

# Learning

## Lecture 8

# Examples of learning

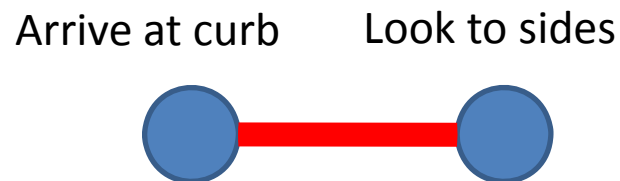
- Learning to play the piano
- Learning to win in chess
- Learning to cross the street safely
- Learning to catch prey
- Learning theoretical neuroscience
- The brain developing its representation of objects

# Three types of learning

- Supervised learning: using feedback, learn to produce correct output given new input
- Reinforcement learning: learn to act in a way that maximizes future rewards
- Unsupervised learning: finding structure in data

# Hebbian learning

- How neural networks change in response to input: activity-dependent synaptic plasticity
- Neural basis of learning and memory
- “Neurons that fire together, wire together”



- Can be supervised or unsupervised

Presynaptic activity:  $\mathbf{u}$

Postsynaptic activity:  $v$

Weights:  $\mathbf{w}$

Simple Hebb rule:  $\tau \frac{d\mathbf{w}}{dt} = v\mathbf{u}$

Averaged Hebb rule:  $\tau \frac{d\mathbf{w}}{dt} = \langle v\mathbf{u} \rangle$

average over ensemble of inputs

Unstable! Constraint needed to prevent weights from growing indefinitely.

Simple case:  $v = \mathbf{w} \cdot \mathbf{u}$


Averaged Hebb rule:  $\tau \frac{d\mathbf{w}}{dt} = \langle \mathbf{u}\mathbf{u} \rangle \mathbf{w} = \mathbf{Q} \cdot \mathbf{w}$

Correlation-based rule

With subtractive normalization and constraint  
→ can predict ocular dominance columns



# Supervised Hebbian learning

$$\tau \frac{d\mathbf{w}}{dt} = \langle v_s \mathbf{u}_s \rangle$$


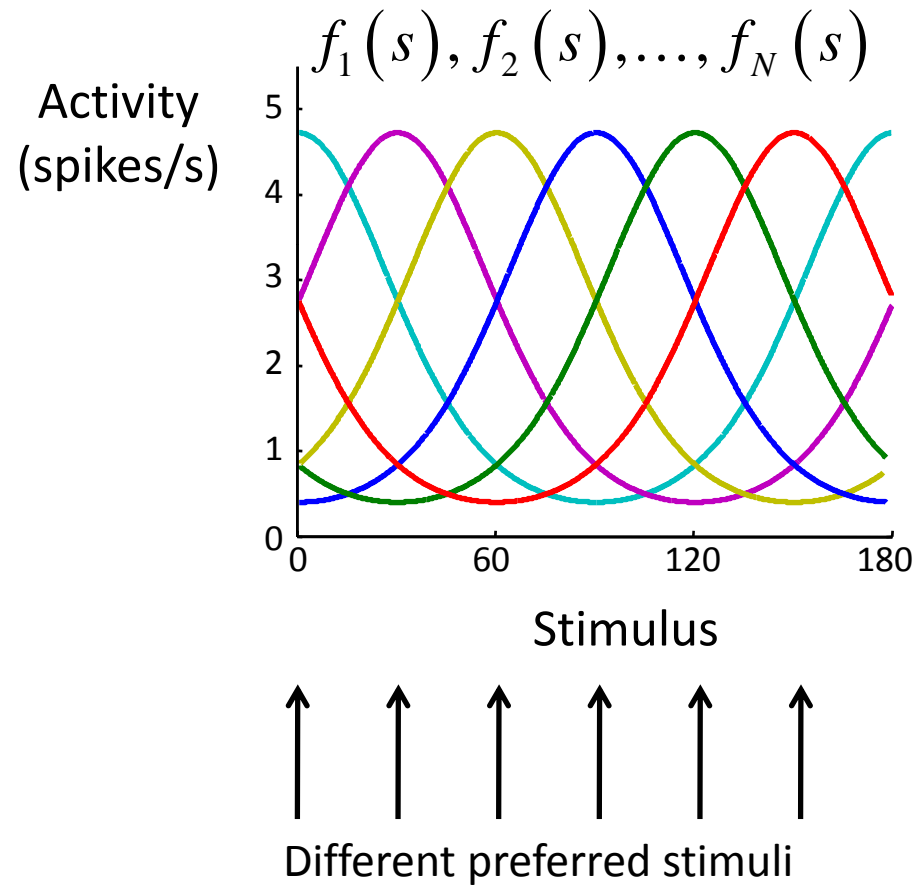
Paired samples

With decay:  $\tau \frac{d\mathbf{w}}{dt} = \langle v_s \mathbf{u}_s \rangle - \alpha \mathbf{w}$

Steady state:  $\mathbf{w} = \frac{1}{\alpha} \langle v_s \mathbf{u}_s \rangle$

Weights proportional to input-output cross-correlation.

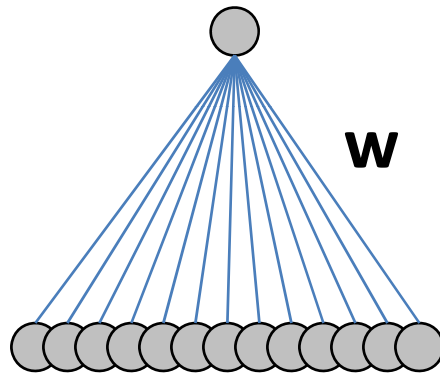
# Neural networks for function approximation





# Computing a new function

Output neuron  $v$ , should have tuning  $h(s)$



Tuning curves  $\mathbf{u} = \mathbf{f}(s)$

$$v = \mathbf{w} \cdot \mathbf{u} = \mathbf{w} \cdot \mathbf{f}(s)$$

# Learning rule

$$E = \left\langle \left( h(s_s) - \mathbf{w} \cdot \mathbf{f}(s_s) \right)^2 \right\rangle_{\text{training data}}$$

Gradient descent:  $\mathbf{w} \rightarrow \mathbf{w} - \varepsilon \nabla_{\mathbf{w}} E$

$$\mathbf{w} \rightarrow \mathbf{w} + \varepsilon \left\langle \left( h(s_s) - v(s_s) \right) \mathbf{f}(s) \right\rangle_{\text{training data}}$$

Stochastic gradient descent (delta rule):

$$\mathbf{w} \rightarrow \mathbf{w} + \varepsilon \left( h(s_s) - v(s_s) \right) \mathbf{f}(s)$$

# Representational learning

- How do neurons acquire their response selectivities?
- Natural images are richly structured and highly constrained.
- System learns statistical structure of visual images and builds a model to reproduce structure → *generative model*
- Use this to identify objects in particular images → *recognition model*

# Types of representational learning

- Mixture of Gaussians
- Factor analysis
- Principal component analysis
- Independent component analysis
- Sparse coding
- Helmholtz machine

# Sparse coding

- Inference on retinal images
- Model images as linear superposition of basis functions
- Sparseness: simple representation, minimize interference between different patterns of input, save energy
- Statistically independent: reduces redundancy

# Image model

$$I(x, y) = \sum_i a_i \varphi_i(x, y) + \eta(x, y)$$

The diagram shows the equation  $I(x, y) = \sum_i a_i \varphi_i(x, y) + \eta(x, y)$  with four arrows pointing from labels below to terms in the equation:

- An arrow points from the label "image" to the  $I(x, y)$  term.
- An arrow points from the label "coefficients (output activities)" to the  $a_i$  term.
- An arrow points from the label "Basis functions" to the  $\varphi_i(x, y)$  term.
- An arrow points from the label "Noise with variance  $\sigma^2$ " to the  $\eta(x, y)$  term.

Infinite number of solutions for  $\mathbf{a}$ ,  $\Phi$ , when basis set is *overcomplete* (more output units than input units)

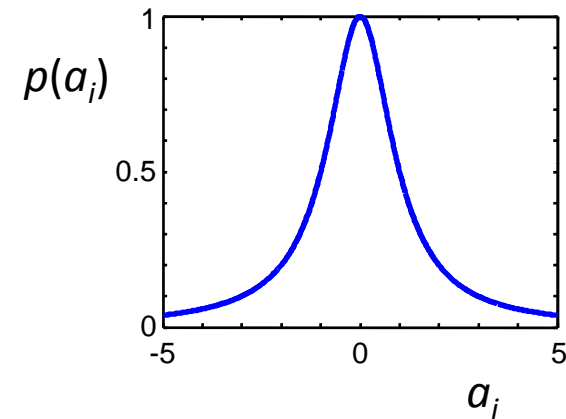
# Questions

- When a given image  $I$  is presented, what should output activity,  $a$ , be? (recognition model)
- Across all images, what is the best choice of basis functions?

# Prior over output activities

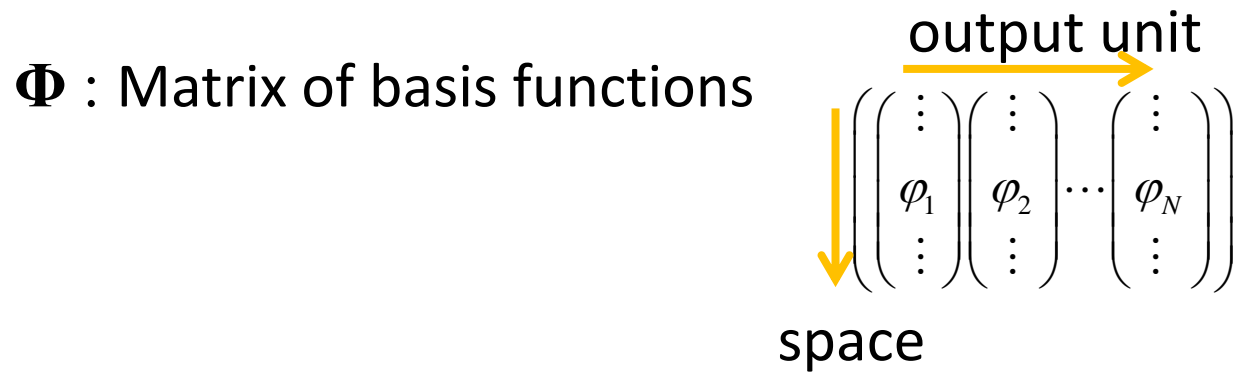
$$p(\mathbf{a}) = \prod_i p(a_i)$$

$$p(a_i) \propto e^{-S(a_i)}$$



Small coefficients are favored: sparseness





Assume  $\Phi$  fixed and known

Posterior over coefficients based on a given image  $\mathbf{I}$ :

$$p(\mathbf{a} | \mathbf{I}, \Phi) \propto p(\mathbf{I} | \mathbf{a}, \Phi) p(\mathbf{a} | \Phi)$$

$$p(\mathbf{I} | \mathbf{a}, \Phi) \propto e^{-\frac{\|\mathbf{I} - \Phi \mathbf{a}\|^2}{2\sigma^2}}$$

$$p(\mathbf{a} | \Phi) \propto \prod_i e^{-S(a_i)}$$

Best possible coefficients:

$$\begin{aligned}\hat{\mathbf{a}} &= \operatorname{argmax}_{\mathbf{a}} p(\mathbf{a} | \mathbf{I}, \Phi) \\ &= \operatorname{argmax}_{\mathbf{a}} \log p(\mathbf{a} | \mathbf{I}, \Phi) \\ &= \operatorname{argmin}_{\mathbf{a}} \left( \frac{\|\mathbf{I} - \Phi \mathbf{a}\|^2}{2\sigma^2} + \sum_i S(a_i) \right)\end{aligned}$$

Learning the coefficients through gradient descent:

$$\begin{aligned}\Delta \mathbf{a} &\propto -\nabla_{\mathbf{a}} \left( \frac{\|\mathbf{I} - \Phi \mathbf{a}\|^2}{2\sigma^2} + \sum_i S(a_i) \right) \\ &= \frac{1}{\sigma^2} \Phi^T \underbrace{(\mathbf{I} - \Phi \mathbf{a})}_{\text{residual image}} - S'(\mathbf{a})\end{aligned}$$

Can be implemented in recurrent neural network

# Learning the basis functions

- So far: fixed  $\Phi$ , learned coefficients  $\mathbf{a}$
- What about different  $\Phi$ ?
- Maximize average log likelihood of parameters (minimize KL distance)

Objective function for learning: log likelihood of model  $\Phi$ :

$$L = \langle \log p(\mathbf{I} | \Phi) \rangle$$

Average over input images

$$p(\mathbf{I} | \Phi) = \int p(\mathbf{I} | \mathbf{a}, \Phi) p(\mathbf{a} | \Phi) d\mathbf{a}$$

Gradient descent on  $\Phi$ :

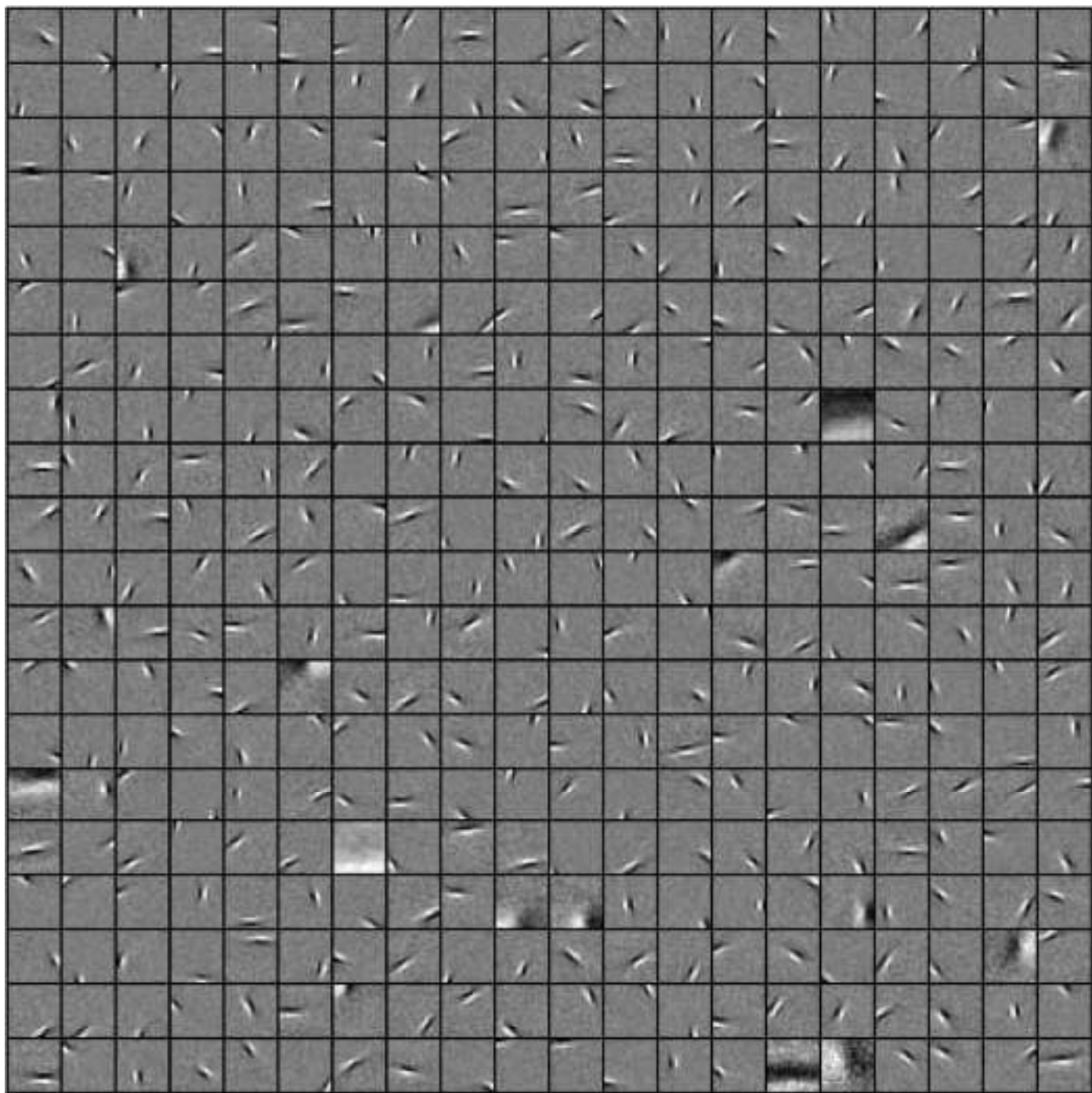
$$\Delta \Phi \propto \frac{\partial L}{\partial \Phi} = \frac{1}{\sigma^2} \left\langle \left\langle (\mathbf{I} - \Phi \mathbf{a}) \mathbf{a}^T \right\rangle_{p(\mathbf{a} | \mathbf{I}, \Phi)} \right\rangle$$

Hebbian learning

Approximate by posterior maximum:

$$\Delta \Phi \propto \left\langle (\mathbf{I} - \Phi \hat{\mathbf{a}}) \hat{\mathbf{a}}^T \right\rangle$$

Constraint needed to prevent growth without bound



## **Exercise**

Reproduce this. Using sparse coding, learn Gabor-like basis functions from any set of photos. Make assumptions where necessary.

**Due April 26 by email**

# Expectation maximization

- Objective function with two parameter sets:

$$F = \langle \log p(\mathbf{a}, \Phi; \mathbf{I}) \rangle$$

(free energy)

- Step 1: fix  $\Phi$ , find  $\mathbf{a}(\mathbf{I})$  (expectation)
- Step 2: fix  $\mathbf{a}$ , optimize  $\Phi$  (maximization)
- Repeat.
- Converges to local maximum

# References

- Olshausen and Field, *Emergence of simple-cell receptive field properties by learning a sparse code for natural images*. Nature 381, 607-9
- Olshausen, *Sparse codes and spikes*. In *Probabilistic models of the brain* (MIT Press, 2002)
- Chapter 10 in Dayan and Abbott, *Theoretical Neuroscience* (MIT Press, 2001)



# Lectures so far

- Neural population coding; how to decode
- Role of correlations in information processing
- Perception as Bayesian inference
- Cue combination
- Bayesian models of behavioral data
- Bayesian model comparison
- Neural implementation of Bayesian inference
- Models of perceptual decision-making
- Representational learning; sparse coding
- **Thursday: optimal inference in higher-level cognition**