A Machine learning | Learning and Memory in the Brain (Les Gastins 2025)

1/ Feedforward network with L-Thiolden layers Z, ..., Ze-1  $Z^{l} = W^{l} \phi(Z^{l-1}) \in \mathbb{R}^{Ne}$   $Z^{1} = W^{1} \times \mathcal{L}$  nonlinearity Output  $f(x;W) = z^{L}(x)$ , input weights &W1,...,W3

XERNO WERNEN

 $(x) \xrightarrow{W^{1}} (x_{1}) \xrightarrow{\phi_{1} W_{2}} (z_{2}) \xrightarrow{\phi_{1} W_{3}} (z_{2}) \xrightarrow{\phi_{1} W_{2}} (z_{2}) = f(x_{1}W)$ input  $(x) \xrightarrow{\psi_{1} W_{2}} (z_{2}) \xrightarrow{\phi_{1} W_{3}} (z_{2}) \xrightarrow{\phi_{1} W_{2}} (z_{2}) = f(x_{1}W)$ 

/ training desired output network output

Defective function C(y, f(x; W)): How much does the network output f(x; W) for given input x differ from desired output y.

E.g.  $C(y,f(x;W)) = \|y - f(x;W)\|^2$  (MSE)

GANASI SEBEST

Given training data {xx, yxz, loss L(w):= 2 ((yx, f(x,w)))

Gradient descent

 $\omega' \leftarrow \omega' - \lambda \frac{\partial \mathcal{L}(\omega)}{\partial \omega_{ij}^{ij}} \Big|_{ij}$ learning rate

\*/ Backpropagation of the error : efficient recursive algo to compute 
$$\frac{\partial \mathcal{L}}{\partial \omega^e}$$
 for all  $l=1,...,L$ 

Needed:  $\frac{\partial \mathcal{L}}{\partial \omega^e} = \frac{\partial \mathcal{L}}{\partial t_i}$  (sum over repeated holices i, i.1...)

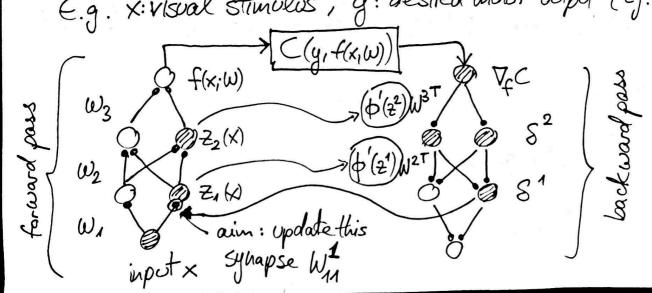
$$= \frac{\partial C}{\partial t_{i_{L}}} W_{i_{L}i_{L-1}}^{L} \phi'(z_{i_{L-1}}^{L-1}) \dots W_{i_{k+1}}^{l+1} \phi'(z_{i}^{l}) \phi(z_{j}^{l-1})$$

or 
$$S^{\ell-1} = \text{diag} \left[ \phi'(z_i^{\ell-1}) \right] (W^{\ell})^T S^{\ell}$$
 with  $S^{\ell} = V_{\ell} C$ 

$$\frac{\partial C}{\partial W^{\ell}} = S^{\ell} \phi(z^{\ell-1})^T \left( \text{and } \frac{\partial C}{\partial w^{\ell}} = S^{\ell} \times^T \right)$$

/ In the Brain: How could a synapse Wij. receive the error signal Si which aletermines how it should be updated to optimize the objective function  $C(y, f(x, \omega))$ ?

E.g. x: visual stimulus, y: desired motor octput (eg. move your orm)



Several issues:

• no nonlinearity in the backward pass:  $Z^{l} = W^{l} \phi(Z^{l-1})$  (forward) but St-1 = d'(zt-1) WITE (backward)

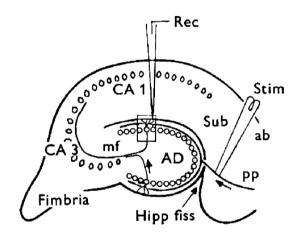
· Either forward and packward pass happen in the same network,
then synapses need to be symmetric Wij = Wii, but:

That's usually not the case and

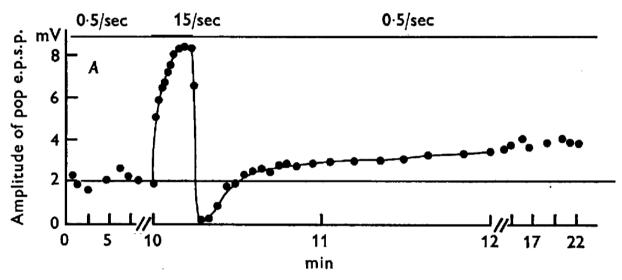
The network would need to precisely time and alternate forward and backward passes

- Or the backward pass happens in a separate network with exactly the same (but reversed) synapses. But then it would also need knowledge of  $\phi'(2^{l-1})$  and still time the alternating forward /backward passes forward /backward passes
- · (Not so clear what exactly the output in the brank would be a

## **Long Term Potentiation (LTP)**



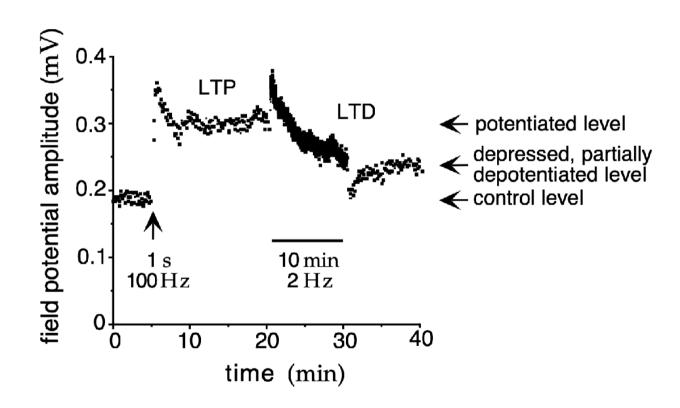
Perforant path fibres (pp) in hippocampus of rabbits anaesthetized with urethane.



Increasing the stimulation frequency from 5/sec to 15/sec during a short period (~2min) increases the neuron's response for a long period

Bliss and Lomo (1973)

## **Long Term Depression**



B) Neuroscience			
· Hebb (1949) conject Contributes to the fi is strengthened (	tores that if ing	oct from neuron A 3, then the synap	often ose A→B
· Experimental evider many stodies after	nce: Bliss & L	20mp (1973) au	nol
I theoretical models based on correlation	: Hebbian lear	ming /plasticity	umus Arina (6
	1 @ "	"firshe rate	."
Swild Vity	1 L	$S = \phi(\omega T)$ "unembrae	ne potential"
	pre post	(later)	e di e

 $\frac{dW_i}{dt} \propto r_i S$ (basic hebbian flamy)

/Problems

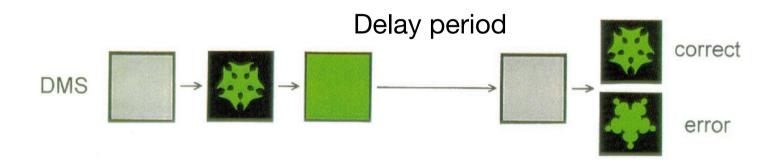
Weights grow unconstraint to the transfer of the standing of the transfer of the

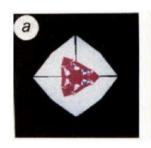
What look to Ash

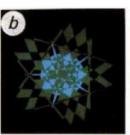
Chscpervised: How to optimize an objective function without acess to any error signal?
 Can this do interesting things? (Yes → Hapfield)

· Oja's modification of Helds rule: Normalize weights + small yelotesy  $w(t+1) = \frac{\omega(t) + \gamma r(t) s(t)}{||\omega(t)| + \gamma r(t) s(t)||} \approx \omega(t) + \gamma (sr - \omega s^2)$   $||\omega(t)| + \gamma r(t) s(t)||$   $|(\omega(t)| + \gamma$ C transfer function: Hembrane pot -> Firsty rate Take away: Weights w converge to be eigenvectors of C:= < rr } (correlation matrix of the input to neurons) with largest eignalue, i.e. a becomes the 1st PC. (principle component). Pcoof: Note that IW(t) 11 = 1 Ht (Def. of the update rule) w(t+1) & w(t) + g 3w(t+1)/0 30(41) = Sr - 2 3 11 w + grs 11  $= sr - \omega \frac{(\omega + \gamma rs)^{T} rs}{|\omega + \gamma rs||} = sr - \omega \omega^{T} rs$  $\Delta \omega(t) = \gamma(sr - \omega s^2) + substitute s = r \omega$  $\langle \Delta \omega t \theta \rangle_{t} = 0 = \gamma \langle \langle r r^{\dagger} \rangle \omega - (\omega^{\dagger} \langle r r^{\dagger} \rangle \omega) \omega$  $\Leftrightarrow$   $C\omega = (\omega T C\omega) \omega$ w is evec. of C with eval wTCw. One can show that the dynamic converges to w st. wTCw is maximal I

## **Delay Match to Sample Tasks**



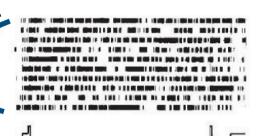


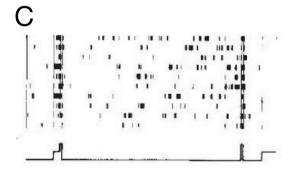


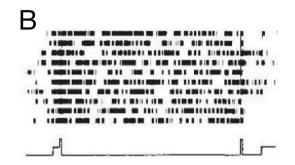


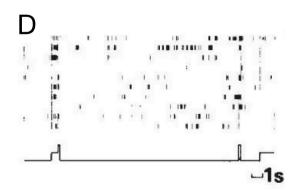


Same neuron, (different trials



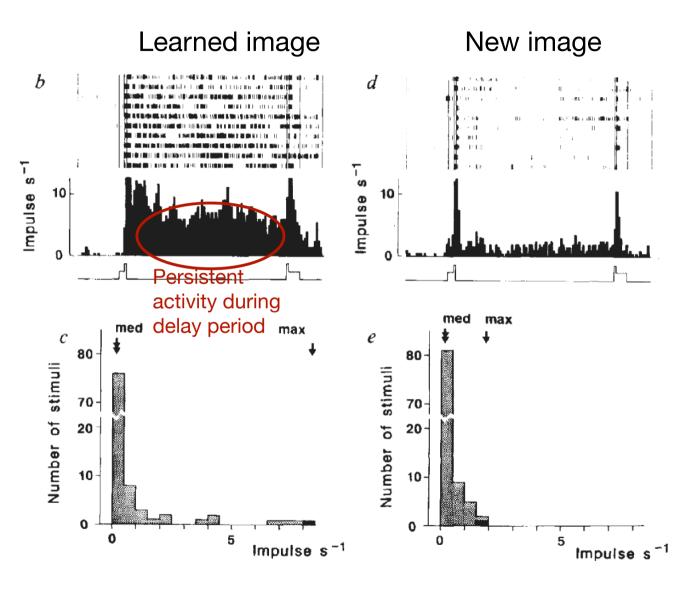






Y. Miyashita, H.S. Chang,

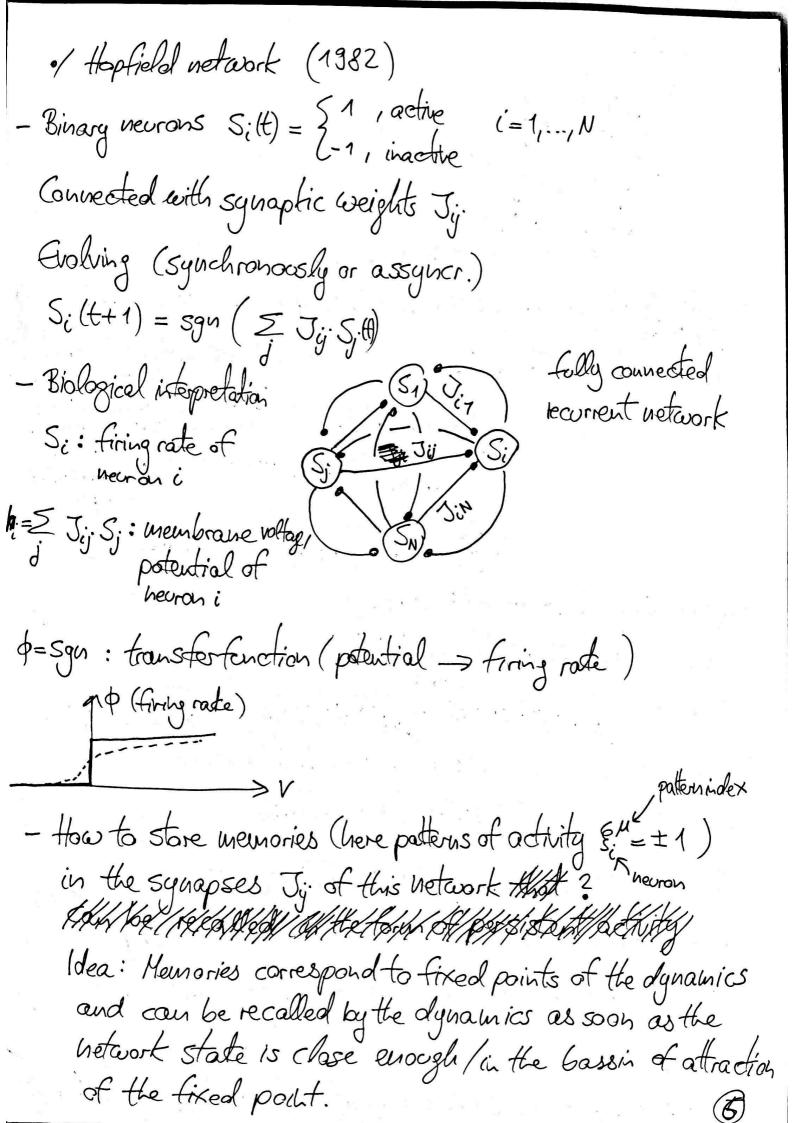
"Neuronal correlate of pictorial short-term memory in the primate temporal cortex" 1988

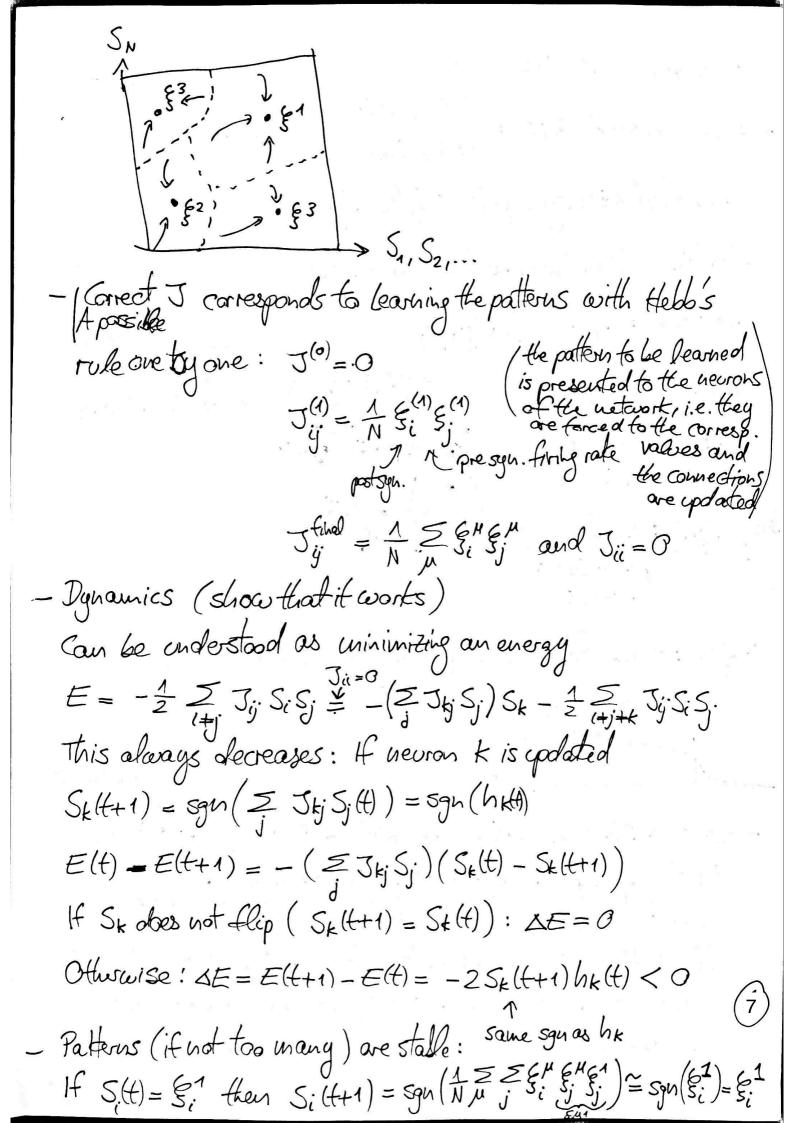


Y Miyashita, "Neuronal correlate of visual associative long-term memory in the primate temporal cortex" (1988)

(b,d) Trial averaged activity of the neuron that is most active for the given image

(c,d) Distribution of firing rates (Impulse/s) during delay period for a given neuron and 100 images





- max number of patterns" that can be stored and perfectly retrieved with N mans Why? Study SNR Assume  $Sit=S_i^1$ . Then  $Si(t+1) = sgn(h_i(t))$  $h_i(t) = \sum_{j \neq i} J_{ij} g_{j}^{1} = g_{i}^{1} + \frac{1}{N} \sum_{j \neq i} \sum_{M=2}^{p} g_{i}^{M} g_{j}^{M} g_{j}^{1}$ Signal = Si uoiso Signal = Si uoiso  $\langle S_i \rangle = 0$  $\langle S_i^2 \rangle \simeq \frac{\gamma}{N} = \chi^{\prime\prime} \rightarrow SNR = \frac{1}{\sqrt{\chi}}$  $g = Prob \left( \text{Noise } S_i < \text{Signal} = 1 \right) = \int_{-\sqrt{1/\kappa'}}^{\infty} \frac{dx}{\sqrt{2\pi}} e^{-x/2} \left( S_i \text{ is approx Gaussian} \right)$ with std  $\sqrt{\kappa}$ (write  $S_i = \sqrt{\alpha} \times \text{ with } \times_{\sim} \mathcal{N}(0,1)$ , then  $\Re(S < 1) = \Re(X < \frac{1}{\sqrt{\alpha}})$  $= \Re\left(X > -\frac{1}{\sqrt{x}}\right) /$  $g \approx 1 - \exp(-\frac{1}{2\kappa})(\sqrt{\alpha} + \mathcal{O}(\alpha^{3/2}))$ Prob(Si<1 \ti)=g" ~ exp(-Nvx exp(-\frac{1}{2x})) € <> 2 (un) : 5 > 0 - if not potect retrieval is required:  $\alpha \sim O(1)$ 

- Critiques of Hopfield's model

· No separation between exitatory and inhibitory neurons (Dale's law: All synapses of a given neuron are either whilitory or exhibitory)

• Electrophys. recordings show rather a sparse cooling of memories (most  $\xi_i^{M} = 0$ )

• Symmetric connections are not biologically plausible (Sompolinsty showed that  $J_{ij} + \gamma_{ij}$  is robust where  $\gamma_{ij} \sim W(0, \Delta)$  with  $\Delta \sim O(1)$ 

· Black-out catastrophe: All paterus are lost if one goes over the capacity. Numerous modifications have been proposed - Either: stop learning by some supervised signal when capacity is reached

- Or: modity st. continual learning is possible where older patterns get forgotten as newer ones are learned.

 $\frac{2\pi}{||S_1||_{L^{\infty}}} = (1+\pi) \tilde{\omega} = 2\pi \int_{\mathbb{R}^2} \frac{2\pi}{||S_1||_{L^{\infty}}} = (1+\pi) \tilde{\omega} = 2\pi \int_{\mathbb{R}^2} \frac{2\pi}{||S_1||_{L$