

6) $p \in P_K$ is a probability vector belonging to a K -dimensional vector space. Therefore any vector which is perpendicular to it, ~~consequently~~ will have its coefficients satisfy,
 $\sum a_j p_j = 0$ as $\langle a, p \rangle = 0$ where $a \in \mathbb{R}^K$ is a vector.

For every vector there exists a normal so such an 'a' exists.

Now if $\hat{x}(i, x) = a_i + \frac{\mathbb{I}_{\{i=1\}} x_1}{p_1}$
 i.e.,

$$\hat{x}(i, x) = \begin{cases} a_i + \frac{x_1}{p_1} & \text{if } i=1 \\ a_i & \text{otherwise} \end{cases}$$

$$\begin{aligned} E(\hat{x}(i, x)) &= \sum p_i \hat{x}(i, x_i) \\ &= \sum_{i=1}^K \left(p_i a_i + \frac{p_i \mathbb{I}_{\{i=1\}} x_1}{p_1} \right) \\ &= \sum_i p_i a_i + \frac{\sum_i p_i \mathbb{I}_{\{i=1\}} x_1}{p_1} \\ &= 0 + \frac{p_1 x_1}{p_1} \\ &= x_1 \end{aligned}$$

\therefore Any a perpendicular to $p \in P_K$ will satisfy the equation.

$$5) R_n^{\text{track}}(\pi, \nu) = \mathbb{E} \left[\sum_{t=1}^T \max_i x_{t,i} - \sum_{t=1}^T x_{t,I_t} \right]$$

Now for any policy, at time t , the sum over probabilities of all actions is 1

$$\therefore \sum_{i \in [K]} p_t(i) = 1$$

As all $p_t(i) \geq 0$ we can say there exists ~~one~~ at least one $i \in [K]$ such that $p_t(i) \leq 1/K$, for any distribution. — ①

Let the sequence of rewards be $\{x_t\}$ such that,

$$x_{t,i} = \begin{cases} 1 & \text{when } i \text{ is such that } p_t(i) = \min_{j \in K} \{p_t(j)\} \\ 0 & \text{otherwise} \end{cases}$$

$$= \mathbb{E} \left[\sum_{t=1}^T \max_i x_{t,i} - \sum_{t=1}^T x_{t,I_t} \right],$$

Now, $\max_i x_{t,i} = 1 \quad \forall t$

$$\therefore \text{We have, } T - \mathbb{E} \sum_{t=1}^T x_{t,I_t}$$

$$= T - \sum_{t=1}^T \mathbb{E}(x_{t,I_t})$$

Now, $\forall x_{t,i} = 0 \quad \forall i$, except one where $p(i) = \min_{j \in K} p(j)$

$$\therefore \mathbb{E}(x_{t,I_t}) = p_t(i^*), \text{ where } i^* = \operatorname{argmin}_{j \in K} p(j)$$

$$= T - \sum_{t=1}^T p_t(i^*) \leq 1$$

$$\geq T - \sum_{t=1}^T \frac{1}{K} \quad (\text{from ①})$$

$$= T - T/K$$

$$= T(1 - 1/K)$$

$$\therefore R_n(\pi, \nu) \geq T(1 - 1/K)$$

4) Now,

$$R_T(\pi, \nu) = \max_{i \in [K]} \sum_{t=1}^T x_{t,i} - \mathbb{E} \left[\sum_{t=1}^T x_{t,I_t} \right]$$

In a deterministic policy, the sequence of actions of the player (i.e., I_t) is fixed. So there might exist an environment which generates a sequence of rewards such that,

$$x_{t,i} = \begin{cases} 0 & i = I_t \\ 1 & i = [K] - I_t \end{cases}$$

$$\begin{aligned} \therefore R_T(\pi, \nu) &= \max_{i \in [K]} \sum_{t=1}^T x_{t,i} - \mathbb{E} \left[\sum_{t=1}^T x_{t,I_t} \right] \\ &= \max_{i \in [K]} \sum_{t=1}^T x_{t,i} \quad [\because x_{t,I_t} = 0 \forall t] \end{aligned}$$

$$\text{Now, } \max_{i \in [K]} \sum_{t=1}^T x_{t,i} \geq \sum_{t=1}^T x_{t,j} \quad \forall j \in [K]$$

$$\therefore K \cdot \max_{i \in [K]} \sum_{t=1}^T x_{t,i} \geq \sum_{j \in [K]} \sum_{t=1}^T x_{t,j} \quad [\text{Summing over all inequalities for } j \in [K]]$$

$$= \sum_{t=1}^T \sum_{j \in [K]} x_{t,j} \quad [\text{Interchanging summation}]$$

$$= \sum_{t=1}^T (K-1)$$

$$\Rightarrow \max_{i \in [K]} \sum_{t=1}^T x_{t,i} \geq \frac{T(K-1)}{K}$$

$$= T \left(1 - \frac{1}{K} \right)$$

$$\therefore R_T(\pi, \nu) \geq T \left(1 - \frac{1}{K} \right) \quad \text{Hence proved.}$$