# Contextual Bandits.
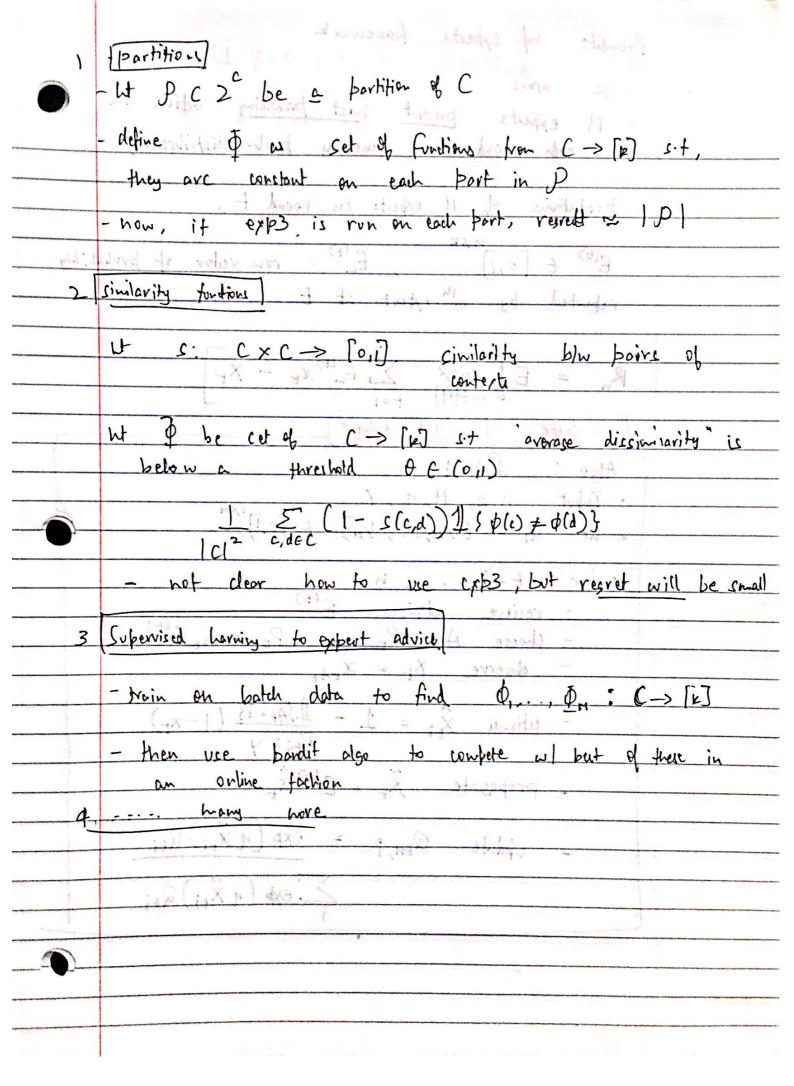
- in many Bandit problems, learner posceces extra/ side information to "predict" q/vality of actions.
- all algorithms + regret def$^n$ thus far ignore these contextual data.
- here we look at __better models__

- Eg: movie recommendation.

Interaction Protocol

- Adversary secretly chooses $(x_t)_{t=1}^n$, $x_t \in [0,1]^k$

- Adversary secretly chooses $(c_t)_{t=1}^n$, $\underline{c_t \in C}$
  <span style="color:red">arbitrary, fixed</span>

- for $t = 1, \ldots, n$

  learner observes $c_t$
  learner select $P_t \in P_{k-1}$, $A_t \sim P_t$

  learner observes $X_t = x_{t A_t}$

__Regret__: 
$$R_n = E\left[ \sum_{c \in C} \max_{i \in [k]} \sum_{t, c_t = c} (x_{ti} - X_t) \right]$$

$$R_{nc} = E\left[ \max_i \sum_{t, c_t = c} (x_{ti} - X_t) \right]$$

if exp3 is used for each context separately,

$$R_{nc} < 2 \sqrt{k \log k \sum_{t=1}^{n} \mathbb{1}\{c_t = c\}}$$

$$R_n < 2 \sum_{c \in C} \sqrt{k \log k \sum_{t} \mathbb{1}\{c_t = c\}}$$

- if $|C| = 1$, then same as adv. bandit

- if all $c \in C$ are equally likely

$$R_n < 2 \sqrt{nk|C| \log k}$$

## Bandits w/ expert advice

- if $|C|$ is large, then exp3 on each context not useful, unless $n$ is enormous.
- however, $C$ is "structured" in real-life
- ex: movie recommendation - users with similar demographics have similar preferences w high likelihood

Let $\Phi$ be set of all functions from $C \to [k]$

$$R_n = E\left[ \max_{\phi \in \Phi} \sum_{t=1}^{n} X_{t\phi(c_t)} - X_t \right]$$

- if $\Phi$ is small, we can get better reward.

1. **Partitions**

- let $P \subseteq 2^C$ be a partition of $C$

- define $\Phi$ as set of functions from $C \to [k]$ s.t. they are constant on each part in $P$

- now, if exp3 is run on each part, regret $\approx |P|$

2. **similarity functions**

let $s: C \times C \to [0,1]$ similarity b/w pairs of contexts

let $\Phi$ be set of $C \to [k]$ s.t 'average dissimilarity' is below a threshold $\theta \in (0,1)$

$$\frac{1}{|C|^2} \sum_{c,d \in C} (1 - s(c,d)) \mathbb{1}\{\phi(c) \neq \phi(d)\}$$

- not clear how to use exp3, but regret will be small

3. **Supervised learning to expert advice**

- train on batch data to find $\phi_1, \ldots, \phi_M : C \to [k]$

- then use bandit algo to compete w/ but of these in an online fashion

4. ...... many more

Bandits w/ experts framework

- $k$ arms
- $M$ experts predict <u>most promising</u> action in each round. (generally prob. distributions)

- predictions of $M$ experts in round $t$,

$$E^{(t)} \in [0,1]^{M \times R}, \quad E_m^{(t)} - \text{row vector of probability}$$
reported by $m^{th}$ expert at $t$

$$R_n = E\left[\max_{m \in [M]} \sum_{t=1}^{n} E_m^{(t)} x_t - X_t\right]$$

[ compete w/ <u>best expert</u> ]

Also: Exp 4
- Input $n, k, M, \eta, Y$
- let $Q_1 = (1/N, \ldots, 1/N) \in [0,1]^{1 \times M}$

- for $t = 1, \ldots, n$
  - recieve advice $E^{(t)}$
  - choose $A_t \sim P_t$, $\quad P_t = Q_t E^{(t)}$
  - observe $X_t = x_{t A_t}$

  - estimate $\hat{X}_{ti} = 1 - \dfrac{\mathbb{1}\{A_t = i\}}{P_{ti} + Y}(1 - x_t)$

  - propagate $\tilde{X}_t = E^{(t)} \hat{X}_t$

  - update $Q_{t+1,i} = \dfrac{\exp(\eta \tilde{X}_{ti}) Q_{ti}}{\sum_{j} \exp(\eta \tilde{X}_{tj}) Q_{tj}}$

**Thm:** Let $Y = 0$, $\eta = \sqrt{2(\log M)/nk}$. Then

$$R_n \leq \sqrt{2nk \log M}$$