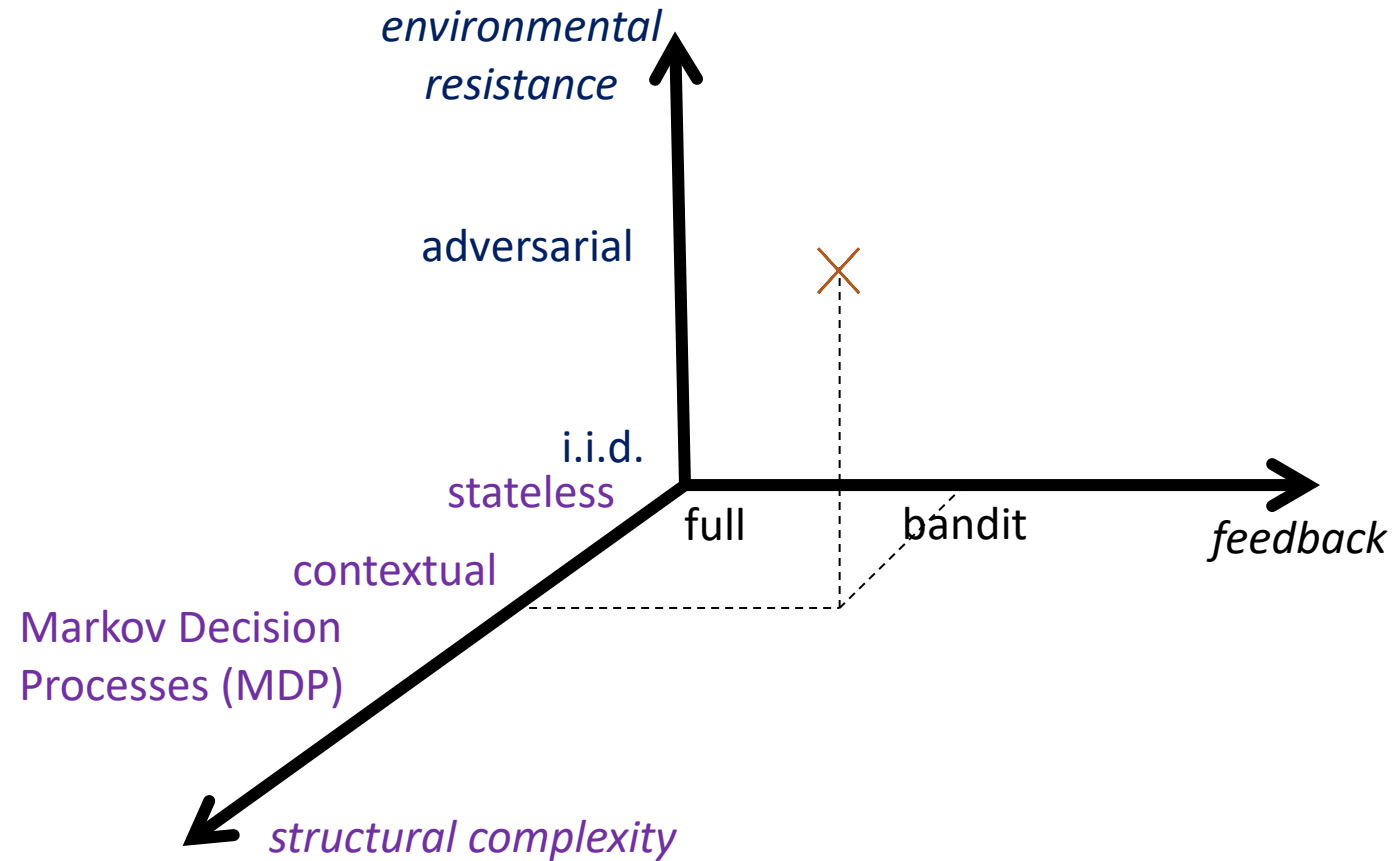


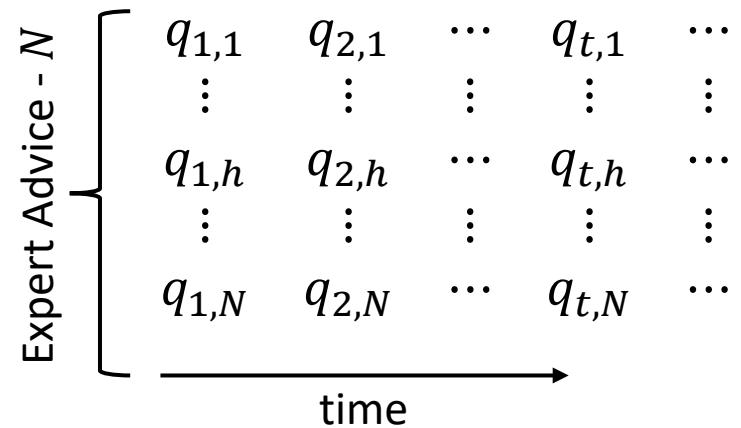
Contextual Bandits

Yevgeny Seldin

Contextual Bandits



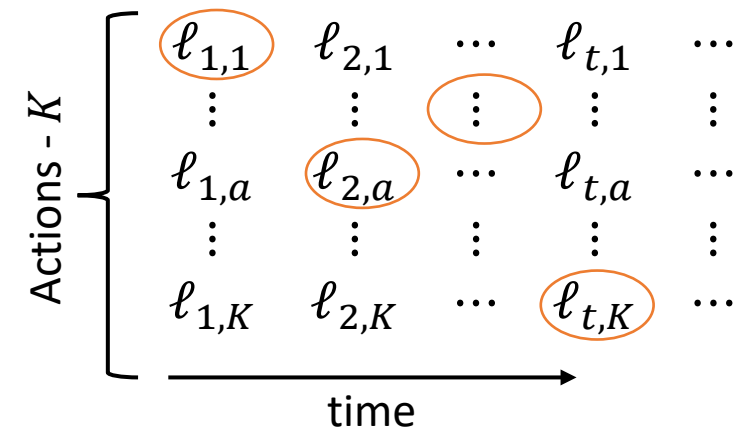
Version #1: Bandits with Expert Advice



Game definition:

- For $t = 1, 2, \dots$
 - Observe advice of N experts $q_{t,1}, \dots, q_{t,N}$
 - where $q_{t,h}$ is a distribution on actions $\{1, \dots, K\}$
 - Play an action A_t
 - Suffer and observe ℓ_{t,A_t}

Performance measure – regret:



$$R_T = \underbrace{\sum_{t=1}^T \ell_{t,A_t}}_{\text{Loss of the algorithm}} - \min_h \underbrace{\sum_{t=1}^T \sum_a q_{t,h}(a) \ell_{t,a}}_{\text{(Expected) loss of expert } h}$$

Deterministic $q_{t,h}$ models a path through loss matrix

Algorithm: EXP4

(Exponential Exploration Exploitation with Expert Advice)

- $\forall h: \tilde{L}_0(h) = 0$
- For $t = 1, 2, \dots$
 - $\forall a: p_t(a) = \sum_h \underbrace{q_{t,h}(a)}_{\text{Advice}} \underbrace{\frac{e^{-\eta_t \tilde{L}_{t-1}(h)}}{\sum_{h'} e^{-\eta_t \tilde{L}_{t-1}(h')}}}_{\text{Weight of expert } h}$
 - $A_t \sim p_t$
 - [Observe ℓ_{t,A_t}]
 - $\forall a: \tilde{\ell}_{t,a} = \frac{\ell_{t,a} \mathbb{1}(A_t=a)}{p_t(a)}$
 - $\forall h: \tilde{\ell}_{t,h} = \sum_a q_{t,h}(a) \tilde{\ell}_{t,a}$
 - $\forall h: \tilde{L}_t(h) = \tilde{L}_{t-1}(h) + \tilde{\ell}_{t,h}$

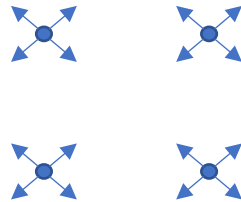
- EXP4 Expected regret upper bound:

$$\mathbb{E}[R_T] \leq \sqrt{2 \underbrace{KT}_{\substack{\text{Price of} \\ \text{bandit} \\ \text{feedback}}} \underbrace{\ln N}_{\substack{\text{Size of the} \\ \text{comparator} \\ \text{class}}}}$$

Version #2: Bandits with side information

Game definition:

- For $t = 1, 2, \dots$
 - Observe side info (state) S_t
 - Play an action A_t
 - Suffer and observe $\ell(A_t, S_t)$



Regret bounds:

- Running EXP4:

$$\mathbb{E}[R_T] = O\left(\sqrt{KT|\mathcal{S}| \ln K}\right)$$

- Lower bound in a nutshell:

- Generate $|\mathcal{S}|$ independent bandit problems
- Take $N = K^{|\mathcal{S}|}$ experts – all possible ways of assigning best actions to states
- Reminder: Lower bound for a single bandit is $\Omega(\sqrt{KT})$
- Each bandit is played $T/|\mathcal{S}|$ times, so its regret is $\Omega\left(\sqrt{K \frac{T}{|\mathcal{S}|}}\right)$

$$\text{Total regret } \mathbb{E}[R_T] = \Omega\left(\underbrace{|\mathcal{S}|}_{\#(\text{bandits})} \underbrace{\sqrt{K \frac{T}{|\mathcal{S}|}}}_{\text{regret of each bandit}}\right) = \Omega\left(\sqrt{KT|\mathcal{S}|}\right)$$

$$\text{Regret: } R_T = \underbrace{\sum_{t=1}^T \ell_t(A_t, S_t)}_{\text{Loss of the algorithm}} - \underbrace{\sum_{s \in \mathcal{S}} \min_a \sum_{t: S_t=s} \ell_t(a, s)}_{\substack{\text{Loss of the best action} \\ \text{in hindsight in state } s}} \\ \text{Total loss assuming the best} \\ \text{action in hindsight in each state}$$

Reduction (assuming expert advice is deterministic – each expert recommends just 1 action):

- #2→#1: Experts \rightarrow all possible mappings $h: \mathcal{S} \rightarrow \{1, \dots, K\}$
 - $N = K^{|\mathcal{S}|}$

- $|\mathcal{S}|$ - structural complexity

Summary – Contextual Bandits

- Bandits with Expert Advice

- EXP4 algorithm

- $p_t(a) = \sum_h \underbrace{q_{t,h}(a)}_{\text{Advice}} \underbrace{\frac{e^{-\eta_t \tilde{L}_{t-1}(h)}}{\sum_{h'} e^{-\eta_t \tilde{L}_{t-1}(h')}}}_{\text{Weight of expert } h}$

- $\mathbb{E}[R_T] \leq \sqrt{2KT \ln N}$

- Bandits with side information

- Reduction to prediction with expert advice

- $\mathbb{E}[R_T] = O\left(\sqrt{KT|\mathcal{S}| \ln K}\right)$

- $\mathbb{E}[R_T] = \Omega\left(\sqrt{KT|\mathcal{S}|}\right)$

