

# Image Classification

# Image Classification



This image by Nikita is  
licensed under CC-BY 2.0

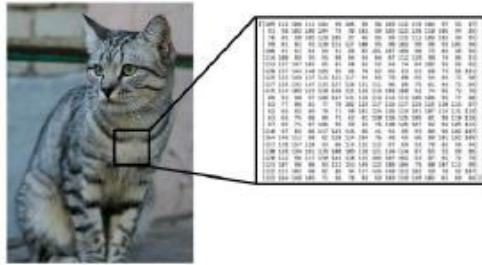
(assume given a set of labels)  
{dog, cat, truck, plane, ...}



cat  
dog  
bird  
deer  
truck

# Challenges of recognition

Viewpoint



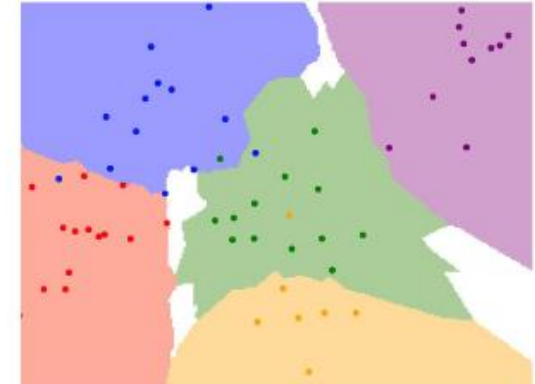
# Data-driven approach, KNN



1-NN classifier



5-NN classifier



# Linear Classifier



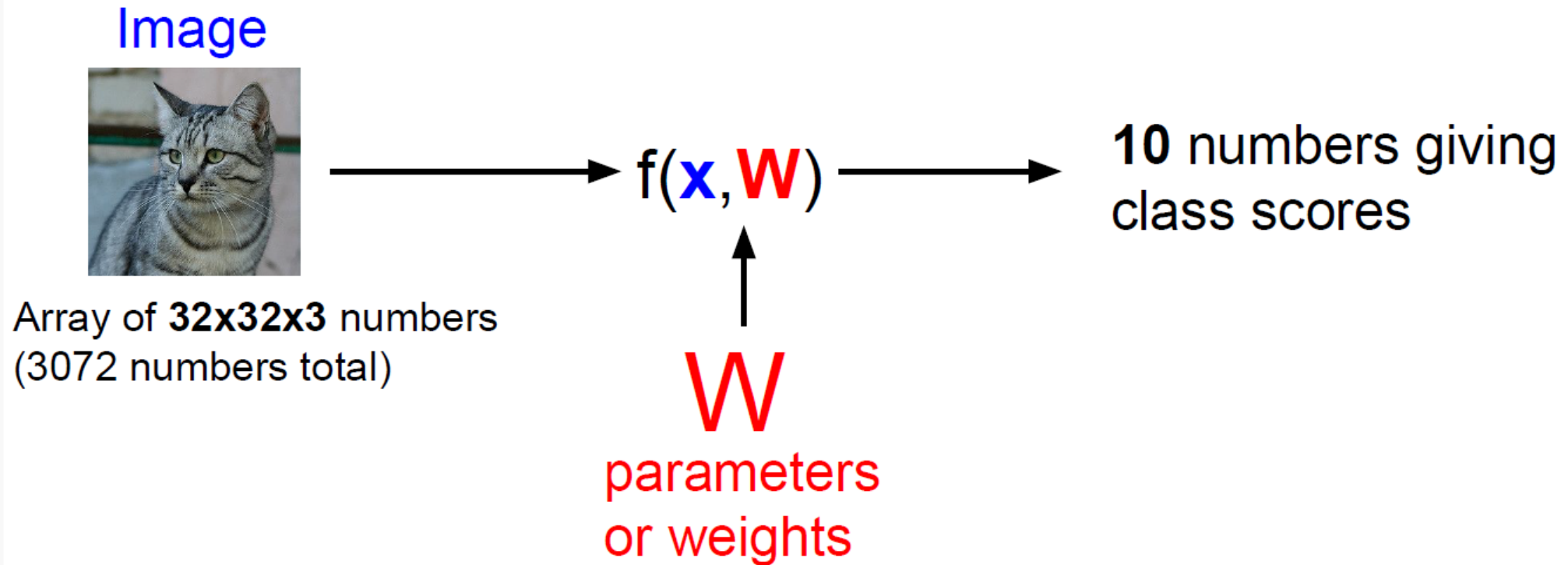
# Neural Network

Linear  
classifiers



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

# Parametric Approach

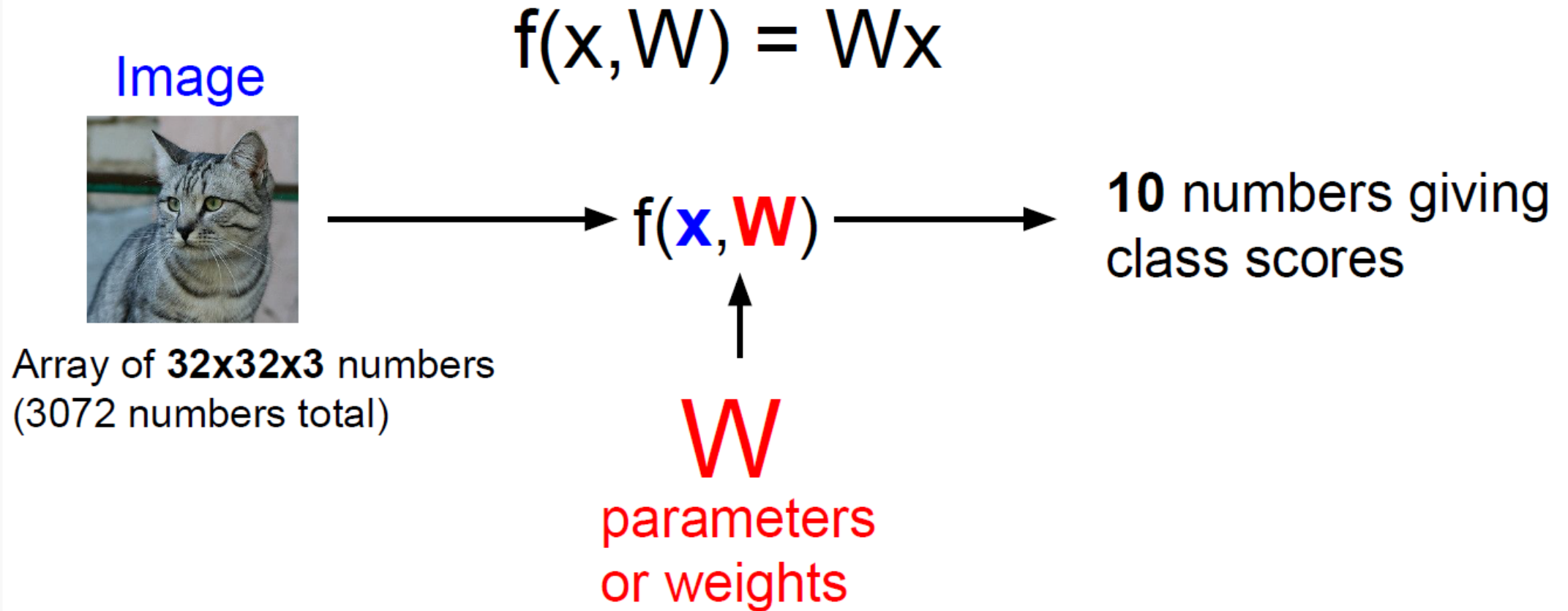


# Parametric Approach

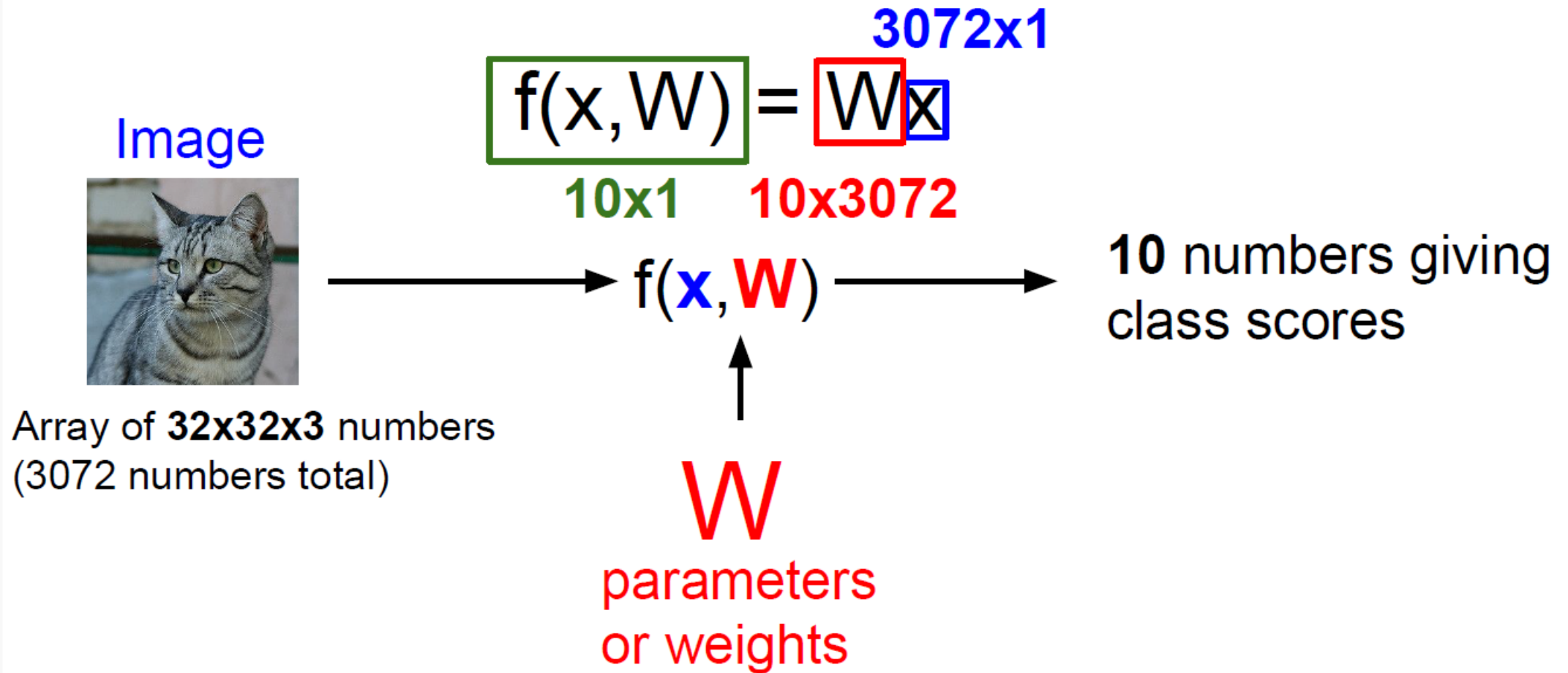
- In image classification, the **parametric approach** refers to the use of a predefined model structure with a fixed set of parameters.
- By summarizing the information into a set of parameters, the parametric approach allows the model to be more efficient and scalable.
- Once the parameters are learned from the training data, the model can quickly make predictions for new, unseen images by applying the learned parameters to the input data.
- This is in contrast to the kNN approach, where the entire dataset needs to be stored and searched for each prediction, making it computationally expensive and less efficient for large datasets.
- We will start out with arguably the simplest possible function, a linear mapping



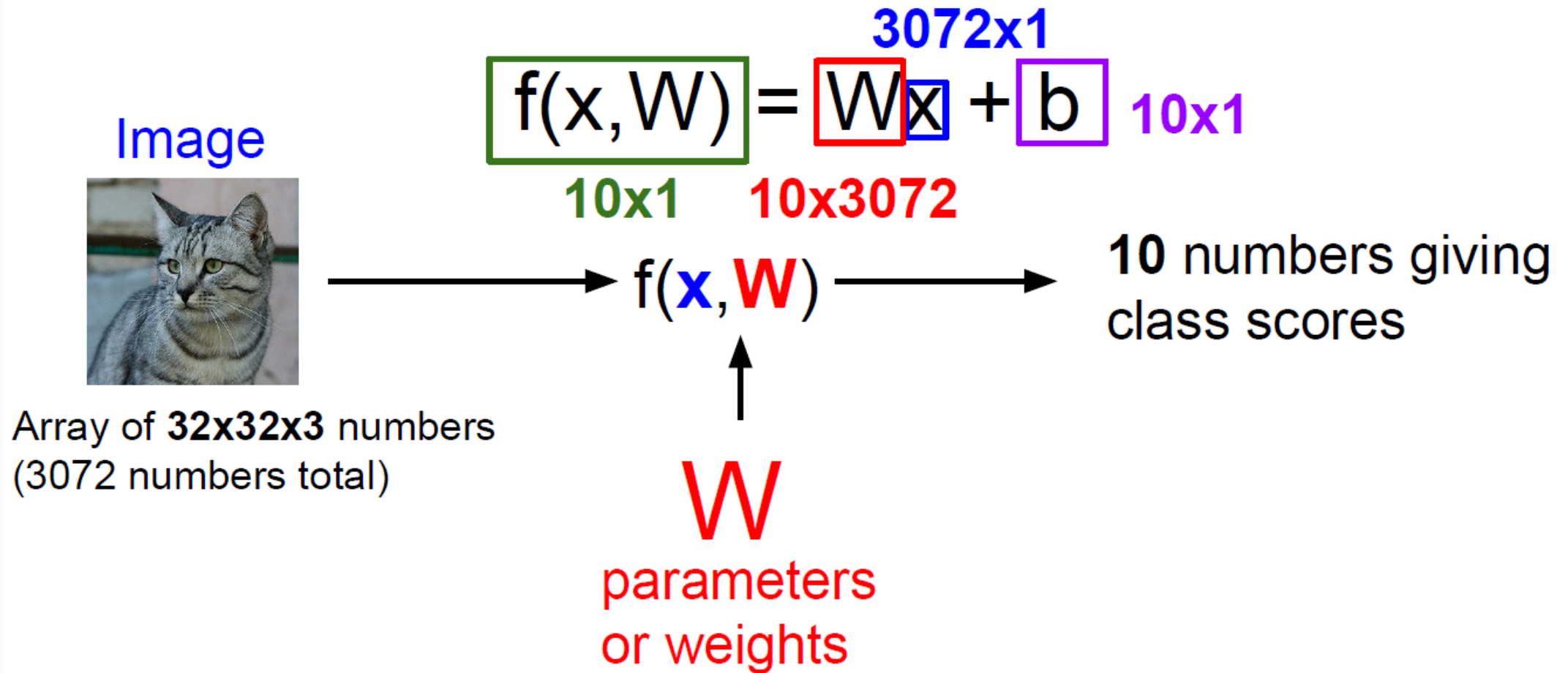
# Parametric Approach: Linear Classifier



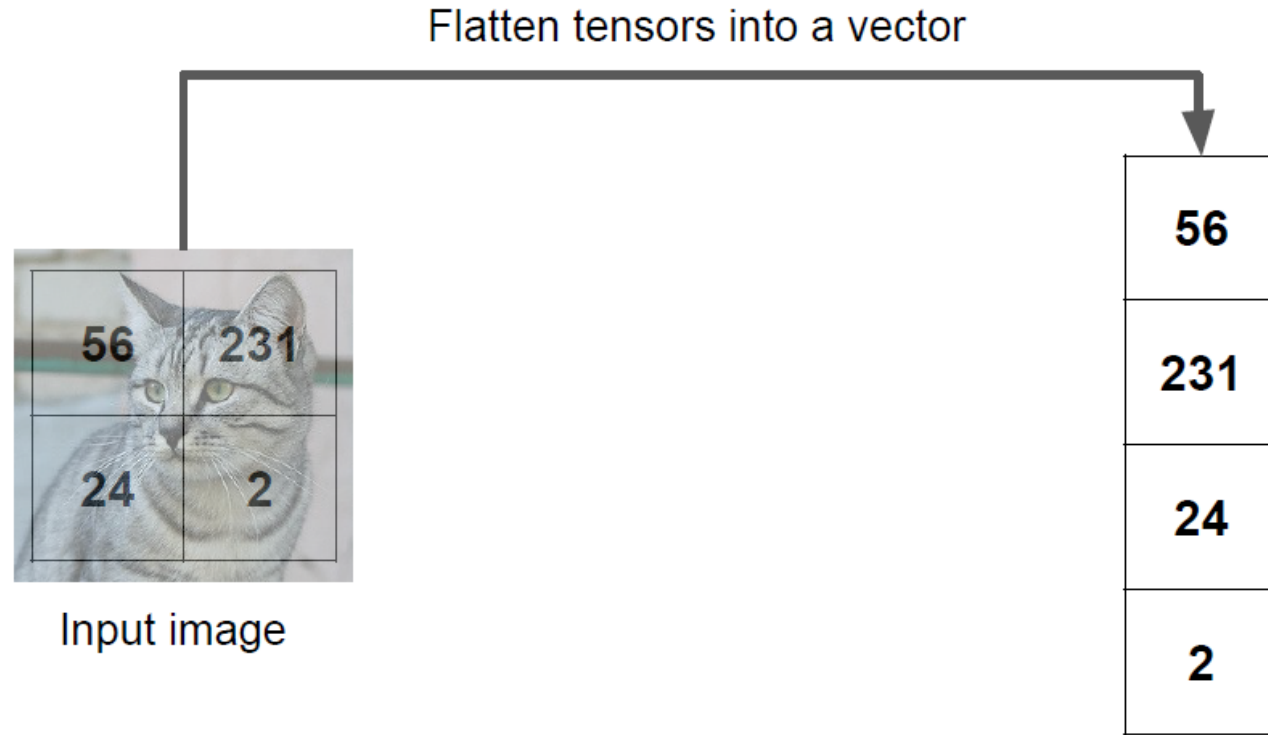
# Parametric Approach: Linear Classifier



# Parametric Approach: Linear Classifier

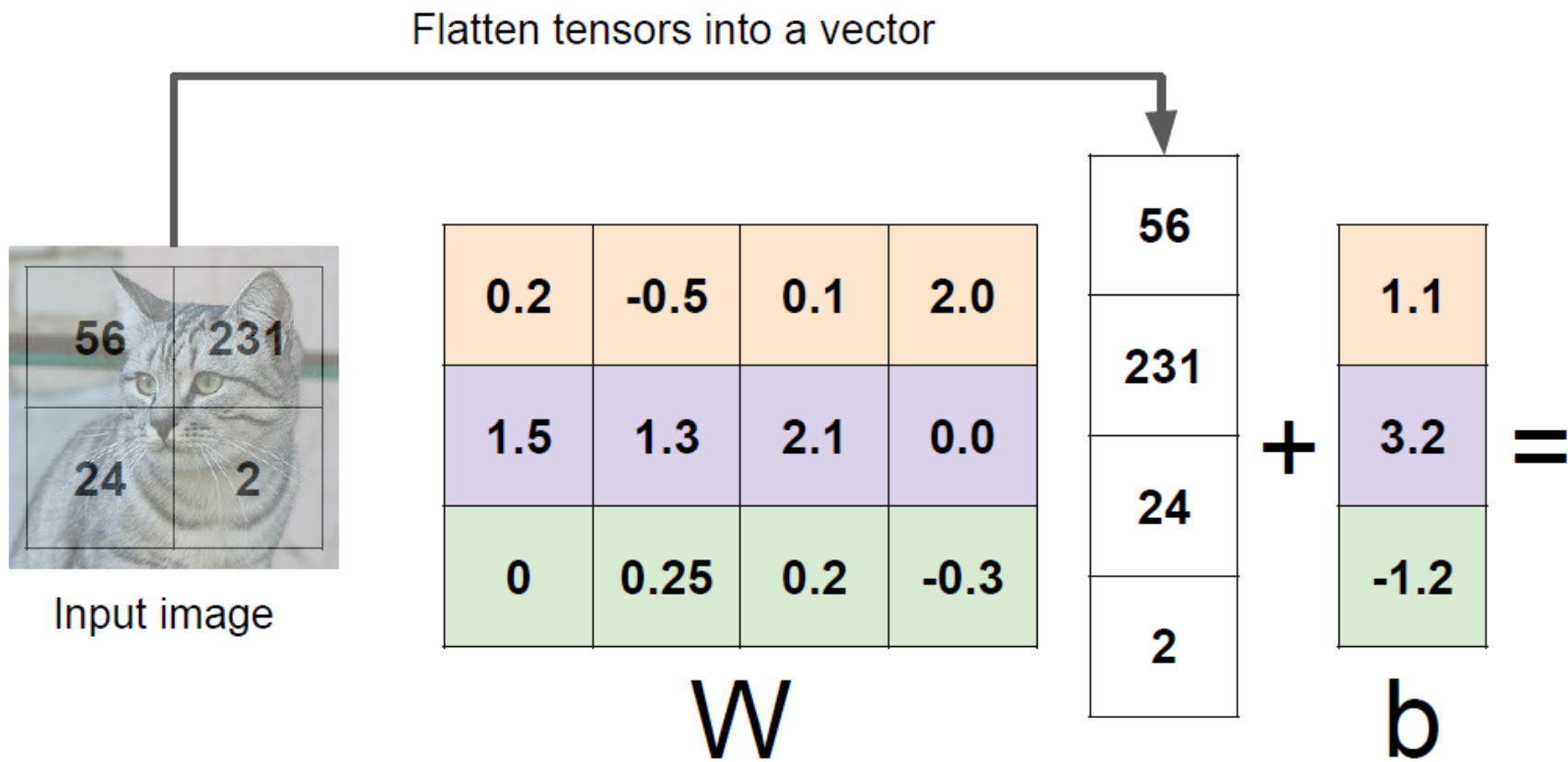


Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



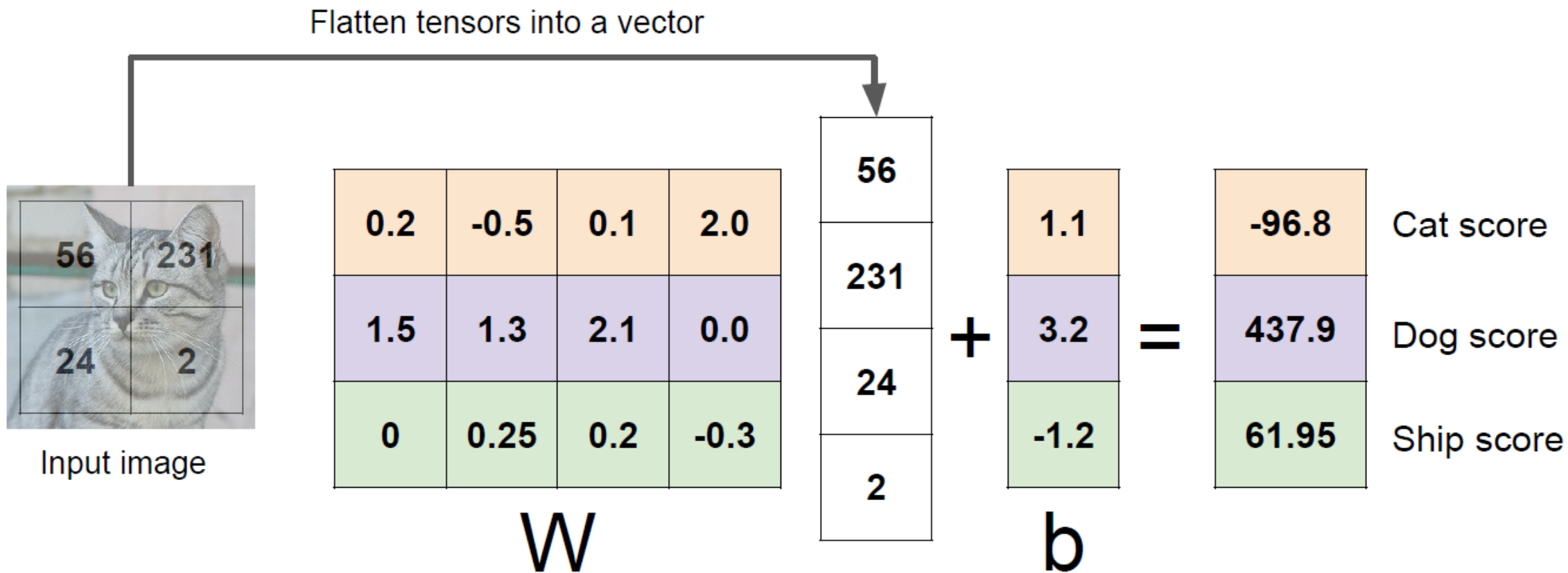
# Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

## Algebraic Viewpoint



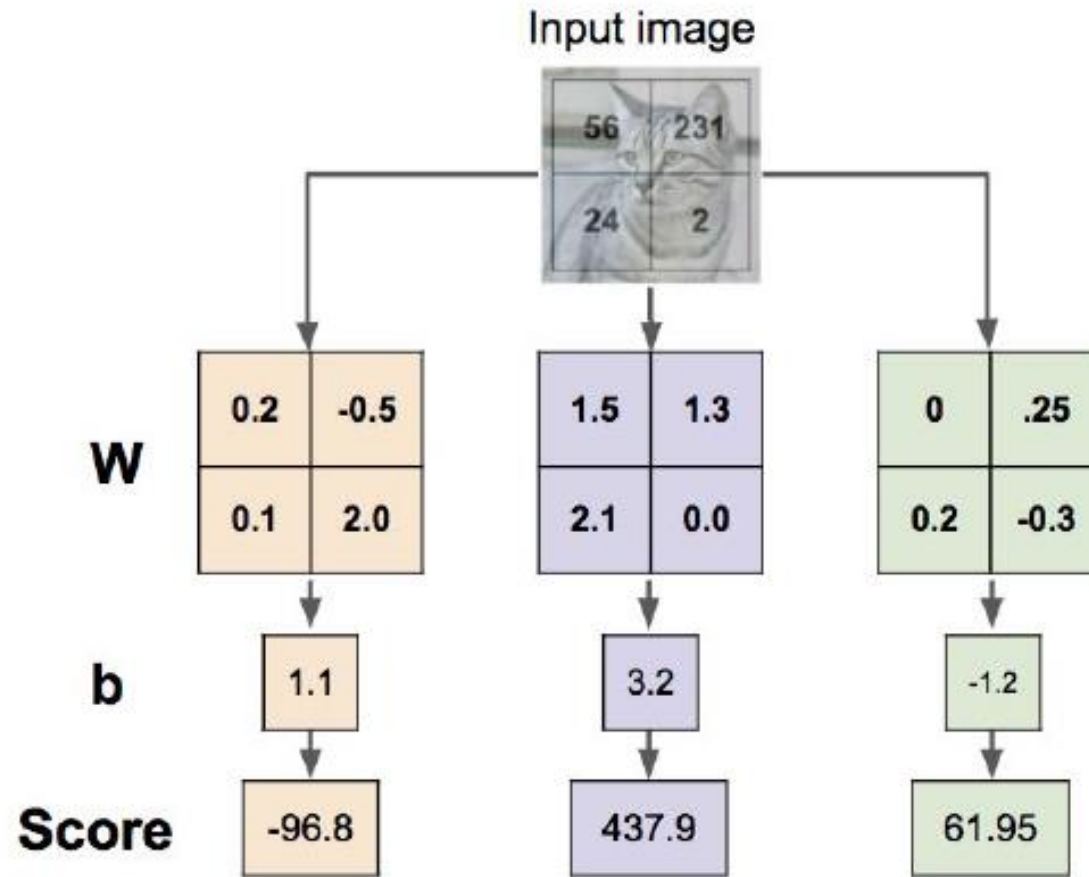
# Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

## Algebraic Viewpoint

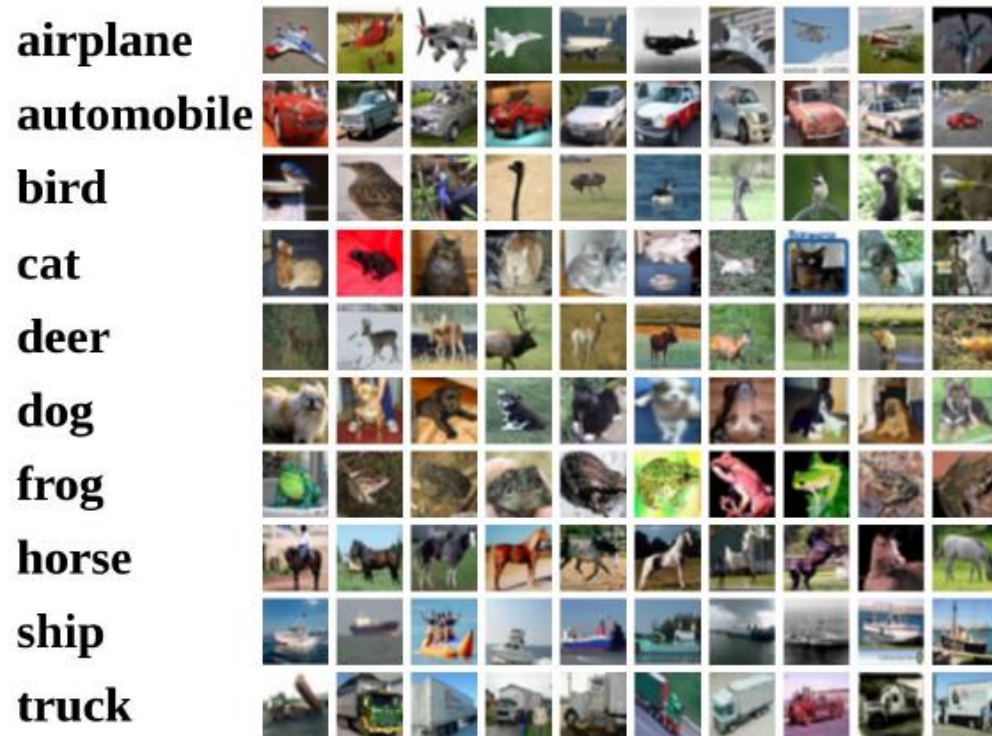




# Interpreting a Linear Classifier: Visual Viewpoint



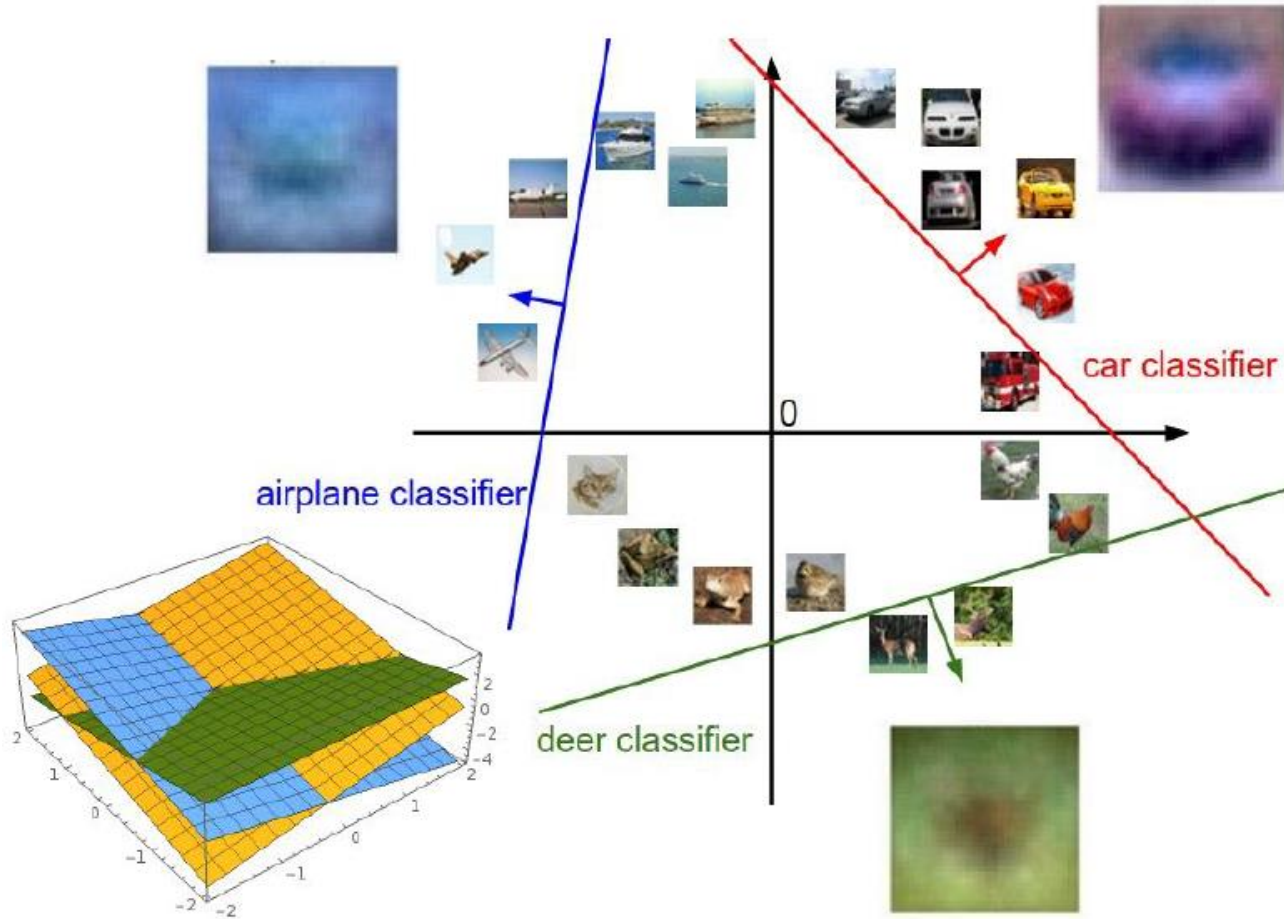
# Interpreting a Linear Classifier: Visual Viewpoint



Example trained weights of  
linear classifier on CIFAR-10



# Interpreting a Linear Classifier: Geometric Viewpoint



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



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

# Linear classifier

- Now we need ...
  - A **score function** that maps the raw data to class scores
  - Define a **loss function** that quantifies the agreement between the predicted scores and the ground truth labels.
  - Come up with a way of efficiently finding the parameters that minimize the loss function. (Optimization)

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>



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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

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

Given a dataset of examples

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

Where  $x_i$  is image and  
 $y_i$  is (integer) label

Loss over the dataset is a  
average of loss over examples:

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



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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

## Multiclass SVM loss:

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

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

the SVM loss has the form:

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

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

## Multiclass SVM loss:

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

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

the SVM loss has the form:

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

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	<b>2.9</b>		

## Multiclass SVM loss:

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

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

the SVM loss has the form:

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



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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	2.9	0	

## Multiclass SVM loss:

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

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

the SVM loss has the form:

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

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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	2.9	0	<b>12.9</b>

## Multiclass SVM loss:

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

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

the SVM loss has the form:

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



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



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	2.9	0	12.9

## Multiclass SVM loss:

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

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

the SVM loss has the form:

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

Loss over full dataset is average:

$$L = \frac{1}{N} \sum_{i=1}^N L_i$$

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



# Softmax classifier (Multinomial Logistic Regression)



cat	<b>3.2</b>
car	5.1
frog	-1.7

# Softmax classifier (Multinomial Logistic Regression)



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

cat	<b>3.2</b>
car	5.1
frog	-1.7

# Softmax classifier (Multinomial Logistic Regression)



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

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

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

Softmax  
Function

cat	<b>3.2</b>
car	5.1
frog	-1.7

# Softmax classifier (Multinomial Logistic Regression)



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

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

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

Probabilities  
must be  $\geq 0$

cat	3.2	exp →	24.5
car	5.1		164.0
frog	-1.7		0.18

unnormalized probabilities

# Softmax classifier (Multinomial Logistic Regression)



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

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

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

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

cat	3.2
car	5.1
frog	-1.7

exp

24.5
164.0
0.18

normalize

0.13
0.87
0.00

unnormalized  
probabilities

probabilities

# Softmax classifier (Multinomial Logistic Regression)



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

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

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

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

cat  
car  
frog

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



# Softmax classifier (Multinomial Logistic Regression)



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

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

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

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

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

cat  
car  
frog

**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

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

# Softmax classifier (Multinomial Logistic Regression)



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

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

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

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

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

cat  
car  
frog

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

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

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

# Softmax classifier (Multinomial Logistic Regression)



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

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

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

Softmax  
Function

Maximize probability of correct class

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

cat            **3.2**

car            **5.1**

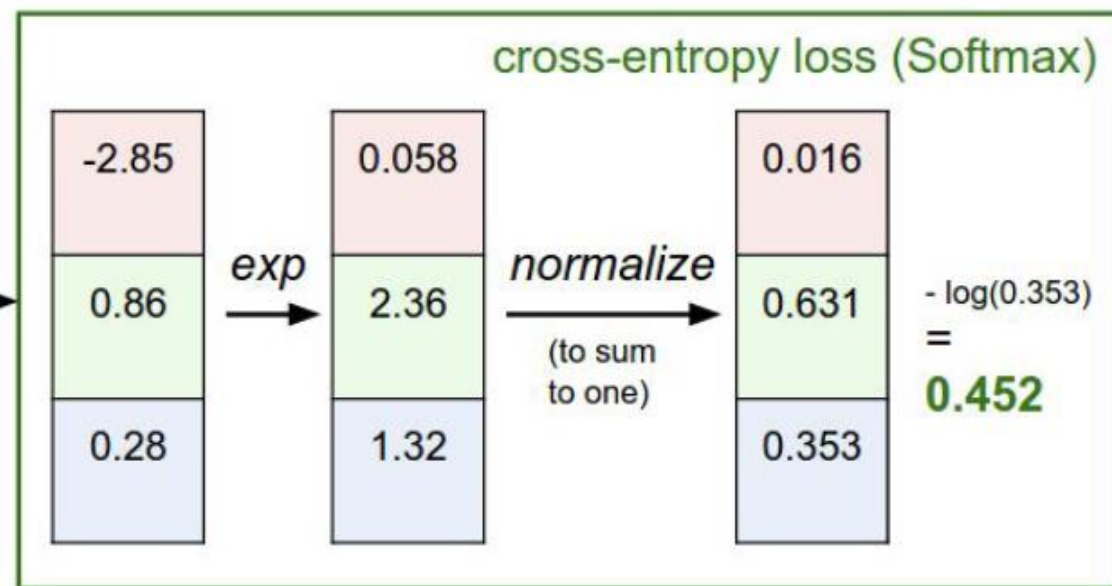
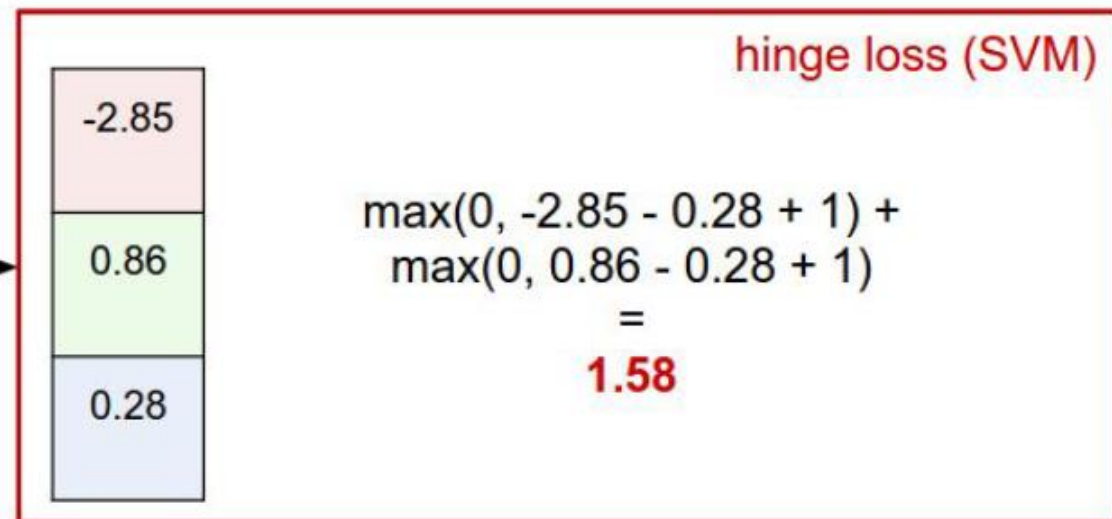
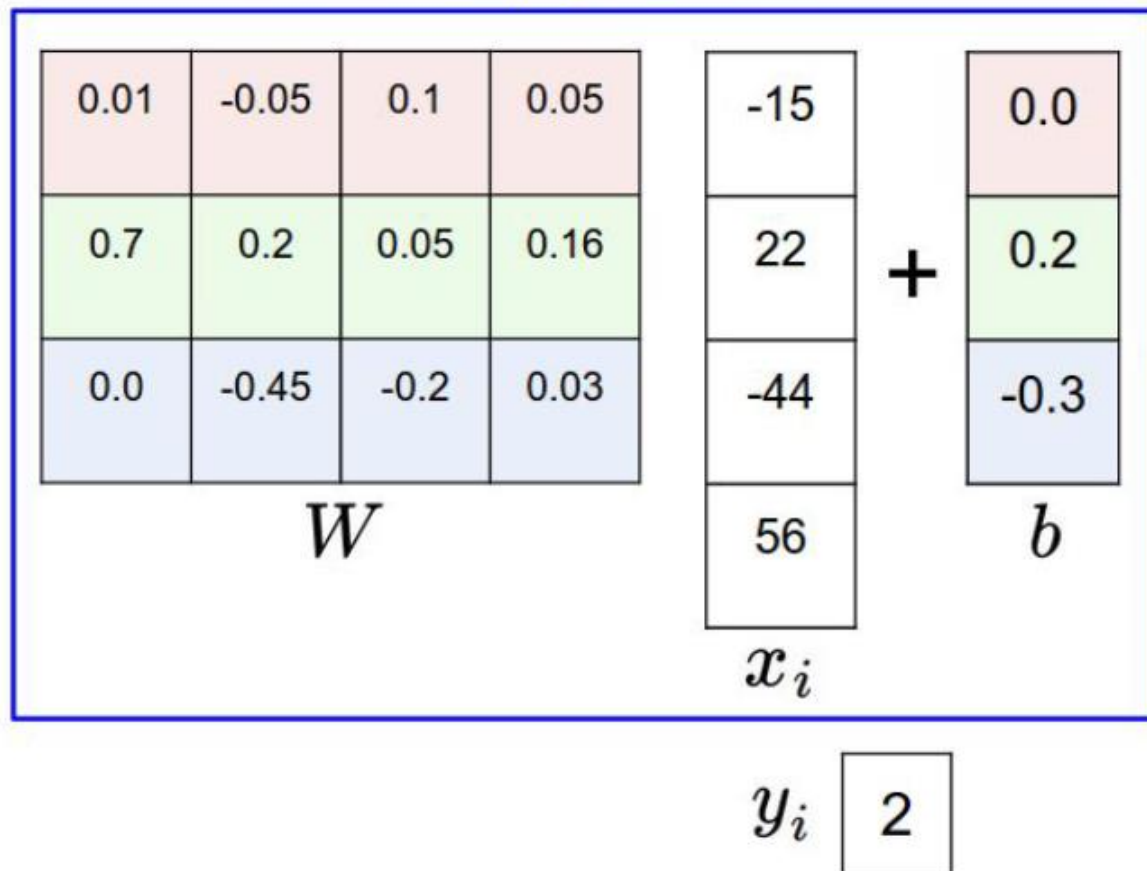
frog           **-1.7**

Putting it all together:

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

# Softmax vs. SVM

matrix multiply + bias offset



# Softmax classifier (Multinomial Logistic Regression)

- In practice, SVM and Softmax are usually comparable. The performance difference between the SVM and Softmax are usually very small, and different people will have different opinions on which classifier works better.
- Compared to the Softmax classifier, the SVM is a more local objective
- the Softmax classifier is never fully happy with the scores it produces: the correct class could always have a higher probability and the incorrect classes always a lower probability and the loss would always get better.
- However, the SVM is happy once the margins are satisfied, and it does not micromanage the exact scores beyond this constraint.



# So Far ...

- We defined a **score function** from image pixels to class scores (in this section, a linear function that depends on weights  $W$  and biases  $b$ ).
- Unlike kNN classifier, the advantage of this **parametric approach** is that once we learn the parameters we can discard the training data. Additionally, the prediction for a new test image is fast since it requires a single matrix multiplication with  $W$ , not an exhaustive comparison to every single training example.
- We defined a **loss function** (we introduced two commonly used losses for linear classifiers: the SVM and the Softmax) that measures how compatible a given set of parameters is with respect to the ground truth labels in the training dataset. We also saw that the loss function was defined in such way that making good predictions on the training data is equivalent to having a small loss.



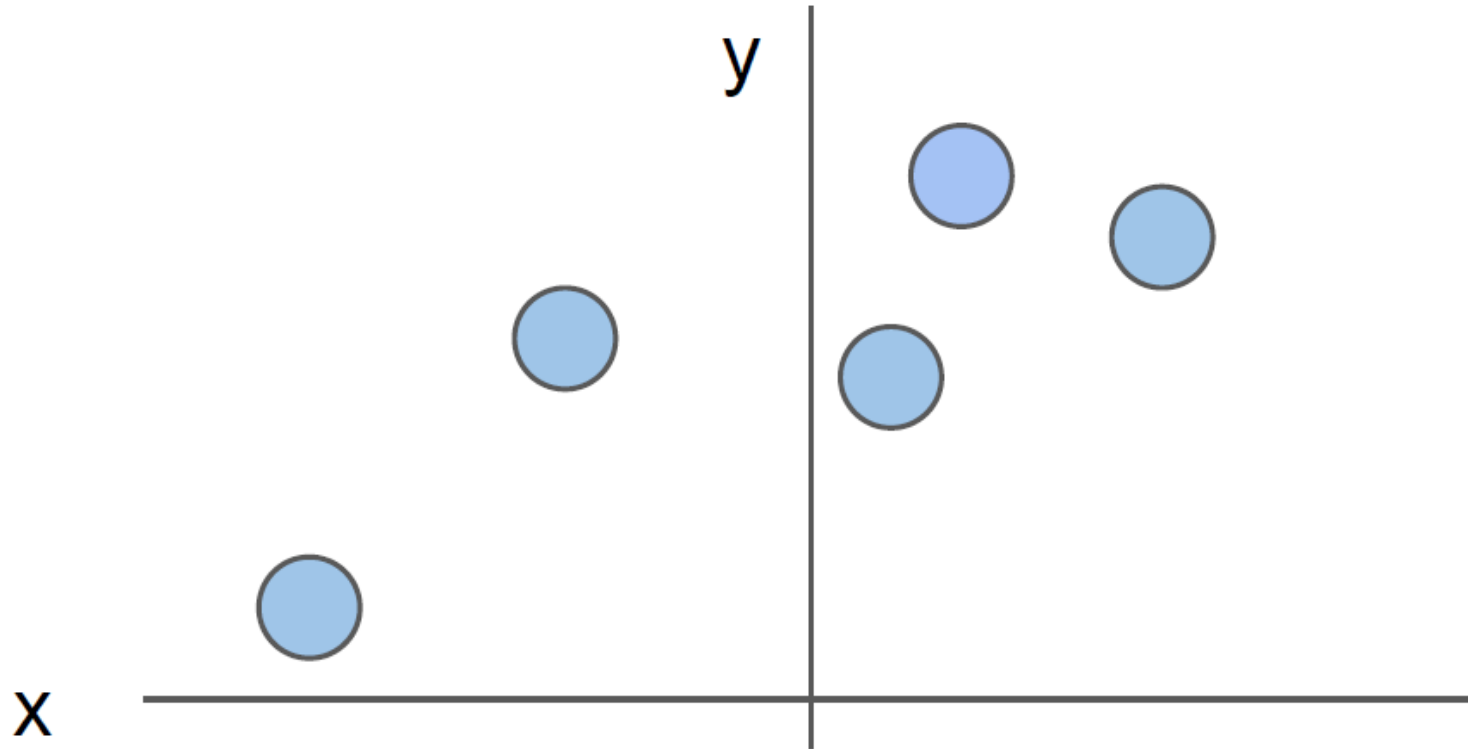
# Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss}} + \underbrace{\lambda R(W)}_{\text{Regularization}}$$

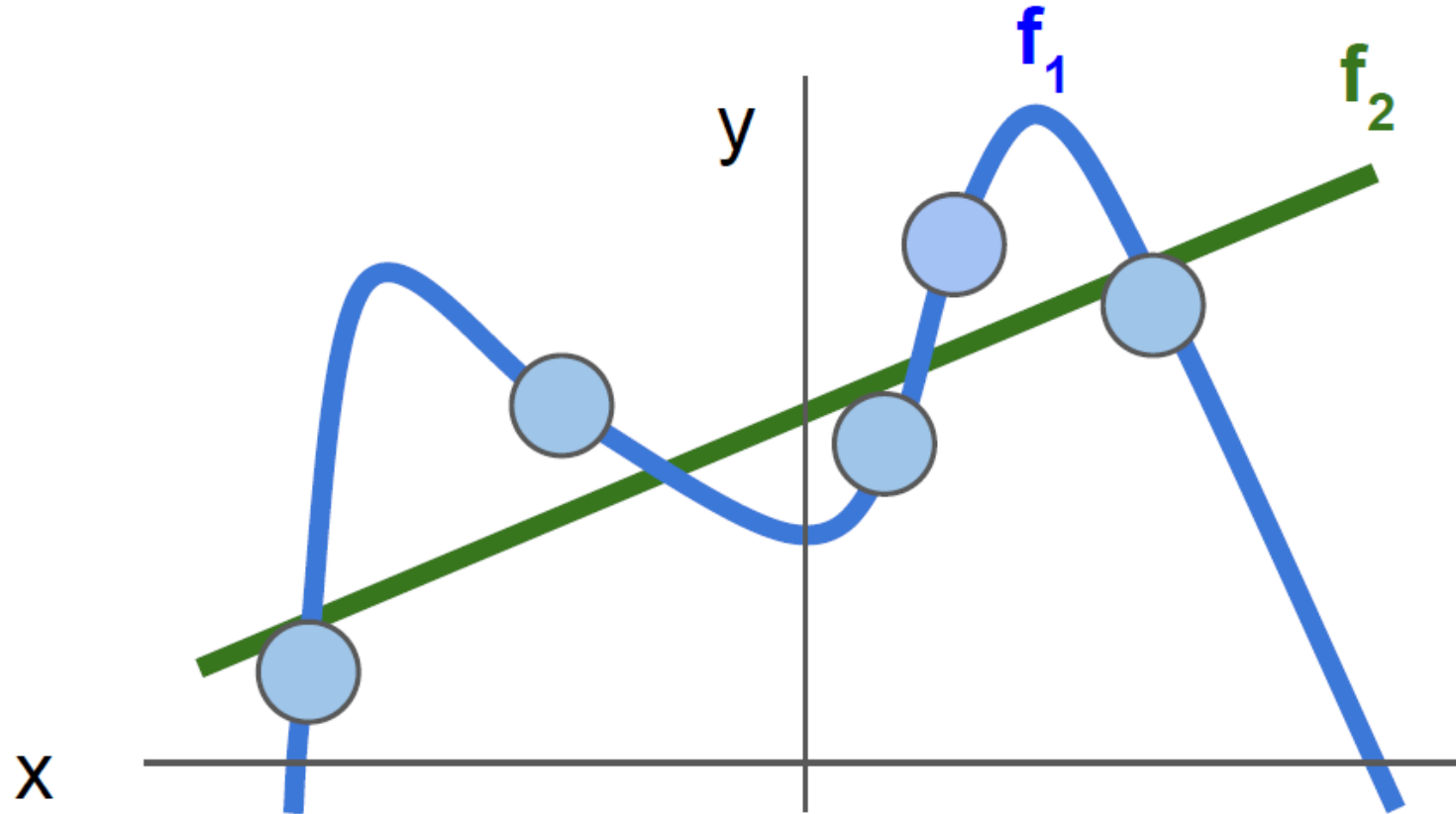
**Data loss:** Model predictions should match training data

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

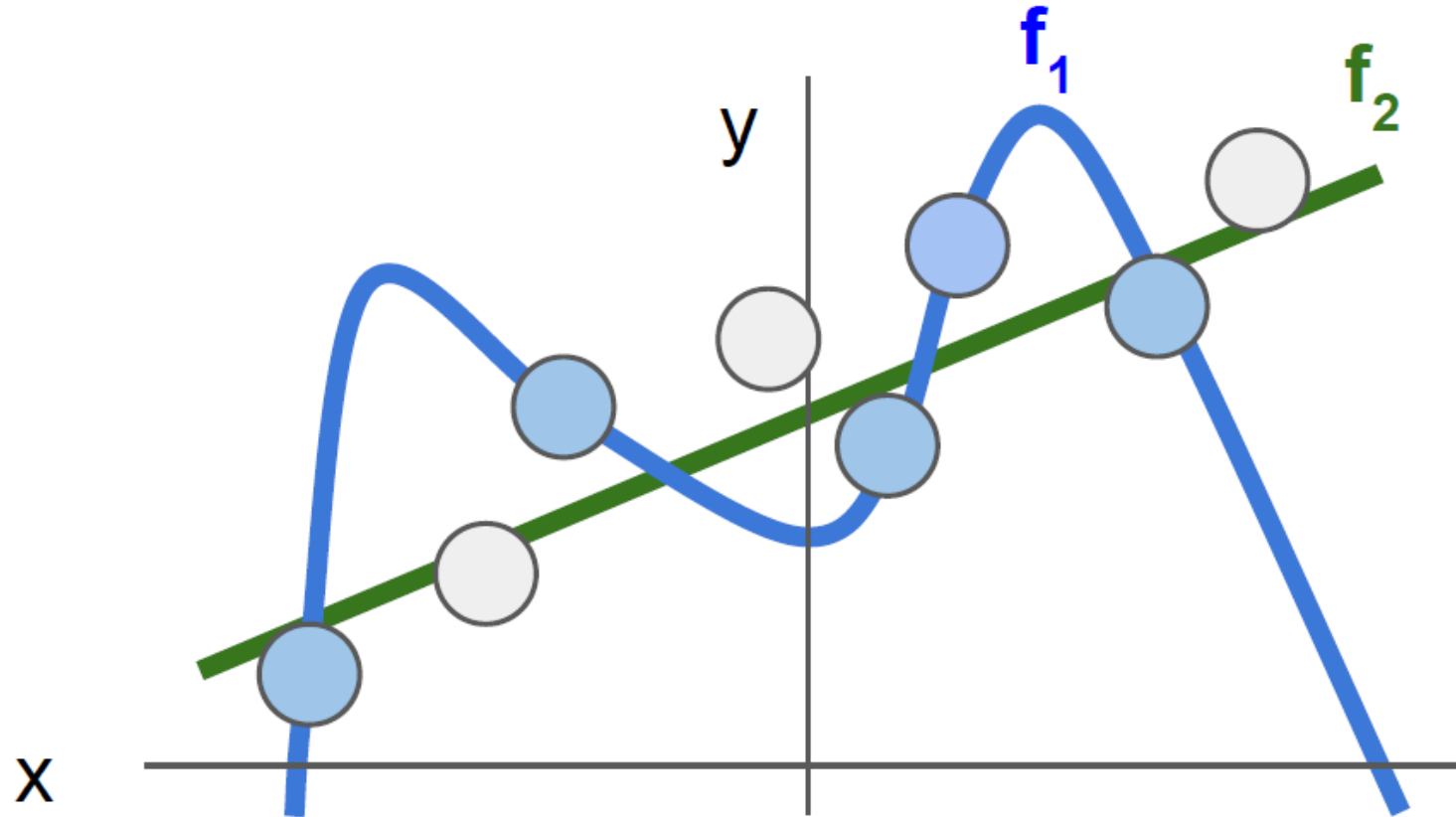
Regularization intuition: toy example training data



Regularization intuition: toy example training data



# Regularization intuition: toy example training data



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

# Regularization

$\lambda$  = regularization strength  
(hyperparameter)

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss}} + \underbrace{\lambda R(W)}_{\text{Regularization}}$$

**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):  $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

## More complex:

Dropout

Batch normalization

Stochastic depth, fractional pooling, etc

# Regularization

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

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

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

L2 Regularization

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

Which of  $w_1$  or  $w_2$  will  
the L2 regularizer prefer?



# Regularization

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

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

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

L2 Regularization

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

Which of  $w_1$  or  $w_2$  will the L2 regularizer prefer?

L2 regularization likes to “spread out” the weights

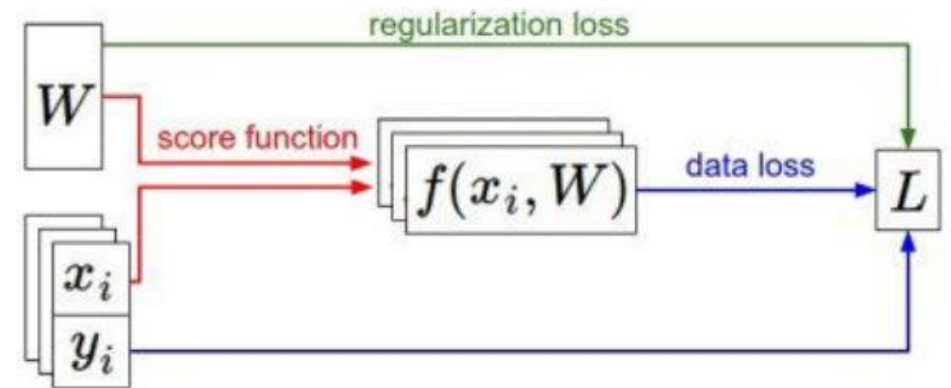
# Recap

- We have some dataset of (x,y)
- We have a **score function**:  $s = f(x; W) \stackrel{\text{e.g.}}{=} Wx$
- We have a **loss function**:

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

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

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



The softmax function is applied to the logits to:

- a) Normalize the scores into a probability distribution
- b) Increase the magnitude of the scores
- c) Decrease the magnitude of the scores
- d) Convert the scores into binary values