# Lecture 10. The Perceptron

COMP90051 Statistical Machine Learning

Lecturer:  Jean Honorio

THE UNIVERSITY OF
MELBOURNE

# This lecture

- ## Perceptron model

  - ∗ Introduction to Artificial Neural Networks

  - ∗ The perceptron model

- ## Perceptron training rule

  - ∗ Stochastic gradient descent

- ## Kernel perceptron

# The Perceptron Model

A building block for artificial neural networks, yet another linear classifier
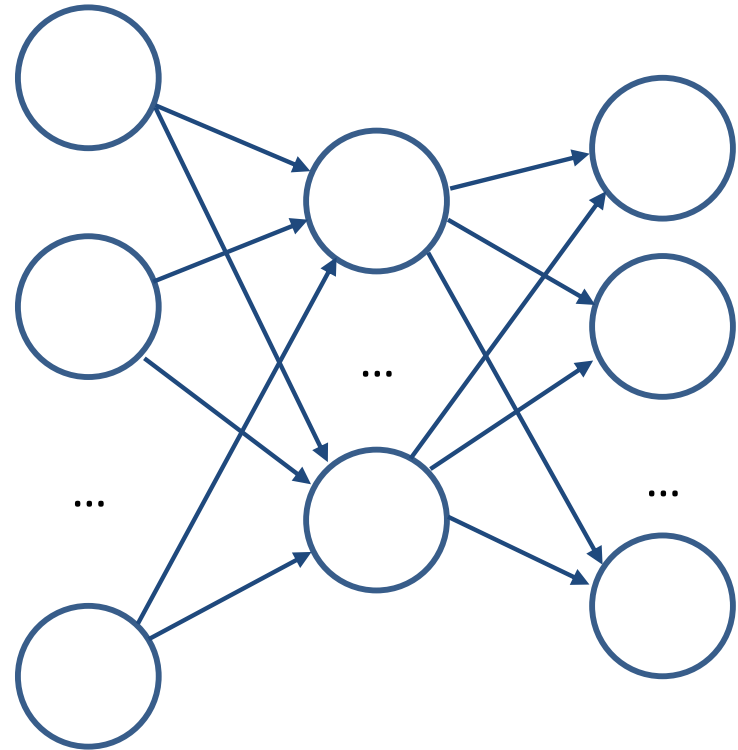
# Biological inspiration

- Humans perform well at many tasks that matter

- Originally neural networks were an attempt to mimic the human brain

photo: Alvesgaspar, Wikimedia Commons, CC3
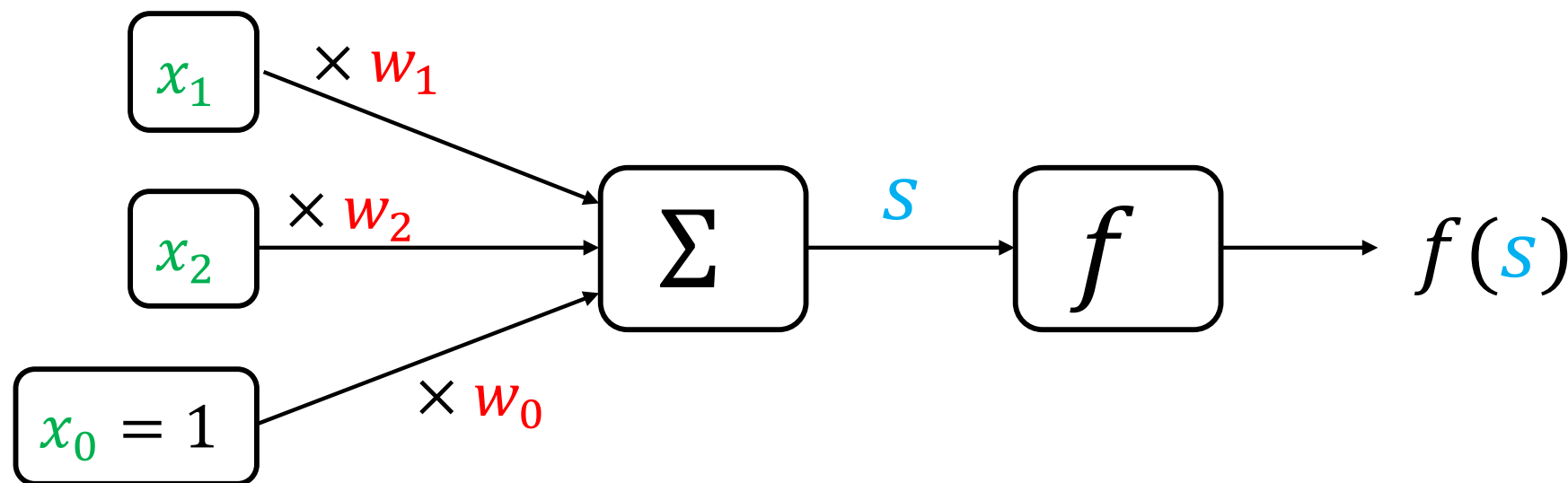
# Artificial neural network

- As a *crude approximation,* the human brain can be thought as a mesh of interconnected processing nodes (neurons) that relay electrical signals

- Artificial neural network is a network of processing elements

- Each element converts inputs to output

- The output is a function (called *activation function*) of a weighted sum of inputs

# Outline

- In order to use an ANN we need (a) to design network topology and (b) adjust weights to given data
  - * In this subject, we will exclusively focus on task (b) for a particular class of networks called feed forward networks

- Training an ANN means adjusting weights for training data given a pre-defined network topology

- First we will turn our attention to an individual network element, before building deeper architectures

# Perceptron model



$x_1 \quad \times w_1$

$x_2 \quad \times w_2$

$x_0 = 1 \quad \times w_0$

$\Sigma \quad s \quad f \quad f(s)$

Compare this model to logistic regression

- $x_1, x_2$ – features/inputs
- $w_1, w_2$ – synaptic weights
- $w_0$ – bias weight
- $f$ – activation function
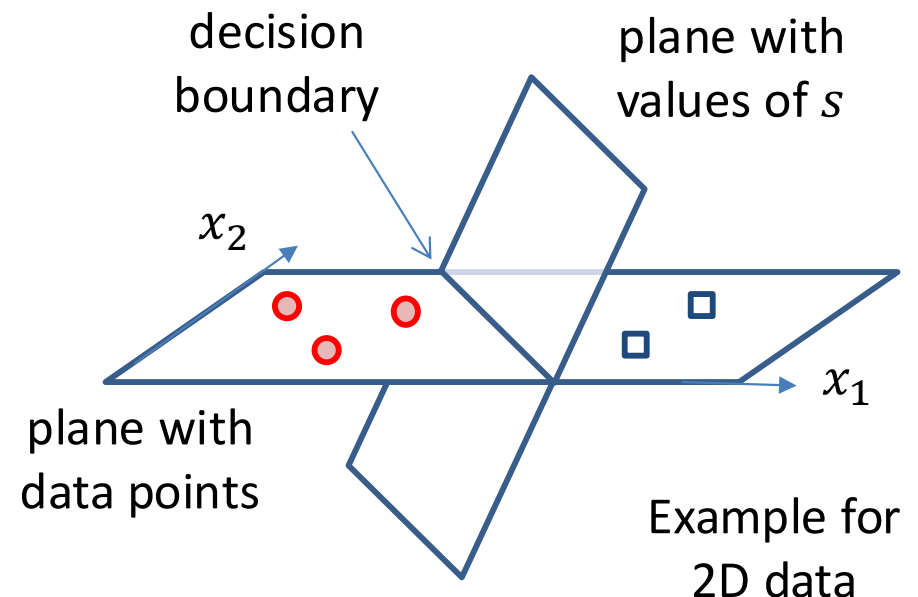
# Perceptron is a linear binary classifier

Perceptron is a binary classifier:

Predict class A if $s \geq 0$
Predict class B if $s < 0$
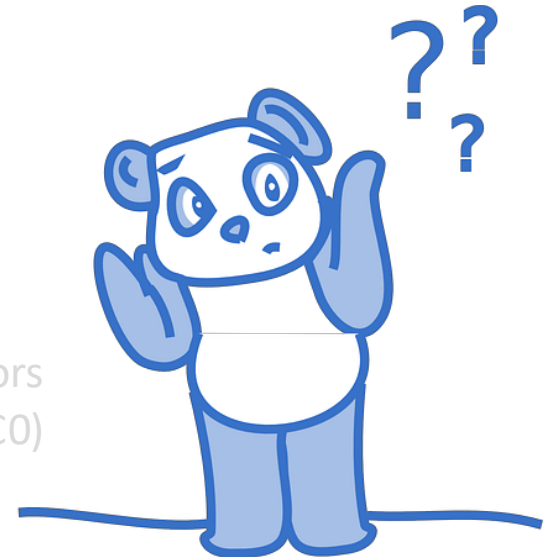where $s = \boldsymbol{w}'\boldsymbol{x} = \sum_{j=0}^{m} w_j x_j$

Perceptron is a <u>linear classifier</u>: $s$ is a linear function of inputs, and the decision boundary is linear

decision boundary

plane with values of $s$

$x_2$

$x_1$

plane with data points

Example for 2D data

Exercise: find weights of a perceptron capable of perfect classification of the following dataset

| $x_1$ | $x_2$ | $y$ |
|-------|-------|---------|
| 0 | 0 | Class B |
| 0 | 1 | Class B |
| 1 | 0 | Class B |
| 1 | 1 | Class A |

Exercise: find weights of a perceptron capable of perfect classification of the following dataset

| $x_1$ | $x_2$ | $y$ | $s$ |
|---|---|---|---|
| 0 | 0 | Class B | -1.5 |
| 0 | 1 | Class B | -0.5 |
| 1 | 0 | Class B | -0.5 |
| 1 | 1 | Class A | 0.5 |

$$w_1 = 1, w_2 = 1, w_0 = -1.5$$
$$s = \boldsymbol{w'x} = w_1 x_1 + w_2 x_2 + w_0$$
$$s = x_1 + x_2 - 1.5$$

# Mini Summary

- Perceptron

  * Introduction to Artificial Neural Networks

  * The perceptron model

Next: Perceptron training

# Perceptron Training Rule

Gateway to stochastic gradient descent.

Convergence guaranteed by convexity.

# Loss function for perceptron

- "Training": finds weights to minimise some loss. Which?

- Our task is binary classification. Encode one class as $+1$ and the other as $-1$. So each training example is now $(\boldsymbol{x}_i, y_i)$, where $y_i$ is either $+1$ or $-1$

- Recall that, in a perceptron, $s_i = \boldsymbol{w}'\boldsymbol{x}_i = \sum_{j=0}^{m} w_j x_{ij}$, and the sign of $s_i$ determines the predicted class: $+1$ if $s_i > 0$, and $-1$ if $s_i < 0$

- Consider a single training example.
  * If $y_i$ and $s_i$ have same sign then the example is classified correctly.
  * If $y_i$ and $s_i$ have different signs, the example is misclassified

# Loss function for perceptron

- The perceptron uses a loss function in which there is no penalty for correctly classified examples, while the penalty (loss) is equal to $s_i$ for misclassified examples*
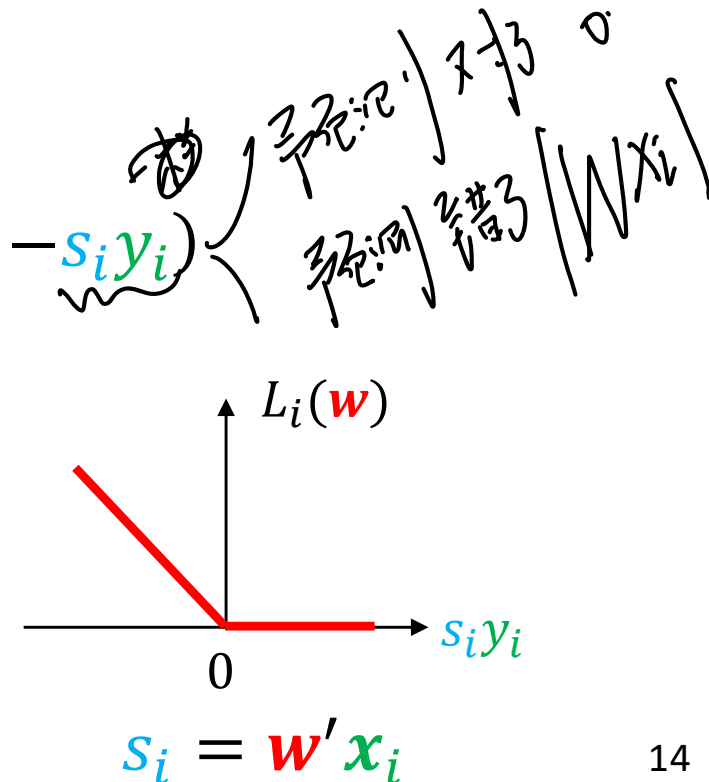
- Formally:
  * $L_i(\boldsymbol{w}) = 0$ if both $s_i, y_i$ have the same sign
  * $L_i(\boldsymbol{w}) = |s_i|$ if both $s_i, y_i$ have different signs

- This can be re-written as $L_i(\boldsymbol{w}) = \max(0, -s_i y_i)$

*该要看 loss的 gradient ★*

* This is similar, but not identical to the SVM's **hinge** loss

$s_i = \boldsymbol{w}'\boldsymbol{x}_i$

14

# Stochastic gradient descent

- Randomly shuffle/split all training examples in $B$ batches

- Choose initial $\boldsymbol{\theta}^{(1)}$

- For $t$ from $1$ to $T$

> Iterations over the entire dataset are called *epochs*

-     For $b$ from $1$ to $B$

-         Do gradient descent update <u>using data from batch $b$</u>

- Advantage of such an approach: computational feasibility for large datasets

true gradient: average over all training samples' gradients (标准)

# Perceptron training algorithm

Choose initial guess $\boldsymbol{w}^{(0)}$, $k = 0$

*(handwritten: Value of weight after k changes)*

For $t$ from 1 to $T$ (epochs)

*(handwritten: 内部 sample 一个个过)*

    For $i$ from 1 to $N$ (training examples)

> Consider example $(\boldsymbol{x}_i, y_i)$
>
> Update*: $\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} - \eta \nabla L_i\big(\boldsymbol{w}^{(k)}\big)$
>
> $k = k + 1$

*(handwritten: learning rate, gradient of loss)*

$$L_i(\boldsymbol{w}) = \max(0, -s_i y_i)$$
$$s_i = \boldsymbol{w}' \boldsymbol{x}_i$$
$\eta$ is learning rate

*(handwritten: 接近90 对应0 一阶导 不连续)*

*There is no derivative when $s_i = 0$, but this case is handled explicitly in the algorithm, see next slides

16

# Perceptron training rule

- We have $\nabla L_i(\boldsymbol{w}) = \boldsymbol{0}$ when $s_i y_i > 0$
  - ✳ We don't need to do update when sample $i$ is correctly classified 预测对了 ⇒ gradient=0

- What is $\nabla L_i(\boldsymbol{w})$ when $s_i y_i < 0$? 预测错了 ⇒ 我挑选 ⇒ $\nabla L_i(w) = -y_i x_i$
  - ✳ We need to update when sample $i$ is misclassified
  - ✳ We have $\nabla L_i(\boldsymbol{w}) = \boldsymbol{x}_i$ when $y_i = -1$ and $s_i > 0$
  - ✳ We have $\nabla L_i(\boldsymbol{w}) = -\boldsymbol{x}_i$ when $y_i = 1$ and $s_i < 0$
  - ✳ Thus $\nabla L_i(\boldsymbol{w}) = -y_i \boldsymbol{x}_i$

  $(-s_i y_i)$

- $L_i(\boldsymbol{w}) = \max(0, -s_i y_i)$

$s_i = \boldsymbol{w}' \boldsymbol{x}_i$

17

# Perceptron training algorithm

Choose initial guess $\boldsymbol{w}^{(0)}$, $k = 0$

For $t$ from 1 to $T$ (epochs)  *go through data set serval times.*

    For $i$ from 1 to $N$ (training examples)  *go through one data point at a time.*

*computed $s_i$*

Compute $s_i = \left(\boldsymbol{w}^{(k)}\right)' \boldsymbol{x}_i$

If $s_i y_i \leq 0$: (sample $i$ misclassified)

$$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} + \eta y_i \boldsymbol{x}_i$$

$k = k + 1$  $-(\eta \nabla L_i(\boldsymbol{w}))$

$-\eta \nabla L_i(\boldsymbol{w})$

$(y_i x_i)$

判达值 错没错

$\nabla L_i(\boldsymbol{w}) = -y_i x_i$

($\eta > 0$ is called *learning rate*)

# Perceptron training algorithm

Choose initial guess $\boldsymbol{w}^{(0)}$, $k = 0$

For $t$ from 1 to $T$ (epochs)

　　For $i$ from 1 to $N$ (training examples)

　　　Compute $s_i = \left(\boldsymbol{w}^{(k)}\right)' \boldsymbol{x}_i$

　　　If $s_i y_i \leq 0$: (sample $i$ misclassified)

　　　　$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} + \eta y_i \boldsymbol{x}_i$

　　　　$k = k + 1$

Strictly speaking it should be $s_i y_i < 0$ but $\leq$ allows handling the case $\boldsymbol{w}^{(k)} = \boldsymbol{0}$

$\boldsymbol{w}^{(k)}$ represents the value of $\boldsymbol{w}$ after $k$ updates (useful for theory). If you implement this, just write: $\boldsymbol{w} = \boldsymbol{w} + \eta y_i \boldsymbol{x}_i$

在 initialization 阶段 $w^0$ 可能从0开始

这也是 什么时候为非0

习惯上从第一次移向做渐逼过程中 $s_i y_i < 0$

倒应信息博将值

若 信息 $s_i y_i < 0 \Rightarrow W$ 需要被更新

初始值

19

# Example 1: Perceptron learning

remove $w_0$, set $\eta = 1$     Basic setup   $w = w + \eta \, y_i x_i$

$x_1$ $\quad w_1$

$x_2$ $\quad w_2$

$$s = w_1 x_1 + w_2 x_2$$

$\Sigma$ $\quad f \quad \rightarrow \begin{cases} -1, & s < 0 \\ 1, & s \geq 0 \end{cases}$

(learning rate $\eta = 1$)

* We drop the sample index $i$ to have a simpler notation.

# Example 1: Perceptron learning

Start with zero weights, Consider sample 1

$w_1 = 0$

$w_2 = 0$

initialize $w$ at zero:

$S_1 = 0$

起始点     $S_1 y_1 = 0 \Rightarrow$ 胡合为 $\Rightarrow$ 更新

$x_2$

□ sample 1

| ● class -1 |
| □ class 1 |

$w = 0$

$x_1$

# Example 1: Perceptron learning

## Update



class -1
class 1

sample 1

$+1$

$\boldsymbol{w}$

$\boldsymbol{x}$

$-1$

$s_i y_i \leq 0$

Since $\boldsymbol{w} = \boldsymbol{0}$, then update (why?)

$$\boldsymbol{w} = \boldsymbol{w} + y\boldsymbol{x} = \boldsymbol{x}$$

$0 + 1 \cdot x = x$

* We drop the sample index $i$ to have a simpler notation.

# Example 1: Perceptron learning

## Consider sample 2



-1 乘し 様S，等出转隆，预测错3)

$S或Y\xi\leq0$ ⇒要update

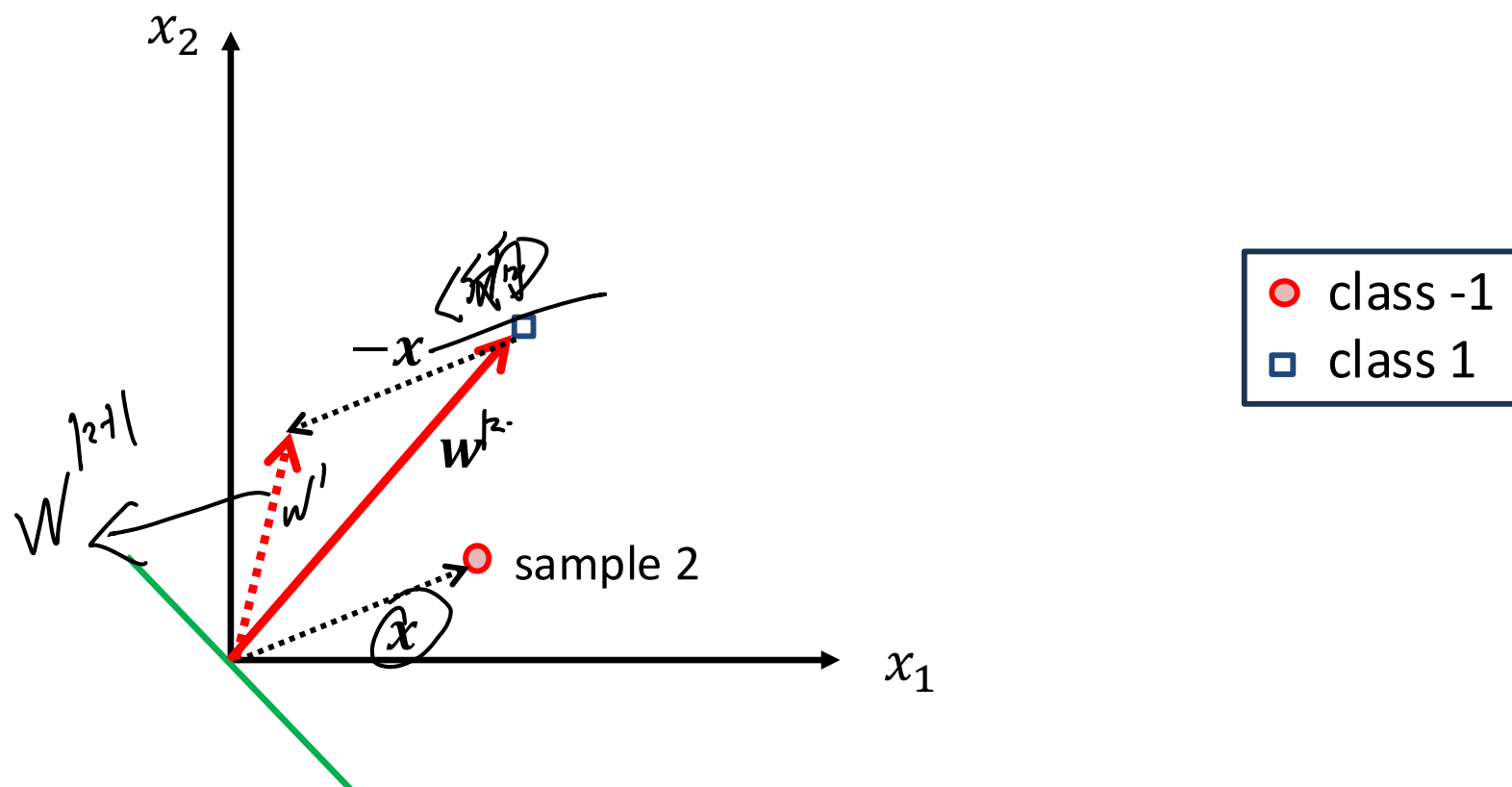basically⇒领工在绿线场的样样座
是负的都会预测错 ⇒ $S=+1$ ⇒ 都要动
W.

class -1
class 1

Sample 2 is in the wrong side of the decision boundary

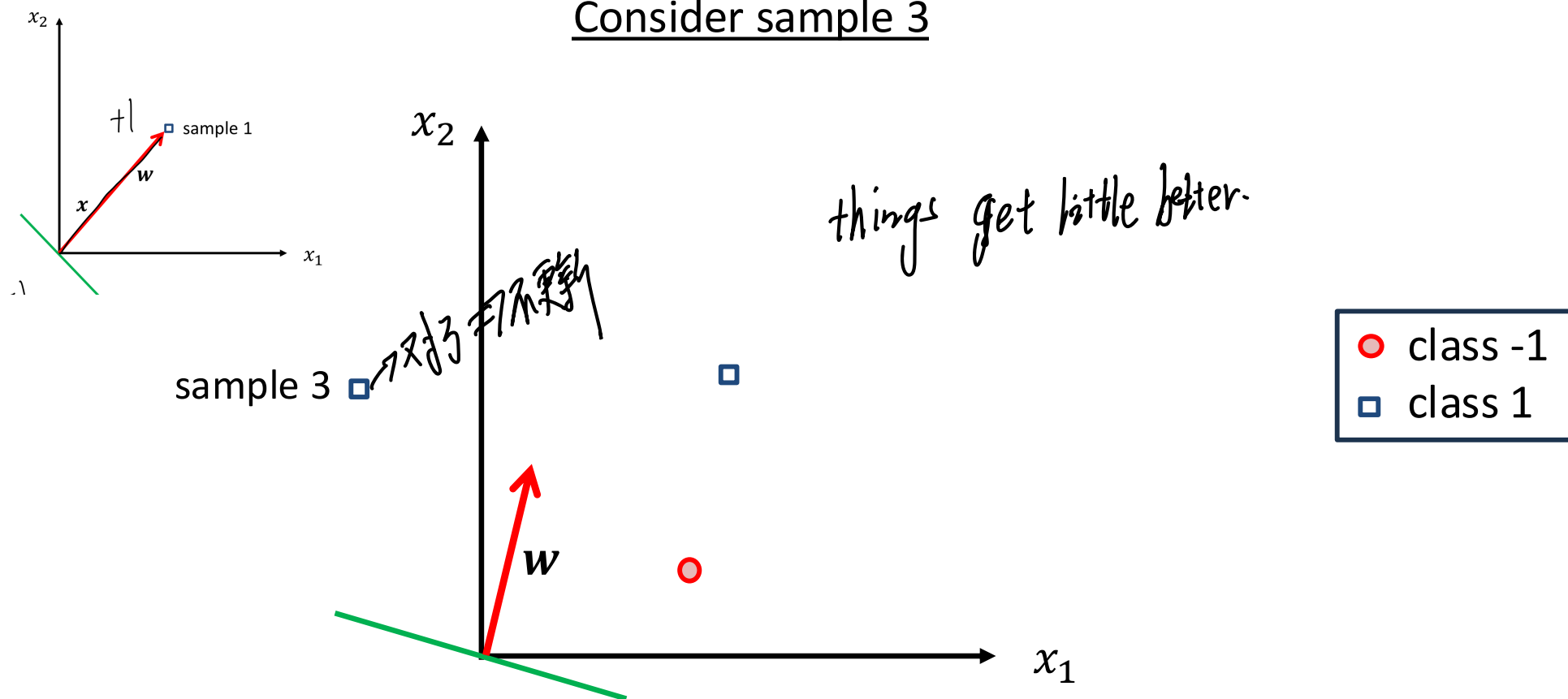# Example 1: Perceptron learning

## Update



Sample 2 is in the wrong side of the decision boundary, then update

$$w = w + yx = w - x$$

* We drop the sample index $i$ to have a simpler notation.
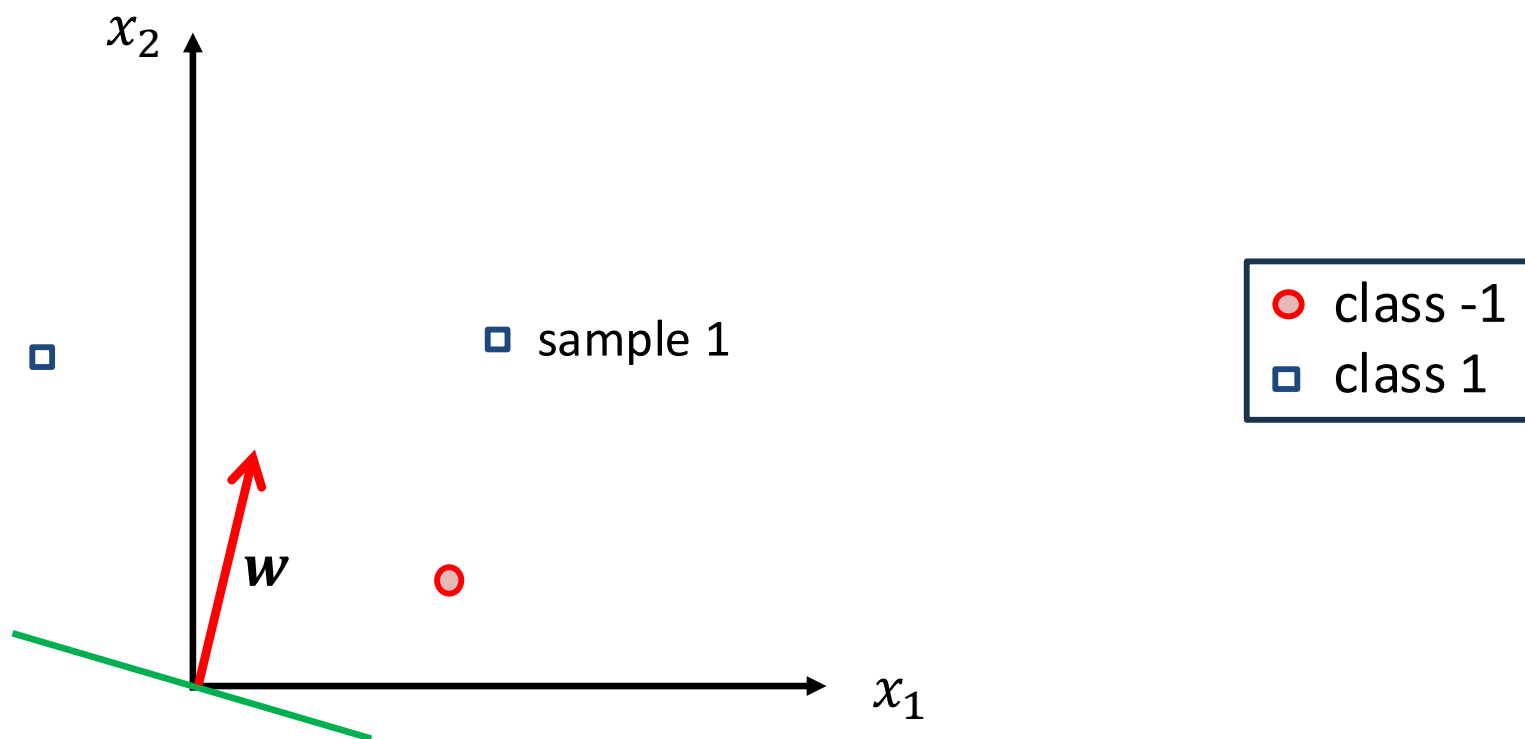
# Example 1: Perceptron learning

## Consider sample 3



things get little better.

| | |
|---|---|
| ⊙ | class -1 |
| □ | class 1 |

Sample 3 is in the correct side of the decision boundary, no update

# Example 1: Perceptron learning

$\frac{x_2^k}{\rho} 2\,batch$ $\overparen{\text{Consider sample 1 again}}$



Sample 1 is in the correct side of the decision boundary, no update

* We drop the sample index $i$ to have a simpler notation.

# Example 1: Perceptron learning

Consider sample 2 again



再更新

class -1
class 1

$x_2$

$w$

$\circ$ sample 2

$x_1$

Sample 2 is in the wrong side of the decision boundary, then update

* We drop the sample index $i$ to have a simpler notation.

# Example 1: Perceptron learning

## Update



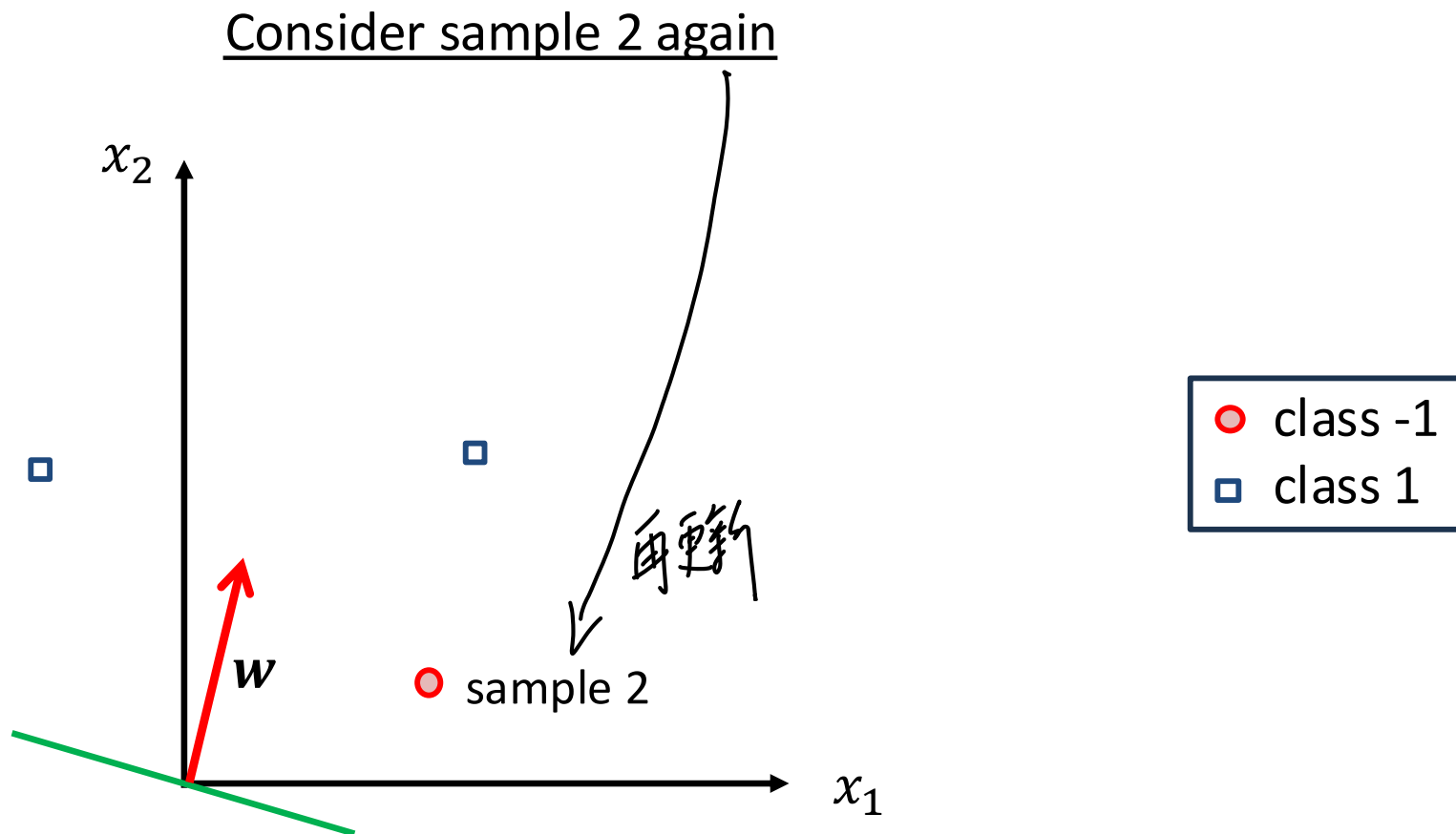Sample 2 is in the wrong side of the decision boundary, then update

$$w = w + yx = w - x$$

$$y = -1$$

* We drop the sample index $i$ to have a simpler notation.

# Example 1: Perceptron learning

## And we continue...

# Perceptron training algorithm

Choose initial guess $\boldsymbol{w}^{(0)}$, $k = 0$

For $t$ from 1 to $T$ (epochs)

把所有点放了遍 最后就停 ⇒ 更新:k 知道

For $i$ from 1 to $N$ (training examples)

Compute $s_i = \left(\boldsymbol{w}^{(k)}\right)' \boldsymbol{x}_i$

If $s_i y_i \leq 0$: (sample $i$ misclassified)

$$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} + \eta y_i \boldsymbol{x}_i$$

$k = k + 1$

*Convergence Theorem*: if the training data is linearly separable, the algorithm is guaranteed to converge to a solution. That is, there exist a finite $k$ such that $L\left(\boldsymbol{w}^{(k)}\right) = 0$

# Perceptron convergence theorem

- Assumptions

  * Linear separability: There exists $\boldsymbol{w}^*$ so that $\frac{y_i(\boldsymbol{w}^*)'\boldsymbol{x}_i}{\|\boldsymbol{w}^*\|} \geq \gamma$ for all training data $i = 1, \dots, N$ and some positive $\gamma$.
  * Bounded data: $\|\boldsymbol{x}_i\| \leq R$ for $i = 1, \dots, N$ and some finite $R$.
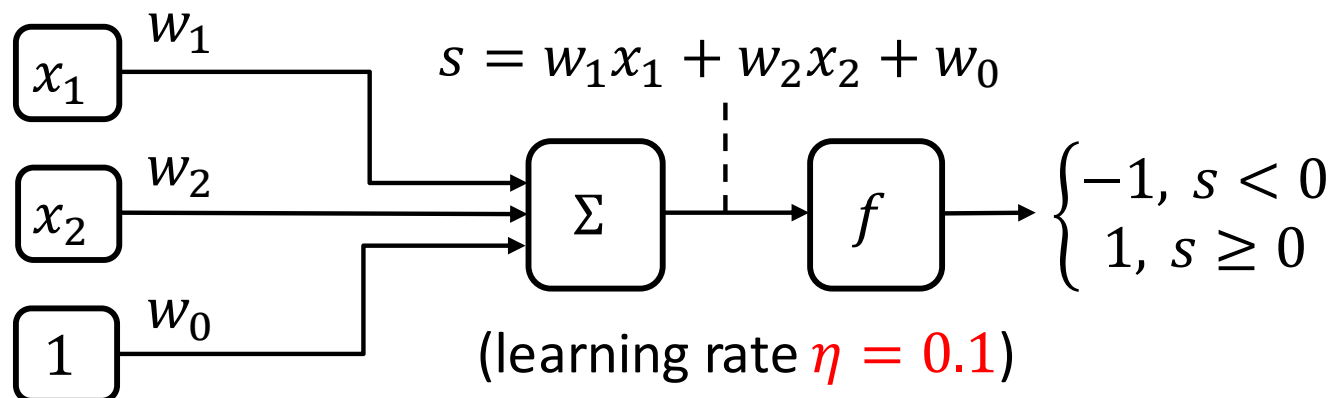
- Proof sketch (for $\eta = 1$)

  * Assumes $\boldsymbol{w}^{(0)} = \boldsymbol{0}$ to
    - Establish that $(\boldsymbol{w}^*)'\boldsymbol{w}^{(k)} \geq k\gamma\|\boldsymbol{w}^*\|$
    - Establish that $\left\|\boldsymbol{w}^{(k)}\right\|^2 \leq kR^2$

  * Note that $1 \geq \cos\left(\boldsymbol{w}^*, \boldsymbol{w}^{(k)}\right) = \frac{(\boldsymbol{w}^*)'\boldsymbol{w}^{(k)}}{\|\boldsymbol{w}^*\|\left\|\boldsymbol{w}^{(k)}\right\|} \geq \frac{k\gamma\|\boldsymbol{w}^*\|}{\|\boldsymbol{w}^*\|\sqrt{k}R}$

  * Take the left-most and right-most equations $1 \geq \frac{k\gamma\|\boldsymbol{w}^*\|}{\|\boldsymbol{w}^*\|\sqrt{k}R}$

  * Rearranging we get $k \leq \frac{R^2}{\gamma^2}$

# Pros and cons of perceptron learning

- If the data is linearly separable, the perceptron training algorithm will converge to a correct solution

  * There is a formal proof ← good!

  * It will converge to some solution (separating boundary), one of infinitely many possible ← bad!

- However, if the data is not linearly separable, the training will fail completely rather than give some approximate solution
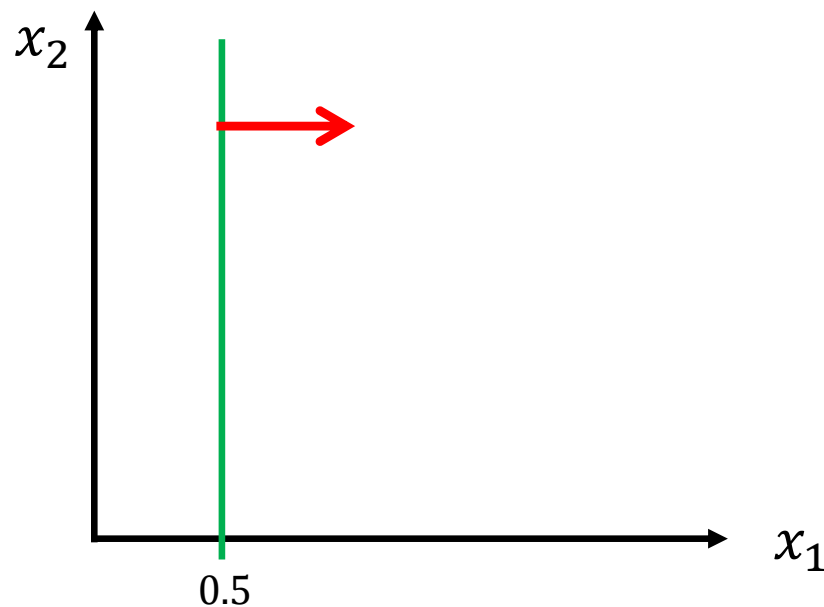
  * Ugly ☹

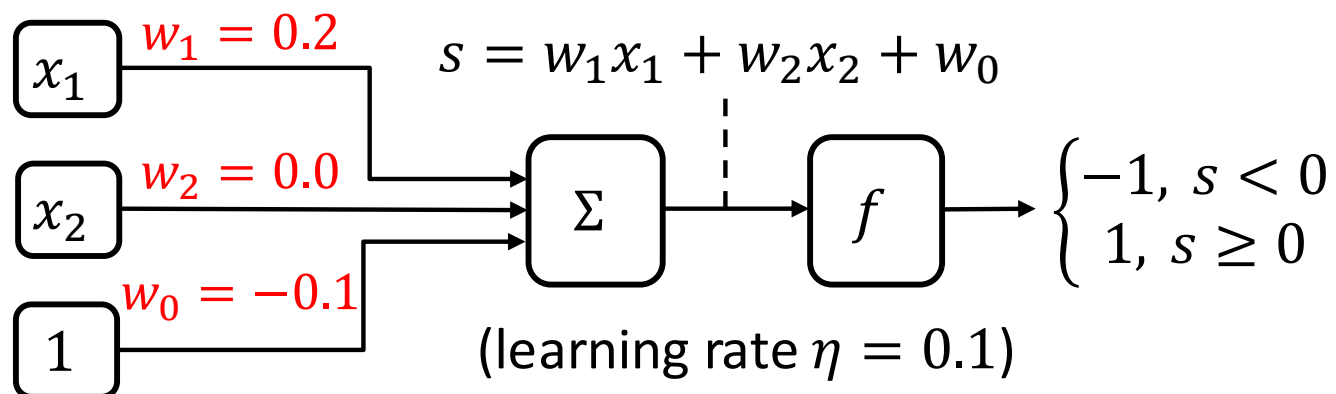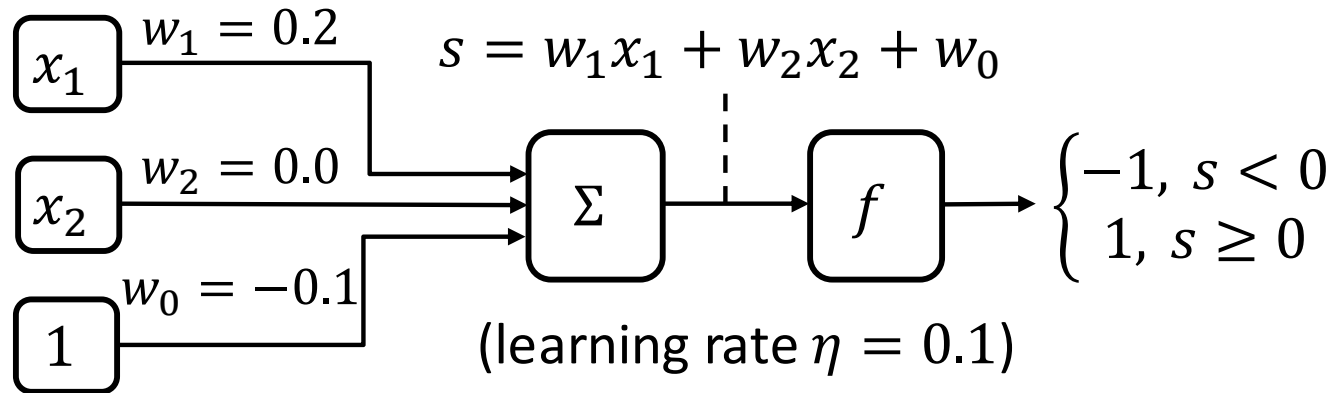# Example 2: Perceptron learning

## Basic setup



$$s = w_1 x_1 + w_2 x_2 + w_0$$

$$\begin{cases} -1, & s < 0 \\ 1, & s \geq 0 \end{cases}$$

(learning rate $\eta = 0.1$)

* We drop the sample index $i$ to have a simpler notation.

# Example 2: Perceptron learning

## Start with random weights



$w_1 = 0.2$

$s = w_1 x_1 + w_2 x_2 + w_0$

$x_1$

$w_2 = 0.0$

$x_2$

$\Sigma$    $f$    $\begin{cases} -1, & s < 0 \\ 1, & s \geq 0 \end{cases}$

$w_0 = -0.1$

$1$

(learning rate $\eta = 0.1$)

$x_2$

$x_1$

$0.5$

\* We drop the sample index $i$ to have a simpler notation.

# Example 2: Perceptron learning

### Consider training example 1



$$s = w_1 x_1 + w_2 x_2 + w_0$$

$$f \rightarrow \begin{cases} -1, & s < 0 \\ 1, & s \geq 0 \end{cases}$$

$w_1 = 0.2$

$w_2 = 0.0$
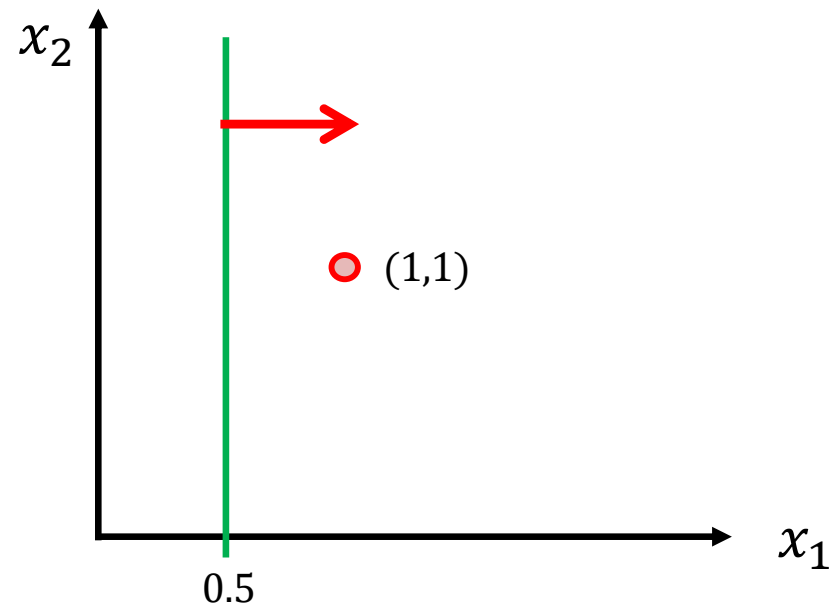
$w_0 = -0.1$

(learning rate $\eta = 0.1$)

○ class -1
□ class 1

$$y(0.2x_1 + 0.0x_2 - 0.1) = -0.1 \leq 0$$

$$w_1 \leftarrow w_1 - \eta x_1 = \quad 0.1$$
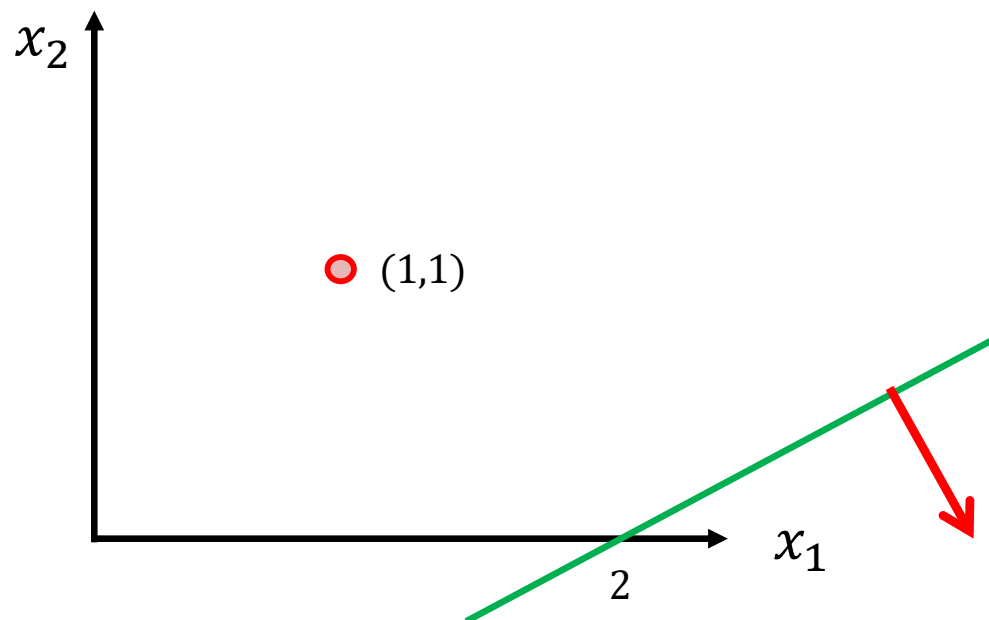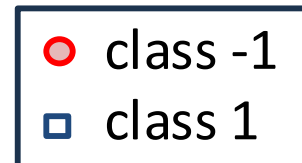$$w_2 \leftarrow w_2 - \eta x_2 = -0.1$$
$$w_0 \leftarrow w_0 - \eta \quad\; = -0.2$$

* We drop the sample index $i$ to have a simpler notation.

# Example 2: Perceptron learning

## Update weights

$x_1$    $w_1 = 0.1$

$x_2$    $w_2 = -0.1$

$1$    $w_0 = -0.2$

$$s = w_1 x_1 + w_2 x_2 + w_0$$

$$\Sigma \longrightarrow f \longrightarrow \begin{cases} -1, & s < 0 \\ 1, & s \geq 0 \end{cases}$$

(learning rate $\eta = 0.1$)

○ class -1
□ class 1

$x_2$

○ (1,1)

2

$x_1$

# Example 2: Perceptron learning

### Consider training example 2

$$w_1 = 0.1$$
$x_1$

$$s = w_1 x_1 + w_2 x_2 + w_0$$

$$w_2 = -0.1$$
$x_2$

$\Sigma$    $f$    $\begin{cases} -1, & s < 0 \\ 1, & s \geq 0 \end{cases}$

$$w_0 = -0.2$$
1

(learning rate $\eta = 0.1$)

○ class -1
□ class 1

$$y(0.1x_1 - 0.1x_2 - 0.2) = -0.1 \leq 0$$

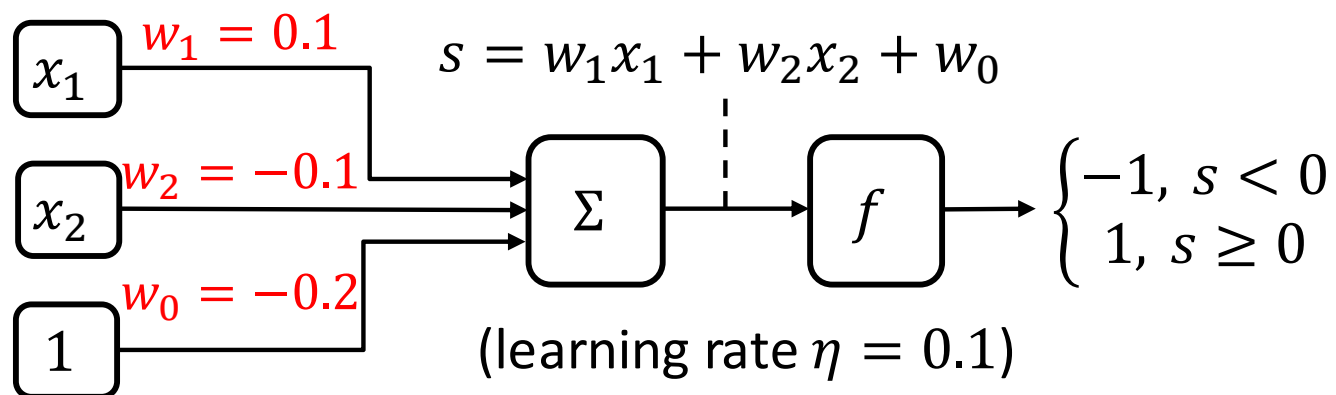$$w_1 \leftarrow w_1 + \eta x_1 = \quad 0.3$$
$$w_2 \leftarrow w_2 + \eta x_2 = \quad 0.0$$
$$w_0 \leftarrow w_0 + \eta \quad\quad = -0.1$$

$x_2$

○    □ (2,1)

2    $x_1$

* We drop the sample index $i$ to have a simpler notation.

# Example 2: Perceptron learning

## Update weights
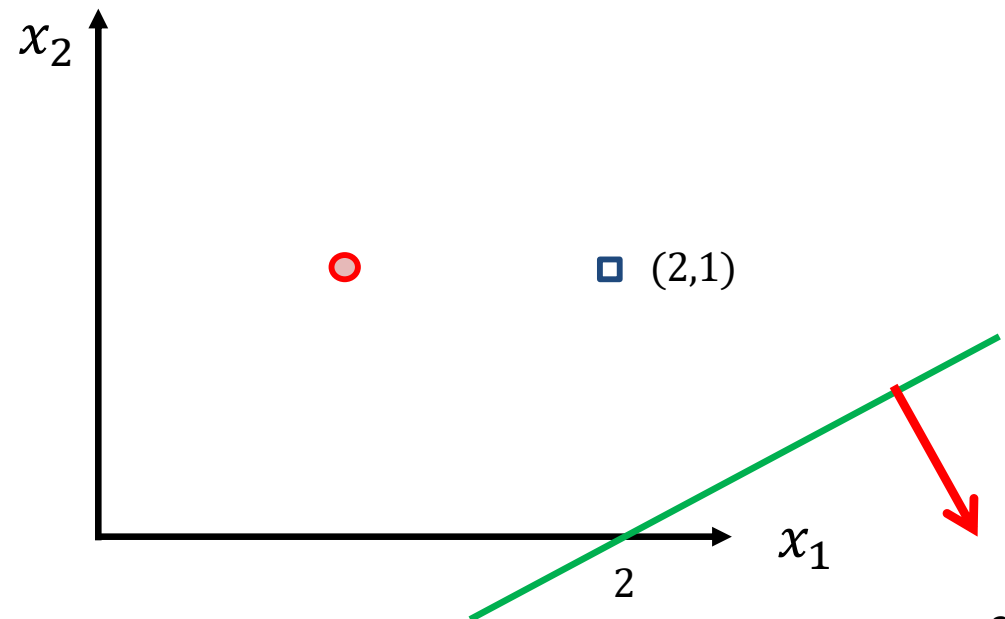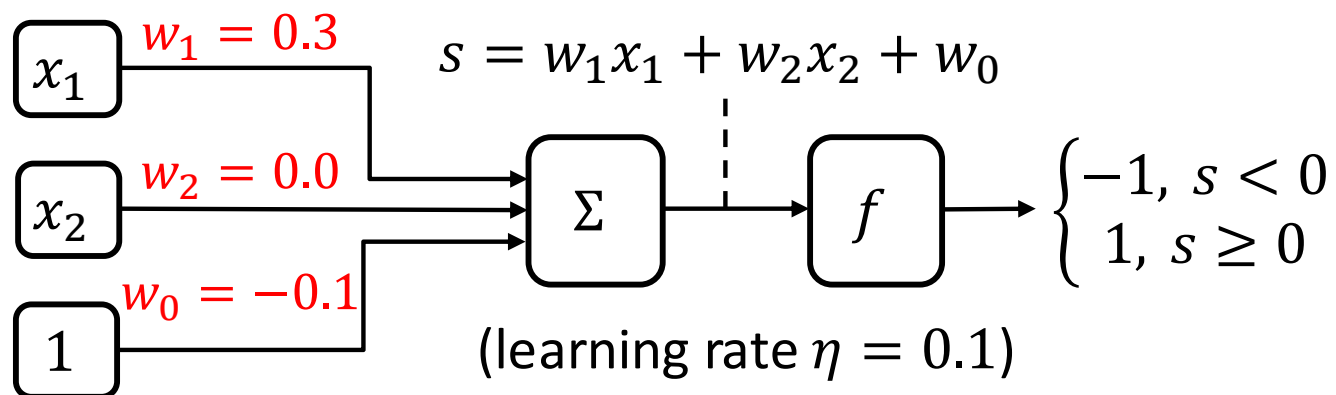


$$w_1 = 0.3$$

$$s = w_1 x_1 + w_2 x_2 + w_0$$

$$w_2 = 0.0$$

$$w_0 = -0.1$$

$$\begin{cases} -1, & s < 0 \\ 1, & s \geq 0 \end{cases}$$

(learning rate $\eta = 0.1$)

○ class -1
□ class 1

$(2,1)$

$1/3$

* We drop the sample index $i$ to have a simpler notation.

# Example 2: Perceptron learning

## Further examples

$x_1$   $w_1 = 0.3$

$x_2$   $w_2 = 0.0$

$1$   $w_0 = -0.1$

$s = w_1 x_1 + w_2 x_2 + w_0$

$\Sigma$   $f$   $\begin{cases} -1, & s < 0 \\ 1, & s \geq 0 \end{cases}$

(learning rate $\eta = 0.1$)

○ class -1
□ class 1

$y(0.3x_1 - 0.0x_2 - 0.1) = 0.35 > 0$
3rd point: correctly classified

4th point: incorrect, update
etc.

$x_2$

4th point

$x_1$

1/3

□ (1.5,0.5)
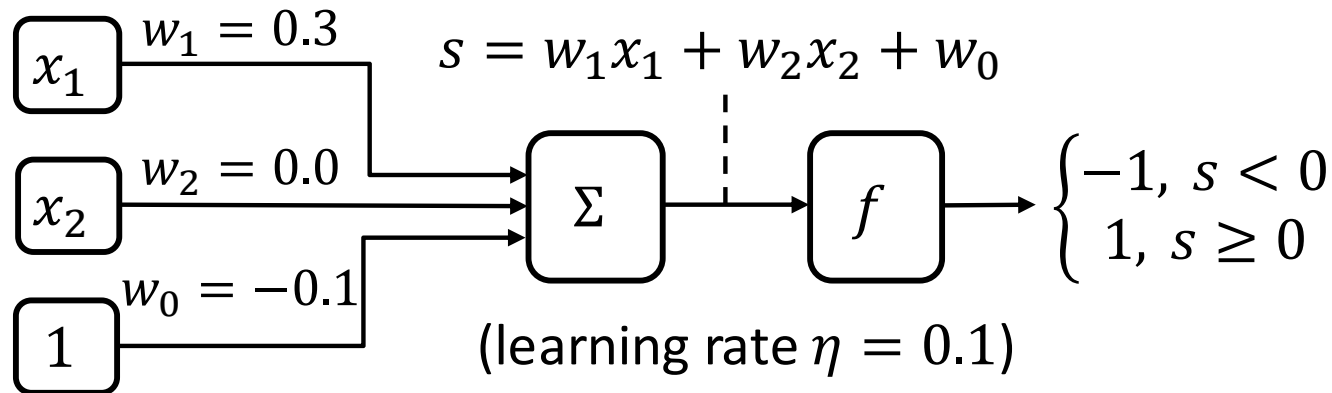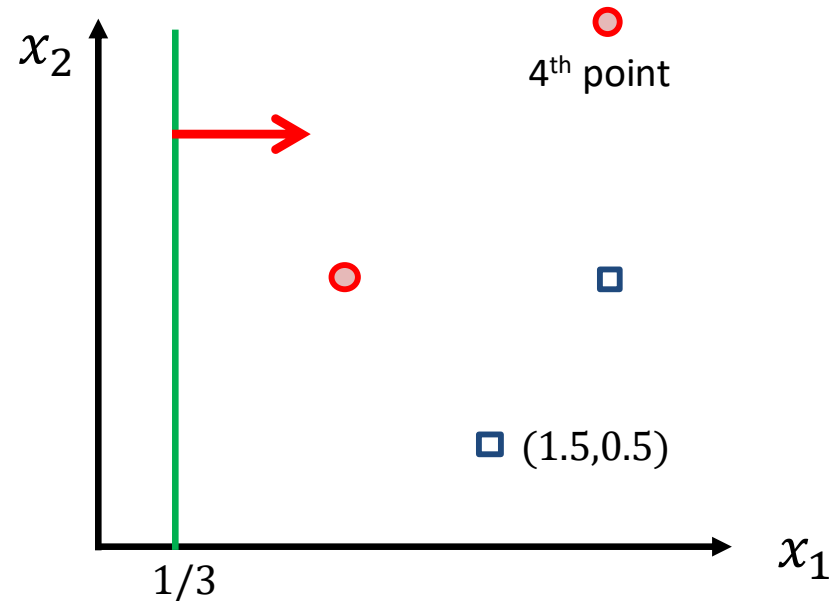
* We drop the sample index $i$ to have a simpler notation.

39

# Example 2: Perceptron learning

## Further examples

$x_1$   $w_1 = \cdots$

$s = w_1 x_1 + w_2 x_2 + w_0$

$x_2$   $w_2 = \cdots$   $\Sigma$   $f$   $\begin{cases} -1, & s < 0 \\ 1, & s \geq 0 \end{cases}$

$1$   $w_0 = \cdots$

(learning rate $\eta = 0.1$)

○ class -1
□ class 1

Eventually, all the data will be correctly classified (provided it is linearly separable)

$x_2$

$x_1$

* We drop the sample index $i$ to have a simpler notation.

40

# Mini Summary

- Perceptron loss function

- Stochastic gradient descent

- Perceptron training rule
    * Perceptron convergence theorem

Next: Kernel perceptron

# Kernel Perceptron

Another example of a kernelizable learning algorithm (like the SVM).

# Perceptron training rule: Recap

Compute $s_i = \left(\boldsymbol{w}^{(k)}\right)' \boldsymbol{x}_i$

If $s_i y_i \le 0$: (sample $i$ misclassified)

$$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} + \eta y_i \boldsymbol{x}_i$$

$$k = k + 1$$

Suppose weights are initially set to $\boldsymbol{w}^{(0)} = \boldsymbol{0}$
Suppose the algorithm misclassifies sample 1, 7, 29, and 1 again

First update:   $\boldsymbol{w}^{(1)} = \eta y_1 \boldsymbol{x}_1$
Second update: $\boldsymbol{w}^{(2)} \doteq \eta y_1 \boldsymbol{x}_1 + \eta y_7 \boldsymbol{x}_7$
Third update:   $\boldsymbol{w}^{(3)} = \eta y_1 \boldsymbol{x}_1 + \eta y_7 \boldsymbol{x}_7 + \eta y_{29} \boldsymbol{x}_{29}$
Third update:   $\boldsymbol{w}^{(4)} = 2\eta y_1 \boldsymbol{x}_1 + \eta y_7 \boldsymbol{x}_7 + \eta y_{29} \boldsymbol{x}_{29}$   etc.

*(handwritten annotations)*
$W^k - \eta \Delta l_i(w) = w^k + \eta x_i y_i$
$\Delta l_i(w) = (l - s_i y_i)$
$\frac{\partial}{\partial w_i}$
$= -x_i y_i$
那两次错了加 两次
那个怎5错3
如口那次错3加哪次

43

# Accumulating updates: Data enters via dot products

- Weights always take the form $\mathbf{w} = \sum_{j=1}^{N} \alpha_j y_j \mathbf{x}_j$,
  where $\boldsymbol{\alpha}$ some coefficients

  $\alpha$: summation of how many $\eta$ I have to assign to that sample when it was misclassified.
  — how many learning rate $\eta$ was add to that.

- Perceptron weights always linear comb. of data!

- Recall that prediction for a new point $\mathbf{x}$ is based on sign of $\mathbf{w}'\mathbf{x}$

- Substituting $\mathbf{w}$ we get $\mathbf{w}'\mathbf{x} = \sum_{j=1}^{N} \alpha_j y_j \mathbf{x}_j'\mathbf{x}$

- The dot product $\mathbf{x}_j'\mathbf{x}$ can be replaced with a kernel ✗

  keeping $\alpha \cdot y$

  $\alpha_j y_j x_j' x$

# Kernelised perceptron training rule

Set $\boldsymbol{\alpha} = \mathbf{0}$

For $t$ from 1 to $T$ (epochs)

    For $i$ from 1 to $N$ (training examples)

第i sample 被分类正确,对应的 α是 0

有错 one time → $+\eta$

错两次 two times → $+2\eta$ = 2

Compute $s_i = \sum_{j=1}^{N} \alpha_j y_j x_j' x_i$ kernel

If $s_i y_i \leq 0$: (sample $i$ misclassified)

对i 猜错一次 就对 α_i 加一次η

$$\alpha_i \leftarrow \alpha_i + \eta$$

($\eta > 0$ is called *learning rate*)

# Kernelised perceptron training rule

Set $\boldsymbol{\alpha} = \mathbf{0}$

For $t$ from 1 to $T$ (epochs)

    For $i$ from 1 to $N$ (training examples)

Compute $s_i = \sum_{j=1}^{N} \alpha_j y_j K(\boldsymbol{x}_j, \boldsymbol{x}_i)$    等标签    kernel→inner product

If $s_i y_i \leq 0$: (sample $i$ misclassified)    说们错了

$\alpha_i \leftarrow \alpha_i + \eta$

$\eta$ learning rate

$\eta =1$ 时，在数分类错误的次数

郭社星

# Mini Summary

- Accumulating weight updates leads to linear combinations of data

- Predictions are dot products with data

- Can replace these with kernel evaluations

- Leads to kernel perceptron with kernel training rule

Next time: Deep learning

$$\Longleftrightarrow \sup_{g \in G} \left( E g(Z) - \frac{1}{n} \sum_{i=1}^{n} (z_i) \right) \le 2R_n(G) + \sqrt{\frac{\log \frac{1}{\delta}}{2n}}$$

$$\Longleftrightarrow \Phi(Z_1, \cdots Z_n) = \sup_{g \in G} \left( E g(Z) - \frac{1}{n} \sum_{i=1}^{n} g(z_i) \right)$$

$$\le E \Phi(Z_1, \cdots, Z_n) \qquad + \frac{\sqrt{\log \frac{1}{\delta}}}{\sqrt{2n}}$$

$$\le 2R_n(G$$

$$\Longleftrightarrow \Phi(Z_1, \cdots Z_n) - E_{Z_1, \cdots, Z_n} \Phi(Z_1, \cdots, Z_n)$$

$$\le \xi \qquad\qquad \xi = \sqrt{\frac{\log \frac{1}{\delta}}{2n}} \quad \delta = \exp(-2n\xi)$$

$$\Longleftrightarrow E g(Z) \le \frac{1}{n} \sum_{i=1}^{n} g(z_i) + 2\hat{R}_s(G) + \sqrt{\frac{\log \frac{2}{\delta}}{2n}}$$

Apply McDiamid to $\hat{R}_s(G)$ as function of $S = (z_1, \cdots, z_n)$

Bounded difference property:

$$\hat{R}_{z_1, \dots, z_n}(\mathcal{G}) - \hat{R}_{z_1, \dots, z_{k-1}, z'_1, \dots}(\mathcal{G})$$

$$= E_6 \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_i \sigma_i g(z_i) - \sup_{g \in \mathcal{G}} \left( \frac{1}{n} \sum_i \sigma_i g(z_i) + \sigma_k g(z'_k) \right) \right]$$

$$\leq E_6 \sup_{g \in \mathcal{G}} \left( \frac{1}{n} \sigma_k g(z_k) - \frac{1}{n} \sigma_k g(z'_k) \right) \leq \frac{2}{n}$$

---

how to apply to binary classification?

Lemma: [Mt Lemma 3.4]

Let $\mathcal{H}$: ~~faith~~ family of functions $h: \mathcal{X} \to \{-1, +1\} =: \mathcal{Y}$

$\quad$ G: loss function

$\quad$ $\mathcal{G} := \{g: \mathcal{Z} \to \{0,1\}, (x,y) \mapsto 1_{h(x) \neq y}: h \in \mathcal{H}\}$

$\qquad\qquad : \mathcal{X} \times \mathcal{Y}$

$\quad$ $z_1 = (x_1, y_1), \dots, z_n = (x_n, y_n)$

$\quad$ Then $\hat{R}_{z_1, \dots, z_n}(\mathcal{G}) = \frac{1}{2} \hat{R}_{x_1, \dots, x_n}(\mathcal{H})$ $\quad$ (+) $\quad$ ☆

$\quad$ Proof see [Mt]

Simple Corollary $R_n(G) = \frac{1}{2} R_n(H)$

Theorem : H family of functions

$P_X$ dist on $X$

$\delta \in (0,1)$

with.h.p $\geq 1-\delta$ over $q_X^n$, for all $h \in H$ simultaneously

$$\begin{cases} R(h) \leq \hat{R}(h) + R_n(H) + \sqrt{\dfrac{\log \frac{1}{\delta}}{2n}} \\[4mm] R(h) \leq \hat{R}(h) + \hat{R}_{x_1,\cdots,x_n}(H) + 3\sqrt{\dfrac{\log \frac{2}{\delta}}{2n}} \end{cases}$$

Proof:

Let $q \in G$ be associated to $h$ i.e $g(x,y) \to \mathbb{1}_{h(x) \neq y}$

Then $R(h) = \mathbb{E} g(z)$ $\qquad \hat{R}(h) = \frac{1}{n}\sum_{i=1}^{n} g(z_i)$

Lemma : $2\hat{R}_{z_1,\cdots,z_n}(G) = \hat{R}_{x_1,\cdots,x_n}(H)$

$2R_n(G) = R_n(H)$

$$\Pi_H(m) = 2^m$$

maximum set of this set of dichotomy.

If $\exists S$ with $|S| = m$ shattered by $H$, i.e. $|\Pi_H(S)| = 2^m$ then $\Pi_H^{(m)} = \max_{\substack{j=2^m \\ S \subseteq X^m}} |\Pi_j(S)|$

$$dvc(H) := \max\left\{ m : \Pi_H(m) = 2^m \right\}$$

Note: ① does not imply all sets of size $\leq d$ are shattered by $H$.

② often easier to compute the growth function $G$ for $RC$

③ Connect to VC dimension & growth function. Via Sauer's Lemma.

Sauer's Lemma: Let $H$ be a class of function $h: X \to \{+1, 1\}$ and let $d := dvc(H)$. Then $\Pi_H(m) \leq \sum \binom{m}{i}$

upper bound for growth function

Proof: (by induction on the sum $m + d$) Define $\Phi := \sum_{i=0}^{d} \binom{m}{i}$

Basic case: $m = 0$ (labelling on an empty set)

$$\Pi_H(0) = 1 \quad \left.\begin{array}{c} \end{array}\right\} \text{ if } d = 0, \text{ no dataset can be shattered}$$

$$\Phi_d(0) = \binom{m}{0} = 1 \quad \left.\begin{array}{c}\end{array}\right.$$

$$\Pi_H(m) \leq 1$$

$$\Phi_0(m) = \binom{m}{0} = 1$$

(24): if $m=1$ (1-elt set) $d=0$

if $m=1$

$$\Pi_H(1) \leq$$

$$\Phi_d(m) = \sum_{i=0}^{d} \binom{m}{i}$$

$$\longrightarrow \Pi_{(H)}(1) \leq 2.$$

inductive step:
  assume for some $m, d \geq 1$, we have $m'+d' \leq m+d$.
  and $\Pi_{H'}(m') \leq \Phi_{d'}(m') = \sum_{i=0}^{d'} \binom{m'}{i}$,

  where $H' = H$ restricted to $m' = m-1$  ($H' \subseteq H$ subset).

  $$d' = d_{VC}(H') = d \text{ or } d-1$$

Consider: labellings induced by $H$ on any set $S = \{\vec{x_1} \cdots \vec{x_m}\}$

w log let $S_1 = \{x_1, \cdots, x_{m-1}\}$.
            $= S \setminus \{x_m\}$

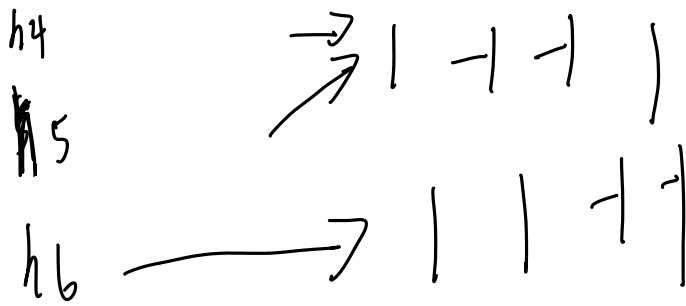Let $H_1$ be the set of hypothesis restricted to $S_1$

$$\Pi_H(S) \qquad \Pi_{H_1}(S_1) \quad \leq (1-\vec{3})$$

$$x_1 \ x_2 \ x_3 \ x_4 \ x_5$$

$h_1$  | 1 1 1 + + $\rightarrow$  -| | | -|
$h_2$  -| | | +| $\rightarrow$  -| | ? )
$h_3$              $\rightarrow$  -| | ? )

h4

h5 $\rightarrow$ | -| -| )

h6 $\longrightarrow$ | | -| -|

$$\underline{H \ 与 \ S \ 不同}$$

H shatters set S, $H_1$ shatter a set S

$H_1$ shatter a set $\Rightarrow$ then so does H.

$$d_{VC}(H_1) \leq d_{VC}(H) = d$$