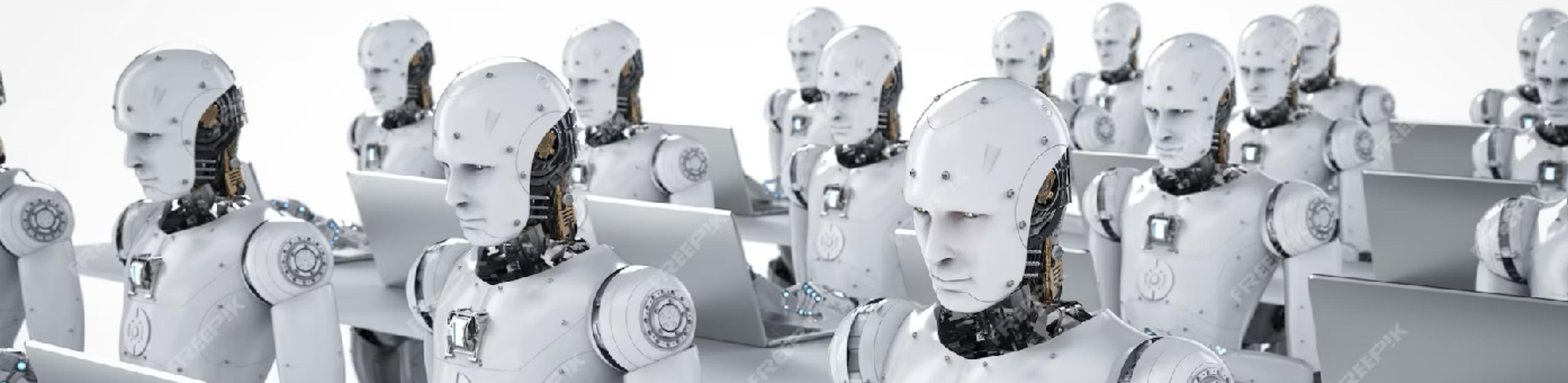


Učení robotů



Under the hood of linear classifier

Linear classifier on RGB images

Pre-requisites: linear algebra,

Karel Zimmermann

Czech Technical University in Prague

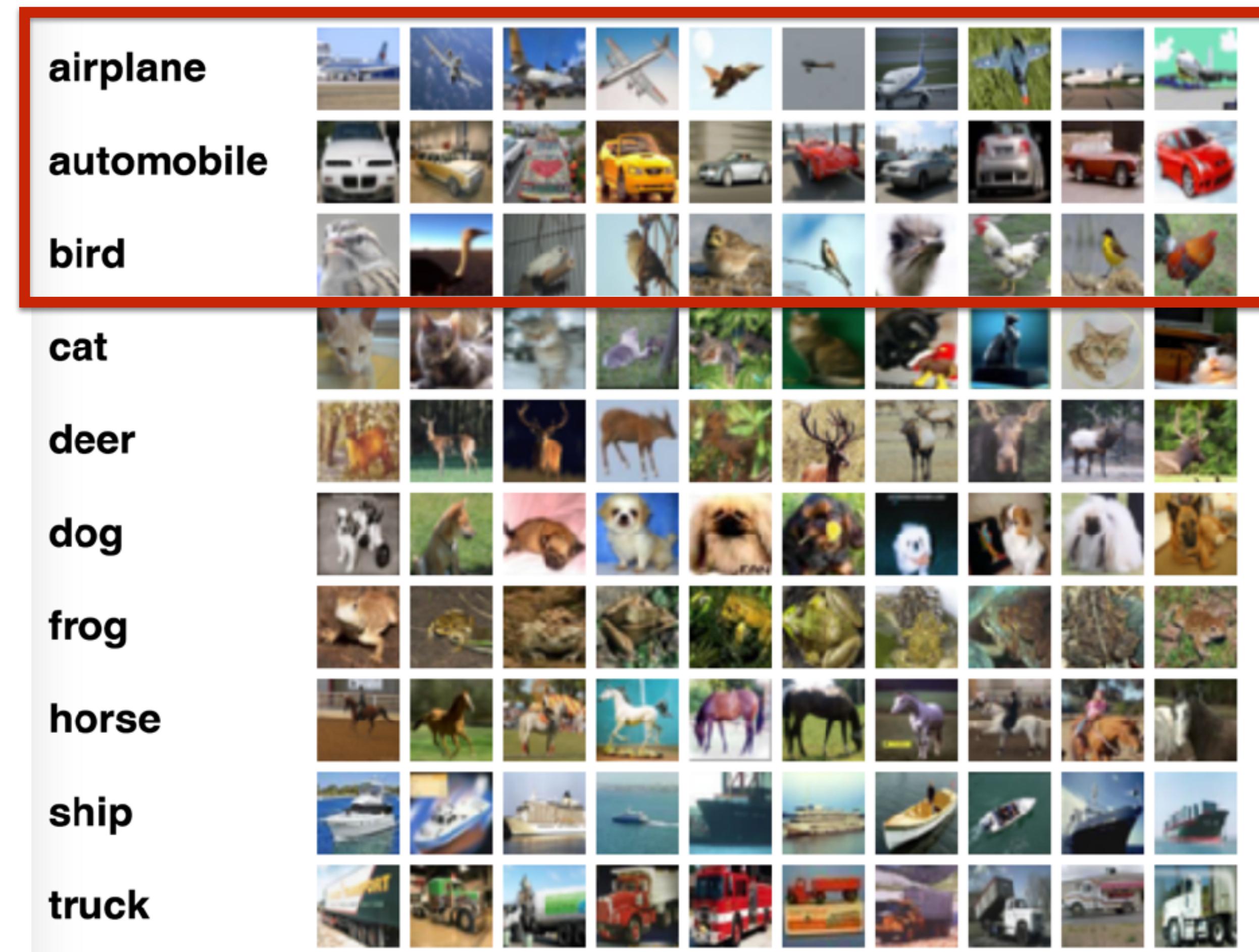
Faculty of Electrical Engineering, Department of Cybernetics



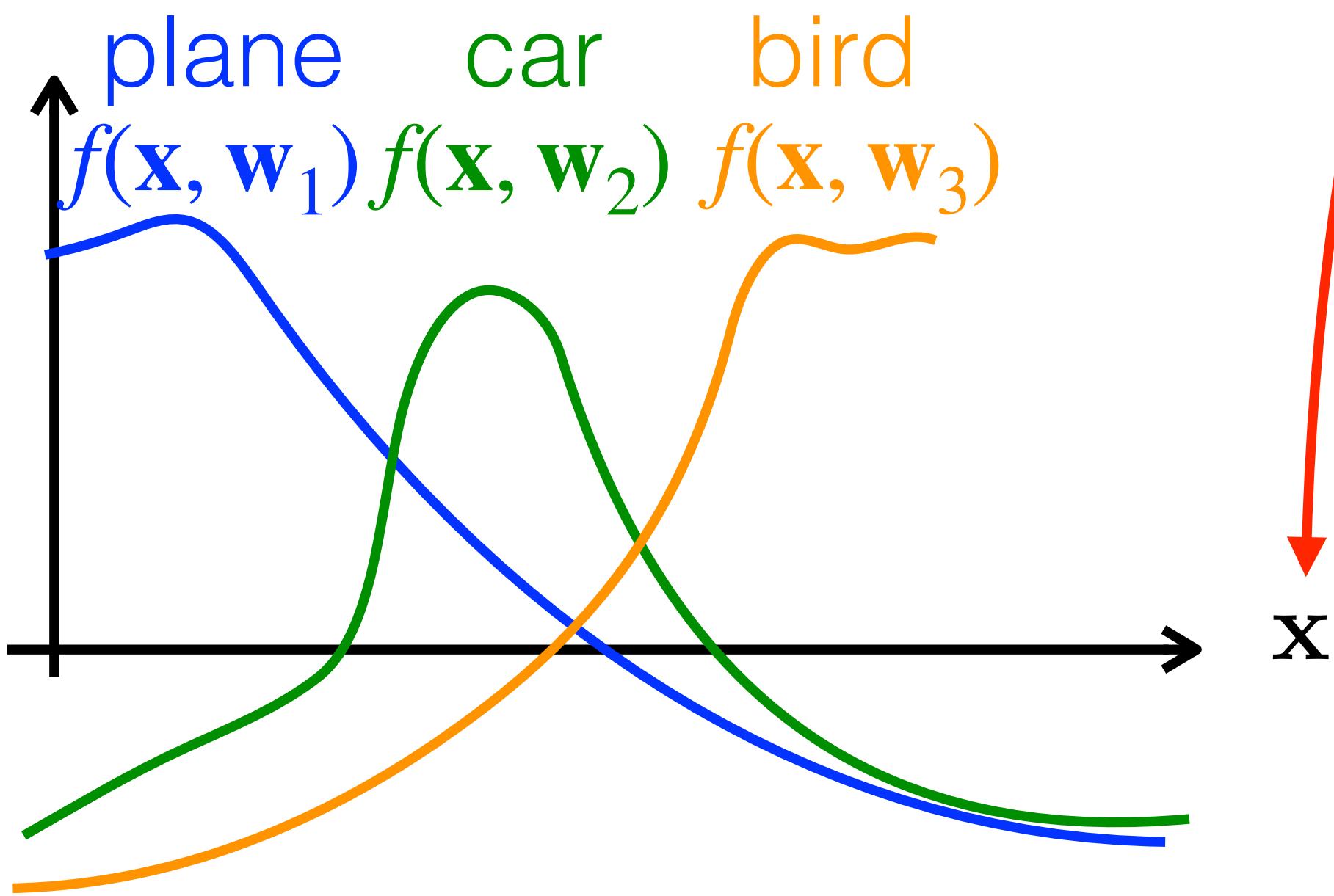
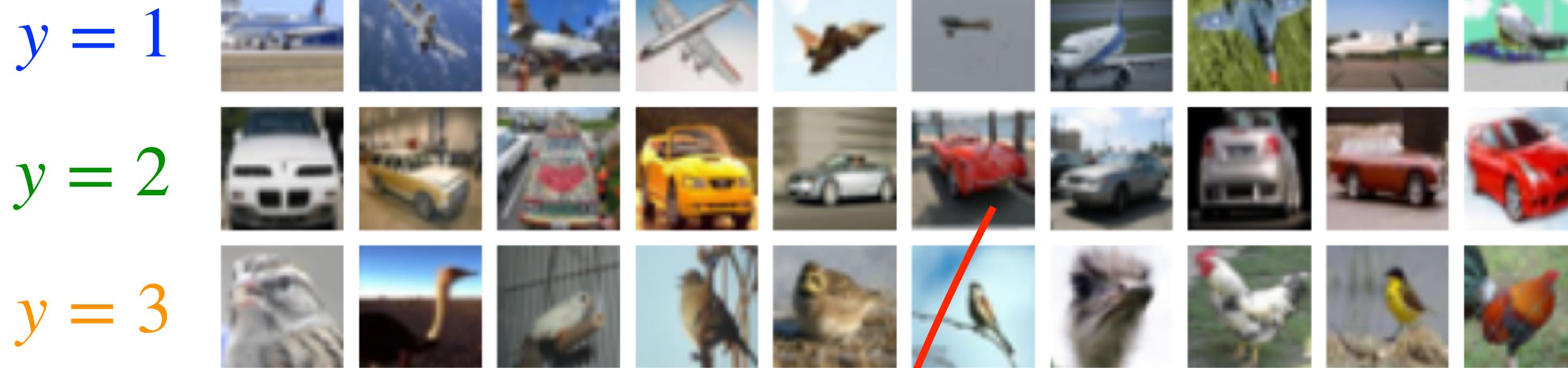
Financováno
Evropskou unií
NextGenerationEU

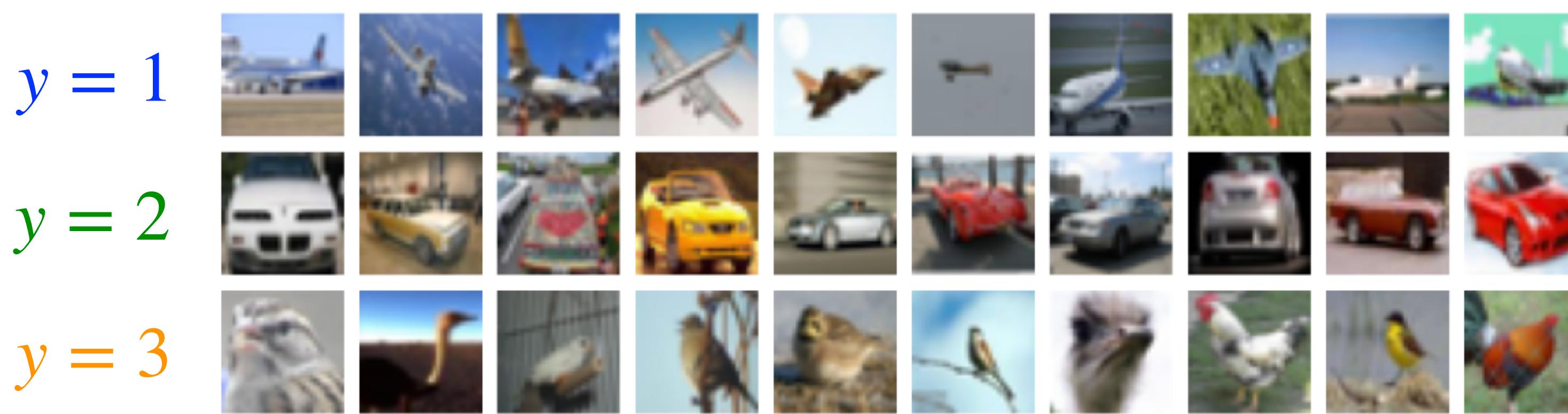


Recognition problem

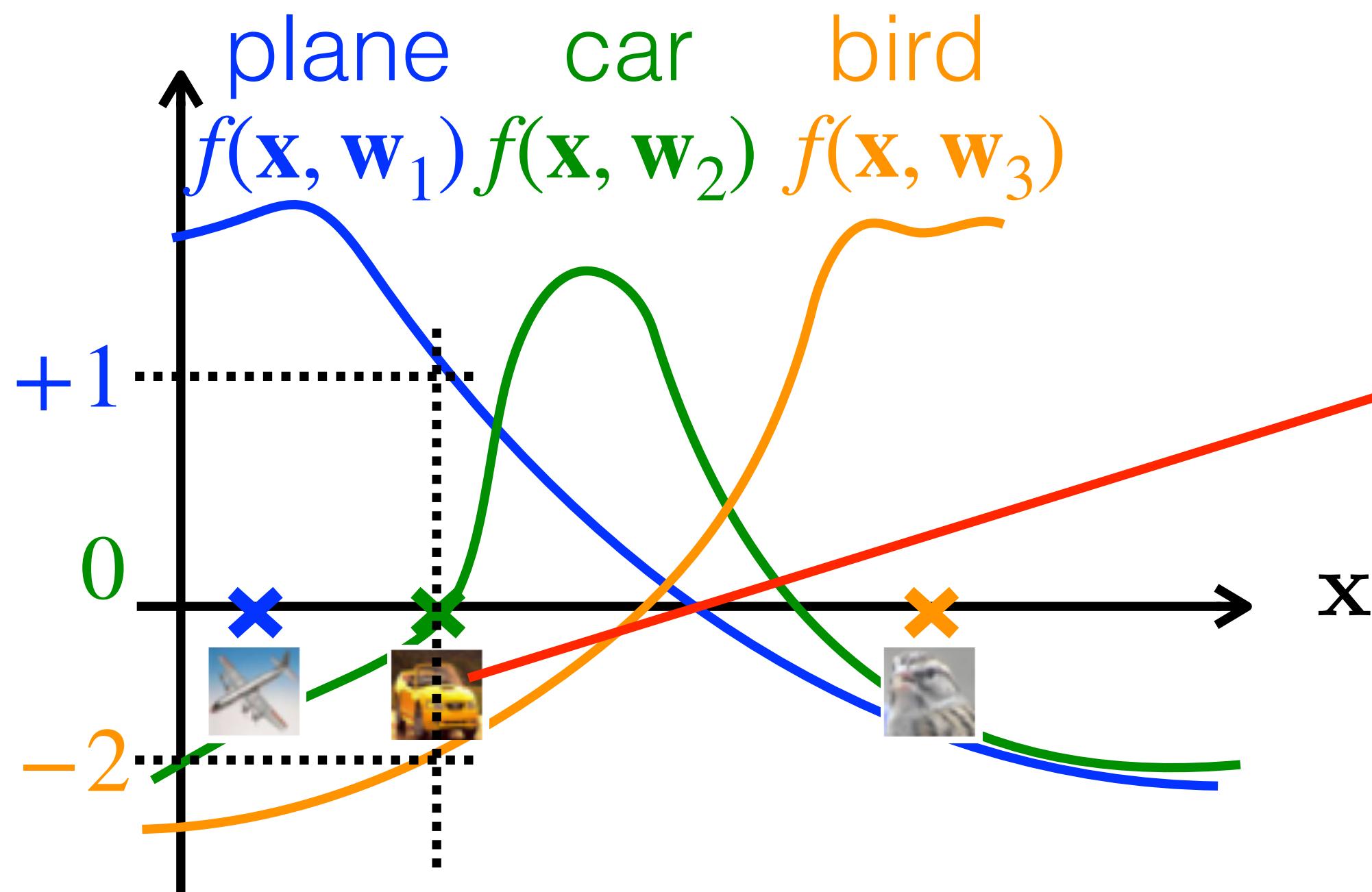


CIFAR-10: classify 32x32 RGB images into 10 categories
<https://www.cs.toronto.edu/~kriz/cifar.html>





$$p(y | \mathbf{x}, \mathbf{W}) = \text{???}$$



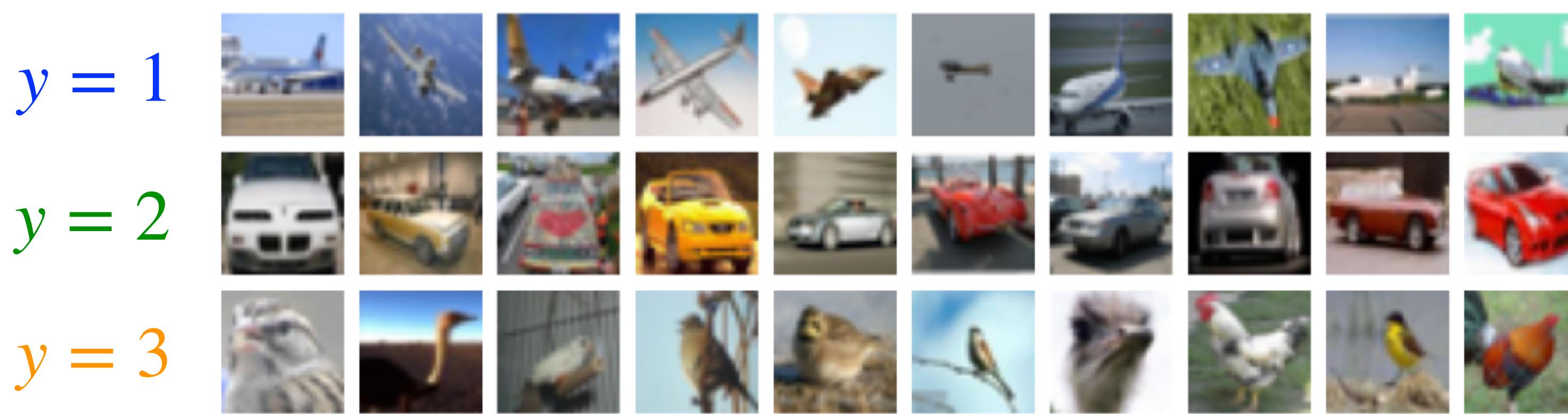
Linear
classifier

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

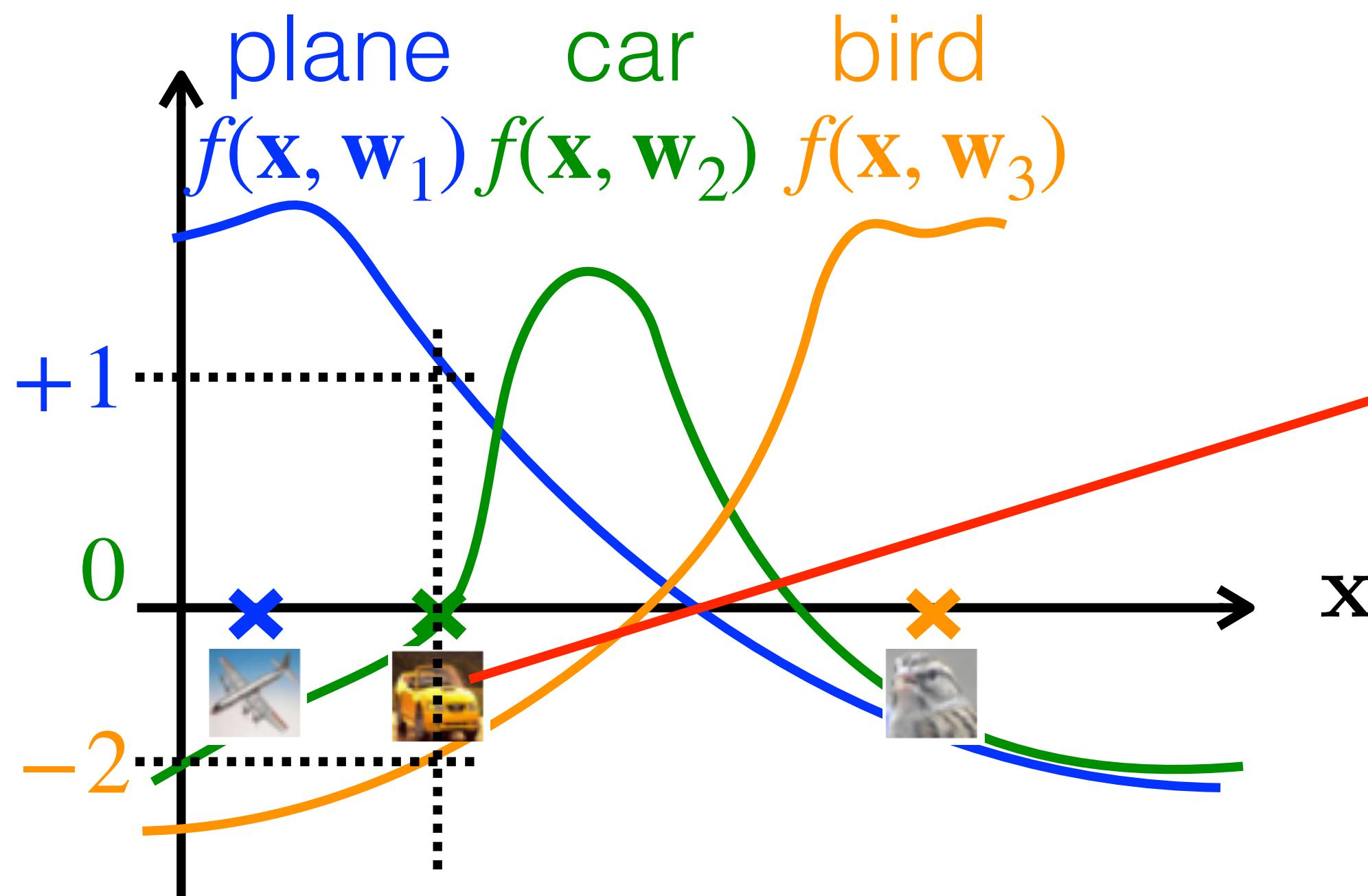
$$\mathbf{f}(\mathbf{x}, \mathbf{W}) = \boxed{\mathbf{W}} \quad \boxed{\mathbf{x}}$$

$$\text{logits} = \begin{bmatrix} +1 \\ 0 \\ -2 \end{bmatrix}$$

How does the linear classifier work?



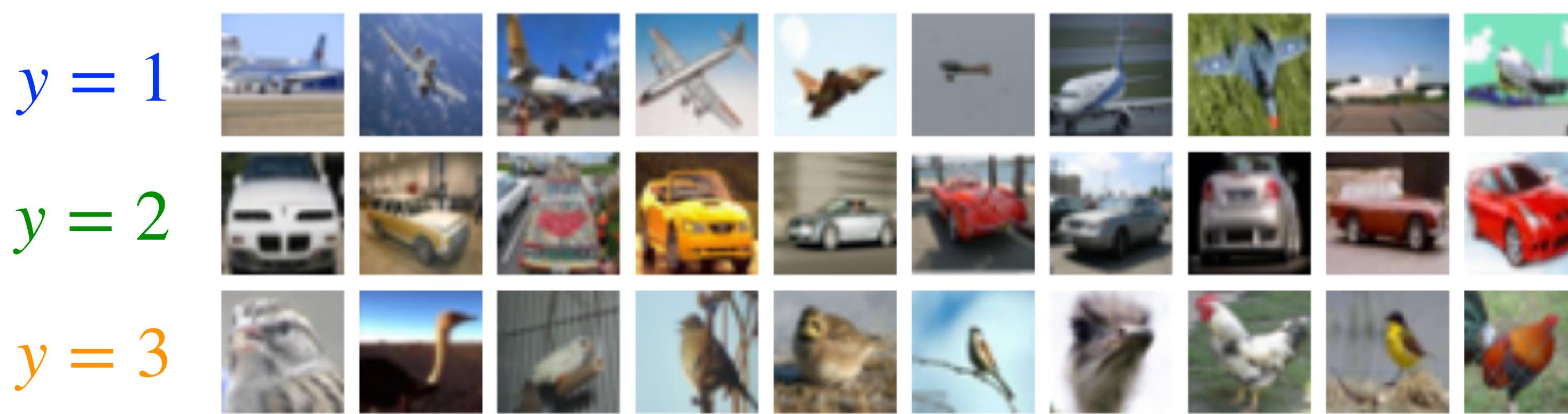
$$p(y|x, \mathbf{W}) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k)} = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) = \mathbf{s}\left(\begin{bmatrix} +1 \\ 0 \\ -2 \end{bmatrix}\right) = \begin{bmatrix} 0.71 \\ 0.26 \\ 0.03 \end{bmatrix}$$



Linear
classifier

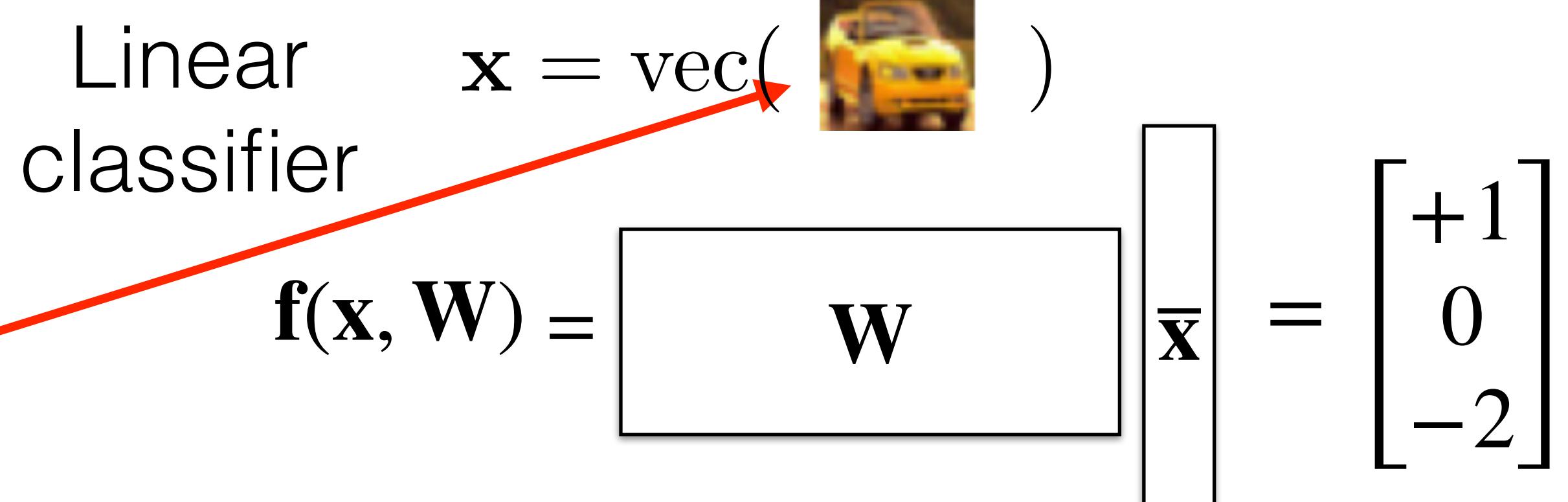
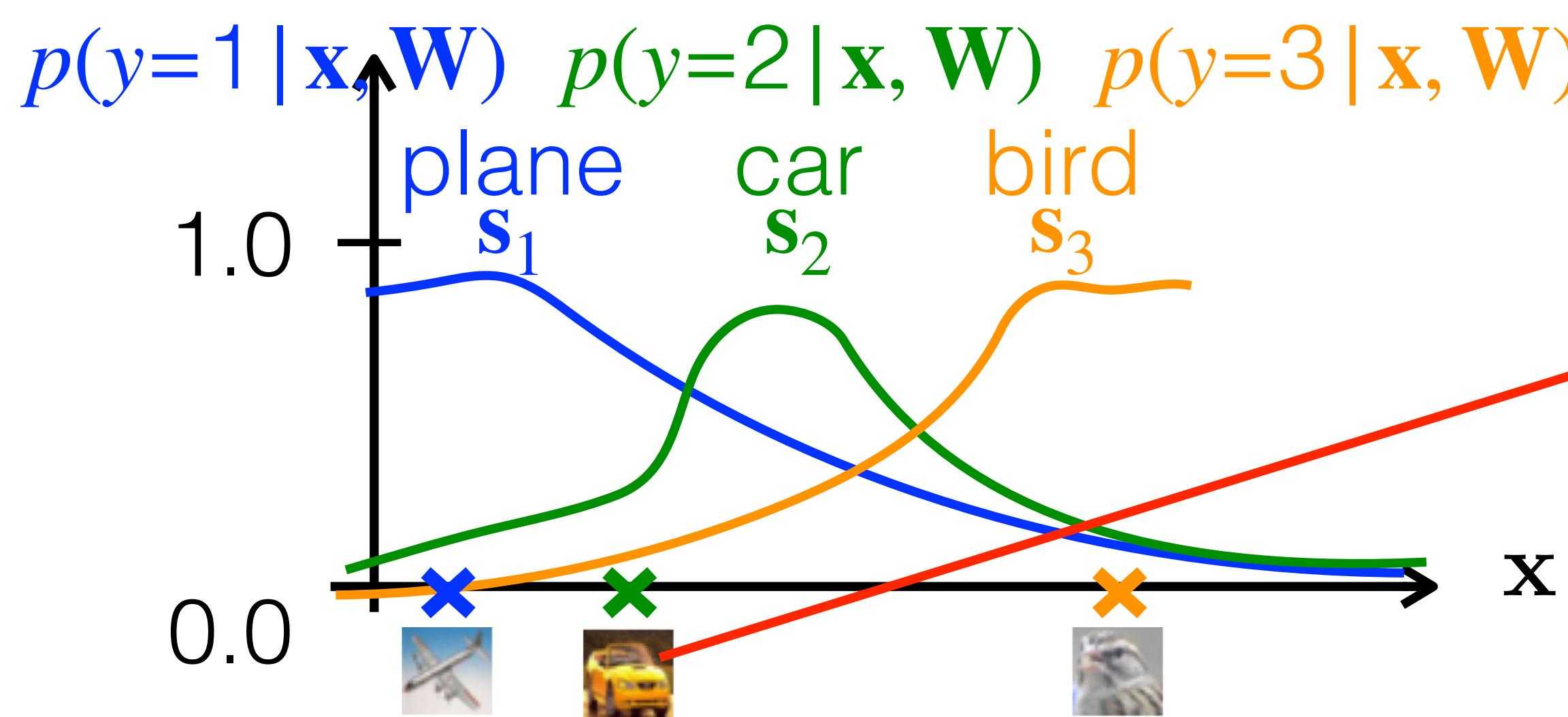
$\mathbf{x} = \text{vec}()$

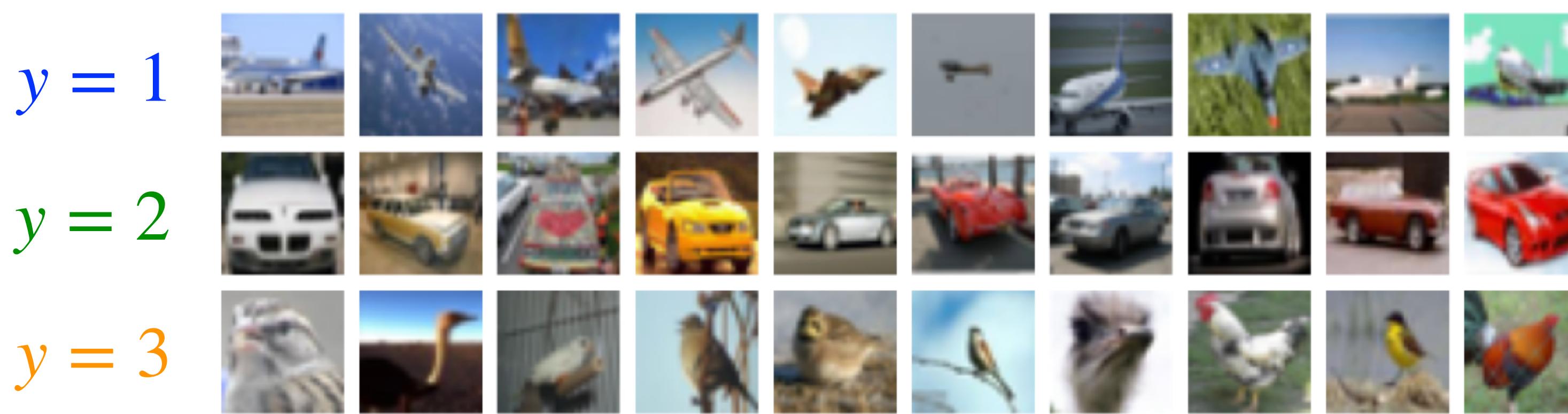
$\mathbf{f}(\mathbf{x}, \mathbf{W}) = \boxed{\mathbf{W}} \quad \boxed{\bar{\mathbf{x}}} = \begin{bmatrix} +1 \\ 0 \\ -2 \end{bmatrix}$



Model probability distribution over classes by softmax function

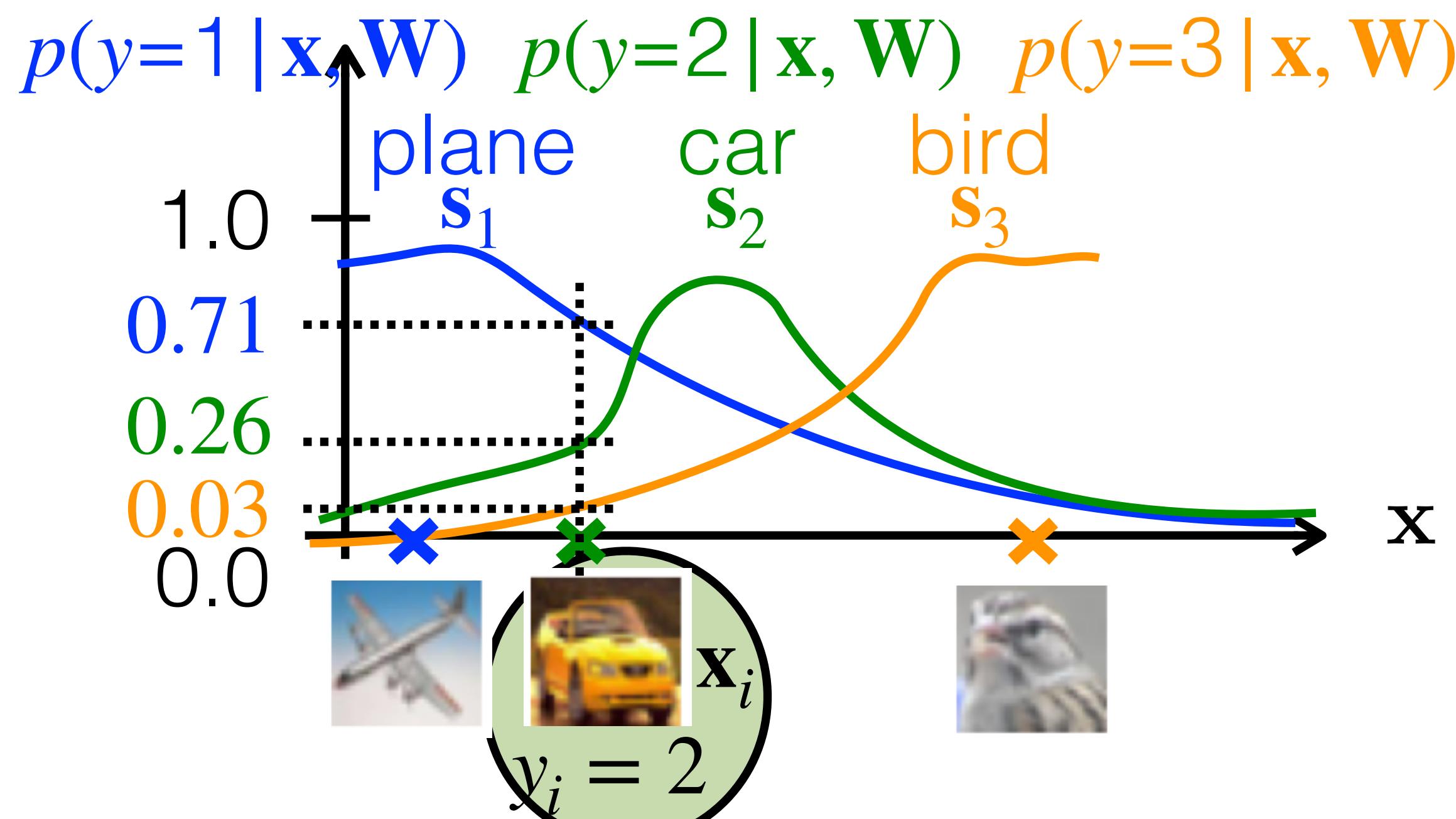
$$p(y|x, \mathbf{W}) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k)} = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) = \mathbf{s}\left(\begin{bmatrix} +1 \\ 0 \\ -2 \end{bmatrix}\right) = \begin{bmatrix} 0.71 \\ 0.26 \\ 0.03 \end{bmatrix}$$





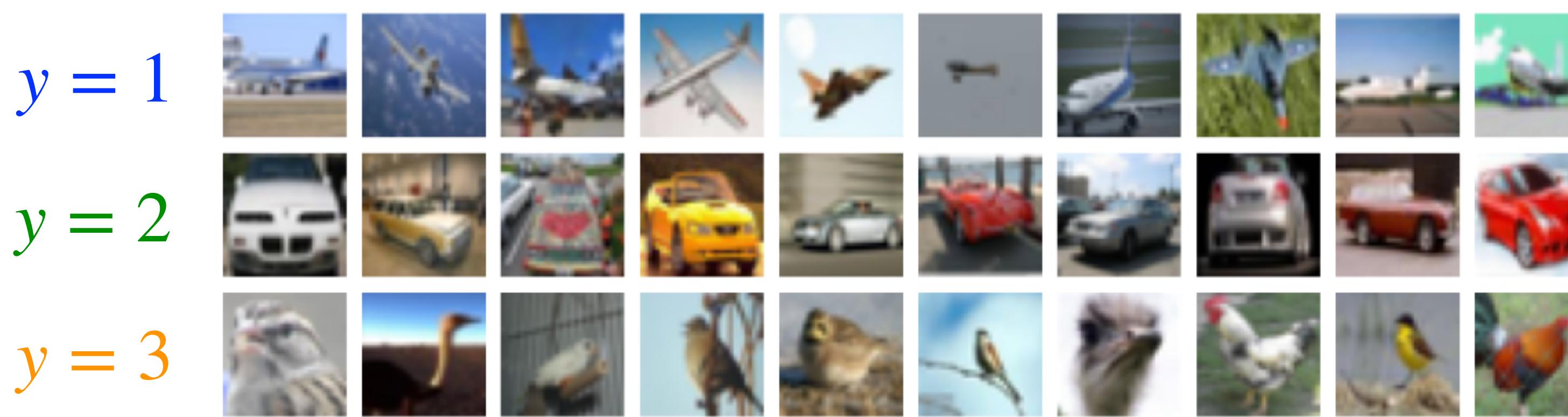
Model probability distribution over classes by softmax function

$$p(y|x, W) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k)} = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) = \mathbf{s}\left(\begin{bmatrix} +1 \\ 0 \\ -2 \end{bmatrix}\right) = \begin{bmatrix} 0.71 \\ 0.26 \\ 0.03 \end{bmatrix}$$



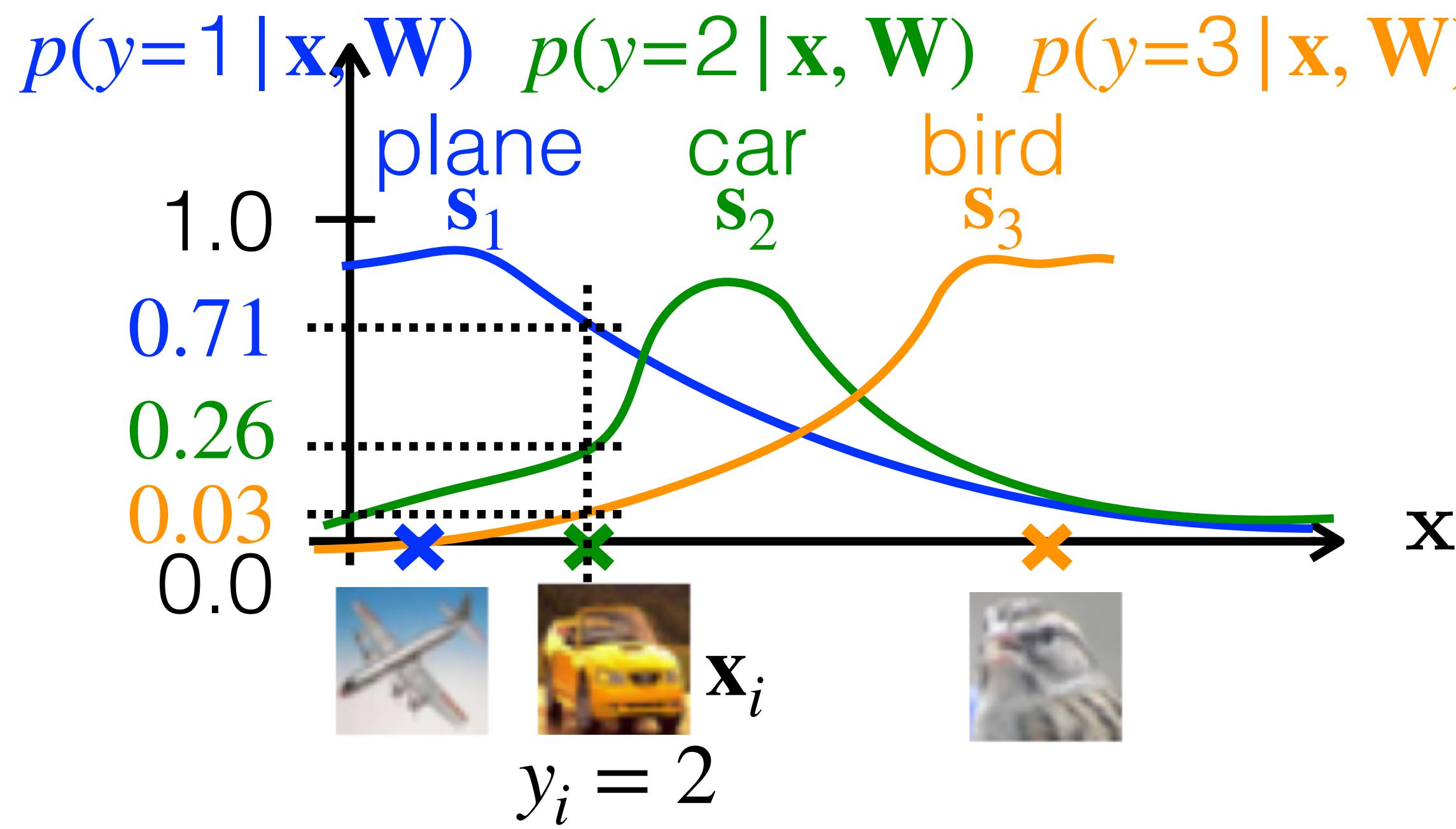
Which image is incorrectly classified?





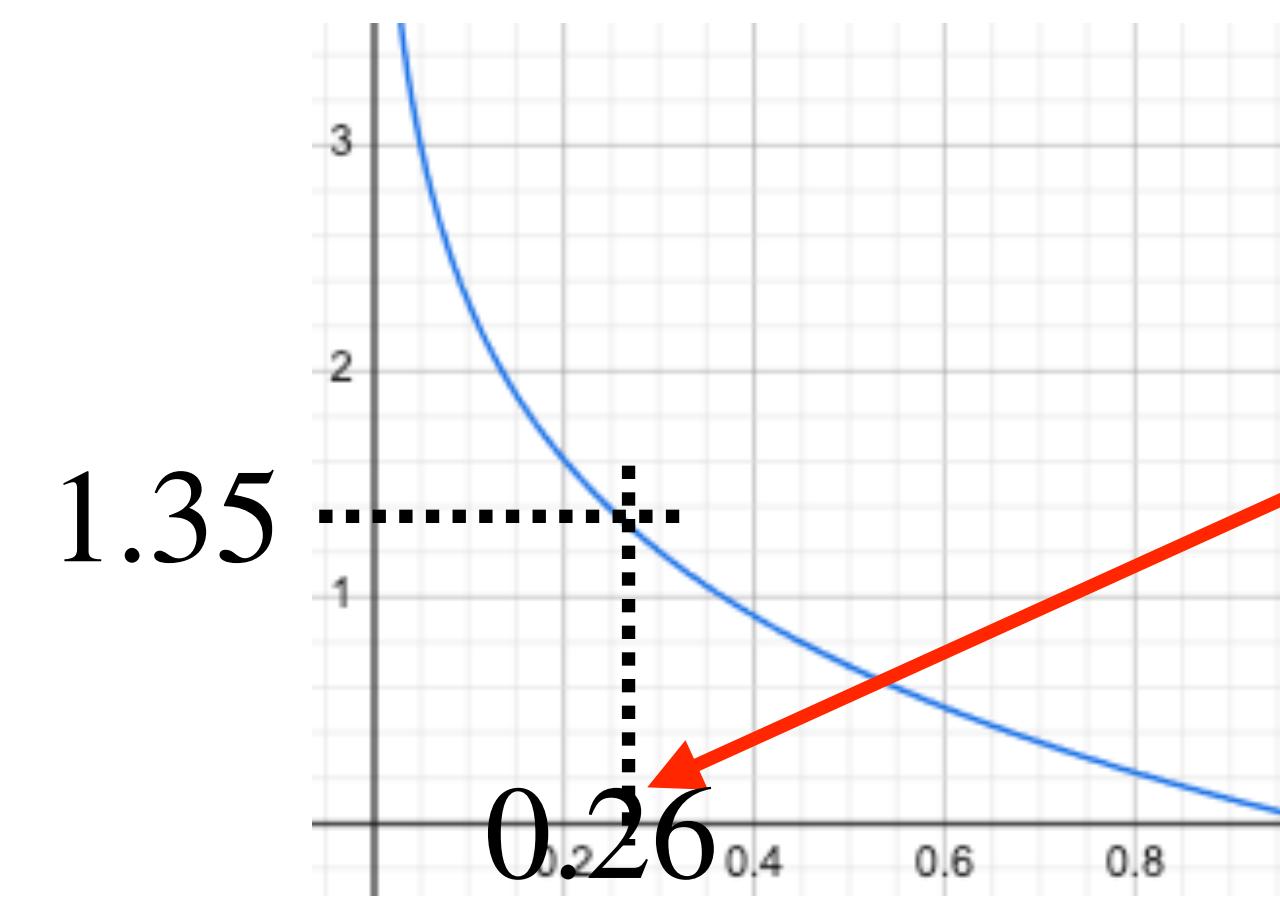
Model probability distribution over classes by softmax function

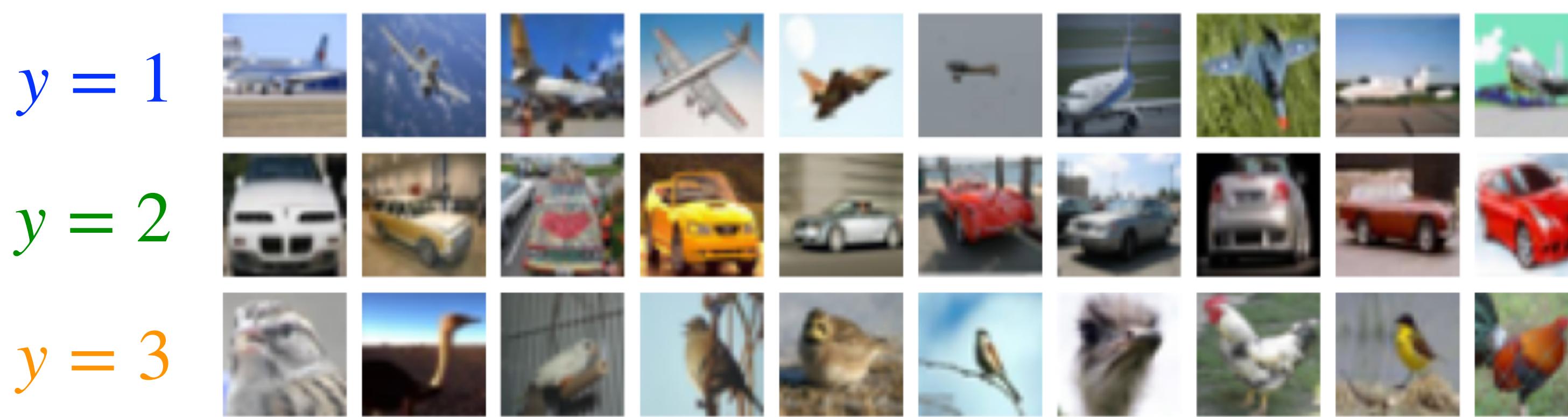
$$p(y|\mathbf{x}, \mathbf{W}) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k)} = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) = \mathbf{s}\left(\begin{bmatrix} +1 \\ 0 \\ -2 \end{bmatrix}\right) = \begin{bmatrix} 0.71 \\ 0.26 \\ 0.03 \end{bmatrix}$$



Loss function:

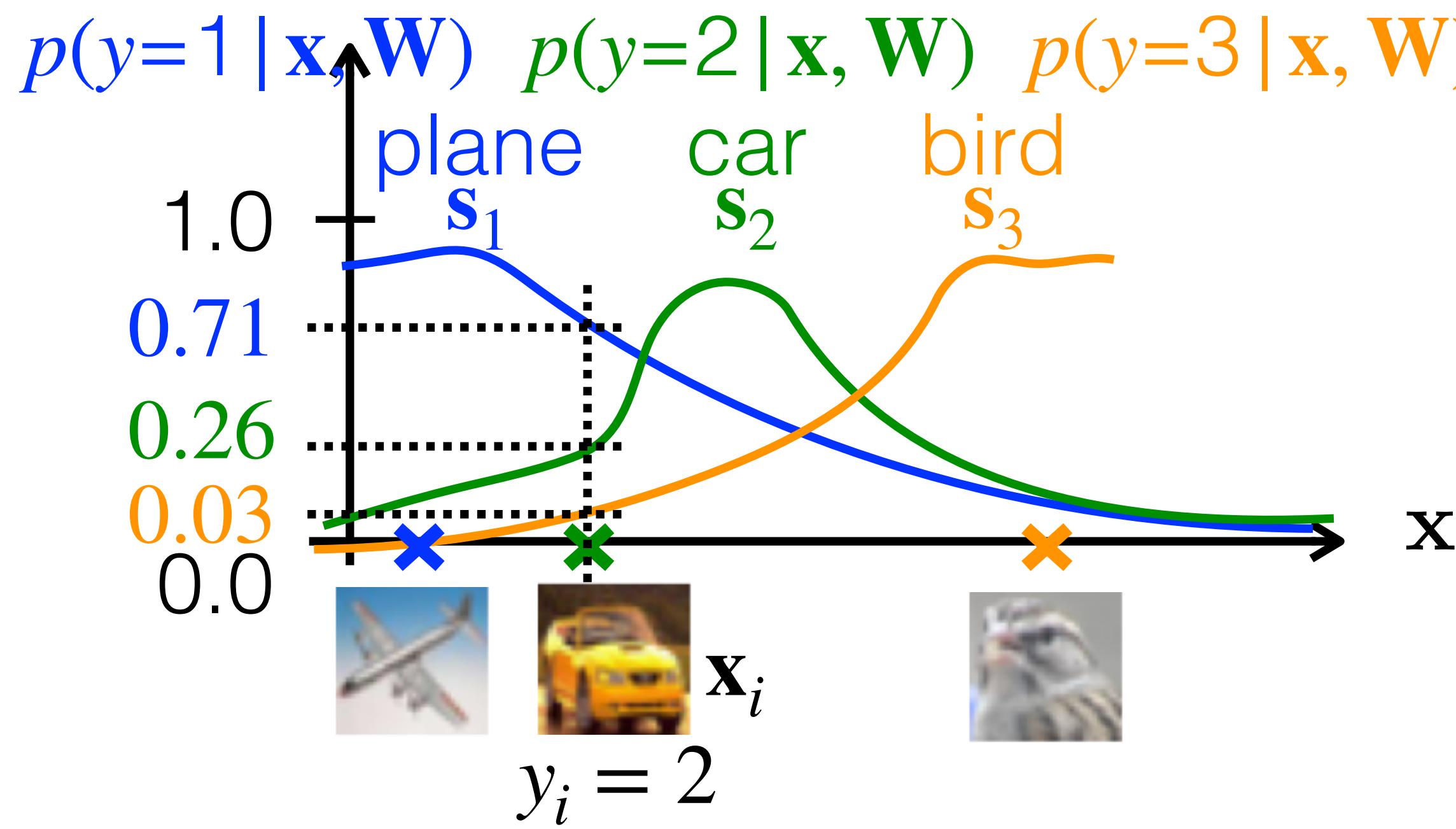
$$\begin{aligned} \mathcal{L}(\mathbf{W}) &= -\log p(y = y_i | \mathbf{x}_i, \mathbf{W}) = -\log \mathbf{s}_{y_i}(\mathbf{W} \bar{\mathbf{x}}_i) \\ &= -\log(0.26) = 1.35 \end{aligned}$$





Model probability distribution over classes by softmax function

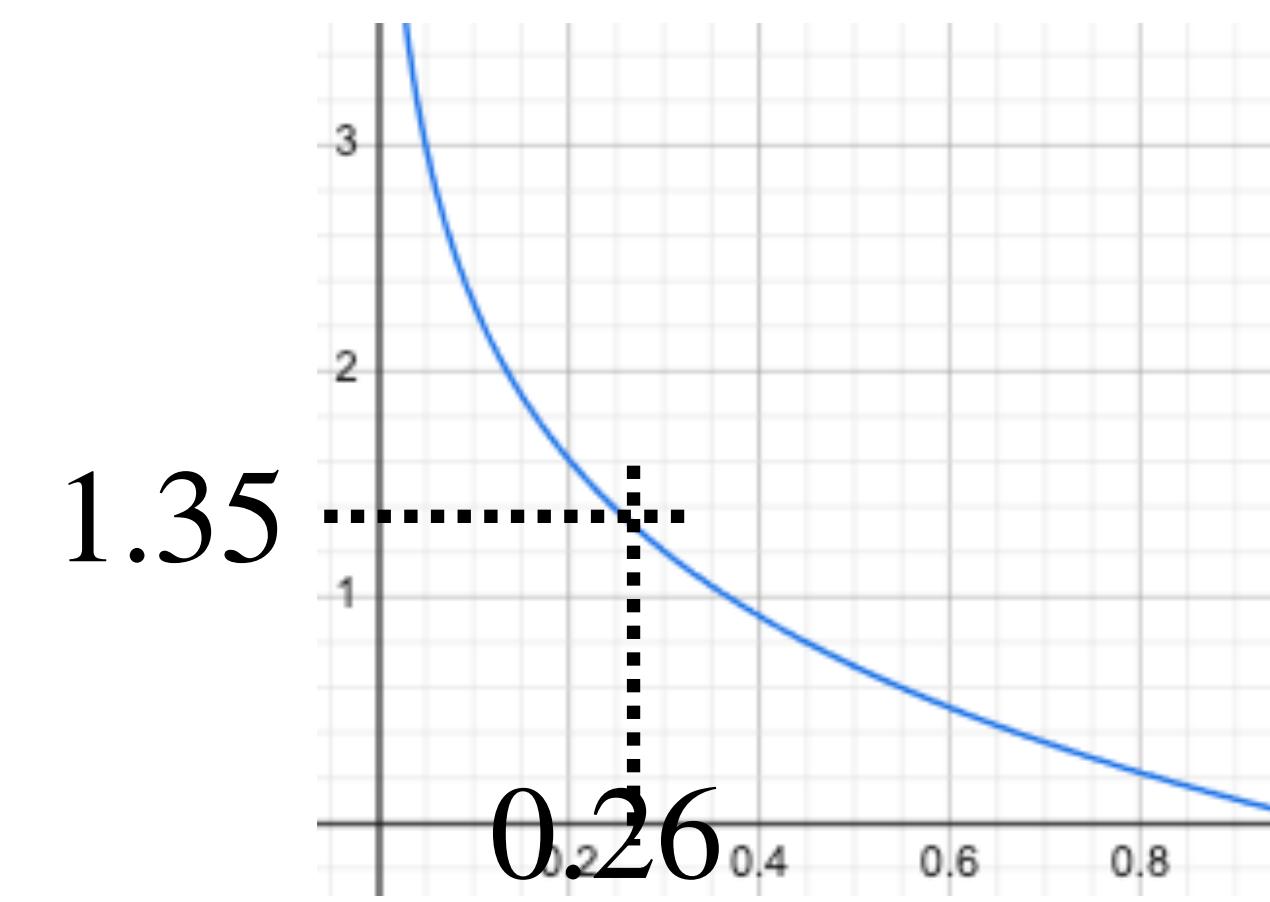
$$p(y|\mathbf{x}, \mathbf{W}) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k)} = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) = \mathbf{s}\left(\begin{bmatrix} +1 \\ 0 \\ -2 \end{bmatrix}\right) = \begin{bmatrix} 0.71 \\ 0.26 \\ 0.03 \end{bmatrix}$$

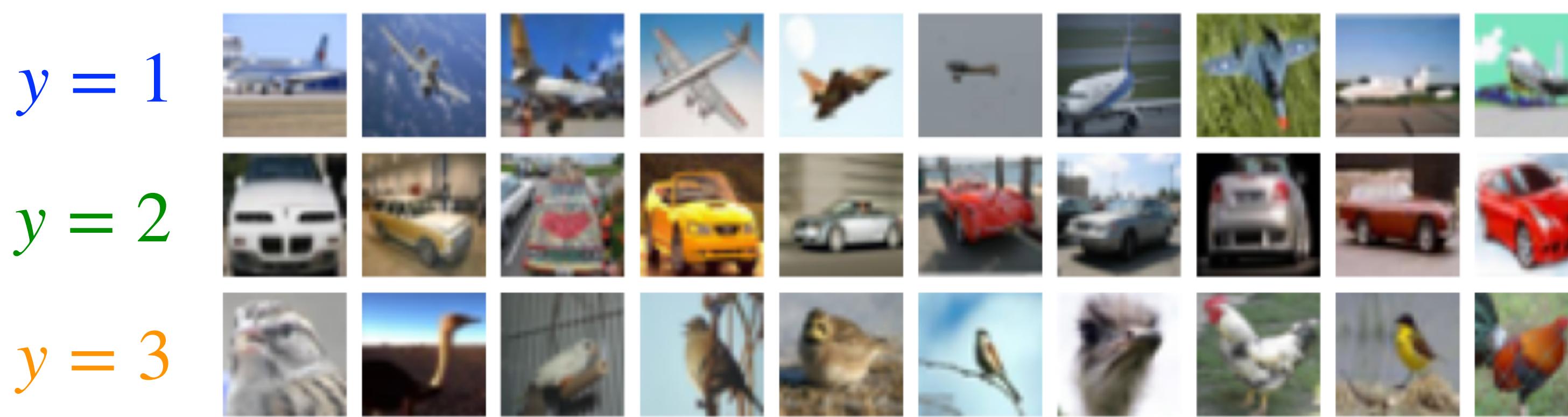


Loss function:

$$\mathcal{L}(\mathbf{W}) = -\log p(y = y_i | \mathbf{x}_i, \mathbf{W}) = -\log \mathbf{s}_{y_i}(\mathbf{W} \mathbf{x}_i)$$

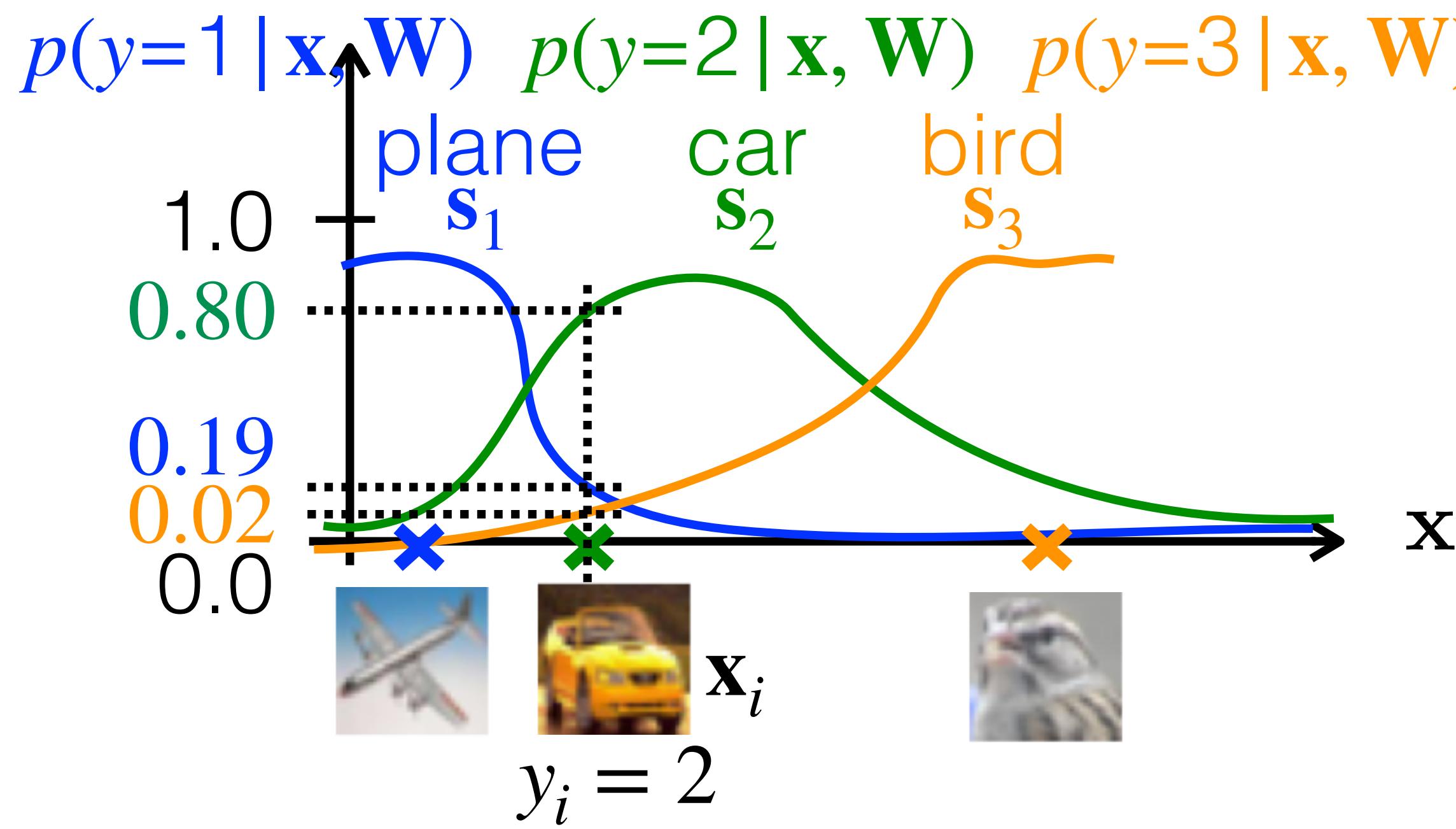
Learning:
 $\arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W})$





Model probability distribution over classes by softmax function

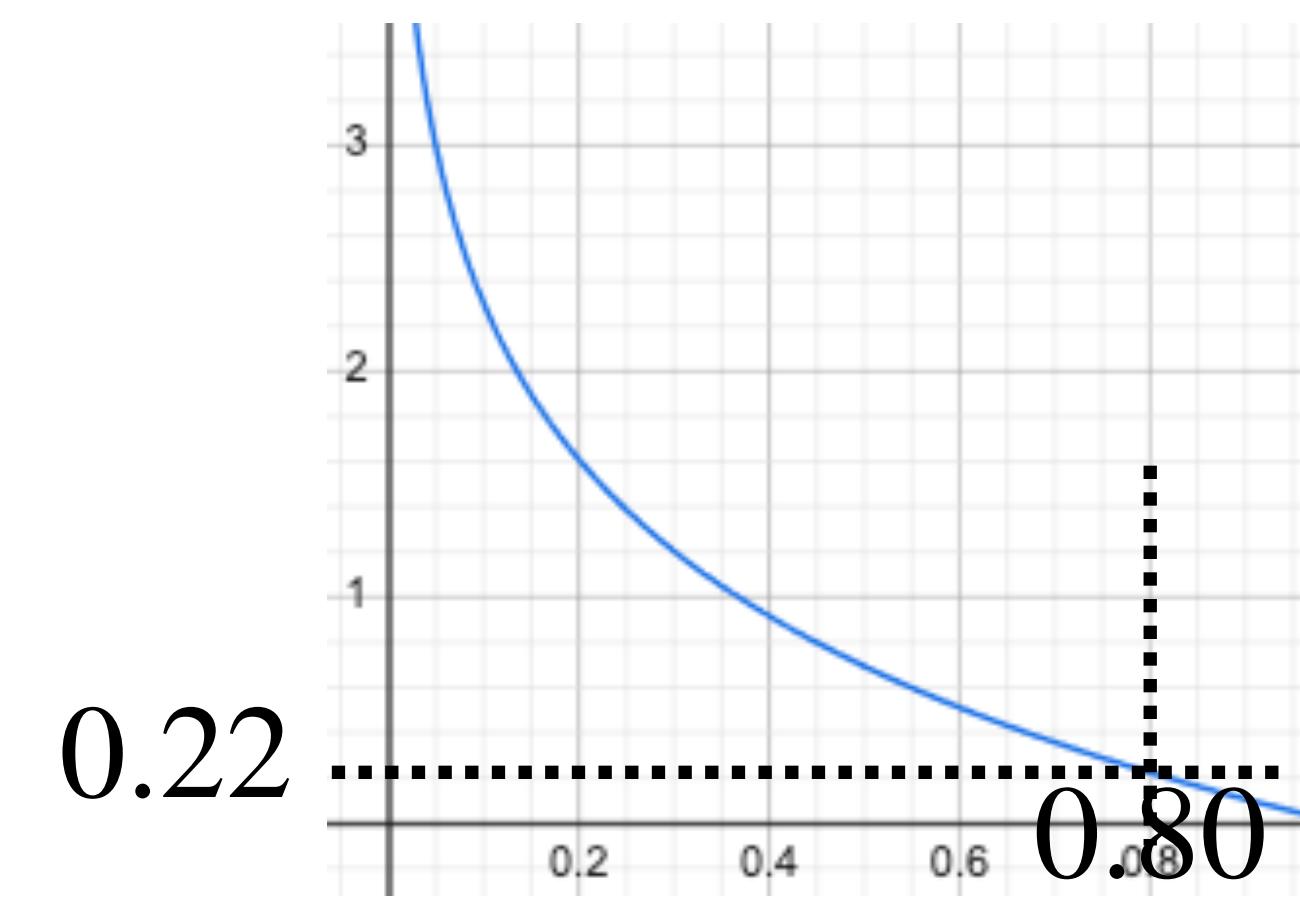
$$p(y|\mathbf{x}, \mathbf{W}) = \frac{\begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \\ \exp(f(\mathbf{x}, \mathbf{w}_3)) \end{bmatrix}}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k)} = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) = \mathbf{s}\left(\begin{bmatrix} +1 \\ 0 \\ -2 \end{bmatrix}\right) = \begin{bmatrix} 0.71 \\ 0.26 \\ 0.03 \end{bmatrix}$$



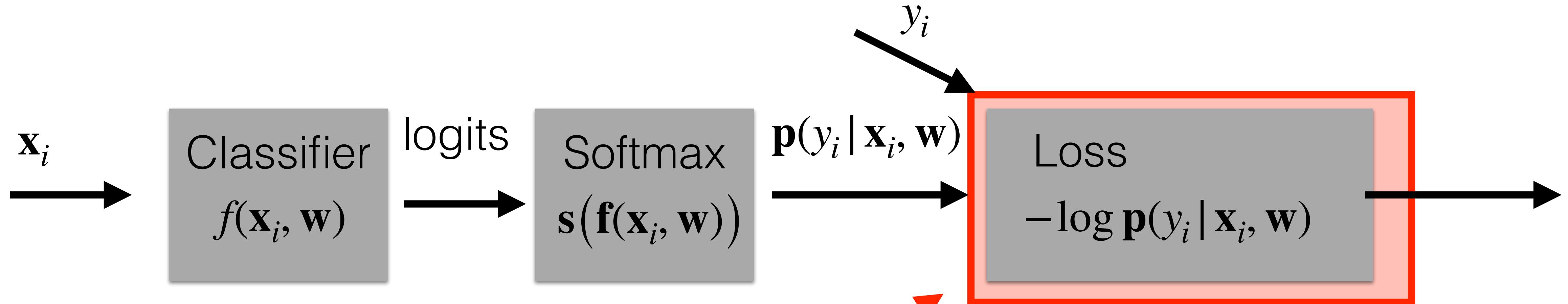
Loss function:

$$\mathcal{L}(\mathbf{W}) = -\log p(y = y_i | \mathbf{x}_i, \mathbf{W}) = -\log s_{y_i}(\mathbf{W} \mathbf{x}_i)$$

Learning:
 $\arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W})$



Various implementation and names



Input: probabilities

(Multinomial) Logistic loss,
Cross-entropy loss, ...

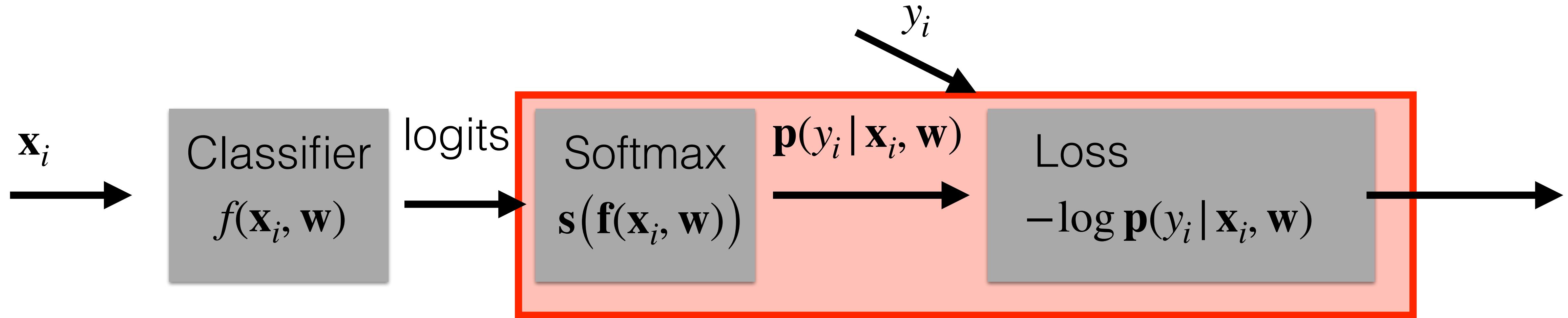
$$\mathcal{L}(w) = CE(y_i, \mathbf{p}(y_i | \mathbf{x}_i, w))$$

Pytorch: [BCELoss](#)

TensorFlow: [log_loss](#)

Caffe: [Multinomial Logistic Loss Layer](#)

Various implementation and names



Input: probabilities

(Multinomial) Logistic loss,
Cross-entropy loss, ...

$$\mathcal{L}(\mathbf{w}) = CE(y_i, \mathbf{p}(y_i | \mathbf{x}_i, \mathbf{w}))$$

Pytorch: [NLLLoss](#)

TensorFlow: [log_loss](#)

Caffe: [Multinomial Logistic Loss Layer](#)

Input: logits

Softmax loss,
Categorical Cross-entropy loss, ...

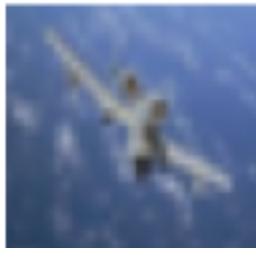
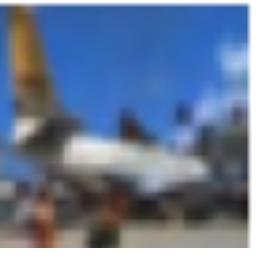
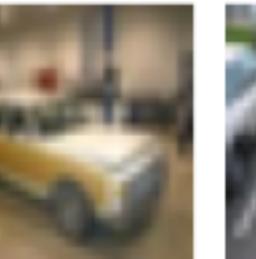
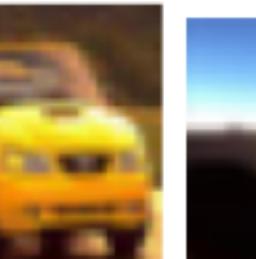
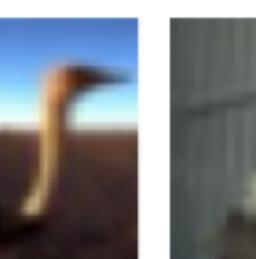
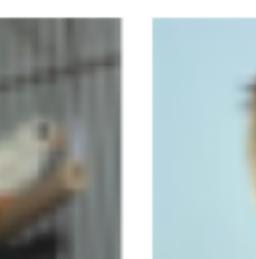
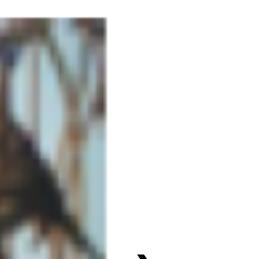
$$\mathcal{L}(\mathbf{w}) = CE(y_i, f(\mathbf{x}_i, \mathbf{w}))$$

Pytorch: [CrossEntropyLoss](#)

TensorFlow: [softmax_cross_entropy](#)

Caffe: [SoftmaxWithLoss Layer](#)

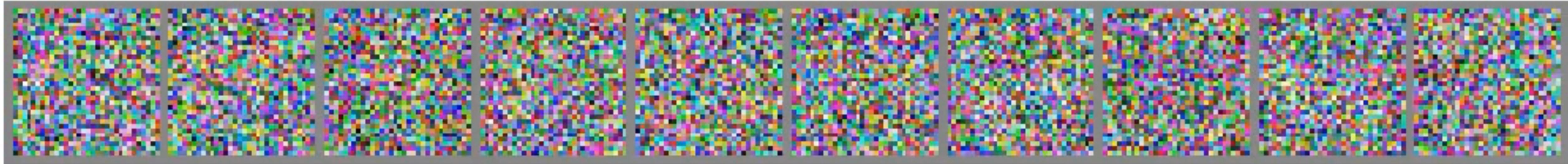
```

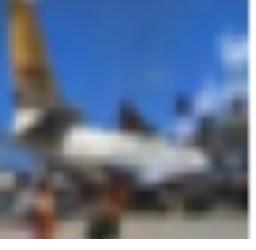
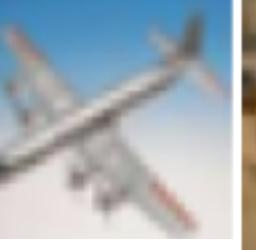
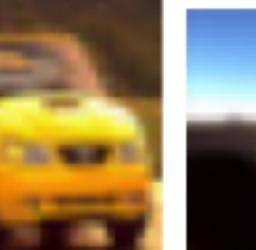
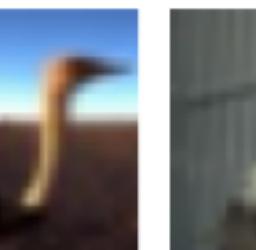
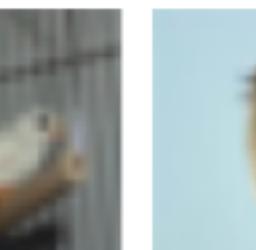
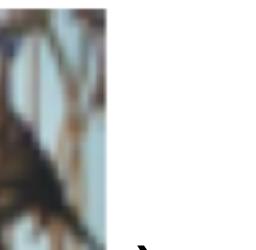
def train(  1  1  1  2  2  2  3  3  3 ):
     $\mathbf{x}_i = \text{vec}(\begin{matrix} \text{airplane} \end{matrix})$ 

 $\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}) = \arg \min_{\mathbf{W}} \sum_i -\log s_{y_i}(\mathbf{W} \bar{\mathbf{x}}_i)$  ... gradient optimization
    return  $\mathbf{W}^*$ 

```

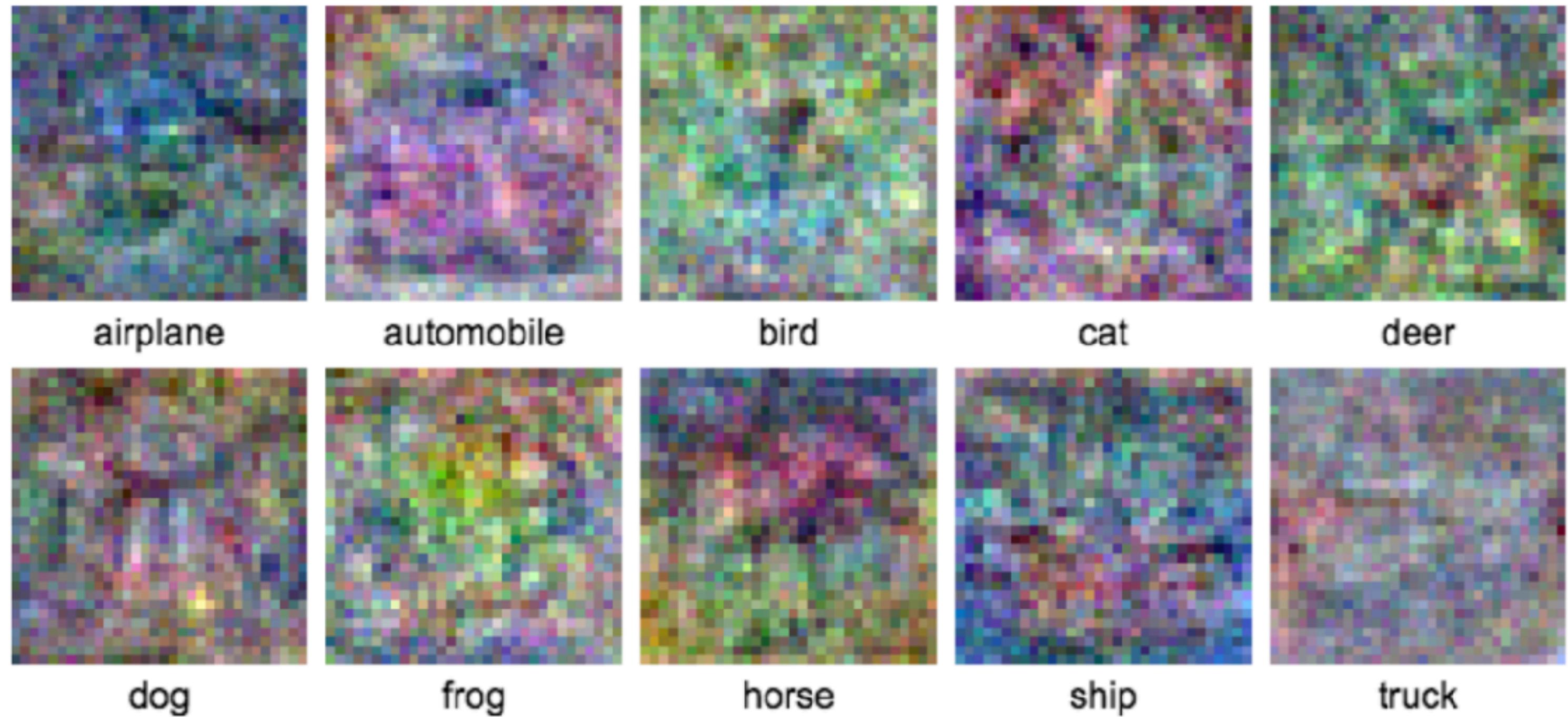
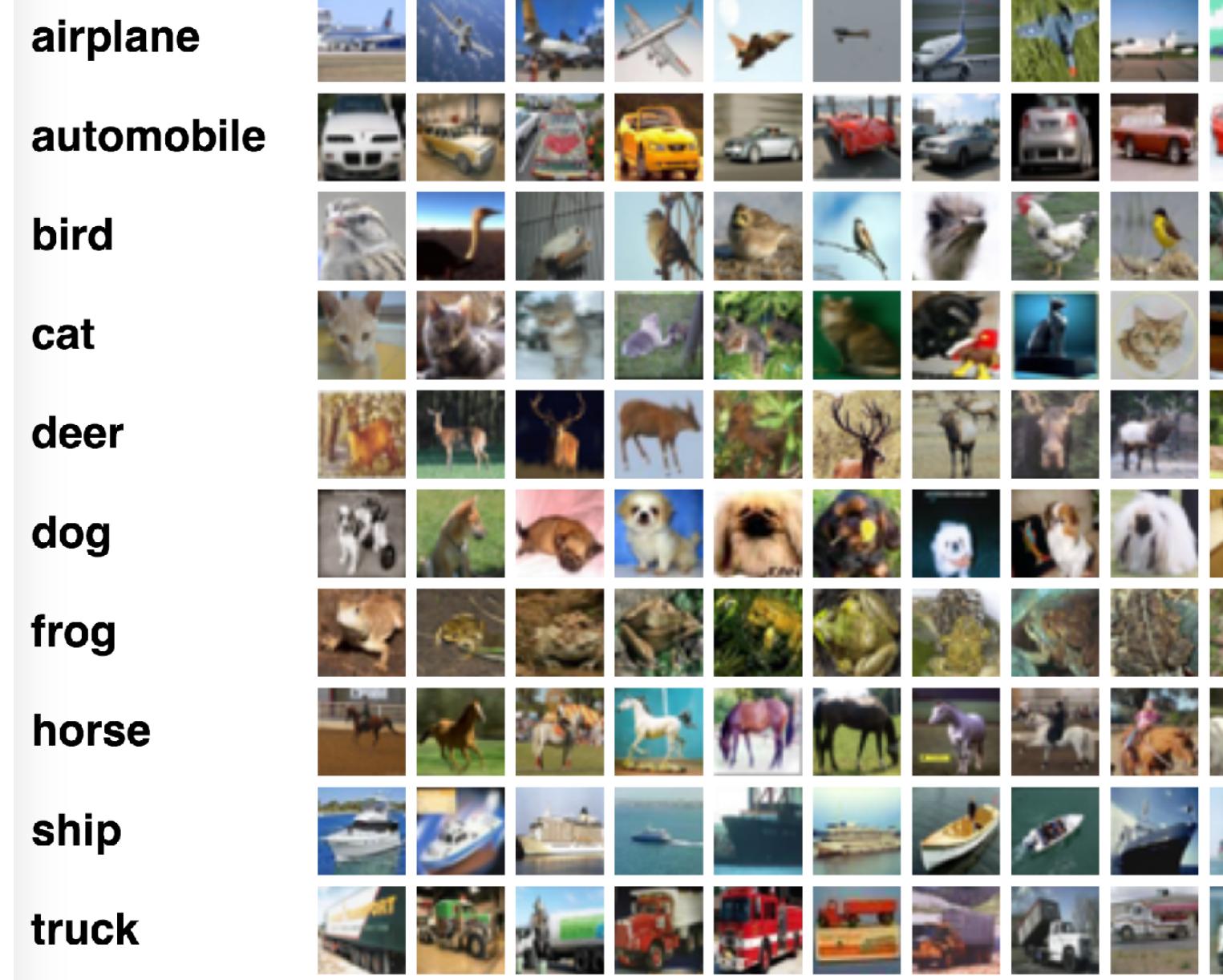
reshape(\mathbf{w}_1) reshape(\mathbf{w}_2) reshape(\mathbf{w}_3) reshape(\mathbf{w}_4) reshape(\mathbf{w}_5) reshape(\mathbf{w}_6) reshape(\mathbf{w}_7) reshape(\mathbf{w}_8) reshape(\mathbf{w}_9) reshape(\mathbf{w}_{10})

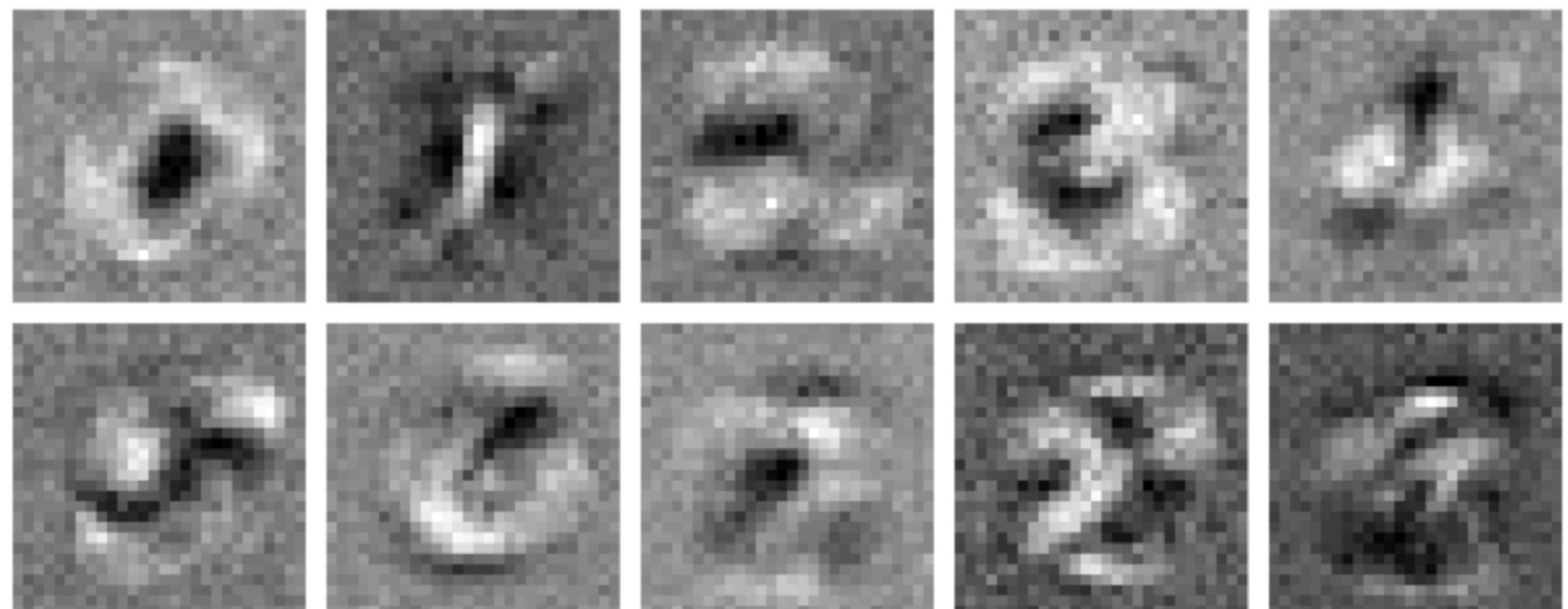
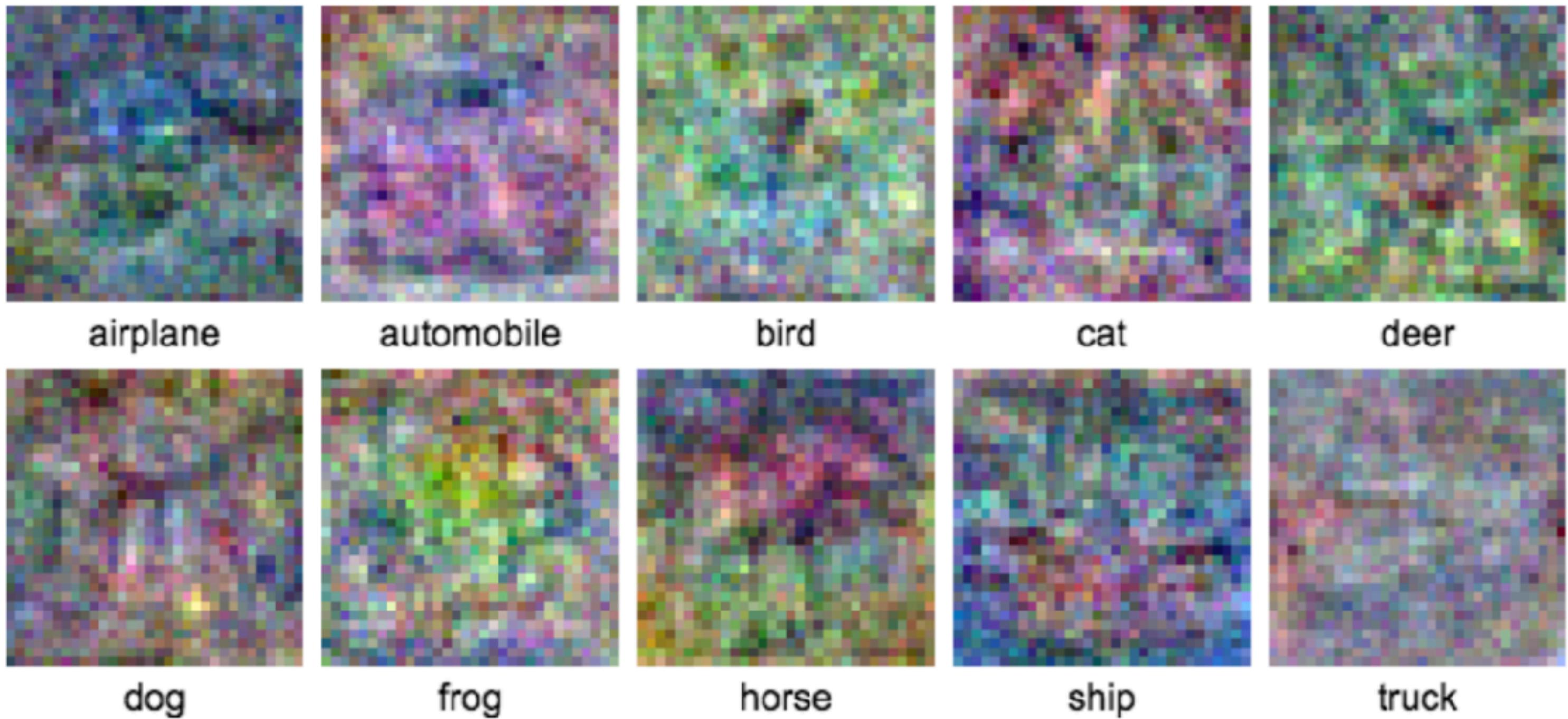
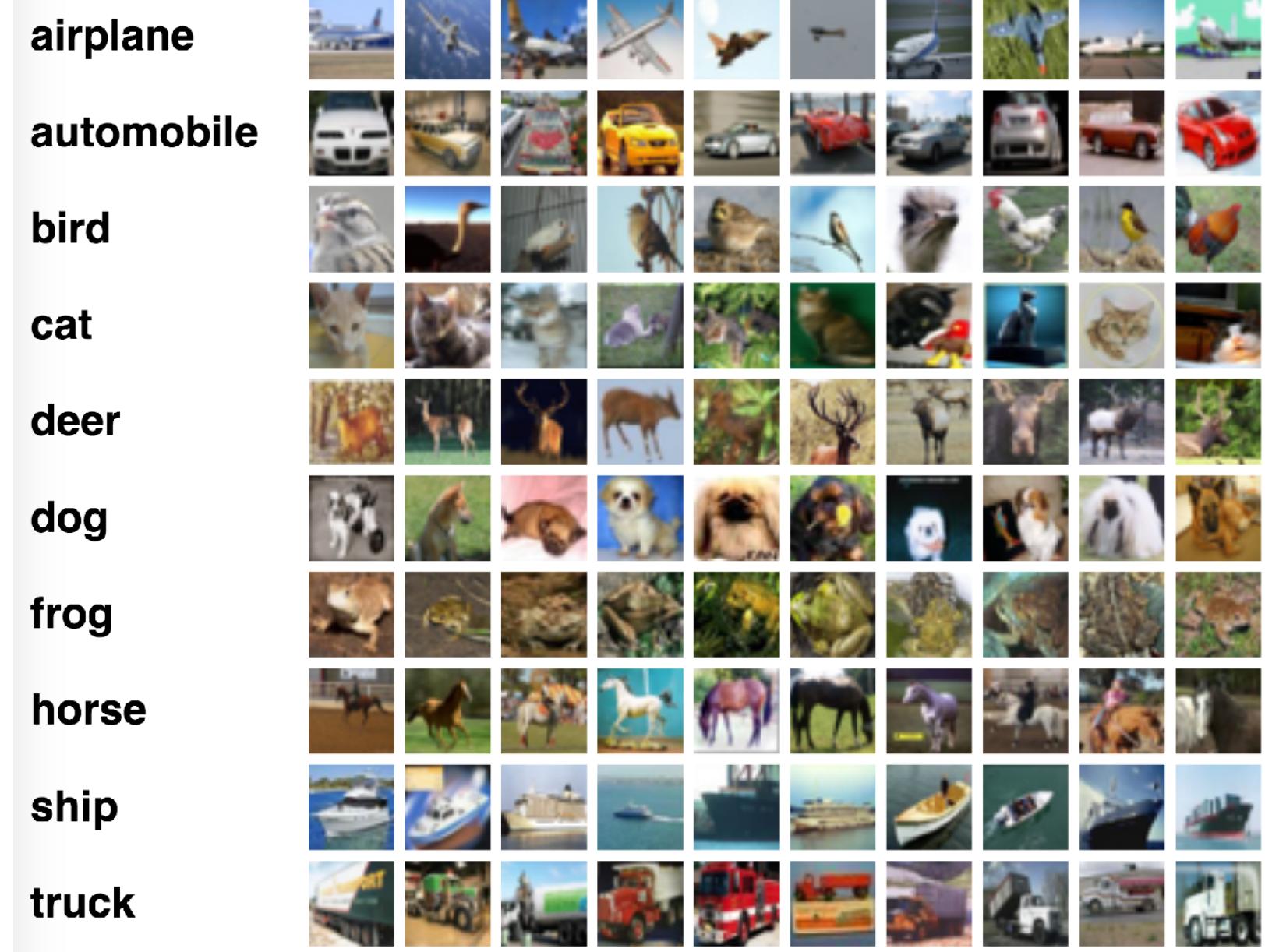


```
def train(  1  1  1  2  2  2  3  3  3 ):  

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

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}) = \arg \min_{\mathbf{W}} \sum_i -\log s_{y_i}(\mathbf{W} \bar{\mathbf{x}}_i) \dots \text{gradient optimization}$$
  
return  $\mathbf{W}^*$ 
```



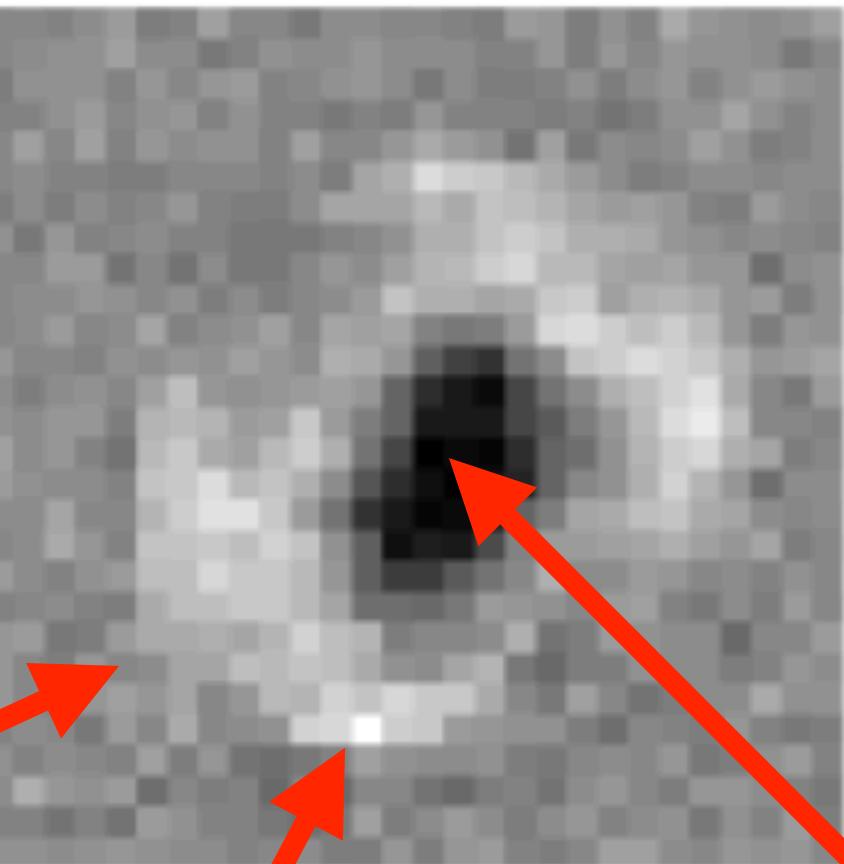


Σ

pixels

 \mathbf{x}  $*$

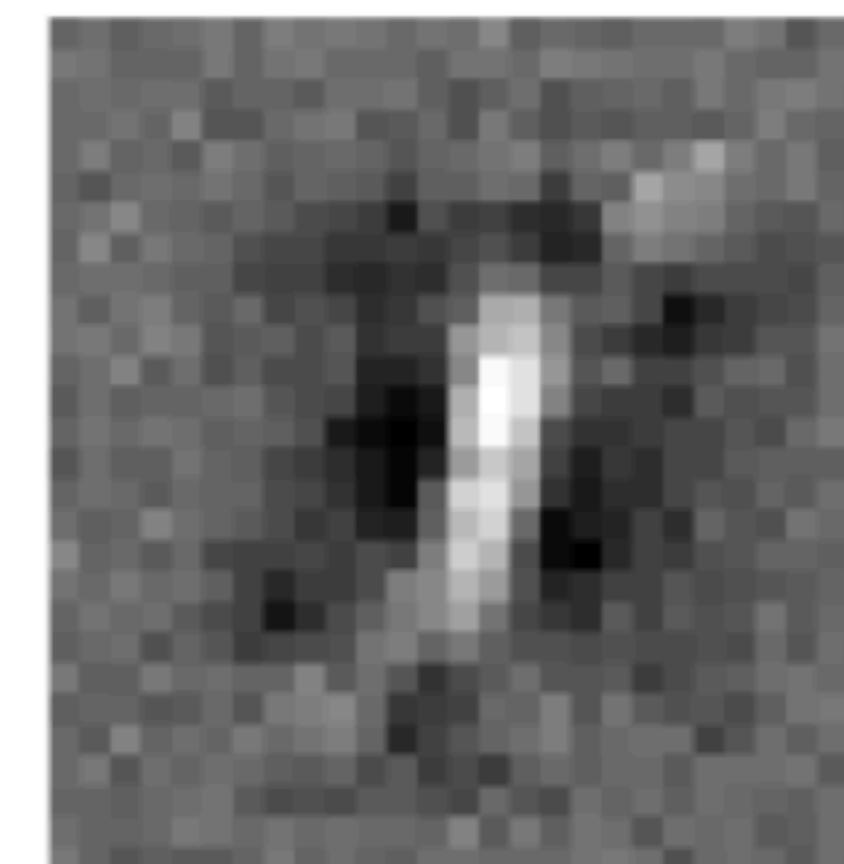
zeros

 \mathbf{w}_1  $=$

big number

 \square positive \blacksquare negative Σ

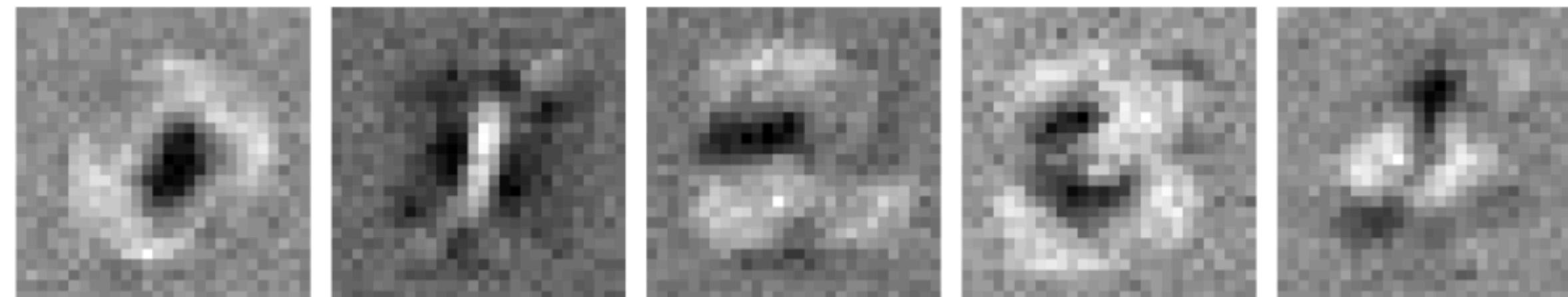
pixels

 \mathbf{x}  $*$  $=$

small number

Dataset Label Images Learned weights of linear classifier Error

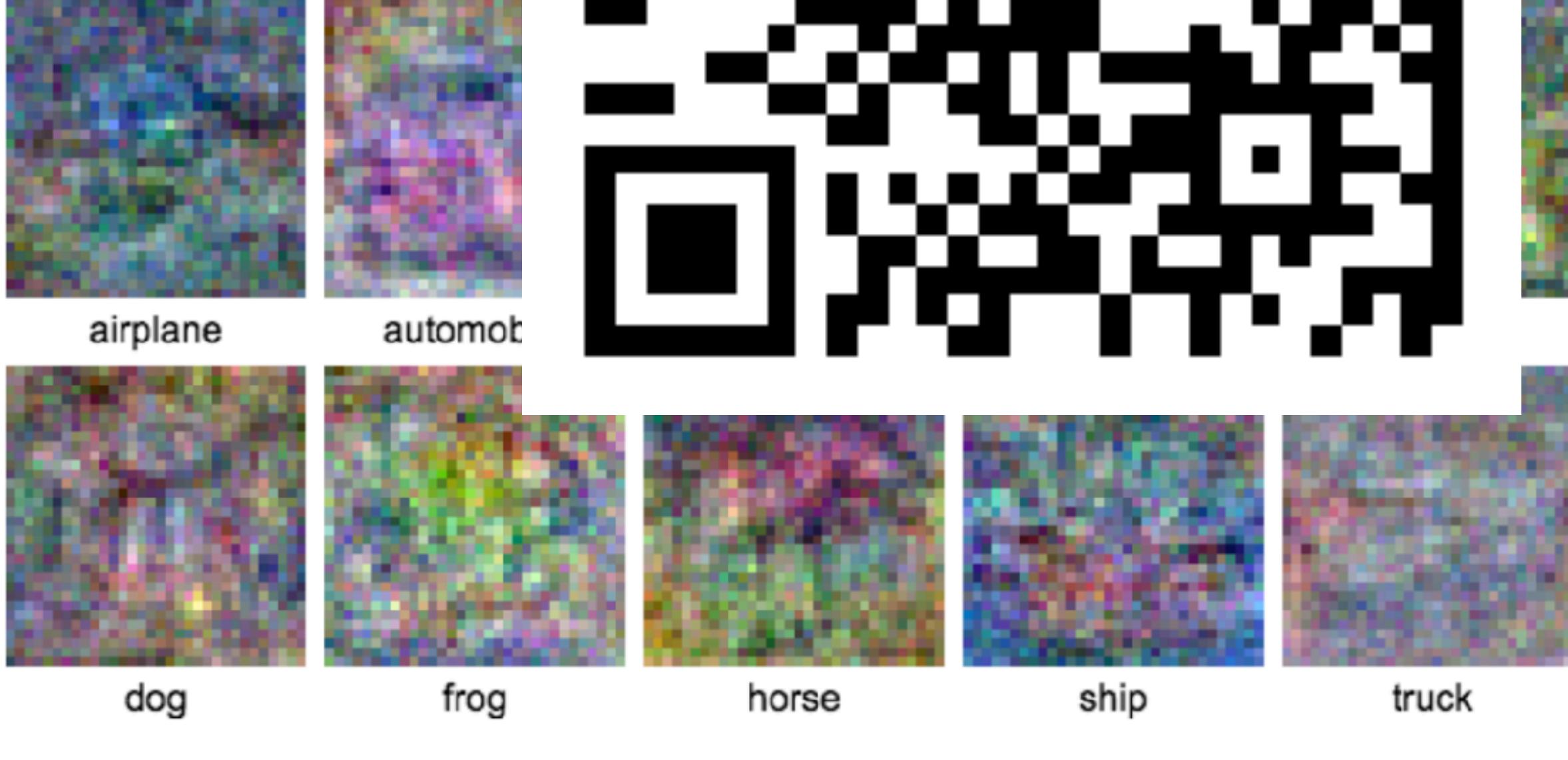
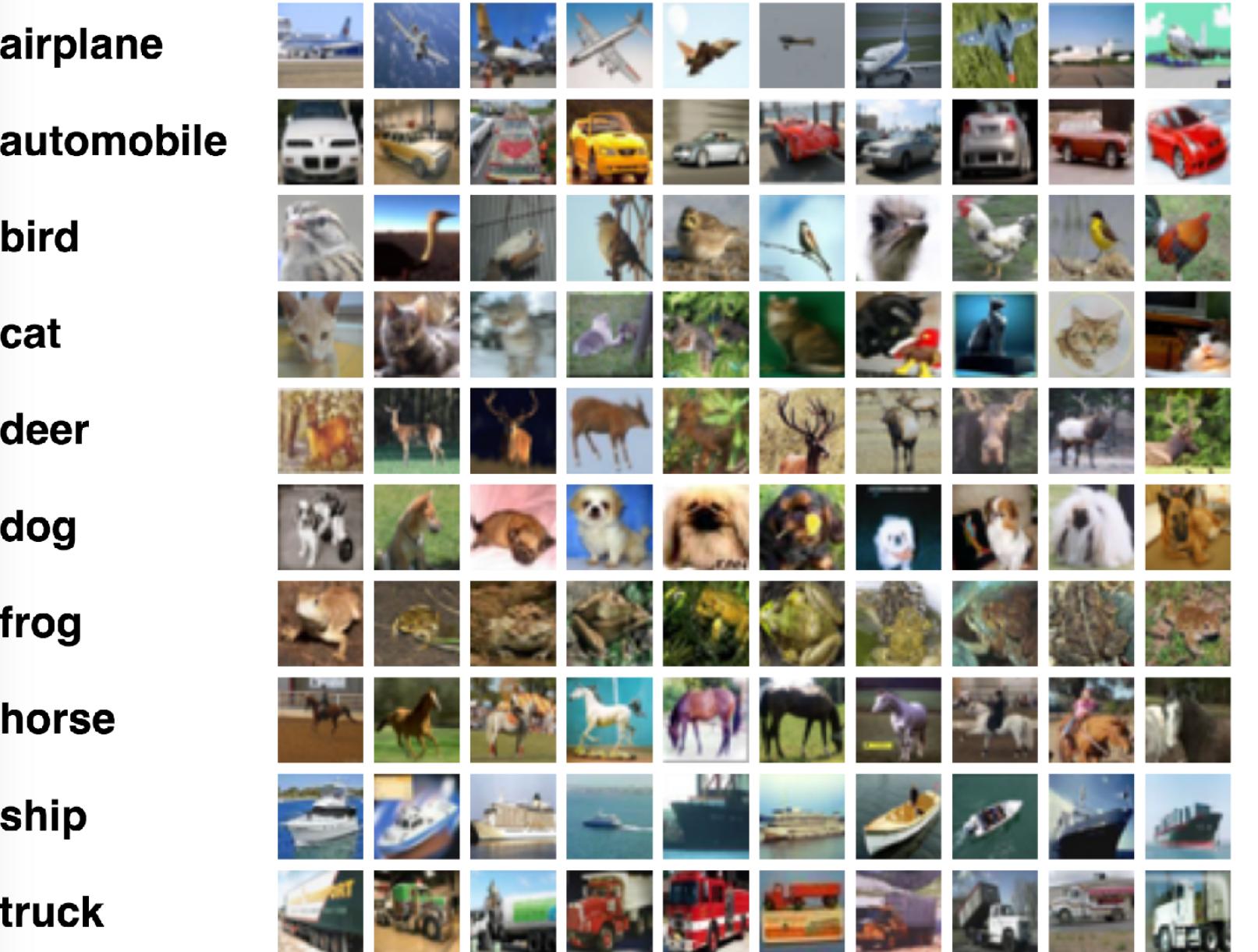
MNIST



???



CIFAR-10

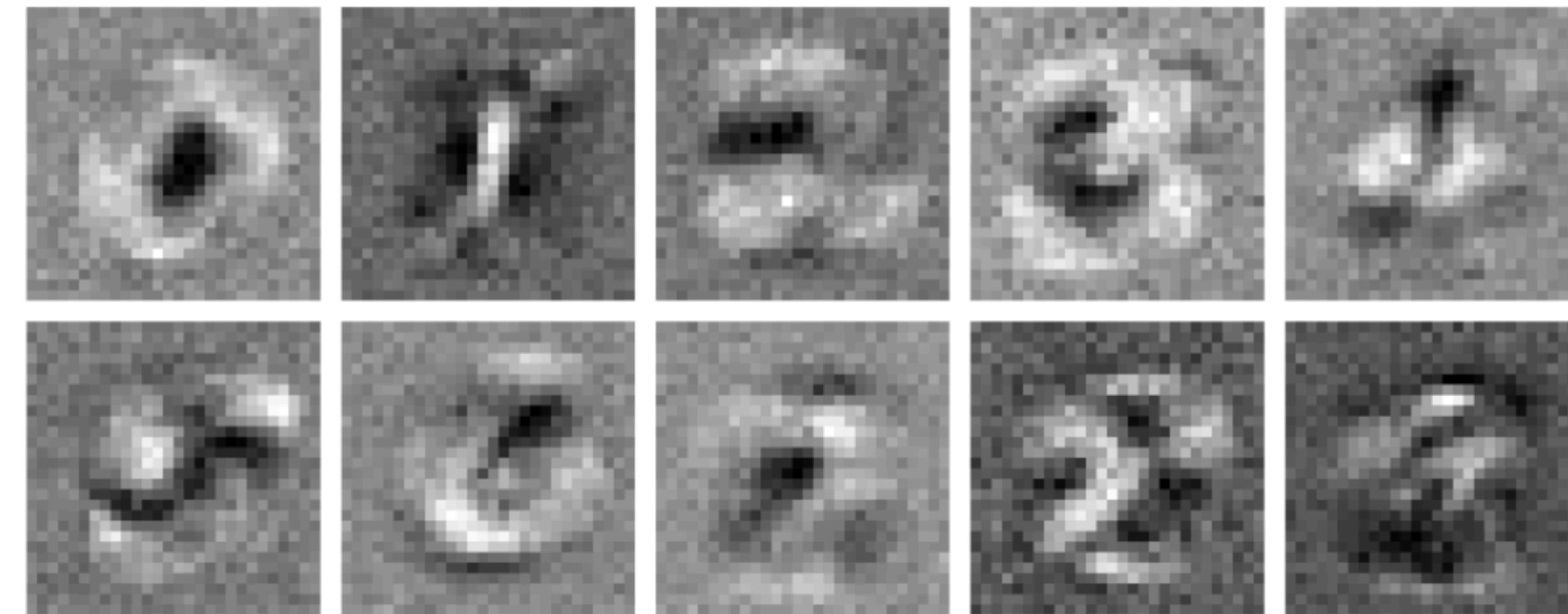


???

Dataset Label Images Learned weights of linear classifier Error

MNIST

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

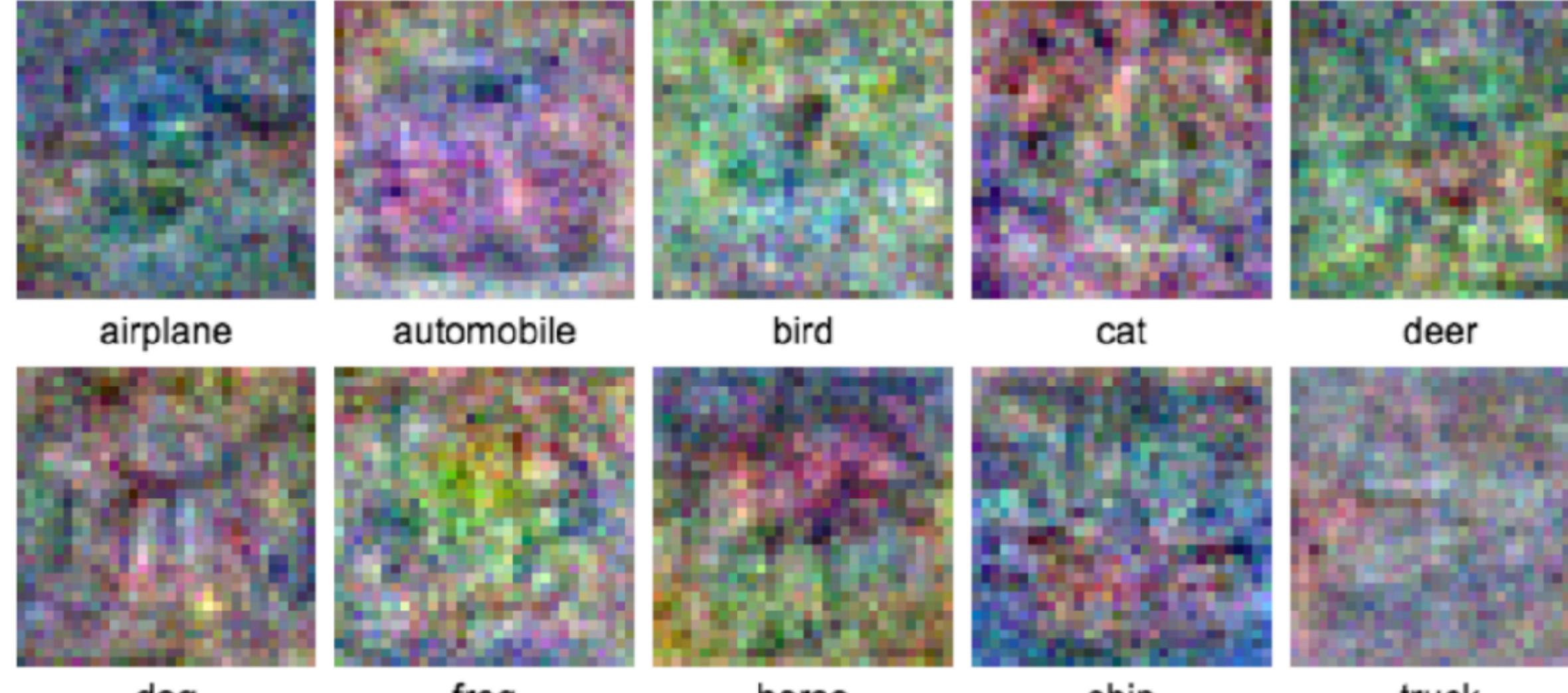
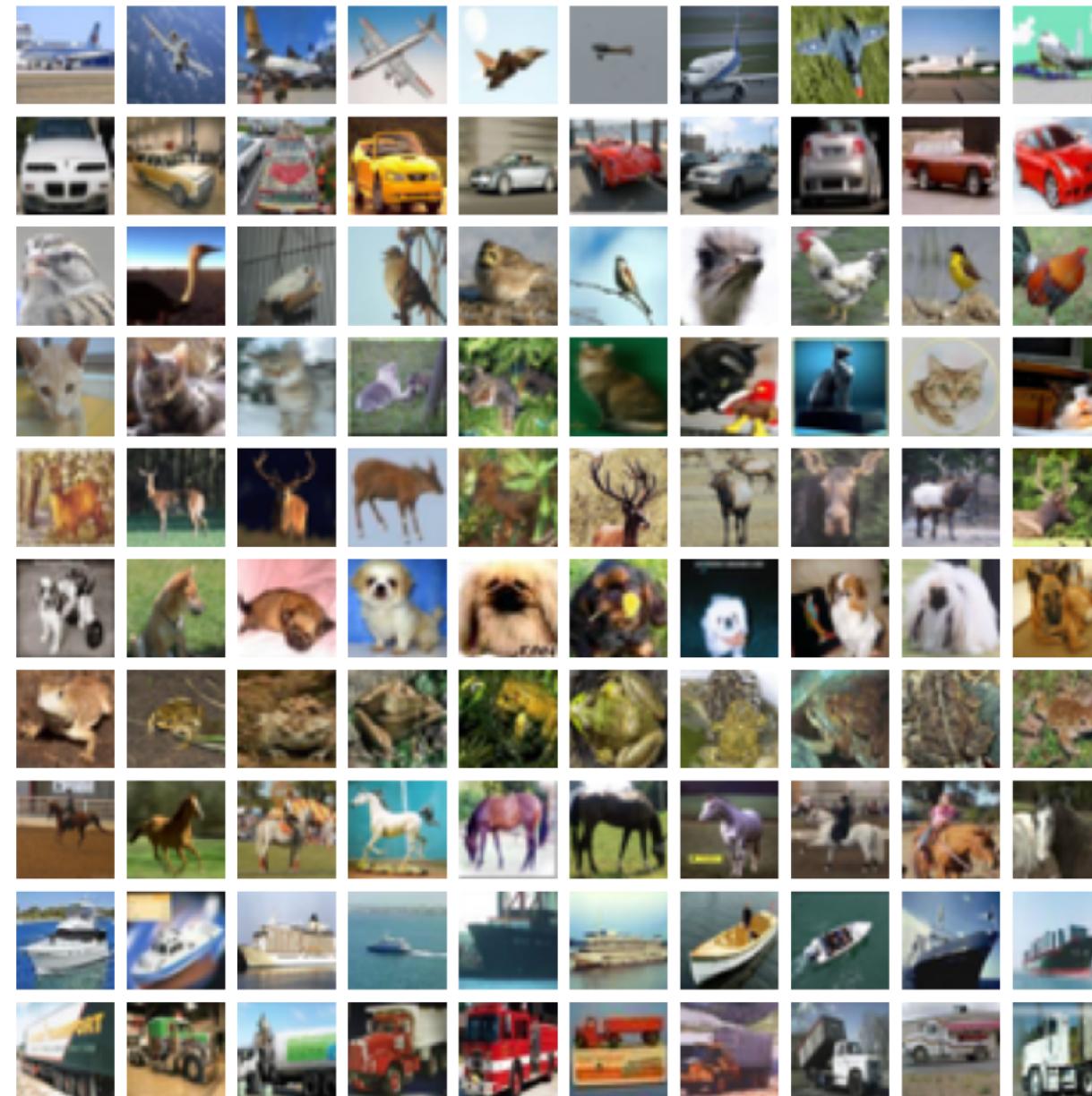


8%

Why is it simple?

CIFAR-10

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck



63%

Why is it hard?

Why is it hard?

Huge within-class variability!

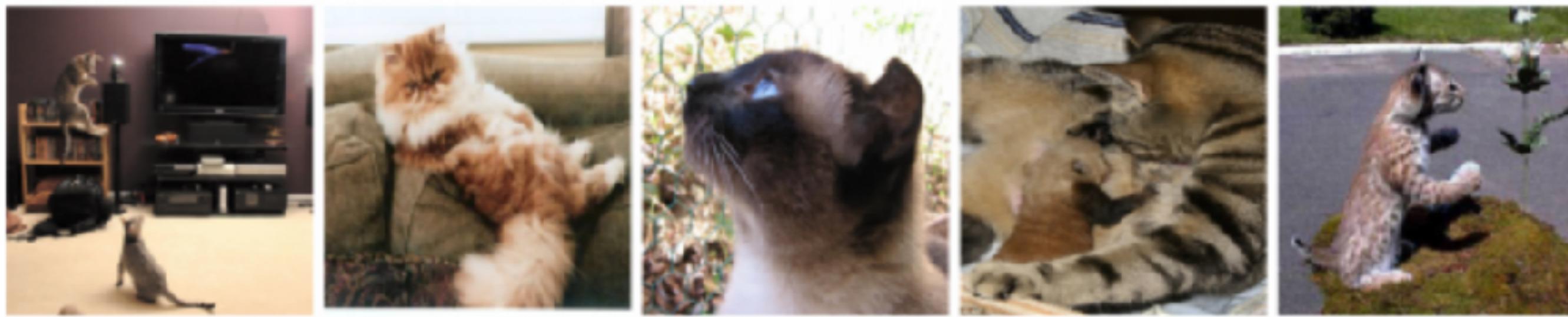
bird



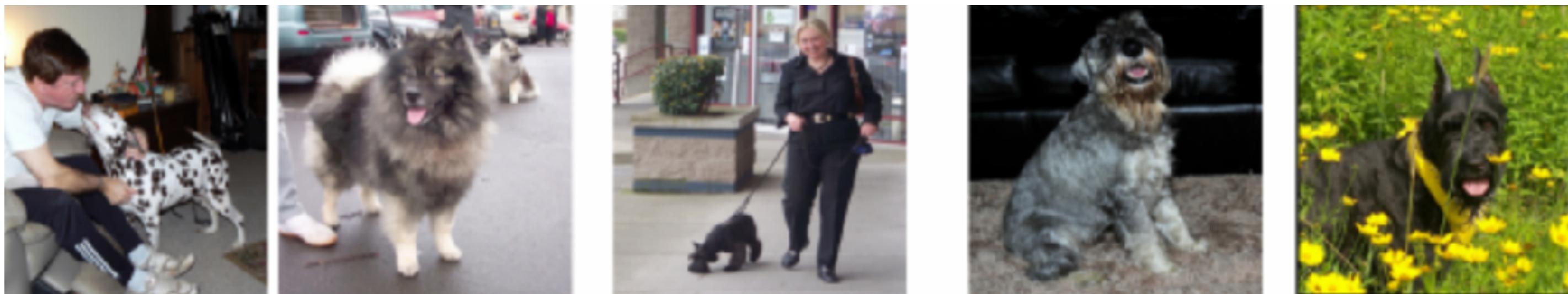
Due to:
???

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

cat



dog



Why is it hard?

Huge within-class variability!



Due to:

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

Why is it hard?

Huge among-class similarity!



Why is it hard?

Huge among-class similarity!

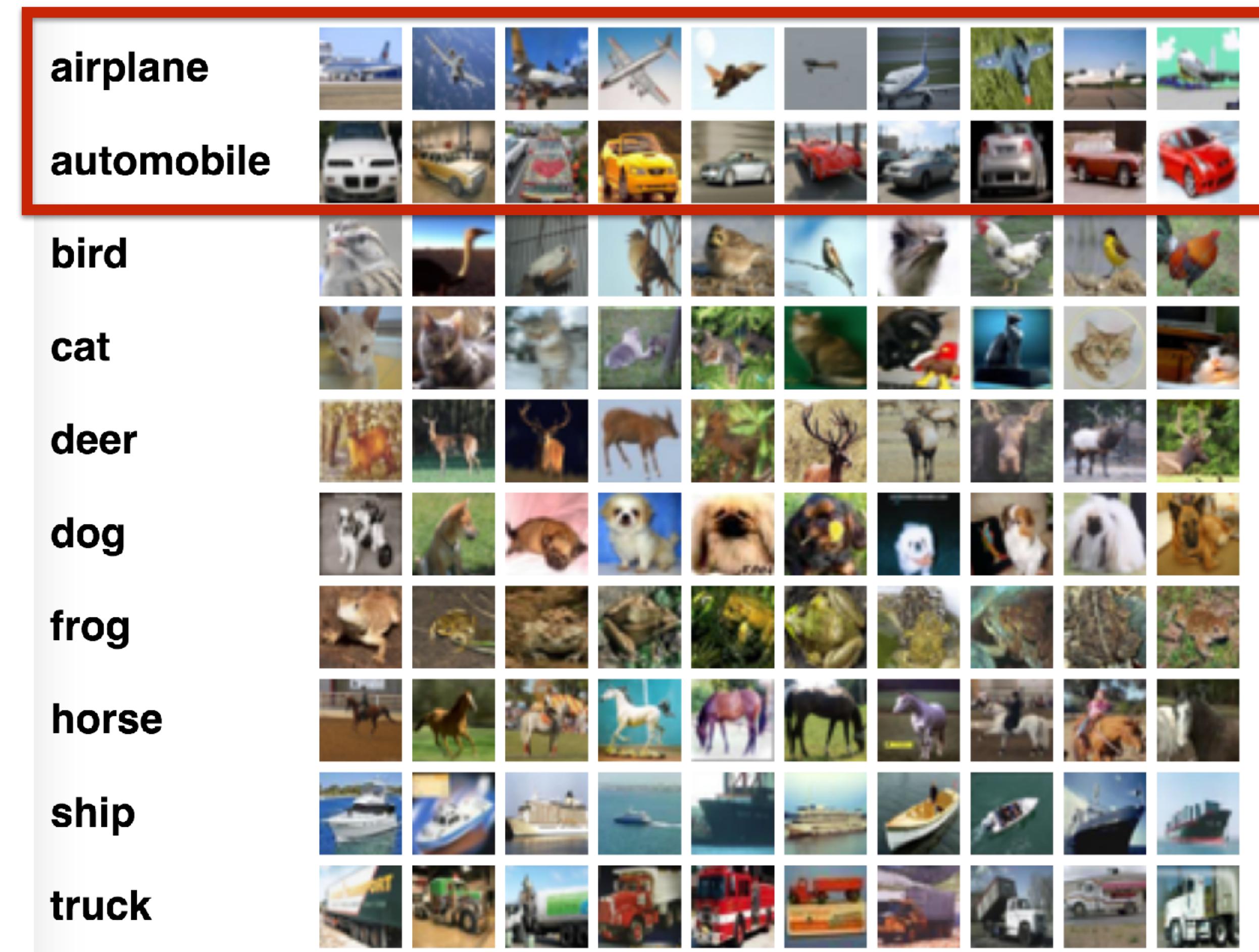


Why is it hard?

Huge among-class similarity!



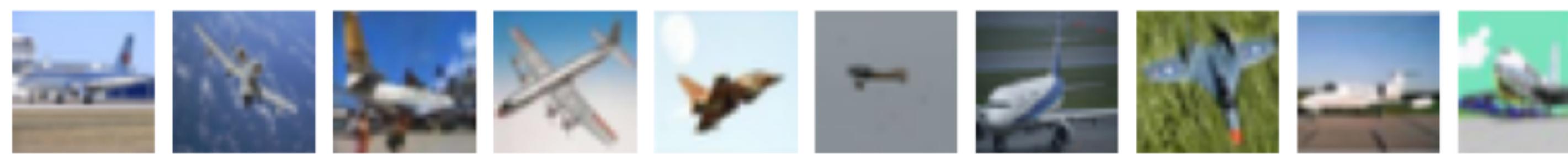
Recognition problem



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

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

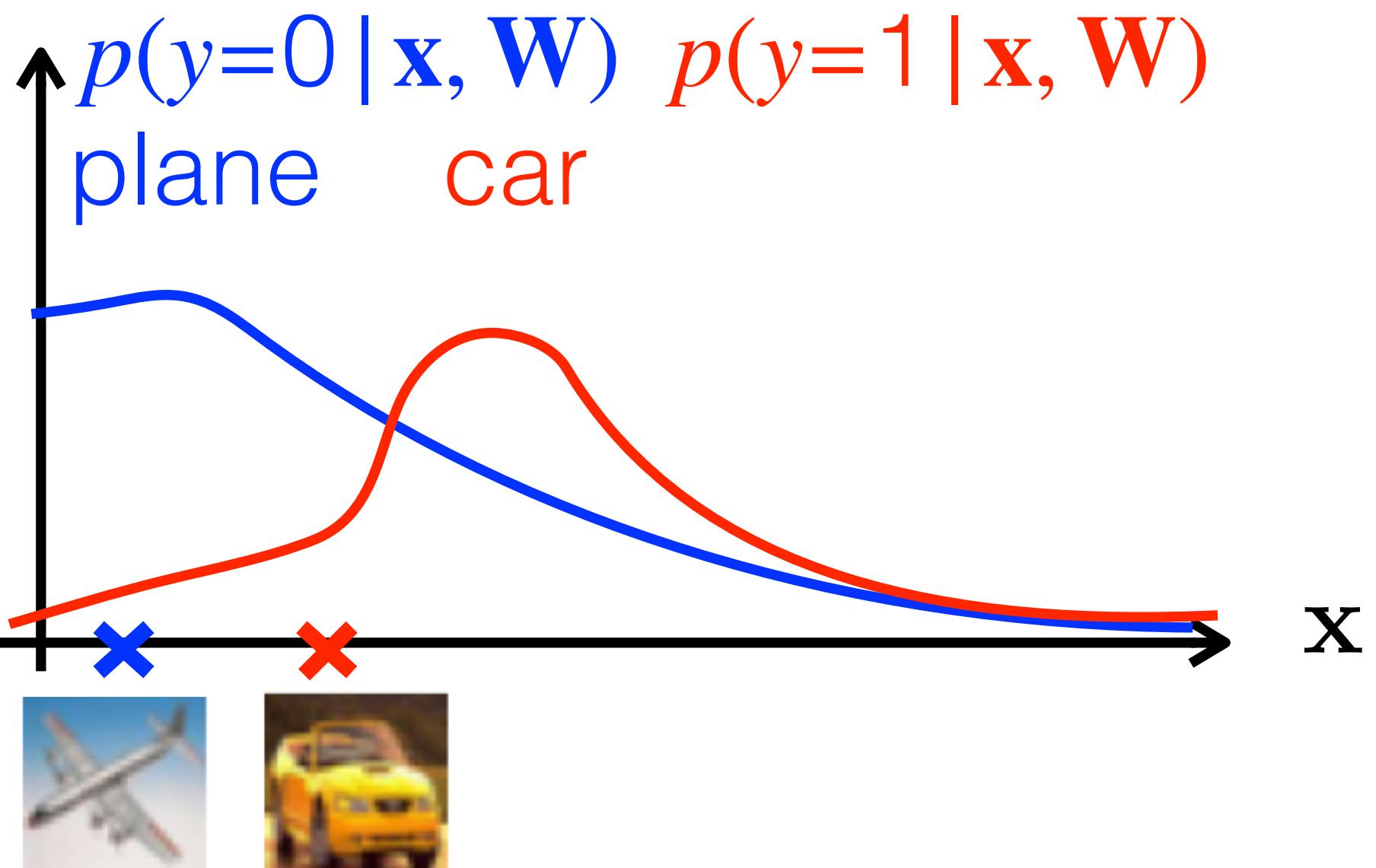
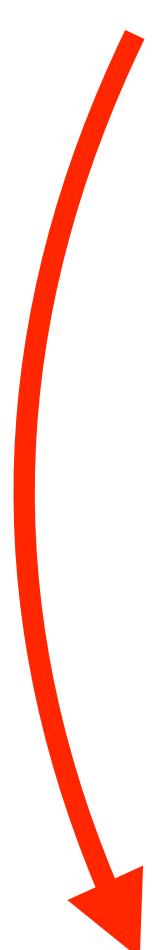
$$y = 0$$

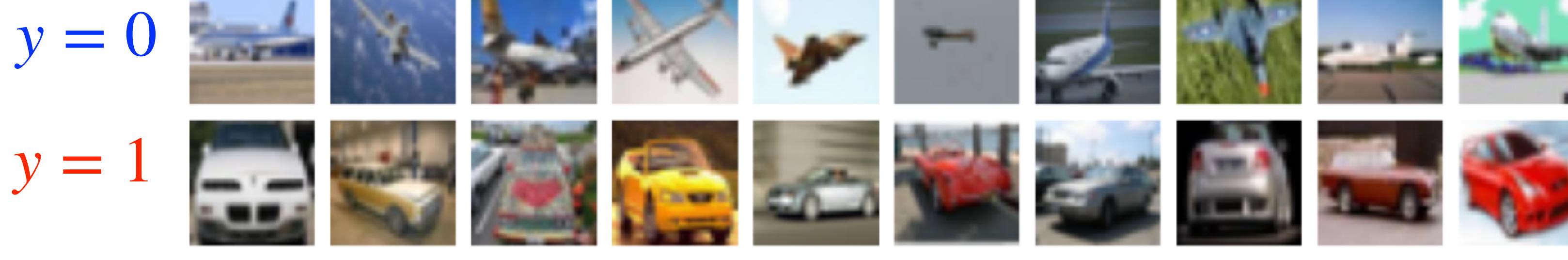


$$y = 1$$



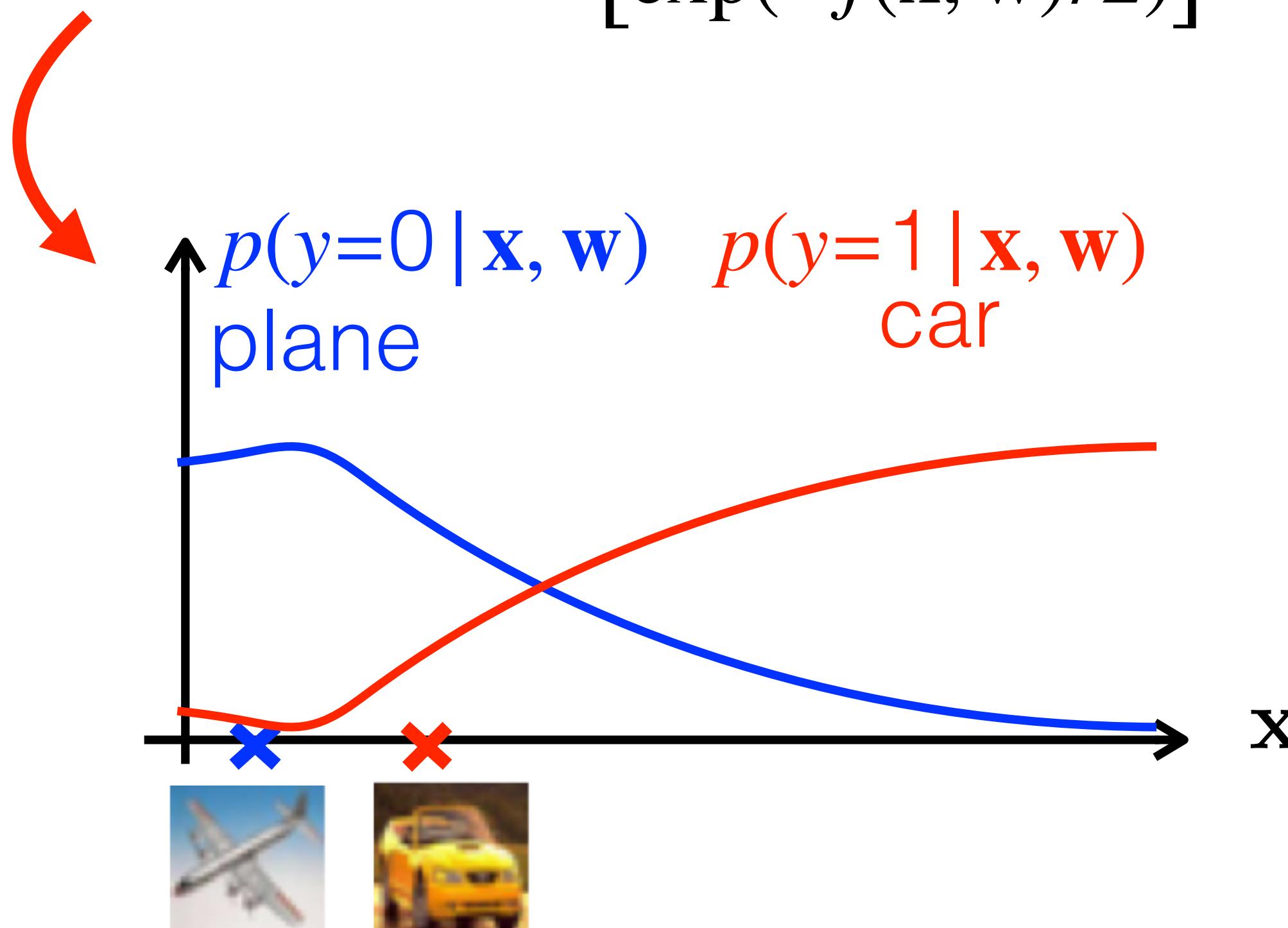
$$p(y | \mathbf{x}, \mathbf{W}) = \frac{\exp(f(\mathbf{x}, \mathbf{w}_1))}{\exp(f(\mathbf{x}, \mathbf{w}_2))} / \sum_k \exp(f(\mathbf{x}, \mathbf{w}_k)) = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) \in \mathbb{R}^2$$





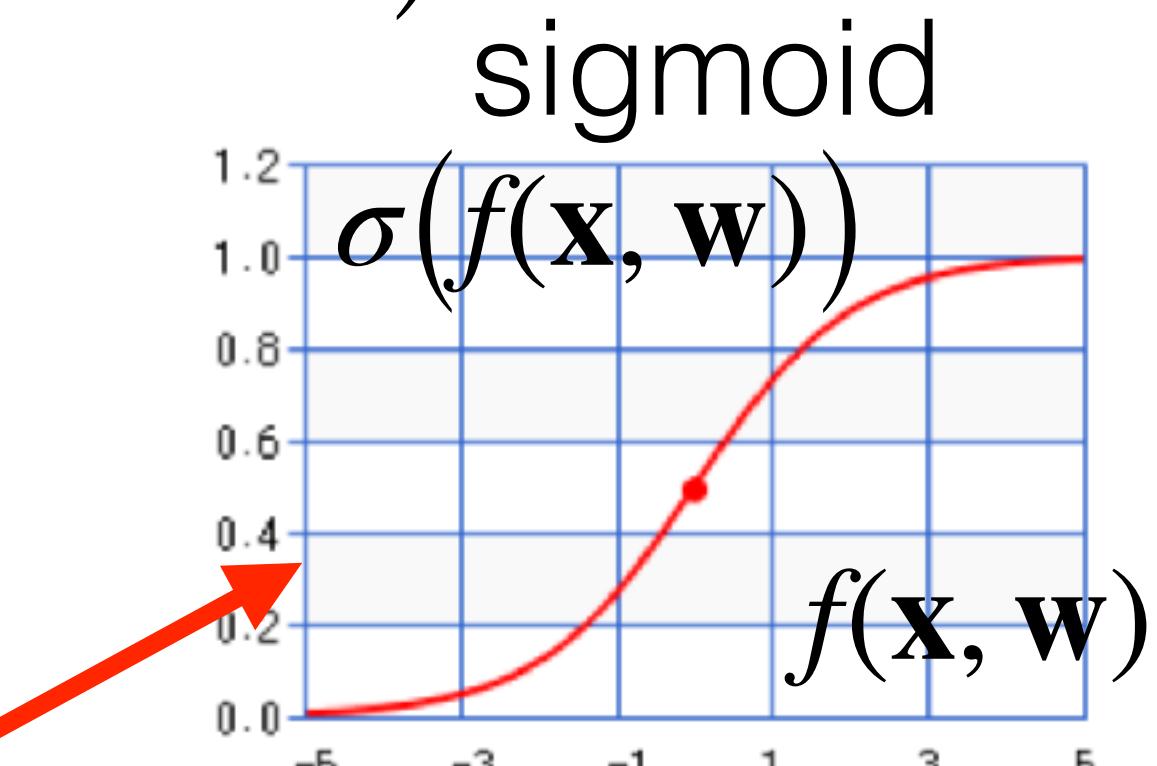
$$p(y | \mathbf{x}, \mathbf{W}) = \begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \end{bmatrix} / \sum_k \exp(f(\mathbf{x}, \mathbf{w}_k)) = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) \in \mathbb{R}^2$$

$$p(y | \mathbf{x}, \mathbf{w}) = \begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w})/2) \\ \exp(-f(\mathbf{x}, \mathbf{w})/2) \end{bmatrix} / \left(\exp(f(\mathbf{x}, \mathbf{w})/2) + \exp(-f(\mathbf{x}, \mathbf{w})/2) \right)$$



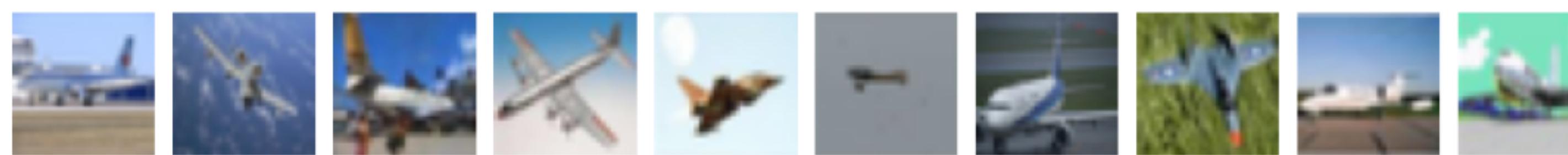
$$= \begin{bmatrix} \frac{1}{1 + \frac{\exp(-f(\mathbf{x}, \mathbf{w})/2)}{\exp(f(\mathbf{x}, \mathbf{w})/2)}} \\ 1 - \frac{1}{1 + \frac{\exp(-f(\mathbf{x}, \mathbf{w})/2)}{\exp(f(\mathbf{x}, \mathbf{w})/2)}} \end{bmatrix}$$

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



$$= \begin{bmatrix} \sigma(f(\mathbf{x}, \mathbf{w})) \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) \end{bmatrix}$$

$$y = 0$$

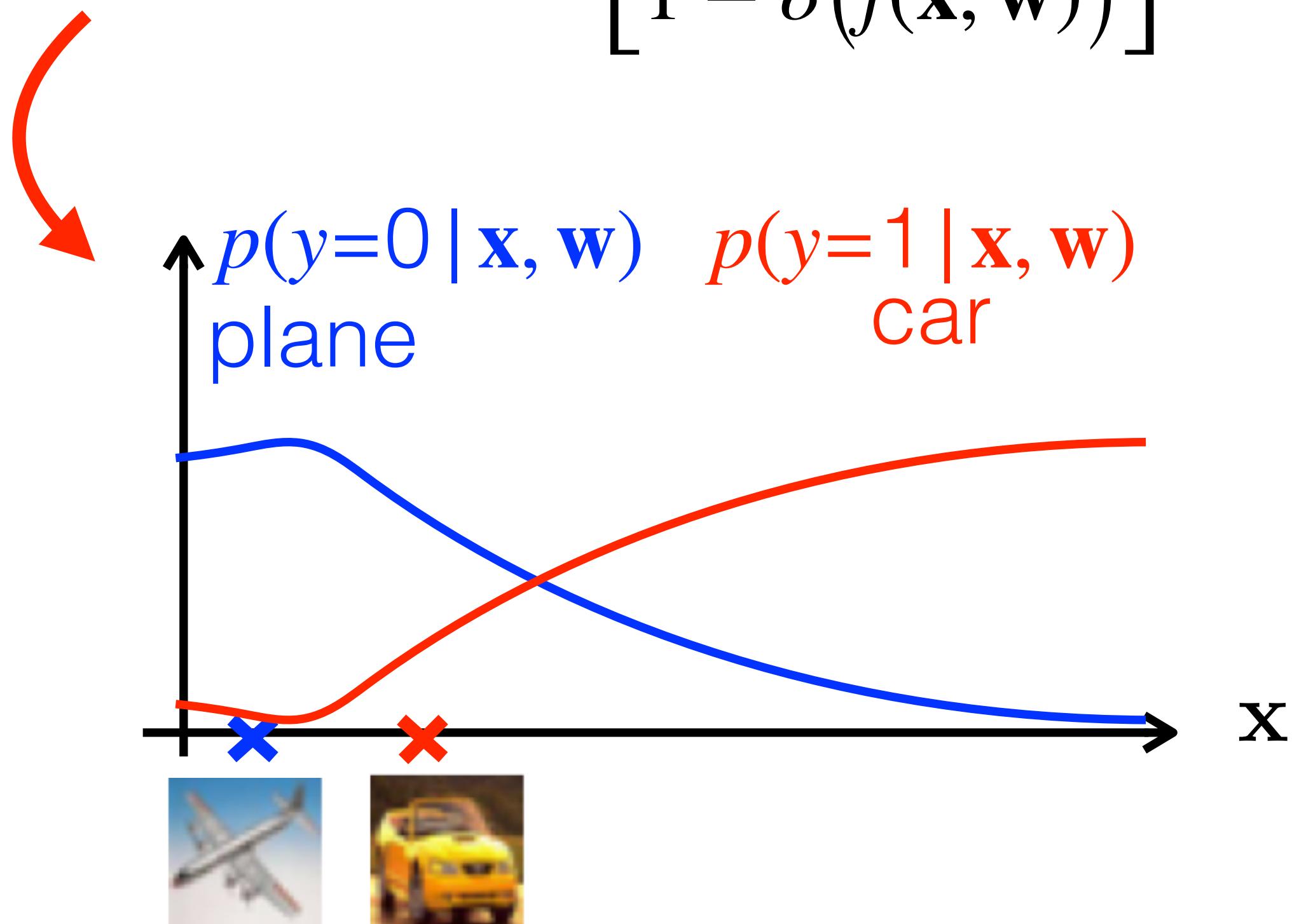


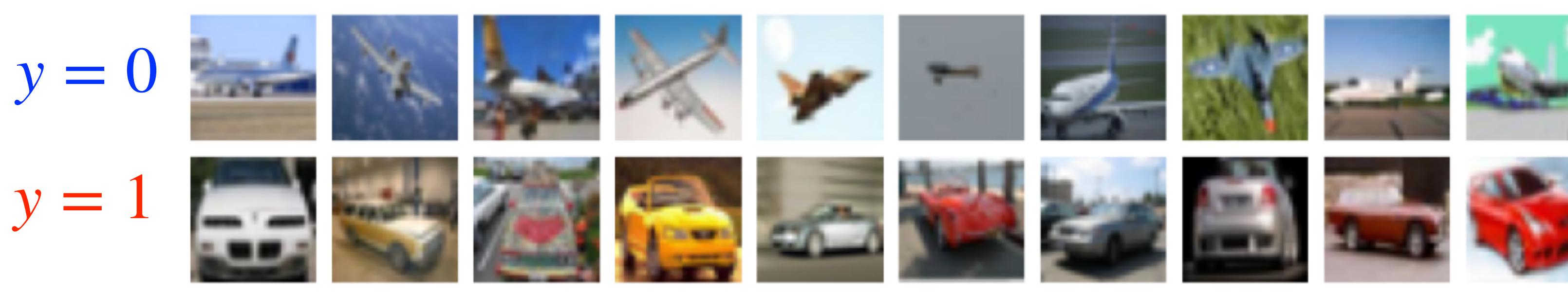
$$y = 1$$



$$\mathbf{p}(y \mid \mathbf{x}, \mathbf{W}) = \begin{bmatrix} \exp(f(\mathbf{x}, \mathbf{w}_1)) \\ \exp(f(\mathbf{x}, \mathbf{w}_2)) \end{bmatrix} / \sum_k \exp(f(\mathbf{x}, \mathbf{w}_k)) = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) \in \mathbb{R}^2$$

$$\mathbf{p}(y \mid \mathbf{x}, \mathbf{w}) = \begin{bmatrix} \sigma(f(\mathbf{x}, \mathbf{w})) \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) \end{bmatrix} \in \mathbb{R}^2$$

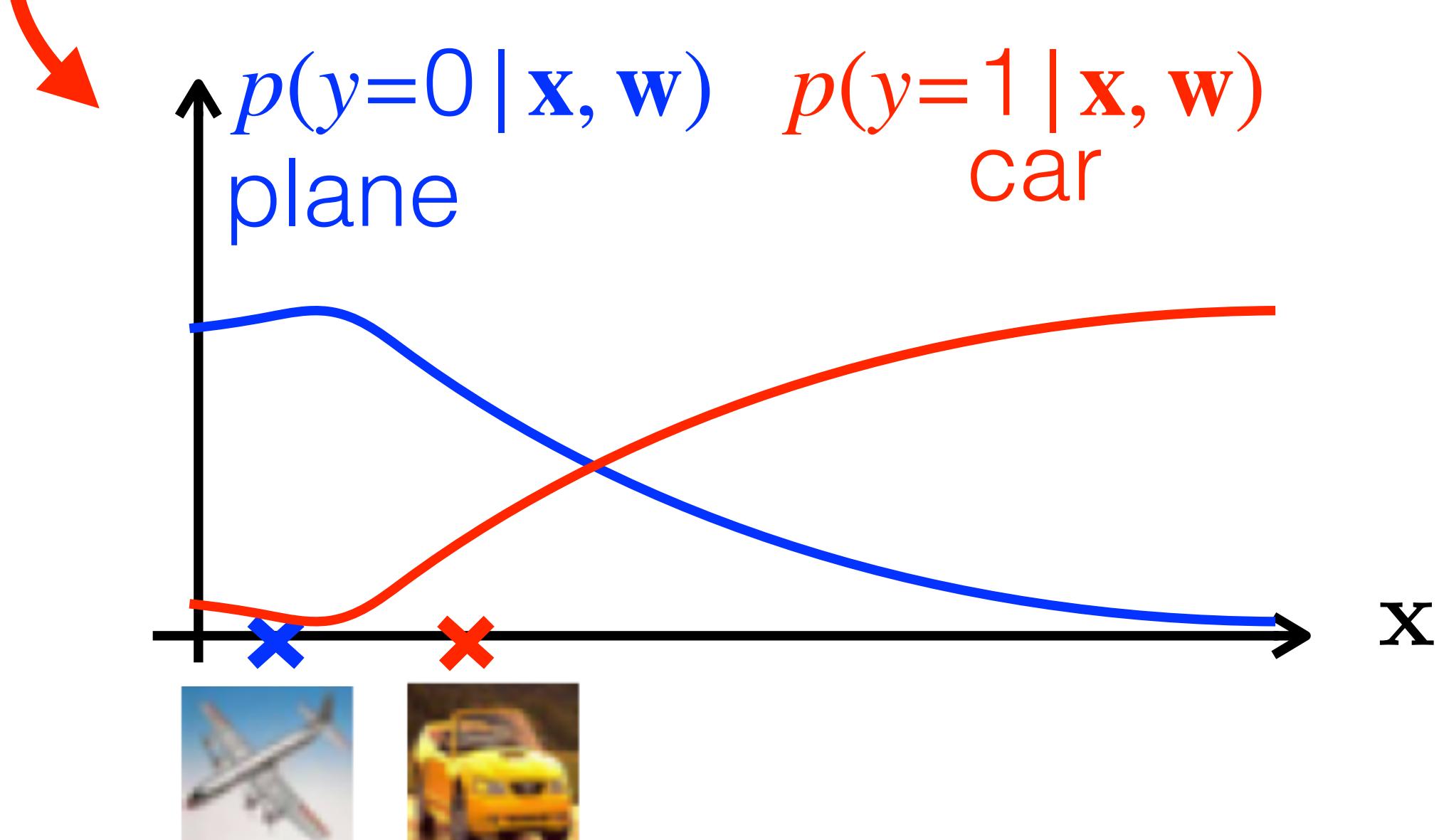




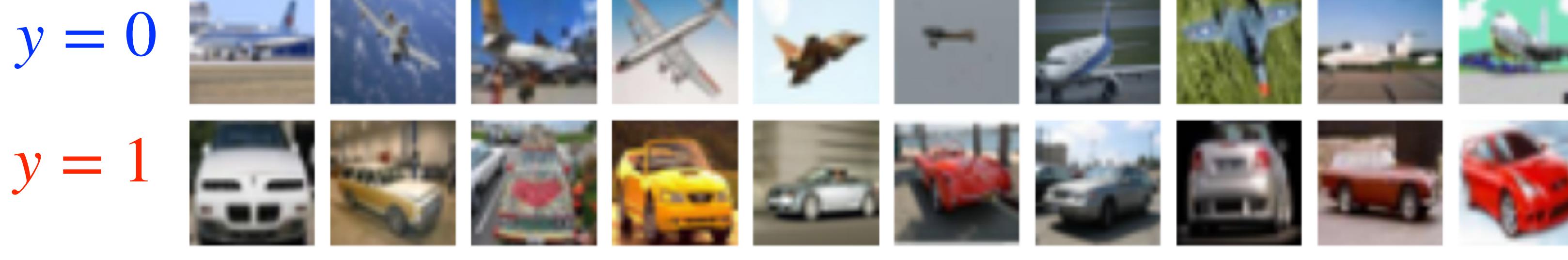
$$p(y | \mathbf{x}, \mathbf{W}) = \frac{\exp(f(\mathbf{x}, \mathbf{w}_1))}{\sum_k \exp(f(\mathbf{x}, \mathbf{w}_k))} = \mathbf{s}(\mathbf{f}(\mathbf{x}, \mathbf{W})) \in \mathbb{R}^2$$

Red arrow pointing from the equation above to this one:

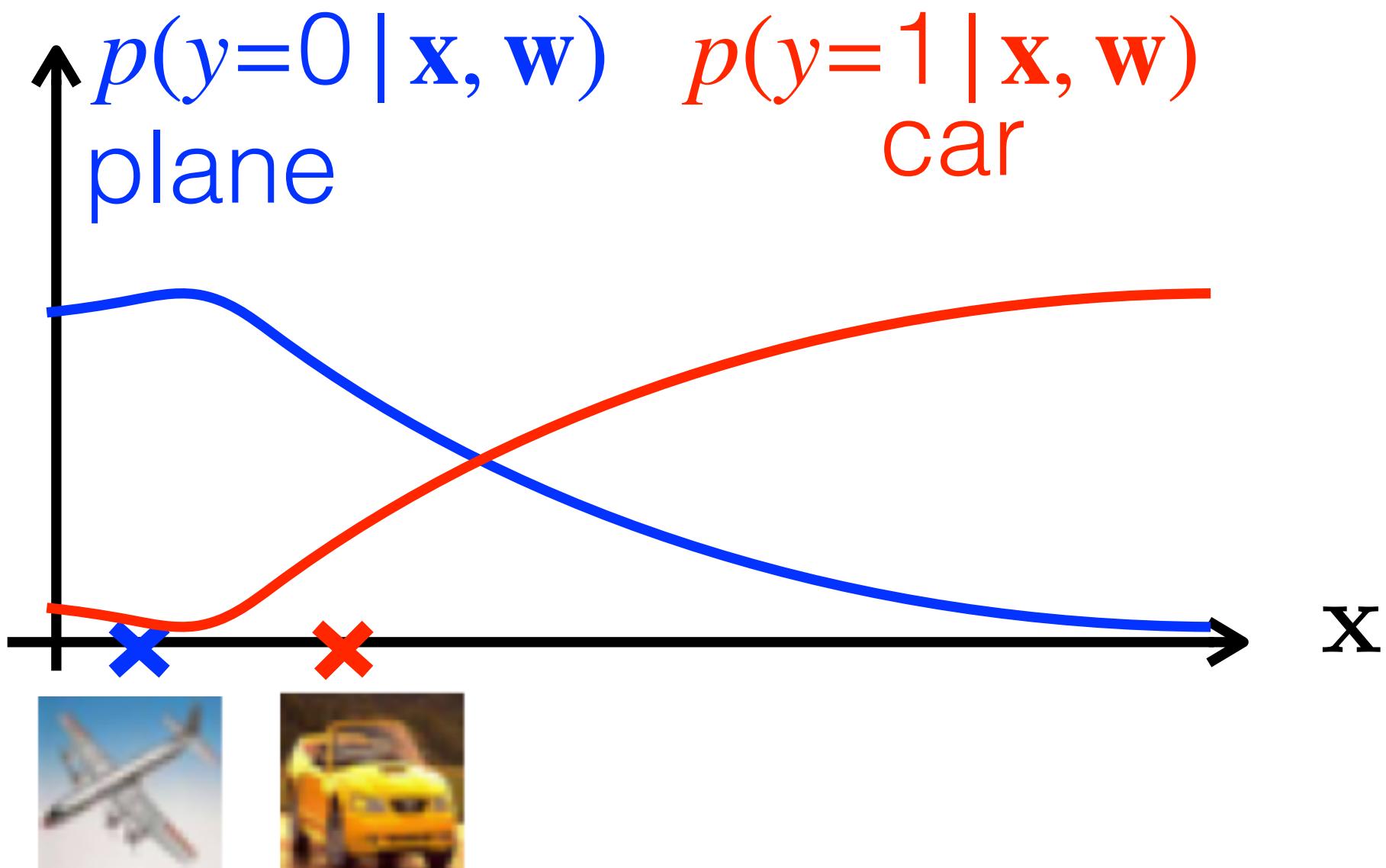
$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = 1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = 0 \end{cases} \in \mathbb{R}^1$$

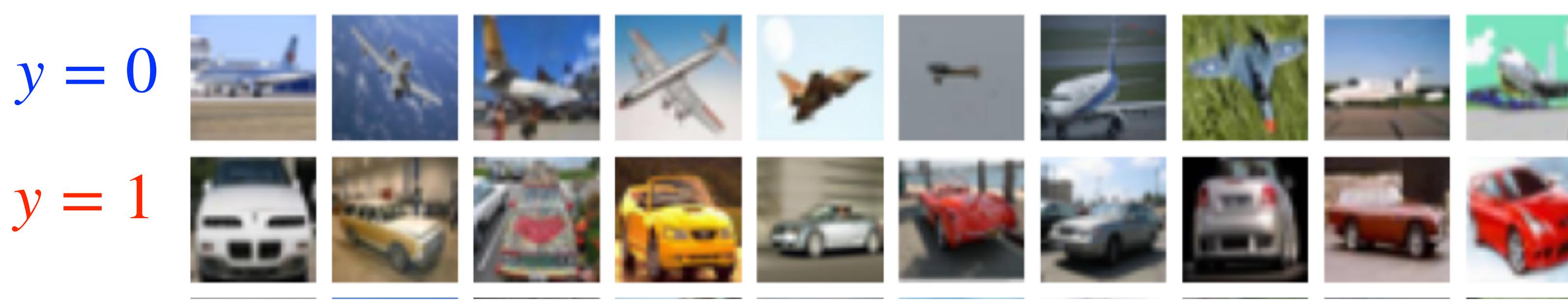


$$p(y | \mathbf{x}, \mathbf{w}) = \begin{bmatrix} \sigma(f(\mathbf{x}, \mathbf{w})) \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) \end{bmatrix} \in \mathbb{R}^2$$

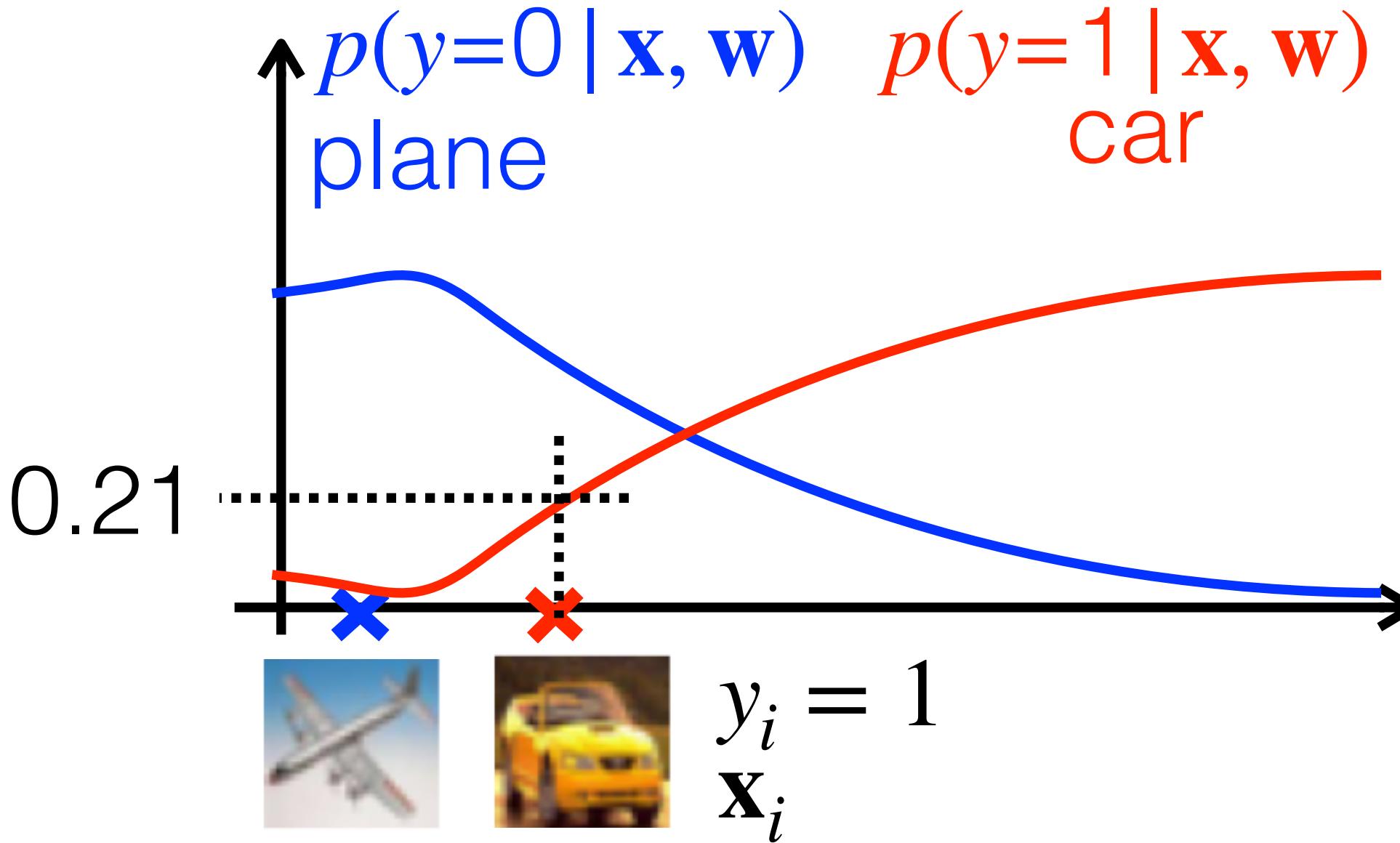


$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = 1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = 0 \end{cases} = y \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w})))$$



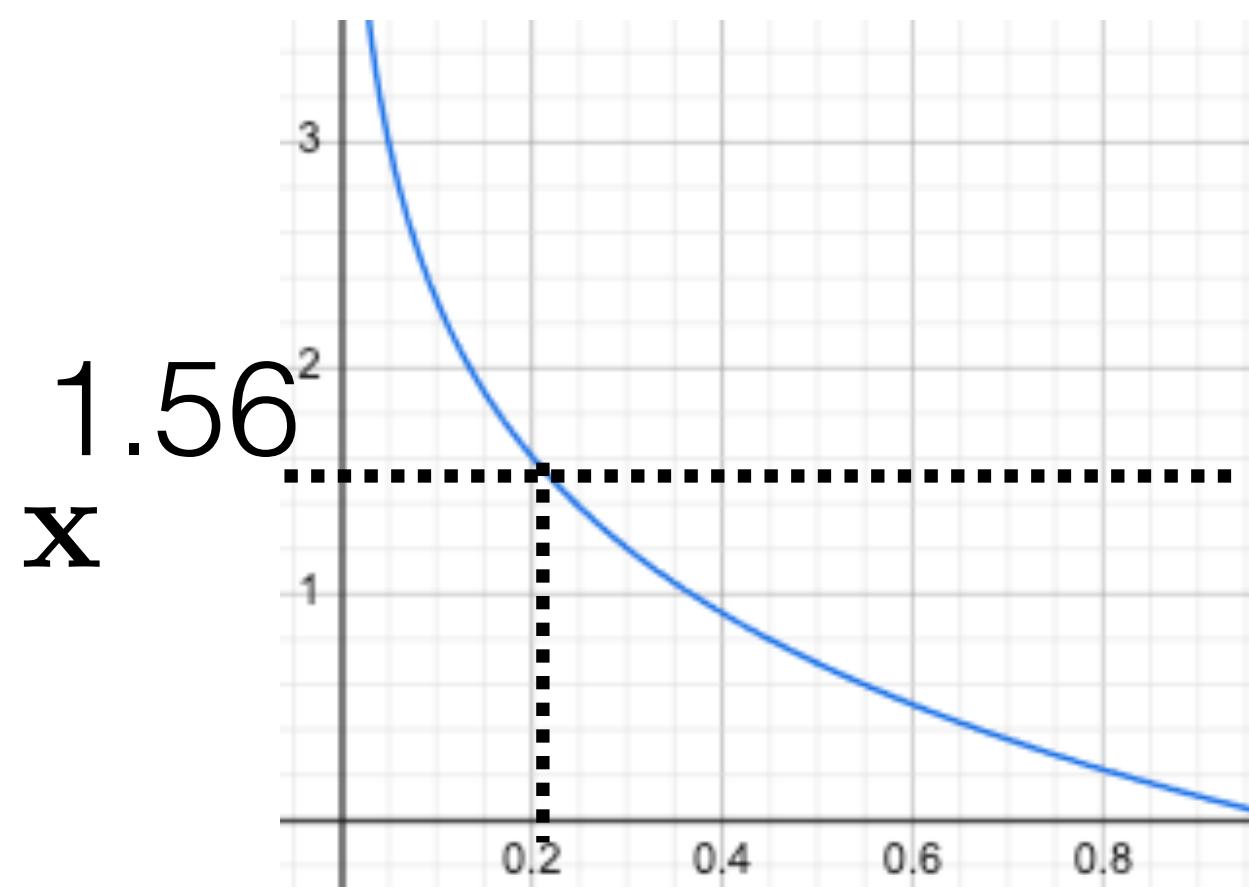


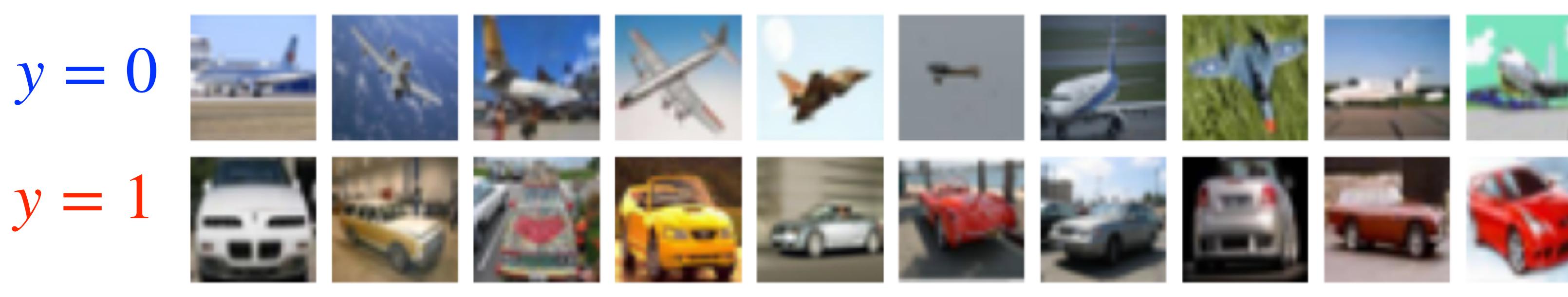
$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = 1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = 0 \end{cases} = y_i \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y_i) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w})))$$



Loss function (subst. car \times):

$$\begin{aligned} \mathcal{L}(\mathbf{w}) &= -\log p(y_i | \mathbf{x}_i, \mathbf{w}) \\ &= -\log(y_i \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y_i) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w})))) \\ &= -\log(1 \cdot 0.21 + (1 - 1) \cdot 0.21) = 1.56 \end{aligned}$$





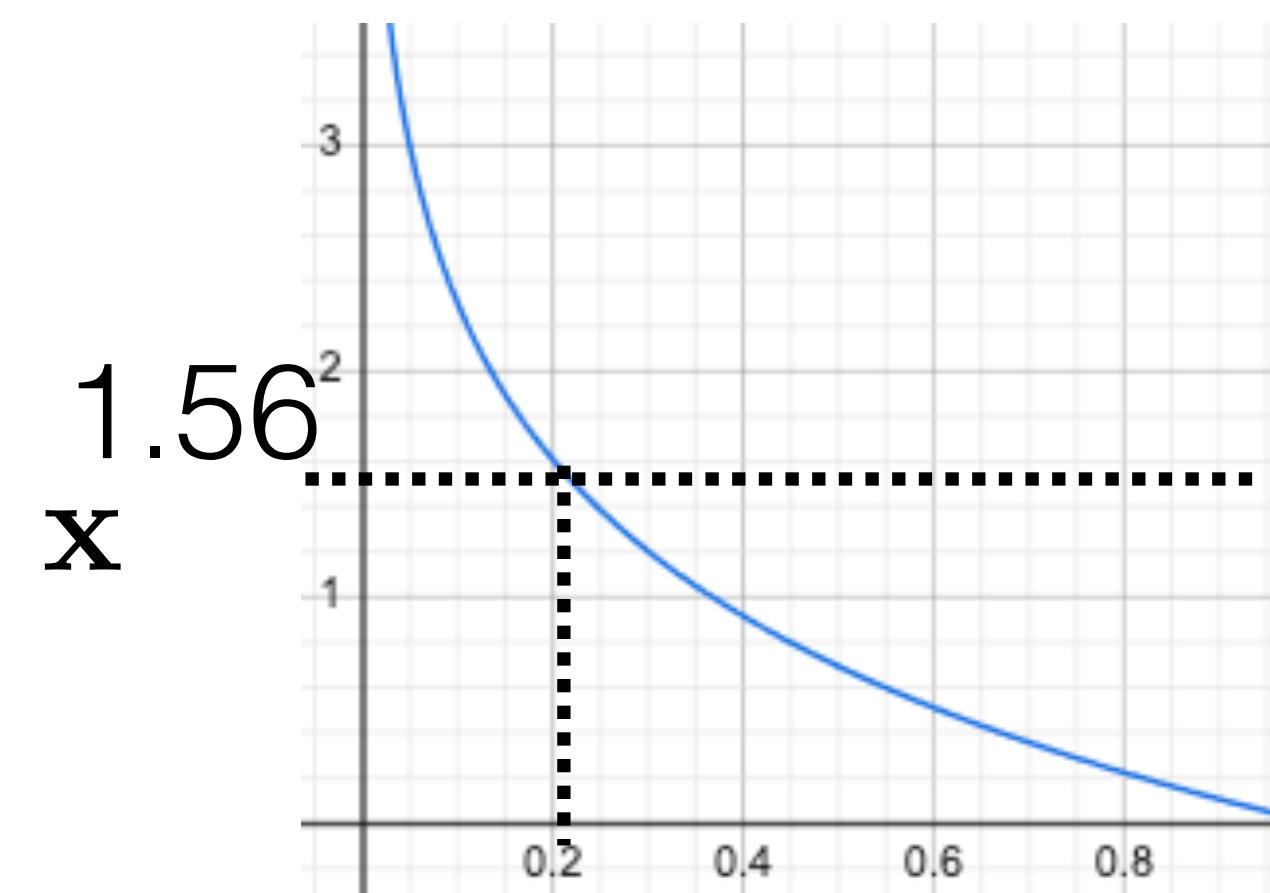
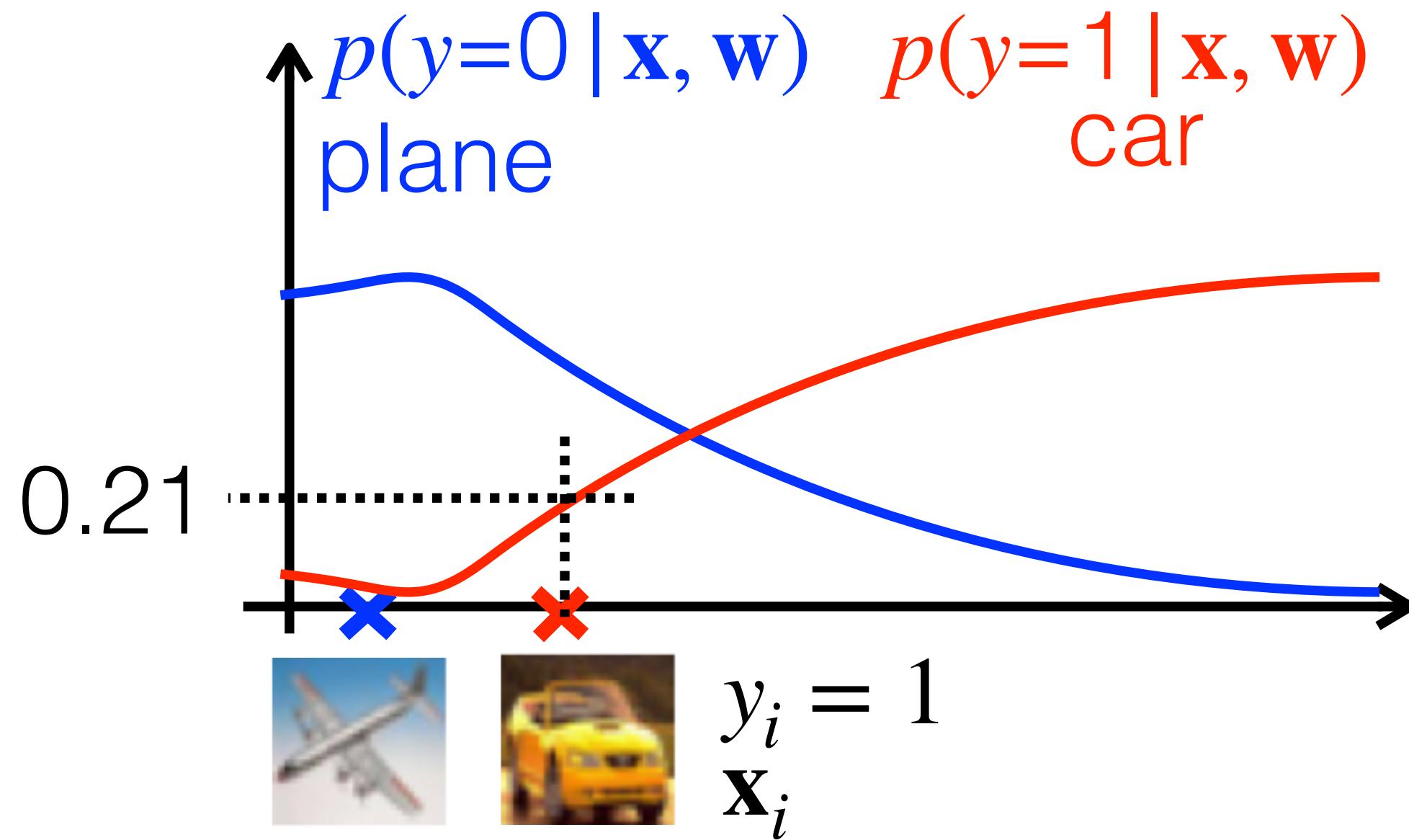
$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = 1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = 0 \end{cases} = y \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w})))$$

Loss function (subst. car \times):

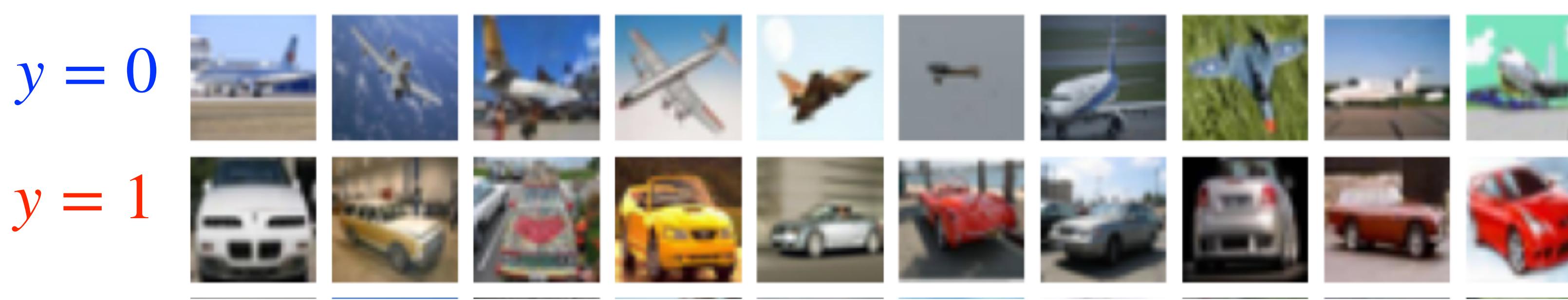
$$\mathcal{L}(\mathbf{w}) = -\log p(y_i | \mathbf{x}_i, \mathbf{w})$$

$$= -\log(y_i \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y_i) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w}))))$$

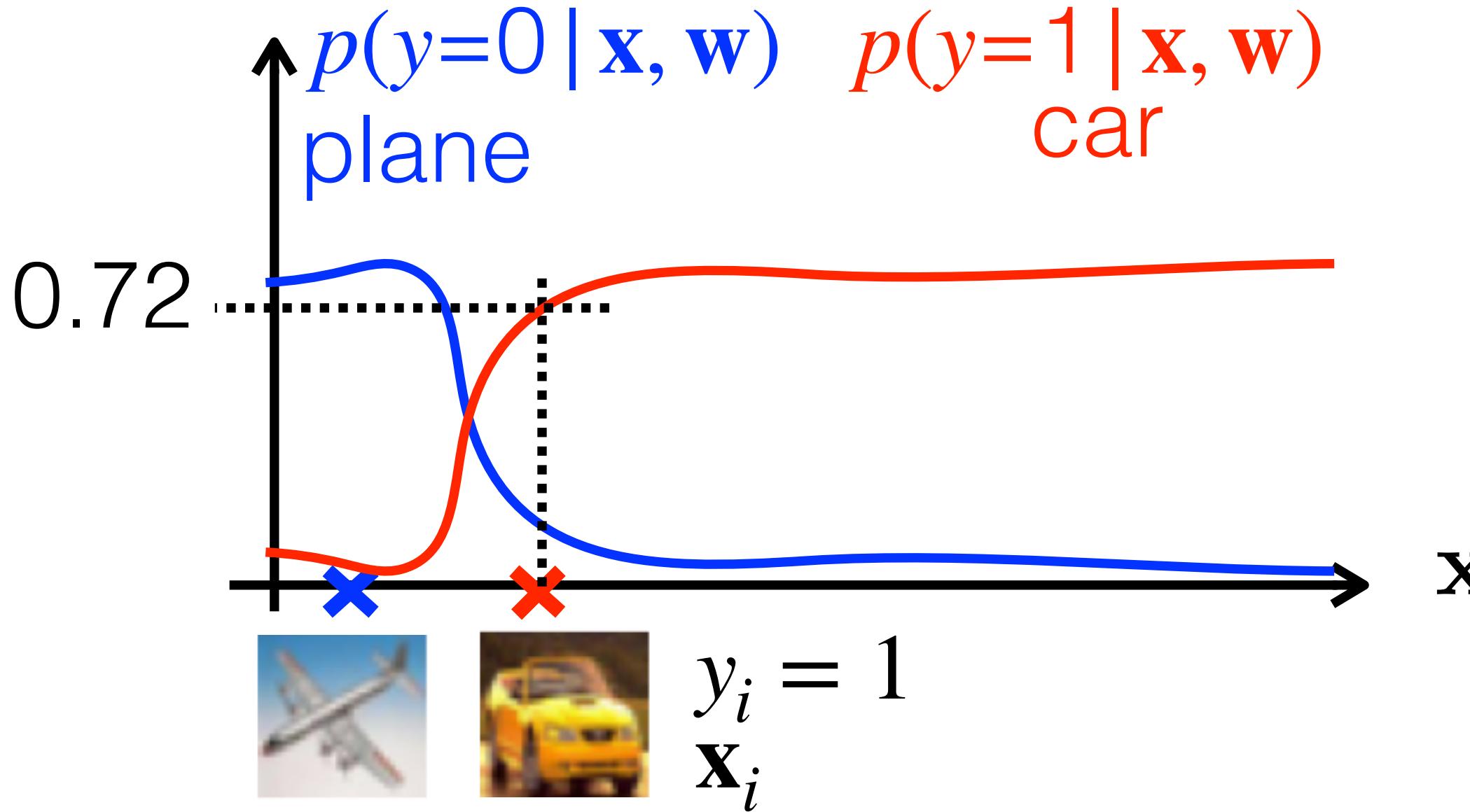
$$= -\log(1 \cdot 0.21 + (1 - 1) \cdot 0.21) = 1.56$$



Learning:
 $\arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w})$

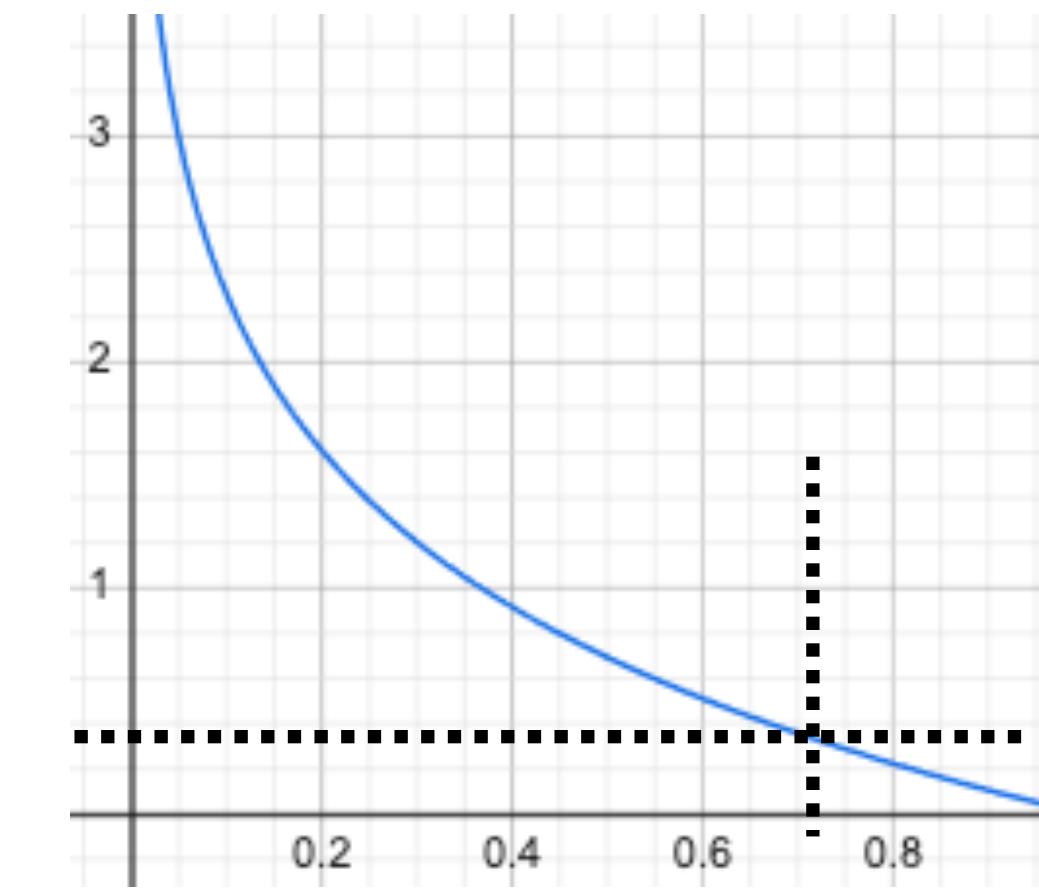


$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = 1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = 0 \end{cases} = y \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w})))$$

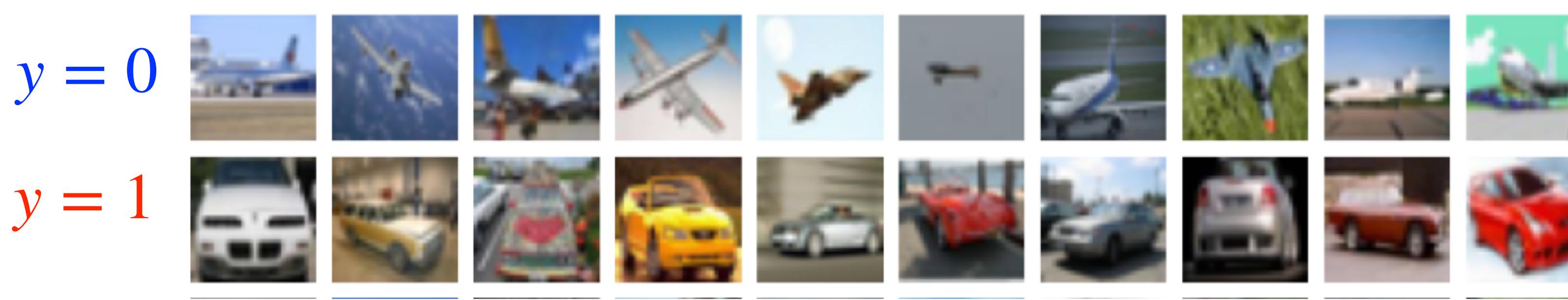


Loss function (subst. car \times):

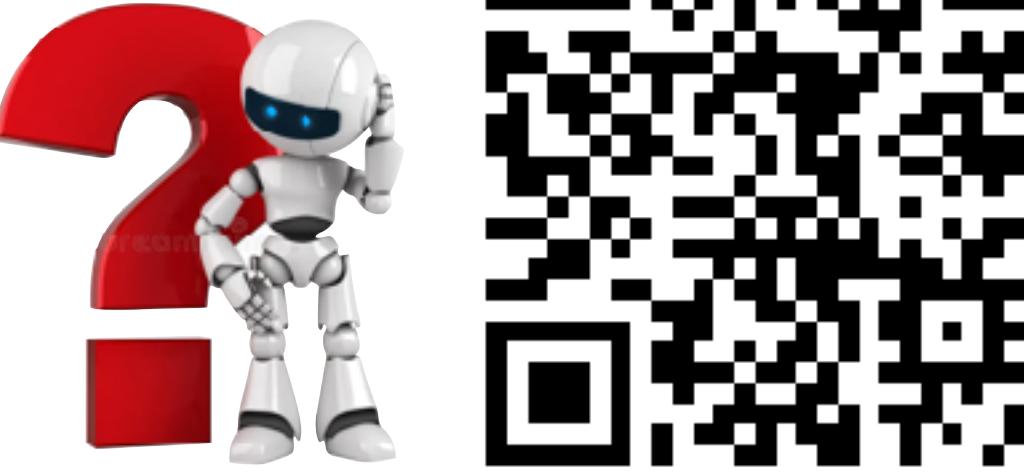
$$\begin{aligned} \mathcal{L}(\mathbf{w}) &= -\log p(y_i | \mathbf{x}_i, \mathbf{w}) \\ &= -\log(y_i \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y_i) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w})))) \\ &= -\log(1 \cdot 0.21 + (1 - 1) \cdot 0.21) = 1.56 \end{aligned}$$



Learning:
 $\arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w})$



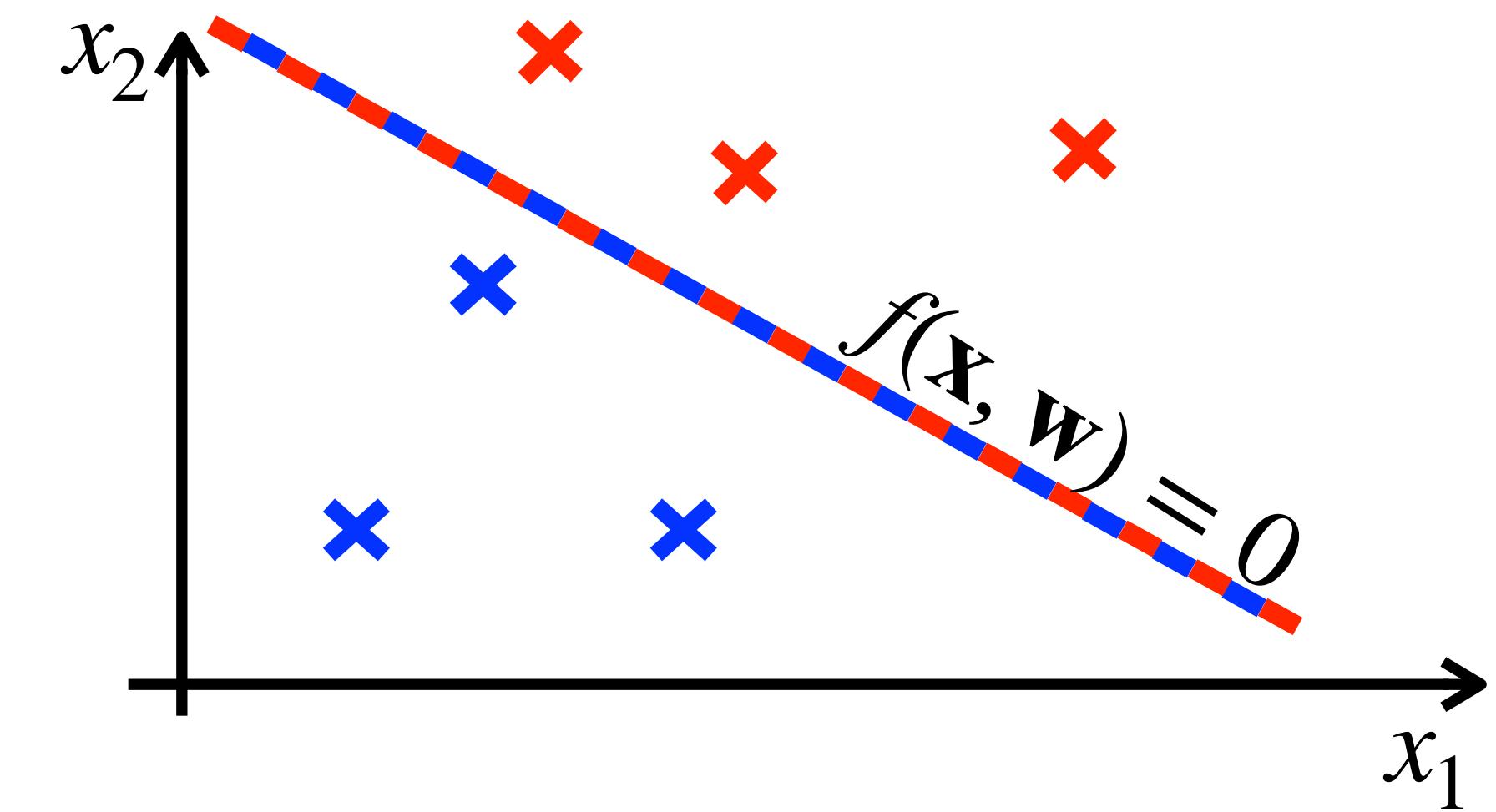
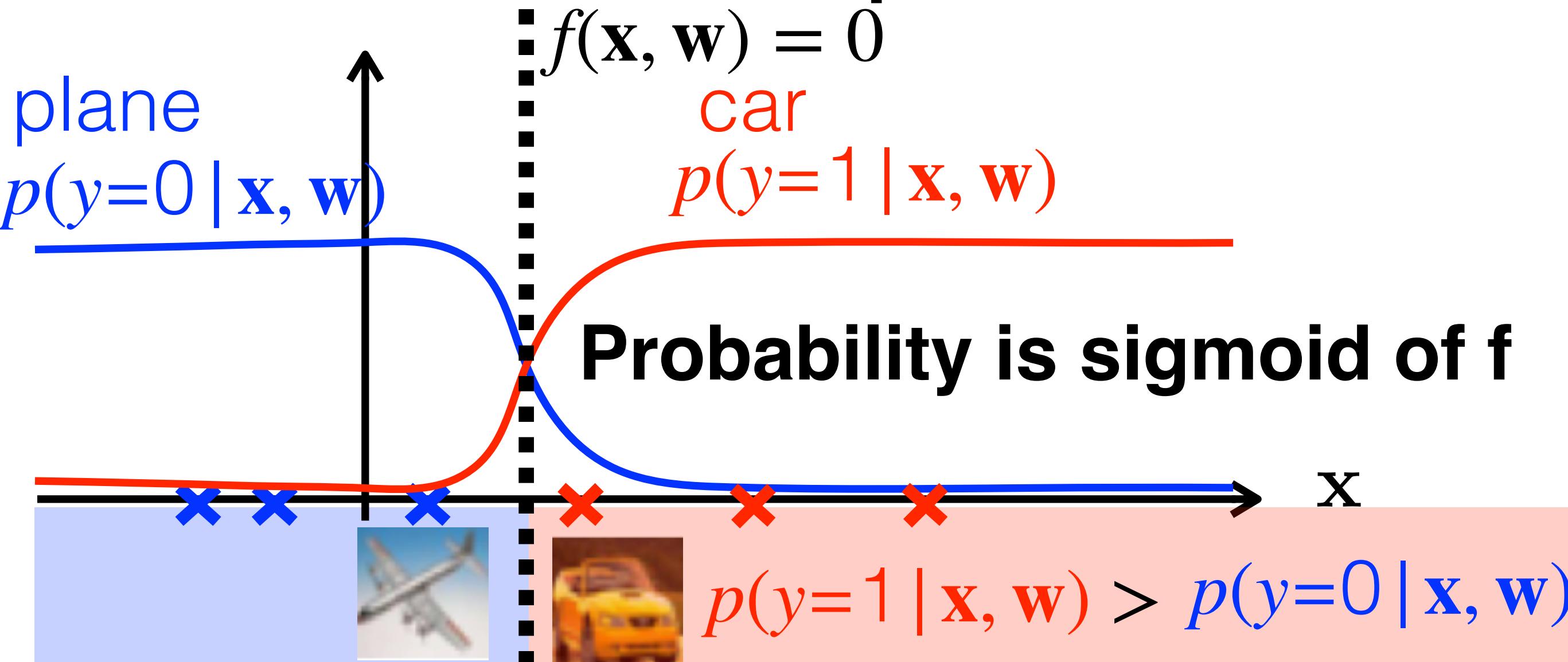
$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = 1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = 0 \end{cases} = y \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w})))$$

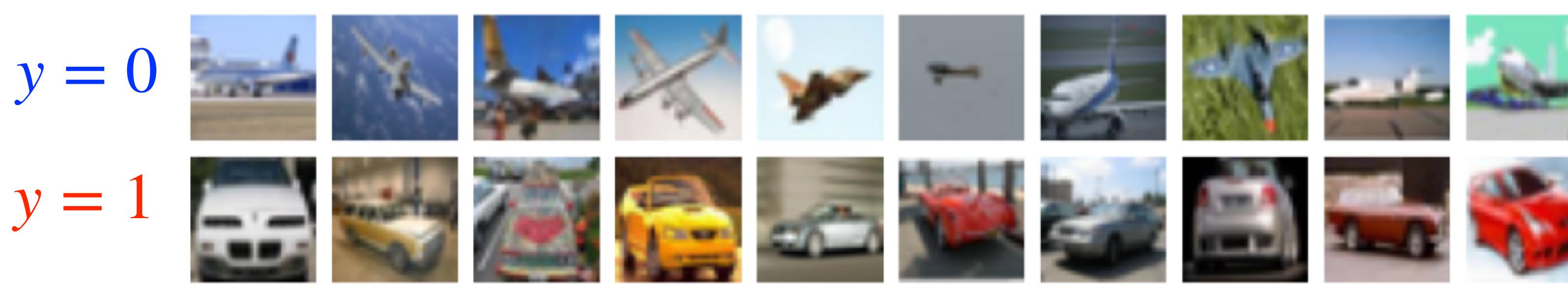


1D linear classifier

2D linear classifier

Which value of f corresponds to this threshold?

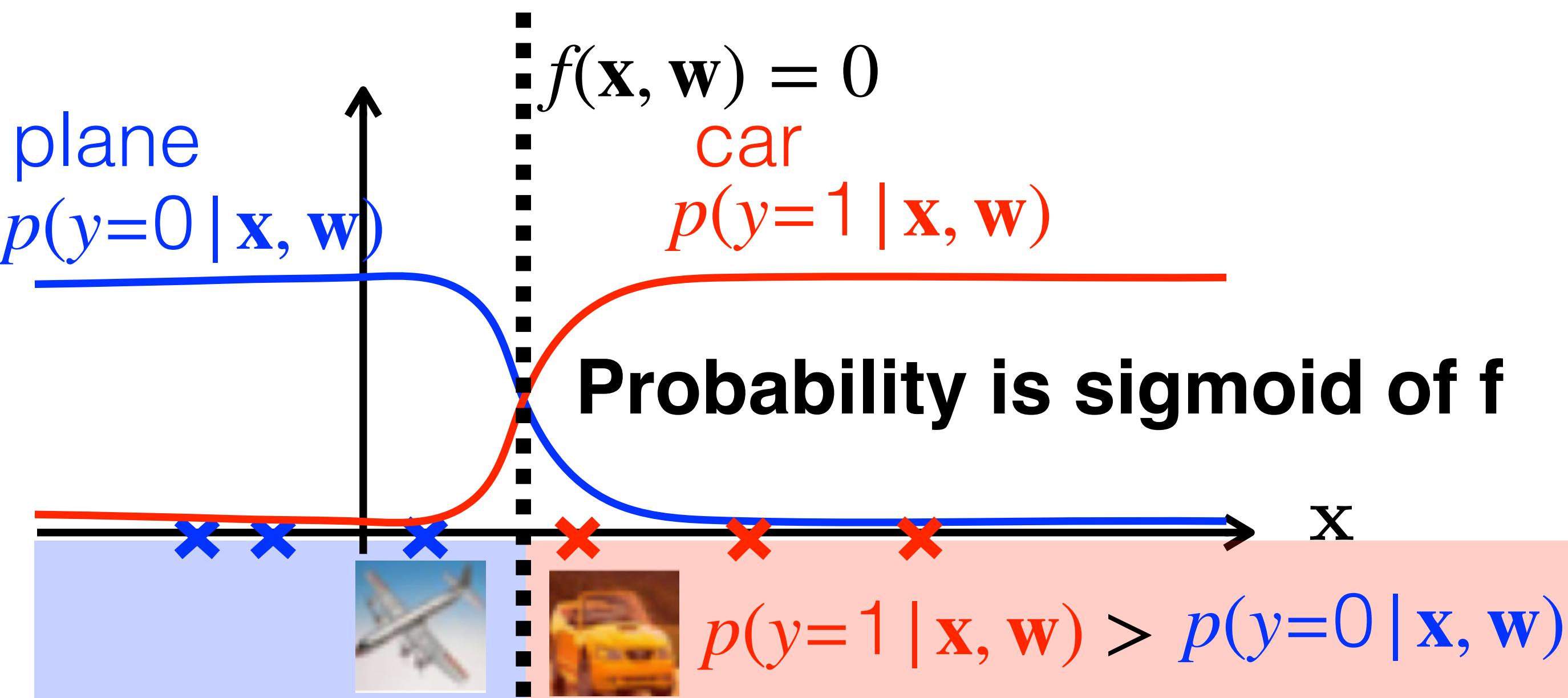




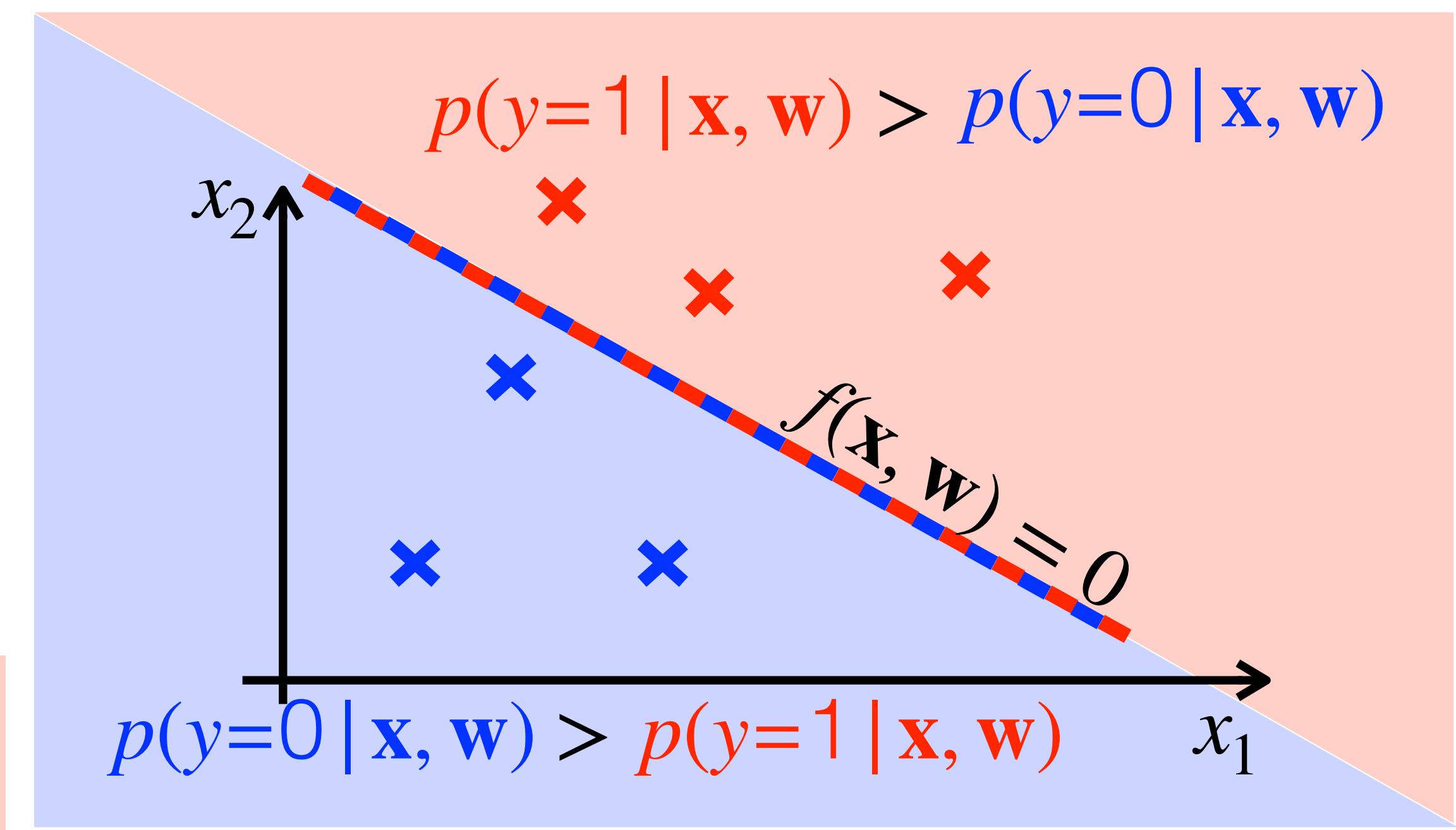
$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = 1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = 0 \end{cases} = y \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w})))$$

Let's implement it !!!

1D linear classifier



2D linear classifier



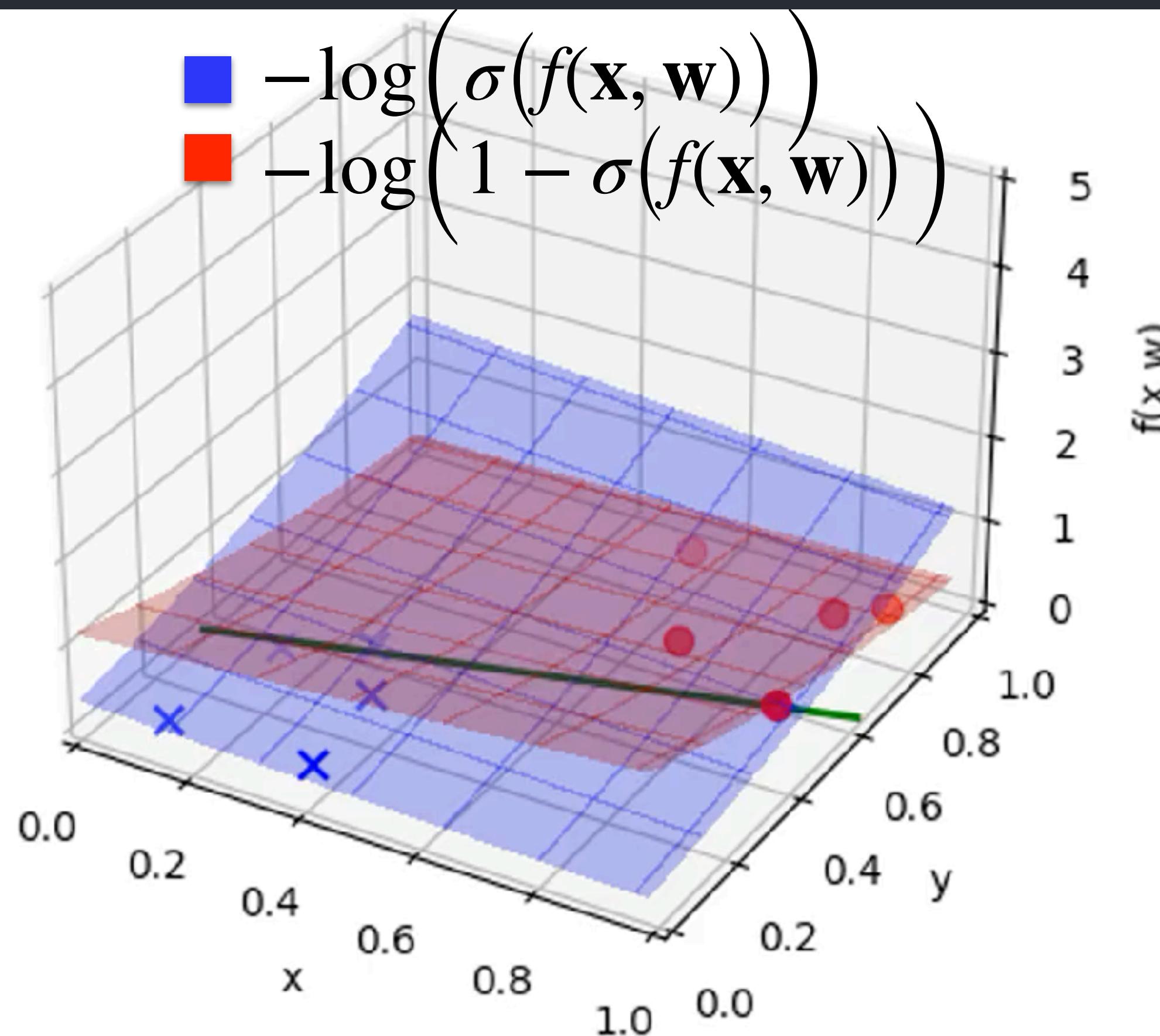
Numpy implementation

```
for i in range(30):
    f = w[0] * x[:, 0] + w[1] * x[:, 1] + w[2]
    p = sigmoid(u)
    loss = (-np.log(p)*y + -np.log(1-p)*(1-y)).sum()
    grad = -1/p * sigmoid(f) * (1-sigmoid(f)) * ...
w = w - 0.1 * grad
```

$$\mathcal{L}(\mathbf{W}) = -\log(y \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w}))))$$

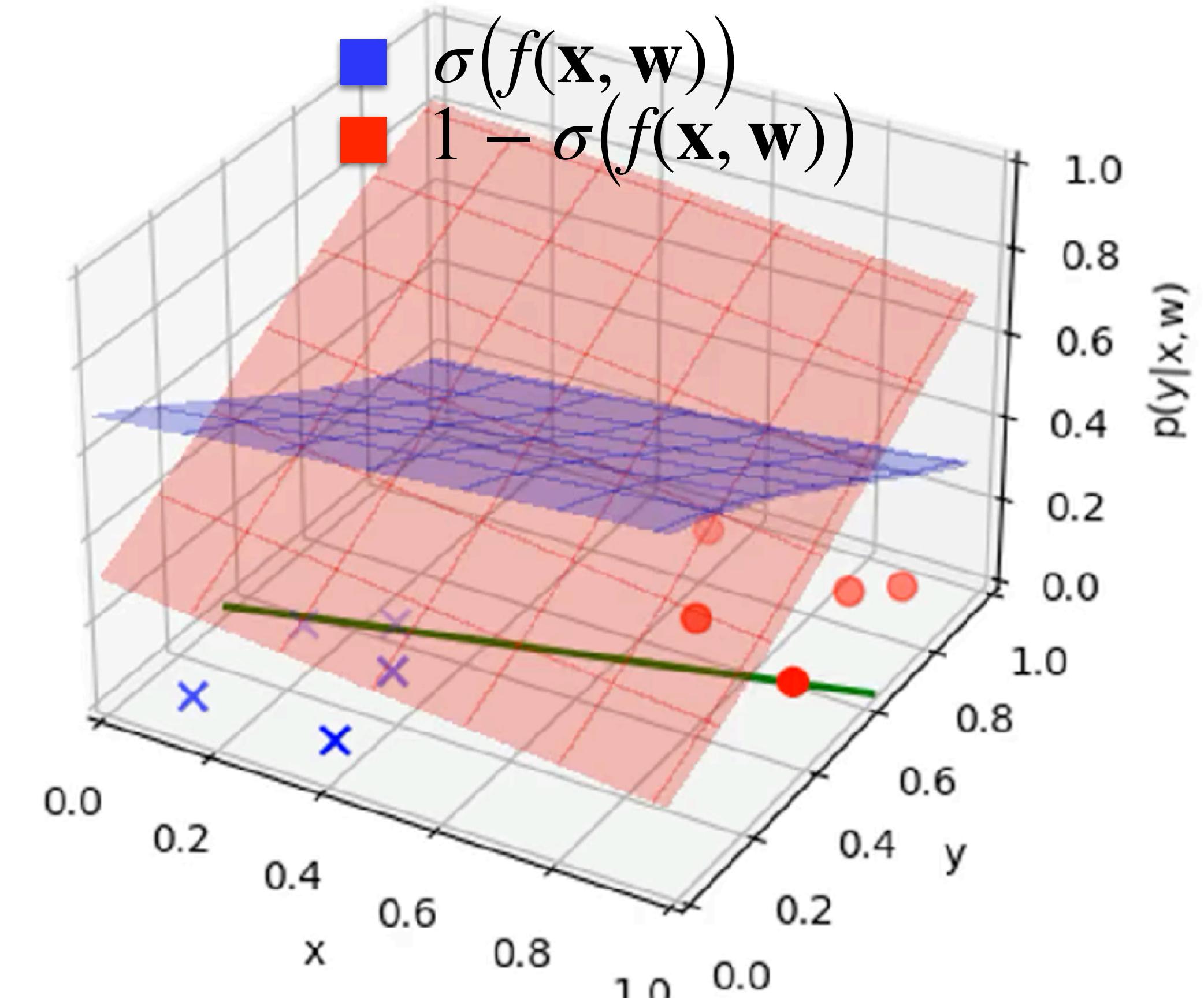
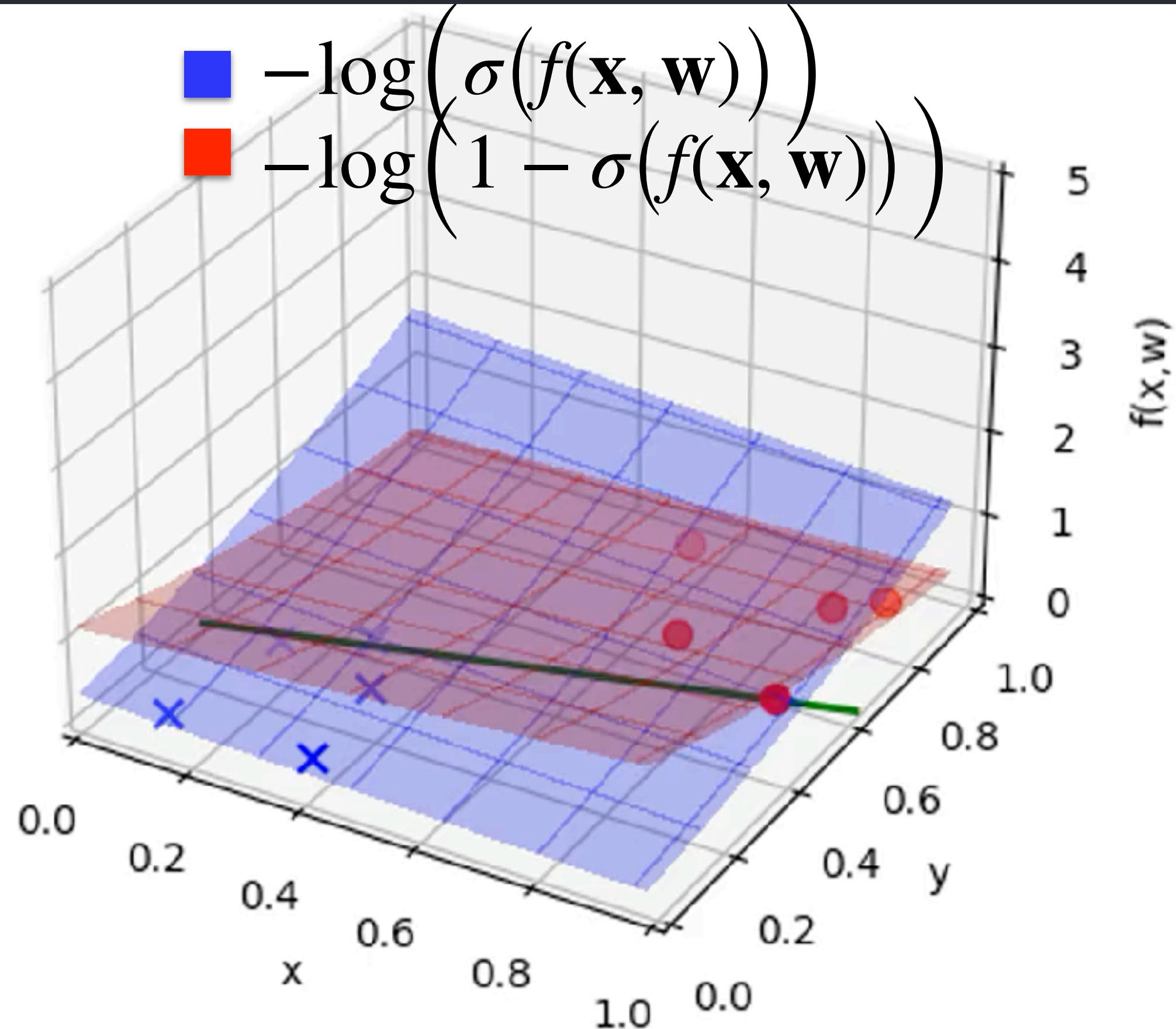
Numpy implementation

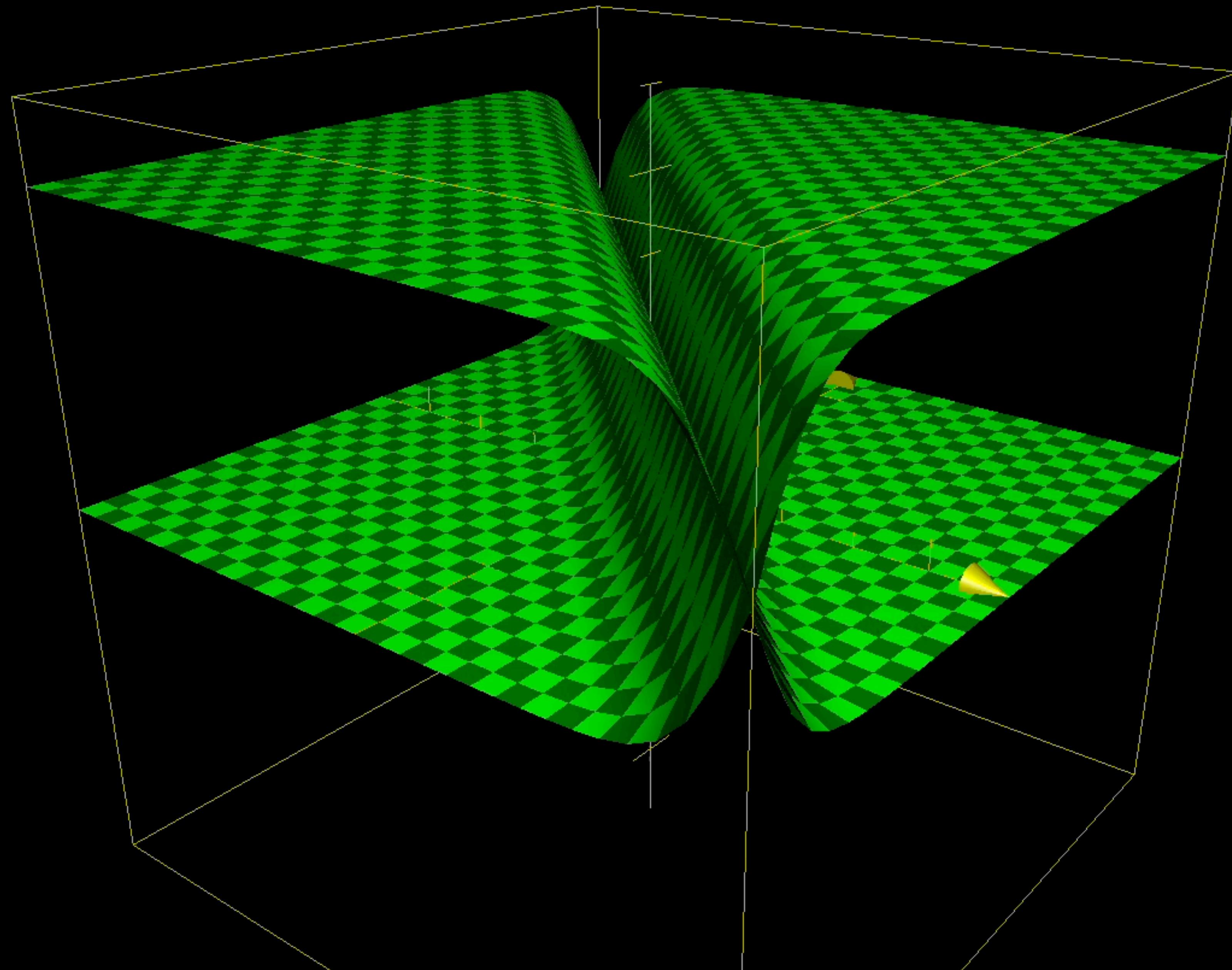
```
for i in range(30):
    f = w[0] * x[:, 0] + w[1] * x[:, 1] + w[2]
    p = sigmoid(u)
    loss = (-np.log(p)*y + -np.log(1-p)*(1-y)).sum()
    grad = -1/p * sigmoid(f) * (1-sigmoid(f)) * ...
w = w - 0.1 * grad
```

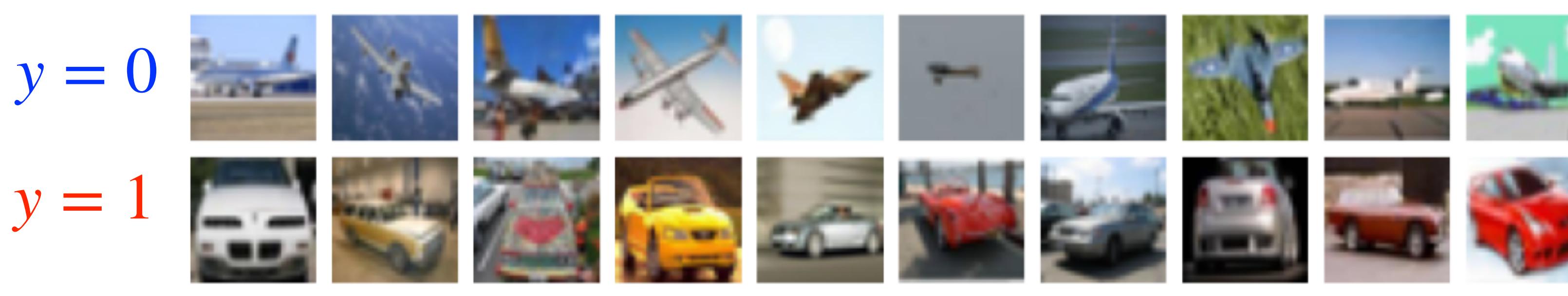


Numpy implementation

```
for i in range(30):
    f = w[0] * x[:, 0] + w[1] * x[:, 1] + w[2]
    p = sigmoid(u)
    loss = (-np.log(p)*y + -np.log(1-p)*(1-y)).sum()
    grad = -1/p * sigmoid(f) * (1-sigmoid(f)) * ...
w = w - 0.1 * grad
```



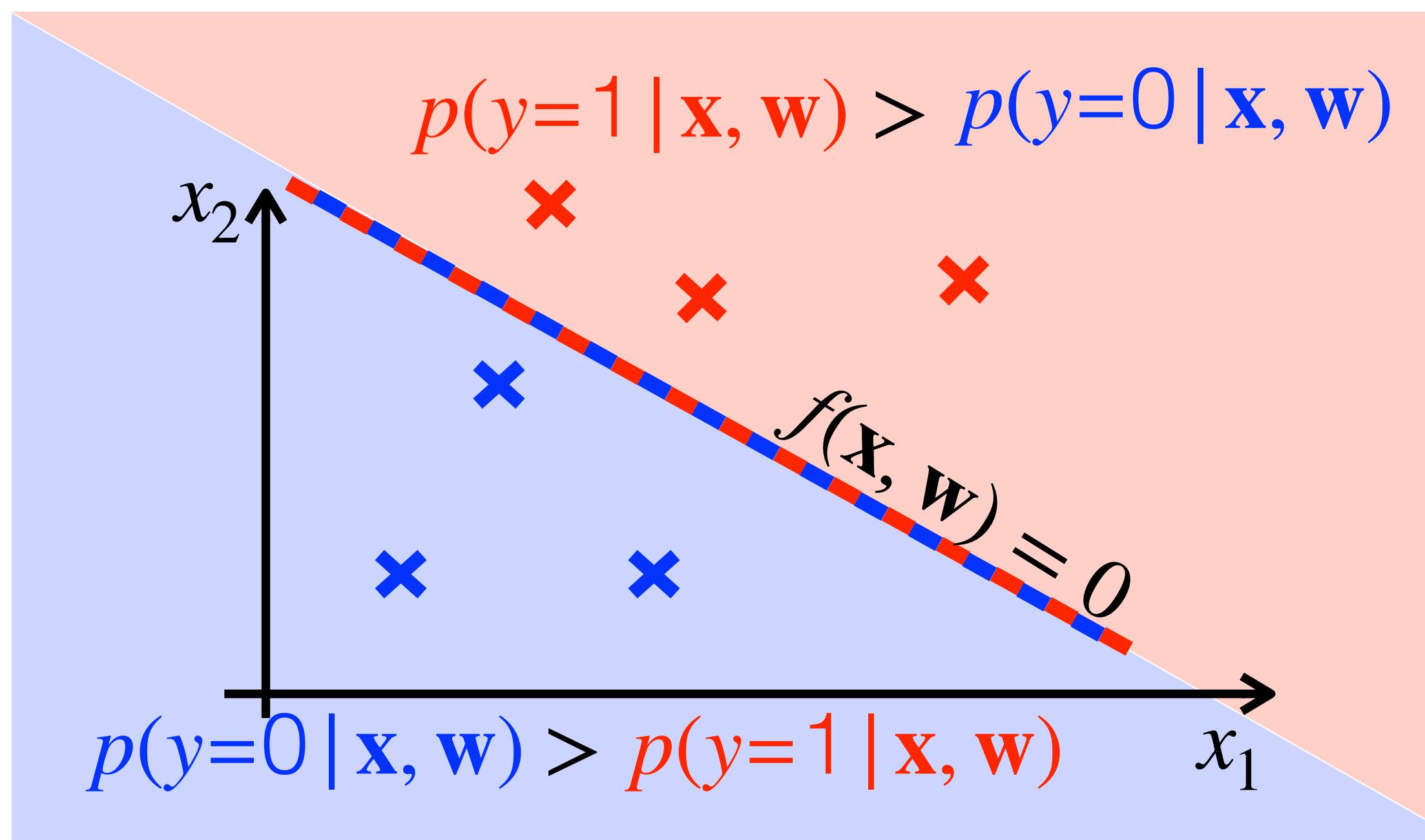




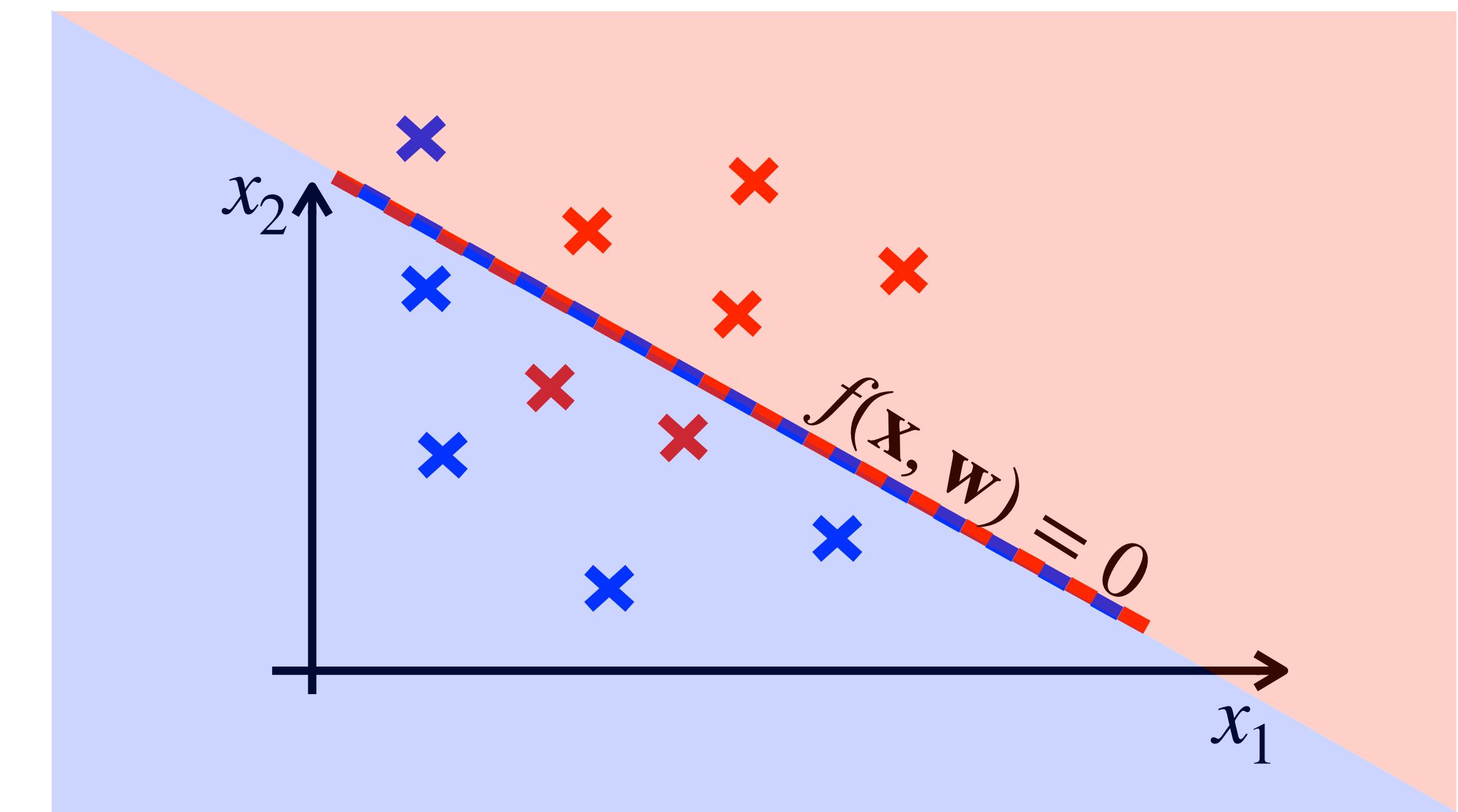
$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = 1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = 0 \end{cases} = y_i \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y_i) \cdot (1 - \sigma(f(\mathbf{x}, \mathbf{w})))$$

Draw training data that cannot be separated by linear classifier????

2D linear classifier



2D linear classifier

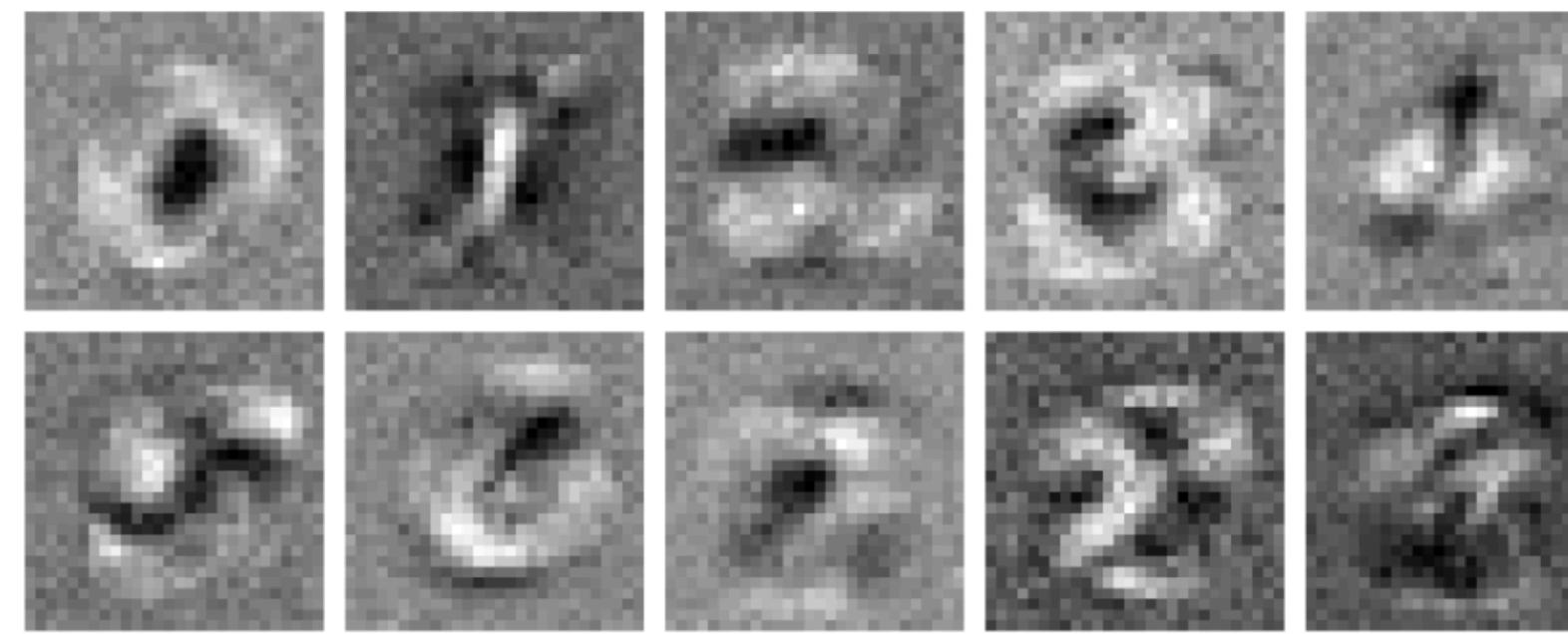


Label

Images

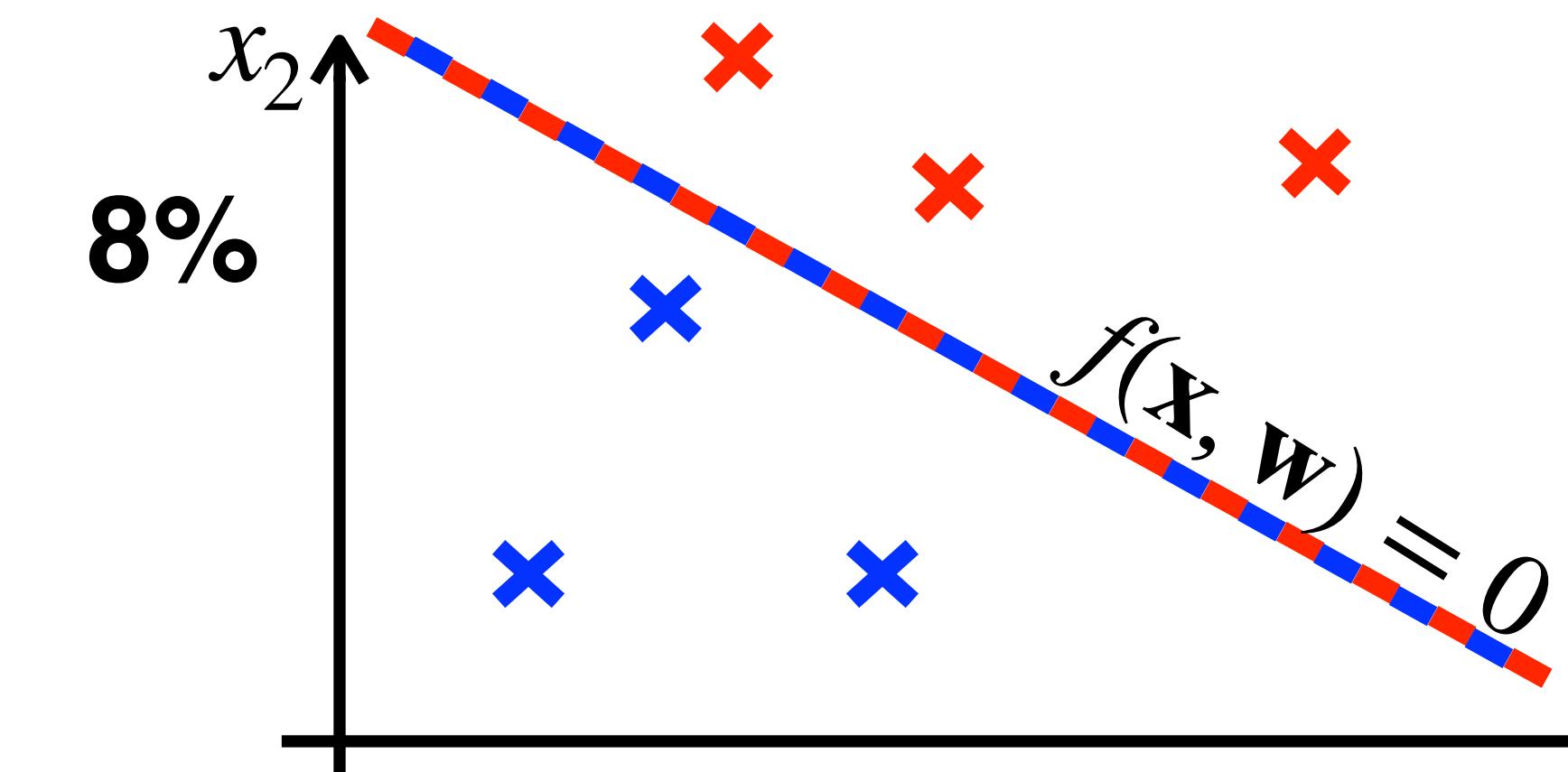


Learned weights



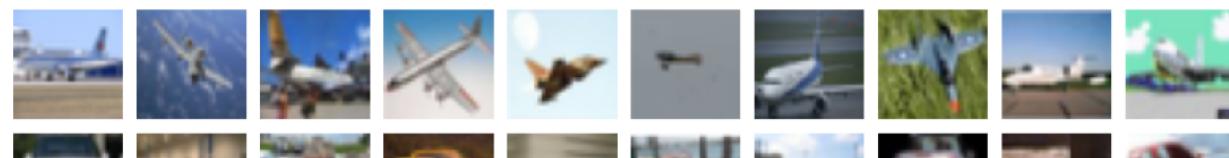
Error

2D linear classifier



2D linear classifier

airplane



automobile



bird



cat



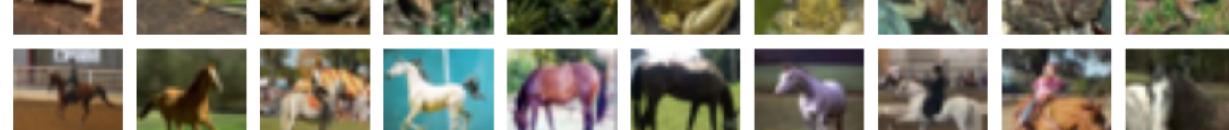
deer



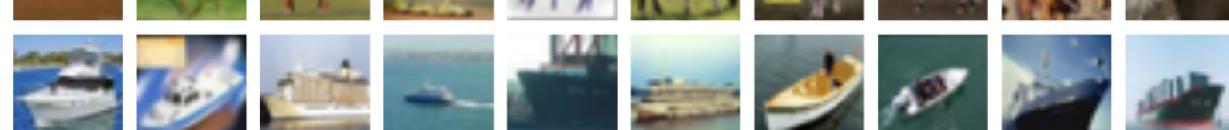
dog



frog



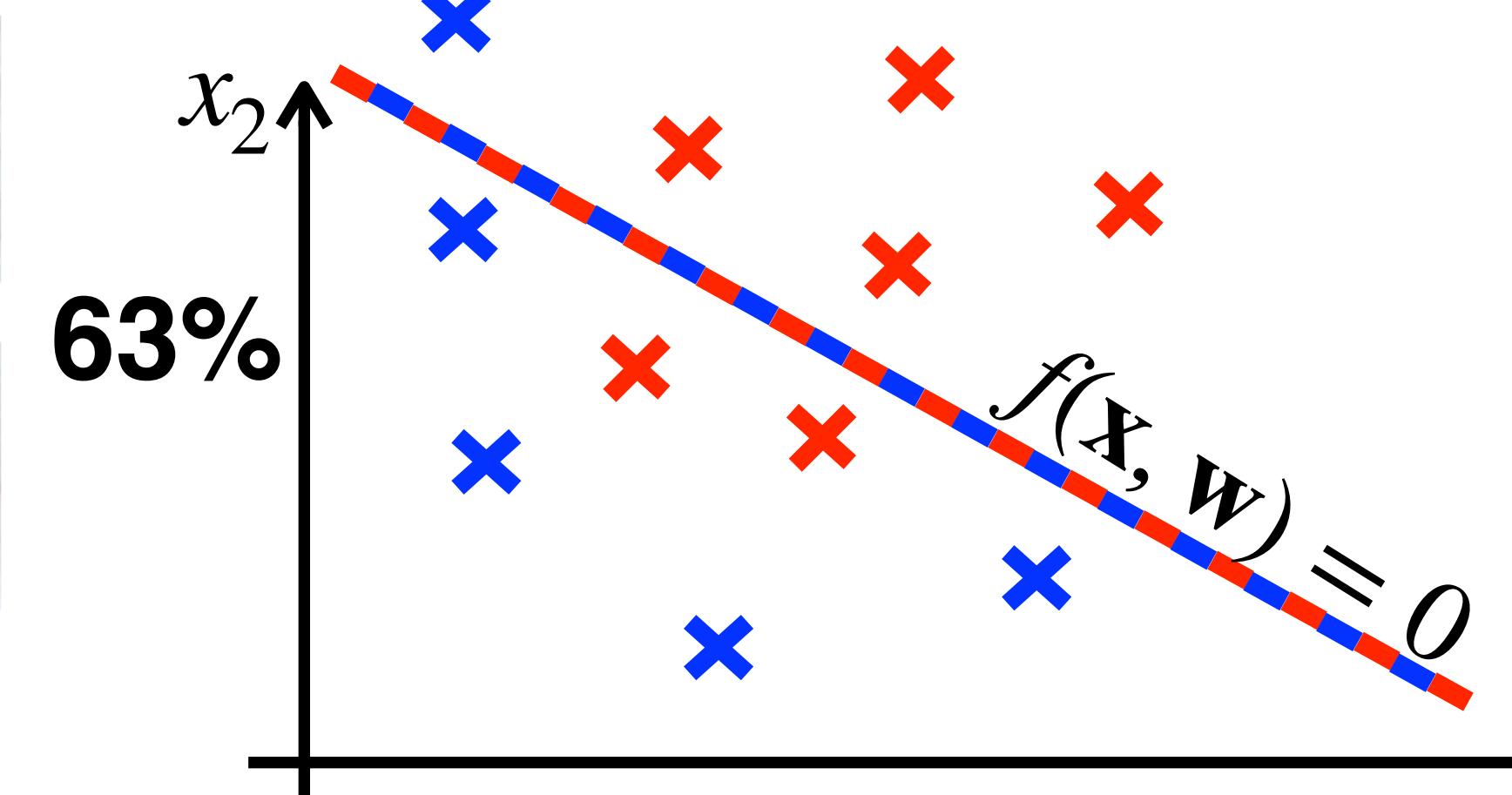
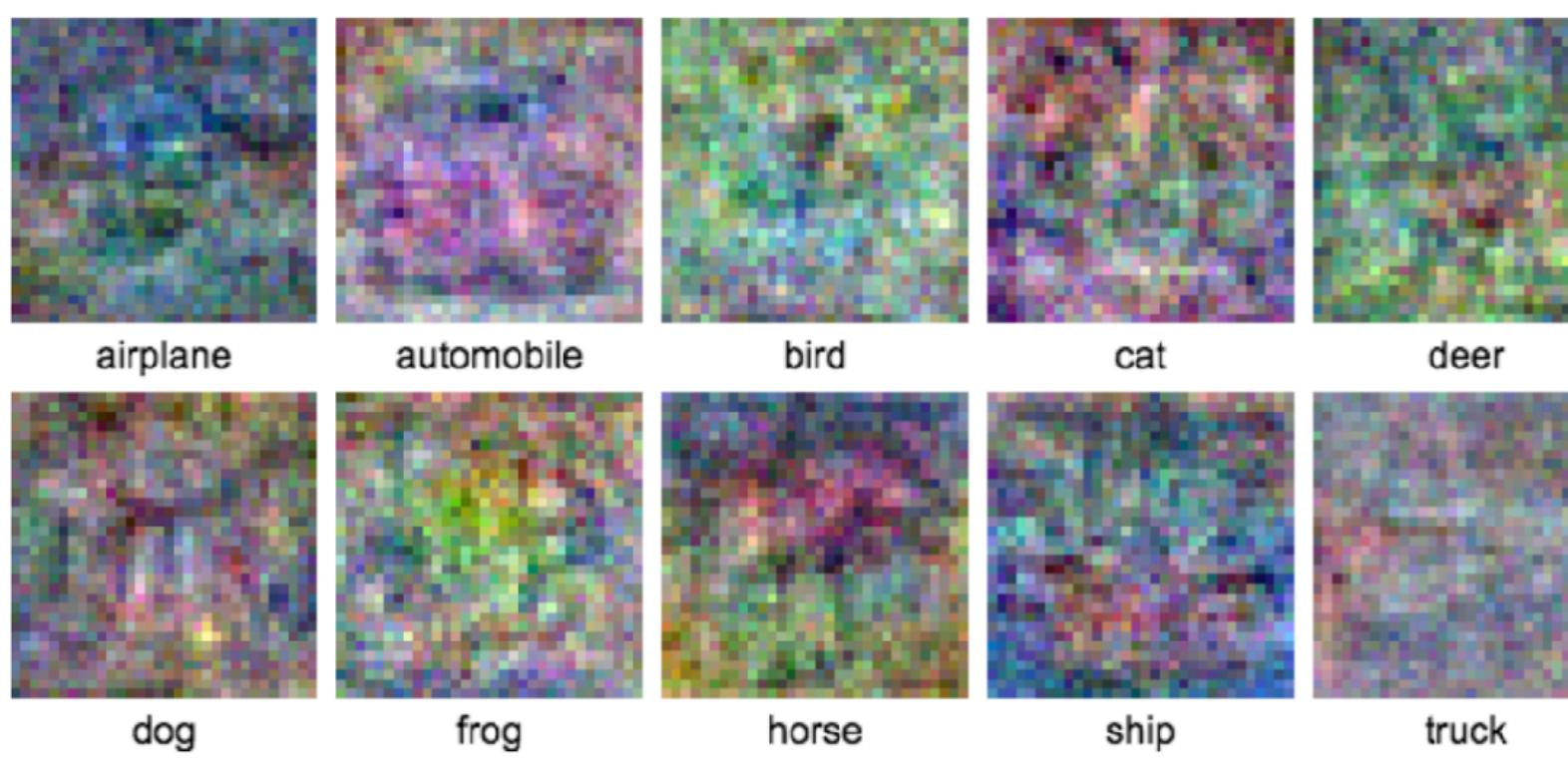
horse



ship

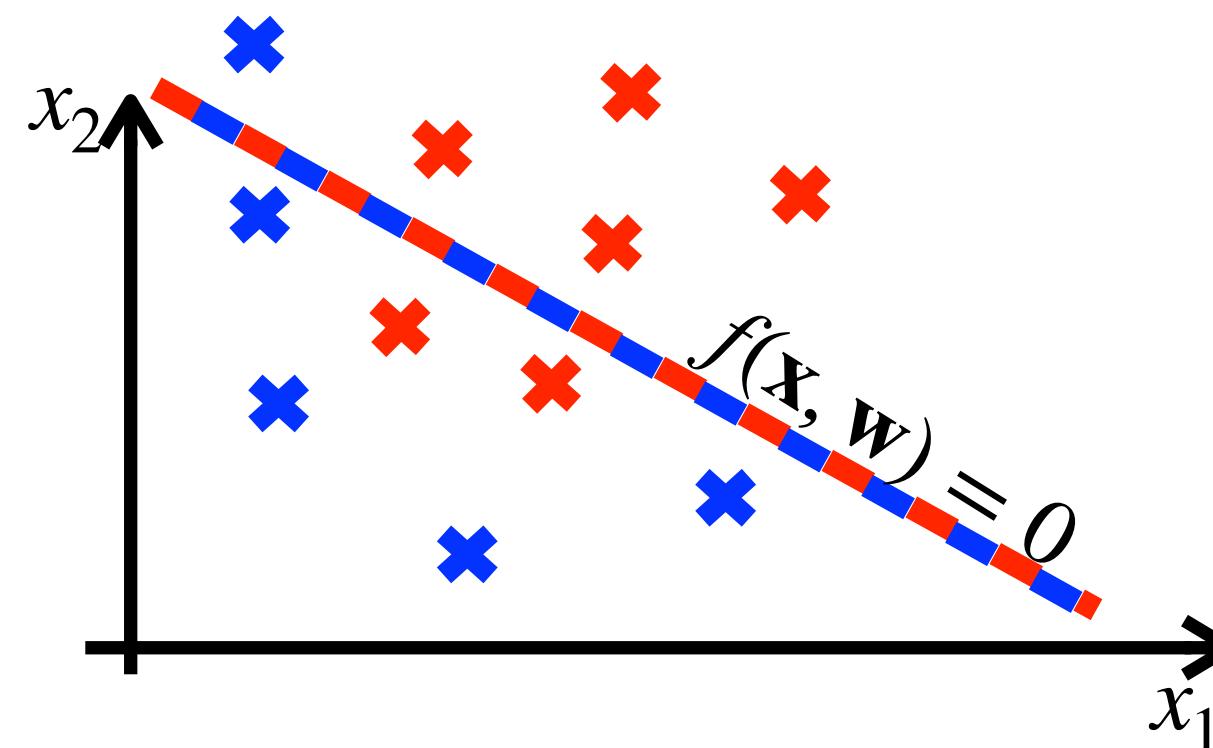


truck



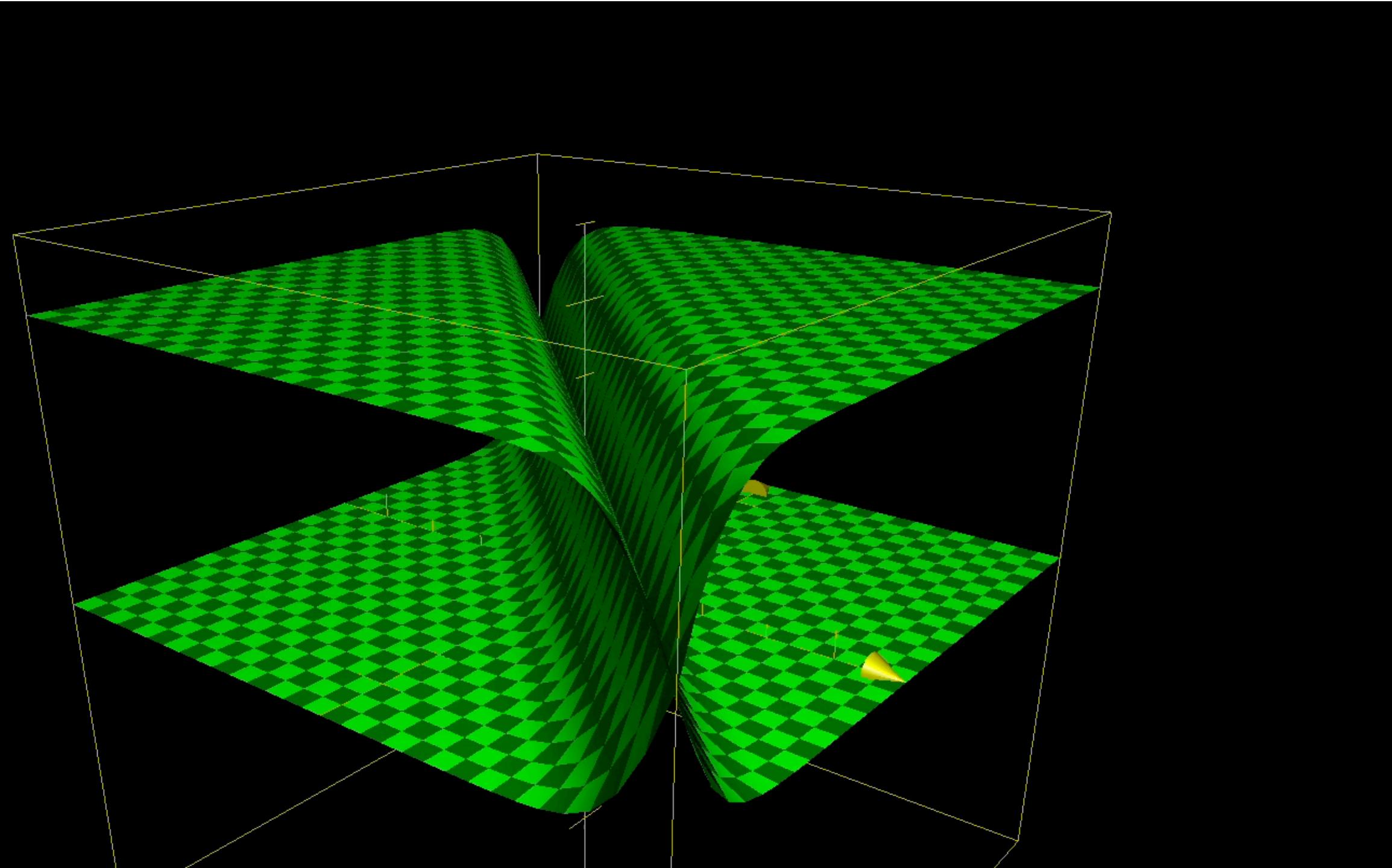
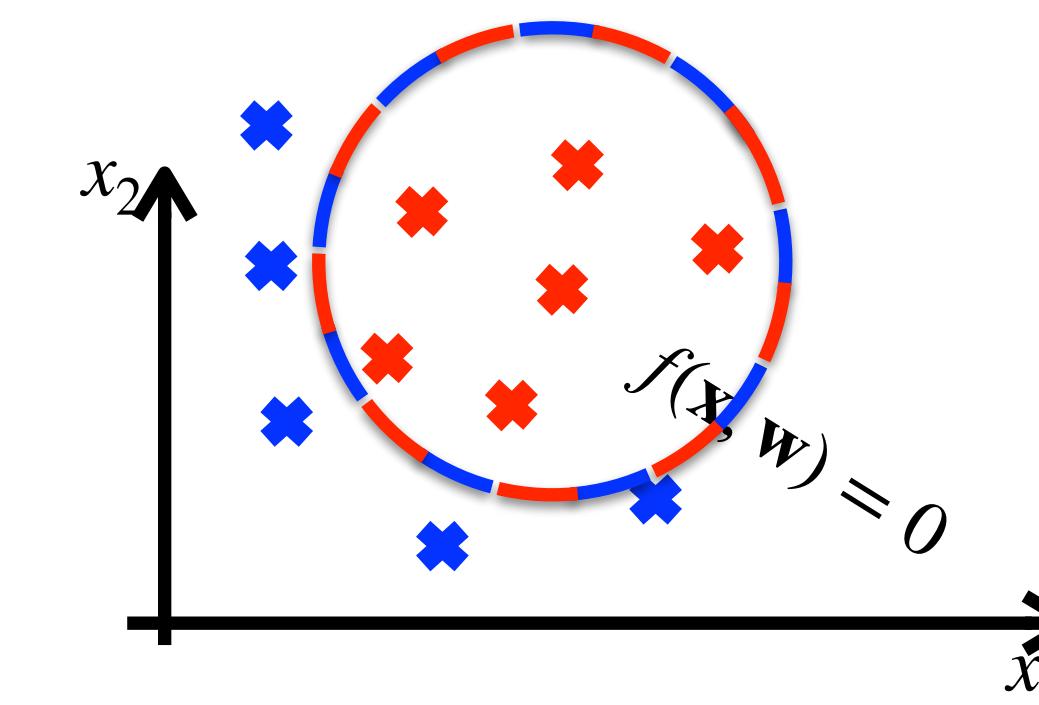
2D linear classifier

$$f([x_1, x_2], \mathbf{w}) = w_1 x_1 + w_2 x_2 + w_0$$



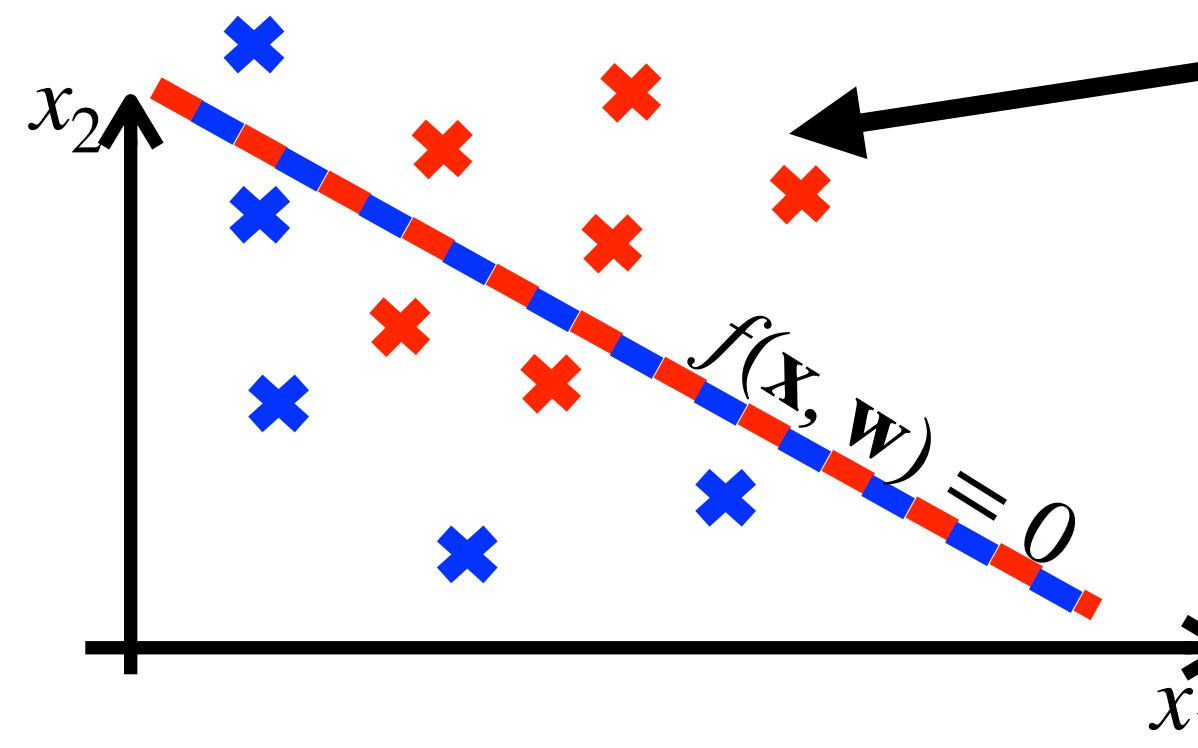
2D non-linear classifier

$$\begin{aligned} f([x_1, x_2], \mathbf{w}) &= \\ &= w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2 + w_0 \end{aligned}$$



2D linear classifier

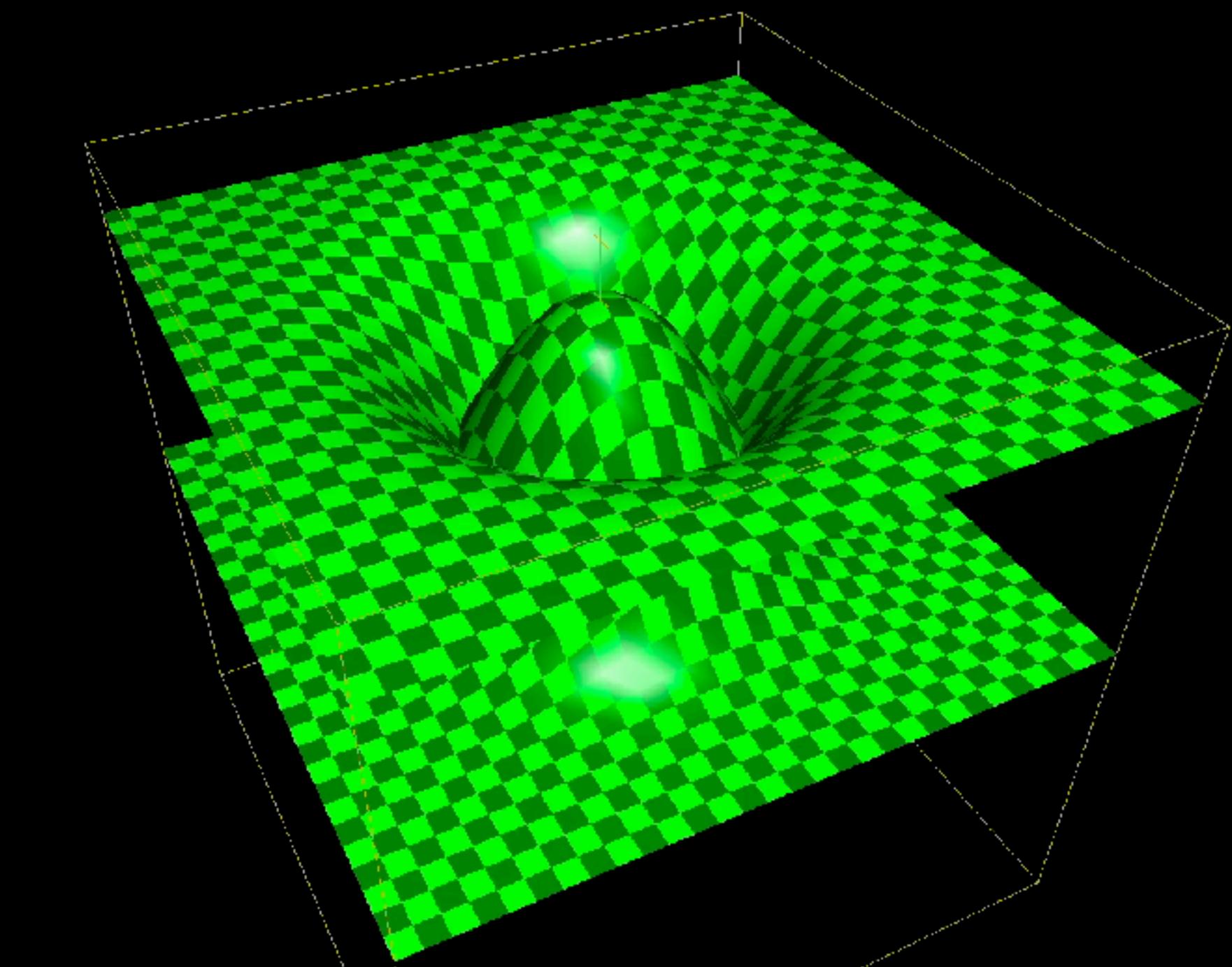
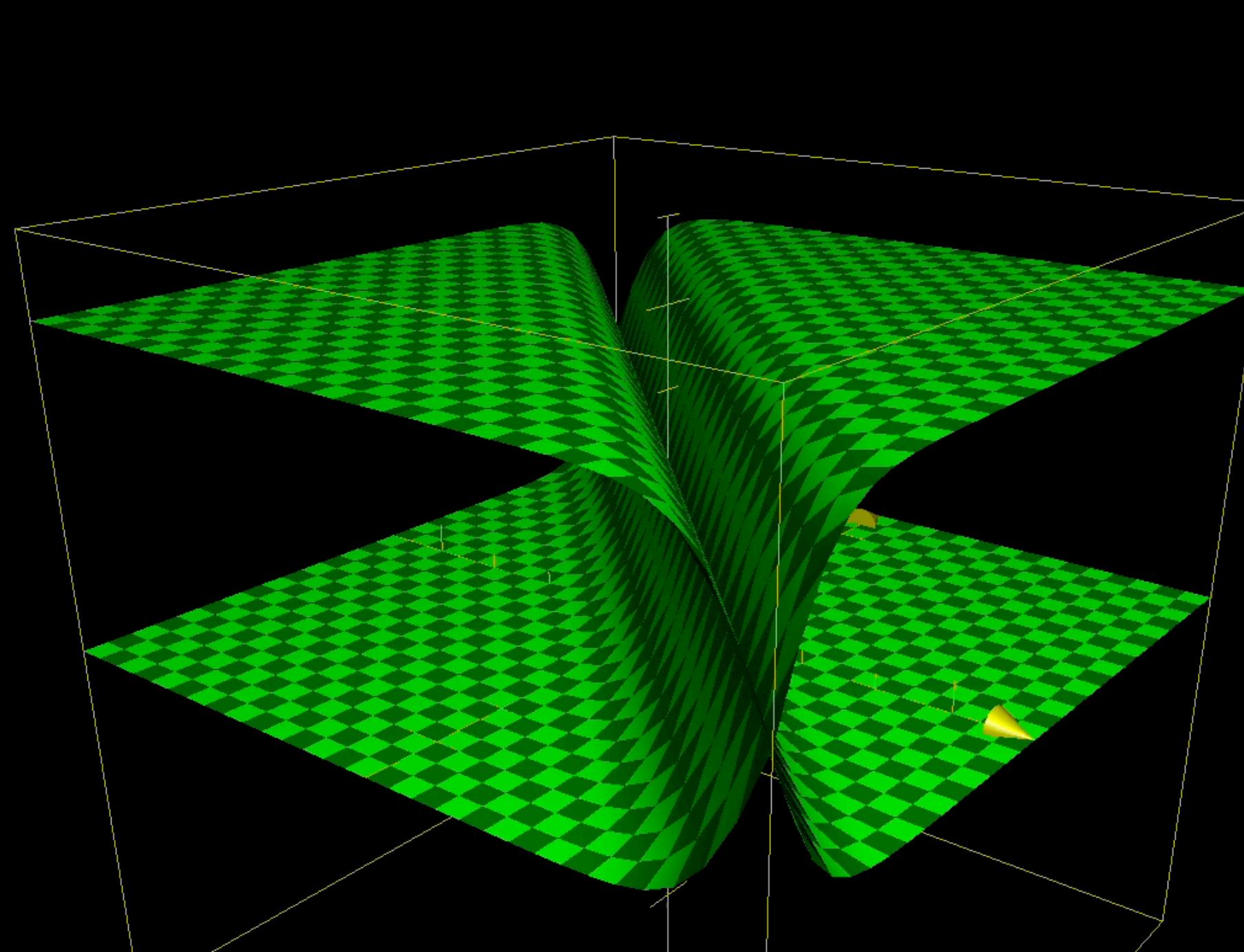
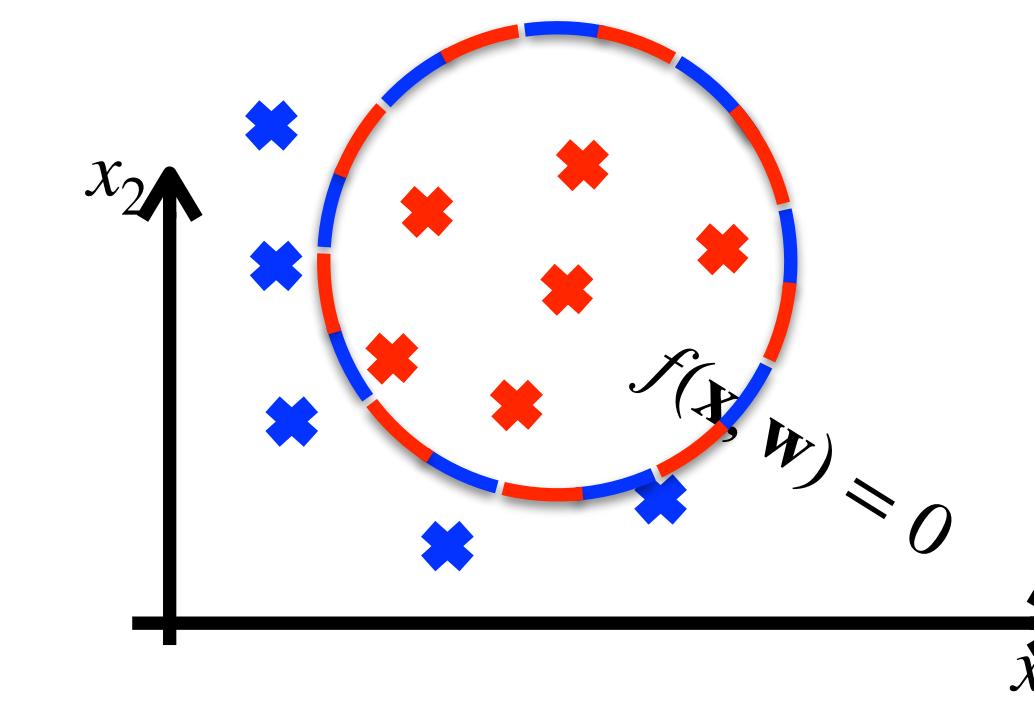
$$f([x_1, x_2], \mathbf{w}) = w_1 x_1 + w_2 x_2 + w_0$$



**Is there a linear classifier
that separates this data?**

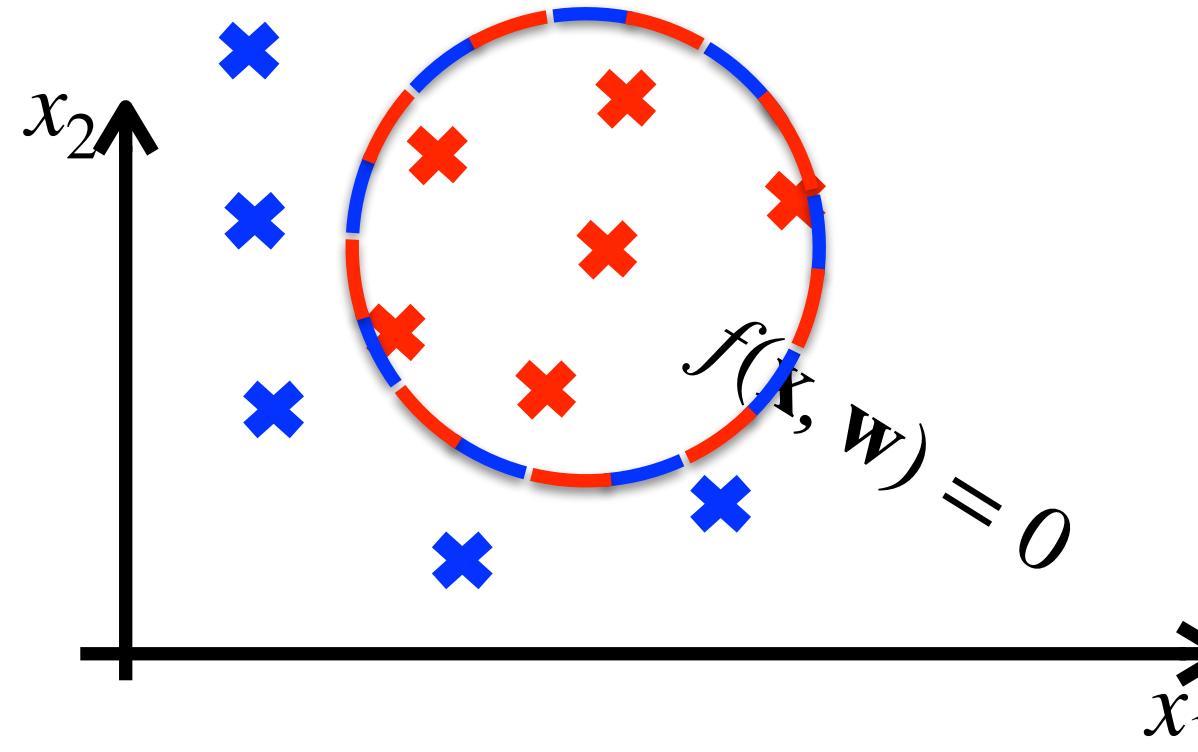
2D non-linear classifier

$$\begin{aligned} f([\mathbf{x}_1, \mathbf{x}_2], \mathbf{w}) = \\ = w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2 + w_0 \end{aligned}$$



5D linear classifier

$$f([x_1, x_2, x_1^2, x_2^2, x_1x_2], \mathbf{w}) = \\ = w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2 + w_0$$

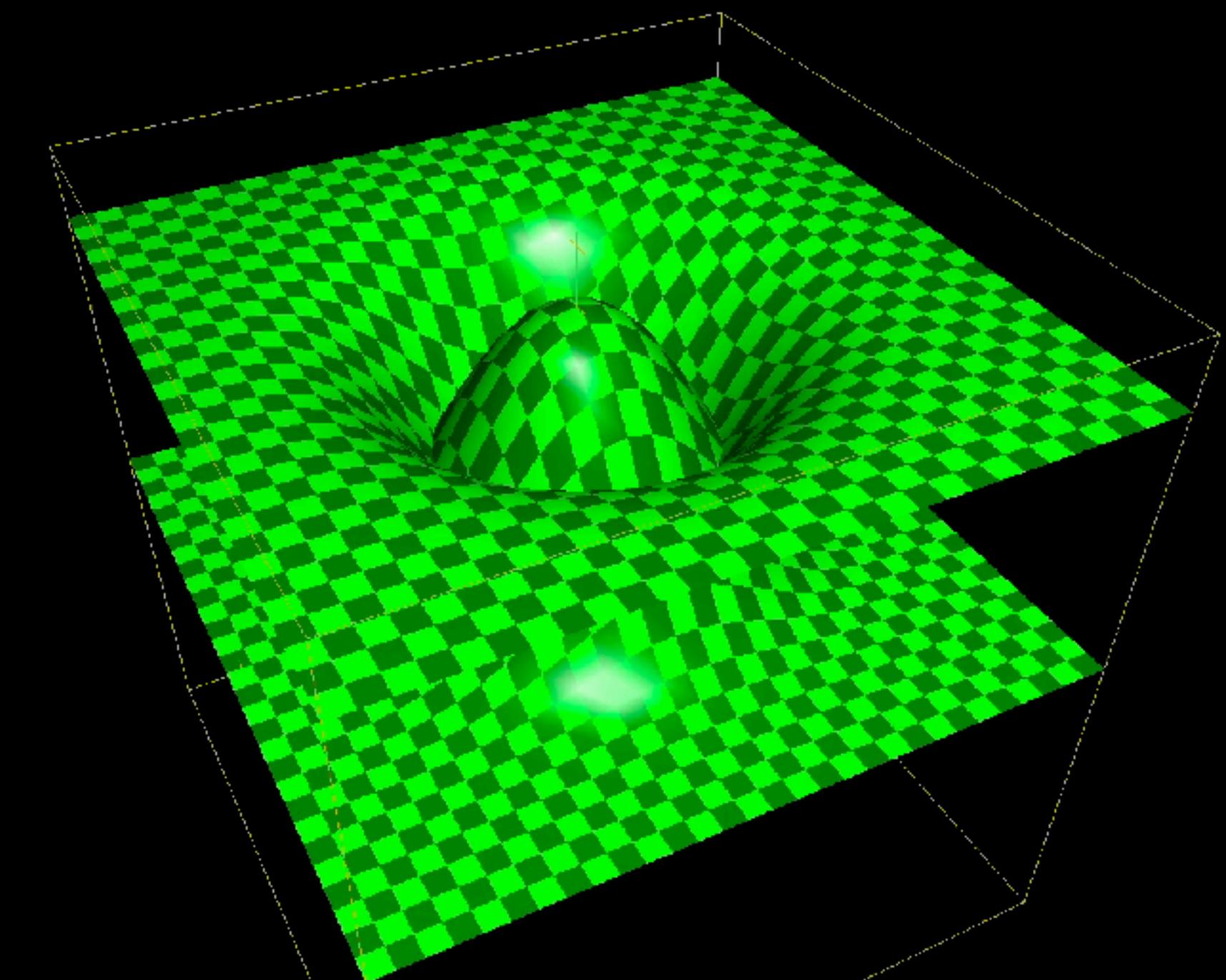
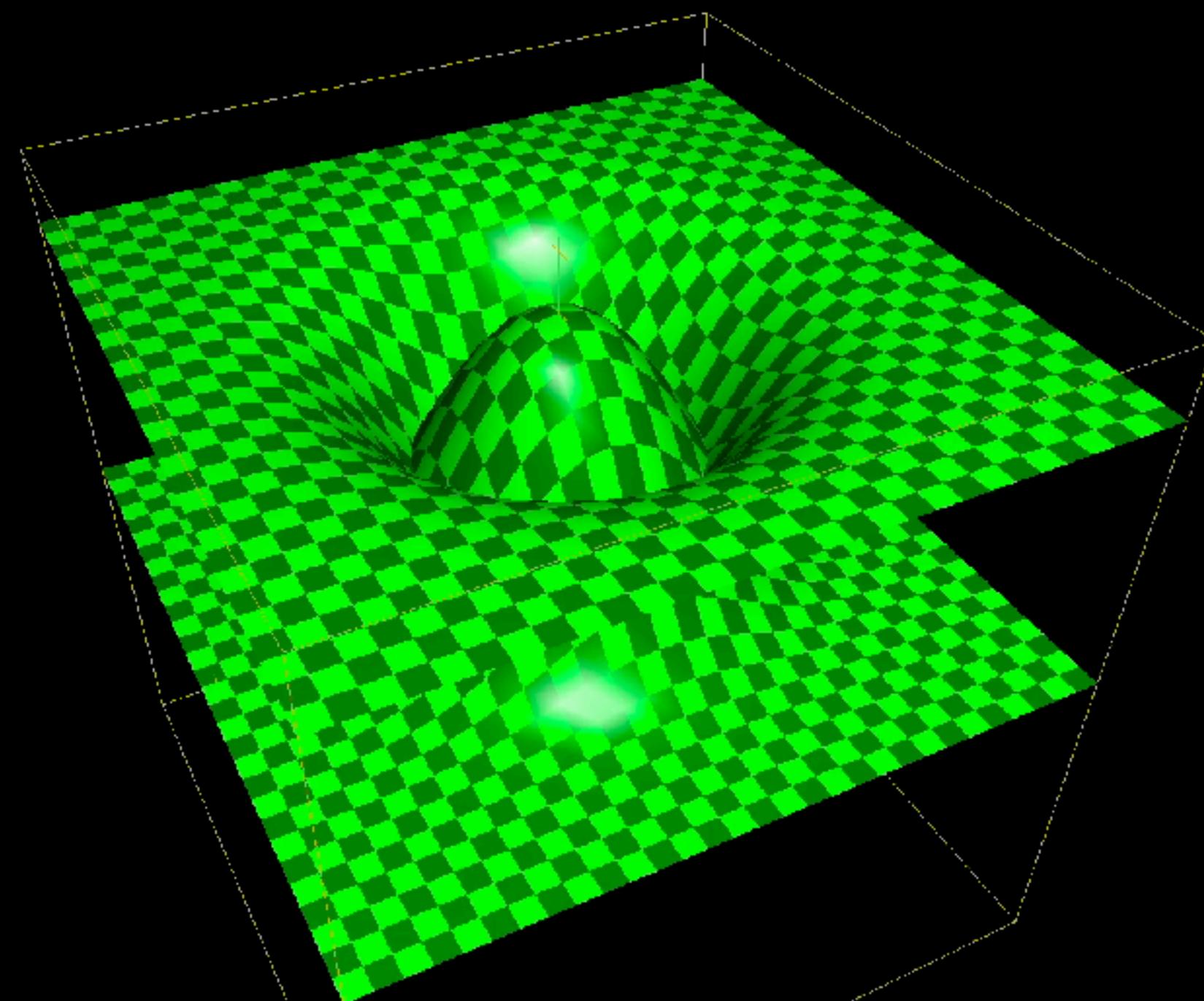
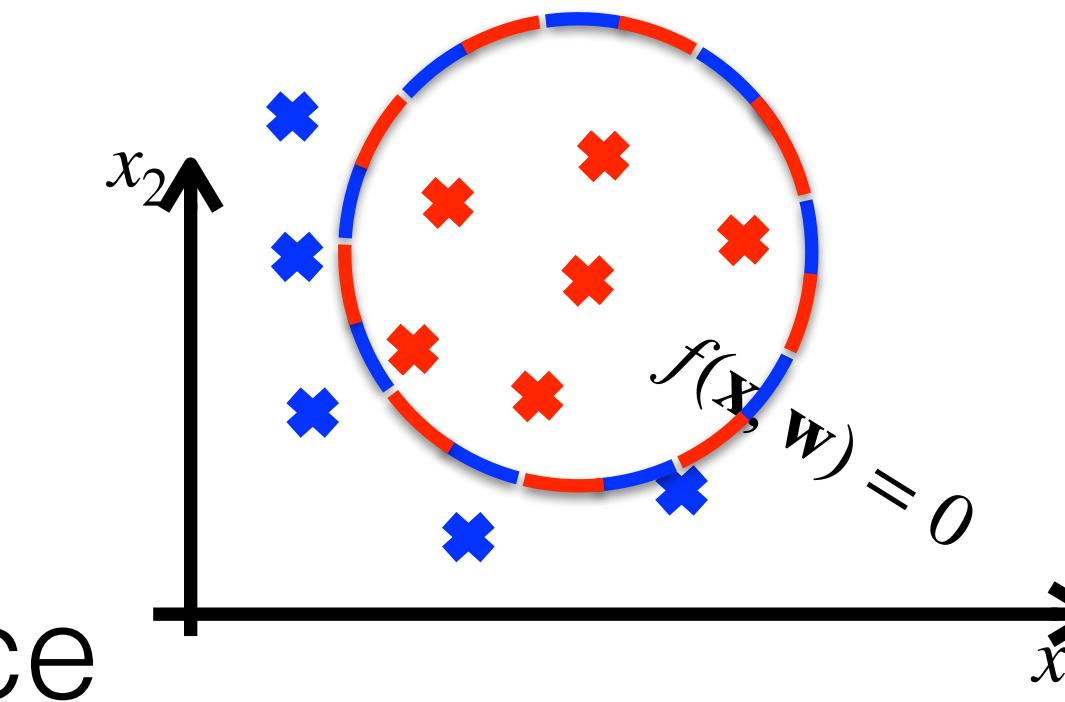


2D non-linear classifier

$$f([x_1, x_2], \mathbf{w}) = \\ = w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2 + w_0$$

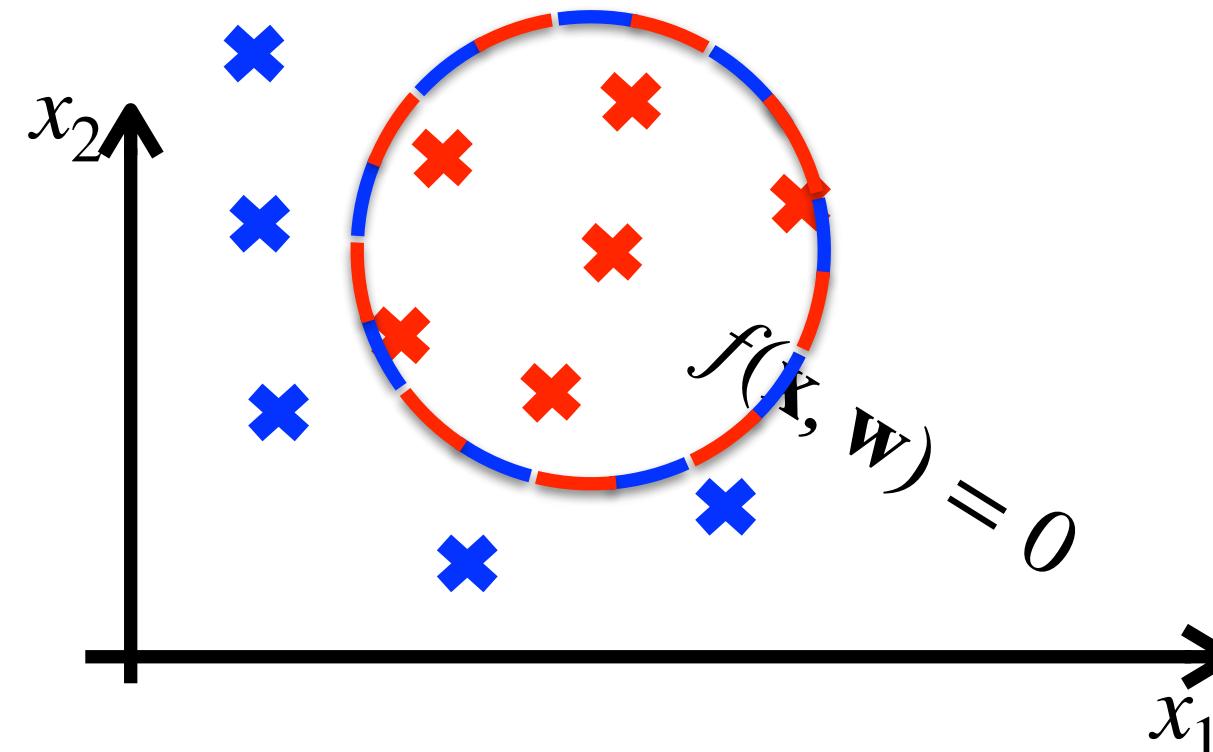
**Is there a linear classifier
that separates this data?**

- Taylor polynomia
- Fourier transform
- Hilbert or Banach space



5D linear classifier

$$f([x_1, x_2, x_1^2, x_2^2, x_1x_2], \mathbf{w}) = \\ = w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2 + w_0$$

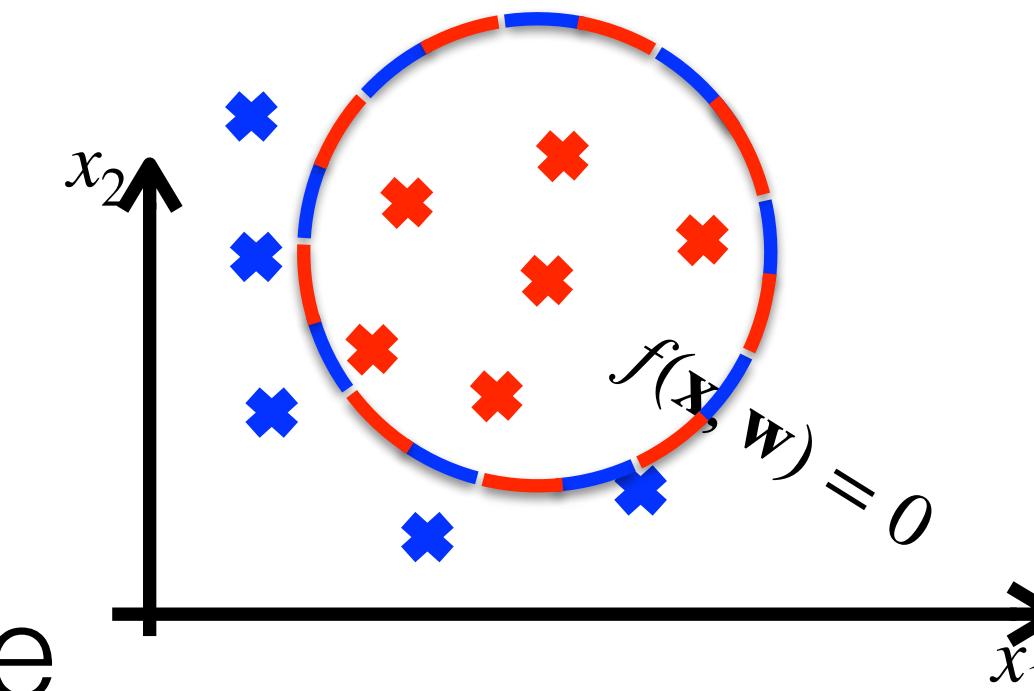


2D non-linear classifier

$$f([x_1, x_2], \mathbf{w}) = \\ = w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2 + w_0$$

**Is there a linear classifier
that separates this data?**

- Taylor polynomia
- Fourier transform
- Hilbert or Banach space



There are two ways you can understand deep learning:

- Learn **non-linear classifier** that separates complicated areas in pixel space
- Learn **non-linear features** that makes classes linearly separable

Equivalence of common binary losses

Binary cross-entropy loss:

$$p(y | \mathbf{x}, \mathbf{w}) = \begin{cases} \sigma(f(\mathbf{x}, \mathbf{w})) & y = 1 \\ 1 - \sigma(f(\mathbf{x}, \mathbf{w})) & y = 0 \end{cases} = y_i \cdot \sigma(f(\mathbf{x}, \mathbf{w})) + (1 - y_i) \cdot \sigma(f(\mathbf{x}, \mathbf{w}))$$

$$\mathcal{L}(\mathbf{W}) = \sum_i -\log p(y_i | \mathbf{x}_i, \mathbf{w}) = \sum_i -\log(y_i \cdot \sigma(f(\mathbf{x}_i, \mathbf{w})) + (1 - y_i) \cdot \sigma(f(\mathbf{x}_i, \mathbf{w})))$$

```
loss = y*(-torch.log(torch.sigmoid(w@x)))  
      + (1-y)*(-torch.log(1-torch.sigmoid(w@x)))  
tensor(3.9876, requires_grad=True)
```

Logistic loss:

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

$$\mathcal{L}(\mathbf{W}) = \sum_i -\log p(y_i | \mathbf{x}_i, \mathbf{w}) = \sum_i -\log(\sigma(y_i f(\mathbf{x}_i, \mathbf{w}))) = \sum_i \log(1 + \exp(-y_i f(\mathbf{x}_i, \mathbf{w})))$$

```
loss = -torch.log(torch.sigmoid(y*(w@x)))  
tensor(3.9876, requires_grad=True)  
loss = torch.log(1+torch.exp(-y*(w@x)))  
tensor(3.9876, requires_grad=True)
```

Competencies required for the test T1

- Classification loss for two-class and K-class classification problem.
- Equivalence of common binary classification losses
- Meaning of weights in linear classifier
- 2D visualization, decision boundary, features
- Drawback of linear classifier (too simple for some problems)