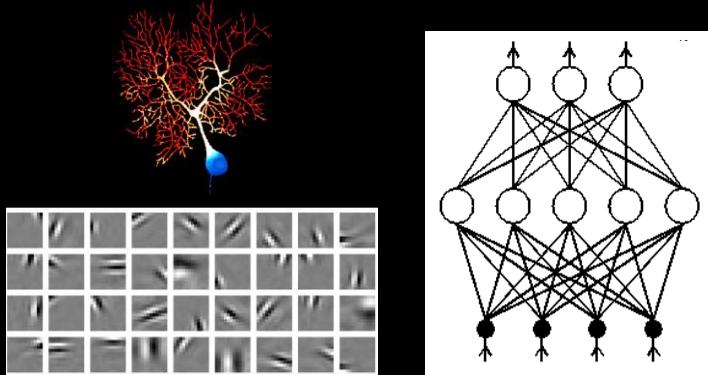


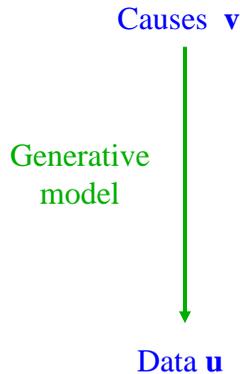
CSE/NB 528
Lecture 13: From Unsupervised to
Reinforcement Learning
(Chapters 8-10)



Today's Agenda: All about Learning

- ◆ Unsupervised Learning
 - ⇒ Sparse Coding
 - ⇒ Predictive Coding
- ◆ Supervised learning
 - ⇒ Perceptrons and Backpropagation
- ◆ Reinforcement Learning
 - ⇒ TD and Actor-Critic learning

Recall from Last Time: Linear Generative Model



Suppose input \mathbf{u} was generated by a linear superposition of causes v_1, v_2, \dots, v_k with basis vectors (or “features”) \mathbf{g}_i

$$\mathbf{u} = \sum_i \mathbf{g}_i v_i + \text{noise} = \mathbf{G}\mathbf{v} + n$$

(Assume noise is Gaussian white noise with mean zero)

Bayesian approach

- Find \mathbf{v} and G that maximize posterior:

$$p[\mathbf{v} | \mathbf{u}; G] = k \cdot p[\mathbf{u} | \mathbf{v}; G] p[\mathbf{v}; G]$$

- Equivalently, find \mathbf{v} and G that maximize log posterior:

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

$$\mathbf{u} = \mathbf{G}\mathbf{v} + n$$

log of Gaussian

$$\log N(\mathbf{u}; \mathbf{G}\mathbf{v}, I)$$

$$= -\frac{1}{2}(\mathbf{u} - \mathbf{G}\mathbf{v})^T (\mathbf{u} - \mathbf{G}\mathbf{v}) + C$$

If v_a independent

$$p[\mathbf{v}; G] = \prod_a p[v_a; G]$$

$$\log p[\mathbf{v}; G] = \sum_a \log p[v_a; G]$$

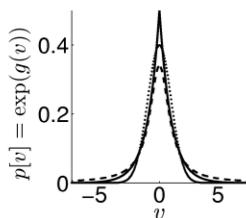
Prior for individual causes (what should this be?)

What do we know about the causes \mathbf{v} ?

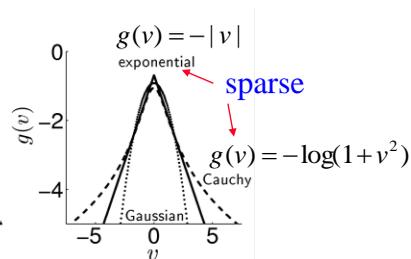
- ◆ Idea: Causes independent: only a few of them will be active for any input
 - ⇒ v_a will be 0 most of the time but high for a few inputs
 - ⇒ Suggests a **sparse distribution** for $p[v_a; G]$: peak at 0 but with heavy tail (also called *super-Gaussian* distribution)

Examples of Prior Distributions for Causes

Possible prior distributions



Log prior



$$p[\mathbf{v}; G] = c \cdot \prod_a \exp(g(v_a))$$

$$\log p[\mathbf{v}; G] = \sum_a g(v_a) + c$$

Finding the optimal \mathbf{v} and G

- ◆ Want to maximize:

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

$$= -\frac{1}{2}(\mathbf{u} - G\mathbf{v})^T (\mathbf{u} - G\mathbf{v}) + \sum_a g(v_a) + K$$

- ◆ Alternate between two steps:

⇒ Maximize F with respect to \mathbf{v} keeping G fixed

◆ How?

⇒ Maximize F with respect to G , given the \mathbf{v} above

◆ How?

Estimating the causes \mathbf{v} for a given input

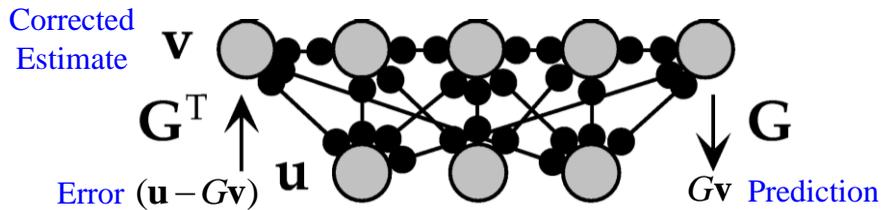
Gradient ascent $\frac{d\mathbf{v}}{dt} \propto \frac{dF}{d\mathbf{v}} = G^T (\mathbf{u} - G\mathbf{v}) + g'(\mathbf{v})$ Derivative of g

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

Reconstruction (prediction) of \mathbf{u}
Error Sparseness constraint Firing rate dynamics (Recurrent network)

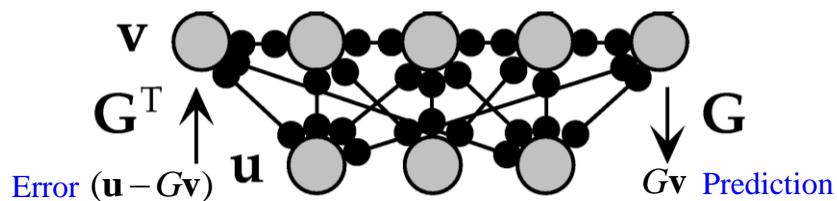
Sparse Coding Network for Estimating \mathbf{v}

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



[Suggests a role for feedback pathways in the cortex (Rao & Ballard, 1999)]

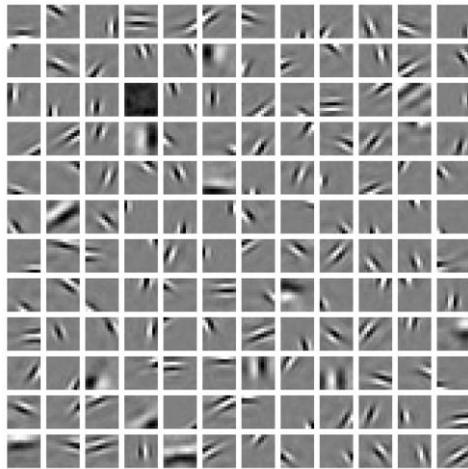
Learning the Synaptic Weights \mathbf{G}



Gradient ascent $\frac{d\mathbf{G}}{dt} \propto \frac{dF}{d\mathbf{G}} = (\mathbf{u} - \mathbf{G}\mathbf{v})\mathbf{v}^T$

Learning rule $\tau_G \frac{d\mathbf{G}}{dt} = (\mathbf{u} - \mathbf{G}\mathbf{v})\mathbf{v}^T$ } Hebbian! (similar to Oja's rule)

Result: Learning G for Natural Images



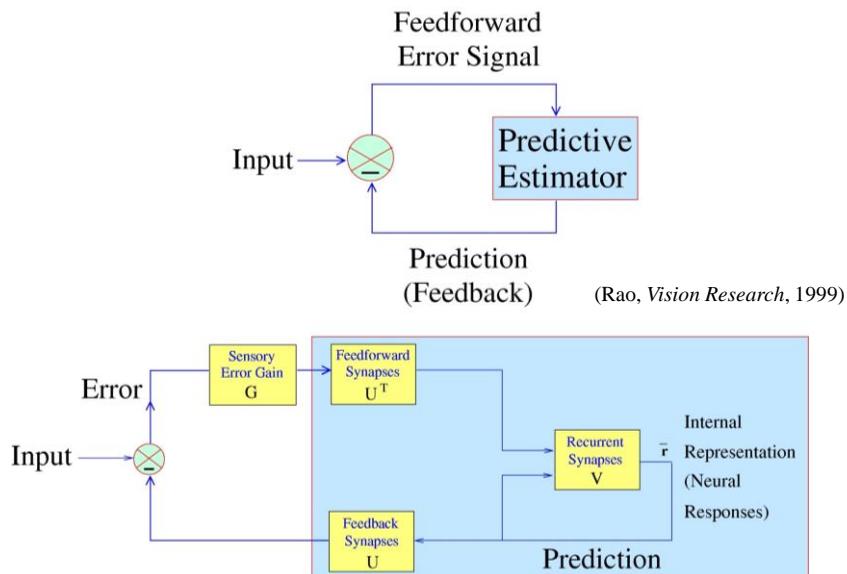
Each square is a column \mathbf{g}_i of G (obtained by collapsing rows of the square into a vector)

Almost all the \mathbf{g}_i represent local edge features

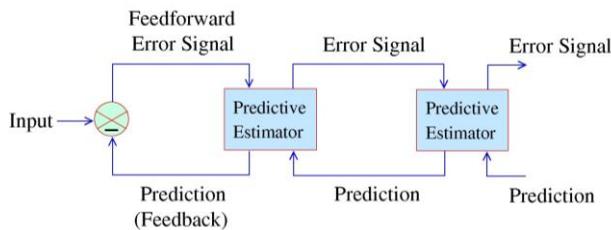
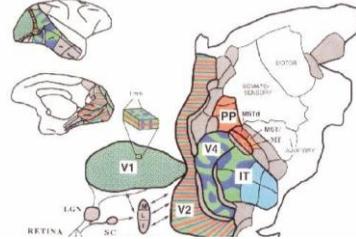
Any image patch \mathbf{u} can be expressed as:

$$\mathbf{u} = \sum_i \mathbf{g}_i v_i = G\mathbf{v}$$

Sparse Coding Network is a special case of Predictive Coding Networks



Predictive Coding Model of Visual Cortex

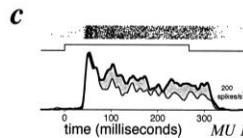
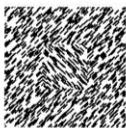
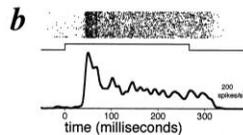


R. Rao, 528: Lecture 13

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

Predictive coding model explains contextual effects

Monkey Primary Visual Cortex

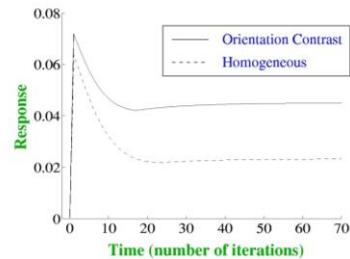


Homogeneous
Random Texture



Orientation Contrast

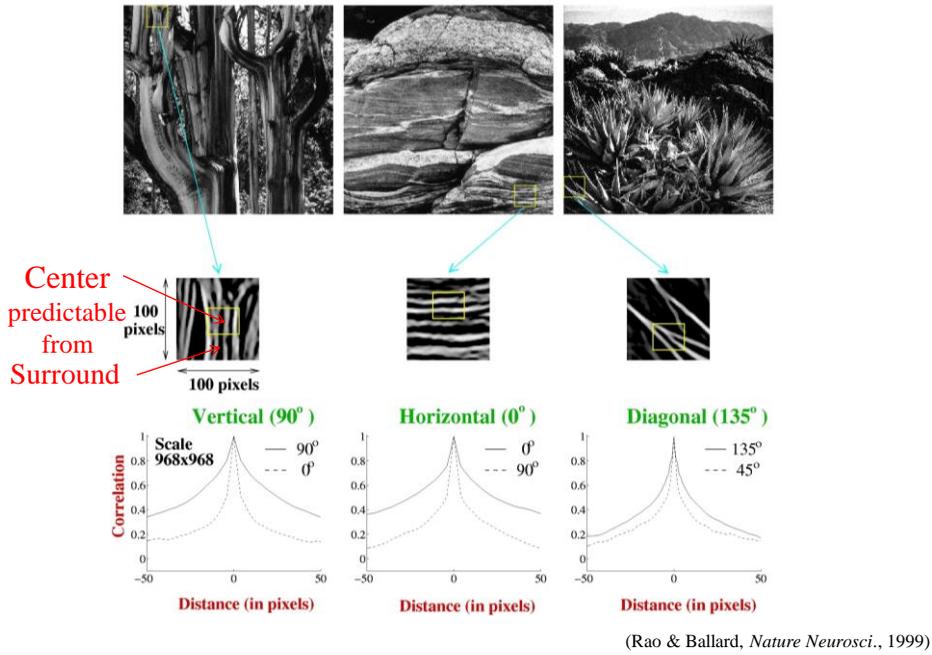
Model



(Zipser et al., *J. Neurosci.*, 1996)

Increased activity for non-homogenous input interpreted as prediction error (i.e., anomalous input): center is not predicted by surrounding context.

Natural Images as a Source of Contextual Effects



What if your data comes with not just inputs but also outputs?

Enter...Supervised Learning

Supervised Learning

◆ Two Primary Tasks

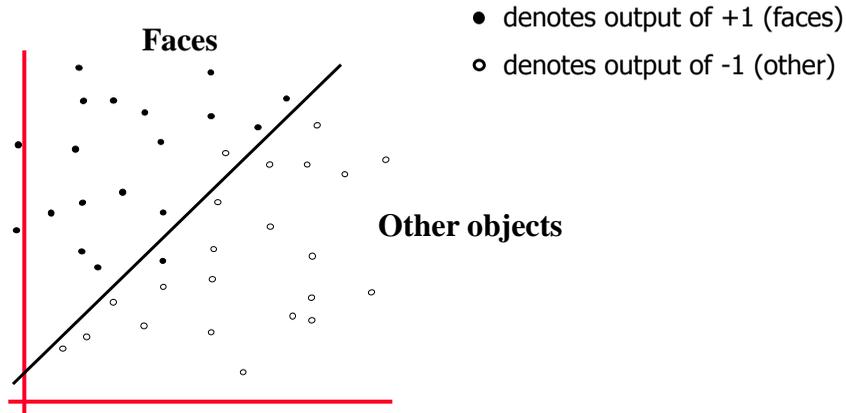
1. Classification

- ◆ Inputs u_1, u_2, \dots and discrete classes C_1, C_2, \dots, C_k
- ◆ Training examples: $(u_1, C_2), (u_2, C_7)$, etc.
- ◆ Learn the mapping from an arbitrary input to its class
- ◆ Example: Inputs = images, output classes = face, not a face

2. Regression

- ◆ Inputs u_1, u_2, \dots and continuous outputs v_1, v_2, \dots
- ◆ Training examples: (input, desired output) pairs
- ◆ Learn to map an arbitrary input to its corresponding output
- ◆ Example: Highway driving
Input = road image, output = steering angle

The Classification Problem



Idea: Find a separating hyperplane (line in this case)

Neurons as Classifiers: The “Perceptron”

- ◆ Artificial neuron:

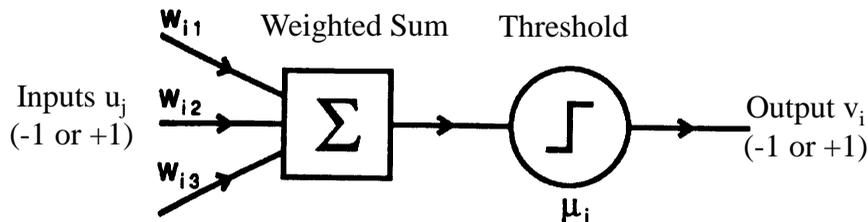
- ⇒ m binary inputs (-1 or 1) and 1 output (-1 or 1)

- ⇒ Synaptic weights w_{ij}

- ⇒ Threshold μ_i

$$v_i = \Theta\left(\sum_j w_{ij}u_j - \mu_i\right)$$

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



What does a Perceptron compute?

- ◆ Consider a single-layer perceptron

- ⇒ Weighted sum forms a *linear hyperplane (line, plane, ...)*

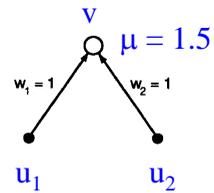
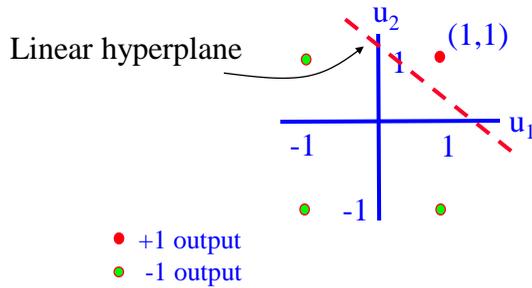
$$\sum_j w_{ij}u_j - \mu_i = 0$$

- ⇒ Everything *on one side* of hyperplane is in **class 1** (output = +1) and everything *on other side* is **class 2** (output = -1)

- ⇒ Any function that is linearly separable can be computed by a perceptron

Linear Separability

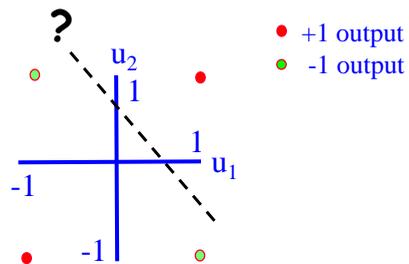
- ◆ Example: AND function is linearly separable
 - ⇒ $a \text{ AND } b = 1$ if and only if $a = 1$ and $b = 1$



Perceptron for AND

What about the XOR function?

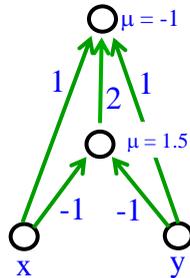
u_1	u_2	XOR
-1	-1	+1
1	-1	-1
-1	1	-1
1	1	+1



Can a straight line separate the +1 outputs from the -1 outputs?

Multilayer Perceptrons

- ◆ Removes limitations of single-layer networks
 - ⇨ Can solve XOR
- ◆ An example of a two-layer perceptron that computes XOR



- ◆ Output is +1 if and only if $x + y + 2\Theta(-x - y - 1.5) > -1$

R. Rao, 528: Lecture 13 (Inputs x and y can be +1 or -1)

23

What if you want to approximate a *continuous* function (i.e., regression)?

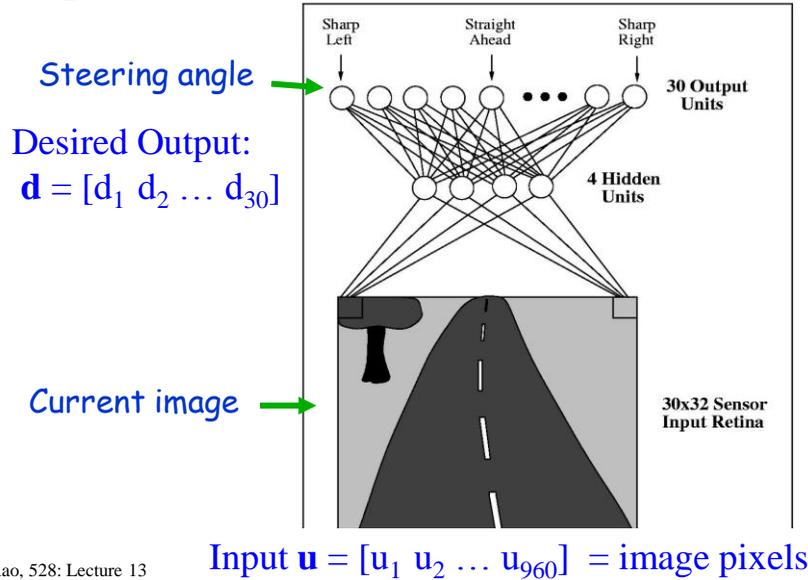


Can a network learn to drive?

R. Rao, 528: Lecture 13

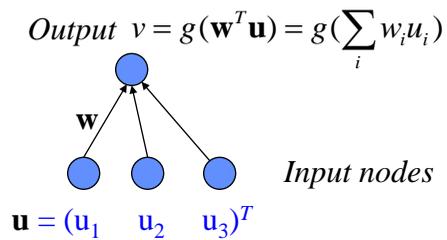
24

Example Network



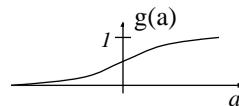
R. Rao, 528: Lecture 13

Sigmoid Networks



Sigmoid output function:

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



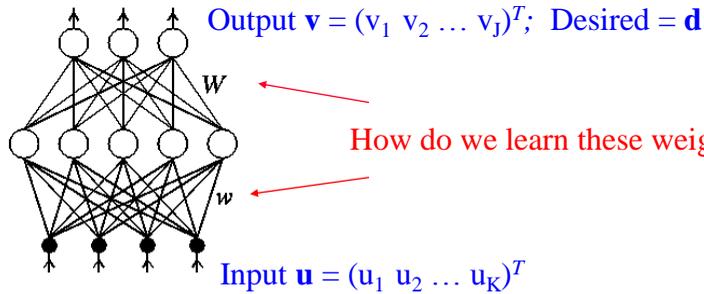
Sigmoid is a non-linear “squashing” function: Squashes input to be between 0 and 1. Parameter β controls the slope.

R. Rao, 528: Lecture 13

26

Multilayer Sigmoid Networks

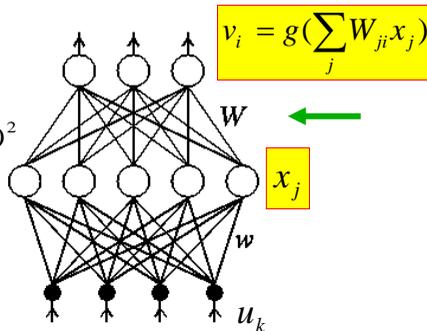
$$v_i = g\left(\sum_j W_{ji} g\left(\sum_k w_{kj} u_k\right)\right)$$



Backpropagation Learning: Uppermost layer

Minimize output error:

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



Learning rule for hidden-output weights \mathbf{W} :

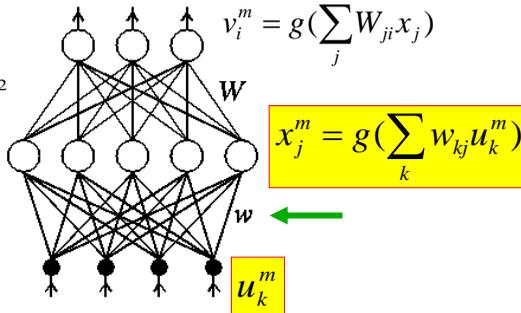
$$W_{ji} \rightarrow W_{ji} - \varepsilon \frac{dE}{dW_{ji}} \quad \{\text{gradient descent}\}$$

$$\frac{dE}{dW_{ji}} = -(d_i - v_i) g'\left(\sum_j W_{ji} x_j\right) x_j \quad \{\text{delta rule}\}$$

Backpropagation: Inner layer (chain rule)

Minimize output error:

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



Learning rule for input-hidden weights w :

$$w_{kj} \rightarrow w_{kj} - \varepsilon \frac{dE}{dw_{kj}} \quad \text{But: } \frac{dE}{dw_{kj}} = \frac{dE}{dx_j} \cdot \frac{dx_j}{dw_{kj}} \quad \{\text{chain rule}\}$$

$$\frac{dE}{dw_{kj}} = \left[- \sum_{m,i} (d_i^m - v_i^m) g'(\sum_j W_{ji} x_j^m) W_{ji} \right] \cdot \left[g'(\sum_k w_{kj} u_k^m) u_k^m \right]$$

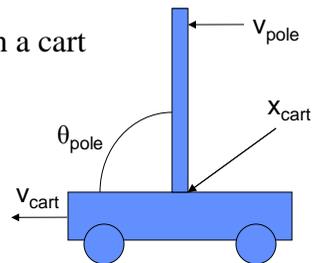
29

Demos: Pole Balancing and Backing up a Truck

(courtesy of Keith Grochow, CSE 599)

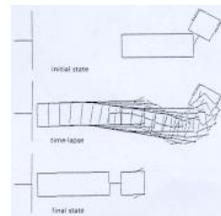
- Neural network learns to balance a pole on a cart

- System:
 - 4 state variables: $x_{\text{cart}}, v_{\text{cart}}, \theta_{\text{pole}}, v_{\text{pole}}$
 - 1 input: Force on cart
- Backprop Network:
 - Input: State variables
 - Output: New force on cart



- NN learns to back a truck into a loading dock

- System (Nyugen and Widrow, 1989):
 - State variables: $x_{\text{cab}}, y_{\text{cab}}, \theta_{\text{cab}}$
 - 1 input: new θ_{steering}
- Backprop Network:
 - Input: State variables
 - Output: Steering angle θ_{steering}



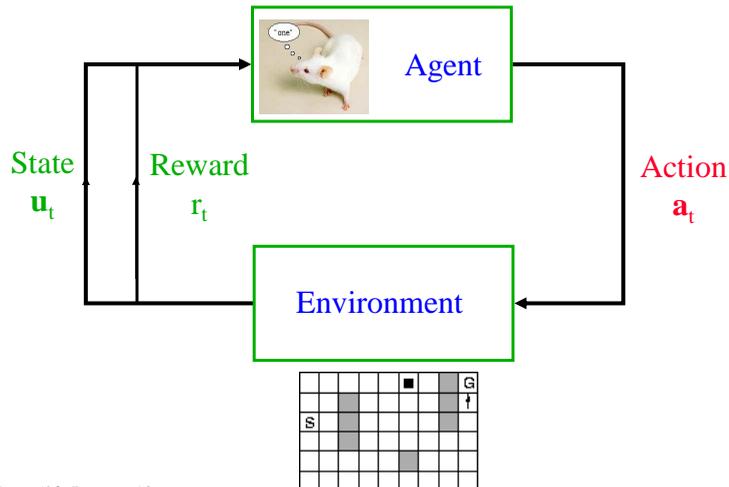
Humans (and animals in general) don't get exact supervisory signals (commands for muscles) for learning to talk, walk, ride a bicycle, play the piano, drive, etc.

We learn by trial-and-error
(with hints from others)

Might get "rewards and punishments" along the way

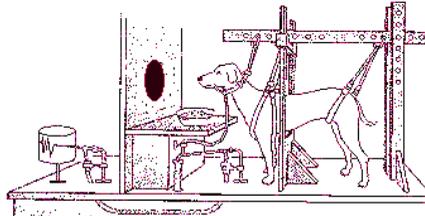
Enter...Reinforcement Learning

The Reinforcement Learning "Agent"



Early Results: Pavlov and his Dog

- ◆ Classical (Pavlovian) conditioning experiments
- ◆ Training: Bell → Food
- ◆ After: Bell → Salivate
- ◆ Conditioned stimulus (bell) predicts future reward (food)



(<http://employees.csbsju.edu/tcreed/pb/pdoganim.html>)

Predicting Delayed Rewards

- ◆ Reward is typically delivered at the end (when you know whether you succeeded or not)
- ◆ Time: $0 \leq t \leq T$ with stimulus $u(t)$ and reward $r(t)$ at each time step t (Note: $r(t)$ can be zero at some time points)
- ◆ Key Idea: Make the output $v(t)$ predict *total expected future reward* starting from time t

$$v(t) \approx \left\langle \sum_{\tau=0}^{T-t} r(t+\tau) \right\rangle$$

Learning to Predict Delayed Rewards

- Use a set of modifiable weights $w(t)$ and *predict based on all past stimuli* $u(t)$:

$$v(t) = \sum_{\tau=0}^t w(\tau)u(t-\tau)$$

- Would like to find the weights (or filter) $w(\tau)$ that minimize:

$$\left(\sum_{\tau=0}^{T-t} r(t+\tau) - v(t) \right)^2 \quad \text{(Can we minimize this using gradient descent and delta rule?)}$$

Yes, BUT...not yet available are the future rewards



Temporal Difference (TD) Learning

- Key Idea:** Rewrite squared error to get rid of future terms:

$$\left(\sum_{\tau=0}^{T-t} r(t+\tau) - v(t) \right)^2 = \left(r(t) + \sum_{\tau=0}^{T-t-1} r(t+1+\tau) - v(t) \right)^2$$

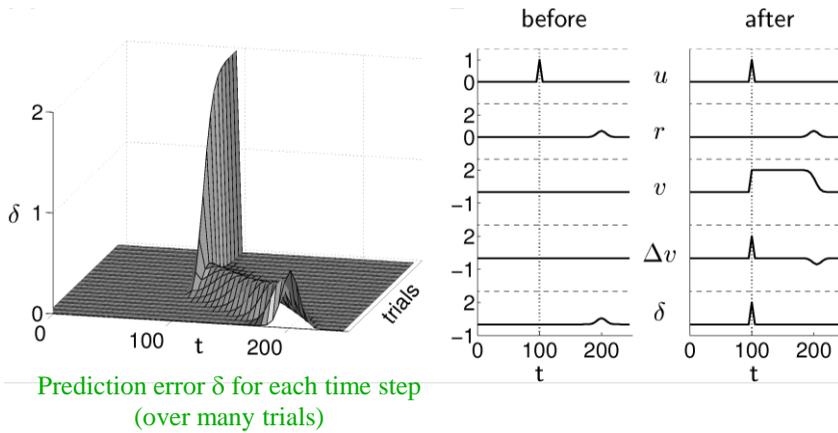
$$\approx (r(t) + v(t+1) - v(t))^2 \quad \text{Minimize this using gradient descent!}$$

- Temporal Difference (TD) Learning:**

$$w(\tau) \rightarrow w(\tau) + \varepsilon \left[\underbrace{r(t) + v(t+1)}_{\text{Expected future reward}} - \underbrace{v(t)}_{\text{Prediction}} \right] u(t-\tau)$$

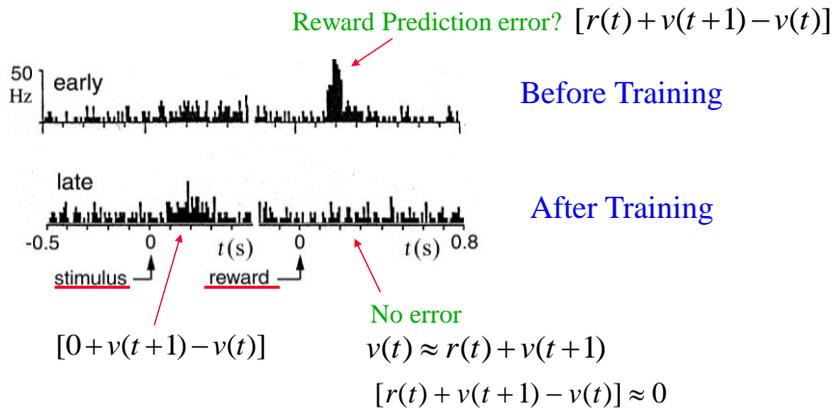
Predicting Delayed Reward: TD Learning

Stimulus at $t = 100$ and reward at $t = 200$



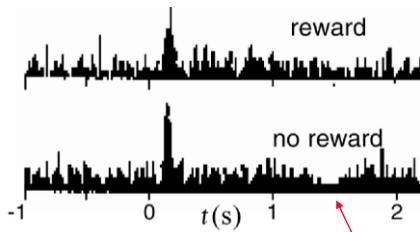
Possible Reward Prediction Error Signal in the Primate Brain

Dopaminergic cells in Ventral Tegmental Area (VTA)



More Evidence for Prediction Error Signals

Dopaminergic cells in VTA



Negative error

$$r(t) = 0, v(t+1) = 0$$

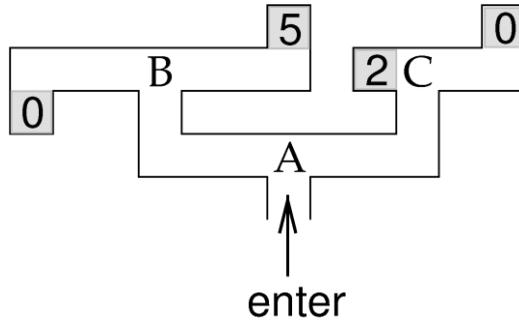
$$[r(t) + v(t+1) - v(t)] = -v(t)$$

Reward predicted
but not delivered

That's great, but how does
all that math help me get
food in a maze?



Selecting Actions when Reward is Delayed

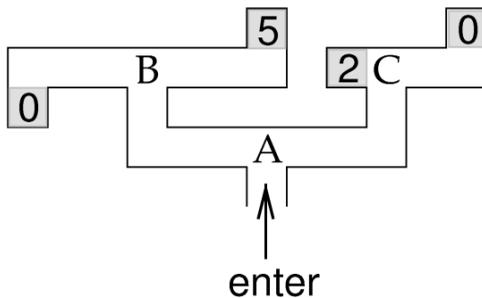


States: A, B, or C

Possible actions at any state: Left (L) or Right (R)

If you randomly choose to go L or R (random “policy”), what is the **expected value v** of each state?

Policy Evaluation



For random policy:

$$v(B) = \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 5 = 2.5$$

$$v(C) = \frac{1}{2} \cdot 2 + \frac{1}{2} \cdot 0 = 1$$

$$v(A) = \frac{1}{2} \cdot v(B) + \frac{1}{2} \cdot v(C) = 1.75$$

(Location, action) \Rightarrow new location

$(u, a) \Rightarrow u'$

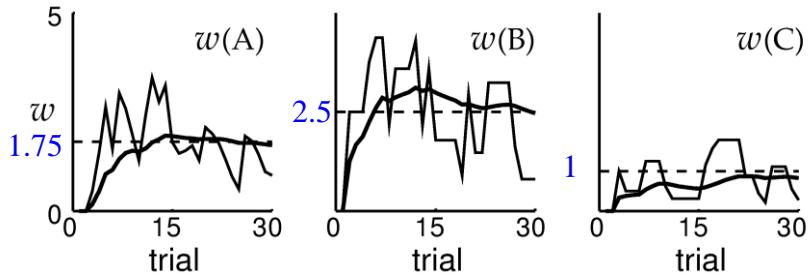
Let value of location

$v(u) =$ weight $w(u)$

$$w(u) \rightarrow w(u) + \epsilon [r_a(u) + v(u') - v(u)]$$

Can learn value of locations using TD learning:

Maze Value Learning for Random Policy

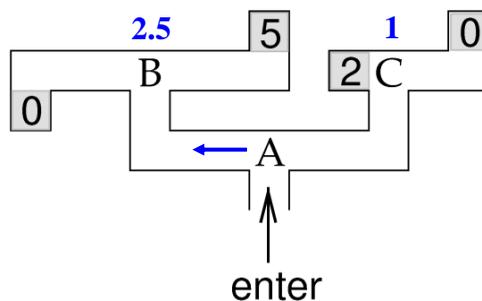


(For all three,
 $\epsilon = 0.5$)

Once I know the values, I can pick the action
that leads to the higher valued state!



Selecting Actions based on Values



Values act as
surrogate immediate
rewards \rightarrow Locally
optimal choice leads
to globally optimal
policy (for “Markov”
environments)
Related to *Dynamic
Programming* in CS
(see appendix in text)

Actor-Critic Learning

- ◆ Two separate components: Actor (maintains policy) and Critic (maintains value of each state)

1. Critic Learning (“Policy Evaluation”):

Value of state $u = v(u) = w(u)$

$$w(u) \rightarrow w(u) + \varepsilon [r_a(u) + v(u') - v(u)] \quad (\text{same as TD rule})$$

2. Actor Learning (“Policy Improvement”):

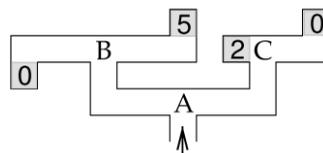
$$P(a;u) = \frac{\exp(\beta Q_a(u))}{\sum_b \exp(\beta Q_b(u))} \quad \begin{array}{l} \text{Use this to select an action } a \\ \text{at state } u \end{array}$$

For all a' :

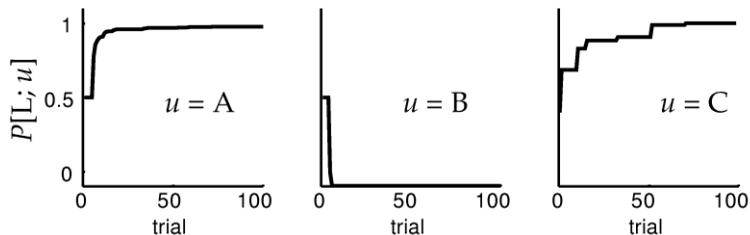
$$Q_{a'}(u) \rightarrow Q_{a'}(u) + \varepsilon [r_a(u) + v(u') - v(u)] (\delta_{aa'} - P(a';u))$$

3. Interleave 1 and 2

Actor-Critic Learning in the Maze Task

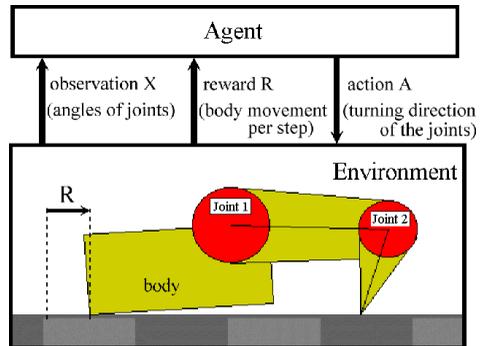


Probability of going Left at a location



Demo of Reinforcement Learning in a Robot

(from <http://sysplan.nams.kyushu-u.ac.jp/gen/papers/JavaDemoML97/robodemo.html>)



Things to do:

Finish homework 3
Work on group project

Thanks, dopamine!

