

Lecture 18-1

Linear classifiers and neural networks

1950s Age of the Perceptron

1957 The Perceptron (Rosenblatt)

1969 Perceptrons (Minsky, Papert)

1980s Age of the Neural Network

1986 Back propagation (Hinton)

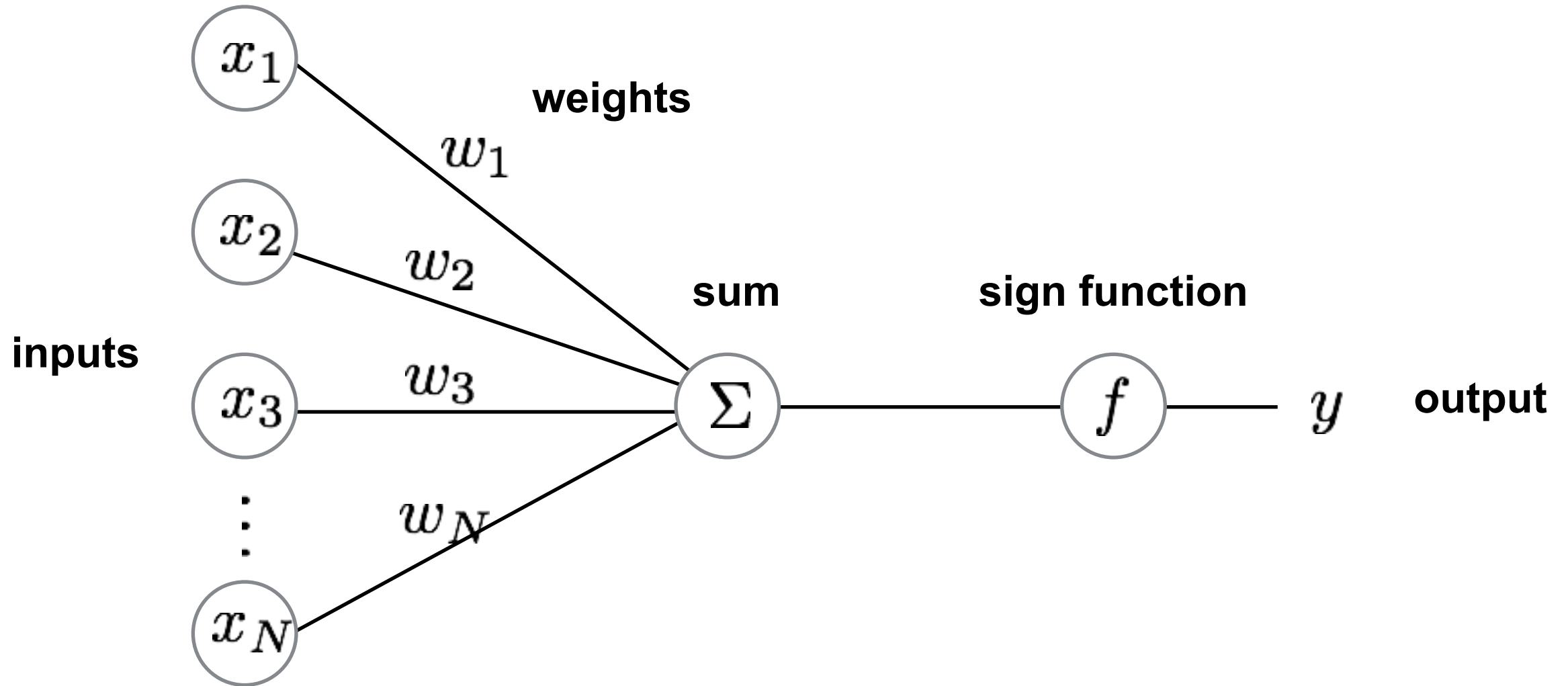
1990s Age of the Graphical Model

2000s Age of the Support Vector Machine

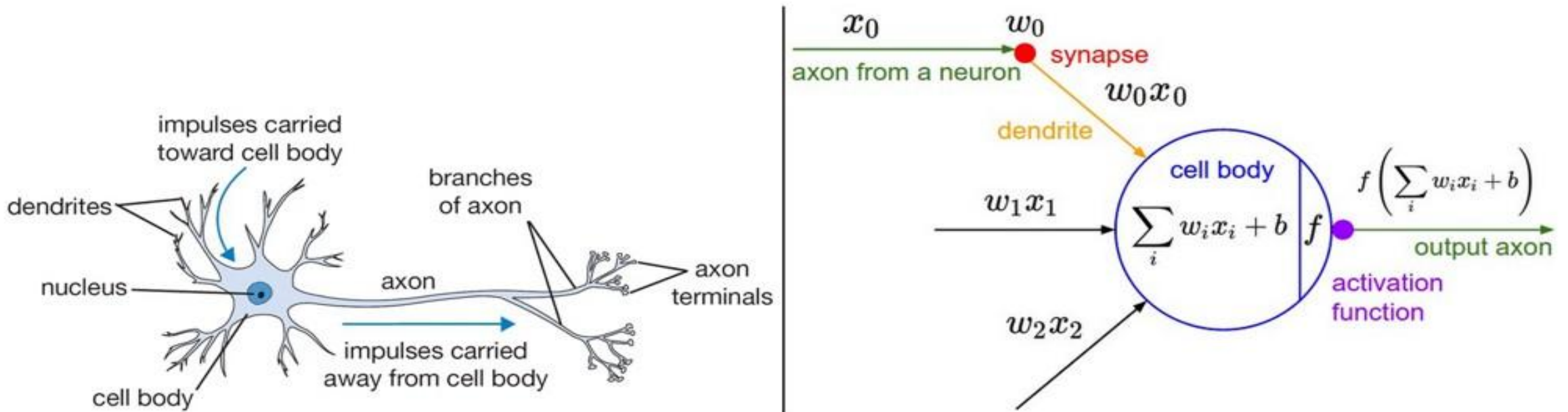
2010s Age of the Deep Network

deep learning = known algorithms + computing power + big data

Perceptron



Aside: Inspiration from Biology

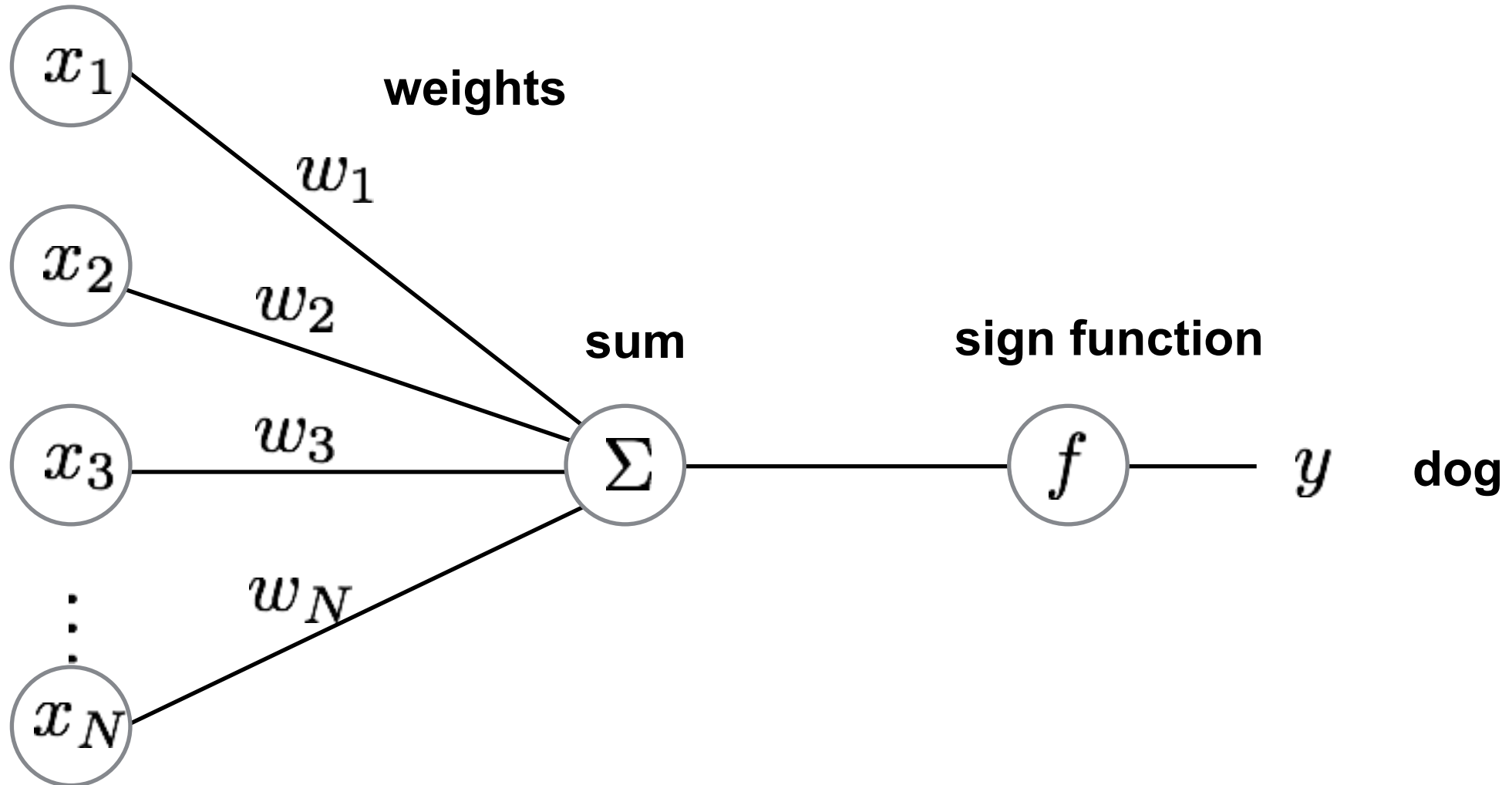
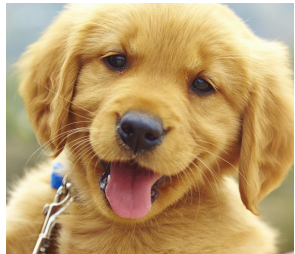


A cartoon drawing of a biological neuron (left) and its mathematical model (right).

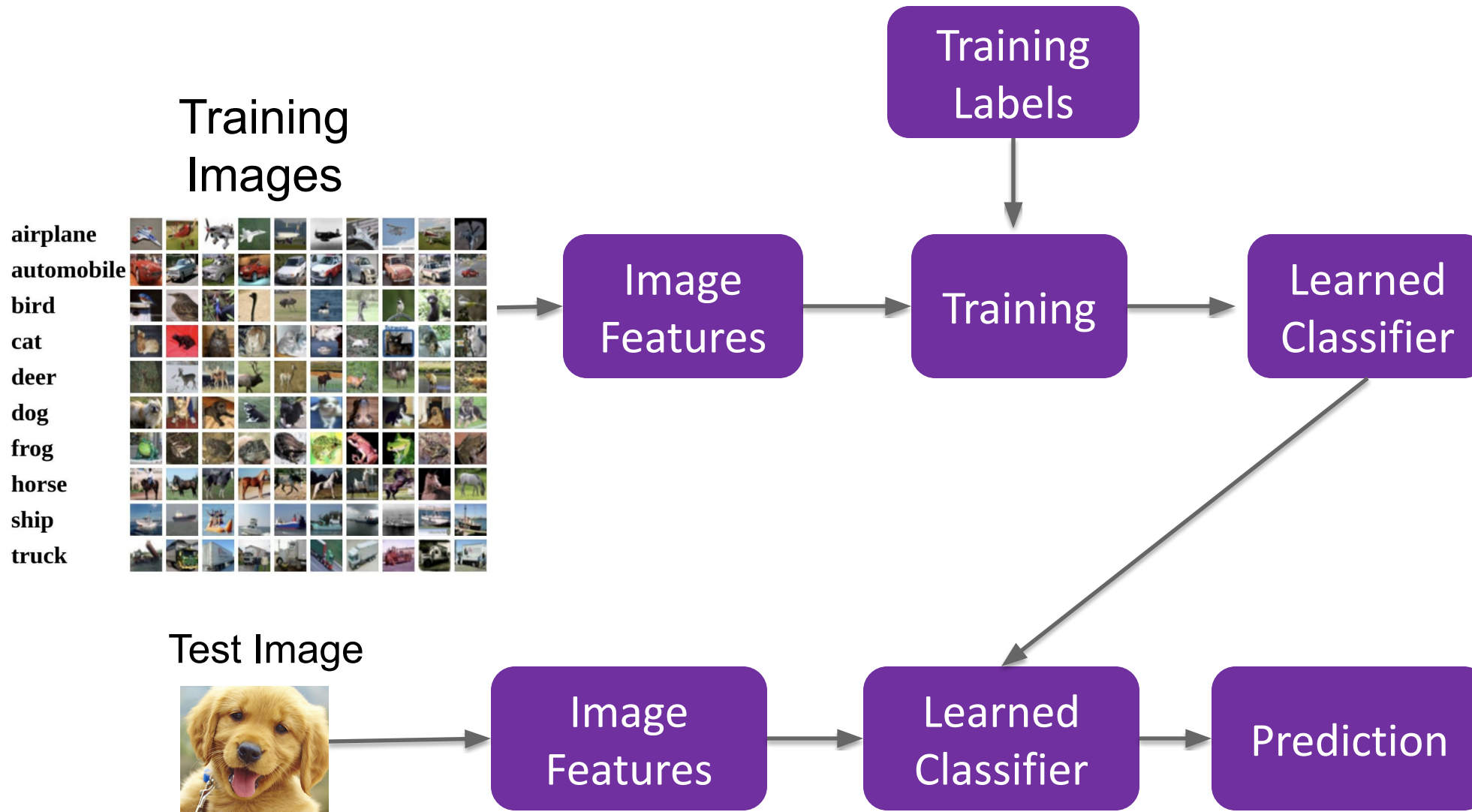
Neural nets/perceptrons are loosely inspired by biology.

But they are NOT how the brain works, or even how neurons work.

Perceptron: for image classification



Let's revisit our simple recognition pipeline to explain where perceptrons fit in



Recall: we can featurize images into a vector

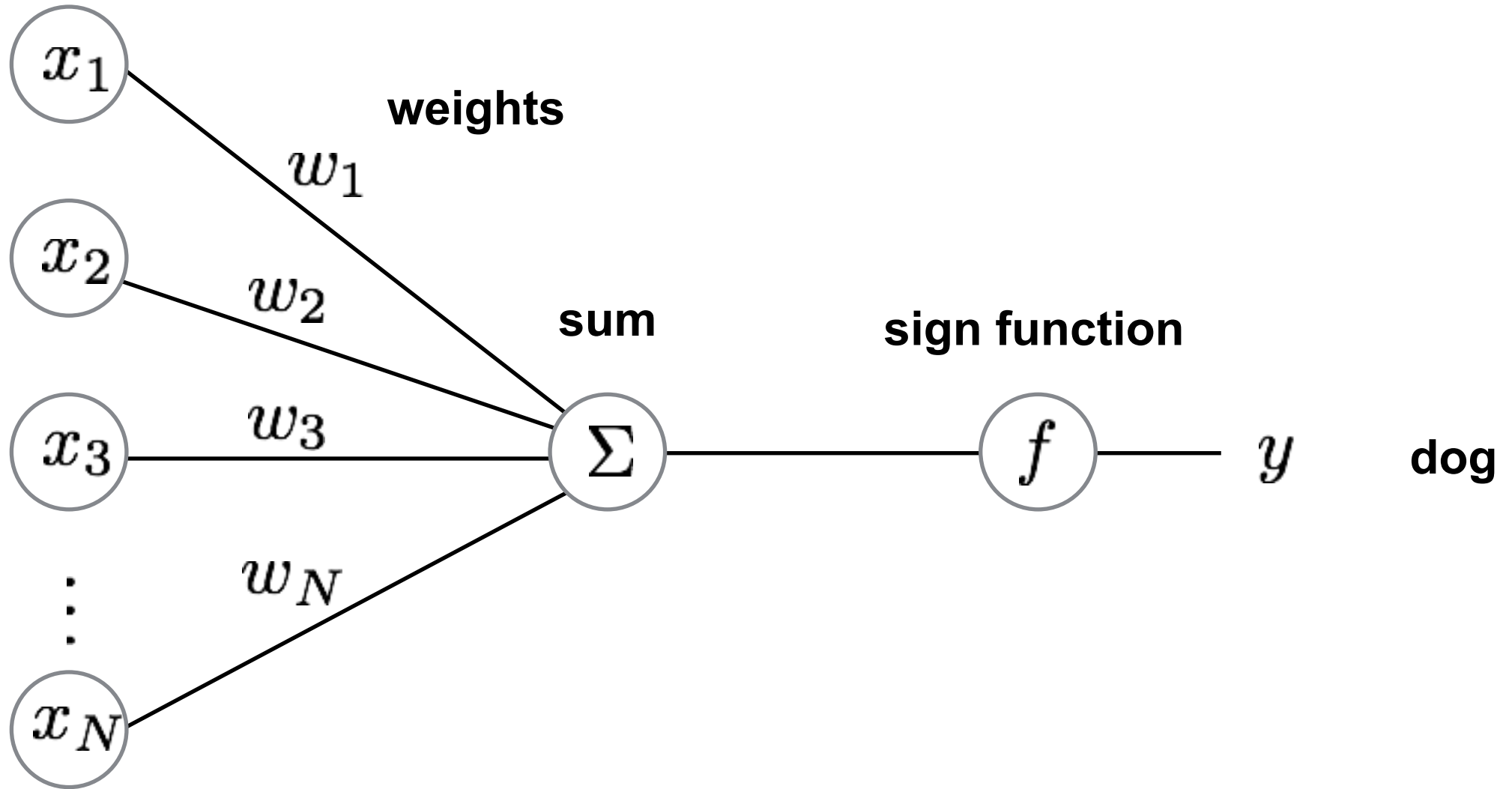


Raw pixels
Raw pixels + (x,y)
PCA
BoW
BoW + spatial pyramids

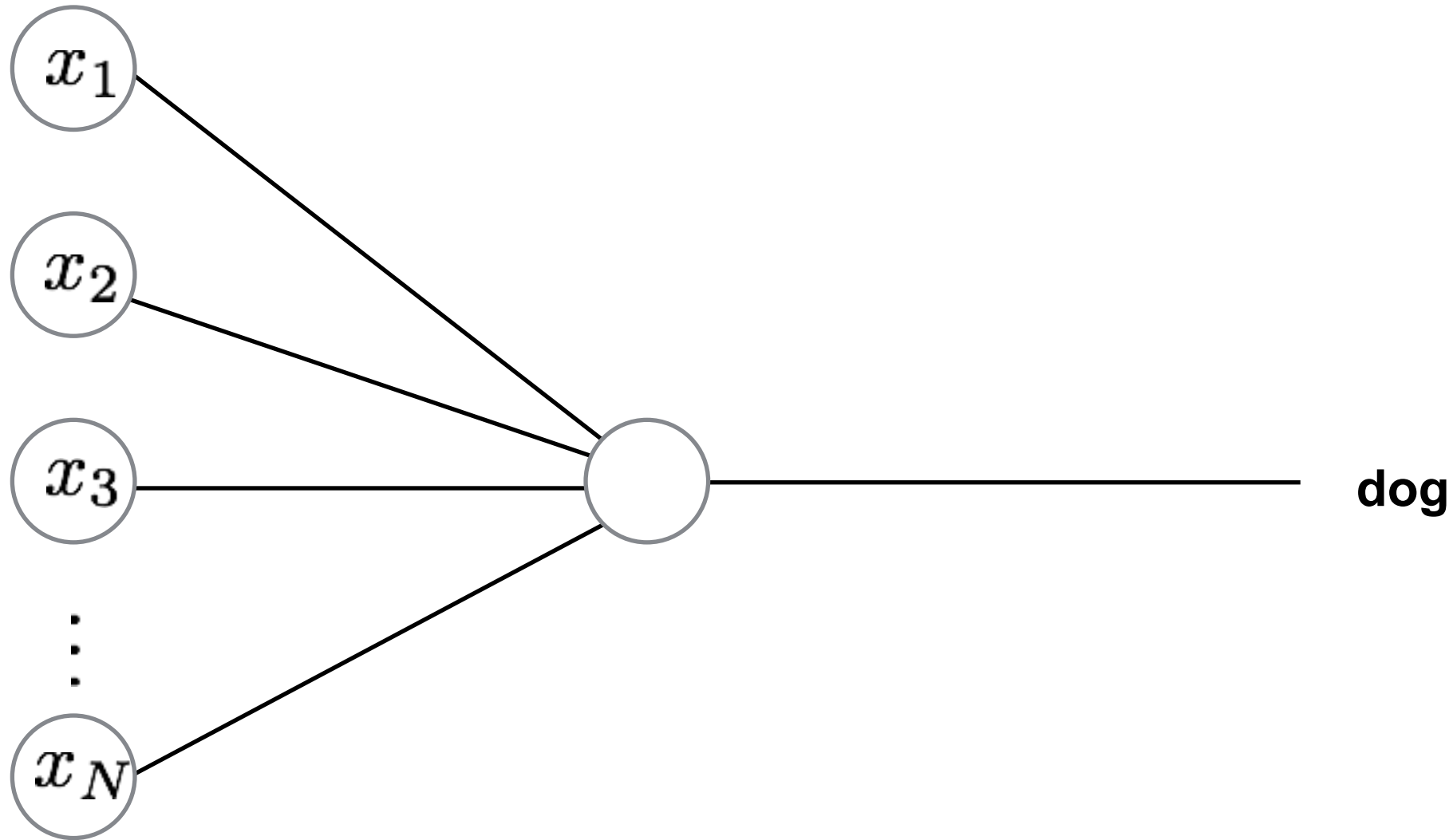
Image
Vector

x_1
 x_2
 x_3
...
...
...
...
...
...
...
 x_n

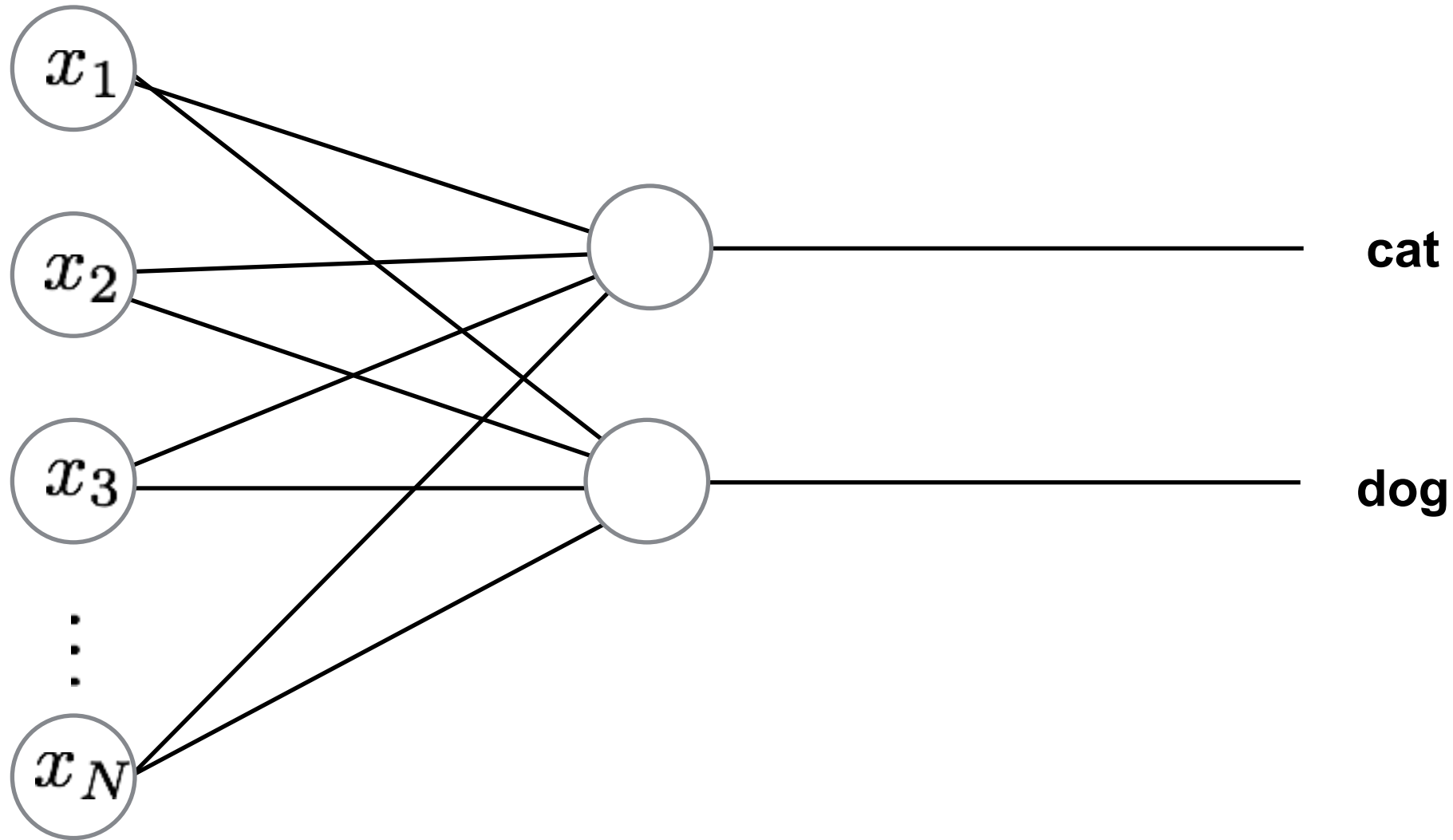
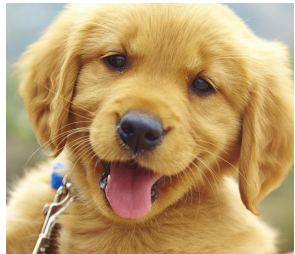
Perceptrons are a simple transformation that converts feature vectors into recognition scores

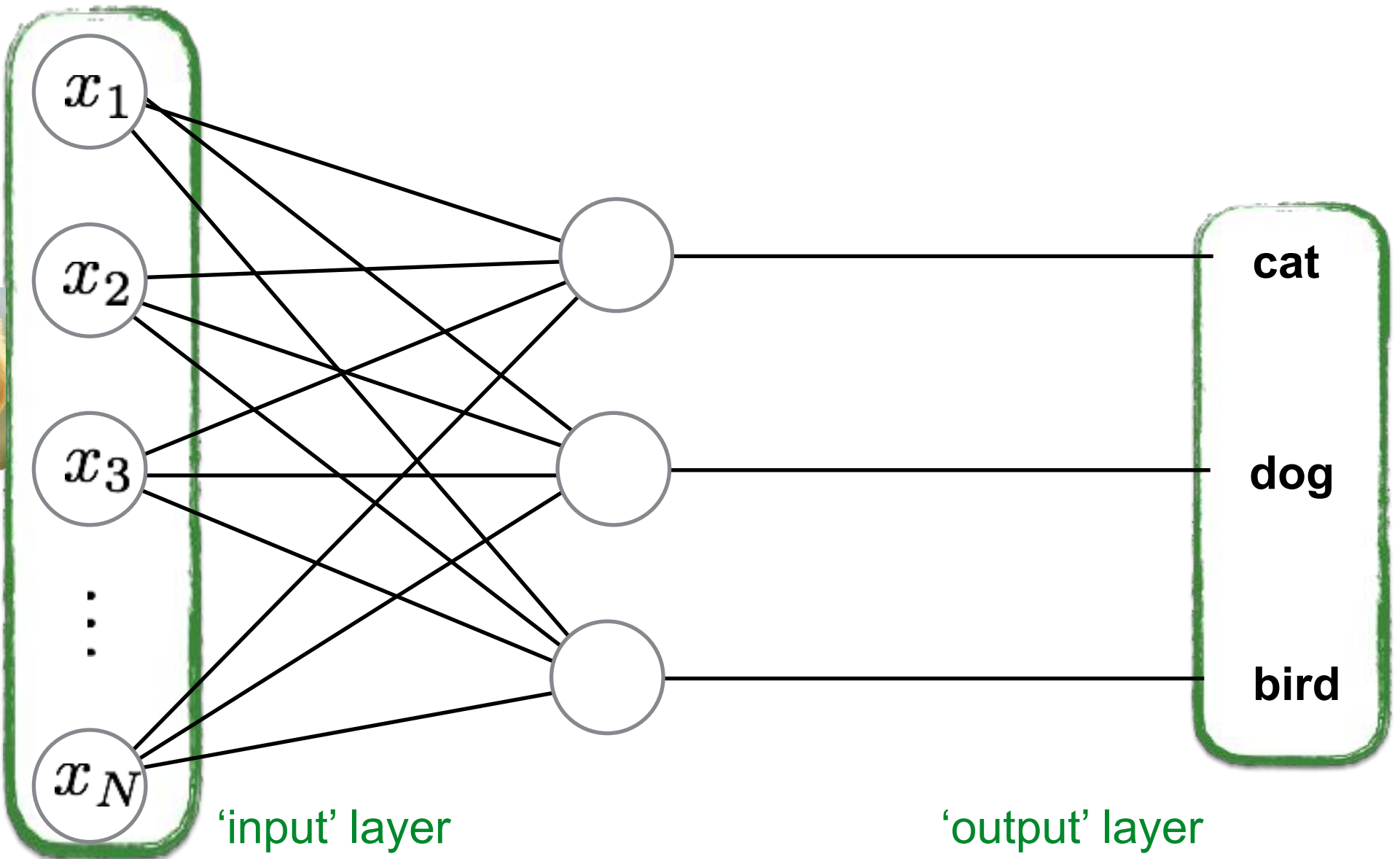


Perceptron: simplified view with **one** perceptron
(produces 1 score for one category)

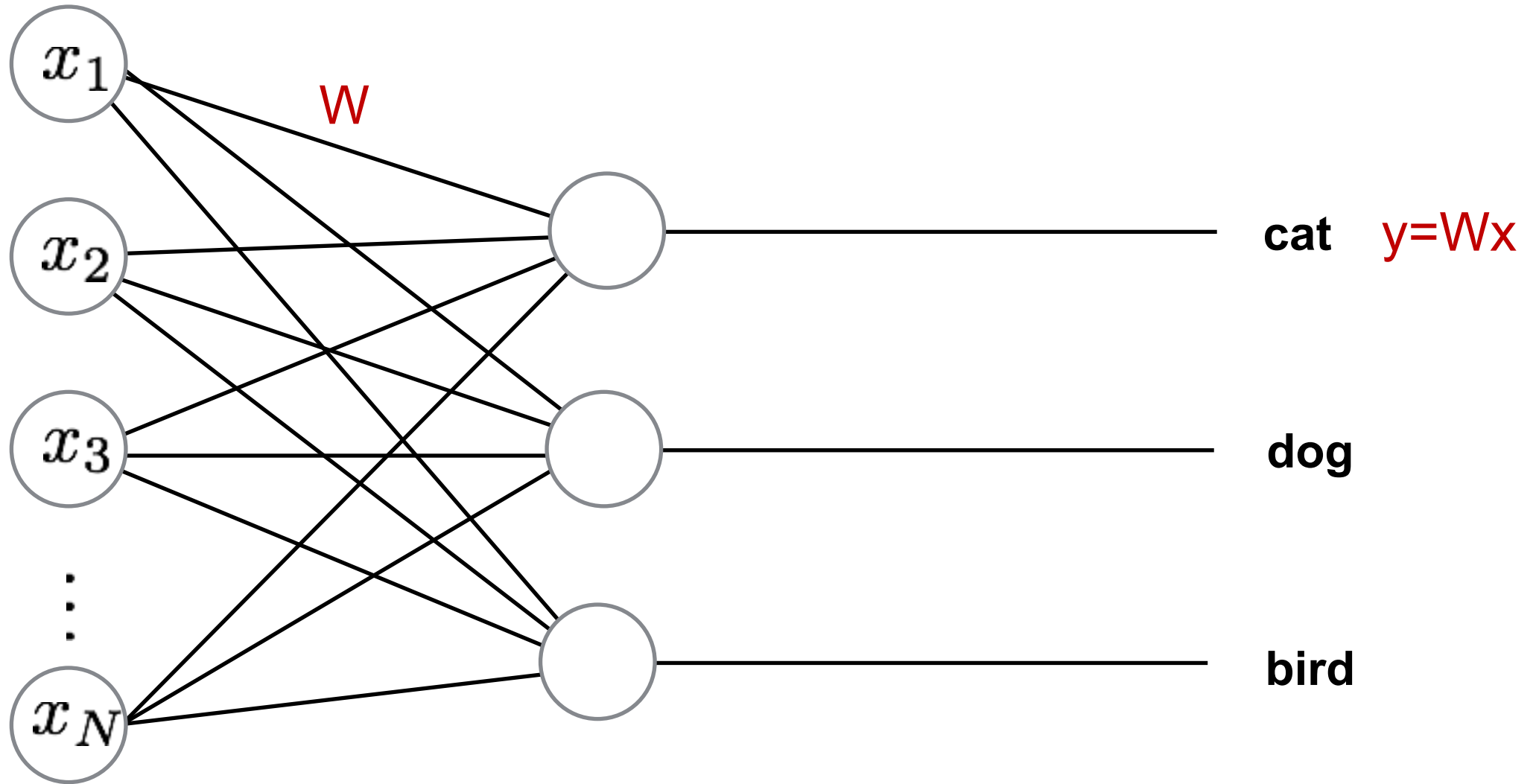


Perceptron: simplified view with **two** perceptrons
(produces 2 scores with 2 categories)

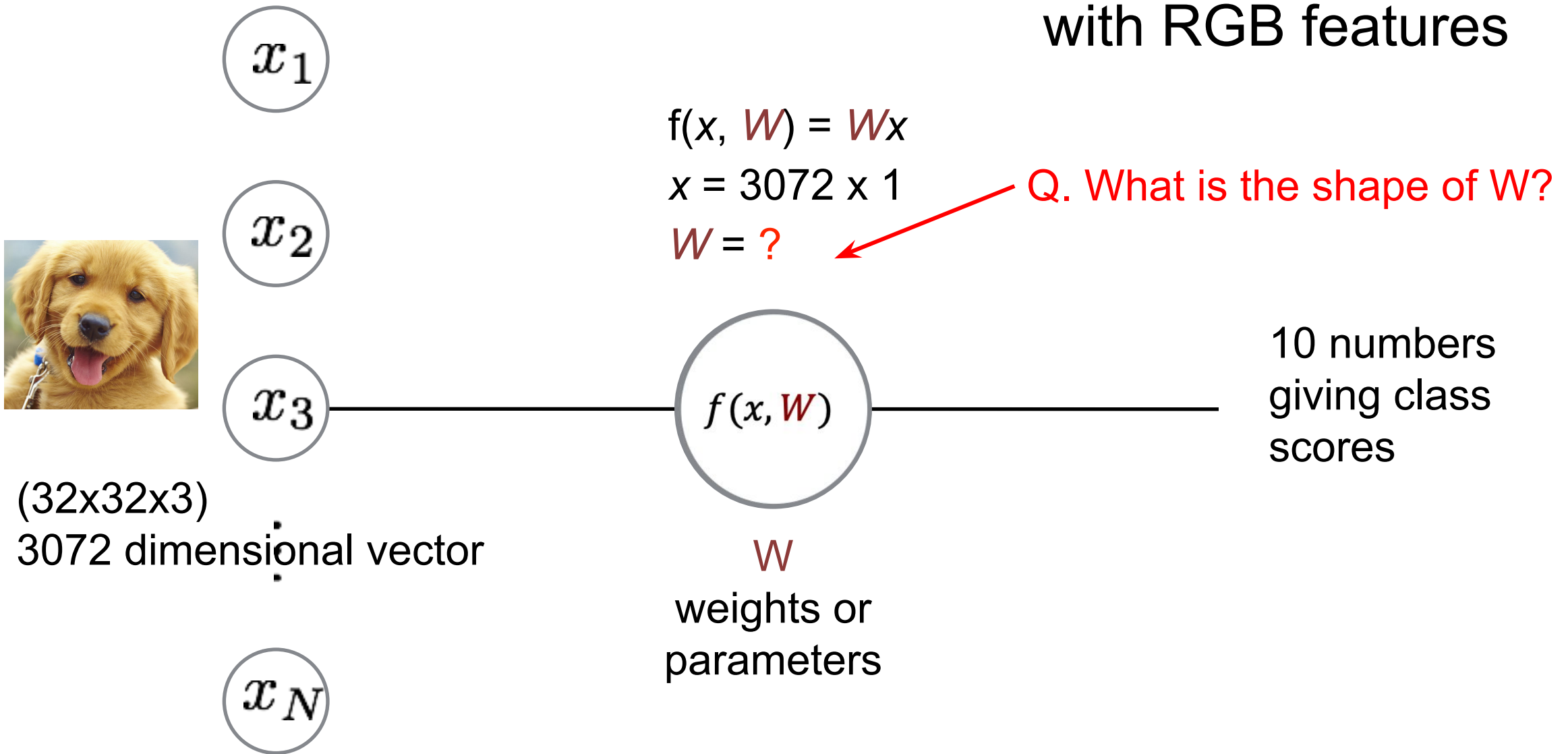




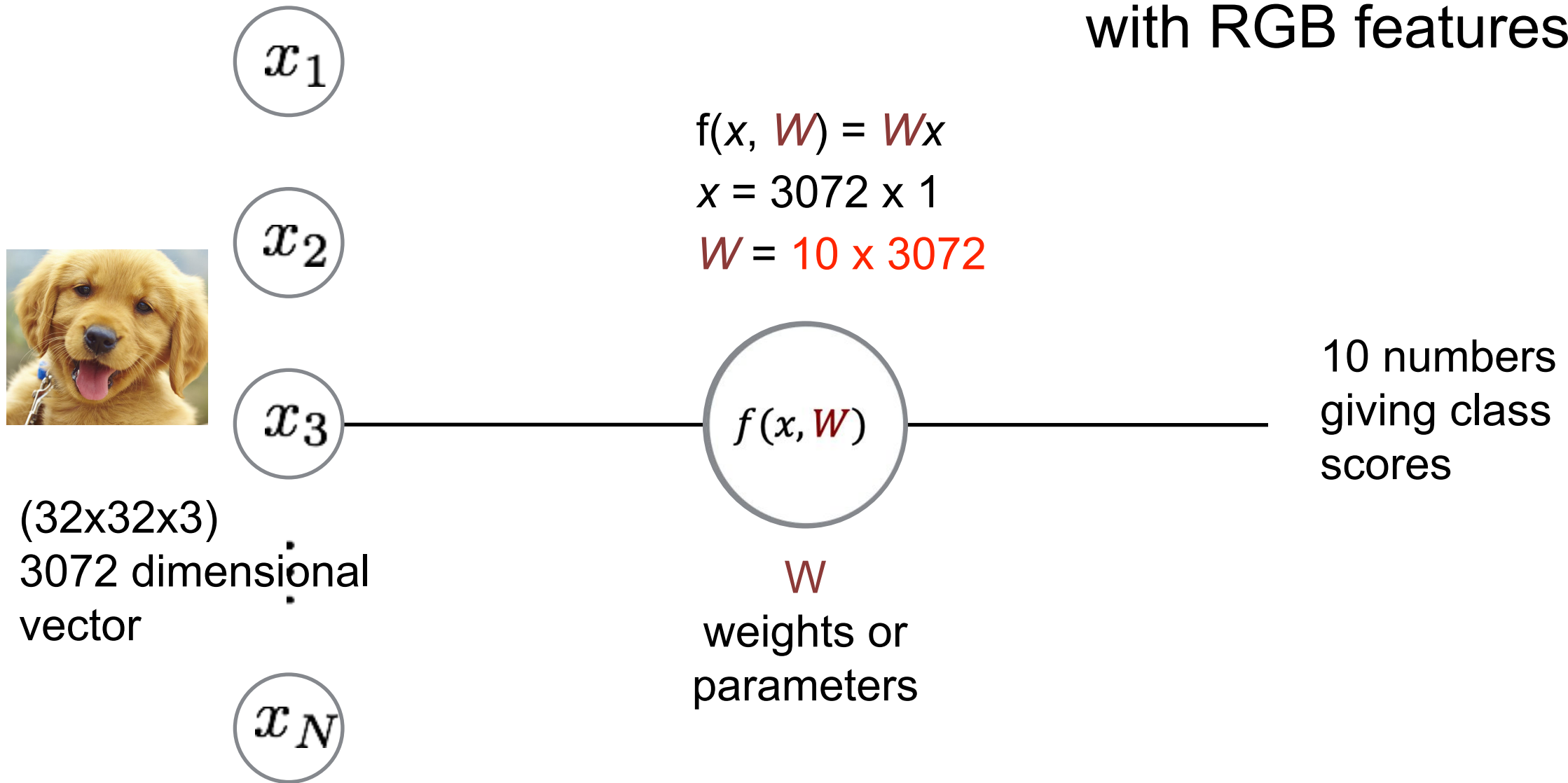
Linear classifier is a set of perceptrons produces one score for every category



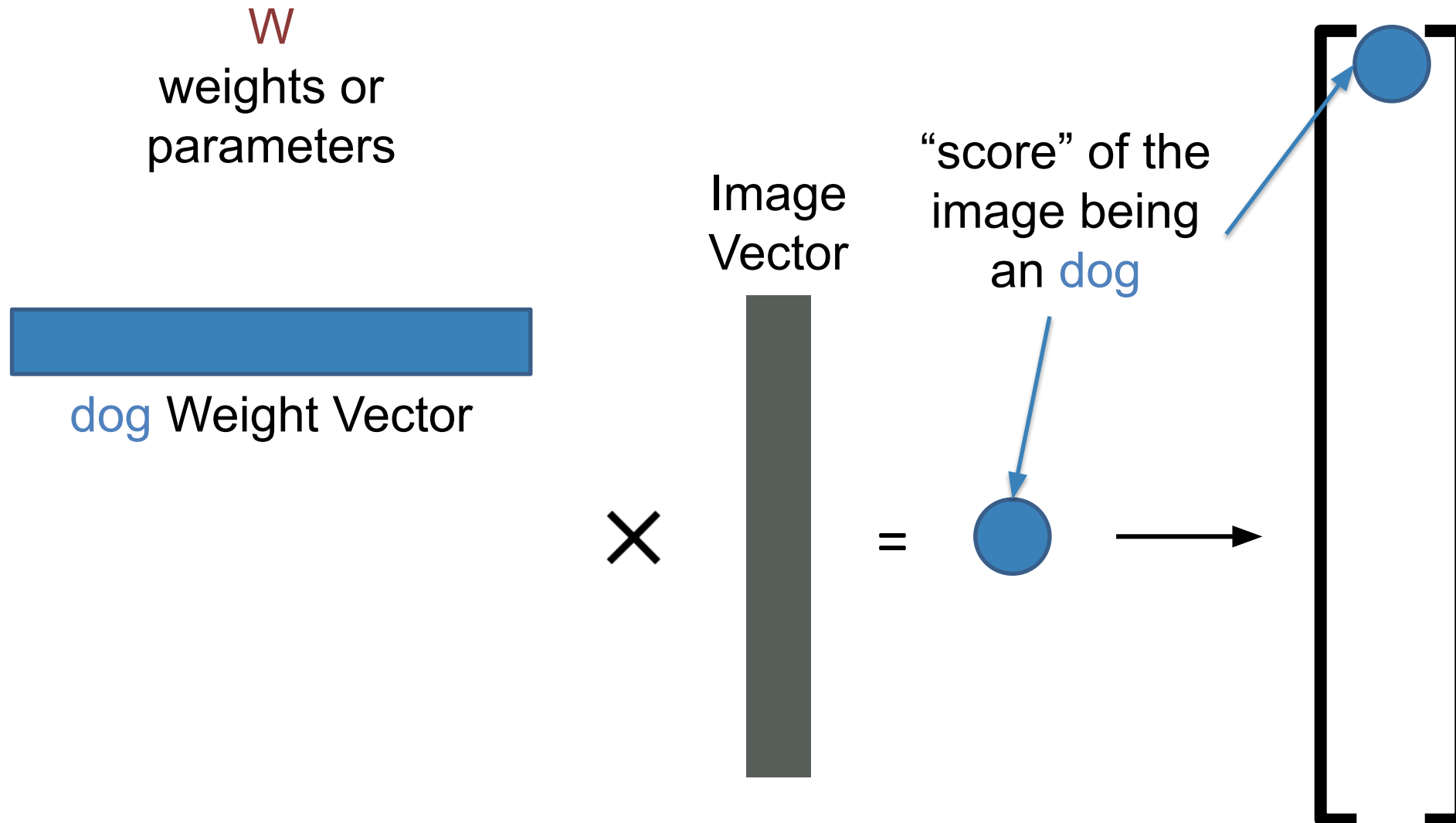
Linear classifier: mathematical formulation with RGB features



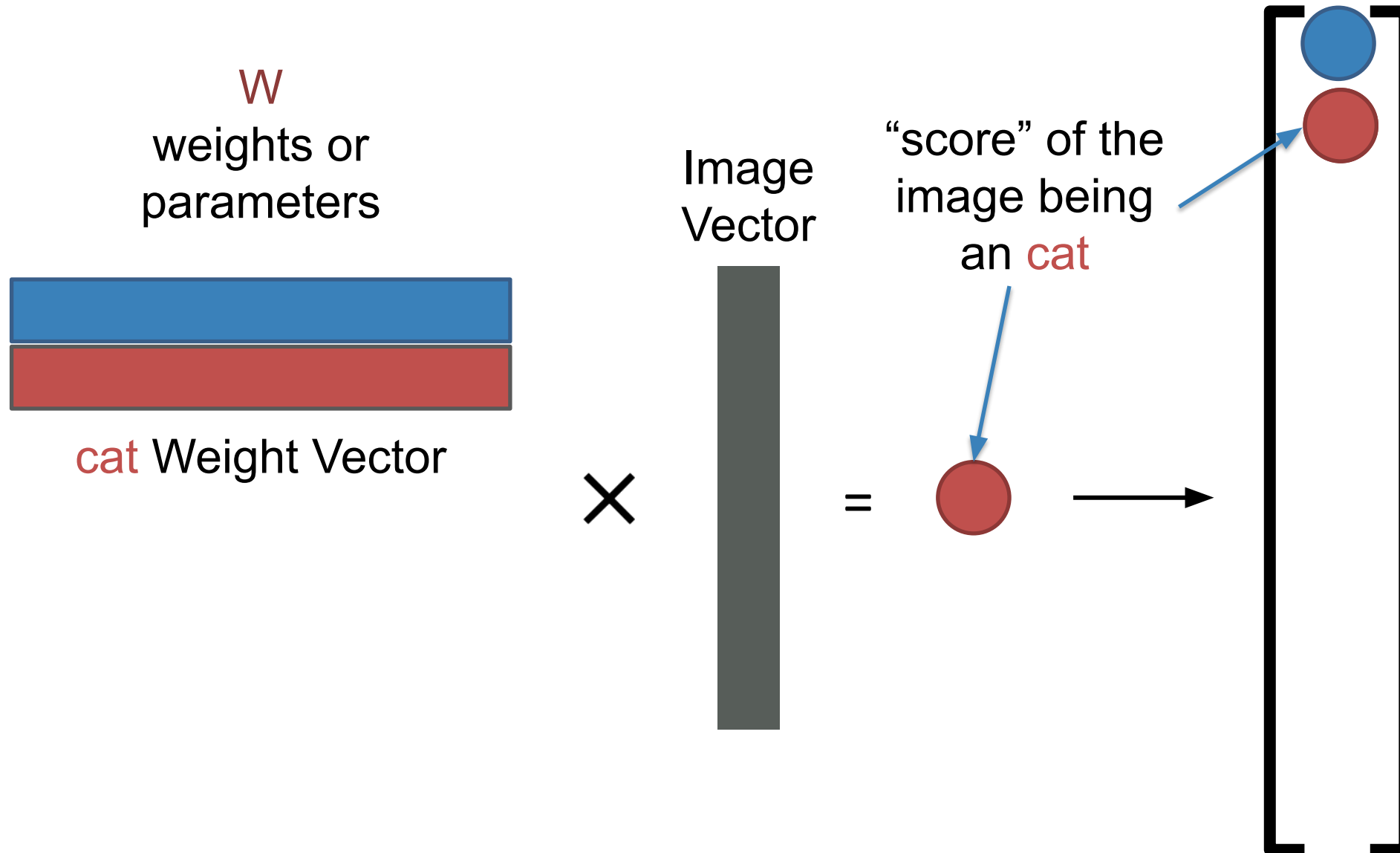
Linear classifier: mathematical formulation with RGB features



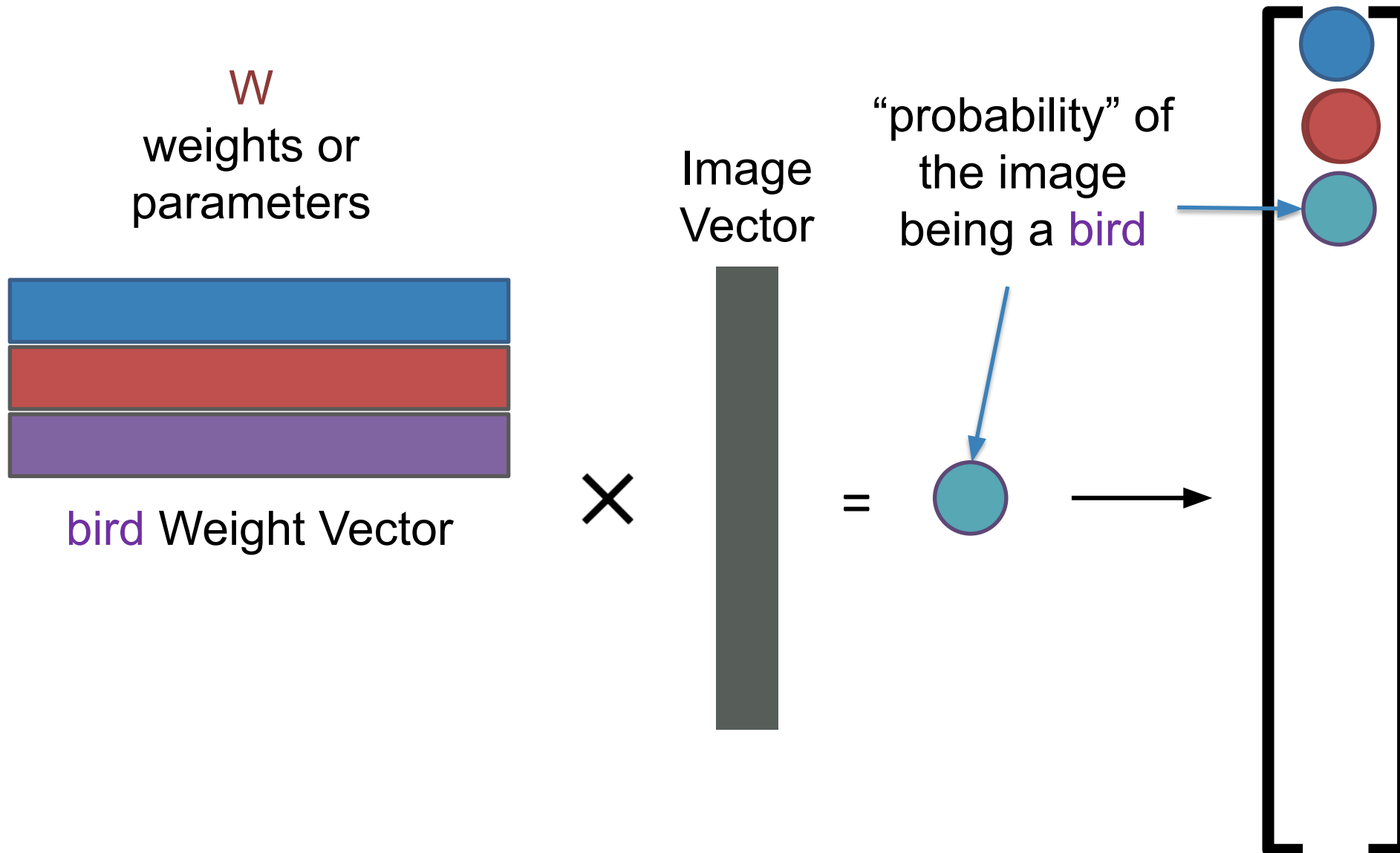
Linear classifier: function visualized



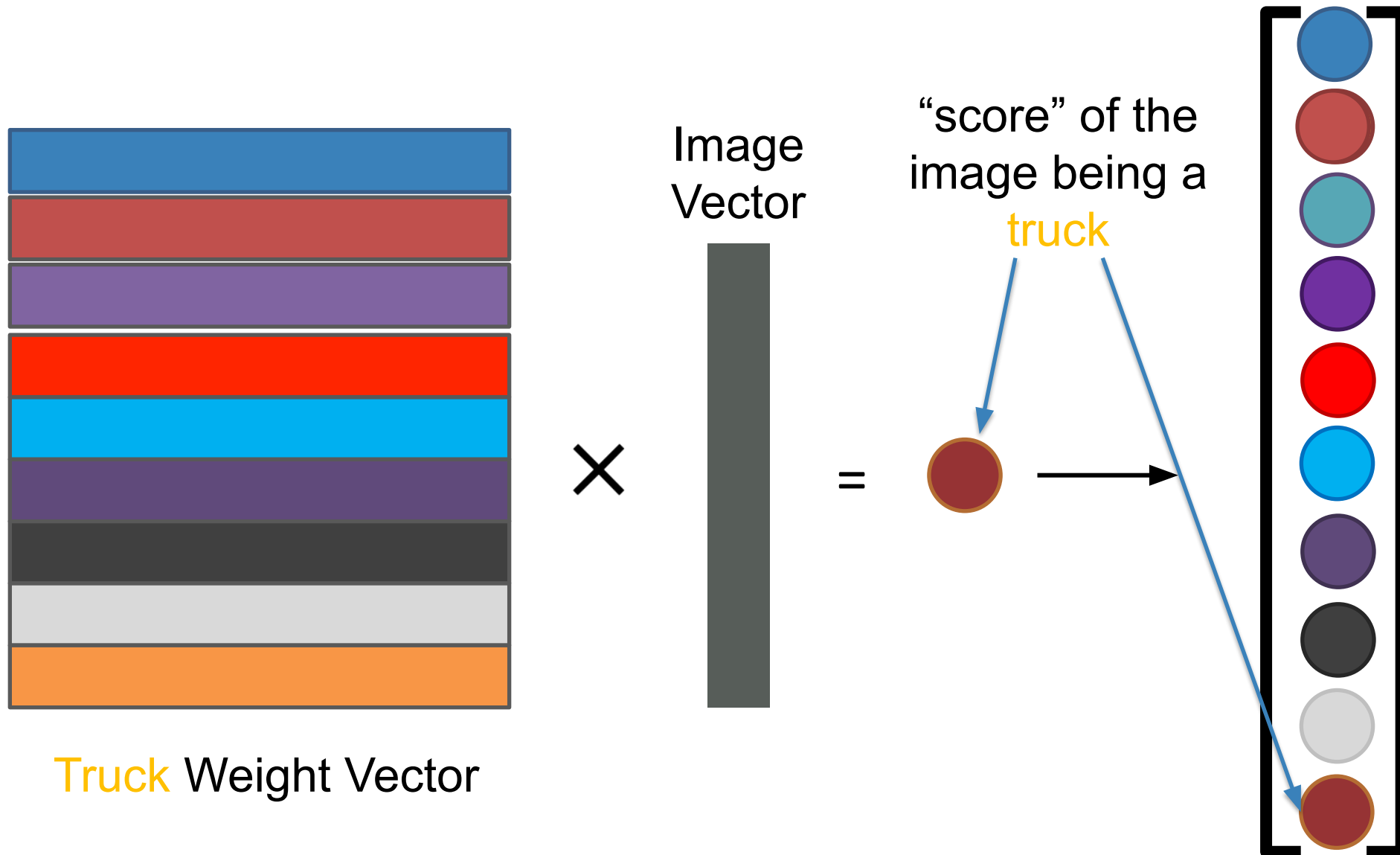
Linear classifier: function visualized



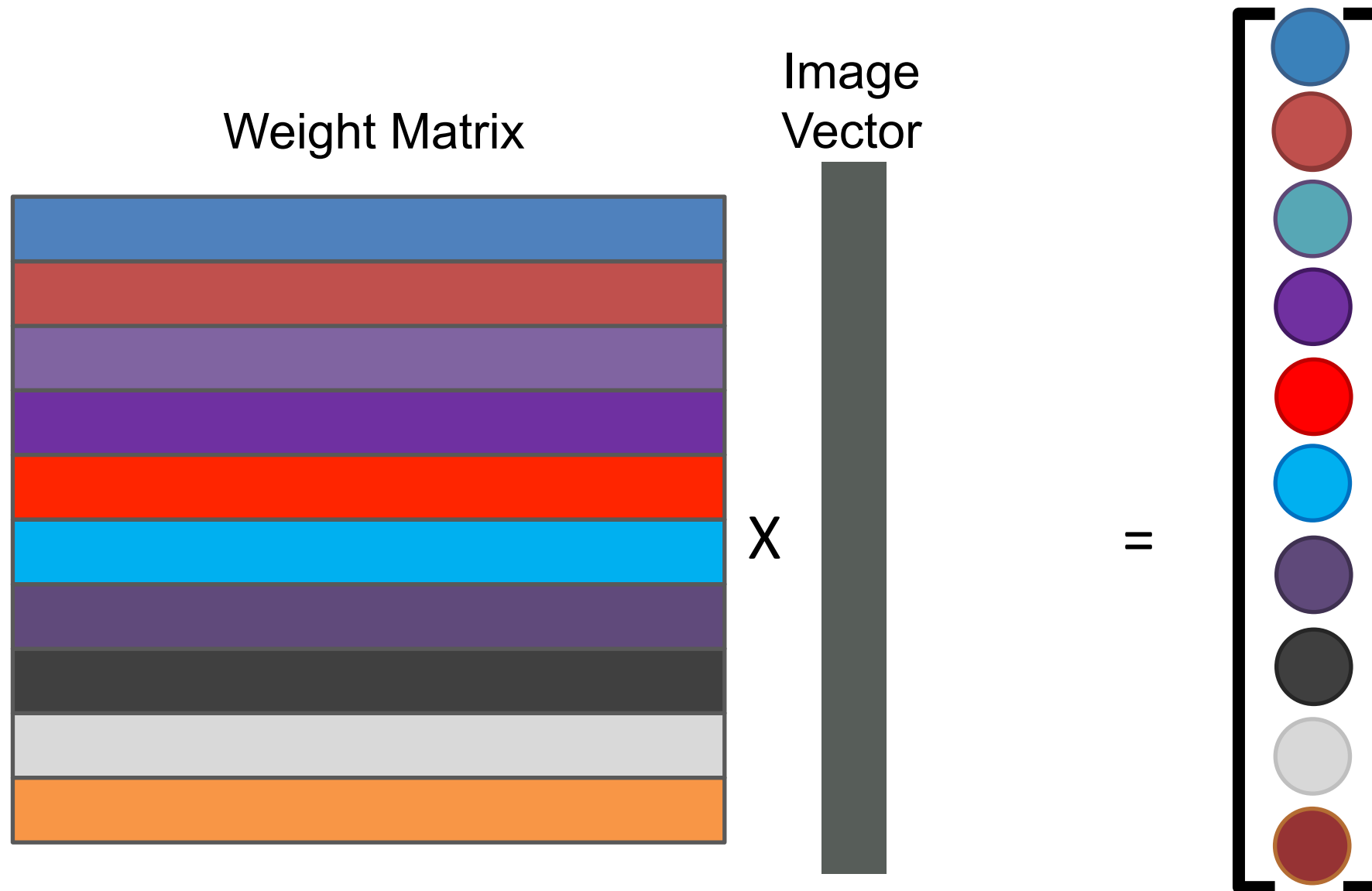
Linear classifier: function visualized



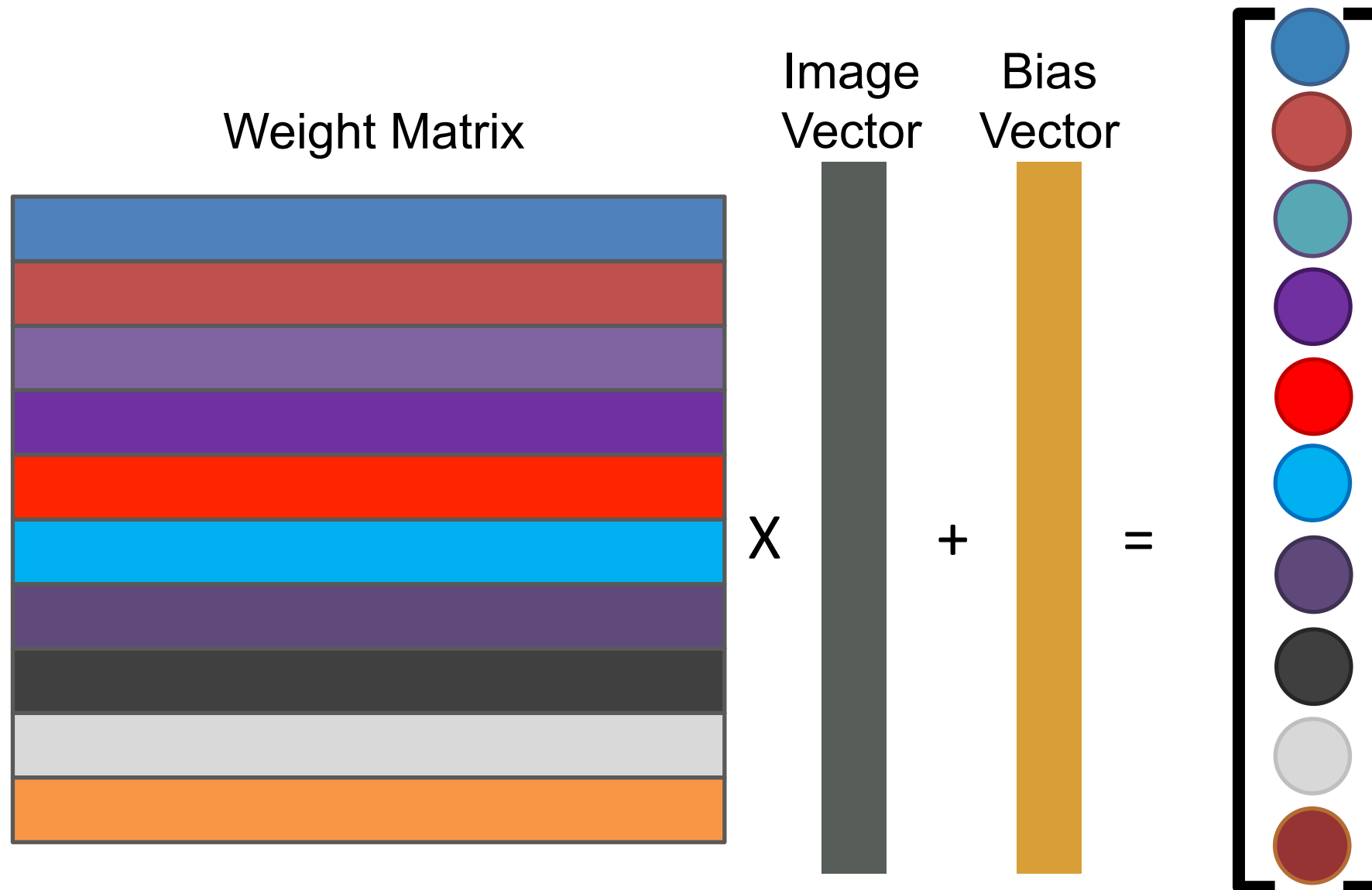
Linear classifier: function visualized



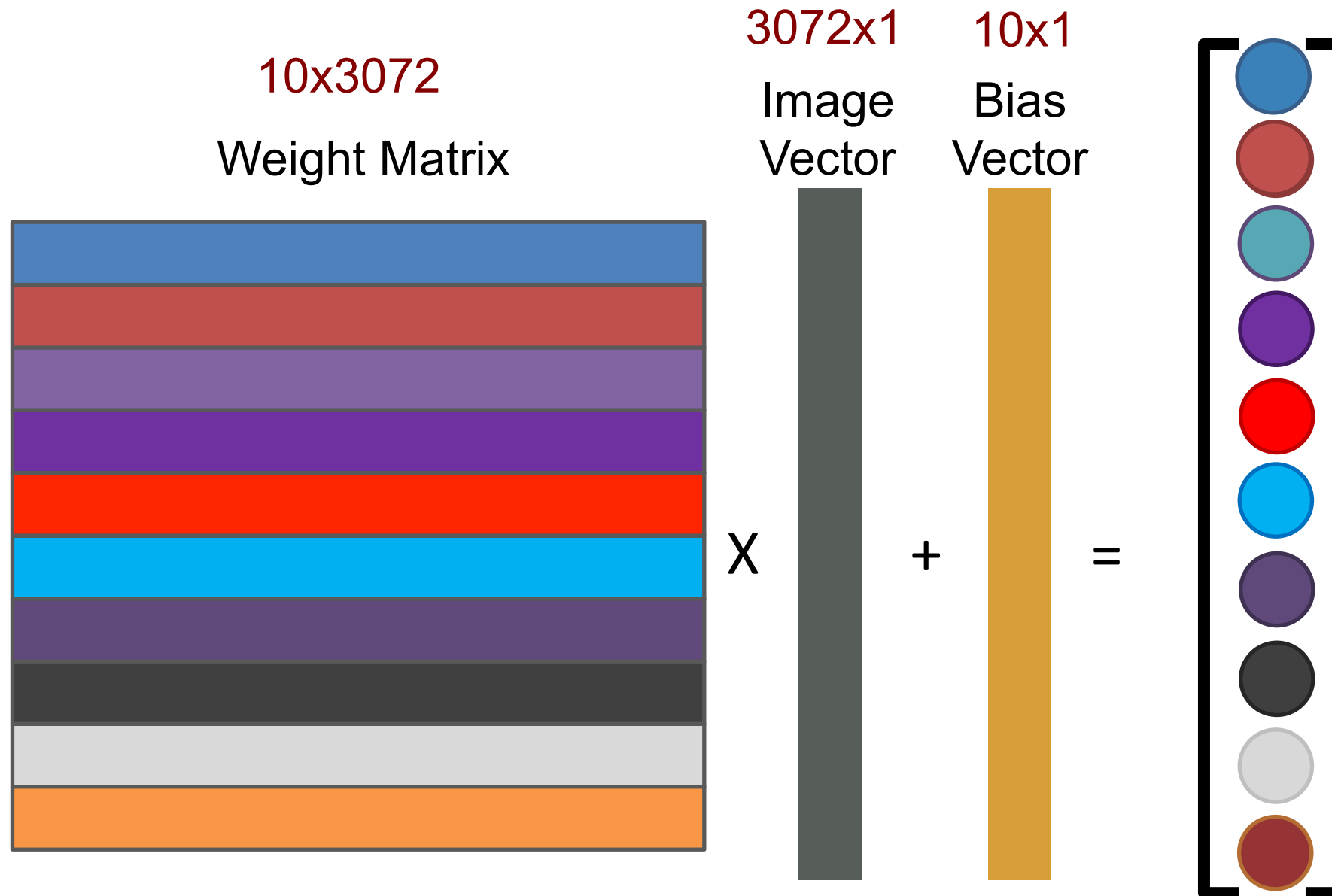
Linear classifier: **function visualized**



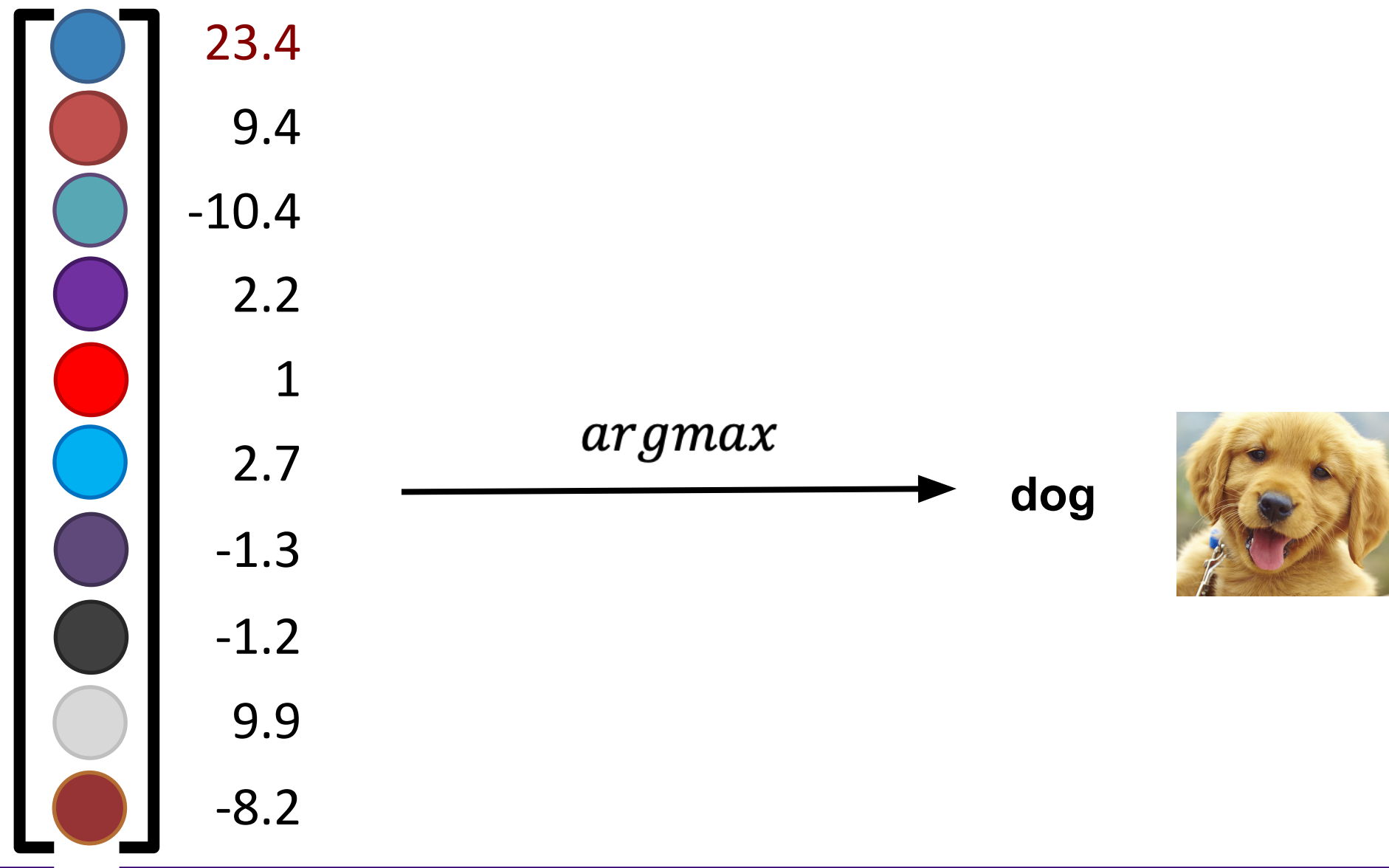
Linear classifier: **bias vector**



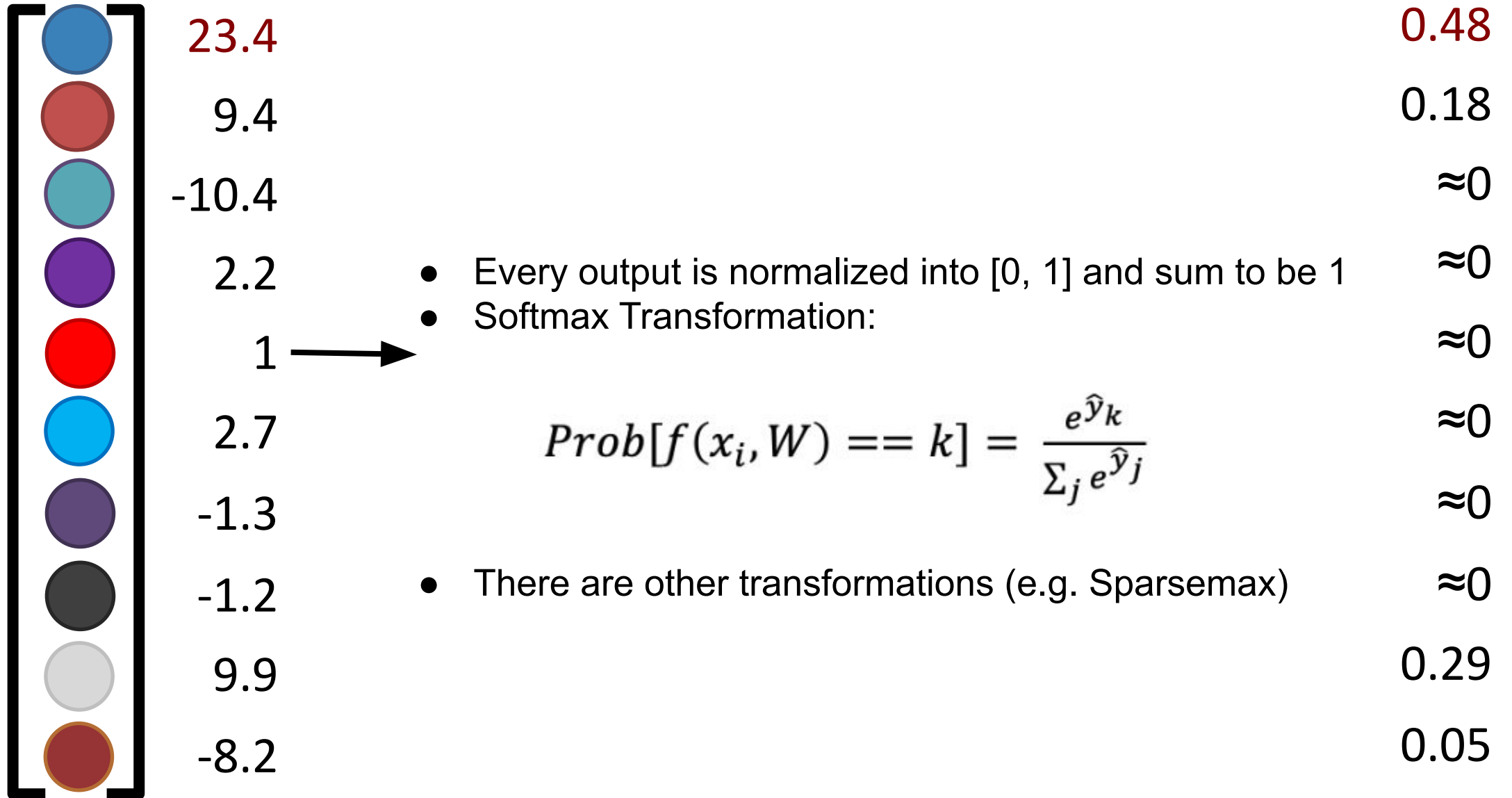
Linear classifier: size



Linear classifier: Making a classification



Interpret the Output as Probability



Interpreting the weights

- Assume our weights are trained on the CIFAR 10 dataset with **raw pixels**:

airplane



automobile



bird



cat



deer



dog



frog



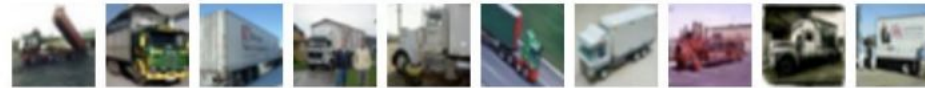
horse



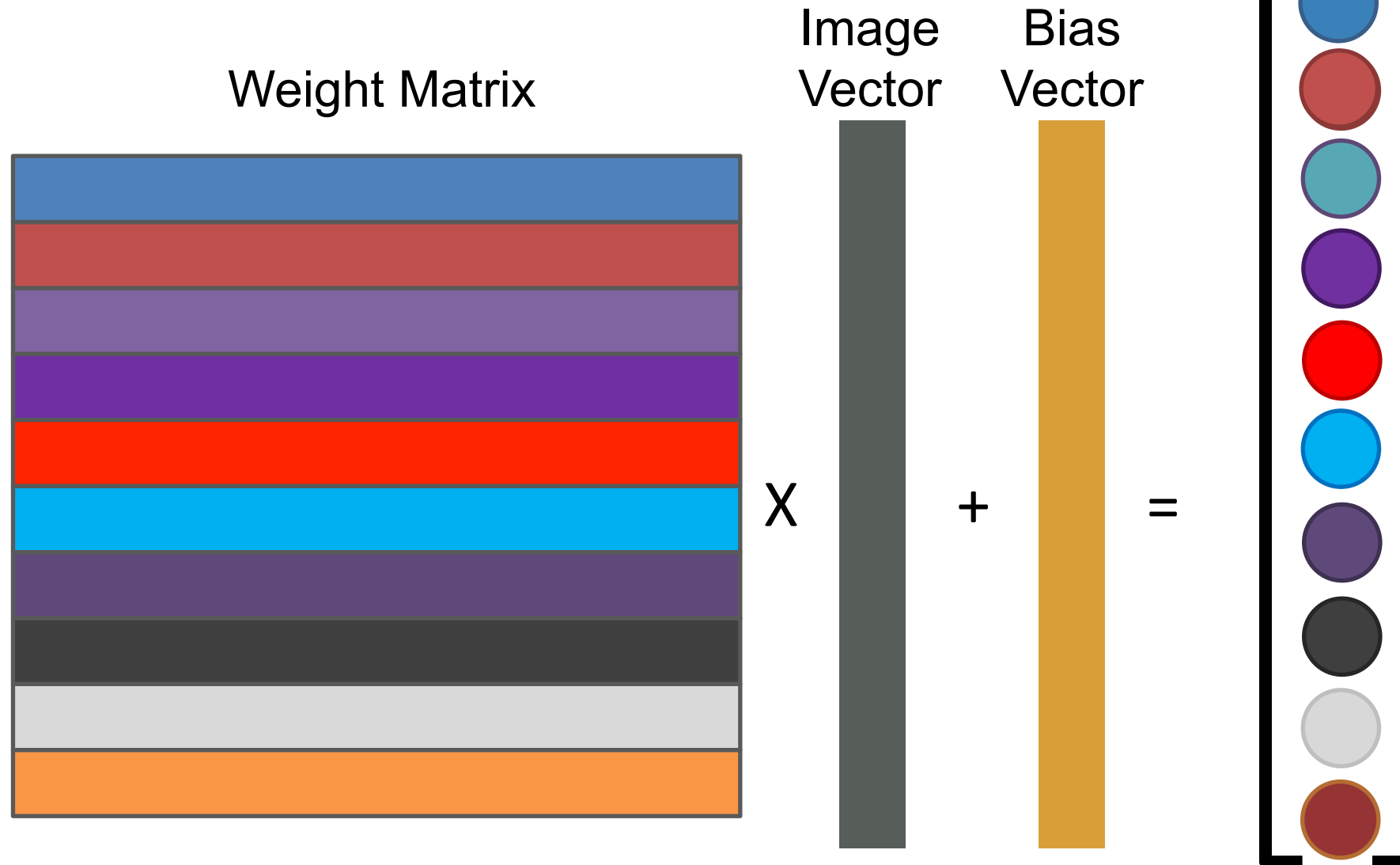
ship



truck

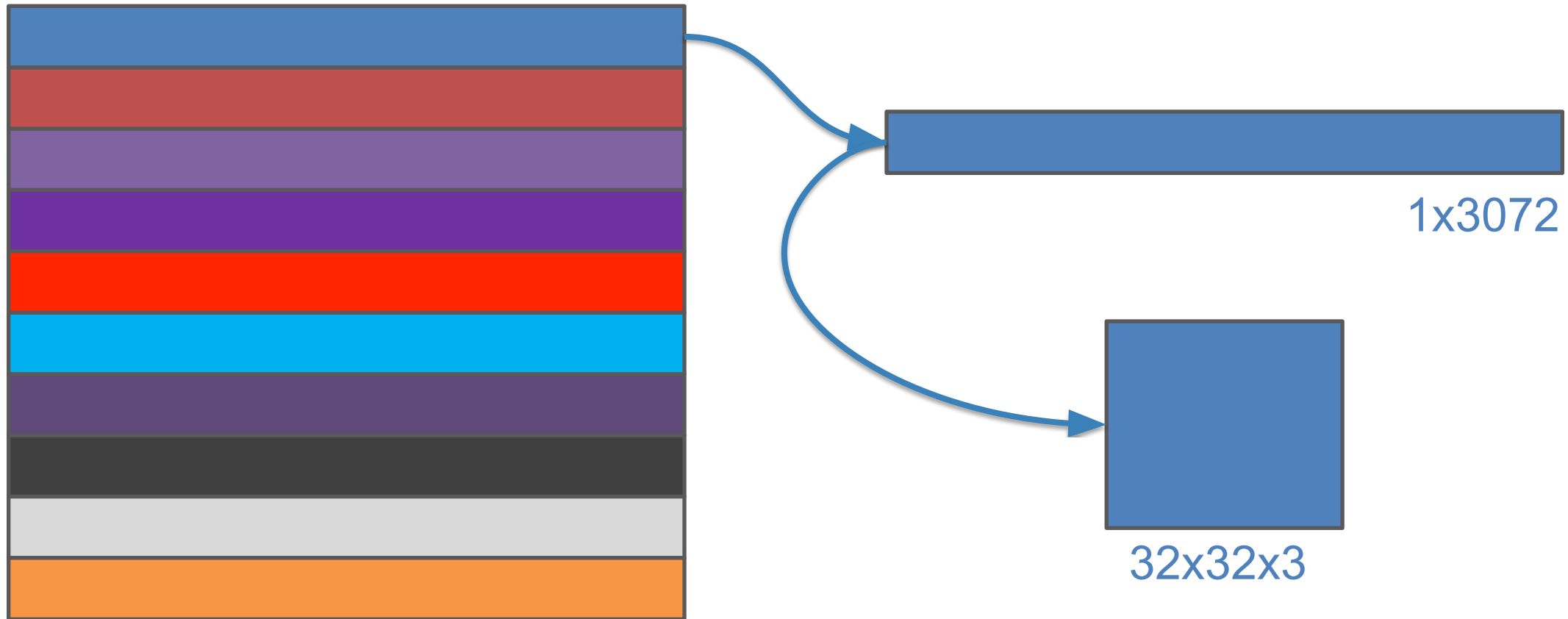


Interpreting the weights **as templates**



Interpreting the weights **as templates**

We can reshape the vector back in to the shape of an image

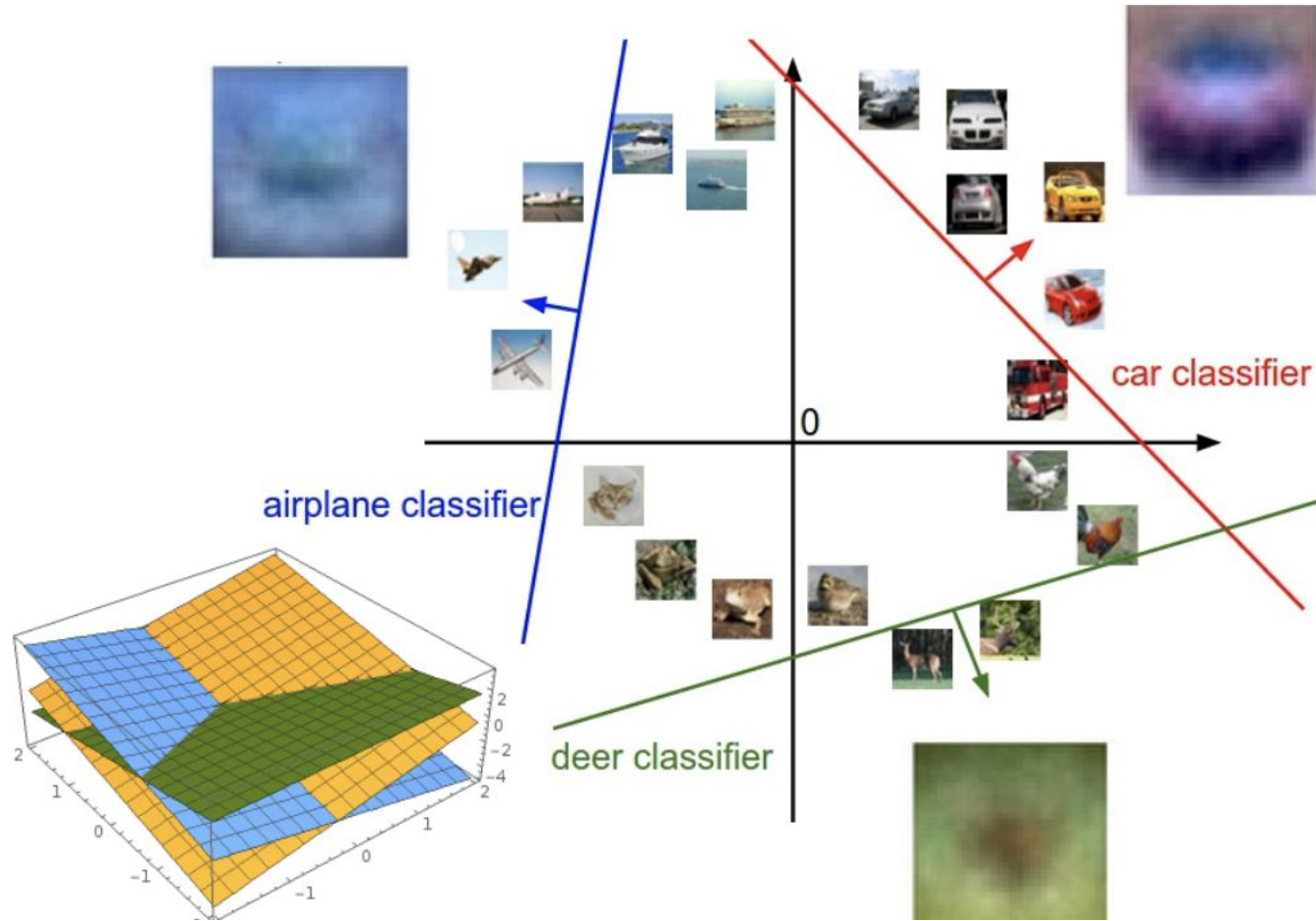


Let's visualize what the **templates** look like

We can reshape the row back to the shape of an image



Interpreting the weights **geometrically**



- Assume the image vectors are in 2D space to make it easier to visualize.

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

Today's agenda

- Perceptron
- Linear classifier
- **Loss function**
- Gradient descent and backpropagation
- Neural networks

Training linear classifiers

We need to learn how to **pick the weights** in the first place.

Formally, we need to find **W** such that

$$\min_{\mathbf{W}} \text{Loss}(\mathbf{y}, \hat{\mathbf{y}})$$

Where \mathbf{y} is the true label, $\hat{\mathbf{y}}$ is the model's predicted label.

All we have to do is **define a loss function!**

- When the classifier predicts **correctly**, the loss should be low
- When the classifier makes **mistakes**, the loss should be high

Properties of a loss function

Given several training examples: $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

and a perceptron: $\hat{y} = wx$

where x is image and y is (integer) label (0 for dog, 1 for cat, etc)

Loss over the entire dataset is an average of loss over examples

$$L = \frac{1}{N} \sum_{i=1}^N L_i(y_i, \hat{y}_i)$$

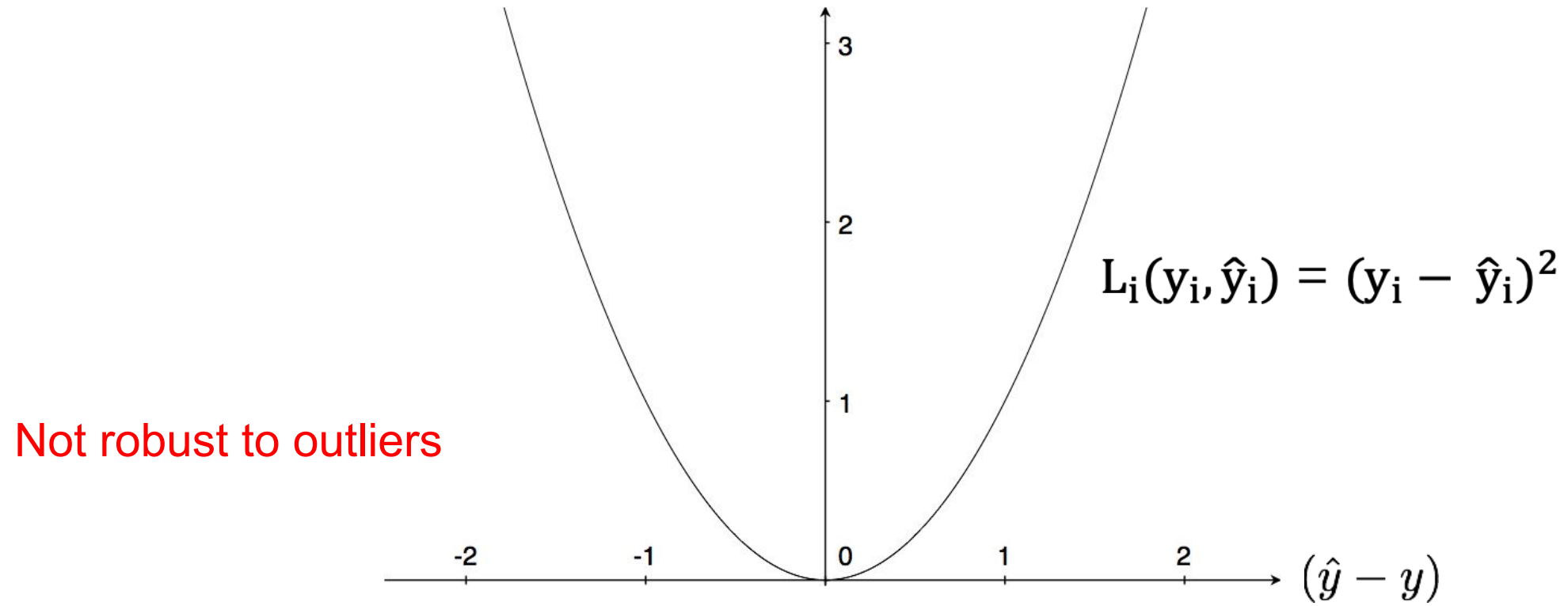
How do we choose the loss function L_i ?

YOU get to chose the loss function!

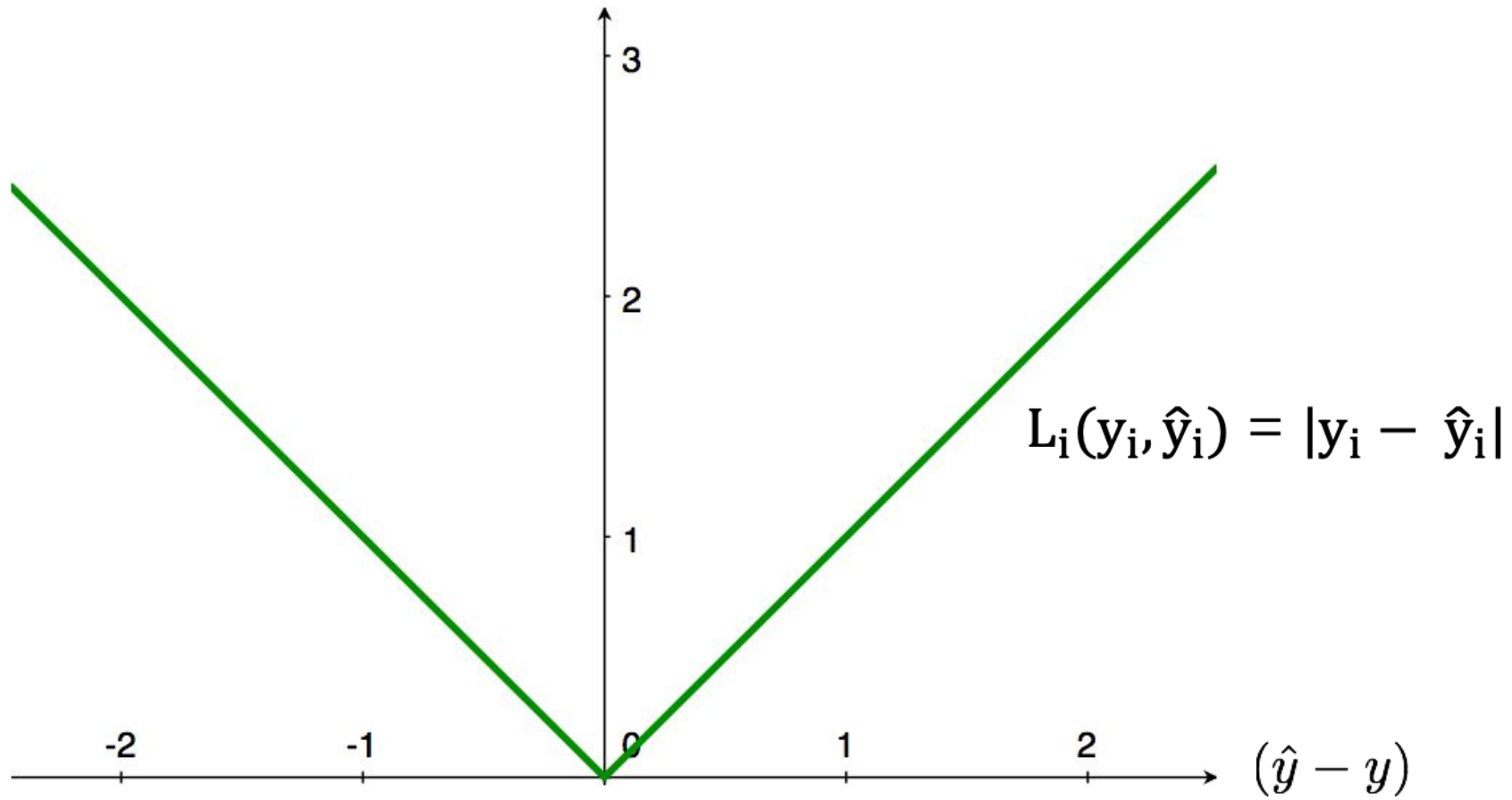
(some are better than others depending on what you want to do)

Squared Error (L2)

(a popular loss function) ((why?))

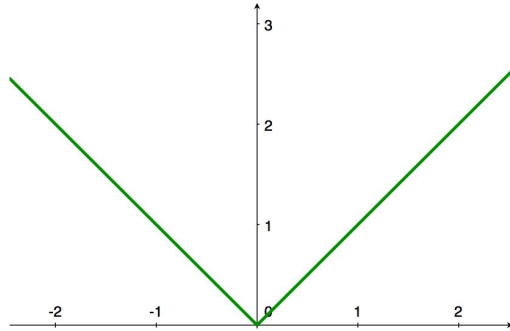


L1 loss



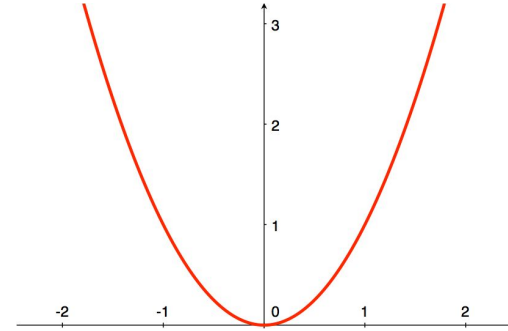
L1 Loss

$$L_i(y_i, \hat{y}_i) = |y_i - \hat{y}_i|$$



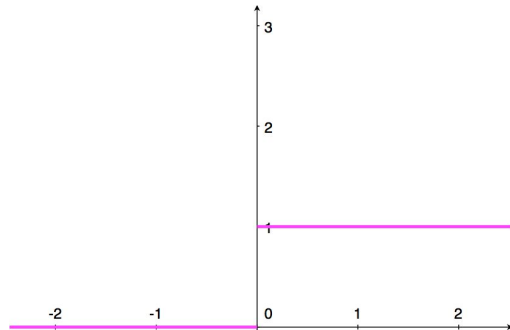
L2 Loss

$$L_i(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2$$



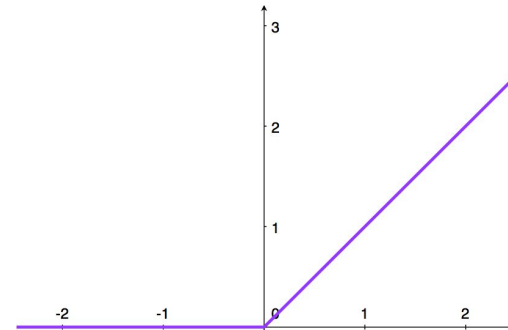
Zero-One Loss

$$L_i(y_i, \hat{y}_i) = 1 \text{ if } y_i \neq \hat{y}_i$$



Hinge Loss (only if y ranges [0,1])

$$L_i(y_i, \hat{y}_i) = \max(0, 1 - y_i \hat{y}_i)$$



Softmax Classifier (Multinomial Logistic Regression)

- **Recall:** we can treat the outputs of a model as probabilities for each class
- common way of measuring distance between probability distributions is **Kullback-Leibler (KL) divergence:**

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

- where P is the ground truth distribution and Q is the model's output score distribution

Softmax Classifier (Multinomial Logistic Regression)

KL divergence:

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

In our case, P is only non-zero for correct class

For example, consider the case we only have 3 classes:



$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

dog
cat
bird

correct outputs

Softmax Classifier (Multinomial Logistic Regression)

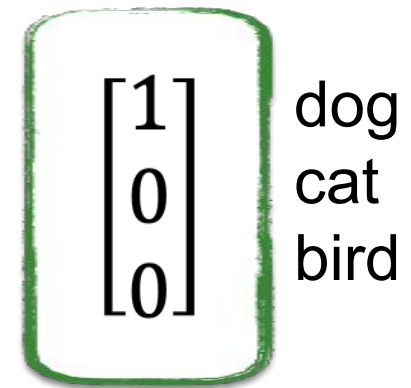
KL divergence:

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

$$= -\log Q(y) \text{ when } y = \text{dog}$$

$$= -\log \text{Prob}[f(x_i, W) = y_i]$$

(It's also called Cross Entropy)



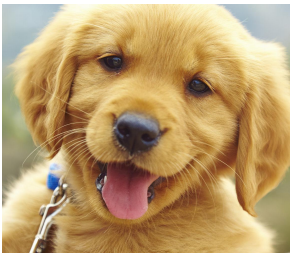
correct outputs

Softmax Classifier (Multinomial Logistic Regression)

$$L_i = -\log \text{Prob}[f(x_i, W) == y_i]$$

We need a mechanism to convert or normalize the output into probability range [0, 1]

Recall: **SOFTMAX**: $\text{Prob}[f(x_i, W) == k] = \frac{e^{\hat{y}_k}}{\sum_j e^{\hat{y}_j}}$



$$\begin{bmatrix} 3.2 \\ 5.1 \\ -1.7 \end{bmatrix}$$

model outputs

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

dog
cat
bird

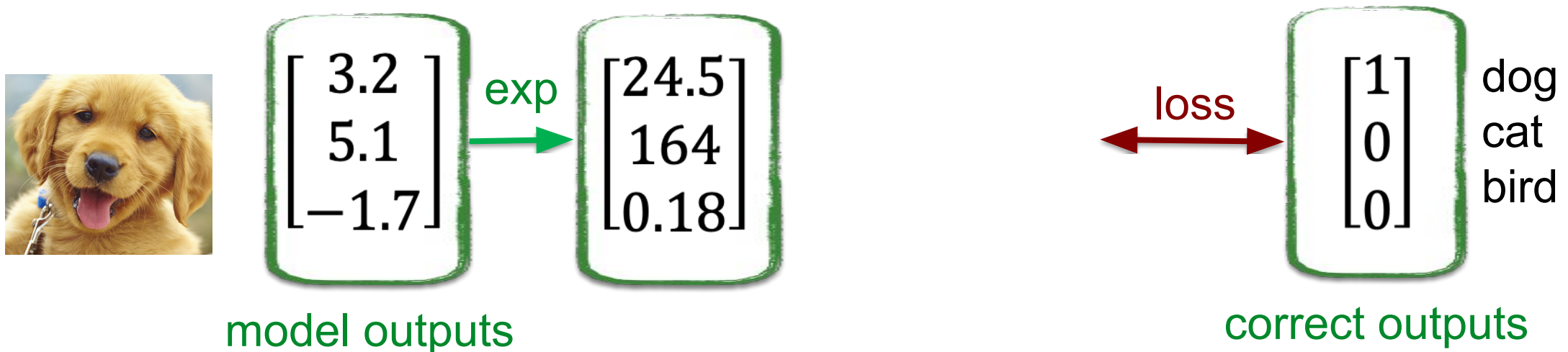
correct outputs

Softmax Classifier (Multinomial Logistic Regression)

$$L_i = -\log \text{Prob}[f(x_i, W) == y_i]$$

We need a mechanism to convert or normalize the output into probability range [0, 1]

Recall: **SOFTMAX**: $\text{Prob}[f(x_i, W) == k] = \frac{e^{\hat{y}_k}}{\sum_j e^{\hat{y}_j}}$

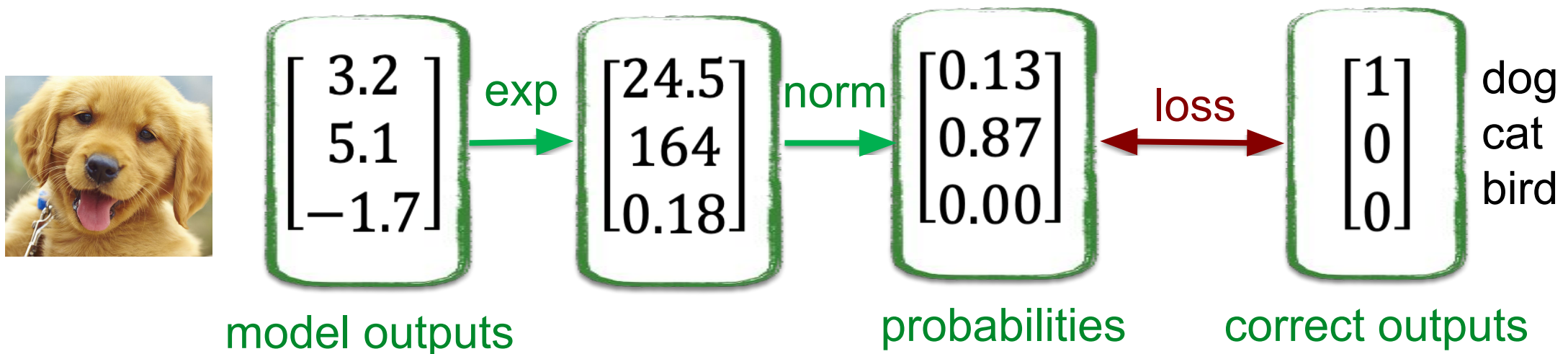


Softmax Classifier (Multinomial Logistic Regression)

$$L_i = -\log \text{Prob}[f(x_i, W) == y_i]$$

We need a mechanism to convert or normalize the output into probability range [0, 1]

Recall: **SOFTMAX**: $\text{Prob}[f(x_i, W) == k] = \frac{e^{\hat{y}_k}}{\sum_j e^{\hat{y}_j}}$

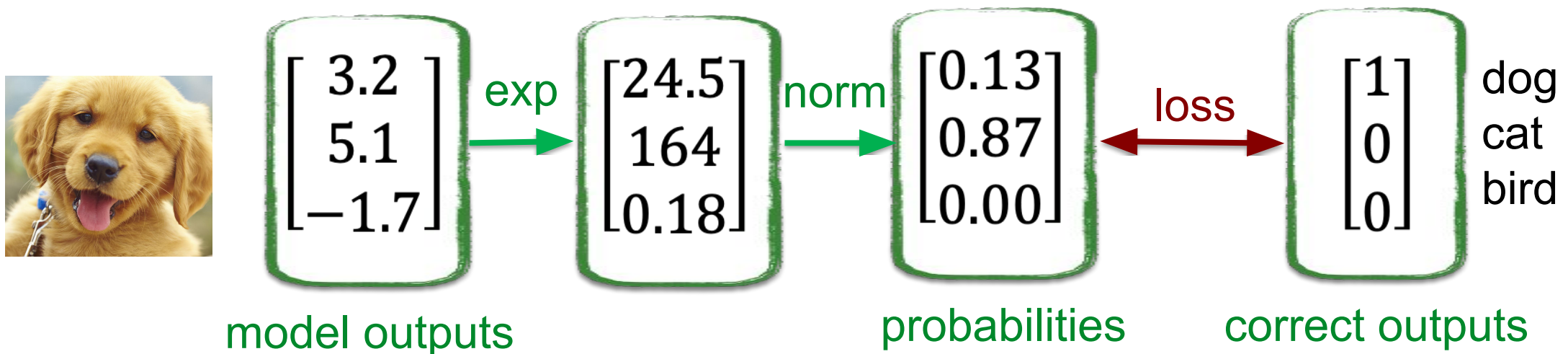


Softmax Classifier (Multinomial Logistic Regression)

$$L_i = -\log \text{Prob}[f(x_i, W) == y_i]$$

In this case, **what is the loss:**

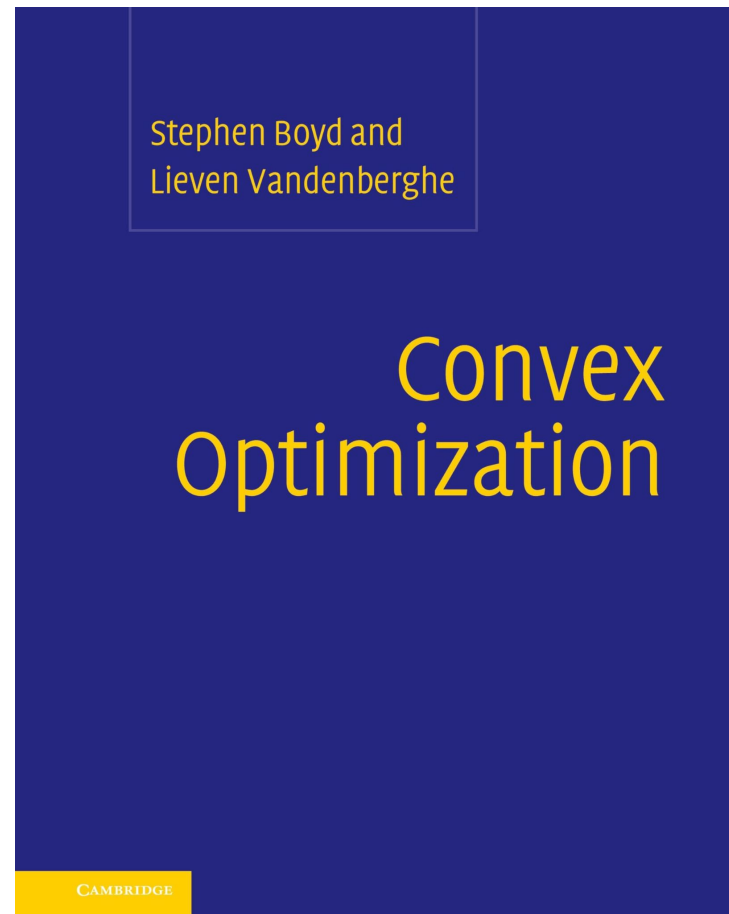
$$L_i = -\log(0.13) = 2.04$$



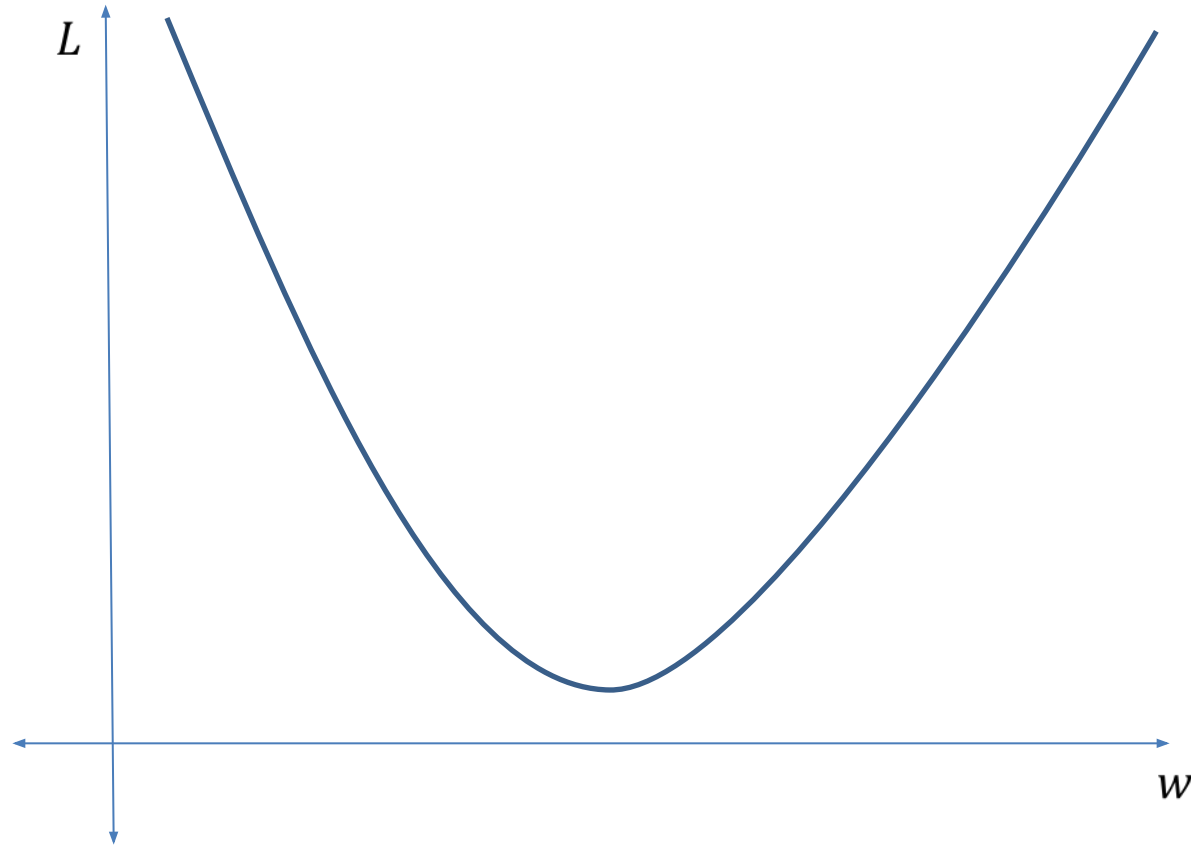
Today's agenda

- Perceptron
- Linear classifier
- Loss function
- Gradient descent and backpropagation
- Neural networks

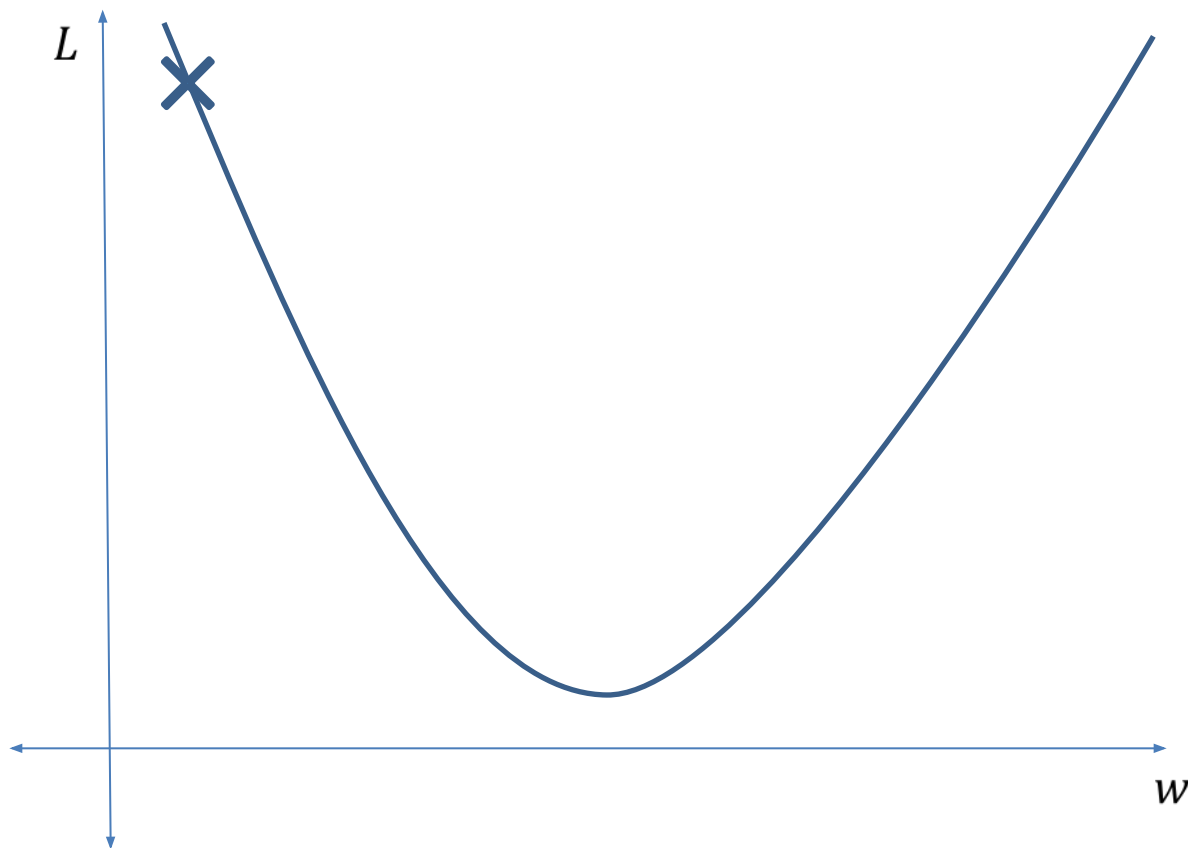
How do we find the weights that minimize the loss?



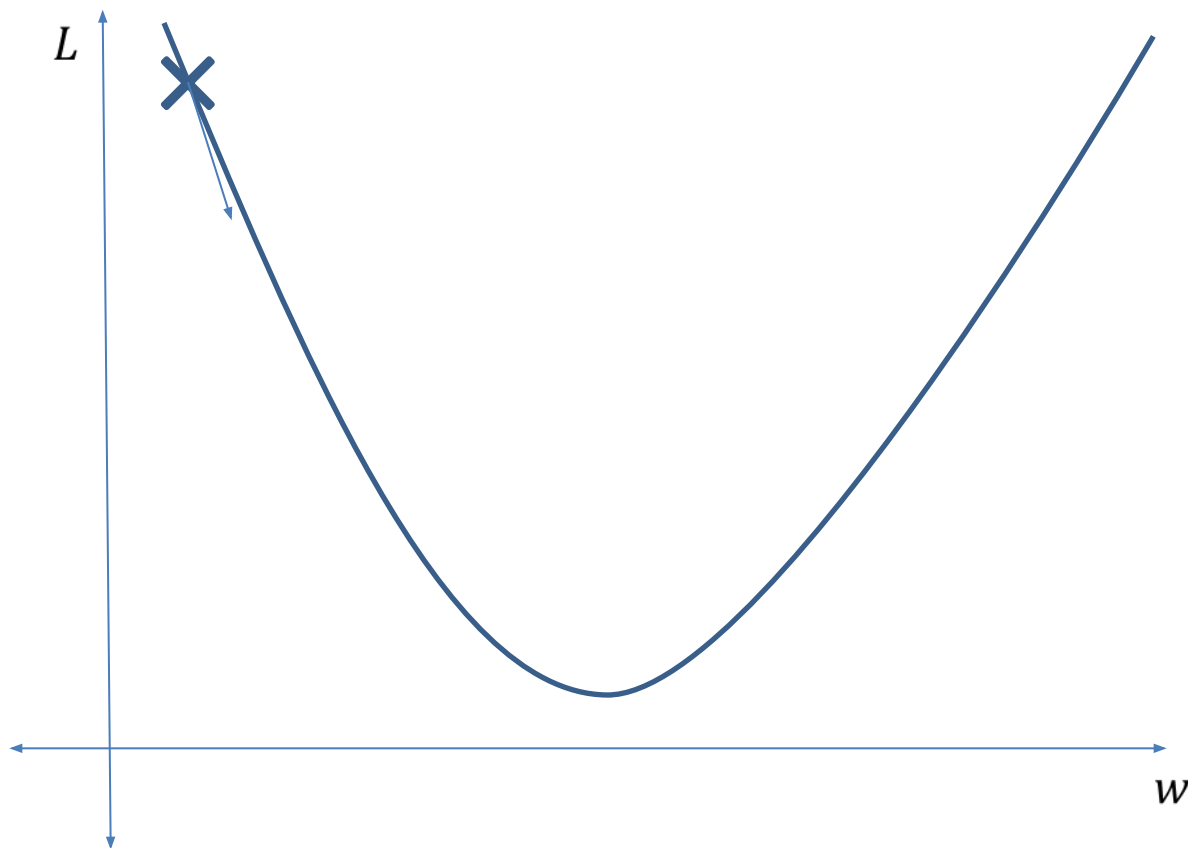
Gradient descent visualized: Minimizing loss



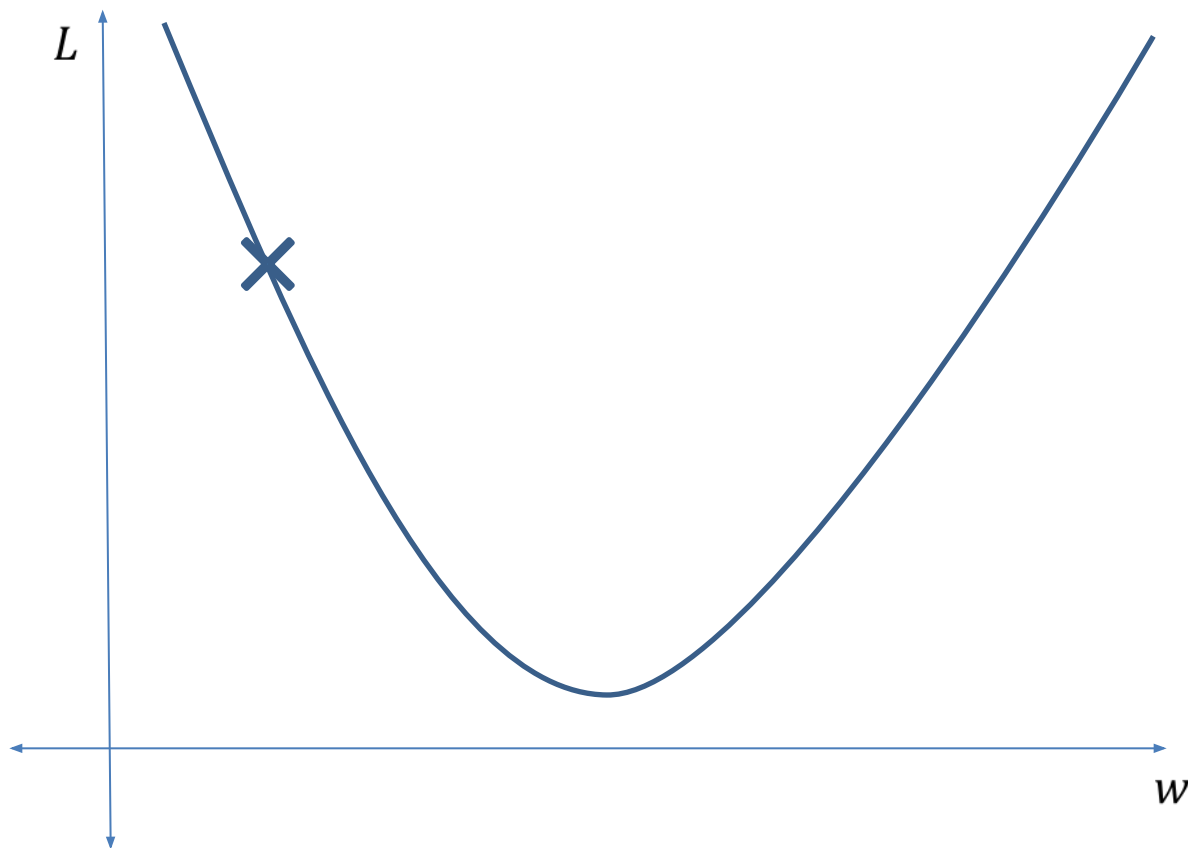
Gradient descent visualized: Minimizing loss



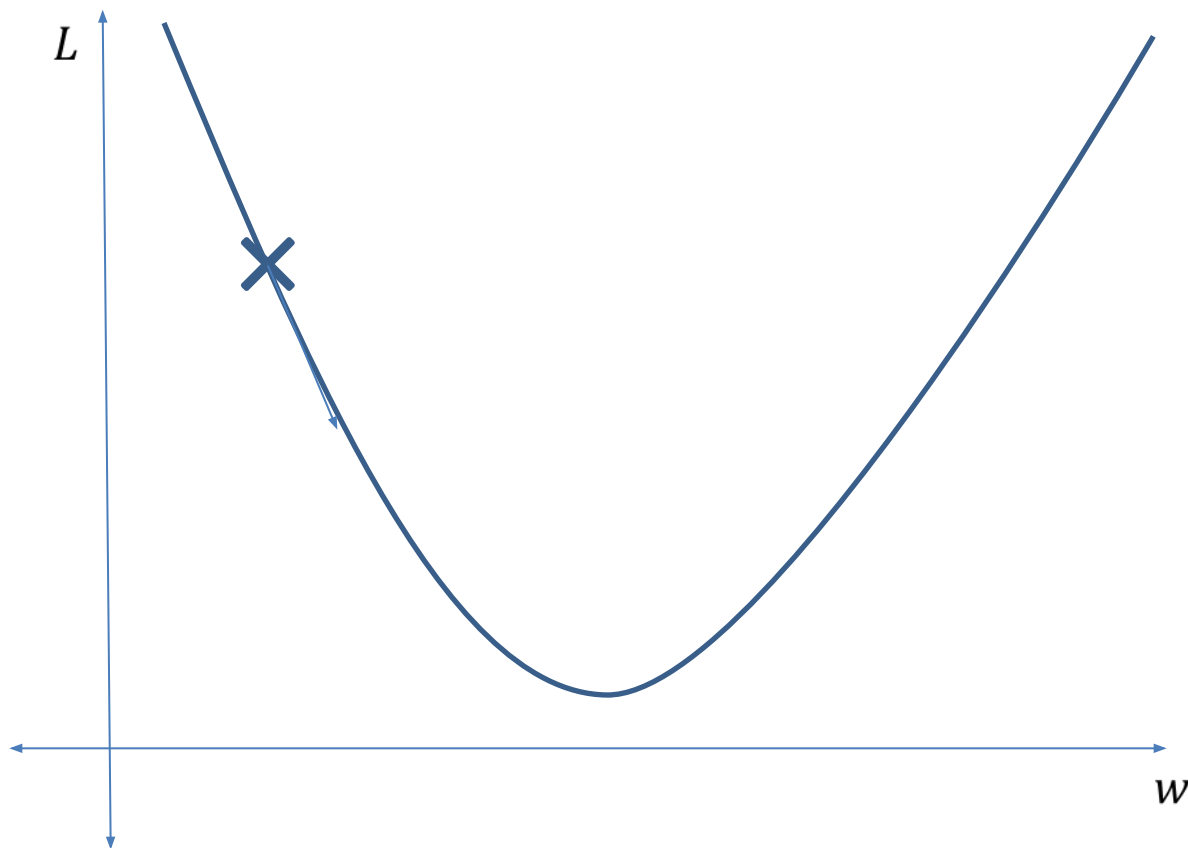
Gradient descent visualized: Minimizing loss



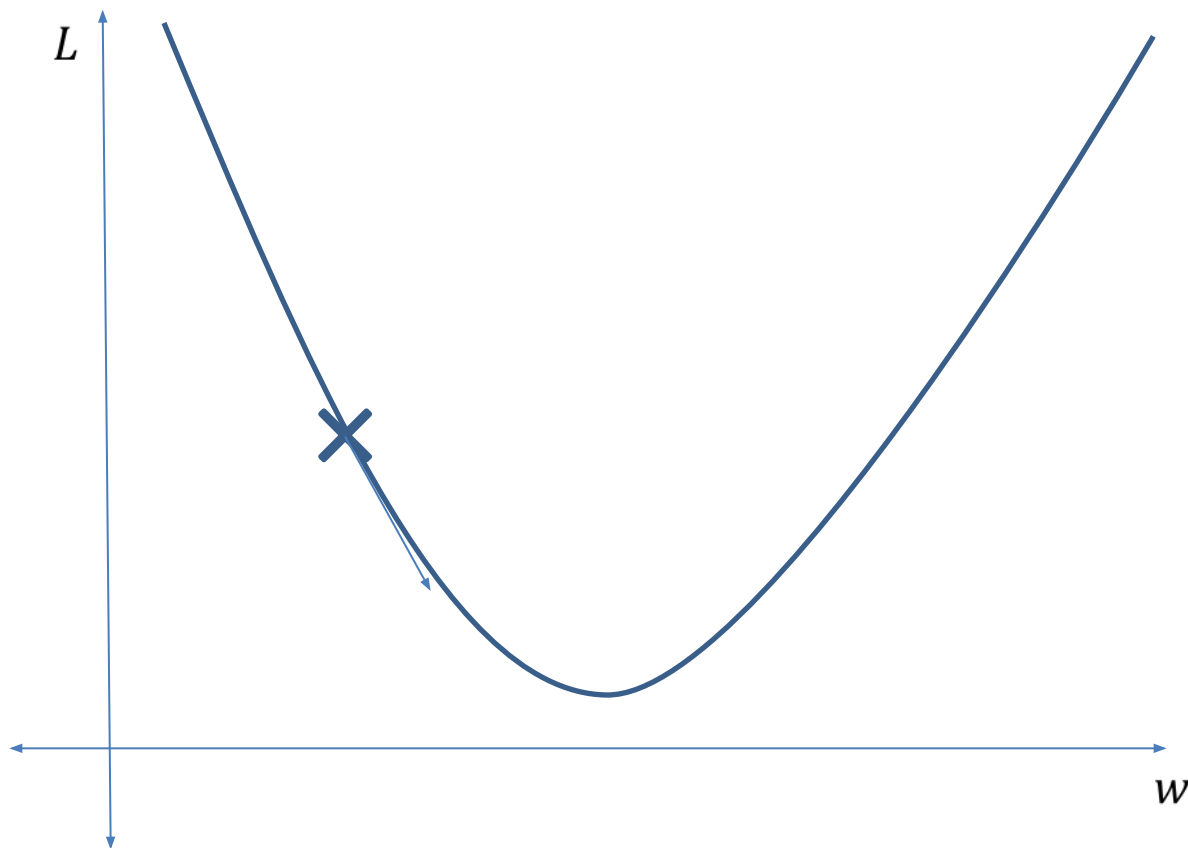
Gradient descent visualized: Minimizing loss



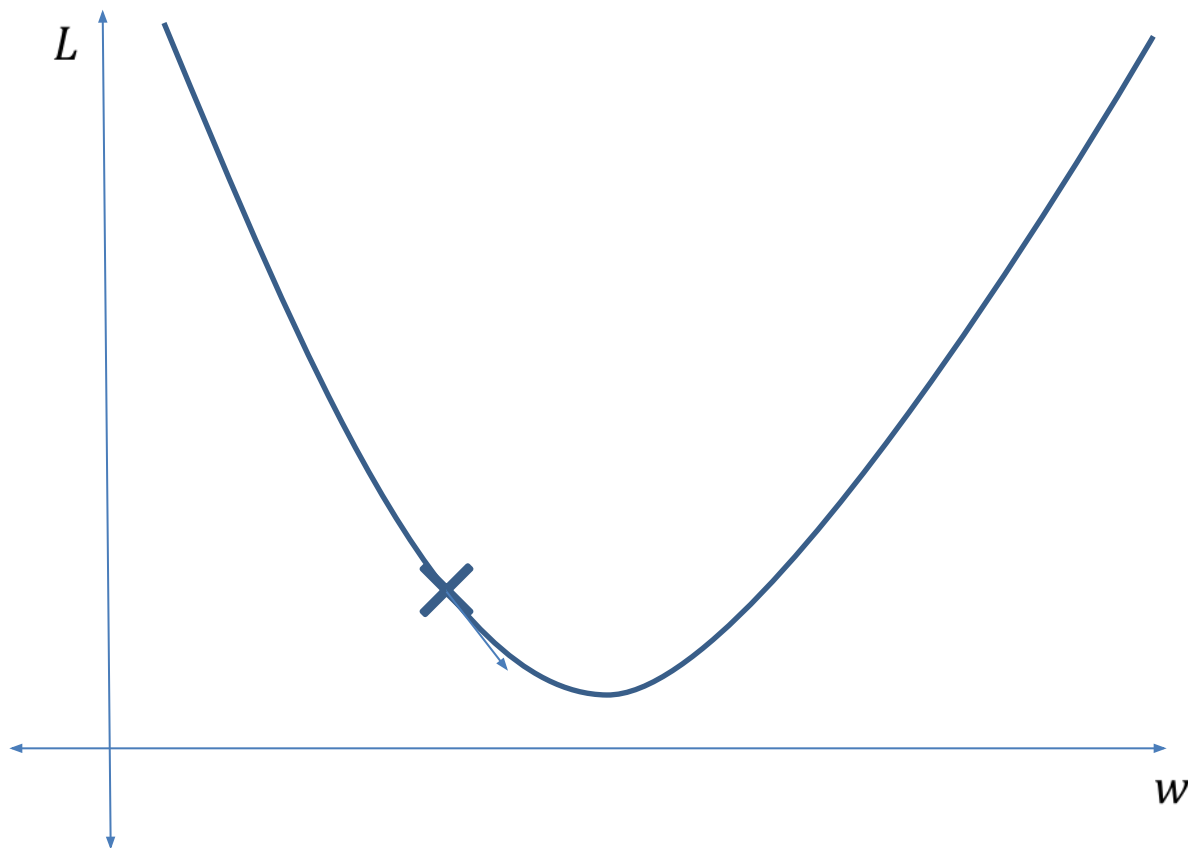
Gradient descent visualized: Minimizing loss



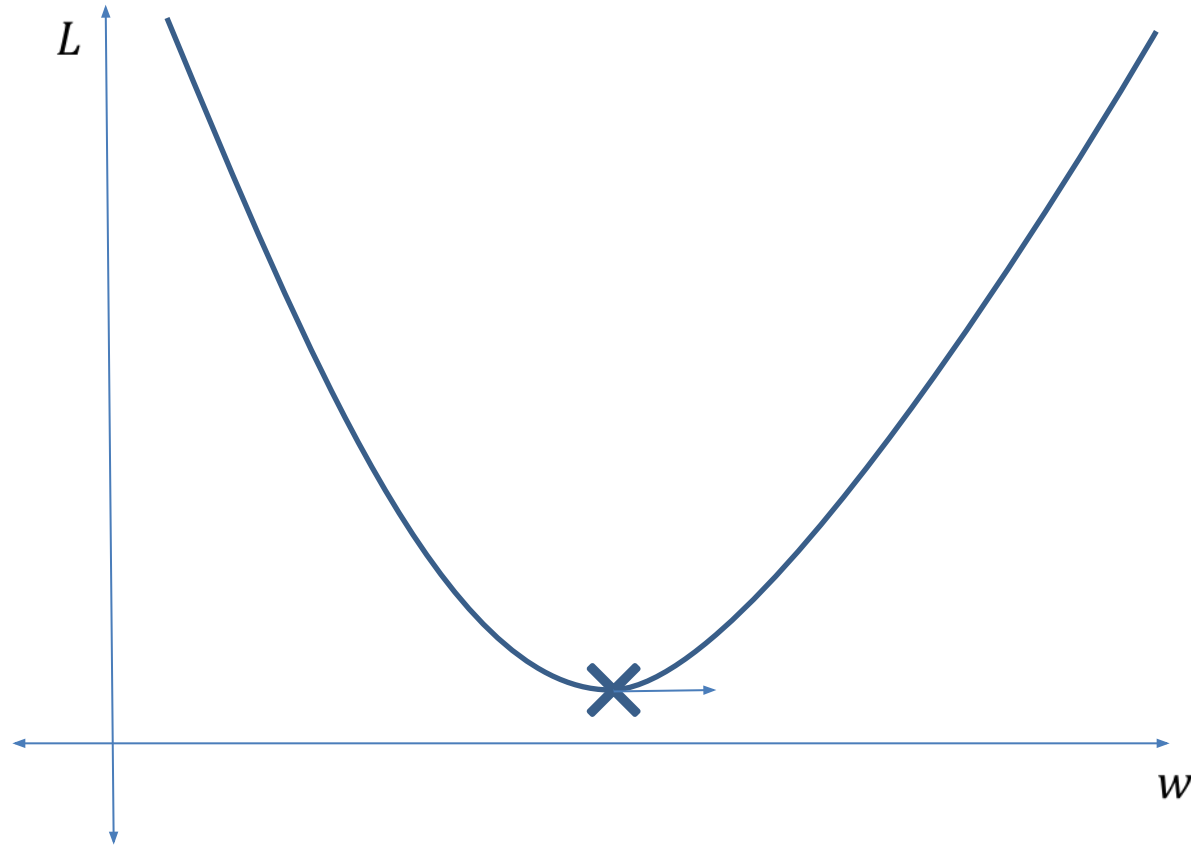
Gradient descent visualized: Minimizing loss



Gradient descent visualized: Minimizing loss



Gradient descent visualized: Minimizing loss



Gradient Descent Pseudocode

```
for _ in {0,...,num_epochs}:
```

```
     $L = 0$ 
```

```
    for  $x_i, y_i$  in data:
```

```
         $\hat{y}_i = f(x_i, \mathbf{W})$ 
```

```
         $L += L_i(y_i, \hat{y}_i)$ 
```

```
     $\frac{dL}{d\mathbf{W}} = ???$ 
```

```
     $\mathbf{W} := \mathbf{W} - \alpha \frac{dL}{d\mathbf{W}}$ 
```

Small step x Gradient

Gradient Descent Pseudocode

```
for _ in {0,...,num_epochs}:
```

```
    L = 0
```

```
    for  $x_i, y_i$  in data:
```

$$\hat{y}_i = f(x_i, \mathbf{W})$$

$$L += L_i(y_i, \hat{y}_i)$$

$$\frac{dL}{d\mathbf{W}} = ???$$

$$\mathbf{W} := \mathbf{W} - \alpha \frac{dL}{d\mathbf{W}}$$

Small step x Gradient

Exercise on linear classification:

$$\hat{\mathbf{y}} = \mathbf{W}\mathbf{x}$$

$$\text{Loss} = -\hat{y}_k + \log \sum_j e^{\hat{y}_j}$$

$$\frac{dL}{d\mathbf{W}} = \begin{bmatrix} \frac{e^{\hat{y}_0}}{\sum_j e^{\hat{y}_j}} \\ \dots \\ -1 + \frac{e^{\hat{y}_k}}{\sum_j e^{\hat{y}_j}} \\ \dots \\ \frac{e^{\hat{y}_{3071}}}{\sum_j e^{\hat{y}_j}} \end{bmatrix} \mathbf{x}$$

Partial derivative of loss to update weights

Given training data point (\mathbf{x}, \mathbf{y}) , the linear classifier formula is: $\hat{\mathbf{y}} = \mathbf{W}\mathbf{x}$

Let's assume that the correct label is class \mathbf{k} , implying $\mathbf{y}=\mathbf{k}$

$$\text{Loss} = -\hat{y}_k + \log \sum_j e^{\hat{y}_j}$$

Now, we want to update the weights \mathbf{W} by calculating the direction in which to change the weights to reduce the loss: $\frac{dL}{dW} = \frac{dL}{d\hat{\mathbf{y}}} \frac{d\hat{\mathbf{y}}}{dW}$

we know that $\frac{d\hat{\mathbf{y}}}{dW} = \mathbf{x}$, but what about $\frac{dL}{d\hat{\mathbf{y}}}$?

Partial derivative of loss to update weights

$$L = -\hat{y}_k + \log \sum_j e^{\hat{y}_j}$$

To calculate $\frac{dL}{d\hat{y}}$, we need to consider two cases:

Case 1:

$$\frac{dL}{d\hat{y}_k} = -1 + \frac{e^{\hat{y}_k}}{\sum_j e^{\hat{y}_j}}$$

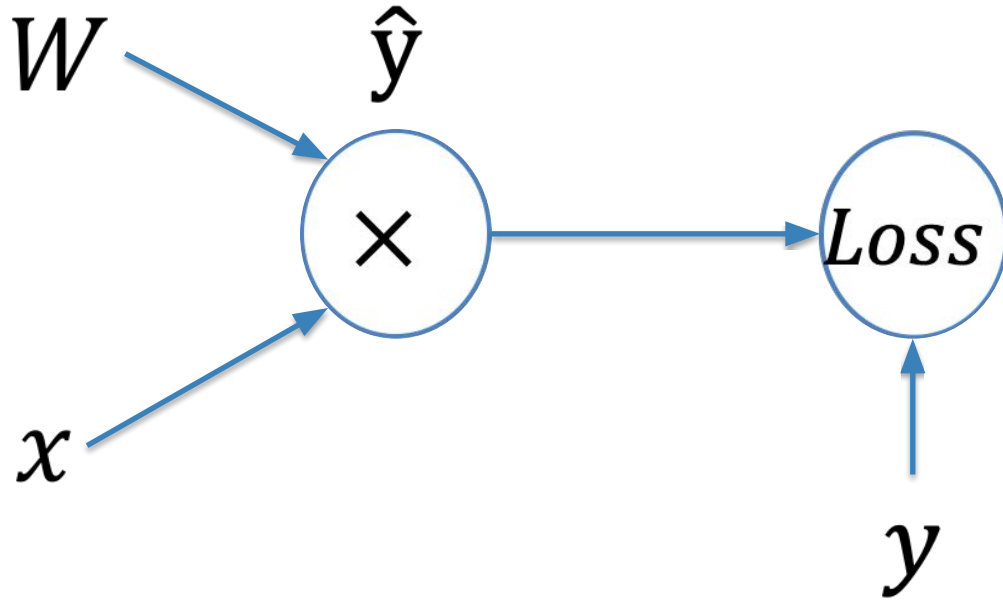
Case 2:

$$\frac{dL}{d\hat{y}_{l \neq k}} = \frac{e^{\hat{y}_l}}{\sum_j e^{\hat{y}_j}}$$

Exercise: Partial derivative of loss to update weights

$$\frac{dL}{dW} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dW}$$
$$\frac{dL}{dW} = \begin{bmatrix} \frac{e^{\hat{y}_0}}{\sum_j e^{\hat{y}_j}} \\ \dots \\ -1 + \frac{e^{\hat{y}_k}}{\sum_j e^{\hat{y}_j}} \\ \dots \\ \frac{e^{\hat{y}_{3071}}}{\sum_j e^{\hat{y}_j}} \end{bmatrix} x$$

Backprop – another way of computing gradients

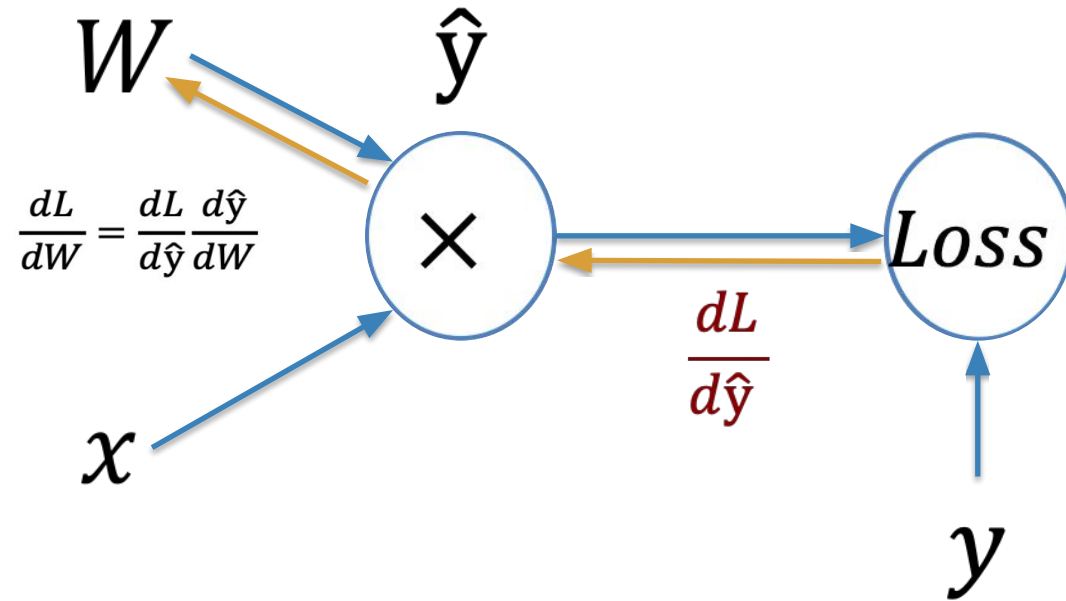


$$\hat{y} = Wx$$
$$L = \text{Loss}(\hat{y}, y)$$

$$\frac{dL}{dW} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dW}$$

Calculating the gradient is hard, but we can use the chain rule to make it simpler

Backprop – a way of computing gradients



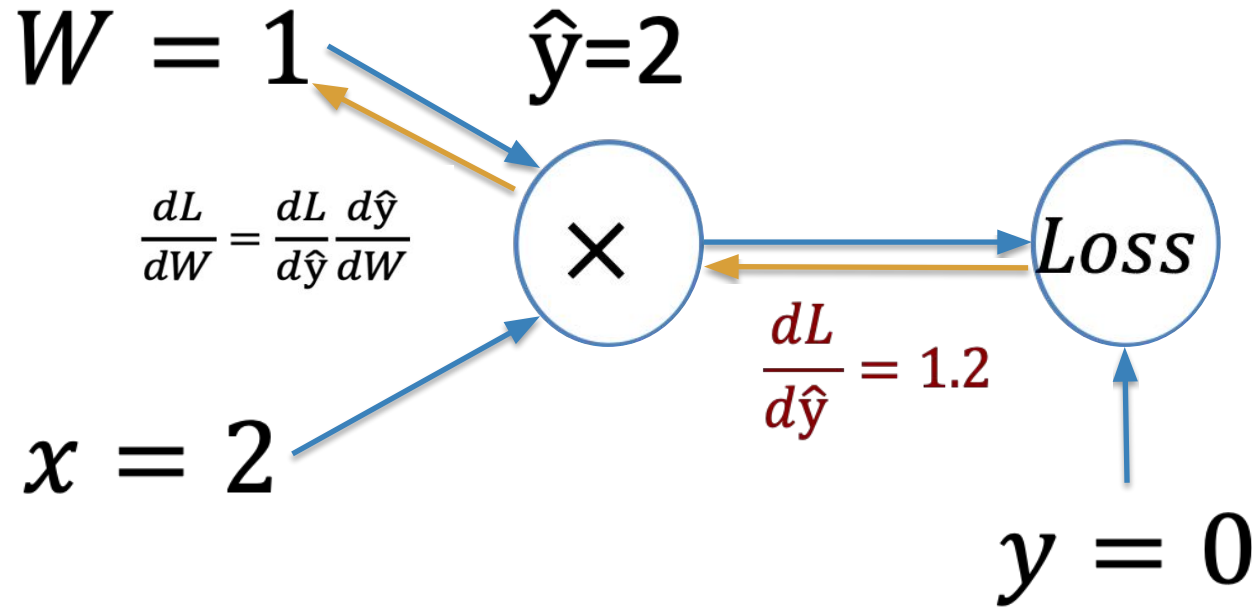
$$\hat{y} = Wx$$
$$L = \text{Loss}(\hat{y}, y)$$

$$\frac{dL}{dW} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dW}$$

Key Insight:

- visualize the **computation as a graph flow**
- Compute the **forward pass** to calculate the loss.
- Compute all **gradients** for each **pair of nodes backwards**

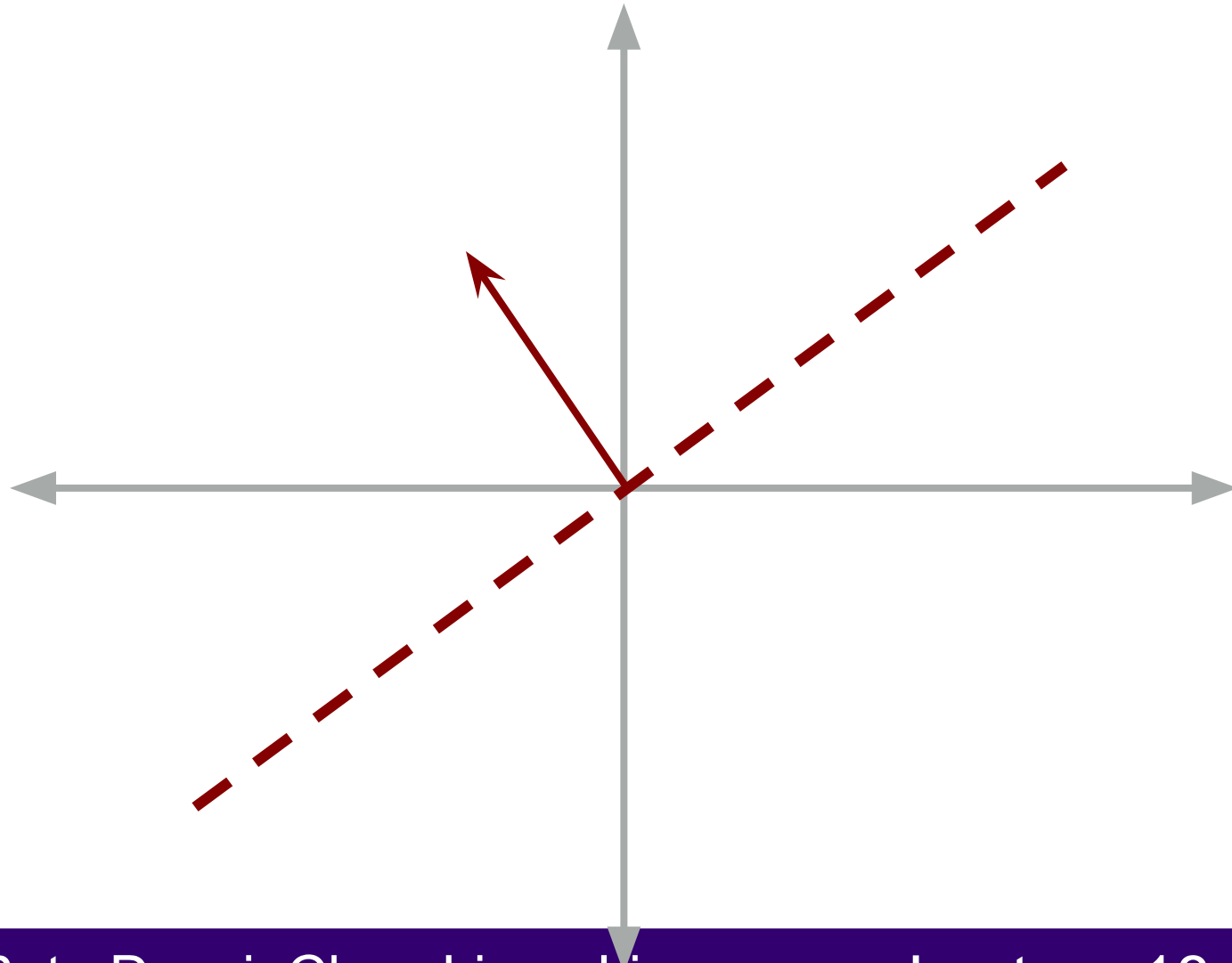
Backprop example in 1D:



We know the chain rule

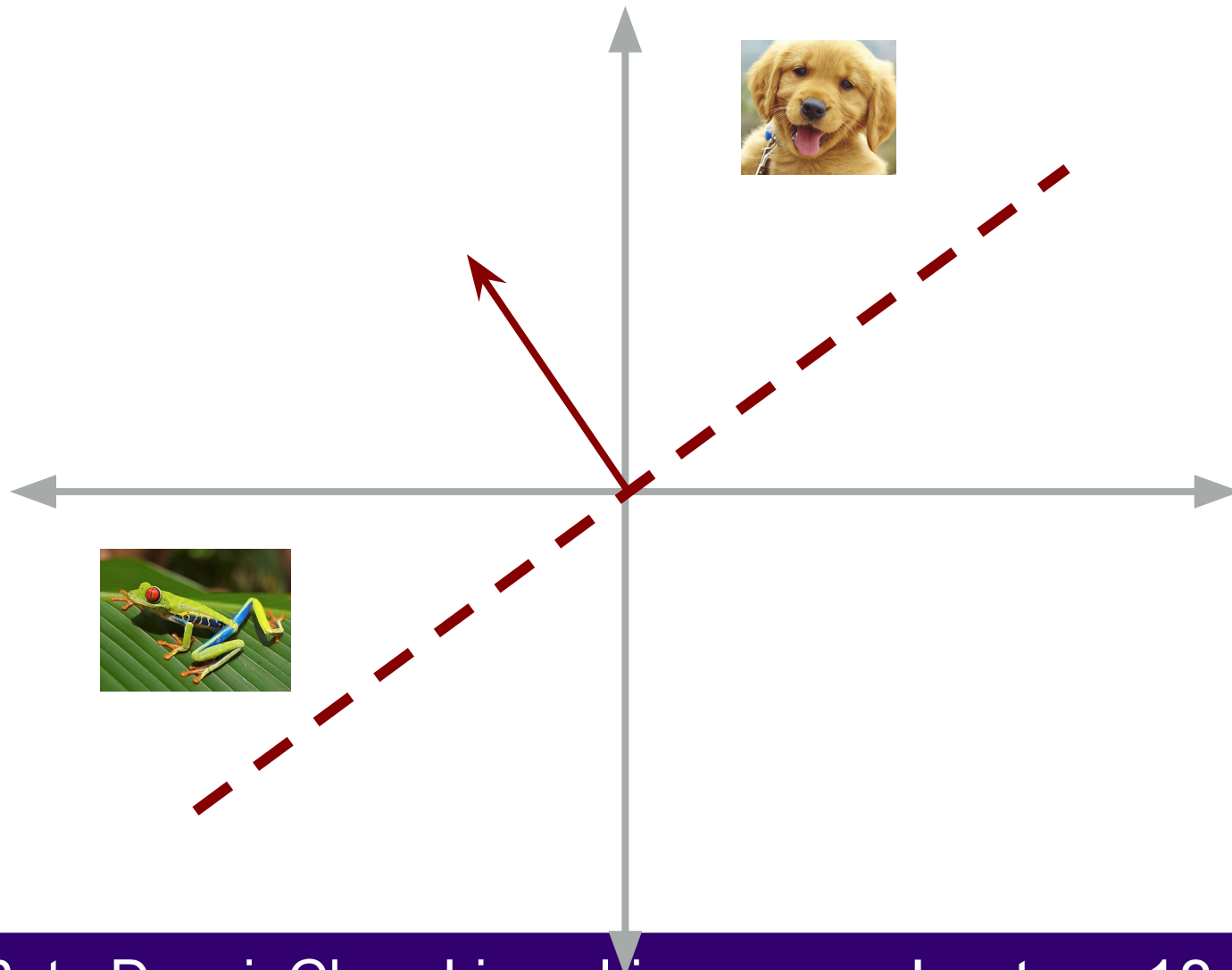
$$\begin{aligned}\frac{dL}{dW} &= \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dW} \\ &= \frac{dL}{d\hat{y}} x \\ &= 1.2x \\ &= 1.2 \times 2 \\ &= 2.4\end{aligned}$$

Interpreting the weights **geometrically**



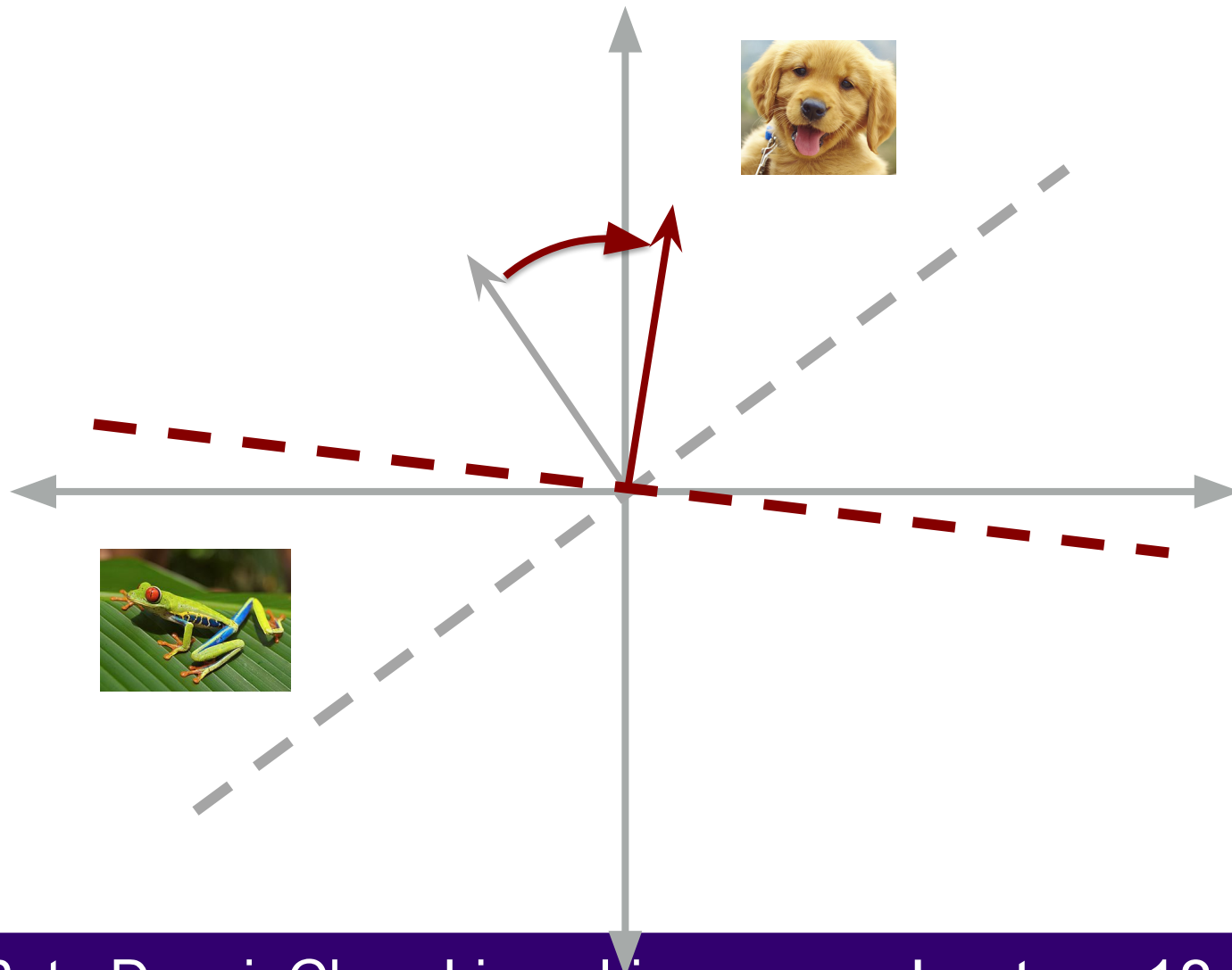
- Assume the image vectors are in 2D space to make it easier to visualize.
- Let's start with one class: **dog**.
- Initialize the weights randomly

Interpreting the weights **geometrically**



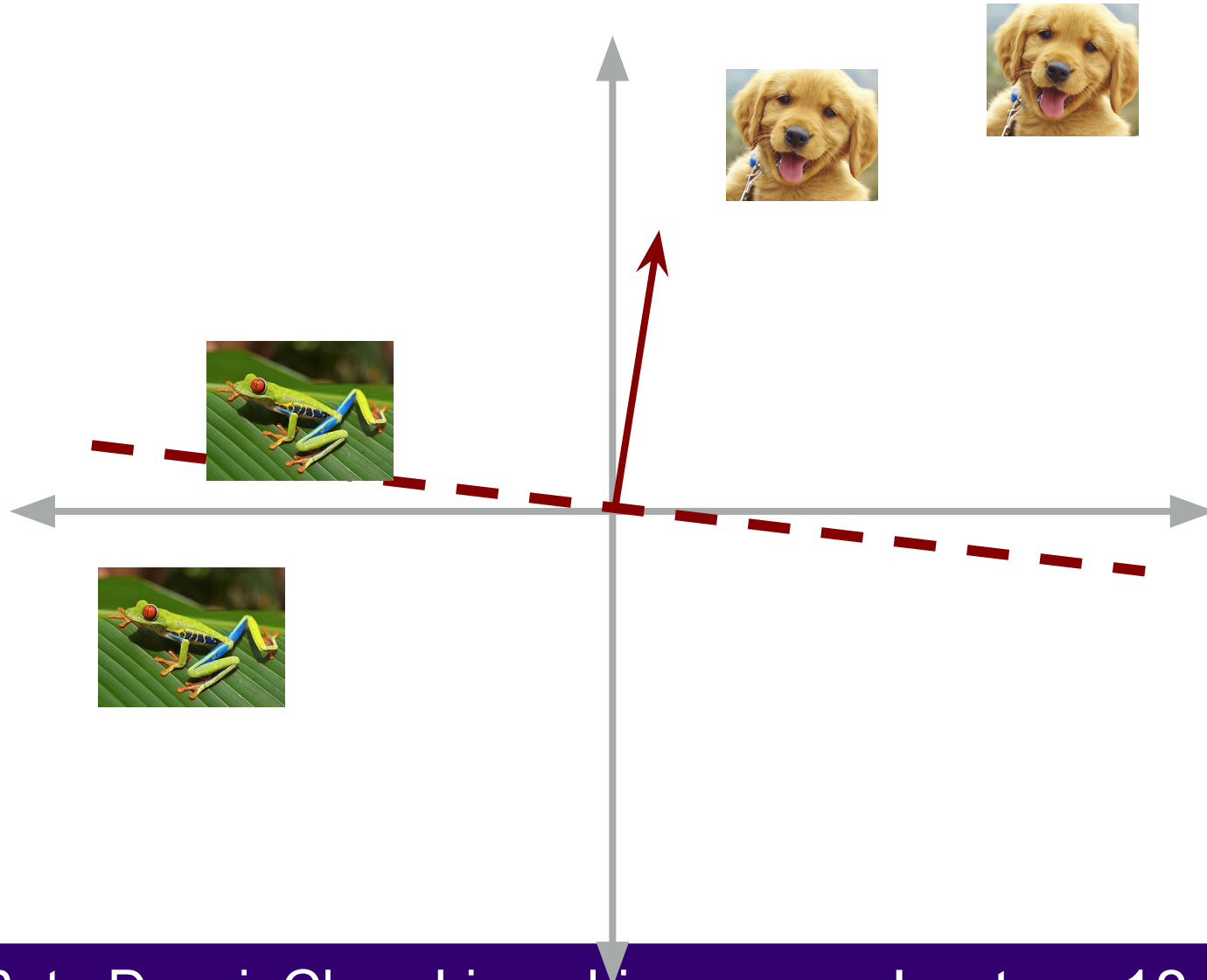
- Assume the image vectors are in 2D space to make it easier to visualize.
- Let's start with one class: **dog**.
- Initialize the weights randomly
- Now let's **add two data points**

Interpreting the weights **geometrically**



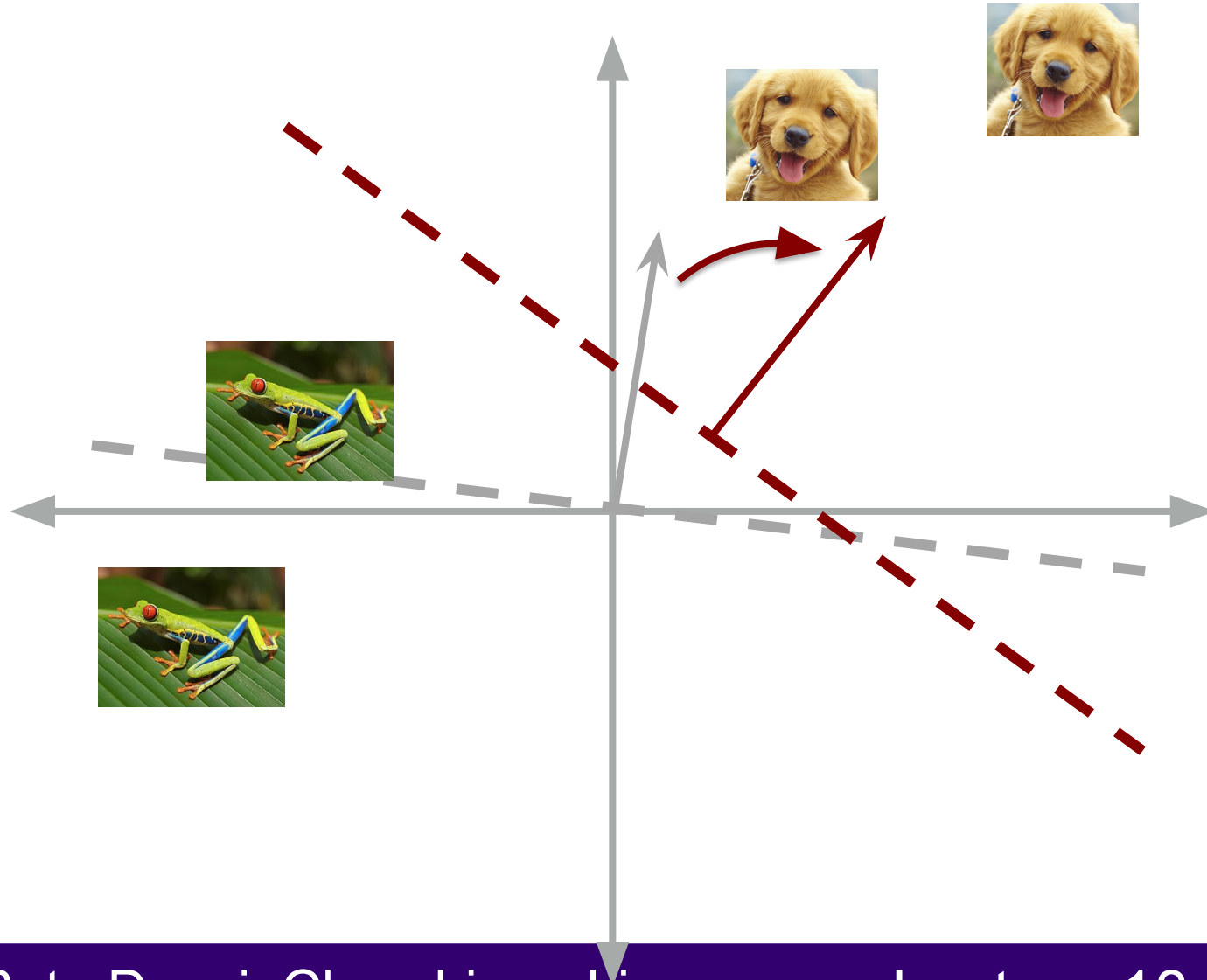
- Assume the image vectors are in 2D space to make it easier to visualize.
- Let's start with one class: **dog**.
- Initialize the weights randomly
- Now let's add two data points
- **Update the weights**

Interpreting the weights **geometrically**



- Assume the image vectors are in 2D space to make it easier to visualize.
- Let's start with one class: **dog**.
- Initialize the weights randomly
- Now let's **add two more data points**
- Update the weights

Interpreting the weights **geometrically**



- Assume the image vectors are in 2D space to make it easier to visualize.
- Let's start with one class: **dog**.
- Initialize the weights randomly
- Now let's add two more data points
- **Update the weights**

Interpreting the weights **geometrically**

