

Lecture #7: Intro to neural networks

Goals/Agenda

- From logistic regression to neural networks (NNs)
 - Why use NNs?
 - How to train NNs?
 - Intuition behind backpropagation
 - NNs for materials science

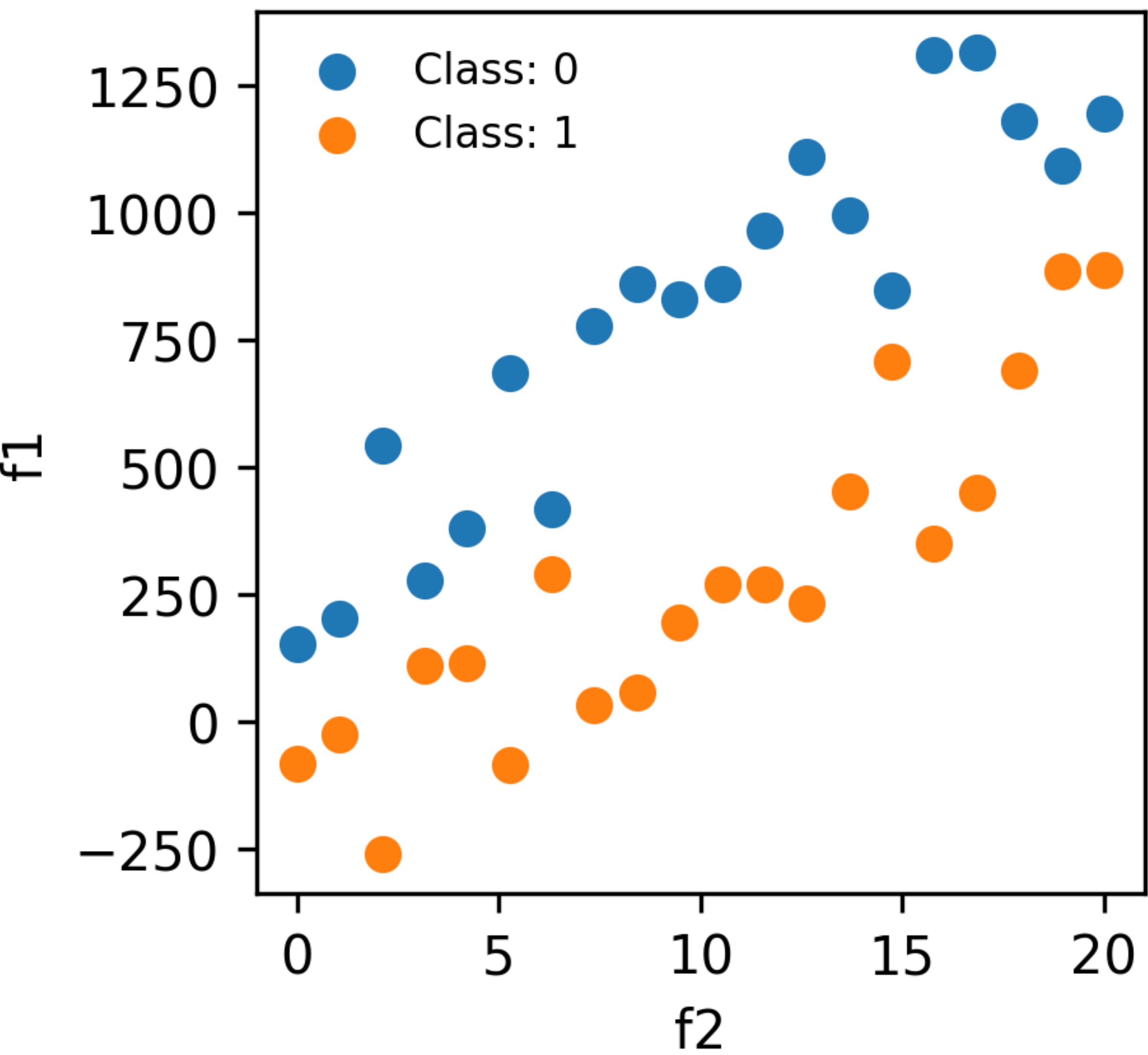
Part #1: Logistic regression

Data samples:

$$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_i, y_i), \dots, (\mathbf{x}_n, y_n)\}$$

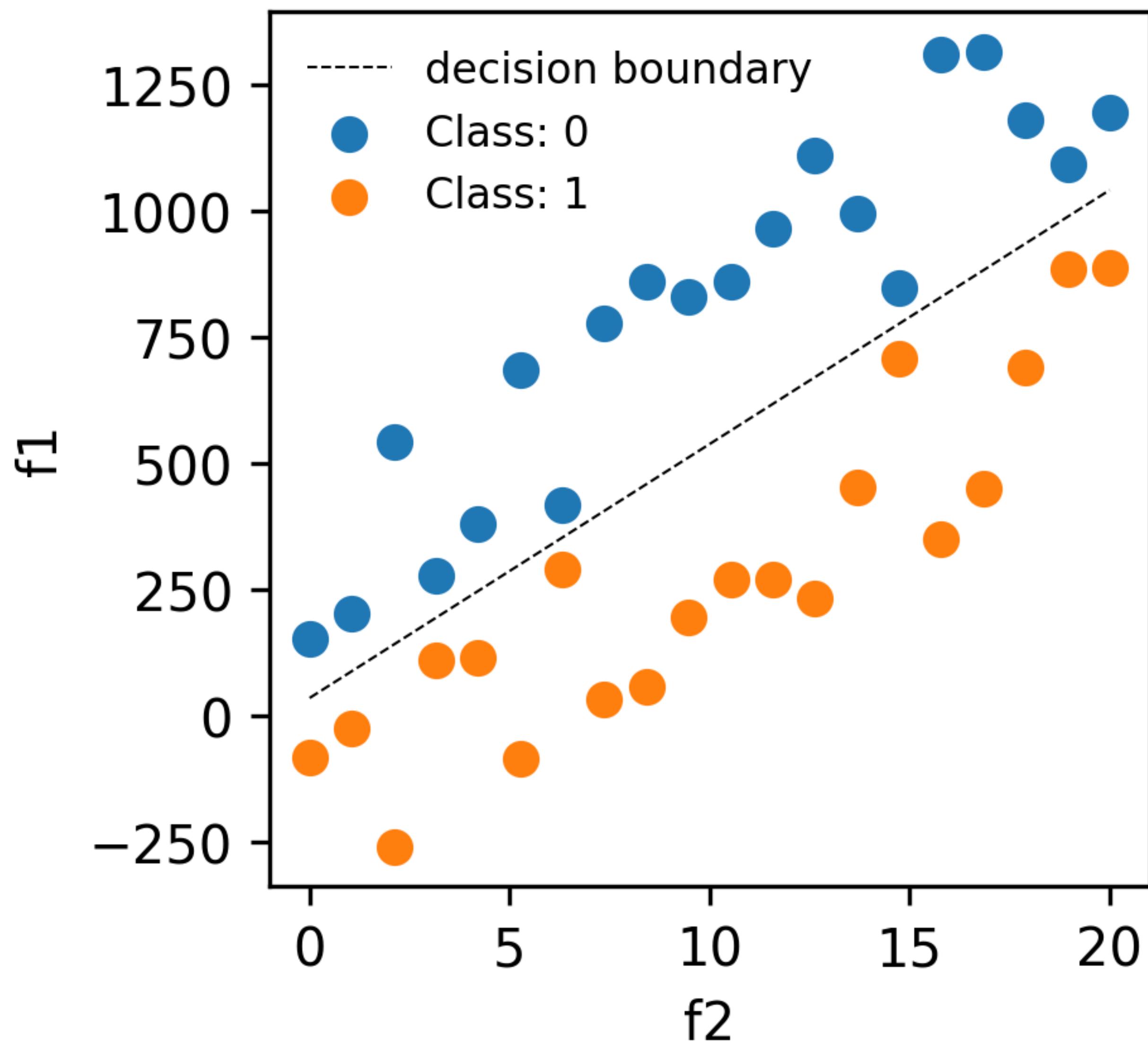
y_i - class of a sample (0 or 1)

$\mathbf{x}_i = (x_i^1, \dots, x_i^d)$ - $1 \times d$ vector of
input features



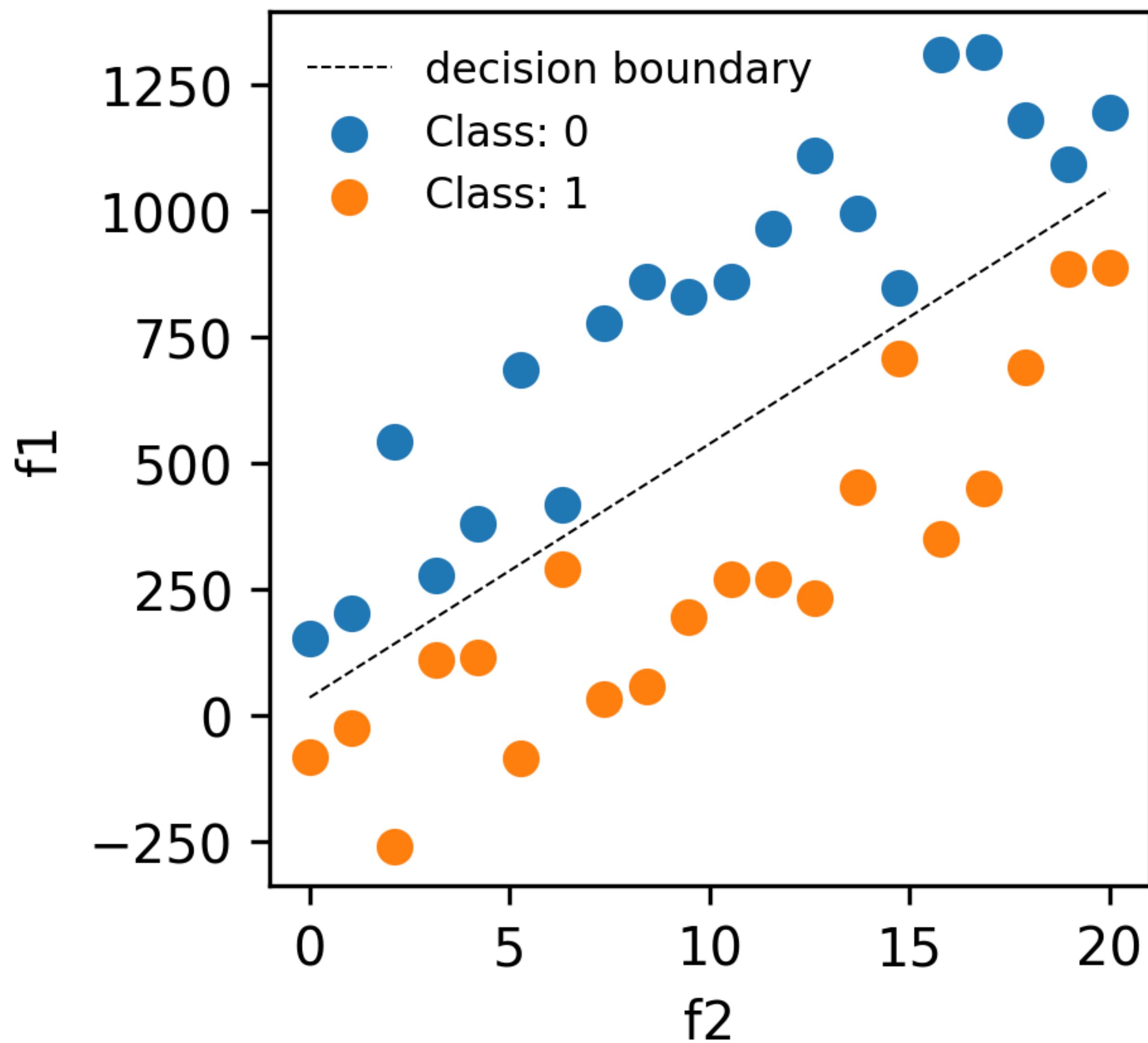
We can draw a line that separates
these two groups

Classes are linearly separable



This line is called a decision boundary

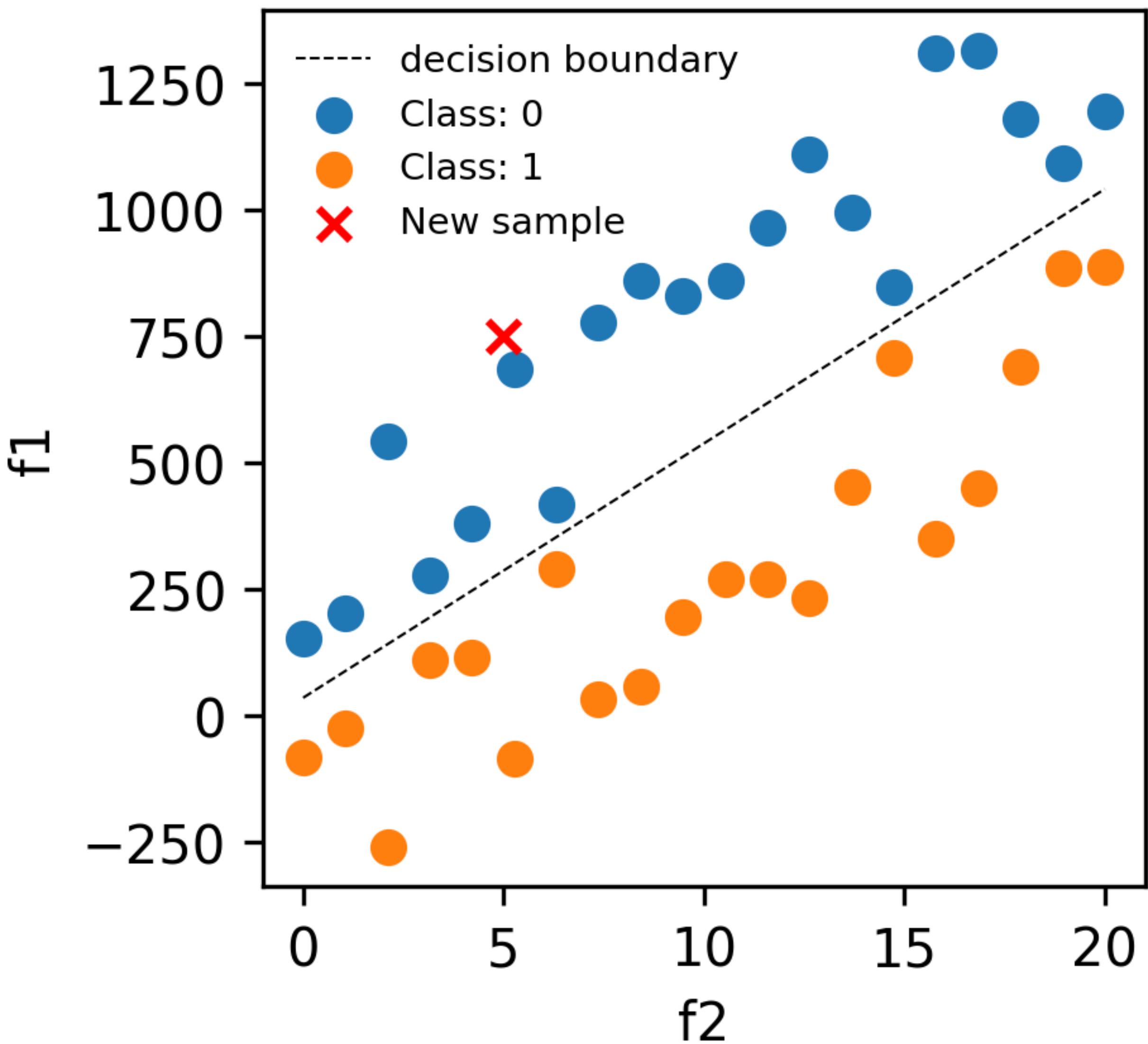
$$w_1 \cdot x_1 + w_2 \cdot x_2 + b = 0$$



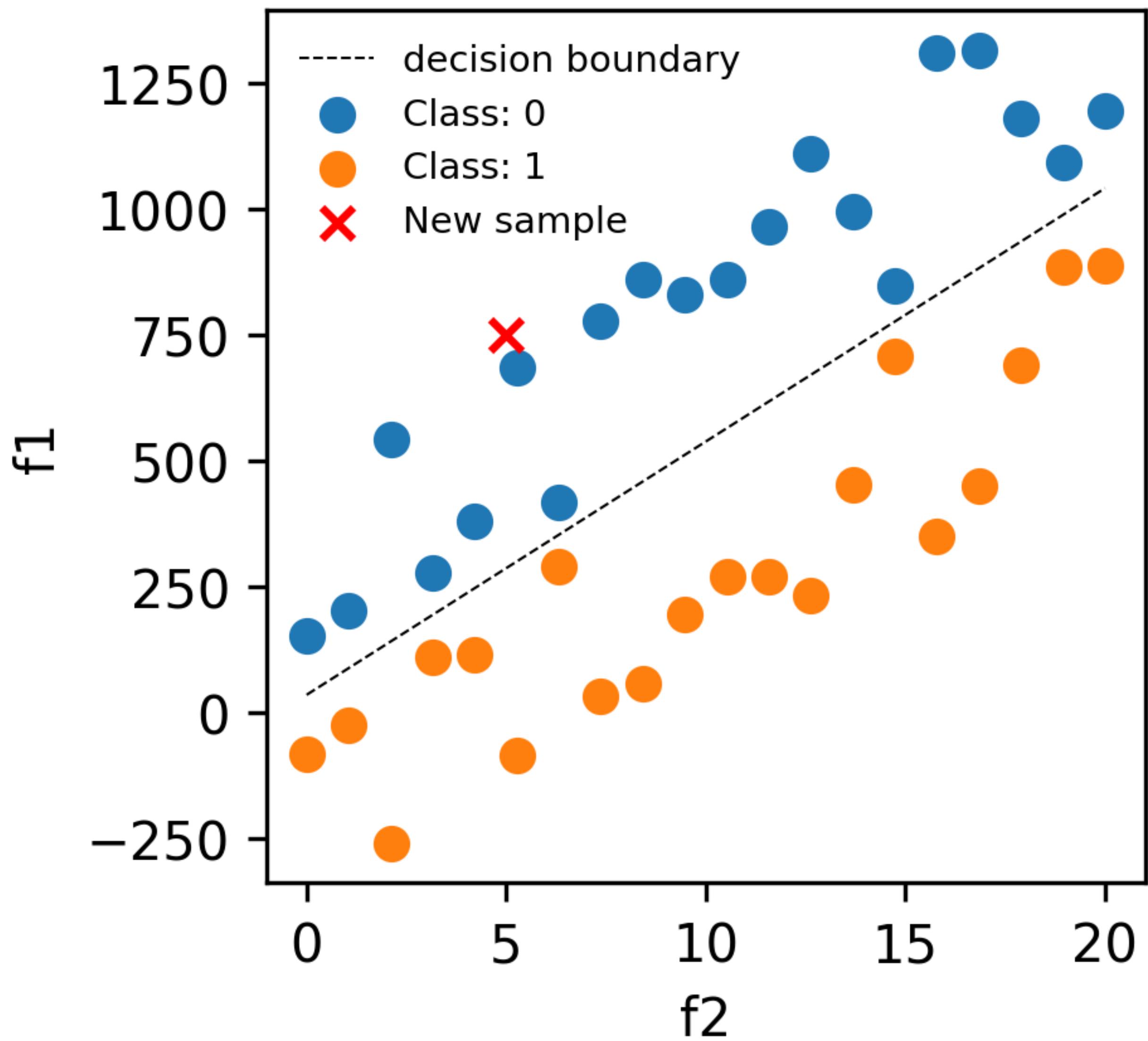
Let z be the score metric that helps answer the question "does this new sample belong to class 1?"

$$z = w_1 \cdot x_1 + w_2 \cdot x_2 + b$$

Note: $z = 0$ is a decision boundary



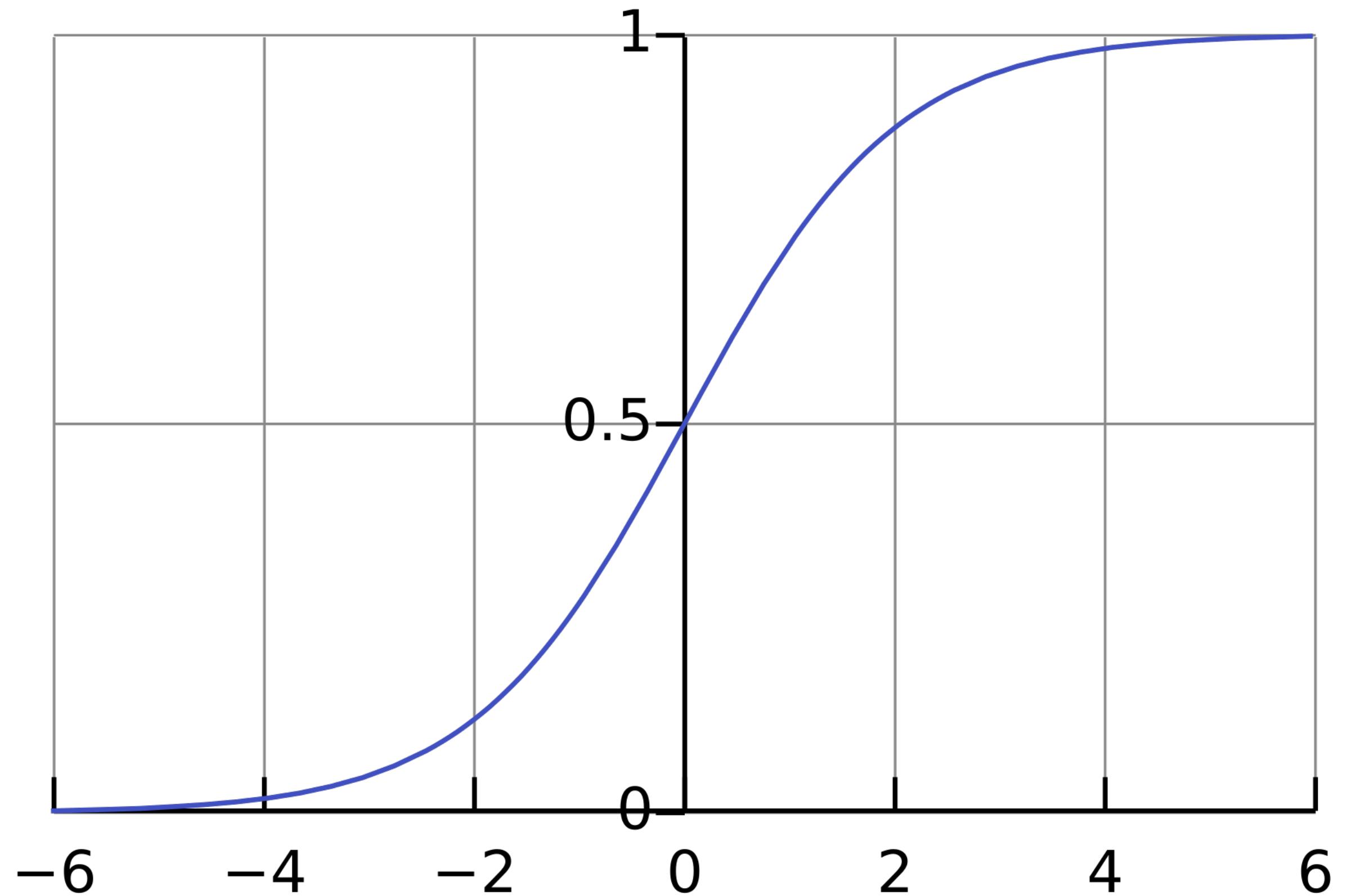
- In principle, z can take any value between $-\infty$ and ∞
- We can map this value to the range of $[0, 1]$ using some non-linear function
- And interpret the result as a probability



Logistic function

$$f(z) = \frac{1}{1 + \exp(-z)}, \text{ where } z = \mathbf{x}\mathbf{W} + b$$

- Tells us how likely the point is to be in Class 1
- i.e. the probability



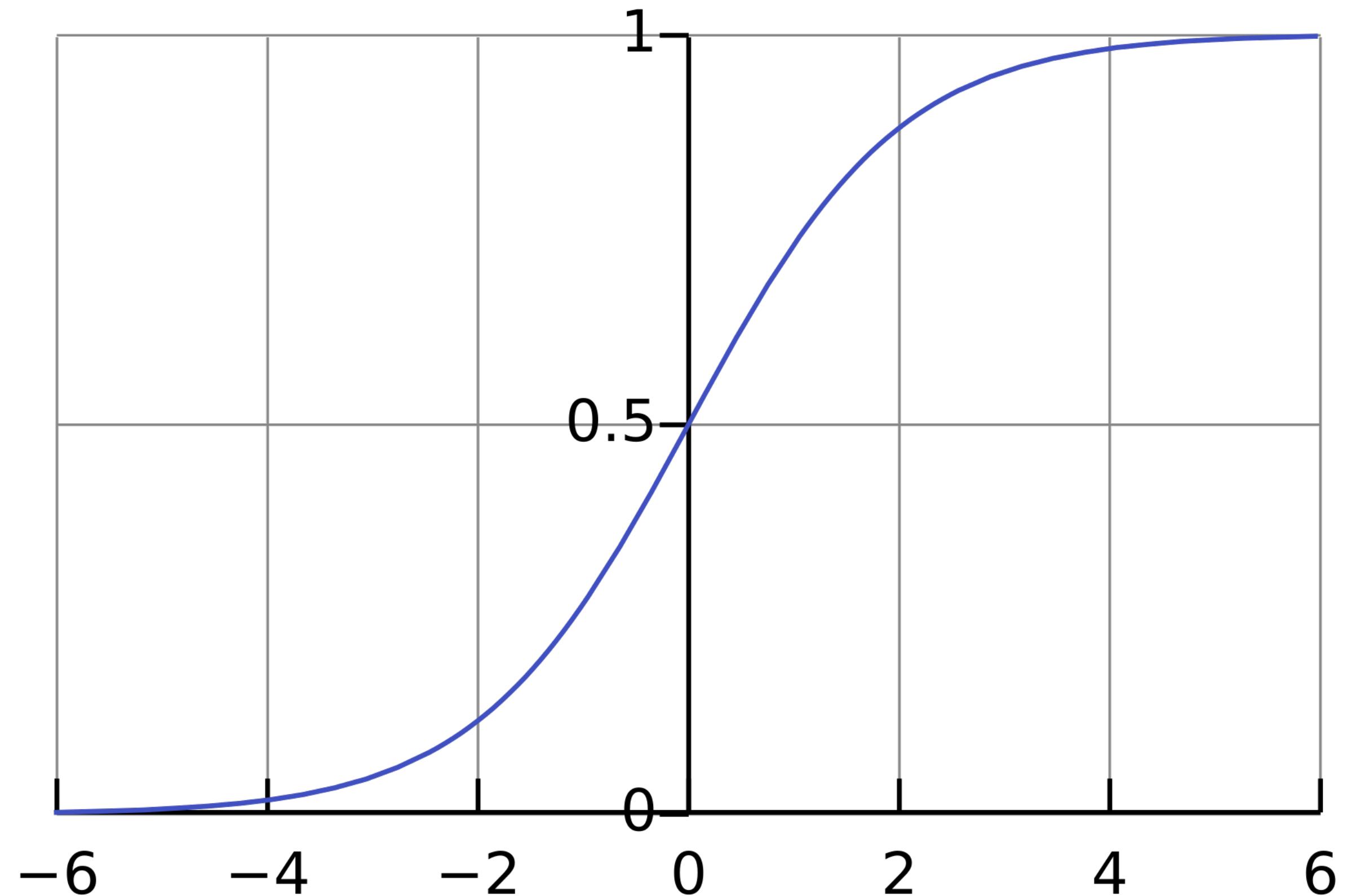
Logistic function

$$f(z) = \frac{1}{1 + \exp(-z)}, \text{ where } z = \mathbf{x}\mathbf{W} + b$$

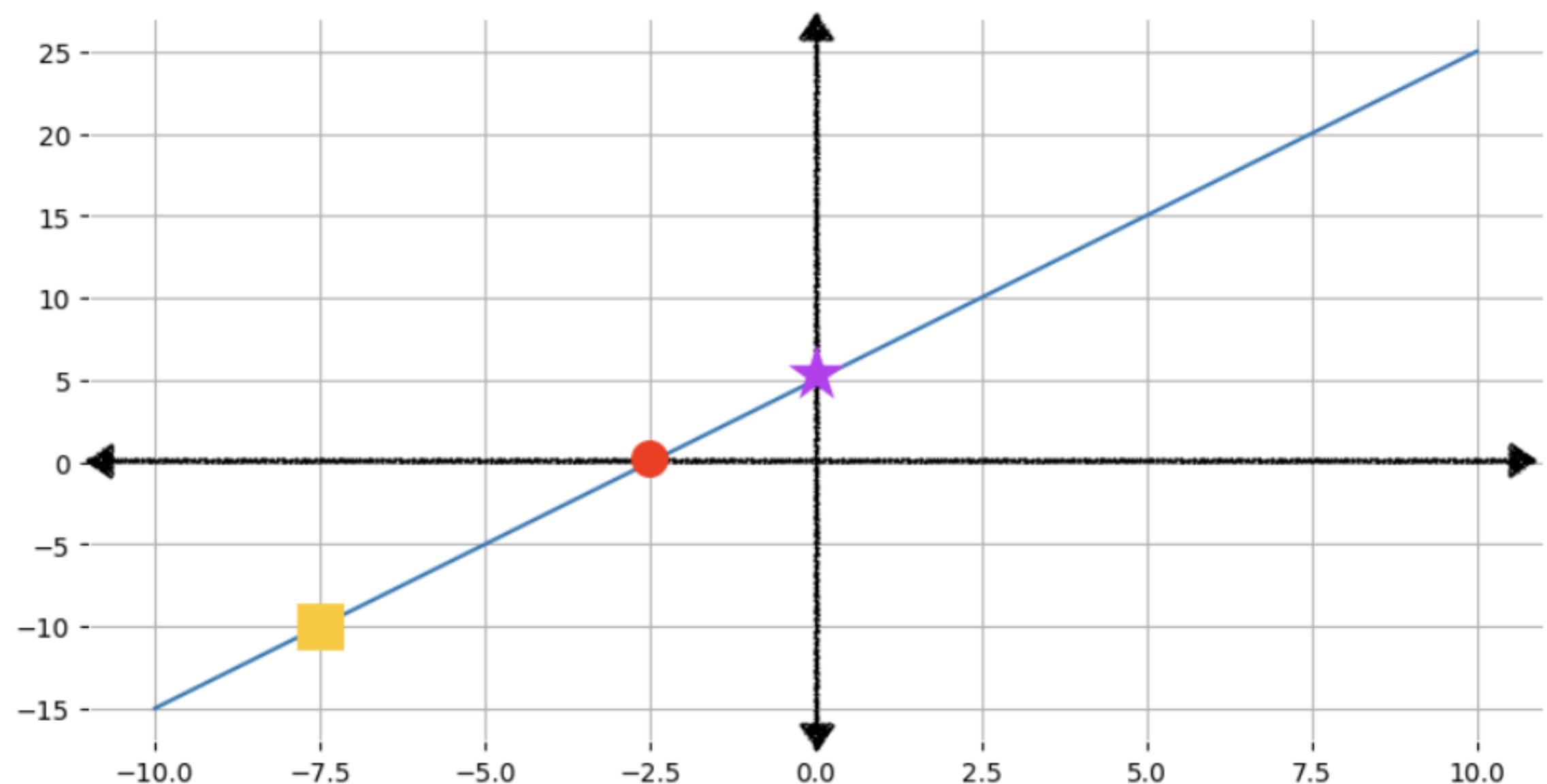
Post-processing step -> get class prediction

$$t(f(z)) = \begin{cases} 1 & \text{if } f(z) \geq 0.5 \\ 0 & \text{if } f(z) < 0.5 \end{cases}$$

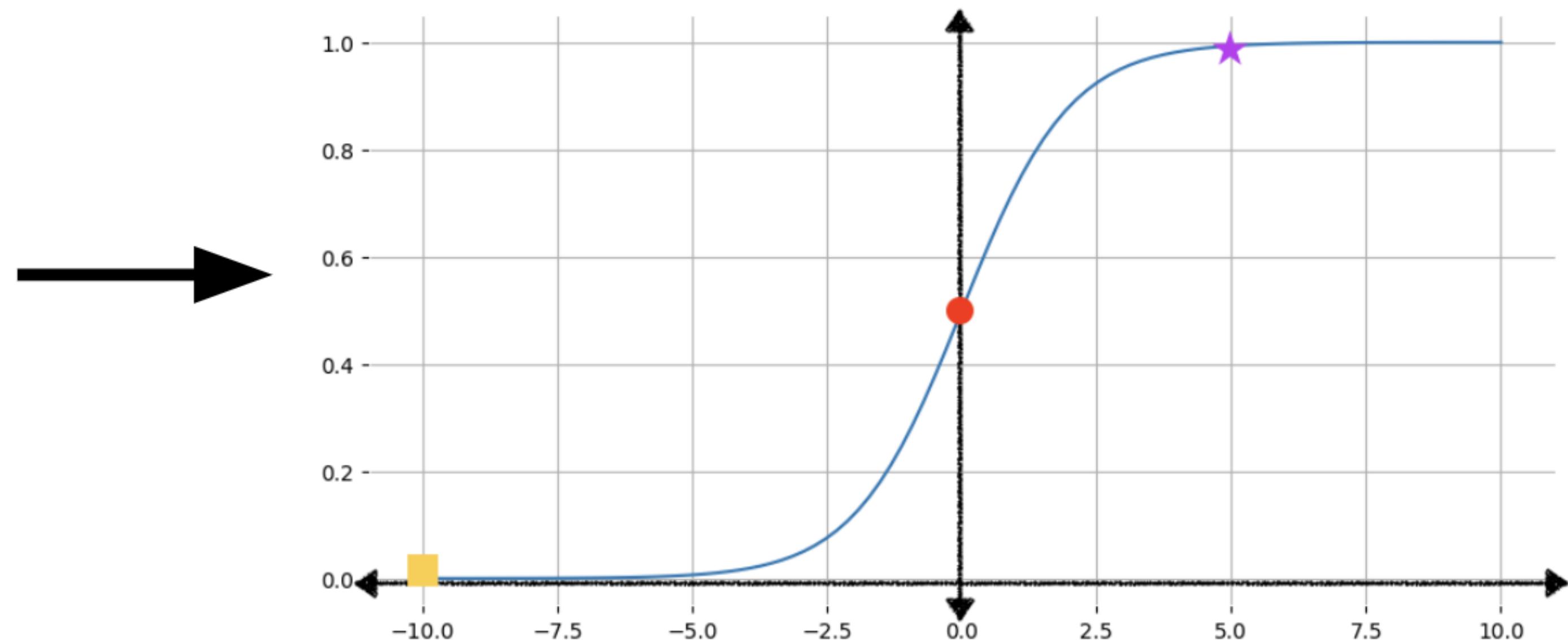
(You can use other than 0.5 threshold)



$$z = 2x + 5$$



$$y' = 1 / (1 + e^{-z})$$



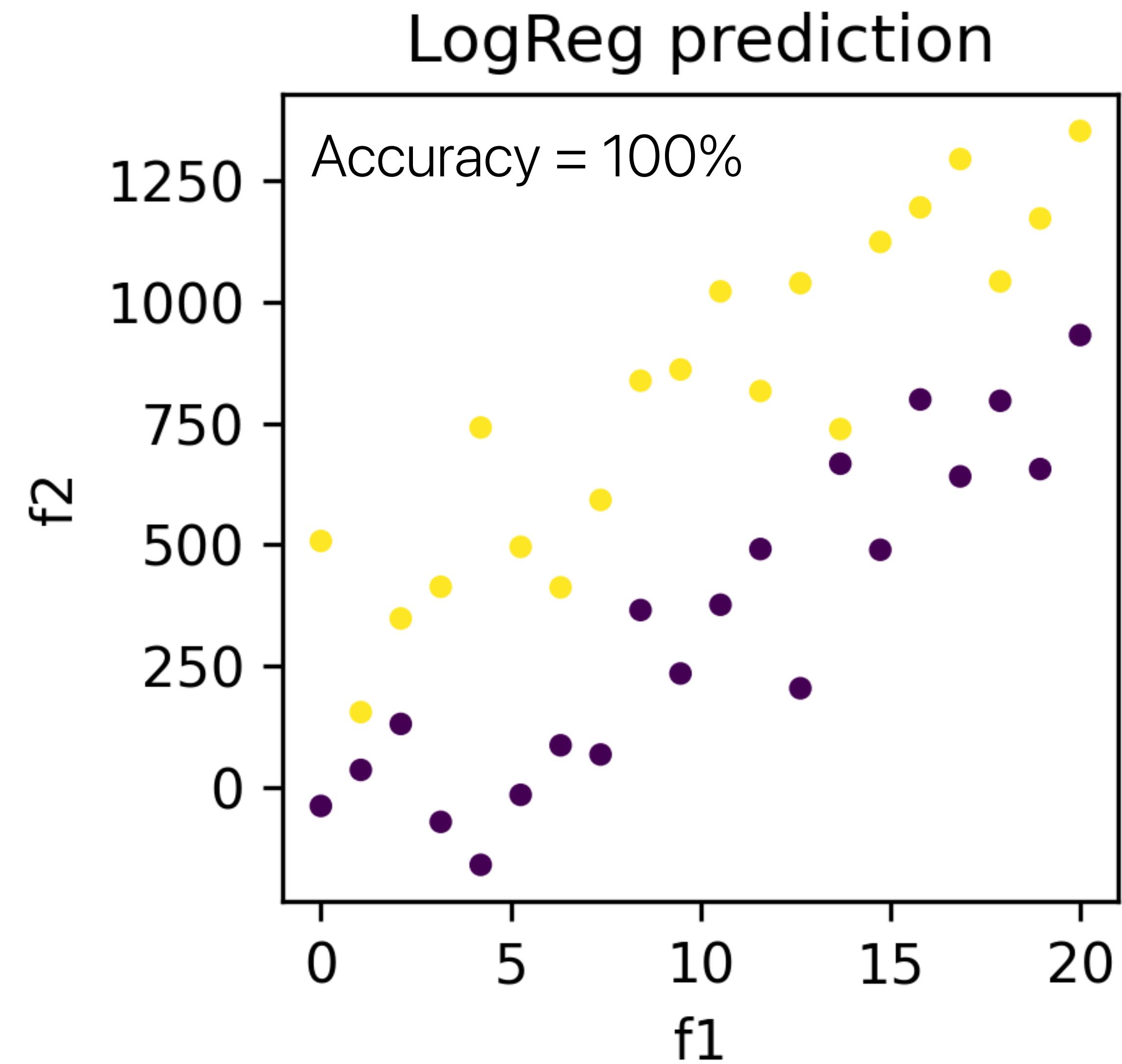
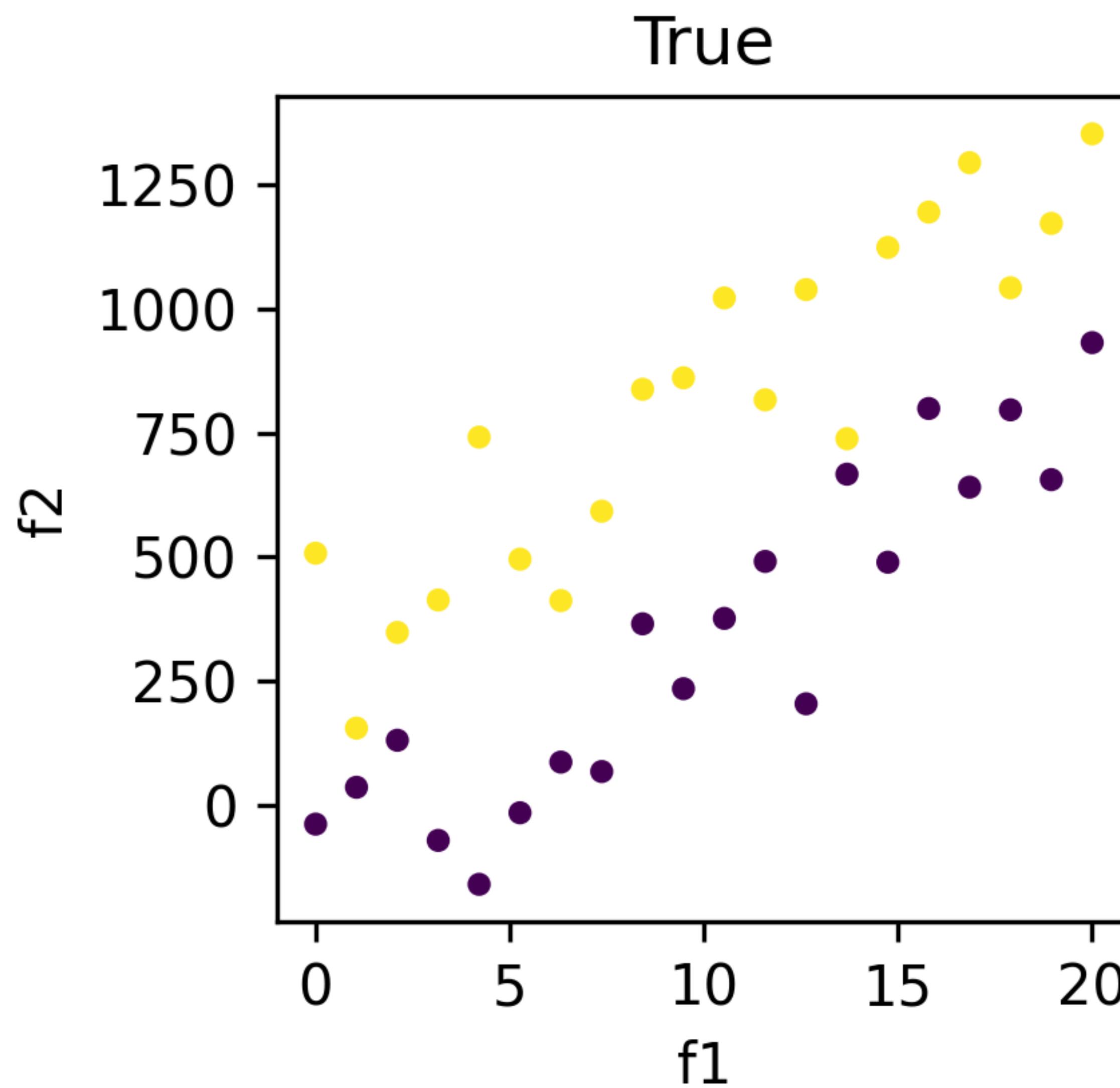
Loss: Cross-entropy loss (Log loss)

$$J(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

"...the logarithm of the magnitude of the change, rather than just the distance from data to prediction."

Metrics: Accuracy, recall, precision, and related metrics

$$\text{Accuracy} = \frac{\text{correct classifications}}{\text{total classifications}} = \frac{TP + TN}{TP + TN + FP + FN}$$



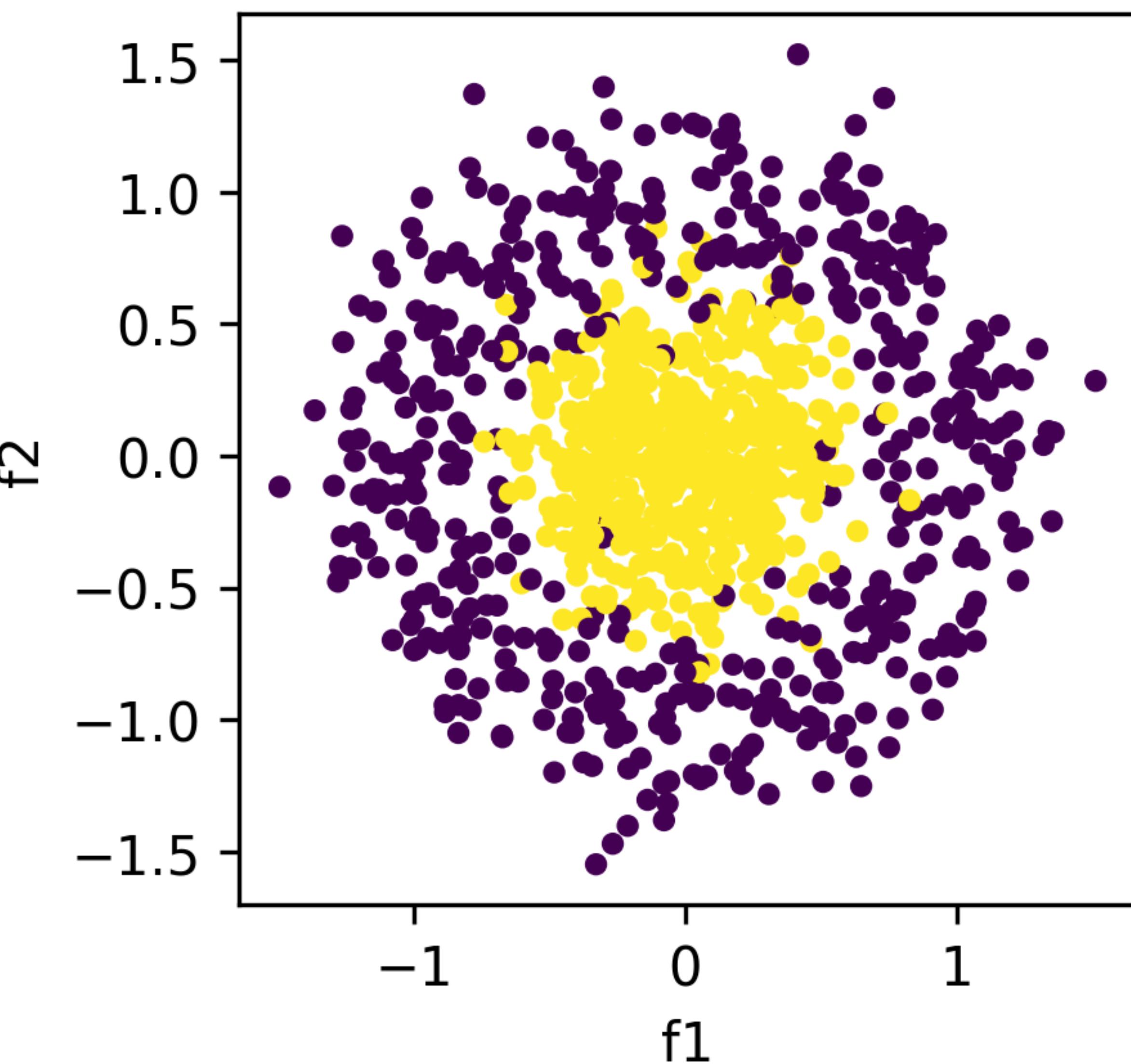
Part #2: Neural nets

Data samples:

$$\{(x_1, y_1), \dots, (x_i, y_i), \dots, (x_n, y_n)\}$$

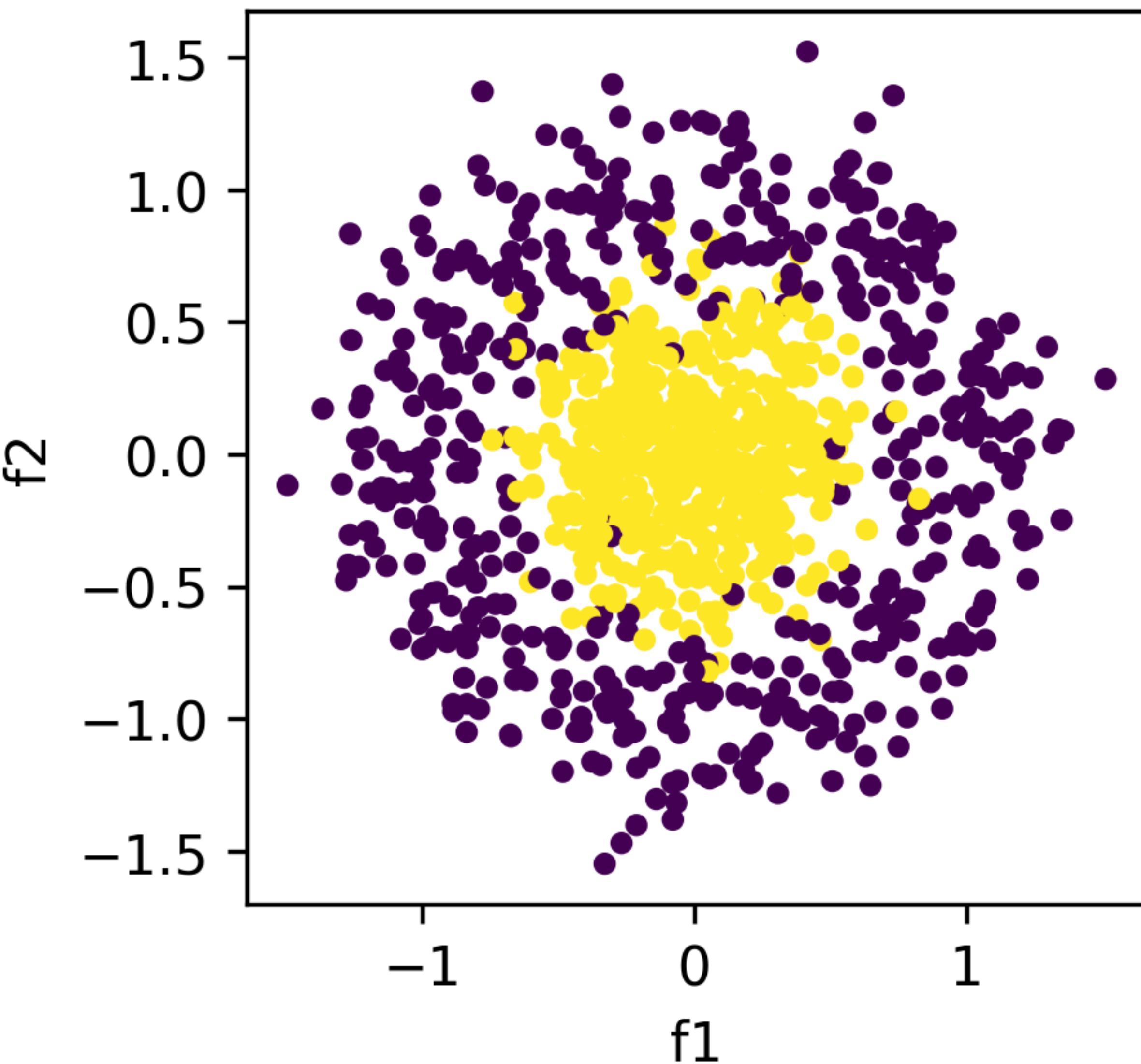
y_i - class of a sample (0 or 1)

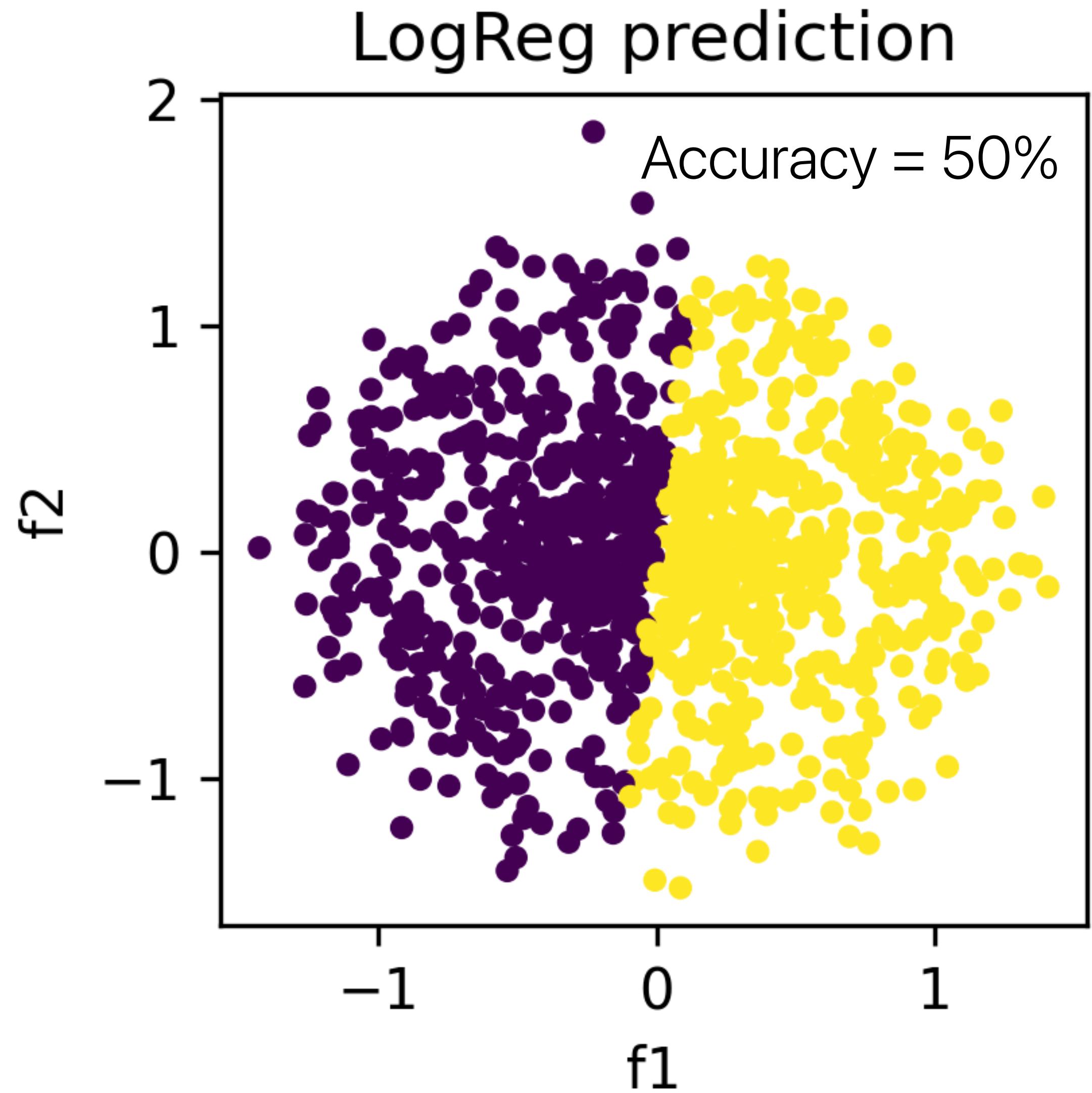
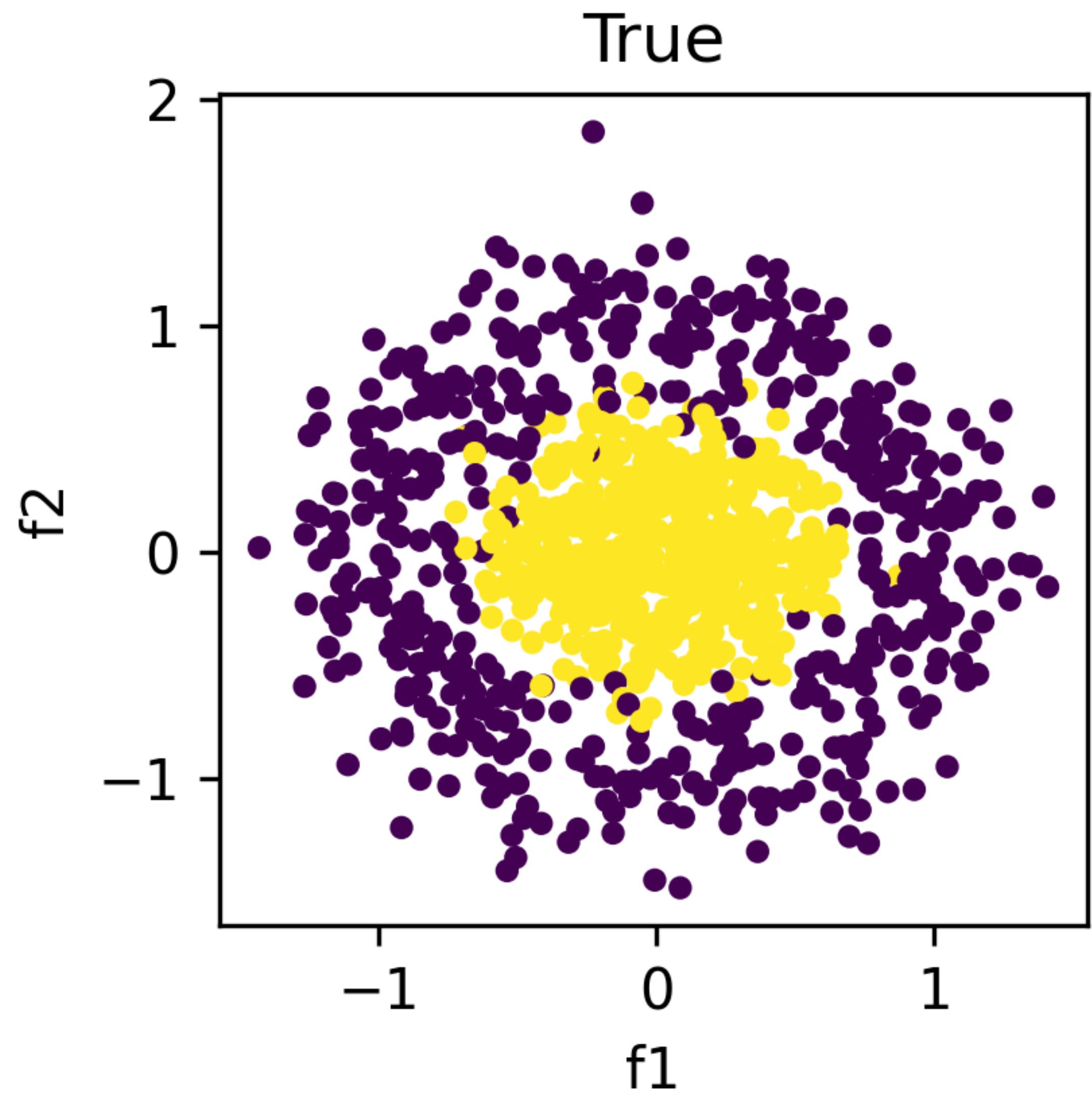
$x_i = (x_i^1, \dots, x_i^d)$ - $1 \times d$ vector of
input features



Can we draw a line to separate the classes?

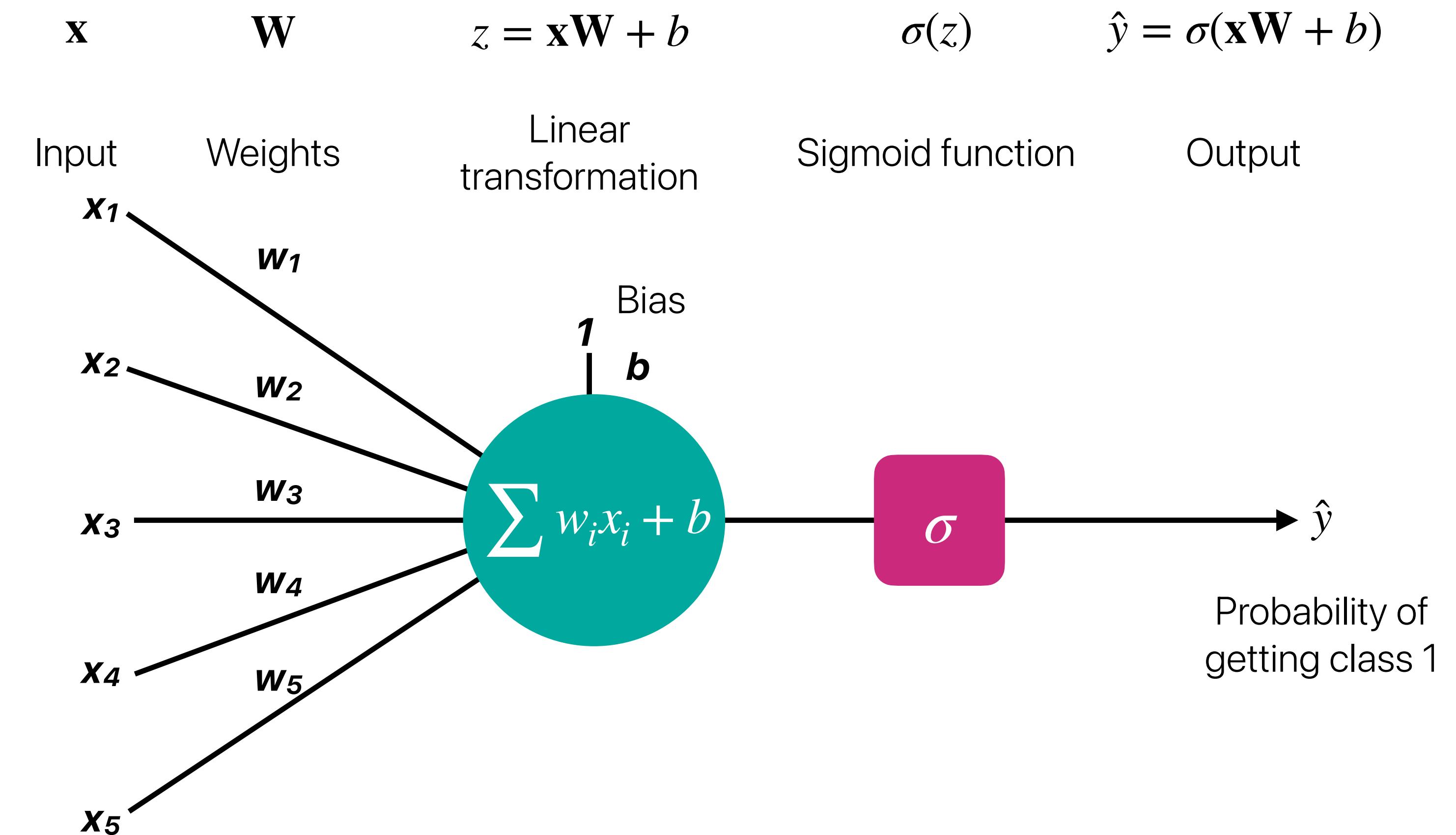
Classes are linearly non-separable



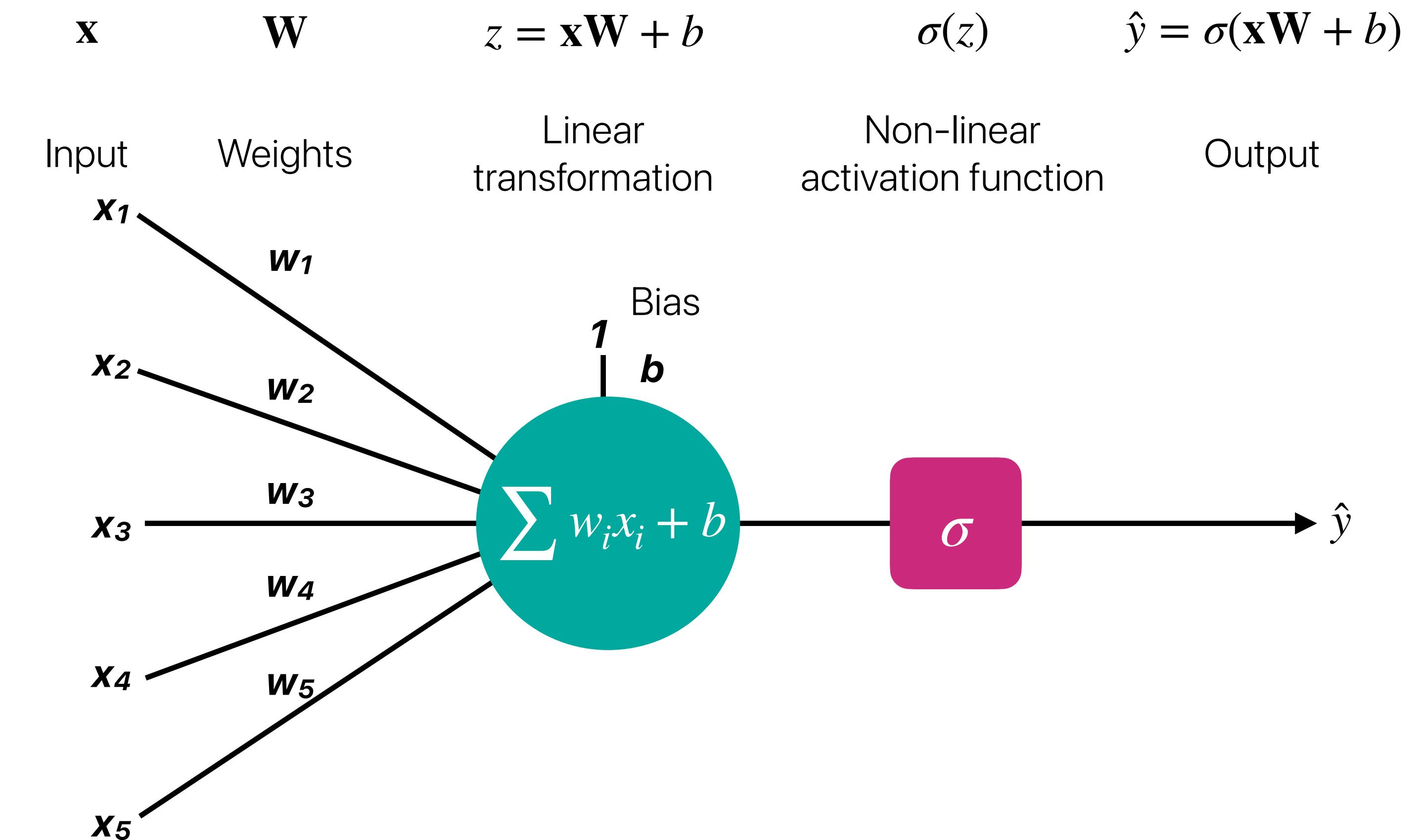


That's why we need neural nets

Let's represent a LogReg in a
following way



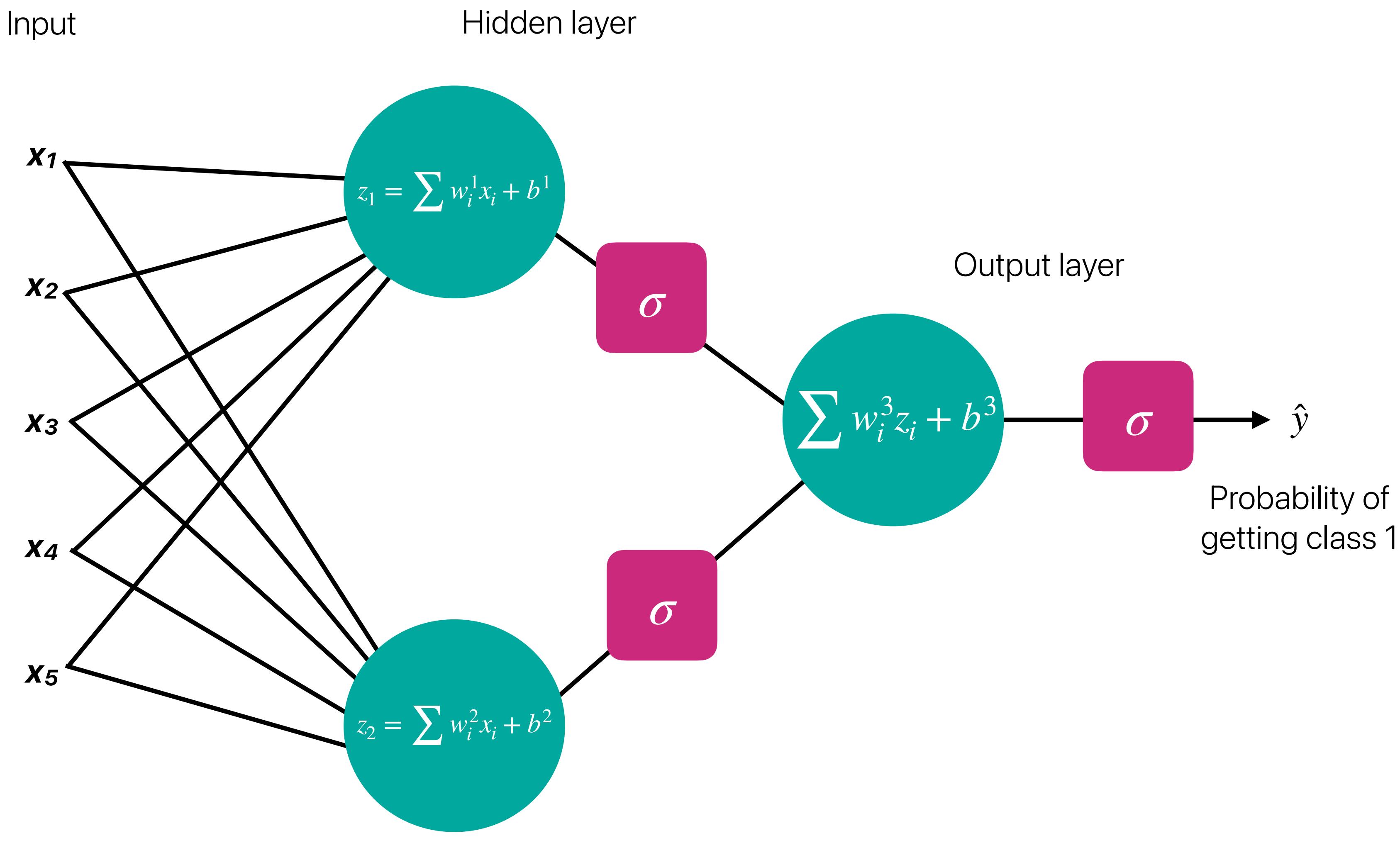
This unit is called a perceptron



One hidden layer

multilayer perceptron

- Input layer
 - Features
- Hidden layer
 - Weighted sum
 - Activation
- Output layer
 - Weighted sum
 - Sigmoid
- Post-processing
 - Convert probabilities to labels



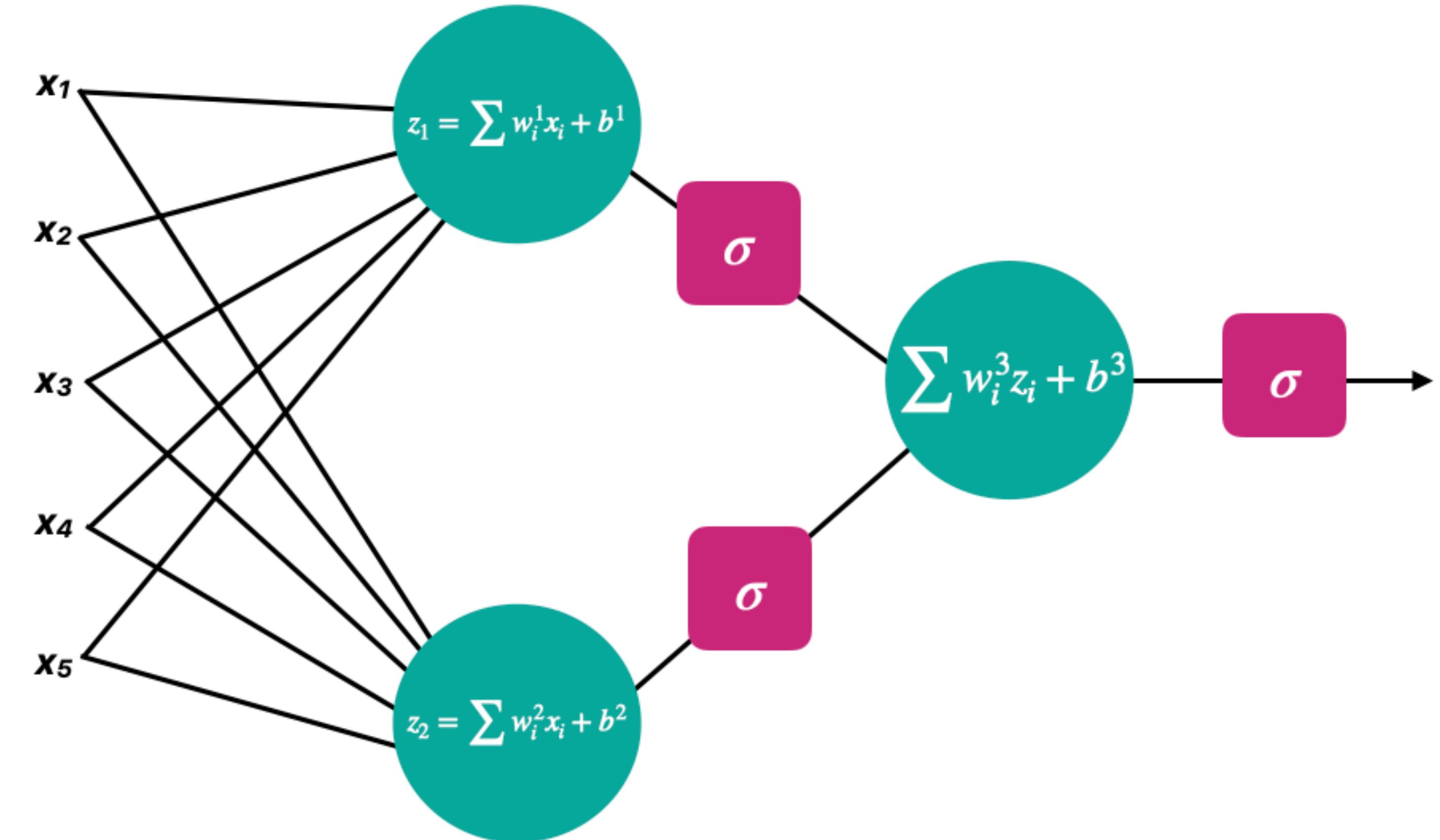
Matrix form

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1k} \\ w_{21} & w_{22} & \cdots & w_{2k} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

$$\mathbf{z} = W\vec{x} + \vec{b}$$

$$h = \sigma(W\vec{x} + \vec{b})$$

*“The output of each layer is the product of its weight matrix and the input vector plus its bias vector, all wrapped in a non-linear activation function.”



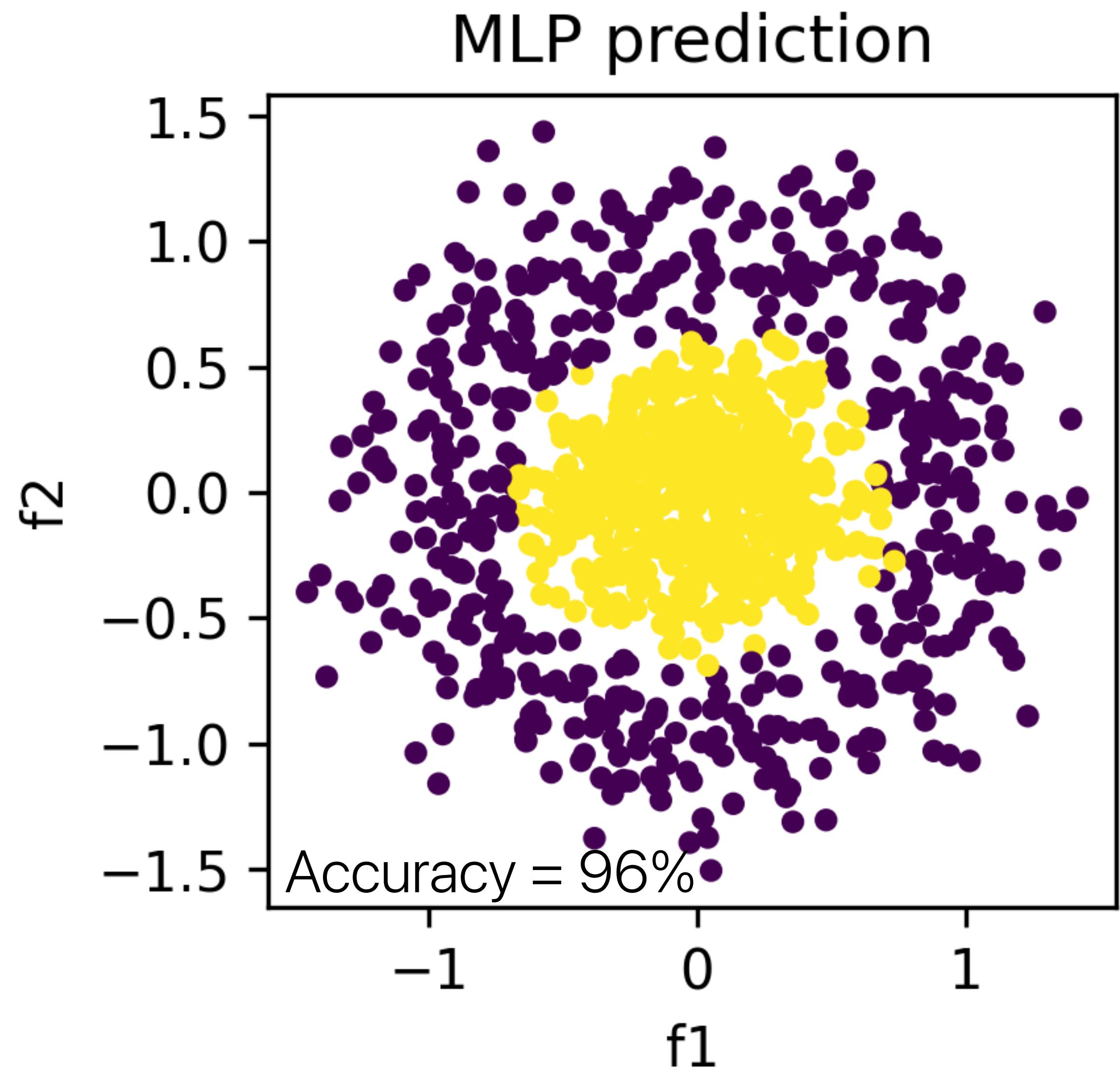
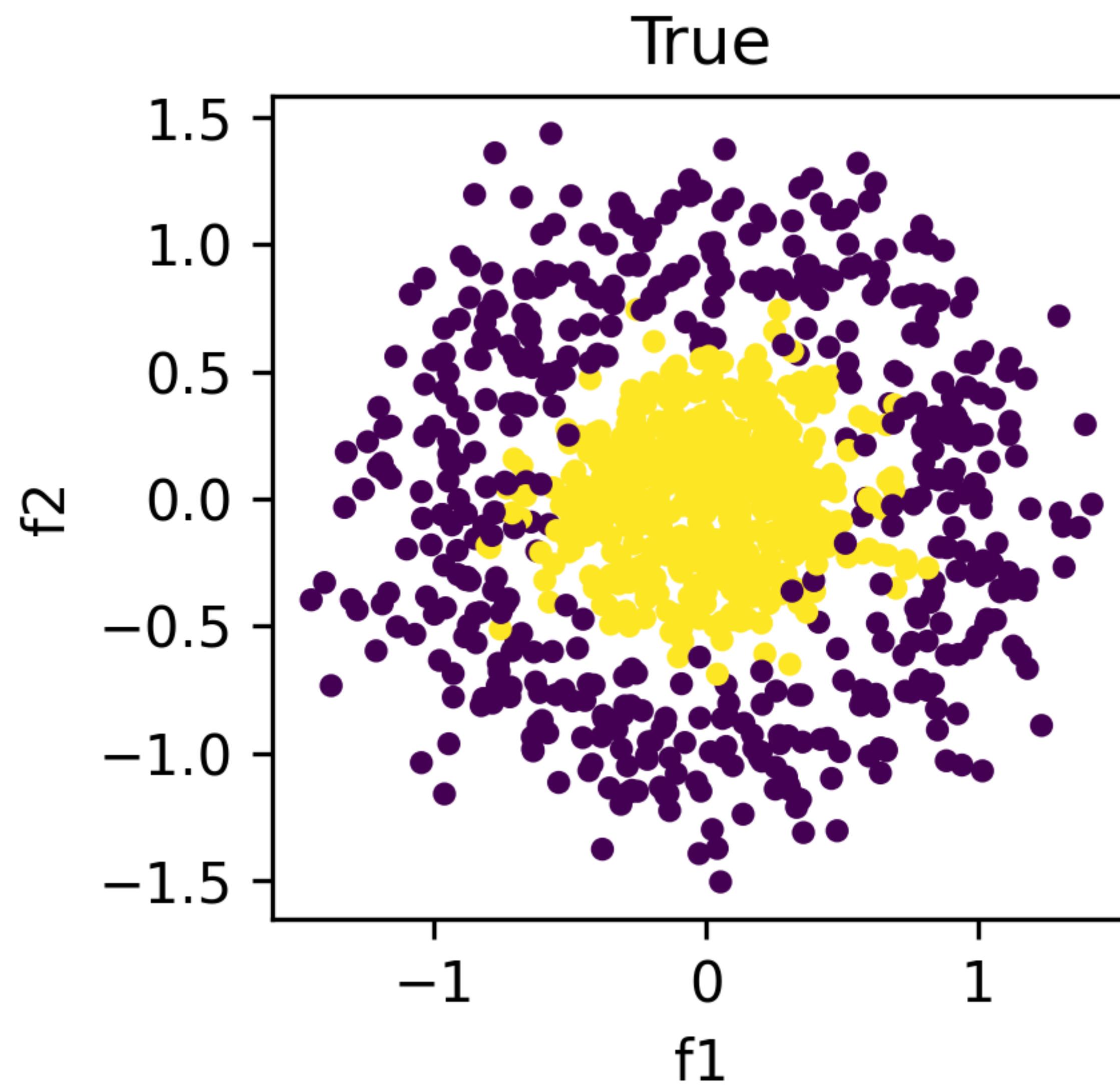
Non-linear coupling between the layers is the key

...that enables neural networks to be universal function approximators and learn complex, non-linear patterns.



Neural Network Activation Functions: a small subset!

ReLU $\max(0, x)$	GELU $\frac{x}{2} \left(1 + \tanh \left(\sqrt{\frac{2}{\pi}} (x + ax^3) \right) \right)$	PReLU $\max(0, x)$
ELU $\begin{cases} x & \text{if } x > 0 \\ \alpha(x \exp x - 1) & \text{if } x < 0 \end{cases}$	Swish $\frac{x}{1 + \exp -x}$	SELU $\alpha(\max(0, x) + \min(0, \beta(\exp x - 1)))$
SoftPlus $\frac{1}{\beta} \log(1 + \exp(\beta x))$	Mish $x \tanh \left(\frac{1}{\beta} \log(1 + \exp(\beta x)) \right)$	RReLU $\begin{cases} x & \text{if } x \geq 0 \\ ax & \text{if } x < 0 \text{ with } a \sim \mathcal{R}(l, u) \end{cases}$
HardSwish $\begin{cases} 0 & \text{if } x \leq -3 \\ x & \text{if } x \geq 3 \\ x(x+3)/6 & \text{otherwise} \end{cases}$	Sigmoid $\frac{1}{1 + \exp(-x)}$	SoftSign $\frac{x}{1 + x }$
Tanh $\tanh(x)$	Hard tanh $\begin{cases} a & \text{if } x \geq a \\ b & \text{if } x \leq b \\ x & \text{otherwise} \end{cases}$	Hard Sigmoid $\begin{cases} 0 & \text{if } x \leq -3 \\ 1 & \text{if } x > 3 \\ x/6 + 1/2 & \text{otherwise} \end{cases}$
Tanh Shrink $x - \tanh(x)$	Soft Shrink $\begin{cases} x - \lambda & \text{if } x > \lambda \\ x + \lambda & \text{if } x < -\lambda \\ 0 & \text{otherwise} \end{cases}$	Hard Shrink $\begin{cases} x & \text{if } x > \lambda \\ x & \text{if } x < -\lambda \\ 0 & \text{otherwise} \end{cases}$



[Universal Approximation Theorem](#) states that certain models, most notably feedforward neural networks with a single hidden layer, can approximate any continuous function to an arbitrary degree of accuracy, given enough neurons in the hidden layer.

How to train NNs: Backpropagation and chain rule

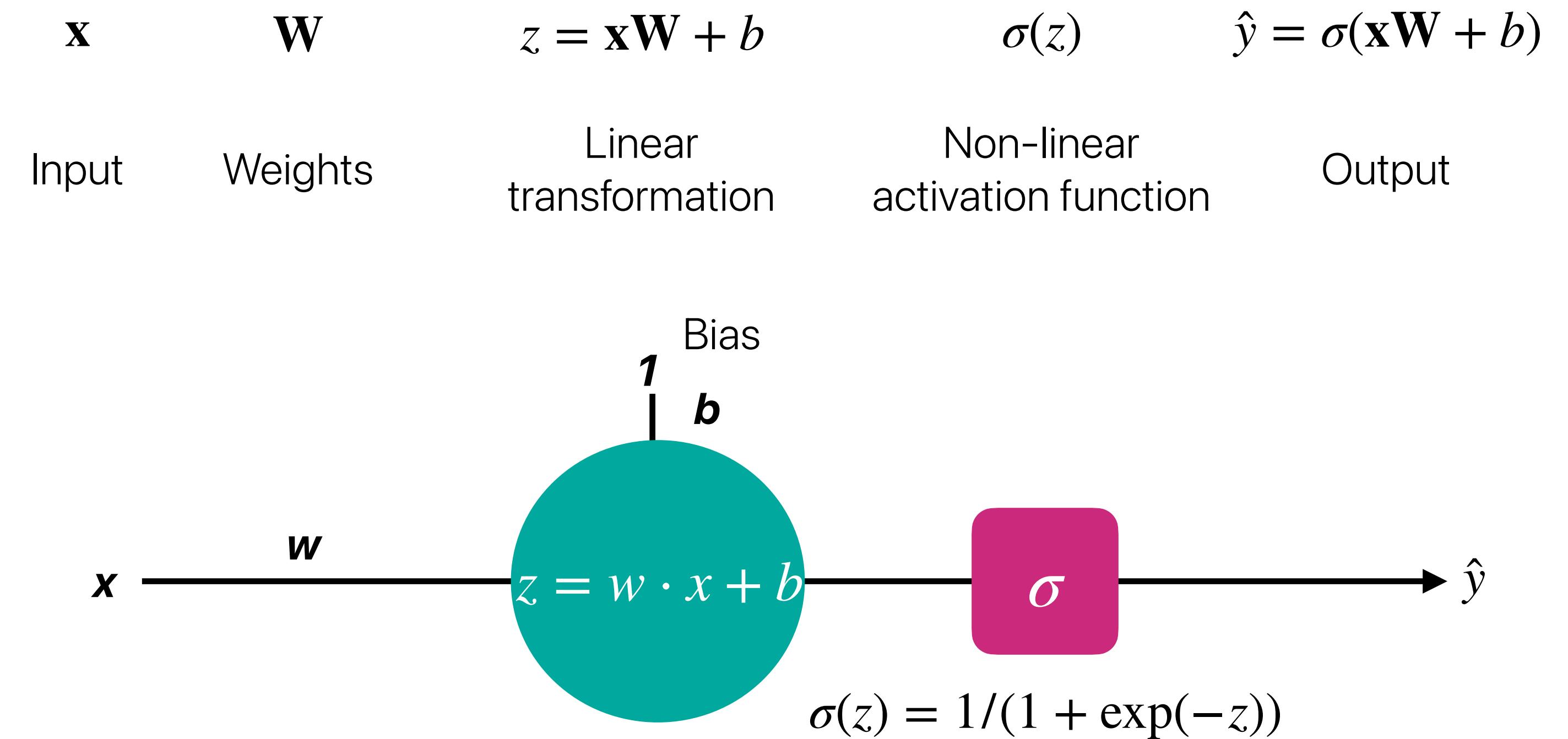
Consider a simple example with a scalar feature

We want to find the best parameters of our model
using gradient descent, i.e. update weights as follows:

$$w = w - \alpha \cdot \frac{\partial E}{\partial w}$$

$$b = b - \alpha \cdot \frac{\partial E}{\partial b}$$

Where $E = \frac{1}{2}(\hat{y} - y)^2$ is the MSE loss function



Chain rule reminder:

$$y = f(g(x))$$

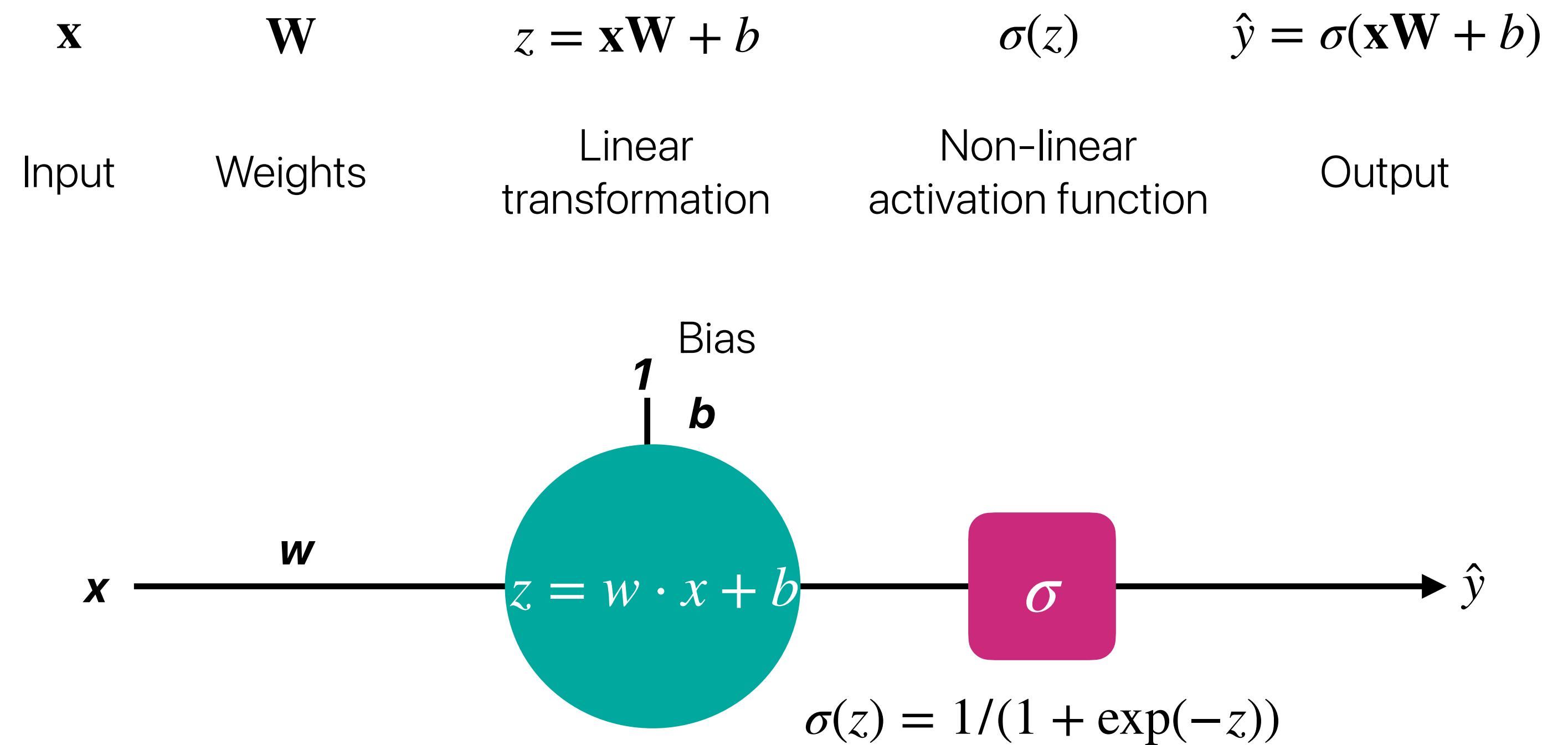
$$\frac{dy}{dx} = \frac{df}{dg} \cdot \frac{dg}{dx}$$

According to the chain rule

$$E = \frac{1}{2}(\hat{y}(z(w)) - y)^2$$

$$\frac{\partial E}{\partial w} = \frac{\partial E(\hat{y}(z(w)))}{\partial w}$$

$$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w}$$



$$\sigma(z) = 1/(1 + \exp(-z))$$

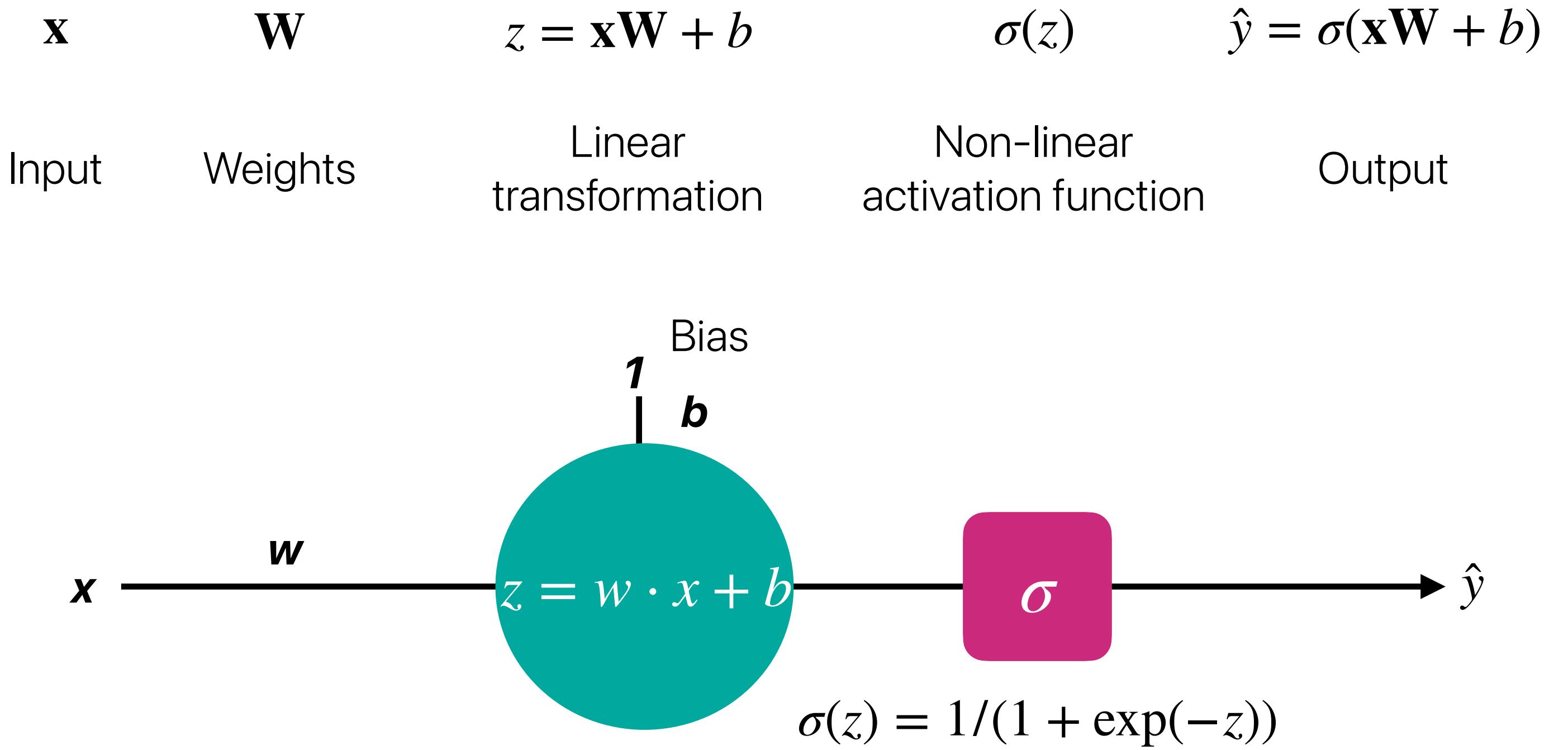
We are propagating gradients back from the last
(output) layer to the first hidden layer

$$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$\frac{\partial E}{\partial y} = \hat{y} - y$$

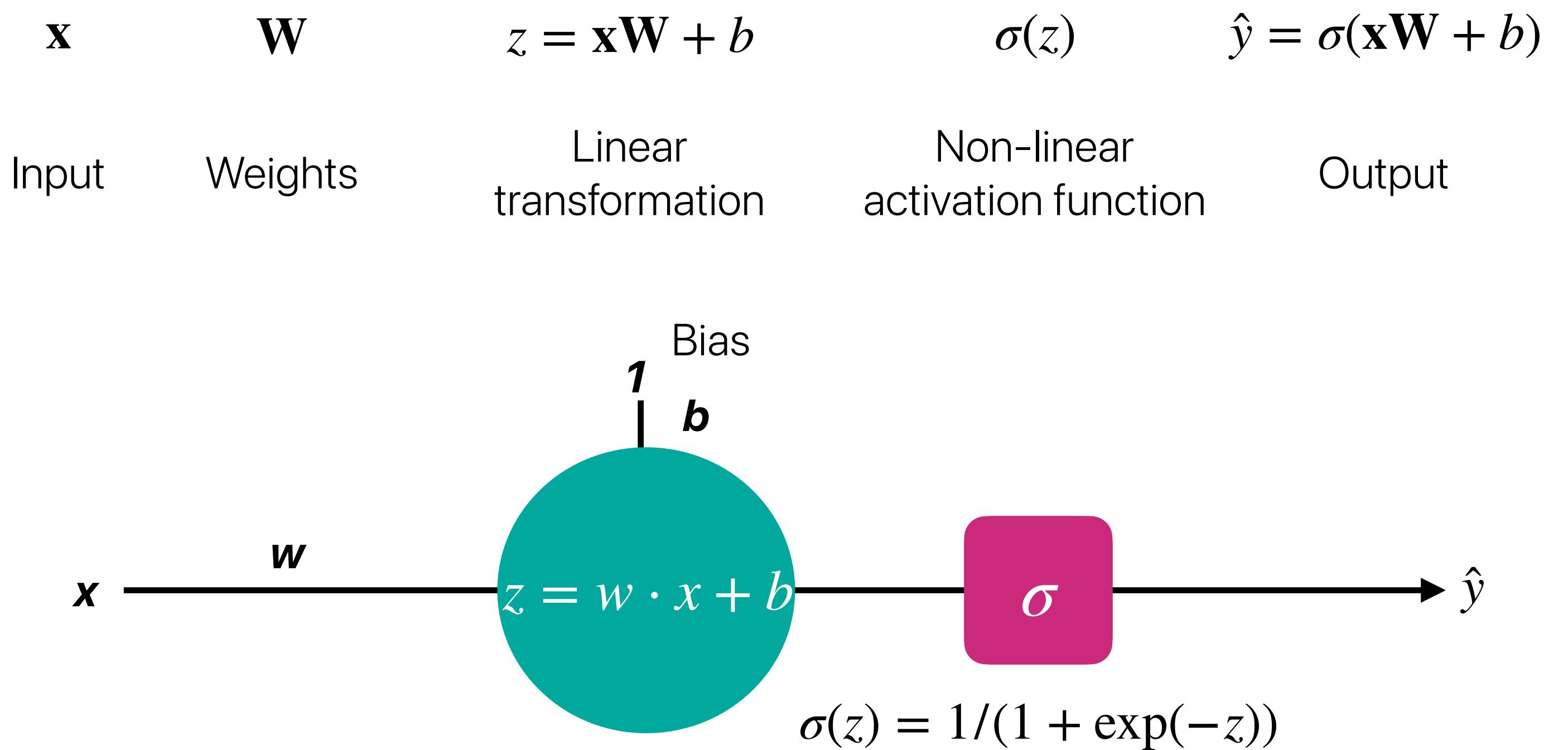
$$\frac{\partial y}{\partial z} = \frac{\partial(\frac{1}{1 + \exp(-z)})}{\partial z} = -\frac{\exp(-z)}{1 + \exp(-z)} = \sigma(z) \cdot (1 - \sigma(z))$$

$$\frac{\partial z}{\partial w} = x$$



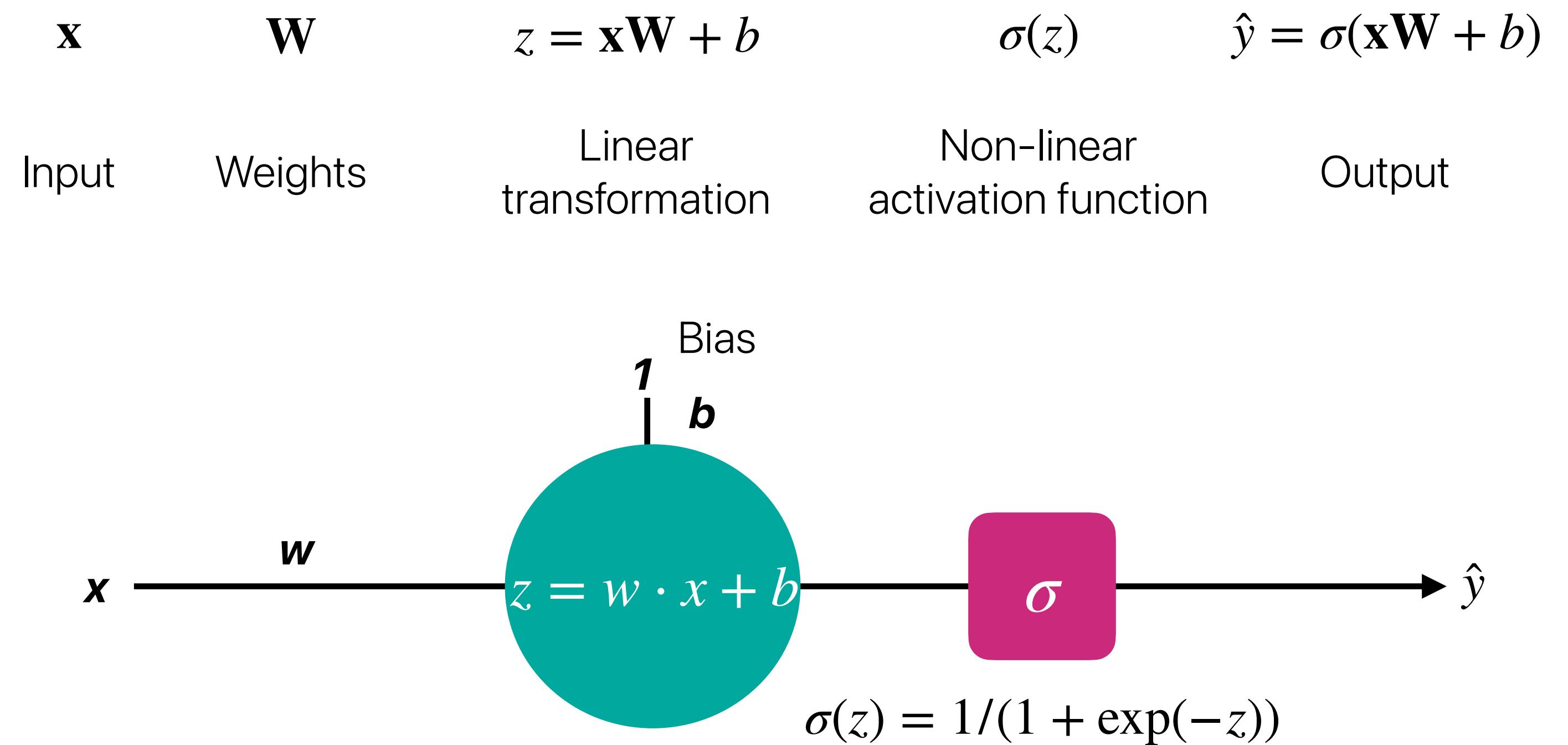
$$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$\frac{\partial E}{\partial w} = (\hat{y} - y) \cdot \sigma(z)(1 - \sigma(z)) \cdot x$$



Similarly

$$\frac{\partial E}{\partial b} = (\hat{y} - y) \cdot \sigma(z)(1 - \sigma(z))$$



Let's rewrite MLP as follows:

$$\text{MLP}_N(\mathbf{x}) = \sigma_N \left(\mathbf{W}_N \cdots \sigma_2 \left(\mathbf{W}_2 \cdot \sigma_1 \left(\mathbf{W}_1 \cdot \mathbf{x} + \mathbf{b}_1 \right) + \mathbf{b}_2 \right) \cdots + \mathbf{b}_N \right)$$

- just a fancy (composite) function, isn't?

Or even simpler:

$$\text{MLP}_N(\mathbf{x}) = L_N(L_{N-1}(\cdots L_2(L_1(\mathbf{X}, \mathbf{W}_1^*), \mathbf{W}_2^*) \cdots), \mathbf{W}_n^*)$$

(* symbol specifies that W matrix includes bias term)

$$\text{MLP}_N(\mathbf{x}) = L_N(L_{N-1}(\dots L_2(L_1(\mathbf{X}, \mathbf{W}_1^*), \mathbf{W}_2^*)\dots), \mathbf{W}_n^*)$$

\mathbf{x} - input features, \mathbf{y} - true labels

Let:

$$\mathbf{h}_0 = \mathbf{x}$$

$$\mathbf{h}_k = \sigma_k \left(\mathbf{W}_k \cdot \mathbf{h}_{k-1} + \mathbf{b}_k \right) \quad \text{for } k = 1, 2, \dots, L$$

$$\hat{\mathbf{y}} = \mathbf{h}_L$$

Feed forward:

$$\mathbf{h}_0 \rightarrow \mathbf{h}_1 \rightarrow \dots \rightarrow \mathbf{h}_l$$

Backward propagation:

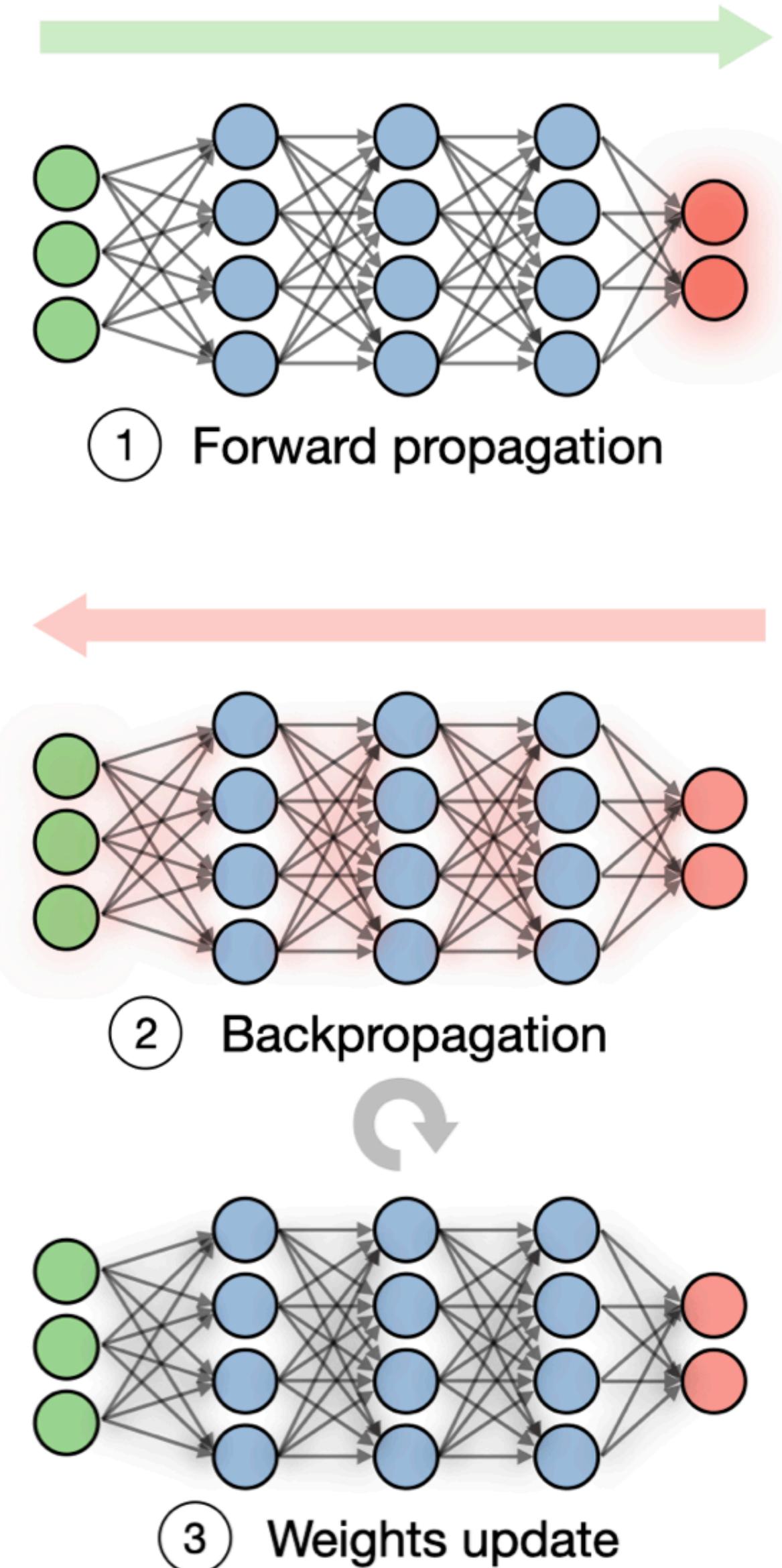
$$\frac{\partial E}{\partial \mathbf{W}_k} = \frac{\partial E}{\partial \mathbf{h}_N} \cdot \frac{\partial \mathbf{h}_N}{\partial \mathbf{h}_{N-1}} \dots \cdot \frac{\partial \mathbf{h}_{k+1}}{\partial \mathbf{h}_k} \cdot \frac{\partial \mathbf{h}_k}{\partial \mathbf{W}_k}$$

For a more detailed explanation consider:

- [Backpropagation](#) lecture by Prof. Roger Grosse
- [Deep learning handout](#) lecture by Prof. Evgeny Burnaev
- [Neural networks handout](#) lecture by Prof. Evgeny Burnaev
- <https://en.wikipedia.org/wiki/Backpropagation>
- <https://www.qwertee.io/blog/an-introduction-to-backpropagation/>

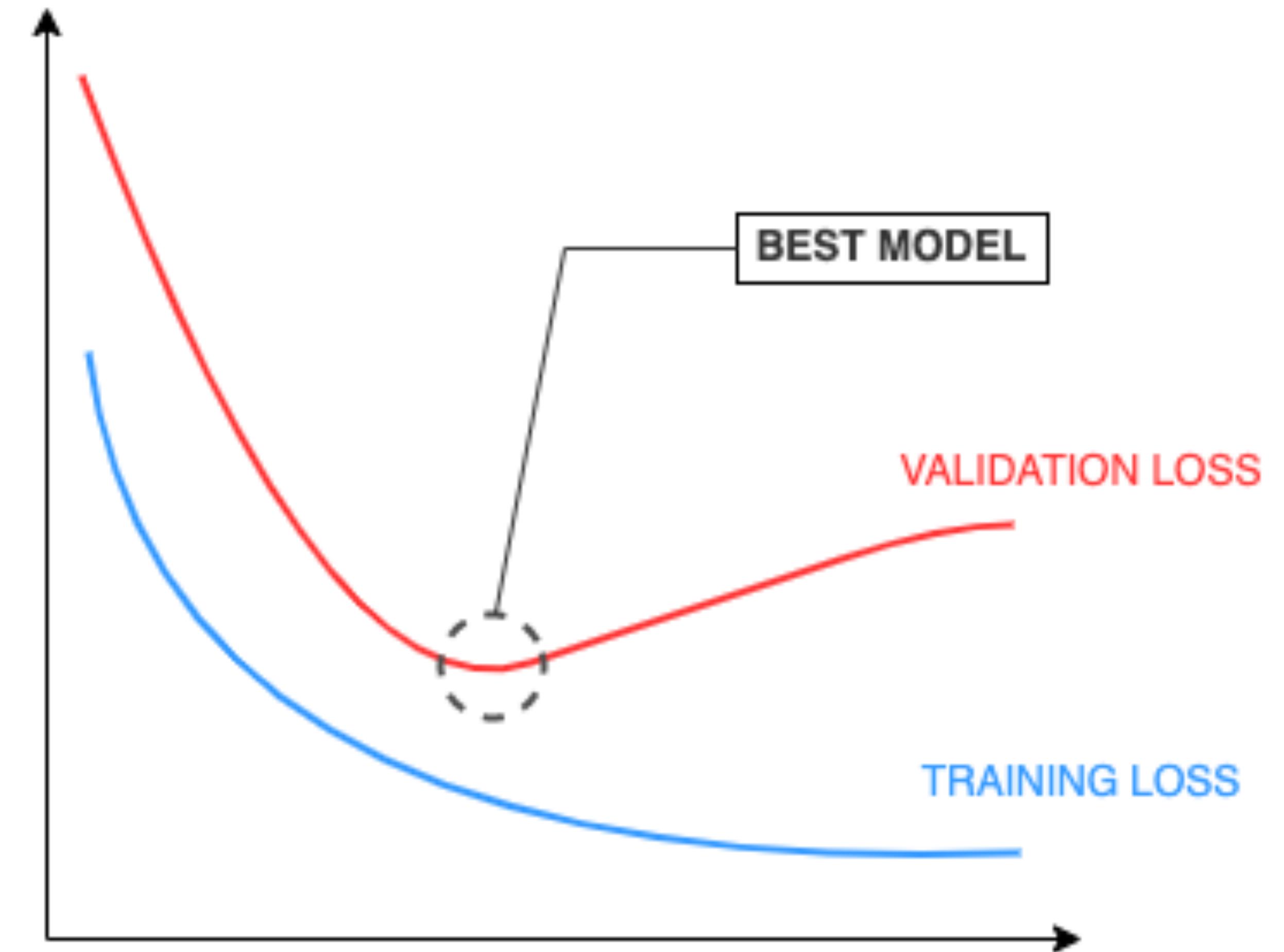
(Approximate) Training loop

1. Randomly initialize the model's parameters
2. For a specified number of epochs:
 - A. Shuffle the training dataset $\{\mathbf{x}_i, y_i\}$
 - B. Split training dataset into non-overlapping mini-batches
 - C. For each mini-batch in the dataset
 - I. Forward pass: Compute predictions for the batch $\hat{y} = \text{MLP}(x)$
 - II. Compute the average loss over the batch: $\text{loss} = E(\hat{y}, y)$
 - III. Compute gradients via backpropagation
 - IV. Update the weights using the gradients: $w_{ij}^k = w_{ij}^k - \alpha \cdot \frac{\partial E}{\partial w_{ij}}$



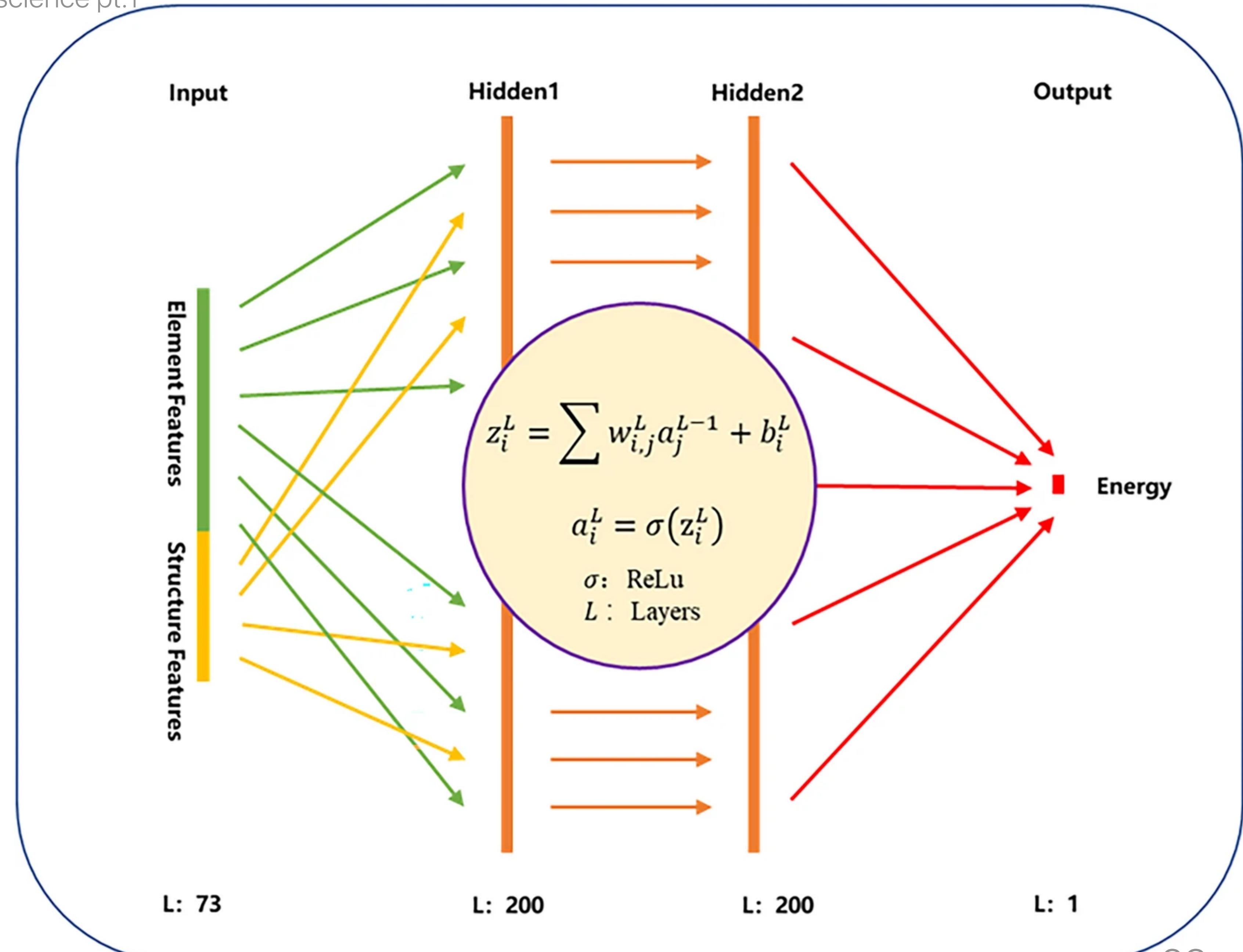
When to stop training procedure?

- Model parameters are updated using training data
- Validation data monitors loss on unseen data to determine when to stop
- Best model is then evaluated on test data



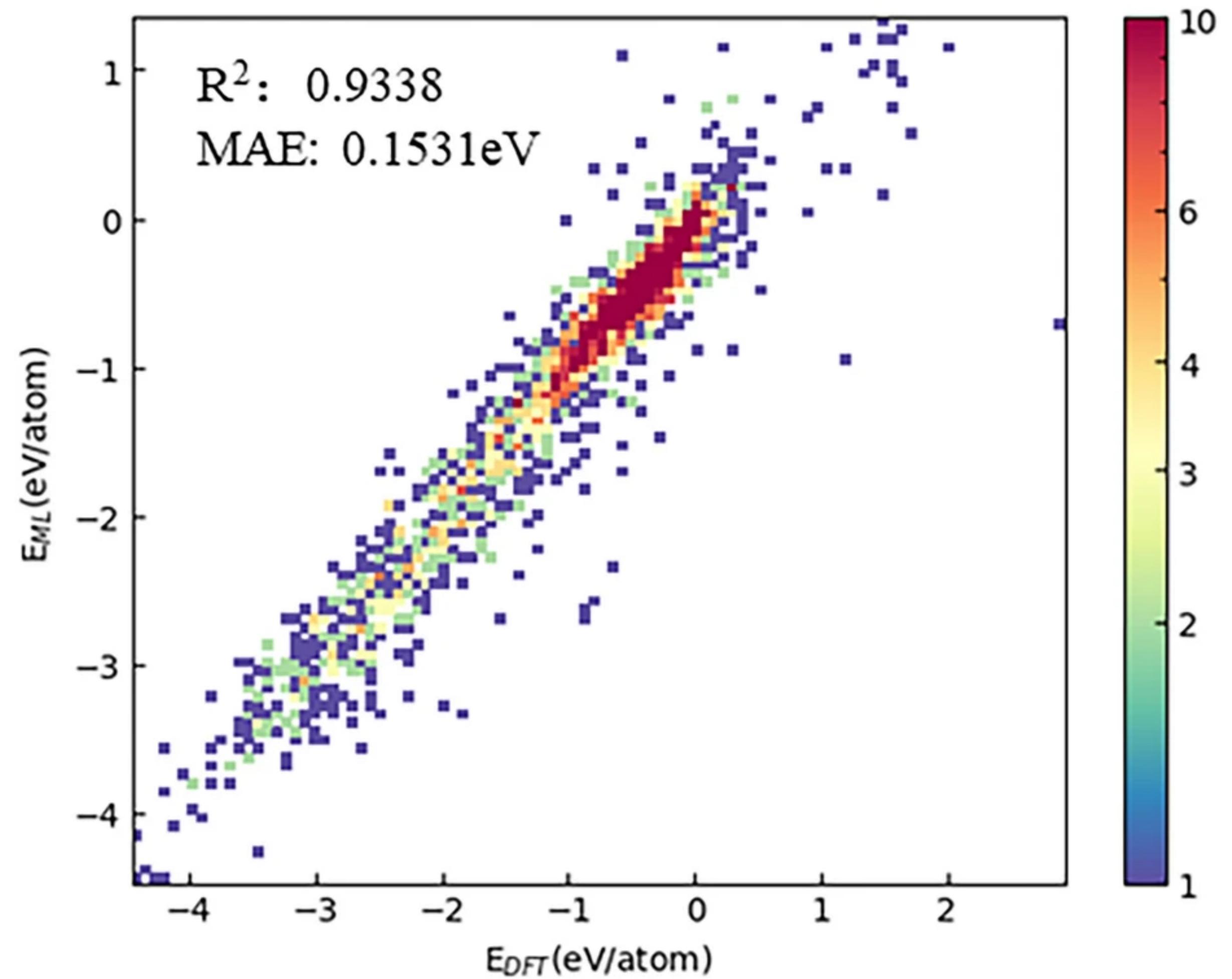
MLP for formation energy

- Input:
 - Aggregated element features
 - Aggregated structure features
- Output:
 - Formation energy per atom



MLP for formation energy

- Input:
 - Aggregated element features
 - Aggregated structure features
- Output:
 - Formation energy per atom



MLP pros/cons

- Pros
 - Universal approximator
 - Versatility -> easy to implement for various tasks (regression/classification)
- Cons
 - Requires large datasets
 - In theory, MLP can capture more complex non-linear relationships compared to e.g. RF
 - In practice, well tuned RF may outperform MLP
 - Black box -> no interpretability
 - Computational complexity

- The Multilayer Perceptron (MLP) is one class of neural network
- These days, Graph Neural Networks (GNNs) are the most popular nets used in computational materials science
- We will discuss GNNs in the next lecture

Resources

- [https://uvadlc-notebooks.readthedocs.io/en/latest/tutorial_notebooks/tutorial2/
Introduction_to_PyTorch.html](https://uvadlc-notebooks.readthedocs.io/en/latest/tutorial_notebooks/tutorial2/Introduction_to_PyTorch.html)
- <https://pabloinsente.github.io/the-multilayer-perceptron>