

* Access Models

→ full probabilistic description.

costs $\cdot P_{ss}^a$

→ Can't store

→ Information theoretically heavy

→ Deterministic Generative Model

$$f(x, u) \rightarrow x'$$

Set the random seed.
to recreate same scenario

→ Generative model.

Programmatic access

→ Reset Model.

Encode $f(x, u)$ sequences
and anytime reset to x_0 .

→ Trace Model.

Can't ever go back.

Approximate Dynamic Programming:-

(1) Too expensive to compute and store
value iteration $V(x, t)$

(2) Need weaker access to model
↳ generative model.

Fitted Value Iteration :- (FVI)

Regression Algorithm as oracle.

$$D = \{x_i, y_i\}_{i=1 \dots m}$$

↑
vector of features

{ representation of state space }

output, R
(value)

$$f: X \rightarrow Y$$

Ex: $\frac{1}{N} \sum_{i=1}^N (f(x_i) - y_i)^2$ is small

→ Squared error, sensitive to outliers

→ Overfitting

↓

How to fix?

by

→ minimize $E_{r(x,y)} [(f(x) - y)^2] \leq \epsilon$

and ~~also~~ Regularize

What Regression Algos?

Linear Regression.

Locally Weighted Regression.

Gaussian Processes. [or other kernel machine]

Regression Tree.

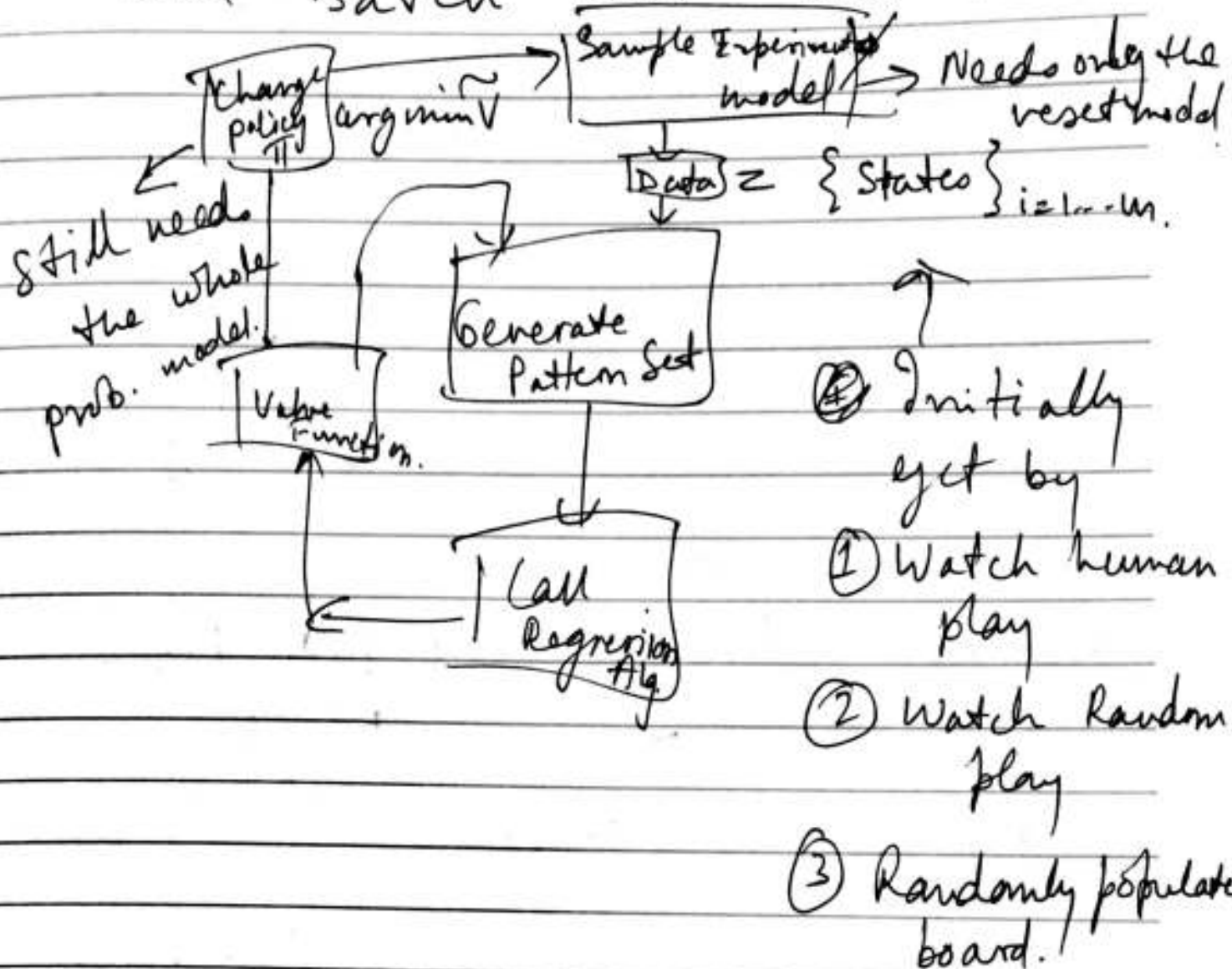
Random Forests.

Neural Networks.

Splines.

$$V_k(s) = \min_a \text{cost}(s, a) + \gamma \sum_{s'} P(s'|s, a) V_{k+1}(s')$$

* General Batch ADP Data Flow:-

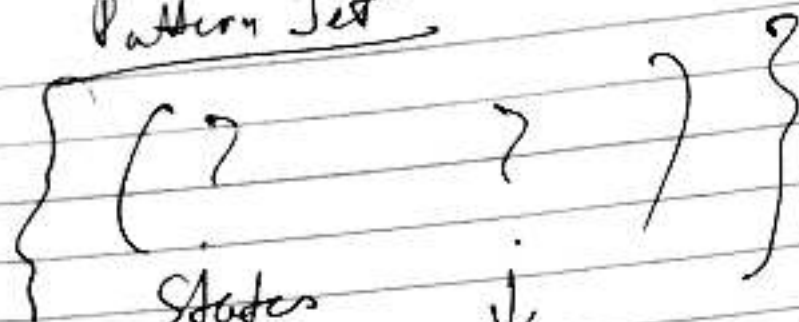


- Initially get by
- ① Watch human play
 - ② Watch Random play
 - ③ Randomly populate board.

But keep in mind that regular Value Iteration also has that problem.

Some states are not reachable.

Pattern Set



States
(x)

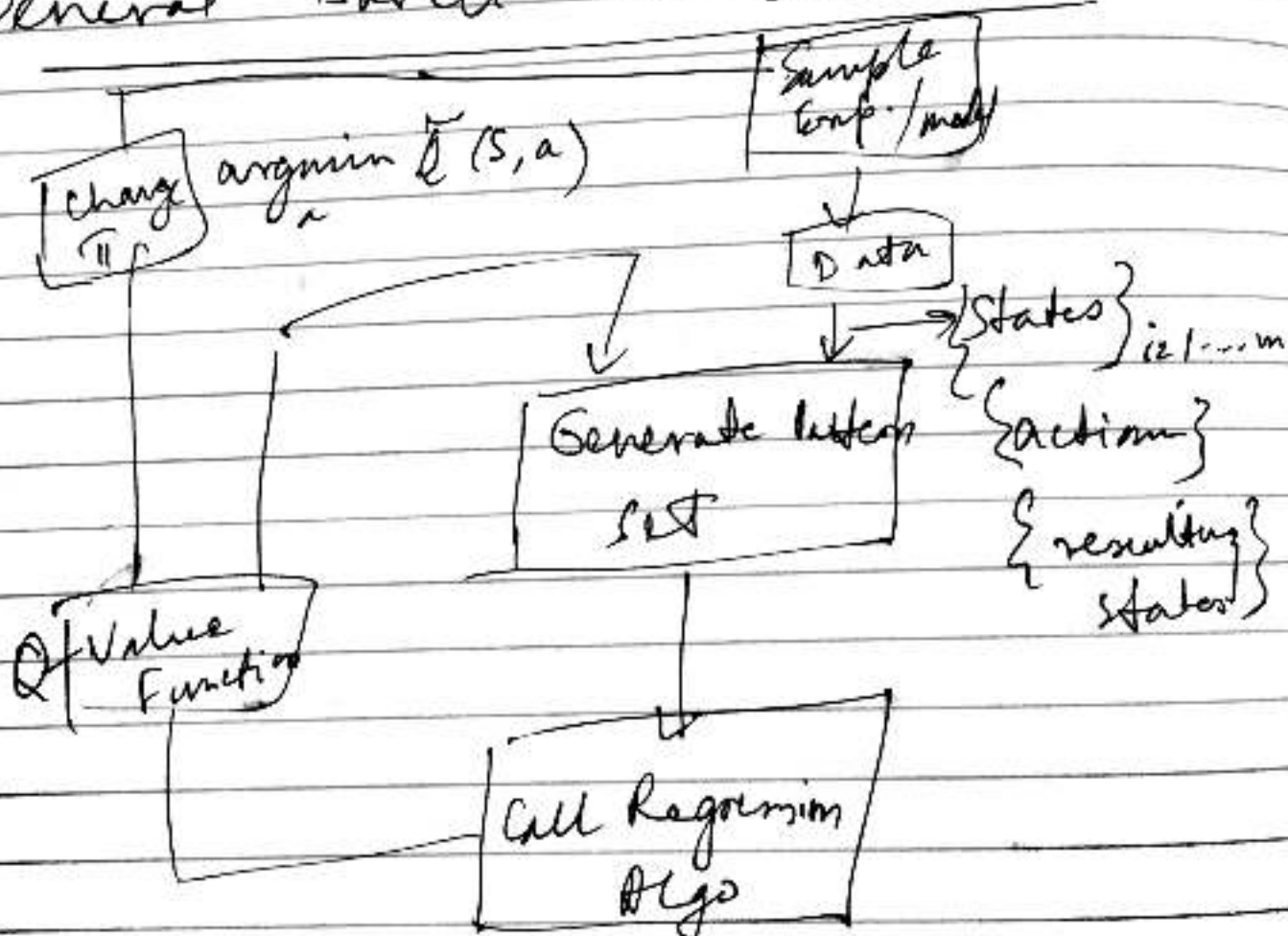
$$\min_a \left[\text{cost}(s, a) + \sum_{s'} p(s'|s, a) \cdot V_{k+1}(s') \right]$$

to get V^* value.

This is set to
0 for the
1st iteration

(Generative model)

General Batch ADP Data Flow:-



Fitted Q-Iteration (FQI)

Quality Function.

Action-Value Function.

$Q(s, a) \rightarrow$ Value you get if
 $V(s)$ you choose a at s
and then follow the
optimal policy

$$* Q(s, a) \geq V(s)$$

$$* V(s) = \min_a Q(s, a)$$

$$* a^* = \operatorname{argmin}_a Q(s, a)$$

$$Q(s, a, t) = \text{cost}(s, a) + \min_{a'} Q(s', a', t+1)$$

Storing this
requires more
space.

$$= \text{cost}(s, a) + \sum_{(s', a')} \min_{a'} Q(s', a', t)$$

Pattern Search

$$\left\{ (s, a), \text{cost}(s, a) + \gamma \min_{a'} Q^{\pi}(s', a') \right\}$$

resulting state
↑
←

Just the
samples
input.

No model.

Fitted Policy Iteration (FPI)

Recap of Policy Iteration

Policy Iteration

(a) Evaluate Policy

Calculate $V^{\pi}(s)$

$$V^{\pi}(s) = \text{cost}(s, \pi(s)) + \gamma \mathbb{E}_{p(s'|s, \pi(s))} [V^{\pi}(s')]$$

(b) Improve Policy

$$Q^{\pi}(s) = \text{cost}(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} [Q^{\pi}(s', \pi(s'))]$$

$$\pi' \leftarrow \underset{a}{\text{argmin}} Q^{\pi}(s, a)$$

(Note: Don't need to use model)

Pattern Set

$$\{(s, a), \text{cost}(s, a) + \gamma Q^{\pi}(s', \pi(s'))\}$$

(Data Flow Diagram is same as FDI)

Advantages (to FDI)

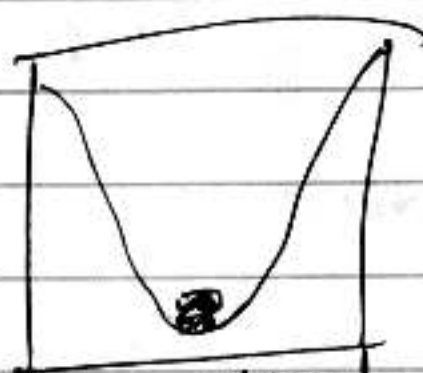
* ^{→ w/o error blow up.} ~~Stable~~ than FDI. (as it is linear)
in the Evaluate Policy equation.

* Toy Examples

(Boyan and Moore)

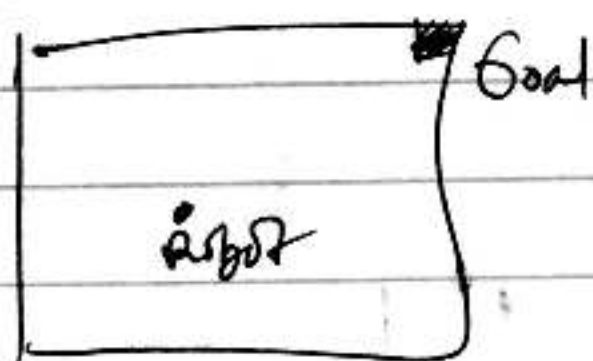
→ Cart on the Hill

problem



→ under actuated
Non-minimum
phase problem

Empty grid with goal at one corner.



You have to go
the wrong way to
go the right
way.

Convergence Results?

— Relatively weak
"The Averagers"
(Gordon et al.)

— "On policy" temporal difference ones will
have stronger results.

→ Hang on, still coming.

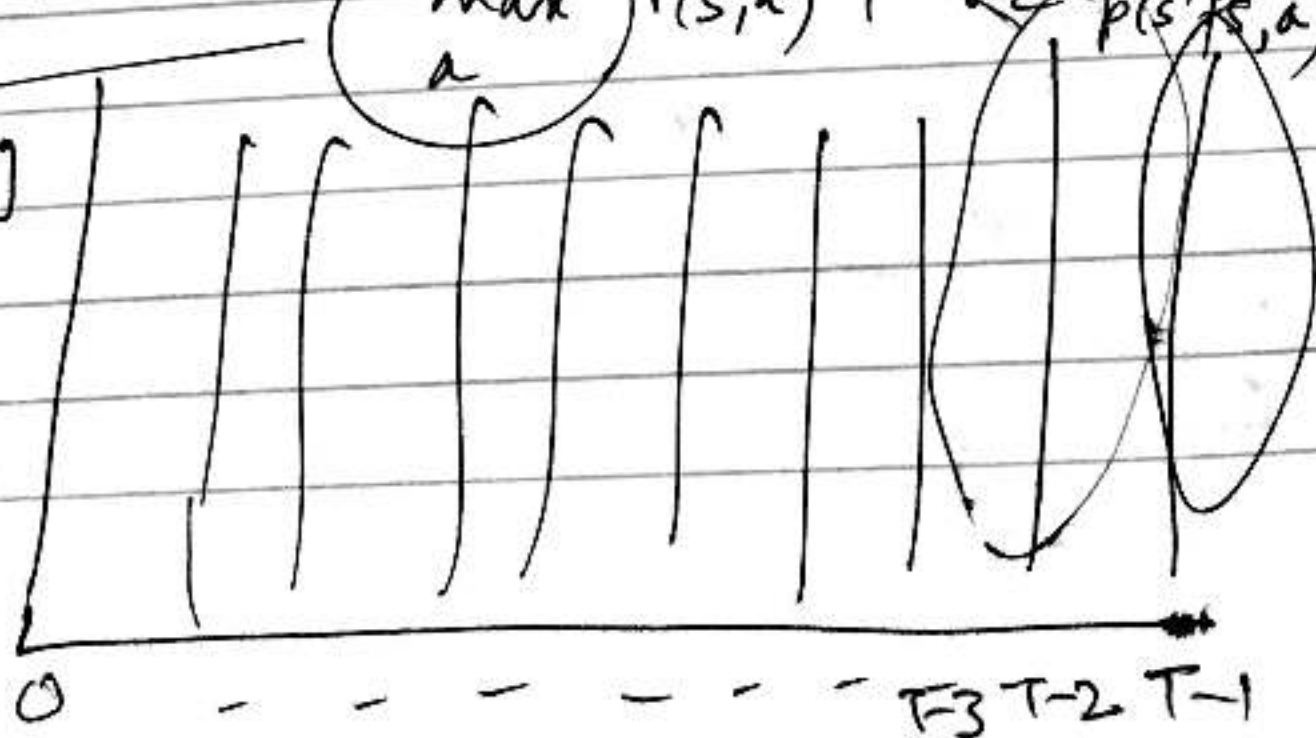
— Error Analysis?

Horrendous.

Why does FQI, ~~TD~~ go horribly wrong
sometimes?



Amplifying
error τ .
S



Problem

① where should sample set be?

② max in value iteration amplifies errors.

In policy iteration, the inner loop has ^{Can Fix!} no ~~max~~ max or min term and hence is more stable.

Recap

Fitted Dynamic Programming $\left[\begin{array}{l} \text{FPI} \\ \text{FSI} \\ \vdots \end{array} \right]$

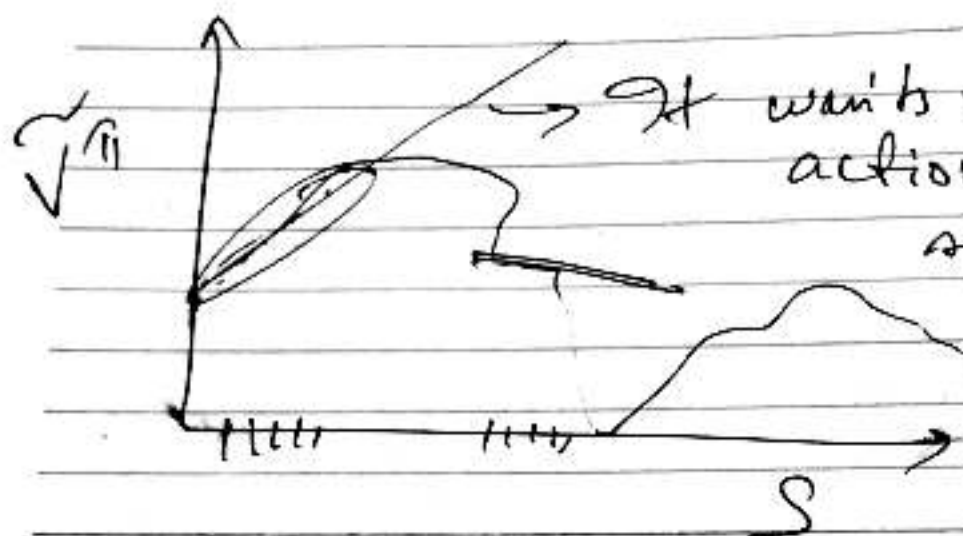
→ ① Bootstrapping increases error (feedback) we can fix this

② where should the sample set be?

For problem ①

the problem in FPI case

happens in the policy improvement
ex. (in the outer loop)



It wants to change the actions for every state such that the next state you land in has a higher value.

But you don't have values in that area.

Fix

High level Idea:-

cache policies ~~not~~

value functions.

We do this in both FPI and FSS

* An optimal policy has the property that no matter what you do now, hereafter

you must act optimally.

Alg 2 to fix

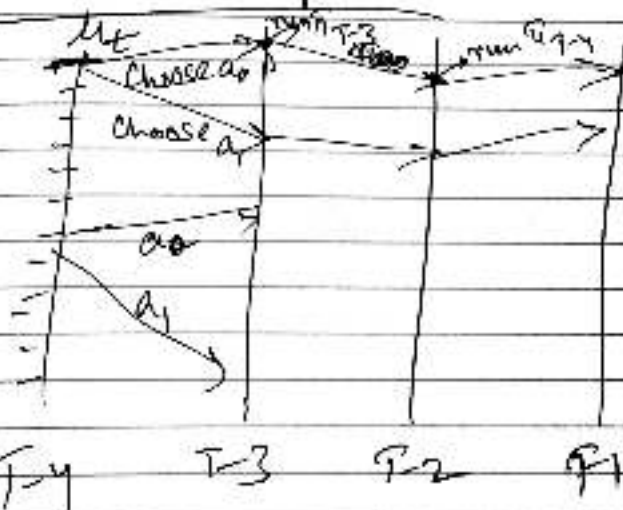
Policy Search by Dynamic Prog. ^{may not be practical as tied to finite horizon.}

for each $t \in T-1, \dots, 0$:

$d_t = \text{GenerateSamples}(\pi_{\{t+1, \dots, T-1\}}, d_t)$

$\pi_t = \text{Learn}(d_t)$

Generate Samples



$\tilde{Q}^{\pi}(a_0, s, T-4)$

$\tilde{Q}^{\pi}(a_1, s, T-4)$

Choose the \tilde{Q} as the higher reward

Learn (d_t)
① We can say that this is a
classification problem.

Find which state at $T-1$ maps
to which action ~~at~~.

Weight the classification problem as.

$$Q^{\pi}(a_1, s, T-1) - \underline{Q^{\pi}(a_0, s, T-1)}$$

For multiple
weights make
this the best
action.

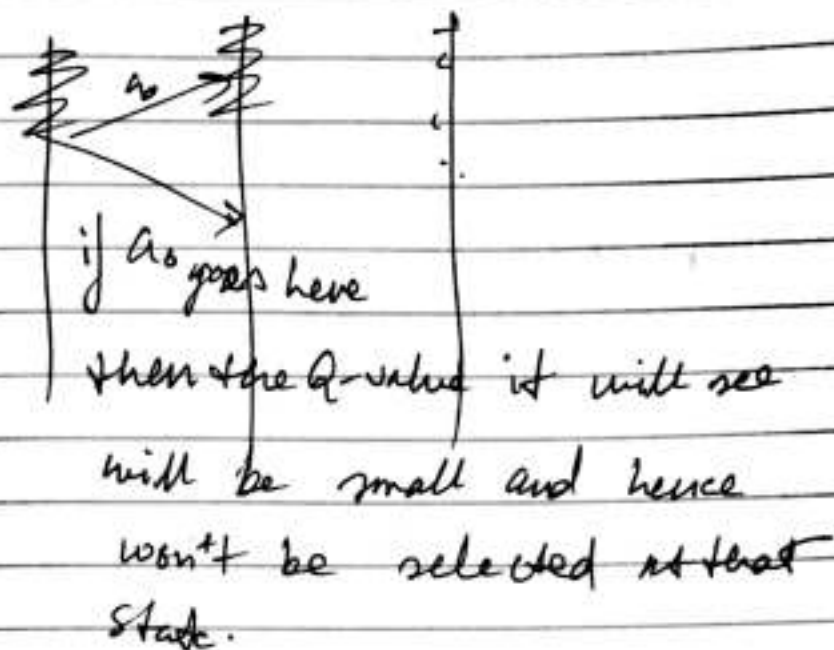
② Regression on Q 's

$$a(s, t) = \arg \max_a \tilde{Q}(s, a, t)$$

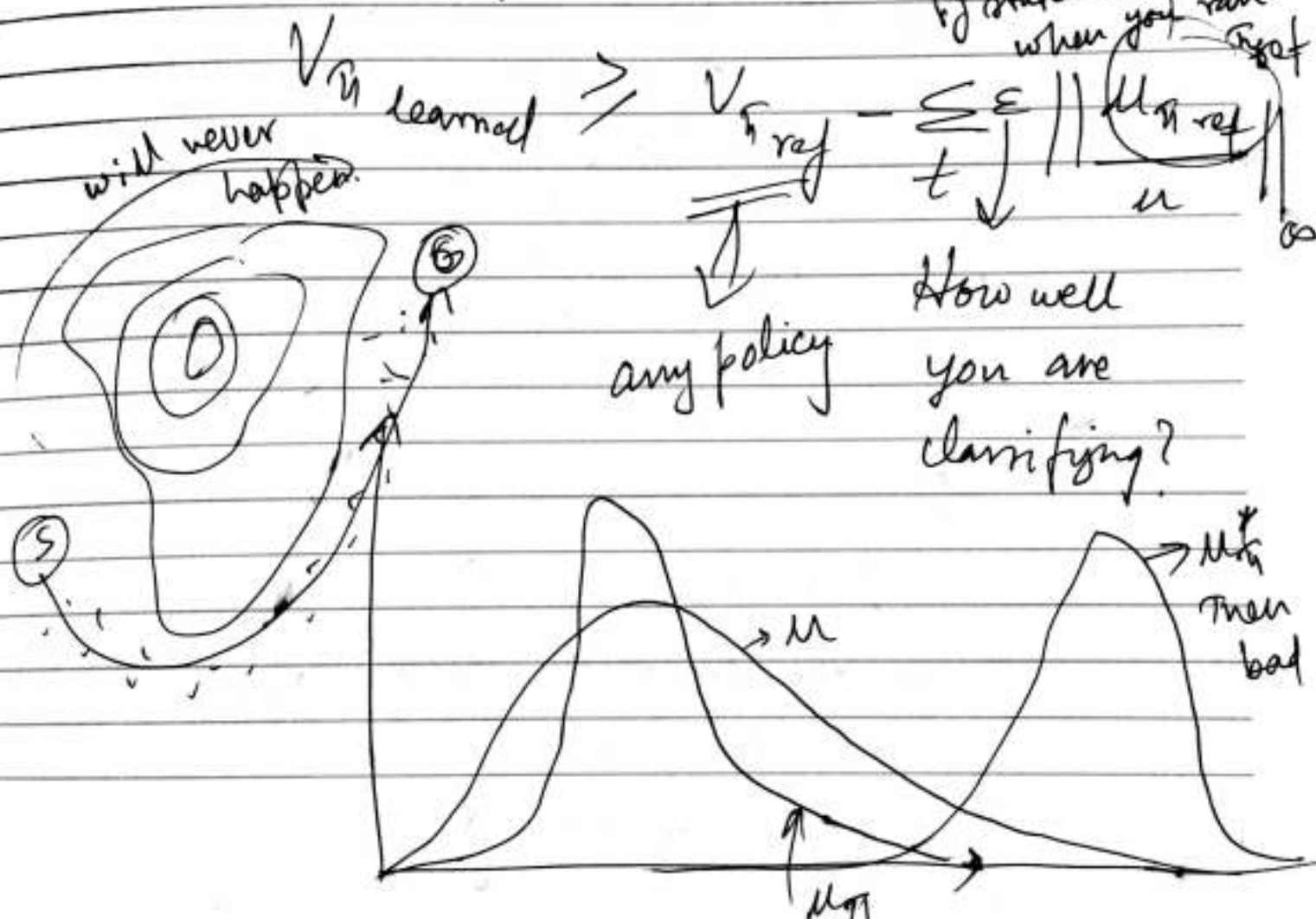
* This is different than FVI, FQI,

* This is more expensive as we have to do
 T times more work $[O(T^2)]$
but no bootstrapping - stable.

* Errors in bootstrapping are ~~low~~ propagated but rollouts are not.



* Performance guarantee



Algo 2 to fix

Conservative Policy Iteration

Evaluate Policy

Improve Policy

$$\pi' \leftarrow \arg \max_a \tilde{Q}^{\pi}(s, a)$$

Conservative

$$\pi_{\text{new}} \leftarrow \alpha \pi' + (1 - \alpha) \pi$$

↑
mixing parameter.

$$\text{where } \alpha \approx O\left(\frac{1}{T}\right)$$

→ $\tilde{Q}^{\pi} = ?$

Sample from π

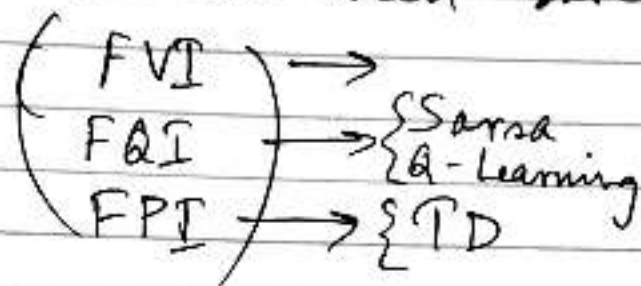
Run the current policy π_{new}
for some # of steps.



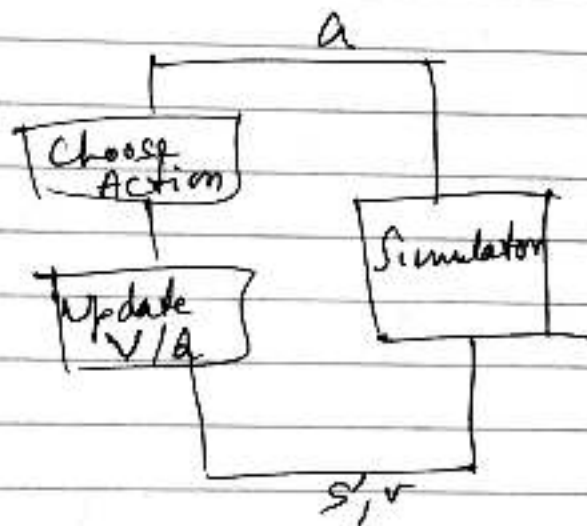
with $(1-\gamma)$ probability you stop and try out actions. γ is the discount factor

Temporal Difference Methods:-

Upto now we have seen batch methods



Overview



* PD: policy evaluation.

given π , estimate $V_{\pi}(s)$, $\forall s$

Goal
Find st. $V_{\pi}(s) = R(s, \pi(s)) + \gamma \mathbb{E}_{s'} [V_{\pi}(s')]$