

# Machine Learning for Cognitive Sciences: Principles and Applications

**Simona Cocco (Physics Departement of ENS)**  
[simona.cocco@phys.ens.fr](mailto:simona.cocco@phys.ens.fr) CM &TD

**Vito Dichio (Physics Departement of ENS)**  
[vito.dichio@phys.ens.fr](mailto:vito.dichio@phys.ens.fr) TD

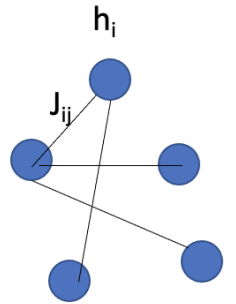
Week 10

# Recall of previous lesson

## Network inference

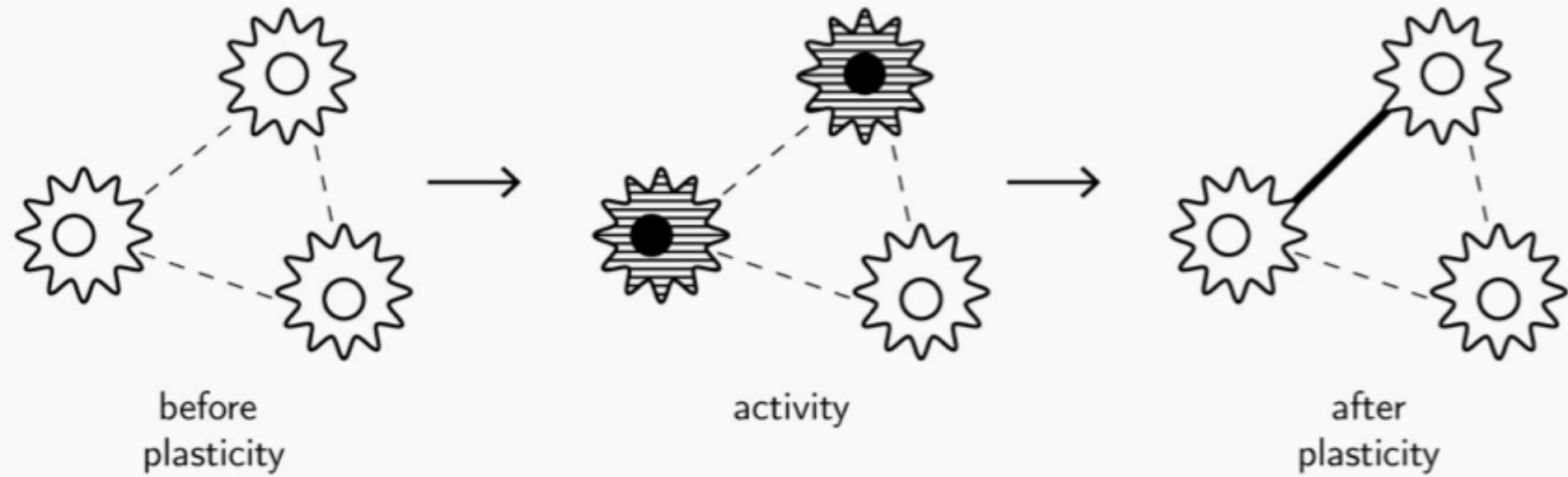
Different approximations and heuristics to infer couplings in an efficient way : Multivariate Gaussian variables, exact maximization of Log-Likelihood

**Hopfield Model::** Model for auto-associative memory: patterns/memories are stored in the coupling matrix.



# Hebbian learning through synaptic Potentiation

*Cells which fire together wire together*



- Already suggested by Ramon Y Cajal in 1894: memory results from the change of couplings between neurons, rather than from the appearance of new neurons

Hebb (1949)

- Cell assemblies hypothesis (Hebb 1949 ): Information is represented , replayed and stored trough a group of neurons firing together

# Hopfield model for auto-associative memories

One can build up a coupling matrix storing  $P$  features or 'patterns' as energy minima of the model (associative memories of the network):

Hebbe rule 
$$J_{ij} = \sum_{\mu=1}^P w_i^{\mu} w_j^{\mu}$$

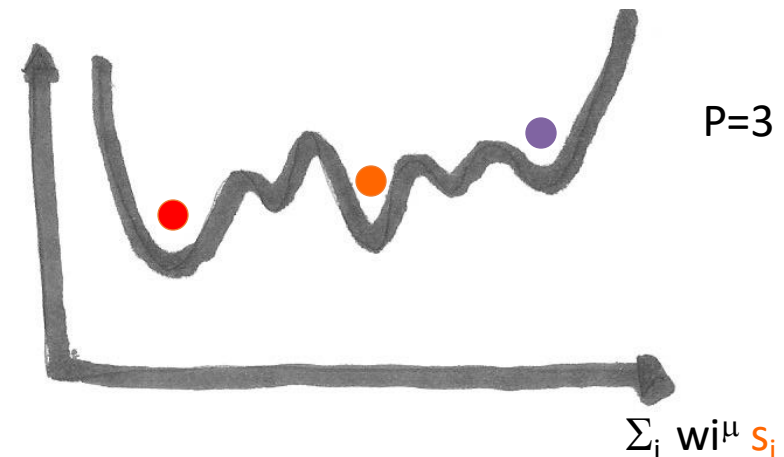
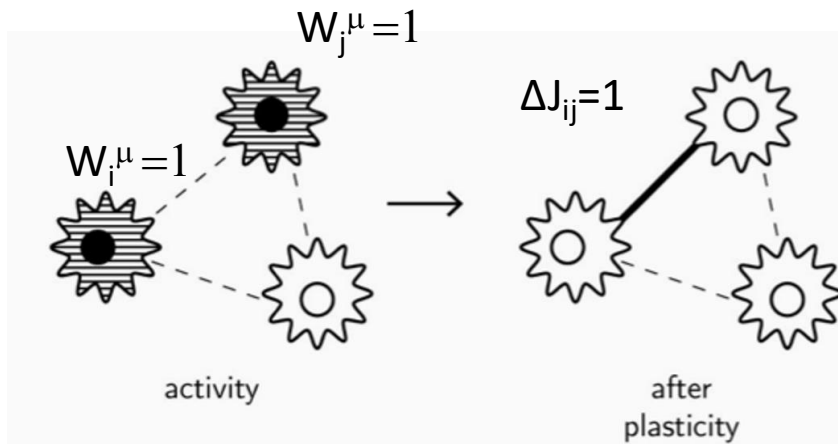
$$w_i^{\mu} \in [-1, 1]$$

$$s_i \in [-1, 1]$$

Configurations which look like the patterns are minima of the energy

$$E(\mathbf{s}) = -\frac{1}{2} \sum_{i < j} J_{ij} s_i s_j$$

$$E(\mathbf{s}) = -\frac{1}{2} \sum_{\mu=1}^P \left( \sum_i w_i^{\mu} s_i \right)^2$$



# Hopfield model for auto-associative memories

One can build up a coupling matrix storing  $M$  features or 'patterns' as energy minima of the model (associative memories of the network):

Hebb rule  $J_{ij} = \sum_{\mu=1}^P w_i^{\mu} w_j^{\mu}$

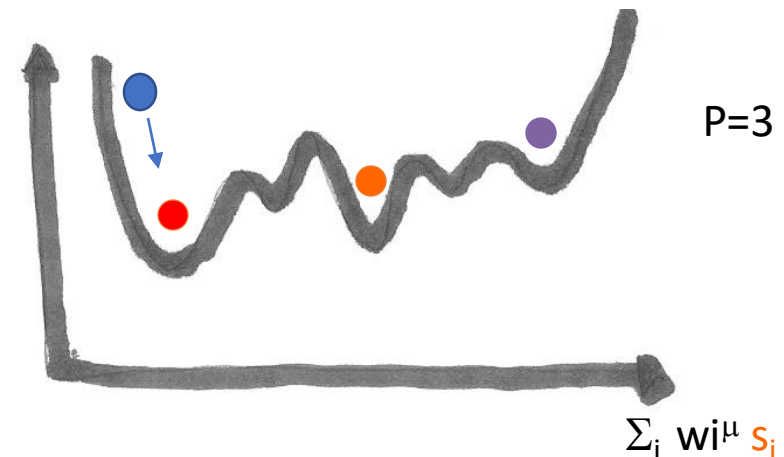
$w_i^{\mu} \in [-1, 1]$   
 $s_i \in [-1, 1]$

Configurations which look like the patterns are minima of the energy:

$$E(\mathbf{s}) = -\frac{1}{2} \sum_{i < j} J_{ij} s_i s_j$$

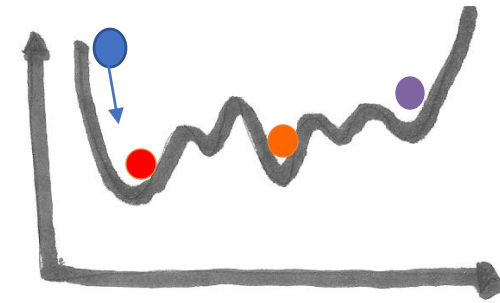
$$E(\mathbf{s}) = -\frac{1}{2} \sum_{\mu=1}^P \left( \sum_i w_i^{\mu} s_i \right)^2$$

This network has properties of associative memory: dynamics of the system starting on  $\bullet$  fall into the nearest minima  $\bullet$



# Hopfield model for auto-associative memories

$$s_i^{t+1} = \text{sign} \left( \sum_{j \neq i} J_{i,j} s_j^t \right)$$



$$\sum_i w_i^\mu s_i$$

Starting from a pattern  $s^t = \mathbf{w}^{\mu'}$  (or nearby a pattern),  
the pattern is a stable fixed point of the dynamics above if the number of stored patterns  
 $P < L$  (number of neurons)

Proof:

$$s_i^{t+1} = \text{sign} \left( \sum_{\mu=1}^P w_i^\mu \sum_{j \neq i} w_j^\mu s_j^t \right)$$

Starting pattern

Other patterns

$$s_i^{t+1} = \text{sign} \left( w_i^{\mu'} (L - 1) + \sum_{\mu \neq \mu'} \sum_{j \neq i} w_i^\mu w_j^\mu w_j^{\mu'} \right)$$

The second term is a random variables with

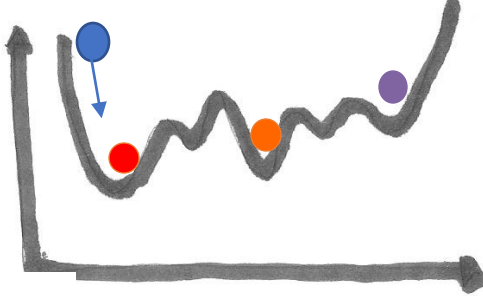
z: Sum of (P-1)(L-1) variables [-1,1]

$$\langle z \rangle = 0 \quad \sigma_z^2 = (P - 1)(L - 1)$$

# Hopfield model for auto-associative memories

$$s_i^{t+1} = \text{sign} \left( w_i^{\mu'} (L - 1) + z \right)$$

Starting pattern
Other patterns



$\langle z \rangle = 0 \quad \sigma_z^2 = (P - 1)(L - 1)$

$\sum_i w_i^{\mu} s_i$

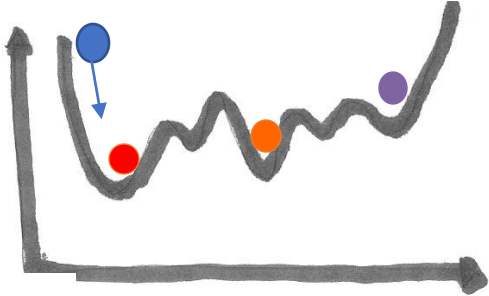
Starting from a pattern  $s^t = \mathbf{w}^{\mu'}$  (or nearby a pattern),  
 the pattern is a stable fixed point of the dynamics above if the number of stored patterns  
 $P < L$  (number of neurons)

# Hopfield model for auto-associative memories

$$s_i^{t+1} = \text{sign} \left( w_i^{\mu'} (L - 1) + z \right)$$

$\langle z \rangle = 0 \quad \sigma_z^2 = (P - 1)(L - 1)$

$\sum_i w_i^{\mu} s_i$



Starting from a pattern  $s^t = \mathbf{w}^{\mu'}$  (or nearby a pattern), the pattern is a stable fixed point of the dynamics above if the number of stored patterns  $P < L$  (number of neurons)

Proof:

Probability of an error:  $s_i^{t+1} \neq w_i^{\mu'}$

$$P_{\text{error}} \approx \int_{-\infty}^{\infty} \frac{dz}{\sqrt{L/P}} \frac{1}{2\pi LP} e^{-z^2/(2LP)} = \int_{-\infty}^{\infty} \frac{dz}{\sqrt{L/P}} \frac{1}{2\pi} e^{-z^2/2}$$

Ratio of number of synapses  
(on each neuron) and pattern to be stored

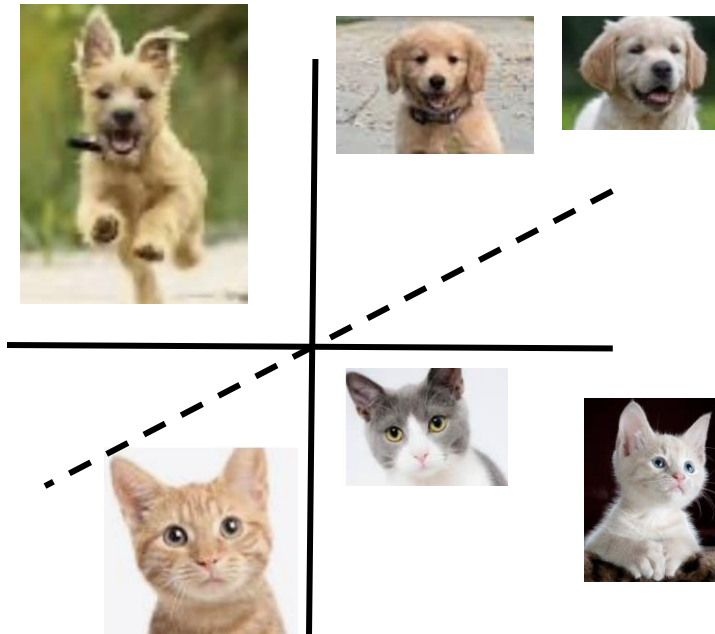
Example:  $P_{\text{error}} = 0.001$  for  $P < 0.1 L$



# Perceptron Algorithm

Up to now we have seen **Unsupervised** learning, where data are generic configuration of the system we want to model.

Now we will study **Supervised** learning where a label is associated to the data.



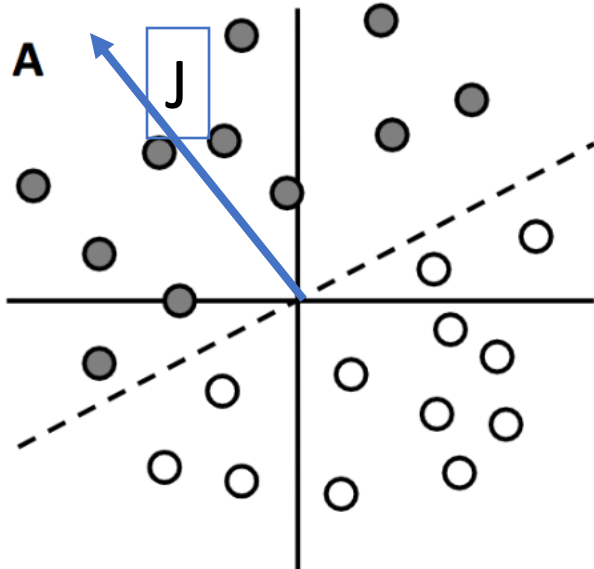
If the two classes are linearly separable

The Perceptron algorithm learn the rule to separate the class and

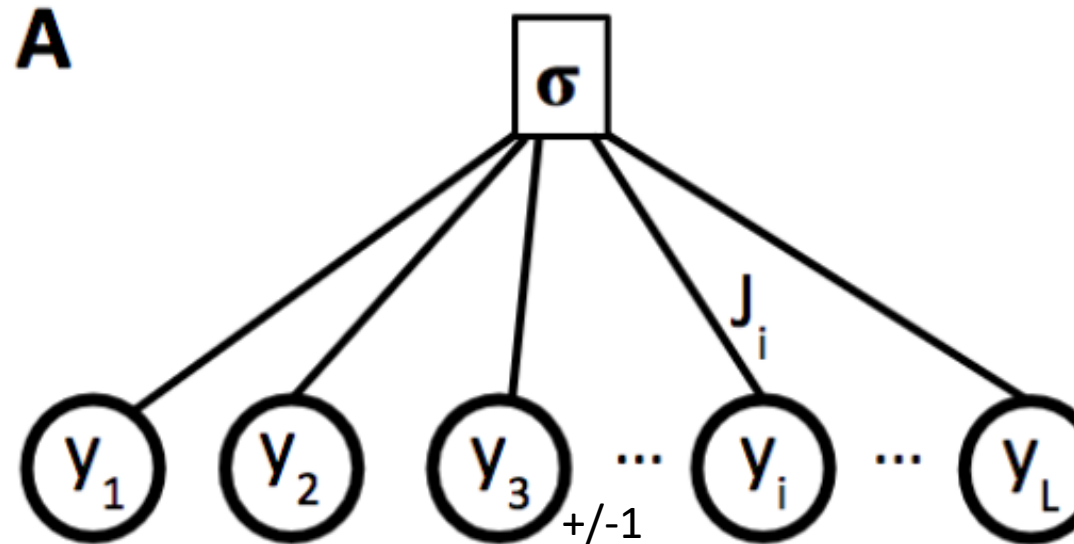
Learn to generalise the rule for the label.

# Perceptron Algorithm

The perceptron learns to generalise the rule for the label. It is also called Support Vector Machine.



Learn a vector  $\mathbf{J}$  which has a scalar product positif with a class and negative with the other



$$\sigma = \text{sign} \left( \sum_{i=1}^L J_i y_i - \theta \right)$$

Inspired from the behavior of a real neurons  
 $\Theta$  is an additional term to bias the output  $\sigma$ .

# Perceptron Algorithm

$$\sigma = \text{sign} \left( \sum_{i=1}^L J_i y_i - \theta \right)$$

To treat the bias  $\Theta$  it is convenient to add a component to the weight vector  $J_{L+1} = -\theta$  and to each input vector:  $y_{L+1} = 1$

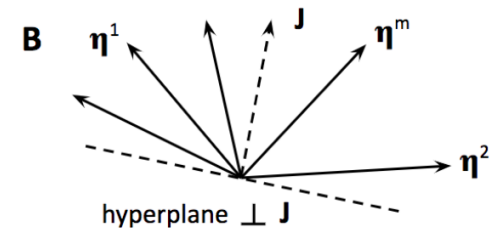
In a classification problem we have  $M$  examples ( $m=1..M$ ) and we look for the weight vector such that these examples are correctly classified:

$$\sigma^m = \text{sign}(\mathbf{J} \cdot \mathbf{y}^m), \quad \forall m = 1, \dots, M.$$

By defining:  $\boldsymbol{\eta}^m = \sigma^m \mathbf{y}^m$ .

We obtain the condition:

$$\mathbf{J} \cdot \boldsymbol{\eta}^m > 0, \quad \forall m = 1, \dots, M.$$



# Perceptron Algorithm

$$\sigma = \text{sign} \left( \sum_{i=1}^L J_i y_i - \theta \right)$$

To treat the bias  $\Theta$  it is convenient to add a component to the weight vector  $J_{L+1} = -\theta$  and to each input vector:  $y_{L+1} = 1$

In a classification problem we have  $M$  examples ( $m=1..M$ ) and we look for the weight vector such that these examples are correctly classified:

$$\sigma^m = \text{sign}(\mathbf{J} \cdot \mathbf{y}^m), \quad \forall m = 1, \dots, M.$$

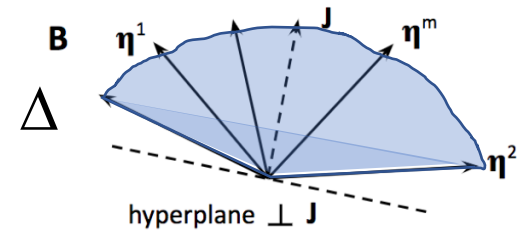
By defining:  $\boldsymbol{\eta}^m = \sigma^m \mathbf{y}^m$ .

We obtain the condition:

$$\mathbf{J} \cdot \boldsymbol{\eta}^m > 0, \quad \forall m = 1, \dots, M.$$

To make the perceptron-based classification more robust we write:

$$\mathbf{J} \cdot \boldsymbol{\eta}^m > \Delta, \quad \forall m = 1, \dots, M, \quad \Delta \text{ is called the stability or margin}$$



# Perceptron Algorithm

We suppose that a normalized weight vector  $\mathbf{J}^*$  satisfying the perceptron condition for any example exists. We define the stability as

$$\min_m \mathbf{J}^* \cdot \boldsymbol{\eta}^m = \Delta^* > 0$$

Under the normalization condition  $|\mathbf{J}^*| = 1$ .

A simple iterative algorithm guaranteed to output the vector  $\mathbf{J}^*$  is the following:

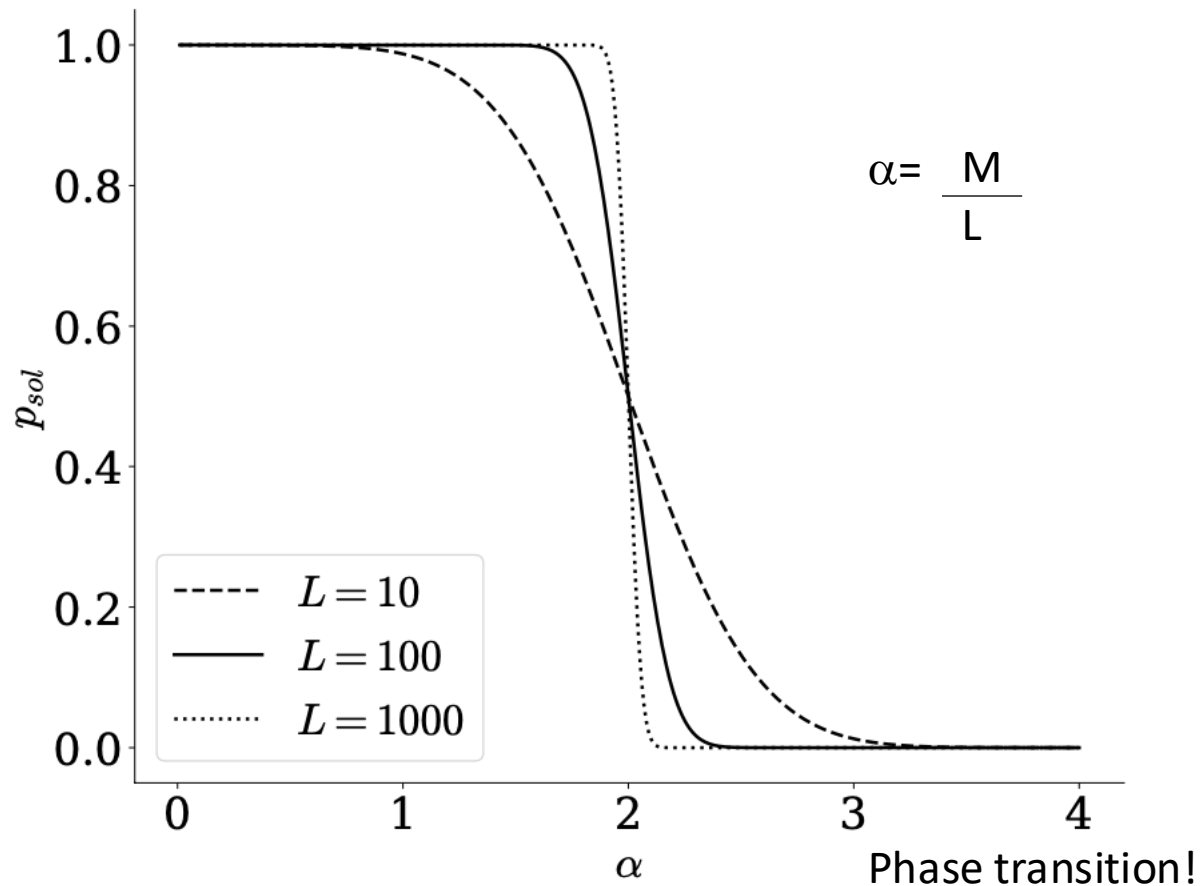
## Perceptron learning algorithm

- Set  $\mathbf{J}^{(0)} = \mathbf{0}$ ,  $t = 0$ , and let  $c$  be a strictly positive real number.
- Repeat as long as  $\min_m \mathbf{J}^{(t)} \cdot \boldsymbol{\eta}^m < c$ 
  - Define  $m_t = \operatorname{argmin}_m \mathbf{J}^{(t)} \cdot \boldsymbol{\eta}^m$ ;
  - Update  $\mathbf{J}^{(t+1)} \leftarrow \mathbf{J}^{(t)} + \boldsymbol{\eta}^{m_t}$ ,  $t \leftarrow t + 1$ .
- Output  $\mathbf{J} = \mathbf{J}^{(t)} / |\mathbf{J}^{(t)}|$ .

At each step  $t$ , the algorithm looks for the least stable pattern, i.e. having the smaller scalar product with  $\mathbf{J}^{(t)}$ , and pushes the weight vector in its direction at step  $t+1$ .

It stops when the stability with all patterns is larger than  $c$ .

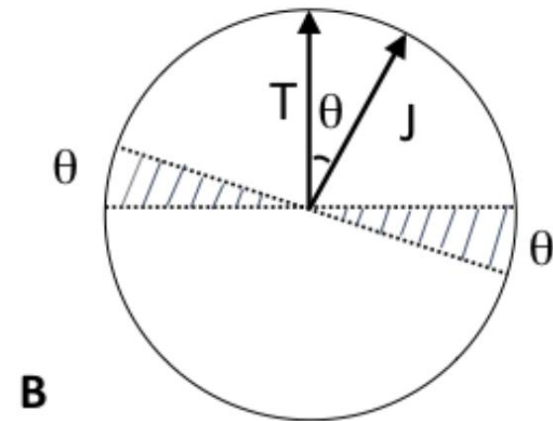
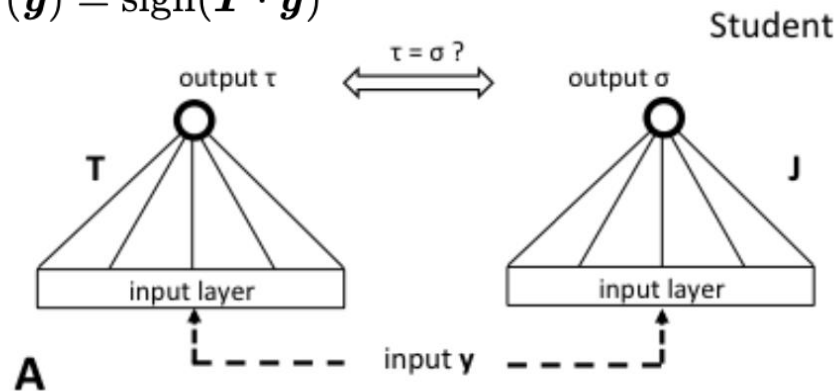
# Critical Capacity of the Perceptron



$$p_{sol}(\alpha) \equiv \lim_{L \rightarrow \infty} p_{sol}(M = \alpha L, L) = \begin{cases} 1 & \text{if } \alpha < 2, \\ 0 & \text{if } \alpha > 2. \end{cases}$$

# Generalization and Teacher Student Framework

$$\tau(\mathbf{y}) = \text{sign}(\mathbf{T} \cdot \mathbf{y})$$



Generalization error

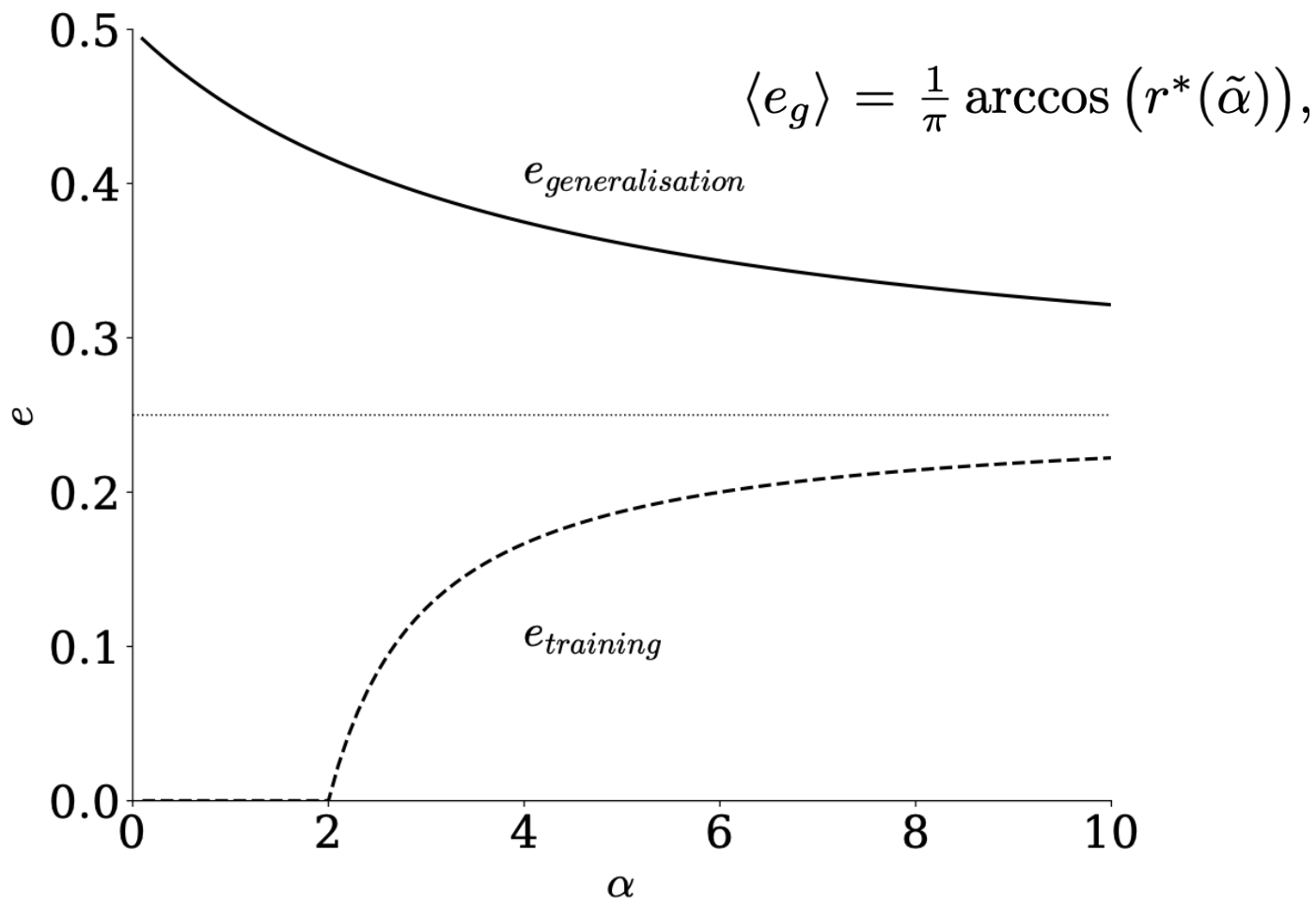
After training: test on a new example

$$e_g = (2\theta)/(2\pi)$$

$$e_g(\mathbf{J}|\mathbf{T}) = \text{Proba}((\mathbf{J} \cdot \mathbf{y}) \times (\mathbf{T} \cdot \mathbf{y}) < 0)$$

## Generalization error

The student weight vector  $J$  gets more and more aligned with the teacher weight vector  $T$  as the number of data increases





# Take Home Message

**Hopfield Model:** associative memory for  $P < L$

The perceptron as a classifier ,

**Algorithm to find the weight vector** if it does exists

- Algorithm to find the weight vectors in the linear case

**Does Weight Vector exists?**

- **Critical capacity**  $\alpha_c = M/L = 2$
- Use statistical mechanics to compute generalization error as a function of  $\alpha$   
It goes from

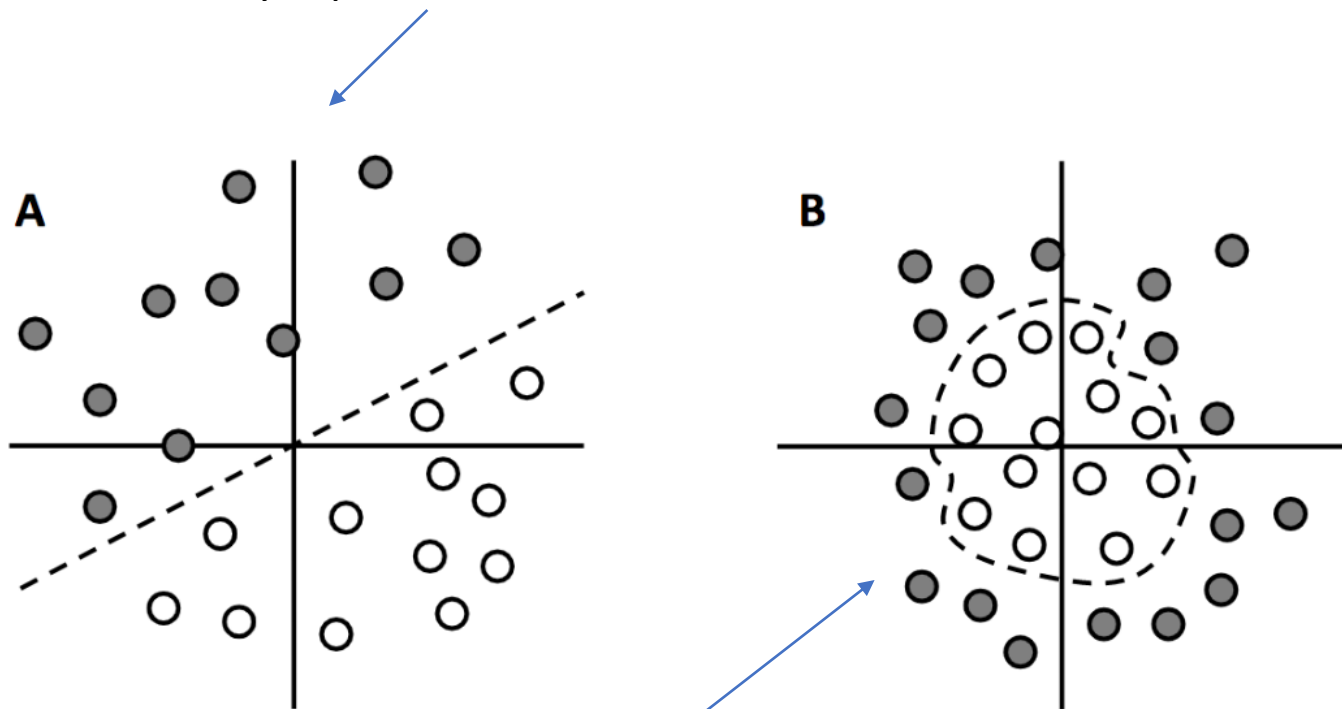
**Generalization Error**

Can be computed in teacher-student framework, decays as the number of examples increases

## Supplementary Slides

# Non linear classification and Kernel Method

The version of the perceptron we have studied above allows one to classify data points that are linearly separable



But sometimes they are not:

# Non linear classification and Kernel Method

To handle non-linear classification tasks one has to replace the decision for classification:

$$\sigma = \text{sign}(\mathbf{J} \cdot \mathbf{y})$$

by

$$\sigma = \text{sign}(\tilde{\mathbf{J}} \cdot \tilde{\phi}(\mathbf{y}))$$

Where  $\Phi$  is a non-linear function.

Eg.

$$\tilde{\phi}(\mathbf{y}) = (y_1, \dots, y_L, y_1 y_2, y_1 y_3, \dots, y_{L-1} y_L)$$

$$\tilde{\mathbf{J}} = (J_1, \dots, J_L, J_{1,2}, J_{1,3}, \dots, J_{L-1,L}).$$

$$\sigma = \text{sign}\left(\sum_i J_i y_i + \sum_{i < j} J_{i,j} y_i y_j\right)$$

Going to higher dimension  $D \gg L$  may make data linearly classifiable.