

Feb 16, 2017 How does an RL agent select actions?

① purely random

↳ con: can't exploit.

② greedy (take the best all the time)

↳ con: doesn't explore

③ ϵ -greedy ($\epsilon \in [0, 1]$ but small)

↳ take the best action with probability $\gamma(1-\epsilon)$
and a random action with probability ϵ .



con: Taking the best frequently at the beginning
is not correct.

④ Softmax action selection

$$p(a) = \frac{e^{Q(s,a)/\tau}}{\sum_b e^{Q(s,b)/\tau}}$$

(
a = specific action
b = running variable
1 to # of actions

estimated value of a in s

How should an RL agent store the rewards?

The average of the first k rewards $Q_k = r_1 + r_2 + \dots + r_k / k$

Can we do this
incrementally?

requires a lot of memory

↳ "running sum"

$$Q_{k+1} = Q_k + \frac{1}{k+1} (r_{k+1} - Q_k)$$

New estimate = old estimate + stepsize (Target - old estimate)

How can an RL agent follow the trajectory?

Q_k to be a simple average then this may
work for a stationary problem. But for non-stationary
problems, it's better to use: $Q_{k+1} = Q_k + \alpha [r_{k+1} - Q_k]$

probability that a^* is the best action when you are in state s

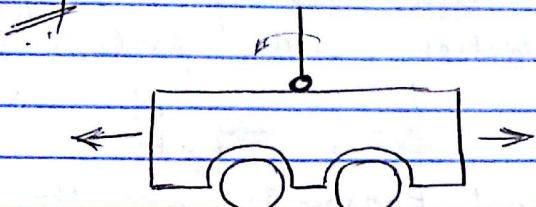
What should an RL agent learn?

↳ RL has to learn a policy π

$\pi(s, a) = \text{probability that } a_i = a \text{ when } s_i = s$
but π changes as a result of more experience.

How can an RL agent function? episodic or never ending?

Csg.



task: avoid failure beyond a certain angle.

We can understand this as an episodic problem.

Each episode ends with a failure.

reward = +1 for each step before failure.

return = # of steps before failure.

Continuing Task: with discounted return.

reward = -1 upon failure
0 otherwise.

return = $-\gamma^k$ for k steps before failure.

In either case, return will be maximized by avoiding failure for as long as possible.

RL Implementation: Tabular

↳ Q-Learning

↳ One step Q-Learning

Q Learning

① Initialize $Q(s, a)$ randomly

② Repeat (for each episode)

2a) initialize s

2b) Repeat (for each step of the episode)

Reinforcement learning is with you in the field.
Neural Networks are not so.

(Chose an action a from set of states S using the policy π from the estimated matrix Q .)

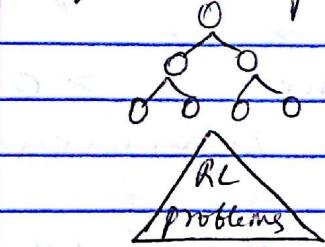
- 2b1) Chose a from S using π from Q
- 2b2) Take action a . observe r, π, S'
- 2b3) Update the Q values.

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

$$2b4) S \leftarrow S'$$

2c) Until S is terminal.

Dynamic Programming



Monte Carlo

used for subsets of entire problem. no compress
no dimensionality reduction.

Temporal differencing

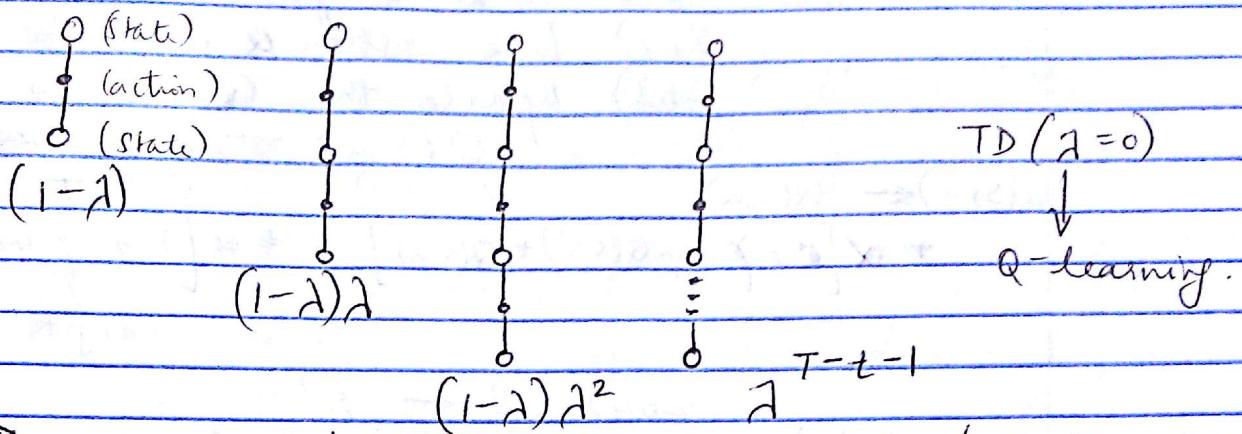
Class of algo works with a compact understanding of what the system is about.

$TD(\lambda)$ is a class of temporal differencing RL methods
(averaging all n -step backups
 \hookrightarrow weight by λ^{n-1} $\lambda \in [0, 1]$
 \hookrightarrow λ -return.

Backup using λ -return:

$$\Delta V_t(s_t) = \alpha [R_t^\lambda - V_t(s_t)]$$

TD (λ) λ -return:



TD ($\lambda = 0$)

\downarrow
Q-learning.

- * TD methods don't require a model, but experience
(Q-Learning is a special case of TD)
- (Q-Learning) is the most simplest RL algorithm
(very short sighted.)
- if you want far sighted, use some $\lambda > 0$.
- * TD methods can be fully incremental.

Q	a_1	a_2	a_3	a_4	the table is initialized with random numbers.
s_1					
s_2					
$\rightarrow s_3$	2.1	1.3	0.9	(3.7)	choose this if you're greedy.
s_4					every cell $Q(s, a)$ has to be visited multiple times.

$\epsilon \rightarrow 0$

$\epsilon = 0.00001$ (never 0)

\Rightarrow always learn. never finish learning

Compare:

(Som vs tmean)
(auto encodes vs nn)

Midterm:

① Tutorial Questions: 25%

② Lecture Questions: 75%

2a) Quiz like $\leq 20\%$ of 75%

2b) general questions: explain / describe $\leq 20\%$

2c) philosophical questions: 20% (e.g. Turing vs Chinese)

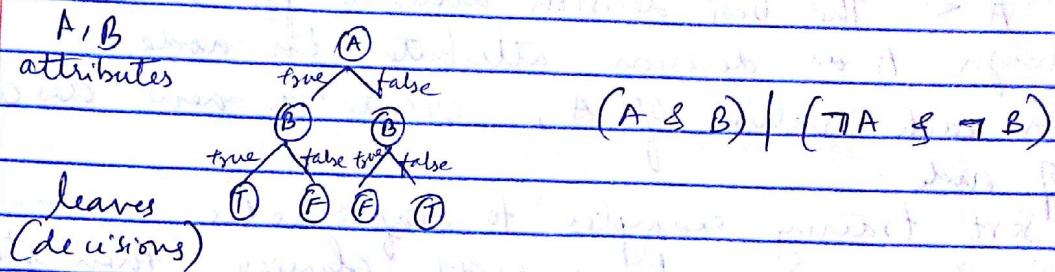
2d) comparative questions: 20% (this vs that)

2e) Math questions $\approx 20\%$ of 75%
(entirely about Back propagation &
Gradient descent)

(derivative of sigmoid function
etc.)

Feb 28th, 2017

Decision Trees (DT)

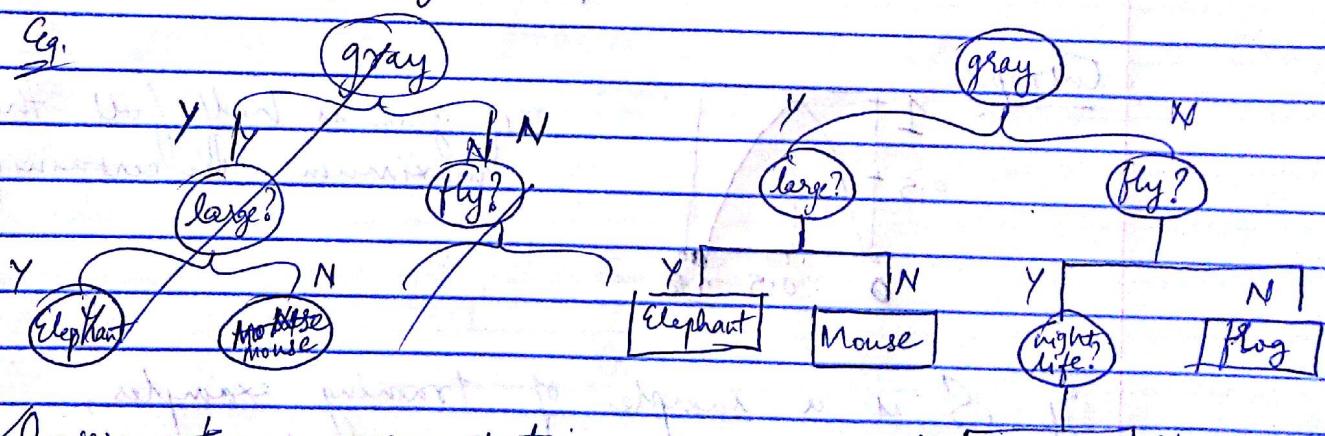


Can we use trees (ADT) for learning to recognize scenes/
objects?

"Attributes" of trees (not features)

↳ more commonly used with "trees"

Eg.



Decision tree representation

- Attributes : internal nodes
- Branches : values of attributes
- leaves : classification

When to use DTs:

- ↳ instances are describable by (attribute, value) pair.
- ↳ Target function is discrete.
- ↳ Disjunctive hypothesis may be necessary
- ↳ noisy data.

Eg. for applications

- (1) Credit risk analysis
- (2) Medical diagnosis (computer aided assistance to clinical experts)
- (3) Modelling scheduling references.

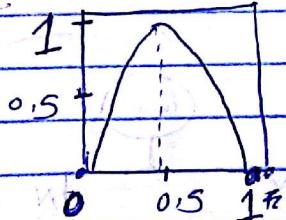
⇒ General scheme for building a Decision Tree

- (1) $A \leftarrow$ the best decision attribute for the next node.
- (2) Assign A as decision attribute for node.
- (3) For each value of A , create a new descendant of node.
- (4) Sort training examples to leaf nodes
- (5) If training examples perfectly classify, then stop, otherwise iterate over new leaf nodes.

Repeat (1) to (5).

What does it mean for an attribute to be the best?

Entropy.



If glass is half full, there is maximum uncertainty.

If S is a sample of training examples,
 P_{\oplus} is proportion of \oplus samples
and
 P_{\ominus} is proportion of \ominus samples.

Entropy H measures the impurity of S .

$$H(S) = -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$$

The expected # of bits needed to encode class (\oplus or \ominus) of randomly drawn member of S (shortest-length condition)

Information Theory.

optimal length code assigns $-\log_2 p$ bits to a message having the probability p .

Entropy \rightarrow how bad an attribute is

? \rightarrow how good an attribute is.

\hookrightarrow i.e., how much would I gain from adding the attribute A to my tree?

Gain (S, A) = expected reduction in entropy due to (Sample) \leftarrow (attribute) sorting on A .

$$\text{Gain } (S, A) = H(S) - \sum_{v \in \text{Value}(A)} \frac{|S_v|}{|S|} H(S_v)$$

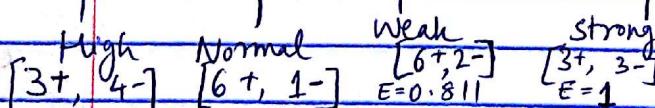
Refer to $|...|$ = cardinality of set, not the absolute value
e.g. on LEARN.

For playing Tennis problem, which attribute is the best?

Sample: 9 positives (Yes, play tennis)
5 negatives (No)

$$S = [9+, 5-]$$

$$E = 0.940$$



$$E = 0.985 \quad E = 0.592$$

$$\text{Gain } (S, "Wind")$$

$$= 0.940 - \left(\frac{8}{14} \right) 0.811 - \left(\frac{6}{14} \right) 1$$

$$\text{Gain } (S, "Humidity") =$$

$$= 0.940 - \left(\frac{7}{14} \right) 0.985 - \left(\frac{7}{14} \right) 0.592$$

$$= 0.048$$

$= 0.151 \rightarrow$ information gained if we add "humidity" to tree.

\therefore humidity has most discrimination power. Keep it @ the top of decision tree

Main Question: Small or Large tree?

⇒ Occam's Razor

A short hypothesis that fits the data is unlikely to be coincidence.
∴ we prepare small trees.

Overfitting

Training data : error_{train}(h)

hypothesis = solutions

Distribution of data : error_D(h)

Hypothesis $h \in H$ overfits the training data

if there is $h' \in H$ such that
 $\text{error}_{\text{train}}(h) \leq \text{error}_{\text{train}}(h')$

but $e(h) > e(h')$.

for entire data D
(the one that overfits is bigger)

