# Online Supplement for

# Online Contextual Learning with Perishable Resources Allocation

Xin Pan, Jie Song, Jingtong Zhao, Van-Anh Truong

## A Proof for Lemma 1

*Proof.* Fix any realization of demand represented as a sequence  $\mathbb{D} = \{(i_d, t_d), D\}$  (d = 1, 2, ..., D)where  $i_d$  and  $t_d$  are the customer type and arrival time of the *d*-th arrival, respectively, and D is a random variable representing the total number of arrivals. Let  $p_{i_d j k}(t_d) \in [0, 1]$  be the assignment probability of resource (j, k) to the *d*-th arrival. Note that the optimal offline decisions must satisfy

$$\sum_{d=1}^{D} p_{idjk}(t_d) \le c_j \quad j = 1, 2, ..., J; k = 1, 2, ..., K,$$
$$\sum_{j=1}^{J} \sum_{k=1}^{K} p_{ijk}(t_d) \le \mathbb{1}_{i_d=i}, d = 1, 2, ..., D, i = 1, 2, ..., I$$
$$p_{i_djk}(t_d) = 0, \ \forall t_d = [0, s_{jk} - W) \cup (s_{jk}, T].$$

Since the inequalities hold for any realization  $\{(i_d, t_d)\}$  (d = 1, 2, ..., D), then it is obvious that the constraints will also be true after taking expectations on both sides. Therefore, the expected optimal offline solution is feasible in problem (1), so that  $V^{OFF} \ge V^*$ .

# **B** Proofs for Section 4

#### B.1 Proof for Lemma 2

*Proof.* For the first part, because  $\mathbf{r}_{\tau} = \mathbf{Z}_{\tau}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , we have

$$\begin{split} \hat{\boldsymbol{\beta}} &= (\mathbf{Z}_{\tau}^{\mathrm{T}} \mathbf{Z}_{\tau})^{-1} \mathbf{Z}_{\tau}^{\mathrm{T}} \mathbf{r}_{\tau} \\ &= (\mathbf{Z}_{\tau}^{\top} \mathbf{Z}_{\tau})^{-1} \mathbf{Z}_{\tau}^{\top} (\mathbf{Z}_{\tau} \boldsymbol{\beta} + \boldsymbol{\epsilon}) \\ &= \boldsymbol{\beta} + (\mathbf{Z}_{\tau}^{\top} \mathbf{Z}_{\tau})^{-1} \mathbf{Z}_{\tau}^{\top} \boldsymbol{\epsilon}. \end{split}$$

Hence  $E[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ , and

$$Var(\hat{\boldsymbol{\beta}}) = E[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^{\top}]$$
  
=  $E[((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}\mathbf{Z}_{\tau}^{\top}\boldsymbol{\epsilon})((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}\mathbf{Z}_{\tau}^{\top}\boldsymbol{\epsilon})^{\top}]$   
=  $(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}\mathbf{Z}_{\tau}^{\top}E[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^{\top}]\mathbf{Z}_{\tau}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}$   
=  $\sigma^{2}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}.$ 

The second part follows immediately from the first part.

### B.2 Proof for Lemma 3

*Proof.* According to Lemma 2,  $\mathbf{z}\hat{\boldsymbol{\beta}}$  follows the normal distribution, and the variance is known. So we can construct a confidence interval around  $\mathbf{z}\hat{\boldsymbol{\beta}}$ 

$$\mathbf{z}\hat{\boldsymbol{\beta}} \pm \mathbf{z}_{\frac{\alpha}{2}} \frac{\sigma\sqrt{\mathbf{z}^{\top}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}\mathbf{z})}}{\sqrt{\tau}}$$
(1)

with confidence level  $1 - \alpha$ , where  $\mathbf{z}_{\frac{\alpha}{2}}$  is the standard normal distribution coefficient.

Next we show how the width of confidence interval changes with the scaling parameter. As a property of eigenvalues, we have

$$\mathbf{z}^{\top} (\mathbf{Z}_{\tau}^{\top} \mathbf{Z}_{\tau})^{-1} \mathbf{z} \le \lambda_{max} \mathbf{z}^{\top} \mathbf{z} \le tr((\mathbf{Z}_{\tau}^{\top} \mathbf{Z}_{\tau})^{-1}) \mathbf{z}^{\top} \mathbf{z},$$
(2)

where  $\lambda_{max}$  is the maximum eigenvalue of  $(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}$ . By the Sherman-Morrison formula, after

adding one more sample  $\mathbf{z}'$  into  $\mathbf{Z}_{\tau}$ , the trace of the inverse matrix  $(\mathbf{Z}^{\top}\mathbf{Z})_{\tau}^{-1}$  becomes

$$tr((\mathbf{Z}_{\tau+1}^{\top}\mathbf{Z}_{\tau+1})^{-1}) = tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}) - \frac{tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}\mathbf{z}'\mathbf{z}'^{\top}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1})}{1 + \mathbf{z}'^{\top}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}\mathbf{z}'}$$

$$\leq tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}) - \frac{tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}\mathbf{z}'\mathbf{z}'^{\top}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1})}{1 + tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1})\mathbf{z}'^{\top}\mathbf{z}'}$$

$$\leq tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}) - \frac{tr(\mathbf{z}'\mathbf{z}'^{\top})}{tr^{2}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})(1 + tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1})\mathbf{z}'^{\top}\mathbf{z}')},$$
(3)

where the first inequality follows (2), and the second inequality follows from  $tr(AB) \leq tr(A)tr(B)$ so that  $tr(\mathbf{z}'\mathbf{z}'^{\top}) = tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}\mathbf{z}'\mathbf{z}'^{\top}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})) \leq tr(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1}\mathbf{z}'\mathbf{z}'^{\top}(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1})$  $tr(\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})$ . Suppose  $a \leq \mathbf{z}'^{\top}\mathbf{z}' \leq b$  for all possible  $\mathbf{z}'$ , then it is obvious that  $\tau a \leq tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})) \leq$  $\tau b$ . Let  $x_{\tau} = tr((\mathbf{Z}_{\tau}^{\top}\mathbf{Z}_{\tau})^{-1})$ , then (3) becomes

$$x_{\tau} \le x_{\tau-1} - \frac{a}{(\tau-1)^2 b^2 (1+bx_{\tau-1})}.$$
(4)

We suppose that  $\tau \in [\tau_0, \infty)$  where  $\tau = \tau_0$  is the first time that the sample matrix  $\mathbf{Z}$  becomes a full rank matrix. It is obvious that sequence  $[x_{\tau}]$  is decreasing, so that  $x_{\tau} \leq x_{\tau_0}$  for all  $\tau \in [\tau_0, \infty)$ . Then  $\mathbf{z}_{\frac{\alpha}{2}} \frac{\sigma \sqrt{\mathbf{z}^{\top} (\mathbf{Z}^{\top} \mathbf{Z})^{-1} \mathbf{z})}}{\sqrt{\tau}} \leq \frac{\eta_1}{\sqrt{\tau}}$  where  $\eta_1 = \mathbf{z}_{\frac{\alpha}{2}} \sigma \sqrt{x_{\tau_0} b}$ , so that

$$P\{||r - \tilde{r}||_{\infty} \le \frac{\eta_1}{\sqrt{\tau}}\} \ge 1 - \alpha.$$
(5)

#### B.3 Proof for Lemma 4

Proof. From Chebyshev's inequality, we have

$$P\left(|\tau - E[\tau]| \ge \sqrt{\frac{Var[\tau]}{\delta_{\tau}}}\right) \le \delta_{\tau}.$$

 $\operatorname{So}$ 

$$P(\tau \le E[\tau] - \sqrt{\frac{Var[\tau]}{\delta_{\tau}}}) \le P(|\tau - E[\tau]| \ge \sqrt{\frac{Var[\tau]}{\delta_{\tau}}}) \le \delta_{\tau}.$$

Let  $\tau_{jk}$  be the number of all customer types admitted to the kth resource of type j under the exploration subroutine. Then  $E[\tau] = K_0 \sum_{j=1}^N E[\tau_{jk}]$  considering recurrent arrivals.  $\tau_{jk}$  is a truncated Poisson random variable with rate  $\lambda_j = \sum_{i=1}^M D_{ij}$ , and it cannot exceed the resource capacity  $c_j$ . Therefore,

$$\begin{split} E[\tau_{jk}] &= \sum_{l=1}^{c_j} l \frac{\lambda_j^l e^{-\lambda_j}}{l!} + \sum_{l=c_j+1}^{\infty} c_j \frac{\lambda_j^l e^{-\lambda_j}}{l!} \\ &= \lambda_j \sum_{l=1}^{c_j} \frac{\lambda_j^{l-1} e^{-\lambda_j}}{(l-1)!} + \sum_{l=c_j+1}^{\infty} c_j \frac{\lambda_j^l e^{-\lambda_j}}{l!} \\ &\geq \lambda_j \sum_{l=0}^{c_j-1} \frac{\lambda_j^l e^{-\lambda_j}}{l!} \\ &= \lambda_j P(\text{Poisson}(\lambda_j) < c_j) \\ &= \lambda_j (1 - P(\text{Poisson}(\lambda_j) \ge c_j)) \\ &\geq \lambda_j (1 - \frac{\lambda_j}{c_j}), \end{split}$$

where we have used the Markov inequality in the last step. Hence,  $E[\tau] \ge K_0 \sum_{j=1}^N \lambda_j (1 - \frac{\lambda_j}{c_j})$ .

Since  $\tau$  is the sum of those truncated Poisson random variables, its variance is smaller than that of the sum of those un-truncated Poisson random variables. So  $Var[\tau] \leq K_0 \sum_{j=1}^N \lambda_j$ . Let  $\mu = K_0 \sum_{j=1}^N \sum_{i=1}^M D_{ij} (1 - \frac{\sum_{i=1}^M D_{ij}}{c_j}) - \sqrt{\frac{K_0 \sum_{j=1}^N \sum_{i=1}^M D_{ij}}{\delta_{\tau}}} \leq E[\tau] - \sqrt{\frac{Var[\tau]}{\delta_{\tau}}}$ . Then  $P(\tau \leq \mu) \leq P(\tau \leq E[\tau] - \sqrt{\frac{Var[\tau]}{\delta_{\tau}}}) \leq \delta_{\tau}$ .

#### B.4 Proof for Theorem 1

*Proof.* Let  $F_{MN+1}$  denote the event that the observation matrix **Z** has full rank  $M \times N + 1$ . Then by Lemma 3 and Lemma 4 we have

$$\begin{split} P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}} |\tau, F_{MN+1}\} &= P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}} |\tau \geq \mu, F_{MN+1}\} P(\tau \geq \mu) \\ &+ P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}} |\tau \leq \mu, F_{MN+1}\} P(\tau \leq \mu) \\ &\leq P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}} |\tau \geq \mu, F_{MN+1}\} + P(\tau \leq \mu) \\ &\leq \alpha + \delta_{\tau}. \end{split}$$

## B.5 Proof for Lemma 5

*Proof.* Note that  $1 - x \le e^{-x}$  for  $0 \le x \le 1$ , therefore

$$1 - \frac{d_{ij}^*}{c_j e} \le e^{-\frac{d_{ij}^*}{c_j e}}.$$
 (6)

Thus  $(1 - \frac{d_{ij}^*}{c_j e})^{K_0' c_j} \le e^{-\frac{d_{ij}^* K_0'}{e}}$ , and

$$P(F_{MN+1}) \ge \prod_{j=1}^{N} (1 - \sum_{i=1}^{M} [(1 - \frac{d_{ij}^*}{c_j e})^{K'_0 c_j}]) p_{delay}$$

$$\tag{7}$$

$$\geq \prod_{j=1}^{N} (1 - \sum_{i=1}^{M} e^{-\frac{d_{ij}^* K_0'}{e}}) p_{delay}$$
(8)

For an exploration phase of length  $K_0 \ge -\frac{e}{\min_{i,j} d_{ij}^*} \ln \frac{1 - (\frac{p_f}{p_{delay}})^{\frac{1}{N}}}{M} - \frac{\log(1 - p_{delay})}{D_L}$ , we have

$$P(F_{MN+1}) \ge \prod_{j=1}^{N} (1 - \sum_{i=1}^{M} e^{\frac{d_{ij}^*}{\min_{i,j} d_{ij}^*} \ln \frac{1 - (\frac{p_f}{p_{delay}})^{\frac{1}{N}}}{M}}) p_{delay}$$
$$\ge \prod_{j=1}^{N} (1 - \sum_{i=1}^{M} \frac{1 - (\frac{p_f}{p_{delay}})^{\frac{1}{N}}}{M}}{M}) p_{delay}$$
$$= p_f.$$

н		

#### B.6 Proof for Theorem 2

*Proof.* From Theorem 1, we get:

$$\begin{split} P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}}\} &= P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}}|F_{MN+1}\}P(F_{MN+1}) + \\ P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}}|\bar{F}_{MN+1}\}P(\bar{F}_{MN+1}) \\ &\leq P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}}|F_{MN+1}\}P(F_{MN+1}) + P(\bar{F}_{MN+1}) \\ &= P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}}|F_{MN+1}\}P(F_{MN+1}) + 1 - P(F_{MN+1}) \\ &= (P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}}|F_{MN+1}\} - 1)P(F_{MN+1}) + 1 \\ &\leq (P\{||r - \tilde{r}||_{\infty} \geq \frac{\eta_{1}}{\sqrt{\mu}}|F_{MN+1}\} - 1)p_{f} + 1 \\ &\leq (\alpha + \delta_{\tau} - 1)p_{f} + 1 \end{split}$$

# C Proofs for Section 6

#### C.1 Proof for lemma 6

*Proof.* We obtain the total regret by summing up the regret in the exploration phase and that in the exploitation phase as follows

$$Reg^{TP} = Reg^{TP}(0, T_0) + Reg^{TP}(T_0, T).$$
(9)

The first term is the regret incurred in the exploration phase, and it is smaller than  $K_0||c||_1$ since each single resource can only incur regret at most 1. Along with regret of exploitation phase, we have

$$Reg^{TP}(0,T) \le K_0 ||c||_1 + \frac{3}{2}(K - K_0) ||c||_1 \eta_1 (\eta_2 K_0 - \eta_3 \sqrt{K_0})^{-0.5}$$
$$\le K_0 ||c||_1 + \frac{3}{2}(K - K_0) ||c||_1 \eta_1 (\eta_2 - \eta_3)^{-0.5} K_0^{-0.5}$$

#### C.2 Proof for Lemma 7

*Proof.* According to the expression  $Reg(K_0) = K_0 ||c||_1 + \eta_4 ||c||_1 (K - K_0) K_0^{-0.5}$ , we have

$$\frac{\partial Reg(K_0)}{\partial K_0} = a_1 - a_2 K_0^{-1.5} + a_3 K_0^{-0.5},\tag{10}$$

where  $a_1 = ||c||_1$ ,  $a_2 = 0.5||c||_1\eta_4 K$ , and  $a_3 = -0.5||c||_1\eta_4$ . Therefore, to find the root of the above equation, we have to solve a cubic polynomial

$$a_1 x^3 + a_3 x^2 - a_2 = 0, (11)$$

where  $x = K_0^{0.5}$ . Since the discriminant of the equation  $\Delta = -4a_3^3(-a_2) - 27a_1^2a_2^2 < 0$ , the equation has only one real root and two conjugate non-real roots.

To solve the cubic polynomial, we transform the equation into another standard form

$$x^3 + px + q = 0, (12)$$

where  $p = \frac{-a_3^2}{3a_1^2}$  and  $q = \frac{-27a_1^2a_2+2a_3^2}{27a_1^3}$ . Note that  $\frac{\partial(a_1x^3+a_3x^2-a_2)}{\partial x} = 3a_1x^2 + 2a_3x > 0$  for  $x \ge 0$ , because when K is a large number there are  $a_1 > 0 > a_3$  and  $|a_1| \gg |a_3|$ . So there is only one solution for (11) when x > 0, which is

$$x = \sqrt[3]{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}} + \sqrt[3]{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}.$$
(13)

Since p = O(1) and q = O(K), when K is a large number we have  $|q| \gg |p|$ . Hence, we have an approximate form of the solution, namely  $x = \sqrt[3]{-q} = \sqrt[3]{\frac{27a_1^2a_2 - 2a_3^3}{27a_1^3}}$ . Furthermore, we have

$$K_0^* = x^2 \approx \left(\frac{a_2}{a_1}\right)^{\frac{2}{3}} = \left(0.5\eta_4 K\right)^{\frac{2}{3}}.$$
(14)

# D Proof for Theorem 4

*Proof.* According to Theorem 2, with the current available data  $\mu_0^{\omega-1}$ , if an exploration phase  $K_0^{\omega}$  is selected, the estimation error will be

$$||r - \tilde{r}||_{\infty} \le \frac{\eta_1}{\sqrt{\mu^{\omega} + \mu_0^{\omega - 1}}} = \eta_1 (\eta_2 K_0^{\omega} - \eta_3 \sqrt{K_0^{\omega}} + \mu_0^{\omega - 1})^{-0.5} \le \eta_1 ((\eta_2 - \eta_3) K_0^{\omega} + \mu_0^{\omega - 1})^{-0.5},$$
(15)

where  $\mu^{\omega} = \eta_2 K_0^{\omega} - \eta_3 \sqrt{K_0^{\omega}}$ . So we can again proceed as before and now we have

$$Reg^{\omega}(K_0^{\omega}) = K_0^{\omega}||c||_1 + ||c||_1\eta_4(K - K_0^{\omega})(K_0^{\omega} + \frac{\mu_0^{\omega-1}}{\eta_2 - \eta_3})^{-0.5},$$

where  $\eta_4 = \frac{3}{2}\eta_1(\eta_2 - \eta_3)^{-0.5}$ . In the following let  $X_0^{\omega-1} = \frac{\mu_0^{\omega-1}}{\eta_2 - \eta_3}$ . So  $X_0^{\omega-1}$  is the total number of previous data in terms of the number of cycles.

According to the expression of  $Reg^{\omega}(K_0^{\omega})$ , we have

$$\frac{\partial Reg^{\omega}(K_0^{\omega})}{\partial K_0^{\omega}} = a_1 - a_2(K_0^{\omega} + X_0^{\omega-1})^{-1.5} + a_3(K_0^{\omega} + X_0^{\omega-1})^{-0.5},$$
(16)

where  $a_1 = ||c||_1$ ,  $a_2 = 0.5||c||_1\eta_4(K + X_0^{\omega-1})$ , and  $a_3 = -0.5||c||_1\eta_4$ . Therefore, to find the root of the above equation, we have to solve a cubic polynomial

$$a_1 x^3 + a_3 x^2 - a_2 = 0, (17)$$

where  $x = (K_0^{\omega} + X_0^{\omega-1})^{0.5}$ . Since the discriminant of the equation  $\Delta = -4a_3^3(-a_2) - 27a_1^2a_2^2 < 0$ , the equation has only one real roots.

To solve the cubic polynomial, we transform the equation into another standard form

$$x^3 + px + q = 0, (18)$$

where  $p = \frac{-a_3^2}{3a_1^2}$  and  $q = \frac{-27a_1^2a_2 + 2a_3^3}{27a_1^3}$ . The solution is

$$x = \sqrt[3]{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}} + \sqrt[3]{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}.$$
(19)

Since p = O(1) and  $q = -O(X_0^{\omega-1})$ , when  $X_0^{\omega-1} \to \infty$  we have  $|q| \gg |p|$ . Hence, we have an approximate form of the solution, namely  $x \approx \sqrt[3]{-q} = \sqrt[3]{\frac{27a_1^2a_2 - 2a_3^3}{27a_1^3}}$ . Furthermore, since  $a_1 = O(1), a_2 = O(X_0^{\omega - 1})$  and  $a_3 = O(1)$ , when  $X_0^{\omega - 1} \to \infty$ , we have

$$K_0^{\omega} = x^2 - X_0^{\omega - 1} \approx \left(\frac{a_2}{a_1}\right)^{\frac{2}{3}} - X_0^{\omega - 1} = \left(\frac{\eta_4(K + X_0^{\omega - 1})}{2}\right)^{\frac{2}{3}} - X_0^{\omega - 1}.$$
 (20)

Inserting  $K_0^{\omega}$  into the regret bound and considering  $K \to \infty$ , we have

$$Reg^{\omega} = O(\left[\eta_4(K + X_0^{\omega-1})\right]^{\frac{2}{3}} ||c||_1 - X_0^{\omega-1}||c||_1).$$
(21)

It is clear that within each epoch, the regret  $Reg^{\omega}$  has the same order as the exploration length  $K_0^{\omega}$ , i.e.,  $Reg^{\omega} = O(K_0^{\omega})$ . Let  $Reg^{MP} = \sum_{\omega=1}^{\Omega} Reg^{\omega}$ . Then  $Reg^{MP} = O(X_0^{\Omega})$ , for  $X_0^{\Omega} = \sum_{\omega=1}^{\Omega} K_0^{\omega}$ . Note that  $X_0^{\omega-1}$  is increasing with  $\omega$  according to the definition, and when  $X_0^{\omega-1}$  is large enough,  $K_0^{\omega}$  goes to 0. This implies that  $X_0^{\omega-1}$  will converge. Note that  $X_0^1 = O(K^{\frac{2}{3}})$ , and it also implies that  $X_0^{\omega-1} = O(K^{\frac{2}{3}})$  for all  $\omega = 3, ..., \Omega + 1$  since if  $X_0^{\omega-1}$  exceeds  $O(K^{\frac{2}{3}})$ ,  $K_0^{\omega}$  would become negative when  $K \to \infty$ . Therefore, when  $K \to \infty$ , we have  $Reg^{MP} = O(X_0^{\Omega}) = O(K^{\frac{2}{3}})$ . Let  $K_{total} = \Omega K$  be the total number of cycles for all horizons, then we have  $Reg^{MP} = O(K_{total}^{\frac{2}{3}})$ .

# E Table of average regret of t-test

	Mean of average regret			Std of average regret			Confidence interval width			T-test result	
K	TP	Greedy	$\epsilon$ Greedy	TP	Greedy	$\epsilon$ Greedy	TP	Greedy	$\epsilon$ Greedy	Greedy	$\epsilon$ Greedy
11		exp.	- Conceany		exp.			exp.	- Conceany	exp.	c orecuy
5	68.2	83.0	73.2	33.4	18.1	11.4	12.0	6.5	4.1	0.72	0.78
10	64.0	88.0	75.0	32.5	24.6	7.8	11.6	8.8	2.8	1.48	1.80
15	63.3	84.0	79.3	6.4	14.1	5.2	2.3	5.0	1.9	5.67	10.65
20	61.5	84.5	76.0	7.1	13.0	5.4	2.5	4.7	1.9	5.37	8.94
25	62.4	82.4	73.6	6.7	7.0	5.2	2.4	2.5	1.9	6.34	7.23
30	62.7	80.3	74.7	5.5	6.3	6.1	2.0	2.3	2.2	7.84	7.98
35	63.1	80.6	74.6	6.1	6.1	5.3	2.2	2.2	1.9	7.27	7.79
40	62.8	79.0	73.3	6.5	5.3	5.2	2.3	1.9	1.9	6.84	6.90
45	63.8	79.8	73.8	5.6	6.4	4.3	2.0	2.3	1.5	6.48	7.80
50	62.0	79.0	72.8	6.1	6.3	5.0	2.2	2.3	1.8	6.77	7.52
55	61.8	77.5	72.9	6.2	5.6	4.7	2.2	2.0	1.7	7.27	7.84
60	63.0	78.3	73.7	6.0	5.1	4.3	2.1	1.8	1.5	7.41	7.93
65	61.5	77.7	72.5	5.4	5.8	4.4	1.9	2.1	1.6	7.55	8.53
70	61.4	78.6	72.9	6.0	4.6	3.8	2.1	1.6	1.4	8.32	8.84
75	61.3	77.3	72.0	6.0	3.0	4.0	2.2	1.1	1.4	8.71	8.11
80	62.5	77.5	72.5	5.5	3.5	4.7	2.0	1.2	1.7	8.46	7.63
85	62.4	77.6	71.8	5.7	4.4	4.9	2.0	1.6	1.8	7.13	6.82
90	61.1	76.7	72.2	5.4	4.4	4.2	1.9	1.6	1.5	8.78	8.88
95	62.1	77.9	72.6	6.0	5.3	4.2	2.1	1.9	1.5	7.20	7.86
100	61.0	77.0	72.0	6.2	2.8	3.7	2.1	1.0	1.3	8.86	8.38
105	61.0	77.1	72.0 72.4	6.0	3.4	3.7 4.4	2.2 2.1	1.0 1.2	1.5 1.6	9.07	8.44
105	61.8	78.2	72.4	5.4	3.4 4.8	4.4 4.1	1.9	1.2 1.7	1.0 1.5	9.07 8.29	8.84
115	60.9		72.2	$5.4 \\ 5.9$	$4.0 \\ 5.6$	4.1 3.0	1.9 2.1	2.0	1.5	7.64	0.04 9.42
		77.4									
120	61.7	77.5	71.7	5.9	4.6	3.3	2.1	1.6	1.2	7.30	8.06
125	60.8	76.8	71.2	5.5	2.2	3.3	2.0	0.8	1.2	9.68	8.95
130	61.5	77.7	72.3	5.7	3.7	2.8	2.1	1.3	1.0	8.62	9.23
135	61.5	77.0	71.9	5.6	3.9	2.5	2.0	1.4	0.9	8.31	9.19
140	61.4	77.1	72.1	5.8	4.8	4.2	2.1	1.7	1.5	7.83	8.19
145	61.4	77.2	72.4	5.5	4.0	2.2	2.0	1.4	0.8	8.85	10.14
150	61.3	77.3	72.0	5.7	4.1	3.3	2.0	1.5	1.2	8.39	8.92
155	61.3	76.8	71.6	5.7	3.5	4.0	2.0	1.3	1.4	8.40	8.11
160	61.9	77.5	72.5	5.7	2.7	2.8	2.1	1.0	1.0	9.17	9.09
165	61.2	77.0	72.1	5.8	3.6	3.0	2.1	1.3	1.1	8.73	9.17
170	61.2	77.1	71.8	5.9	2.9	3.4	2.1	1.0	1.2	8.89	8.53
175	60.6	76.6	71.4	5.9	2.7	3.8	2.1	1.0	1.3	9.11	8.46
180	61.1	77.2	72.2	5.9	3.6	2.4	2.1	1.3	0.9	8.73	9.49
185	61.1	76.8	71.9	6.3	3.7	2.2	2.3	1.3	0.8	8.09	8.84
190	60.5	76.8	71.6	5.7	2.3	2.1	2.0	0.8	0.8	9.85	9.98
195	61.5	77.4	72.3	5.4	3.4	2.0	1.9	1.2	0.7	9.25	10.28
200	61.5	77.0	72.0	5.5	3.5	2.5	2.0	1.3	0.9	8.79	9.52
205	61.0	77.1	71.7	5.7	3.3	2.3	2.0	1.2	0.8	8.98	9.61
210	61.4	77.1	72.4	5.3	3.5	1.9	1.9	1.3	0.7	9.46	10.70
215	61.4	77.2	72.1	5.4	4.4	2.1	1.9	1.6	0.7	8.37	10.06
220	60.9	76.8	71.8	5.8	2.3	1.1	2.1	0.8	0.4	9.54	10.08
225	60.9	76.9	71.6	5.8	4.0	2.1	2.1	1.4	0.8	8.33	9.48
230	61.3	77.4	72.2	5.7	3.7	1.6	2.0	1.3	0.6	8.83	10.13
235	61.3	77.0	71.9	5.3	3.0	1.7	1.9	1.1	0.6	9.60	10.53
240	61.3	77.1	71.7	5.5	0.8	1.1	2.0	0.3	0.4	10.19	10.12
245	61.2	77.1	72.2	5.4	2.8	1.6	1.9	1.0	0.6	9.89	10.66
250	61.2	77.2	72.0	5.4	3.3	1.9	1.9	1.2	0.7	9.37	10.35

Table 1: Summary of average regret and results of t-test.