

# VIR2: MLE for classification

Karel Zimmermann

<http://cmp.felk.cvut.cz/~zimmerk/>



Vision for Robotics and Autonomous Systems

<https://cyber.felk.cvut.cz/vras/>



Center for Machine Perception

<https://cmp.felk.cvut.cz>

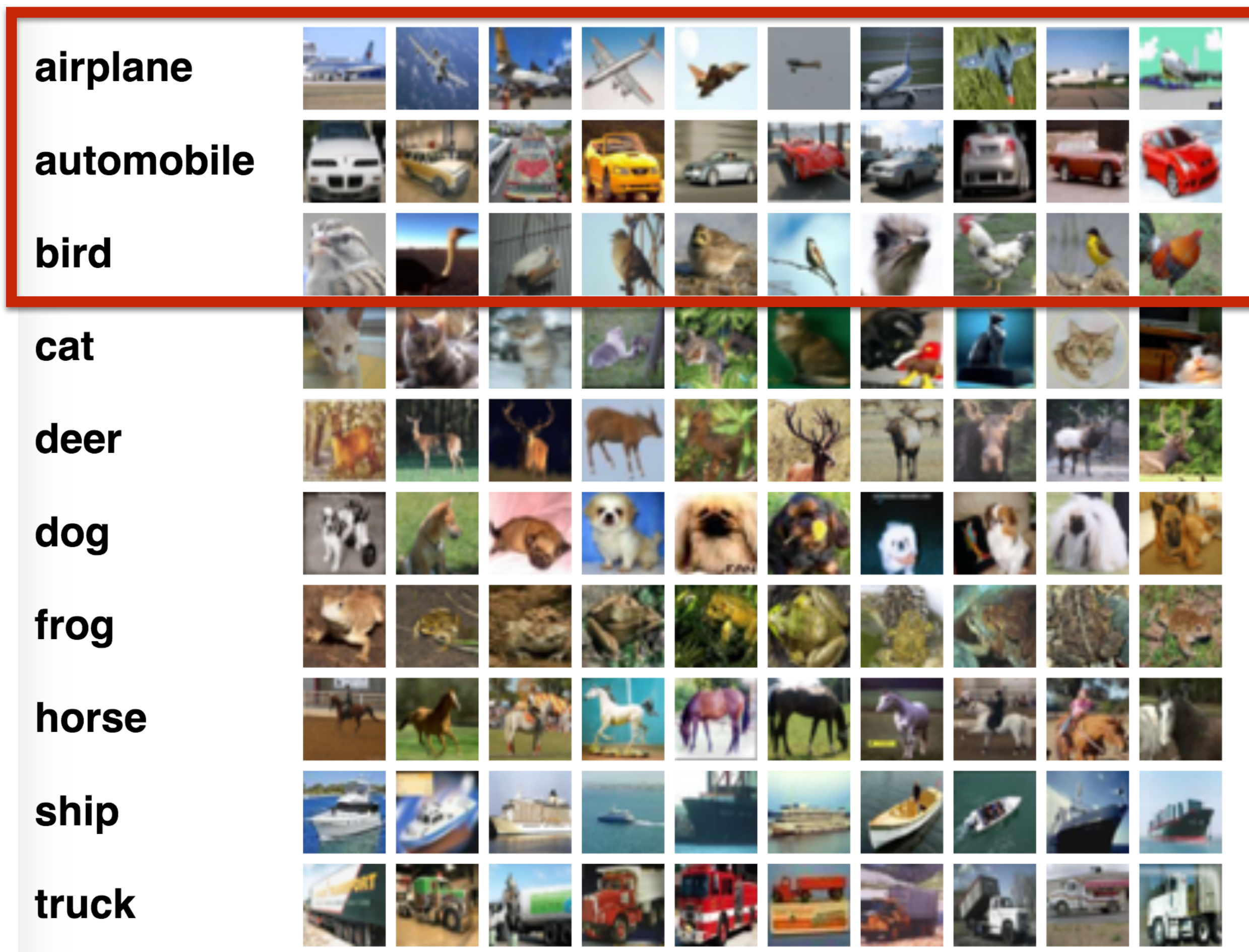


Department for Cybernetics  
Faculty of Electrical Engineering  
Czech Technical University in Prague

# Outline

- Pre-requisites: linear algebra, Bayes rule
- MAP/ML estimation, prior and overfitting
- Linear regression
- Linear classification

# Recognition problem



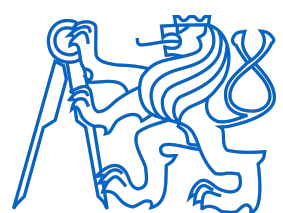
Why is it hard?

CIFAR-10: classify 32x32 RGB images into 10 categories

<https://www.cs.toronto.edu/~kriz/cifar.html>

Czech Technical University in Prague

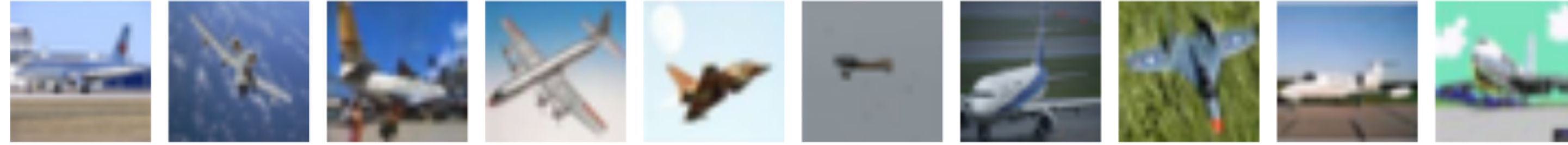
Faculty of Electrical Engineering, Department of Cybernetics



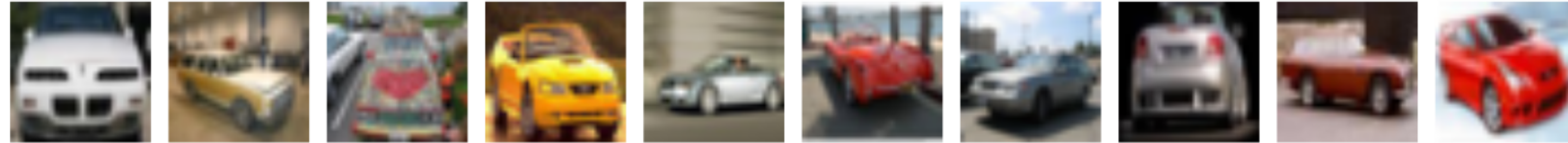
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

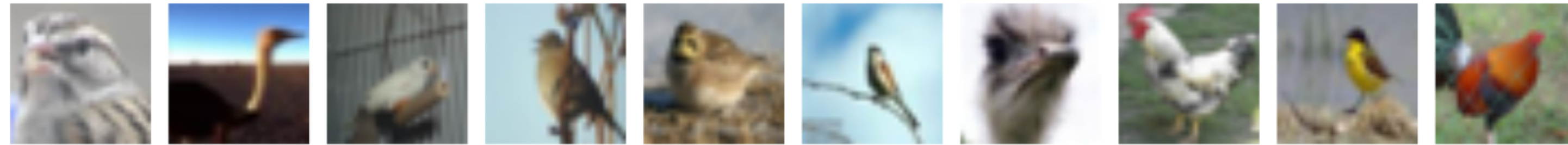
1



2



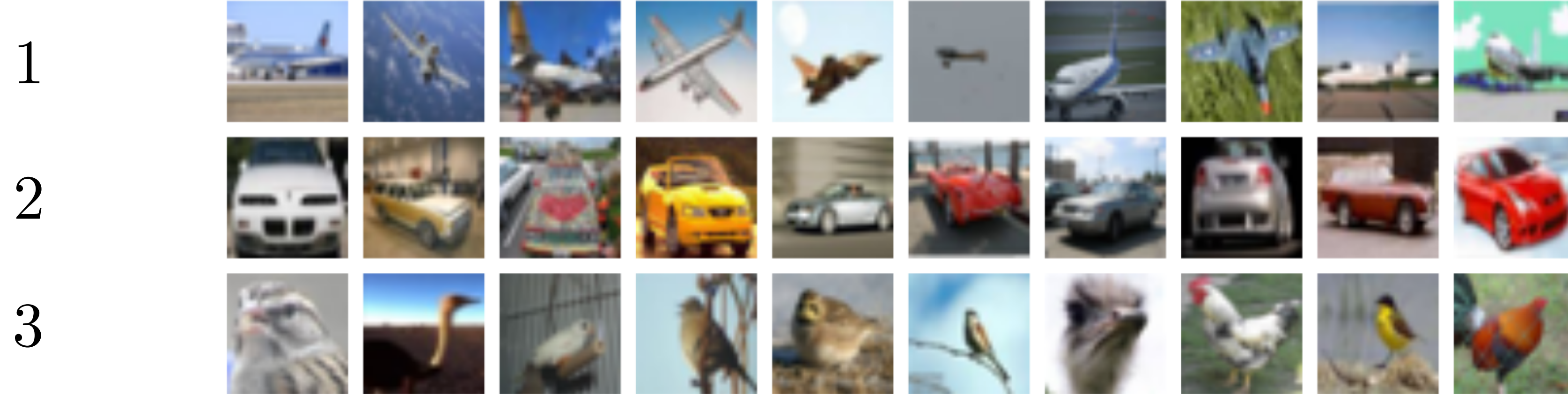
3



**Three-class** recognition problem:

Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )



**Three-class** recognition problem:

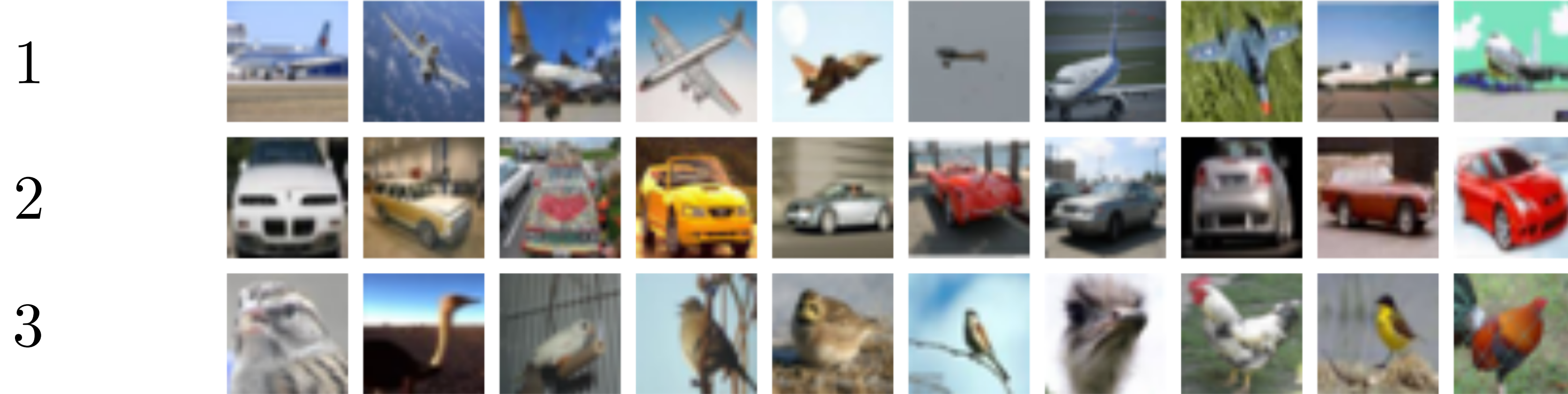
```
def classify():
```

```
    ???
```

```
    return  $\mathbf{p}$ 
```

Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

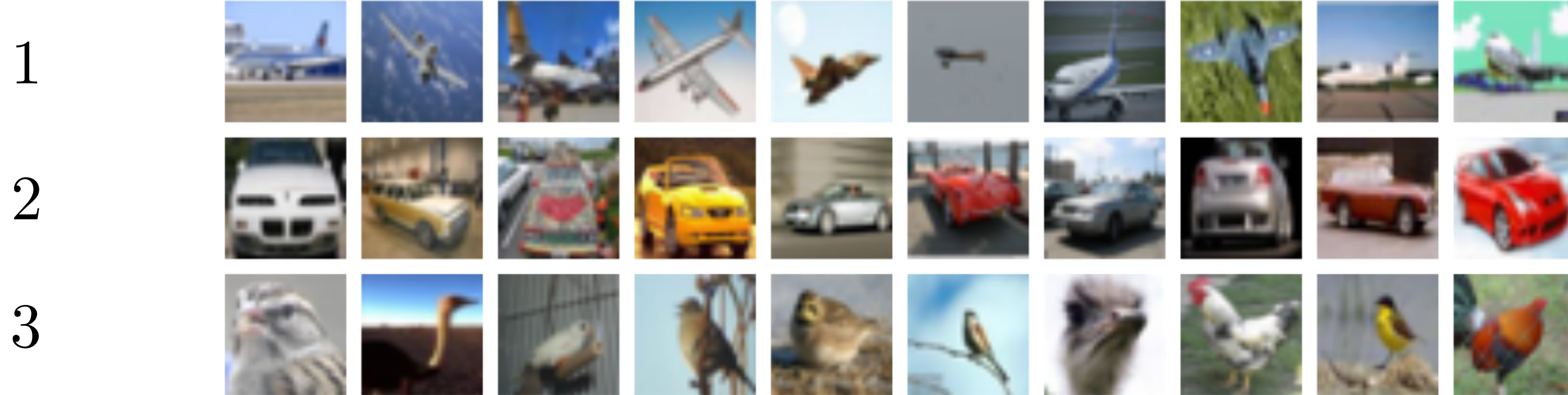


Model probability distribution over classes by softmax function

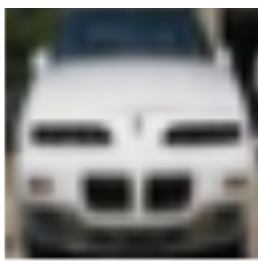
$$p(y|\mathbf{x}, W) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k))} = \mathbf{s}(\mathbf{f}(\mathbf{x}, W))$$

Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )



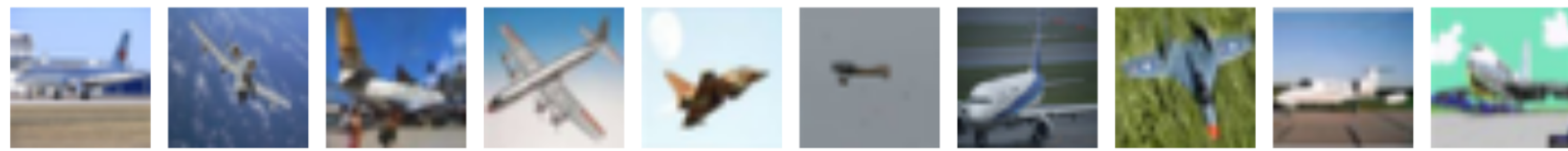
Example: Linear classifier:

```
def classify(  
     $\mathbf{p} = \mathbf{s}(W \bar{\mathbf{x}})$   
    return  $\mathbf{p}$ 
```

$$W \bar{\mathbf{x}} = \begin{bmatrix} -2 \\ +1 \\ 0 \end{bmatrix}$$

$$\mathbf{s} \left( \begin{bmatrix} -2 \\ +1 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} 0.03 \\ 0.71 \\ 0.26 \end{bmatrix}$$

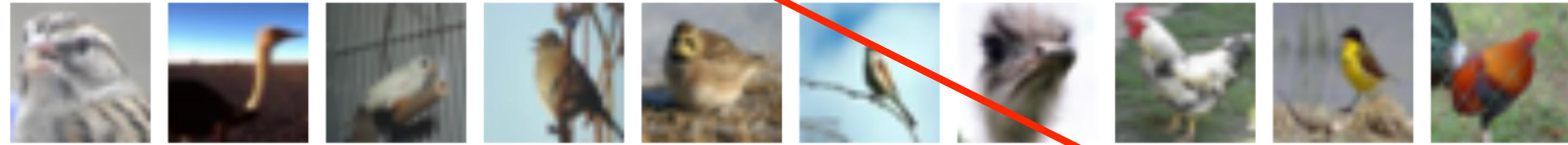
1



2

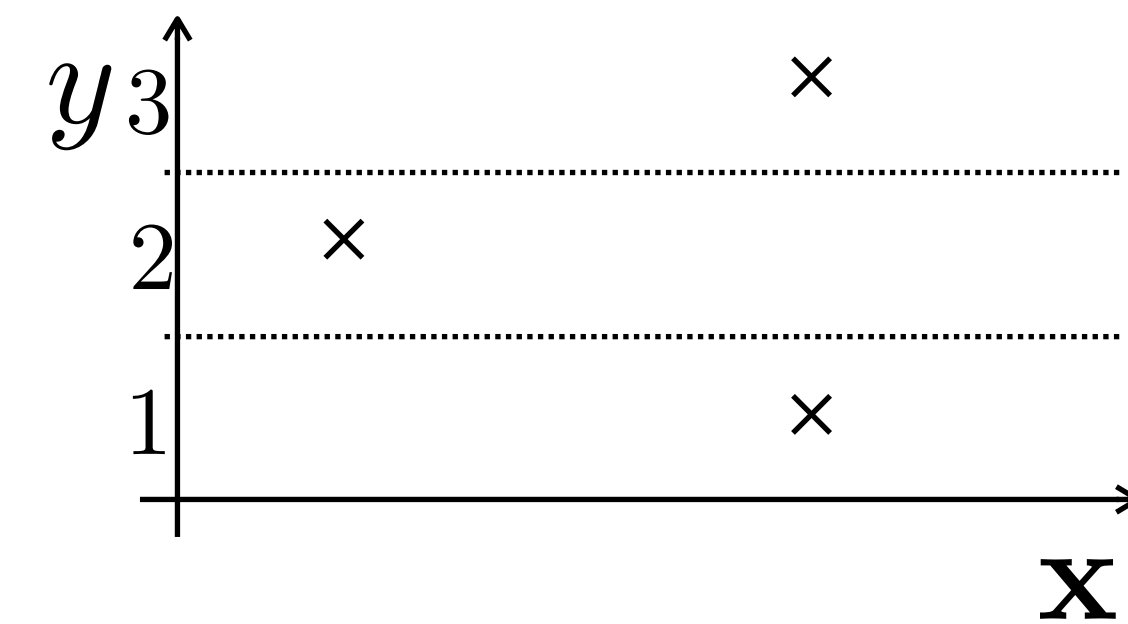


3



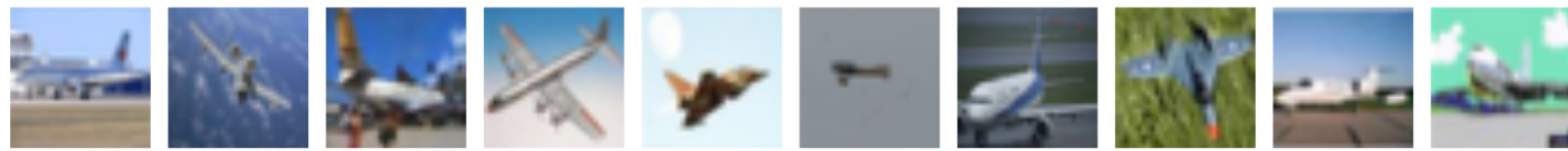
- **Classification** (probability modeled by soft-max function):

$$p(y|\mathbf{x}, W) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k))} = \mathbf{s}(\mathbf{f}(\mathbf{x}, W))$$





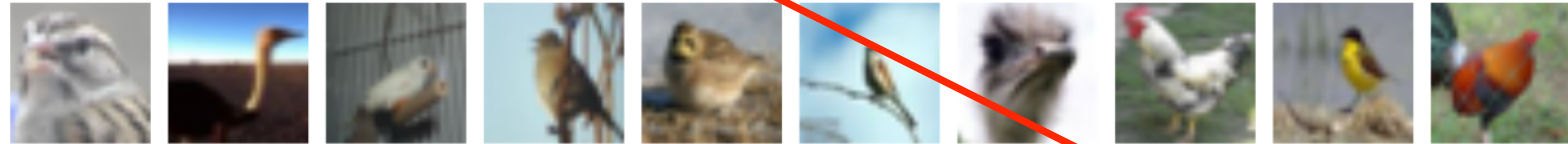
1



2

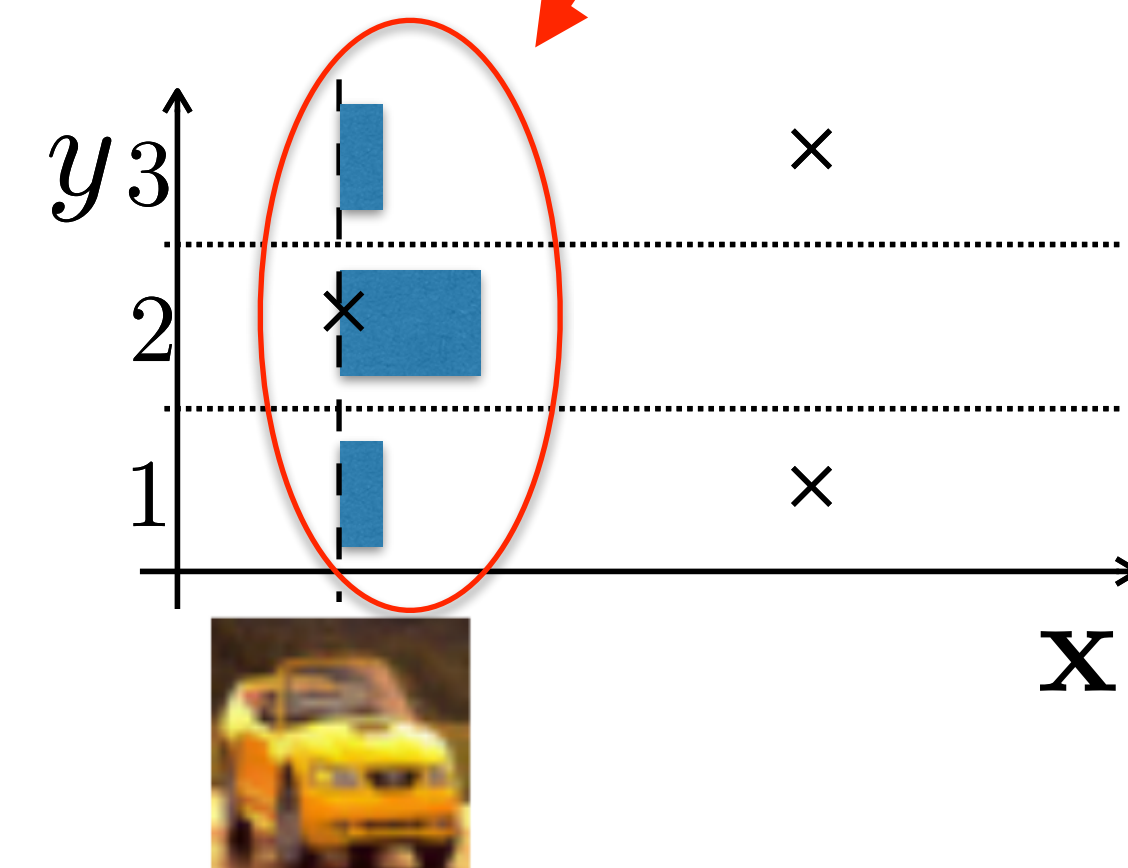


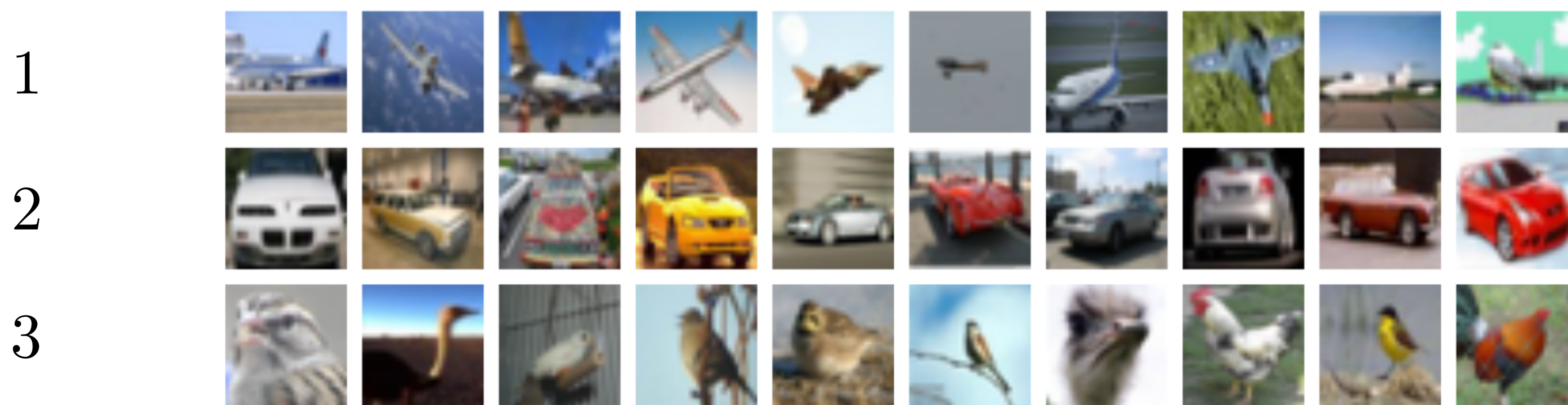
3



- **Classification** (probability modeled by soft-max function):

$$p(y|\mathbf{x}, W) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k))} = \mathbf{s}(\mathbf{f}(\mathbf{x}, W))$$



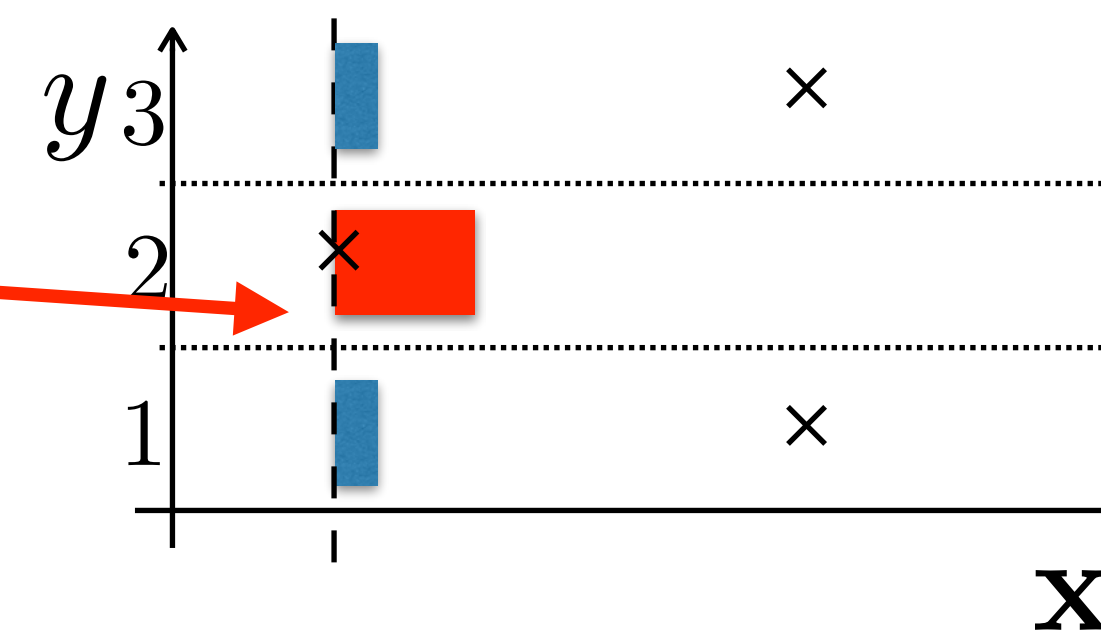


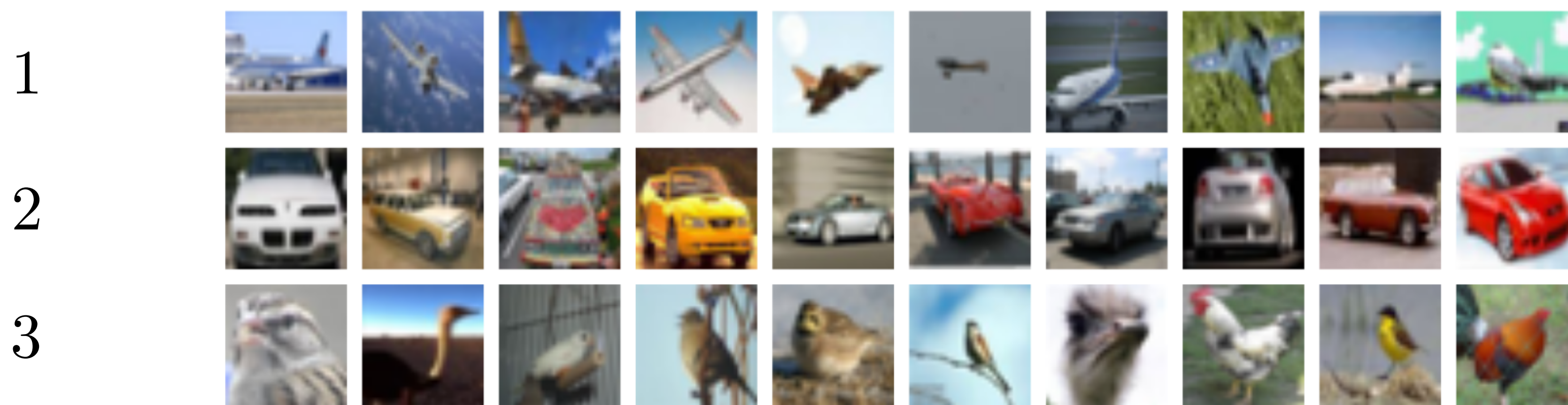
- **Classification** (probability modeled by soft-max function):

$$p(y|\mathbf{x}, W) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k))} = \mathbf{s}(\mathbf{f}(\mathbf{x}, W))$$

- Probability of observing  $y_i$  when measuring  $\mathbf{x}_i$  is

$$p(y_i|\mathbf{x}_i, W) = \mathbf{s}_{y_i}(\mathbf{f}(\mathbf{x}_i, W))$$





- **Classification** (probability modeled by soft-max function):

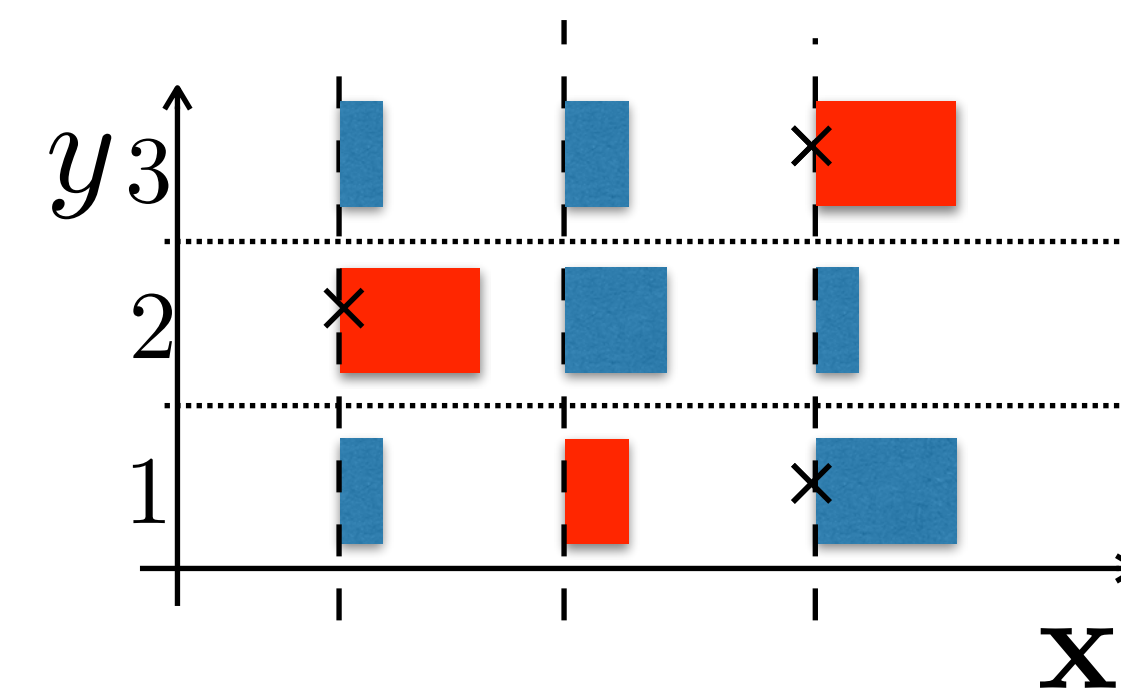
$$p(y|\mathbf{x}, W) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k))} = \mathbf{s}(\mathbf{f}(\mathbf{x}, W))$$

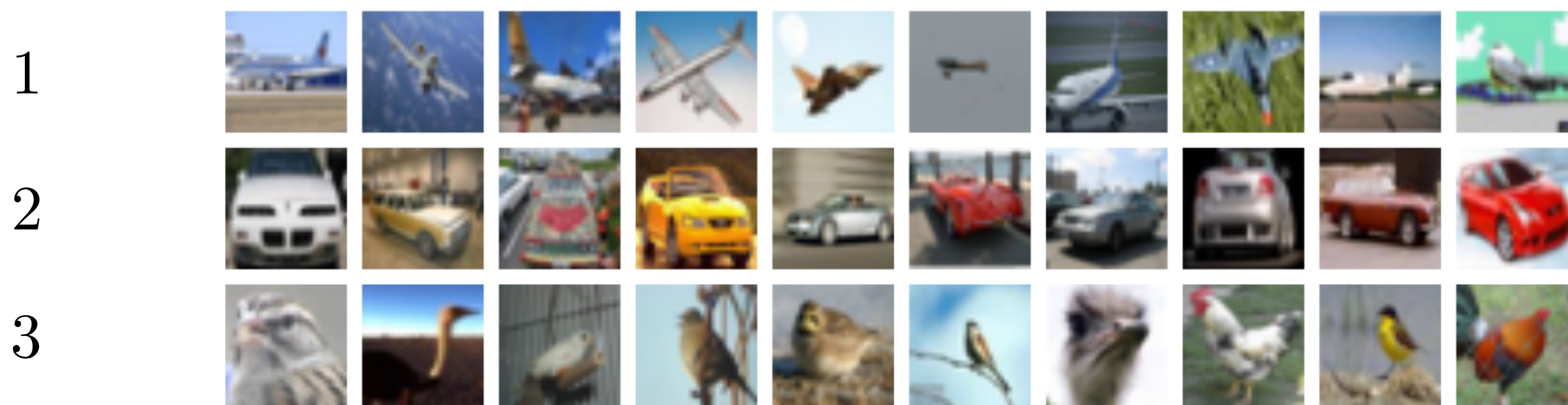
- Probability of observing  $y_i$  when measuring  $\mathbf{x}_i$  is

$$p(y_i|\mathbf{x}_i, W) = \mathbf{s}_{y_i}(\mathbf{f}(\mathbf{x}_i, W))$$

- **Training:** MLE estimate of  $W$

$$W^* = \arg \min_W \sum_i -\log \mathbf{s}_{y_i}(\mathbf{f}(\mathbf{x}_i, W))$$





- **Classification** (probability modeled by soft-max function):

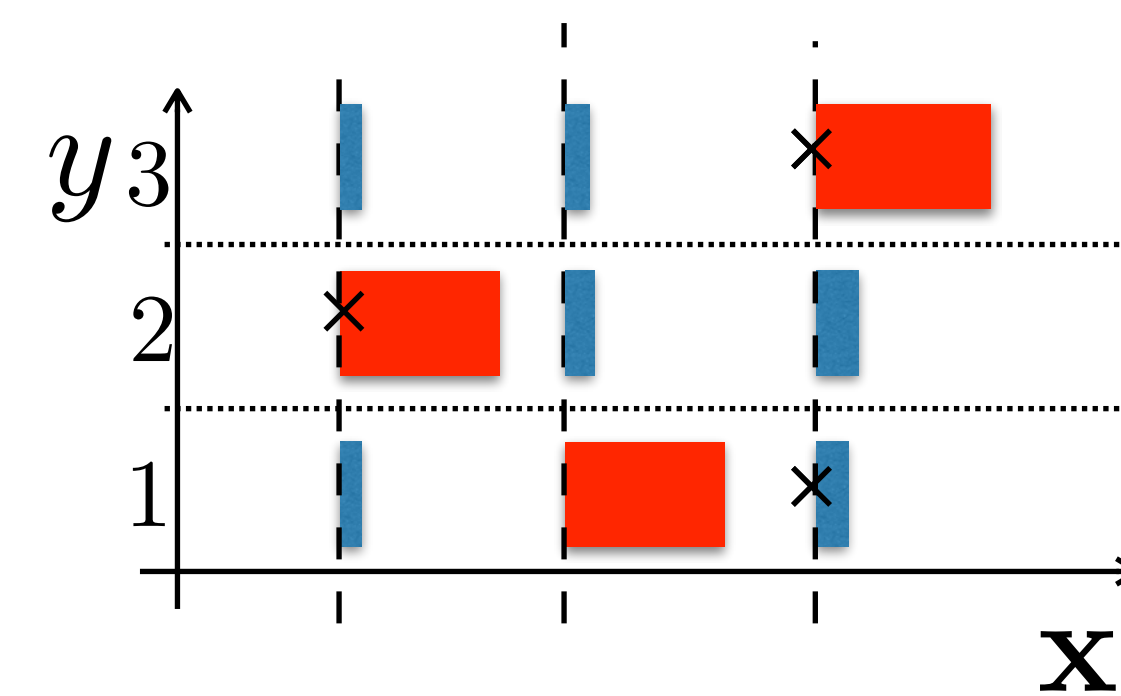
$$p(y|\mathbf{x}, W) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k))} = \mathbf{s}(\mathbf{f}(\mathbf{x}, W))$$

- Probability of observing  $y_i$  when measuring  $\mathbf{x}_i$  is

$$p(y_i|\mathbf{x}_i, W) = \mathbf{s}_{y_i}(\mathbf{f}(\mathbf{x}_i, W))$$

- **Training:** MLE estimate of  $W$

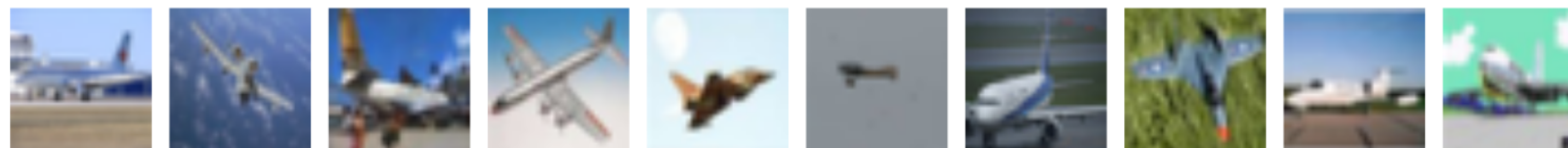
$$W^* = \arg \min_W \sum_i -\log \mathbf{s}_{y_i}(\mathbf{f}(\mathbf{x}_i, W))$$



Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

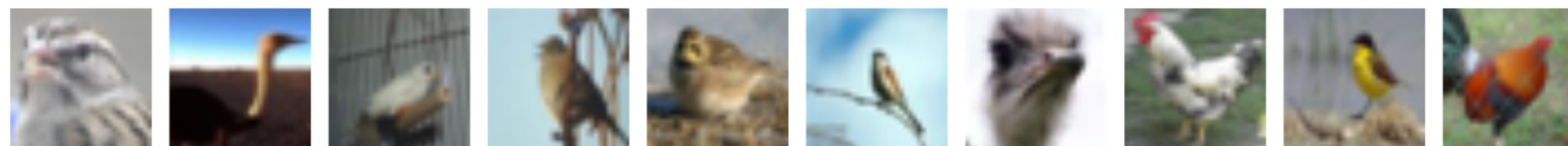
1



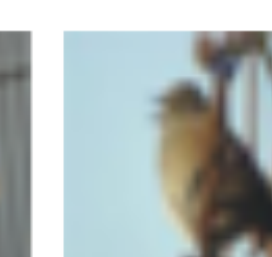
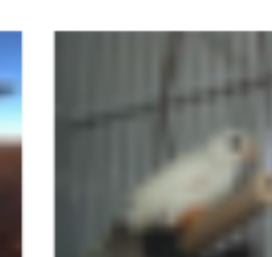
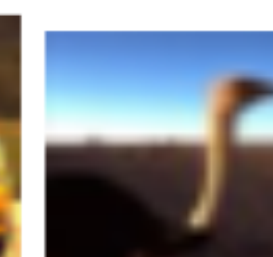
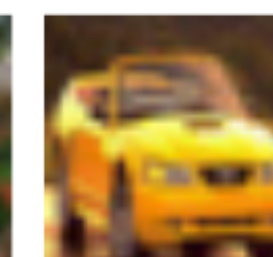
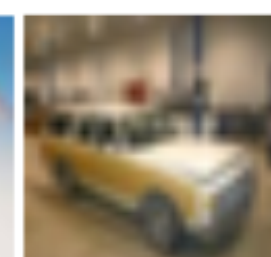
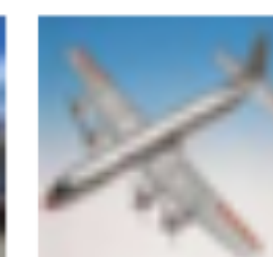
2



3



```
def train(
```



1

1

1

2

2

2

3

3

3

):

```
 $\mathbf{x}_i = \text{vec}(\text{img})$ 
```

```
 $W^* = \arg \min_W \sum_i -\log s_{y_i}(W \bar{\mathbf{x}}_i)$ 
```

```
return  $W^*$ 
```



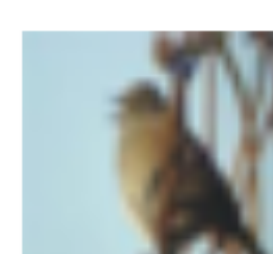
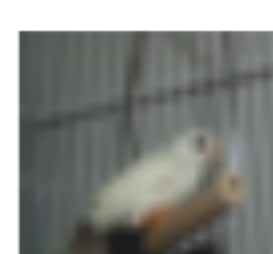
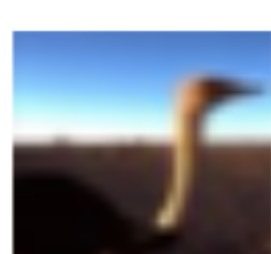
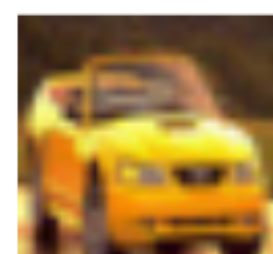
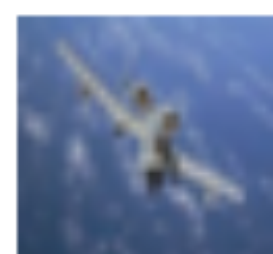
label  
 $y_i = 2$

$$\mathbf{s}(W \bar{\mathbf{x}}_i) = \begin{bmatrix} 0.03 \\ 0.71 \\ 0.26 \end{bmatrix}$$

Car classified as car yields small loss and gradient

$$\Rightarrow -\log \mathbf{s}_{y_i}(W \bar{\mathbf{x}}_i) = -\log(0.71) = 0.15$$

def train(



1

1

1

2

2

2

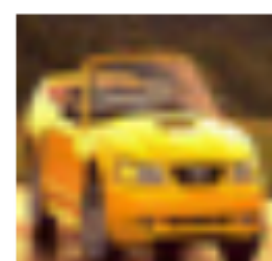
3

3

3

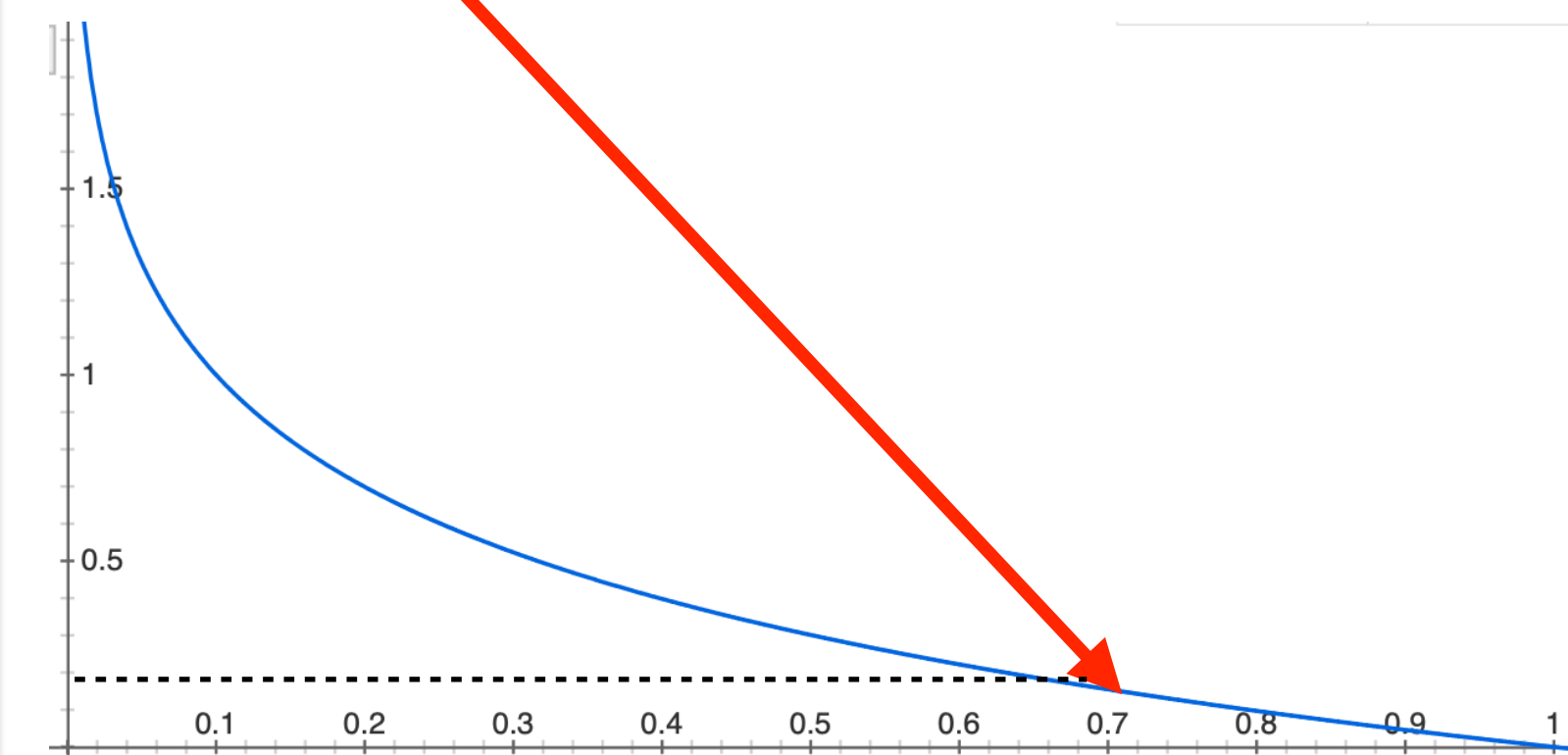
):

$$\mathbf{x}_i = \text{vec}(\text{img})$$



$$W^* = \arg \min_W \sum_i -\log \mathbf{s}_{y_i}(W \bar{\mathbf{x}}_i)$$

return  $W^*$





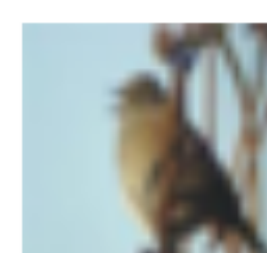
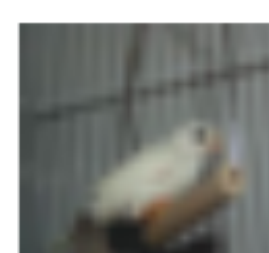
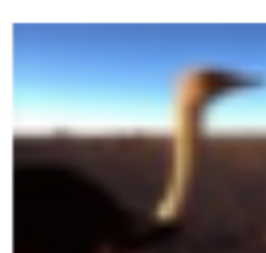
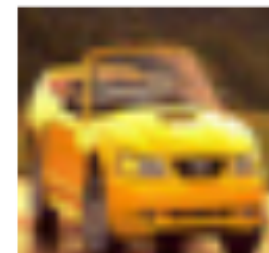
label  
 $y_i = 1$

$$\mathbf{s}(W \bar{\mathbf{x}}_i) = \begin{bmatrix} 0.03 \\ 0.57 \\ 0.40 \end{bmatrix}$$

Plane classified as car yields huge loss and gradient

$$\Rightarrow -\log \mathbf{s}_{y_i}(W \bar{\mathbf{x}}_i) = -\log(0.03) = 1.52$$

def train(



1

1

1

2

2

2

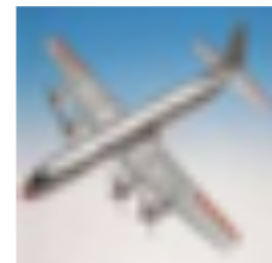
3

3

3

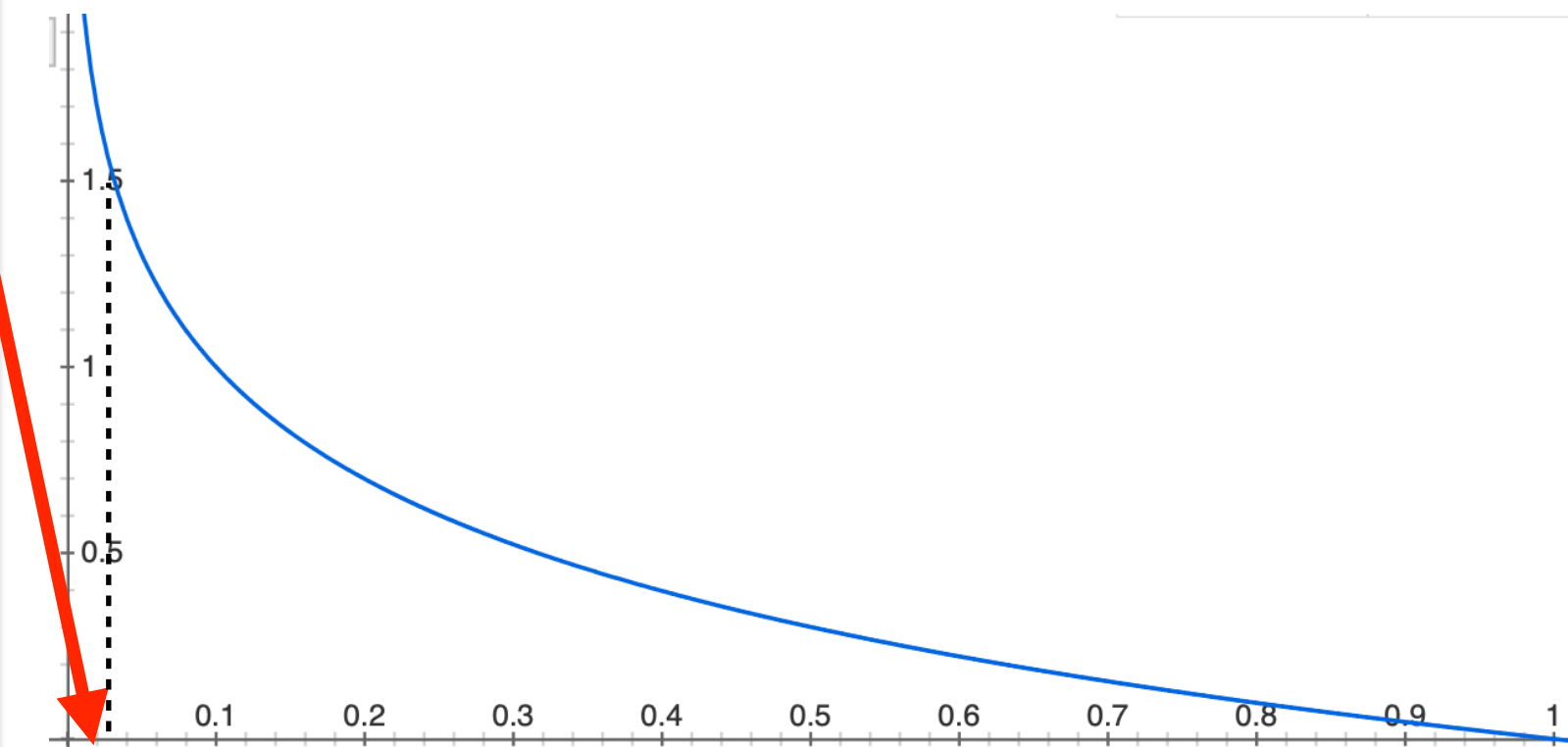
):

$$\mathbf{x}_i = \text{vec}(\text{img})$$



$$W^* = \arg \min_W \sum_i -\log \mathbf{s}_{y_i}(W \bar{\mathbf{x}}_i)$$

return  $W^*$



```
def train(
    
    
    
    
    
    
    
    
    
    ):

```

1

1

1

2

2


2

3

3

3

):




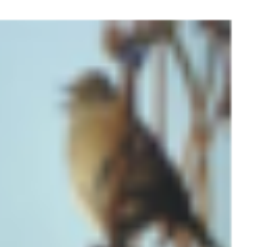
```
 $\mathbf{x}_i = \text{vec}(\text{$ )
```


$$W^* = \arg \min_W \sum_i -\log s_{y_i}(W \bar{\mathbf{x}}_i) = \sum_i \ell(W, \bar{\mathbf{x}}_i, y_i)$$

```
return  $W^*$ 
```

$$\mathcal{L}(\mathbf{w})$$

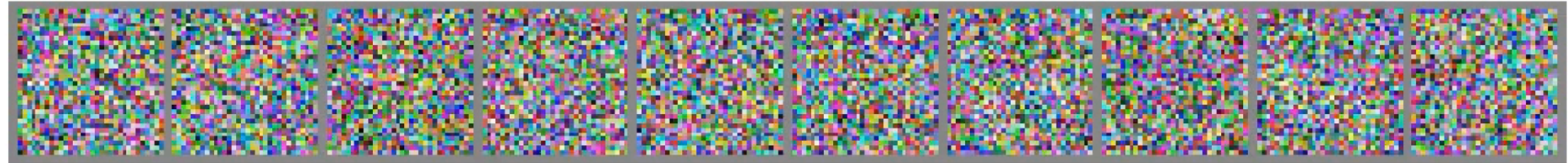








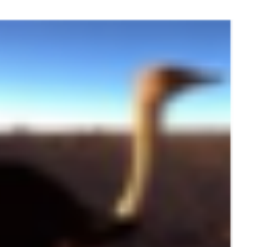
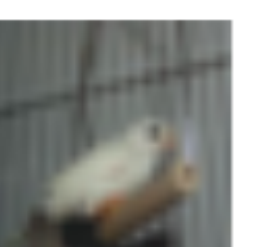
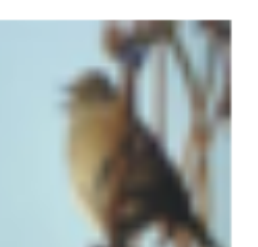

```
def train(          ):
```

$\mathbf{x}_i = \text{vec}(\text{  })$

$W^* = \arg \min_W \mathcal{L}(W)$  ... gradient optimization

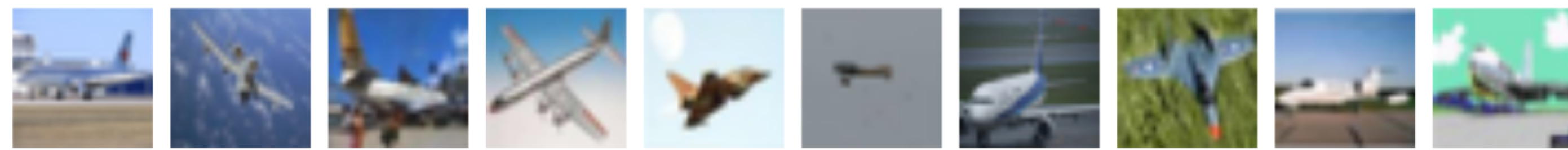
return  $W^*$



```
def train(          ):  
     $\mathbf{x}_i = \text{vec}(\text{  })$   
     $W^* = \arg \min_W \mathcal{L}(W)$  ... gradient optimization  
    return  $W^*$ 
```

3	4	2	1	9	5	6	2	1	8
8	9	1	2	5	0	0	6	6	4
6	7	0	1	6	3	6	3	7	0
3	7	7	9	4	6	6	1	8	2
2	9	3	4	3	9	8	7	2	5
1	5	9	8	3	6	5	7	2	3
9	3	1	9	1	5	8	0	8	4

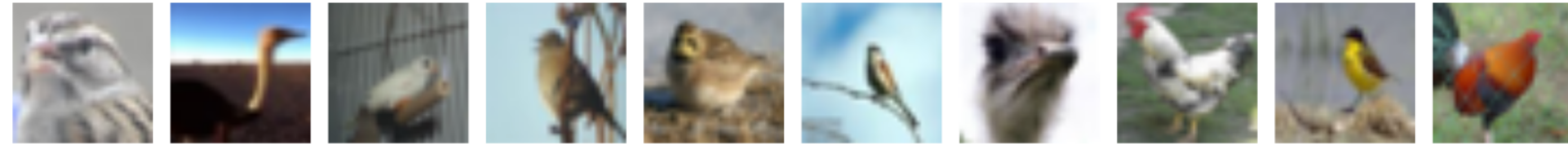
**airplane**



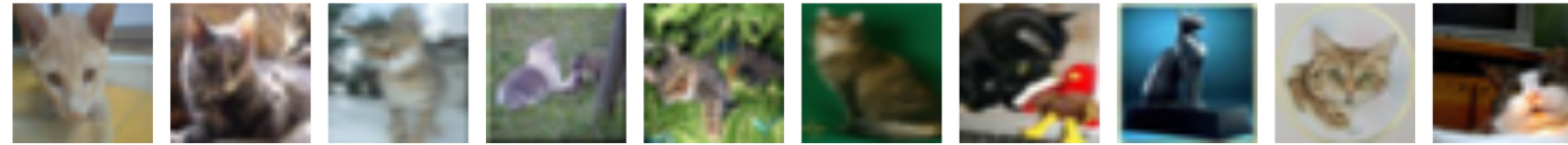
**automobile**



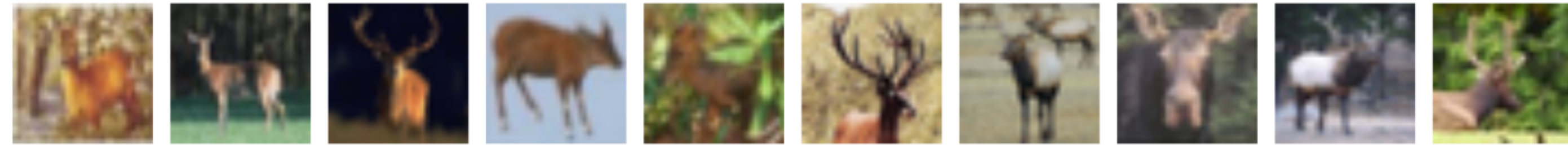
**bird**



**cat**



**deer**



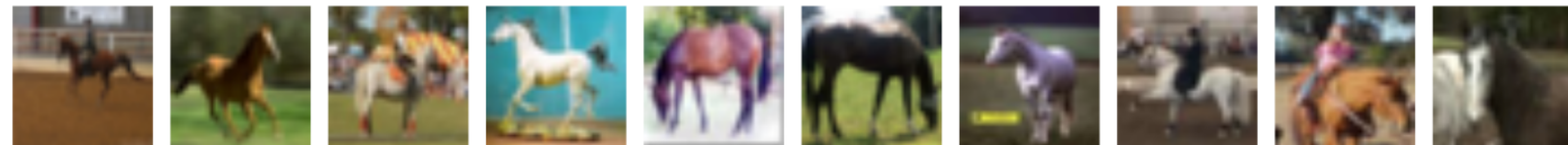
**dog**



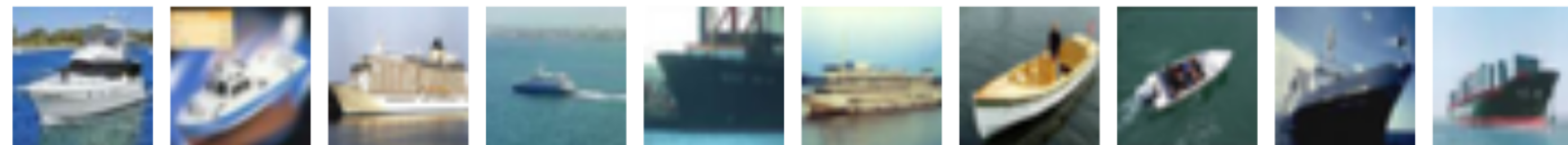
**frog**



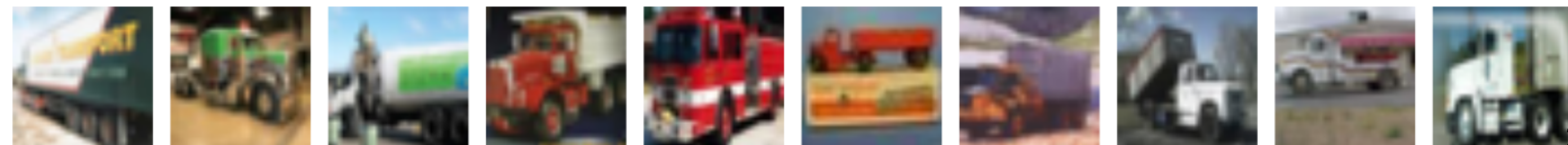
**horse**



**ship**



**truck**

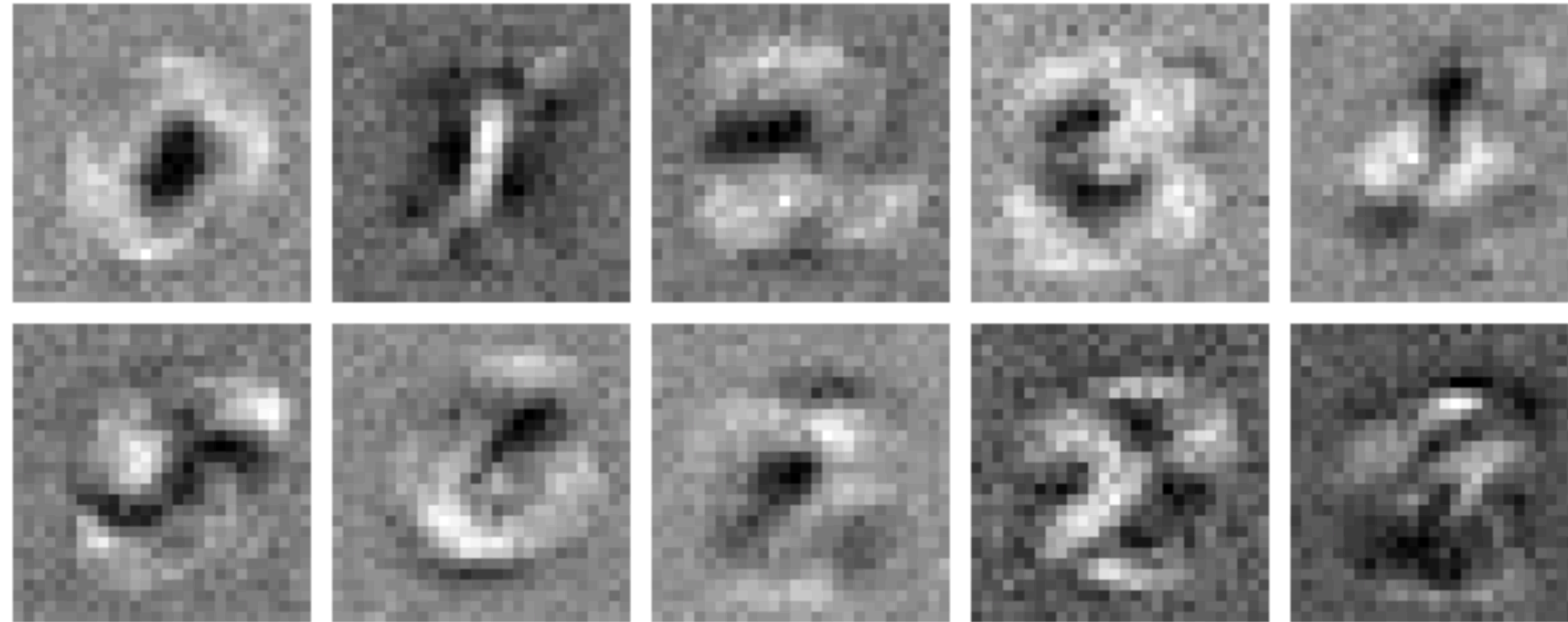


Dataset

Learned weights of linear classifier

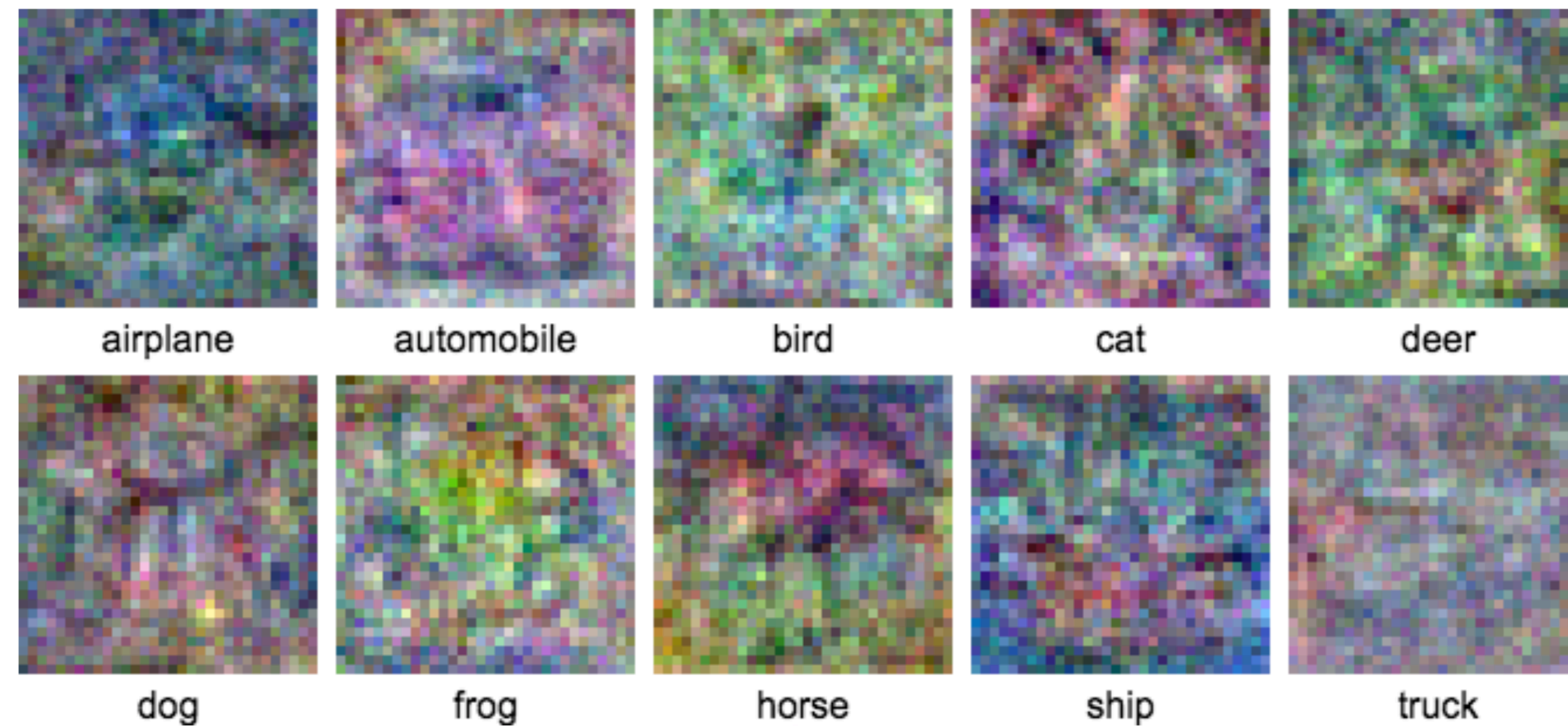
Error

MNIST



8%

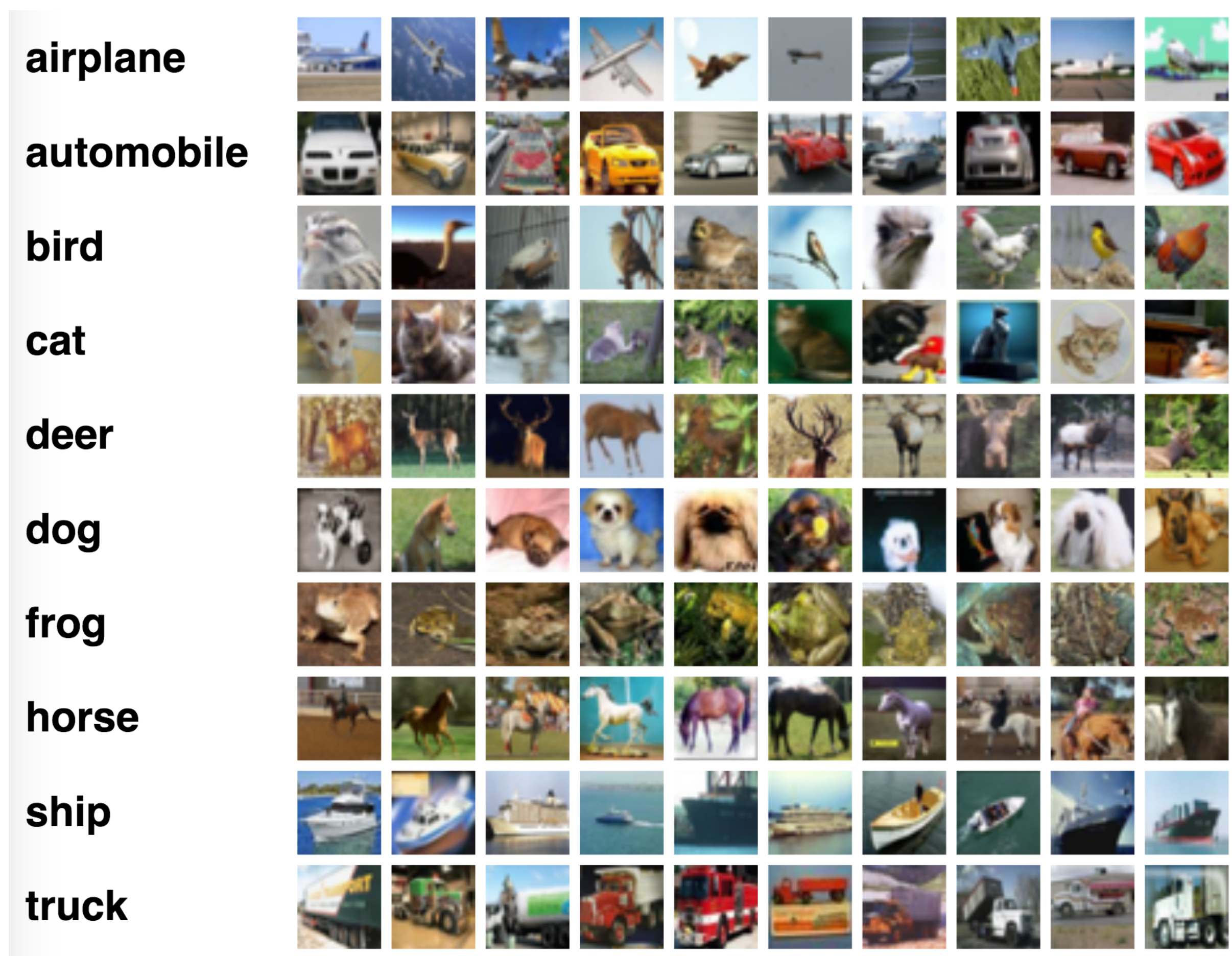
CIFAR-10



63%

<https://benchmarks.ai>

# Recognition problem



Why is it hard?

CIFAR-10: classify 32x32 RGB images into 10 categories  
<https://www.cs.toronto.edu/~kriz/cifar.html>

# Recognition problem

## Huge within-class variability!

Why it is hard?

- Viewpoint
- Occlusion
- Illumination
- Pose
- Type
- Context



Why is it hard?

## Huge (i) within-class variability

bird



cat



dog





Why is it hard?

**Huge (i) within-class variability (ii) among class similarity!**



Czech Technical University in Prague  
Faculty of Electrical Engineering, Department of  
Cybernetics

Why is it hard?

**Huge (i) within-class variability (ii) among class similarity!**



Czech Technical University in Prague  
Faculty of Electrical Engineering, Department of  
Cybernetics

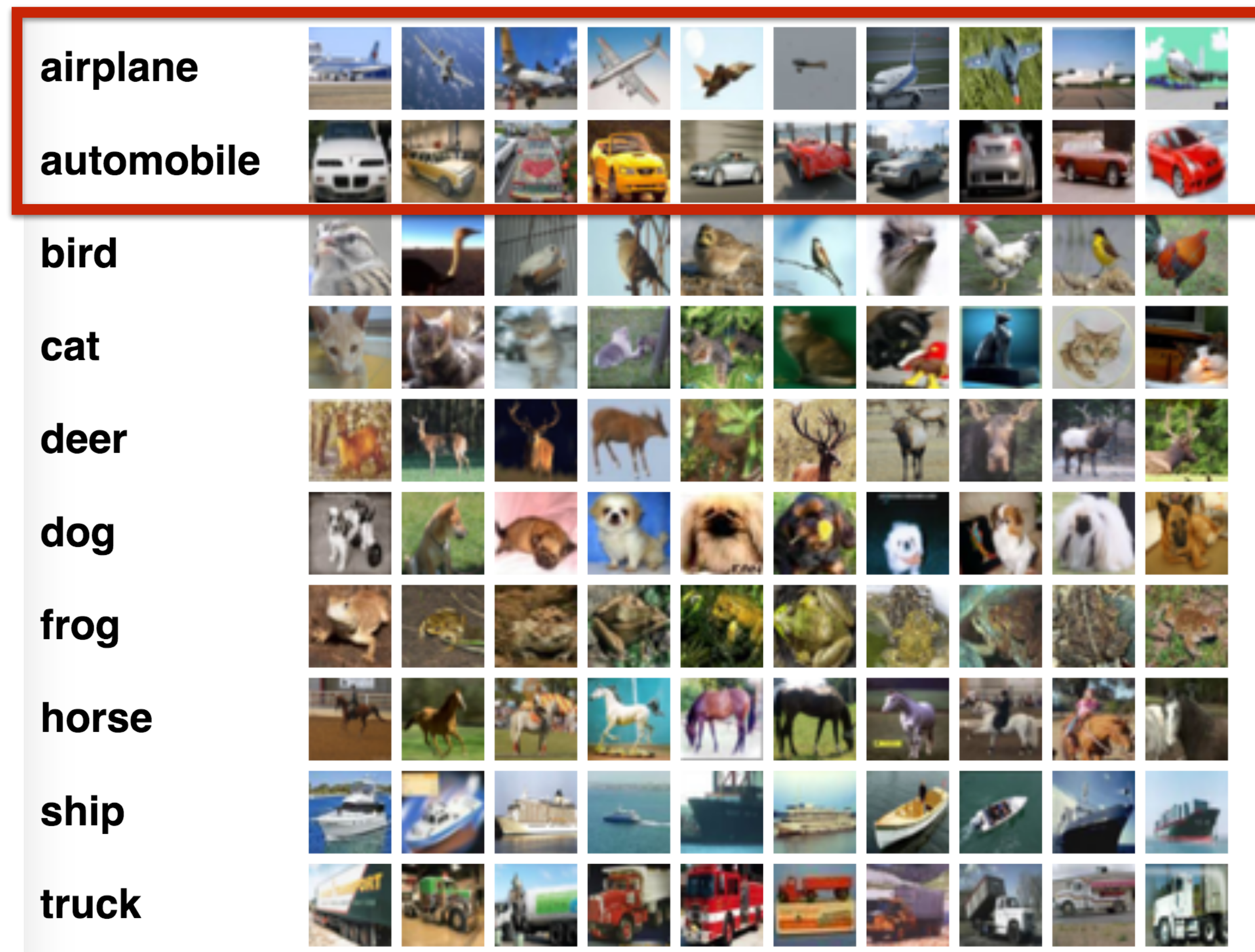
Why is it hard?

**Huge (i) within-class variability (ii) among class similarity!**



Czech Technical University in Prague  
Faculty of Electrical Engineering, Department of  
Cybernetics

# Recognition problem



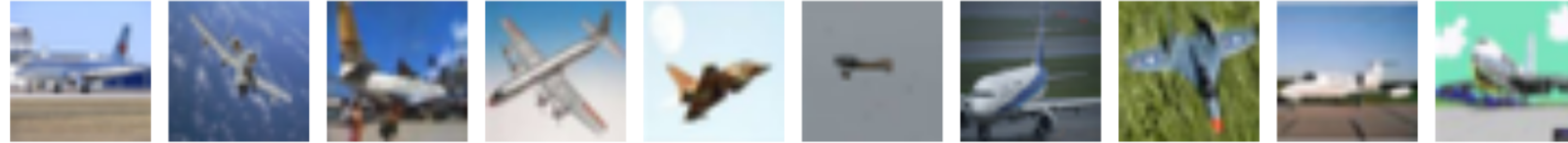
CIFAR-10: classify 32x32 RGB images into 10 categories

<https://www.cs.toronto.edu/~kriz/cifar.html>

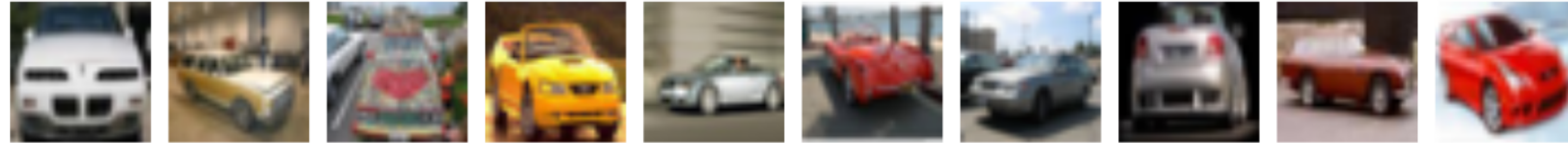
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

**airplane**



**automobile**



Two-class recognition problem: classify airplane/automobile

```
def classify():
```

```
    ???
```

```
    return p
```

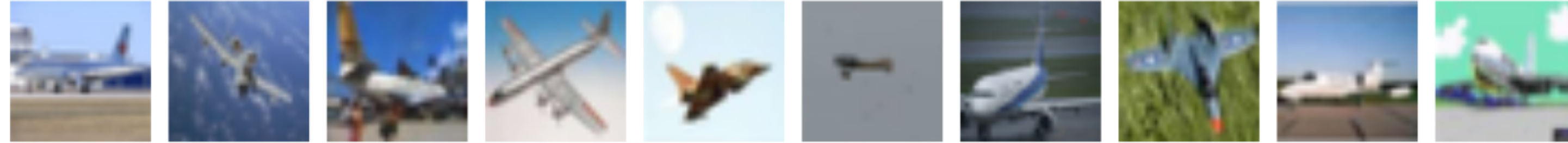
Probability of image being from the class airplane

How to model it?

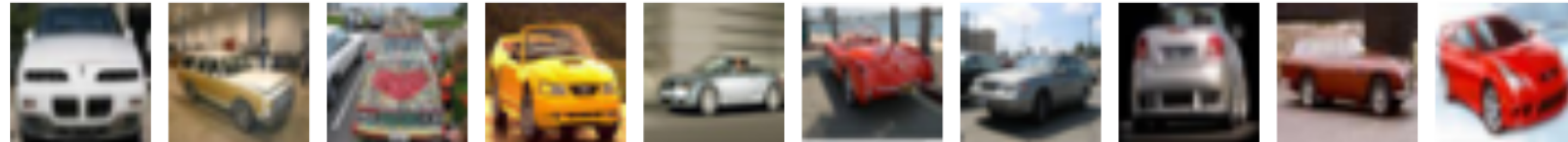
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



-1



## Classification

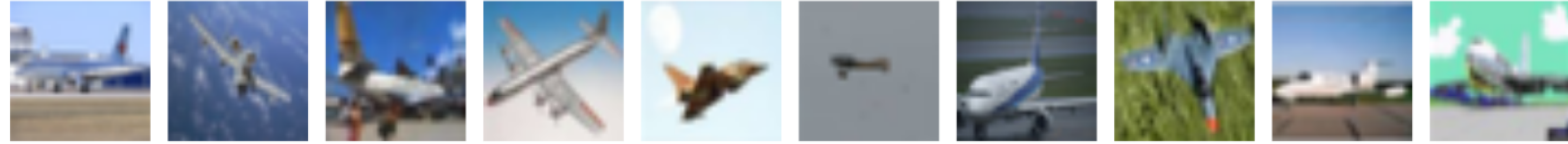
We model probability of image  $\mathbf{x}$  being label +1 or -1 as

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$

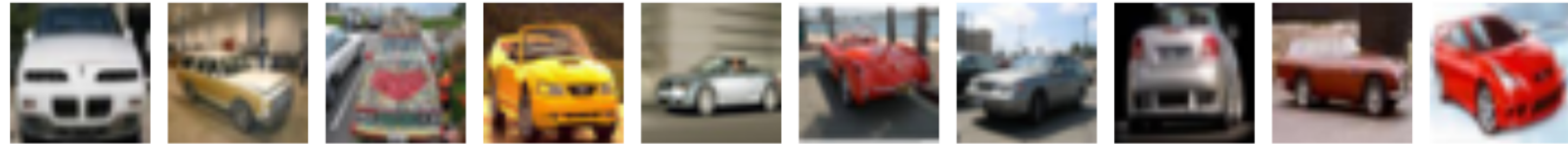
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



-1



## Classification

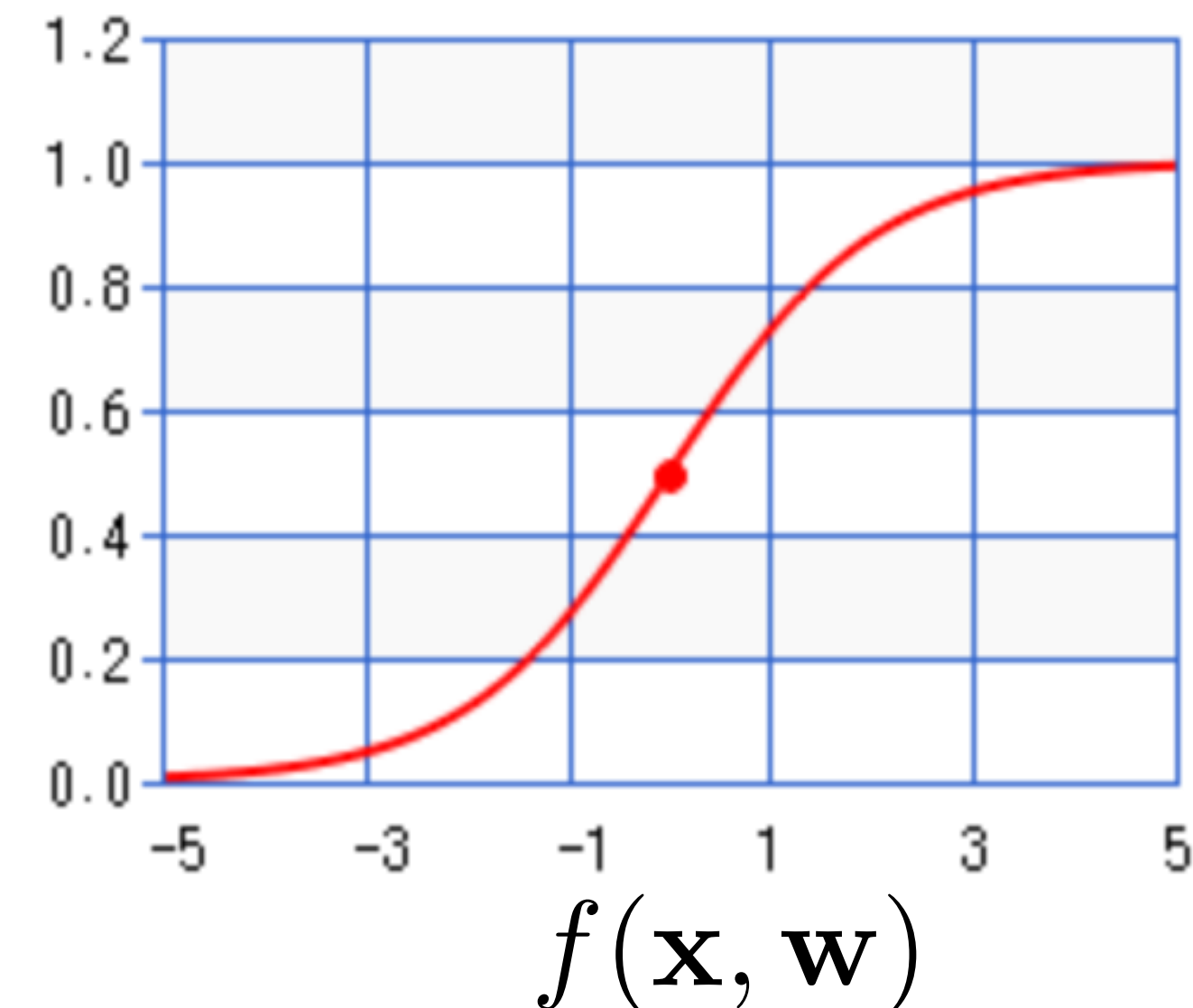
We model probability of image  $\mathbf{x}$  being label +1 or -1 as

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$

where

$$\sigma(f(\mathbf{x}, \mathbf{w})) = \frac{1}{1 + \exp(-f(\mathbf{x}, \mathbf{w}))}$$

is sigmoid function.

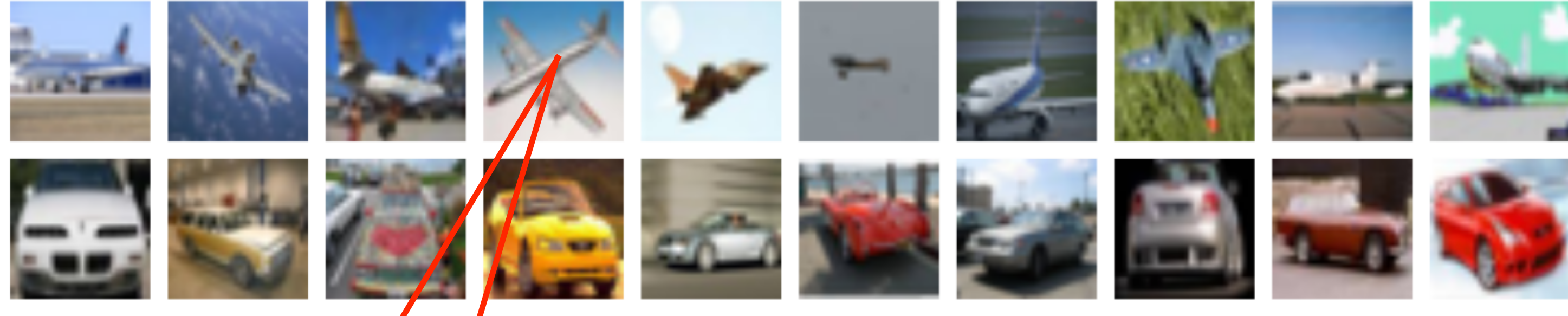


Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1

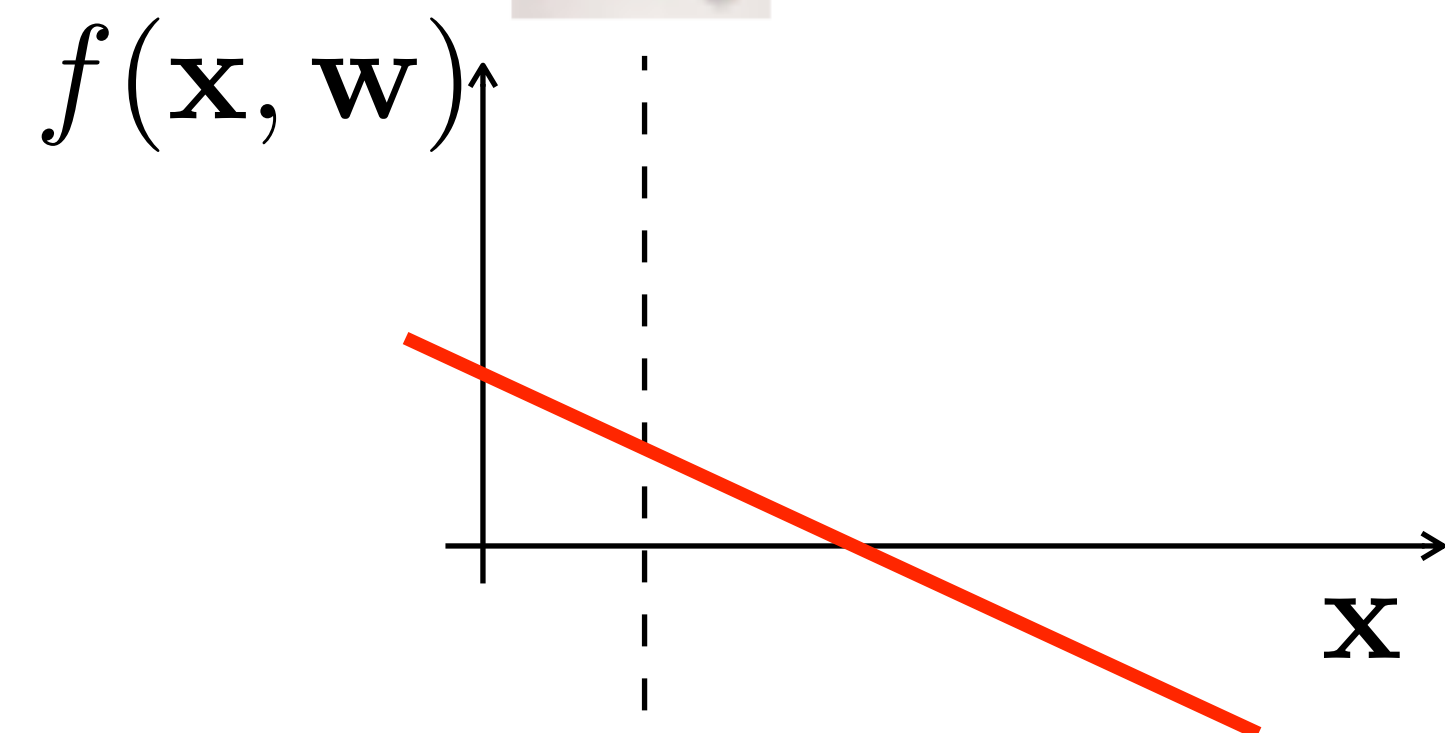
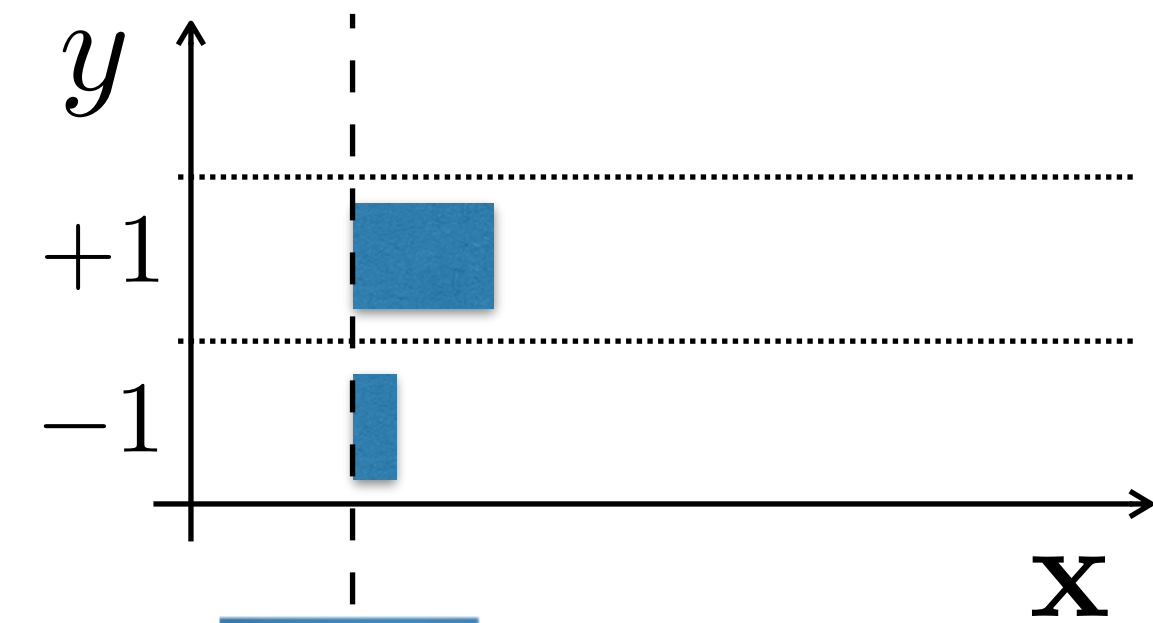
-1



### Classification

We model probability of image  $\mathbf{x}$  being label +1 or -1 as

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$

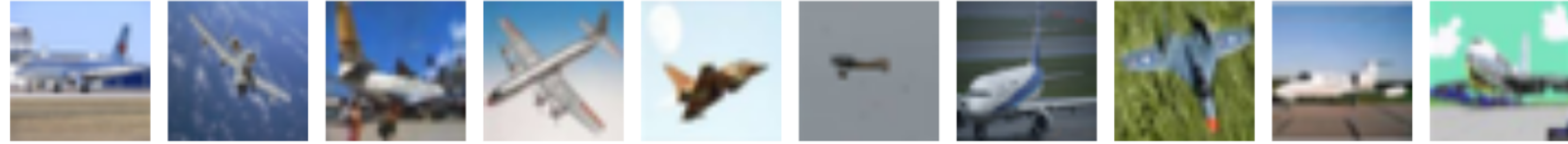




Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



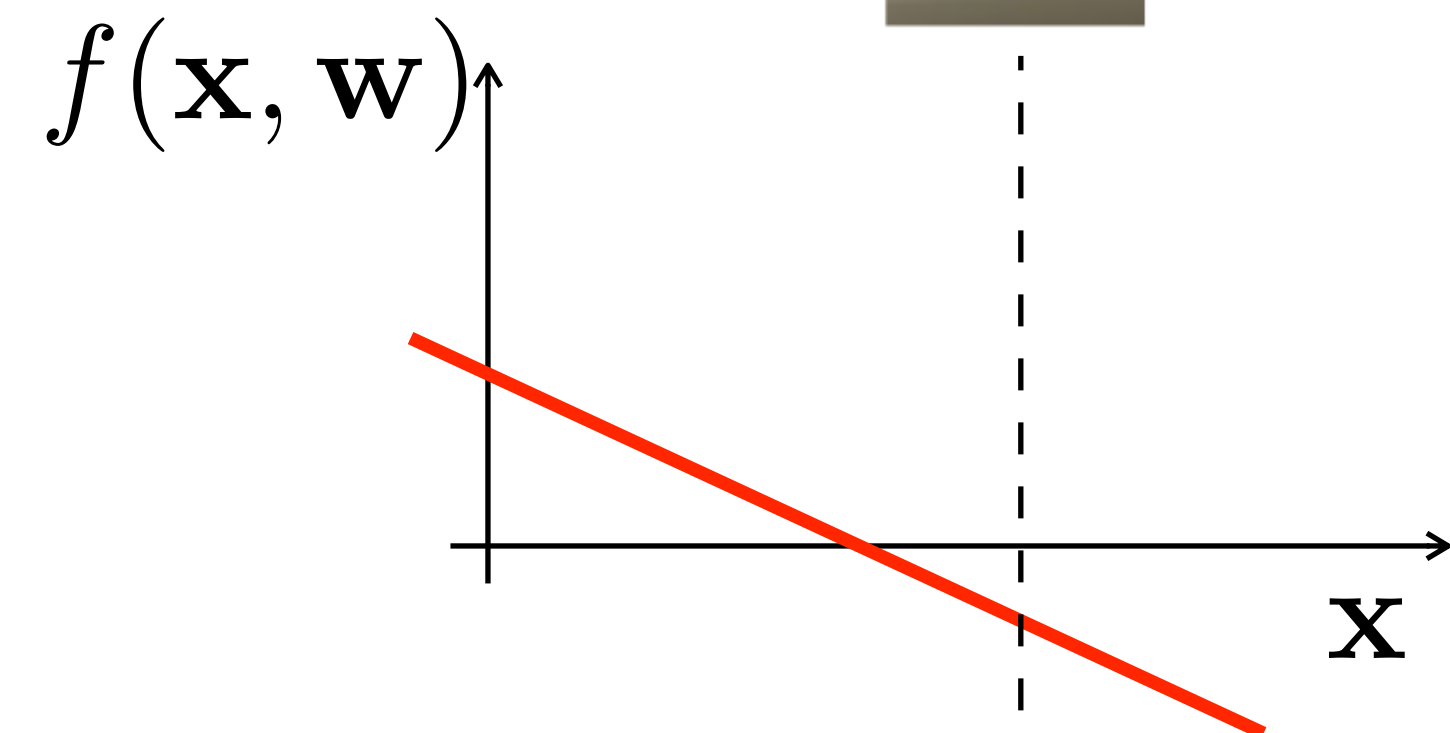
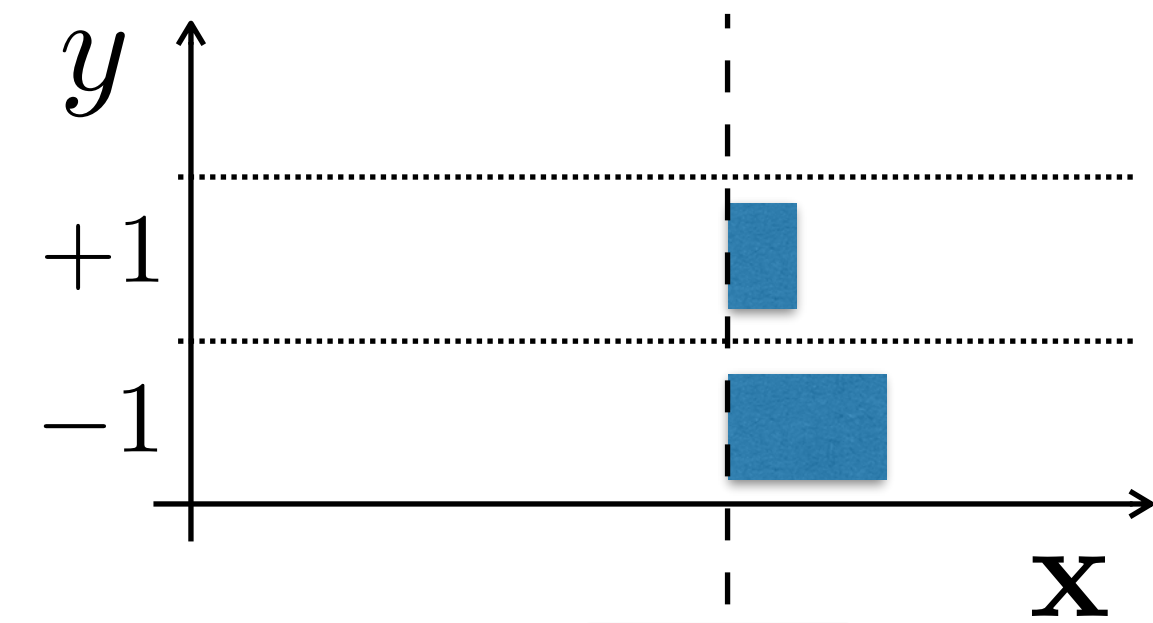
-1



### Classification

We model probability of image  $\mathbf{x}$  being label +1 or -1 as

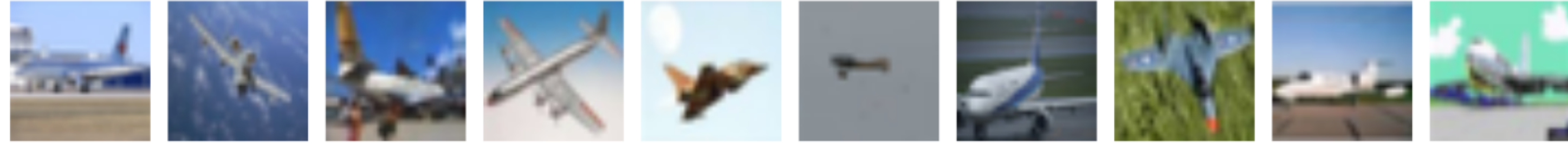
$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



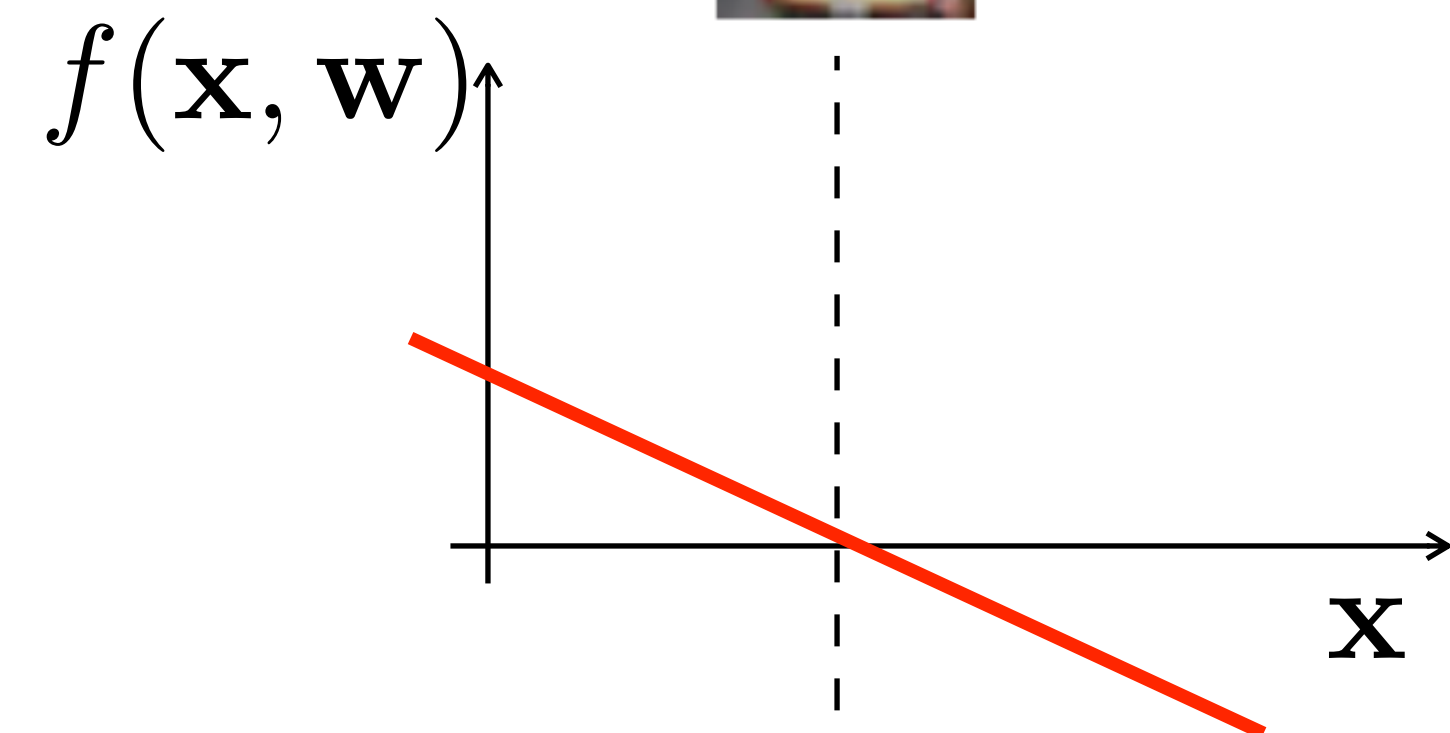
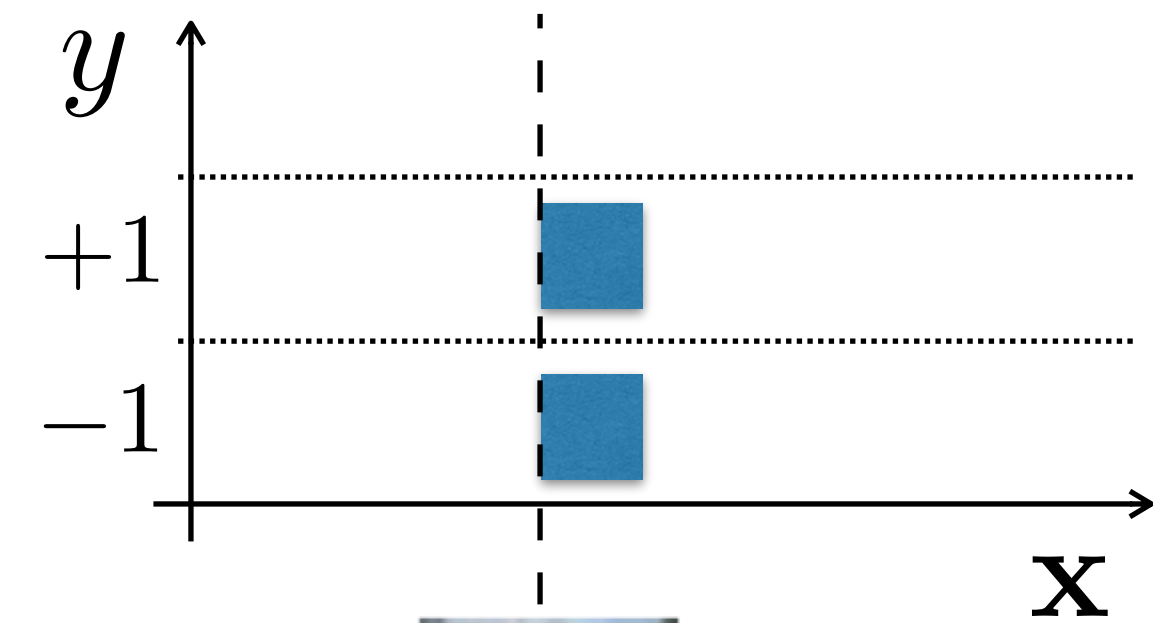
-1



### Classification

We model probability of image  $\mathbf{x}$  being label +1 or -1 as

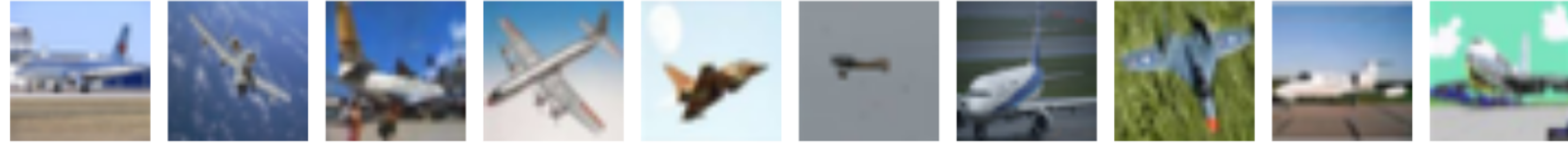
$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



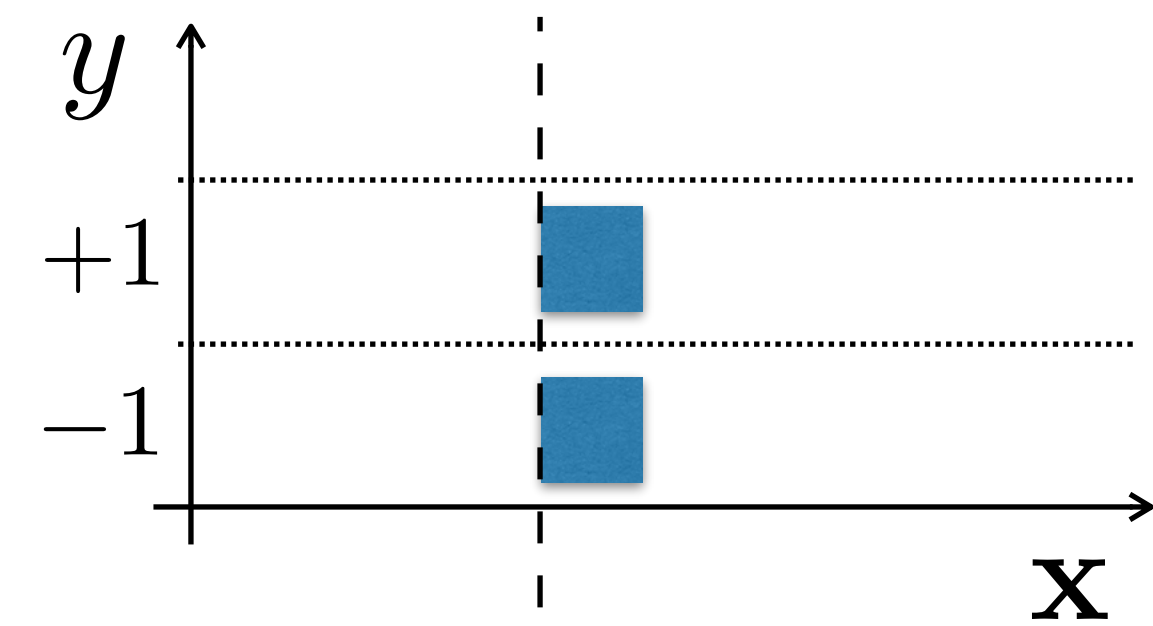
-1



## Classification

We model probability of image  $\mathbf{x}$  being label +1 or -1 as

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



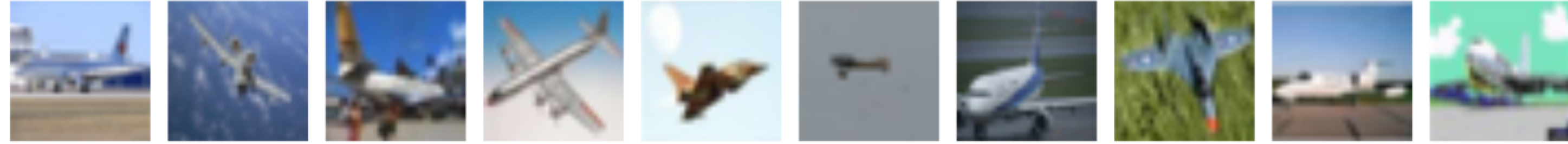
Linear classifier model probability of being from class +1 as  $p = \sigma(\mathbf{w}^\top \bar{\mathbf{x}})$

What is dimensionality of  $\mathbf{x}$  and  $\mathbf{w}$ ?

Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



-1



## Classification

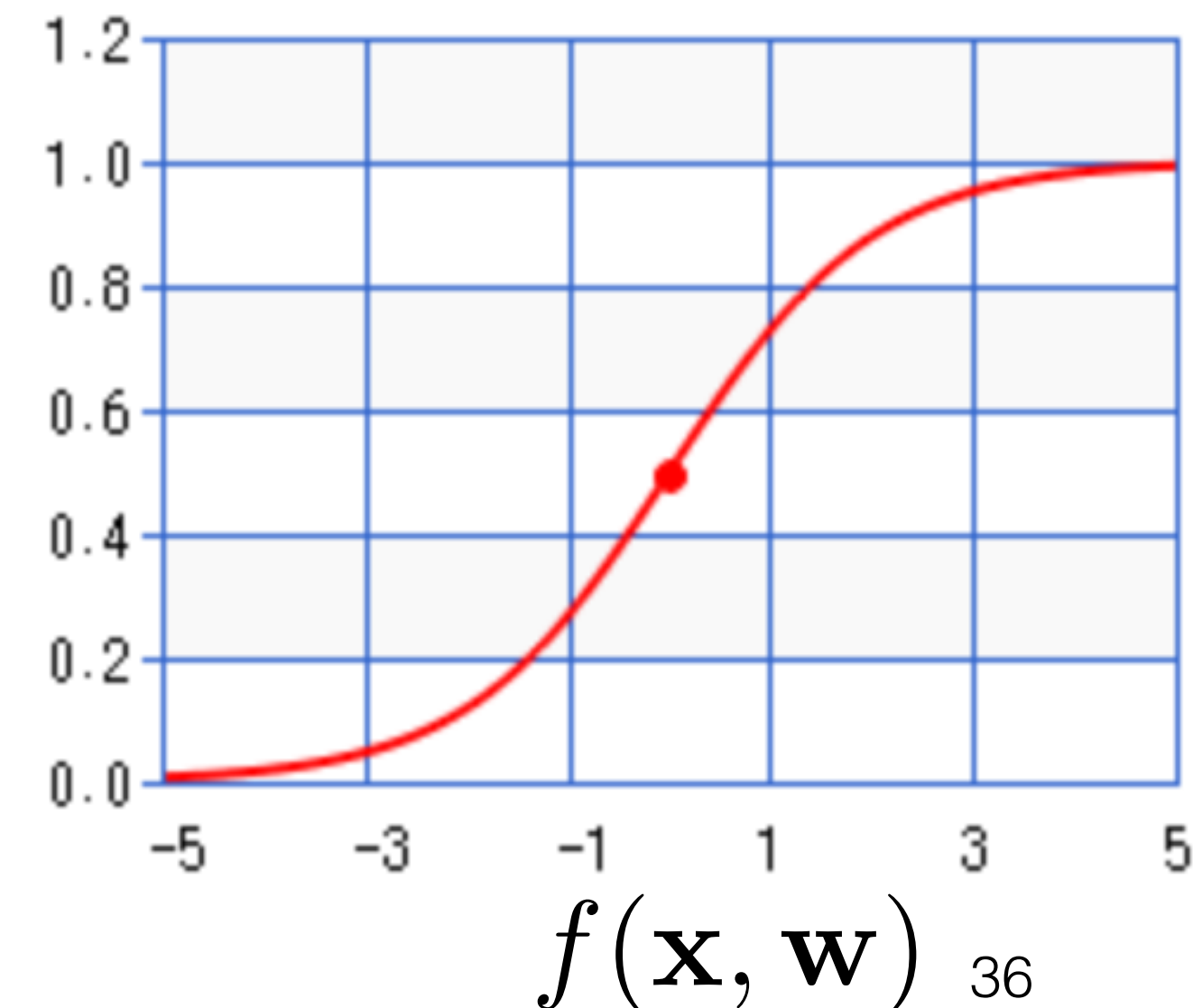
We model probability of image  $\mathbf{x}$  being label +1 or -1 as

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$

where

$$\sigma(f(\mathbf{x}, \mathbf{w})) = \frac{1}{1 + \exp(-f(\mathbf{x}, \mathbf{w}))}$$

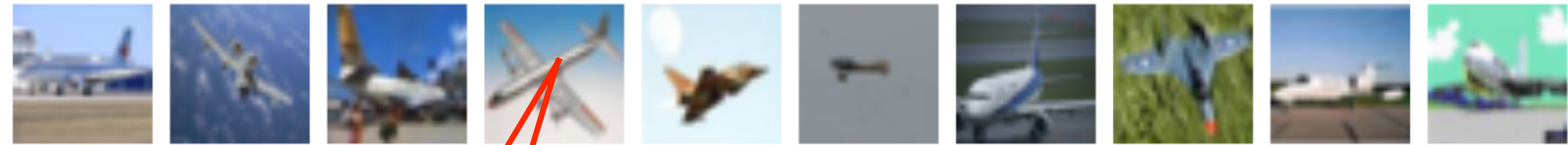
is sigmoid function.



Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



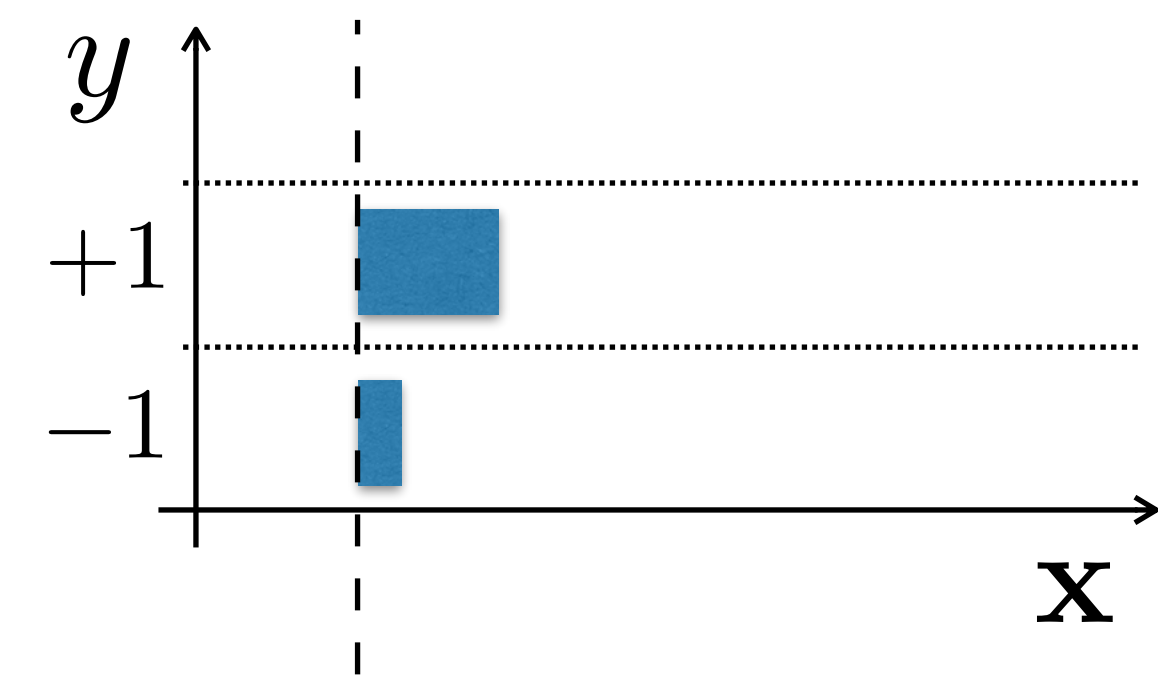
-1



### Classification

We model probability of image  $\mathbf{x}$  being label +1 or -1 as

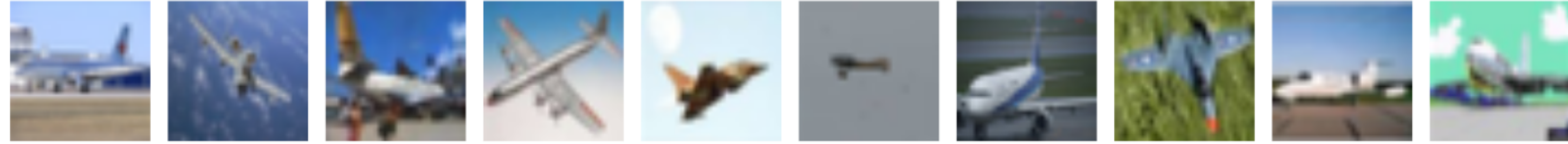
$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



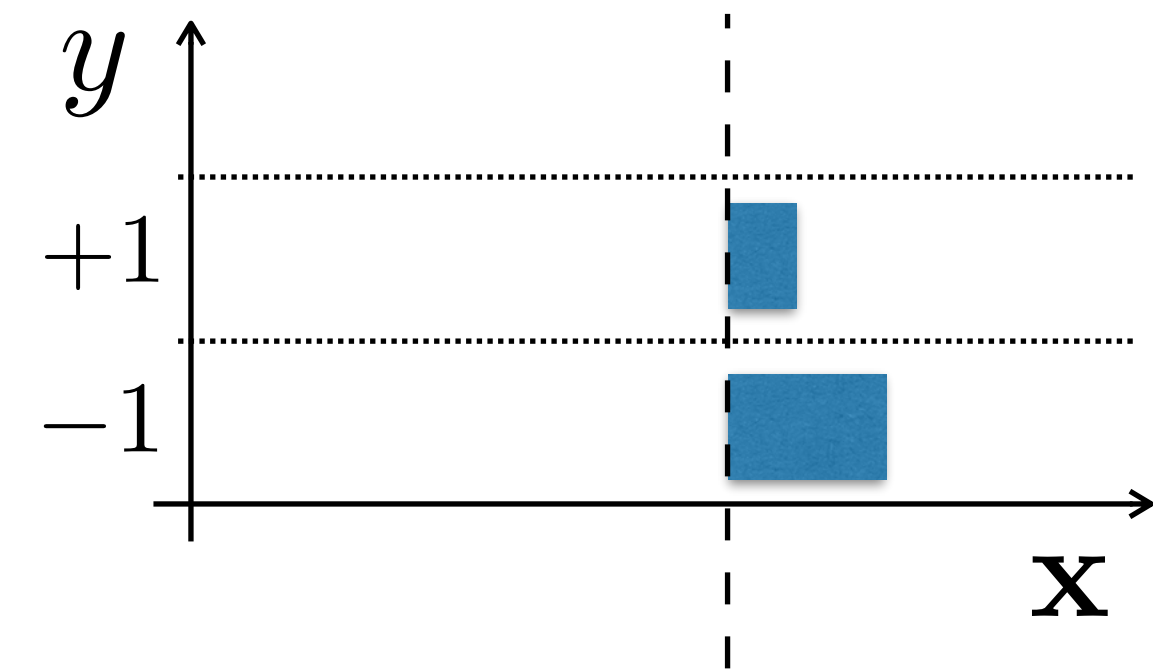
-1



### Classification

We model probability of image  $\mathbf{x}$  being label +1 or -1 as

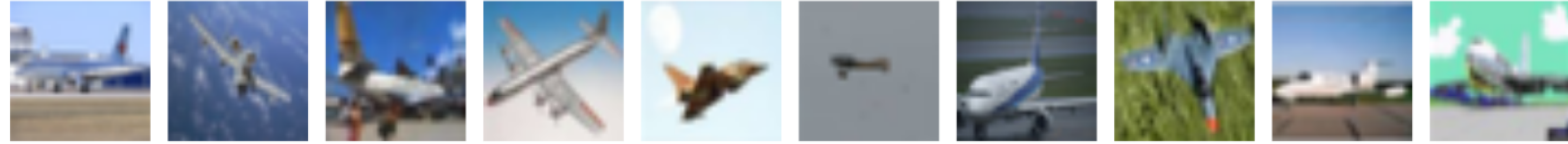
$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



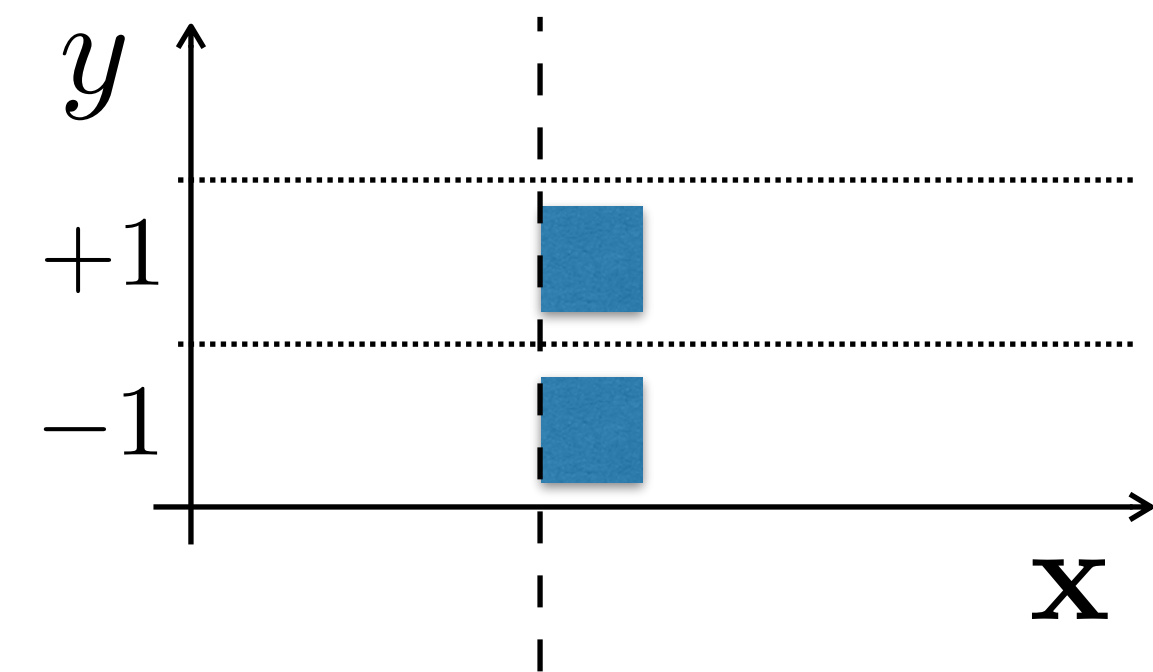
-1



### Classification

We model probability of image  $\mathbf{x}$  being label +1 or -1 as

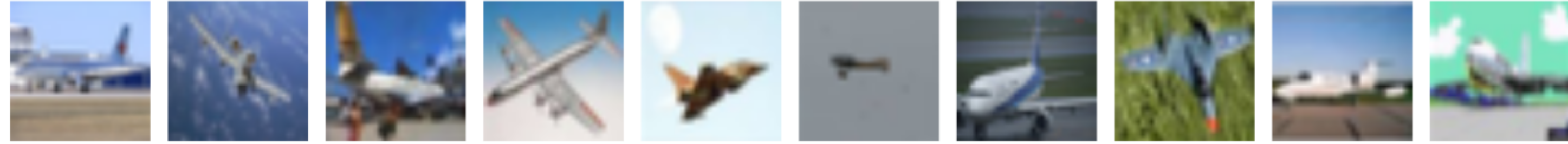
$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



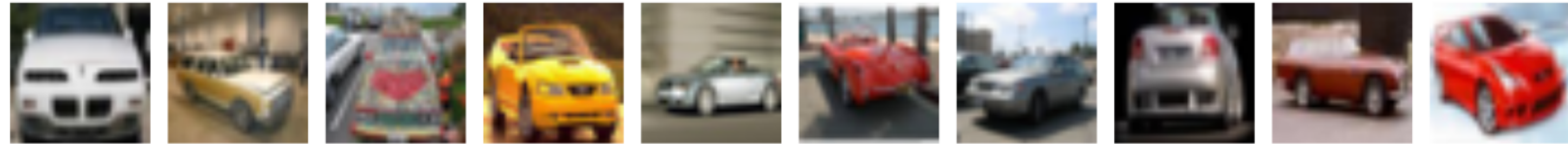
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



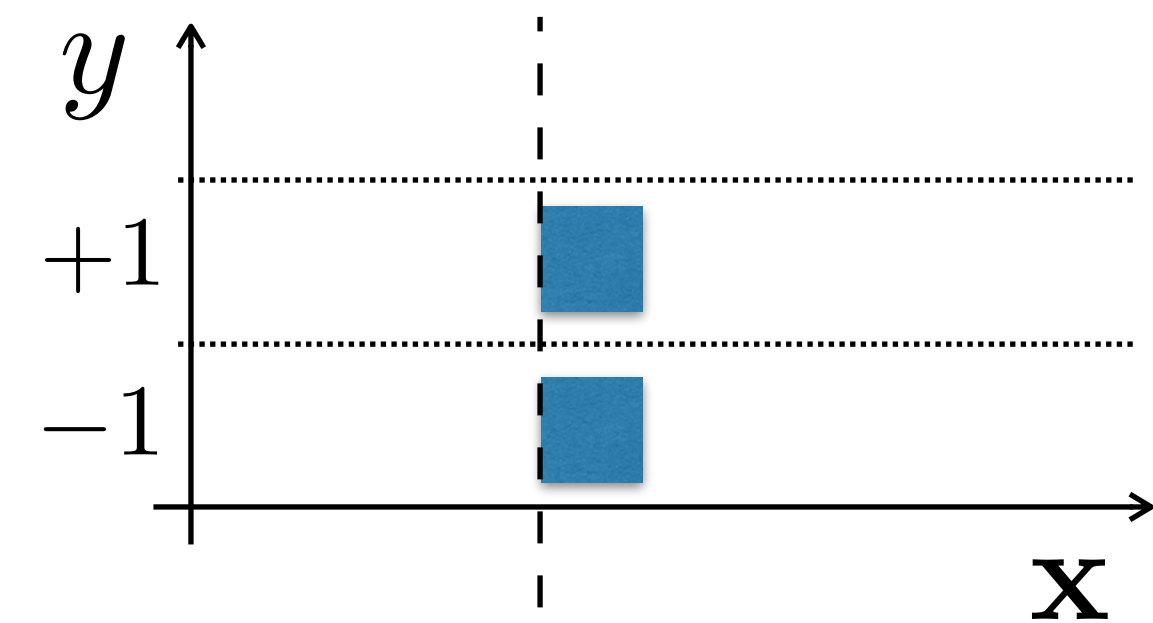
-1



## Classification

We model probability of image  $\mathbf{x}$  being label +1 or -1 as

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



Linear classifier models probability of image being from class +1 as  $p = \sigma(\mathbf{w}^\top \bar{\mathbf{x}})$



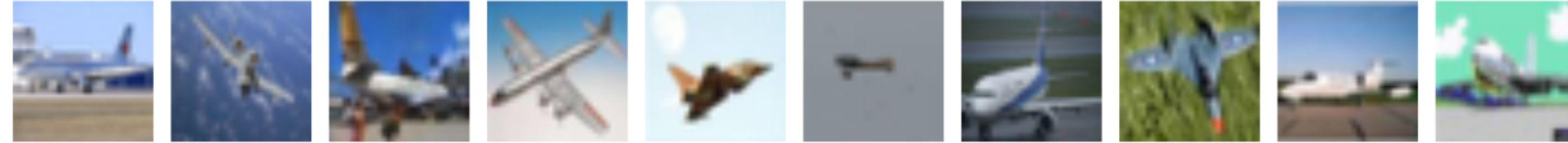
What is dimensionality of  $\mathbf{x}$  and  $\mathbf{w}$ ?



Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



-1



## Classification

Example: Linear classifier

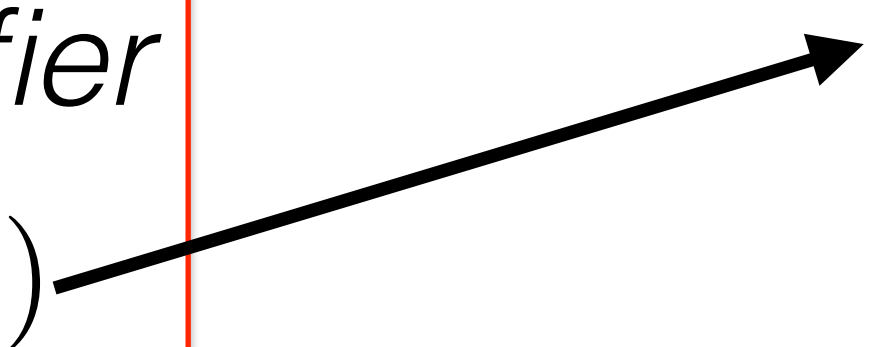
```
def classify():
```

```
    # Linear classifier
```

```
     $\mathbf{x} = \text{vec}(\img alt="small car image" data-bbox="285 601 336 691"/>$ )
```

```
     $p = \sigma(\mathbf{w}^\top \bar{\mathbf{x}})$ 
```

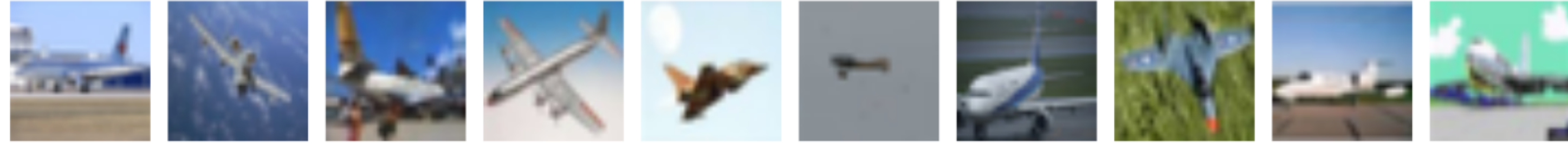
```
    return  $p$ 
```

$$\mathbf{w}^\top \bar{\mathbf{x}} = 2.5 \text{ score}$$


Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1





-1



## Classification

Example: Linear classifier

```
def classify():  
    # Linear classifier  
     $\mathbf{x} = \text{vec}(\text{$ )  
     $p = \sigma(\mathbf{w}^\top \bar{\mathbf{x}})$   
    return  $p$ 
```

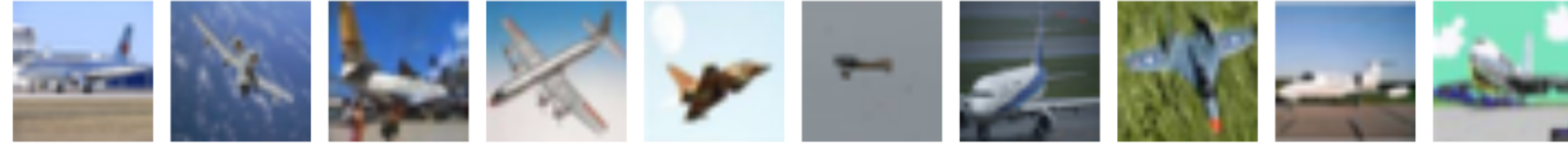
$$\mathbf{w}^\top \bar{\mathbf{x}} = 2.5 \text{ score}$$



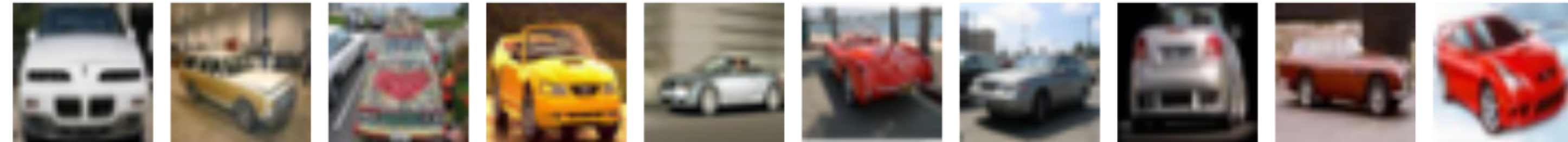
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1




-1

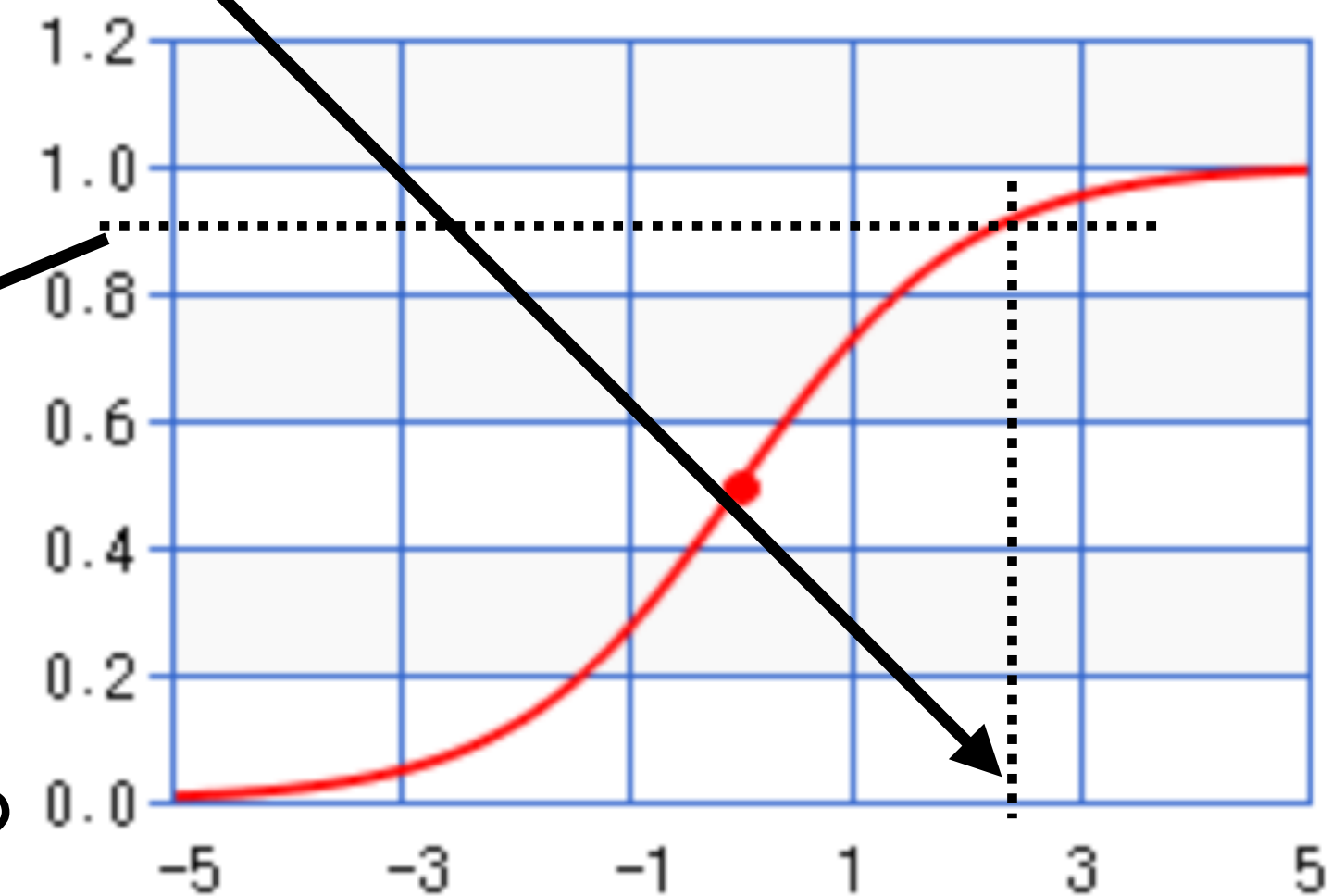


## Classification

Example: Linear classifier

```
def classify(  
     $p = \sigma(\mathbf{w}^\top \bar{\mathbf{x}})$   
    return  $p = \sigma(2.5) = 0.92$ 
```

$$\mathbf{w}^\top \bar{\mathbf{x}} = 2.5 \text{ score}$$



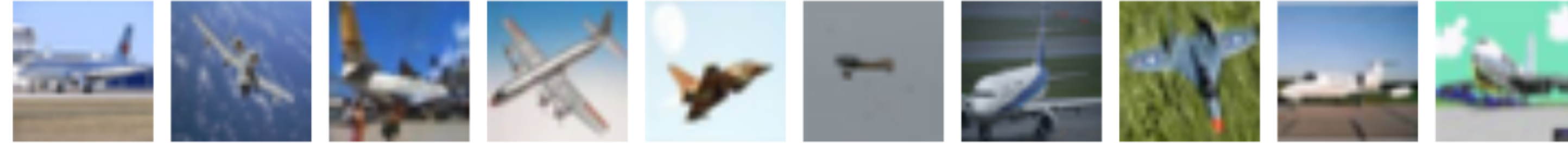
is it a good classifier?

Show python code

Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



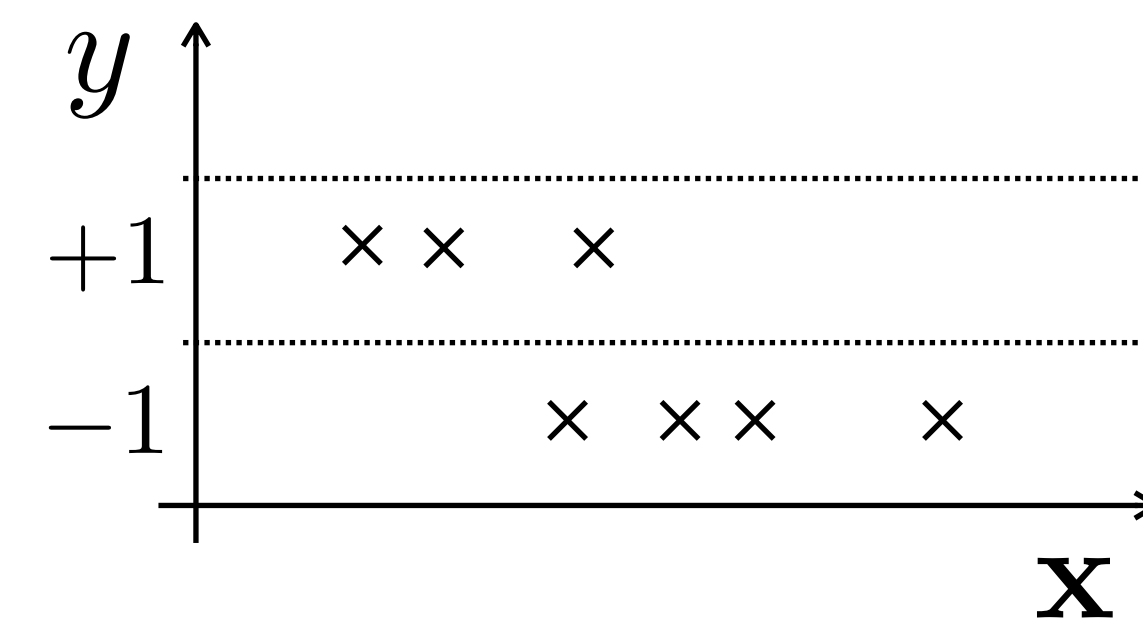
-1



### Training

Training = search for unknown parameters  $\mathbf{w}$   
which fits a given data

Training data

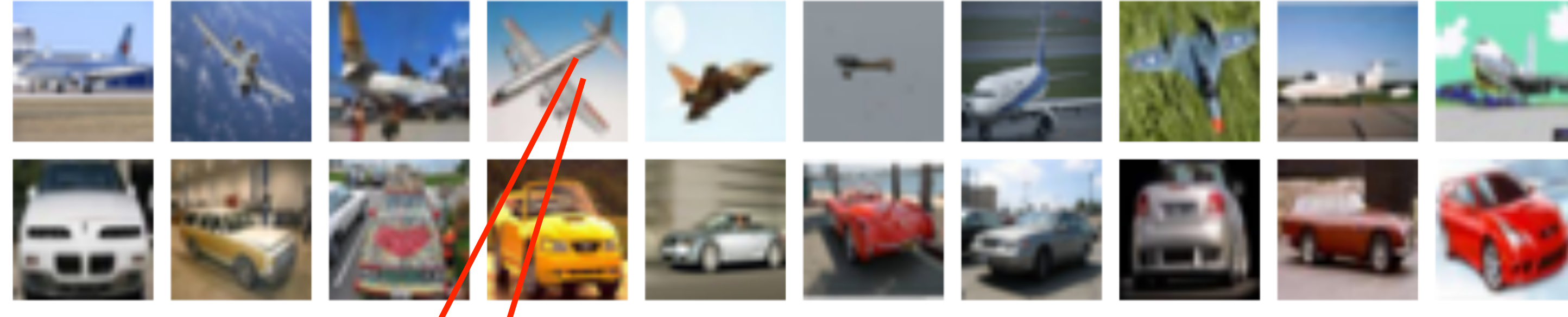


Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1

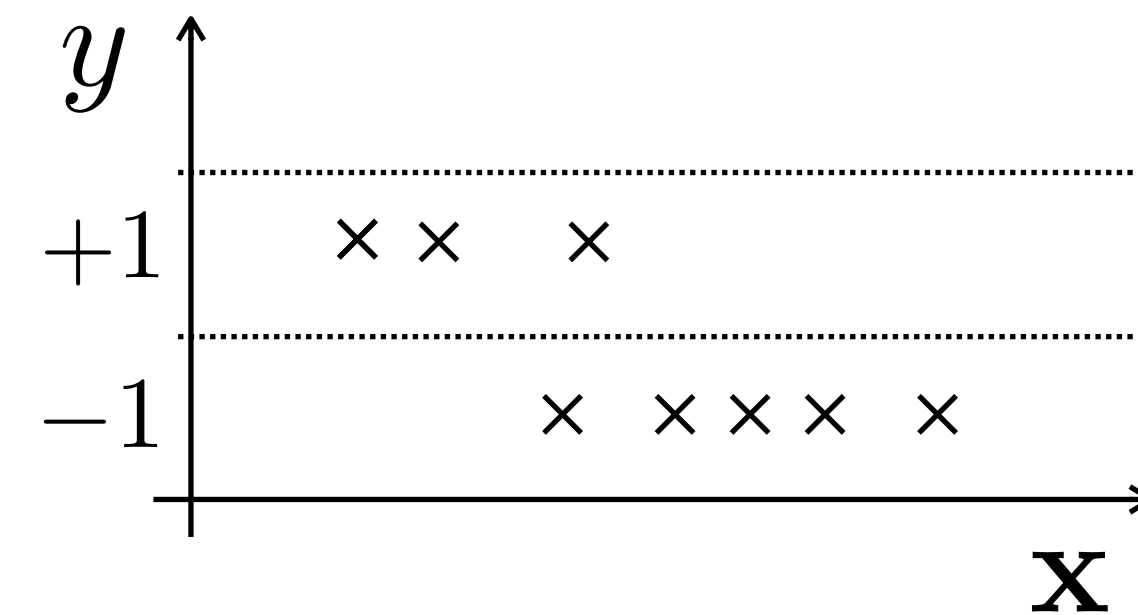
-1



### Training

Training = search for unknown parameters  $\mathbf{w}$   
which fits a given data

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$

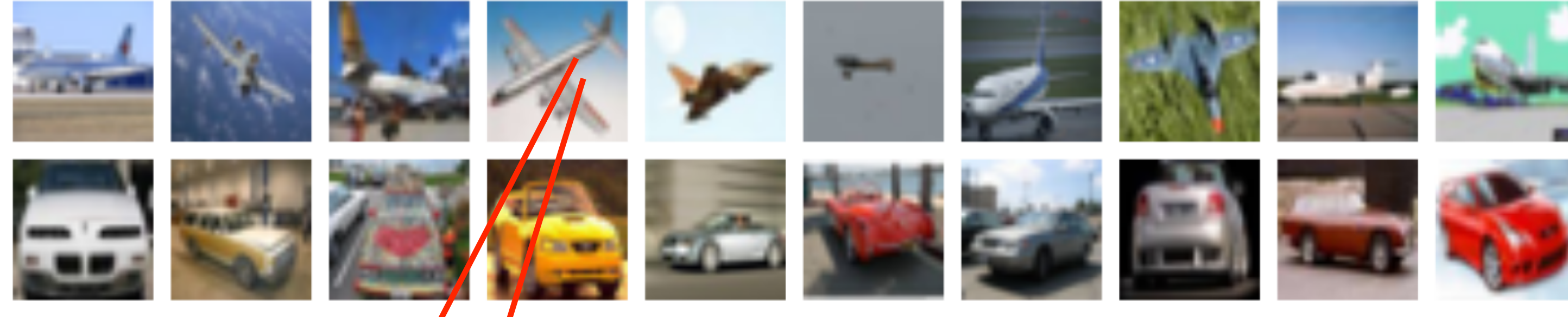


Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1

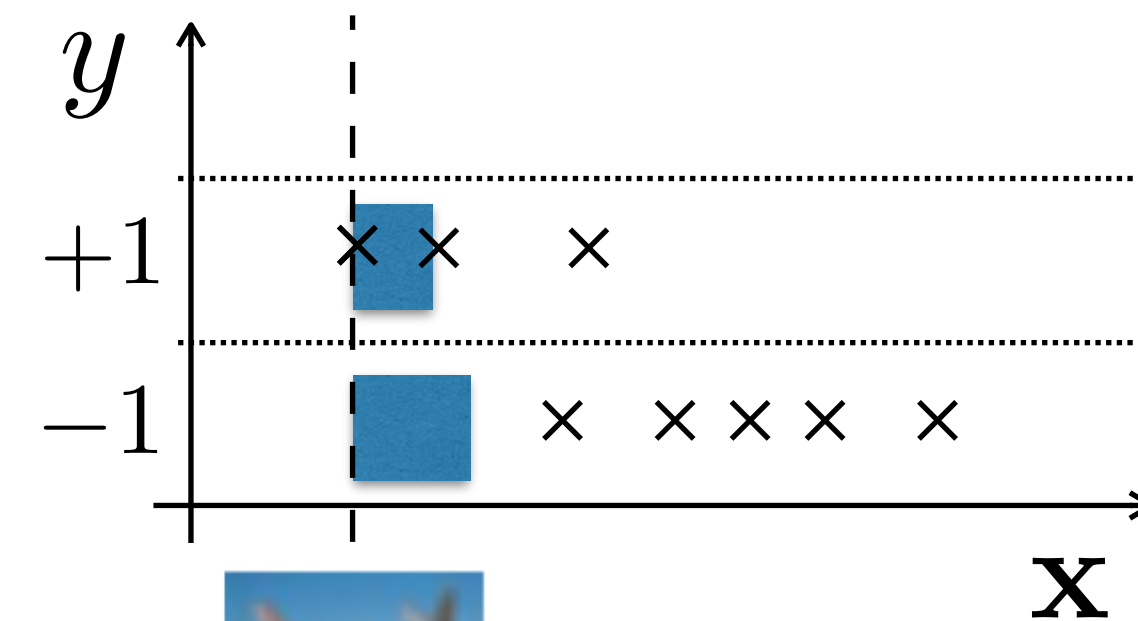
-1



### Training

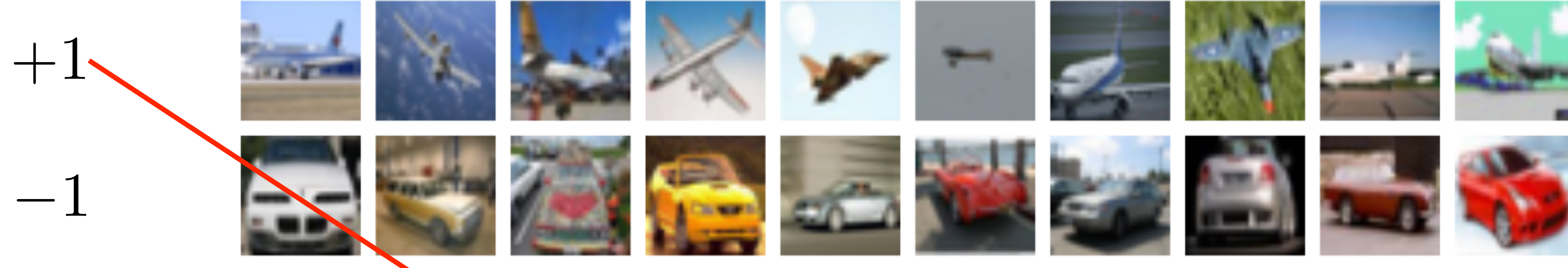
Training = search for unknown parameters  $\mathbf{w}$  which fits a given data

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



Labels ( $y_i$ )

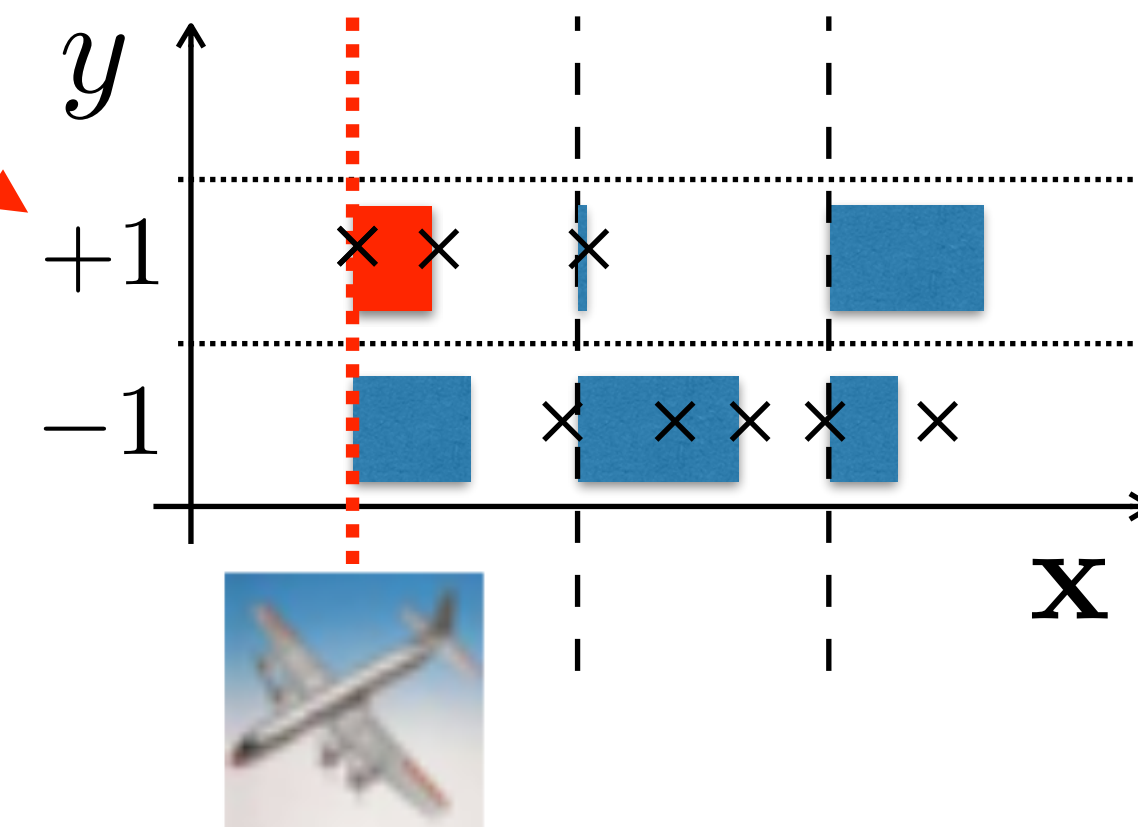
RGB images ( $\mathbf{x}_i$ )



### Training

Training = search for unknown parameters  $\mathbf{w}$  which fits a given data

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



Probability of observing  $y_i$  when measuring  $\mathbf{x}_i$  is

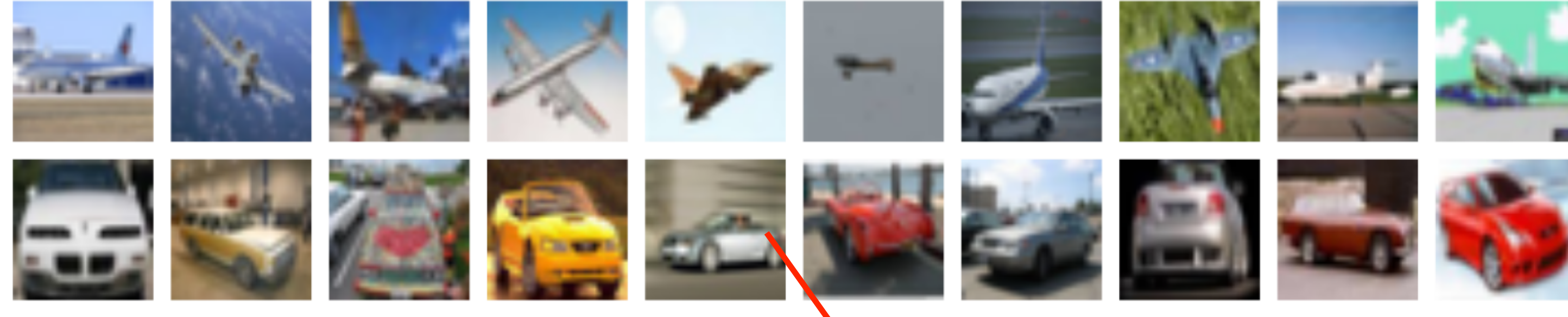


Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1

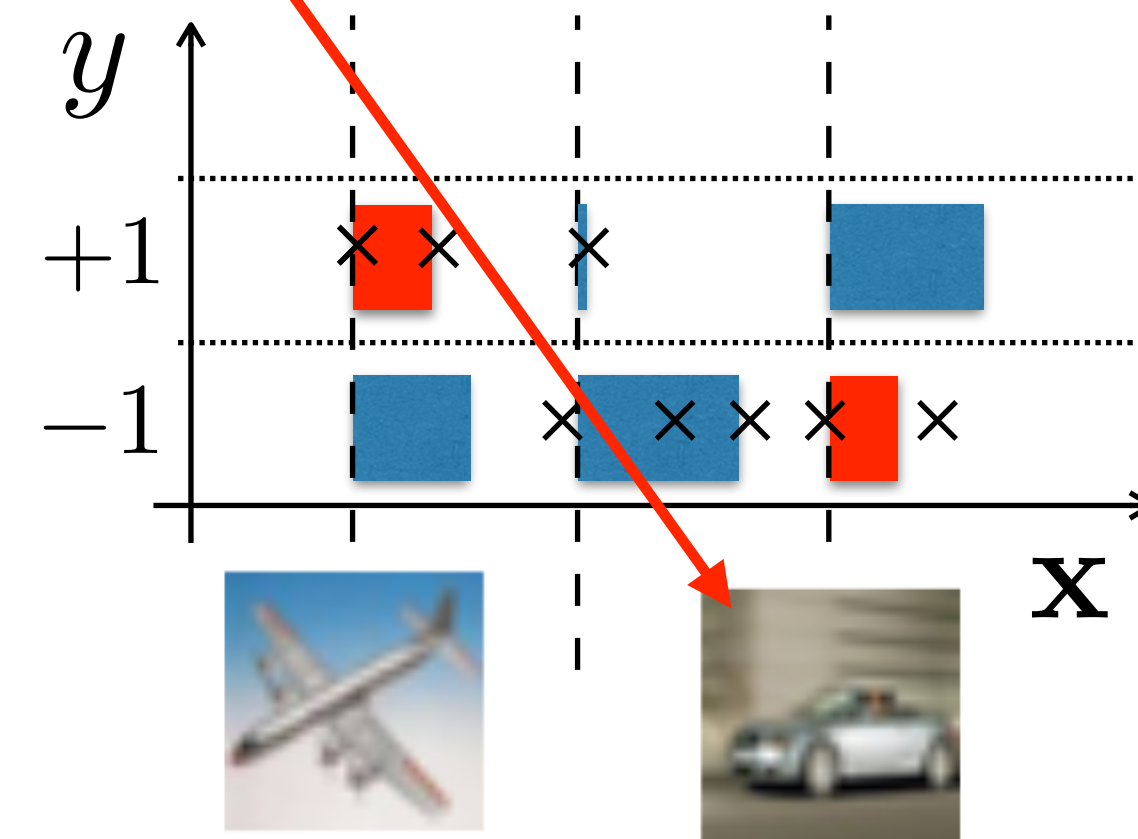
-1



### Training

Training = search for unknown parameters  $\mathbf{w}$  which fits a given data

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



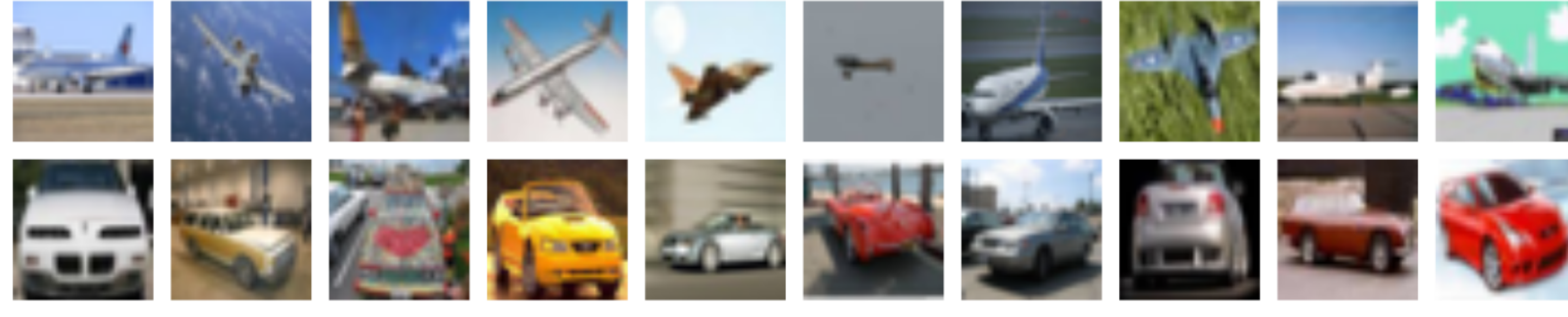
Probability of observing  $y_i$  when measuring  $\mathbf{x}_i$  is

Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1

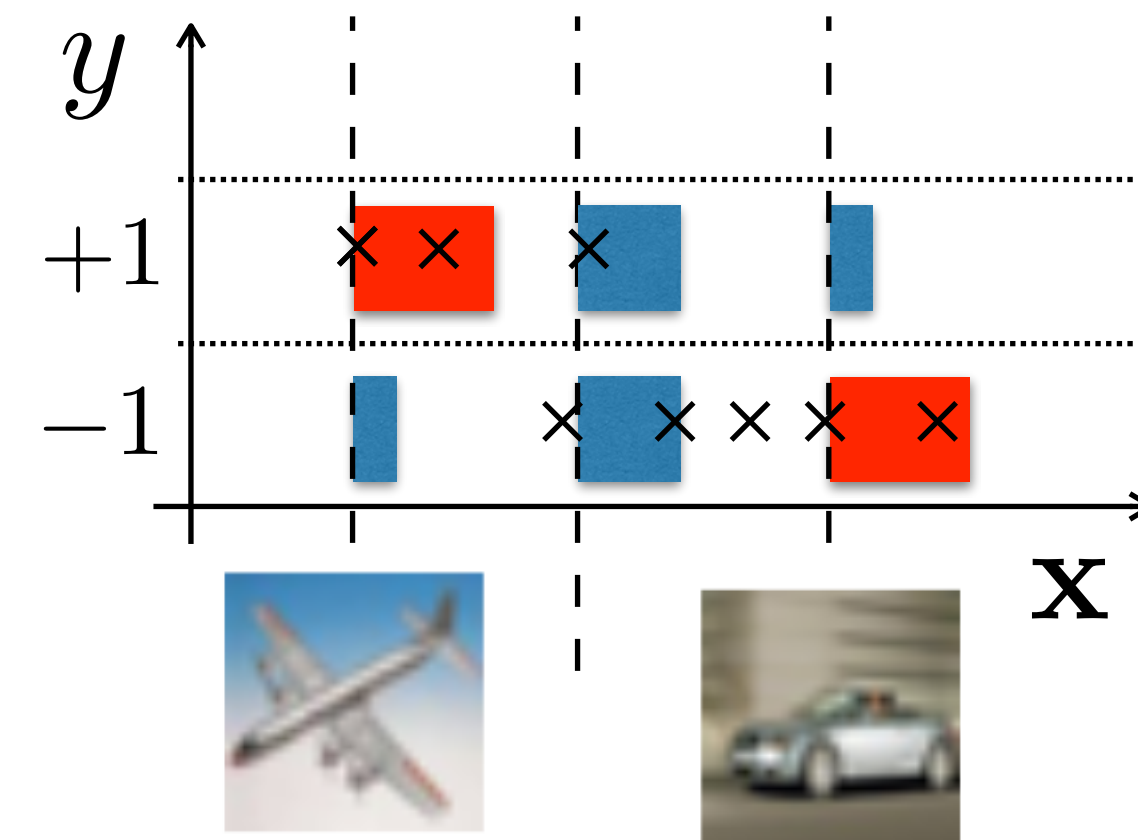
-1



### Training

Training = search for unknown parameters  $\mathbf{w}$  which fits a given data

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$



Probability of observing  $y_i$  when measuring  $\mathbf{x}_i$  is

## Two-class classification problem

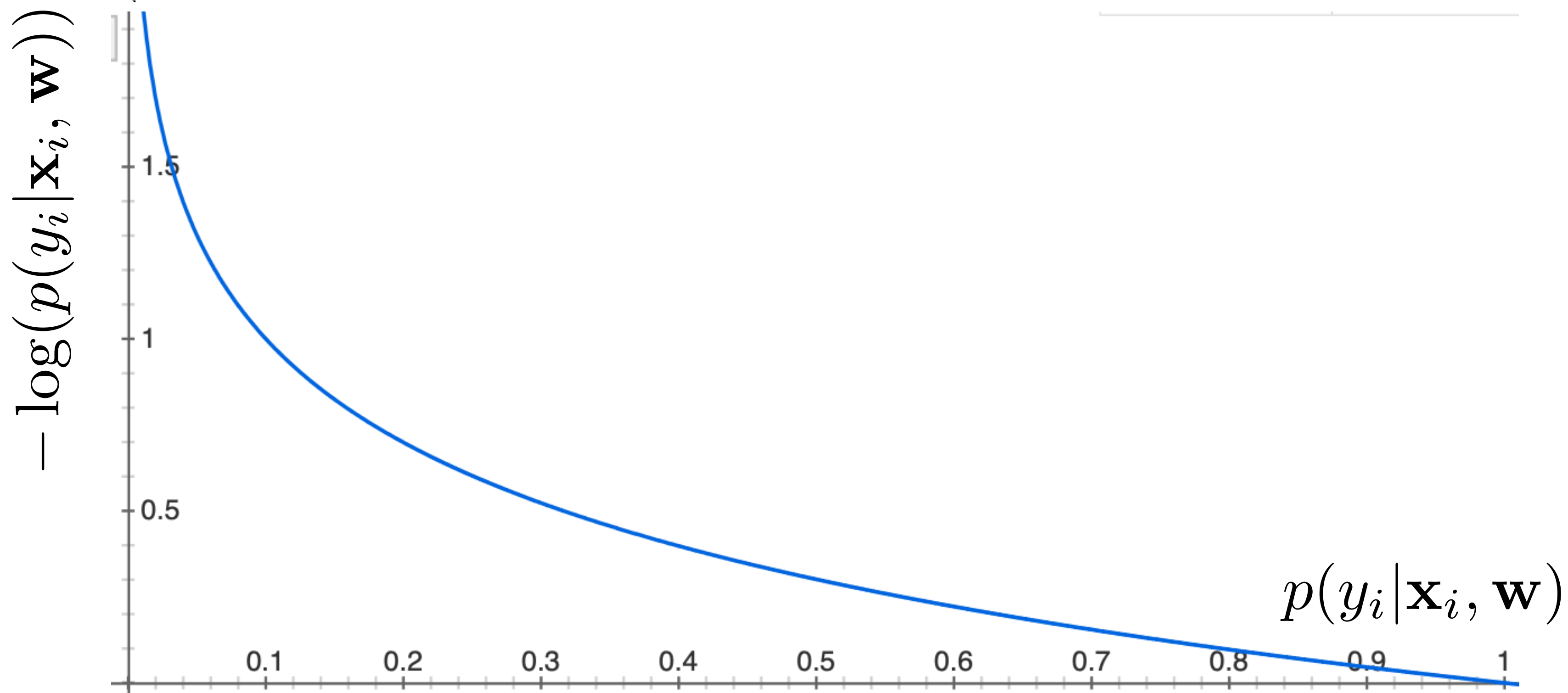
$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases} = \sigma(y \cdot f(\mathbf{x}, \mathbf{w}))$$

# Two-class classification problem

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases} = \sigma(y \cdot f(\mathbf{x}, \mathbf{w}))$$

MLE:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right)$$



## Equivalence of common binary classification losses

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases} = \sigma(y \cdot f(\mathbf{x}, \mathbf{w}))$$

MLE:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right)$$
$$= -\log [\sigma(y_i f(\mathbf{x}_i, \mathbf{w}))]$$

# Equivalence of common binary classification losses

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases} = \sigma(y \cdot f(\mathbf{x}, \mathbf{w}))$$

MLE:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right)$$
$$= -\log [\sigma(y_i f(\mathbf{x}_i, \mathbf{w}))]$$



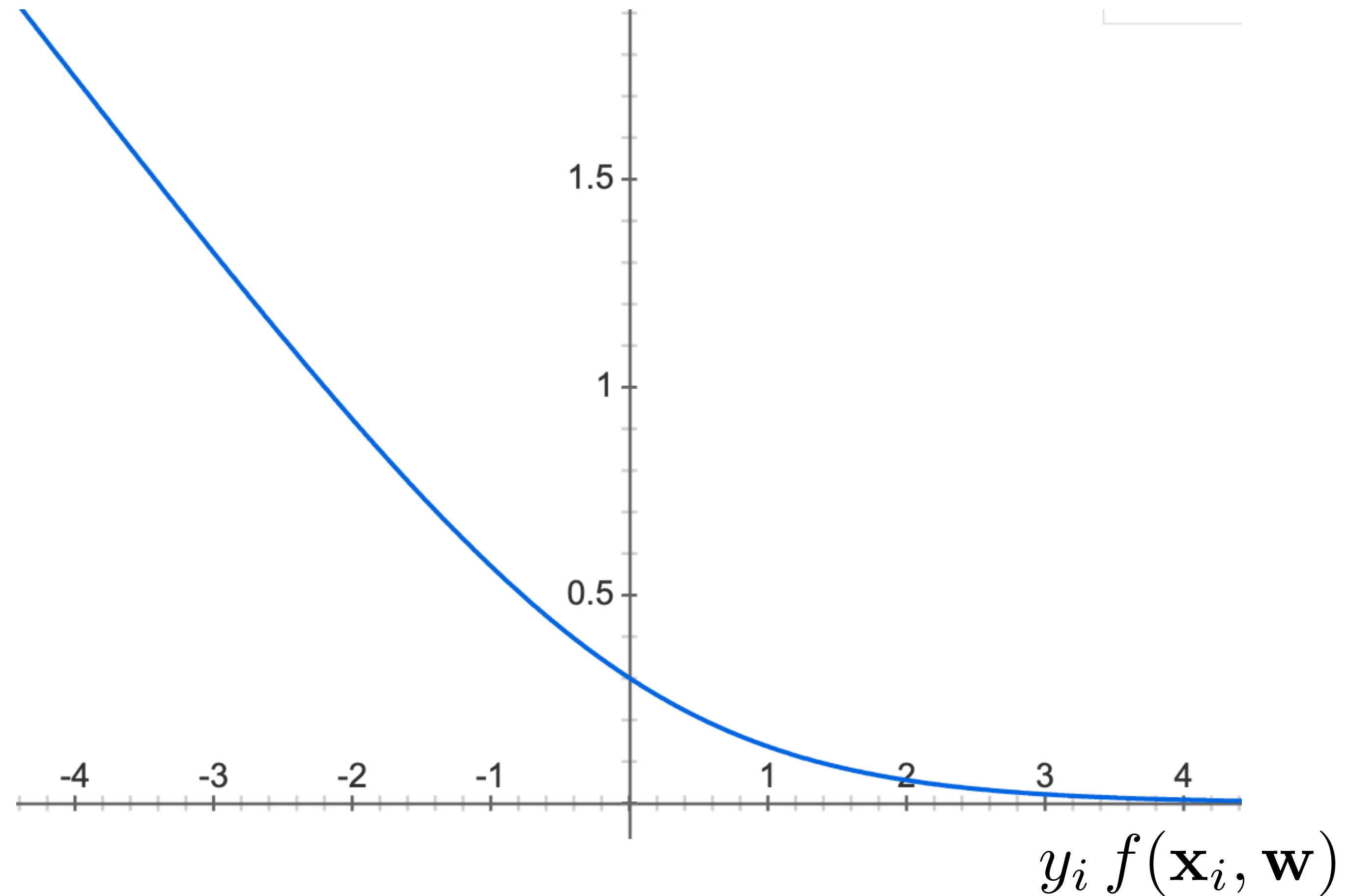
Logistic loss

$$\log [1 + \exp(-y_i f(\mathbf{x}_i, \mathbf{w}))]$$

# Equivalence of common binary classification losses



$$\text{Logistic loss}$$
$$\log [1 + \exp(-y_i f(\mathbf{x}_i, \mathbf{w}))]$$



# Equivalence of common binary classification losses

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases} = \sigma(y \cdot f(\mathbf{x}, \mathbf{w}))$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right)$$
$$= -\log [\sigma(y_i f(\mathbf{x}_i, \mathbf{w}))]$$

$\Leftrightarrow$

Logistic loss

$$\log [1 + \exp(-y_i f(\mathbf{x}_i, \mathbf{w}))]$$



# Equivalence of common binary classification losses

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases} = \sigma(y \cdot f(\mathbf{x}, \mathbf{w}))$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right)$$
$$= -\log [\sigma(y_i f(\mathbf{x}_i, \mathbf{w}))]$$

$\Leftrightarrow$   
Logistic loss  
 $\log [1 + \exp(-y_i f(\mathbf{x}_i, \mathbf{w}))]$

# Equivalence of common binary classification losses

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right)$$
$$= -\log [\sigma(y_i f(\mathbf{x}_i, \mathbf{w}))]$$

$\Leftrightarrow$

Logistic loss

$$\log [1 + \exp(-y_i f(\mathbf{x}_i, \mathbf{w}))]$$

$\Leftrightarrow$

Cross-entropy loss

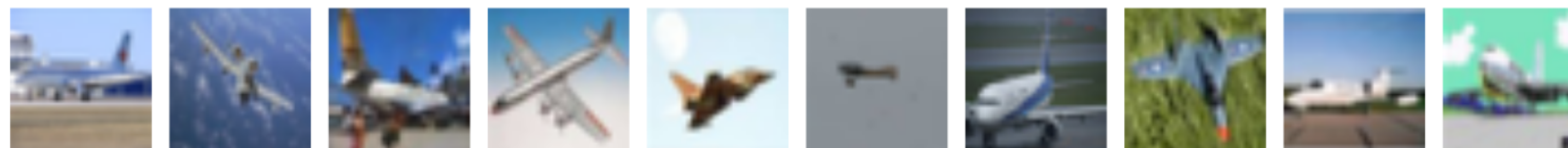
$$-\left[ \bar{y}_i \cdot \log(\sigma(f(\mathbf{x}_i, \mathbf{w}))) + (1 - \bar{y}_i) \cdot \log(1 - \sigma(f(\mathbf{x}_i, \mathbf{w}))) \right]$$

$$\bar{y}_i = \frac{y_i + 1}{2} \in \{0, 1\}$$

Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1







-1



## Training

Example: Training linear classifier

```
def train(       =
```

+1 +1 +1 -1 -1 -1 ):

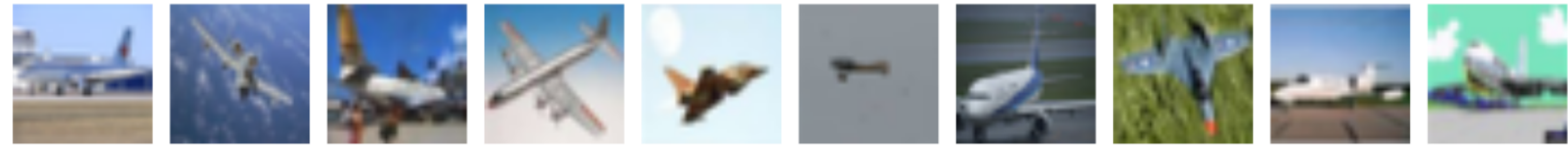
$\mathbf{x}_i = \text{vec}(  ) \quad \forall_i$

```
return  $\mathbf{w}^*$ 
```

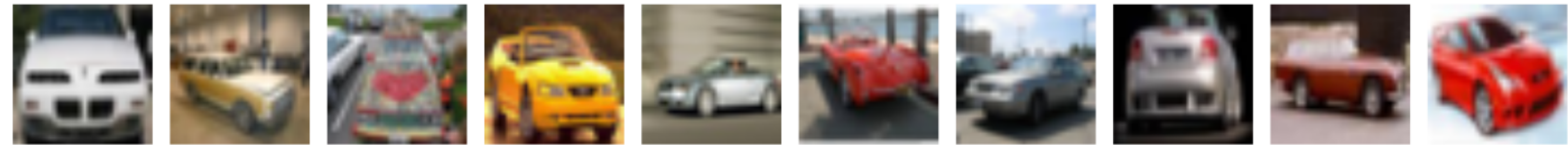
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1










-1



## Training

Example: Training linear classifier

```
def train(       =  
          +1 +1 +1 -1 -1 -1 ):  
     $\mathbf{x}_i = \text{vec}(\text{  }) \quad \forall_i$   
     $\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_i \log [1 + \exp(-y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)]$   
    return  $\mathbf{w}^*$ 
```

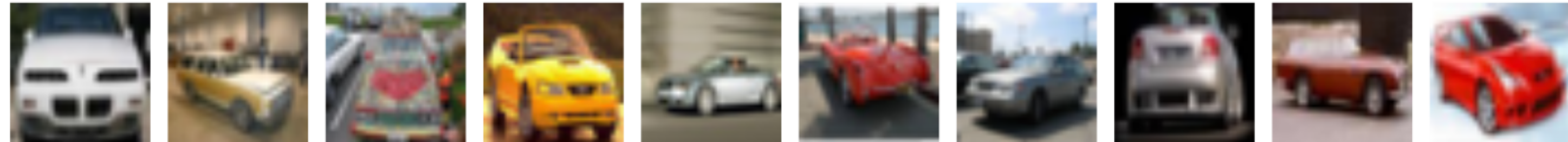
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1



-1



### Training

Example: Training linear classifier

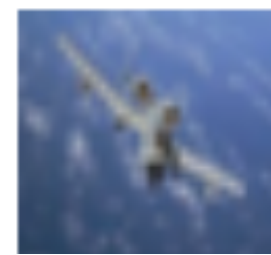
score

label

-2.5

$y_i = -1$

```
def train(
```

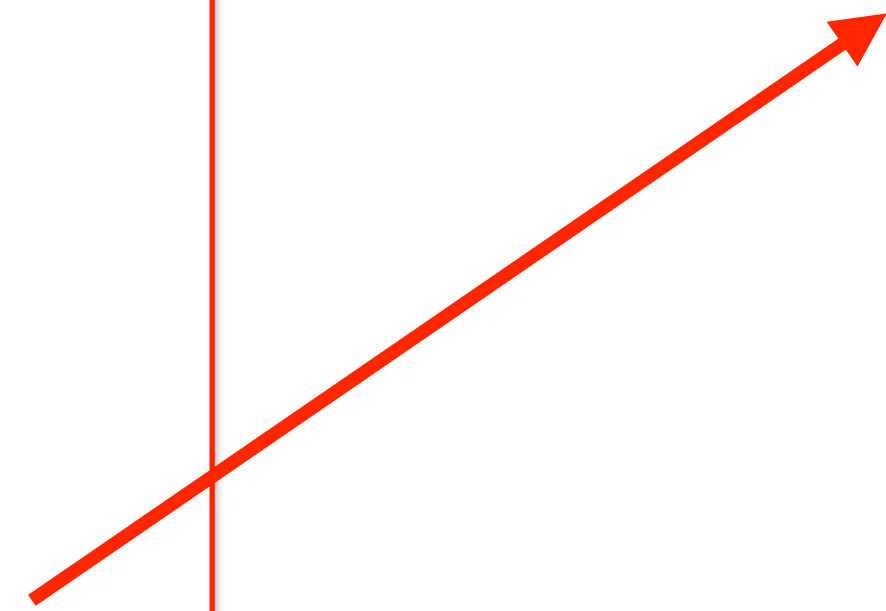


```
    +1 +1 +1 -1 -1 -1 ):
```

$$\mathbf{x}_i = \text{vec}(\text{img}_i) \quad \forall_i$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_i \log [1 + \exp(-y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)]$$

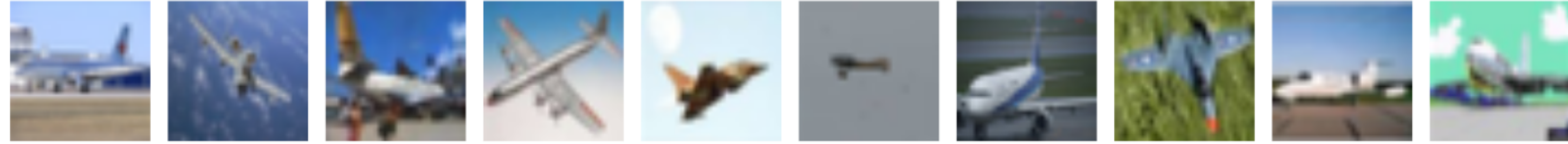
```
return  $\mathbf{w}^*$ 
```



Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1










-1



### Training

Example: Training linear classifier

```
def train(       =  
         +1 +1 +1 -1 -1 -1 ):  
     $\mathbf{x}_i = \text{vec}(\text{  }) \quad \forall_i$   
     $\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_i \log [1 + \exp(-y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)]$   
    return  $\mathbf{w}^*$ 
```

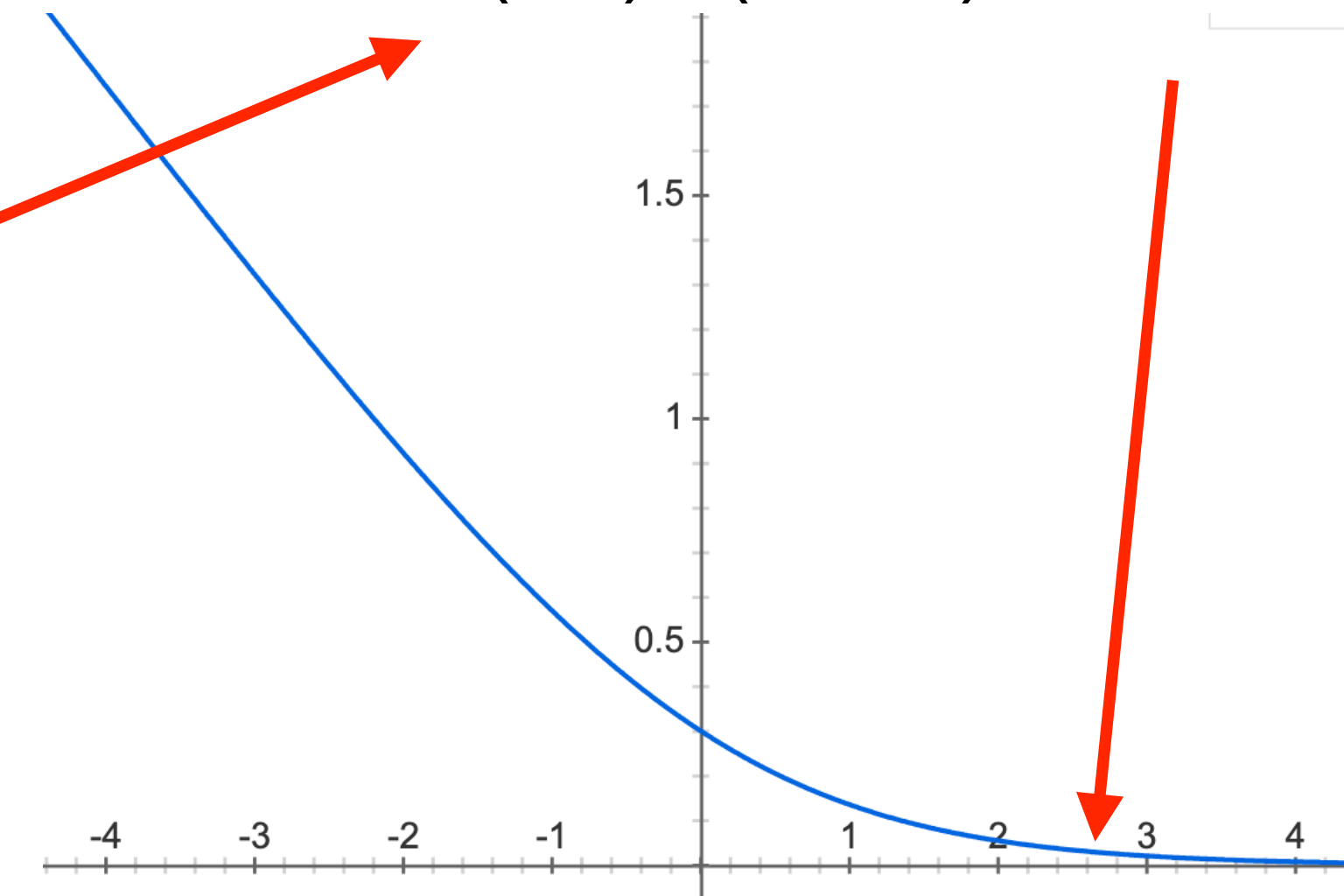
score

label

-2.5

$y_i = -1$

$$(-1) \times (-2.5) = 2.5$$

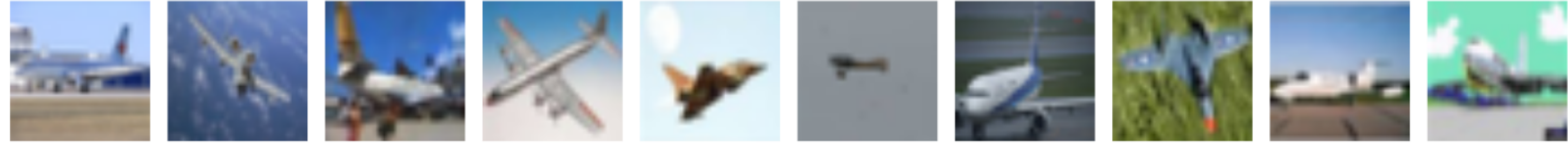


$y_i f(\mathbf{x}_i, \mathbf{w})$

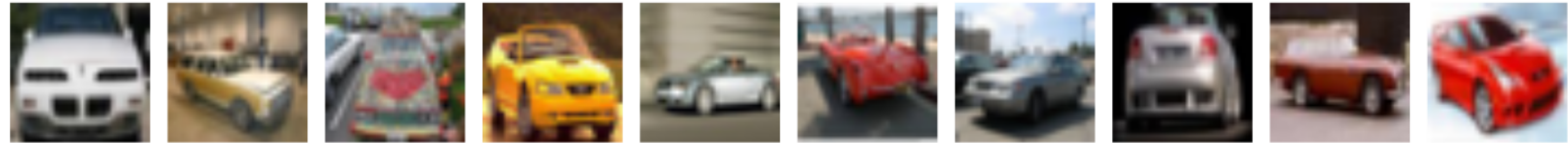
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1










-1



### Training

Example: Training linear classifier

```
def train(       =  
         +1 +1 +1 -1 -1 -1 ):  
     $\mathbf{x}_i = \text{vec}(\text{  }) \quad \forall_i$   
     $\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_i \log [1 + \exp(-y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)]$   
    return  $\mathbf{w}^*$ 
```

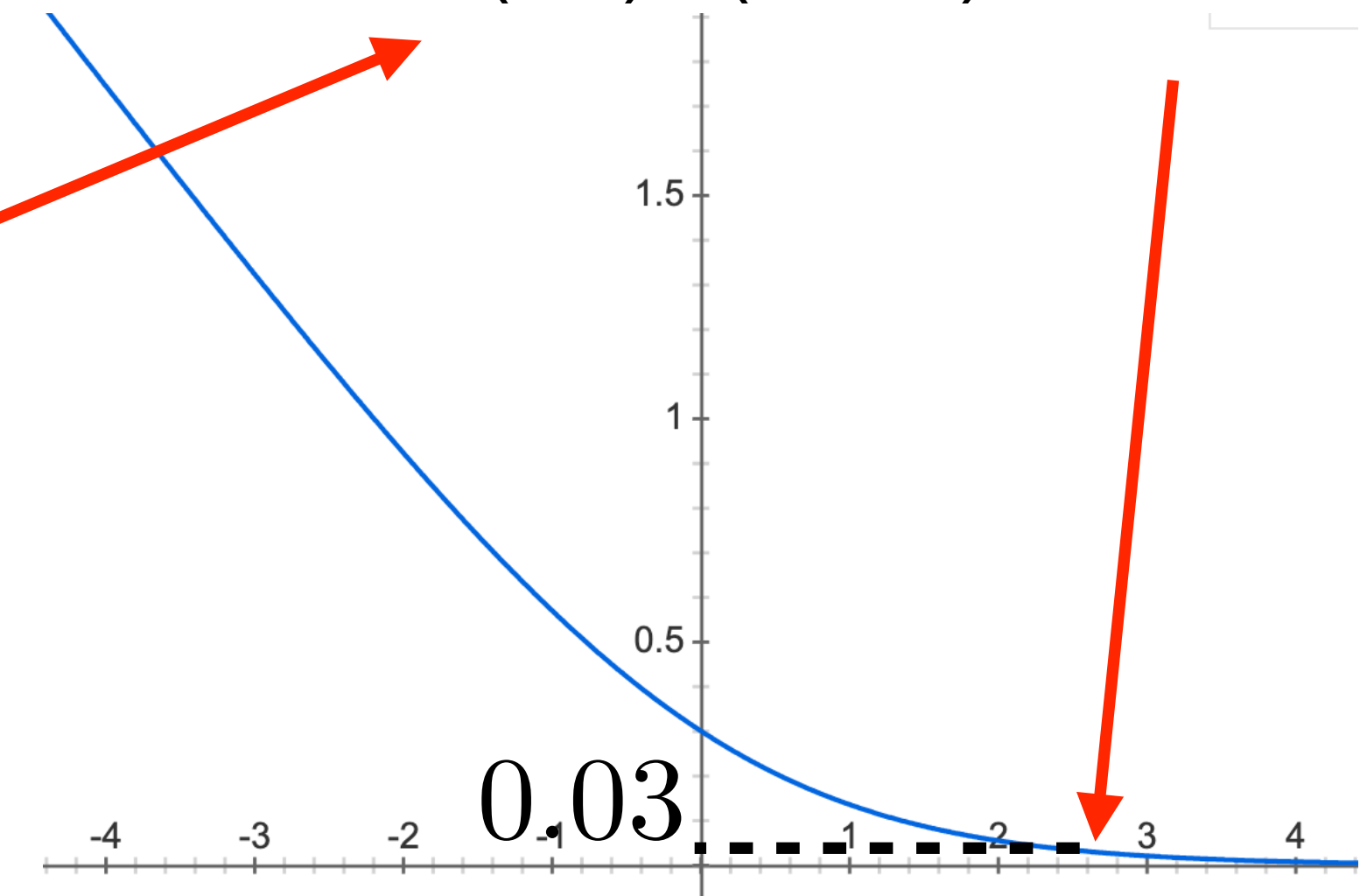
score

label

-2.5

$y_i = -1$

$$(-1) \times (-2.5) = 2.5$$

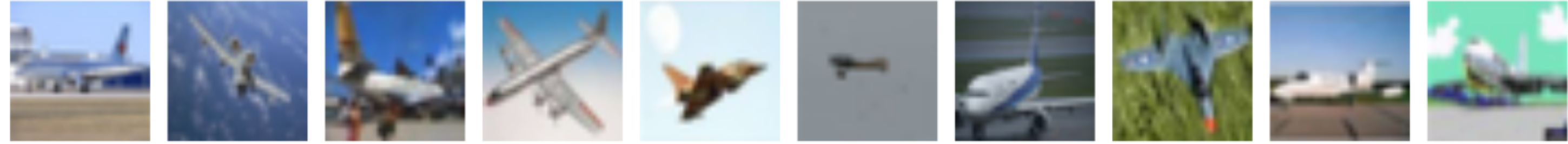


$y_i f(\mathbf{x}_i, \mathbf{w})$

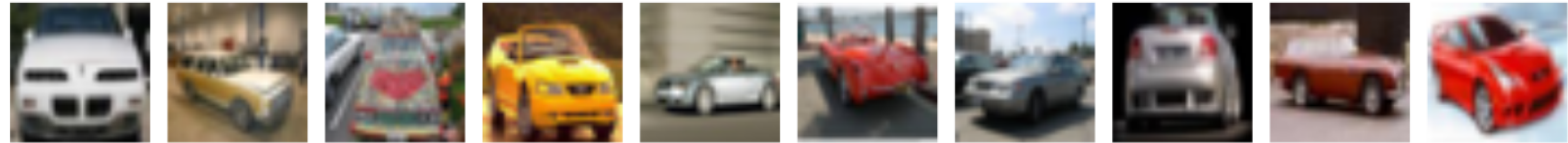
Labels ( $y_i$ )

RGB images ( $\mathbf{x}_i$ )

+1










-1



### Training

Example: Training linear classifier

```
def train(       =  
          +1 +1 +1 -1 -1 -1 ):  
     $\mathbf{x}_i = \text{vec}(\text{  }) \quad \forall_i$   
     $\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_i \log [1 + \exp(-y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)]$   
    return  $\mathbf{w}^*$ 
```

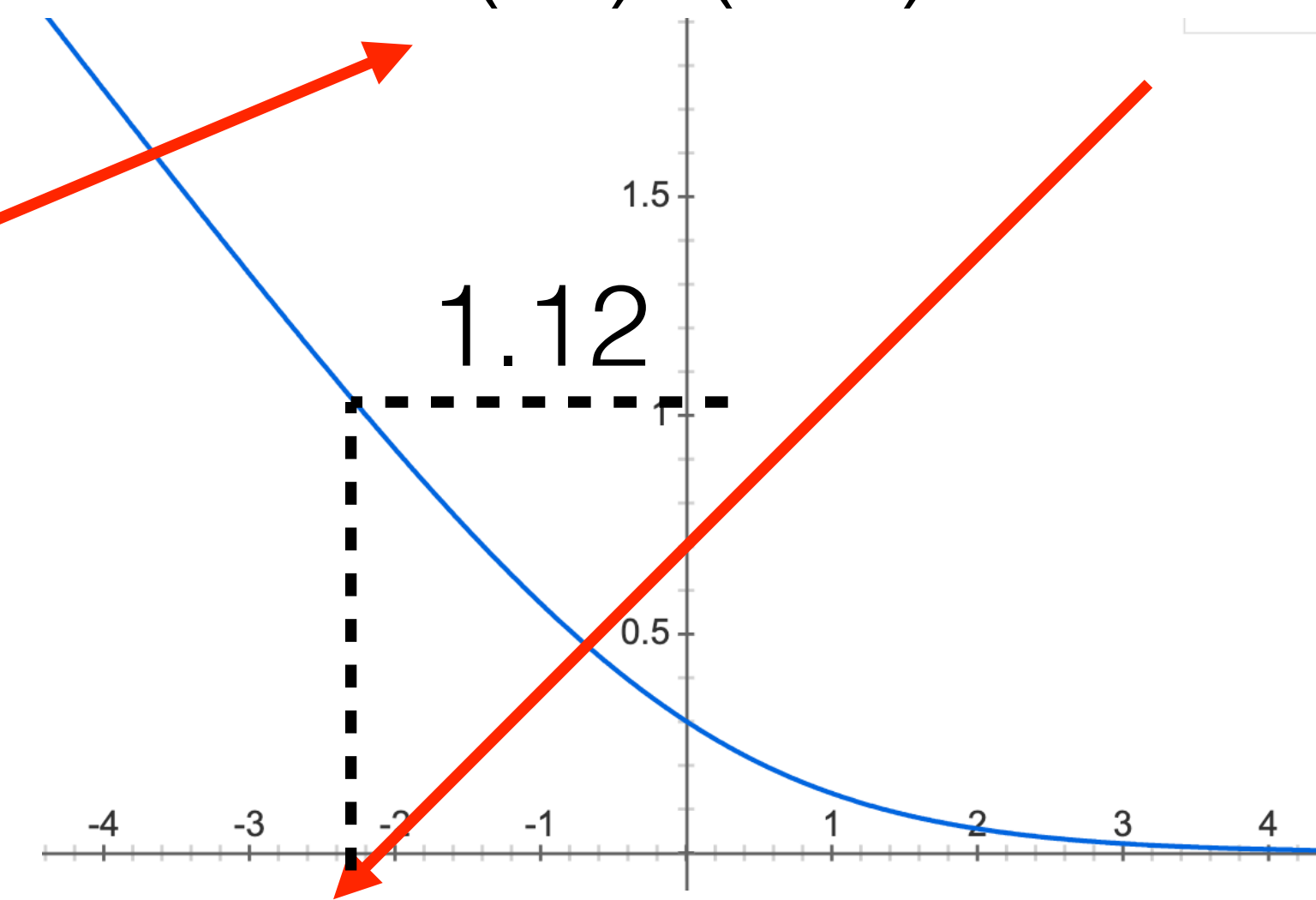
score

label

2.5

$y_i = -1$

$$(-1) \times (2.5) = -2.5$$



$$y_i f(\mathbf{x}_i, \mathbf{w})$$



$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_i \log [1 + \exp(-y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)]$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \underbrace{\sum_i \log [1 + \exp(-y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)]}_{\mathcal{L}(\mathbf{w})}$$

- There is no closed-form solution
- Gradient optimization

$$\mathbf{w} = \mathbf{w} - \alpha \left[ \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} \right]^\top \text{ where } \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = ?$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_i \log [1 + \exp(-y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)]$$
$$\mathcal{L}(\mathbf{w})$$

- There is no closed-form solution
- Gradient optimization

$$\mathbf{w} = \mathbf{w} - \alpha \left[ \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} \right]^\top \text{ where } \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \sum_i \frac{-y_i \bar{\mathbf{x}}_i^\top}{1 + \exp(y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)}$$

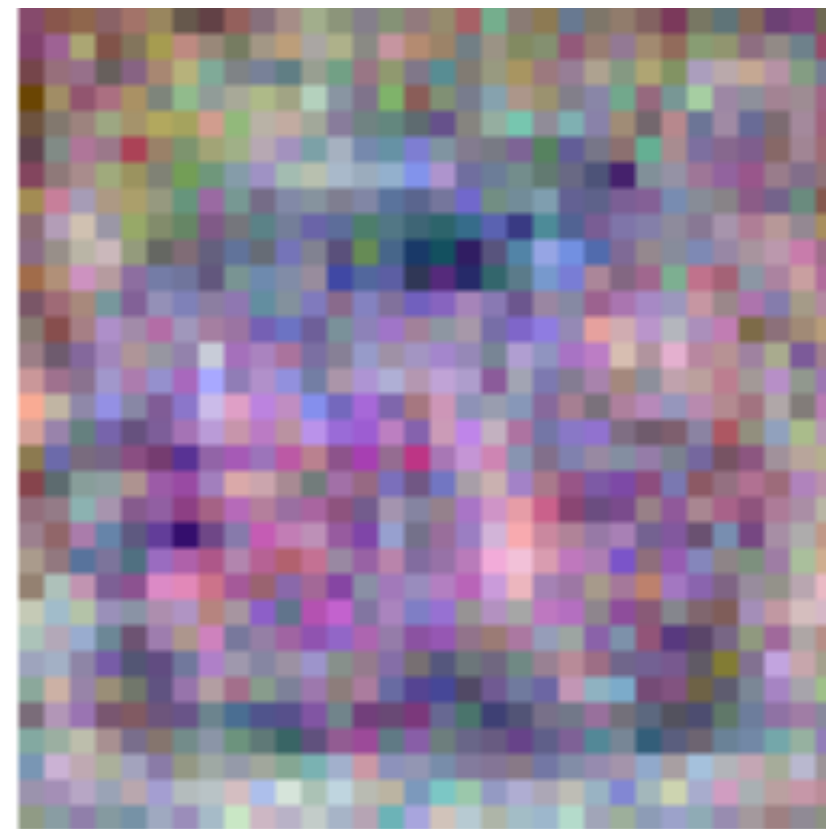
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_i \log [1 + \exp(-y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)]$$

$$\mathcal{L}(\mathbf{w})$$

- There is no closed-form solution
- Gradient optimization

$$\mathbf{w} = \mathbf{w} - \alpha \left[ \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} \right]^\top \text{ where } \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \sum_i \frac{-y_i \bar{\mathbf{x}}_i^\top}{1 + \exp(y_i \mathbf{w}^\top \bar{\mathbf{x}}_i)}$$

Learned weights  
as a template:



automobile

Show python code

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function

prior/regulariser

- **Classification:**  $p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$
- Choice of  $f(\mathbf{x}, \mathbf{w})$  is crucial

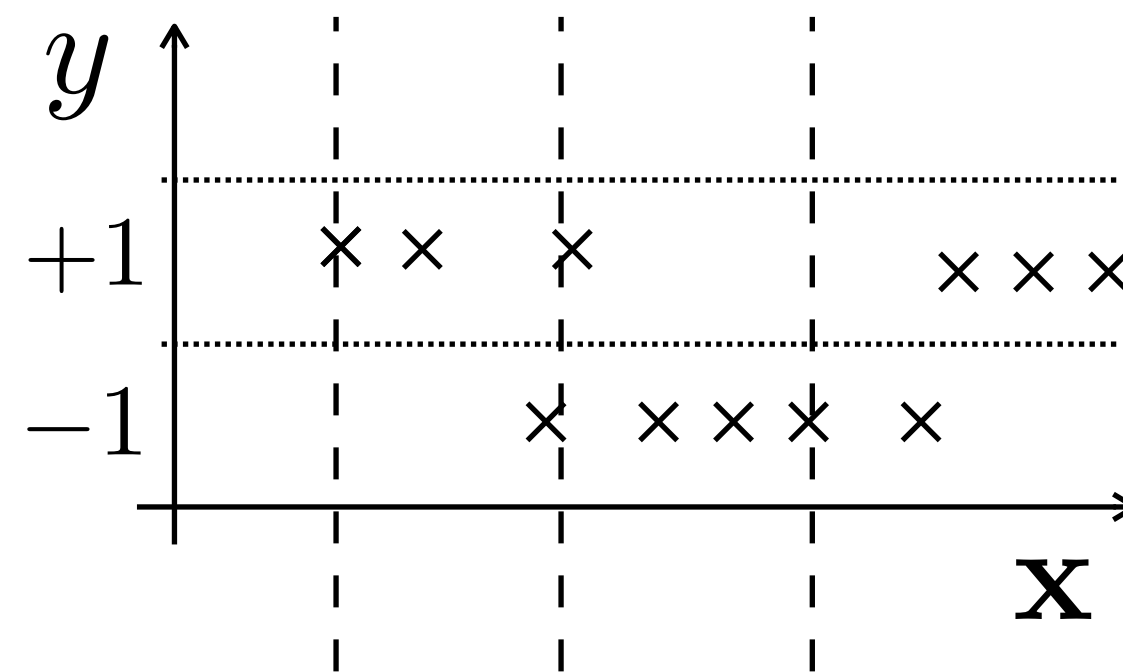
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function

prior/regulariser

- **Classification:**  $p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$
- Linear  $f(\mathbf{x}, \mathbf{w})$  cannot generate wild decision boundary

1D example:



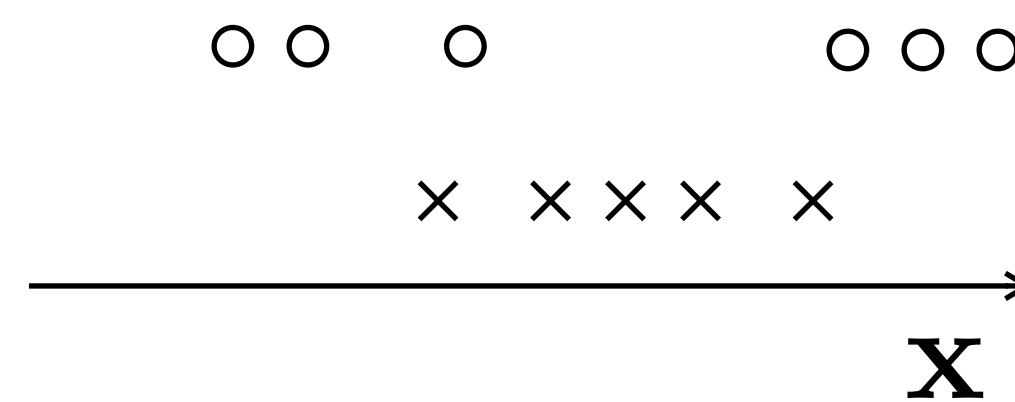
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function

prior/regulariser

- **Classification:**  $p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$
- Linear  $f(\mathbf{x}, \mathbf{w})$  cannot generate wild decision boundary

1D example:





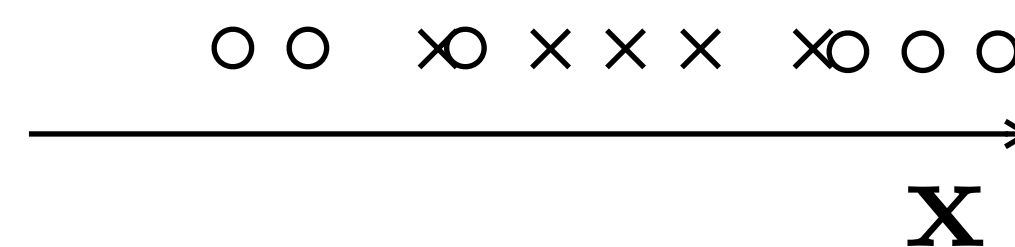
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function

prior/regulariser

- **Classification:**  $p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$
- Linear  $f(\mathbf{x}, \mathbf{w})$  cannot generate wild decision boundary

1D example:



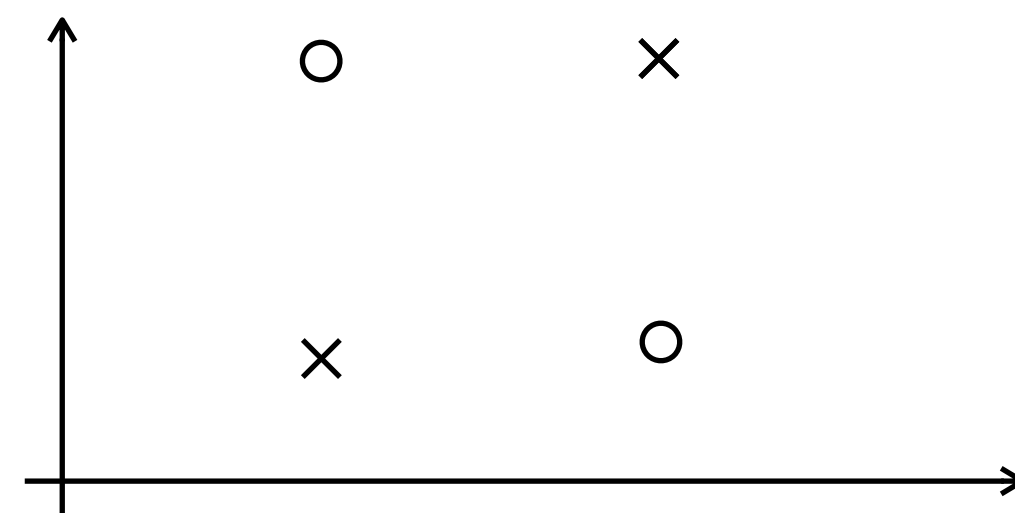
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function

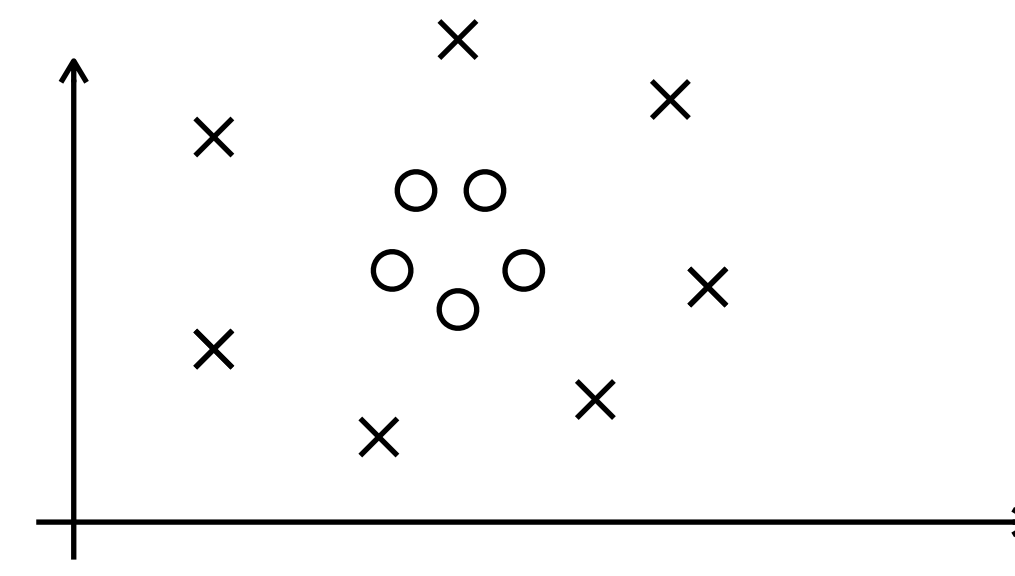
prior/regulariser

- **Classification:**  $p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$
- Linear  $f(\mathbf{x}, \mathbf{w})$  cannot generate wild decision boundary

2D example:



XOR



circle

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function

prior/regulariser

- **Classification:**  $p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$
- Wild  $f(\mathbf{x}, \mathbf{w})$  with high-dimensional  $\mathbf{w}$  suffers from the curse of dimensionality and overfitting



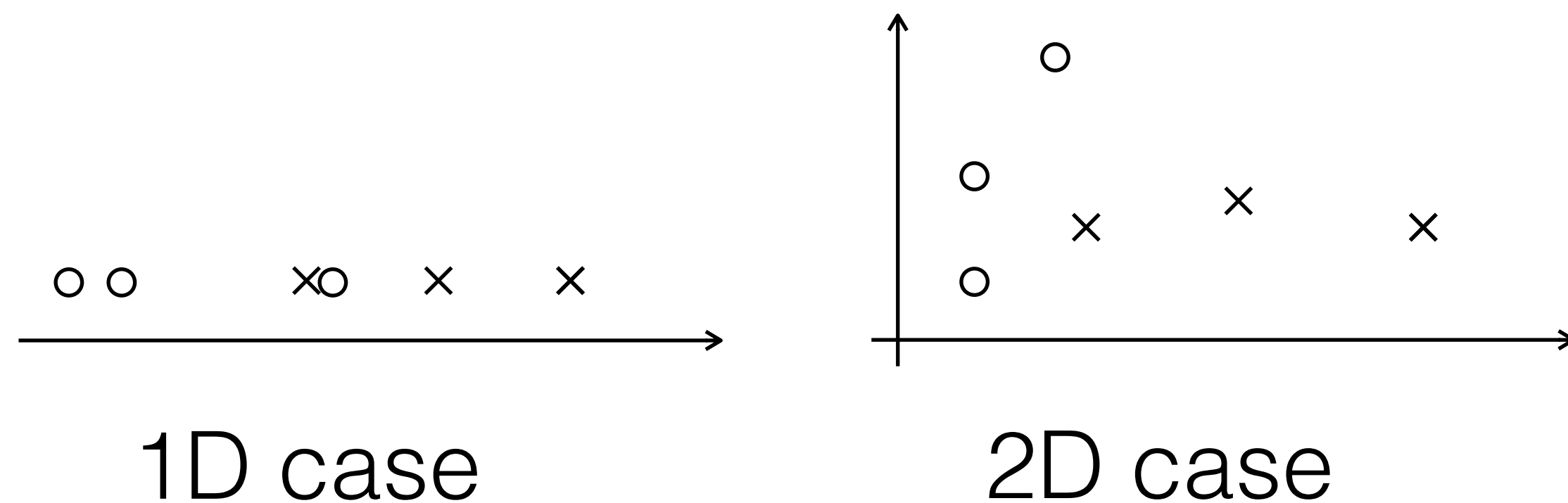
1D case

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function

prior/regulariser

- **Classification:**  $p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$
- Wild  $f(\mathbf{x}, \mathbf{w})$  with high-dimensional  $\mathbf{w}$  suffers from the curse of dimensionality and overfitting

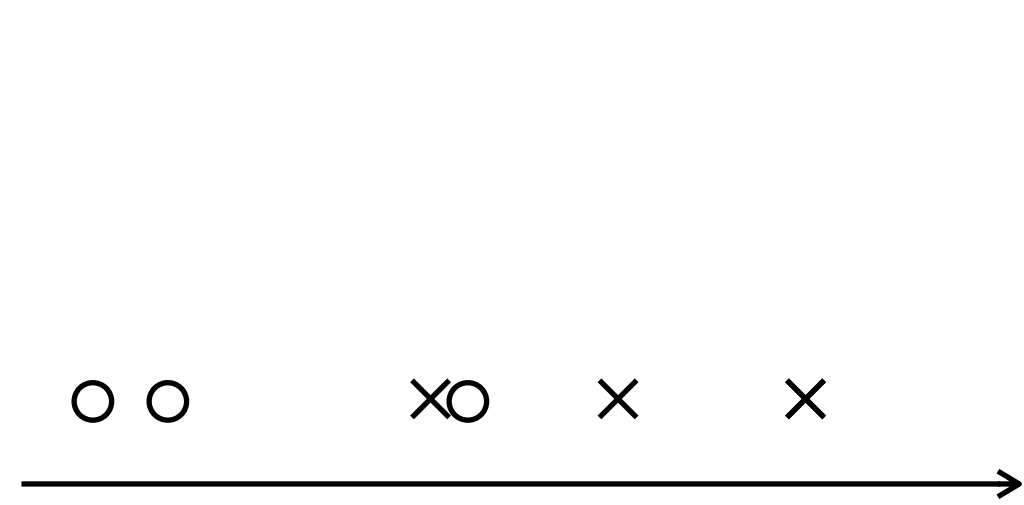


$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

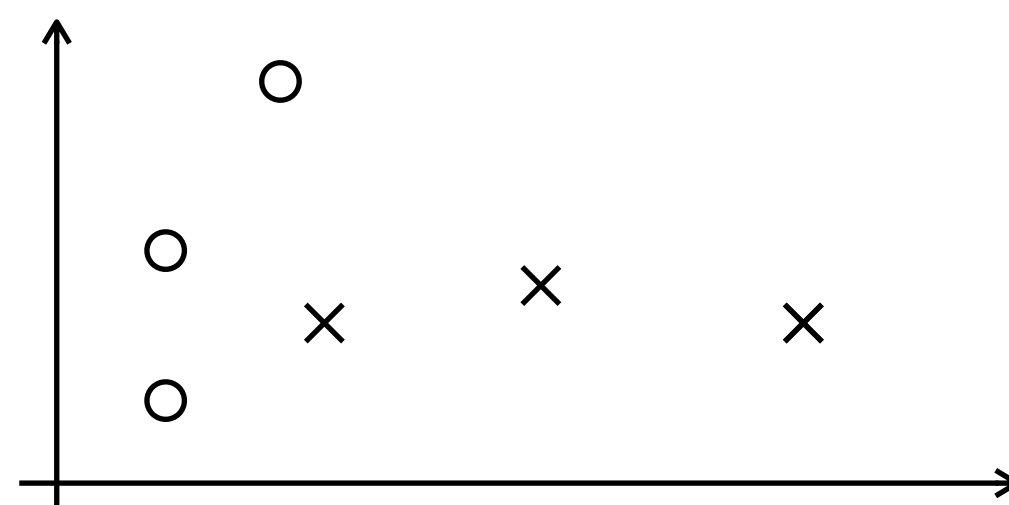
loss function

prior/regulariser

- **Classification:**  $p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$
- Wild  $f(\mathbf{x}, \mathbf{w})$  with high-dimensional  $\mathbf{w}$  suffers from the curse of dimensionality and overfitting



1D case



2D case

???

CIFAR case

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function

prior/regulariser

- **Classification:**  $p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = +1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = -1 \end{cases}$
- Wild  $f(\mathbf{x}, \mathbf{w})$  with high-dimensional  $\mathbf{w}$  suffers from the curse of dimensionality and overfitting
- We exploit prior  $p(\mathbf{w})$  to restrict the wildness of  $f(\mathbf{x}, \mathbf{w})$ 
  - L2 regulariser  $p(\mathbf{w}) = \mathcal{N}_{\mathbf{w}}(0, \sigma^2) \Rightarrow \|\mathbf{w}\|_2^2$
  - L1 regulariser, L1+L2 regulariser (elastic net)
  - prior on  $f(\mathbf{x}, \mathbf{w})$  structure (e.g. consists of convolutions)
  - batch normalization

## Conclusions

- Explained regression and linear classifier as MAP/ML estimator
- Discussed under/overfitting and regularisations
- Next lesson will go deeper

### **Competencies required for the test T1**

- Derive MAP/ML estimate for two-class and K-class classification problem.
- Compute logistic-loss and cross-entropy-loss
- Understand when classifier score has high/low values.
- Understand when loss has high/low values.