



DeepRob

Lecture 3

Linear Classifiers

University of Michigan and University of Minnesota



Project 0

- Instructions and code available on the website
- Here: <https://rpm-lab.github.io/CSCI5980-Spr23-DeepRob/projects/project0/>
- **Due tonight! January 24th, 11:59 PM CT**

Project 1

- Instructions and code will be available on the website **today**.
- Classification using K-Nearest Neighbors and Linear Models

Calendar	
Week 1	
Jan 17:	LEC 1 Course Introduction PROJECT 0 DUE UMich discussion link - Intro to Python, Pytorch, and Colab
Jan 19:	LEC 2 Image Classification
Week 2	
Jan 24:	LEC 3 Linear Classifiers PROJECT 0 DUE PROJECT 1 DUE
Jan 28:	LEC 4 Regularization + Optimization
Week 3	
Jan 31:	LEC 5 Neural Networks
Feb 02:	LEC 6 Backpropagation

 We're here!

Gradescope Quizzes

- Quiz links will be published at 7am on the day of lecture.
 - **This will start from next lecture on 01/26**
- The quiz will close before that day's lecture time i.e. 2:30pm
- Time limit of 15 min once quiz is opened
- Covers material from previous lectures and graded projects

Recap: Image Classification—A Core Computer Vision Task

Input: image



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

Chocolate Pretzels

Granola Bar

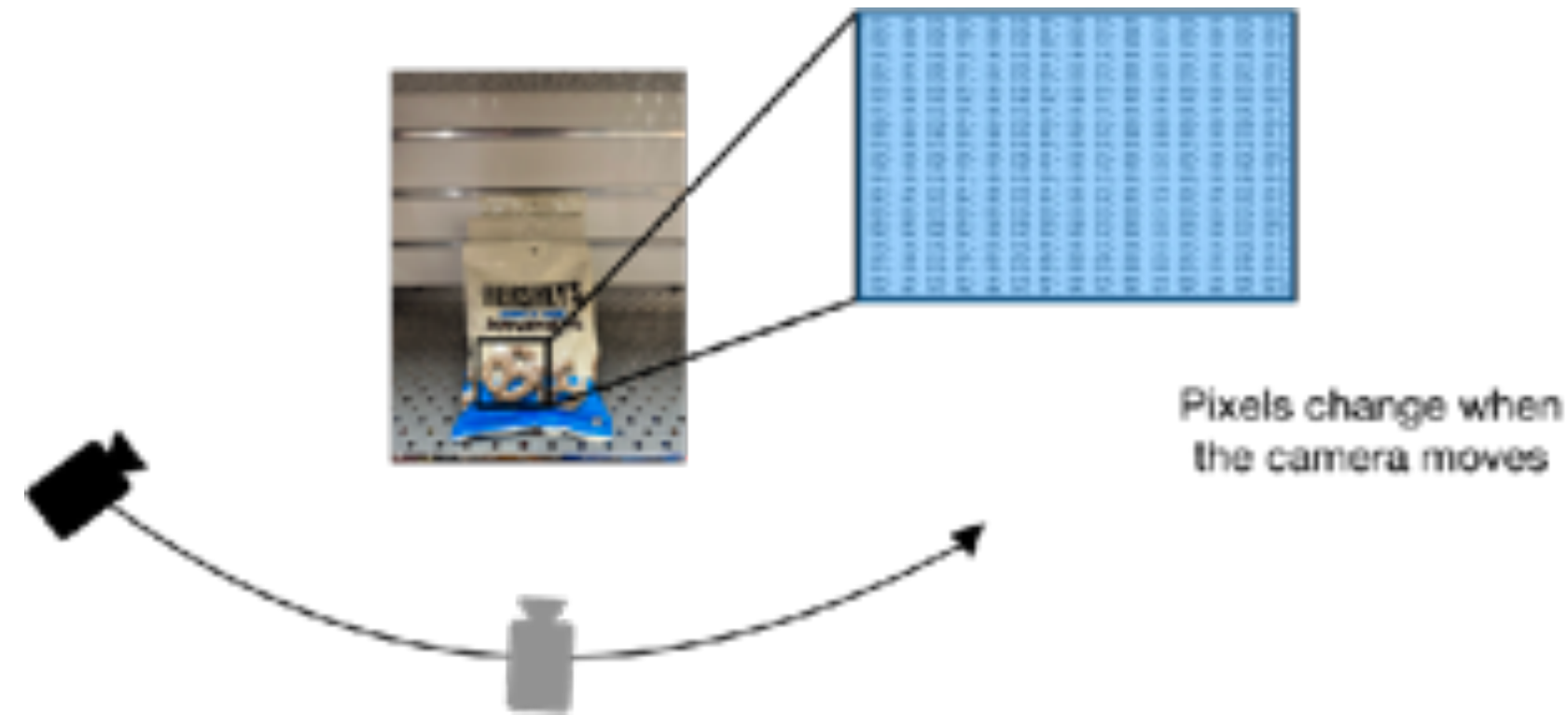
Potato Chips

Water Bottle

Popcorn

Recap: Image Classification Challenges

Viewpoint Variation & Semantic Gap



Illumination Changes



Intraclass Variation



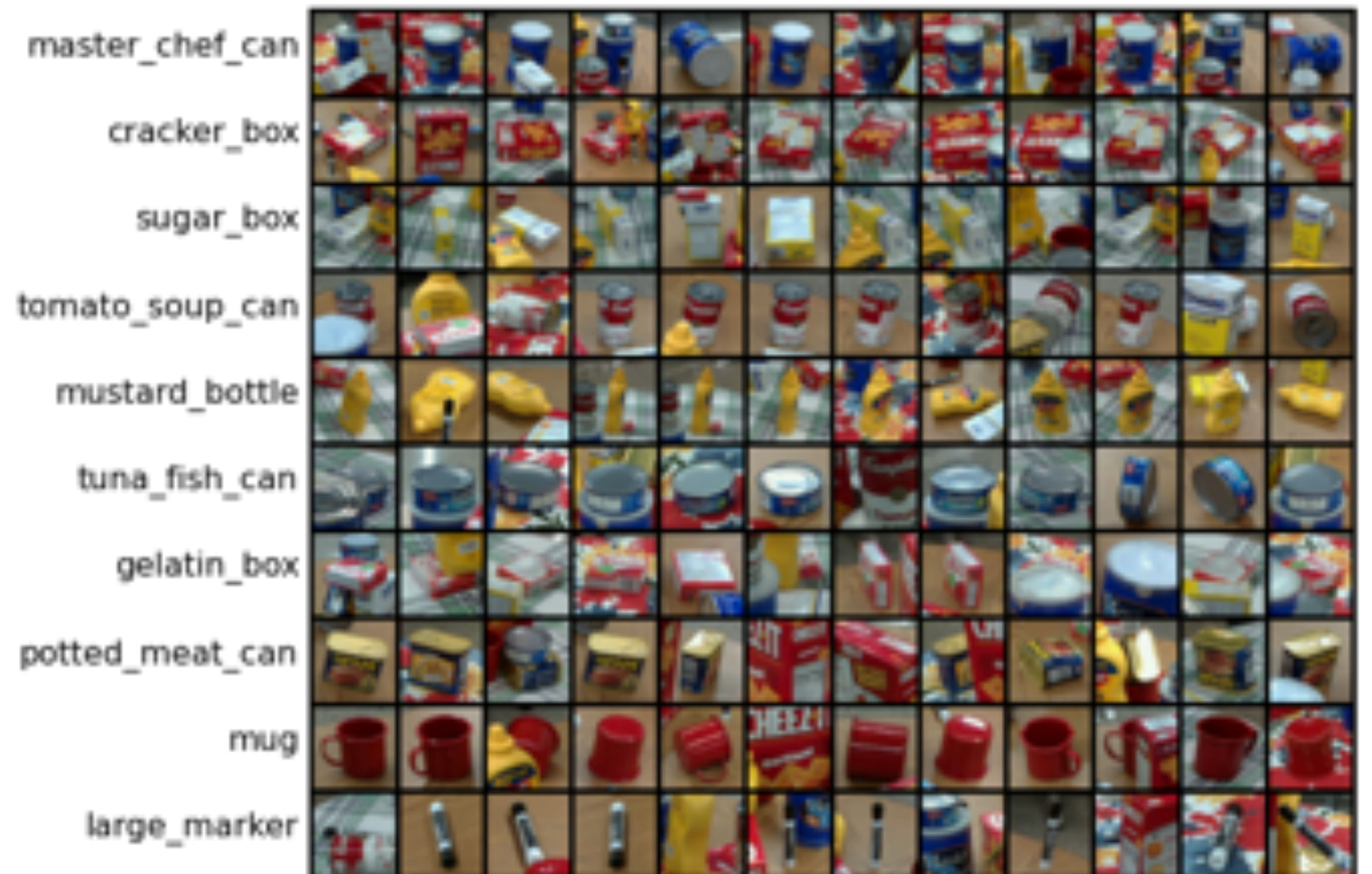
Recap: Machine Learning—Data-Driven Approach

1. Collect a dataset of images and labels
2. Use Machine Learning to train a classifier
3. Evaluate the classifier on new images

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

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

Example training set





Linear Classifiers



Building Block of Neural Networks

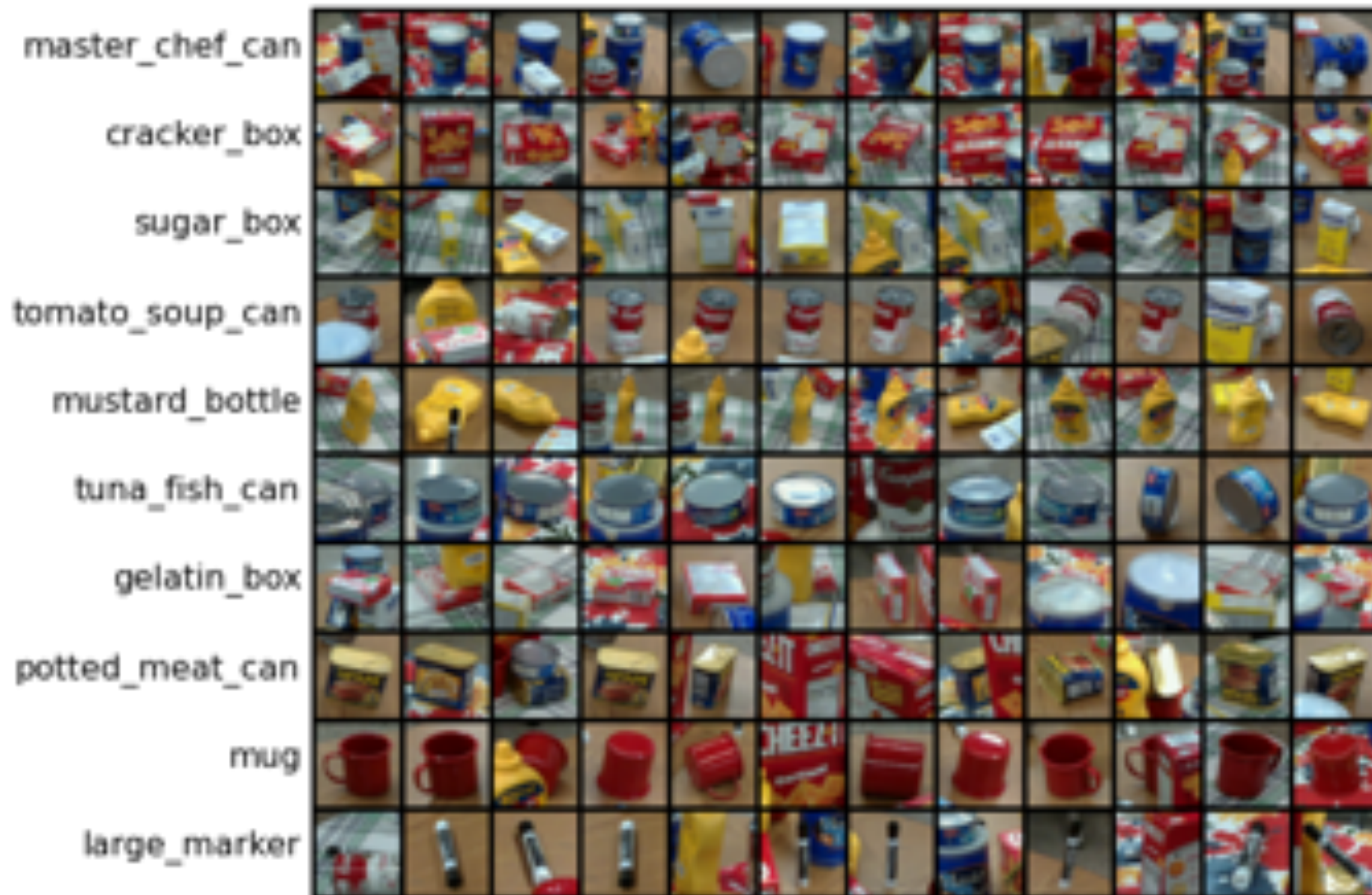
Linear
classifiers



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

Recall PROPS

Progress Robot Object Perception Samples Dataset



10 classes

32x32 RGB images

50k training images (5k per class)

10k test images (1k per class)

A video link will be posted to the website today discussing about PROPS dataset that you will use for P1

Chen et al., "ProgressLabeller: Visual Data Stream Annotation for Training Object-Centric 3D Perception", IROS, 2022.

Parametric Approach

Image



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



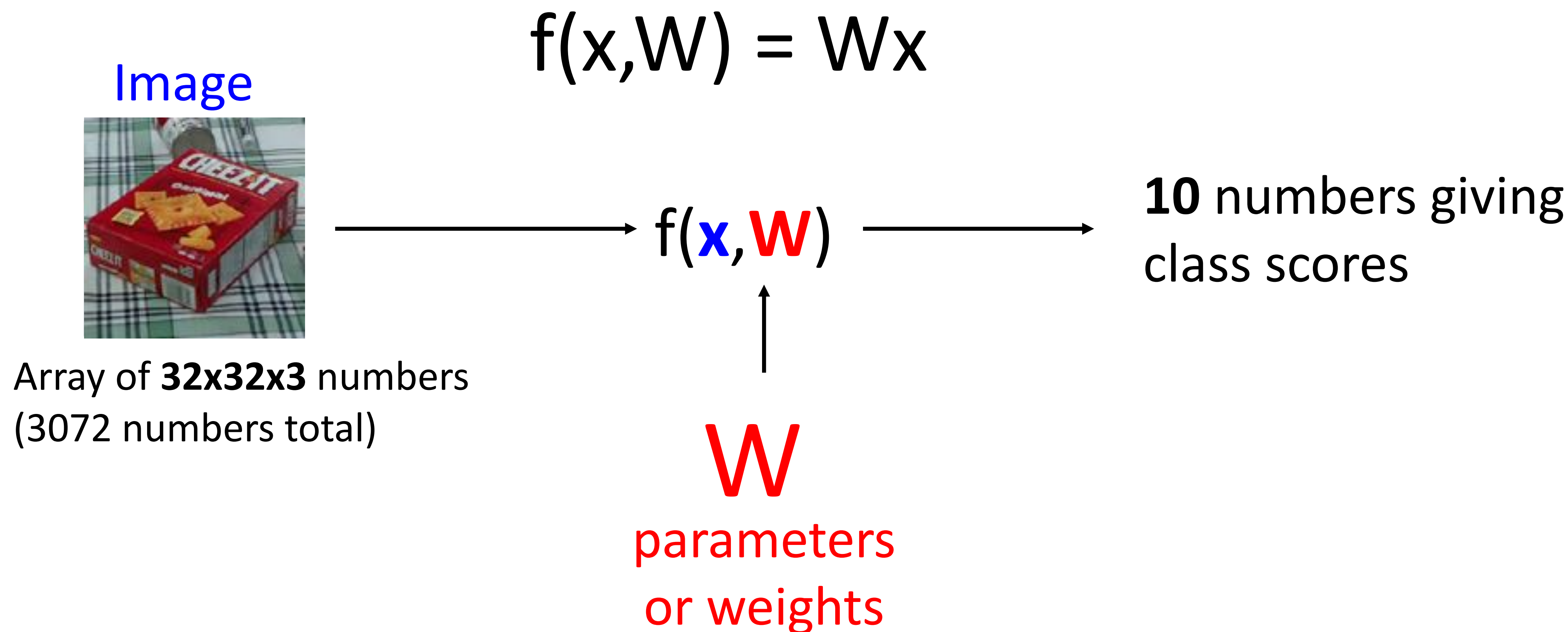
10 numbers giving
class scores



W

parameters
or weights

Parametric Approach—Linear Classifier



Parametric Approach—Linear Classifier

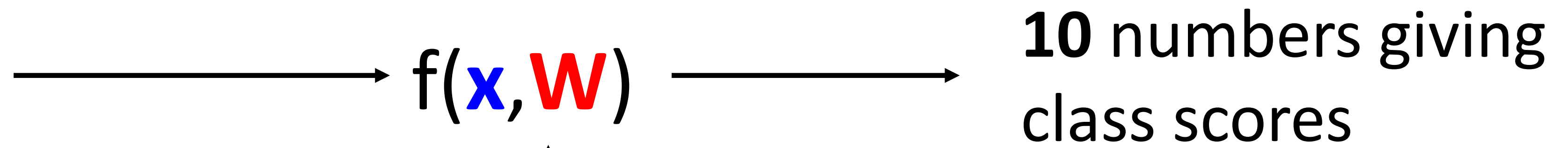
Image



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

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

(10,) (10, 3072) (3072,)



W

parameters
or weights



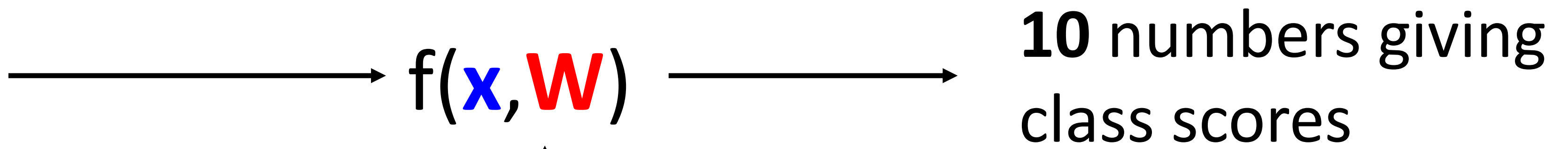
Parametric Approach—Linear Classifier

Image



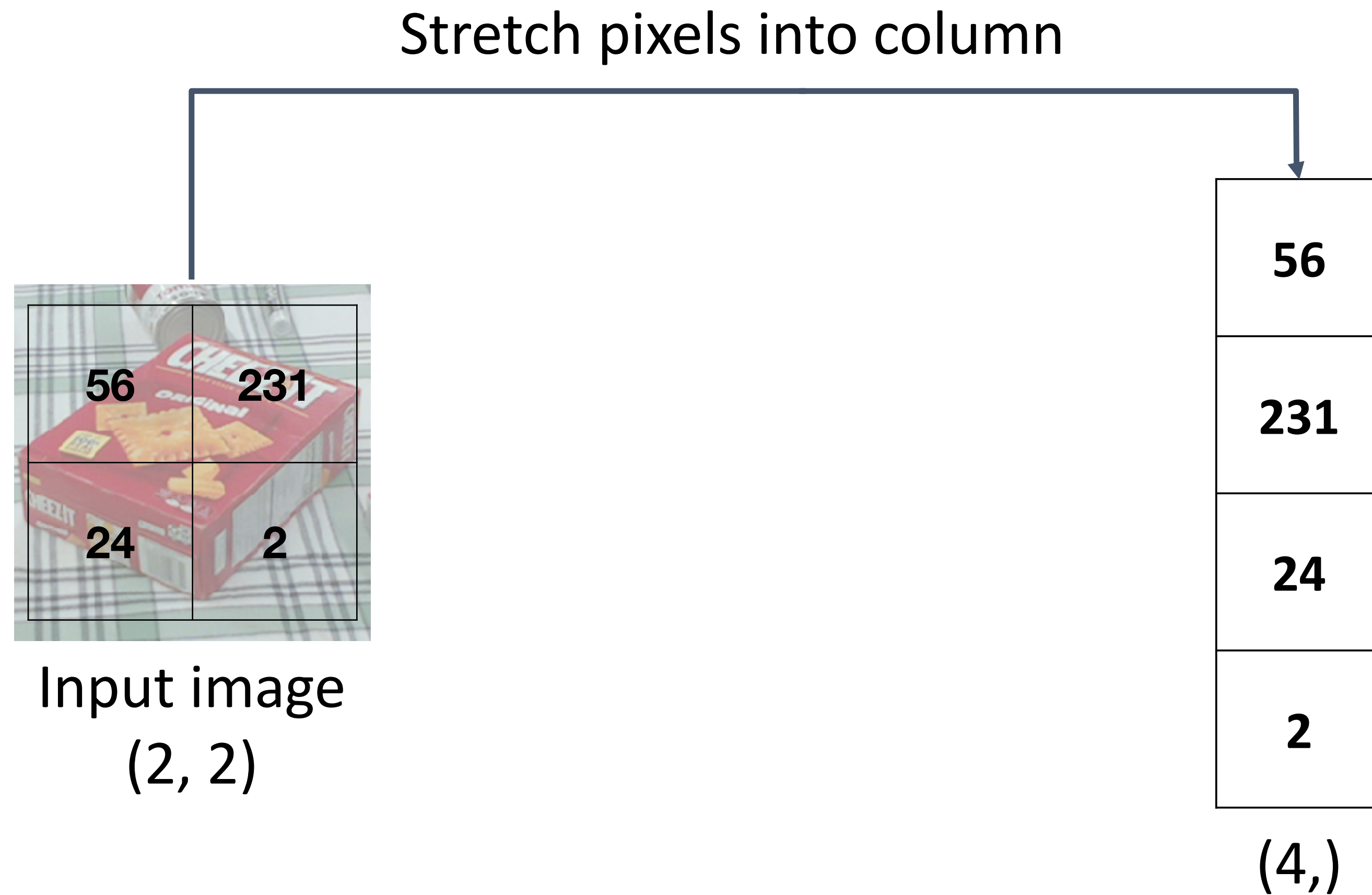
Array of **32x32x3** numbers
(3072 numbers total)

$$\boxed{f(x,W)}_{(10,)} = \boxed{W}_{(10, 3072)} \boxed{x}_{(3072,)} + \boxed{b}_{(10,)}$$



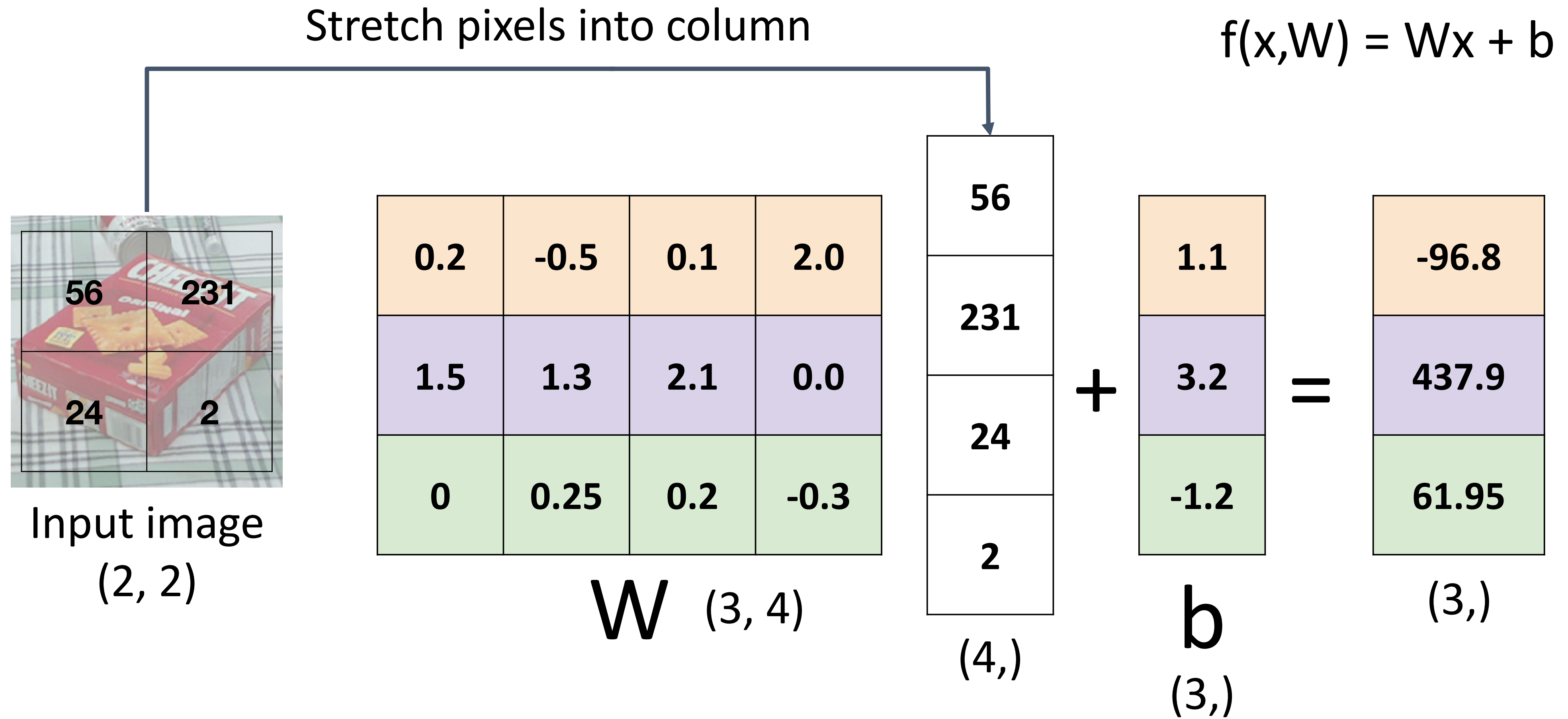
W
parameters
or weights

Example for 2x2 Image, 3 classes (crackers/mug/sugar)

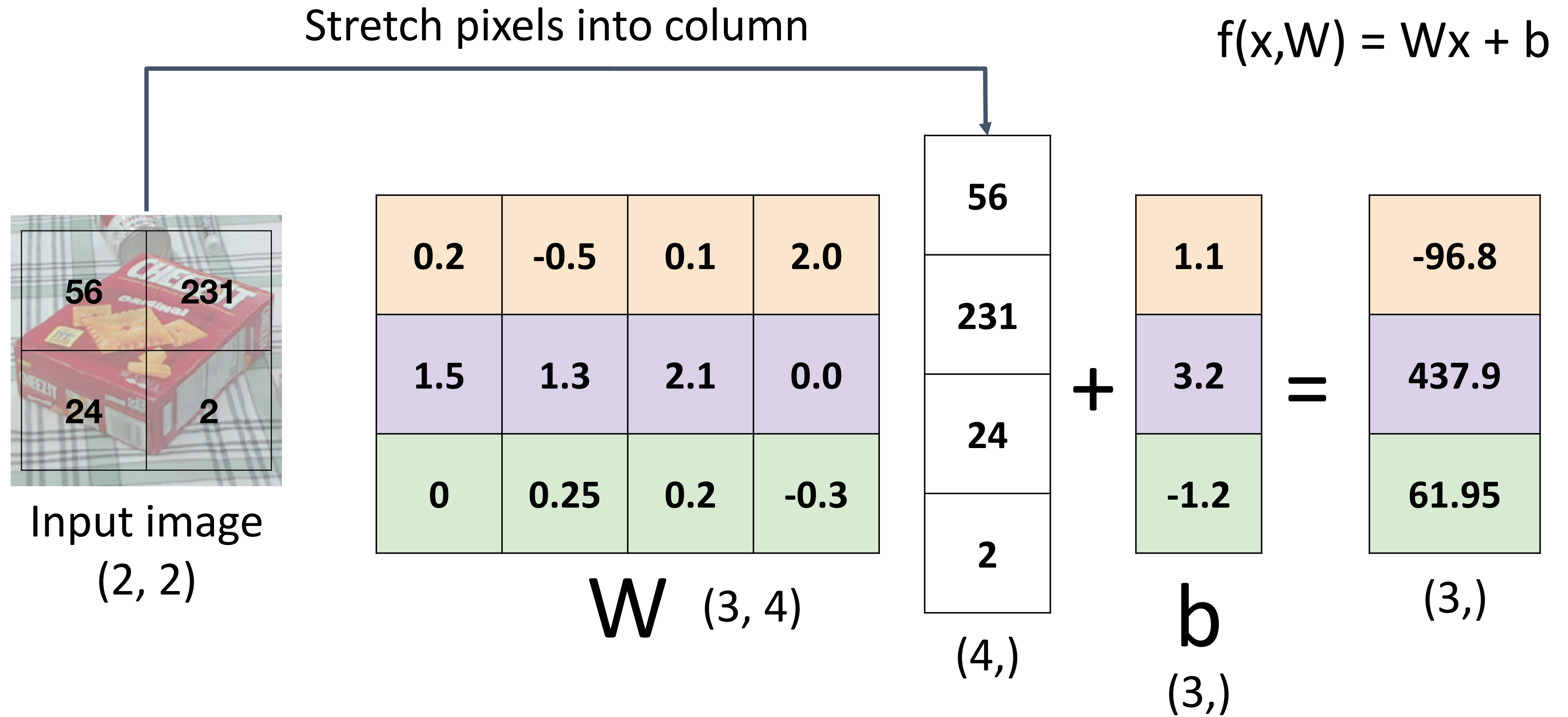


$$f(x,W) = Wx + b$$

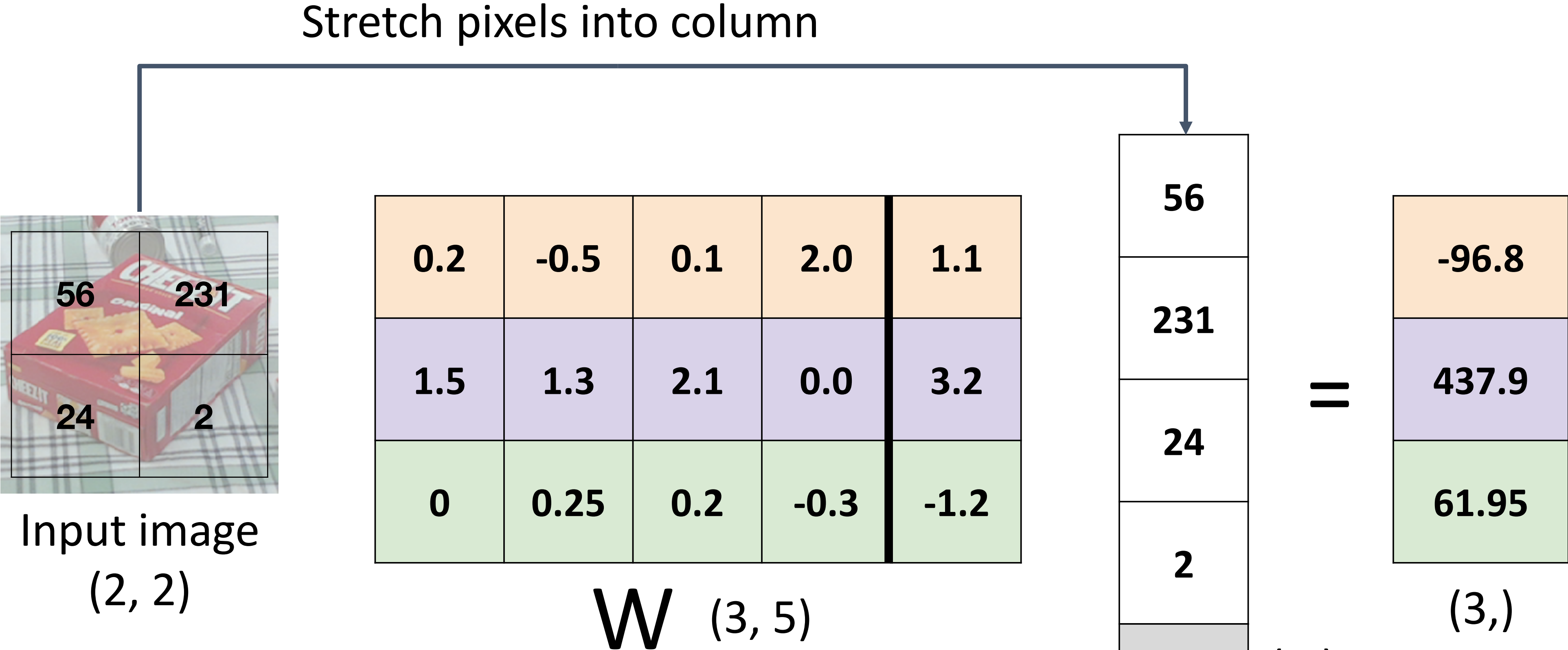
Example for 2x2 Image, 3 classes (crackers/mug/sugar)



Linear Classifier—Algebraic Viewpoint



Linear Classifier—Bias Trick



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



Linear Classifier—Predictions are Linear

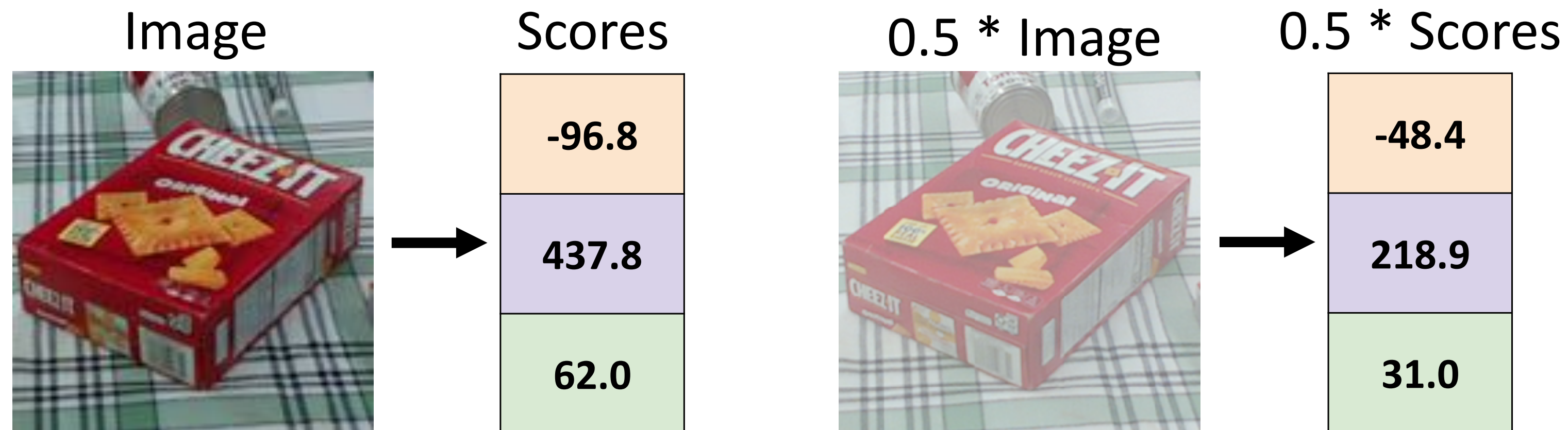
$$f(x, W) = Wx \quad (\text{ignore bias})$$

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

Linear Classifier—Predictions are Linear

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

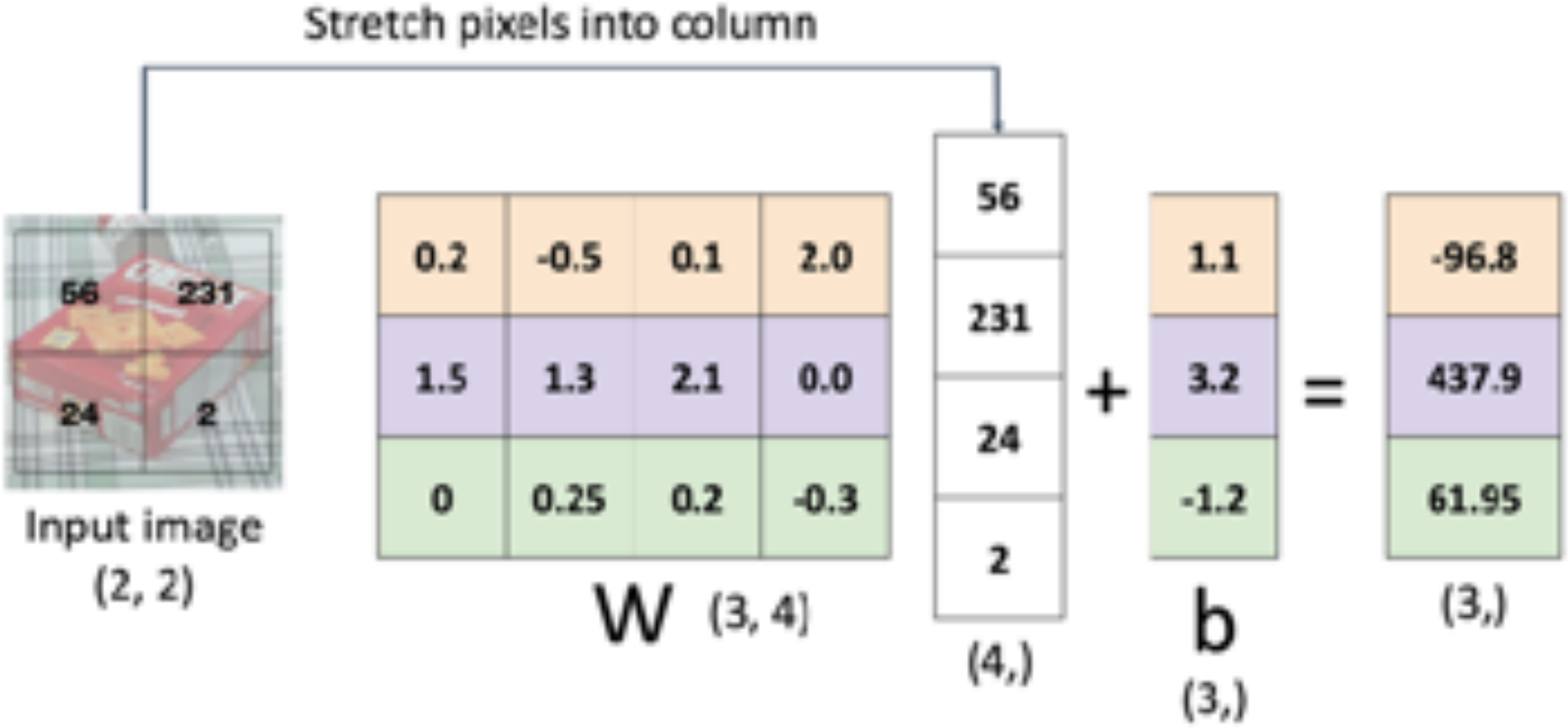
$$f(cx, W) = W(cx) = c * f(x, W)$$



Interpreting a Linear Classifier

Algebraic Viewpoint

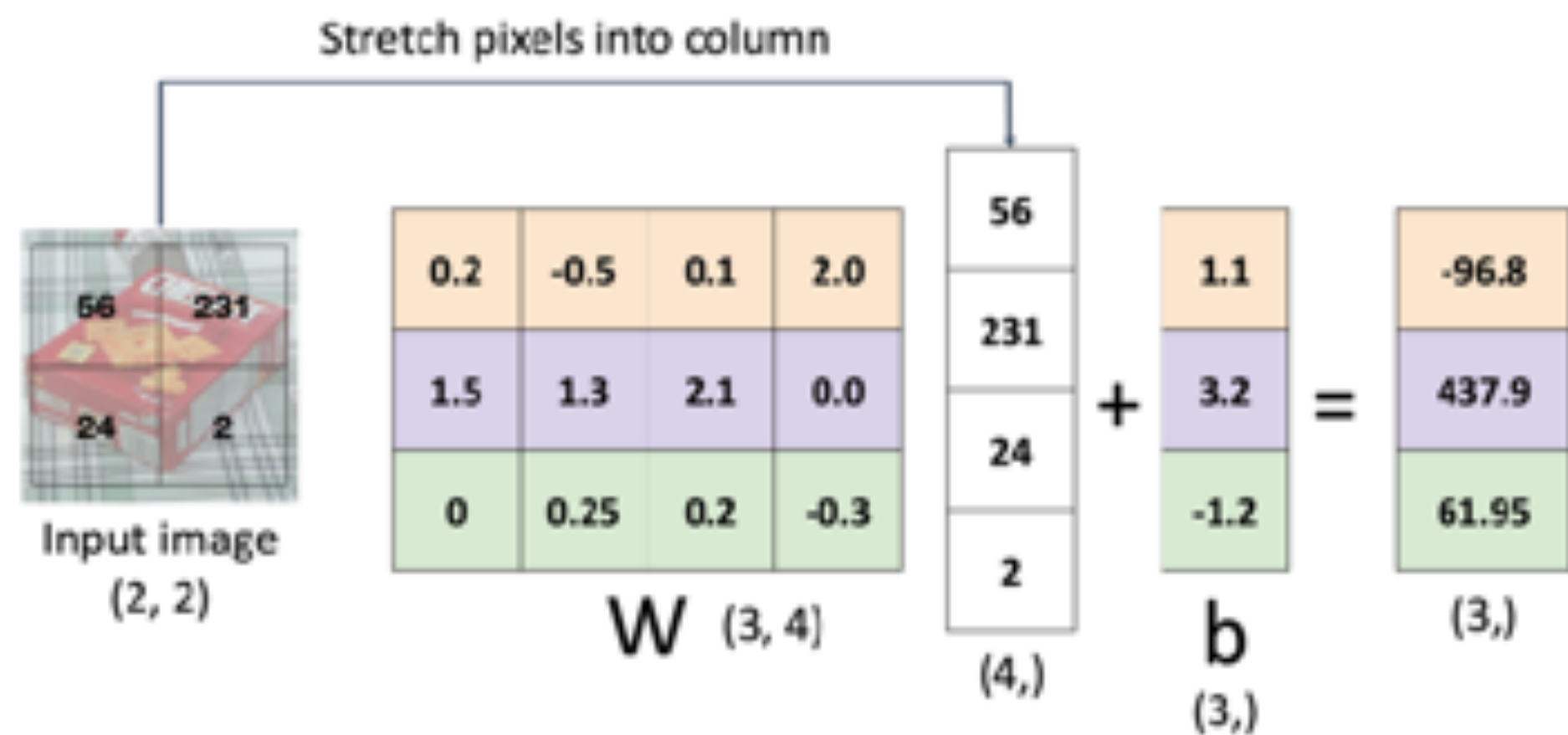
$$f(x,W) = Wx + b$$



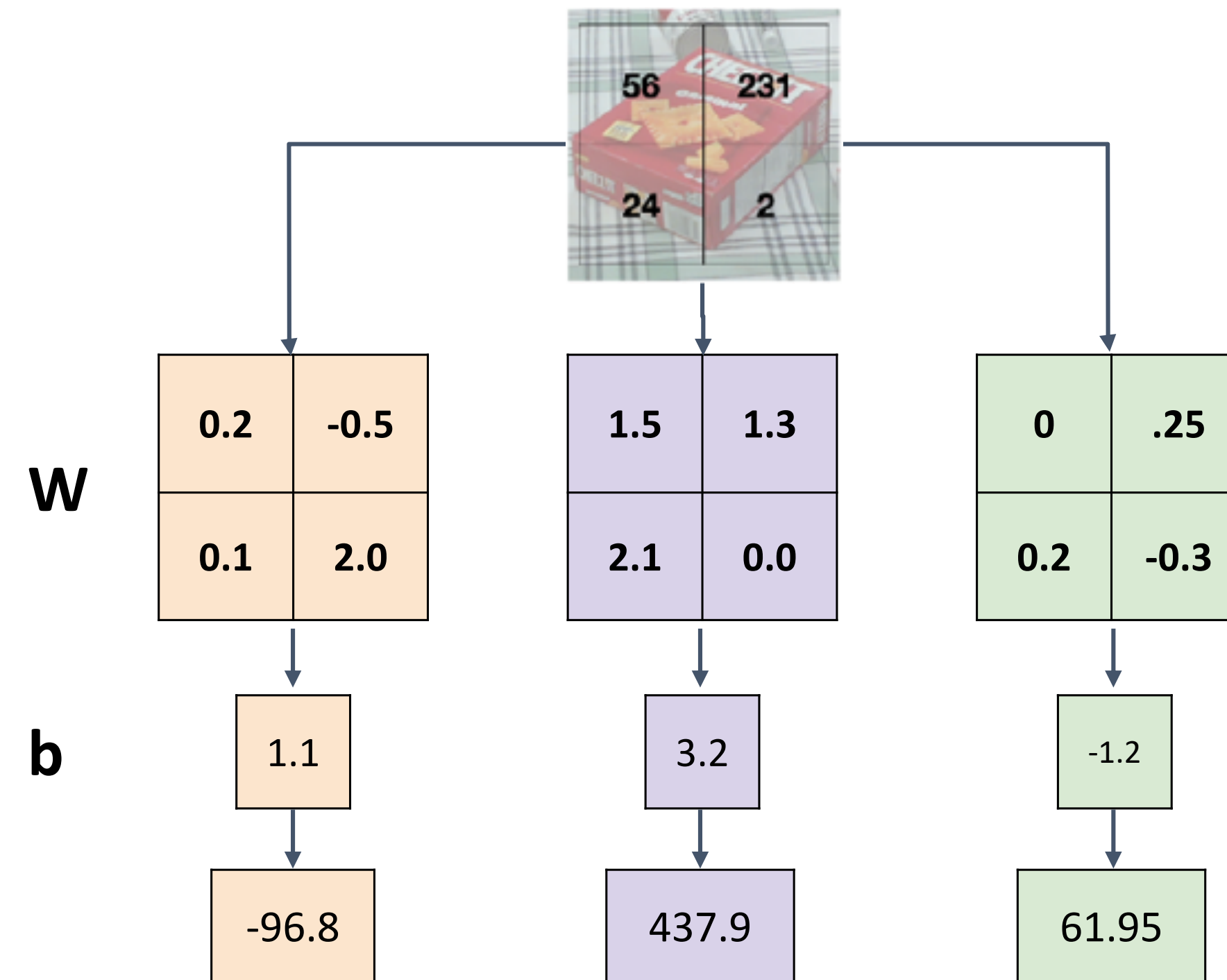
Interpreting a Linear Classifier

Algebraic Viewpoint

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

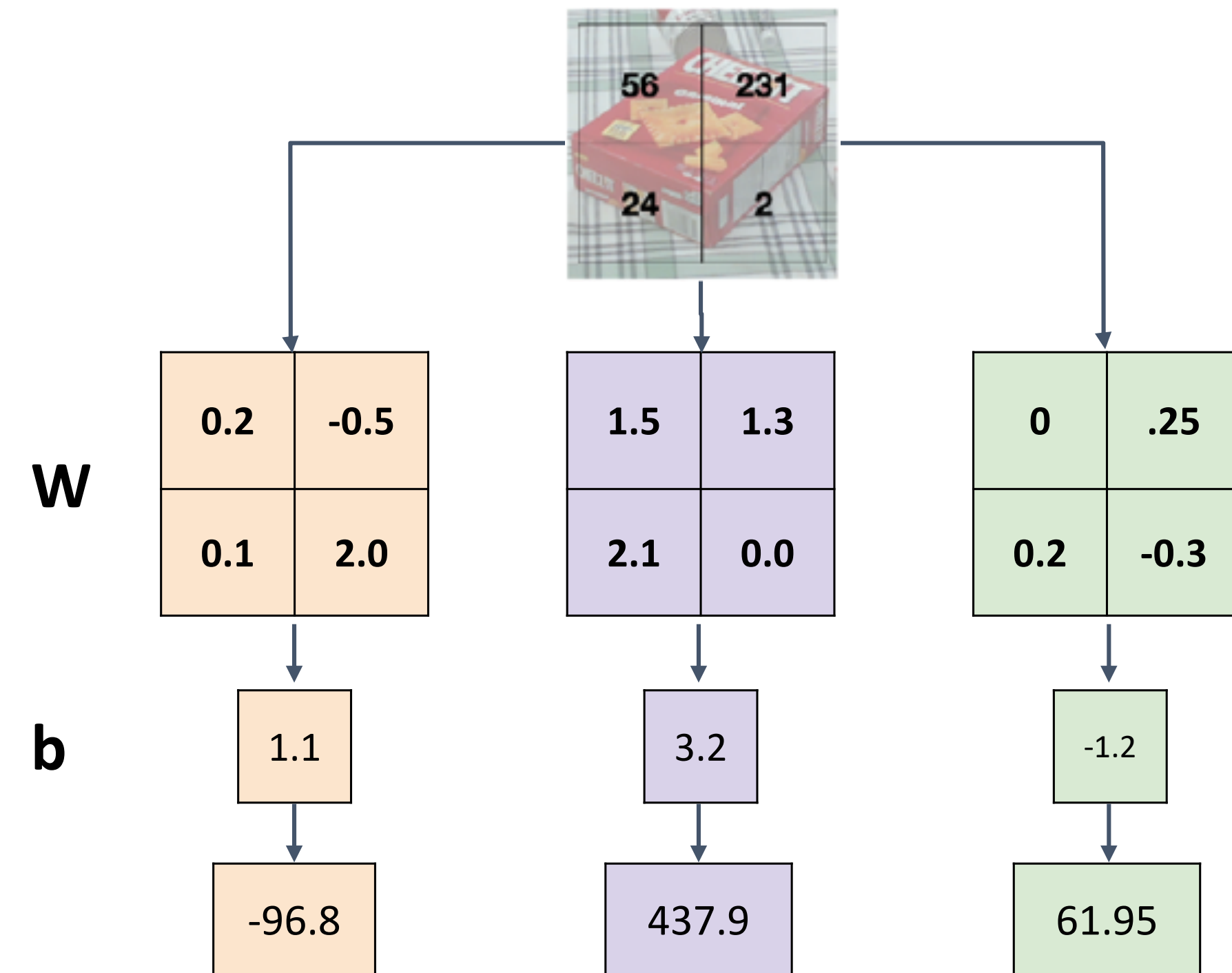
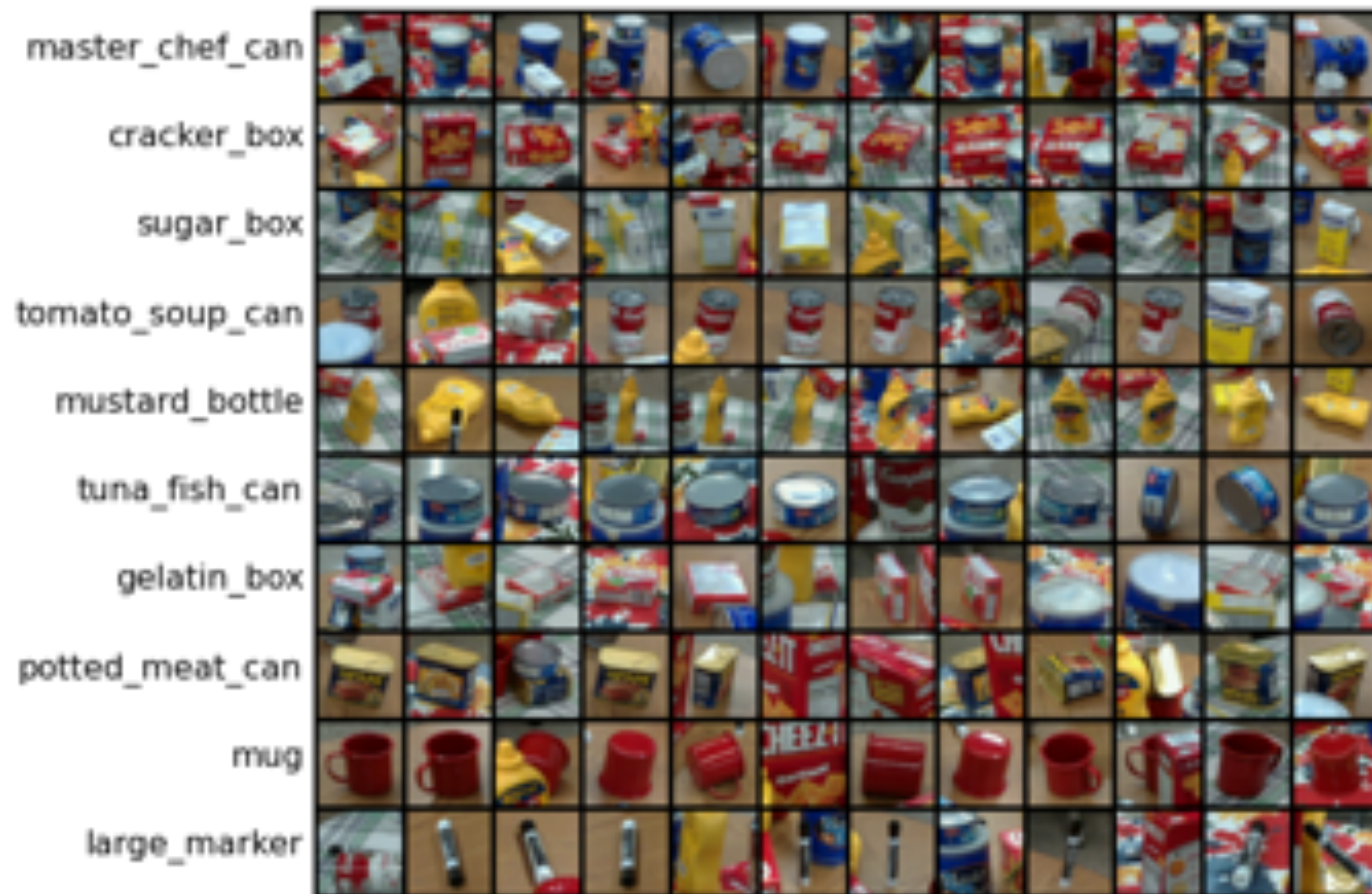


Instead of stretching pixels into columns, we can equivalently stretch rows of W into images!



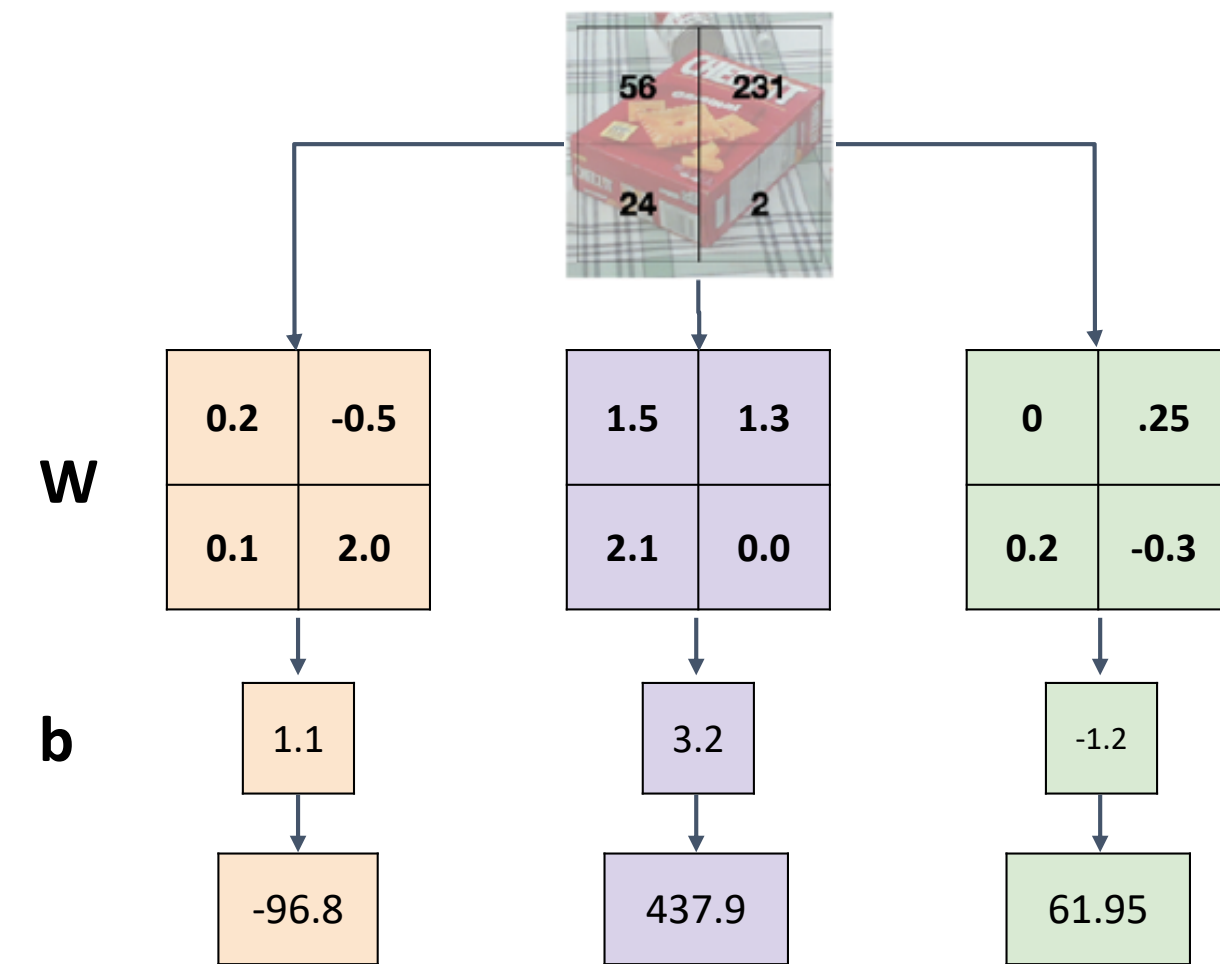
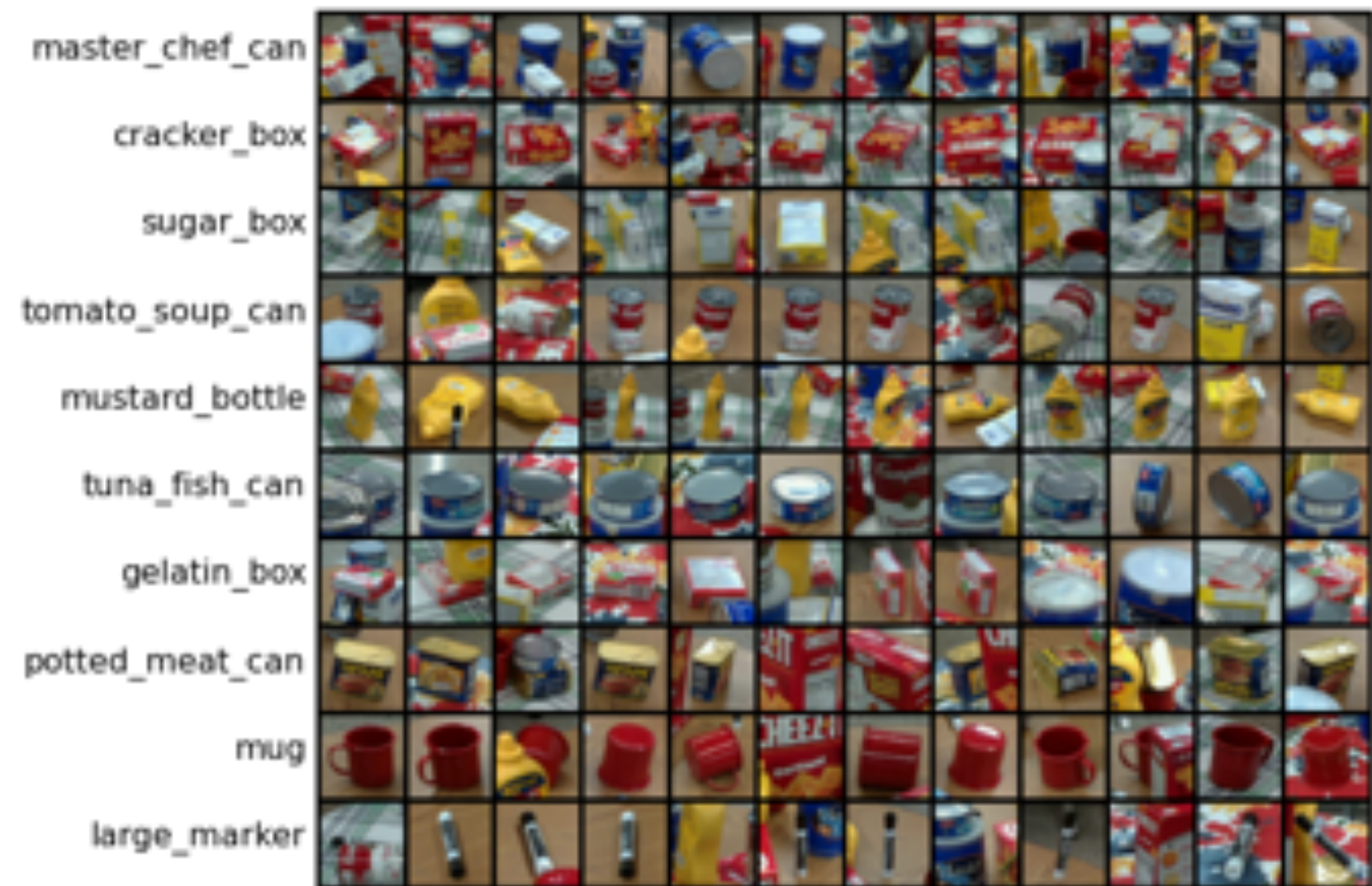
Interpreting a Linear Classifier

Instead of stretching pixels into columns, we can equivalently stretch rows of W into images!



Interpreting a Linear Classifier

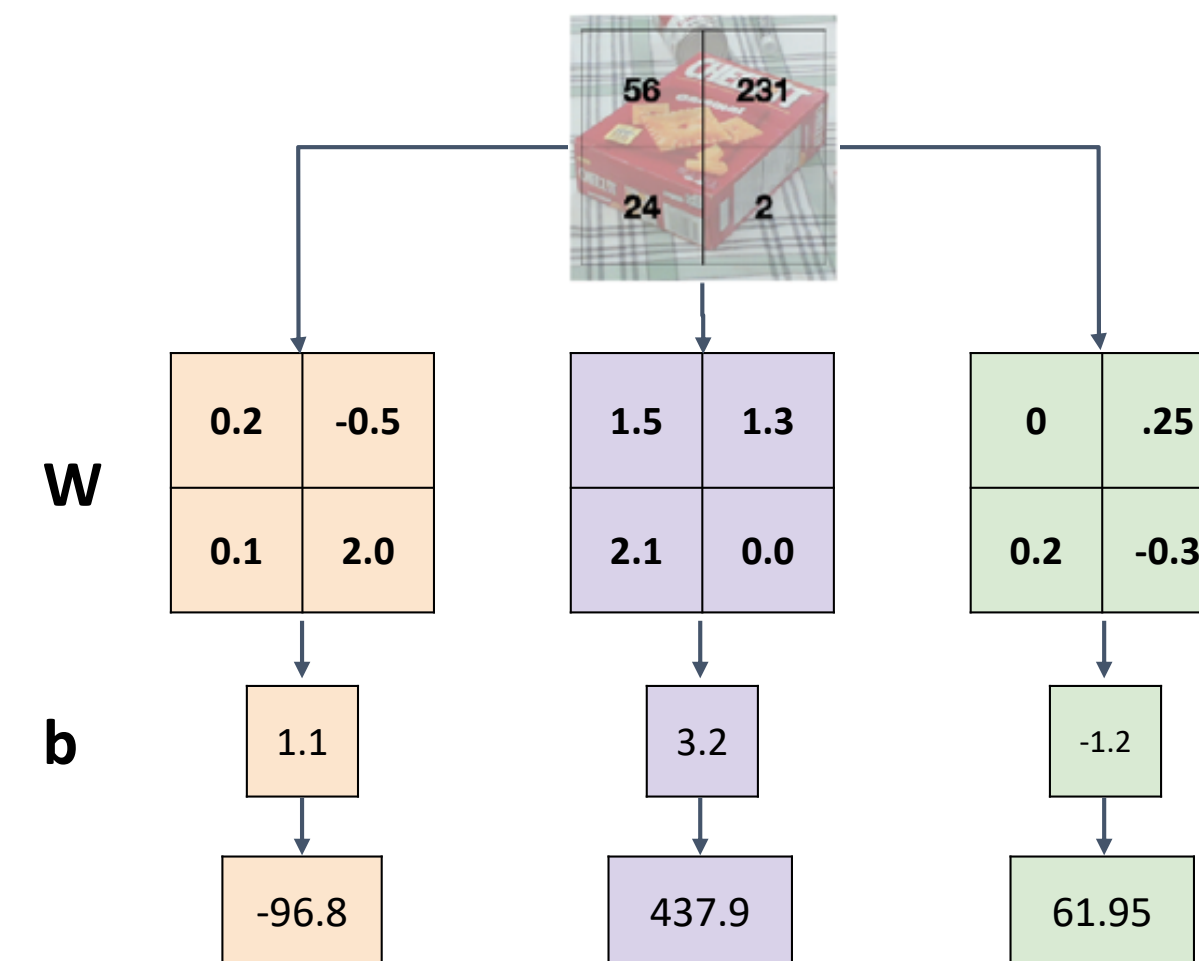
Instead of stretching pixels into columns, we can equivalently stretch rows of W into images!



Interpreting a Linear Classifier—Visual Viewpoint

Linear classifier has one “template” per category

Instead of stretching pixels into columns, we can equivalently stretch rows of W into images!



Interpreting a Linear Classifier—Visual Viewpoint

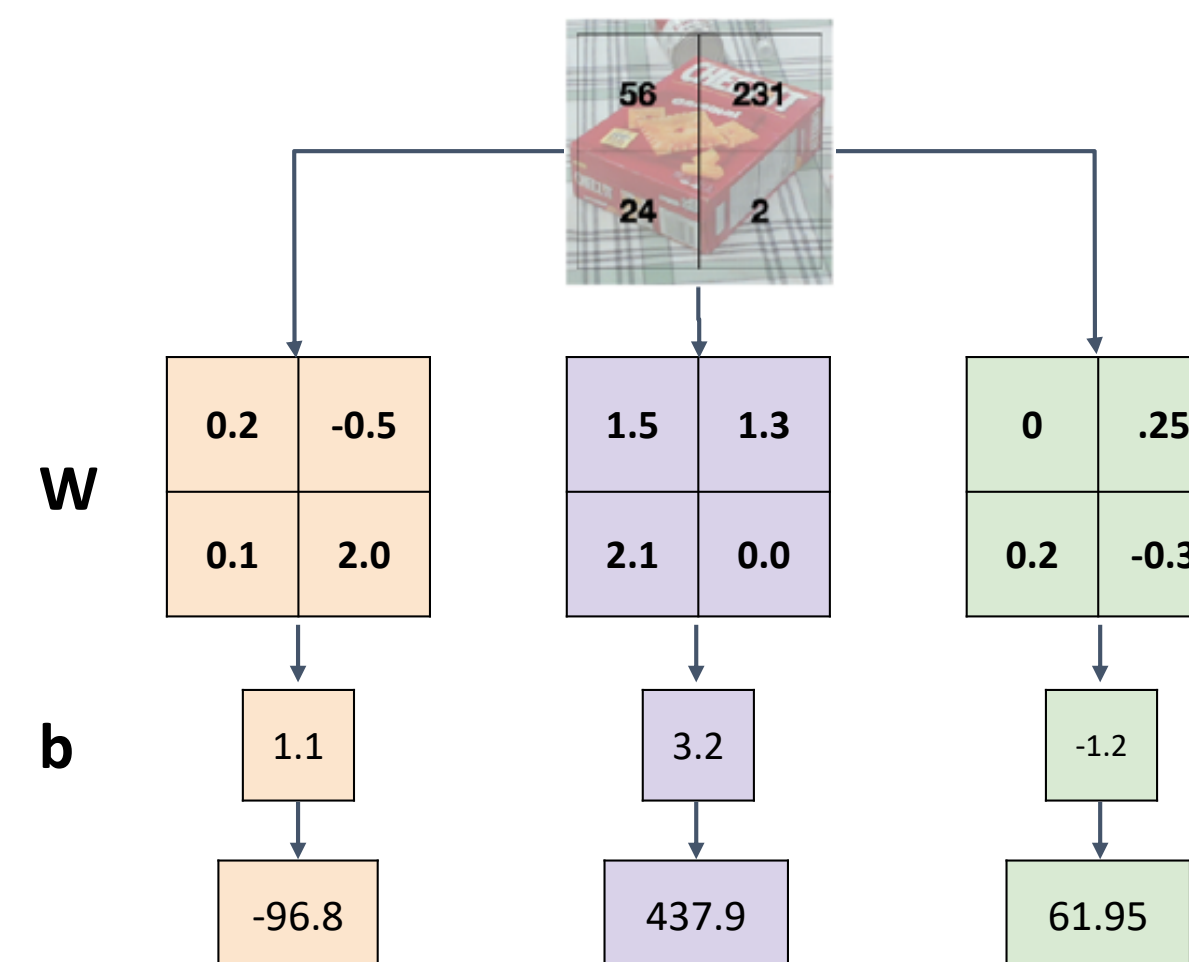
Linear classifier has one “template” per category

A single template cannot capture multiple modes of the data

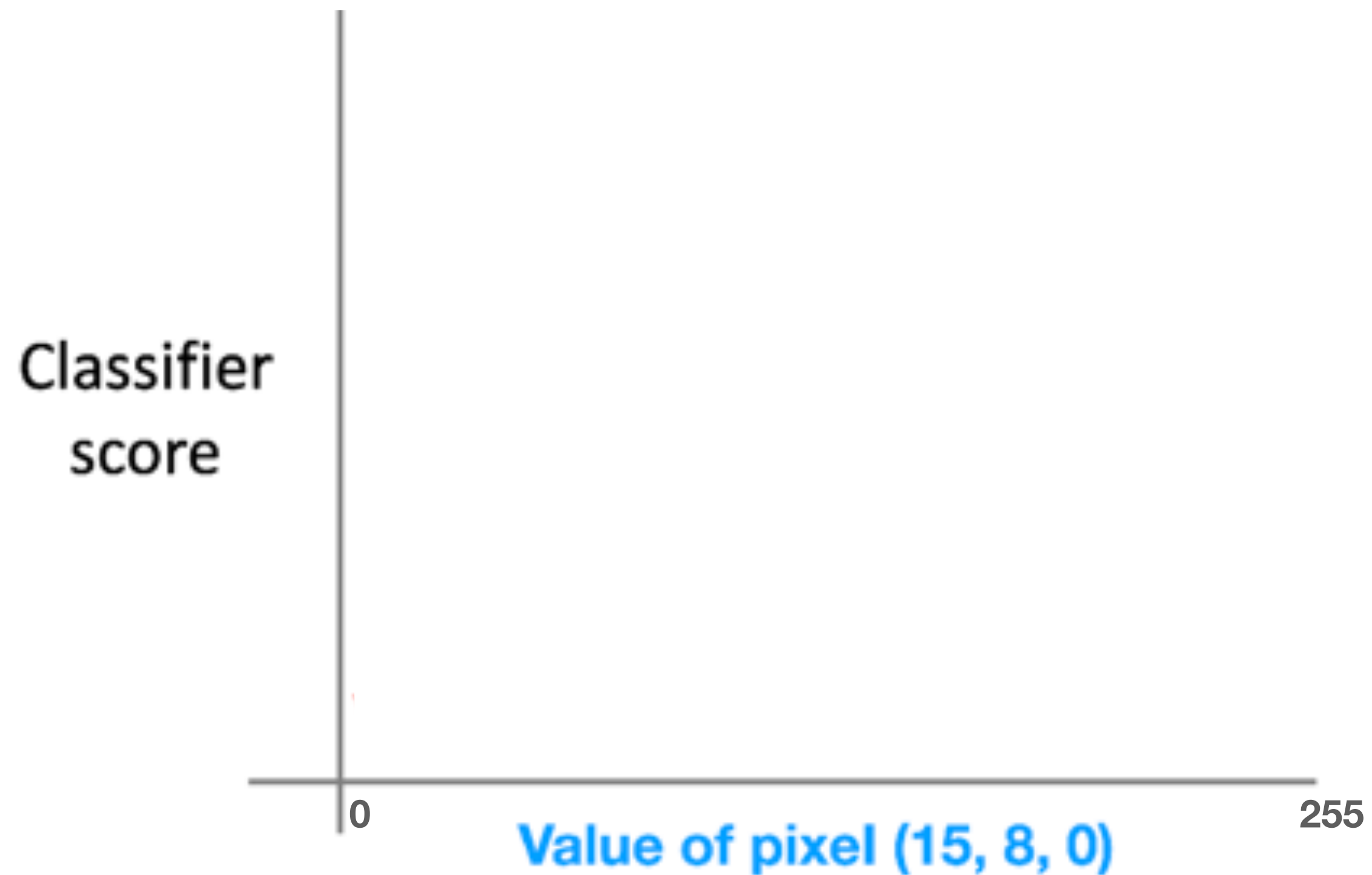
e.g. mustard bottles can rotate



Instead of stretching pixels into columns, we can equivalently stretch rows of W into images!



Interpreting a Linear Classifier—Geometric Viewpoint

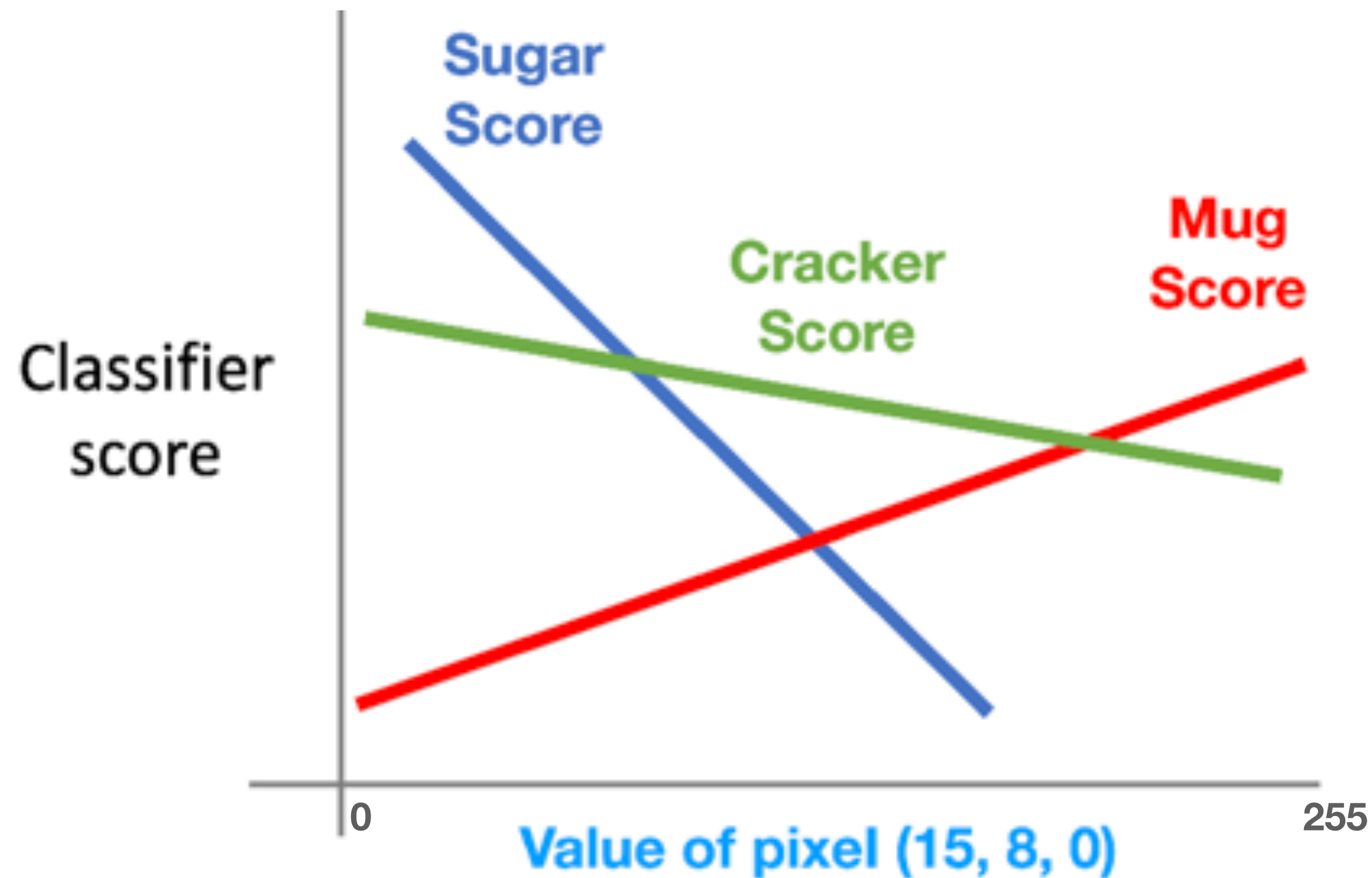


$$f(x, W) = Wx + b$$



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

Interpreting a Linear Classifier—Geometric Viewpoint

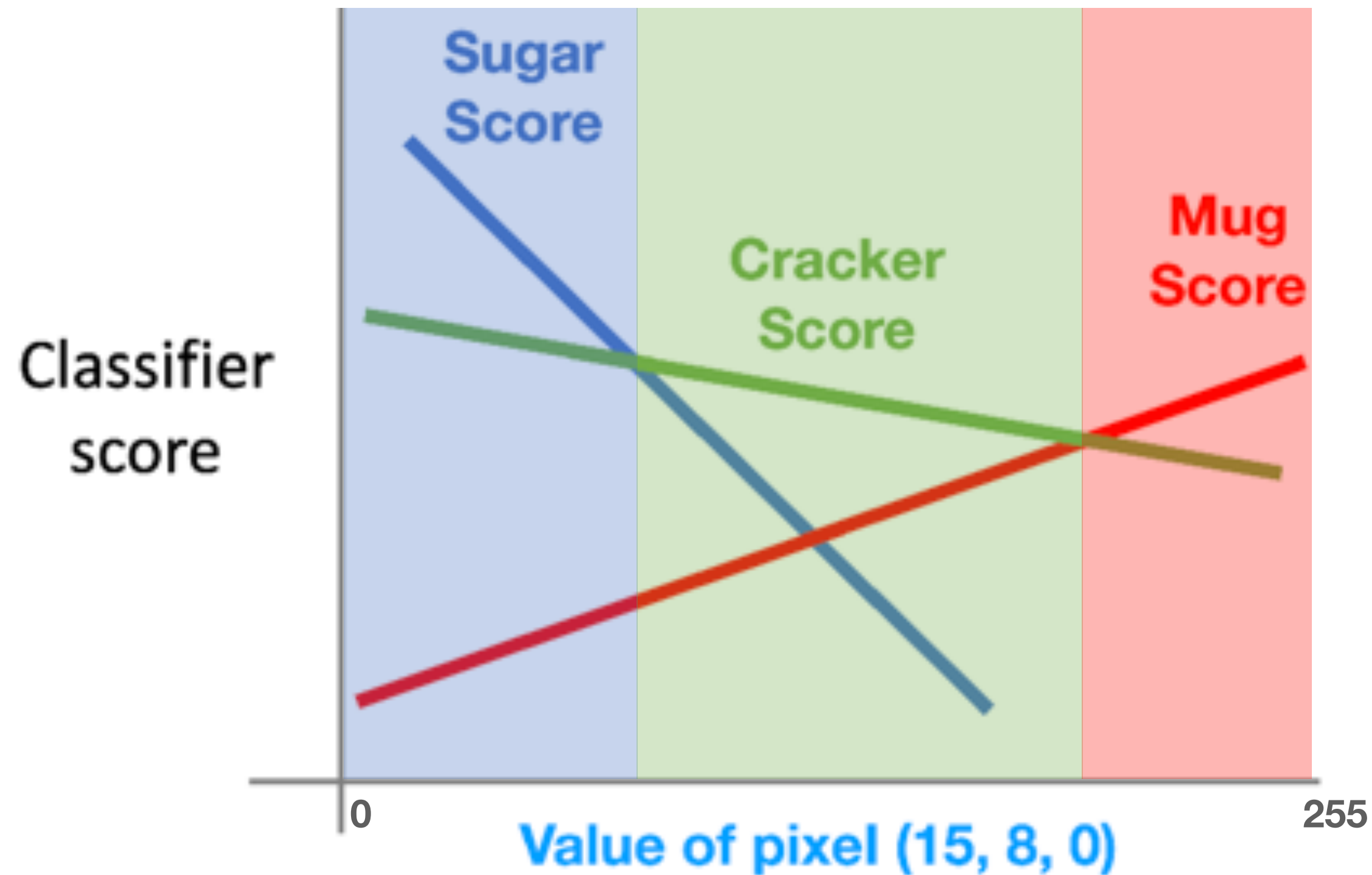


$$f(x, W) = Wx + b$$



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

Interpreting a Linear Classifier—Geometric Viewpoint

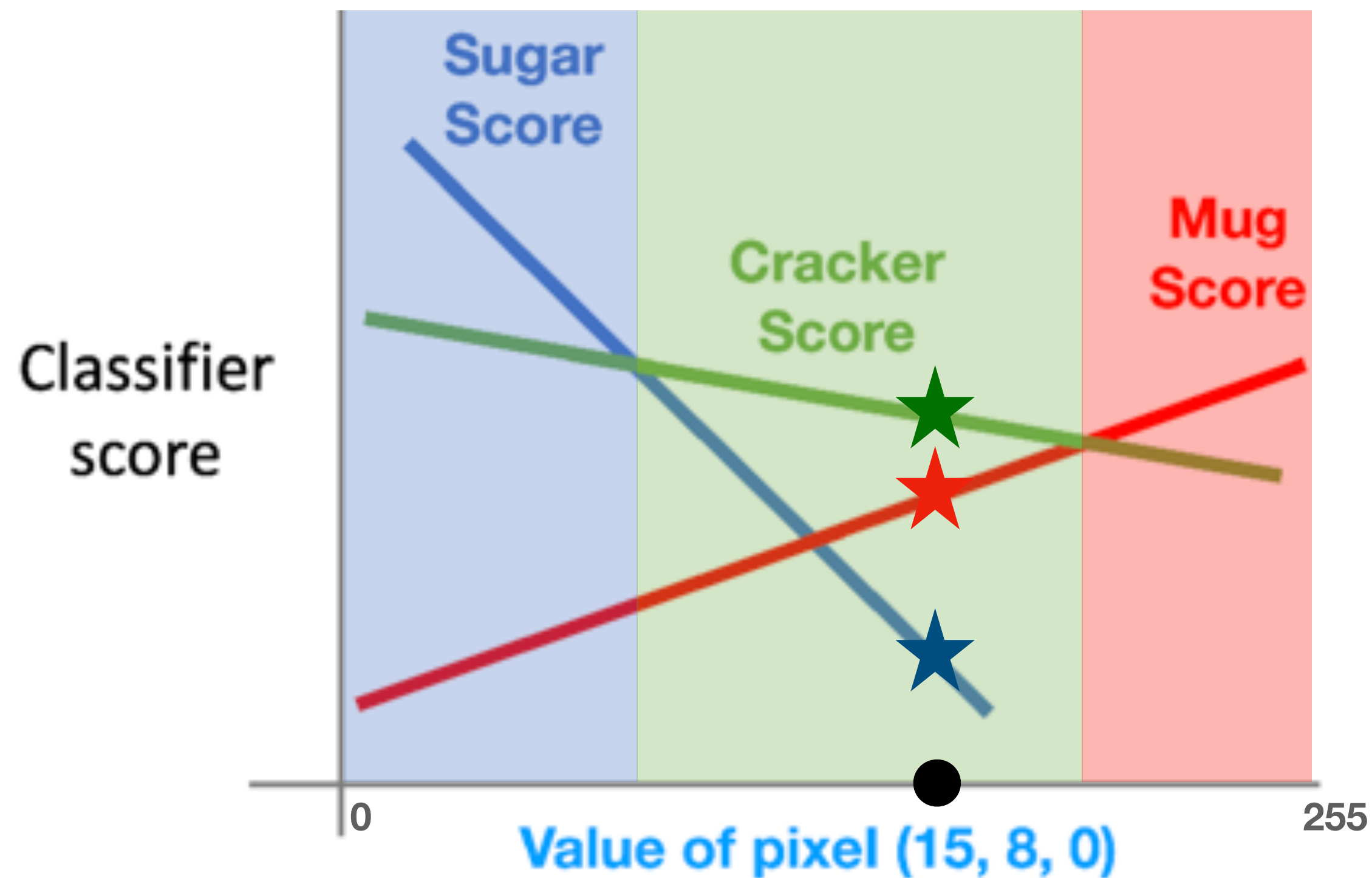


$$f(x, W) = Wx + b$$



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

Interpreting a Linear Classifier—Geometric Viewpoint

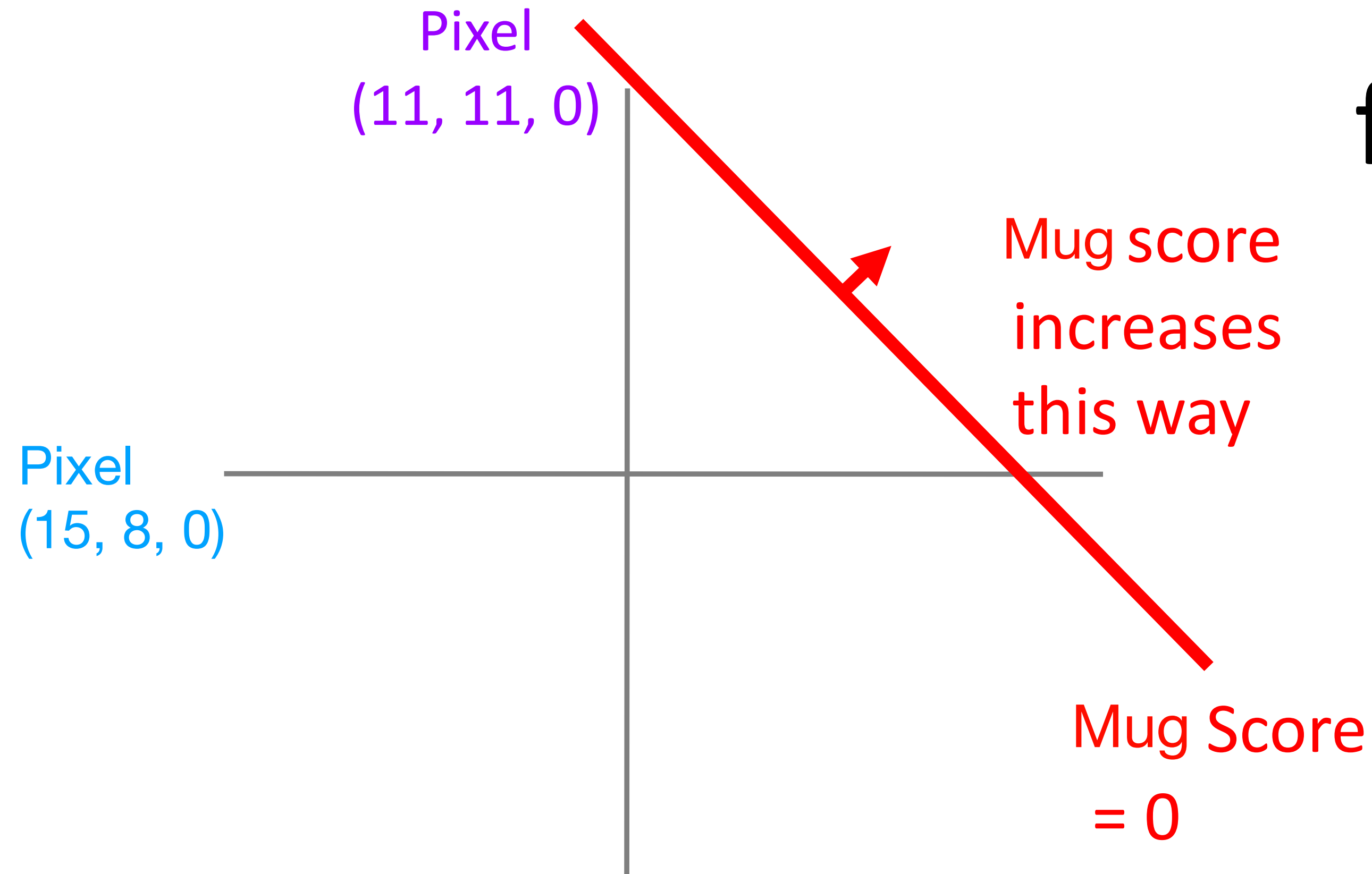


$$f(x, W) = Wx + b$$



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

Interpreting a Linear Classifier—Geometric Viewpoint

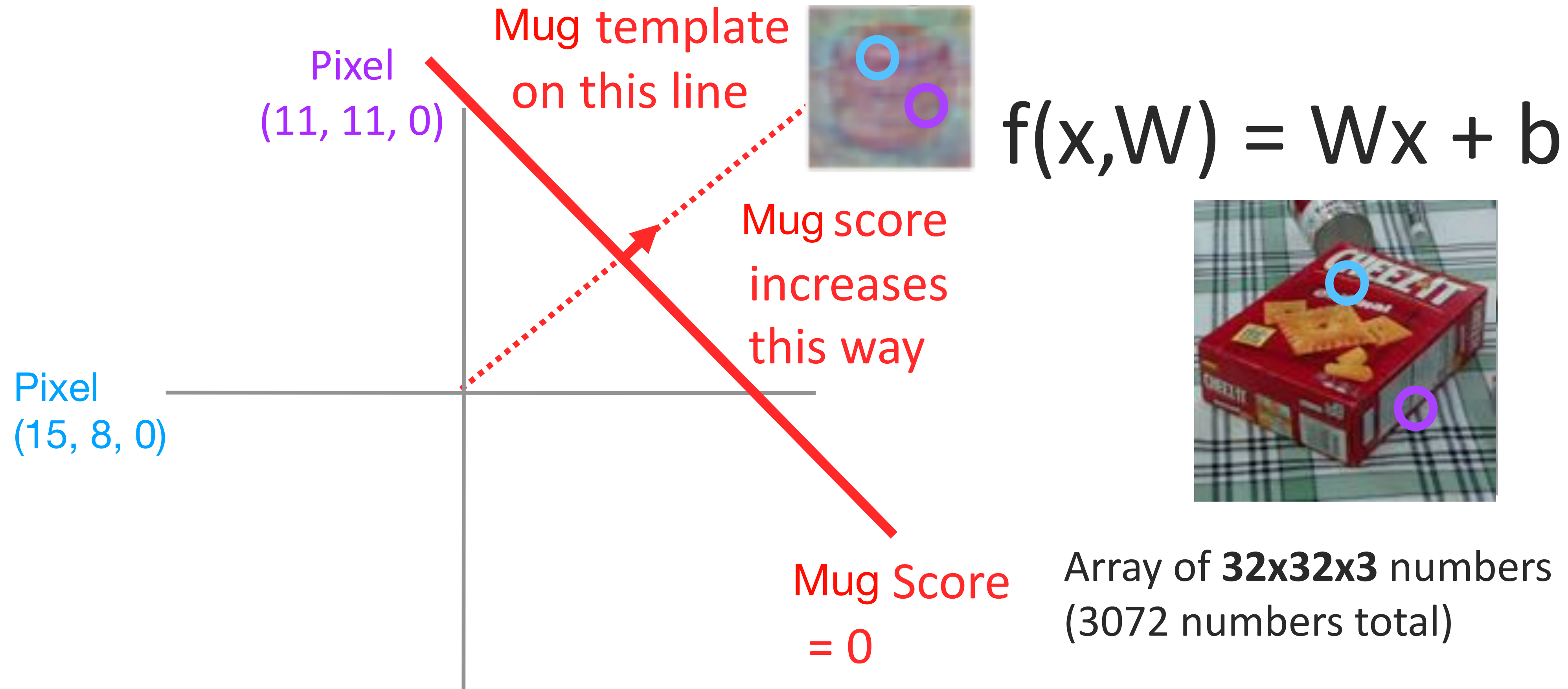


$$f(x, W) = Wx + b$$

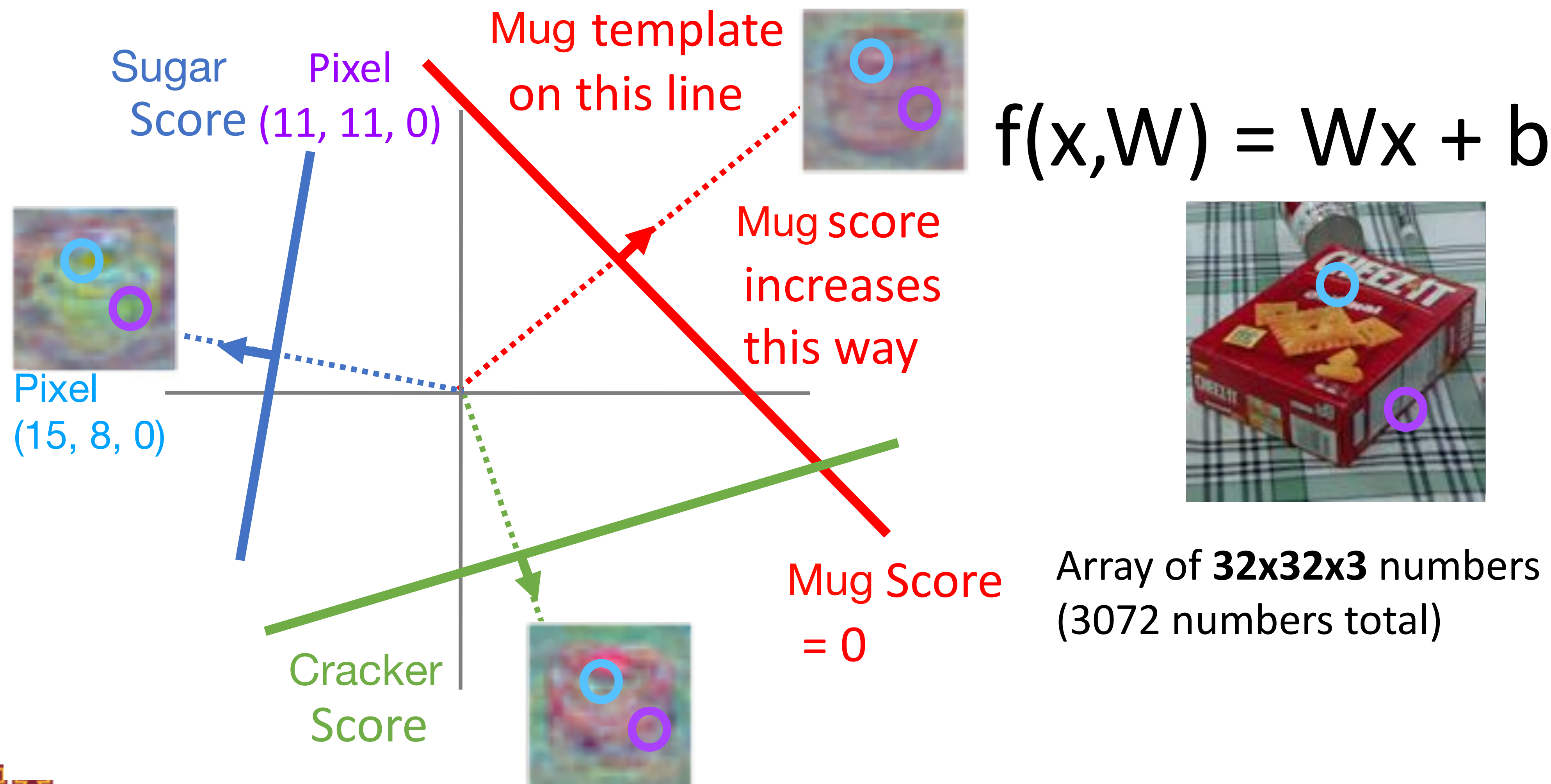


Array of **32x32x3** numbers
(3072 numbers total)

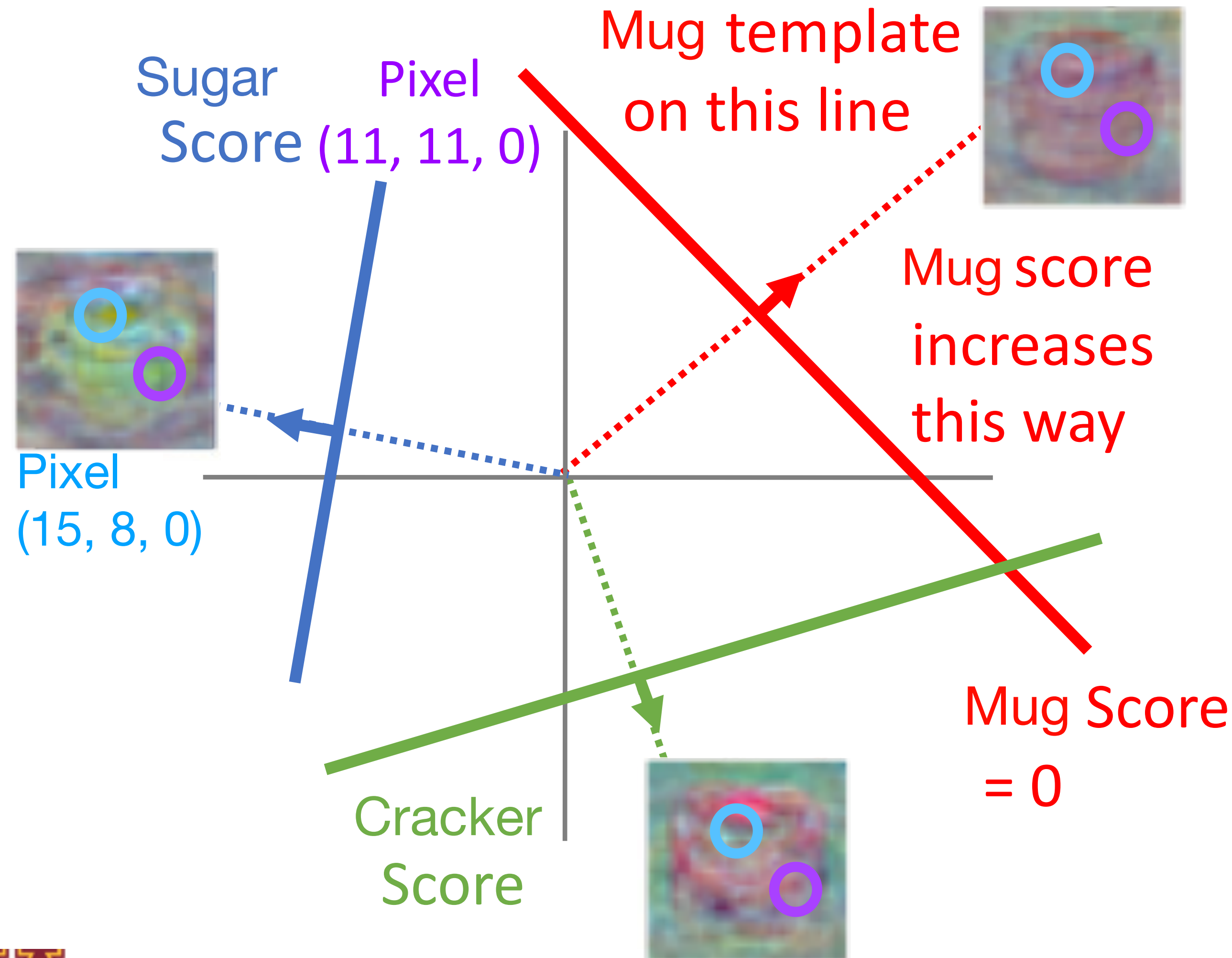
Interpreting a Linear Classifier—Geometric Viewpoint



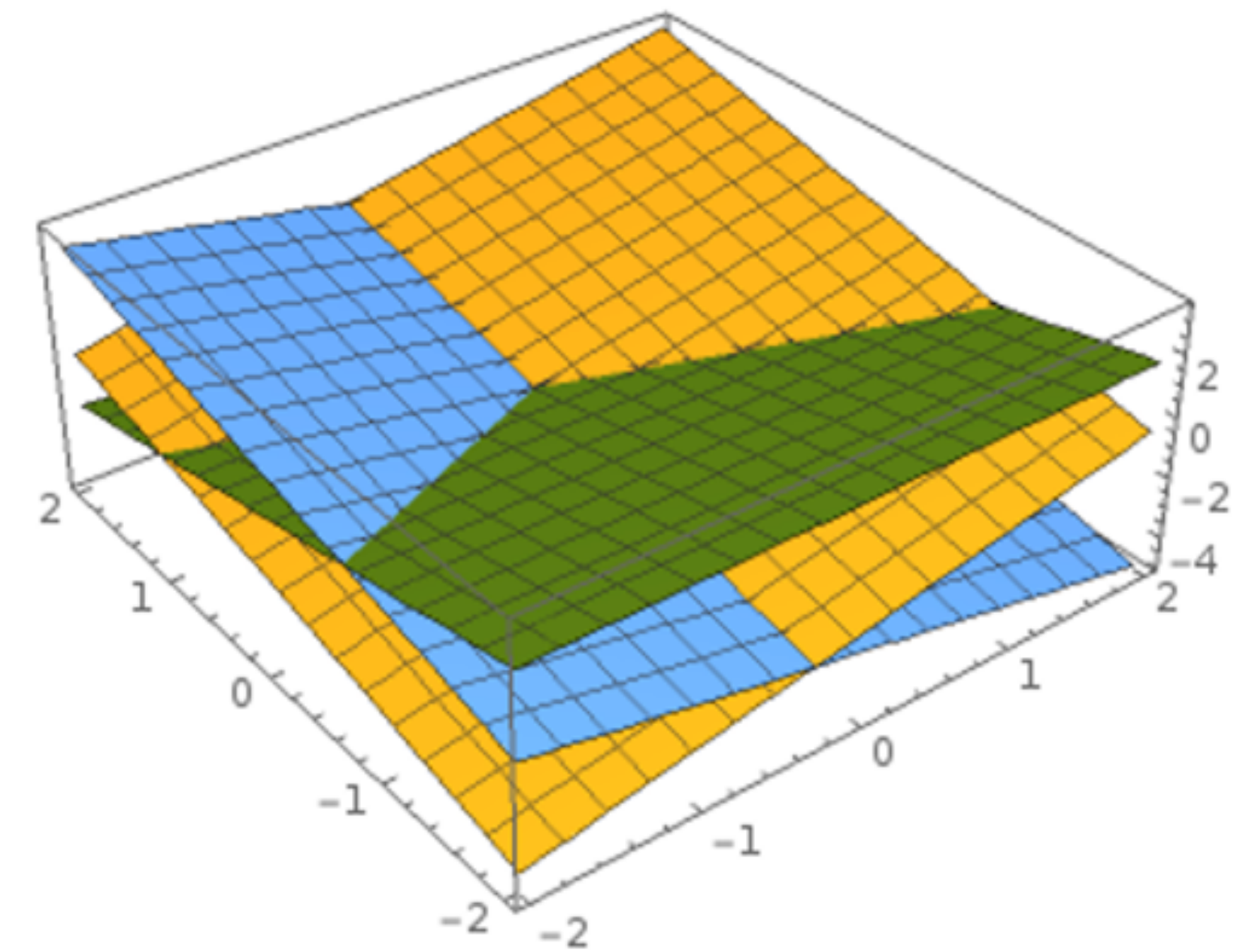
Interpreting a Linear Classifier—Geometric Viewpoint



Interpreting a Linear Classifier—Geometric Viewpoint



Hyperplanes carving up a high-dimensional space



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



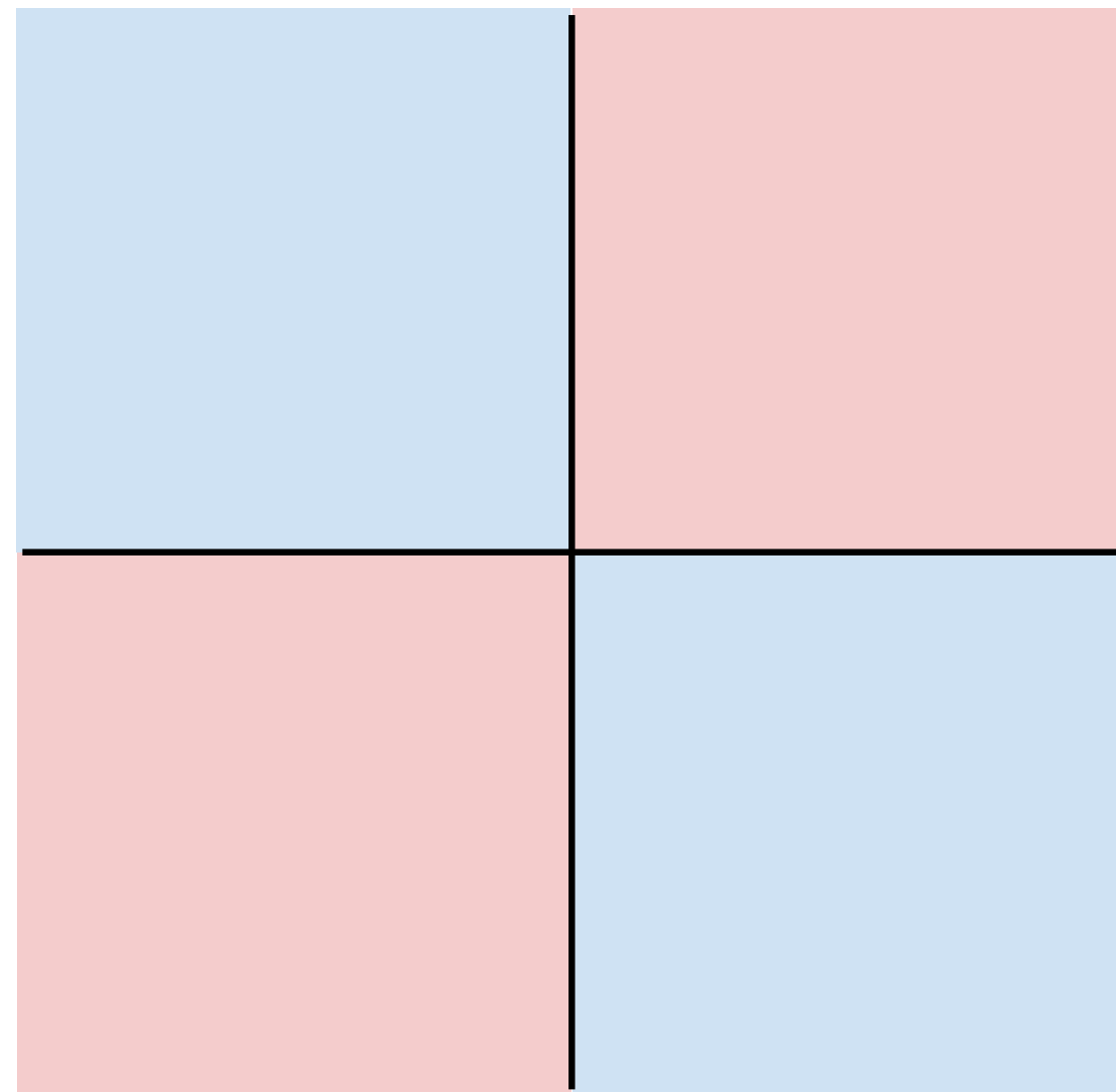
Hard Cases for a Linear Classifier

Class 1:

First and third quadrants

Class 2:

Second and fourth quadrants



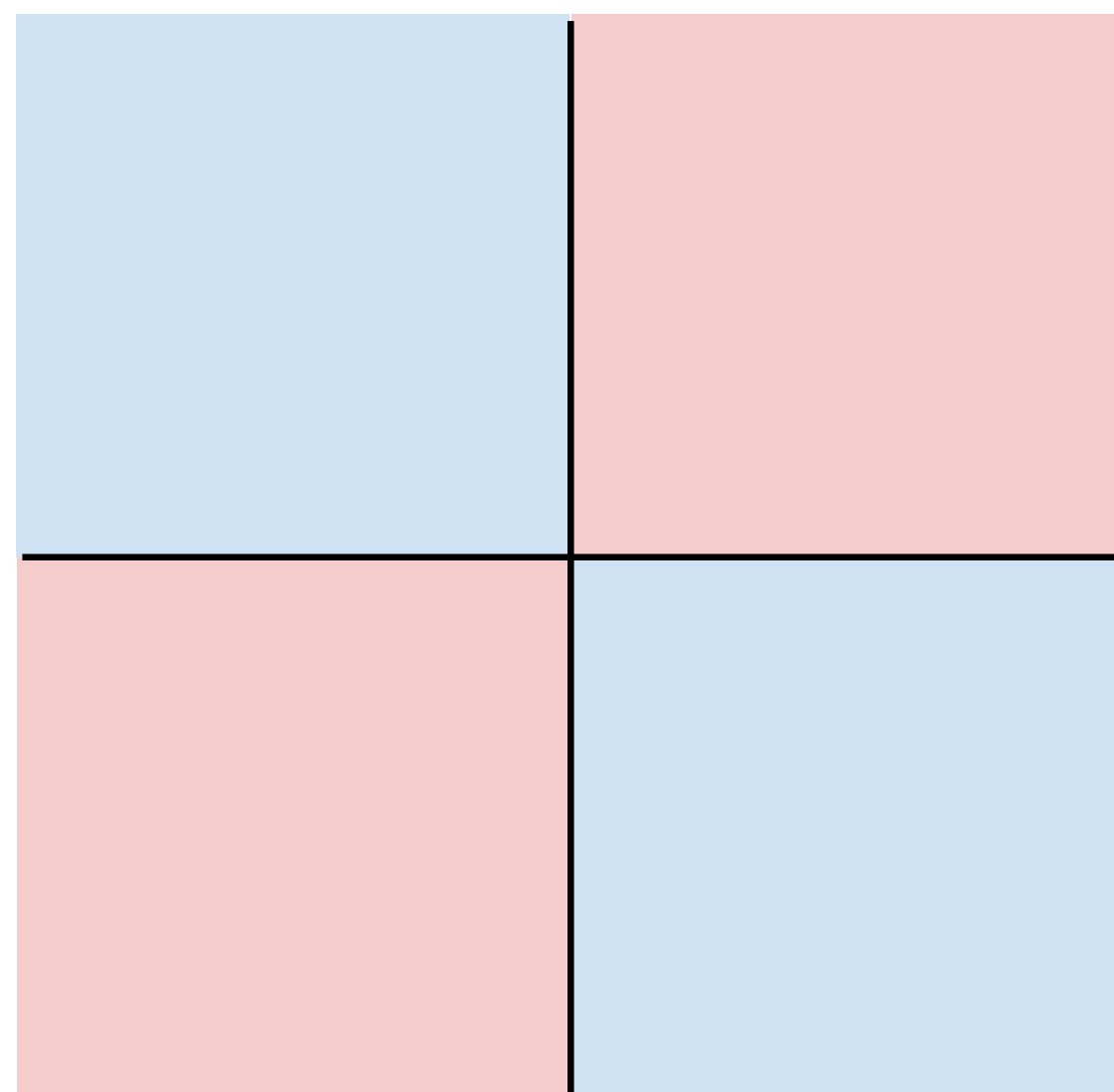
Hard Cases for a Linear Classifier

Class 1:

First and third quadrants

Class 2:

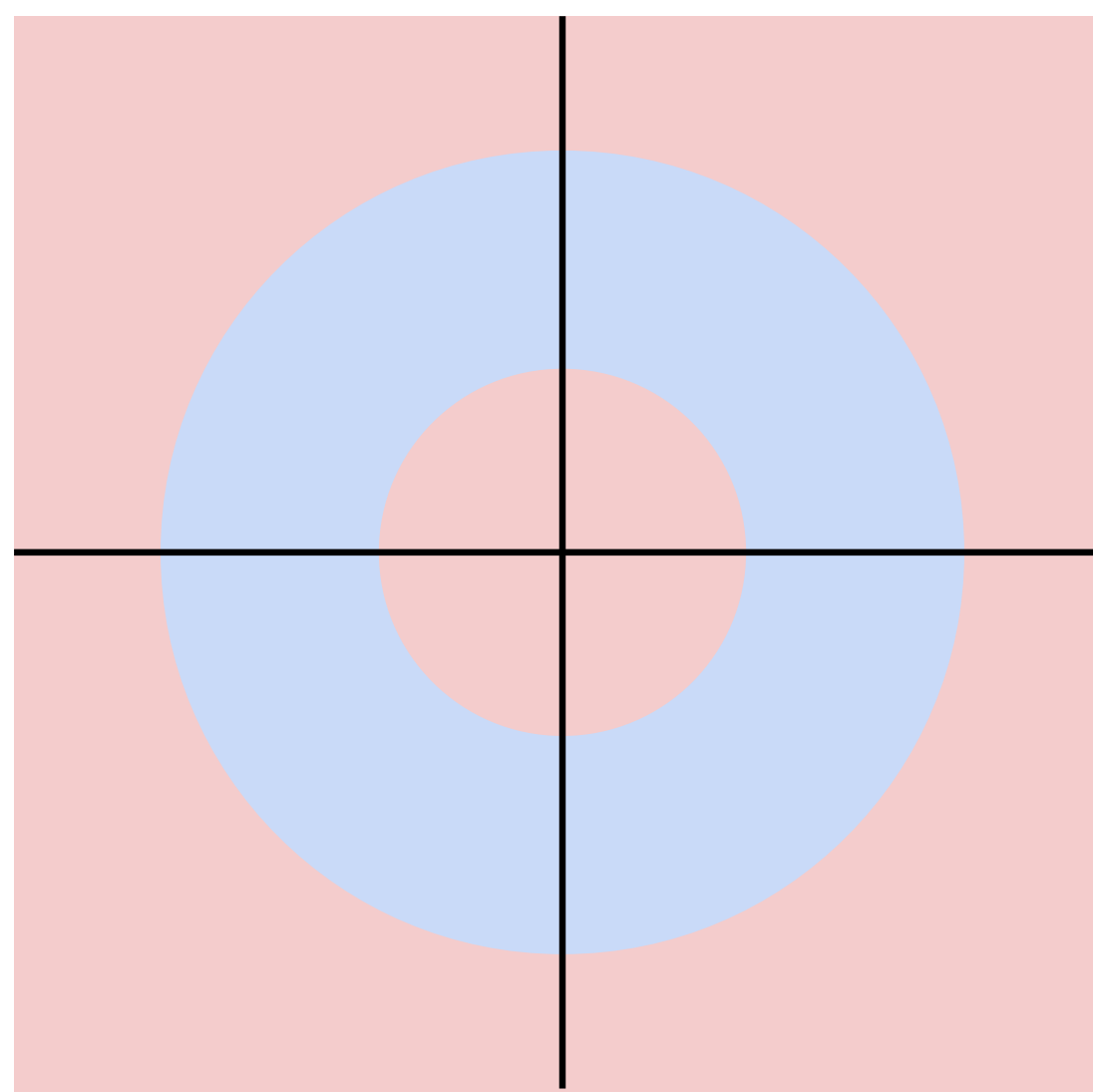
Second and fourth quadrants

**Class 1:**

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

Class 2:

Everything else



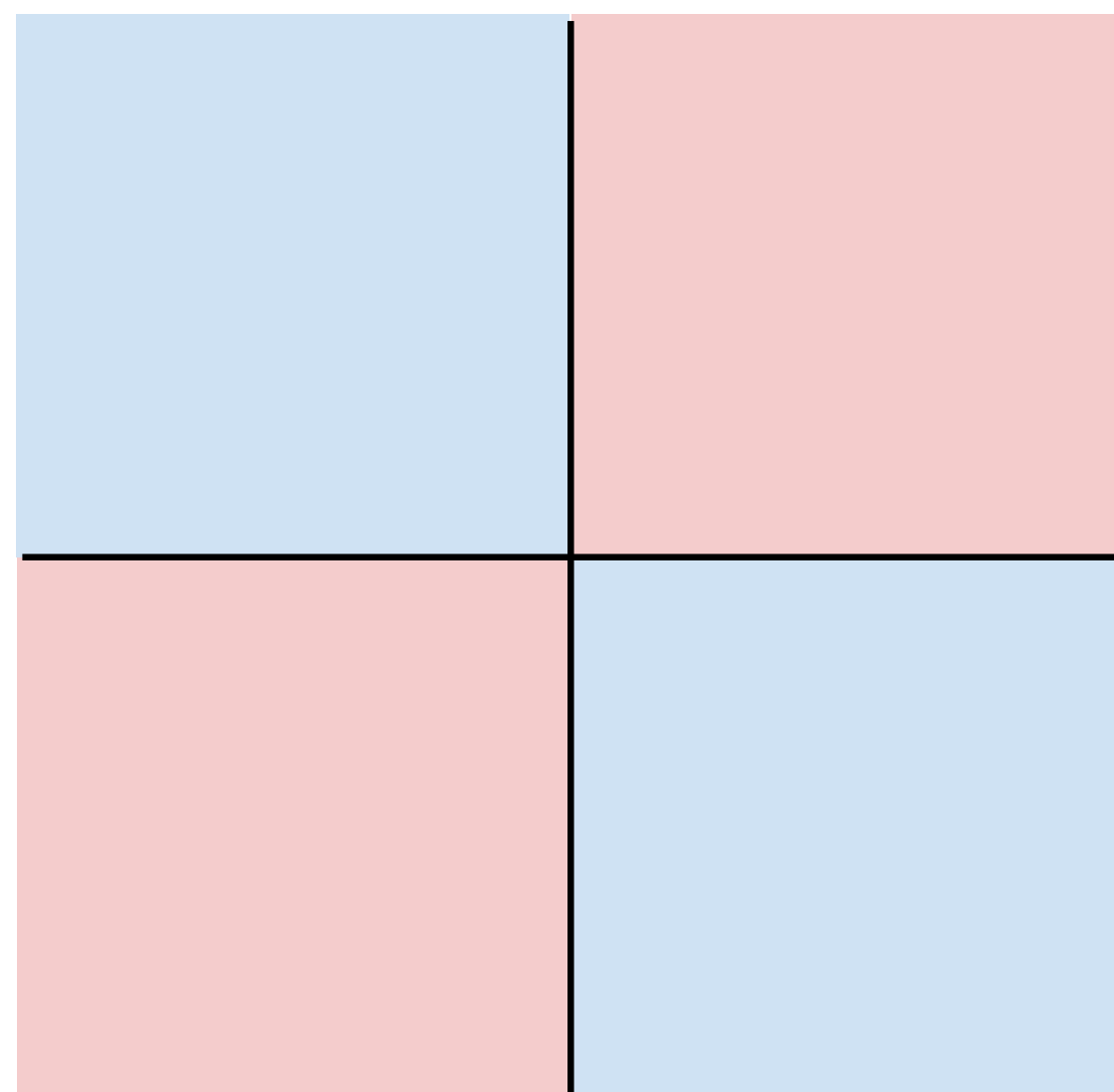
Hard Cases for a Linear Classifier

Class 1:

First and third quadrants

Class 2:

Second and fourth quadrants

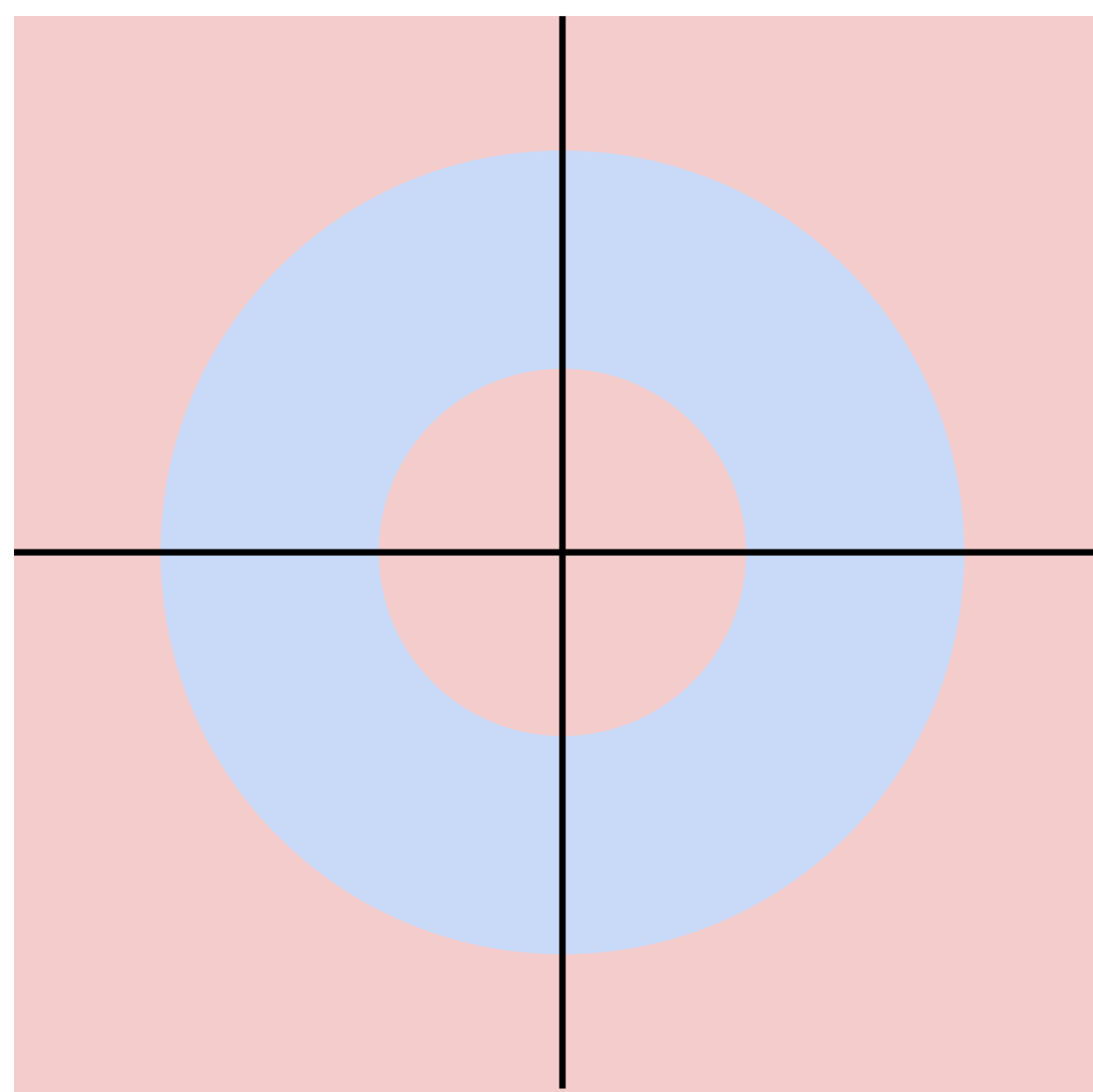


Class 1:

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

Class 2:

Everything else

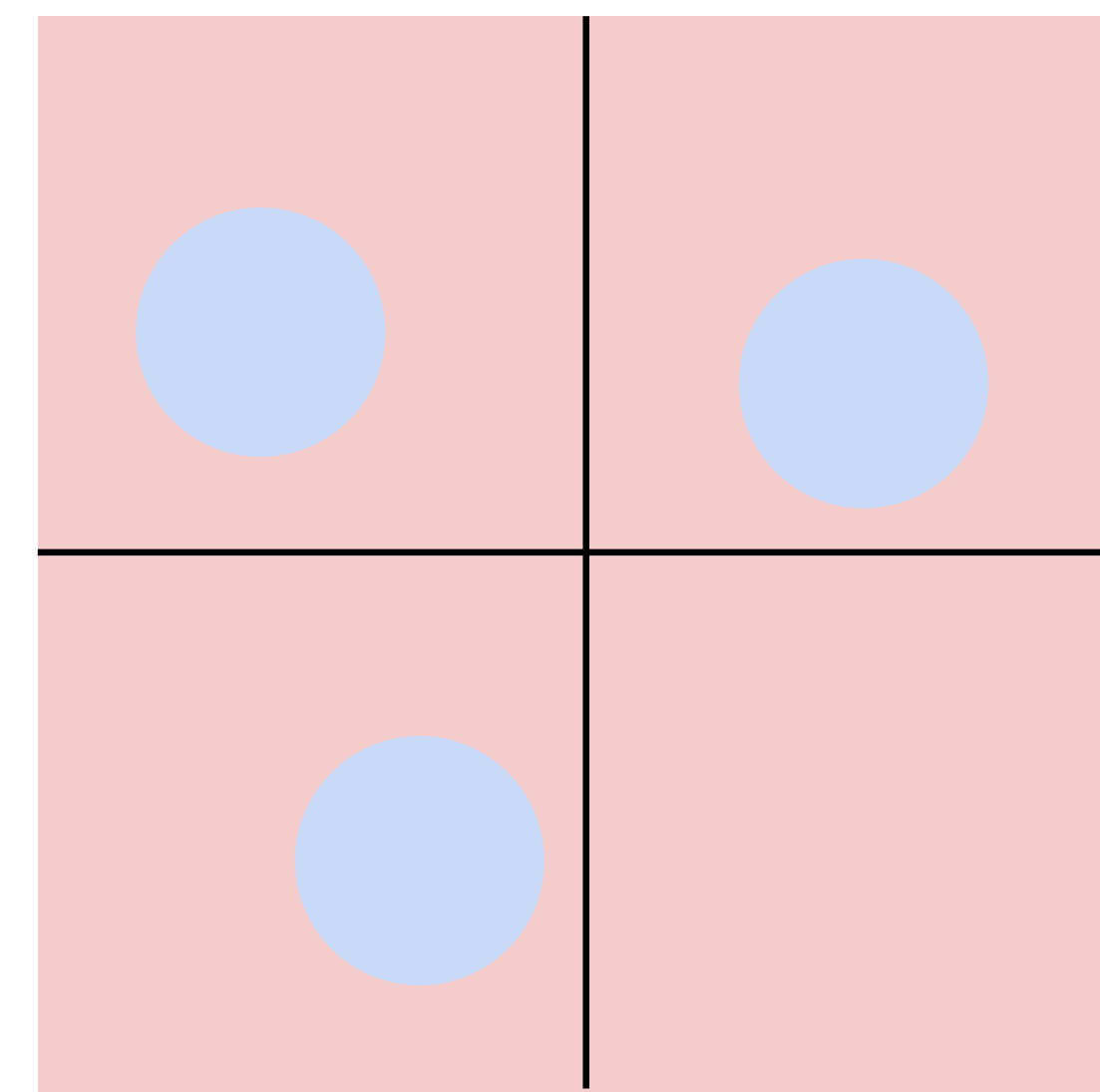


Class 1:

Three modes

Class 2:

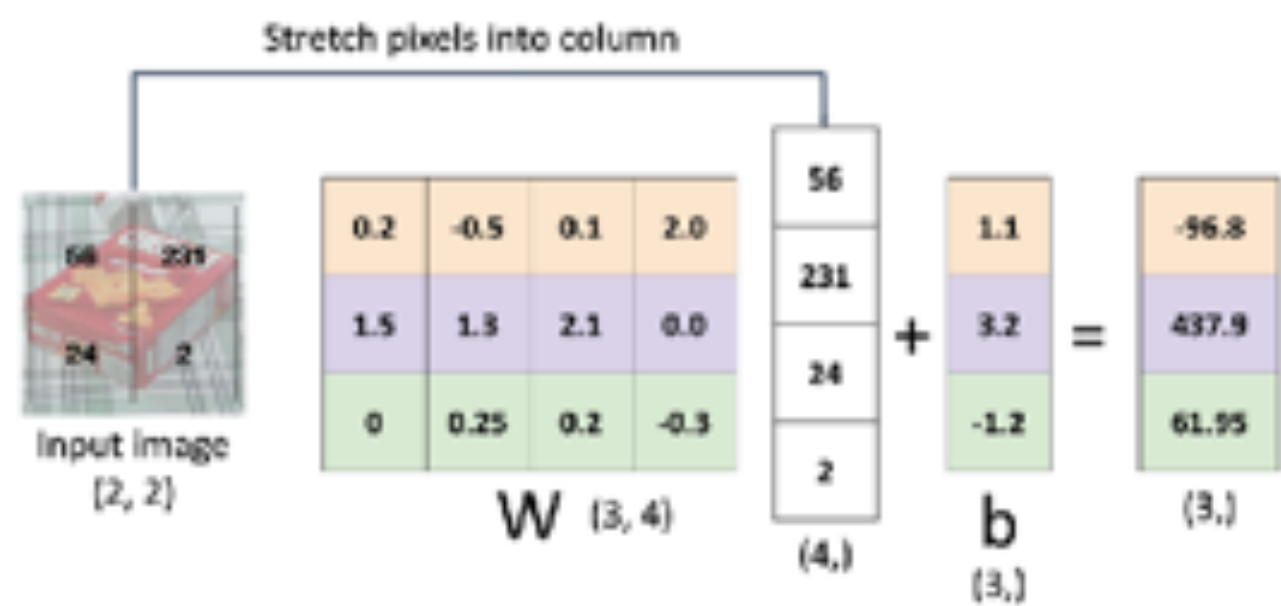
Everything else



Linear Classifier—Three Viewpoints

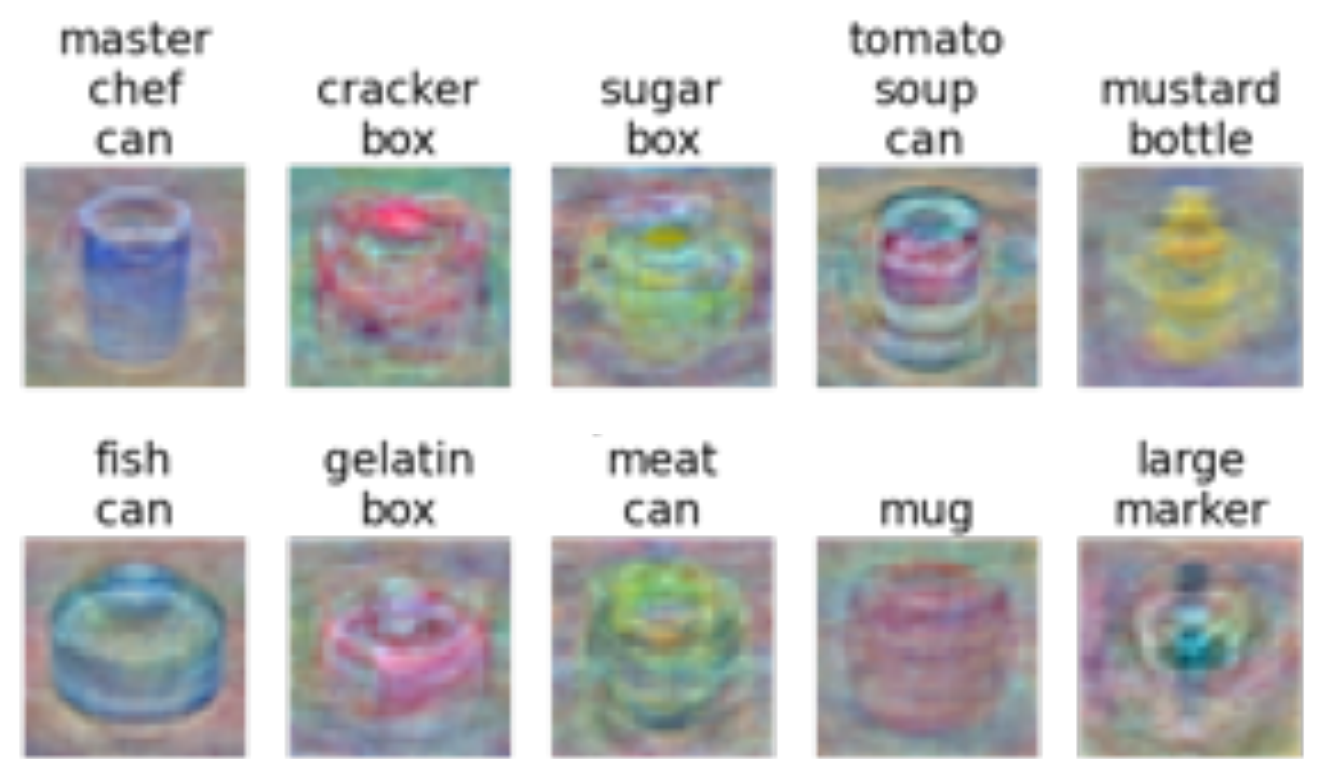
Algebraic Viewpoint

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



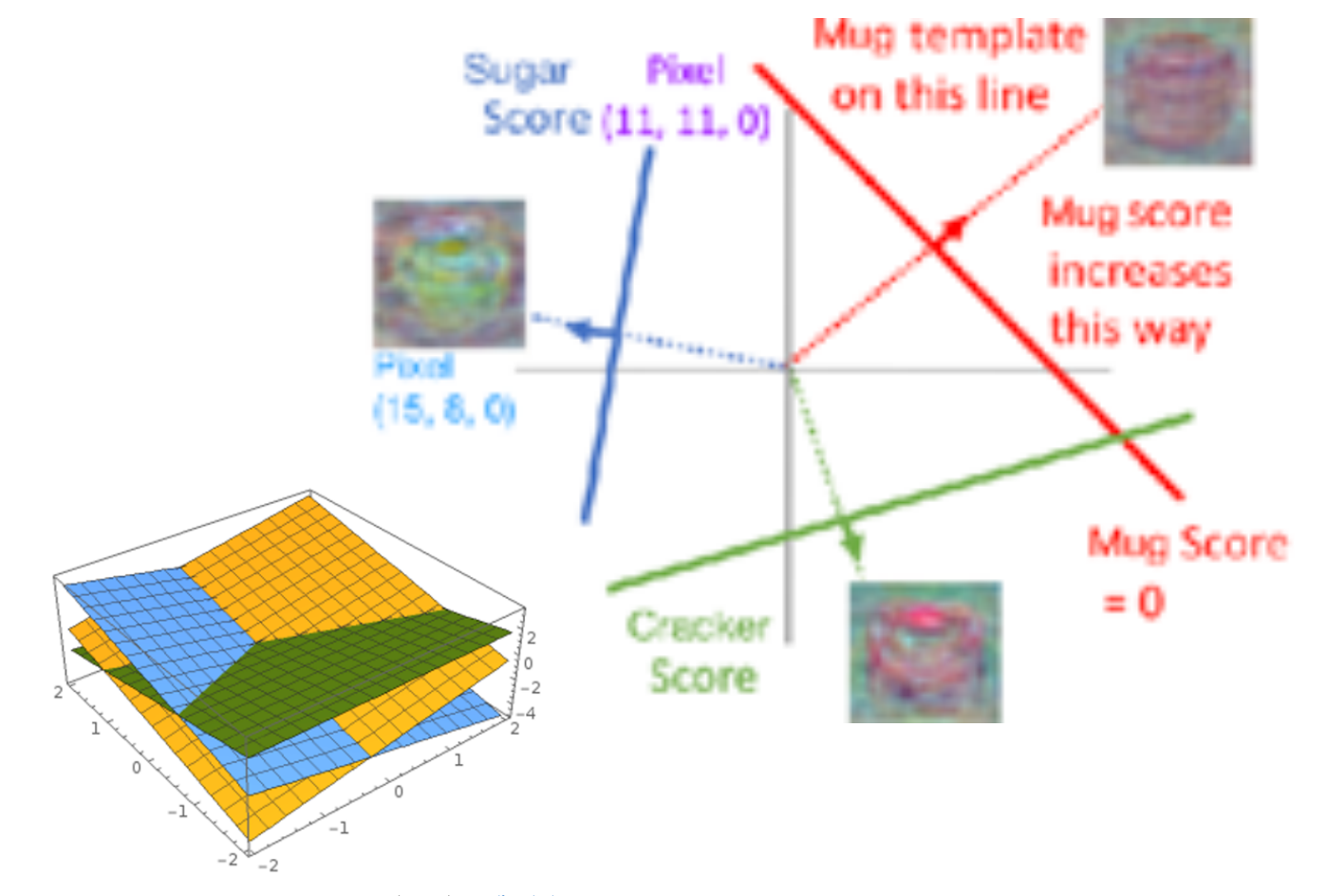
Visual Viewpoint

One template per class



Geometric Viewpoint

Hyperplanes cutting up space



So far—Defined a Score Function



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

master chef can	-3.45	-0.51	3.42
mug	-8.87	6.04	4.64
tomato soup can	0.09	5.31	2.65
cracker box	2.9	-4.22	5.1
mustard bottle	4.48	-4.19	2.64
tuna fish can	8.02	3.58	5.55
sugar box	3.78	4.49	-4.34
gelatin box	1.06	-4.37	-1.5
potted meat can	-0.36	-2.09	-4.79
large marker	-0.72	-2.93	6.14

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

But how can we actually choose a good W ?

So far—Choosing a Good W



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

master chef can	-3.45	-0.51	3.42
mug	-8.87	6.04	4.64
tomato soup can	0.09	5.31	2.65
cracker box	2.9	-4.22	5.1
mustard bottle	4.48	-4.19	2.64
tuna fish can	8.02	3.58	5.55
sugar box	3.78	4.49	-4.34
gelatin box	1.06	-4.37	-1.5
potted meat can	-0.36	-2.09	-4.79
large marker	-0.72	-2.93	6.14

TODO:

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** measures 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** measures 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** measures 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 an image and

y_i is a (discrete) label



Loss Function

A **loss function** measures 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 an image and

y_i is a (discrete) label

Loss for a single example is

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



Loss Function

A **loss function** measures 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 an image and

y_i is a (discrete) 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:

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



Cross-Entropy Loss

Multinomial Logistic Regression



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

cracker **3.2**

mug **5.1**

sugar **-1.7**

Cross-Entropy Loss

Multinomial Logistic Regression



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

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

cracker **3.2**

mug **5.1**

sugar **-1.7**

Cross-Entropy Loss

Multinomial Logistic Regression



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

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

cracker **3.2**

mug 5.1

sugar -1.7

Unnormalized log-probabilities (logits)

Cross-Entropy Loss

Multinomial Logistic Regression



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

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

Probabilities
must be ≥ 0

cracker

3.2

exp(·)

24.5

mug

5.1

164.0

sugar

-1.7

0.18

Unnormalized log-
probabilities (logits)

Unnormalized
probabilities

Cross-Entropy Loss

Multinomial Logistic Regression



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

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

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

Probabilities must be ≥ 0

Probabilities must sum to 1

cracker
mug
sugar

cracker	3.2
mug	5.1
sugar	-1.7

exp(·)

cracker	24.5
mug	164.0
sugar	0.18

normalize

cracker	0.13
mug	0.87
sugar	0.00

Unnormalized log-probabilities (logits)

Unnormalized probabilities

Probabilities



Cross-Entropy Loss

Multinomial Logistic Regression



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

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

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

Probabilities must be ≥ 0

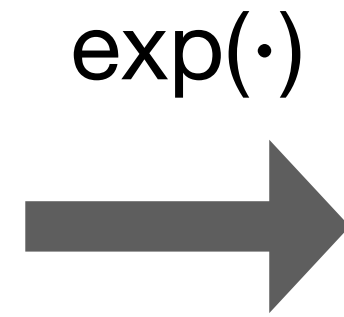
Probabilities must sum to 1

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

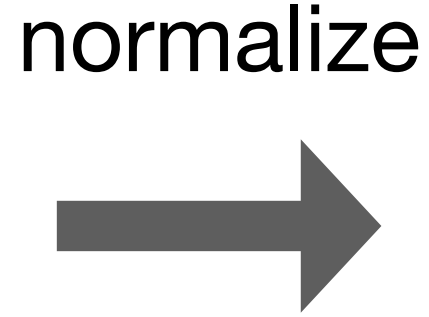
$$L_i = -\log(0.13) = 2.04$$

cracker
mug
sugar

cracker	3.2
mug	5.1
sugar	-1.7



cracker	24.5
mug	164.0
sugar	0.18



cracker	0.13
mug	0.87
sugar	0.00

Unnormalized log-probabilities (logits)

Unnormalized probabilities

Probabilities



Cross-Entropy Loss

Multinomial Logistic Regression



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

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

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

Probabilities must be ≥ 0

Probabilities must sum to 1

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

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

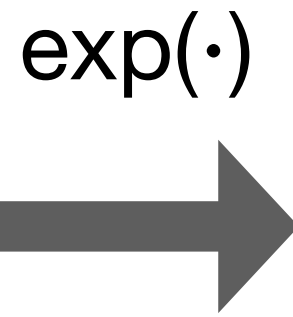
Maximum Likelihood Estimation

Choose weights to maximize the likelihood of the observed data (see CSCI 5521)

cracker
mug
sugar

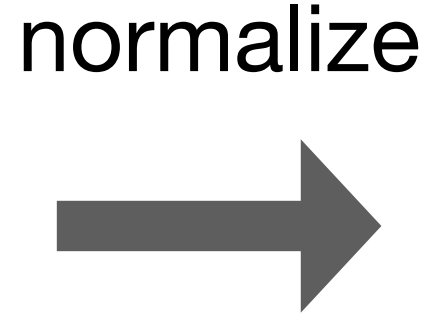
cracker	3.2
mug	5.1
sugar	-1.7

Unnormalized log-probabilities (logits)



cracker	24.5
mug	164.0
sugar	0.18

Unnormalized probabilities



cracker	0.13
mug	0.87
sugar	0.00

Probabilities

Cross-Entropy Loss

Multinomial Logistic Regression



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

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

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

Probabilities must be ≥ 0

Probabilities must sum to 1

cracker
mug
sugar

cracker	3.2
mug	5.1
sugar	-1.7

exp(.)

cracker	24.5
mug	164.0
sugar	0.18

normalize

cracker	0.13
mug	0.87
sugar	0.00

compare

cracker	1.00
mug	0.00
sugar	0.00

Unnormalized log-probabilities (logits)

Unnormalized probabilities

Probabilities

Correct probabilities



Cross-Entropy Loss

Multinomial Logistic Regression



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

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

Probabilities must be ≥ 0

Probabilities must sum to 1

cracker	3.2
mug	5.1
sugar	-1.7

Unnormalized log-probabilities (logits)

exp(·)

cracker	24.5
mug	164.0
sugar	0.18

Unnormalized probabilities

normalize

cracker	0.13
mug	0.87
sugar	0.00

Probabilities

compare

Kullback-Leibler divergence

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

cracker	1.00
mug	0.00
sugar	0.00

Correct probabilities



Cross-Entropy Loss

Multinomial Logistic Regression



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

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

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

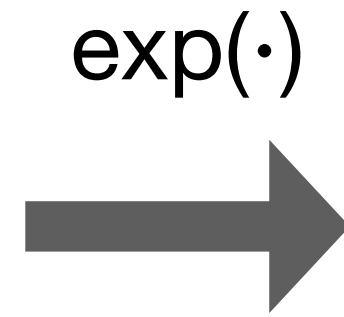
Probabilities must be ≥ 0

Probabilities must sum to 1

cracker
mug
sugar

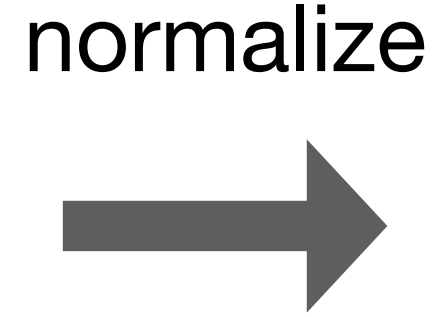
cracker	3.2
mug	5.1
sugar	-1.7

Unnormalized log-probabilities (logits)



cracker	24.5
mug	164.0
sugar	0.18

Unnormalized probabilities



cracker	0.13
mug	0.87
sugar	0.00

Probabilities



cracker	1.00
mug	0.00
sugar	0.00

Correct probabilities

Cross Entropy

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



Cross-Entropy Loss

Multinomial Logistic Regression



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

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

Maximize probability of correct class

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

Putting it all together

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

cracker **3.2**
 mug 5.1
 sugar -1.7

Cross-Entropy Loss

Multinomial Logistic Regression



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

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

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 ?

cracker	3.2
mug	5.1
sugar	-1.7

Cross-Entropy Loss

Multinomial Logistic Regression



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

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

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: $+\infty$

cracker **3.2**
mug **5.1**
sugar **-1.7**

Cross-Entropy Loss

Multinomial Logistic Regression



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

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

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 small random values, what is the loss?

cracker **3.2**
mug **5.1**
sugar **-1.7**

Cross-Entropy Loss

Multinomial Logistic Regression



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

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

Maximize probability of correct class

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

Putting it all together

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

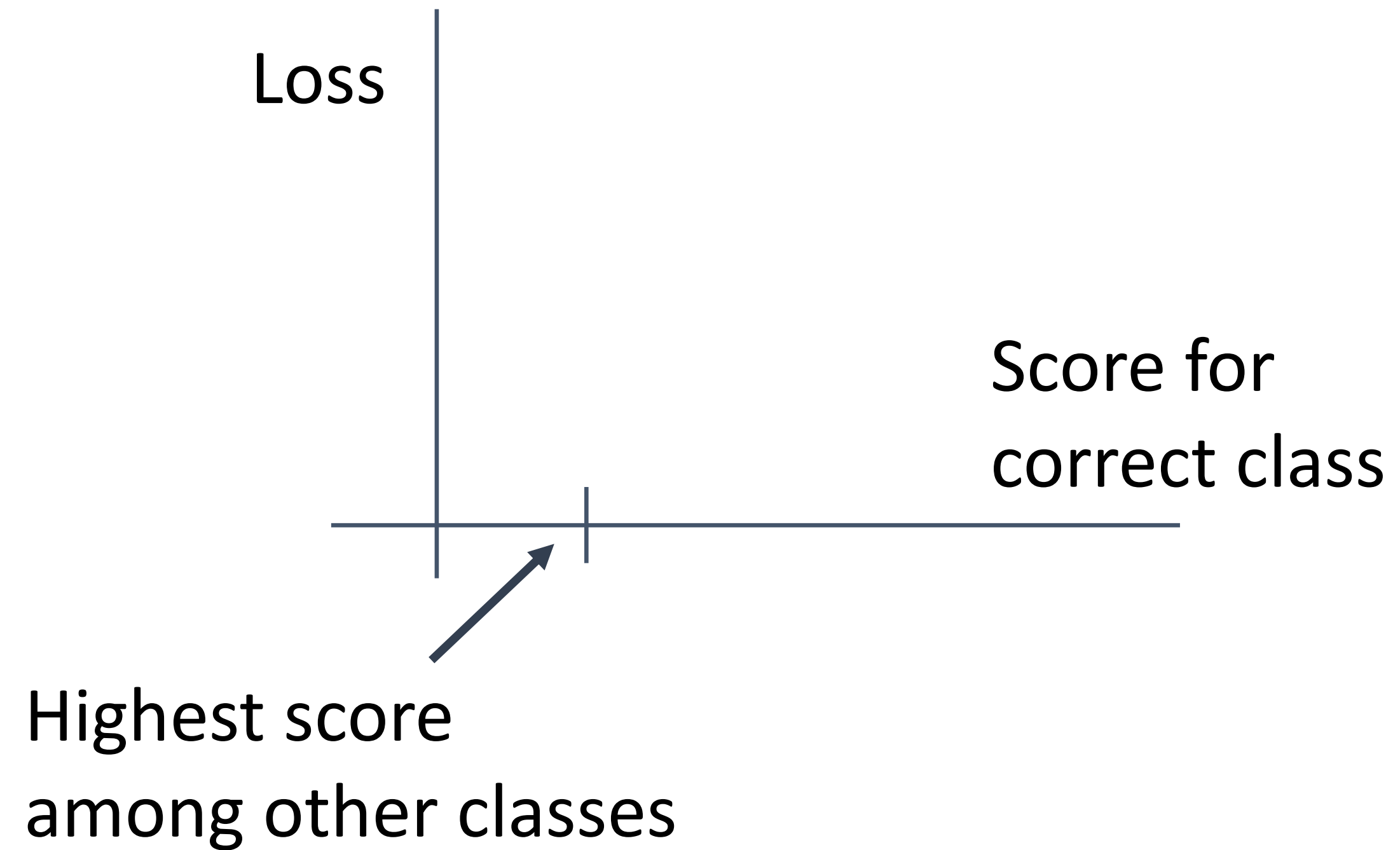
Q: If all scores are small random values, what is the loss?

A: $-\log\left(\frac{1}{C}\right)$

$$\log\left(\frac{1}{10}\right) \approx 2.3$$

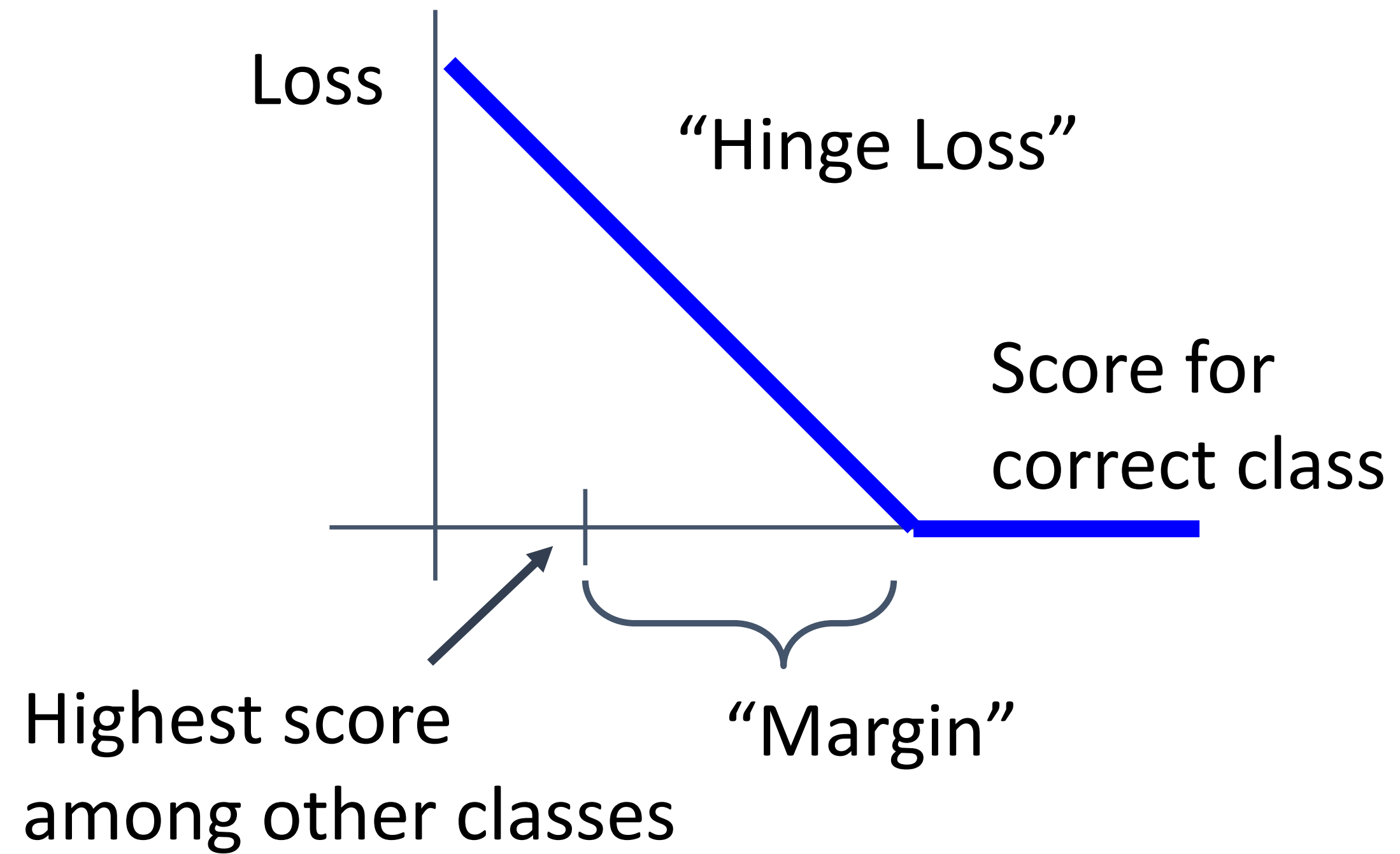
Multiclass SVM Loss

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



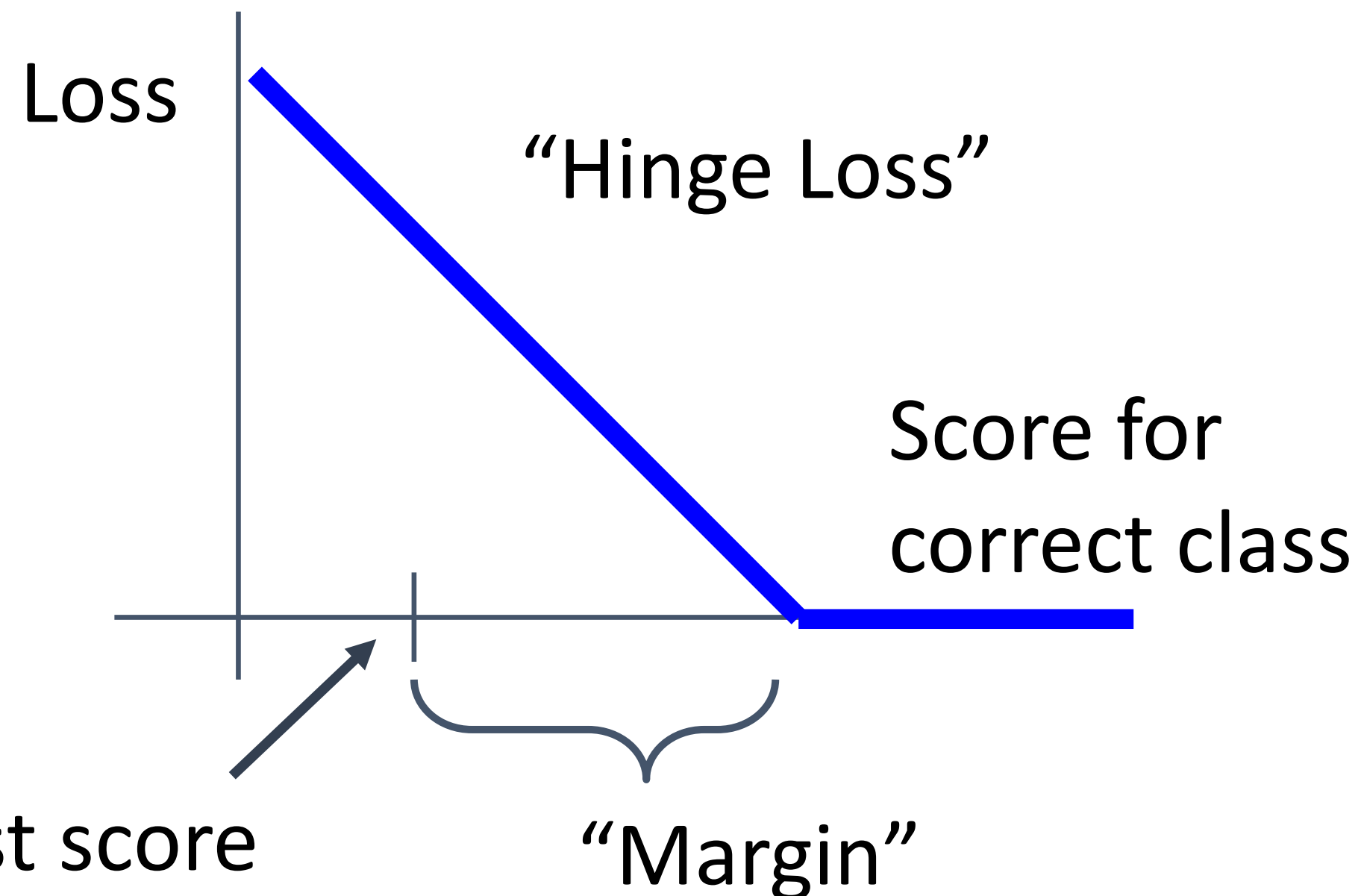
Multiclass SVM Loss

“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”



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



cracker	3.2	1.3	2.2
mug	5.1	4.9	2.5
sugar	-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



cracker	3.2	1.3	2.2
mug	5.1	4.9	2.5
sugar	-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:

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

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

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

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

$$= 2.9 + 0$$

$$= 2.9$$

Multiclass SVM Loss



cracker	3.2	1.3	2.2
mug	5.1	4.9	2.5
sugar	-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:

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

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

Multiclass SVM Loss



cracker	3.2	1.3	2.2
mug	5.1	4.9	2.5
sugar	-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



cracker	3.2	1.3	2.2
mug	5.1	4.9	2.5
sugar	-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:

$$L = (2.9 + 0.0 + 12.9) / 3 \\ = 5.27$$

Multiclass SVM Loss



cracker	3.2	1.3	2.2
mug	5.1	4.9	2.5
sugar	-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:

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

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

Multiclass SVM Loss



cracker	3.2	1.3	2.2
mug	5.1	4.9	2.5
sugar	-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:

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

Q2: What are the min
and max possible loss?

Multiclass SVM Loss



cracker	3.2	1.3	2.2
mug	5.1	4.9	2.5
sugar	-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:

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

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

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

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and $y_i = 0$

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

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 for 1st and 3rd example
SVM loss will change for the 2nd

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

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

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and $y_i = 0$

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

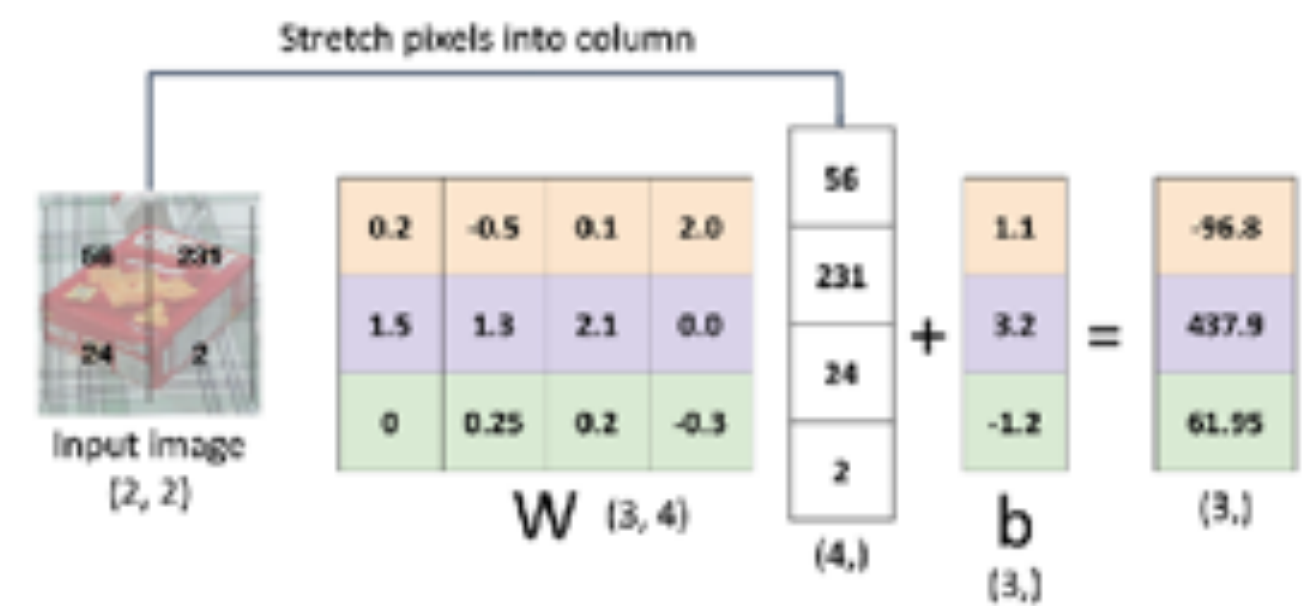
Q: What happens to each loss if I double the score of the correct class from 10 to 20?

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

Recap—Three Ways to Interpret Linear Classifiers

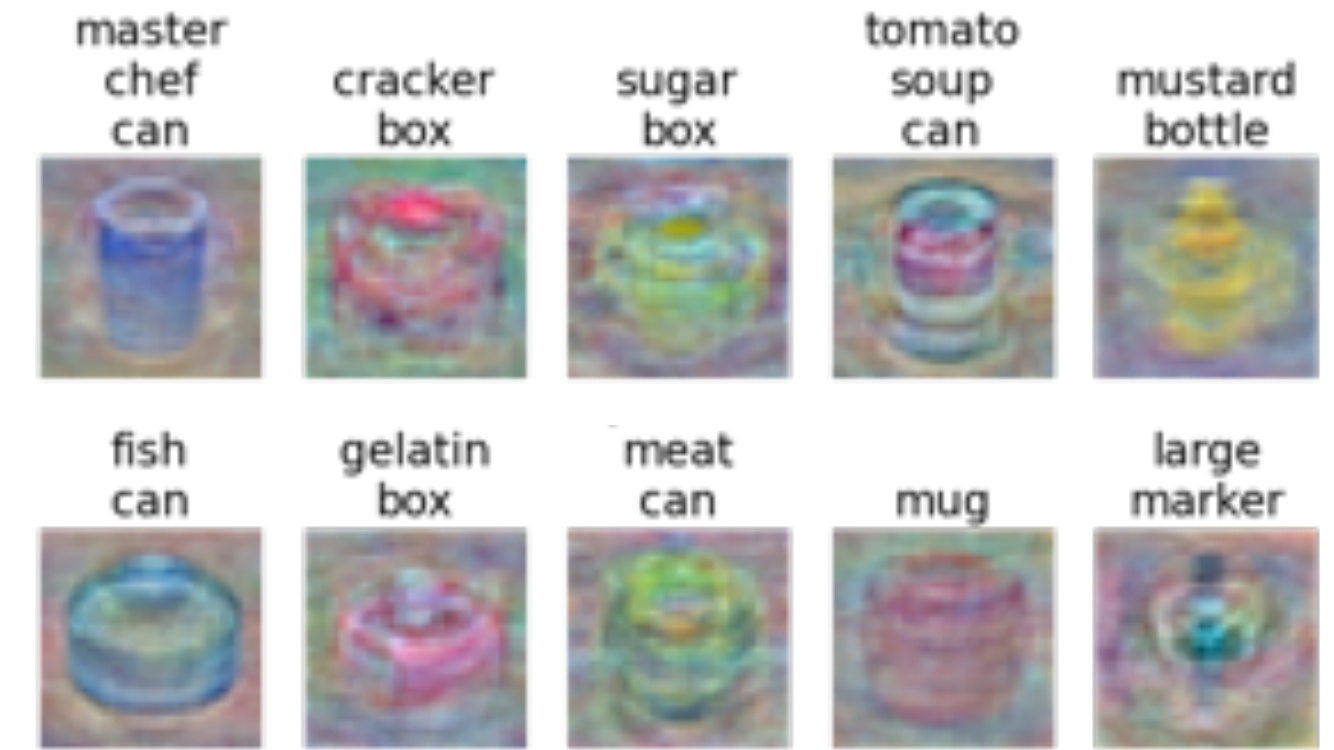
Algebraic Viewpoint

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



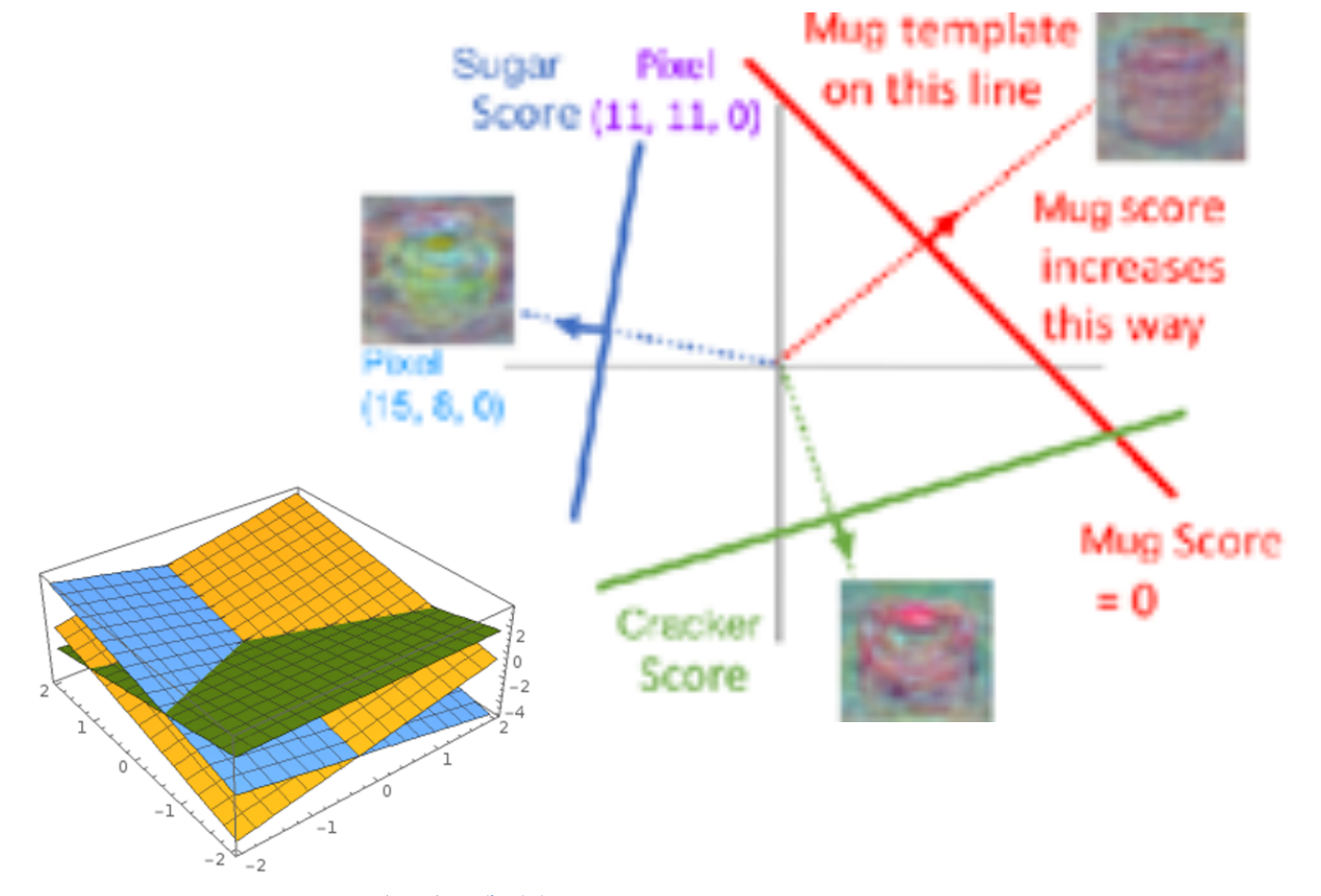
Visual Viewpoint

One template per class



Geometric Viewpoint

Hyperplanes cutting up space



Recap—Loss Functions Quantify Preferences

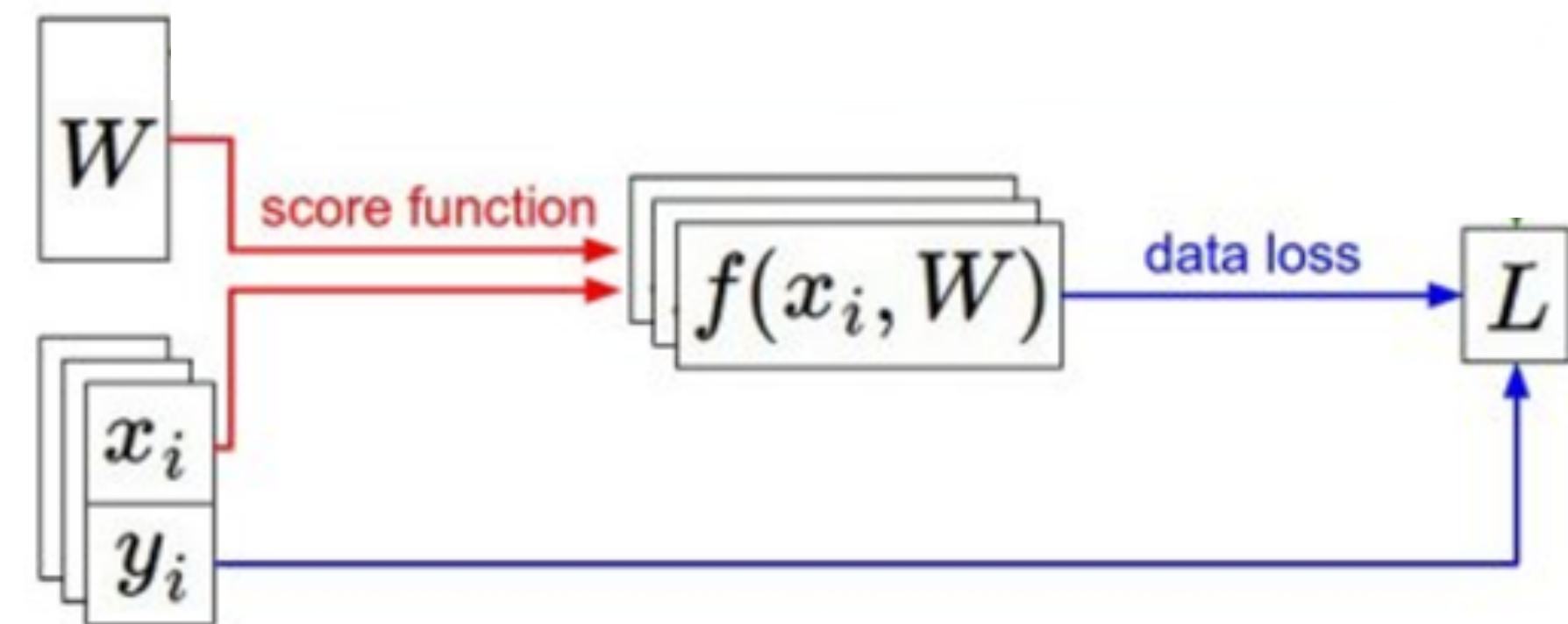
- We have some dataset of (x, y)
- We have a **score function**:
- We have a **loss function**:

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

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

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

Linear classifier



Recap—Loss Functions Quantify Preferences

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

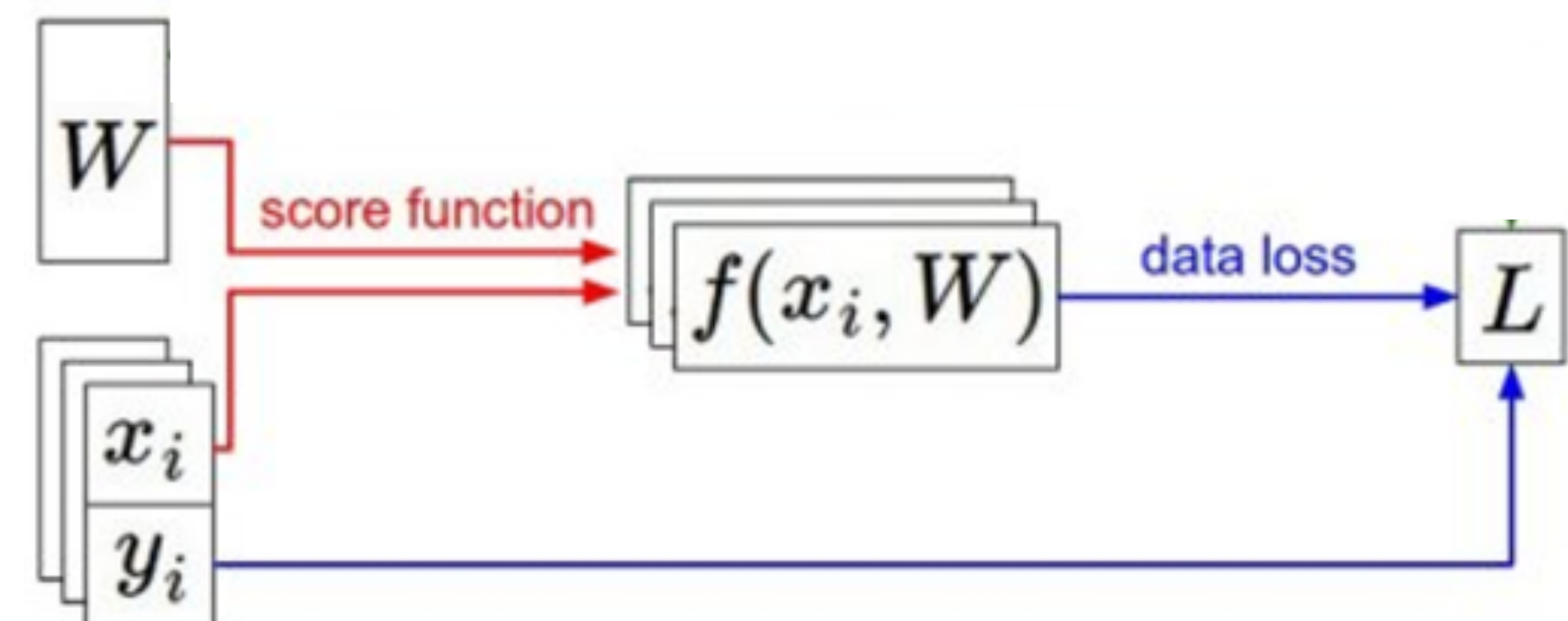
Softmax: $L_i = -\log \left(\frac{\exp(s_{y_i})}{\sum_j \exp(s_j)} \right)$

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

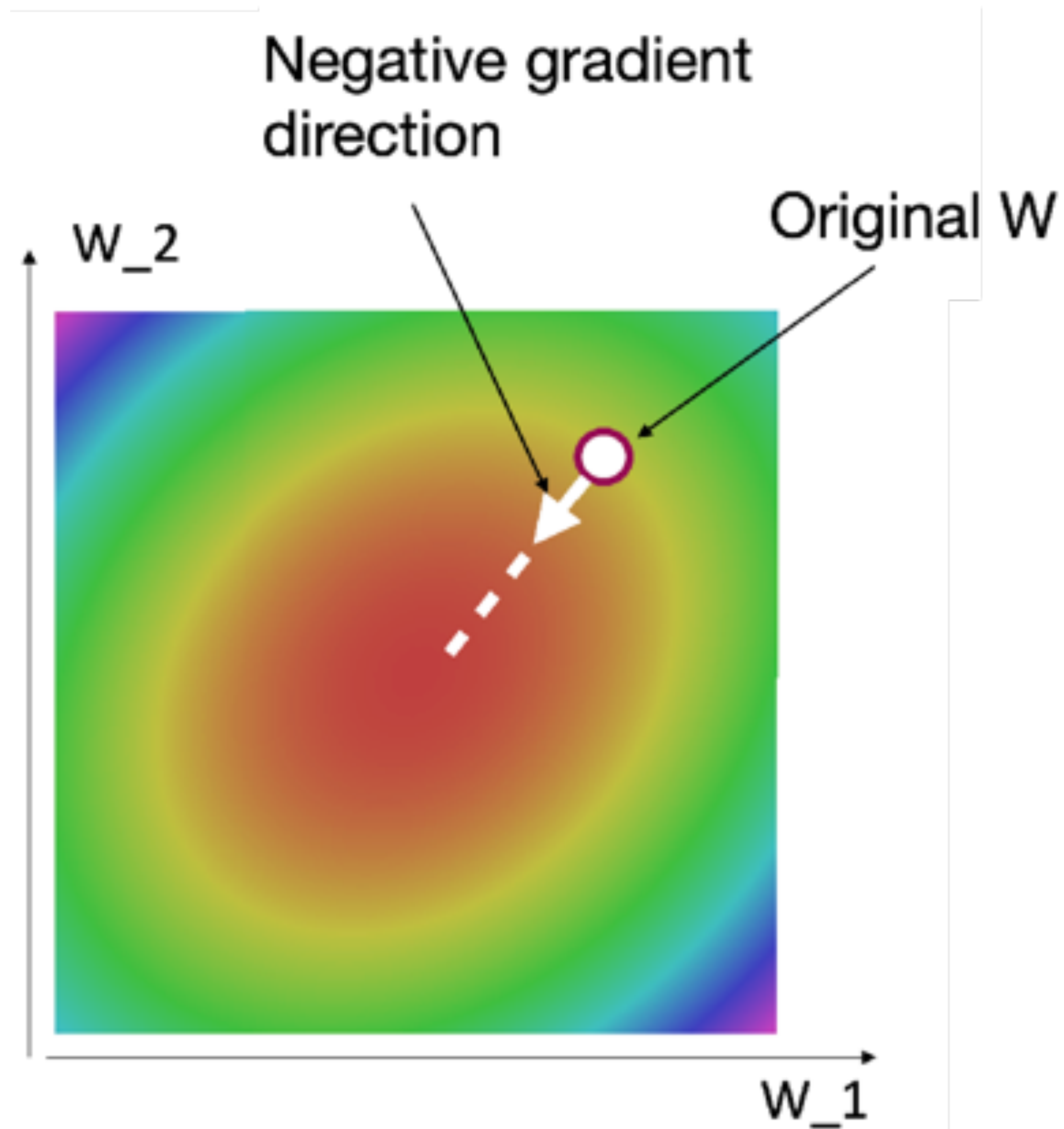
Q: How do we find the best W, b ?

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

Linear classifier



Next time: Regularization + Optimization





DeepRob

Lecture 3
Linear Classifiers
University of Michigan and University of Minnesota

