



Department of Industrial and Systems Engineering

UNIVERSITY OF WISCONSIN-MADISON



Problem Statement

Optimal team formation is crucial in domains like work preferences for working with others, which influence team performance. Key challenges include:

- **Unknown Preferences:** Individual preferences are not always known or clearly defined.

cial
ing

N

ent

places and sports. Each individual has team performance. The main challenges are not directly observable.

Optimizing

Ka



Overview: Using by updating beliefs are estimates of shown in Figure

Team Assignments by F

asper Veje Jakobsen, kjakobsen@wisd.edu

¹Department of Industrial and Systems Engineering



Learning Preferences: Kalman Filter

g a Kalman filter, the preferences between individuals are dynamic based on feedback received after team interactions, using $\hat{\sigma}$ the process noise and feedback noise. An example of the lea

1.



Preferences

μ^1



Simulation

nically estimated
 w and $\hat{\sigma}_f$ which
rning process is

Setup: The simulation evaluates team a
periods, and $s = 3$ teams. Preference
($\sigma_f = 0.1$). A **Random Assignment** strat
weights β are compared. Results are ave
of performances shown in Figure 3 and



on Setup and Results

assignment strategies using $n = 10$ individuals, $m = 100$ tasks evolve with noise ($\sigma_w = 0.1$), and feedback is noisy. The **Greedy Strategy** and **Optimal Assignments** with varying exploration were averaged over 10 simulations, with the mean and variance

Figure 4

- **Dynamic Evolution:** Preferences evolve continuously.
- **Complexity:** Large action space (all possible team assignments given preferences and uncertainties).

Objective: Develop a framework to assign individuals to

- **Maximizing** overall **team performance**.
- Balancing **exploration** (learning preferences) and **exploitation**



Mathematical Notation

Individuals: $I = \{1, 2, \dots, n\}$: Set of n individuals.

Team Assignments: At each period p , individuals are grouped into teams

$$\mathcal{T}^p = \{T_1^p, T_2^p, \dots, T_{s_p}^p\}, \quad T_j^p \cap T_k^p = \emptyset$$

Preferences: The true (unobserved) preference $w_{i,i}^p$ evolves over time



y towards all individuals.

signments) and continuous state space

teams dynamically while:

exploitation (using known preferences).

Optimization

grouped into mutually exclusive teams:

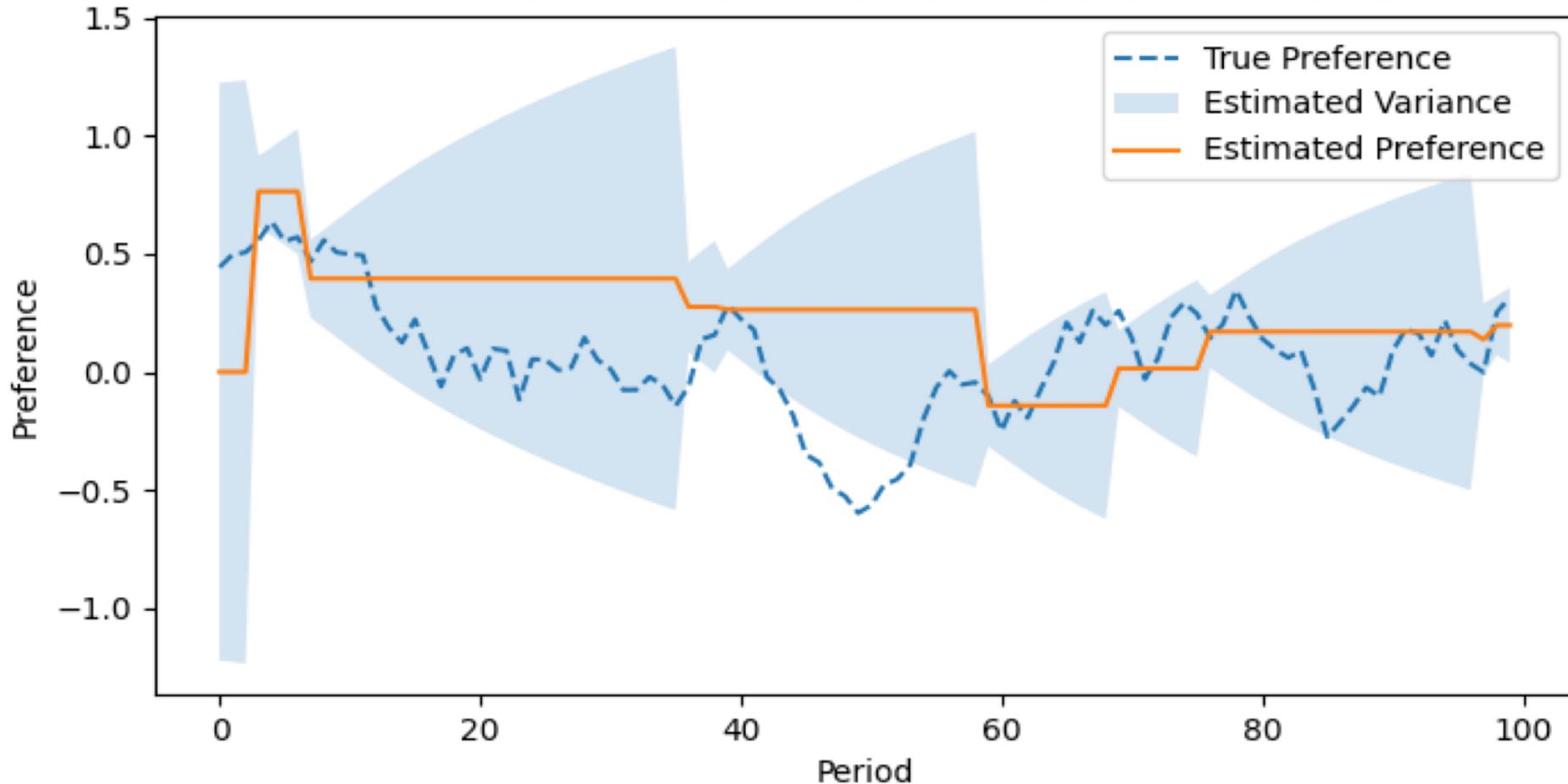
$$= \emptyset, \quad \bigcup_{j=1}^{s_p} T_j^p = I.$$

yes continuously:

Figure 1.



Preference Estimates for Individual 55 and Individual 33



Estimated preferences μ_{ij} compared to true preferences w_{ij} . σ_{ij} reflects un

Optimal Assignments

Metrics: The reward measures team performance by summing up the individual assignments. The preference distance measures the deviation from the true preferences (w_{ij}) using the L_1 norm.

certainty.

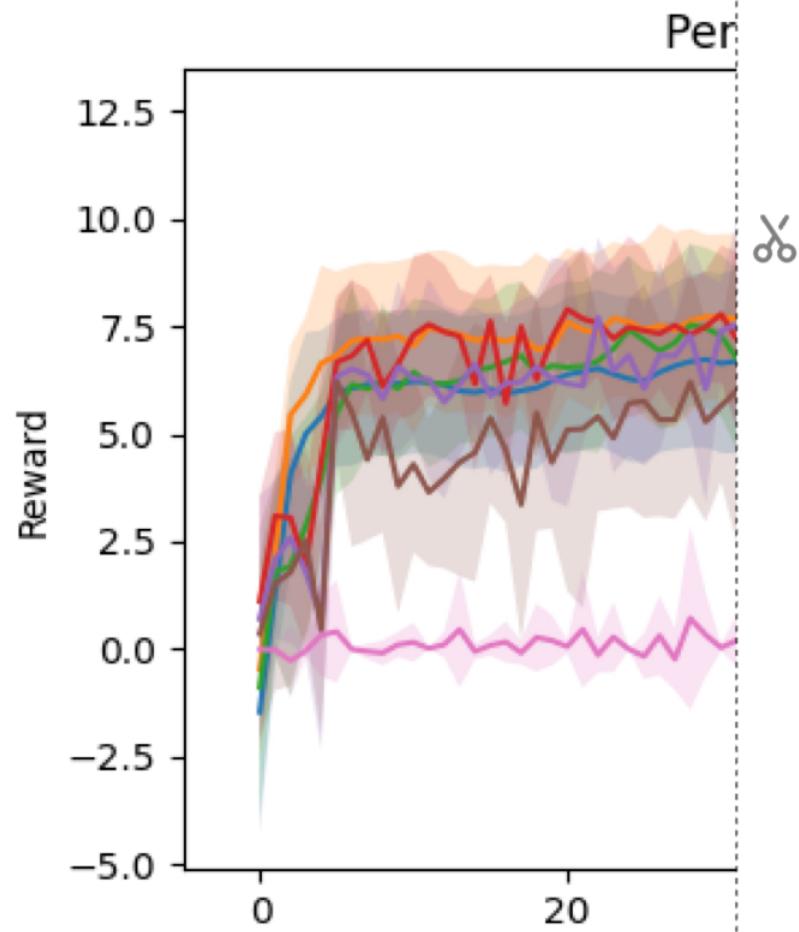
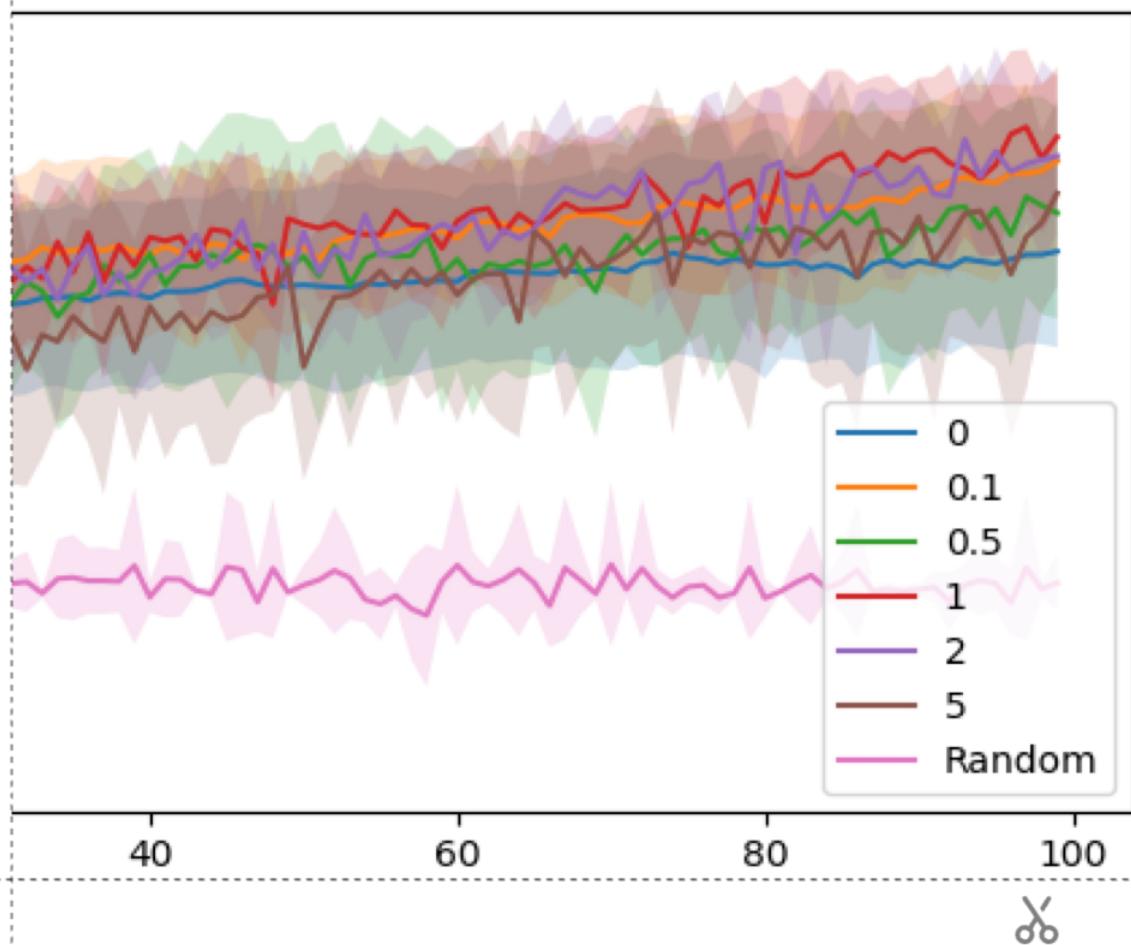


Figure 7.

performance over periods, reflecting the effectiveness of evaluates the convergence of estimated preferences (μ_{ij}) ℓ_2 -norm, indicating the learning accuracy of the model.

Performance of Different Actors



$$w_{ij}^0 = N(0, \sigma_p^2), \quad w_{ij}^{p+1} = w_{ij}^p + \varepsilon_{ij}^p$$

Feedback: Sparse feedback F_{ij}^p is collected for individual i

$$F_{ij}^p = w_{ij}^p + \eta_{ij}^p, \quad \eta_{ij}^p \sim \text{Unif}(0, 1)$$

DP Formulation



State Space: The state at each period represents the belief probability distribution:

$$\hat{s}_p = \left\{ P(w_{ij}^p \mid \text{history up to } p) \right\}_{i,j}$$

This belief can be fully characterized by the current estimates, where μ_{ij}^p is the estimated preference and σ_{ij}^p is the standard deviation.

$$s_n = \{\mu_{ij}^p, \sigma_{ij}^p \mid i, j \in \mathcal{N}\}$$



$$, \varepsilon_{ij}^p \sim N(0, \sigma_w^2)$$

Is in the same team:

$$N(0, \sigma_f^2)$$

on

belief about preferences, expressed as a

$$) \mid i, j \in I \} .$$

States of preferences and their uncertainty:

$$I \},$$

Objective: Assign
(learning preference)

Upper Confidence

into a score S_{ij}^p

Teams are assigned

MILP: Used as a
at each period.

Future Directions:
efficiently handle

n individuals to teams to maximize overall performance by balances) and **exploitation** (using known preferences).

nce Bound (UCB) Strategy: Combines mean preference μ_{ij}^p and using an exploration weight $\beta \geq 0$:

$$S_{ij}^p = \mu_{ij}^p + \beta \cdot \sigma_{ij}^p.$$

ned to maximize the sum of scores:

$$a_p = \arg \max_{a \in \mathcal{A}} \sum_{T \in a} \sum_{i, j \in T, i \neq j} S_{ij}^p.$$

baseline, MILP directly optimizes team assignments by maximizin

ns: Explore scalable methods like multi-agent reinforcement learning for large action spaces and dynamic team assignment scenarios.

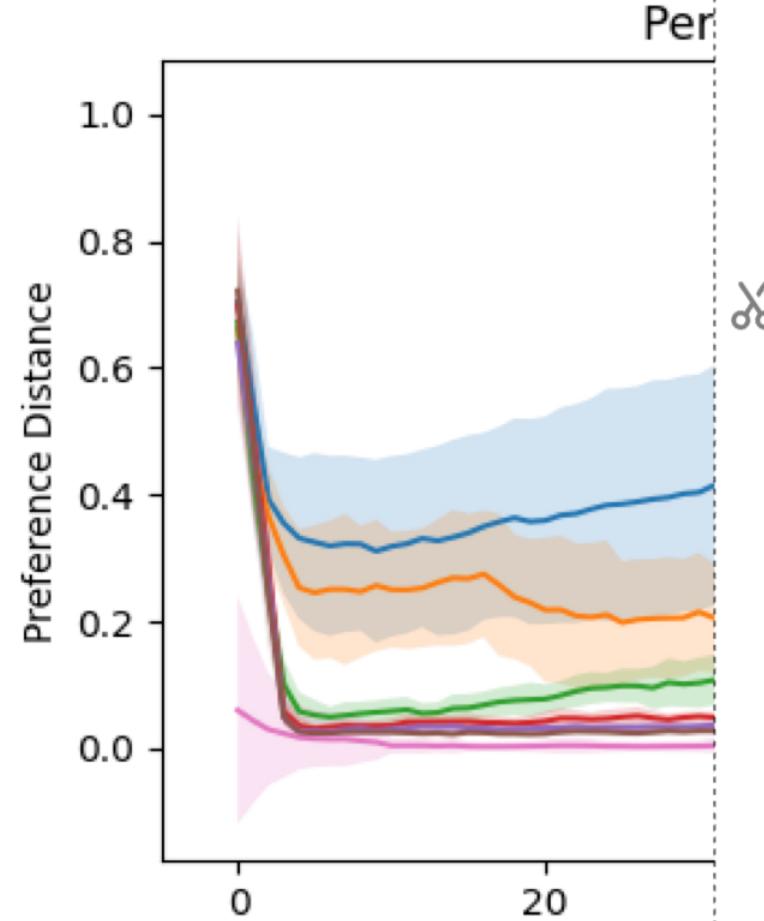
cing exploration

uncertainty σ_{ij}^p

g the score (S_{ij}^p)

rning (MARL) to

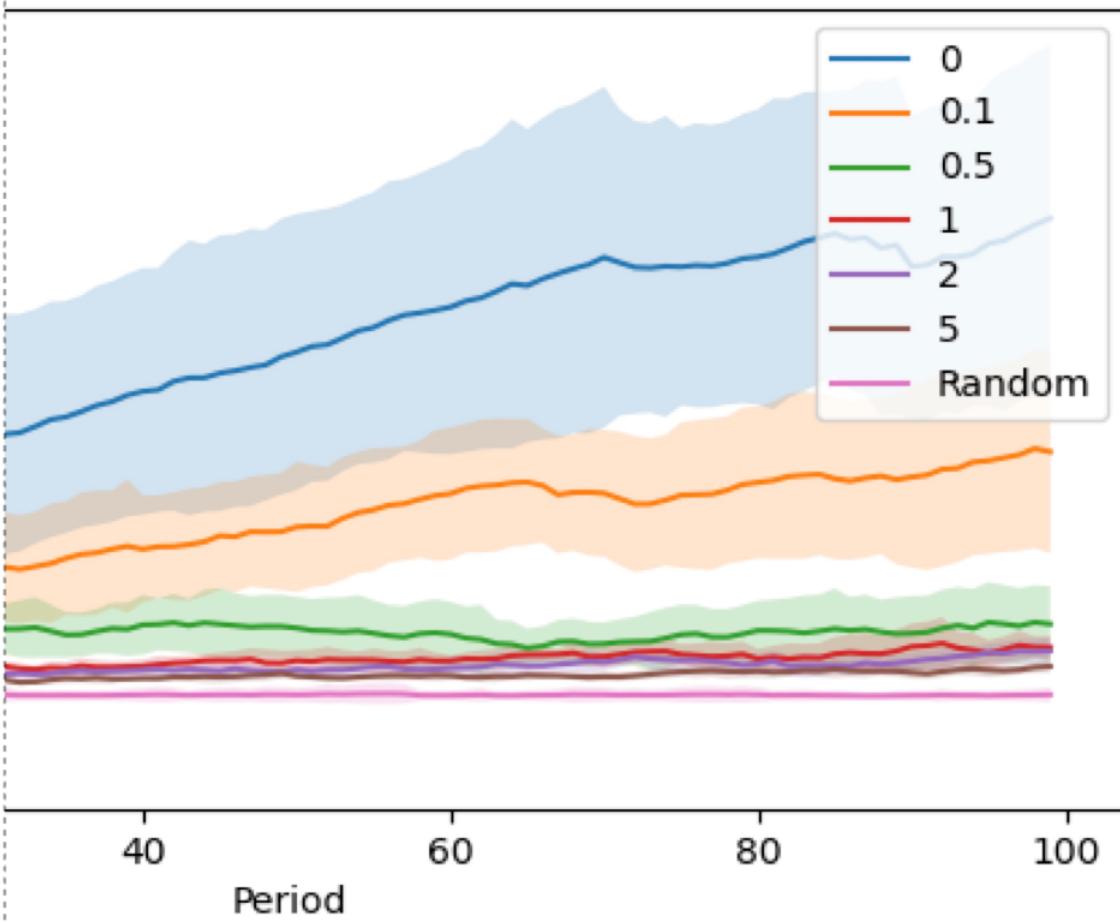
Figure 3. Performance comparison across



Period

strategies. Higher rewards indicate better team assignments.

Performance of Different Actors



Action Space: The action $a_p \in \mathcal{A}$ is a valid team assignment.

Transition Function: Beliefs about the preferences evolve

$$s_{p+1} = f(s_p, F^p), \quad \text{where } F^p \text{ incorporates feedback.}$$

Reward Function: The reward reflects team performance.

$$r(s_p, a_p) = \sum_{T \in a_p} \left(\sum_{i, j \in T, i < j} \dots \right)$$

Objective: Find the policy π^* that maximizes the expected reward.

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{p=1}^m r(s_p) \right]$$

ent for the period.

e based on feedback:

back on evolving preferences.

e, measured by the beliefs:

$$\mu_{ij}^p \Big)_{i \neq j}$$

ed cumulative reward:

$$r, a_p) \Big] .$$

The integrated f
signments, and t

Initi



Integrated Framework

framework, which uses a *learner* for preference estimation, an *actor* for action selection, and an *environment* for updates and feedback is shown in Figure 2.

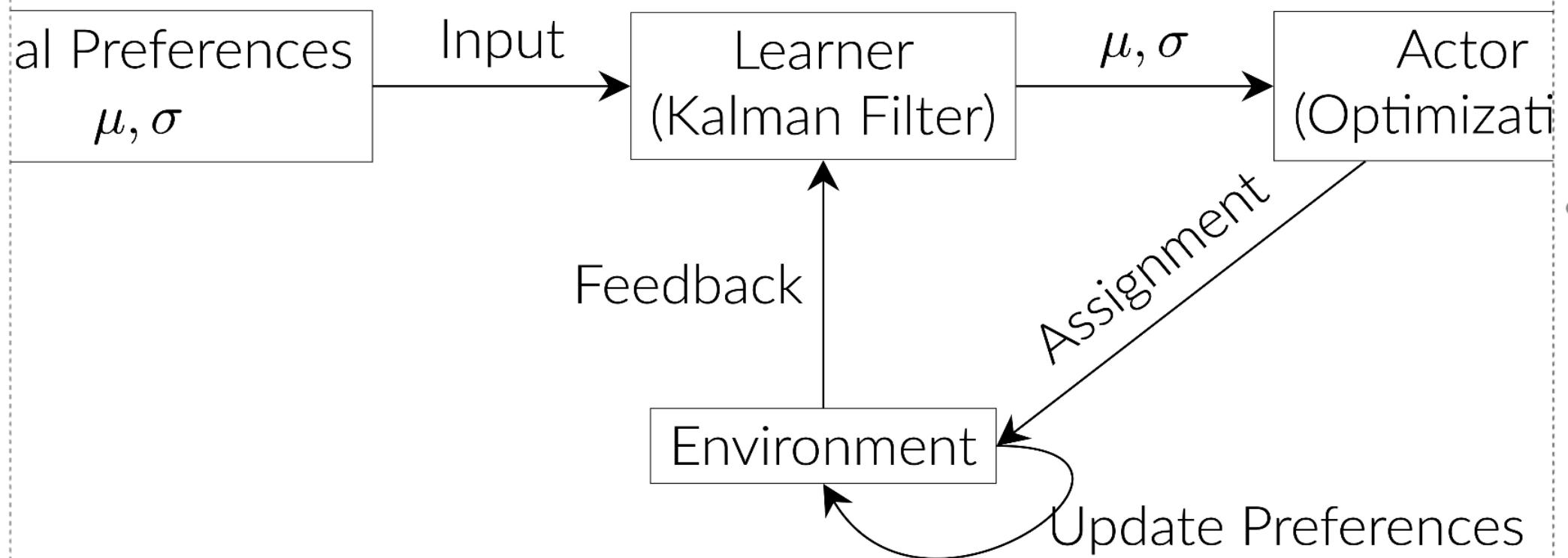


Figure 4. Preference distance between estimate accuracy.

tor for team as-

on)

Discuss

The results demonstrate that incorporating mean preferences) significantly enhance The Kalman filter effectively learns individual as the number of interactions increases.

The optimization-based approach balance complexity limits its scalability to larger settings reinforcement learning or heuristic-base

This framework provides a foundation for in corporate, academic, and sports enviro

d and true preferences. Lower values show improved learning

Solution and Conclusion

ng both **exploration** (via uncertainty) and **exploitation** (via
s team performance compared to random assignments.
idual preferences over time, with convergence improving

ces learning and performance, but computational com-
ngs. Future work will explore **scalable methods** such as
d optimization to handle larger action spaces.

or dynamic team assignment, with potential applications
nments where collaboration dynamics evolve over time.

Department of Industrial and Systems Engineering



Figure 2. Inte



✂

gration of Learner, Actor, and Feedback Loop with Preference Update in the

University of Wisconsin-Madison



Environment.



ISYE 723: Dynamic Programming