

LEC019 MAB I

Multi-armed bandit prob.



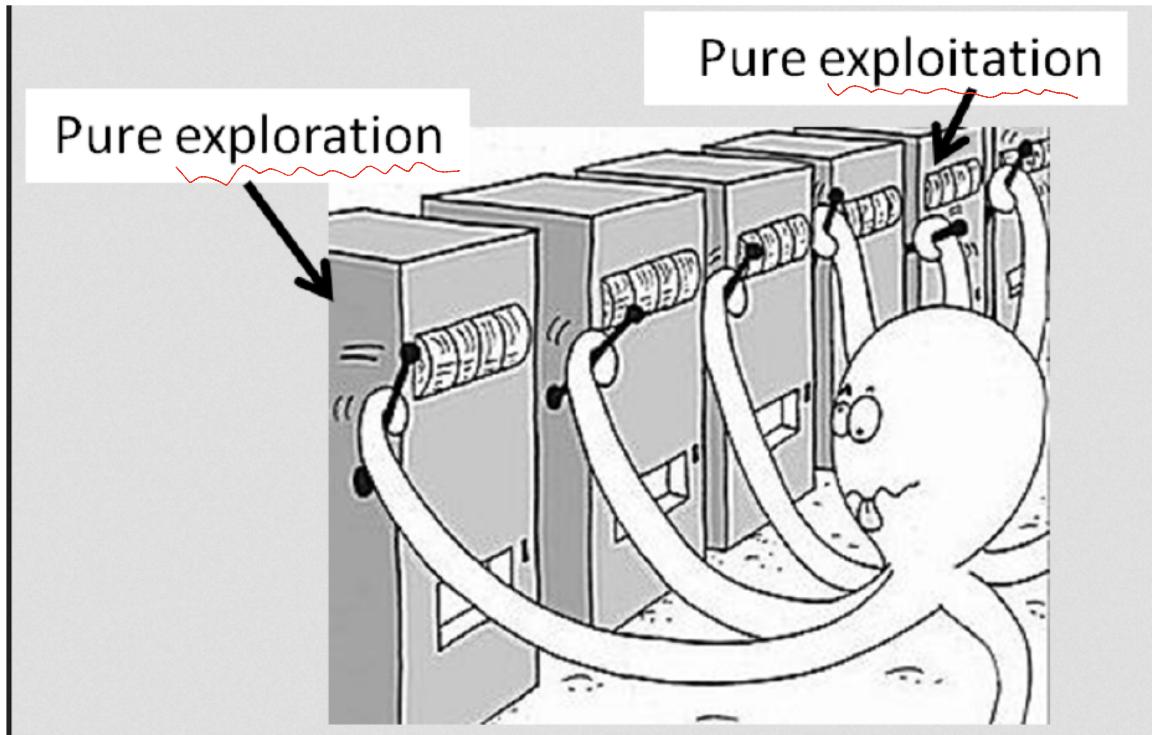
Basic form
of RL.

VG441 SS2020

Cong Shi
Industrial & Operations Engineering
University of Michigan

Basic RL for combining learning and decisions

- Multi-Armed Bandit Problem



MDP

MAB

- Different machine generates different random rewards
- Gambler decides which slot machine to play with each token
- Maximize reward (\$\$)



違った機械
の賞金が違う

Online decision-making: learning while doing

- Online decision-making involves a fundamental choice:
 - Exploration: Gather more information
 - Exploitation: Make the best decision given current information

Trade
- off



- The best long-term strategy may involve short-term sacrifices

Example: Insufficient Exploration

1	2	3	4	5	6	7	8	
---	---	---	---	---	---	---	---	--



Example: Insufficient Exploration

1	2	3	4	5	6	7	8	
---	---	---	---	---	---	---	---	--



	\$0		\$0	\$0				
--	-----	--	-----	-----	--	--	--	--



\$5		\$5			\$5	\$5	\$5	...
-----	--	-----	--	--	-----	-----	-----	-----

Example: Insufficient Exploration

1	2	3	4	5	6	7	8	
---	---	---	---	---	---	---	---	--



	\$0		\$0	\$0				
--	-----	--	-----	-----	--	--	--	--



\$5	\$5	\$5	\$5	\$5	\$5	\$5	\$5	...
-----	-----	-----	-----	-----	-----	-----	-----	-----

It turns out  always pays \$5/round

Example: Insufficient Exploration

1	2	3	4	5	6	7	8	
	\$0		\$0	\$0				
	\$5	\$5	\$5	\$5	\$5	\$5	\$5	...

It turns out  always pays **\$5**/round

事实上，

 pays **\$100** a quarter of the time
(\$25/round on average)

Example: Insufficient Exploration



1	2	3	4	5	6	7	8	
\$100	\$0	\$0	\$0	\$0	\$100	\$0	\$100	



\$5	\$5	\$5	\$5	\$5	\$5	\$5	\$5	...
-----	-----	-----	-----	-----	-----	-----	-----	-----

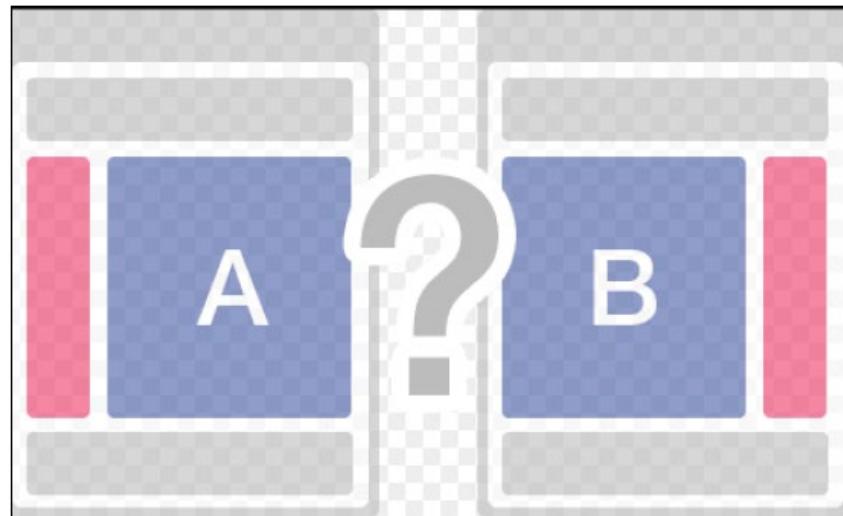
It turns out  always pays **\$5**/round

 pays **\$100** a quarter of the time
(\$25/round on average)

A/B Testing

多項，則 A/B/C/D/E Testing

- Exploration: Gather more information about which design is better
- Exploration: Show the best design to the customer



Revenue Management

= pricing decision

- Retailers are interested in finding an optimal (pricing) policy to max revenue
- Unknown relationship between price and customer's purchasing decision (demand distribution)
 - Exploration: Gather more information about customers behavior using different prices
 - Exploitation: Make the best price based on the current information

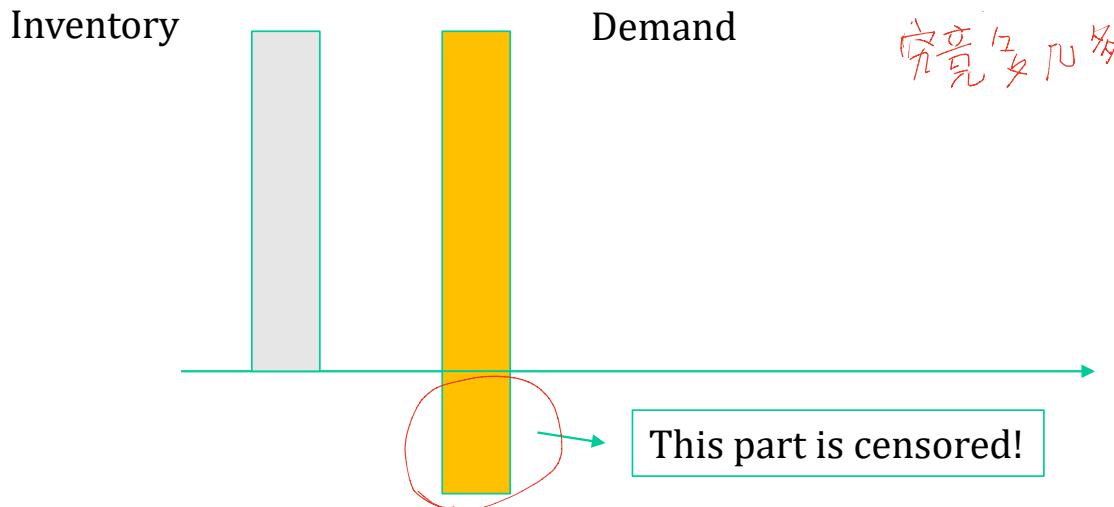
D(P) unknown

demand
in terms of
price.



Inventory Management

- Retailers are interested in finding an optimal (ordering) policy to min cost
- Unknown demand distribution (can only observe sales – censored demand)
 - Exploration: Order more to find out about true demand distribution
 - Exploitation: Order just right to minimize the cost



Other Applications

exploration vs exploitation



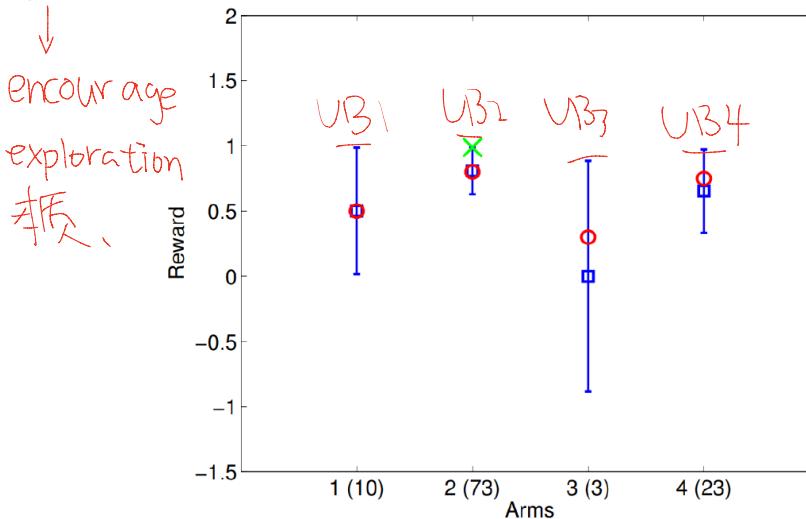
- Clinical trials 选 drug
- Recommender systems 推荐系统 eg. Netflix, 抖音看一天的
算法工程师.
- Advertising: what ad to put on a web-page?
- Auctions
- Financial portfolio design
- Crowdsourcing

in practice
 $\epsilon = 0.05$

Many algorithms for MAB

- ϵ -greedy algorithm 每次 {select the current best
some random} $1 - \epsilon$
- Upper confidence bound (UCB)
 - Add confidence bonus to the estimated mean
 - If the estimator is reliable, add less; if not, add more

Select the arm with highest UB



optimal } UCB empirical avg
 TS frequentist
 Bayes theorem
sub-optimal : ϵ -greedy Bayes theorem

$$i_t = \arg \max_{\text{Empirical mean}} \left[\hat{\mu}_i + \sqrt{\frac{c \log t}{n_i}} \right]$$

□ blue = empirical mean
○ red = true mean

- Thompson sampling (TS)
 - Bayesian setup with a prior distribution over reward parameters
 - Choose the auction that maximizes the expected reward under posterior

posterior \propto prior \times likelihood.

例 1

(large
if
under-
explored)

Online Network RM using TS

REVENUE management.

- ~\$300B industry with ~10% annual growth over the last 5 years
 - IBISWorldUS Industry Report; excludes online sales of traditionally brick & mortar stores



- Online retailers have additional information as compared to brick & mortar retailers, e.g. real-time customer purchase decisions (buy / no buy)
 - How can we use this information to develop a more effective revenue management strategy?

Setting

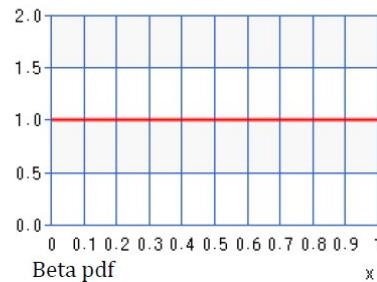
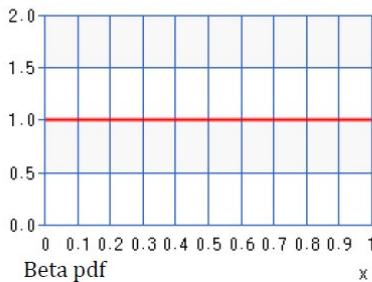
- Finite selling horizon of T periods
 - One customer arrives per period
 - Sequentially observe customer purchase decisions
- Finite set of prices; i -th price denoted by $\underline{p_i}$
 - mean demand
 - $\text{at price } \underline{p_i}$
- Unknown mean demand per price (“purchase probability”) $\underline{d_i}$
- Given unlimited inventory and known demand, select price with highest revenue = $\underline{p_i} \times \underline{d_i}$
 - expected
 - $(\text{之後會有 finite inventory})$
- Challenges: unknown demand
- Exploration vs. Exploitation Tradeoff

- Retailer decides...
 - Which price to offer to a customer
 - How many times to offer each price
 - In what order to offer prices to customers
- Learns demand at each price to max revenue



RM-MAB

→ ~~最简单的设置~~ Simple bb setting
 Customer behavior Bernoulli { \u25b2 buy \u25bc not buy }



$\hat{d}_1 \sim \text{Beta}(1,1)$ $\hat{d}_2 \sim \text{Beta}(1,1)$
 True (unknown) $d_1 = 0.6$ True (unknown) $d_2 = 0.3$

$\text{Beta}(1,1)$
 = uniform[0,1]

- Customer arrives
- Retailer samples θ_1 and θ_2 from current distributional estimation of d_1 and d_2
- Retailer offers price that maximizes $p_i \theta_i$
- Customer makes purchase decision (according to d_i)
- Retailer observes purchase decision and updates demand estimation

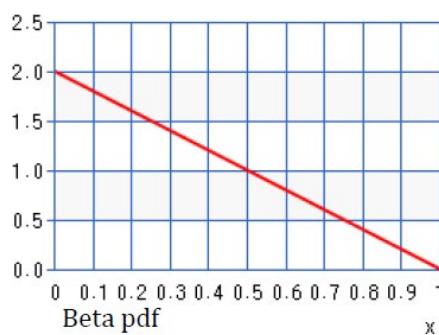
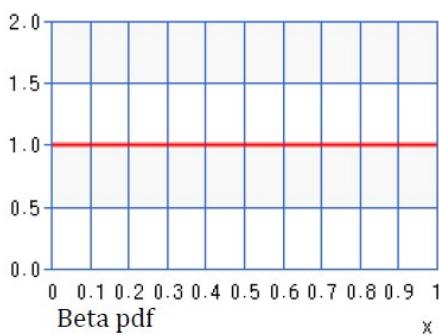
RM-MAB

Posterior = Prior \times likelihood
 Beta Beta Bern...

$$\rightarrow E[\text{Beta}(\alpha, \beta)] = \frac{\alpha}{\alpha + \beta}, \quad \text{Beta is conjugate prior for Bernoulli likelihood function}$$



pdf



$$\theta_1 = 0.41, \theta_2 = 0.83$$

$$p_2 \theta_2 > p_1 \theta_1$$

$$39.9 \times 0.83 > 29.9 \times 0.41$$



$\hat{d}_1 \sim \text{Beta}(1,1)$
 True (unknown) $d_1 = 0.6$

$\hat{d}_2 \sim \text{Beta}(1,1)$
 True (unknown) $d_2 = 0.3$

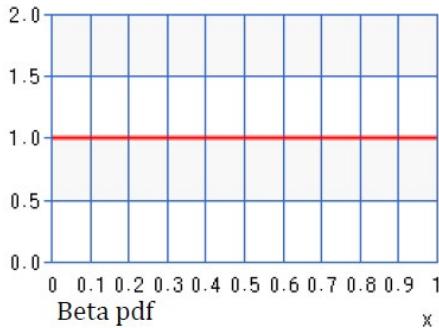
update
 $\hat{d}_2 \sim \text{Beta}(1, 1 + 1)$
 True (unknown) $d_2 = 0.3$

Customer does not buy item 2

good ex. on "Bayes rule"

e.g. if buy
 $(+, +) \rightarrow (+, -)$

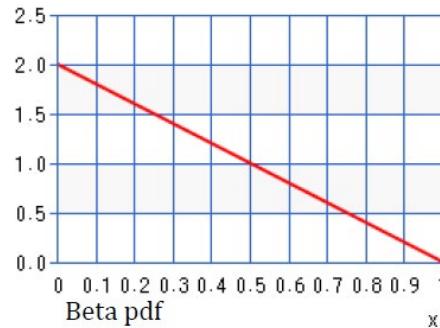
RM-MAB



$\hat{d}_1 \sim \text{Beta}(1,1)$
True (unknown) $d_1 = 0.6$



$\hat{d}_1 \sim \text{Beta}(1+1,1)$
True (unknown) $d_1 = 0.6$



$\hat{d}_2 \sim \text{Beta}(1,2)$
True (unknown) $d_2 = 0.3$

update

$$\theta_1 = 0.93, \theta_2 = 0.12$$

$$29.9 \times 0.93 > 39.9 \times 0.12$$



Customer buys item 1

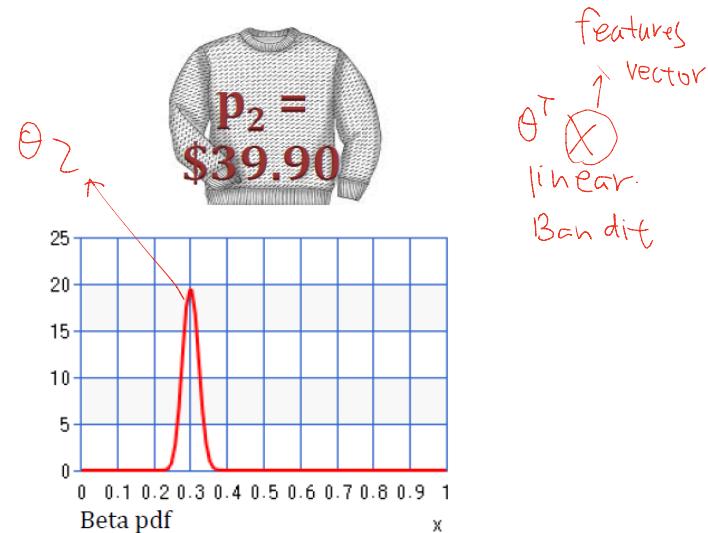
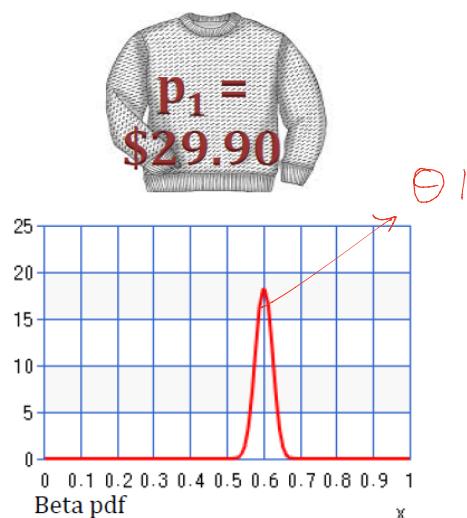
\rightarrow Bernoulli

RM-MAB: 2 Price Example

Amazon 用 IS 来 pricing strategy. (但现实中顾客有 features)

As each price is offered more times...

- Beta pdf converges to reflect true mean demand
- Will choose optimal price with high probability



$\hat{d}_1 \sim \text{Beta}(1 + \# \text{ "buy"}, 1 + \# \text{ "no buy"})$
True (unknown) $d_1 = 0.6$

$\hat{d}_2 \sim \text{Beta}(1 + \# \text{ "buy"}, 1 + \# \text{ "no buy"})$
True (unknown) $d_2 = 0.3$

Advantages of Thompson Sampling

- Empirical and theoretical results show it's a highly competitive algorithm for unlimited inventory
 - Easy to implement and understand
 - Non-parametric
 - Continuous exploration & exploitation
"Learning while doing"
- online (linear time)
near (linear)

How do we incorporate inventory constraints?

Key Tradeoffs: ~~時間很大~~

- Exploration vs. Exploitation
- Explore at the cost of running out of inventory

{ concentration
anti ~

RM-with inventory constraint

1. Customer arrives
2. Retailer samples θ_1 and θ_2
3. Retailer solves a deterministic LP to identify the optimal fraction of remaining customers to offer p_1 and p_2 , using
 - θ_1 and θ_2
 - Remaining unsold inventory & customers
4. Retailer offers price p_i with probability based on fraction found in Step 3
5. Customer makes purchase decision
6. Retailer observes decision and updates \hat{d}_i

(linear program
//)

RM-with inventory constraint

solve LP,
at time t randomized policy.

$\underline{x_i}$ = fraction of remaining customers $(T-t)$ to offer price $\underline{p_i}$
 (剩余)

$$\max_{x_1, x_2} \sum_{T-t} p_1 \theta_1 x_1 + p_2 \theta_2 x_2$$

$$s.t. \underline{x_1 + x_2 \leq 1}$$

$$(T-t)(\theta_1 x_1 + \theta_2 x_2) \leq Inv(t)$$

remaining time horizon \times expected demand

$$x_1, x_2 \geq 0$$

maximize revenue over remaining customers

fraction of remaining customers ≤ 1

expected inventory sold is upper-bounded by remaining inventory

\leq remaining inventory

"Bandits with Knapsack".
 $Ax \leq b$ 之类 constrain.

现实中可以加入
threshold 临界值

RM-with inventory constraint

希望越小越好

clairvoyant

$$\text{Regret} = \underbrace{\mathbb{E}[\text{Revenue of Optimal Policy with Known Demand}]}_{\text{希望越小越好}} - \mathbb{E}[\text{Revenue of Algorithm}]$$

$$\leq \mathbb{E}[\text{Upper Bound on Optimal Policy} - \mathbb{E}[\text{Revenue of Algorithm}]]$$

$$f(T) = O(T)$$

Theorem

Suppose the LP of the underlying true demand (i.e. benchmark) is nondegenerate. Then, for the modified Thompson Sampling with Inventory Algorithm,

$$\text{Regret}(T) \leq O(\sqrt{T} \log T \log \log T) = \tilde{O}(\sqrt{T})$$

主要用 T
 $\log \log (T)$

不希望 $\text{Regret}(T) = O(T)$ bad (随机都不能 $O(T)$)

$O(T^{2/3}), O(\sqrt{T}), O(\log T)$ 就很好

\tilde{O} → hide \log terms.

這是 lower bound (optimal)

regret (T) $\geq \Omega(\sqrt{T})$ for any algorithms.

以上內容是 phd/graduate level, 不會考了.