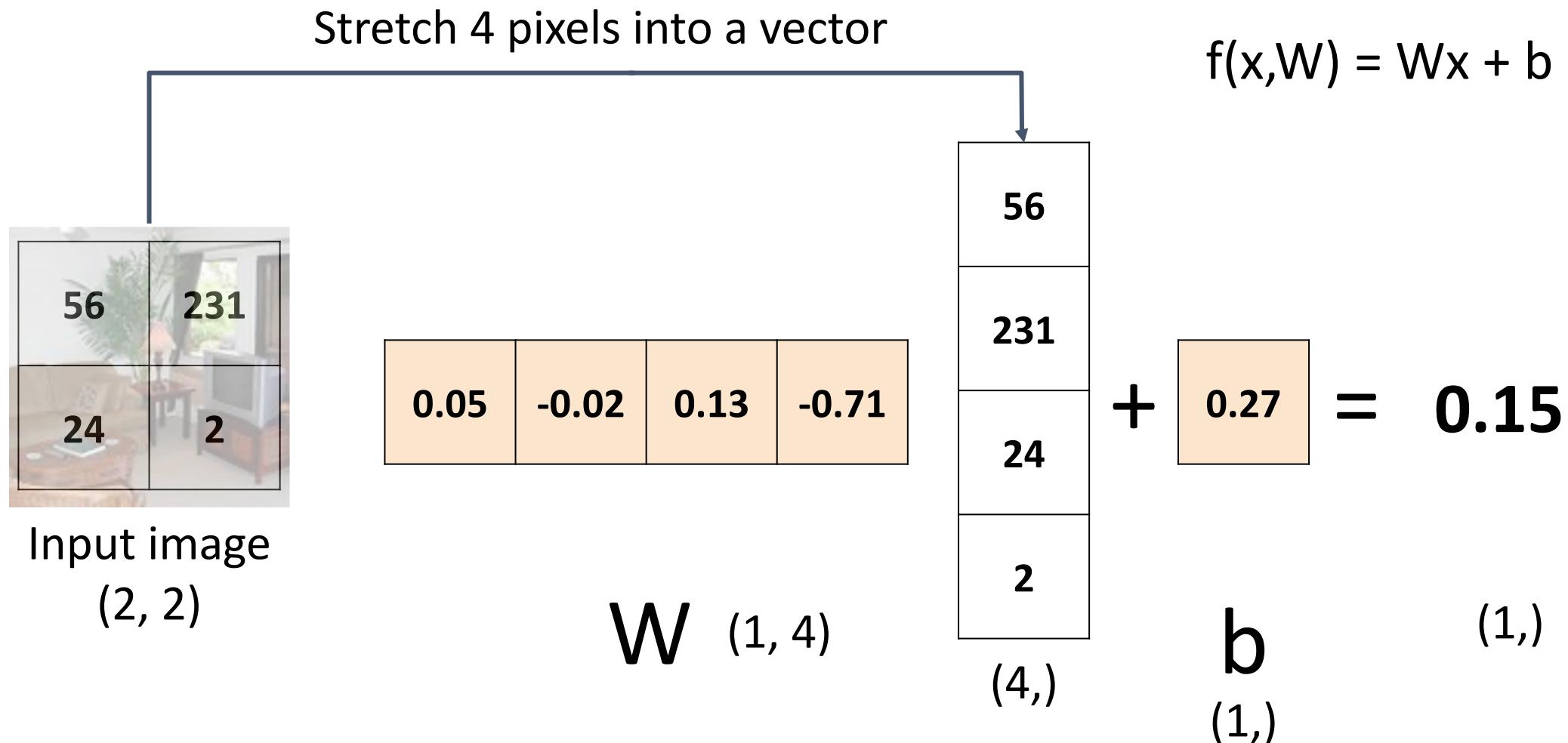
Array of **32x32x3** numbers
(3072 numbers total)

W
parameters
or weights

Example for 2x2 image, outdoor-ness prediction

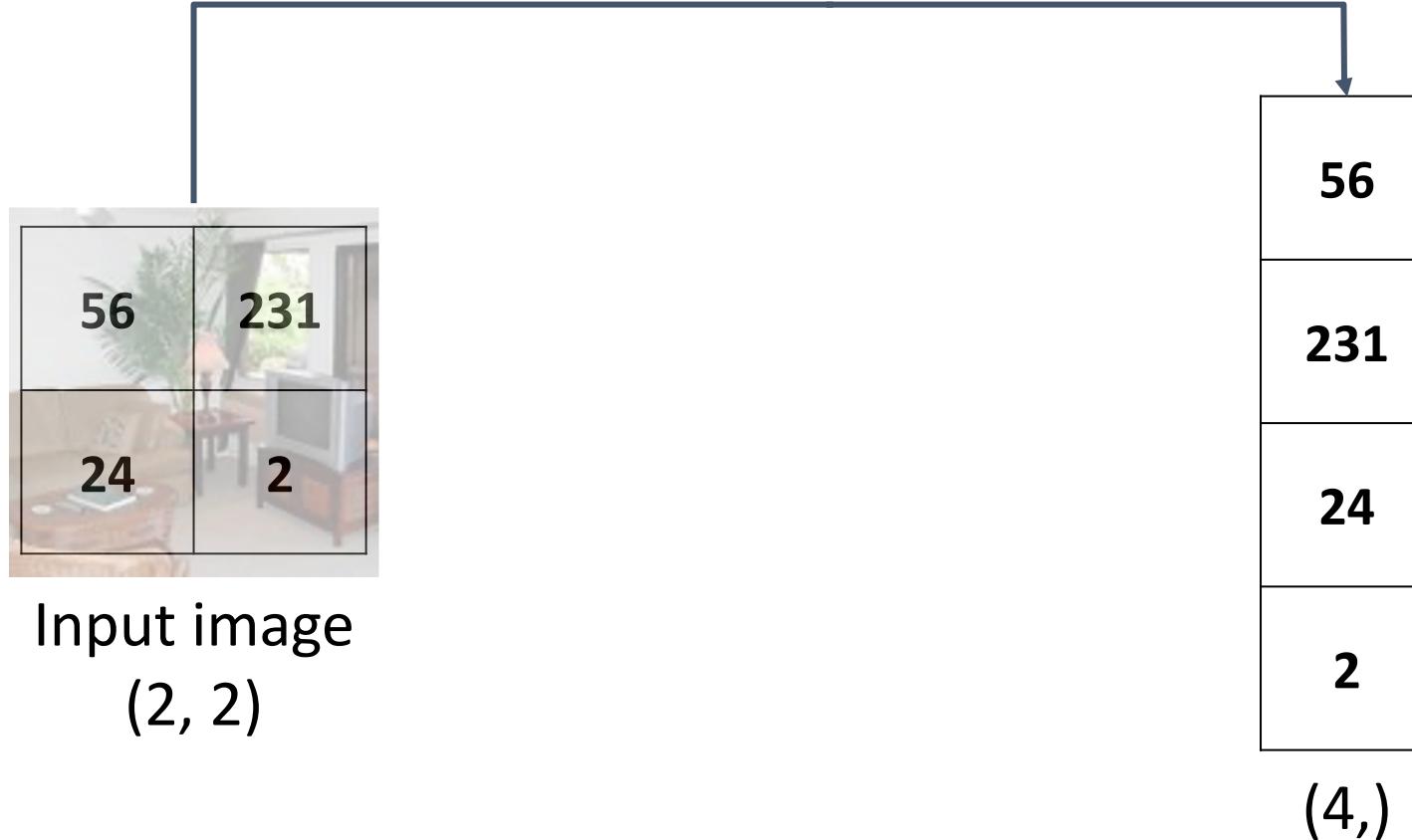


Example for 2x2 image, outdoor-ness prediction



Example for 2x2 image, 3 classes (livingroom/highway/mountain)

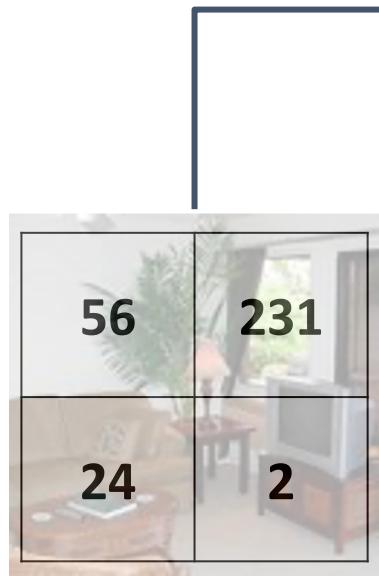
Stretch 4 pixels into a vector



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

Example for 2x2 image, 3 classes (livingroom/highway/mountain)

Stretch 4 pixels into a vector



Input image
(2, 2)

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

W (3, 4)

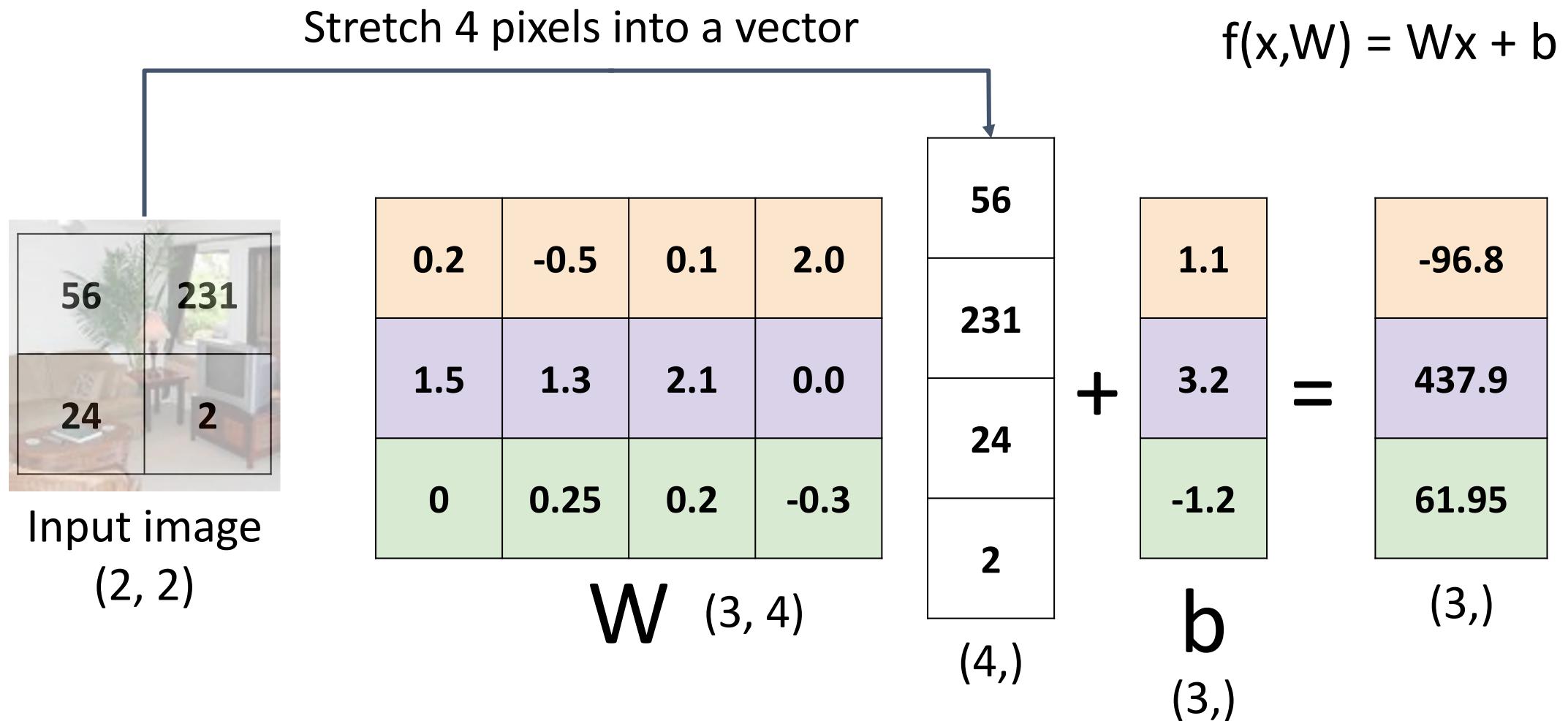
56
231
24
2

(4,)

$$f(x, W) = Wx + b$$
$$\begin{array}{c} + \\ \hline \end{array} \quad \begin{array}{c} b \\ (3,) \end{array} \quad = \quad \begin{array}{c} -96.8 \\ 437.9 \\ 61.95 \end{array}$$

The diagram shows the computation of the output vector. It starts with the input vector (4,) on the left, followed by a plus sign, then the bias vector b with dimensions (3,), and finally an equals sign followed by the resulting output vector (3,). The output vector contains values -96.8, 437.9, and 61.95.

Linear Classifier: Linear Transformation



Linear Classifier: Bias Trick

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

Stretch pixels into column



Input image
(2, 2)

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

W (3, 5)



= (5,)

-96.8
437.9
61.95

(3,)

Linear Classifier: Predictions are Linear!

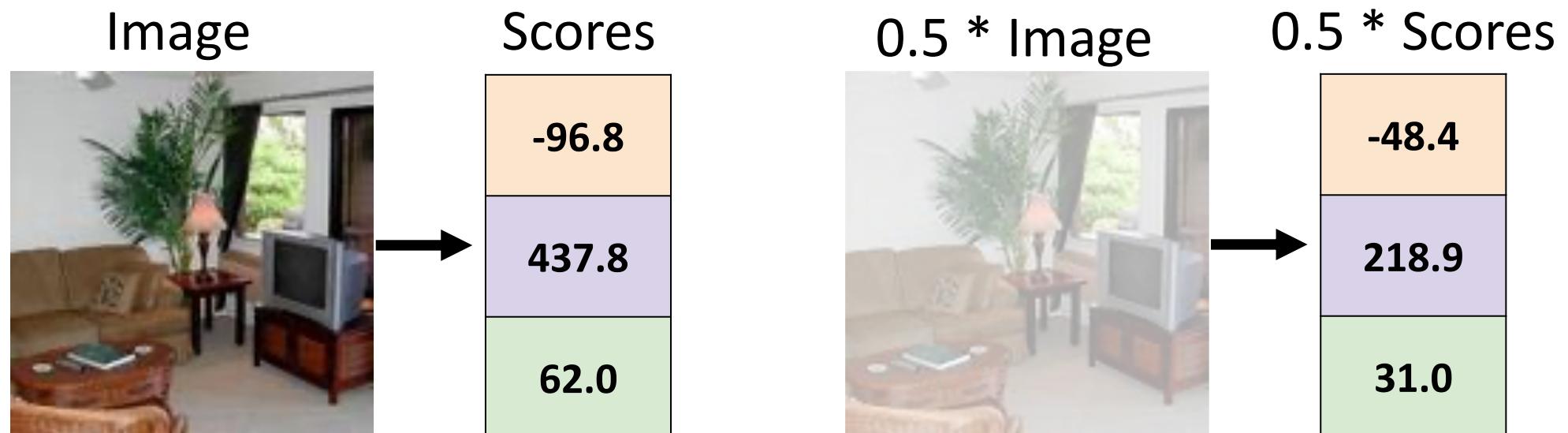
$$f(x, W) = Wx \quad (\text{ignore bias})$$

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

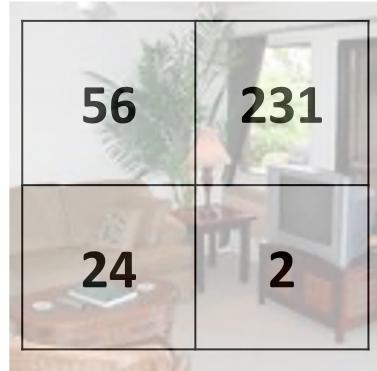
Linear Classifier: Predictions are Linear!

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

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

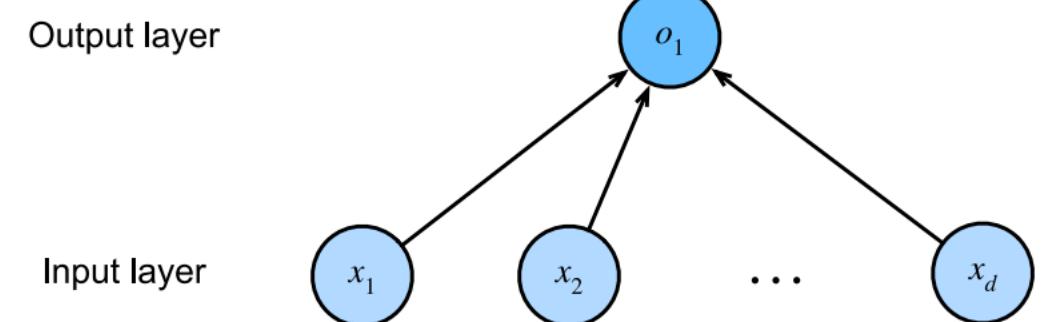


Linear Classifier is a single-layer neural network



Outdoor-ness

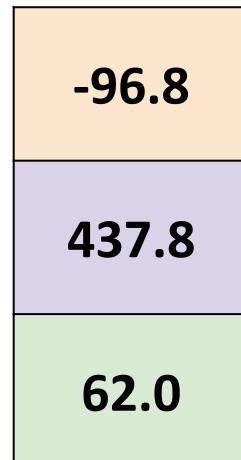
→ 0.15



3 Scene Classification

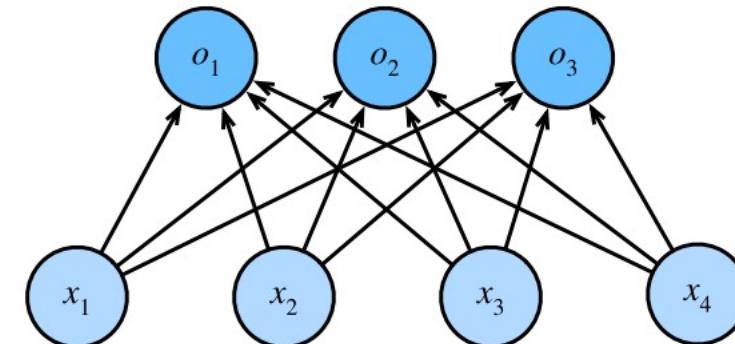


→



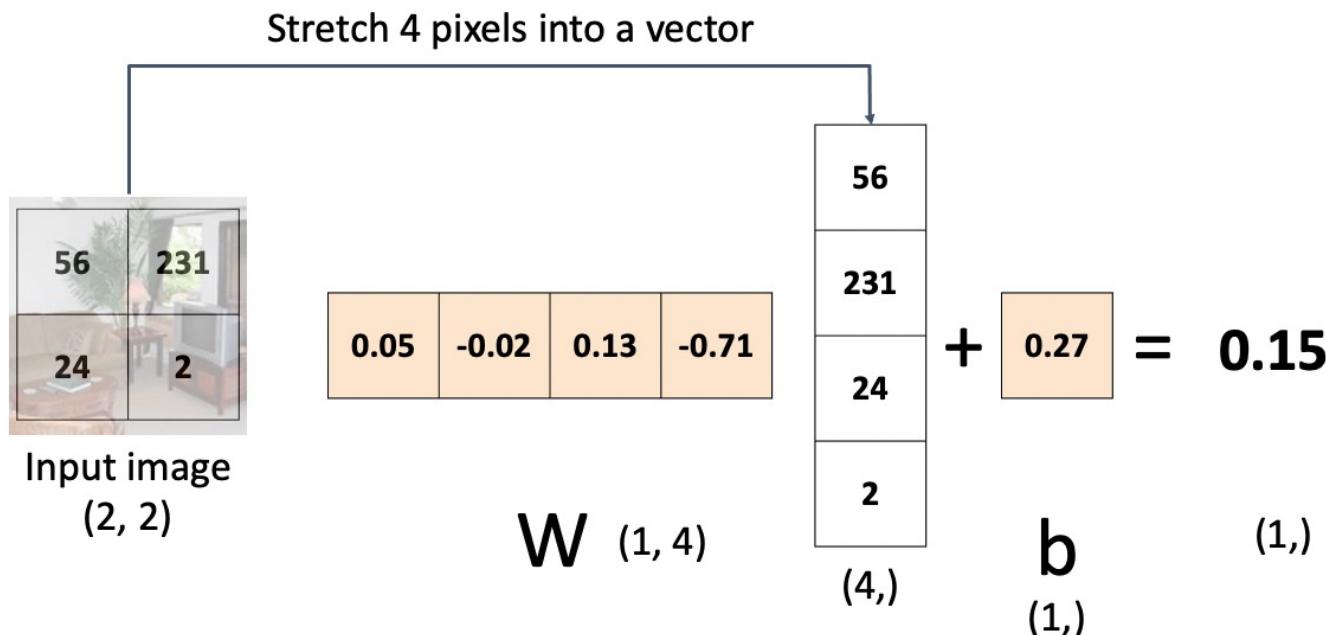
Output layer

Input layer



Interpreting a Linear Classifier

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



Interpreting a Linear Classifier

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

Instead of stretching pixels into a vector we can equivalently stretch rows of W into image!

Stretch 4 pixels into a vector

Input image $(2, 2)$

$W \quad (1, 4)$

$b \quad (1, 1)$

0.15

Matrix W and vector b are shown below:

$W = \begin{bmatrix} 0.05 & -0.02 & 0.13 & -0.71 \end{bmatrix}$

$b = \begin{bmatrix} 0.27 \end{bmatrix}$



W

$$\begin{bmatrix} 0.05 & -0.02 \\ 0.13 & -0.71 \end{bmatrix}$$

b

$$0.27$$

$$0.15$$

Interpreting a Linear Classifier

10 outdoor scenes



10 indoor scenes



More blue
at the top

More dark at
the bottom



Average images for
outdoor and indoor



W

b

$$\begin{matrix} 0.05 & -0.02 \\ 0.13 & -0.71 \end{matrix}$$

$$0.27$$

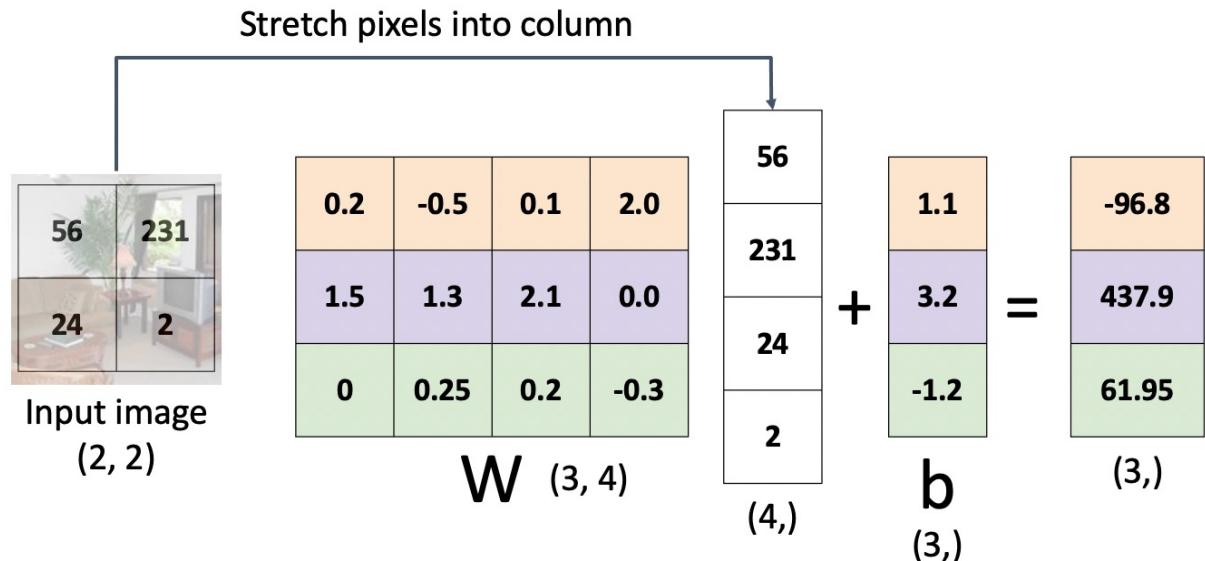
$$0.15$$

56	231
24	2

Interpreting a Linear Classifier

Similar to scene classifier

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



Interpreting a Linear Classifier

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

Stretch pixels into column

56	231
24	2

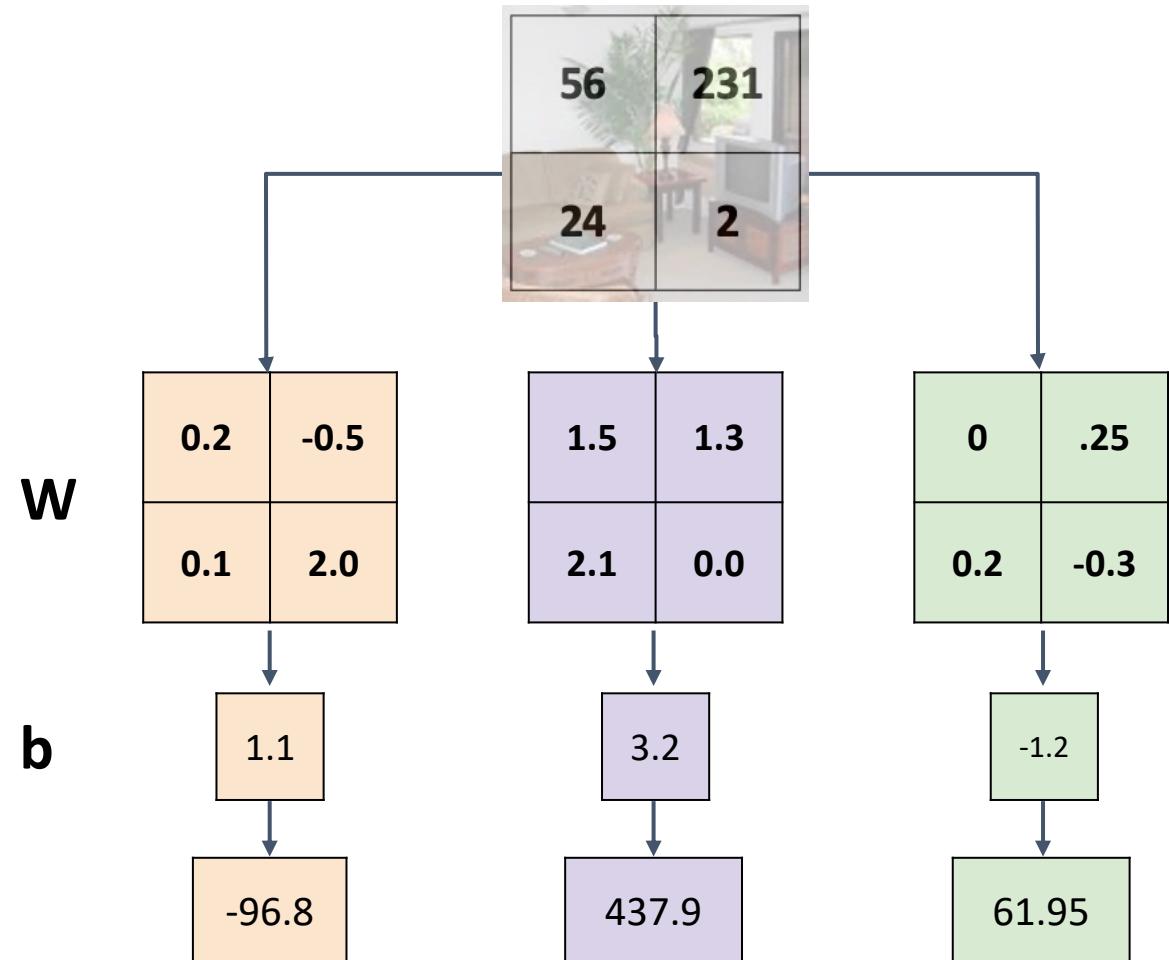
Input image (2, 2)

$W \quad (3, 4)$

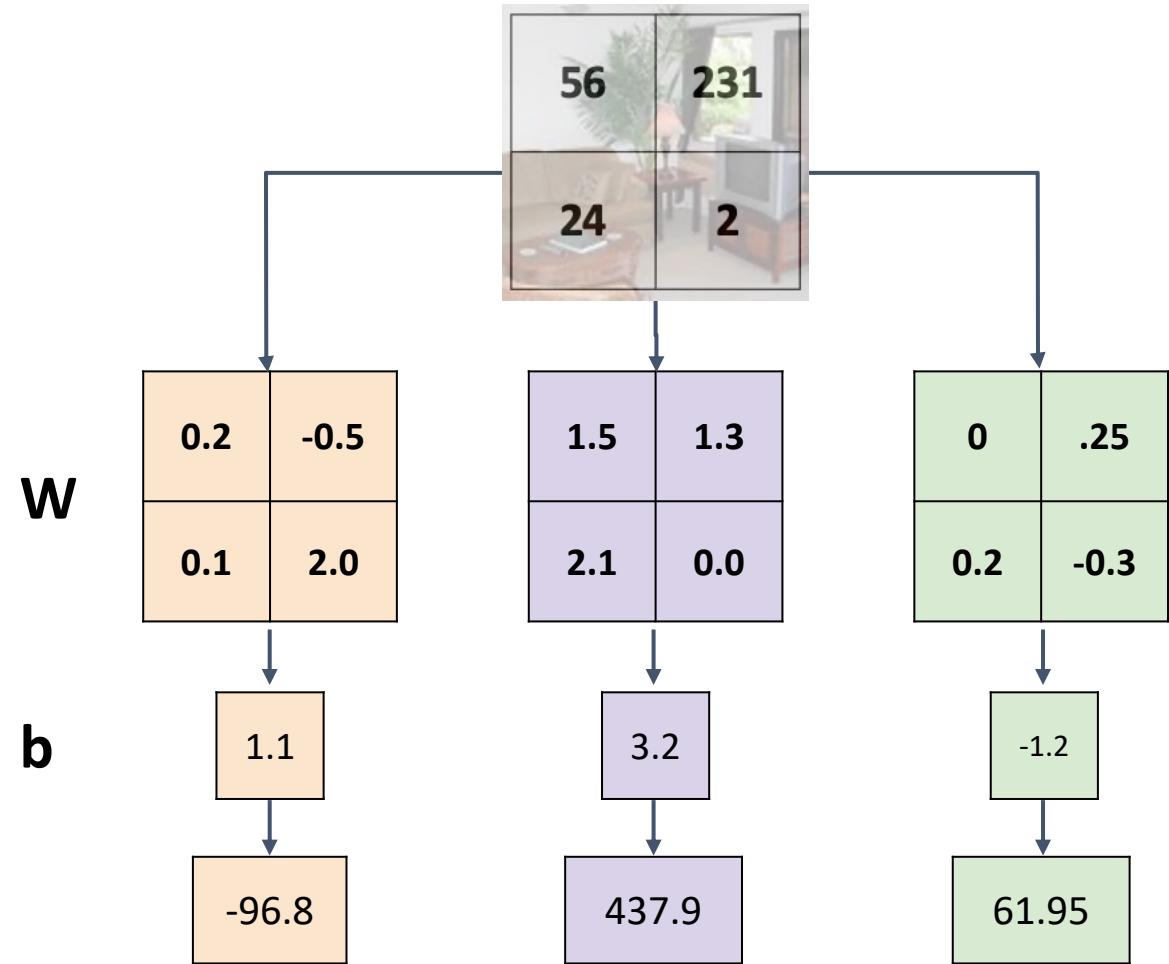
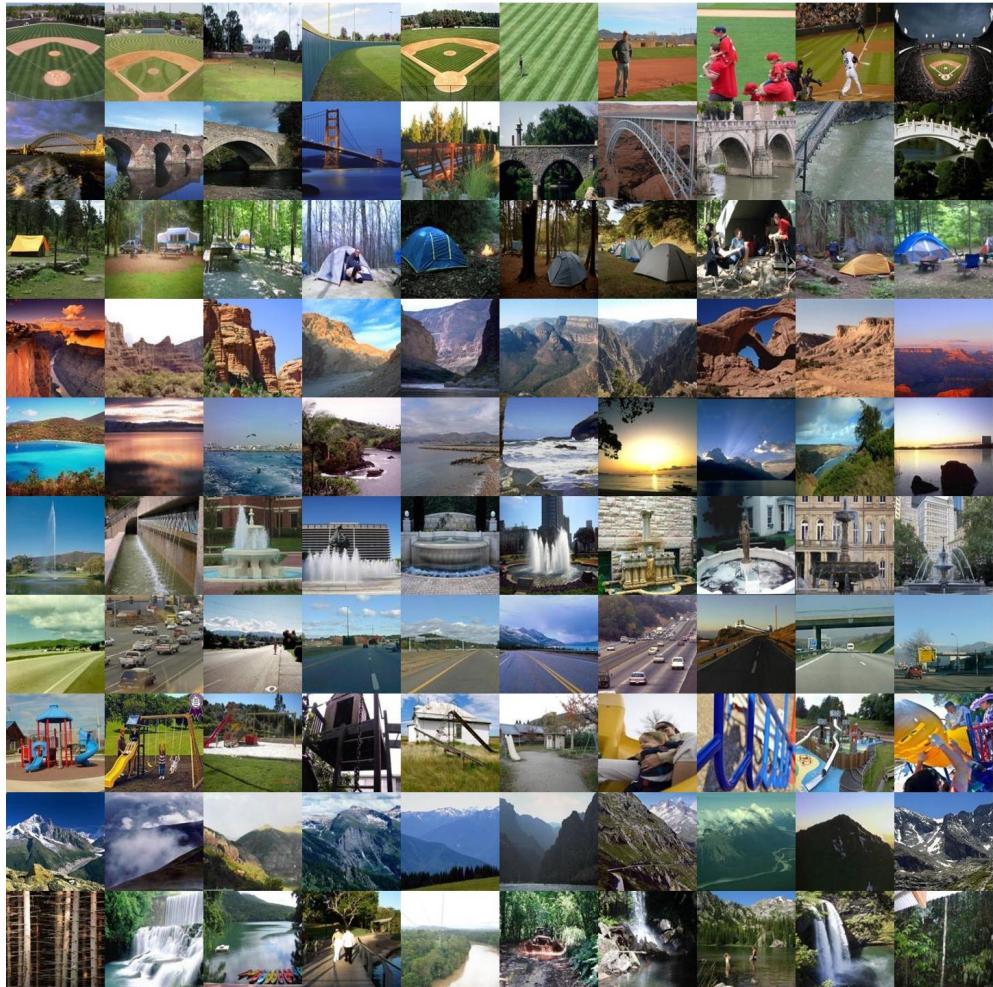
$$\begin{matrix} 0.2 & -0.5 & 0.1 & 2.0 \\ 1.5 & 1.3 & 2.1 & 0.0 \\ 0 & 0.25 & 0.2 & -0.3 \end{matrix} + \begin{matrix} 56 \\ 231 \\ 24 \\ 2 \end{matrix} = \begin{matrix} 1.1 \\ 3.2 \\ -1.2 \end{matrix} \quad \begin{matrix} -96.8 \\ 437.9 \\ 61.95 \end{matrix}$$

$b \quad (3,)$

(4,)



Interpreting a Linear Classifier



Average Images for Categories of TinyPlaces

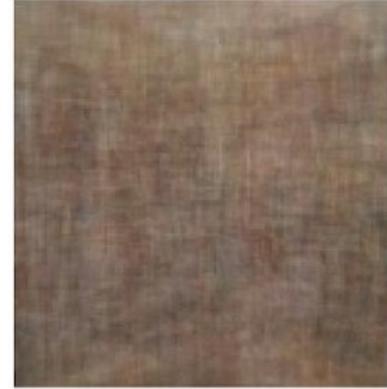
bathroom



bedroom



bookstore



classroom



diningroom



fountain



highway



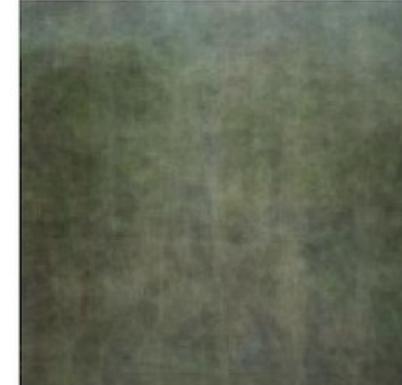
playground



mountain



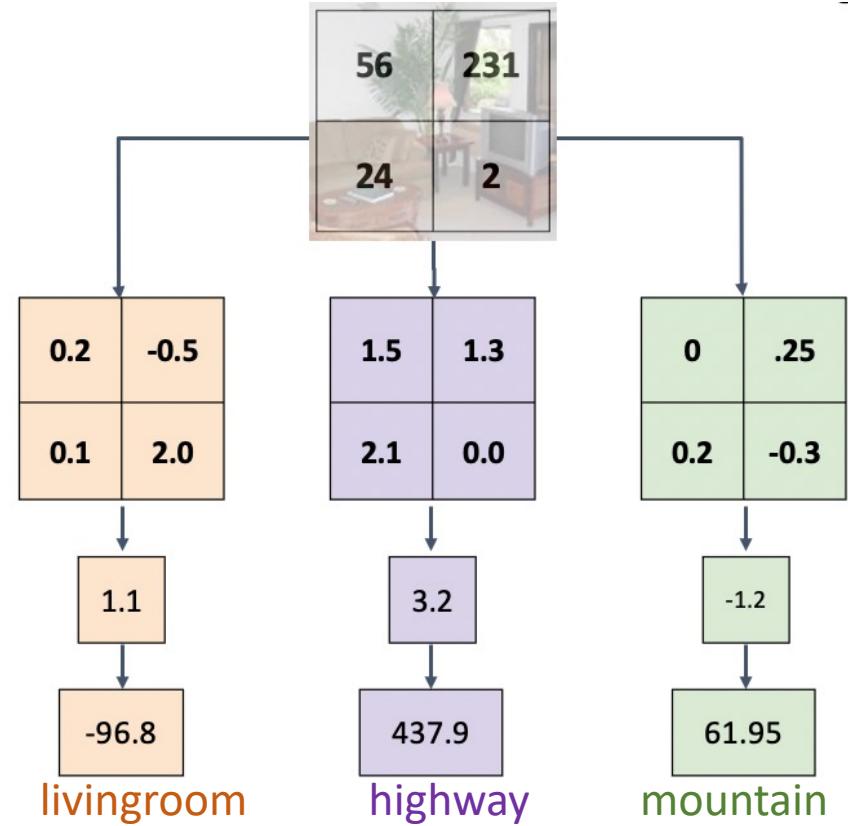
rainforest



Interpreting a Linear Classifier

Linear classifier has one “template” per category, then compute the correlation.

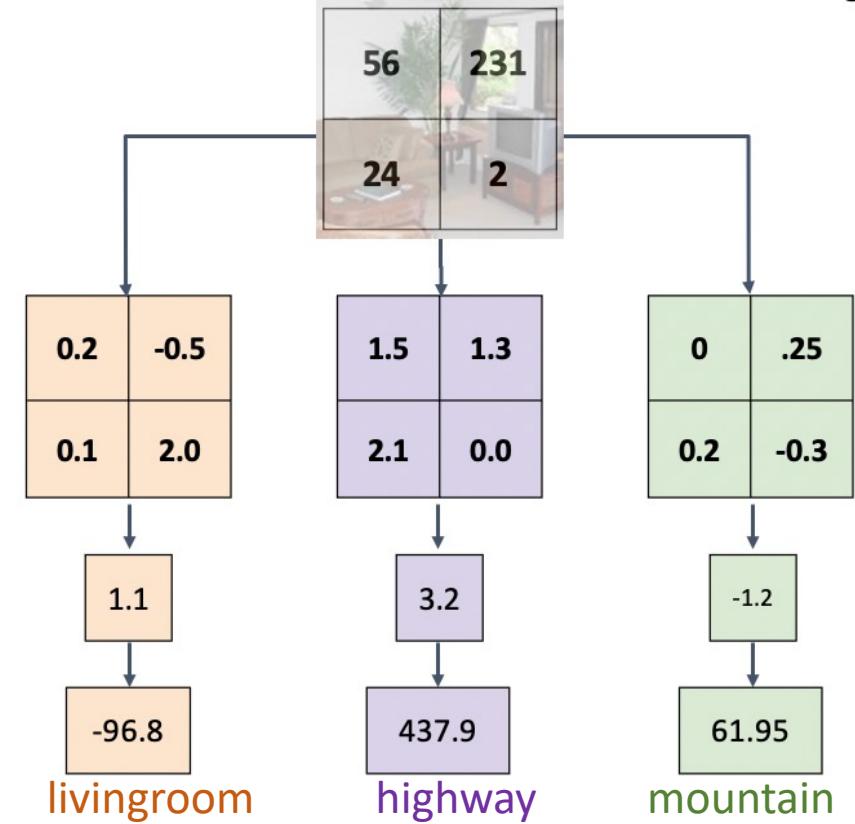
Drawback: A single template cannot capture multiple modes of the data



Understanding the Linear Classifier

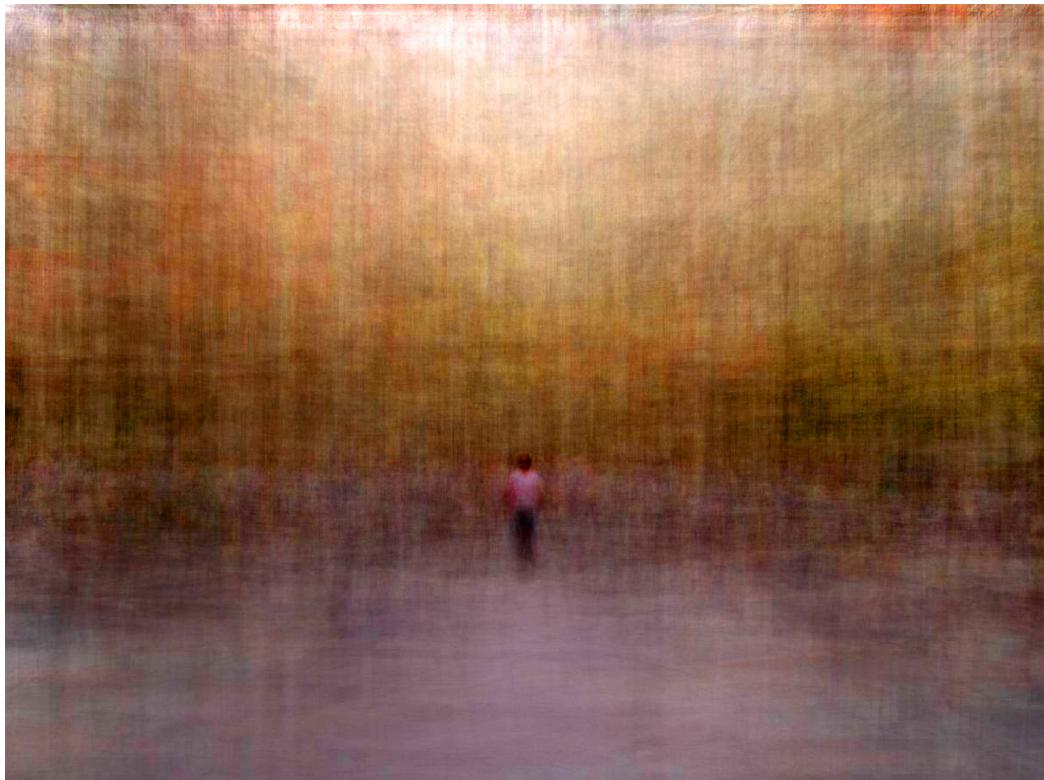
Drawback: A single template cannot capture multiple modes of the data

Potential improvement: go beyond pixel template



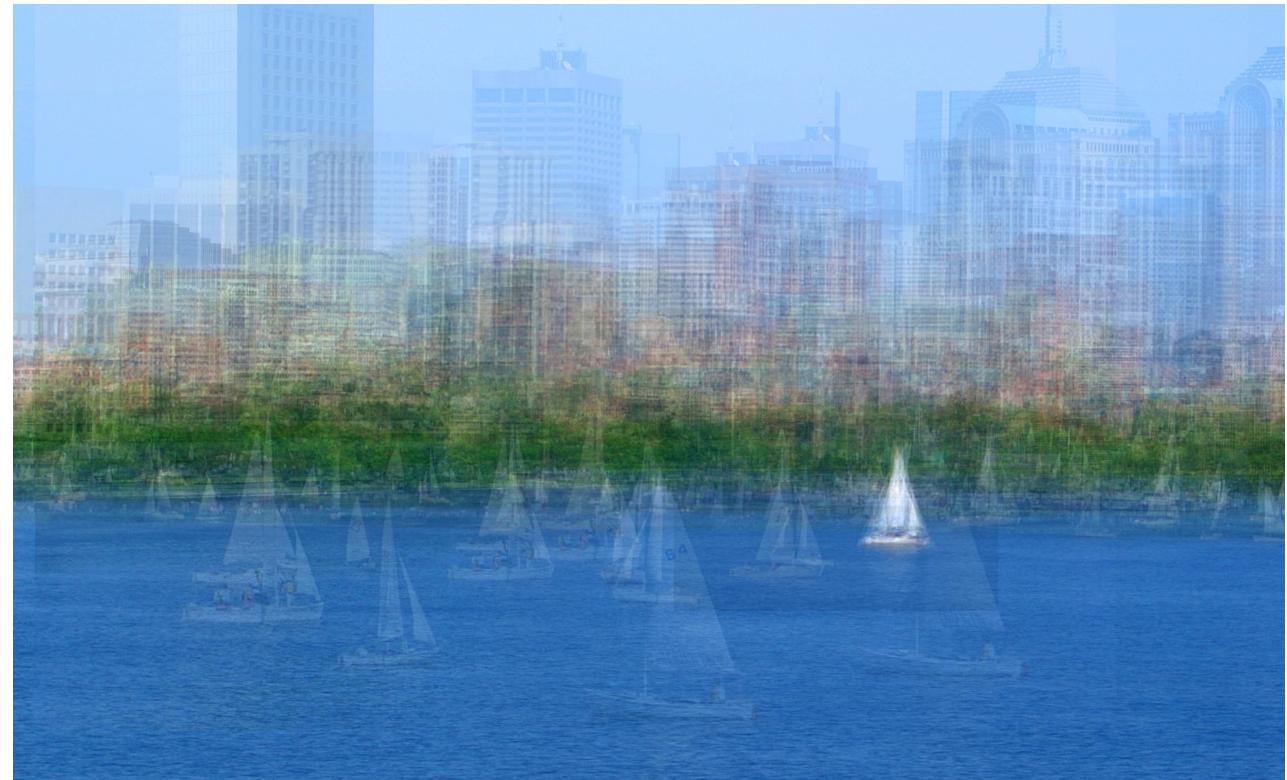
Average Image to explore regularities in data

- Side project from Prof. Antonio Torralba (my former PhD advisor)



People in Valladolid, Spain

<https://people.csail.mit.edu/torralba/gallery/>

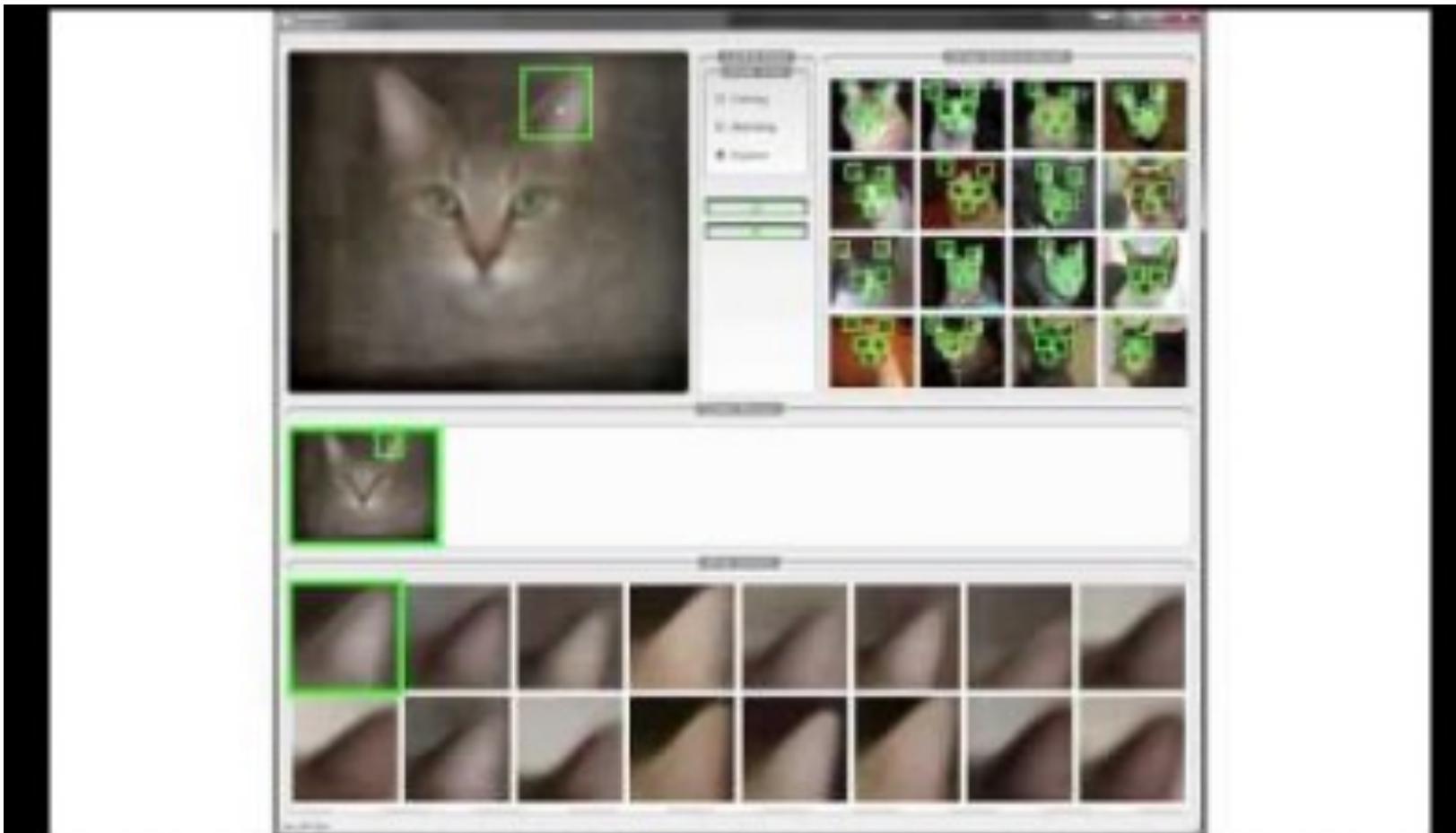


Sailboats in Charles river (spring). Pictures aligned on a sailboat on the right.

Lecture 3 - 34

Average Image to explore regularities in data

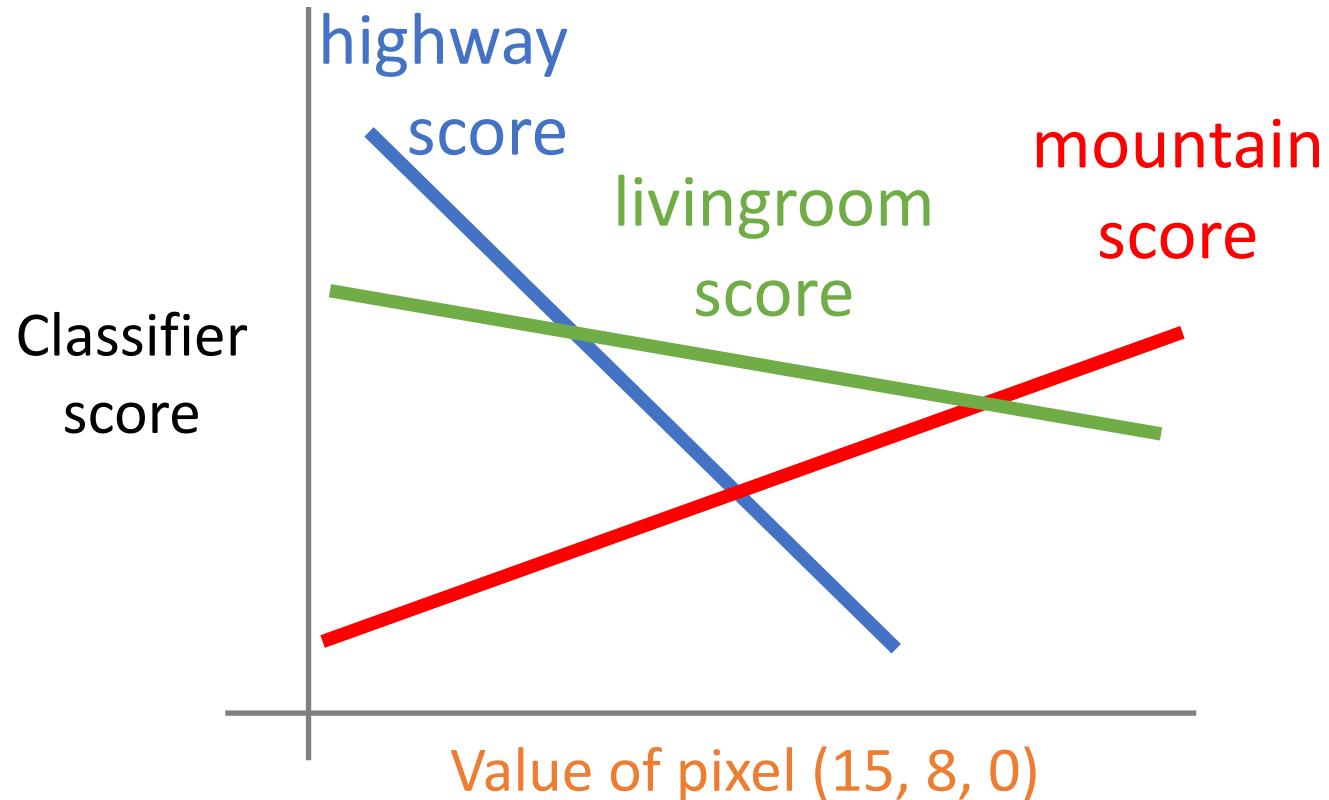
- AverageExplorer from Zhu, Lee, Efros. SIGGRAPH'14



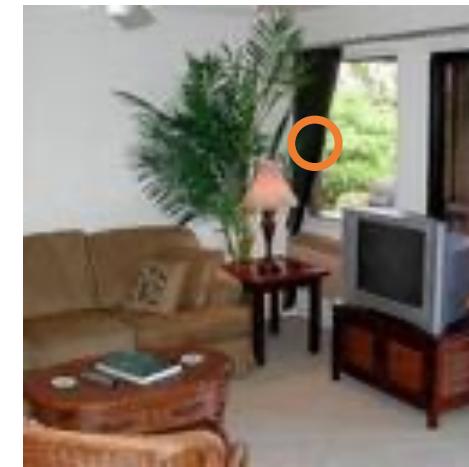
<https://www.cs.cmu.edu/~junyanz/projects/averageExplorer>

L

Interpreting a Linear Classifier: Geometric Viewpoint



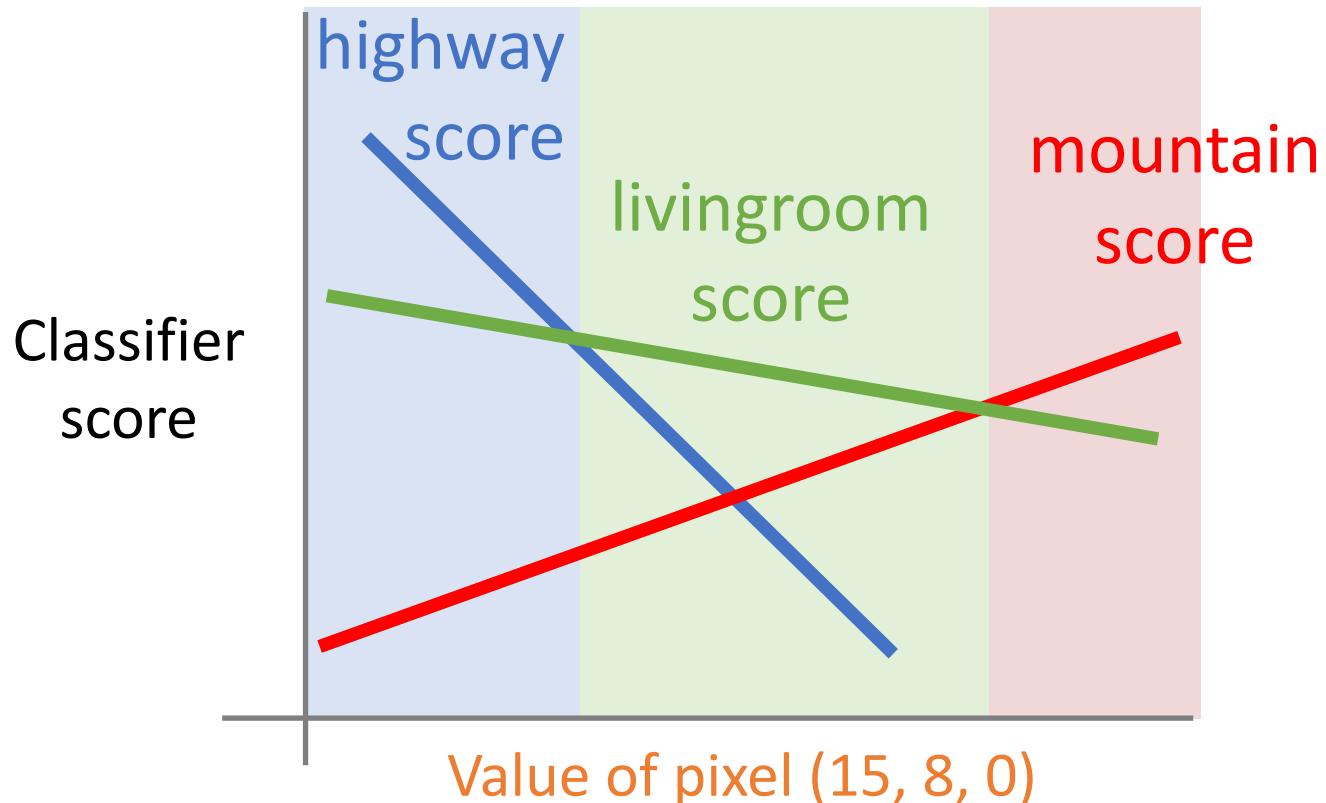
$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + b$$



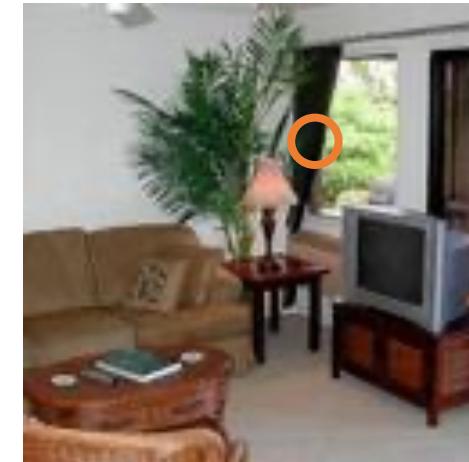
Array of **32x32x3** numbers
(3072 numbers total)

Interpreting a Linear Classifier: Geometric Viewpoint

Decision Regions

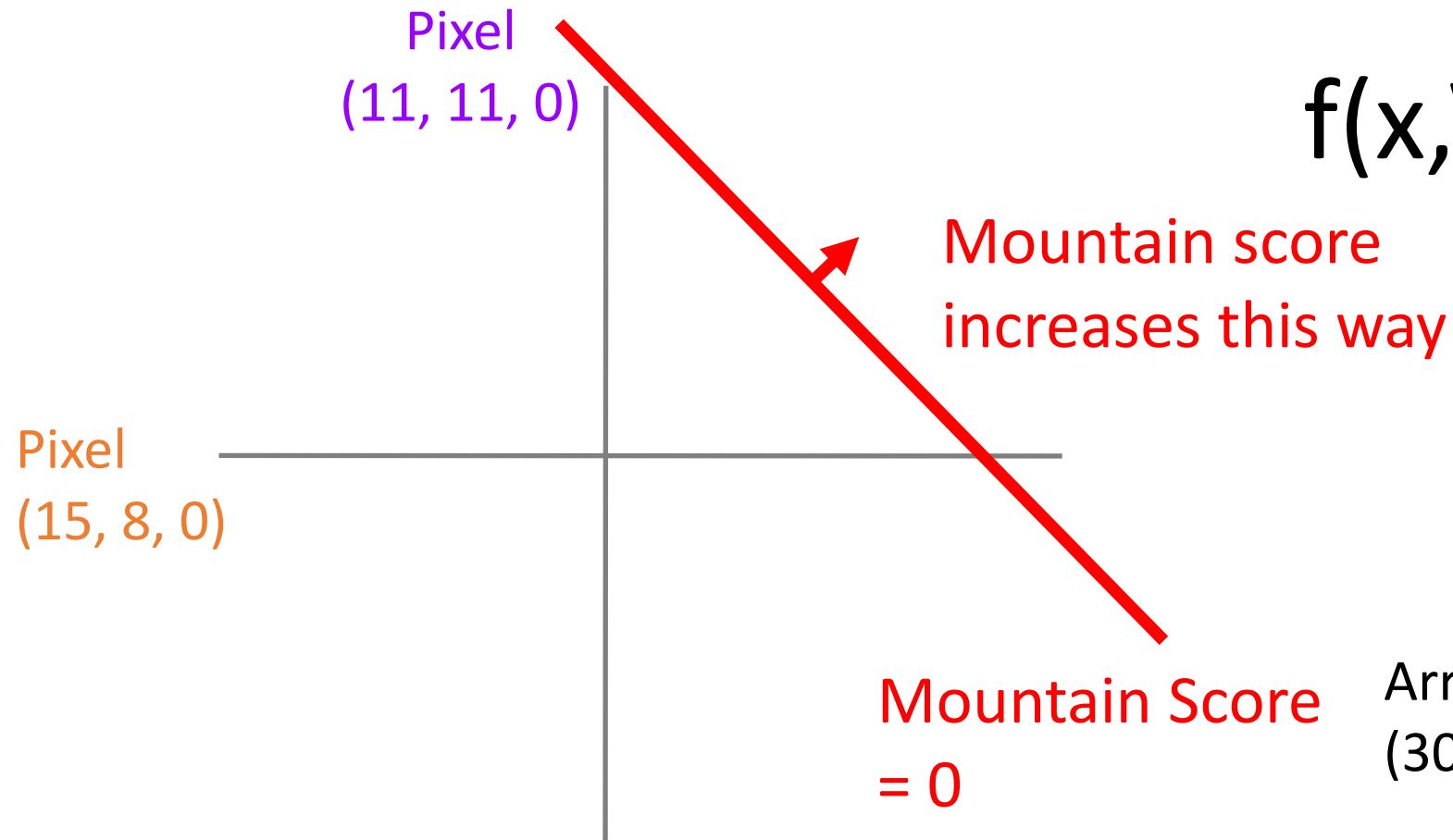


$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$

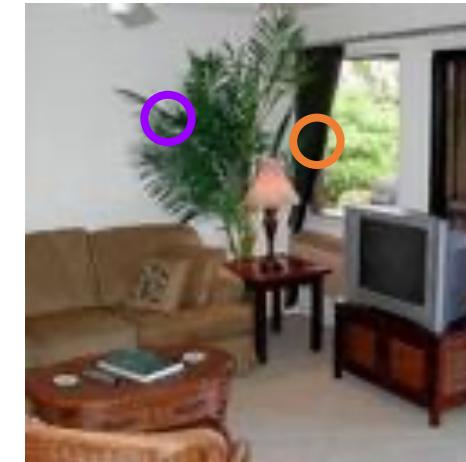


Array of **32x32x3** numbers
(3072 numbers total)

Interpreting a Linear Classifier: Geometric Viewpoint

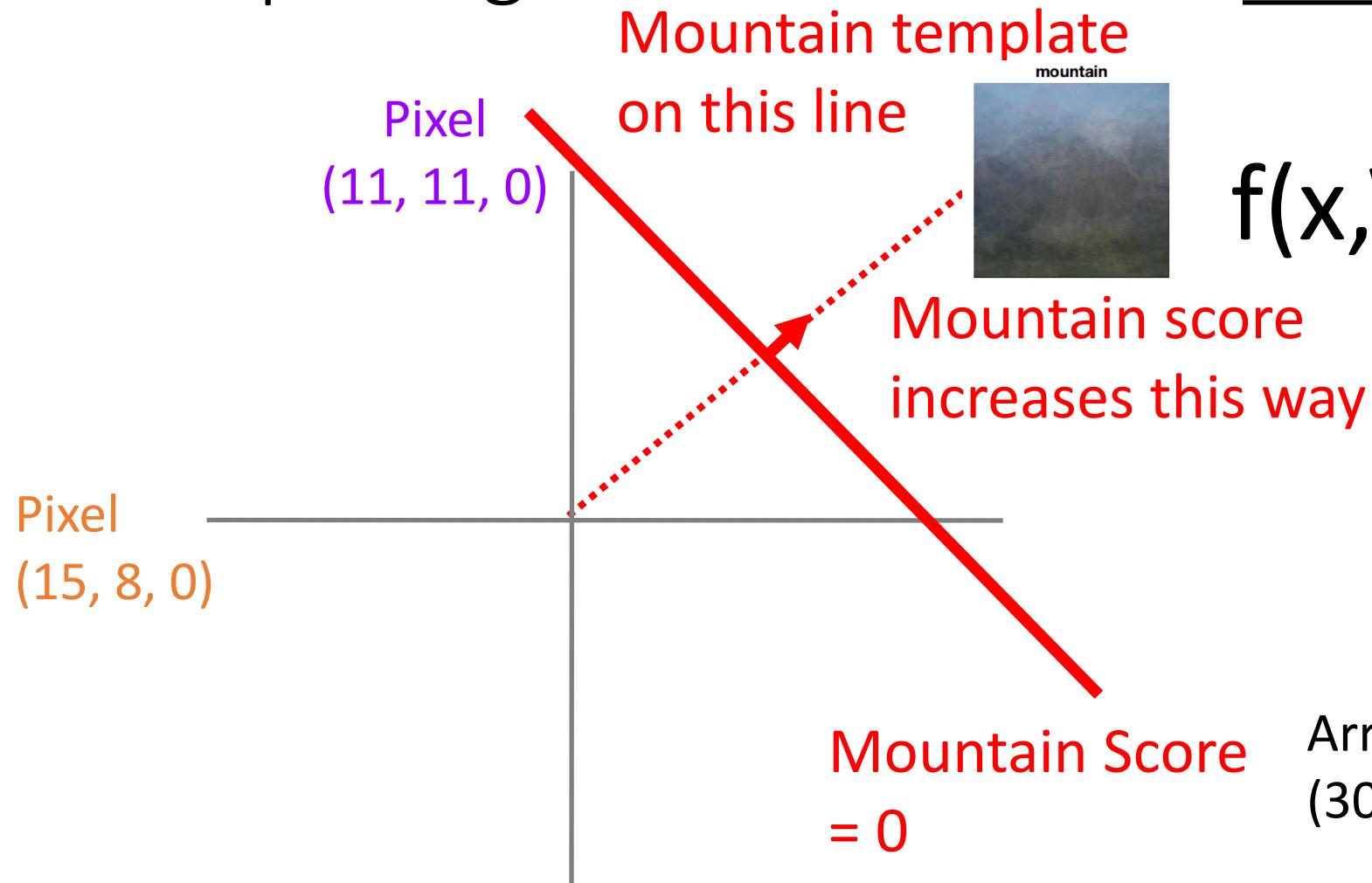


$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + b$$

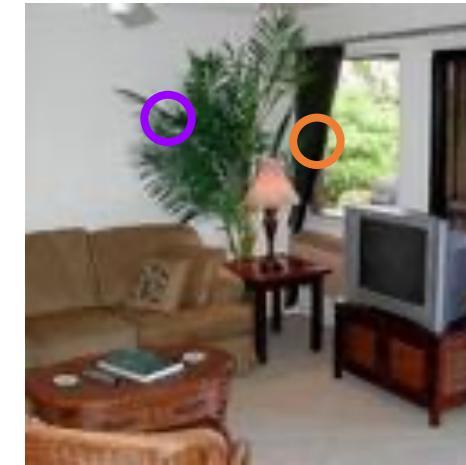


Array of **32x32x3** numbers
(3072 numbers total)

Interpreting a Linear Classifier: Geometric Viewpoint

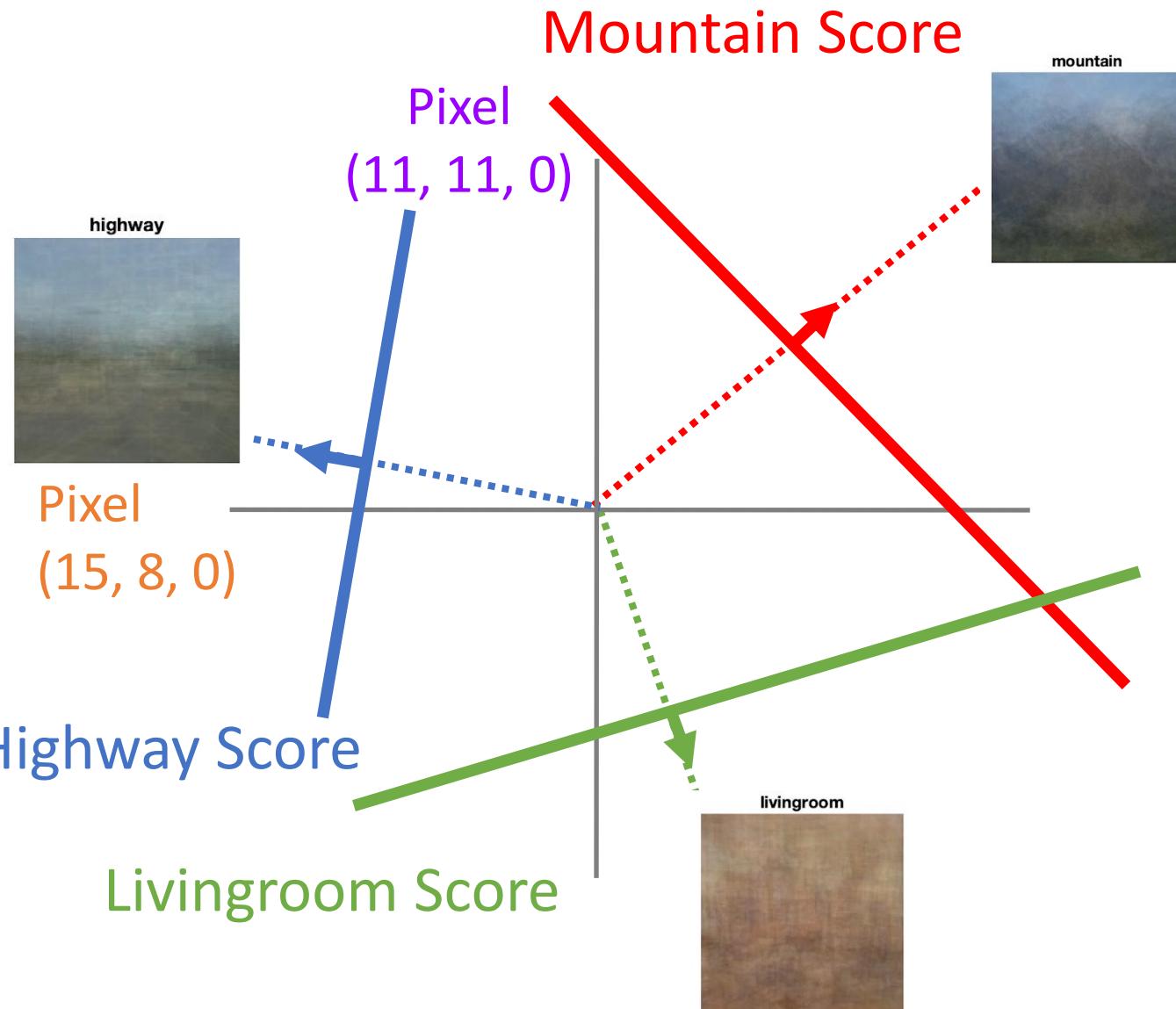


$$f(x, W) = Wx + b$$

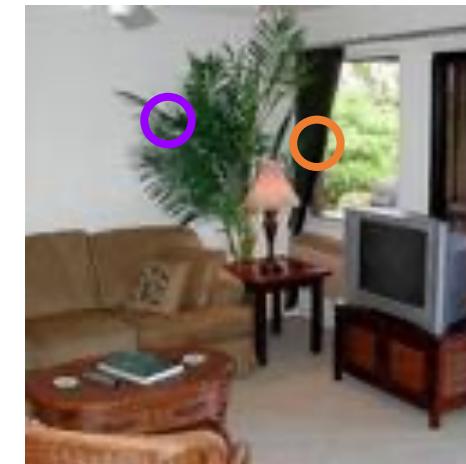


Array of **32x32x3** numbers
(3072 numbers total)

Interpreting a Linear Classifier: Geometric Viewpoint

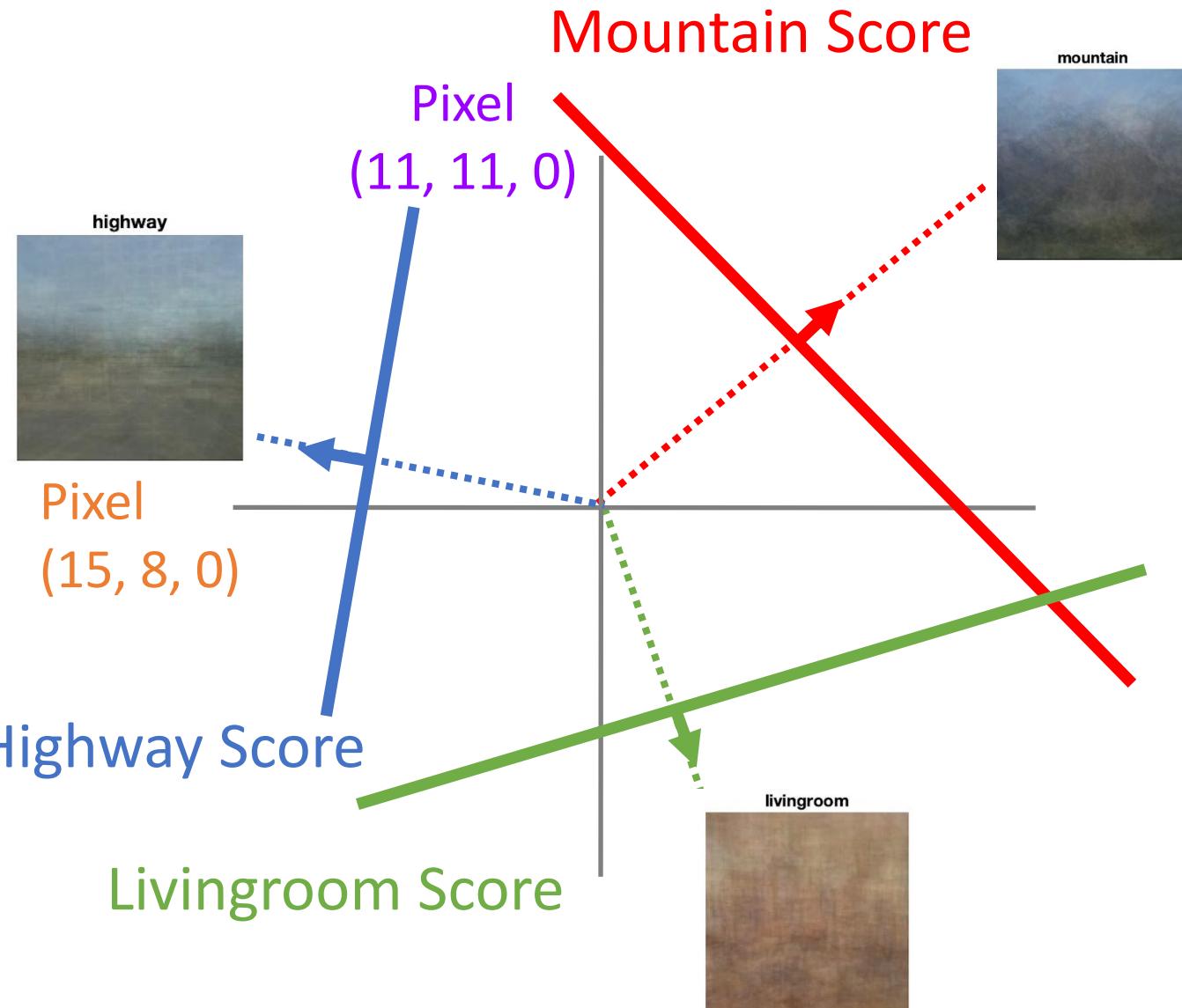


$$f(x, W) = Wx + b$$

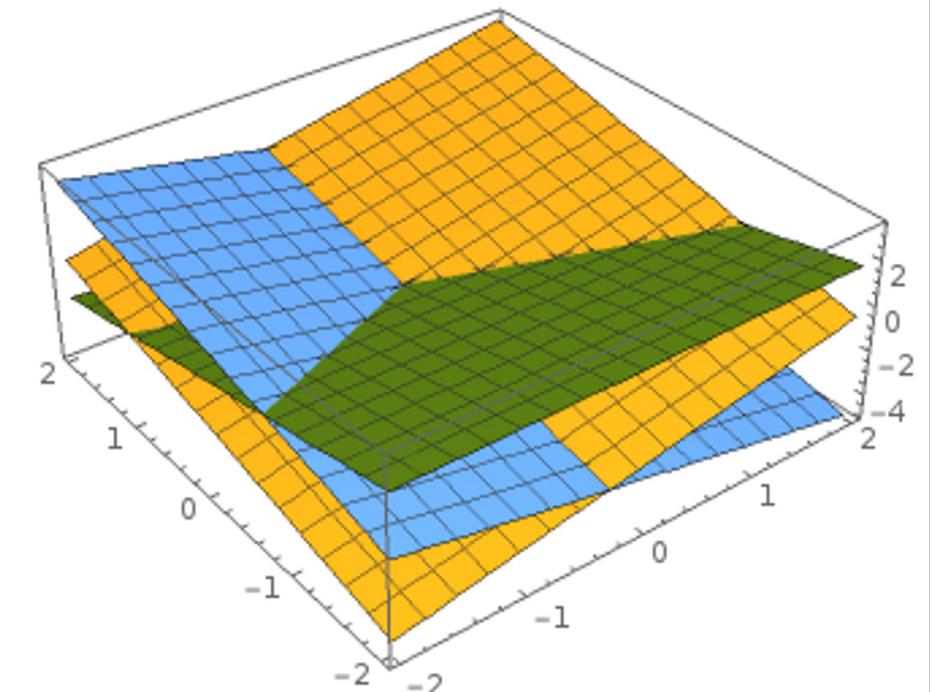


Array of **32x32x3** numbers
(3072 numbers total)

Interpreting a Linear Classifier: Geometric Viewpoint



Hyperplanes carving up a high-dimensional space



Plot created using [Wolfram Cloud](#)

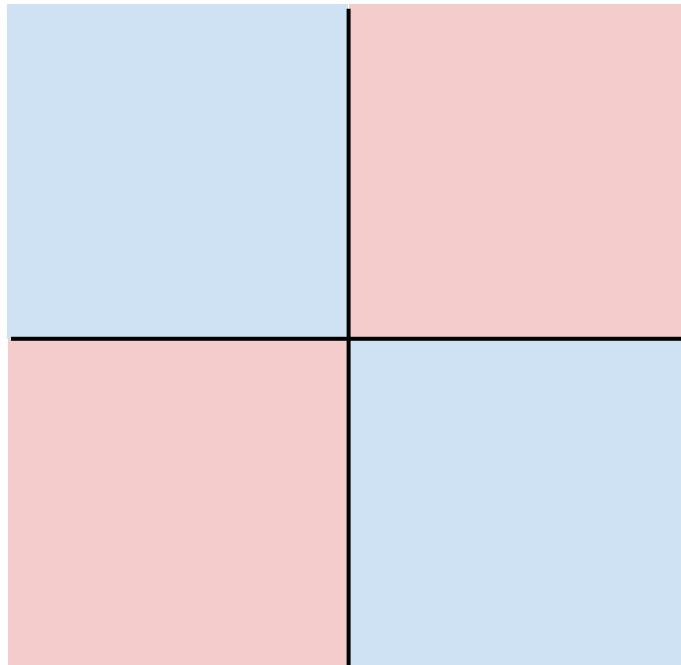
Hard Cases for a Linear Classifier

Class 1:

First and third quadrants

Class 2:

Second and fourth quadrants

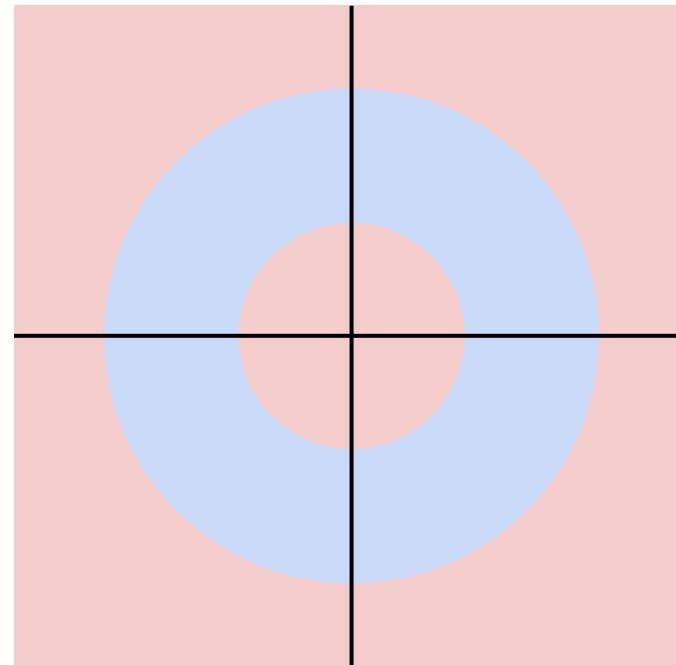


Class 1:

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

Class 2:

Everything else

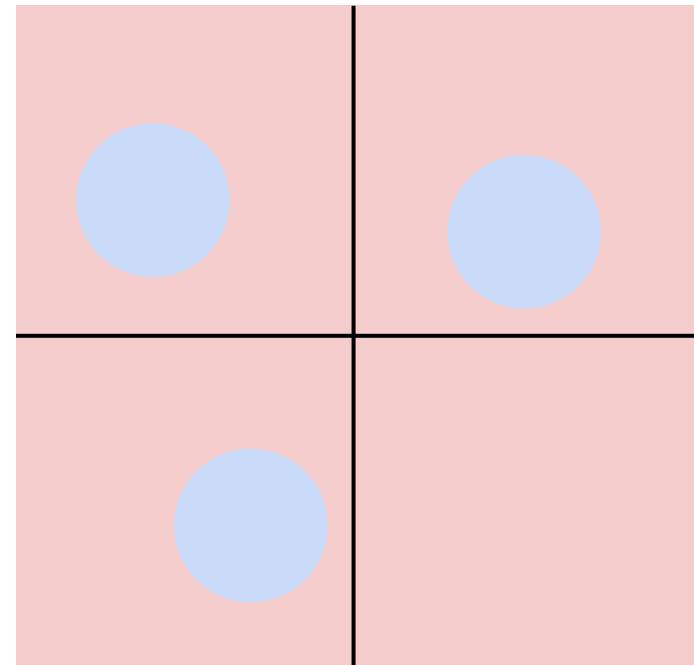


Class 1:

Three modes

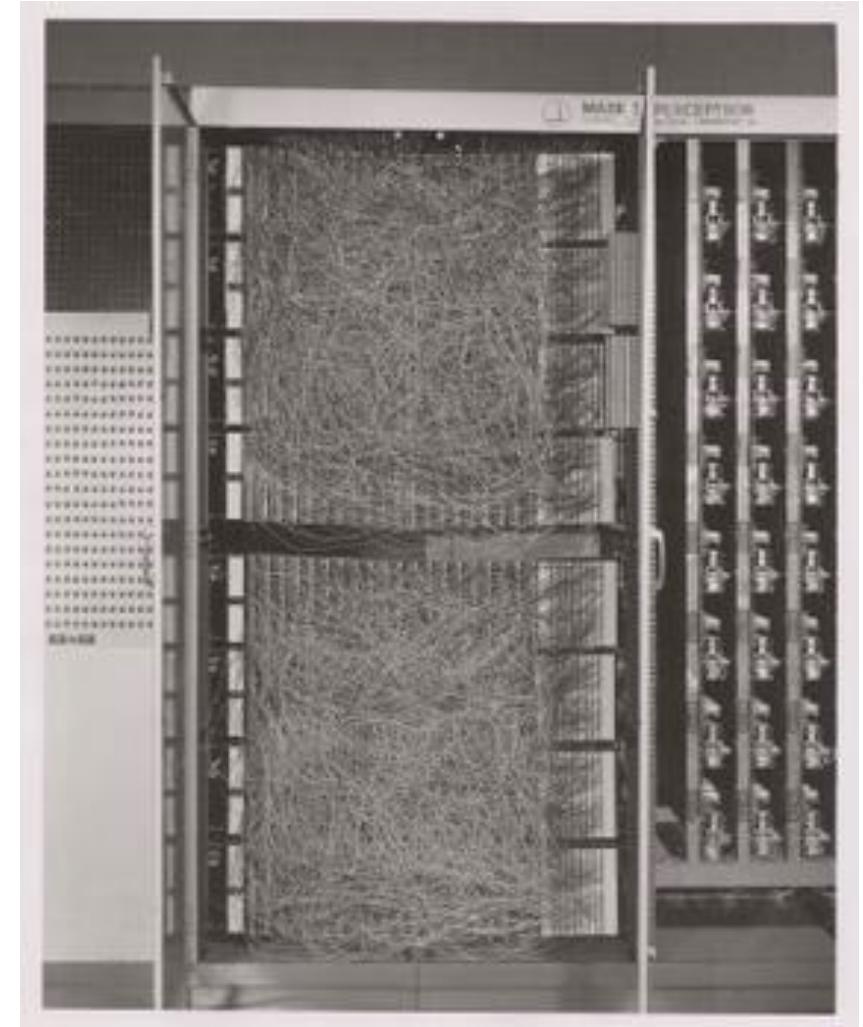
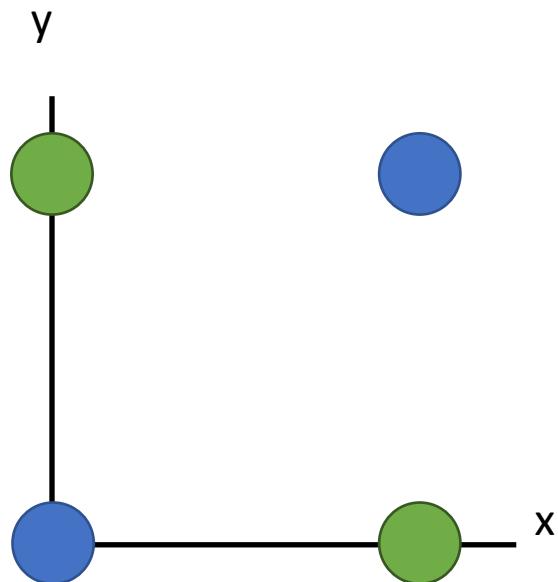
Class 2:

Everything else



Recall: Perceptron couldn't learn XOR

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



A live demo on training classifier

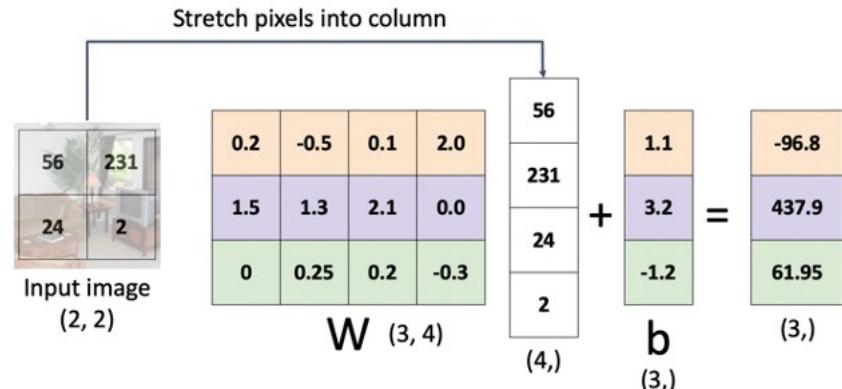
- <https://playground.tensorflow.org/>

Linear NN/Classifier: Three Viewpoints

Equation Viewpoint

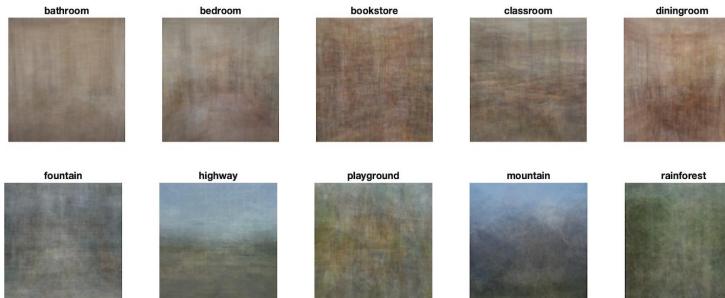
$$f(x, W) = Wx$$

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



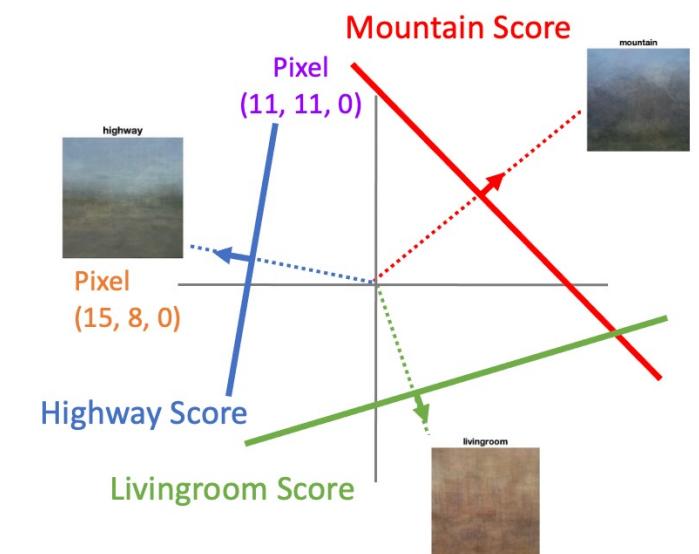
Visualization Viewpoint

One template per class



Geometric Viewpoint

Hyperplanes cutting up space



So Far: Defined a linear score function

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



Wrong prediction1

bathroom	-1.09	1.42	2.33
bedroom	0.03	0.29	-0.89
bookstore	0.55	0.20	0.10
classroom	1.10	1.59	-0.54
dining_room	1.54	-0.80	0.30
food_court	0.09	0.70	-0.60
kitchen	-1.49	0.84	0.49
lobby	-0.74	-0.24	0.74
living_room	2.01	0.22	1.71
office	2.35	-1.17	-0.19
baseball_field	-0.62	-1.15	-2.14
bridge	0.75	0.10	-0.84
campsite	-0.19	0.72	1.35
canyon	0.89	2.59	-1.07
coast	-0.76	-0.67	0.96
fountain	-1.40	2.76	0.12
highway	-1.42	-0.08	1.44
playground	0.49	-1.93	-1.96
mountain	-0.18	-0.44	-0.20
rainforest	-0.20	-1.79	-1.21

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

But how can we have a good W ?

Choosing a good W

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



bathroom	-1.09	1.42	2.33
bedroom	0.03	0.29	-0.89
bookstore	0.55	0.20	0.10
classroom	1.10	1.59	-0.54
dining_room	1.54	-0.80	0.30
food_court	0.09	0.70	-0.60
kitchen	-1.49	0.84	0.49
lobby	-0.74	-0.24	0.74
living_room	2.01	0.22	1.71
office	2.35	-1.17	-0.19
baseball_field	-0.62	-1.15	-2.14
bridge	0.75	0.10	-0.84
campsite	-0.19	0.72	1.35
canyon	0.89	2.59	-1.07
coast	-0.76	-0.67	0.96
fountain	-1.40	2.76	0.12
highway	-1.42	-0.08	1.44
playground	0.49	-1.93	-1.96
mountain	-0.18	-0.44	-0.20
rainforest	-0.20	-1.79	-1.21

Objective:

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

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

Given a dataset of examples

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

Where x_i is image and
 y_i is (integer) label

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

Given a dataset of examples

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

Where x_i is image and
 y_i is (integer) label

Loss for a single example is

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

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

Given a dataset of examples

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

Where x_i is image and
 y_i is (integer) label

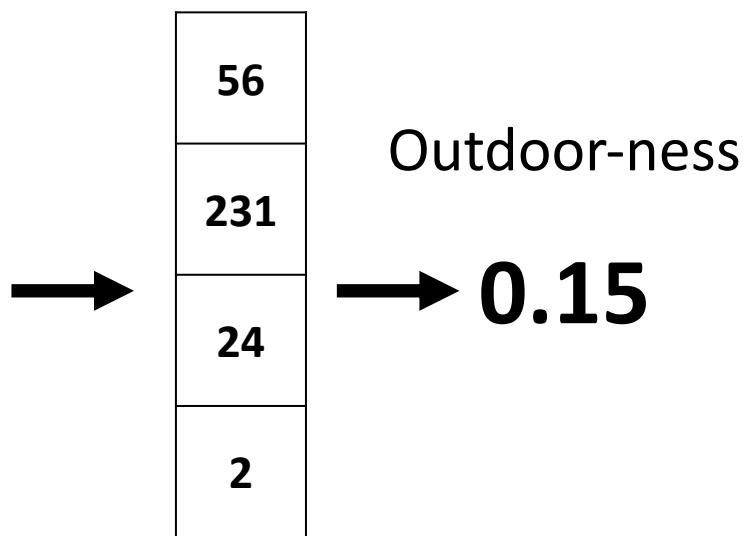
Loss for a single example is

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

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

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

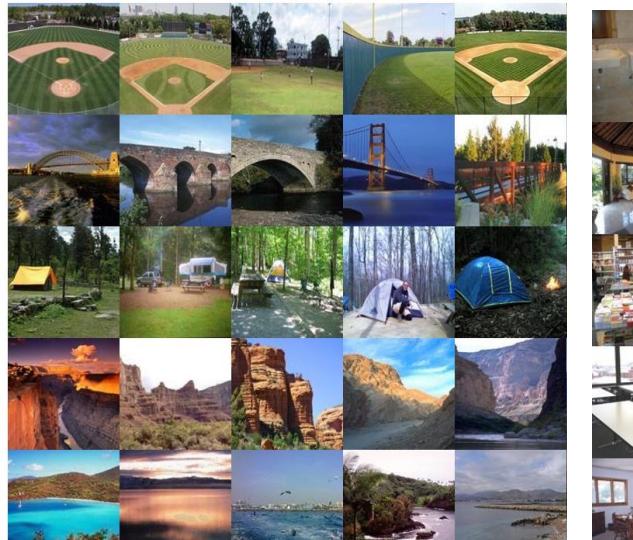
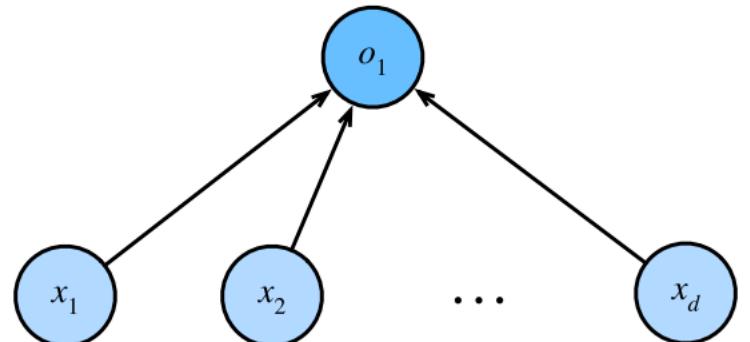
Linear Regression Loss



Output layer

Input layer

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



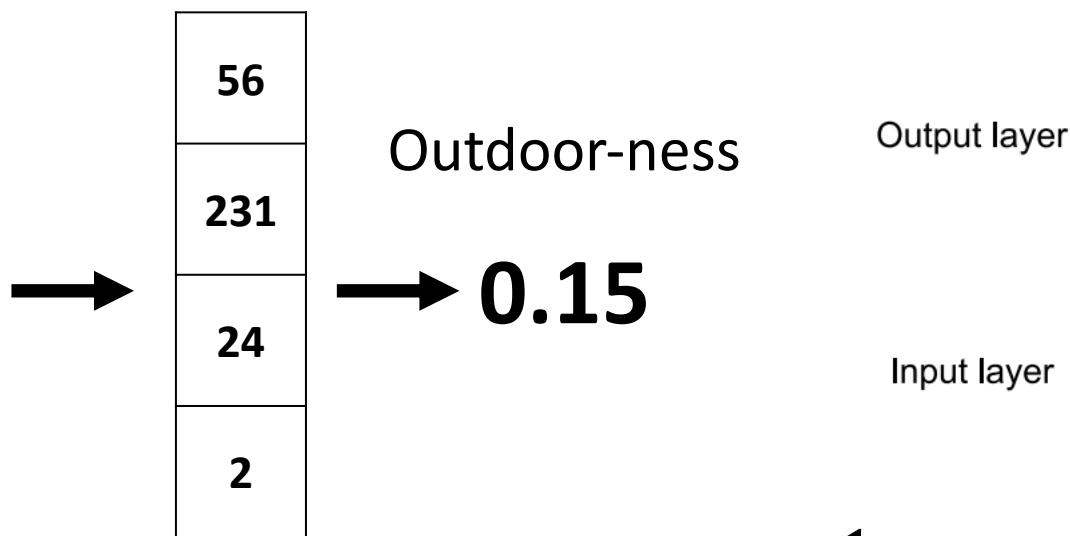
Outdoor images



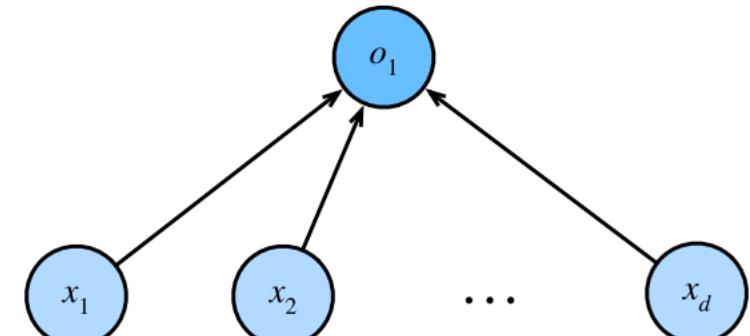
Indoor images

- Set the output for indoor images as 0 and for outdoor images as 1
- Categorical label is turned into a scalar value

Linear Regression Loss



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



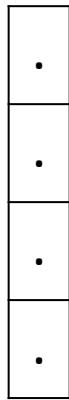
Squared loss for one sample: $L_i(W) = \frac{1}{2}(f(x_i, W) - y_i)^2$

Squared loss for all the training samples:

$$L(W) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2}(f(x_i, W) - y_i)^2$$

Learning objective: $W^* = \operatorname{argmin}_W L(W)$

Linear Regression Loss



Model output

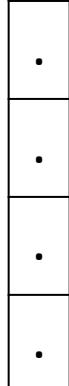
→ **0.15**

Loss

$$\frac{1}{2}(0.15 - 1)^2$$

Ground truth

1



→ **0.05**

$$\frac{1}{2}(0.05 - 0)^2$$

0

Analytic Solution for Linear Regression Loss

- You will have a practice in the Assignment 1 for outdoor and indoor classification with a linear regression loss

$$L(W) = \frac{1}{2n} \sum_{i=1}^n (y_i - Wx_i)^2$$

$$W^* = (X^T X)^{-1} X^T y, \text{ where } X \text{ is the batch of all the data}$$

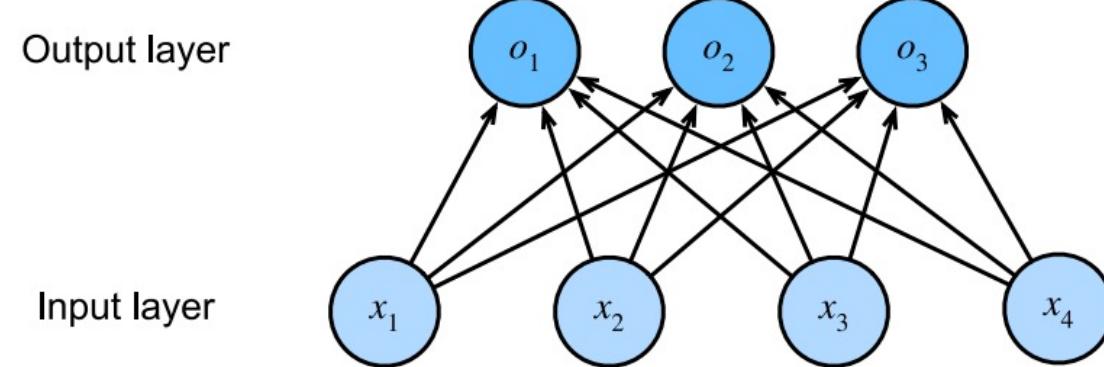
In which situation this close-form solution cannot be used directly?

Cross-Entropy Loss (Multinomial Logistic Regression)

3 scene classification



Bridge **3.2**



Mountain 5.1

Coast -1.7

Cross-Entropy Loss (Multinomial Logistic Regression)

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



Bridge **3.2**

Mountain 5.1

Coast -1.7

Cross-Entropy Loss (Multinomial Logistic Regression)

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



$$s = f(x_i; W) \quad P(Y = k | X = x_i) = \frac{\exp(s_k)}{\sum_j \exp(s_j)} \quad \text{Softmax function}$$

Bridge **3.2**

Mountain 5.1

Coast -1.7

Cross-Entropy Loss (Multinomial Logistic Regression)

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



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

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

Softmax
function

Bridge

3.2

5.1

Mountain

Coast

-1.7

Unnormalized log-
probabilities / logits

Cross-Entropy Loss (Multinomial Logistic Regression)



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

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

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

Softmax
function

Probabilities
must be ≥ 0

Bridge	3.2
Mountain	5.1
Coast	-1.7

Unnormalized log-
probabilities / logits

exp →

24.5
164.0
0.18

unnormalized
probabilities

Cross-Entropy Loss (Multinomial Logistic Regression)



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

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

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

Softmax
function

Probabilities
must be ≥ 0

Probabilities
must sum to 1

Bridge
Mountain
Coast

3.2
5.1
-1.7

Unnormalized log-
probabilities / logits

\exp

24.5
164.0
0.18

unnormalized
probabilities

normalize

0.13
0.87
0.00

probabilities

Cross-Entropy Loss (Multinomial Logistic Regression)



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

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

Probabilities
must be ≥ 0

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

Softmax
function

Probabilities
must sum to 1

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

Bridge
Mountain
Coast

3.2
5.1
-1.7

Unnormalized log-probabilities / logits

\exp

24.5
164.0
0.18

unnormalized probabilities

normalize

0.13
0.87
0.00

probabilities

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

Cross-Entropy Loss (Multinomial Logistic Regression)



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

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

Probabilities
must be ≥ 0

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

Softmax
function

Probabilities
must sum to 1

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

Bridge
Mountain
Coast

3.2
5.1
-1.7

Unnormalized log-probabilities / logits

\exp

24.5
164.0
0.18

unnormalized probabilities

normalize

0.13
0.87
0.00

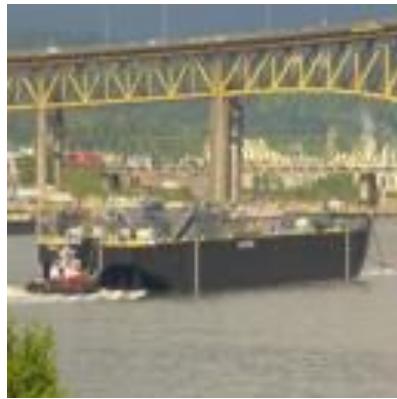
probabilities

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

Maximum Likelihood Estimation
Choose weights to maximize the likelihood of the observed data

Cross-Entropy Loss (Multinomial Logistic Regression)

One-hot encoding for the categorical label of each sample



Bridge	1	1	0	0
--------	---	---	---	---

Mountain	0	0	1	0
----------	---	---	---	---

Coast	0	0	0	1
-------	---	---	---	---

Cross-Entropy Loss (Multinomial Logistic Regression)



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

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

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

Softmax
function

Probabilities
must be ≥ 0

Probabilities
must sum to 1

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

Bridge

3.2

\rightarrow
 \exp

5.1

Mountain

-1.7

Unnormalized log-
probabilities / logits

24.5

\rightarrow
normalize

164.0

unnormalized
probabilities

0.13

0.87

0.00

probabilities

Compare \leftarrow

1.00

0.00

0.00

Correct
probs

Cross-Entropy Loss (Multinomial Logistic Regression)



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

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

Probabilities
must be ≥ 0

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

Softmax
function

Probabilities
must sum to 1

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

Bridge

3.2

\rightarrow
 \exp

5.1

24.5

normalize
 \rightarrow

164.0

0.13

Compare \leftarrow

-1.7

0.18

Unnormalized log-
probabilities / logits

unnormalized
probabilities

0.87

0.00

probabilities

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

1.00

0.00

0.00

Correct
probs

Cross-Entropy Loss (Multinomial Logistic Regression)



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

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

Probabilities
must be ≥ 0

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

Softmax
function

Probabilities
must sum to 1

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

Bridge

3.2
5.1
-1.7

Unnormalized log-
probabilities / logits

\exp

24.5
164.0
0.18

unnormalized
probabilities

normalize

0.13
0.87
0.00

probabilities

$$H(P, Q) = -\sum p(x) \log q(x)$$
$$H(P) + D_{KL}(P || Q)$$

Correct
probs

Then $H(P)=0$, and $D_{KL}(P || Q)$ will have just one term for the correct class

Cross-Entropy Loss (Multinomial Logistic Regression)



Bridge

3.2

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

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

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

Softmax
function

Maximize probability of correct class

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

Putting it all together:

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

Mountain 5.1

Coast -1.7

Cross-Entropy Loss (Multinomial Logistic Regression)



Bridge

3.2

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

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

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

Softmax
function

Maximize probability of correct class

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

Putting it all together:

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

Mountain

5.1

Coast

-1.7

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

Cross-Entropy Loss (Multinomial Logistic Regression)



Bridge	3.2
Mountain	5.1
Coast	-1.7

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

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

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

Softmax
function

Maximize probability of correct class

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

Putting it all together:

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

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

A: min 0, max +infinity

Cross-Entropy Loss (Multinomial Logistic Regression)



Bridge **3.2**

Mountain 5.1

Coast -1.7

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

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

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

Softmax
function

Maximize probability of correct class

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

Putting it all together:

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

Q: If all scores are uniformly random, what is the loss?

Cross-Entropy Loss (Multinomial Logistic Regression)



Bridge **3.2**

Mountain 5.1

Coast -1.7

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

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

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

Softmax
function

Maximize probability of correct class

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

Putting it all together:

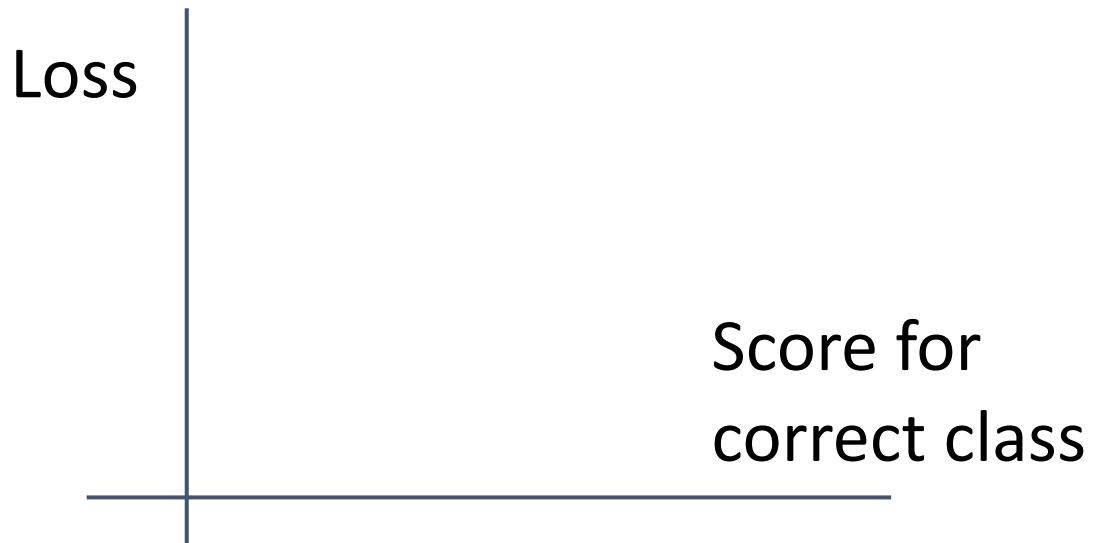
$$L_i = -\log \left(\frac{\exp(s_{y_i})}{\sum_j \exp(s_j)} \right)$$

Q: If all scores are uniformly random, what is the loss?

A: $-\log(1/C)$
 $\log(10) \approx 2.3$

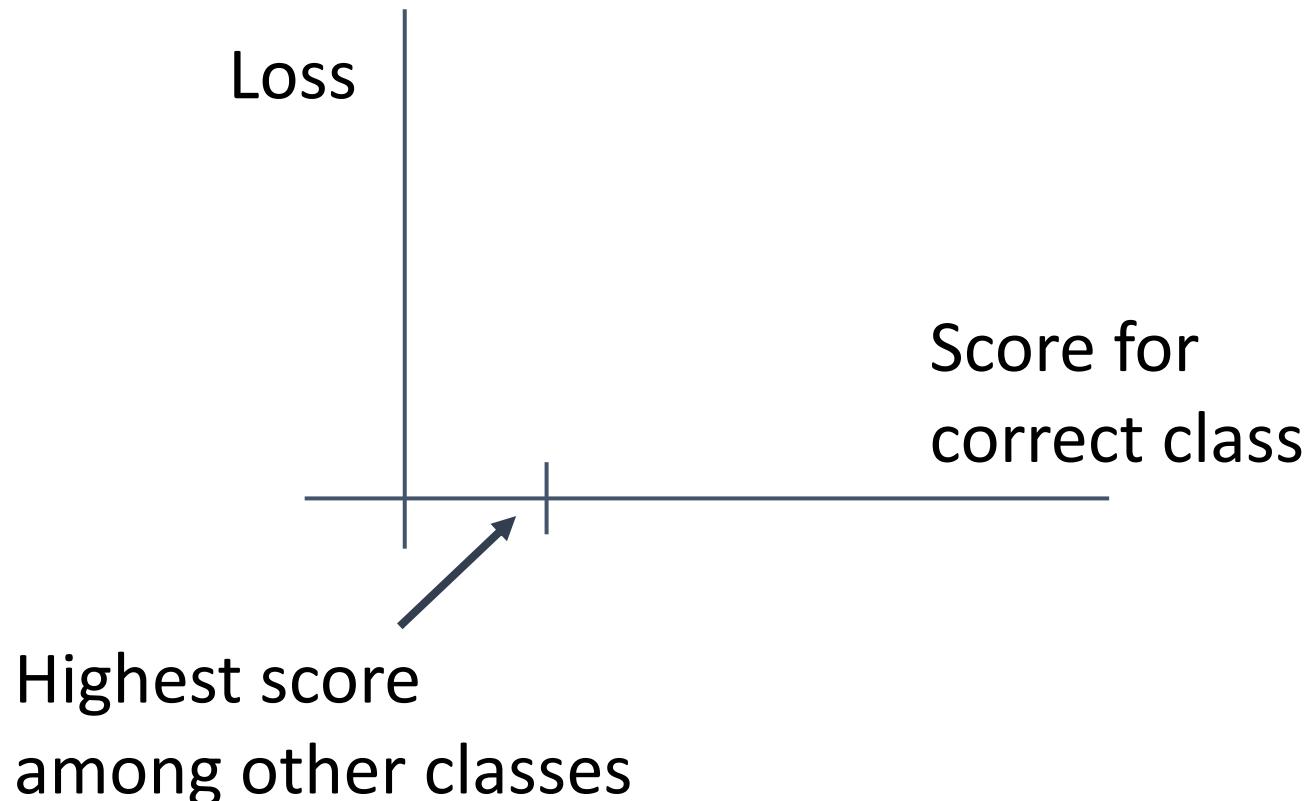
Multiclass SVM Loss

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



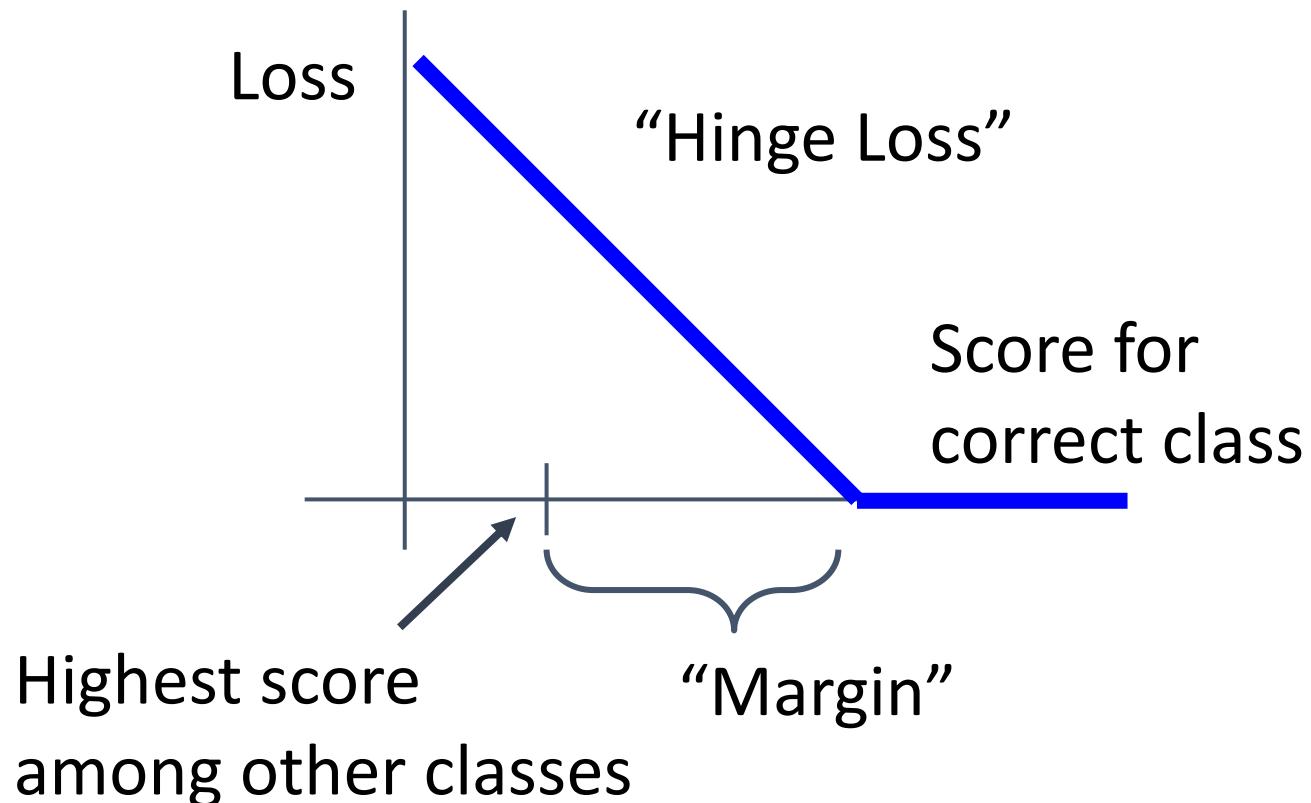
Multiclass SVM Loss

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



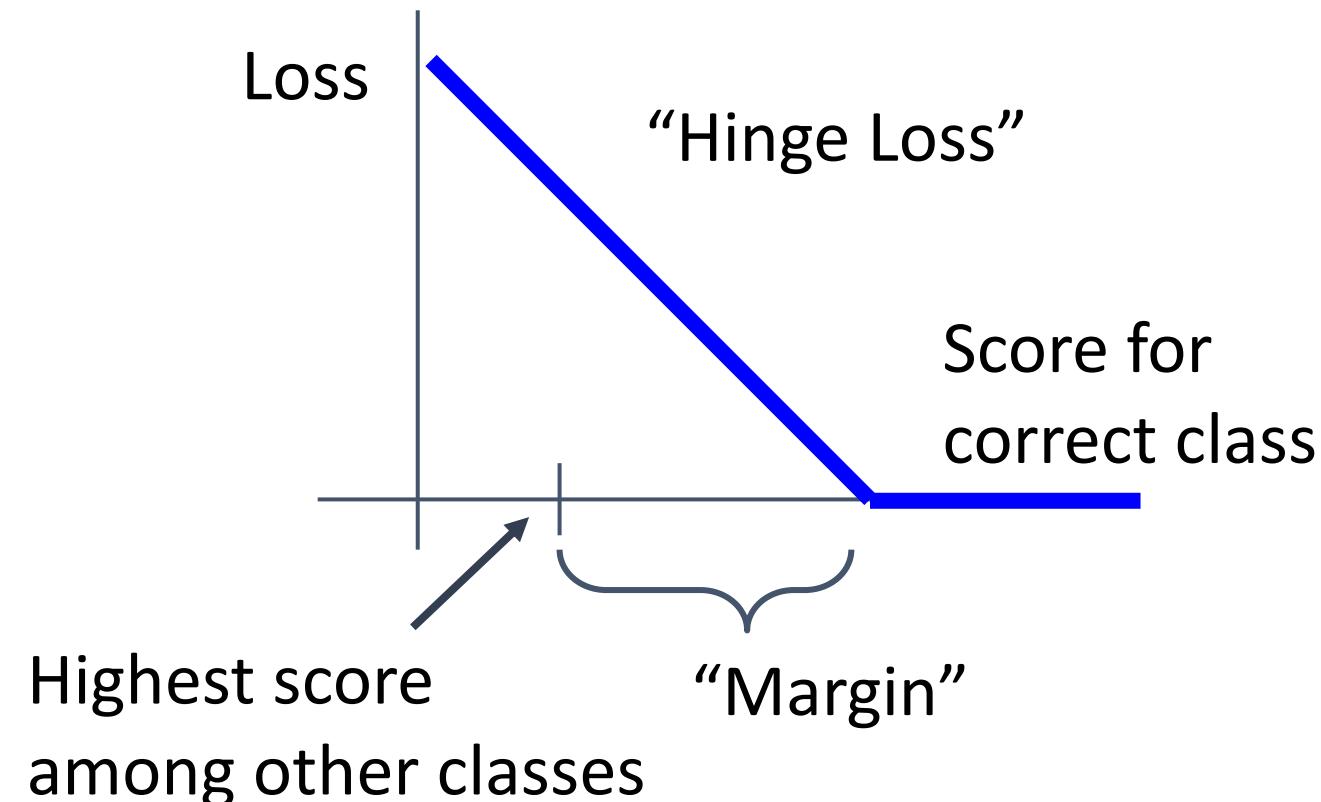
Multiclass SVM Loss

“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



Bridge	3.2	1.3	2.2
Mountain	5.1	4.9	2.5
Coast	-1.7	2.0	-3.1

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



Bridge	3.2	1.3	2.2
Mountain	5.1	4.9	2.5
Coast	-1.7	2.0	-3.1
Loss	2.9		

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:

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

Multiclass SVM Loss



Bridge	3.2	1.3	2.2
Mountain	5.1	4.9	2.5
Coast	-1.7	2.0	-3.1
Loss	2.9	0	

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:

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

Multiclass SVM Loss



Bridge	3.2	1.3	2.2
Mountain	5.1	4.9	2.5
Coast	-1.7	2.0	-3.1
Loss	2.9	0	12.9

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:

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

Multiclass SVM Loss



Bridge	3.2	1.3	2.2
Mountain	5.1	4.9	2.5
Coast	-1.7	2.0	-3.1
Loss	2.9	0	12.9

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:

$$\begin{aligned} L &= (2.9 + 0.0 + 12.9) / 3 \\ &= 5.27 \end{aligned}$$

Cross-Entropy vs SVM Loss

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

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and

$$y_i = 0$$

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

Cross-Entropy vs SVM Loss

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

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and

$$y_i = 0$$

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

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

Cross-Entropy vs SVM Loss

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

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and

$$y_i = 0$$

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

Cross-Entropy vs SVM Loss

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

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and

$$y_i = 0$$

Q: 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

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

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and

$$y_i = 0$$

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

Cross-Entropy vs SVM Loss

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

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and

$$y_i = 0$$

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

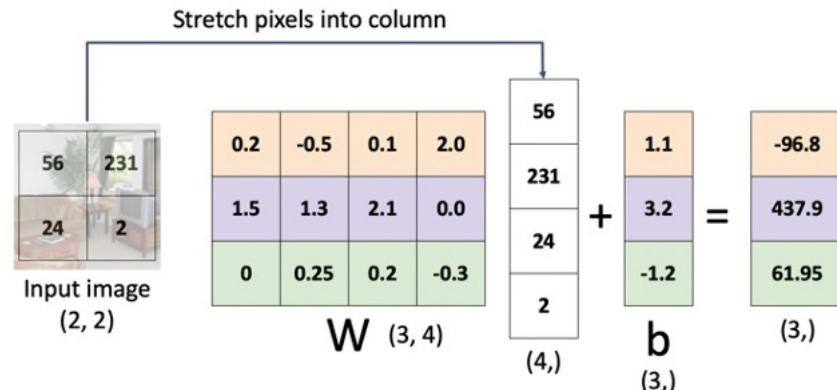
A: Cross-entropy loss will decrease,
SVM loss still 0

Recap: Linear NN/Classifier: Three Viewpoints

Equation Viewpoint

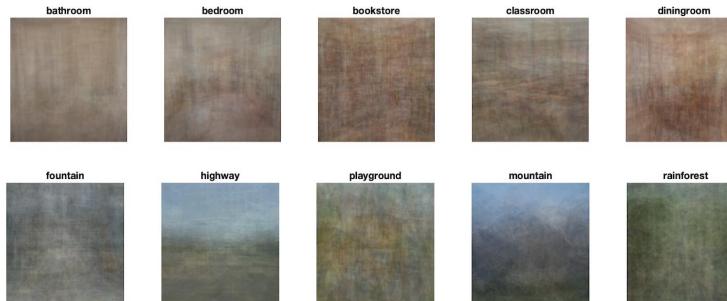
$$f(x, W) = Wx$$

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



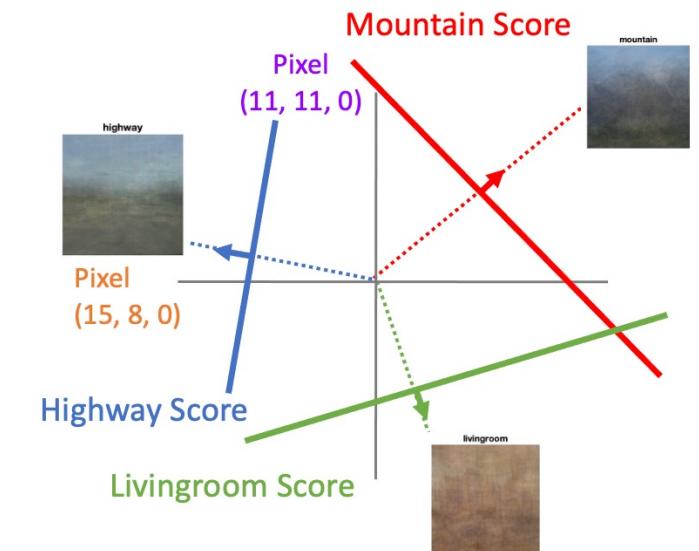
Visualization Viewpoint

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**: $s = f(x; W, b) = Wx + b$
- We have a **loss function**:

Linear regression loss:

$$L_i(W) = \frac{1}{2} (s_i - y_i)^2$$

Softmax loss:

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

Linear classifier

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

Analytical solution:

$$w^* = (X^T X)^{-1} X^T y$$

SVM loss:

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

But after we compute the loss,
how do we find the best W ?

Reading and Practice

- Textbook: D2L: <https://d2l.ai/d2l-en-pytorch.pdf>
- Go through Chapter 3. Linear Neural Networks
- Go through the notebooks for that chapter: https://github.com/d2l-ai/d2l-pytorch-colab/tree/master/chapter_linear-classification

```
def show_images(imgs, num_rows, num_cols, titles=None, scale=1.5): #save
    """Plot a list of images."""
    figsize = (num_cols * scale, num_rows * scale)
    _, axes = d2l.plt.subplots(num_rows, num_cols, figsize=figsize)
    axes = axes.flatten()
    for i, (ax, img) in enumerate(zip(axes, imgs)):
        if torch.is_tensor(img):
            # Tensor Image
            ax.imshow(img.numpy())
        else:
            # PIL Image
            ax.imshow(img)
        ax.axes.get_xaxis().set_visible(False)
        ax.axes.get_yaxis().set_visible(False)
        if titles:
            ax.set_title(titles[i])
    return axes
```

Here are the images and their corresponding labels (in text) for the first few examples in the training dataset.

```
X, y = next(iter(data.DataLoader(mnist_train, batch_size=18)))
show_images(X.reshape(18, 28, 28), 2, 9, titles=get_fashion_mnist_labels(y));
```

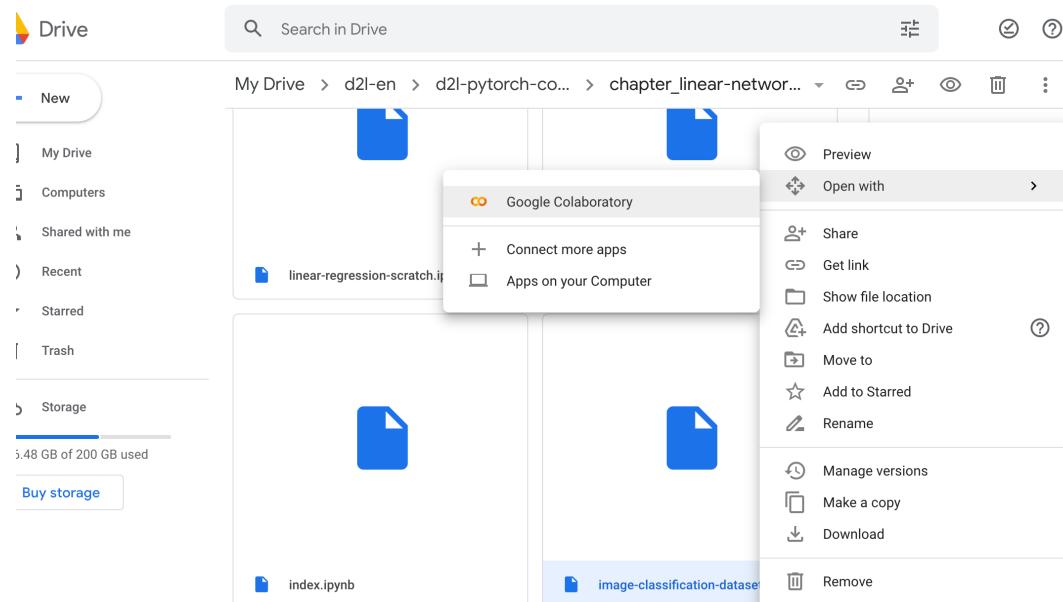


Two ways of running D2L Jupyter notebook

- Run on local machine/desktop/labtop
- https://d2l.ai/chapter_installation/index.html

Two ways of running D2L Jupyter notebook

- Run one Google Colab:
 - Download <https://github.com/d2l-ai/d2l-pytorch-colab>
 - Upload the notebook files to your Google Drive
 - In Google Drive d2l-pytorch-colab folder, click the notebook, open with Google Colab



Next time:
Optimization Algorithms

