

Nonlinearity

Lecture 03 — CS 577 Deep Learning

Instructor: Yutong Wang

Computer Science
Illinois Institute of Technology

September 4, 2024

Administrative matter

- For the course project, you can form your own groups (2-4 people) or have the groups be assigned to you randomly.

Notations

Let $i = 1, \dots, N$ (the sample index)

- Training samples $\mathbf{x}^{(i)} \in \underline{\mathcal{X}} \subseteq \mathbb{R}^d$ 
- Labels $y^{(i)} \in \mathcal{Y} = \mathbb{R}$
- $f(\cdot; \boldsymbol{\theta}) : \mathcal{X} \rightarrow \mathcal{Y}$

Linearity

Definition. $f(\cdot; \theta)$ is *linear* if there exists some C such that

$$\begin{aligned} f(\mathbf{x} + \mathbf{x}'; \theta) - C &= \frac{f(\mathbf{x}; \theta) + f(\mathbf{x}'; \theta)}{\lambda f(\lambda \mathbf{x}; \theta) - \lambda C} - C \\ f(\lambda \mathbf{x}; \theta) - C &= \end{aligned}$$

Note: See [GBC16, p. 5.1.4] regarding “affine” vs “linear”

Linearity

Definition. $f(\cdot; \theta)$ is *linear* if there exists some C such that

$$f(\mathbf{x} + \mathbf{x}'; \theta) - C = \mathbf{w}^\top (\mathbf{x} + \mathbf{x}') + b$$

$$f(\lambda \mathbf{x}; \theta) - C = \underbrace{\mathbf{w}^\top \mathbf{x} + \mathbf{w}^\top \mathbf{x}' + b + b - C}$$

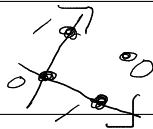
Example (linear regression).

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^\top \mathbf{x} + b \quad f(\mathbf{x}; \theta) + f(\mathbf{x}'; \theta) - C$$

Note: Limitations of linearity

Do: `lec03-in-class-ex1-xor.ipynb`

Discuss: does “learning” occur?



Set to
 $C = b$

Linearity

Definition. $f(\cdot; \theta)$ is *linear* if there exists some C such that

$$f(\mathbf{x} + \mathbf{x}'; \theta) - C =$$

$$f(\lambda \mathbf{x}; \theta) - C =$$

Example (linear regression).

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^\top \mathbf{x} + b$$

Note: Limitations of linearity

Do: [lec03-in-class-ex1-xor.ipynb](#)

Discuss: does “learning” occur?

It seems like

$w \rightarrow 0$.

What does this mean?

y doesn't depend on x

Not at all

$f_{\text{true}}(x_1, x_2) = x_1$, exclusive x_2

Linearity

Definition. $f(\cdot; \theta)$ is *linear* if there exists some C such that

$$f(\mathbf{x} + \mathbf{x}'; \theta) - C =$$

$$f(\lambda \mathbf{x}; \theta) - C =$$

Example (linear regression).

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^\top \mathbf{x} + b$$

Note: Limitations of linearity

Do: `lec03-in-class-ex1-xor.ipynb`

Discuss: does “learning” occur?

Empirical observation

Beyond linearity: feature map

Definition. $f(\cdot; \theta)$ is *linear* if there exists some C such that

$$f(\mathbf{x} + \mathbf{x}'; \theta) - C =$$

$$f(\lambda \mathbf{x}; \theta) - C =$$

Example (linear regression with a (fixed) feature map). $\phi: \mathbb{R}^1 \rightarrow \mathbb{R}^D$

$$f(x; \mathbf{w}, b) = \mathbf{w}^\top \underbrace{\phi(x)}_{\tilde{x}} + b$$

degree d polynomial feature map $\tilde{x} = \text{engineered features}$

$$\phi(x) = [1, x, x^2, \dots, x^d]$$

For 1-D data only

Can we apply it to XOR?

Beyond linearity: feature map

```
def polynomial_feature_map(x, degree):  
    return [x**d for d in range(degree+1)]  
  
def polynomial_feature_map(x_array, degree):  
    return np.array([polynomial_feature_map(x_array[i], degree) for i in range(len(x_array))])  
  
deg = 37  
  
# has a single HYPERPARAMETER -> deg  
Xtilde = polynomial_feature_map(x, deg)
```

outer func

degree = 2

$\{ ax^2 + bx + c : a, b, c \in \mathbb{R} \}$

poly in 1 unknown

Beyond linearity: feature map

```
def polynomial_feature_map(x, degree):  
    return [x**d for d in range(degree+1)]  
  
def polynomial_feature_map(x_array, degree):  
    return np.array([polynomial_feature_map(x_array[i], degree) for i in range(len(x_array))])  
  
deg = 37  
  
# has a single HYPERPARAMETER -> deg  
Xtilde = polynomial_feature_map(x, deg)
```

$$\begin{aligned} \text{Poly deg 2} &= ax^2 + bx + c \\ &= \underbrace{\begin{bmatrix} a & b & c \end{bmatrix}}_w \begin{bmatrix} x^2 \\ x \\ 1 \end{bmatrix} = \phi(x) \end{aligned}$$

Beyond linearity: feature map

```
def polynomial_feature_map_(x, degree):  
    return [x**d for d in range(degree+1)]  
  
def polynomial_feature_map(x_array, degree):  
    return np.array([polynomial_feature_map_(x_array[i], degree) for i in range(len(x_array))])  
  
deg = 37  
  
# has a single HYPERPARAMETER -> deg  
Xtilde = polynomial_feature_map(x, deg)
```

poly of deg 2 in 2 free variable

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad a x_2^2 + b x_1^2 + c x_1 x_2 + d x_1 + e x_2 + f \cdot 1$$

Beyond linearity: feature map

Definition. $f(\cdot; \theta)$ is *linear* if there exists some C such that

$$f(\mathbf{x} + \mathbf{x}'; \theta) - C =$$

$$f(\lambda \mathbf{x}; \theta) - C =$$

Example (linear regression with a (fixed) feature map). $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$

Note: Implement the feature map for degree 2 polynomials in 2 free variables.

Can capture
XOR

with
deg 2
poly
feat
map

$$d = 2$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$D = ?$$

all degree 2 polynomial feature

Linearity

Example (linear regression with a (varying) feature map). $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$

$$f(\mathbf{x}; \mathbf{w}^{(2)}, b^{(2)}) = \mathbf{w}^{(2)\top} \phi(\mathbf{x}) + b^{(2)}$$

$$\phi(\mathbf{x}; \mathbf{w}^{(1)}, b^{(1)}) = \underbrace{\mathbf{w}^{(1)\top}}_{\text{feature map with tunable parameter}} \mathbf{x} + b^{(1)}$$

feature
map
with

tunable parameter

$$\in \mathbb{R}^d \times D \quad \Bigg| \quad \in \mathbb{R}^D$$

Linearity

Example (linear regression with a (varying) feature map). $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$

$$\phi(\mathbf{x}) = \mathbf{w}^{(1)\top} \mathbf{x} + b^{(1)}$$

Still linear

with disguise \rightarrow

$$= (\mathbf{w}^{(1)} \mathbf{w}^{(2)})^\top \mathbf{x} + \mathbf{w}^{(2)\top} b^{(1)} + b^{(2)}$$

Linearity

Example (linear regression with a (varying) feature map). $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^\top \phi(\mathbf{x}) + b$$

$$\phi(\mathbf{x}) =$$

Composition of matrix mul
+ bias

= mat mul + bias

Linearity

2-layer neural network with “linear” activation. $\mathbf{W}^{(1)} \in \mathbb{R}^{d \times n_1}$ Mat Mul + bias

$$f^{(1)}(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}) = \phi(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)})$$

Matmul + bias \rightarrow

$$f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}) = \mathbf{w}^{(2)\top} \mathbf{h} + b^{(2)}$$

hidden nodes

$$f(\mathbf{x}; \mathbf{w}^{(2)}, b^{(2)}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}) := f^{(2)}(f^{(1)}(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}); \mathbf{w}^{(2)}, b^{(2)})$$

still linear

Note: Is this still linear?

Linearity

2-layer neural network with “linear” activation.

$$f\left(\mathbf{x}; \mathbf{w}^{(2)}, b^{(2)}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}\right) := \mathbf{w}^{(2)\top} (\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}$$

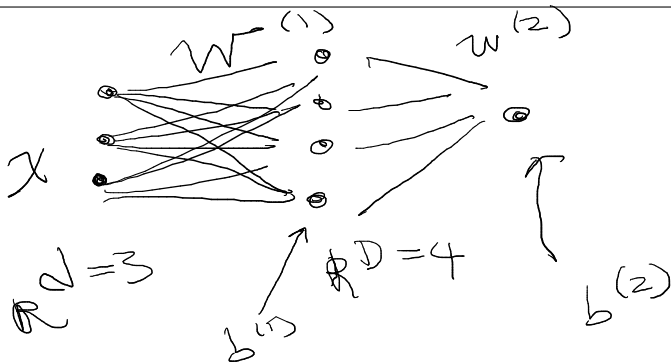
Note: Is this still linear?

Linearity

2-layer neural network with “linear” activation.

$$f(\mathbf{x}; \mathbf{w}^{(2)}, b^{(2)}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}) := \mathbf{w}^{(2)\top} (\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}$$

Note: Draw the network architecture diagram



Linearity

$$\text{mat mul } g(\text{mat mul} + \text{bias}) + \text{bias}$$

2-layer neural network with “non-linear” activation $g : \mathbb{R} \rightarrow \mathbb{R}$.

$$f^{(1)}(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}) = g(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)})$$

$$f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}) = \mathbf{w}^{(2)\top} \mathbf{x} + b^{(2)}$$

$$f(\mathbf{x}; \mathbf{w}^{(2)}, b^{(2)}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}) := f^{(2)}(f^{(1)}(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}); \mathbf{w}^{(2)}, b^{(2)})$$

Activation function

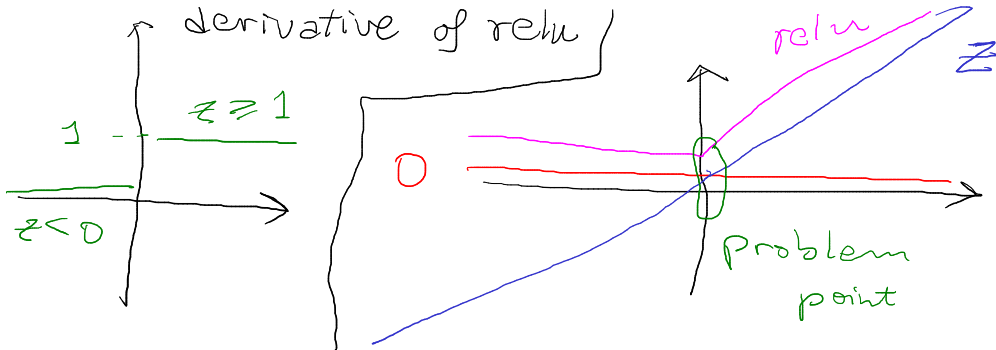
Rectified linear unit or “relu”

$$g(z): \mathbb{R} \rightarrow \mathbb{R}$$

$$\text{relu}(z) := \max\{0, z\}$$

$$\text{relu}'(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

Note: Plot “relu” and its derivative



Linearity

2-layer neural network with “non-linear” activation $g : \mathbb{R} \rightarrow \mathbb{R}$.

$$f(\mathbf{x}; \mathbf{w}^{(2)}, b^{(2)}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}) := \mathbf{w}^{(2)\top} g(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}$$

Note: Is this still linear?

No!

6:21 PM

We will see through
example

Calculation of the gradient

2-layer NN
with hidden width
 n_1

$$f(\mathbf{x}; \mathbf{w}^{(2)}, b^{(2)}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}) := \mathbf{w}^{(2)\top} g(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}$$

- 1-dimensional data $d=1$
- $\mathbf{x}^{(i)} = x^{(i)}$ (no bold)
- $\mathbf{W}^{(1)} \in \mathbb{R}^{d \times n_1}$

$$\mathbf{W}^{(1)} = [W_{n_1}^{(1)}, \dots, W_1^{(1)}]$$
$$\in \mathbb{R}^{d \times n_1}$$

$$\underbrace{\mathbf{w}^{(2)\top} g(\underbrace{\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}}_{n_1})}_{\text{entrywise scalar}} + b^{(2)} \quad \text{still scalar}$$

Calculation of the gradient

$$\mathbf{w}^{(2)}, \mathbf{w}^{(1)} \in \mathbb{R}^{n_1}$$

$$f(\mathbf{x}; \mathbf{w}^{(2)}, b^{(2)}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}) := \mathbf{w}^{(2)\top} g(\underbrace{\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}}_{\text{lower case } \mathbf{w}^{(1)} \text{ or } x}) + b^{(2)}$$

- 1-dimensional data $d = 1$
- $\mathbf{x}^{(i)} = x^{(i)}$ (no bold)
- $\mathbf{W}^{(1)} \in \mathbb{R}^{d \times n_1}$

x just a number

lower case

$\mathbf{w}^{(1)}$
or x

rename

$$\mathbf{W}^{(1)} = \begin{bmatrix} w_1^{(1)} & \dots & w_{n_1}^{(1)} \end{bmatrix}$$

$$\begin{bmatrix} w_1^{(1)} \\ \vdots \\ w_{n_1}^{(1)} \end{bmatrix} x$$

Derivative with respect to a generic point

Note: Calculate the derivative when $d = 1$

MSE

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N J_i(\boldsymbol{\theta})$$

Sum of sample-wise MSE

where

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - \mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)})^2$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}} J_i(\boldsymbol{\theta})$$

Focus on a single i .

$$\begin{bmatrix} w_1^{(1)} \\ \vdots \\ w_n^{(1)} \end{bmatrix}$$

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \left\{ \nabla_{w^{(2)}} J(\boldsymbol{\theta}) \right. \\ &\quad \left. \nabla_{b^{(2)}} J(\boldsymbol{\theta}), \dots \right\} \end{aligned}$$

Derivative with respect to a generic point

$$J_i(\theta) := (y^{(i)} - \mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)})^2 = (y^{(i)} - z^{(i)})^2$$

Note: Calculate the derivative

$$\frac{\partial J_i(\theta)}{\partial \mathbf{w}^{(2)}} = 2(y^{(i)} - z^{(i)}) (-1) \frac{\partial z^{(i)}}{\partial \mathbf{w}^{(2)}}$$

$\parallel z^{(i)}$ "model output"

b.c it's

$$-z_i$$

Fact

$$\frac{\partial (a^\top b)}{\partial a} = b$$

Requires
b not
dependent
on a

Derivative with respect to a generic point

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - \underbrace{\mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)}}_{z_i})^2$$

Note: Calculate the derivative

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(2)}} =$$

z_i

$$\frac{\partial z_i}{\partial w^{(2)}} = g(w^{(1)} x^{(i)} + b^{(1)})$$

Derivative with respect to a generic point

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - \mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)})^2$$

Note: Calculate the derivative

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial b^{(2)}} = -2 (y^{(i)} - z^{(i)}) \frac{\partial z^{(i)}}{\partial b^{(2)}}$$

⏟
11
1

Derivative with respect to a generic point

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - \overbrace{\mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)})}^{z^{(i)}} + b^{(2)})^2$$

Note: Calculate the derivative

Hint: \odot denotes element-wise product between vectors

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(1)}} =$$

$$\mathbf{w}^{(1)} = \begin{bmatrix} w_1^{(1)} \\ \vdots \\ w_{n_1}^{(1)} \end{bmatrix}$$

$$\frac{\partial}{\partial w_v^{(1)}} J_i(\boldsymbol{\theta}) = -2(y^{(i)} - z^{(i)}) \frac{\partial z^{(i)}}{\partial w_v^{(1)}}$$

Derivative with respect to a generic point

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - \mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)})^2$$

↓ entrywise

Note: Calculate the derivative

Hint: \odot denotes element-wise product between vectors

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(1)}} = \frac{\partial z^{(i)}}{\partial w_1^{(1)}} = \frac{\partial}{\partial w_1^{(1)}} \left(w_1^{(2)} g(w_1^{(1)} x^{(i)} + b_1^{(1)}) + \dots + w_{n_i}^{(2)} g(w_{n_i}^{(1)} x^{(i)} + b_{n_i}^{(1)}) + b^{(2)} \right)$$

only part that depends $w_1^{(1)}$

Derivative with respect to a generic point

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - \mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)})^2$$

Note: Calculate the derivative

Hint: \odot denotes element-wise product between vectors

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(1)}} =$$

$$\frac{\partial z^{(i)}}{\partial w_1^{(1)}} = \frac{\partial}{\partial w_1^{(1)}} \left(\underbrace{w_1^{(2)} g(w_1^{(1)} x^{(i)} + b_1^{(1)})}_{\text{only part}} \rightarrow \left. \begin{array}{l} \text{const} \\ \text{wrt} \\ w_1^{(1)} \end{array} \right\} \right)$$
$$= w_1^{(2)} g'(w_1^{(1)} x^{(i)} + b_1^{(1)}) x^{(i)}$$

Derivative with respect to a generic point

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - \mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)})^2$$

Note: Calculate the derivative

Hint: \odot denotes element-wise product between vectors

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(1)}} = \frac{\partial z^{(1)}}{\partial \mathbf{w}^{(1)}} = \frac{\partial}{\partial \mathbf{w}^{(1)}} \left(\mathbf{w}^{(2)} \odot g(\mathbf{w}^{(1)} x^{(1)} + \mathbf{b}^{(1)}) \right)$$

only part

$$- \mathbf{w}^{(2)} g'(\mathbf{w}^{(1)} x^{(1)} + \mathbf{b}^{(1)}) x^{(1)}$$

Derivative with respect to a generic point

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - \mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)})^2$$

Note: Calculate the derivative

Hint: \odot denotes element-wise product between vectors

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(1)}} =$$

Sanity check

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(1)}} = \mathbf{w}^{(2)} \odot g'(\underbrace{\mathbf{w}^{(1)\top} \mathbf{x}^{(i)} + b^{(1)}}_{n_i}) \mathbf{x}^{(i)}$$

$\dim(\cdot) = n_i$

$\mathbf{w}^{(1)} + \frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(1)}}$

Derivative with respect to a generic point

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - \mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)})^2$$

Note: Calculate the derivative

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(1)}} =$$

$$-2(y^{(i)} - z^{(i)}) \underbrace{\frac{\partial z_i}{\partial b^{(1)}}}_{\text{chain rule}}$$

$$w^{(2)} \odot g'(\mathbf{w}^{(1)} x^{(i)} + b^{(1)})$$

1-layer neural network

$$J_i(\boldsymbol{\theta}) := (y^{(i)} - (y^{(i)} - z^{(i)}))^2 \quad \text{where} \quad z_i = \mathbf{w}^{(2)\top} g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}) + b^{(2)}$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(2)}} = -2(y^{(i)} - z^{(i)}) g(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)})$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial b^{(2)}} = -2(y^{(i)} - z^{(i)})$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{w}^{(1)}} = -2(y^{(i)} - z^{(i)}) (\mathbf{w}^{(2)} \odot g'(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)})) x^{(i)}$$

$$\frac{\partial J_i(\boldsymbol{\theta})}{\partial \mathbf{b}^{(1)}} = -2(y^{(i)} - z^{(i)}) (\mathbf{w}^{(2)} \odot g'(\mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}))$$

1-layer neural network

Given params
calculate grad

$$J_i(\theta) := (y^{(i)} - (y^{(i)} - z^{(i)}))^2 \quad \text{where} \quad z_i = \mathbf{w}^{(2)\top} g(\mathbf{h}^{(i)}) + b^{(2)} \quad \text{and} \quad \mathbf{h}^{(i)} = \mathbf{w}^{(1)} x^{(i)} + \mathbf{b}^{(1)}$$



$$\frac{\partial J_i(\theta)}{\partial \mathbf{w}^{(2)}} = -2(y^{(i)} - z^{(i)})g(\mathbf{h}^{(i)})$$

$$\frac{\partial J_i(\theta)}{\partial b^{(2)}} = -2(y^{(i)} - z^{(i)})$$

(4)

$$\frac{\partial J_i(\theta)}{\partial \mathbf{w}^{(1)}} = -2(y^{(i)} - z^{(i)}) (\mathbf{w}^{(2)} \odot g'(\mathbf{h}^{(i)})) x^{(i)}$$

$$\frac{\partial J_i(\theta)}{\partial \mathbf{b}^{(1)}} = -2(y^{(i)} - z^{(i)}) (\mathbf{w}^{(2)} \odot g'(\mathbf{h}^{(i)}))$$

(2)

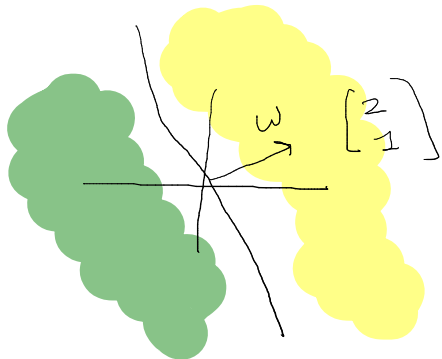
$$\nabla J = \frac{\sum_{i=1}^N \nabla J_i}{N}$$

Note: Do: [lec03-in-class-ex2-relu-net.ipynb](#)

Binary linear classifier

$\mathbf{w} \in \mathbb{R}^d$ with classifier given by

$$f(\mathbf{x}; \mathbf{w}) := \underline{\text{sign}(\mathbf{w}^\top \mathbf{x})} \in \{\pm 1\} \quad (1)$$



Perceptron

$\mathbf{w} \in \mathbb{R}^d$ with classifier given by

$$f(\mathbf{x}; \mathbf{w}) := \text{sign}(\mathbf{w}^\top \mathbf{x}) \in \{\pm 1\}$$

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{w} = \mathbf{0}$.

Output: \mathbf{w}

Perceptron

$\mathbf{w} \in \mathbb{R}^d$ with classifier given by

$$f(\mathbf{x}; \mathbf{w}) := \text{sign}(\mathbf{w}^\top \mathbf{x}) \in \{\pm 1\}$$

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{w} = \mathbf{0}$.
2. For $t = 1, 2, \dots, T$

Output: \mathbf{w}

Perceptron

$\mathbf{w} \in \mathbb{R}^d$ with classifier given by

$$f(\mathbf{x}; \mathbf{w}) := \text{sign}(\mathbf{w}^\top \mathbf{x}) \in \{\pm 1\}$$

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{w} = \mathbf{0}$.
2. For $t = 1, 2, \dots, T$
 - 2.1 If $y^{(t)} \mathbf{w}^\top \mathbf{x}^{(t)} > 0$, then $\mathbf{w} \leftarrow \mathbf{w}$,

Output: \mathbf{w}

Perceptron

$\mathbf{w} \in \mathbb{R}^d$ with classifier given by

$$f(\mathbf{x}; \mathbf{w}) := \text{sign}(\mathbf{w}^\top \mathbf{x}) \in \{\pm 1\}$$

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{w} = \mathbf{0}$.
2. For $t = 1, 2, \dots, T$
 - 2.1 If $y^{(t)} \mathbf{w}^\top \mathbf{x}^{(t)} > 0$, then $\mathbf{w} \leftarrow \mathbf{w}$,
 - 2.2 Else, then $\mathbf{w} \leftarrow \mathbf{w} + y^{(t)} \mathbf{x}^{(t)}$.

Output: \mathbf{w}

Perceptron: an example

$\mathbf{w} \in \mathbb{R}^d$ with classifier given by

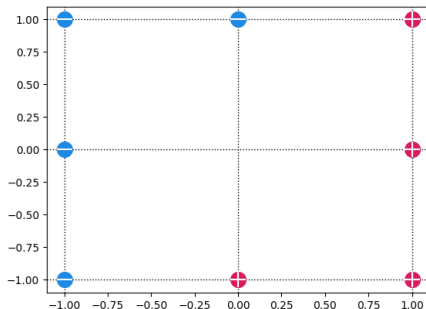
$$f(\mathbf{x}; \mathbf{w}) := \text{sign}(\mathbf{w}^\top \mathbf{x}) \in \{\pm 1\}$$

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{w} = \mathbf{0}$.
2. For $t = 1, 2, \dots, T$
 - 2.1 If $y^{(t)} \mathbf{w}^\top \mathbf{x}^{(t)} > 0$, then $\mathbf{w} \leftarrow \mathbf{w}$,
 - 2.2 Else, then $\mathbf{w} \leftarrow \mathbf{w} + y^{(t)} \mathbf{x}^{(t)}$.

Output: \mathbf{w}



Perceptron: is this SGD?

$\mathbf{w} \in \mathbb{R}^d$ with classifier given by

$$f(\mathbf{x}; \mathbf{w}) := \text{sign}(\mathbf{w}^\top \mathbf{x}) \in \{\pm 1\}$$

$$L(f(\mathbf{x}; \theta), y^{(i)}) = \max(0, -y^{(i)} f(\mathbf{x}; \theta))$$

choice $z^{(i)}$

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{w} = \mathbf{0}$.
2. For $t = 1, 2, \dots, T$
 - 2.1 If $y^{(t)} \mathbf{w}^\top \mathbf{x}^{(t)} > 0$, then $\mathbf{w} \leftarrow \mathbf{w}$,
 - 2.2 Else, then $\mathbf{w} \leftarrow \mathbf{w} + y^{(t)} \mathbf{x}^{(t)}$.

Output: \mathbf{w}

- Compute gradient

$$\mathbf{g} \leftarrow \nabla_{\theta} \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}, \theta), y^{(i)})$$

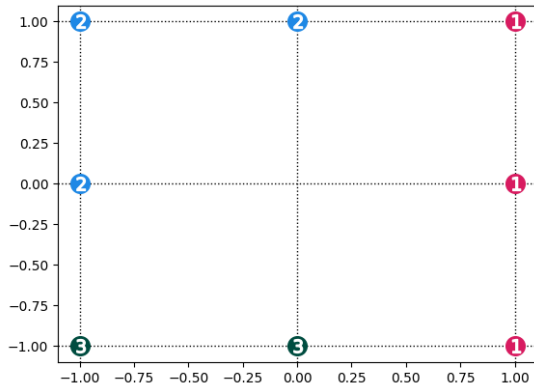
- Compute update $\theta \leftarrow \theta - \epsilon_k \mathbf{g}$

$$\mathbf{w} \leftarrow \mathbf{w} + y^{(t)} \mathbf{x}^{(t)}$$

Multiclass linear classifier (first attempt)

$\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] \in \mathbb{R}^{d \times K}$ with classifier given by

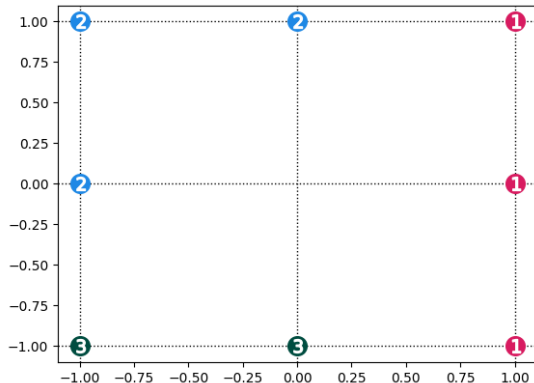
$$f(\mathbf{x}; \mathbf{W}) := \operatorname{argmax}_{\hat{y}=1,\dots,K} \mathbf{w}_{\hat{y}}^\top \mathbf{x} \in \{1, \dots, K\} \quad (2)$$



Multiclass linear classifier (second attempt)

$\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] \in \mathbb{R}^{d \times K}$ with classifier given by

$$f(\mathbf{x}; \mathbf{W}) := \operatorname{argmax}_{\hat{y}=1,\dots,K} \quad \mathbf{w}_{\hat{y}}^\top \mathbf{x} \in \{1, \dots, K\} \quad (3)$$



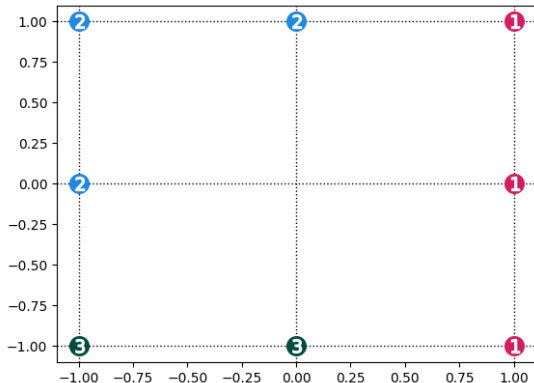
Multiclass perceptron

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

Output: \mathbf{W}



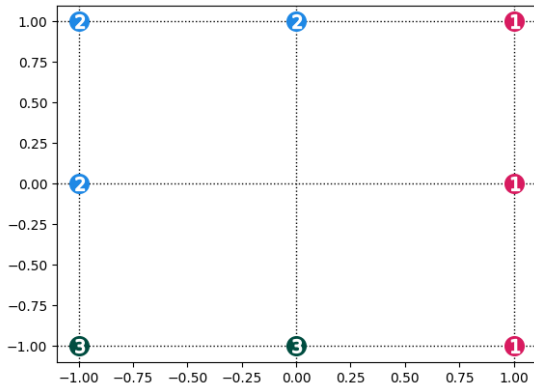
Multiclass perceptron

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.
2. For $t = 1, 2, \dots, T$

Output: \mathbf{W}



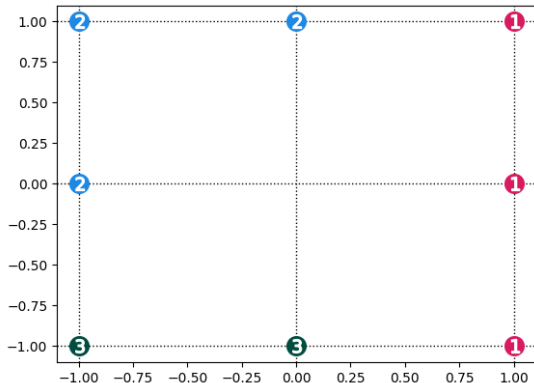
Multiclass perceptron

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.
2. For $t = 1, 2, \dots, T$
 - 2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

Output: \mathbf{W}



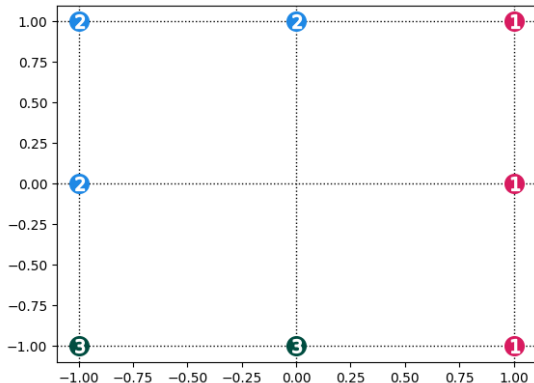
Multiclass perceptron

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.
2. For $t = 1, 2, \dots, T$
 - 2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$
 - 2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

Output: \mathbf{W}



Multiclass perceptron

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

2. For $t = 1, 2, \dots, T$

2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

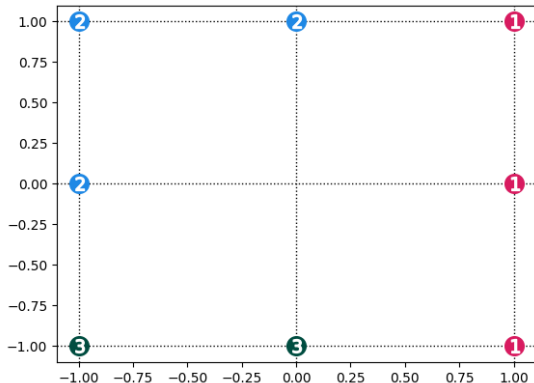
2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

2.3 Else, then

$$\mathbf{w}_{\hat{y}^{(t)}} \leftarrow \mathbf{w}_{\hat{y}^{(t)}} - \mathbf{x}^{(t)}$$

$$\mathbf{w}_{y^{(t)}} \leftarrow \mathbf{w}_{y^{(t)}} + \mathbf{x}^{(t)}$$

Output: \mathbf{W}



Multiclass perceptron (calculations)

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

2. For $t = 1, 2, \dots, T$

2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

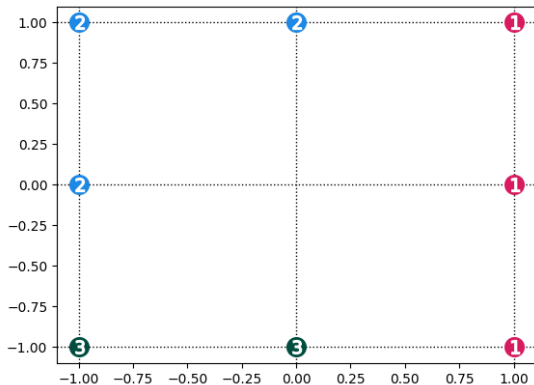
2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

2.3 Else, then

$$\mathbf{w}_{\hat{y}^{(t)}} \leftarrow \mathbf{w}_{\hat{y}^{(t)}} - \mathbf{x}^{(t)}$$

$$\mathbf{w}_{y^{(t)}} \leftarrow \mathbf{w}_{y^{(t)}} + \mathbf{x}^{(t)}$$

Output: \mathbf{W}



Multiclass perceptron (calculations)

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

2. For $t = 1, 2, \dots, T$

2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

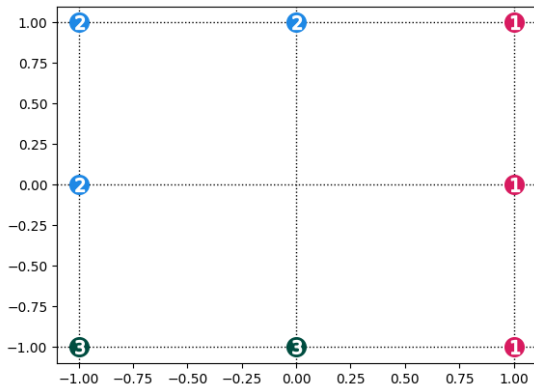
2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

2.3 Else, then

$$\mathbf{w}_{\hat{y}^{(t)}} \leftarrow \mathbf{w}_{\hat{y}^{(t)}} - \mathbf{x}^{(t)}$$

$$\mathbf{w}_{y^{(t)}} \leftarrow \mathbf{w}_{y^{(t)}} + \mathbf{x}^{(t)}$$

Output: \mathbf{W}



Multiclass perceptron (calculations)

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

2. For $t = 1, 2, \dots, T$

2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

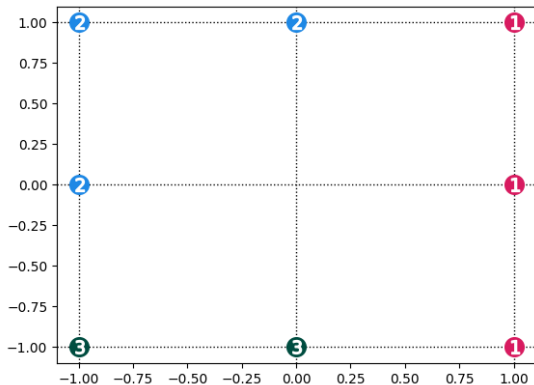
2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

2.3 Else, then

$$\mathbf{w}_{\hat{y}^{(t)}} \leftarrow \mathbf{w}_{\hat{y}^{(t)}} - \mathbf{x}^{(t)}$$

$$\mathbf{w}_{y^{(t)}} \leftarrow \mathbf{w}_{y^{(t)}} + \mathbf{x}^{(t)}$$

Output: \mathbf{W}



Multiclass perceptron (calculations)

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

2. For $t = 1, 2, \dots, T$

2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

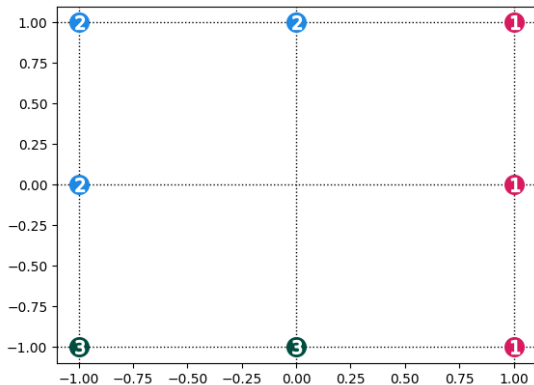
2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

2.3 Else, then

$$\mathbf{w}_{\hat{y}^{(t)}} \leftarrow \mathbf{w}_{\hat{y}^{(t)}} - \mathbf{x}^{(t)}$$

$$\mathbf{w}_{y^{(t)}} \leftarrow \mathbf{w}_{y^{(t)}} + \mathbf{x}^{(t)}$$

Output: \mathbf{W}



Multiclass perceptron (calculations)

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

2. For $t = 1, 2, \dots, T$

2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

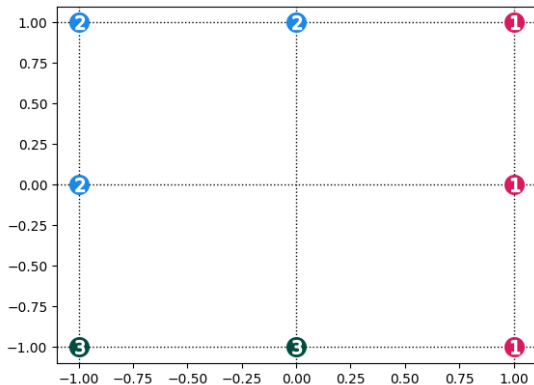
2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

2.3 Else, then

$$\mathbf{w}_{\hat{y}^{(t)}} \leftarrow \mathbf{w}_{\hat{y}^{(t)}} - \mathbf{x}^{(t)}$$

$$\mathbf{w}_{y^{(t)}} \leftarrow \mathbf{w}_{y^{(t)}} + \mathbf{x}^{(t)}$$

Output: \mathbf{W}



Multiclass perceptron (calculations)

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

2. For $t = 1, 2, \dots, T$

2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

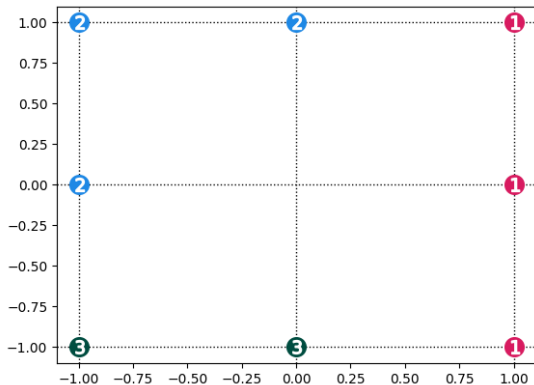
2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

2.3 Else, then

$$\mathbf{w}_{\hat{y}^{(t)}} \leftarrow \mathbf{w}_{\hat{y}^{(t)}} - \mathbf{x}^{(t)}$$

$$\mathbf{w}_{y^{(t)}} \leftarrow \mathbf{w}_{y^{(t)}} + \mathbf{x}^{(t)}$$

Output: \mathbf{W}



Multiclass perceptron (calculations)

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

2. For $t = 1, 2, \dots, T$

2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

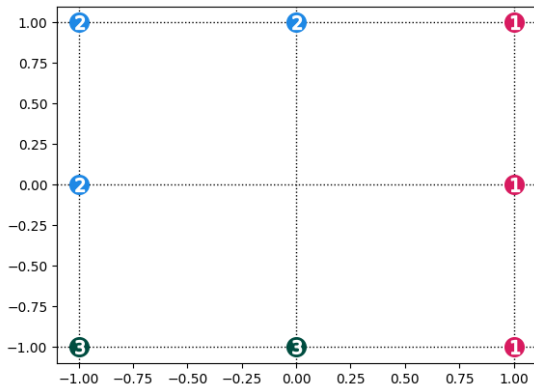
2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

2.3 Else, then

$$\mathbf{w}_{\hat{y}^{(t)}} \leftarrow \mathbf{w}_{\hat{y}^{(t)}} - \mathbf{x}^{(t)}$$

$$\mathbf{w}_{y^{(t)}} \leftarrow \mathbf{w}_{y^{(t)}} + \mathbf{x}^{(t)}$$

Output: \mathbf{W}



Multiclass perceptron (calculations)

Perceptron update

Input: $(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$, time T

1. Initialize $\mathbf{W} = [\mathbf{w}_1 \ \cdots \ \mathbf{w}_K] = \mathbf{0}$.

2. For $t = 1, 2, \dots, T$

2.1 $\hat{y}^{(t)} \leftarrow \operatorname{argmax}_{\hat{y}} \mathbf{w}_{\hat{y}}^\top \mathbf{x}^{(t)}$

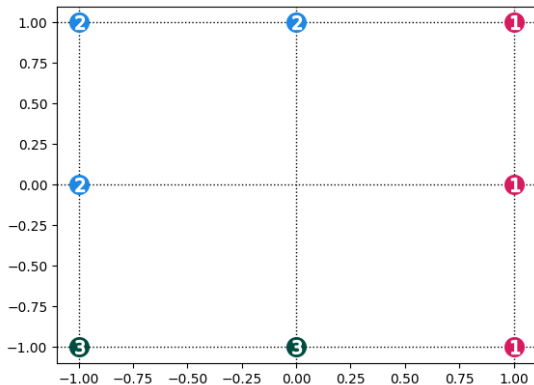
2.2 If $\hat{y}^{(t)} = y^{(t)}$, then $\mathbf{W} \leftarrow \mathbf{W}$,

2.3 Else, then

$$\mathbf{w}_{\hat{y}^{(t)}} \leftarrow \mathbf{w}_{\hat{y}^{(t)}} - \mathbf{x}^{(t)}$$

$$\mathbf{w}_{y^{(t)}} \leftarrow \mathbf{w}_{y^{(t)}} + \mathbf{x}^{(t)}$$

Output: \mathbf{W}



References I

- [GBC16] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT press, 2016.