Rewards over batches

- REINFORCE
- REINFORCE with baseline
- G(PO)MDP
- Optimal policy

Discounted reward vs Batch