# CS440/ECE448 Lecture 10: Perceptron

Mark Hasegawa-Johnson, 2/2024

# Outline

- Linear Classifiers
- Gradient descent
- One-hot vectors and the perceptron loss function
- Perceptron learning algorithm

# Linear classifier: Notation

- The observation $\boldsymbol{x}^T = [x_1, \ldots, x_d]$ is a real-valued vector ($d$ is the number of feature dimensions)

- The class label $y \in \mathcal{Y}$ is drawn from some finite set of class labels.

- Usually the output vocabulary, $\mathcal{Y}$, is some set of strings. For convenience, though, we usually map the class labels to a sequence of integers, $\mathcal{Y} = \{1, \ldots, v\}$, where $v$ is the vocabulary size

# Linear classifier: Definition

A linear classifier is defined by

$$f(\boldsymbol{x}) = \operatorname{argmax} \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b}$$

where:

$$\boldsymbol{W}\boldsymbol{x} + \boldsymbol{b} = \begin{bmatrix} w_{1,1} & \cdots & w_{1,d} \\ \vdots & \ddots & \vdots \\ w_{v,1} & \cdots & w_{v,d} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} + \begin{bmatrix} b_1 \\ \vdots \\ b_v \end{bmatrix} = \begin{bmatrix} w_1^T x + b_1 \\ \vdots \\ w_v^T x + b_v \end{bmatrix}$$

$\boldsymbol{w}_k, b_k$ are the **weight vector** and **bias** corresponding to **class k**, and the argmax function finds the element of the vector $wx$ with the largest value.
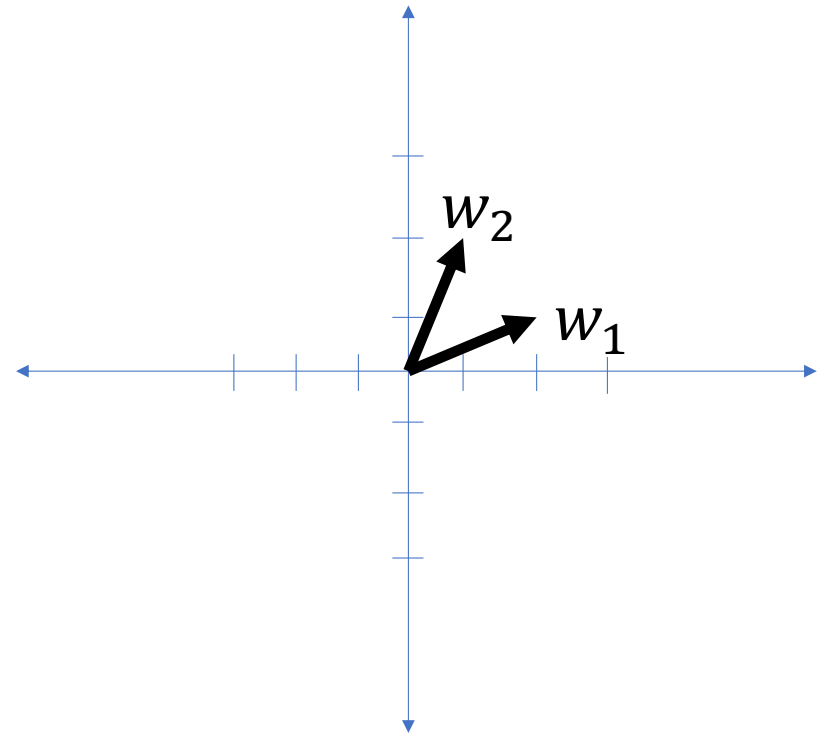
There are a total of $v(d+1)$ trainable parameters: the elements of the matrix $w$.

# Example

Consider a two-class classification problem, with
$$\boldsymbol{w}_1^T = \begin{bmatrix} w_{1,1}, w_{1,2} \end{bmatrix} = [2,1]$$
$$\boldsymbol{w}_2^T = \begin{bmatrix} w_{2,1}, w_{2,2} \end{bmatrix} = [1,2]$$

# Example

Notice that in the two-class case, the equation

$$f(x) = \operatorname{argmax} \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b}$$

Simplifies to

$$f(\boldsymbol{x}) = \begin{cases} 1 & \boldsymbol{w}_1^T\boldsymbol{x} + b_1 > \boldsymbol{w}_2^T\boldsymbol{x} + b_2 \\ 2 & \boldsymbol{w}_1^T\boldsymbol{x} + b_1 < \boldsymbol{w}_2^T\boldsymbol{x} + b_2 \end{cases}$$

The class boundary is the line whose equation is

$$(\boldsymbol{w}_2 - \boldsymbol{w}_1)^T x + (b_2 - b_1) = 0$$



$f(\boldsymbol{x}) = 2$

$\boldsymbol{w}_2$

$\boldsymbol{w}_1$

$f(\boldsymbol{x}) = 1$

# Multi-class linear classifier

In a general multi-class linear classifier,

$$f(\boldsymbol{x}) = \operatorname{argmax} \boldsymbol{W}\boldsymbol{x} + \boldsymbol{b}$$

The boundary between class $k$ and class $l$ is the line (or plane, or hyperplane) given by the equation

$$(\boldsymbol{w}_k - \boldsymbol{w}_l)^T \boldsymbol{x} + (b_k - b_l) = 0$$
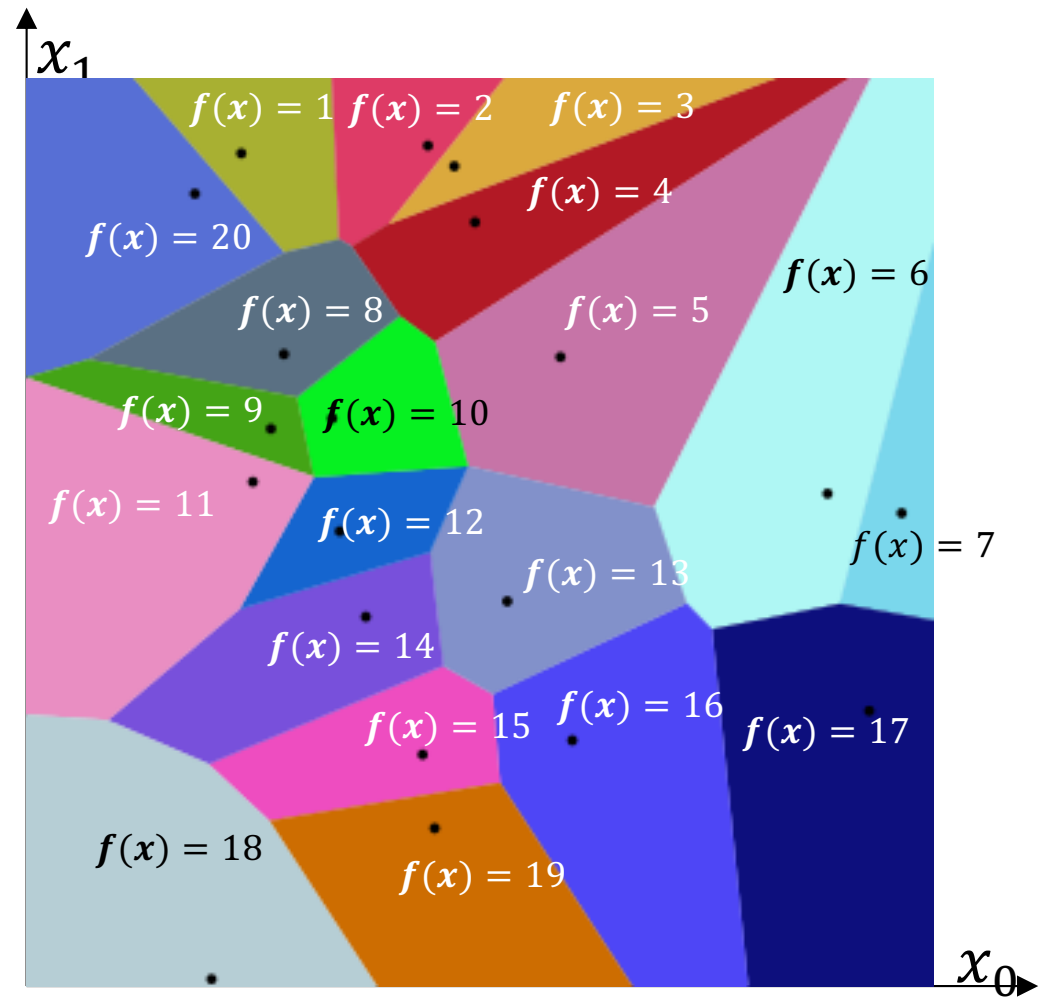
# Voronoi regions

The classification regions in a linear classifier are called Voronoi regions.

A **Voronoi region** is a region that is

- Convex (if $u$ and $v$ are points in the region, then every point on the line segment $\overline{uv}$ connecting them is also in the region)
- Bounded by piece-wise linear boundaries



$x_1$

$f(x) = 1$ $f(x) = 2$ $f(x) = 3$

$f(x) = 4$

$f(x) = 20$

$f(x) = 6$

$f(x) = 8$ $f(x) = 5$

$f(x) = 9$ $f(x) = 10$

$f(x) = 11$

$f(x) = 12$

$f(x) = 7$

$f(x) = 13$

$f(x) = 14$

$f(x) = 16$

$f(x) = 15$ $f(x) = 17$

$f(x) = 18$ $f(x) = 19$

$x_0$

# Outline

- Linear Classifiers
- Gradient descent
- One-hot vectors and the perceptron loss function
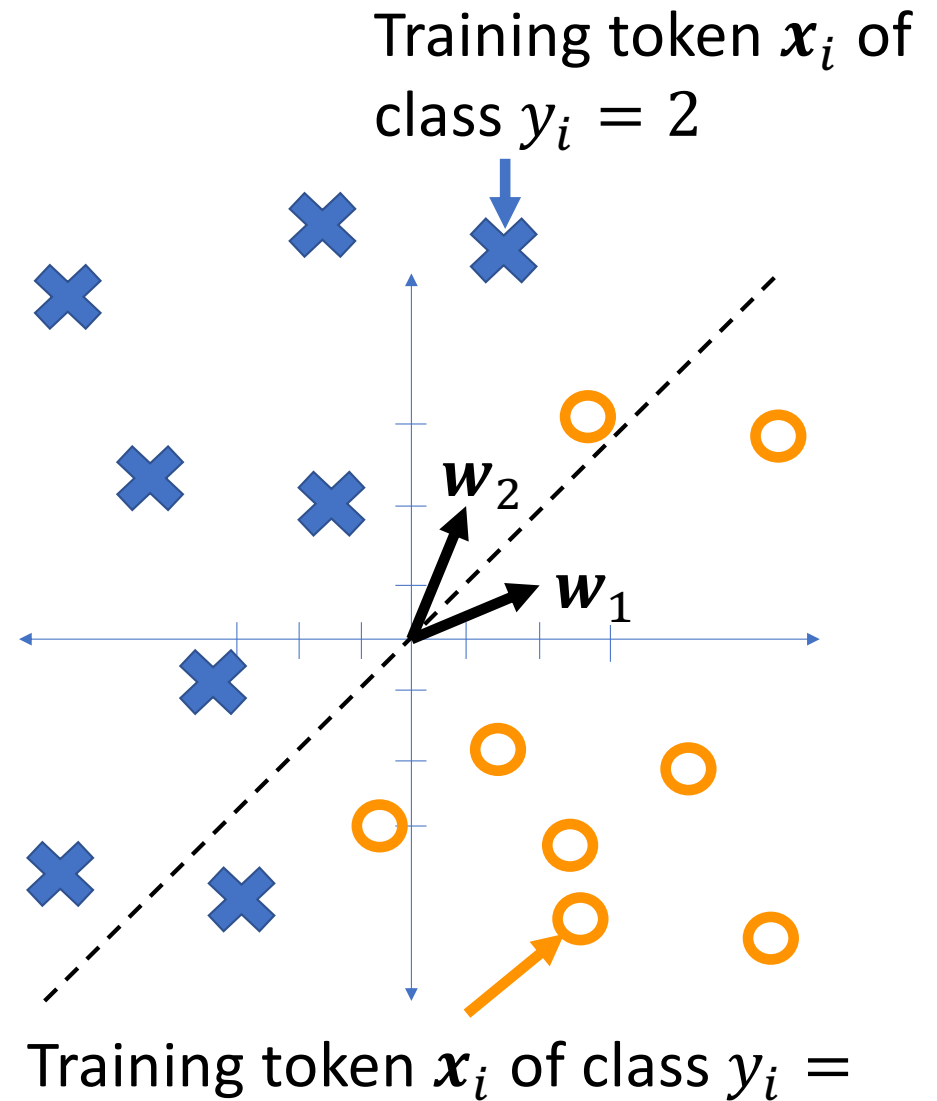- Perceptron learning algorithm

# Gradient descent

Suppose we have training tokens $(x_i, y_i)$, and we have some initial class vectors $w_1$ and $w_2$. We want to update them as

$$\boldsymbol{w}_1 \leftarrow \boldsymbol{w}_1 - \eta \frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_1}$$

$$\boldsymbol{w}_2 \leftarrow \boldsymbol{w}_2 - \eta \frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_2}$$

…where $\mathcal{L}$ is some loss function. What loss function makes sense?

Training token $\boldsymbol{x}_i$ of class $y_i = 2$

$\boldsymbol{w}_2$

$\boldsymbol{w}_1$

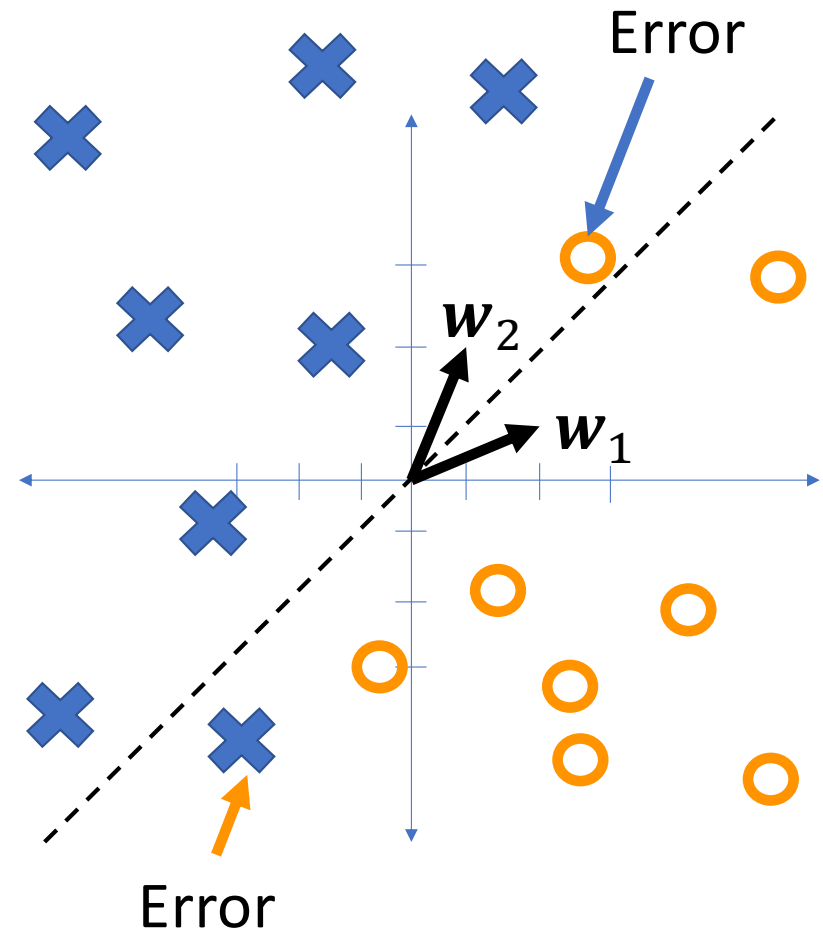Training token $\boldsymbol{x}_i$ of class $y_i =$

# Zero-one loss function

The most obvious loss function for a classifier is its classification error rate,

$$\mathcal{L} = \frac{1}{n}\sum_{i=1}^{n} \ell(f(\boldsymbol{x}_i), y_i)$$

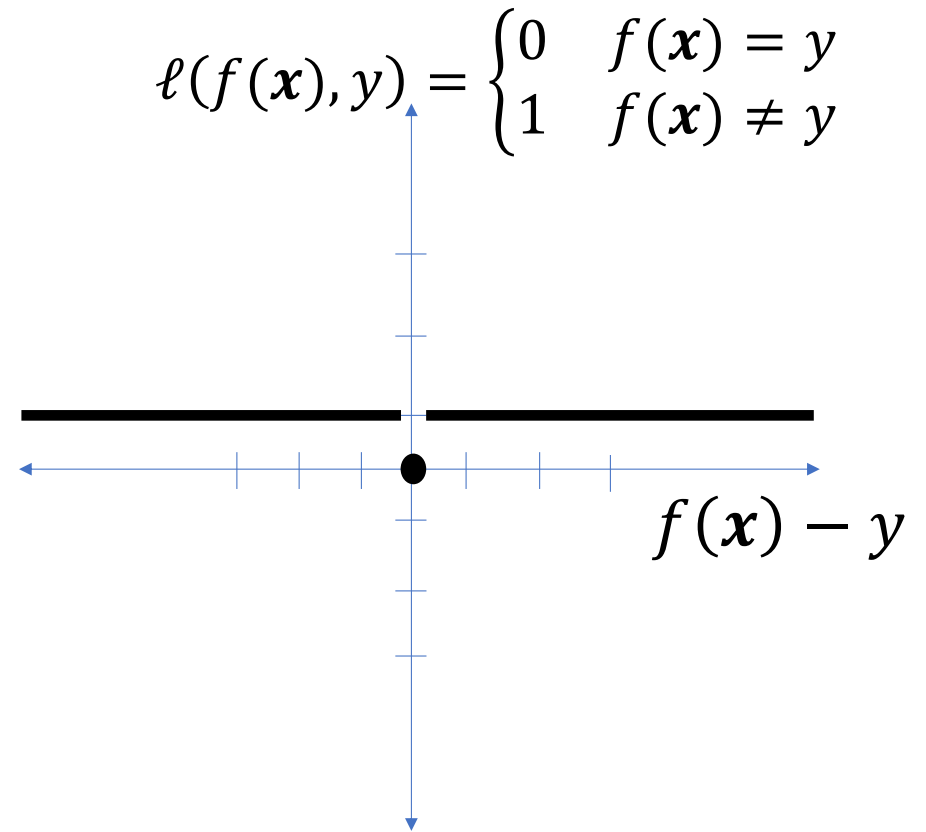Where $\ell(\hat{y}, y)$ is the zero-one loss function,

$$\ell(f(\boldsymbol{x}), y) = \begin{cases} 0 & f(\boldsymbol{x}) = y \\ 1 & f(\boldsymbol{x}) \neq y \end{cases}$$



Error

$\boldsymbol{w}_2$

$\boldsymbol{w}_1$

Error

# Non-differentiable!

The problem with the zero-one loss function is that it's not differentiable:

$$\frac{\partial \ell(f(\boldsymbol{x}), y)}{\partial f(\boldsymbol{x})} = \begin{cases} 0 & f(\boldsymbol{x}) \neq y \\ +\infty & f(\boldsymbol{x}) = y^+ \\ -\infty & f(\boldsymbol{x}) = y^- \end{cases}$$

$$\ell(f(\boldsymbol{x}), y) = \begin{cases} 0 & f(\boldsymbol{x}) = y \\ 1 & f(\boldsymbol{x}) \neq y \end{cases}$$

$$f(\boldsymbol{x}) - y$$

# Outline

- Linear Classifiers: multi-class and 2-class
- Gradient descent
- One-hot vectors and the perceptron loss function
- Perceptron learning algorithm

# One-hot vectors

A **<u>one-hot vector</u>** is a binary vector in which all elements are 0 except for a single element that's equal to 1.
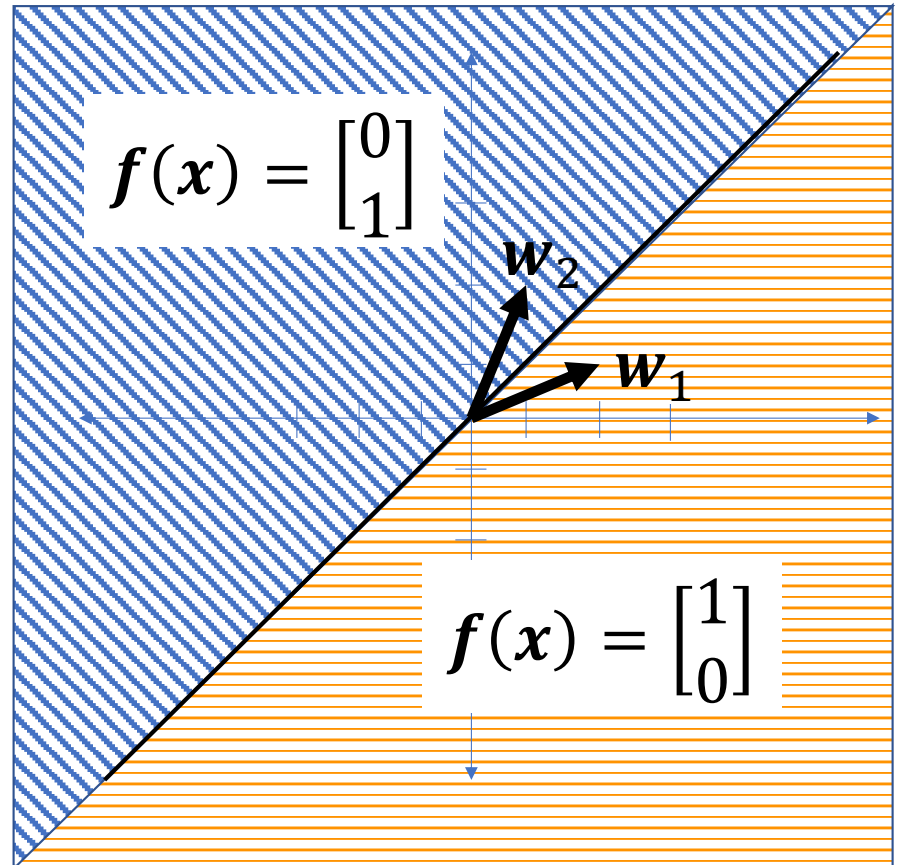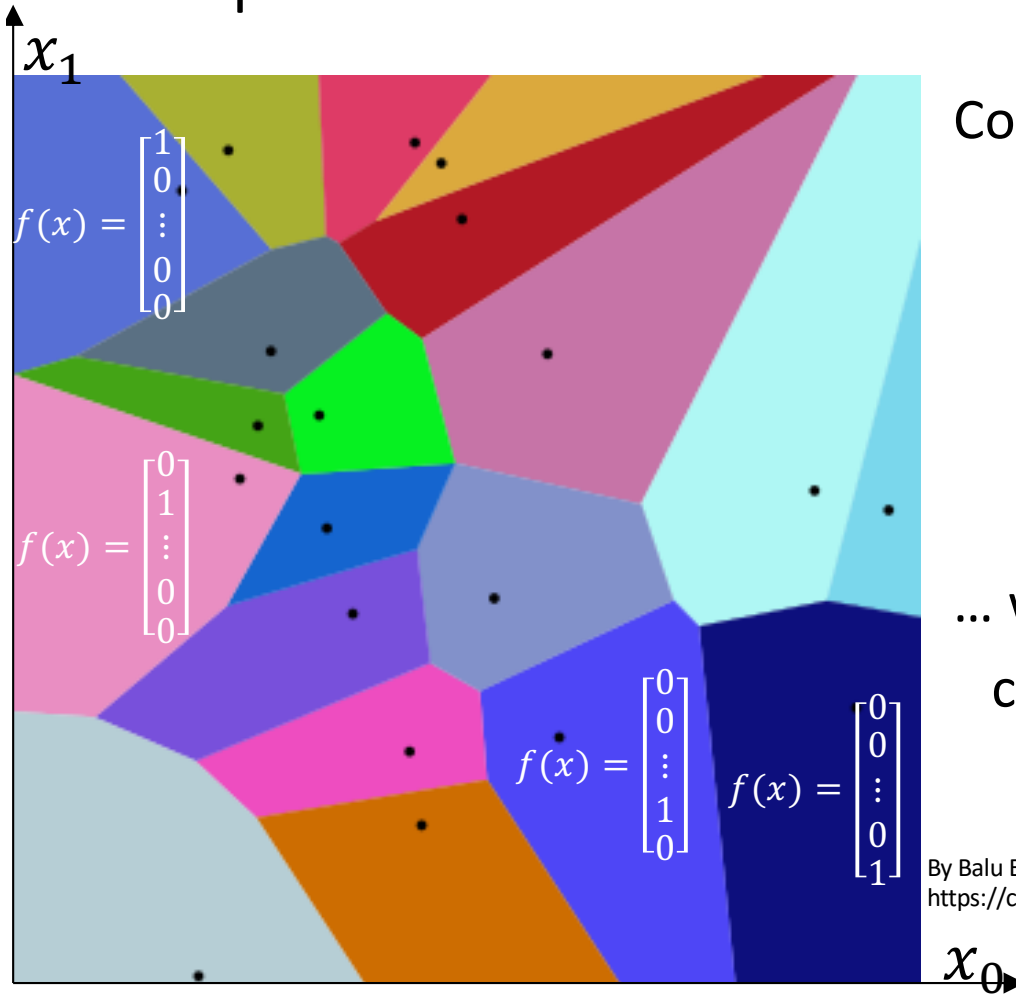
# Example: Binary classifier

Consider the classifier

$$f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{\text{argmax } Wx=1} \\ \mathbb{1}_{\text{argmax } Wx=2} \end{bmatrix}$$

…where $\mathbb{1}_P$ is called the "indicator function," and it means:

$$\mathbb{1}_P = \begin{cases} 1 & P \text{ is true} \\ 0 & P \text{ is false} \end{cases}$$

$$f(x) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$w_2$

$w_1$

$$f(x) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

# Example: Multi-Class



Consider the classifier

$$f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_v(x) \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{\text{argmax } Wx=1} \\ \vdots \\ \mathbb{1}_{\text{argmax } Wx=v} \end{bmatrix}$$

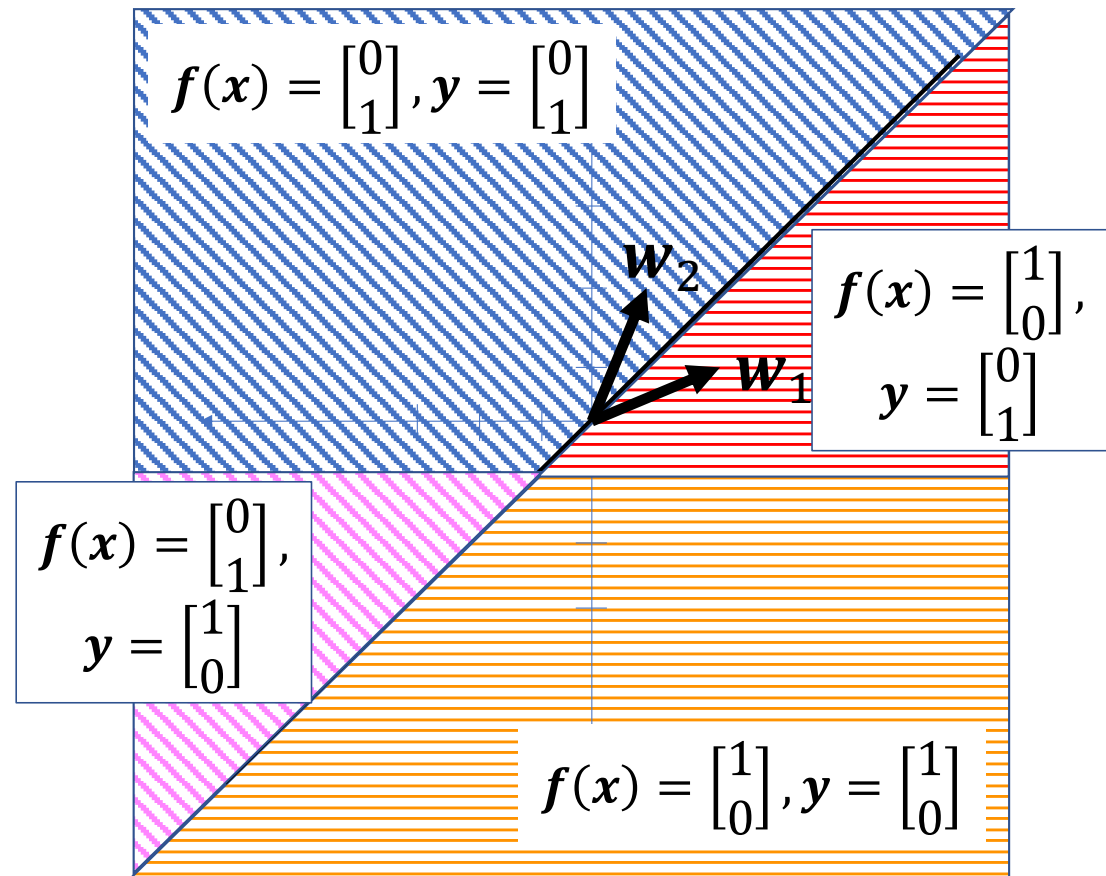… with 20 classes.  Then some of the classifications might look like this.

# One-hot ground truth

We can also use one-hot vectors to describe the ground truth.

Let's call the one-hot vector $\boldsymbol{y}$, and the integer label $y$, thus

$$\boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{y=1} \\ \mathbb{1}_{y=2} \end{bmatrix}$$

Ground truth might differ from classifier output. For example, they might be as shown here:

$$f(\boldsymbol{x}) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \boldsymbol{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$f(\boldsymbol{x}) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$\boldsymbol{w}_2$

$\boldsymbol{w}_1$

$$f(\boldsymbol{x}) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \boldsymbol{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$f(\boldsymbol{x}) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

# Counting errors using one-hot vectors

- An error occurs if $f(x) \neq y$.
- So, to determine whether an error has occurred, we could just check:

$$f(x) - y = \begin{cases} \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} & \text{no error occurred} \\ \text{anything else} & \text{an error occurred} \end{cases}$$

# The perceptron loss

Instead of a one-zero loss, the perceptron uses a weird loss function that gives great results when differentiated.  The perceptron loss function is:

$$\ell(\boldsymbol{x}, \boldsymbol{y}) = (\boldsymbol{f}(\boldsymbol{x}) - \boldsymbol{y})^T (\boldsymbol{W}\boldsymbol{x} + \boldsymbol{b})$$

$$= [f_1(\boldsymbol{x}) - y_1, \quad \cdots, \quad f_v(\boldsymbol{x}) - y_v] \left( \begin{bmatrix} w_{1,1} & \cdots & w_{1,d} \\ \vdots & \ddots & \vdots \\ w_{v,1} & \cdots & w_{v,d} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} + \begin{bmatrix} b_1 \\ \vdots \\ b_v \end{bmatrix} \right)$$

$$= \sum_{k=1}^{v} (f_k(\boldsymbol{x}) - y_k)\left(\boldsymbol{w}_k^T \boldsymbol{x} + b_k\right)$$

# The perceptron loss

$$\ell(\boldsymbol{x}, \boldsymbol{y}) = \sum_{k=1}^{v} (f_k(\boldsymbol{x}) - y_k)(\boldsymbol{w}_k^T \boldsymbol{x} + b_k)$$

Notice that:

$$(f_k(\boldsymbol{x}) - y_k) = \begin{cases} +1 & f_k(\boldsymbol{x}) = 1, y_k = 0 \\ -1 & f_k(\boldsymbol{x}) = 0, y_k = 1 \\ 0 & \text{otherwise} \end{cases}$$

# The perceptron loss

So what the loss really means is:

$$\ell(x, y) = \left(w_{\hat{y}}^T x + b_{\hat{y}}\right) - \left(w_y^T x + b_y\right)$$

Where:

- $y$ is the correct class label for this training token

- $\hat{y} = \underset{k}{\mathrm{argmax}}\left(w_k^T x + b_k\right)$ is the classifier output

- $\ell(x, y) > 0$ if $\hat{y} \neq y$

- $\ell(x, y) = 0$ if $\hat{y} = y$

# Outline

- Linear Classifiers: multi-class and 2-class
- Gradient descent
- One-hot vectors and the perceptron loss function
- **Perceptron learning algorithm**

# Gradient of the perceptron loss

$$\ell(\boldsymbol{x}, \boldsymbol{y}) = \left(\boldsymbol{w}_{\hat{y}}^T \boldsymbol{x} + b_{\hat{y}}\right) - \left(\boldsymbol{w}_y^T \boldsymbol{x} + b_y\right)$$

Its derivative is:

$$\frac{\partial \ell(\boldsymbol{x}, \boldsymbol{y})}{\partial \boldsymbol{w}_k} = \begin{cases} x & k = \hat{y} \\ -x & k = y \\ 0 & \text{otherwise} \end{cases}$$

# The perceptron learning algorithm

1. Compute the classifier output $\hat{y} = \underset{k}{\mathrm{argmax}}\left(\boldsymbol{w}_k^T \boldsymbol{x} + b_k\right)$

2. Update the weight vectors as:

$$\boldsymbol{w}_k \leftarrow \boldsymbol{w}_k - \eta \frac{\partial \ell(\boldsymbol{x}, \boldsymbol{y})}{\partial \boldsymbol{w}_k} = \begin{cases} \boldsymbol{w}_k - \eta \boldsymbol{x} & k = \hat{y} \\ \boldsymbol{w}_k + \eta \boldsymbol{x} & k = y \\ 0 & \text{otherwise} \end{cases}$$

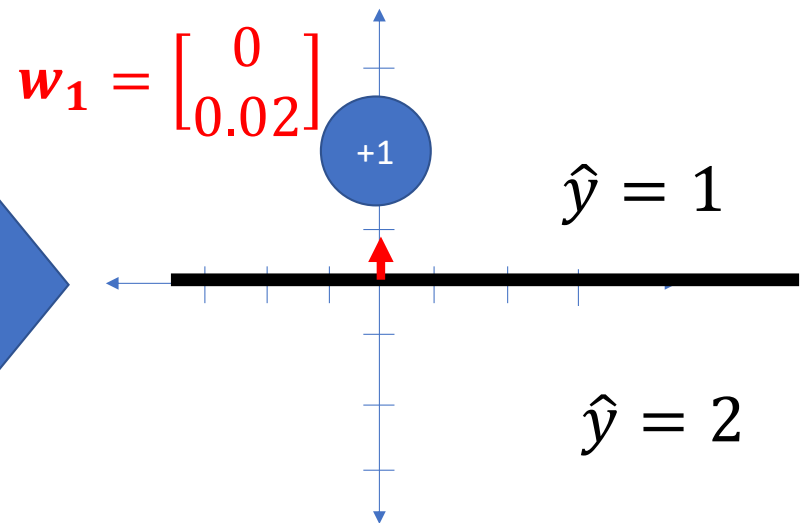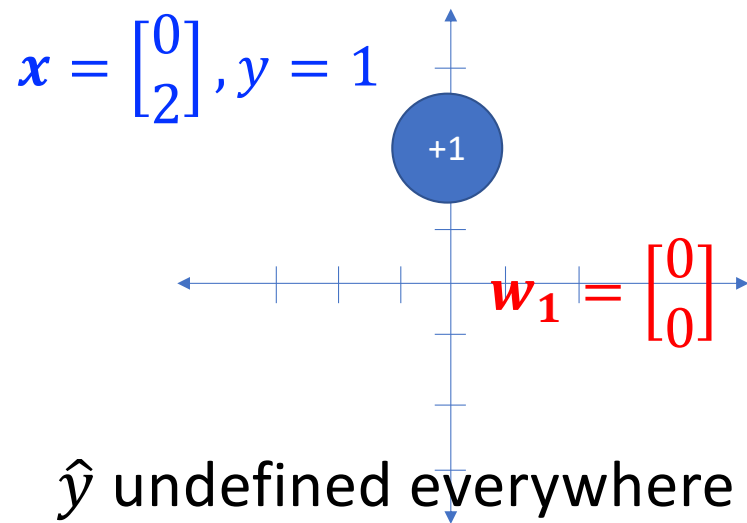where $\eta \approx 0.01$ is the learning rate.

# Example

Start with $w_k = [0,0]^T$ for both classes.

Suppose that $x = [0,2]^T$, with the label $y = 1$.

$\hat{y} = \underset{k}{\mathrm{argmax}}(w_k^T x)$ is undefined, since $w_k^T x = 0$ for both classes, so we only update

$$w_1 \leftarrow w_1 + \eta x = \begin{bmatrix} 0 \\ 0.02 \end{bmatrix}$$



$x = \begin{bmatrix} 0 \\ 2 \end{bmatrix}, y = 1$

$w_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

LEARN!

$\hat{y}$ undefined everywhere

$w_1 = \begin{bmatrix} 0 \\ 0.02 \end{bmatrix}$

$\hat{y} = 1$

$\hat{y} = 2$

# Example
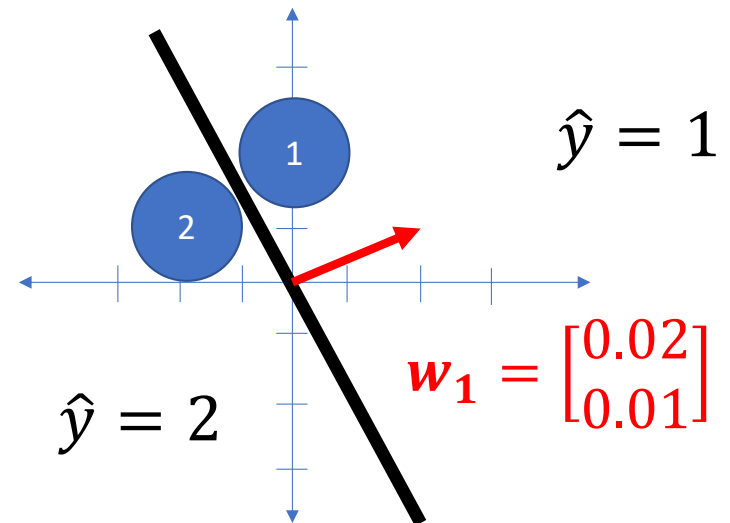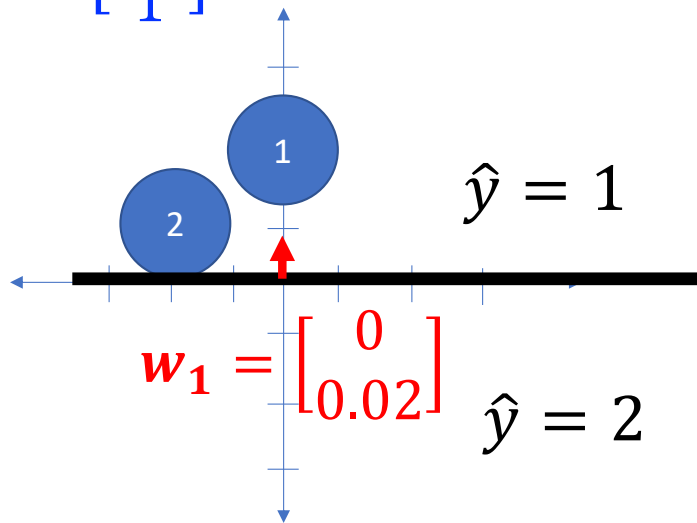
Now $w_1 = [0, 0.02]^T$, but $w_2 = [0, 0]^T$.

Suppose the next $x = [-2, 1]^T$, with the label $y = 2$.

$\hat{y} = \underset{k}{\text{argmax}}(w_k^T x) = 1$ which is wrong, so we update

$$w_1 \leftarrow w_1 - \eta x = \begin{bmatrix} 0.02 \\ 0.01 \end{bmatrix}, \qquad w_2 \leftarrow w_2 + \eta x = \begin{bmatrix} -0.02 \\ 0.01 \end{bmatrix}$$

$x = \begin{bmatrix} -2 \\ 1 \end{bmatrix}, y = 2$



$\hat{y} = 1$

LEARN!

$w_1 = \begin{bmatrix} 0 \\ 0.02 \end{bmatrix}$ $\hat{y} = 2$

$\hat{y} = 1$
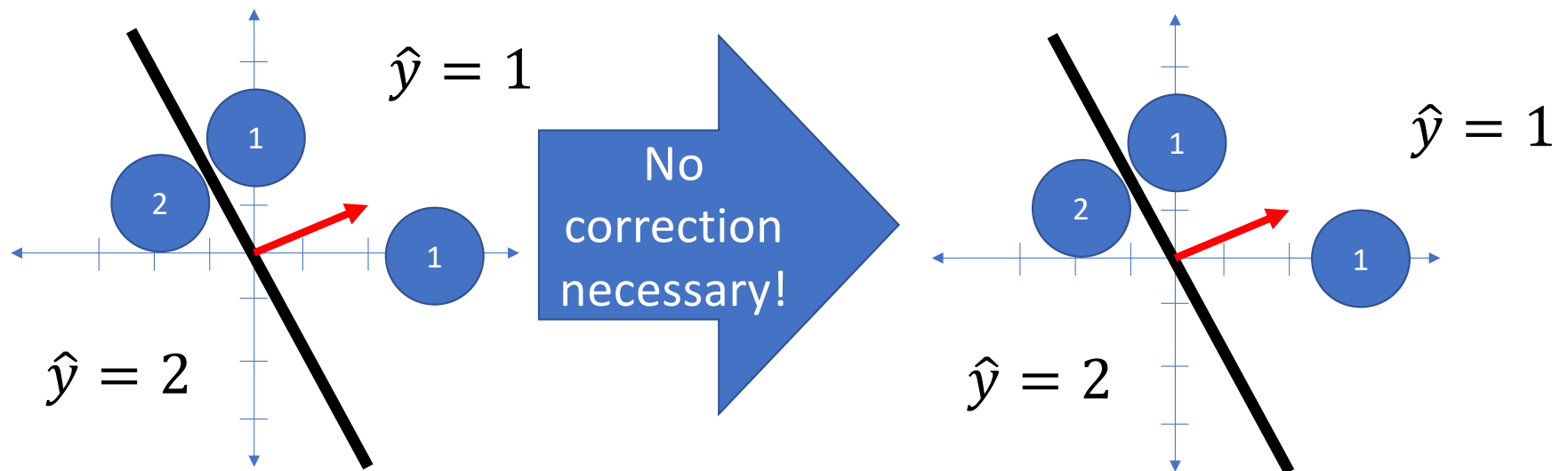
$\hat{y} = 2$ $w_1 = \begin{bmatrix} 0.02 \\ 0.01 \end{bmatrix}$

# Example

Suppose the next token is $x = [3,0]^T$, with the label $y = 1$. Since $\hat{y}$ is right, the weights don't need to be updated:
$$w_k \leftarrow w_k + 0$$

$x = \begin{bmatrix} 3 \\ 1 \end{bmatrix}, y = 1$



$w_k$ unchanged

# The perceptron learning algorithm

1. Compute the classifier output $\hat{y} = \underset{k}{\mathrm{argmax}}\left(\boldsymbol{w}_k^T \boldsymbol{x} + b_k\right)$

2. Update the weight vectors as:

$$\boldsymbol{w}_k \leftarrow \begin{cases} \boldsymbol{w}_k - \eta \boldsymbol{x} & k = \hat{y} \\ \boldsymbol{w}_k + \eta \boldsymbol{x} & k = y \\ 0 & \text{otherwise} \end{cases}$$

where $\eta \approx 0.01$ is the learning rate.

# Try the quiz!

Try the quiz:
https://us.prairielearn.com/pl/course_instance/147925/assessment/2395719

# Special case: two classes

If there are only two classes, then we only need to learn one weight vector, $\boldsymbol{w} = \boldsymbol{w}_1 - \boldsymbol{w}_2$. We can learn it as:

1. Compute the classifier output $\hat{y} = \underset{k}{\text{argmax}}\left(\boldsymbol{w}_k^T \boldsymbol{x} + b_k\right)$

2. Update the weight vectors as:

$$\boldsymbol{w} \leftarrow \begin{cases} \boldsymbol{w} - \eta \boldsymbol{x} & \hat{y} \neq y, y = 2 \\ \boldsymbol{w} + \eta \boldsymbol{x} & \hat{y} \neq y, y = 1 \\ 0 & \hat{y} = y \end{cases}$$

where $\eta \approx 0.01$ is the learning rate. Sometimes we say $y \in \{1, -1\}$ instead of $y \in \{1,2\}$.

# Outline

- Linear Classifiers: $f(x) = \text{argmax } Wx + b$

- Gradient descent: $w_c \leftarrow w_c - \eta \dfrac{\partial \mathcal{L}}{\partial w_c}$

- One-hot vectors: $f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_v(x) \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{\text{argmax } Wx=1} \\ \vdots \\ \mathbb{1}_{\text{argmax } Wx=v} \end{bmatrix}, \ y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{y=1} \\ \mathbb{1}_{y=2} \\ \vdots \end{bmatrix}$

- Perceptron learning algorithm:

$$w_c \leftarrow \begin{cases} w_c - \eta x & c = \hat{y} \\ w_c + \eta x & c = y \\ 0 & \text{otherwise} \end{cases}$$