

# Deterministic Policy Gradient Algorithms

Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., & Riedmiller, M. (2014, June). In ICML.

2015/8/20

D1 金子 貴輝

# 選考理由

- ICML 2014の論文
- 高次元入出力に対応できる強化学習
- DeepMindの中の人々のグループの論文
- (証明のAppendixが見つからず)

# 内容

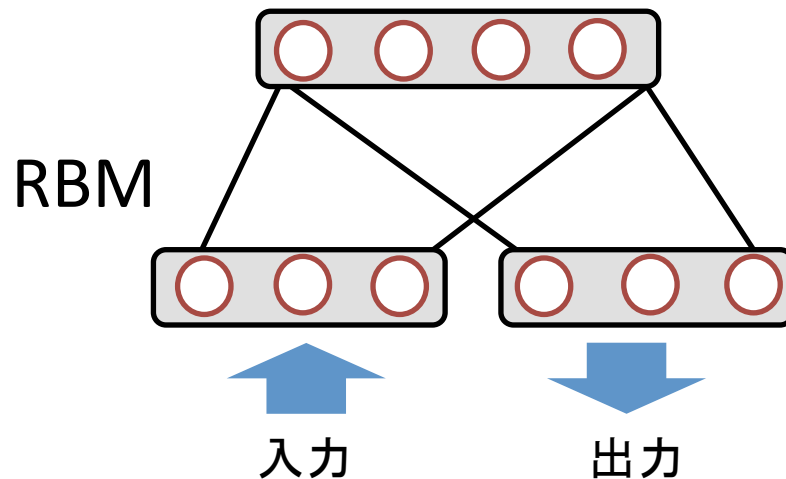
- 背景・課題
  - 高次元入出力の強化学習の課題
- 着想
  - 確率方策＋方策オン → 確定方策＋方策オフ
- 提案手法
  - 確定方策勾配法
  - 方策オフ型確定方策勾配法
  - 互換方策オフ型確定的アクタークリティック(COPDAC)
- 実験
  - 確定方策の有効性の検証タスク
  - 標準的な強化学習タスク
  - 高次元入出力タスク(Octopus Arm)
- 結論

# 内容

- 背景・課題
  - 高次元入出力の強化学習の課題
- 着想
  - 確率方策＋方策オン → 確定方策＋方策オフ
- 提案手法
  - 確定方策勾配法
  - 方策オフ型確定方策勾配法
  - 互換方策オフ型確定的アクタークリティック(COPDAC)
- 実験
  - 確定方策の有効性の検証タスク
  - 標準的な強化学習タスク
  - 高次元入出力タスク(Octopus Arm)
- 結論

# 著者の先行研究

- 2つの強化学習アルゴリズム(ENATDAC, EQNAC)
  - 方策勾配法の一つであるNatural Actor-Critic法を使用
  - 方策分布にRBMを使用し, アルゴリズムをさらに単純化
  - RBMなので**高次元入出力**にも対応できる
    - 実験では入力次元34, 出力次元14のタコ腕タスクも解けた

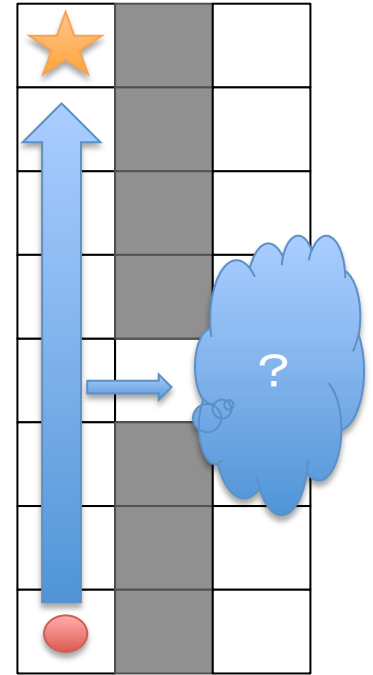




# 先行研究の問題点

従来の方策勾配法では、  
学習が進むに連れて  
方策が確定的になり、  
方策勾配の推定が困難になる

(最適でない行動時の予測が、サンプル数の不足で、不正確になるためだと思われる)



迷路問題での  
確定的な動作例

# 着想

問題点を2つに分ける

- 行動が確定的なときの積分の近似(サンプリング)が不正確
  - **確定的な方策**なら, 積分の必要性が生じない
- 行動が学習とともに確定的になっていく
  - **方策オフ型**にして, 学習の進みと方策の確定さを切り離す



# 内容

- 背景・課題
  - 高次元入出力の強化学習の課題
- 着想
  - 確率方策＋方策オン → 確定方策＋方策オフ
- 提案手法
  - 確定方策勾配法
  - 方策オフ型確定方策勾配法
  - 互換方策オフ型確定的アクタークリティック(COPDAC)
- 実験
  - 確定方策の有効性の検証タスク
  - 標準的な強化学習タスク
  - 高次元入出力タスク(Octopus Arm)
- 結論

# 確定的な方策での方策勾配法

- 従来, 決定方策勾配は存在しないか, 環境モデルを使うときに限られると考えられてきた[Peters, 2010]
- 状態の関数で確定的な方策  $\mu_\theta(s)$  を定義
- この方策で方策勾配法を再定式化

$$\begin{aligned} J(\mu_\theta) &= \int_{\mathcal{S}} \rho^\mu(s) r(s, \mu_\theta(s)) ds \\ &= \mathbb{E}_{s \sim \rho^\mu} [r(s, \mu_\theta(s))] \end{aligned}$$

Deterministic Policy Gradient Theorem

$$\begin{aligned} \nabla_\theta J(\mu_\theta) &= \int_{\mathcal{S}} \rho^\mu(s) \nabla_\theta \mu_\theta(s) \nabla_a Q^\mu(s, a)|_{a=\mu_\theta(s)} ds \\ &= \mathbb{E}_{s \sim \rho^\mu} \left[ \nabla_\theta \mu_\theta(s) \nabla_a Q^\mu(s, a)|_{a=\mu_\theta(s)} \right] \end{aligned}$$

# 方策オフ型確定方策勾配法

- 報酬を最大化する方策 $\mu$ と, 探索のための方策 $\beta$ を分けて, 期待報酬を定義
- 単純に期待値が置き換わる

$$\begin{aligned} J_{\beta}(\mu_{\theta}) &= \int_{\mathcal{S}} \rho^{\beta}(s) V^{\mu}(s) ds & V^{\pi}(s) &= \mathbb{E}[r_1^{\gamma} | S_1 = s; \pi] \\ &= \int_{\mathcal{S}} \rho^{\beta}(s) Q^{\mu}(s, \mu_{\theta}(s)) ds & Q^{\pi}(s, a) &= \mathbb{E}[r_1^{\gamma} | S_1 = s, A_1 = a; \pi] \end{aligned}$$

Off-policy Deterministic Policy Gradient Theorem

$$\begin{aligned} \nabla_{\theta} J_{\beta}(\mu_{\theta}) &\approx \int_{\mathcal{S}} \rho^{\beta}(s) \nabla_{\theta} \mu_{\theta}(a|s) Q^{\mu}(s, a) ds \\ &= \underline{\mathbb{E}_{s \sim \rho^{\beta}}} \left[ \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu}(s, a) \Big|_{a=\mu_{\theta}(s)} \right] \end{aligned}$$

# 提案手法: COPDAC (COPDAC-Q)

- 方策オフ型確定的方策勾配法に基づく
- アクタークリティック法で価値関数を近似

$$Q^\pi(s, a) \approx Q^w(s, a) = \phi(s, a)^\top w$$

- 基底関数  $\phi$  は方策と互換になるように選ぶ (近似精度のため)  
$$\phi(s, a) = a^\top \nabla_{\theta} \mu_{\theta}(s)$$
- 方策と価値関数のパラメータは以下のように更新

$$\begin{aligned}\delta_t &= r_t + \gamma Q^w(s_{t+1}, \mu_{\theta}(s_{t+1})) - Q^w(s_t, a_t) \\ \theta_{t+1} &= \theta_t + \alpha_{\theta} \nabla_{\theta} \mu_{\theta}(s_t) (\nabla_{\theta} \mu_{\theta}(s_t)^\top w_t) \\ w_{t+1} &= w_t + \alpha_w \delta_t \phi(s_t, a_t) \\ v_{t+1} &= v_t + \alpha_v \delta_t \phi(s_t)\end{aligned}$$

※ 状態価値関数のパラメータ  $v$  の扱いはわかりませんでした  
 $\phi(s_t)$  は状態の特徴量

# 内容

- 背景・課題
  - 高次元入出力の強化学習の課題
- 着想
  - 確率方策＋方策オン → 確定方策＋方策オフ
- 提案手法
  - 確定方策勾配法
  - 方策オフ型確定方策勾配法
  - 互換方策オフ型確定的アクタークリティック(COPDAC)
- 実験
  - 確定方策の有効性の検証タスク
  - 標準的な強化学習タスク
  - 高次元入出力タスク(Octopus Arm)
- 結論

# 実験—連続値バンディット

目的: 確率方策勾配と確定方策勾配の比較

環境:

- バンディット: 状態を持たずすぐに報酬が得られるタスク
- 報酬は正しい出力値とのマハラノビス距離 $\times(-1)$
- パフォーマンスはステップごとの平均報酬
- 出力次元は 10, 25, 50

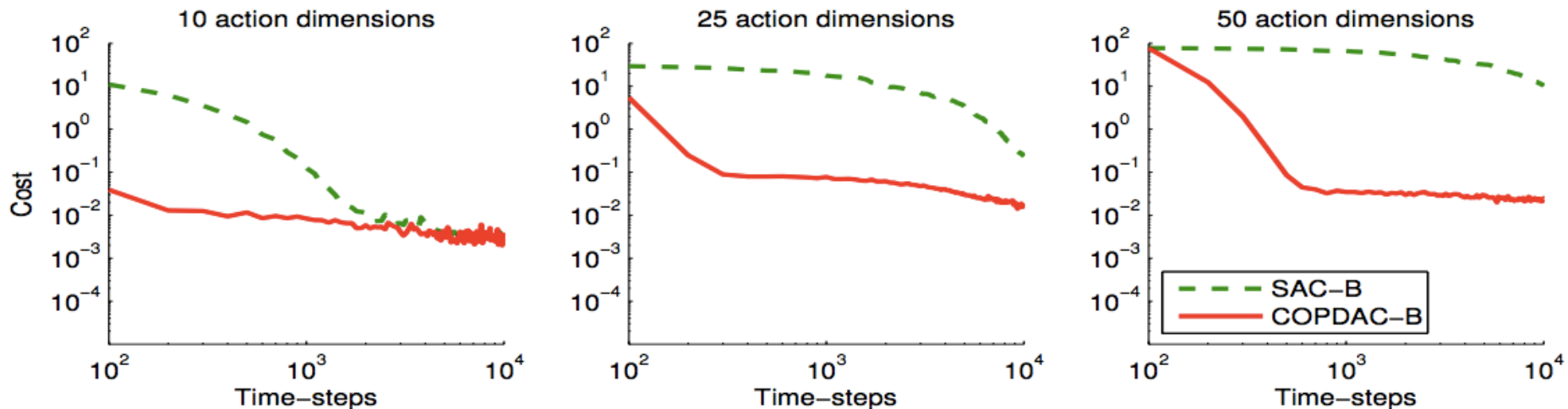
アルゴリズム:

- ガウス分布の方策を使用
- 互換な行動価値関数近似をする
- 方策オン/オフ型による方策の違いは以下

	確率方策(SAC-B)	確定方策(COPDAC-B)
ガウス分布の方策	平均と分散が学習される	分散固定で探索, 平均が学習される

# 実験一連続値バンディット

- 最適なパラメータ, ステップサイズでコストを比較
  - 5試行での平均
- 次元が増えるほど, 確定方策のほうが良くなる



学習が進むほど勾配推定が困難になる確率方策より,  
確定方策なら速く収束する

# 実験一連続値強化学習

目的: 一般的なタスクでの比較

環境:

- 山登りカー, 振り子, 沼地, のベンチマーク
- エピソードは5000ステップ制限
- 割引率0.99(山登りカーと振り子)0.999(沼地)
- 領域外への移動は制限される

アルゴリズム

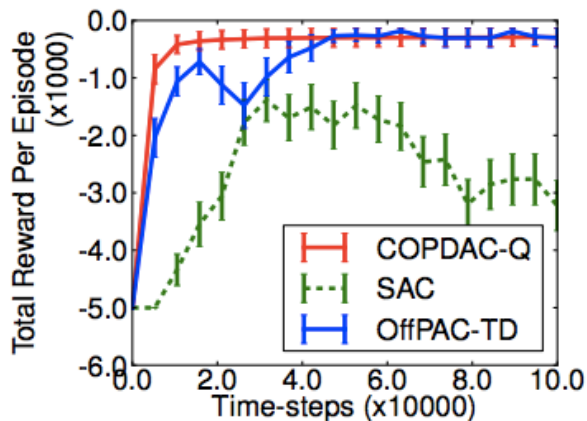
- 特徴量 $\phi$ は状態空間のタイルコーディングで計算される

	方策オン 確率方策	方策オフ 確率方策	方策オフ 確定方策
アルゴリズム	SAC(山登りカーでベストな手法(Degris2012a))	OffPAC	COPDAC-Q
方策	$\mathcal{N}(\theta^\top \phi(s), \exp(y^\top \phi(s)))$	学習方策: 同左 探索方策: 同右	学習方策: $\theta^\top \phi(s)$ 探索方策: $\mathcal{N}(\theta^\top \phi(s), \sigma_\beta^2)$

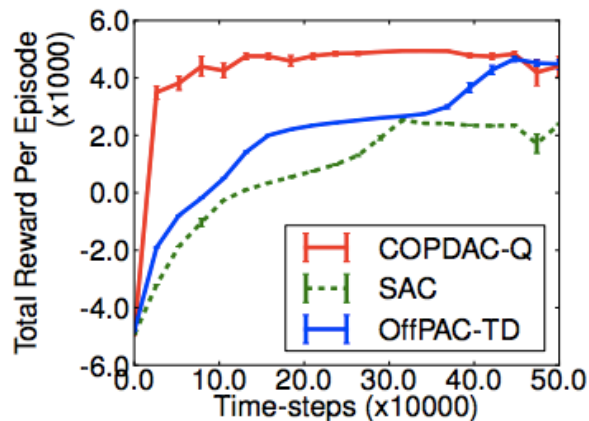


# 実験一連続値強化学習

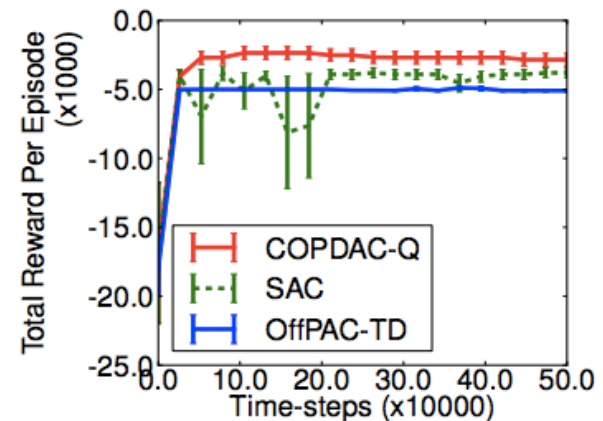
- 最適なパラメータでの比較
  - 30試行での平均
  - エピソード合計報酬の時間変化をグラフ化
- COPDAC-Qがすべての領域で勝っている



(a) Mountain Car



(b) Pendulum

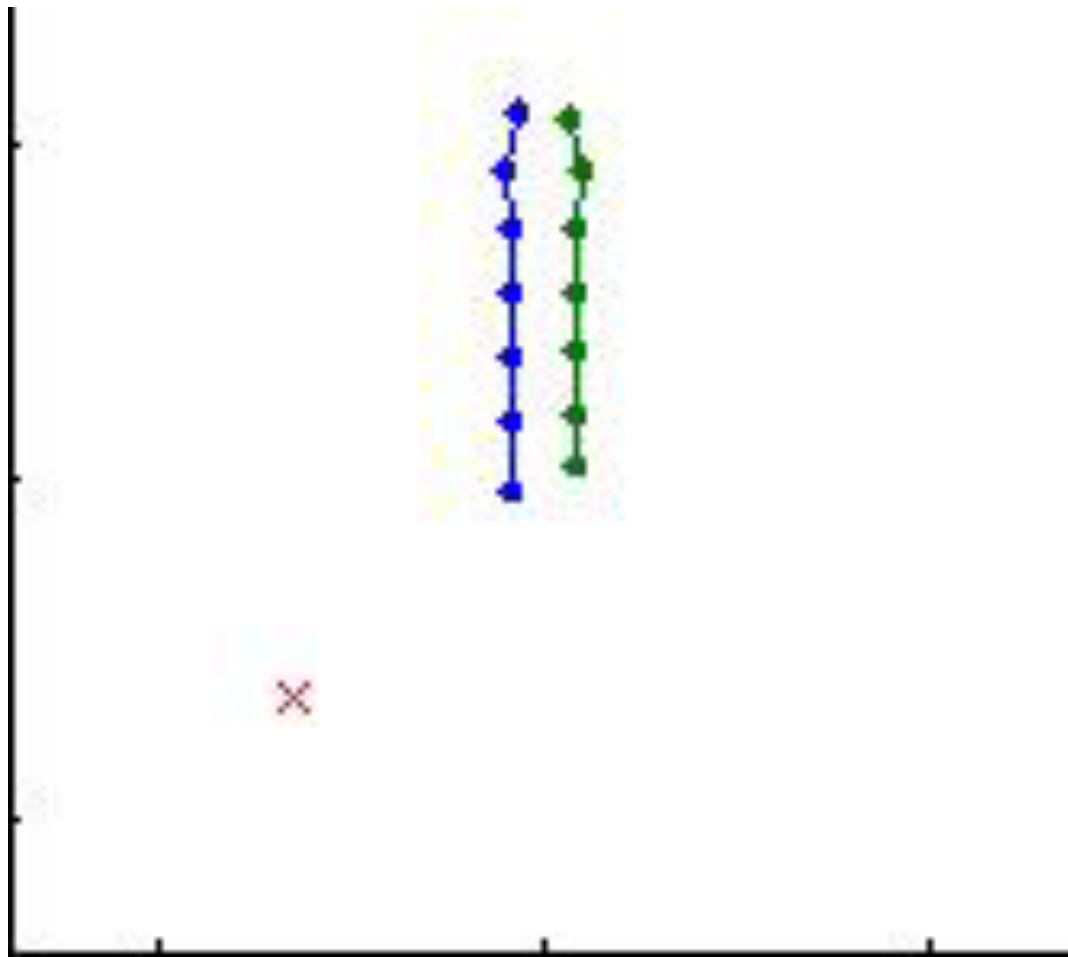


(c) 2D Puddle World

# 実験一タコ腕

- 目標に当たるようシミュレーション上のタコの腕を制御する
  - 6関節に分かれ根本が固定されている
  - 50次元の状態空間
  - 20次元の行動空間(背側/腹側/中央)
  - 報酬は対象と腕の距離で変わる
  - 触れると報酬は+50かつエピソード終了
  - 300ステップで打ち切り
- COPDAC-Qを使用
  - 8ユニットのMLPを方策関数に使用
  - 状態価値関数は40ユニット線形出力のMLPで近似

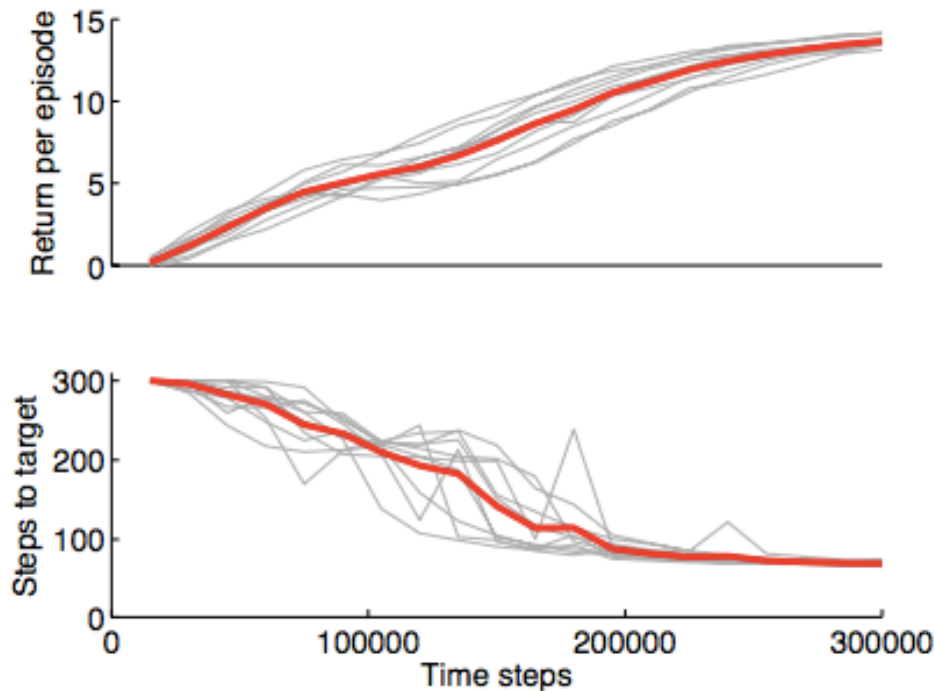
# Octopus Arm Task



# 実験一タコ腕

## 10回の訓練と平均のグラフ

- すべての場合で制御が学習できている



従来手法では行動のマクロを使って単純化したり,  
4関節までの低次元に制限したりしていた



6関節までの学習に成功した

# 結論

- 確定方策勾配法のフレームワークを提案
- 行動空間の積分を避けて確率方策より効率よく勾配推定できる
- 実際に50次元の連続行動空間でのバンディットで圧倒的に勝利
- 50の状態空間, 20の行動空間の難しい課題も解ける