# Iowa State University

## AltGDmin for Policy Evaluation in Offline Reinforcement Learning

Komal Krishna Mogilipalepu

Supervised by Prof. Namrata Vaswani

# Outline

- Offline Reinforcement Learning

- Policy Evaluation for Offline Reinforcement Learning

- Low-Rank Structure on Markov Decision Process

- Offline Dataset Generation

- Low-Rank Matrix Completion using Alternating Minimization

- Policy Evaluation Algorithm with AltGDmin

- Simulation Results of AltGDmin and Altmin

- Simulation Results of Policy Evaluation with Alternating Minimization

- Conclusion, Future Work, and References

- Acknowledgements

➢ Finite discrete Markov decision process (MDP) $\mathcal{M} = (\mathbf{s}, \mathbf{a}, \mathbf{P}, \mathbf{R}, \boldsymbol{\mu}, H)$

- state space $\mathbf{s}$ such that state $s \in \mathbf{s}$,
- action space $\mathbf{a}$ such that action $a \in \mathbf{a}$,
- probability transition kernel $\mathbf{P} = \{\mathbf{P}_t : \mathbf{s} \times \mathbf{a} \times \mathbf{s} \rightarrow [0,1]\}_{t \in [H]}$,
- bounded reward function $\mathbf{R} = \{\mathbf{R}_t : \mathbf{s} \times \mathbf{a} \rightarrow [0,1]\}_{t \in [H]}$,
- initial state distribution $\boldsymbol{\mu}$,
- horizon $H$.

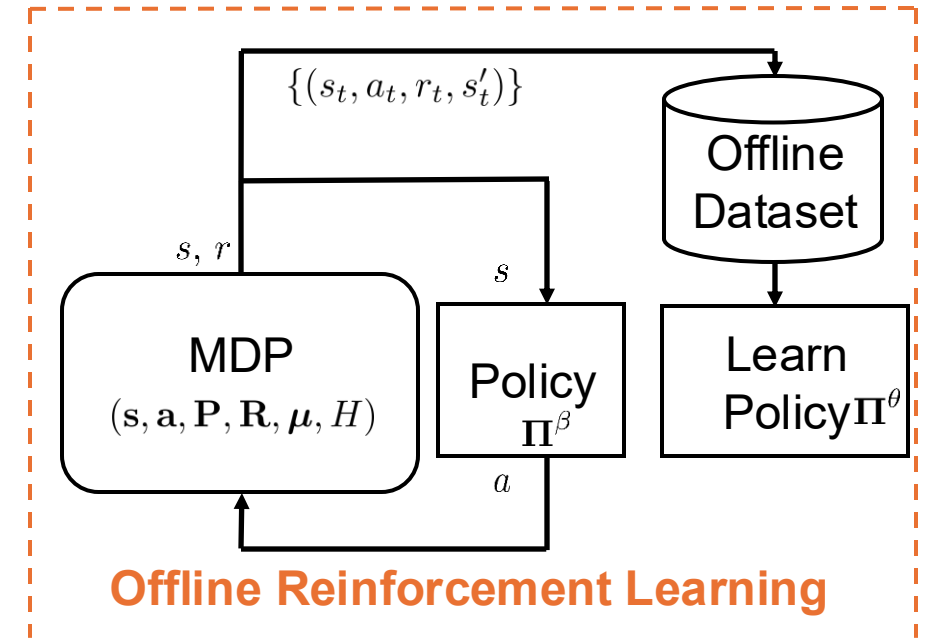➢ Policy $\mathbf{\Pi} : \mathbf{s} \times \mathbf{a} \rightarrow [0,1]$.

➢ Total expected reward with respect to a policy $\mathbf{\Pi}$ is

$$J^{\mathbf{\Pi}} = \mathbb{E}_{\mathbf{\Pi}}[\textstyle\sum_{t=1}^{H} \mathbf{R}(s_t, a_t)|s_1 \sim \boldsymbol{\mu}].$$

➢ State-action value in the $t$-th step according to policy $\mathbf{\Pi}$ is

$$\mathbf{Q}_t^{\mathbf{\Pi}}(s,a) = \mathbb{E}_{\mathbf{\Pi}}[\textstyle\sum_{i=t}^{H} \mathbf{R}_i(s_i, a_i)|s_t = s, a_t = a].$$



$\{(s_t, a_t, r_t, s_t')\}$

Offline Dataset

$s, r$    $s$

MDP $(\mathbf{s}, \mathbf{a}, \mathbf{P}, \mathbf{R}, \boldsymbol{\mu}, H)$

Policy $\mathbf{\Pi}^\beta$

Learn Policy$\mathbf{\Pi}^\theta$

$a$

**Offline Reinforcement Learning**

$\mathbf{\Pi}^\beta$ - unknown behavior policy

$\mathbf{\Pi}^\theta$ - target policy

Learn a target policy that generalize the pattern observed in the offline dataset.

Levine, S., Kumar, A., Tucker, G., & Fu, J. (2020). Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643.*

Offline RL → Policy Evaluation → Policy Deployment

➤ **Distributional shift:**

❖ Evaluating the target policy out of offline data distribution.

➤ **Concentrability coefficient:**

❖ The coverage of the dataset is measured by the concentrability coefficient $\mathcal{C}^{\mathbf{\Pi}} = \max_{s,a} \frac{\mathbf{D}^{\mathbf{\Pi}}(s,a)}{\widehat{\mathbf{D}}(s,a)}$

where, $\mathbf{D}^{\mathbf{\Pi}}(s,a)$ - state-action occupancy measure by a policy $\mathbf{\Pi}$ and

$\widehat{\mathbf{D}}(s,a)$ - empirical measure of a state-action pair observed in the offline data.

❖ Full coverage of the dataset $\mathcal{C}^{\mathbf{\Pi}} < \infty \; \forall \; \mathbf{\Pi}$.

❖ Partial coverage of the dataset $\mathcal{C}^{\mathbf{\Pi}^{\star}} < \infty$ for an optimal policy $\mathbf{\Pi}^{\star}$.

❖ Infinite concentrability coefficient $\mathcal{C}^{\mathbf{\Pi}} = \infty \; \forall \mathbf{\Pi}$.

How to evaluate the target policy with infinite concentrability coefficient?

➢ **Assumption 1:**

The reward matrix $\mathbf{R}_t$ has rank at most $d/2 \; \forall t.$

➢ **Assumption 2:**

The state transition kernel $\mathbf{P}_t$ has following representation

$$\mathbf{P}_t(s, a, s') = \mathbf{u}_t(s', s)^\mathrm{T} \mathbf{w}_t(a)$$

or $\mathbf{P}_t(s, a, s') = \mathbf{u}_t(s)^\mathrm{T} \mathbf{w}_t(s', a) \; \forall \, t, s', s, a,$

where, $\mathbf{u}_t(\cdot)$ and $\mathbf{w}_t(\cdot)$ are unknown functions, which maps the input to a $d/2$ length vector.

---

**Algorithm 1:** Low-Rank Markov decision process

1  $\mathrm{LRMDP}(S, A, H)$
2  $\mathbf{s} = [S], \; \mathbf{a} = [A]$
3  $\boldsymbol{\mu} = [S]/S$
4  reward $\mathbf{R} = (\mathrm{uniform}(0, 1)_{\mathbf{s} \times 1} 1_{1 \times \mathbf{a}}^T)$
5  Initialize state transition kernel $\mathbf{P} = [\;]$
6  **for** $t \leftarrow 1$ **to** $H$ **do**
7  $\quad$ Append $1_{\mathbf{s} \times \mathbf{a} \times \mathbf{s}} * \frac{1}{S}$ to $\mathbf{P}$
8  **end**
9  return $(\mathbf{s}, \mathbf{a}, \mathbf{P}, \mathbf{R}, \boldsymbol{\mu})$

---

$[S] = [1, 2, 3, \ldots, S]$ and $\times$ denotes cartesian product.

➢ **Assumption:** The state-action value matrix $\mathbf{Q}^{\mathbf{\Pi}}$ is at most rank $d$ for any policy $\mathbf{\Pi}$,

$\qquad$ - alleviates finding the unknown functions $\mathbf{u}_t(\cdot)$ and $\mathbf{w}_t(\cdot)$ .

➢ **Objective:**

Estimate total expected reward $J^{\mathbf{\Pi}^\theta}$ for a target policy $\mathbf{\Pi}^\theta$ from the offline dataset that is generated by an unknown behavior policy $\mathbf{\Pi}^\beta$.

➢ Uniform transition model

$\boldsymbol{\mu}_t(s) = 1/n \; \forall \; s \in \mathbf{s}$ here, $n = S = A,$

and $\mathbf{P}_t(\cdot|s,a) = 1/n \; \forall s \in \mathbf{s}, \, a \in \mathbf{a}.$

➢ Uniform policy

- for every $t$ and $\forall \, s \in \mathbf{s}$, target policy $\boldsymbol{\Pi}_t^\theta(\cdot|s)$ sample an action $a$ uniformly from the action subset $\mathbf{a}_t^\theta$ of size $m$, which itself sampled uniformly from action space $\mathbf{a}$ i.e.

$\boldsymbol{\Pi}_t^\theta = \{1, 0\}_{S \times A} \sim \text{Bernoulli}\left(\frac{m}{n}\right) * \frac{1}{m},$

- generate behavior policy from the same target policy model independently i.e.

$\boldsymbol{\Pi}_t^\beta = \{1, 0\}_{S \times A} \sim \text{Bernoulli}\left(\frac{m}{n}\right) * \frac{1}{m}.$

➢ The uniform policy model resulting to an infinite concentrability coefficient

$\mathcal{C}^{\boldsymbol{\Pi}} = \max_{s,a} \frac{\mathbf{D}^{\boldsymbol{\Pi}_t^\theta}(s,a)}{\mathbf{D}^{\boldsymbol{\Pi}_t^\beta}(s,a)} = \infty.$

---

**Algorithm 2:** dataset generation

**Input:** number of states $S$ and actions $A$, number of trajectories $K$, horizon $H$, action subset size $m$.

**Output:** dataset $\mathcal{D}$

1   Initialize dataset $\mathcal{D} = [\,]$
2   $(\mathbf{s}, \mathbf{a}, \mathbf{P}, \mathbf{R}, \boldsymbol{\mu}) = \texttt{MDP}(S, A, H)$
3   **for** $k \leftarrow 1$ **to** $K$ **do**
4      Initialize trajectory $\tau = [\,]$
5      $s_0 \sim \boldsymbol{\mu}(\mathbf{s})$
6      **for** $t \leftarrow 1$ **to** $H$ **do**
7          $a \sim \text{uniform}(\mathbf{a}_t^\theta), \; \text{where} \; |\mathbf{a}_t^\theta| = m, \; \text{and} \; \mathbf{a}_t^\theta \subseteq \mathbf{a}$
8          $r = \mathbf{R}_t(s_0, a)$
9          $s_1 \sim \mathbf{P}_t(\mathbf{s}|s_0, a)$
10        Append $(s_0, a, r, s_1)$ to $\tau$
11        $s_0 \leftarrow s_1$
12      **end**
13      Append trajectory $\tau$ to dataset $\mathcal{D}$
14 **end**

---

➢ Offline data $\mathcal{D} = \{(s_t^k, a_t^k, r_t^k, s_t'^k)\} \; \forall \; t \in H, k \in K.$

➢ Compose a low-rank matrix $\mathbf{X} = \mathbf{UB}$ where, $\mathbf{U} \in \mathbb{R}^{n \times r} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_{r \times r})$ and $\mathbf{B} \in \mathbb{R}^{r \times q} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_{r \times r})$.

➢ Observations $\mathbf{Y} := \mathbf{M} \circ (\mathbf{X} + \mathbf{N})$,

where, $\mathbf{M} \in \{0,1\}^{n \times q}$ is a Bernoulli matrix with probability $p$ and $\mathbf{N} \in \mathbb{R}^{n \times q} \sim \sigma \mathcal{N}(0,1)$ with standard deviation $\sigma$.

➢ Task: retrieve $\mathbf{X}$ from $\mathbf{Y}$.

➢ Optimization problem: $\min\limits_{\mathbf{B}, \mathbf{U}^{\mathrm{T}}\mathbf{U}=\mathbf{I}} \|\mathbf{Y} - \mathbf{M} \circ \mathbf{UB}\|_F^2$.

➢ Metric: subspace distance measure between $\mathbf{U}^{(I)}$ and $\mathbf{U}$,

$\mathrm{SD}(\mathbf{U}^{(I)}, \mathbf{U}) = \|(\mathbf{I} - \mathbf{U}^{(I)}\mathbf{U}^{(I)^{\mathrm{T}}})\mathbf{U}\|_F$.

➢ Altmin implement the least squares instead of the gradient descent for updating $\mathbf{U}$.

➢ We consider $c = 0.1$ for all our experiments.

---

**Algorithm 3:** AltGDmin

**Input:** observations $\mathbf{Y}$, mask $\mathbf{M}$, rank $r$, stepsize $\eta$, and iterations $I$.
**Output:** $\mathbf{U}^{(I)}, \mathbf{B}^{(I)}$

1 Initialize $\mathbf{U}^{(0)}$ by first $r$ left singular vectors of $\mathbf{Y}$
2 **for** $i \leftarrow 1$ *to* $I$ **do**
3   $\quad \mathbf{B}^{(i)}_{:,k} = (\mathbf{U}^{(i-1)^{\mathrm{T}}}_{\mathbf{m}_k,:} \mathbf{U}^{(i-1)}_{\mathbf{m}_k,:})^{-1} \mathbf{U}^{(i-1)^{\mathrm{T}}}_{\mathbf{m}_k,:} \mathbf{Y}_{\mathbf{m}_k,k} \ \forall \ k \in [q]$
4   $\quad \mathbf{U} \leftarrow \mathbf{U}^{(i-1)} - \eta * 2(\mathbf{M} \circ (\mathbf{U}^{(i-1)}\mathbf{B}^{(i)}) - \mathbf{Y})\mathbf{B}^{(i)^{\mathrm{T}}}$
5   $\quad \mathbf{U}^{(i)} \leftarrow \mathrm{QR}(\mathbf{U})$
6 **end**

---

where, step size $\eta = \frac{cp}{\|\mathbf{Y}\|_2^2}$ with constant $c$.

Abbasi, A. A., Moothedath, S., & Vaswani, N. (2023, September). Fast Federated Low Rank Matrix Completion. In *2023 59th Annual Allerton Conference on Communication, Control, and Computing (Allerton)* (pp. 1-6). IEEE.
Jain, P., Netrapalli, P., & Sanghavi, S. (2013, June). Low-rank matrix completion using alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing* (pp. 665-674).

➢ In the $t$-th horizon,

find $\mathbf{N}_t(s,a) = \sum_{k \in [K]} \mathbb{1}_{\{(s_t^k, a_t^k) = (s,a)\}} \ \forall (s,a) \in \mathbf{s} \times \mathbf{a}$,

estimate $\widehat{\mathbf{P}}_t(s'|s,a) = \frac{1}{\mathbf{N}_t(s,a)} \sum_{k \in K} \mathbb{1}_{\{(s_t^k, a_t^k, s_{t+1}^k) = (s,a,s')\}}$,

get $\mathbf{R}_t(s,a) = r(s,a) \ \forall (s,a) \in supp(\mathbf{N}_t)$ from dataset $\mathcal{D}$.

➢ Find, $\mathbf{Z}_t(s,a)$ using $\widehat{\mathbf{P}}_t$ and $\mathbf{R}_t$ $\forall (s,a) \in supp(\mathbf{N}_t)$ through the Bellman update equation, and

$\forall (s,a) \in \mathbf{s} \times \mathbf{a} \setminus supp(\mathbf{N}_t)$ estimate with AltGDmin.

➢ Iterate backward from horizon $H$ to $1$.

➢ Finally, find estimated total expected reward $\widehat{J}$.

➢ True total expected reward

$$J^{\mathbf{\Pi}^\theta} \leftarrow \sum_{s,a} \mathbf{D}_1^{\mathbf{\Pi}^\theta}(s,a) \mathbf{Q}_1^{\mathbf{\Pi}^\theta}(s,a).$$

$$\implies |\widehat{J} - J^{\mathbf{\Pi}^\theta}| = |\langle \mathbf{D}_1^{\mathbf{\Pi}^\theta}, \widehat{\mathbf{Q}}_1^{\mathbf{\Pi}^\theta} - \mathbf{Q}_1^{\mathbf{\Pi}^\theta} \rangle|$$

$$= |\sum_{t=1}^{H} \langle \mathbf{D}_t^{\mathbf{\Pi}^\theta}, \widehat{\mathbf{Q}}_t^{\mathbf{\Pi}^\theta} - \mathbf{Y}_t \rangle| \quad \text{where,} \quad \langle \mathbf{D}, \mathbf{Q} \rangle = \sum_{s,a} \mathbf{D}(s,a) \mathbf{Q}(s,a) \text{ and } \mathbf{Y}_t \text{ is population version of } \mathbf{Z}_t.$$

---

**Algorithm 4:** policy evaluation with AltGDmin

**Input:** dataset $\mathcal{D}$, target policy $\mathbf{\Pi}^\theta$, initial state distribution $\boldsymbol{\mu}_1$, horizon $H$, number of states $S$ and actions $A$, `AltGDmin`, rank $d$, stepsize $\eta$, iterations $I$, and $\mathbf{D}_1^{\mathbf{\Pi}^\theta}$.

**Output:** $\widehat{J}$

1 Initialize $\widehat{\mathbf{Q}}_{H+1}^{\mathbf{\Pi}^\theta} \leftarrow \mathbf{0}_{S \times A}$
2 Number of trajectories $K = \text{len}(\mathcal{D})$
3 **for** $t \leftarrow H$ **to** $1$ **do**
4 $\quad$ Initialize $\mathbf{N}_t \leftarrow \mathbf{0}_{S \times A}$, $\widehat{\mathbf{P}}_t \leftarrow \mathbf{0}_{S \times A \times S}$, and $\mathbf{R}_t \leftarrow \mathbf{0}_{S \times A}$
5 $\quad$ **for** $k \leftarrow 1$ **to** $K$ **do**
6 $\quad\quad$ $(s,a,r,s') \quad = \mathcal{D}(k,t)$
7 $\quad\quad$ $\mathbf{N}_t(s,a) \quad += 1$
8 $\quad\quad$ $\widehat{\mathbf{P}}_t(s,a,s') += 1$
9 $\quad\quad$ $\mathbf{R}_t(s,a) \quad = r$
10 $\quad$ **end**
11 $\quad$ $\widehat{\mathbf{P}}_t(s'|s,a) \leftarrow \widehat{\mathbf{P}}_t(s,a,s')/\mathbf{N}_t(s,a) \ \forall (s,a) \in supp(\mathbf{N}_t)$
12 $\quad$ $\mathbf{Z}_t(s,a) \quad = \mathbf{R}_t(s,a) +$
$\quad\quad \sum_{s',a'} \widehat{\mathbf{P}}_t(s'|s,a) \mathbf{\Pi}_{t+1}^\theta(a'|s') \widehat{\mathbf{Q}}_{t+1}^{\mathbf{\Pi}^\theta}(s',a') \ \forall (s,a) \in supp(\mathbf{N}_t)$
13 $\quad$ $\mathbf{U},\mathbf{B} = \texttt{AltGDmin}(\mathbf{Z}_t, supp(\mathbf{N}_t), d, \eta, T)$
14 $\quad$ $\widehat{\mathbf{Q}}_t^{\mathbf{\Pi}^\theta} \leftarrow \mathbf{U}\mathbf{B}$
15 **end**
16 $\widehat{J} \leftarrow \sum_{s,a} \mathbf{D}_1^{\mathbf{\Pi}^\theta}(s,a) \widehat{\mathbf{Q}}_1^{\mathbf{\Pi}^\theta}(s,a)$

Xi, X., Yu, C., & Chen, Y. Matrix Estimation for Offline Evaluation in Reinforcement Learning with Low-Rank Structure. In *3rd Offline RL Workshop: Offline RL as a"Launchpad"*.

> Observations:

- For a fixed rank, as the number of observations increases the subspace distance measure converge to $10^{-15}$ with fewer iterations.

- For a fixed number of observations, as the rank reduces, the subspace distance measure converge to $10^{-15}$ with fewer iterations.
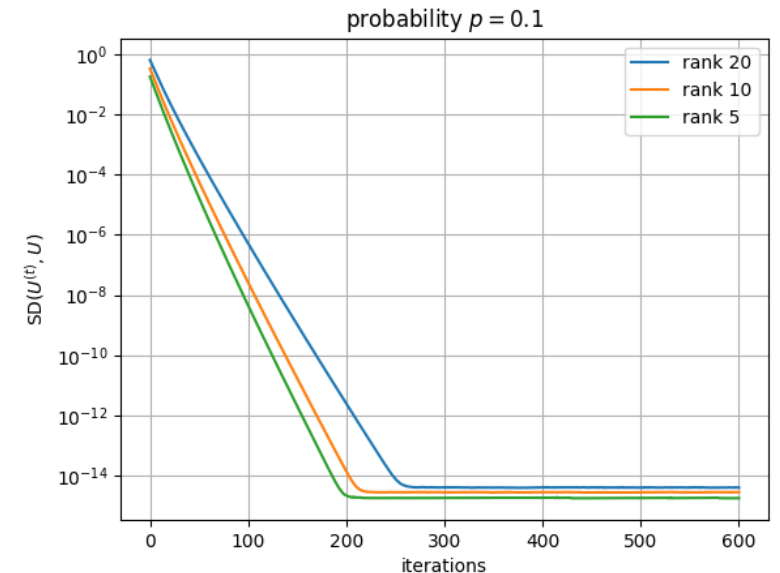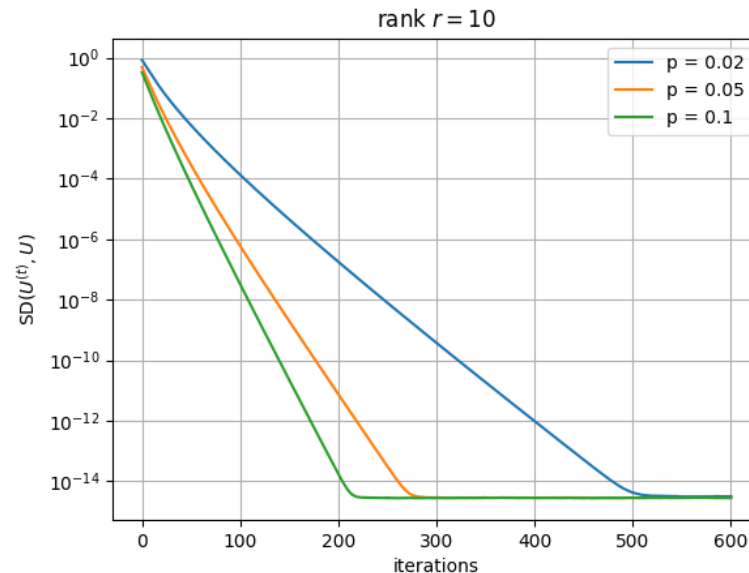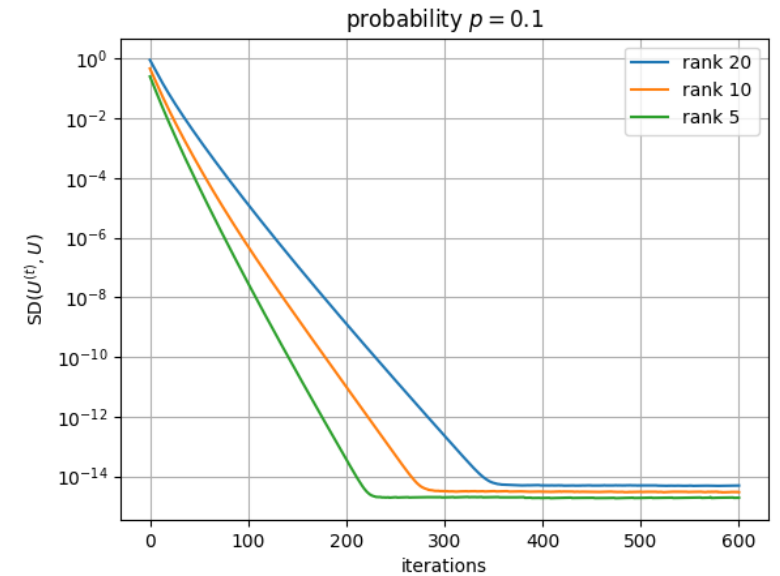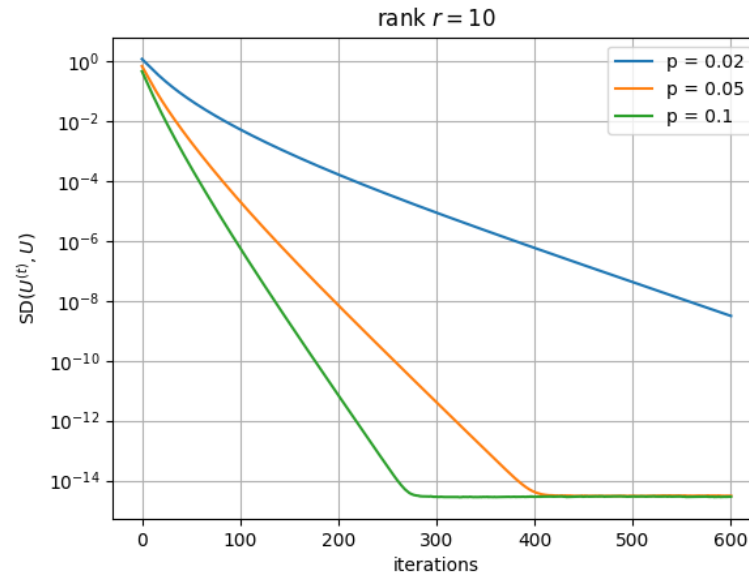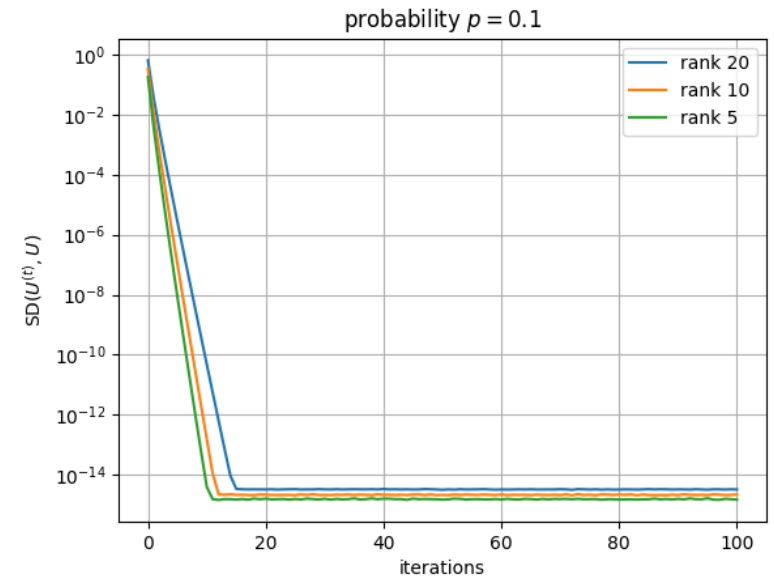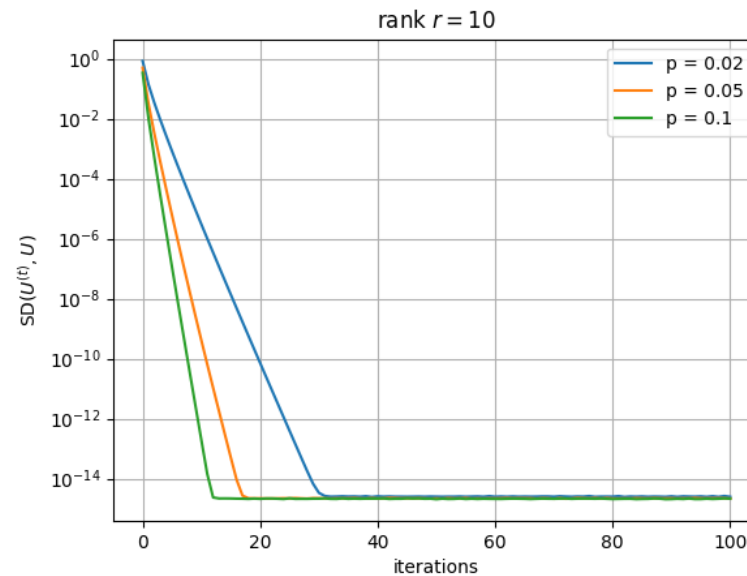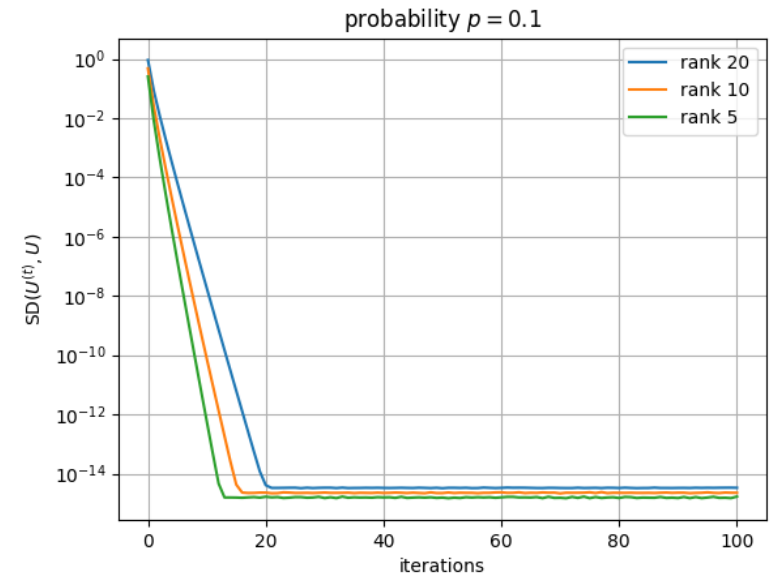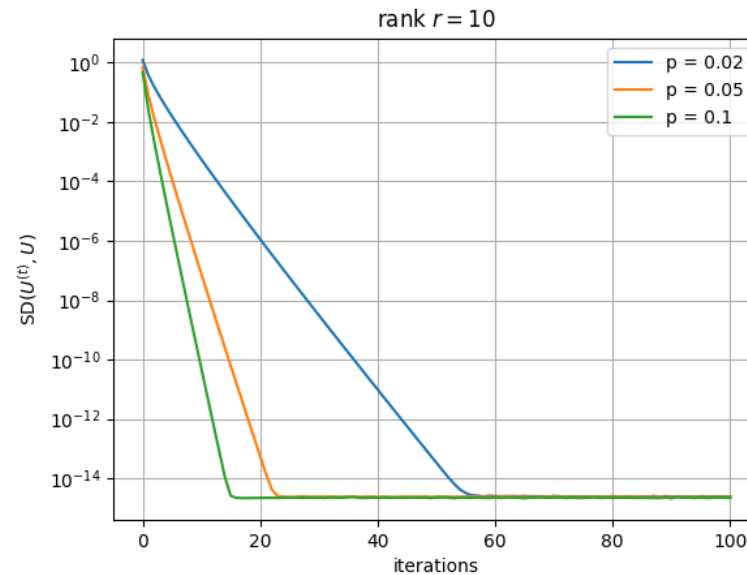

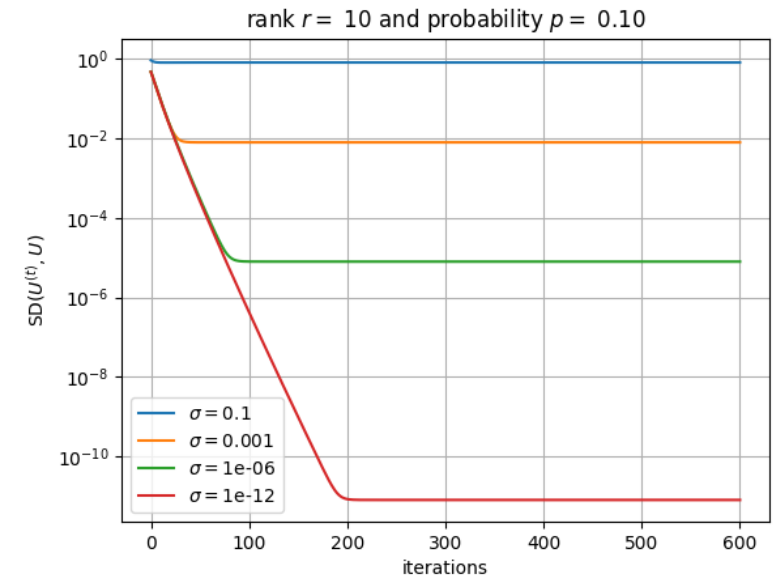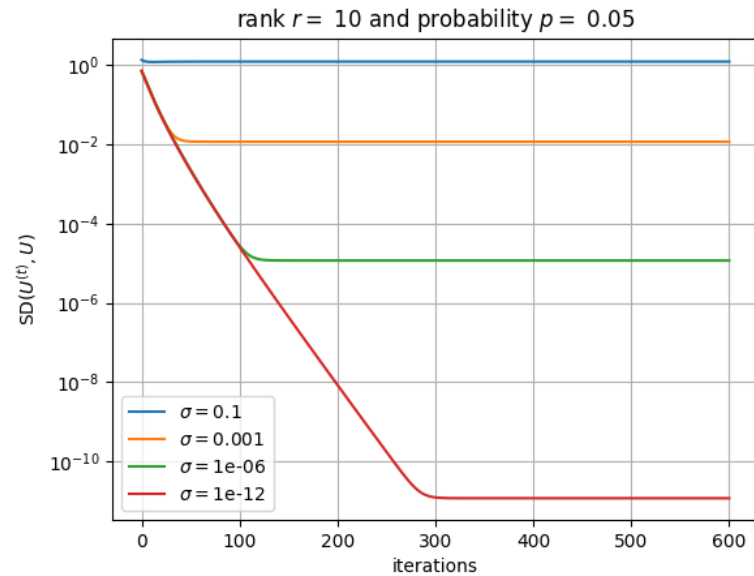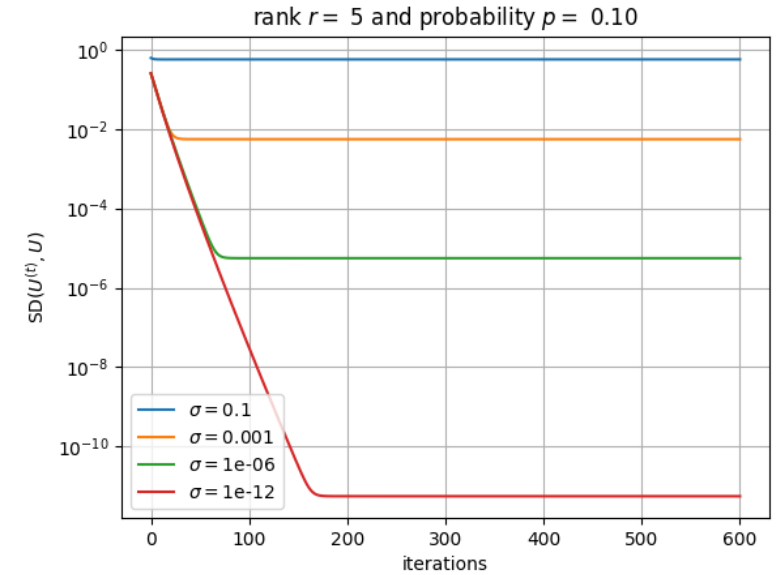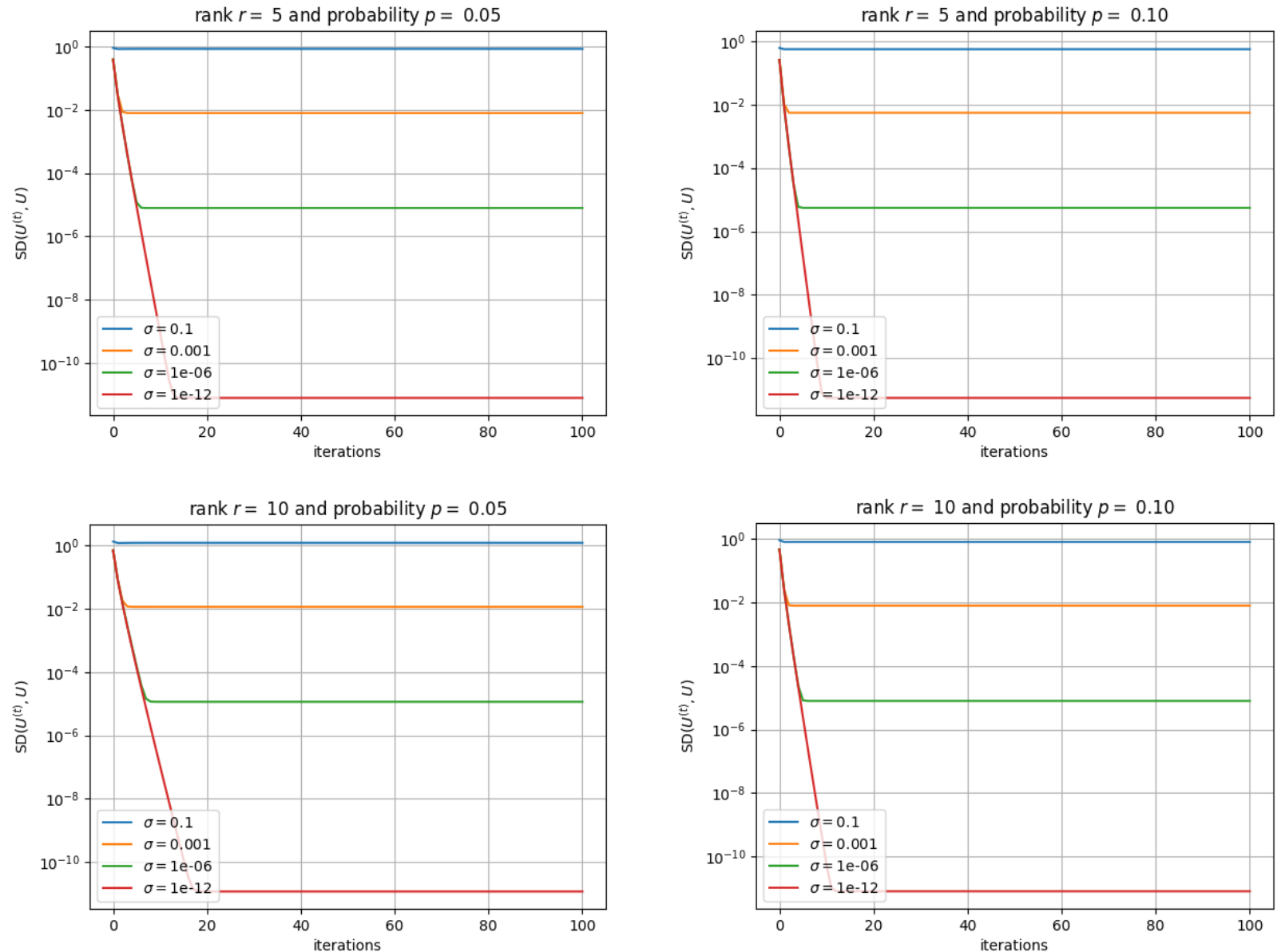
Fig.1: $(3000, 5000)$ and $(5000, 10000)$ low-rank matrix completion using AltGDmin

➢ Observations:

- For a fixed rank, as the number of observations increases the subspace distance measure converge to $10^{-15}$ with fewer iterations.

- For a fixed number of observations, as the rank reduces, the subspace distance measure converge to $10^{-15}$ with fewer iterations.



Fig.2: $(3000, 5000)$ and $(5000, 10000)$ low-rank matrix completion using AltGDmin

> ➢ Observations:

- For a fixed rank, as the number of observations increases the subspace distance measure converge to $10^{-11}$ with fewer iterations.

- For a fixed number of observations, as the rank increases, the subspace distance measure takes more iterations to converge to $10^{-11}$.



Fig.3: $(3000, 5000)$ noisy low-rank matrix completion using AltGDmin

# Simulation Results of Noisy Altmin
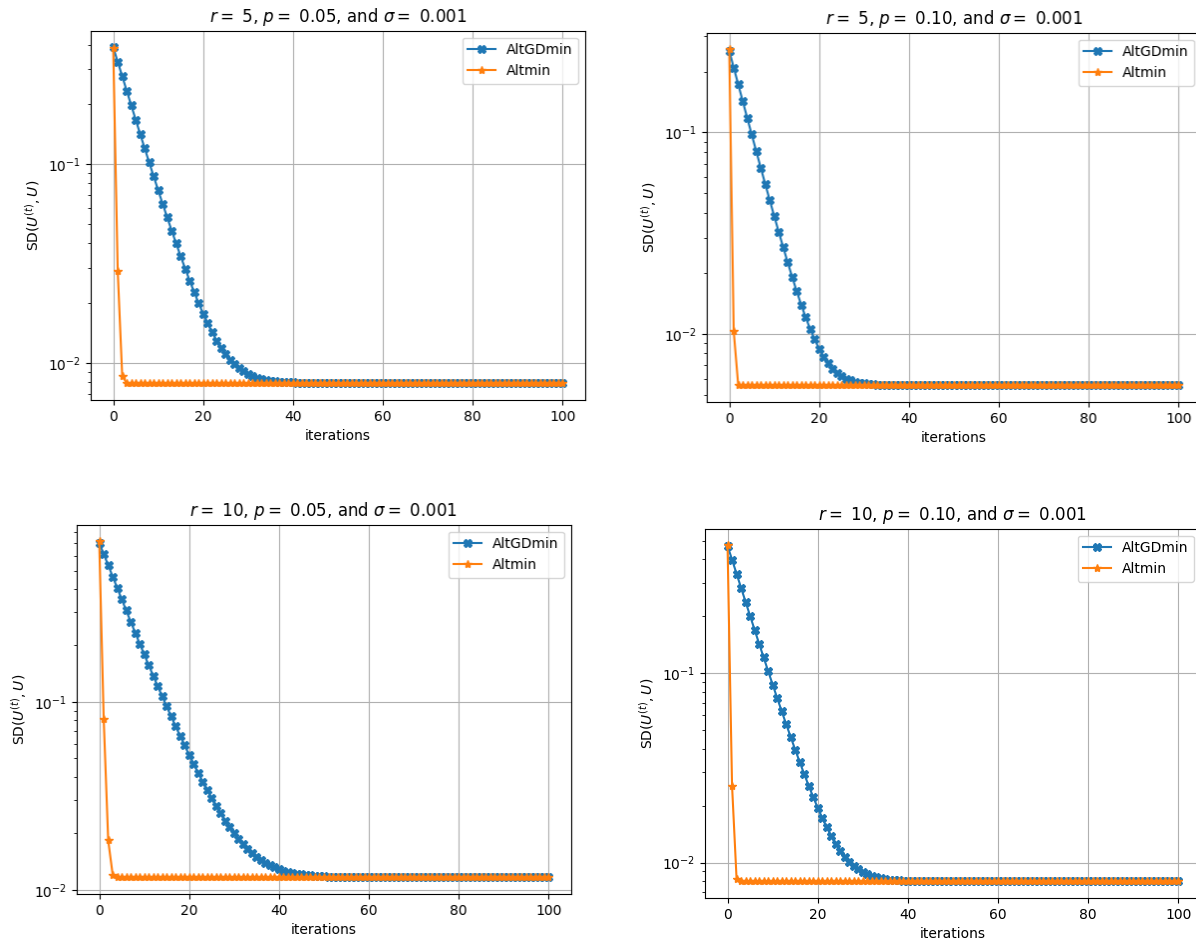
➢ Observations:

- For a fixed rank, as the number of observations increases the subspace distance measure converge to $10^{-11}$ with fewer iterations.

- For a fixed number of observations, as the rank increases, the subspace distance measure takes more iterations to converge to $10^{-11}$.



Fig.4: $(3000, 5000)$ noisy low-rank matrix completion using AltGDmin

# AltGDmin and Altmin Comparison



Fig.6(a): noise standard deviation $\sigma = 0.001$

Fig.6(b): noise standard deviation $\sigma = 1e-06$

Fig.6: Comparing AltGDmin and Altmin for $(3000, 5000)$ noisy LRMC

➢ Altmin takes fewer iterations to converge compared to AltGDmin.

➢ For action subset size $m = 300$ Algorithm 3 do not use matrix completion.

➢ Observations:

- For every value of action subset size $m$ the error increases with horizon $H$.

- For action subset size $m = 300$, as the trajectories $K$ increase from $5000$ to $10000$ the error reduces.

- For action subset size $m \geq n/2$, as the trajectories $K$ increase from $5000$ to $10000$ the error reduces for Altmin.

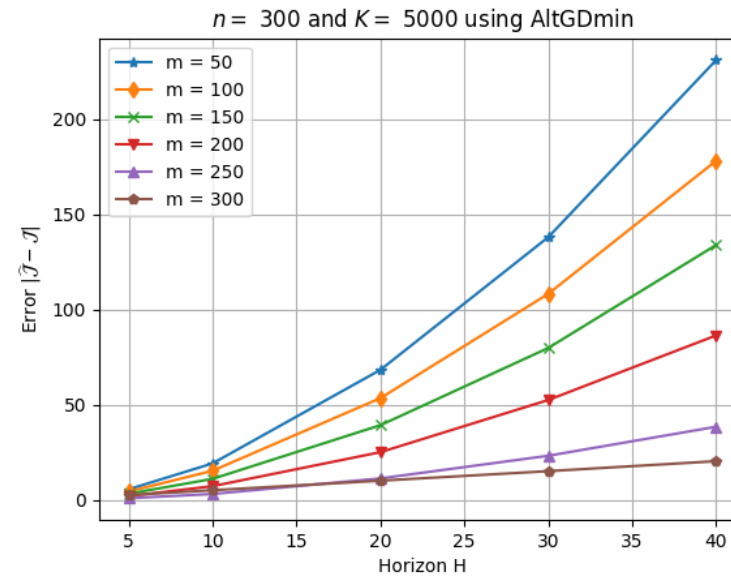- For action subset size $m = 250$, the error is almost increasing linearly with horizon $H$.



Fig.7. Policy Evaluation using AltGDmin and Altmin

➢ Observations:

• For action subset size $m = 300$, as the trajectories $K$ increase from $5000$ to $10000$ the error reduces significantly.

• For action subset size $m \geq n/2$, as the trajectories $K$ increase from $5000$ to $10000$ the error reduces for Altmin.
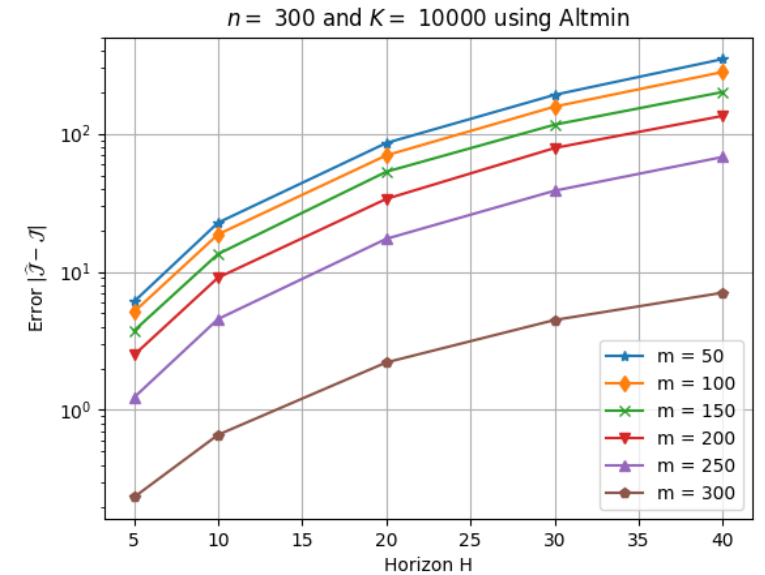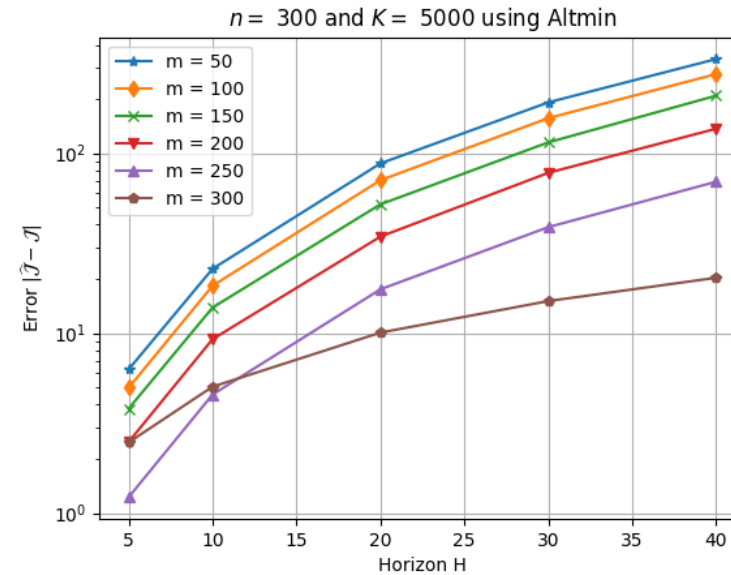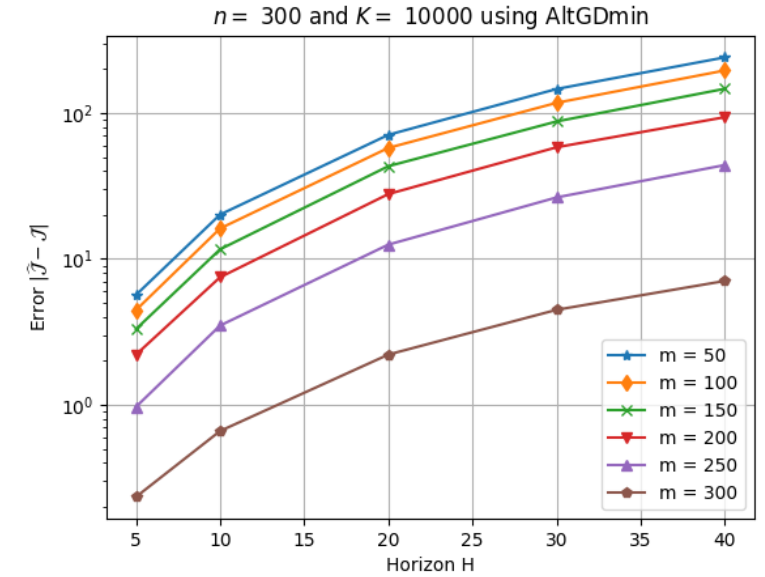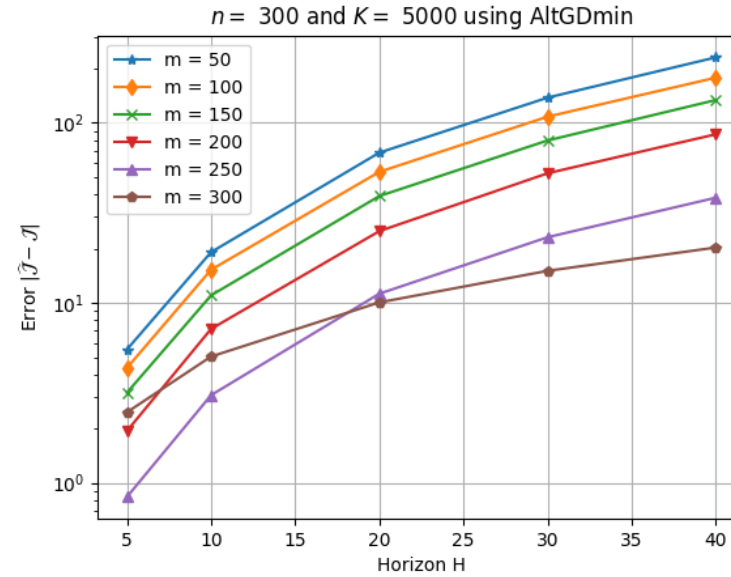


Fig.8. Policy Evaluation using AltGDmin and Altmin

# Comparison of AltGDmin and Altmin for Policy Evaluation
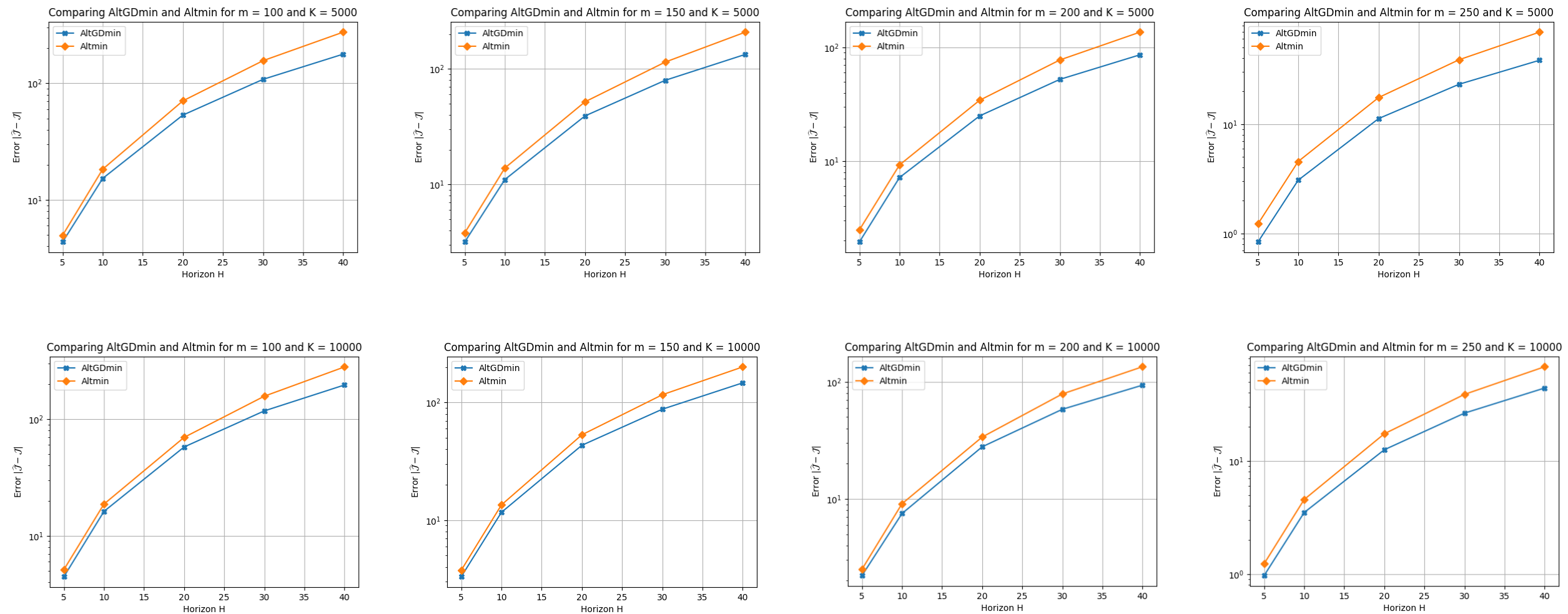


Fig.9. Comparing AltGDmin and Altmin for Policy Evaluation

➢ Observation: AltGDmin based policy evaluation outperform the Altmin one for every action subset size $m$ and trajectory $K$.

# Conclusion

➢ Practically feasible policy evaluation algorithm using alternating minimization for offline reinforcement learning.

➢ Alternating minimization based policy evaluation has a finite sample error bound.

➢ AltGDmin based policy evaluation outperform the Altmin one significantly.

# Future Work

➢ We will provide theoretical guarantees showing that AltGDmin is significantly better than Altmin.

# References

• Abbasi, A. A., Moothedath, S., & Vaswani, N. (2023, September). Fast Federated Low Rank Matrix Completion. In *2023 59th Annual Allerton Conference on Communication, Control, and Computing (Allerton)* (pp. 1-6). IEEE.

• Jain, P., Netrapalli, P., & Sanghavi, S. (2013, June). Low-rank matrix completion using alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing* (pp. 665-674).

• Levine, S., Kumar, A., Tucker, G., & Fu, J. (2020). Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643.*

- Liu, Y., Swaminathan, A., Agarwal, A., & Brunskill, E. (2020). Provably good batch off-policy reinforcement learning without great exploration. *Advances in neural information processing systems*, *33*, 1264-1274.

- Munos, R., & Szepesvári, C. (2008). Finite-Time Bounds for Fitted Value Iteration. *Journal of Machine Learning Research*, *9*(5).

- Nayer, S. S., & Vaswani, N. (2022, June). Fast low rank column-wise compressive sensing. In *2022 IEEE International Symposium on Information Theory (ISIT)* (pp. 3285-3290). IEEE.

- Rashidinejad, P., Zhu, B., Ma, C., Jiao, J., & Russell, S. (2021). Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information Processing Systems*, *34*, 11702-11716.

- Uehara, M., & Sun, W. (2021). Pessimistic model-based offline reinforcement learning under partial coverage. *arXiv preprint arXiv:2107.06226*.

- Xi, X., Yu, C., & Chen, Y. Matrix Estimation for Offline Evaluation in Reinforcement Learning with Low-Rank Structure. In *3rd Offline RL Workshop: Offline RL as a"Launchpad"*.

# Acknowledgements

➢ I am thankful to National Science Foundation (NSF) for supporting on this project.

➢ I am thankful to my advisor Prof. Namrata Vaswani for her consistent support and for teaching the course High Dimensional Probability and Linear Algebra.

➢ I am thankful to the faculty Dr. Shana Moothedath and Dr. Sung Yell Song for being a part of my program of study committee and for the helpful discussions on Reinforcement Learning and Applied Linear Algebra topics.

➢ I am thankful to the faculty Prof. Aditya Ramamoorthy for clarifying all my doubts on the topics of Convex Optimization.

➢ I am thankful to the faculty Prof. Pavan Kumar Aduri, Dr. Hongyang Gao, and Dr. Cheng Huang for teaching the courses Algorithms for Large Datasets, Introduction to Machine Learning, and Deep Learning for Theory and Practice.

➢ I am thankful to the faculty Jenna Bertilson, Ben Godard, and Febriana for teaching Academic Writing and Speaking.

➢ I am thankful to the Electrical and Computer Engineering department for the Teaching Assistant ship.

➢ I am thankful to Ruoyu, Silpa, Ankit, Jiabin, Ahmed, Sifat and Vrindha for helping and supporting.

➢ I am thankful to Dr. Kim, Vicky, Stacy, Shahin, Jason, Lynne, Sara, Tony and ECpE department staff for the administrative support.

➢ I am thankful to Post. Doc Kiran, Aditya Kar, Souradeep, Jaydeep, Nabila, and Aditya for their help and support.

**Thank You!**