

Optimal Task Assignment by Learning Members' Skills (MURI)

Omid Askari

November 4, 2015

1 Problem Setup

Consider we have a group of people and there is a sequence of tasks. Suppose each person has a set of skills and each task requires some of these tasks to a extent.

$people = \{p_1, p_2, \dots, p_n\}$, $tasks = \{t_1, t_2, \dots, t_T\}$

Each person $p = (s_1, s_2, \dots, s_m)$ and task $t = (r_1, r_2, \dots, r_m)$ has m skills set and requirements. We suppose that skill values are discrete and fall in $[0, b]$ and b is a configuration parameter which directly affects problem's complexity. The larger b is, bigger the boundary of learning space would be.

2 Problem Definition

Suppose we know task requirements; however, people's skills are unknown. We need to learn the values of people's skills through applying payoff function. We also consider that for handling each task, merely one person is needed. Additionally, payoff function f for person p and task t is defined as follows

$$f(p, t) = \begin{cases} 1, & s_1 \geq r_1 \wedge s_2 \geq r_2 \wedge \dots \wedge s_m \geq r_m \\ 0, & otherwise \end{cases} \quad (1)$$

In fact, the goal is handling more tasks successfully in the long run.

3 Proposed Method

To maximize the performance, we have to learn people's skills. If we know everybody's expertises, we assign best people to their matching task to a great extent To learn, we propose a simple rule-based learning approach. If agent handles successfully the task, it conveys:

$s_1 \geq r_1 \wedge s_2 \geq r_2 \wedge \dots \wedge s_m \geq r_m$

1 and he or she fails it means:

2 $s_1 < r_1 \vee s_2 < r_2 \vee \dots \vee s_m < r_m$

3 In this method, we can learn, and hence, limit the boundaries for each person's skills with assigning more tasks to him or her.

5 3.1 Rule-based learning example

6 For a simple example, we consider person p with 2 unknown skills, $p = (x, y)$.
7 Then, we assign task $t_1(2, 5)$ which requires known values for each skills as
8 follows,

9 $p(x, y), t_1(2, 5) \Rightarrow \text{success}$, therefore : $x \geq 2 \wedge y \geq 5$

10 Then we assign another task,

11 $p(x, y), t_1(4, 1) \Rightarrow \text{failure}$, therefore : $x \geq 2 \wedge y \geq 5 \wedge (x < 4 \vee y < 1)$

12 Thus with a simple logic we can have,

13 $2 \leq x < 4 \wedge y \geq 5$.

14 Now we have a smaller boundary for x and y and having this knowledge could
15 pave the way for deciding more efficient in task assignment to person p .

16 3.2 Method Steps

17 One goal in this study could be finding a lower bound for the number of task
18 assignment to a single person with m different skills that vary in $[0, b]$ to ensure
19 that more than 90% of his or her skill values are uncovered so far. Even if we
20 consider that pay off function is stochastic, this will still be useful to have a
21 tight boundary for each person's skill.

22 First we start with no knowledge about people's skills and initialize skill
23 values with 0s. To learn actual values there are two important phases in each
24 round. In every round, people either explore or exploit. There is a method to
25 choose one of these phases; however, first we define them descriptively in the
26 following.

27 3.2.1 Exploration

In exploration, we choose one person who is completely unknown and has all
0s in his or her skills. If there is no one left unknown in people, then we will
choose that person has the largest boundary distance (who is more unknown).
To do this, we define a boundary $[x, y]$ such that $x \geq 0$ and $y \leq b$ for each skill
in every person. Boundary distance is defined as

$$dist = \sqrt{\sum_{i=1}^m (x_i - y_i)^2} \quad (2)$$

28 If there are more than one person having the largest distance, then we choose
29 one of them by random.

1 3.2.2 Exploitation

2 In exploitation, we choose one of known people (they have skills values greater
3 than 0) who has the largest skills for the require ones and it means it is the
4 most promising one for handling the task successfully (however in this method
5 we are somehow wasting resources and there is a more intelligent idea which is
6 presented in the 4 section). The selection method for task $t = (r_1, r_2, \dots, r_m)$
7 and if each person such as $p = (s_1, s_2, \dots, s_m)$ is

$$\arg \max_p o(t, p) \quad (3)$$

$$o(t, p) = \begin{cases} \sqrt{\sum_{i=1}^m (s_i - r_i)^2}, & \text{if } s_1 \geq r_1 \wedge s_2 \geq r_2 \wedge \dots \wedge s_m \geq r_m \\ -1, & \text{otherwise} \end{cases}$$

8 3.2.3 Exploration and exploitation trade-off

One of the most substantial part of this setting should be a policy to manage exploration and exploitation. First, we start with a constant probability such as c (which means $0 < c < 1$) for exploration $P_{exploration} = c$ and $P_{exploitation} = 1 - c$ and analyze how this parameter could affect the performance in the long run. To do this, we run the whole setting with different values for c . Second, there are various existing methods in Multi Armed Bandits literature [1, 2, 3] to manage between exploration and exploitation; however, for the sake of simplicity, we will apply the following method which is completely logical,

$$P_{exploration} = \frac{c}{t} \quad (4)$$

9 where c is a constant and t is the time-step variable. In the beginning exploration
10 power is very high; however, since t increases over time; as a consequence, the
11 exploitation power gradually increases and exploration decays. This concept is
12 inspired since in the beginning we do not have enough knowledge and is better
13 to explore more; nonetheless, as time passes, we know more and hence it is
14 better to exploit more and refine the available solutions.

15 4 Extensions

- 16 1. Find a group of people instead of just one person to handle a task.
- 17 2. We suppose there is a prior knowledge regarding each person; however, this
18 prior probability gradually change through time with a Bayesian learning.
3. The payoff could be continuous instead of binary. Also it cannot be euclidean distance for all, thus:

$$f(p, t) = \sqrt{\sum_{i=1}^m u_i}; \quad u_i = \begin{cases} (r_i - s_i)^2, & \text{if } r_i \geq s_i \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

4. The payoff function could be stochastic as follows:

$$f(p, t) = \begin{cases} 1, & \text{if } s_1 + \mathcal{N}(1, \sigma^2) \geq r_1 \wedge s_2 + \mathcal{N}(1, \sigma^2) \geq r_2 \wedge \dots \wedge s_m + \mathcal{N}(1, \sigma^2) \geq r_m \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

- 1 5. We can use some other multi armed bandits algorithms (contextual MAB)
- 2 and for instance Soft Max to facilitate exploration and exploitation trade-
- 3 off [2, 3].
- 4 6. If we consider that tasks are decomposable then we can have a group of
- 5 people handling a task instead of solely one person. Therefore, to solve
- 6 this problem, we will use solutions for multiple knapsack problem [4].
- 7 7. In exploitation phase, it is more logical choosing whom is the most appro-
- 8 priate for the task and not who has the largest skills since we are somehow
- 9 wasting resources and we are using the best person for a easy task. The
- 10 selection method could be changed as

$$\arg \min_p o(t, p) \text{ st } o(t, p) > 0 \quad (7)$$

$$o(t, p) = \begin{cases} \sqrt{\sum_{i=1}^m (s_i - r_i)^2}, & \text{if } s_1 \geq r_1 \wedge s_2 \geq r_2 \wedge \dots \wedge s_m \geq r_m \\ -1, & \text{otherwise} \end{cases}$$

8. We also can suppose that people after handling a task, grow some knowl-
- edge and experience. Consequently, their skills are gradually increased.
- For instance after being successful in a task, skill i will be updated as

$$s_i \leftarrow s_i + \alpha \cdot r_i \quad (8)$$

And if he or she fails then

$$s_i \leftarrow s_i + \beta \cdot r_i \quad (9)$$

- 11 Where α and β both are configurable parameters such that $0 < \alpha < 1$,
- 12 $0 < \beta < 1$, $\alpha > \beta$. Hence, if one task requires one specific skill more than
- 13 other, person will increase that skill more comparing other ones.

14 References

- 15 [1] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of
- 16 the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- 17 [2] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-
- 18 bandit approach to personalized news article recommendation. In *Proceed-*
- 19 *ings of the 19th international conference on World wide web*, pages 661–670.
- 20 ACM, 2010.

- 1 [3] Volodymyr Kuleshov and Doina Precup. Algorithms for multi-armed bandit
2 problems. *arXiv preprint arXiv:1402.6028*, 2014.
- 3 [4] Chandra Chekuri and Sanjeev Khanna. A polynomial time approximation
4 scheme for the multiple knapsack problem. *SIAM Journal on Computing*,
5 35(3):713–728, 2005.