

687 2017-09-19

Note Title

9/19/2017

Pseudocode

```
for episode = 1, 2, ...
  s ~ d_0
  for t = 0, 1, 2, ...
    a = agent.getAction(s)
    s' ~ P(s, a, ·)
    r = R(s, a, s') # assuming deterministic
    agent.train(s, a, r, s')
    if s' == s break
  s = s'
```

Black Box Optimization for Policy Search

- Simple Agent
- Ignore MDP structure
- Phase this as an optimization problem

$$\arg \max_{\pi} J(\pi)$$

- Can access only an estimate $\hat{J}(\pi)$,
e.g. evaluation of π on a sample
of N episodes, average the

discounted reward \rightarrow

$$\hat{J}(\pi) = \frac{1}{N} \sum_{i=1}^N G_i \rightarrow G_i = \sum_{t=0}^{\infty} \gamma^t R_t^i$$

\leftarrow i th episode

A way to represent π : As a table p

	Actions			
	a_1	a_2
States s_1				
s_2				
...				
...				

- # of these parameters is $|S| \times |A|$
- valid if each row sums to 1) and all entries ≥ 0

want a function optimization method that obeys this constraint

\Downarrow

Tubular Soft max Action Selection

- Store π as an $|S| \times |A|$ matrix
- No constraints on the matrix
- Increasing $p(s,a)$ increases $\pi(s,a)$

$$\pi(s,a) = \frac{e^{p(s,a)}}{\sum_{a'} e^{p(s,a')}} \Downarrow$$

This guarantees $(\sum_a \pi(s,a)) = 1$ and $\pi(s,a) \geq 0$

Conventional to call these parameters θ (a vector formed from p)

A couple of alternatives

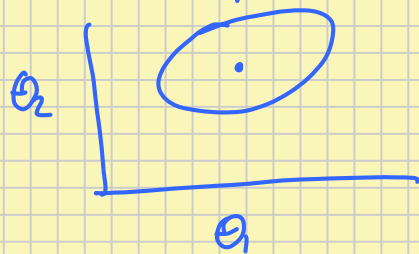
Modeling Diabetes: Policy: Before a meal inject $\frac{BG - BG_{\text{target}}}{\theta_1} + \frac{\text{size of meal}}{\theta_2}$

BG: blood glucose (measured by pump)
BG_{target}: level the doctor believes is good
size of meal = estimated by patient
 θ_1, θ_2 : parameters we may try to learn

Cross-Entropy Method (was used for Tetris) → simple but powerful, general

- ① Generate a random data sample according to a parameterized mechanism
- ② Update parameters based on the samples to do better at the next iteration

- Mechanism = distribution over policies
- Random sample = histories generated by sampled policies



2-D Gaussian

- Draw (θ_1, θ_2) pairs
- Compute $\hat{J}(\pi)$, i.e., $\hat{J}(\theta)$
- Move mean toward higher $\hat{J}(\pi)$
 - Can do weighted sum - regression method
 - K_e : elite population \subseteq K population
 - ↖ best - update mean w/mean over K_e
 - Do likewise with covariance (covariance over K_e)

Pseudo-Code for Cross-Entropy Method (CEM)

Input: θ - mean parameter vector
 Σ - covariance, initially $\sigma^2 I$
 K - population
 K_e - elite population
 N - # of episodes

for $k = 1:K$

$$\theta_k \sim N(\theta, \Sigma)$$

$$\hat{J}_k = \text{evaluate}(\theta_k, N)$$

sort (θ_k, \hat{J}_k) , descending on \hat{J}_k
 $k=1, K$

$$\theta \leftarrow \frac{1}{K_e} \sum_{k=1}^{K_e} \theta_k \quad (\text{top } K_e \theta_k \text{ according to quality estimates } \hat{J}_k)$$

$$\Sigma \leftarrow \frac{1}{K_e} \sum_{k=1}^{K_e} (\theta_k - \theta)(\theta_k - \theta)^T \quad \leftarrow \text{can be numerically unstable so...}$$

$$\Sigma \leftarrow \frac{1}{K_e + \epsilon} \left(\epsilon I + \sum_{k=1}^{K_e} (\theta_k - \theta)(\theta_k - \theta)^T \right)$$

Why called this? Has to do with how to go from one θ distribution to the next while minimizing a certain loss function.

Demo of project things in Visual Studio