

COMP4901K/Math4824B

Machine Learning for Natural Language Processing

Lecture 8: Perceptron, Error-Driven Classification

Instructor: Yangqiu Song

Today

- Algorithms for Discriminative Classification
- Binary classification
 - Perceptron

Binary Classification: examples

- Spam filtering (**spam**, **not spam**)
- Customer service message classification (**urgent** vs. **not urgent**)
- Information retrieval (**relevant**, **not relevant**)
- Sentiment classification (**positive**, **negative**)

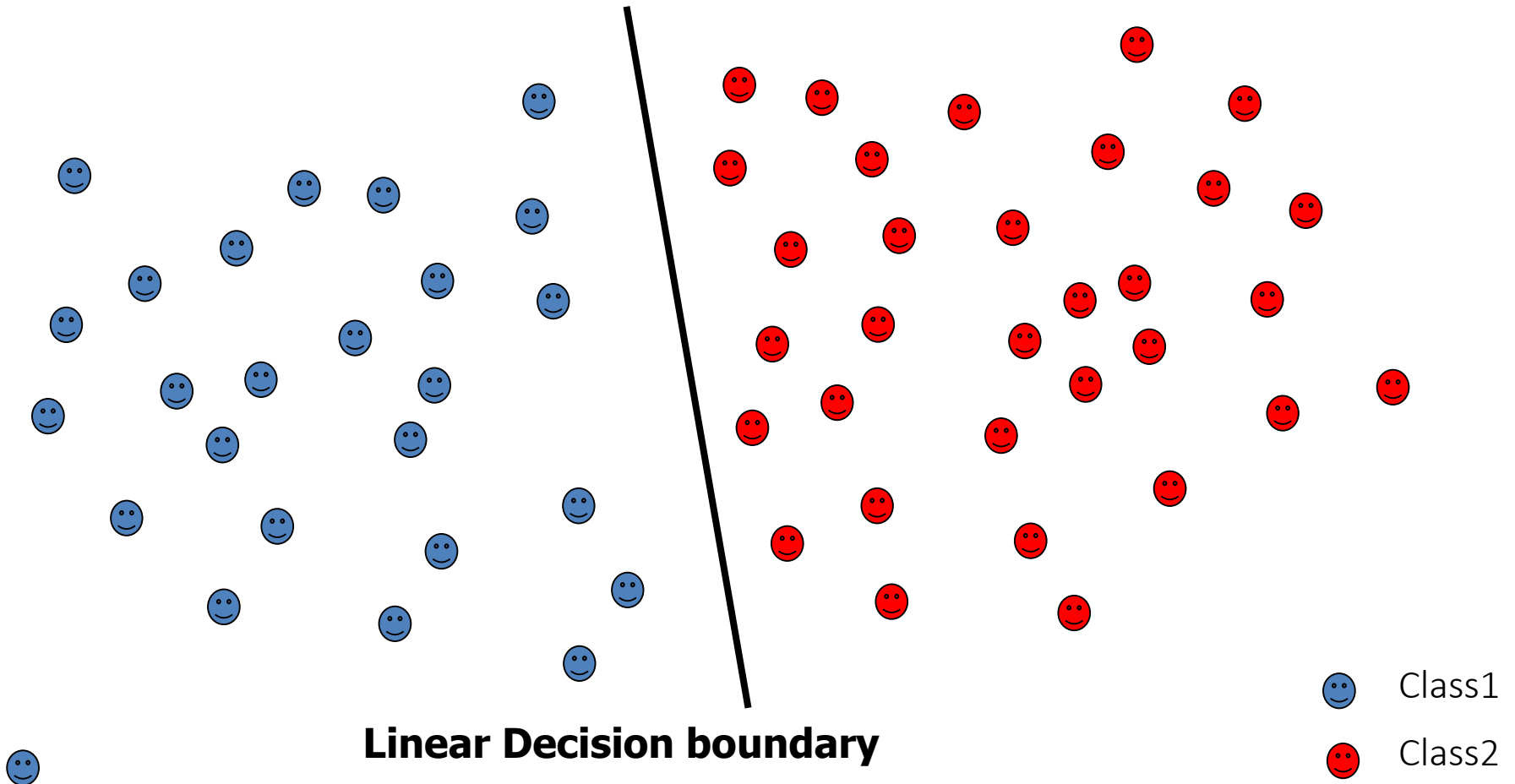
Binary Classification

- **Given:** some data items that belong to a positive (+1 😊) or a negative (-1 😞) class
- **Task:** Train the classifier and predict the class for a new data item
- **Geometrically:** find a separator

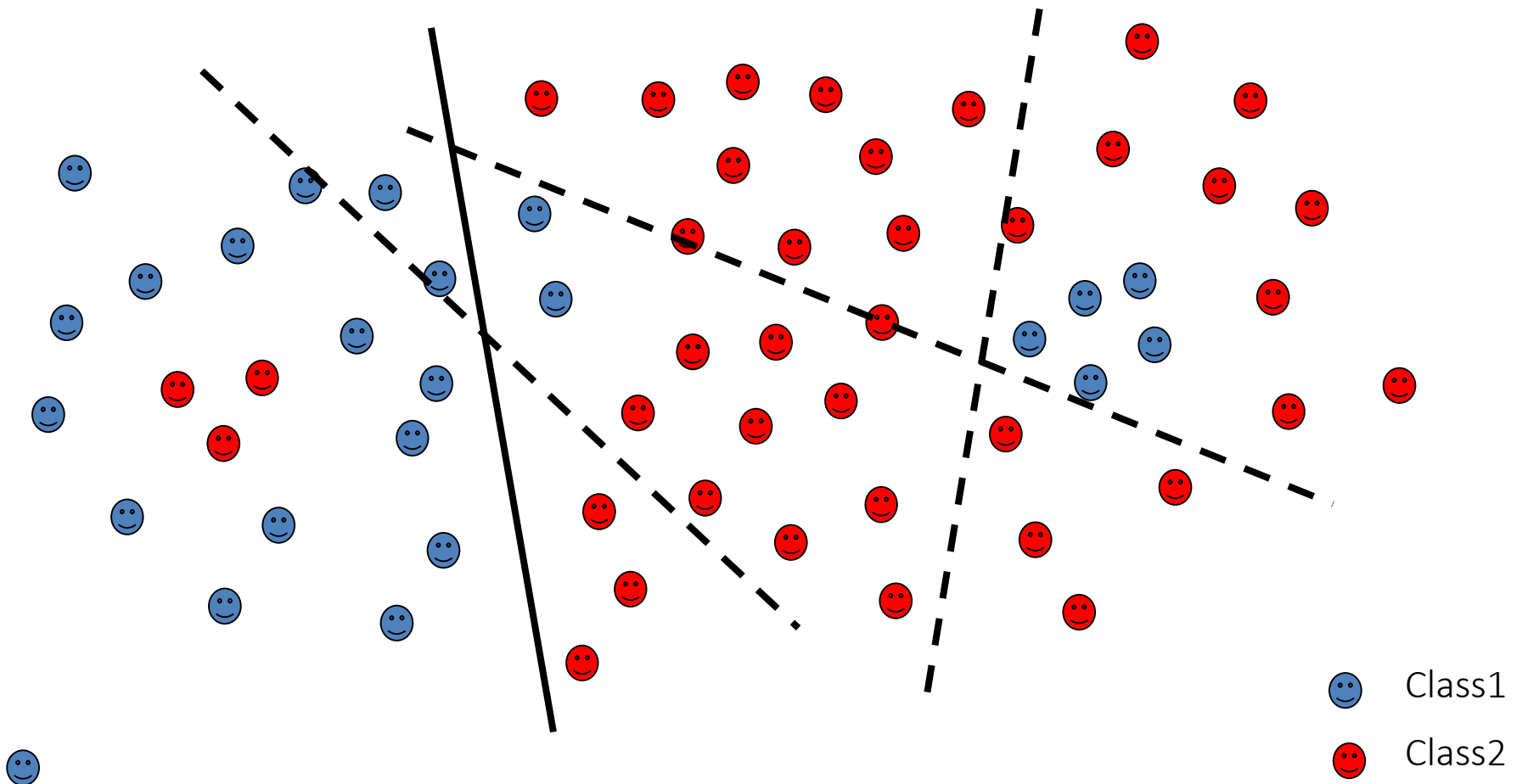
Linear versus Non Linear algorithms

- **Linearly separable data:** if all the data points can be correctly classified by a linear (hyperplanar) decision boundary

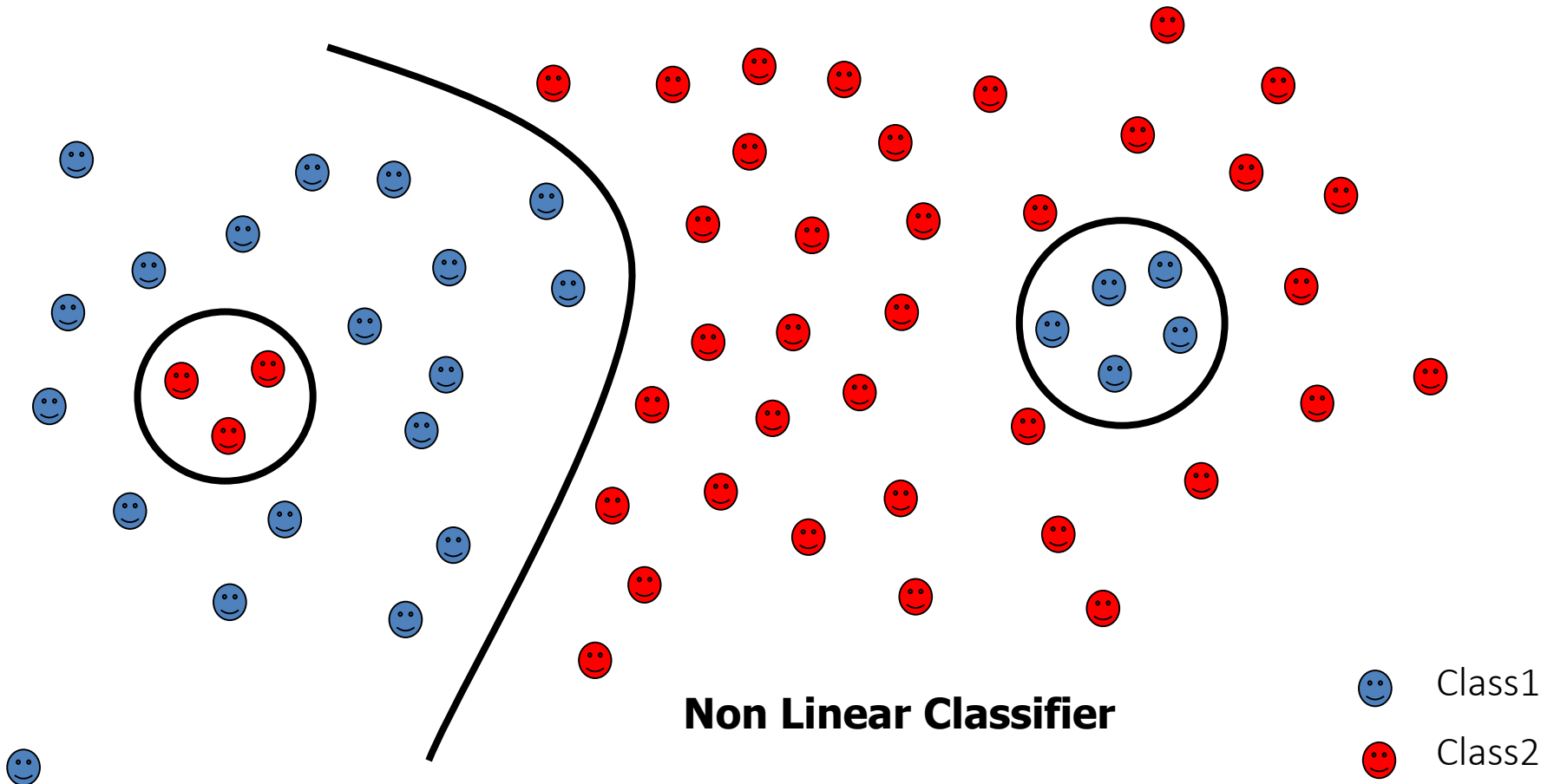
Linearly separable data



Non linearly separable data



Non linearly separable data



Linear versus Non Linear algorithms

- Linear or Non linear separable data?
 - We can find out only empirically
- Linear algorithms (algorithms that find a linear decision boundary)
 - When we think the data is linearly separable
 - Advantages
 - Simpler, less parameters
 - Disadvantages
 - Data (like for texts) is usually not linearly separable
 - Examples: Perceptron, Winnow, SVM

Linear versus Non Linear algorithms

- Non Linear
 - When the data is non linearly separable
 - Advantages
 - More accurate
 - Disadvantages
 - More complicated, more parameters
 - Example: Deep learning

Simple linear algorithms

- Perceptron algorithm
 - Linear
 - Binary classification
 - Online (process data sequentially, one data point at the time)
 - Mistake driven
 - Simple single layer Neural Networks

Linear binary classification

■ Data: $\{(\mathbf{x}_i, y_i)\}_{i=1\dots n}$

- \mathbf{x} in \mathbb{R}^d (\mathbf{x} is a vector in d -dimensional space)

→ feature vector

- y in $\{-1, +1\}$

→ label (class, category)

■ Question:

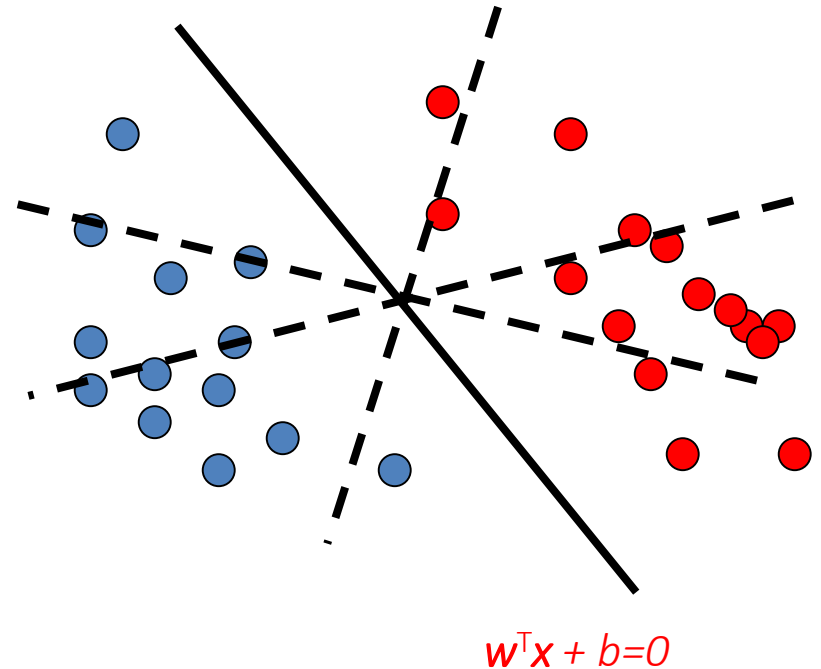
- Design a linear decision boundary: $\mathbf{w}^T \mathbf{x} + b$ (equation of hyperplane) such that the classification rule associated with it has minimal probability of error

- **classification rule**

- $y = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$ which means:
- if $\mathbf{w}^T \mathbf{x} + b > 0$ then $y = +1$
- if $\mathbf{w}^T \mathbf{x} + b < 0$ then $y = -1$

Linear binary classification

- Find a good **hyperplane**
 (\mathbf{w}, b) in \mathbb{R}^{d+1}
that correctly classifies data points as much as possible
- In **online fashion**: one data point at the time, update weights as necessary

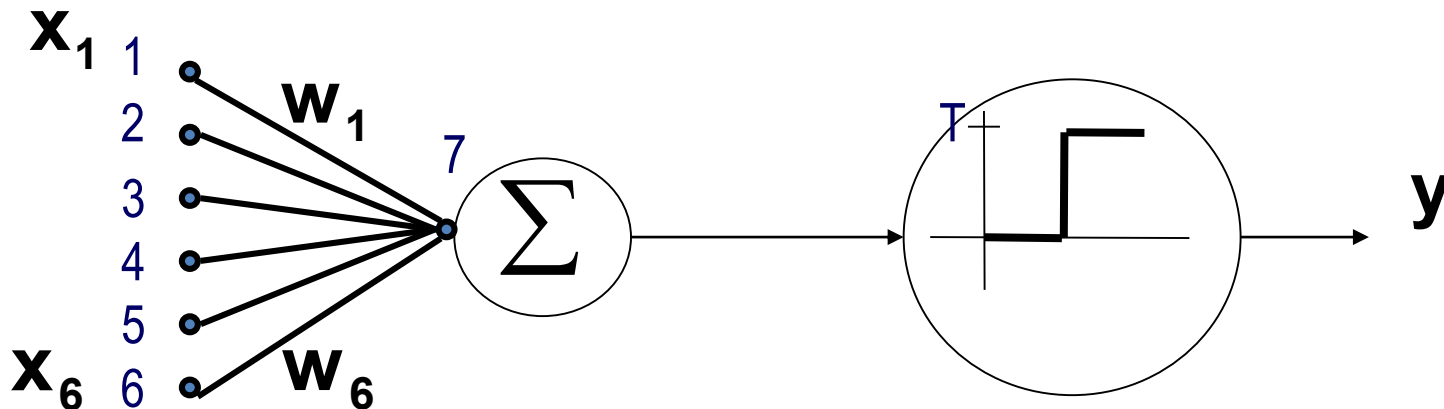


Classification Rule:

$$y = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

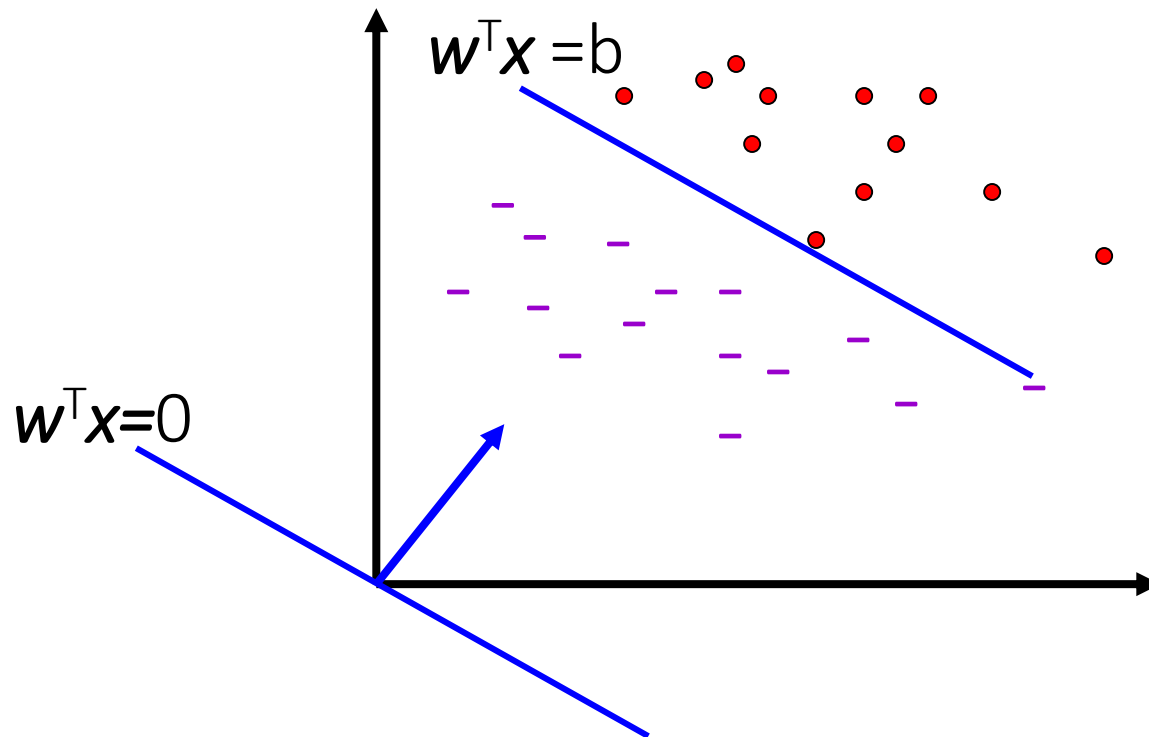
Perceptron learning rule

- On-line, mistake driven algorithm.
- [Rosenblatt \(1959\)](#) suggested that when a target output value is provided for a **single neuron with fixed input**, it can incrementally change weights and learn to produce the output using the Perceptron learning rule
- (Perceptron == Linear Threshold Unit)



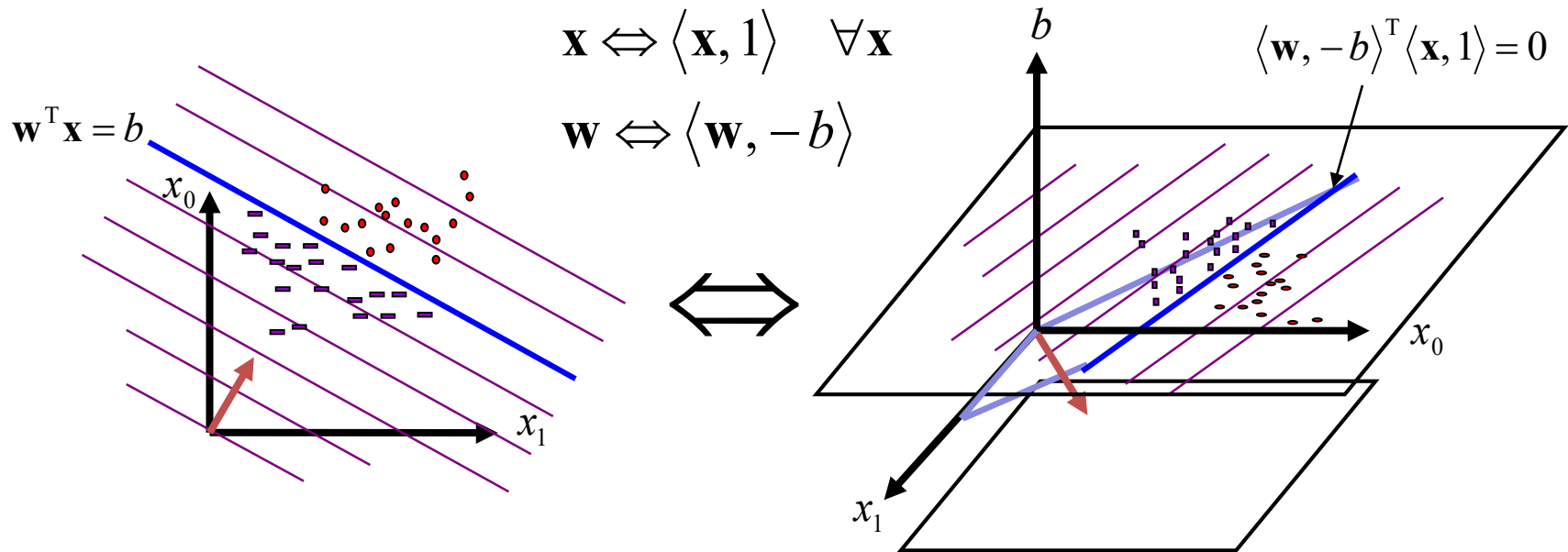
Perceptron learning rule

- We learn $f:\mathbf{X}\rightarrow\{-1,+1\}$ represented as $f=\text{sgn}\{\mathbf{w}^T\mathbf{x}\}$
- Where $\mathbf{X}=\{0,1\}^n$ or $\mathbf{X}=\mathbb{R}^n$ and $\mathbf{w}\in\mathbb{R}^n$
- Given Labeled examples: $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$



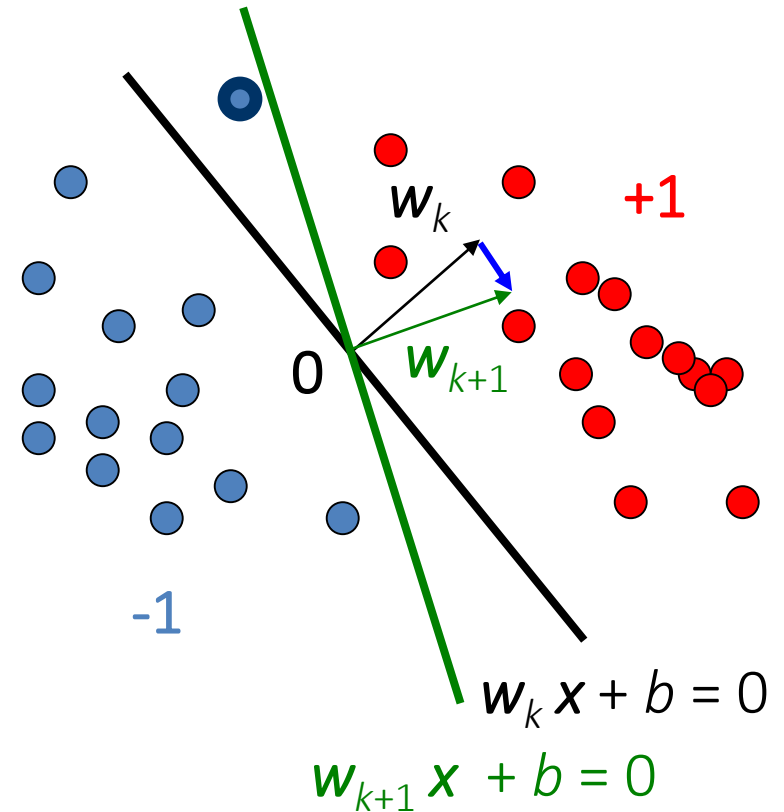
Footnote About the Threshold

- On previous slide, Perceptron has no threshold
- But we don't lose generality:

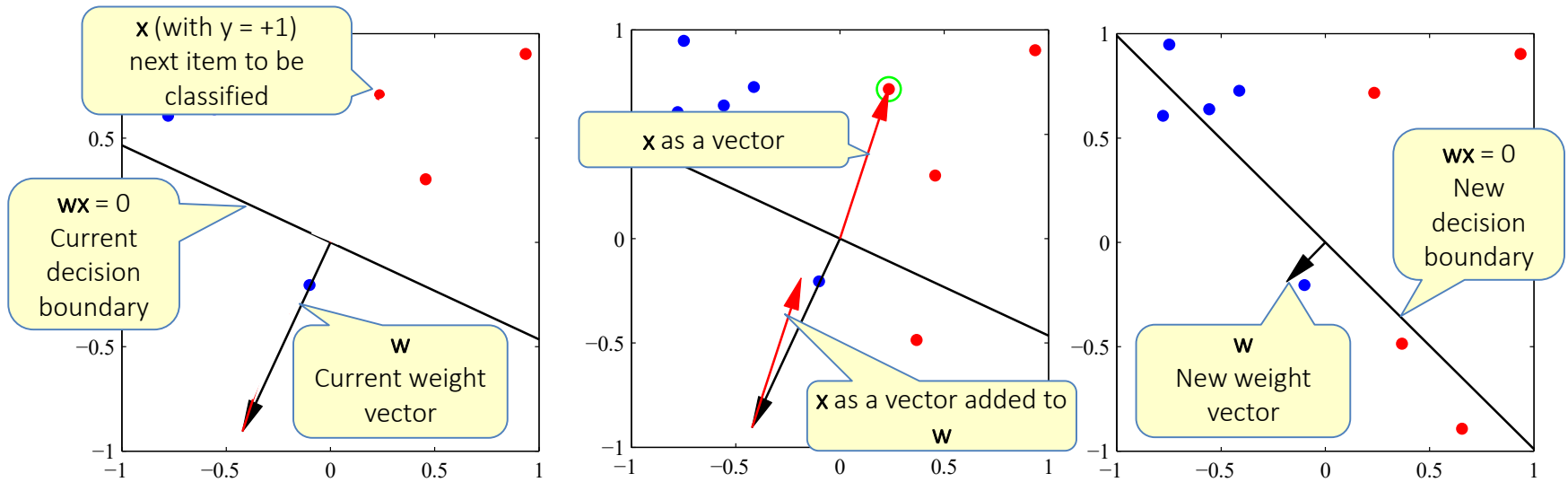


Perceptron algorithm

- Initialize: $\mathbf{w}_1 = \mathbf{0}$ in \mathbb{R}^n
- Updating rule
- For each data point \mathbf{x}
 - Predict the label $y' = \text{sgn}\{\mathbf{w}^T \mathbf{x}\}$
 - if $y' \neq y$, update the weight vector
$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + a y_i \mathbf{x}_i$$
(a - a constant, learning rate)
 - else
$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k$$
- Function: $y' = \text{sgn}\{\mathbf{w}^T \mathbf{x}\}$
 - if $\mathbf{w}^T \mathbf{x} > 0$ return +1
 - else return -1



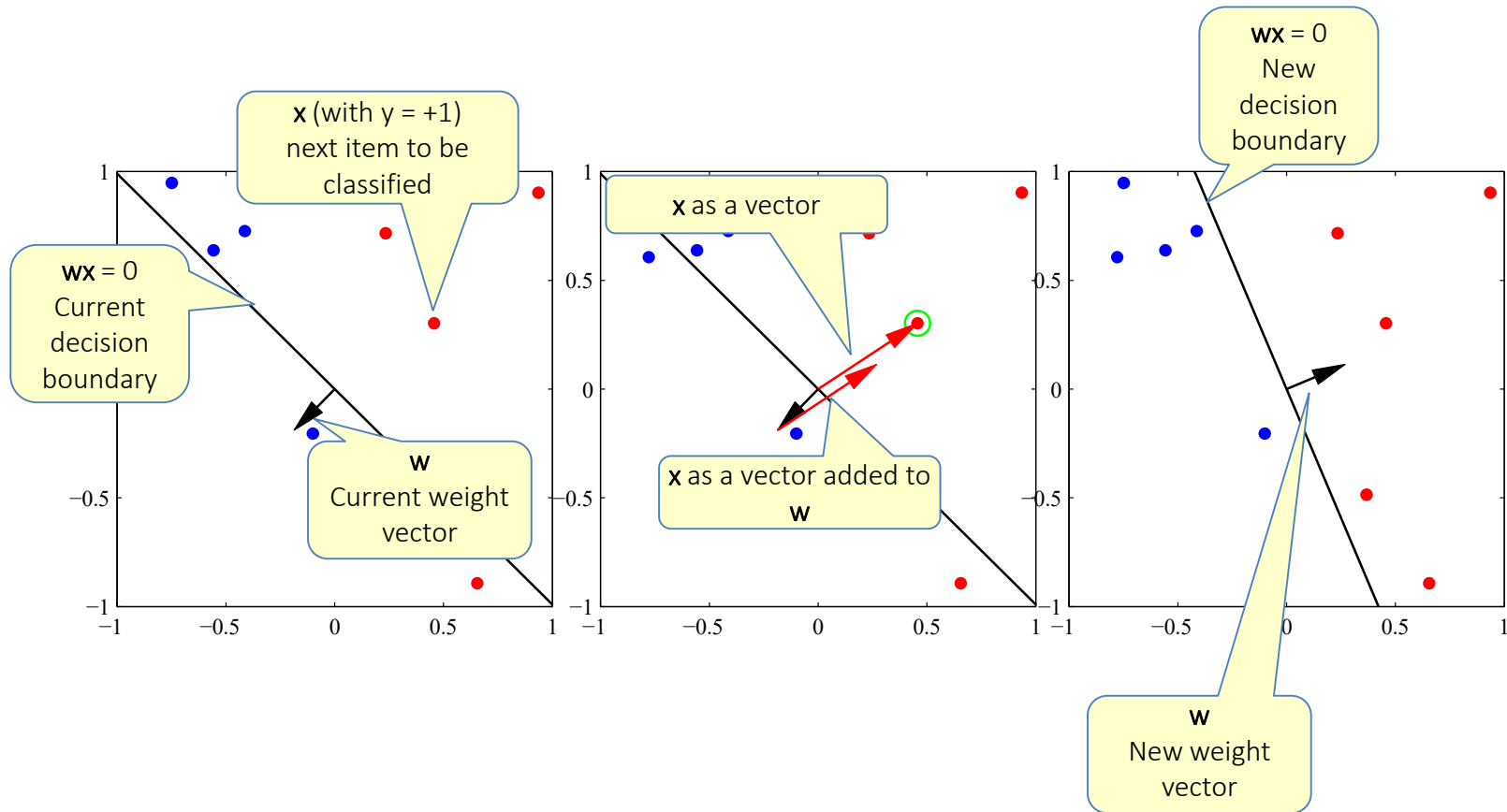
Perceptron in action



(Figures from Bishop 2006)



Perceptron in action



(Figures from Bishop 2006)

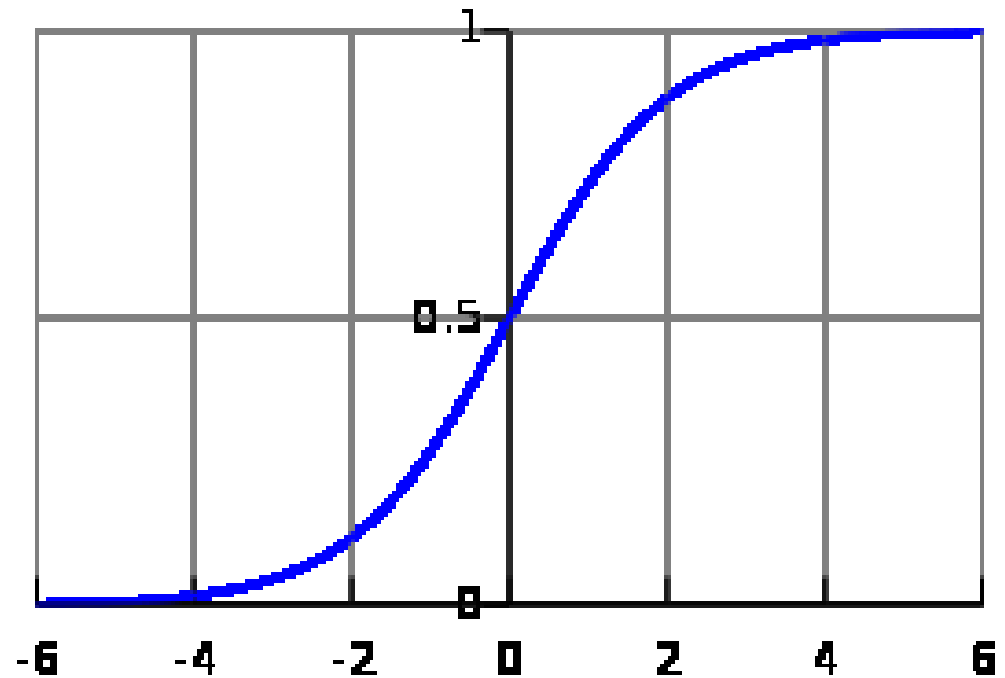
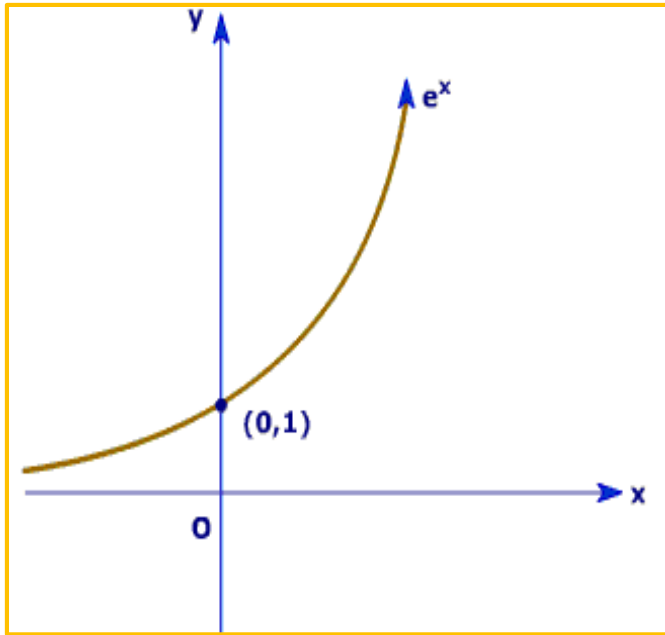
Perceptron learning rule

- If x is Boolean, only weights of **active features** are updated

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \mathbf{x}$$
$$\begin{pmatrix} w^{(1)} + 1 \\ w^{(2)} \\ w^{(3)} - 1 \end{pmatrix} = \begin{pmatrix} w^{(1)} \\ w^{(2)} \\ w^{(3)} \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}$$

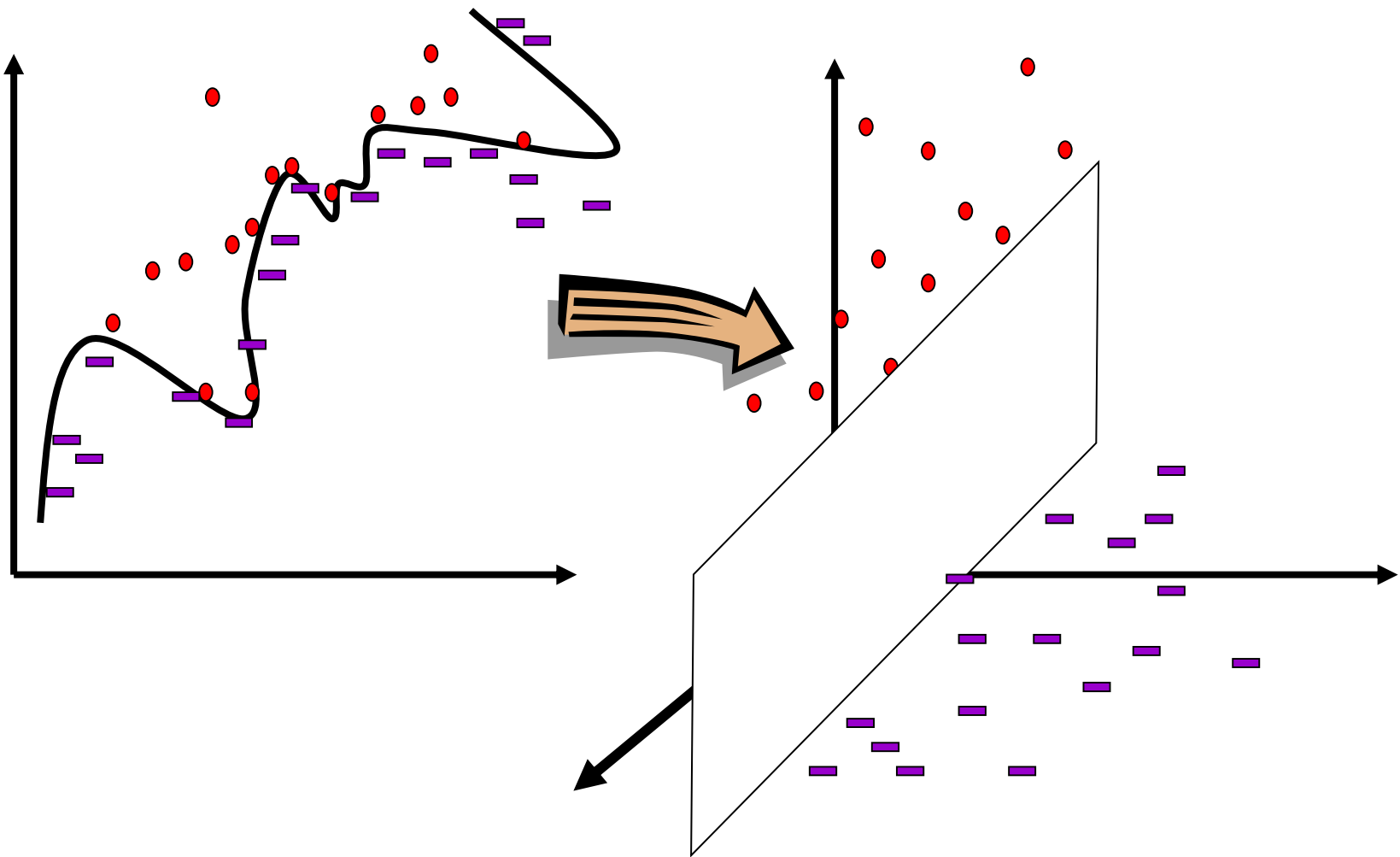
Perceptron learning rule

$$\mathbf{w}^T \mathbf{x} > 0 \text{ is equivalent to } \frac{1}{1 + \exp\{-(\mathbf{w}^T \mathbf{x})\}} > 1/2$$



Perceptron Learnability

- Obviously can't learn what it can't represent Only linearly separable functions
- Minsky and Papert (1969) wrote an influential book demonstrating Perceptron's representational limitations
 - Parity functions can't be learned (XOR)
 - In vision, if patterns are represented with local features, can't represent symmetry, connectivity
- Research on Neural Networks stopped for years
- Rosenblatt himself (1959) asked,
 - *“What pattern recognition problems can be transformed so as to become linearly separable?”*



Perceptron Convergence

- Perceptron Convergence Theorem:
 - If there exist a set of weights that are consistent with the data (i.e., the data is linearly separable), the perceptron learning algorithm will converge
 - How long would it take to converge ?
- Perceptron Cycling Theorem:
 - If the training data is not linearly separable the perceptron learning algorithm will eventually repeat the same set of weights and therefore enter an infinite loop.
 - How to provide robustness, more expressivity ?

Perceptron-Mistake Bound

- Assume that a weight vector $\mathbf{w} \in \mathbb{R}^d$. Upon receiving an example $\mathbf{x} \in \mathbb{R}^d$, we predict according to a linear threshold function.
- Theorem [Novikoff,1963]
 - Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ be a sequence of labeled examples with $\mathbf{x}_i \in \mathbb{R}^d$, $\|\mathbf{x}_i\| \leq r$, and $y_i \in \{-1, +1\}$ for all i .
 - Let $\mathbf{u} \in \mathbb{R}^d$, $\gamma > 0$ be such that $\|\mathbf{u}\| = 1$ and $y_i \cdot \mathbf{u}^T \mathbf{x}_i \geq \gamma$ for all i .
 - Perceptron makes at most $\frac{r^2}{\gamma^2}$ mistakes on this example sequence.
- Assumptions:
 - This theorem assumes that all examples are bounded by some r ; for all \mathbf{x}_i , find the largest one, and r is at least this size.
 - The theorem further assumes that there exists some \mathbf{u} that separates the data.
 - Requiring that $\|\mathbf{u}\| = 1$ is simply a constant that could be arbitrarily scaled.
 - Finally, the theorem assumes that there exists some γ such that the inequality is satisfied.
 - We refer to γ as the complexity parameter: γ is very large, finding a hyperplane is much easier

Proof

- Let \mathbf{w}_k be the hypothesis before the k th mistake.
- Assume that the k th mistake occurs on the input example (\mathbf{x}_i, y_i)

$$\therefore y_i(\mathbf{w}_k^T \mathbf{x}_i) \leq 0$$

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

$$\mathbf{w}_{k+1}^T \mathbf{u} = \mathbf{w}_k^T \mathbf{u} + y_i \mathbf{x}_i^T \mathbf{u} \geq \mathbf{w}_k^T \mathbf{u} + \gamma \quad (\because y_i \mathbf{u}^T \mathbf{x}_i \geq \gamma)$$

$$\geq \mathbf{w}_{k-1}^T \mathbf{u} + 2\gamma$$

...

$$\geq k\gamma$$

$$\begin{aligned} \|\mathbf{w}_{k+1}\|^2 &= \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + \|\mathbf{x}_i\|^2 \\ &\leq \|\mathbf{w}_k\|^2 + r^2 \quad (\because y_i(\mathbf{w}_k^T \mathbf{x}_i) \leq 0) \end{aligned}$$

...

$$\leq kr^2$$

$$\text{Therefore, } \sqrt{k}r \geq \|\mathbf{w}_{k+1}\| \geq \mathbf{w}_{k+1}^T \mathbf{u} \geq k\gamma$$

(Second inequality: $\mathbf{u}^T \mathbf{v} = \|\mathbf{u}\| \cdot \|\mathbf{v}\| \cdot \cos(\mathbf{u}, \mathbf{v}) \leq \|\mathbf{u}\| \cdot \|\mathbf{v}\| = \|\mathbf{v}\|$ and $\|\mathbf{u}\| = 1$)

$$\therefore k \leq \frac{r^2}{\gamma^2}$$

- Further reading

- <https://www.cse.iitb.ac.in/~shivaram/teaching/old/cs344+386-s2017/resources/classnote-1.pdf>

Practical Issues and Extensions

- There are many extensions that can be made to these basic algorithms.
- Some are necessary for them to perform well
 - Regularization (later soon)
- Some are for ease of use and tuning
 - Converting the output of a Perceptron to a conditional probability
$$P(y = +1 \mid x) = [1 + \exp(-a wx)]^{-1}$$
 - Can tune the parameter a
- Multiclass classification

Regularization Via Averaged Perceptron

- An **Averaged Perceptron** Algorithm is motivated by the following considerations:
 - In the mistake bound model: We don't know when we will make the mistakes.
- **Averaged Perceptron** returns a weighted average of a number of earlier hypotheses; the weights are a function of the length of no-mistakes stretch.

Regularization Via Averaged Perceptron

- Training:
[**m**: #(examples); **k**: #(mistakes) = #(hypotheses); **c_i**: consistency count for **w_i**]
- **Input**: a labeled training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$
- Number of epochs T
- **Output**: a list of weighted perceptrons $\{(\mathbf{w}_1, c_1), \dots, (\mathbf{w}_k, c_k)\}$
- Initialize: $k=0$; $\mathbf{w}_1 = 0$, $c_1 = 0$
- Repeat T times:
 - For $i=1, \dots, m$:
 - Compute prediction $y' = \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$
 - If $y' = y$, then $c_k = c_k + 1$
else: $\mathbf{w}_{k+1} = \mathbf{w}_k + \mathbf{w}_i$; $c_{k+1} = 1$; $k = k+1$
- Prediction:
- **Given**: a list of weighted perceptrons $\{(\mathbf{w}_1, c_1), \dots, (\mathbf{w}_k, c_k)\}$; a new example \mathbf{x}
Predict the label(\mathbf{x}) as follows:
$$y(\mathbf{x}) = \text{sign} \left[\sum_{i=1, k} c_i \text{sign}(\mathbf{w}_i^T \mathbf{x}) \right]$$

Perceptron algorithm

- **Online**: can adjust to changing target, over time
- **Advantages**
 - Simple and computationally efficient
 - Guaranteed to learn a linearly separable problem (convergence, global optimum)
- **Limitations**
 - Only linear separations
 - Only converges for linearly separable data
 - Not really “efficient with many features”