2) continuous state space
e-q. qubit example
$$S = \int (O_1 E) / O E(O_1 \pi),$$

 $E \in [0, 27)]$
• cannot explore all states (unconstably many)
=) function approx methods required
similarly, A can be either discrete on
eon timerus:
1) discrete action space $t = \{a_1, ..., a_{|A|}\}$
e.g. noves along Gridworld $\{S, N, E, W\}$
application quantum gakes $\{M_0, M_Y, M_2\}$
note: the number of trajectories τ in
an episodic task of T skps
is $|A|^T$, i.e. exponential
problem: depending on the sparsity of
rewords, exploring a lange
trajectory space care be infeasible

-2-

e-g. a, az ... a_N A works well if dim $t \leq 2$ e-g. if $A = [0, 1]^3$, i.e. the unit cube use linear grid of 10 points per dimension =2 a total of 10³ actions for $A = [0, 1]^3$

b)
$$goal:$$
 find RL algorithms to handle
continuous A
- Policy Greedient Methods (incl. actor-critic)
 $T_0(a/s): t \times S \longrightarrow CO, 1]$
should parametrize a continuous distribution
over action space A
I) unbounded action spaces:
e.g. $t = R$ for $d=1$
singulast possible with distr.: Gaussian
 $T_0(a/s) = \frac{1}{\sqrt{2\pi} \sigma_p^2(s)} exp\left[-\frac{(a - \mu_0(s))}{2\sigma_0^2(s)}^2\right]$
uniquely fixed by its first & second moments:
 $Mo: S \longrightarrow R$ mean
 $\sigma_0: S \longrightarrow [O, \infty)$ variance
 $\frac{Note: \mu(s), \sigma(s)}{\mu(s)}$

-3 -

• every state s has its own policy dilder.
to learn
idea: hind variational approximations to

$$\mu(s), \sigma(s) \longrightarrow \mu_{\theta}(s), \sigma_{\theta}(s)$$

e.g. NN:
 $s \longrightarrow NN_{\theta}$
 $s \longrightarrow NN_{\theta}$
 $s \sigma_{\theta}(s)$
having $\mu_{\theta}(s) \ 8 \sigma_{\theta}(s), sample an action
a from Grassian policy
 $\pi_{\theta}(a|s) = \mathcal{N}(\mu_{\theta}(s), \sigma_{\theta}^{2}(s))$
problem: Grassian distr. has unbounded
support
=) action can take any value in (- ∞, ∞)
what if we need bounded continuous actions?
 \overline{s}
 $bounded$ continuous action spaces:
e.g. $A = \overline{L}0, 17$, or $\overline{L} \in Ca, b7$$

idea : parametrize The by a distribution
with compact support
caveat : we have to be able to sample
actions from it easily:
eq.
$$B(a; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} a^{\alpha-1} (1-\alpha)\beta^{3-1}$$

Beta distribution with params α, β
where $\Gamma(r) = \int_{0}^{\infty} x^{2-1} e^{-x} dx$ Gamma
function
(continuous generalization
of the factorial u!)

$$\frac{bounded \quad support:}{i) \quad CO, 1] \quad hor \quad \alpha_{1}\beta > 1$$

$$ii) \quad (O, 1) \quad hor \quad O \quad \alpha_{1}\beta < 1$$

$$\pi_{O}(\alpha | s) = B(\alpha; \quad \alpha_{O}(s), \beta_{O}(s))$$

$$\pi_{O}(\alpha | s) = B(\alpha; \quad \alpha_{O}(s), \beta_{O}(s))$$

$$\pi_{O}(\alpha | s) = B(\alpha; \quad \alpha_{O}(s), \beta_{O}(s))$$

$$greedy \quad \pi_{O}(\alpha | s) = greedy \quad \pi_{O}(\alpha |$$

• properties:
mean :
$$Ea-B(a; d, \beta)[a] = \frac{d}{d+\beta}$$
 survey
variance: $Vara-B(a; d, \beta)[a] = \frac{d\beta}{(d+\beta)^2(d+\beta+1)}$



way ont : 1) Gaussian mixture models
-> lecture 4, notebook 1
2) normalizing flow models, etc.
Tang & Agrawal arXiv: 1809.10326
- Q-learning with Continuous Actions:
recall: DQN involves the operations:

$$\pi(als) = \begin{cases} 1 & if a = arguax Qy(sca) \\ 0 & otherwise \end{cases}$$

target values:
$$Y_j = V_j + y \max_{a'} Q_p(s',a')$$

-> Q-learning untains a maximization
procedure over actions
problem: max, arguax require separate
optimization procedure in the
case of continuous actions!
-> 3 options to conjute near for
A continuous

$$\frac{\text{Optim 2}}{\text{-sure hunchion classes for Q whichallow to analytically comparts the maxnormalized advantage functions (NAFe) $\text{Q}(s,a) = -\frac{1}{2}(a - \mu_{q}(s))^{t}P_{q}(s)(a - \mu_{q}(s)) + V_{q}(s)$
 $s \rightarrow DNN_{q} \implies P_{q} \in \mathbb{R}^{d}$, $d = dim(A)$
 $y_{q} \in \mathbb{R}$$$

- 9 -

Detion 3
-> use a function approximator to
learn doing the maximitation
-algorithm: Deep Deterministic Policy Goodient
(DDPG)
-> misnomer
recall: notation for deterministic policies

$$\pi: S \rightarrow A$$
, $a = \pi/s$
idea : trein a deterministic policy network
 $Mo(s)$, s.t.
 $Mo(s) \approx angmax Ro(s,a)$
. how do we find the optimal paraens Θ ?
solve: $\Theta \leftarrow argupx R_{D}(s, Mo(s))$
 $Theld tixed$
 $R_{O}(s, Mo(s)) = R_{O}MO dR_{O}$
 $done antomatically in
any auto-diff package, e.g. JAX
-10-$

- ver torget for DQN algorithm
recall: max Q(s,a) = Q(s, augmax Q(s,a))
det

$$= v_{i} + y Q_{p}(s_{i}', augmax Q_{p}(s_{j}', a_{j}'))$$

 $= v_{i} + y Q_{p}(s_{i}', \mu_{0}(s_{j}'))$
- pseudo code (DPP6)
1. take action aj using some policy,
observe (sj, aj, rj, s_{j}'); add it to butter B.
2. sample minibatele from B uniformly
3. comparte $q_{i} = v_{i} + y Q_{p'}(s_{i}', \mu_{0'}(s_{j}'))$
 $= v_{i} + y Q_{p}(s_{j}, a_{j}) (Q_{p}(s_{j}, a_{j}) - q_{i})$
 $= v_{i} + y Q_{p'}(s_{j}, \mu_{0'}(s_{j}))$
 $= v_{i} + y Q_{p'}(s_{j}, \mu_{0'}(s_{j}))$
 $= v_{i} + y Q_{p'}(s_{j}, \mu_{0'}(s_{j}))$
 $= v_{i} + y Q_{p'}(s_{j}, a_{j}) - q_{i})$
 $= v_{i} + y Q_{p}(s_{j}, a_{j}) (Q_{p}(s_{j}, a_{j}) - q_{i})$
 $= 0 + \beta \sum_{i} v_{0} Q_{p}(s_{j}, \mu_{0}(s_{j}))$
 $= u_{i} + v_{i} + v_{i}$