

Apr 28
12:04 pm

RL Objective

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \quad E_{\Sigma \sim \Pi_{\theta}(x)} [r(x)]$$

optimal policy's parameters

$$E_{\Sigma \sim \Pi_{\theta}(x)} [r(x)]$$

↳ This quantity we need to increase
let's call the qty "J"

$$J(\theta) = E_{\Sigma \sim \Pi_{\theta}(x)} [r(x)]$$

↳ This is our Total Reward

↳ To maximize it, we need to find $\frac{\partial J(\theta)}{\partial \theta}$

$$\theta = \theta + \frac{\partial J(\theta)}{\partial \theta}$$

Gradient Ascent

To evaluate,
J(θ)

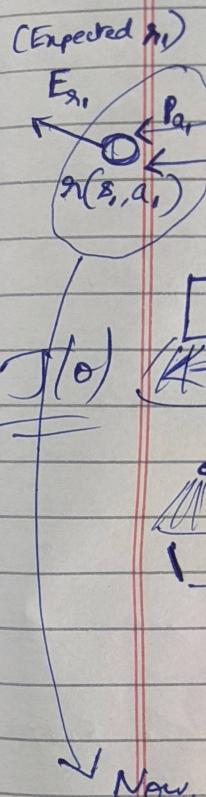
we need to run samples to collect trajectories

$$\underline{J(\theta)} \approx \frac{1}{N} \sum_{i=1}^N r(x_i)$$

as $N \rightarrow \infty$, $J(\theta)$ becomes equal to this quantity!
↳ more the samples, more accurate it is!!

dividing by $\frac{1}{N}$, will give us the probability of trajectory

$$\underline{P(x_i)}$$

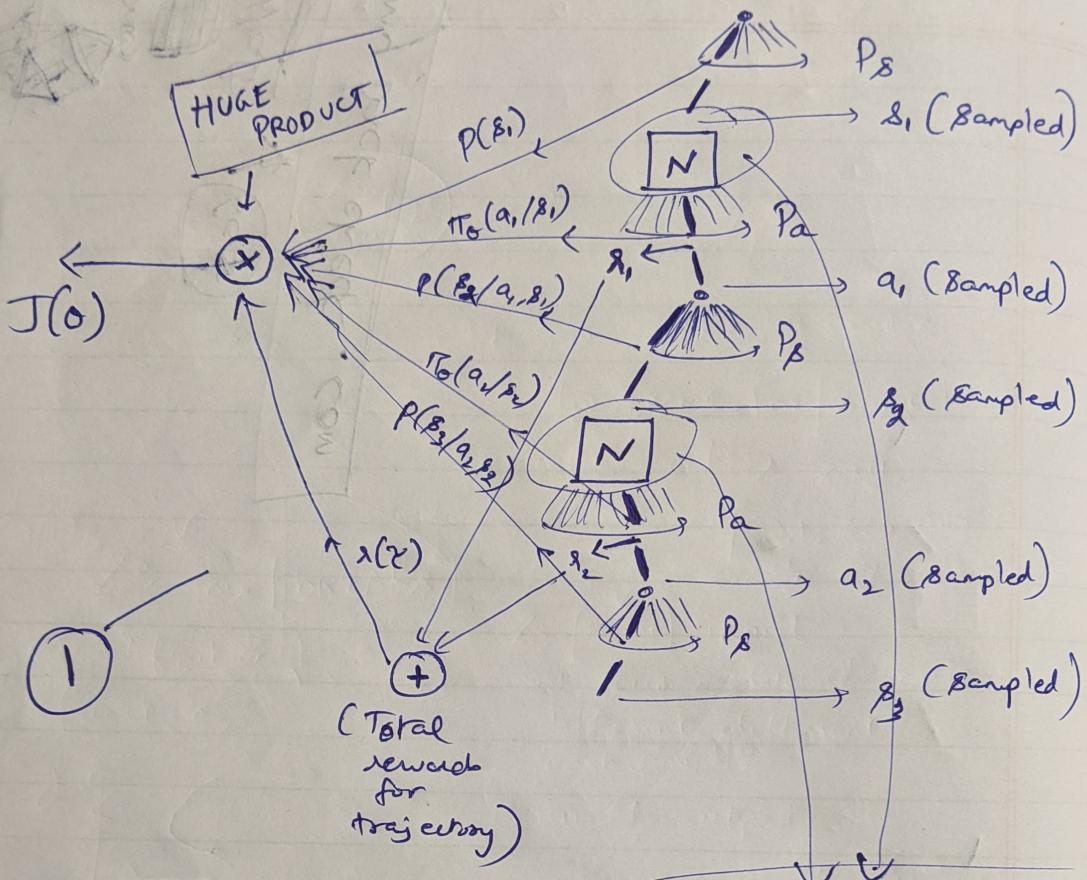


Now, what for the equal

if equal

But equal

In order to calculate the gradient, we can visualize the computation graph & derive it from there!



Probability next letter P

So, here we do the whole

But, include of our

So, if you taking

Now, if you see,
the network is
repeating like RNN
& the weights are shared!

And, just like RNN, it ~~never~~ looks like 1st network's choices affect future network's choices!

So, multiple gradients will flow for one network's choices from all of future + present!

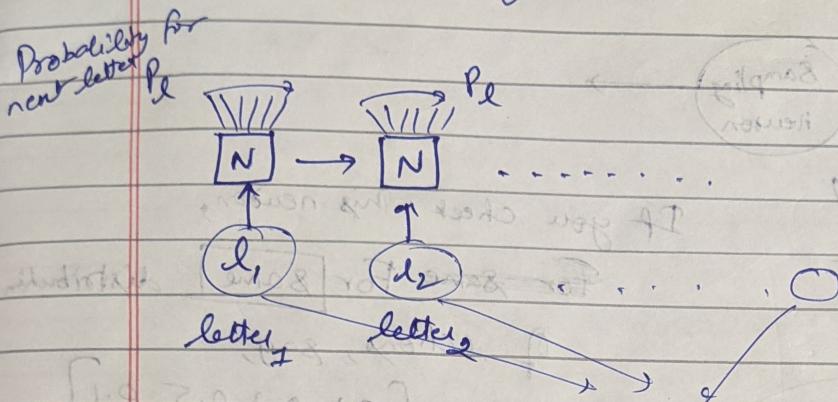
According to this assumption,
it would also look like
you would need ~~to~~ to
know the values of environment
dynamics like $p(s_2 | s_1, a_1)$, $p(s_3 | s_2, a_2)$
~~in order to correctly~~
backpropagate since they
look like they are PART
of the network!

But, That's also NOT "true"!

Let's take an example to understand why!

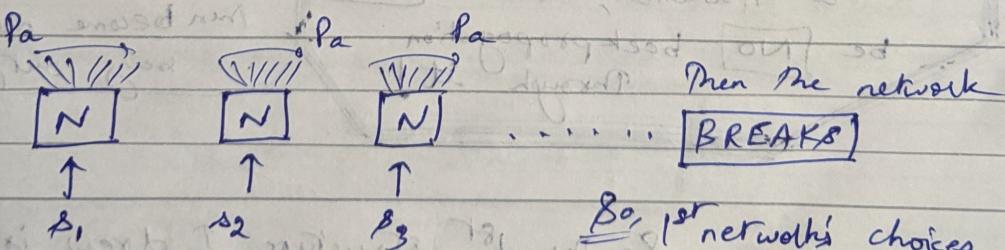
In, [lect 10, 9/CS231n] Winter 2016

Andrij demonstrates a RNN that predicts next letter given ~~current~~ sequence!



If you look @ all the l_i ,
we don't know ← These are actually being [outputted]
the world dynamics! by the world dynamics! we speak
The world dynamics of English language generates these letters!
And what we [see] as input here,
is actually [sampled] from this world distribution probability!

So, in our network, if you incorporate the [same], taking all sampling to be [INPUT], both action + & state!



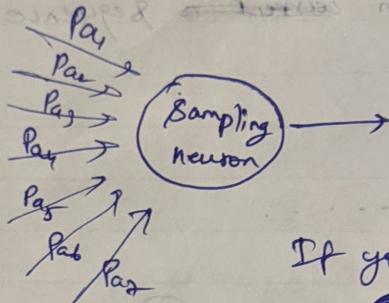
So, 1st network's choices gets [No] feedback or gradient from future, although it plays a part in that!

FNN

The same would have happened, if Andrij's problem!

But, they cleverly introduced connections b/w the networks, do not let the gradient flows

By the way, before we proceed, what is this Sampling step. How do we backpropagate through it?



If you check this neuron,

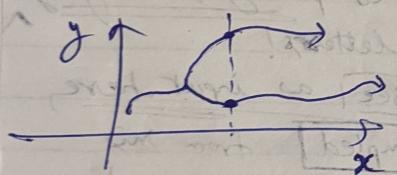
For same distribution of actions, say,

$$[0.1, 0.3, 0.5, 0.1]$$
$$a_1, a_2, a_3, a_4$$

So this BREAKS the definition of a function.

So, same input

the neuron would output different results each time!



For same x , we can't have 2 values of y !

Sampling is not even a function, & there can be NO backpropagation through it!

different labels each time for same input (distribution)

To make it a believed neuron, you can encode the labels in one-hot representation & it would then become numbers being outputted

So, what does EXPECTATION

And this is what we are optimizing!

So, the 1st structure I drew is not even VALID. Also, it's not even a computation graph because in the step of Sampling, there is NO computation involved. The probability distribution values is not taken & passed through a function where multiplication or division happens! In fact, there is NO function that can represent Sampling!

So $\pi(x)$ is INPUT

So, what I drew is a sequence graph, NOT computation graph

It's the order in which the events happened!

It's not a flow of computations!

So, All sampling are INPUT!

But, Q states + actions

if states & actions are input & fixed

Then even, the rewards (functions of states & actions) are FIXED

So, all 3, $(\pi, q, r(\pi, a))$ are fixed, constants & inputs!

So, for a trajectory (τ) , if states, actions & rewards don't change!

Then by changing in policy $\pi(\tau)$

The total reward would still

remain the same right!

since rewards don't change

So, what changes is

EXPECTATION of rewards (not rewards, !)

And this is what we are optimizing!

So, for 2 policies, $\theta_1 + \theta_2$

For the same trajectory τ the rewards are same, $r(\tau)$

But the expectations

of rewards are

different!

So $\pi(\tau)$ is INPUT

and we pick the policy whose expectation is HIGHEST

$$\delta_0 \quad J(\theta) = \left[E_{\sum \pi_\theta(x)} [x(x)] \right]$$

discrete

all x

$\sum x(x) p_\theta(x)$

$\int p_\theta(x) x(x) dx$
all x
 $-\infty \text{ to } \infty$

$$J(\theta) = \sum_{\text{all } x} p_\theta(x) x(x)$$

$$\frac{\partial [J(\theta)]}{\partial \theta} = \frac{\partial}{\partial \theta} \left[\sum p_\theta(x) x(x) \right]$$

$$= \sum_{\text{all } x} \frac{\partial}{\partial \theta} [p_\theta(x) x(x)]$$

constant

$$= \sum_{\text{all } x} x(x) \frac{\partial}{\partial \theta} p_\theta(x)$$

$$p(\theta_1) \times \pi_\theta(a_1 | \theta_1) / p(\theta_2 | \theta_1, a_1) \pi_\theta(a_2 | \theta_2) \dots$$

environment dynamics are **NOT** function of θ, θ_0 , constants!

$$= \sum_{\text{all } x} x(x) \left[p(\theta_1) p(\theta_2 | \theta_1, a_1) p(\theta_3 | \theta_2, a_2) \dots \right] \frac{\partial}{\partial \theta} [\pi_\theta(a_1 | \theta_1) \times \pi_\theta(a_2 | \theta_2) \dots]$$

derivative of products

$$\frac{d}{dx} [y_1 y_2 y_3 \dots y_n]$$

$$= (y'_1 y_2 y_3 \dots y_n + y_1 y'_2 y_3 \dots y_n + \dots + y_1 y_2 y'_3 \dots y_n)$$

Common
can be
computed
ONCE

You will have ' N ' terms,
where in each term, one of the terms
is differentiated!

β_0

$$= \sum_{\text{all } r} r(x) \left[P(\beta_0) p(a_1/a_0, s_1) \dots \right] \frac{\partial}{\partial x} \left[\pi_0(a_1/s_1) \pi_0(a_2/s_2) \dots \right]$$

$$\left(\pi_0(a_1/s_1) \pi_0(a_2/s_2) \dots \pi_0(a_n/s_n) \right)$$

$$+ \left(\pi_0'(a_1/s_1) \pi_0(a_2/s_2) \dots \pi_0(a_n/s_n) \right)$$

$$\pi_0(a_1/s_1) \pi_0(a_2/s_2) \dots \pi_0'(a_n/s_n)$$

To compute this is
difficult to code +

Not efficient! → many of the multiplications
is repeated!

$$\cancel{\pi_0(a_i/s_i)} \cdot \pi_0'(a_i/s_i) \times \pi_0(a_{i+1}/s_{i+1})$$

Duplication of computation

This product is
calculated ' $n-2$ ' times!

If **ONLY** we can compute this
SHARED computation **ONLY ONCE**

β_0

$$\frac{\pi_0'(a_1/s_1)}{\pi_0(a_1/s_1)} \times \left[\pi_0(a_2/s_2) \pi_0(a_3/s_3) \dots \pi_0(a_n/s_n) \right]$$

$$+ \frac{\pi_0'(a_2/s_2)}{\pi_0(a_2/s_2)} \times \left[\pi_0(a_1/s_1) \pi_0(a_3/s_3) \dots \pi_0(a_n/s_n) \right]$$

Common

can be
computed
ONCE

$$+ \frac{\pi_0'(a_n/s_n)}{\pi_0(a_n/s_n)} \times \left[\pi_0(a_1/s_1) \pi_0(a_2/s_2) \dots \pi_0(a_{n-1}/s_{n-1}) \right]$$

$$P_{\theta} = \sum_{a|x} \gamma(x) \left[p(s_1) p(a_1|s_1, a_1) \dots \right] \left[\pi_{\theta}(a_1|s_1) \pi_{\theta}(a_2|s_2) \dots \pi_{\theta}(a_n|s_n) \right]$$

Apr 22
5:08pm

$\boxed{= P_{\theta}(x)}$

$\boxed{\text{in integral}}$

$$\begin{aligned} P_{\theta} &= \sum_{a|x} \gamma(x) P_{\theta}(x) \left[\frac{\partial}{\partial \theta} [\log \pi_{\theta}(a_1|s_1)] \dots \right] \\ &\quad E_{a \sim \pi_{\theta}(x)} \left[\sum_{i=1}^T \frac{\partial}{\partial \theta} (\log \pi_{\theta}(a_i|s_i)) \lambda(x) \right] \end{aligned}$$

$\boxed{\frac{\partial J(\theta)}{\partial \theta}}$

Policy Gradients

Now, if [partial observability] is included

Then,

$$= \sum_{a|x} P_{\theta}(x) \lambda(x)$$

$$p(s_1) p(o_1|s_1) \pi_{\theta}(a_1|s_1) p(a_2|s_1, a_1)$$

$$p(o_2|s_2) \pi_{\theta}(a_2|o_2)$$

[Observations] \rightarrow o ' dynamics are added! \rightarrow These are also CONSTANTS!

\rightarrow it comes out of differentiation!

And formula doesn't CHANGE!