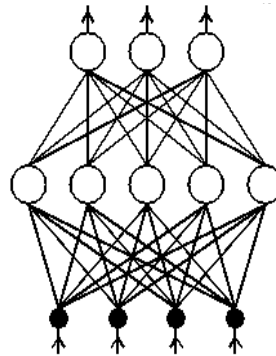
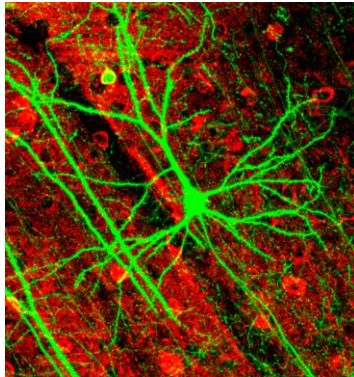


CSE/NEURO  
528  
Lecture 12:  
Unsupervised  
to Supervised  
Learning  
(Chapters 8 & 10)



1

## Last Time: A Linear Model for Natural Images

Input image  $\mathbf{u}$

$$\begin{array}{c}
 \mathbf{g}_1 \qquad \mathbf{g}_2 \qquad \mathbf{g}_3 \\
 \text{[?]} \cdot v_1 + \text{[?]} \cdot v_2 + \text{[?]} \cdot v_3 + \dots
 \end{array}$$

$$\mathbf{u} = \sum_{i=1}^M \mathbf{g}_i v_i + \text{noise}$$

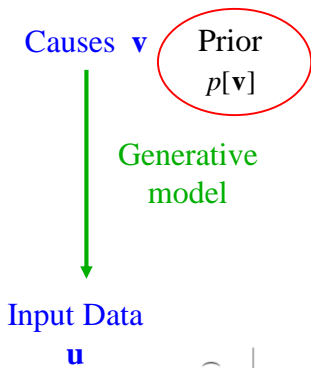
$$= G\mathbf{v} + \text{noise}$$

(Note:  $M$  can be larger than  $N$ )

$G$  = matrix whose columns are  $\mathbf{g}_i$   
 $\mathbf{v}$  = vector whose elements are  $v_i$

2

## Defining the Generative Model: Prior



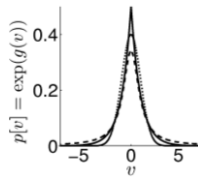
Assuming causes  $v_i$  are independent:

$$p[\mathbf{v}] = \prod_i p[v_i]$$

$$\log p[\mathbf{v}] = \sum_i \log p[v_i; G]$$

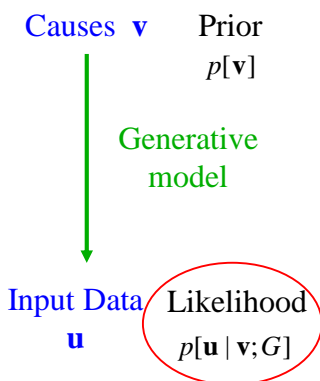
For any input, we want only a few causes  $v_i$  to be active

- $v_i = 0$  most of the time but high for some inputs
- Suggests **sparse distribution** for  $p[v_i]$ : peak at 0 but with heavy tail (also called *super-Gaussian* distribution)



3

## Defining the Generative Model: Likelihood



Linear model:

$$\mathbf{u} = G\mathbf{v} + \text{noise}$$

Assume *noise* is Gaussian white noise:

$$p[\mathbf{u} | \mathbf{v}; G] = \text{Gaussian}(\mathbf{u}; G\mathbf{v}, I)$$

$$\propto \exp\left(-\frac{1}{2} \|\mathbf{u} - G\mathbf{v}\|^2\right)$$

Log likelihood:

$$\log p[\mathbf{u} | \mathbf{v}; G] = -\frac{1}{2} \|\mathbf{u} - G\mathbf{v}\|^2 + C$$

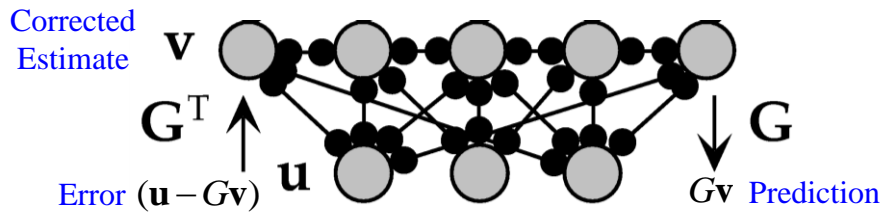
$$F(\mathbf{v}, G) = \log p[\mathbf{u} | \mathbf{v}; G] + \log p[\mathbf{v}; G] + \log k$$

Find  $\mathbf{v}$  and  $G$  that maximize log of posterior probability

$$= -\frac{1}{2} \|\mathbf{u} - G\mathbf{v}\|^2 + \sum_i g(v_i) + K$$

4

## Recurrent Network for Sparse Coding

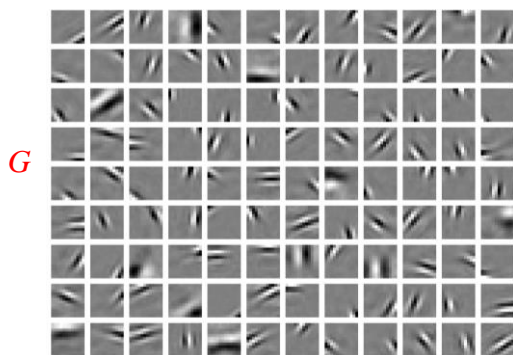


Firing rate dynamics  $\tau \frac{d\mathbf{v}}{dt} = G^T (\mathbf{u} - G\mathbf{v}) + g'(\mathbf{v})$

Learning rule  $\tau_G \frac{dG}{dt} = (\mathbf{u} - G\mathbf{v})\mathbf{v}^T$  *Hebbian* *Sparseness penalty*

5

## Result: Sparse Coding of Natural Images



Each square is a column  $\mathbf{g}_i$  of  $G$  (obtained by collapsing rows of the square into a vector)

The  $\mathbf{g}_i$  look like local edge or bar features similar to receptive fields in primary visual cortex (V1)

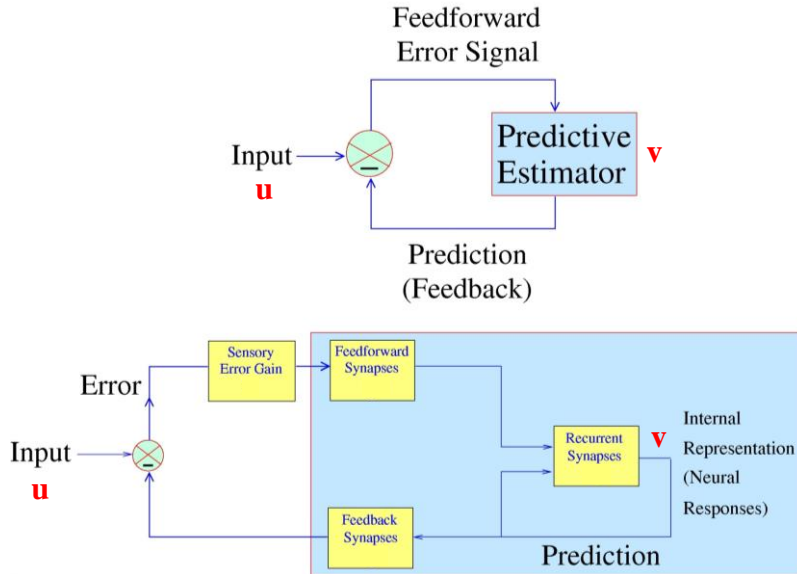
(Olshausen & Field, 1996)

Image patch  $\mathbf{u}$

$$\mathbf{u} = \mathbf{g}_1 \cdot v_1 + \mathbf{g}_2 \cdot v_2 + \mathbf{g}_3 \cdot v_3 + \dots$$

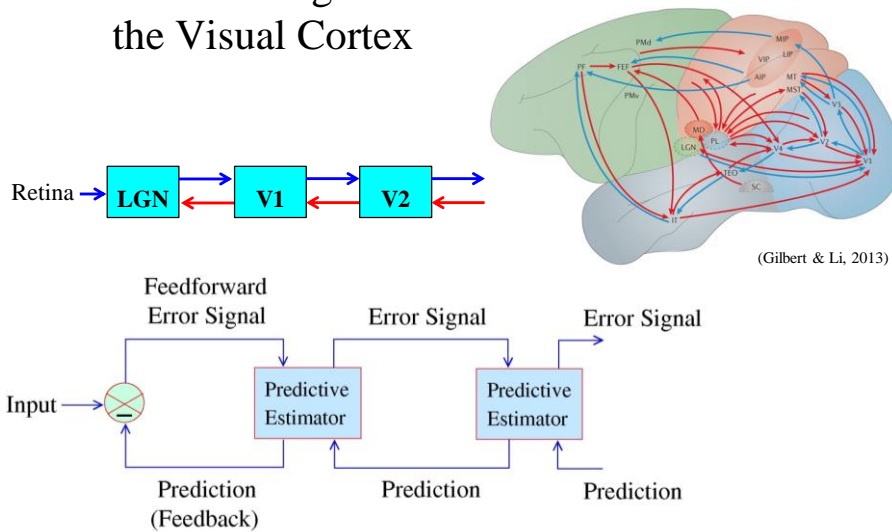
6

## Sparse Coding is a special case of Predictive Coding



For more details, see: (Rao, *Vision Research*, 1999; Rao & Ballard, *Nature Neurosci.*, 1999)

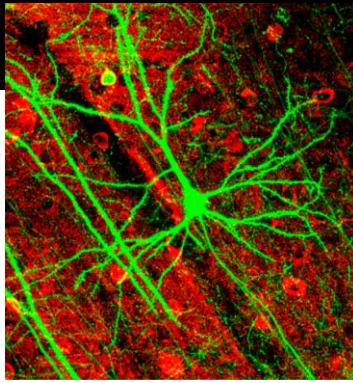
## Predictive Coding Model of the Visual Cortex



(Gilbert & Li, 2013)

(Rao & Ballard, *Nature Neurosci.*, 1999)

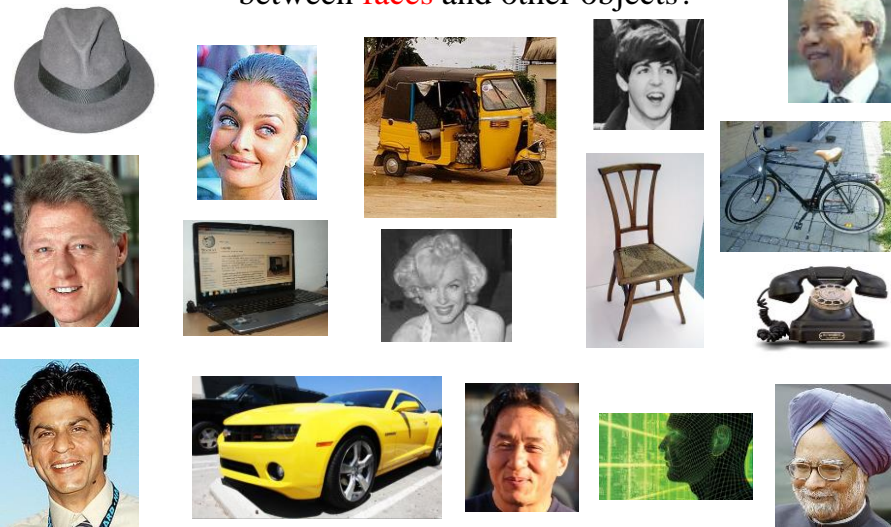
# From Unsupervised to Supervised Learning: Neurons as Classifiers



9

## The Classification Problem

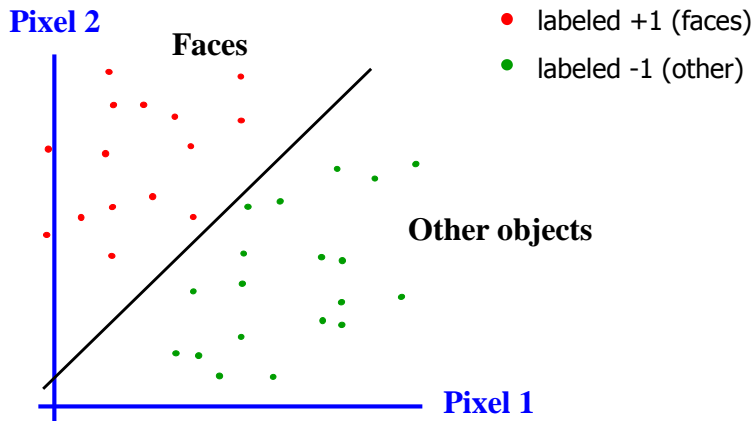
How do we build a classifier to distinguish  
between **faces** and other objects?



10

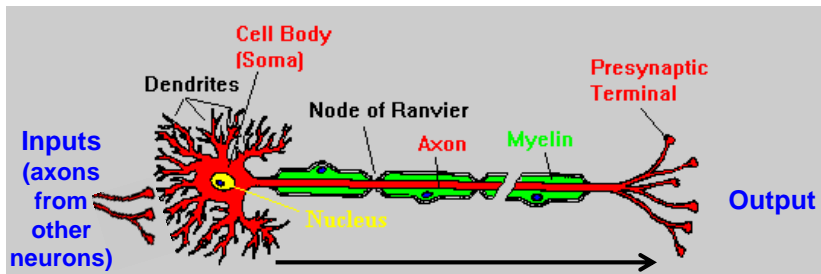
Image Source: Wikimedia Commons

# The Classification Problem



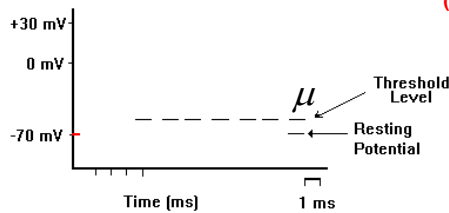
Idea: Find a separating hyperplane (line in this case)  
**Can neurons do that?**

# The Idealized Neuron

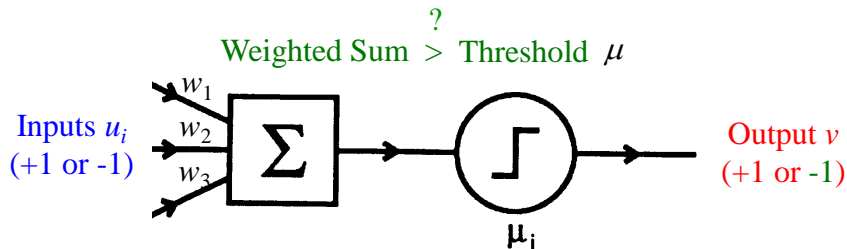


$$u_i \quad w_i \quad \sum_i w_i u_i > \mu?$$

Spike or no spike  
(1 or -1)



## The “Perceptron”



$$v = \Theta\left(\sum_i w_i u_i - \mu\right)$$

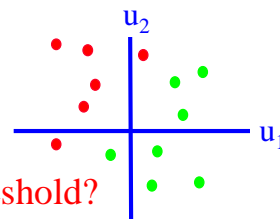
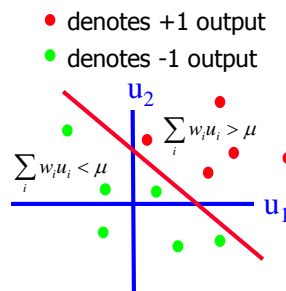
$$\Theta(x) = +1 \text{ if } x > 0 \text{ and } -1 \text{ if } x \leq 0$$

[Introduced by Rosenblatt (1958) building on McCulloch and Pitts (1943)] 13

## What does a Perceptron do?

- ◆ Weighted sum defines a *hyperplane (line, plane, ...)*

$$\sum_i w_i u_i - \mu = 0$$
- ◆ All inputs *on one side* of hyperplane have output = +1 (“class 1”); all inputs *on other side* have output = -1 (“class 2”)
- ◆ Perceptrons can classify!
  - ⇒ Can perform linear classification



How do we learn the weights and threshold?

## Perceptron Learning Rule

---

Given input  $\mathbf{u}$ , output  $v = \Theta(\sum_i w_i u_i - \mu)$ , and **desired output**  $v^d$ :

Adjust  $w_i$  and  $\mu$  according to **output error** ( $v^d - v$ ):

$$\Delta w_i = \varepsilon(v^d - v)u_i$$

$$\begin{array}{cc} +1 & -1 \\ -1 & +1 \end{array}$$

For positive input ( $u_i = +1$ ):

Increases weight if error is positive

Decreases weight if error is negative

(opposite for  $u_i = -1$ )

$$\Delta \mu = -\varepsilon(v^d - v)$$

$$\begin{array}{cc} +1 & -1 \\ -1 & +1 \end{array}$$

Decreases threshold if error is positive

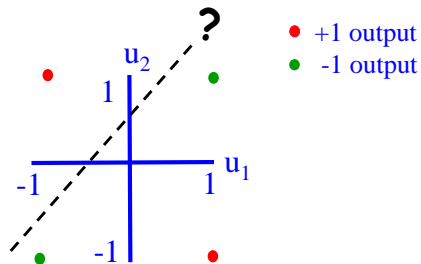
Increases threshold if error is negative

15

## Can Perceptrons learn any function?

---

$u_1$	$u_2$	XOR
-1	-1	-1
1	-1	+1
-1	1	+1
1	1	-1



Perceptrons can only classify **linearly separable** data  
How do we handle linear inseparability?

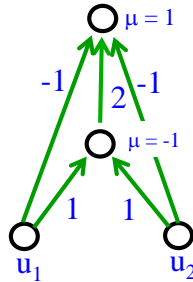
16



## Multilayer Perceptrons

---

- ◆ Can classify linearly inseparable data
  - ⇨ Can solve XOR
- ◆ An example of a two-layer perceptron that computes XOR



(Inputs and outputs are +1 or -1)

17

---

What if you want *continuous* outputs rather than +1/-1 outputs (i.e., regression)?

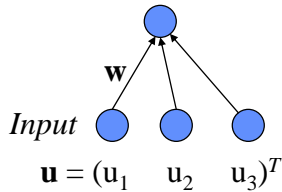


E.g., Teaching a network to drive a truck

18  
Image Source: Wikimedia Commons

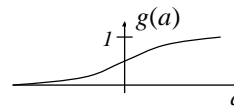
## Continuous Outputs with Sigmoid Networks

Output  $v = g(\mathbf{w}^T \mathbf{u}) = g(\sum_i w_i u_i)$

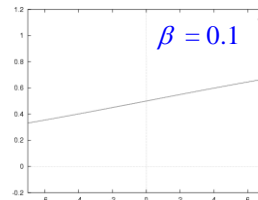
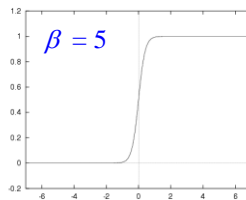


Sigmoid output function:

$$g(a) = \frac{1}{1 + e^{-\beta a}}$$



Parameter  $\beta$  controls the slope

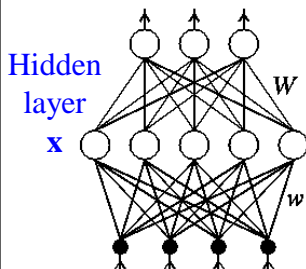


19

## Learning Multilayer Sigmoid Networks

$v_i = g(\sum_j W_{ij} g(\sum_k w_{jk} u_k)) = g(\sum_j W_{ij} x_j)$  Desired output  $\mathbf{d}$  also given

Output  $\mathbf{v} = (v_1 \ v_2 \ \dots \ v_J)^T$



Input  $\mathbf{u} = (u_1 \ u_2 \ \dots \ u_K)^T$

Learn weights that minimize output error:

$$E(\mathbf{W}, \mathbf{w}) = \frac{1}{2} \sum_i (d_i - v_i)^2$$

Use gradient descent!

$$\Delta W_{ij} = -\varepsilon \frac{dE}{dW_{ij}} = \varepsilon \cdot (d_i - v_i) g'(\sum_j W_{ij} x_j) x_j$$

Delta rule

$$\Delta w_{jk} = -\varepsilon \frac{dE}{dw_{jk}} = -\varepsilon \frac{dE}{dx_j} \frac{dx_j}{dw_{jk}} \quad \text{Backpropagation learning rule}$$

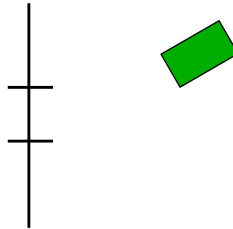
(see Supplementary Materials for details)

## Example: Backing up a Truck (courtesy of Keith Grochow)



### Teaching a Network to back a truck into a loading dock

- Input:  $x$ ,  $y$ ,  $\theta$  of truck
- Output: Steering angle



21

Forget trucks  
doin' supervised  
learning – how do  
I find some food  
in a barn?



Enter... Reinforcement Learning

Image Source: Wikimedia Commons

## Next Class: Reinforcement Learning

---

