

TFG Presentation:

*Toward reducing cumulative bias in
automated decision-making systems
using Multi-Armed Bandits*

Quim De Las Heras Molins (u160402)

quim.delasheras01@estudiant.upf.edu

Juliol 2022

Treball de Fi de Grau

Enginyeria Informàtica

Universitat Pompeu Fabra



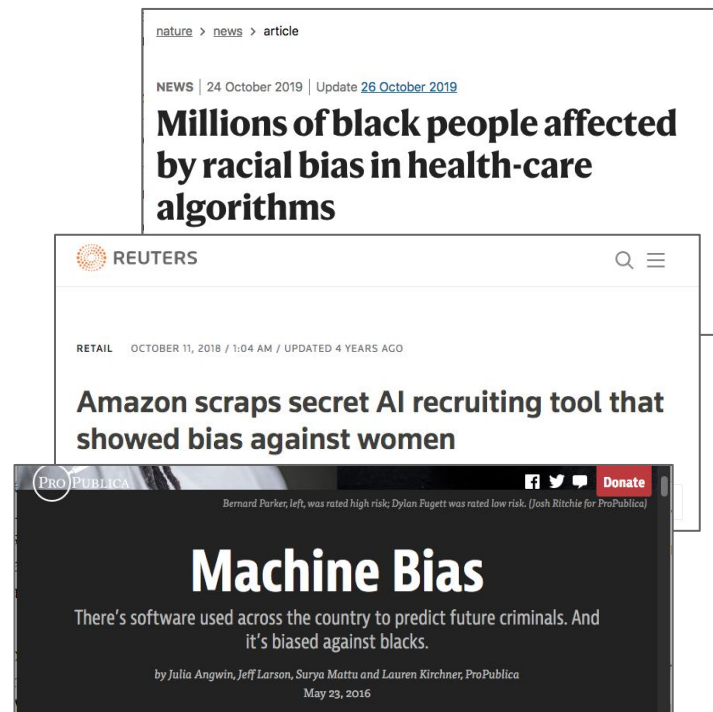
OVERVIEW

- Context and Motivation
- Problem Statement: bias and cumulative bias
- Solution Proposal: BC-UCB
- Results Evaluation
- Conclusions and Future Work

Context and Motivation

[1] [2] [3]

- Automated Decision-Making Systems are on the rise:
 - health care
 - the labor market
 - justice
 - **mortgage lending**
- ↑ Their use becoming more extensive...
- ↑ Higher risk of harmful effects caused by bias.





Context and Motivation

- Types of bias:
 - **Sampling Bias**
 - Negative Set Bias
 - Measurement Bias (Capture, Device or Proxy)
 - **Label Bias**
 - Confounding Bias (Omitted Variable or Proxy)



Context and Motivation

- **Biased** AI models can have harmful **effects** on their **environment**.
- These effects can be **cumulative**.
- Discrimination against **sensitive** groups can be **perpetuated** or **escalated**.
- **Multi-Armed Bandit** problems can be used to model the relationship between an automated-decision making system and its environment.



Context and Motivation

- For a limited amount of rounds T , each round $t \in [1, T]$ the learner will choose one of the k available actions $a_t \in A$ and receive a reward $X_t \in \mathbb{R}$ sampled from distribution P_{a_t} by the environment.

Algorithm 1 - Upper Confidence Bound (UCB)

Input: T and l

Output: $\sum_{t=1}^T X_t$

while $t \leq T$ **do**

 Choose arm $a_t = \operatorname{argmax}_{a \in A} \left[\hat{\mu}_t(a) + l \sqrt{\frac{\log(t)}{N_t(a)}} \right]$

 Observe reward X_t and update $\hat{\mu}_t(a_t)$ and $N_t(a_t)$

end while



Context and Motivation

- The aims of this project are:
 - **Modeling** of an automated decision-making system as a **Multiple Armed Bandit** problem where instances grouped based on **sensitive** variables present arbitrarily different reward distributions that end up being **unfairly optimized** and reveal bias.
 - Study of the **cumulative bias effect** that said model gives rise to when its decisions have an **effect on the environment**, which creates **feedback** loops that expand any existing discriminations **over time** in a manner consistent with real-life scenarios.
 - **Proposing solutions** to MAB problems where **sensitive** features exist within the data which consider not only the exploration-exploitation dilemma but also take into consideration **distances between distributions and arm commitment**.



Problem Statement: bias and cumulative bias

- In our system:
 - Limited amount of resources to be distributed
 - Sensitive features within the data
- Analysis on groups based on intersectional sensitive options.
- Optimizing determinant features should not be on the detriment of sensitive ones.



Problem Statement: bias and cumulative bias

- 717,997 samples from the database by *The Home Mortgage Disclosure Act (HMDA)*.
- US states of Alabama, Arkansas, Georgia, Mississippi, Louisiana and Tennessee, produced during 2020.
- 99 features; a binary target encodes whether the **loan originated or was declined**.

$$k = \prod_{i=1}^s o_i$$

- race = {**W**hite, **B**lack or African American}, ethnicity = {**N**ot **H**ispanic or Latino, **H**ispanic or Latino} and sex = {**M**ale, **F**emale}



Problem Statement: bias and cumulative bias

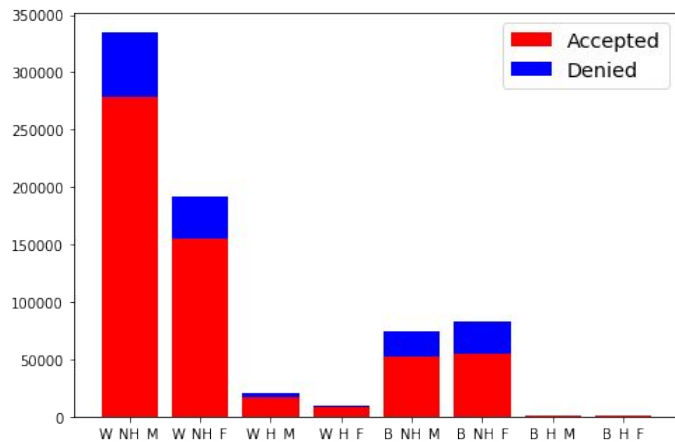


Figure 3.1: Inspection of instance distribution in HMDA database

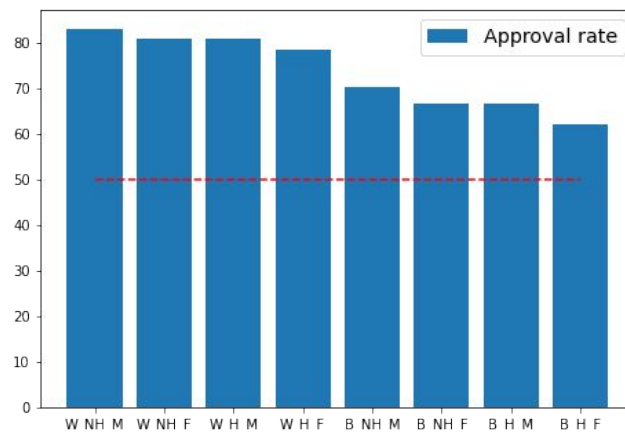


Figure 3.2: Inspection of approval rate spread in HMDA database



Problem Statement: bias and cumulative bias

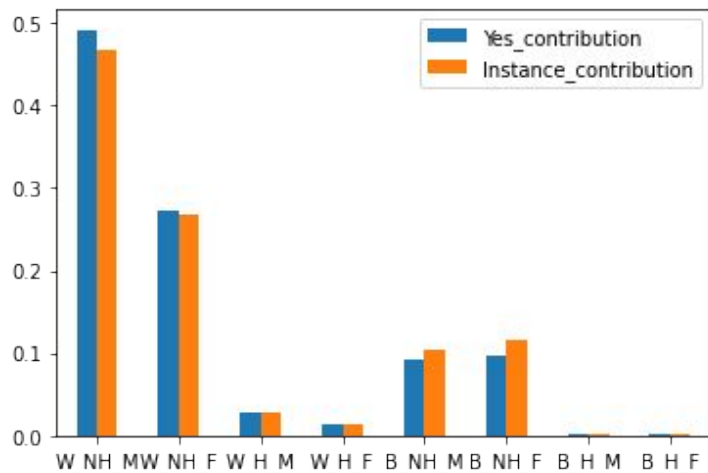


Figure 3.3: Contribution of each demographic to data and to total approvals

$$p_a \in [0.83, 0.807, 0.808, 0.782, 0.702, 0.664, 0.666, 0.619]$$

$$P = (\text{Bernoulli}(p_a) : a \in A)$$

$$w_a \in [0.467, 0.267, 0.028, 0.014, 0.103, 0.116, 0.002, 0.002]$$

$$N = (n_a \sim \text{Binomial}(|C| , w_a) : a \in A)$$

Stationary
Environment



Problem Statement: bias and cumulative bias

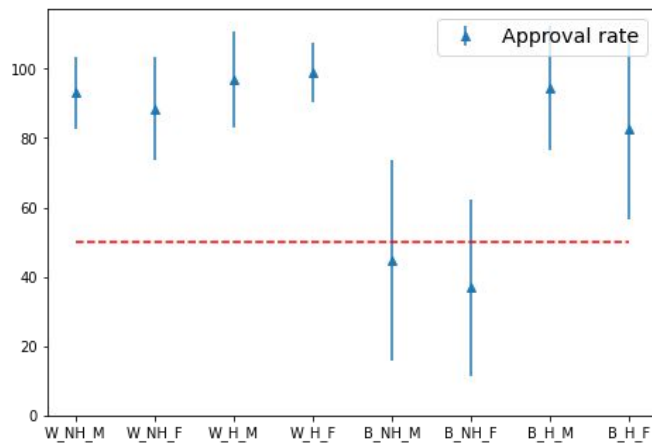


Figure 3.4: Approval rate granted by UCB per each demographic after 100 stationary executions on HMDA-based data

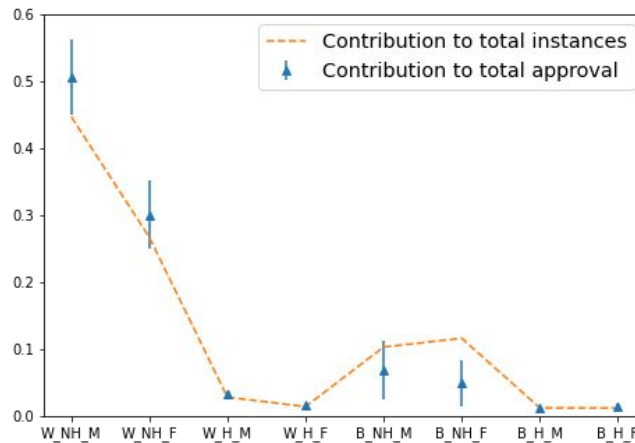


Figure 3.5: Contribution of each demographic to data and to total approvals after 100 stationary executions on HMDA-based data



Problem Statement: bias and cumulative bias

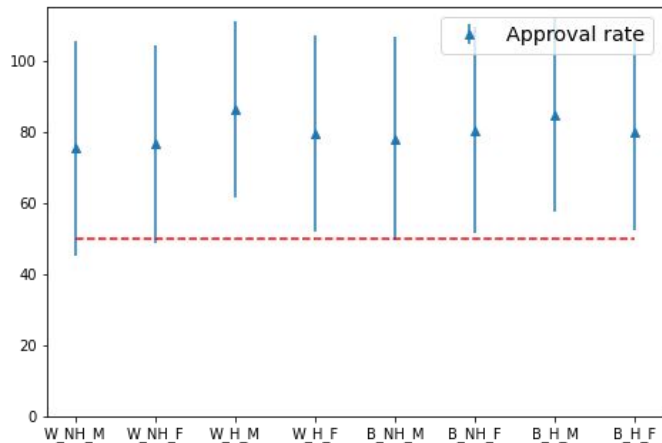


Figure 3.6: Approval rate granted by UCB per each demographic after 100 stationary executions on unbiased data

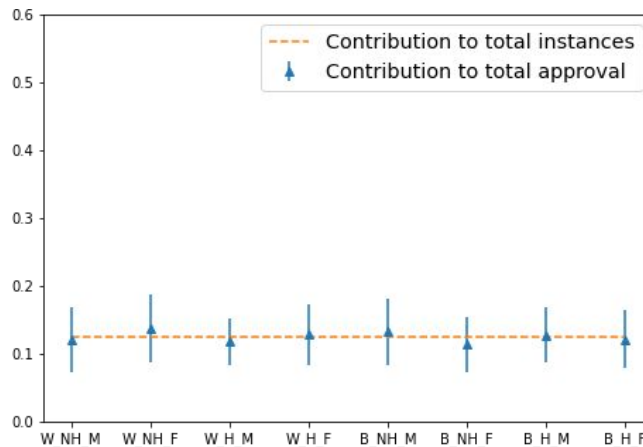
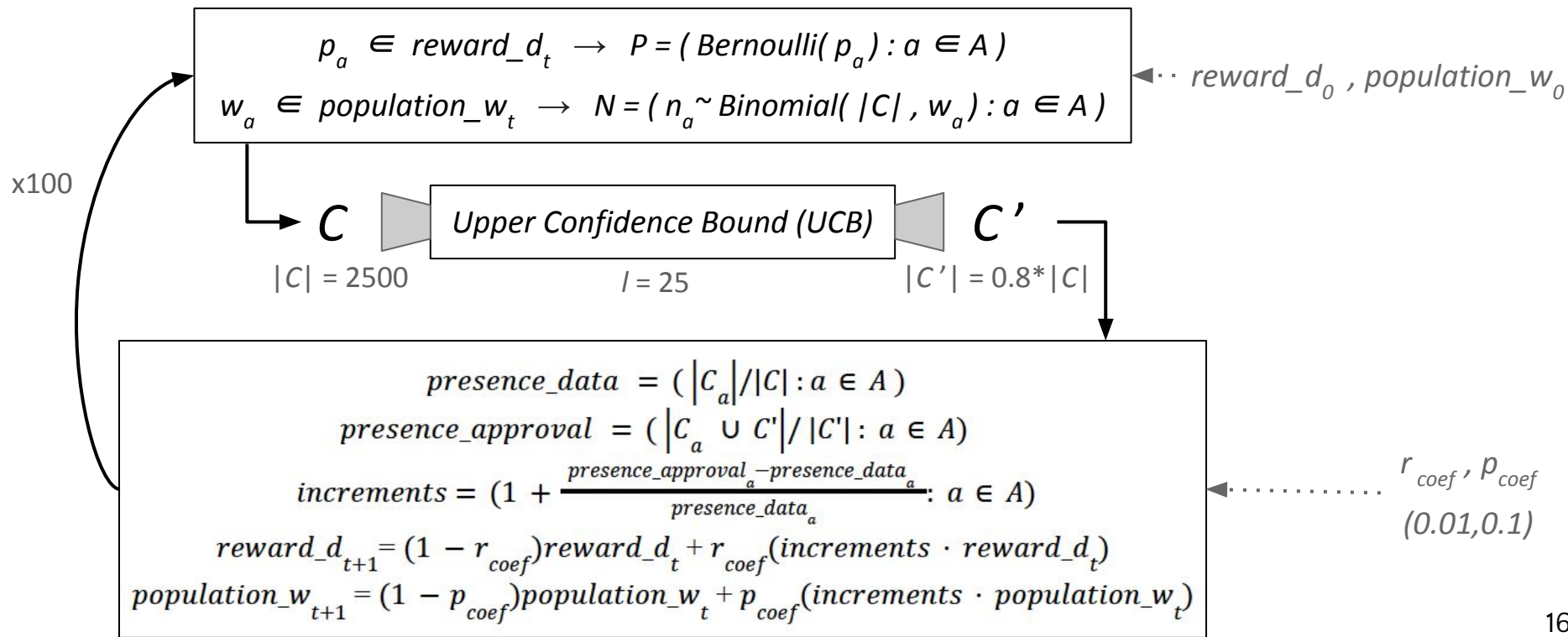


Figure 3.7: Contribution of each demographic to data and to total approvals after 100 stationary executions on unbiased data

Non-Stationary
Environment



Problem Statement: bias and cumulative bias





Problem Statement: bias and cumulative bias

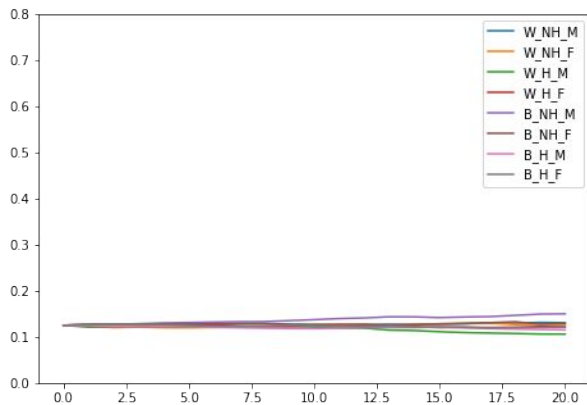


Figure 3.8: Evolution of population weights every 5 epochs in an initially unbiased environment

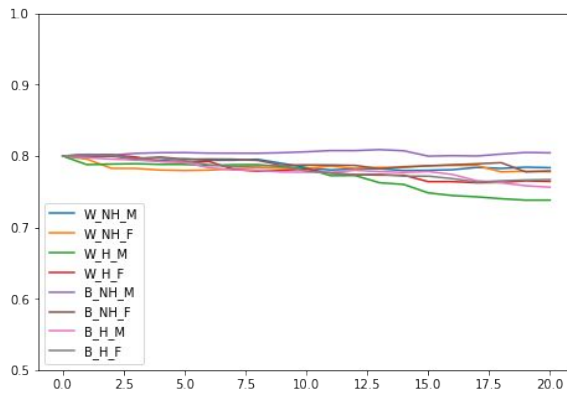


Figure 3.10: Evolution of reward estimates every 5 epochs in an initially unbiased environment

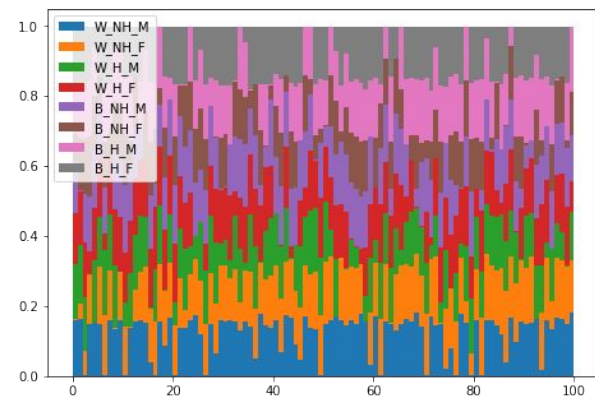


Figure 3.12: Evolution of contribution to total approvals for every epoch in an initially unbiased environment



Problem Statement: bias and cumulative bias

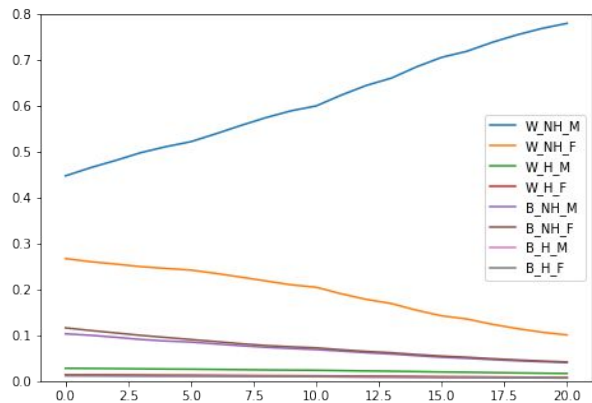


Figure 3.9: Evolution of population weights every 5 epochs in an HMDA-based environment

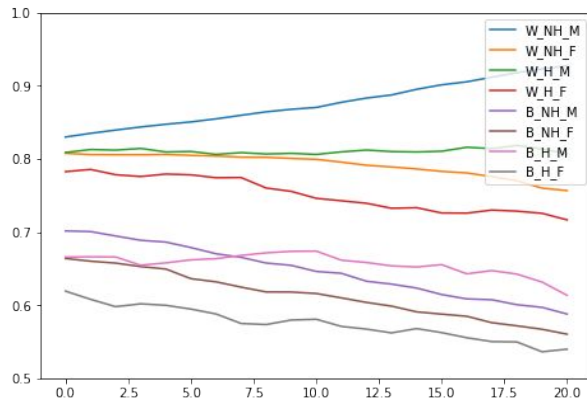


Figure 3.11: Evolution of reward estimates every 5 epochs in an HMDA-based environment

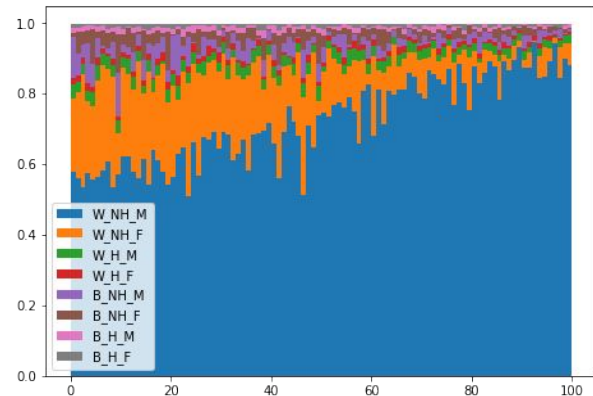


Figure 3.13: Evolution of contribution to total approvals for every epoch in an HMDA-based environment



Solution Proposal: BC-UCB

- Formalization of the four bias metrics:

$$1. \text{ bias_presence} = \sum_{a=1}^k (| \overline{\text{presence_approval}} - \text{presence_data}_a |)$$

$$2. \text{ bias_rates} = \sum_{a=1}^k (| \overline{\text{approval_rates}} - \text{approval_rates}_a |)$$

$$3. \text{ bias_posteriors} = \sum_{a=1}^k (| \overline{\text{posteriors_groups}} - \text{posteriors_groups}_a |)$$

$$4. \text{ bias_rewards} = \sum_{a=1}^k (| \overline{\mu}_t - \hat{\mu}_t(a) |)$$



Solution Proposal: BC-UCB

$$posteriors_groups = (P(c \in C_a | c \in C') : a \in A)$$

$$posteriors_groups = \left(\frac{P(c \in C' | c \in C_a) * P(c \in C_a)}{P(c \in C')} : a \in A \right)$$

$$posteriors_groups = \frac{approval_rates_a * presence_data_a}{0.8} : a \in A)$$



Solution Proposal: BC-UCB

- Bias Constrained UCB (BC-UCB):
 - $bias_metric \in \{bias_presence, bias_rate, bias_posteriors, bias_rewards\}$
 - $bias_term_t = 100 * MinMax(bias_metric(t + 1))$

	Sampling step formula
UCB	$a_t = \underset{a \in A}{argmax} [\hat{\mu}_t(a) + l \sqrt{\frac{\log(t)}{N_t(a)}}]$
BC-UCB	$a_t = \underset{a}{argmax} (1 - \theta) * [\hat{\mu}_t(a) + l \sqrt{\frac{\log(t)}{N_t(a)}}] - \theta * bias_term_t(a)$



Solution Proposal: BC-UCB

Algorithm 2 - Bias Constrained UCB (BC-UCB)

Input: $T, l, \theta, \text{bias_metric}$

Output: $\sum_{t=1}^T X_t, \text{bias_metric}(T)$

while $t \leq T$ **do**

 Compute $\text{bias_term}_t = 100 * \text{MinMax}(\text{bias_metric}(t + 1))$

 Compute $A' = \{a \in A \text{ s.t. } N_t(a) < n_a\}$

 Choose arm $a_t = \underset{a \in A'}{\operatorname{argmax}} ((1 - \theta) * \left[\hat{\mu}_t(a) + l \sqrt{\frac{\log(t)}{N_t(a)}} \right] - \theta * \text{bias_term}_t(a))$

 Observe reward X_t and update $\hat{\mu}_t(a_t)$ and $N_t(a_t)$

end while



Solution Proposal: BC-UCB

- ◉ Navigating the tradeoff between final total reward and bias:

$$Frobenius(v) = v / \sqrt{(\sum_{i=1} v_i^2)}$$

$$\varphi(frew, fbias) = Frobenius(frew) - \beta * Frobenius(fbias)$$

$$\theta_{bias_metric}^* = \underset{(\theta)}{argmax} (\varphi(BC-UCB(\theta, bias_metric)))$$



Results Evaluation

https://github.com/quimHM/QHM_TFG_repository (Last commit: 22nd June)

```
main - QHM_TFG_repository / MAB_sim.ipynb

quimHM Polishing touches Latest commit 196806 12 days ago History

At 1 contributor

2895 Lines (2895 sloc) 935 KB

In [ ]:
#Some of the code and comments based on the MAB posts by Steve Roberts:
#https://towardsdatascience.com/the-upper-confidence-bound-ucb-bandit-algorithm-c85c2b4c13f

In [1]:
# Import modules
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import math
import random
import matplotlib inline
from scipy.stats import binom

In [39]:
from google.colab import drive # Import drive from google colab

ROOT = "/content/drive" # default location for the drive
print(ROOT)
drive.mount(ROOT)

df = pd.read_csv("/content/drive/MyDrive/TFG/hmda/state_AL-GA_actions_taken_1-3.csv")
print(len(df))
extra1 = pd.read_csv("/content/drive/MyDrive/TFG/hmda/state_AR-MS_actions_taken_1-3.csv")
print(len(extra1))
extra2 = pd.read_csv("/content/drive/MyDrive/TFG/hmda/state_TN-LA_actions_taken_1-3.csv")
print(len(extra2))

subsample = pd.concat([df, extra1, extra2])

print(len(subsample.loc[subsample["action_taken"]==1]), len(subsample.loc[subsample["action_taken"]==3]))
subsample = subsample.loc[subsample["action_taken"]].isin([1,3])
subsample["action_taken"] = subsample["action_taken"].replace([3], [0])
print("action:", subsample["action_taken"].unique())
print(len(subsample.loc[subsample["action_taken"]==1]), len(subsample.loc[subsample["action_taken"]==0]))

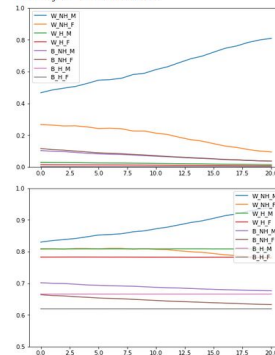
print("race pre:", subsample["derived_race"].unique())
subsample = subsample.loc[subsample["derived_race"]].isin(["White", "Black or African American"])
subsample["derived_race"] = subsample["derived_race"].replace(["White", "Black or African American"], [0,1])
print("race:", subsample["derived_race"].unique())

print("ethnicity pre:", subsample["derived_ethnicity"].unique())
subsample = subsample.loc[subsample["derived_ethnicity"]].isin(["Not Hispanic or Latino", "Hispanic or Latino"])
subsample["derived_ethnicity"] = subsample["derived_ethnicity"].replace(["Not Hispanic or Latino", "Hispanic or Latino"], [0,1])
print("ethnicity:", subsample["derived_ethnicity"].unique())
```

```
In [32]:
r = nonstatic_sol(100, rewards_bias, 0.01, instances_bias, 0.1, 0)

8
Reward distribution: [0.83, 0.8877, 0.8888, 0.7826, 0.7815, 0.6639, 0.6658, 0.6192]
Final estimates: [0.8173, 0.8228, 0.4, 0.8, 0.7243, 0.6429, 0.8, 0.3333]
Presence of each in batch (%): [47.72, 25.44, 2.88, 1.2, 18.44, 15.88, 0.2, 0.24]
Relative approval rate (%): [180.8, 95.84, 188.8, 61.09, 9.43, 26.8, 50.8]
Percentage over total selected (%): [59.62, 28.94, 0.25, 1.5, 8.85, 1.4, 0.45, 0.15]
Total reward: 161380 | and presence bias: 0.3143992883998887
Rate bias: 2.4635158737379366
Posterior bias: 1.2713643178418796
Rewards bias: 2.4598857946355877
Final regret: 2.3983846475762176

99
Reward distribution: [0.9388, 0.7786, 0.8076, 0.7819, 0.6762, 0.6328, 0.6658, 0.6191]
Final estimates: [0.836, 0.7, 0.625, 0.8, 0.8, 0.6, 1.8, 0.8]
Presence of each in batch (%): [88.4, 8.48, 3.32, 8.6, 4.8, 3.92, 0.16, 0.12]
Relative approval rate (%): [198.01, 4.22, 24.24, 6.87, 1.8, 5.1, 188.8, 33.33]
Percentage over total selected (%): [106.45, 0.5, 8.4, 0.85, 0.85, 0.25, 0.2, 0.85]
Total reward: 186488 | and presence bias: 0.3623132433783187
Rate bias: 2.597329882048879
Posterior bias: 1.7191484287851076
Rewards bias: 4.138959398882842
Final regret: 0.722514489513265
```



Stationary
Environment

APPENDICES

Appendix A: Application in a Stationary Environment

a) Presence Bias

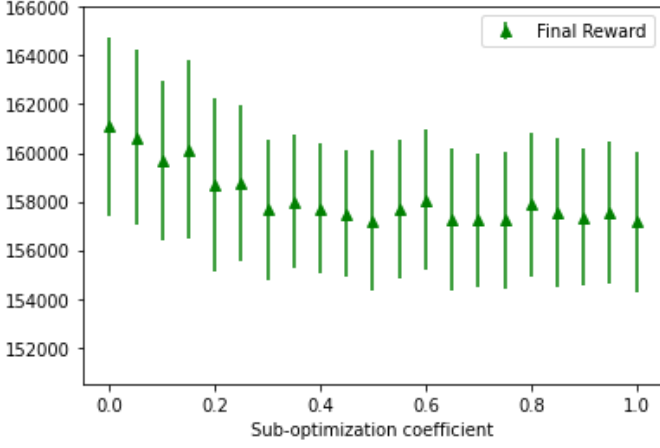


Figure A.1: Final Reward spread depending on θ when $bias_metric = bias_presence$

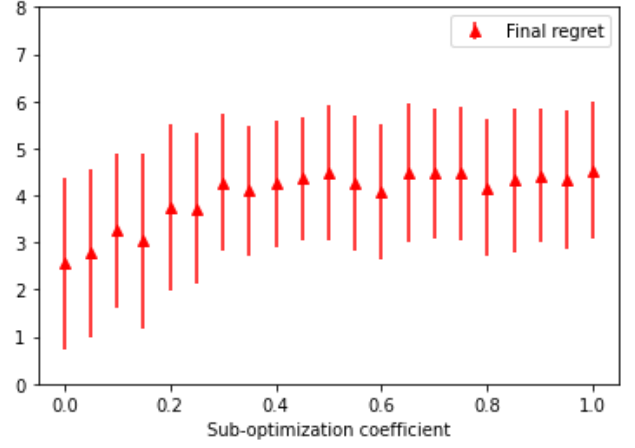


Figure A.2: Final Regret spread depending on θ when $bias_metric = bias_presence$

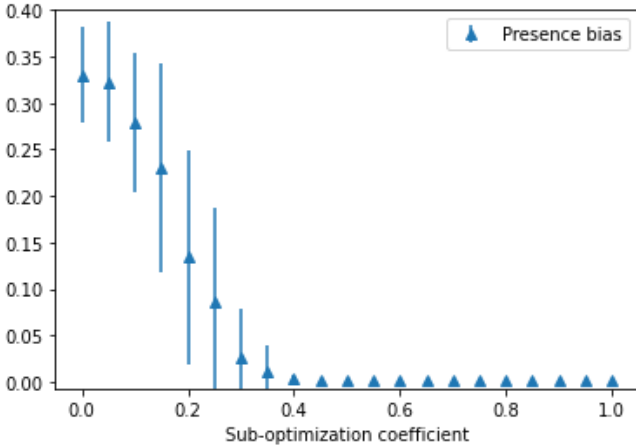


Figure A.3: Final Presence bias spread depending on θ when $bias_metric = bias_presence$

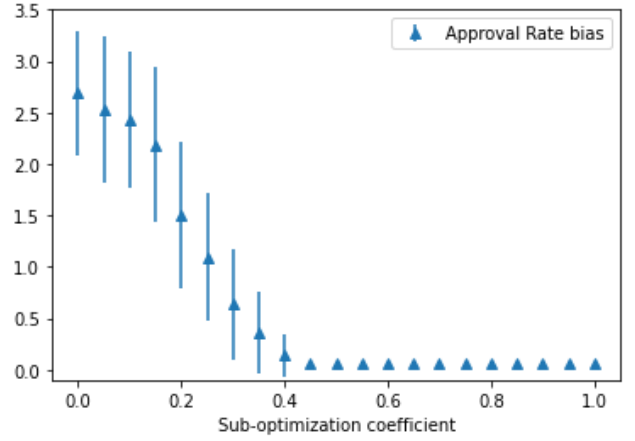


Figure A.4: Final Approval Rate bias spread depending on θ when $bias_metric = bias_presence$

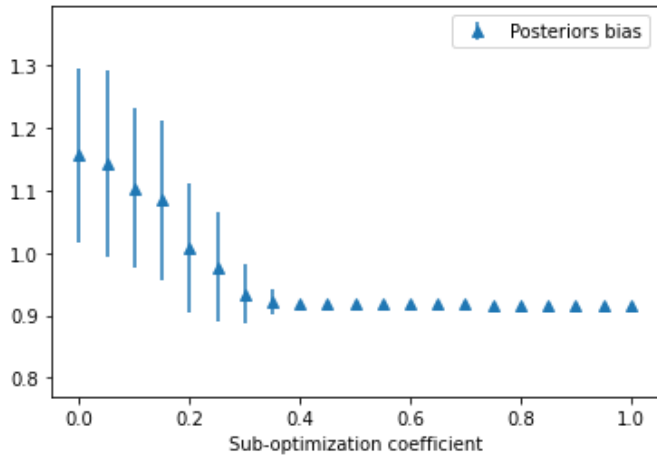


Figure A.5: Final Posteriors bias spread depending on θ when $bias_metric = bias_presence$

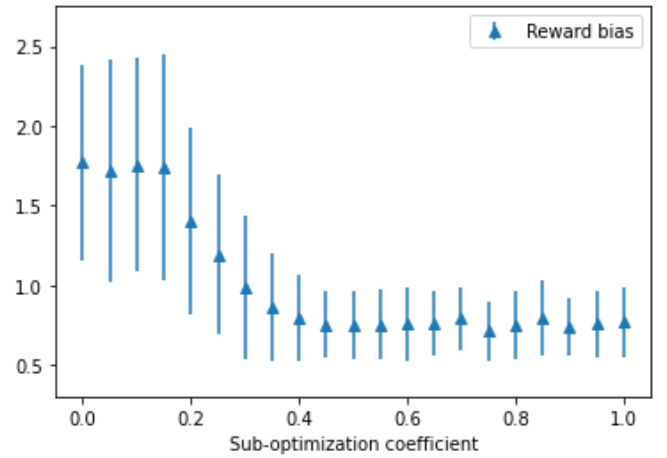


Figure A.6: Final Rewards bias spread depending on θ when $bias_metric = bias_presence$

b) Approval Rate Bias

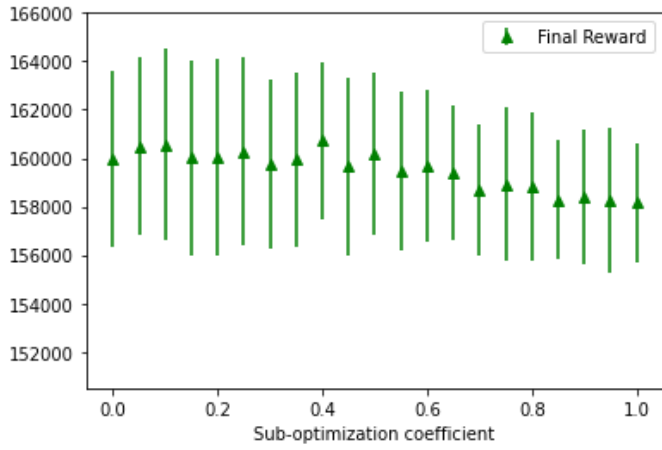


Figure A.7: Final Reward spread depending on θ when $bias_metric = bias_rates$

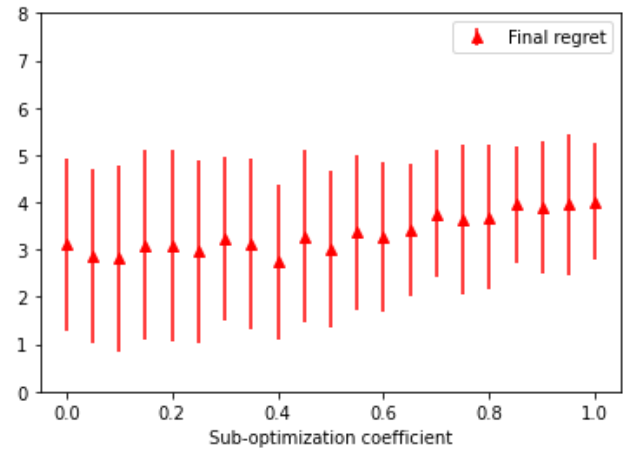


Figure A.8: Final Regret spread depending on θ when $bias_metric = bias_rates$

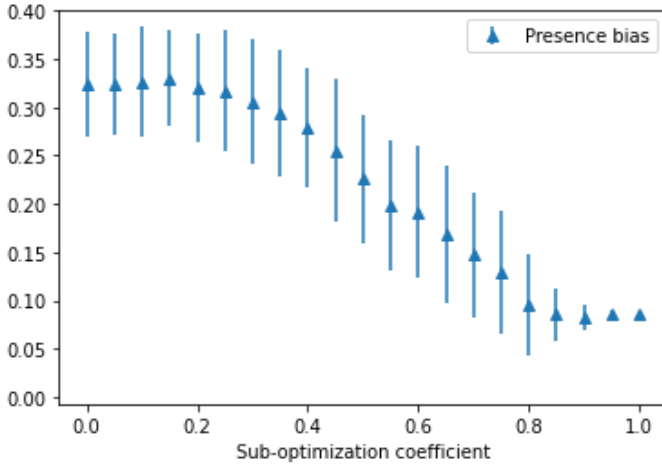


Figure A.9: Final Presence bias spread depending on θ when $bias_metric = bias_rates$

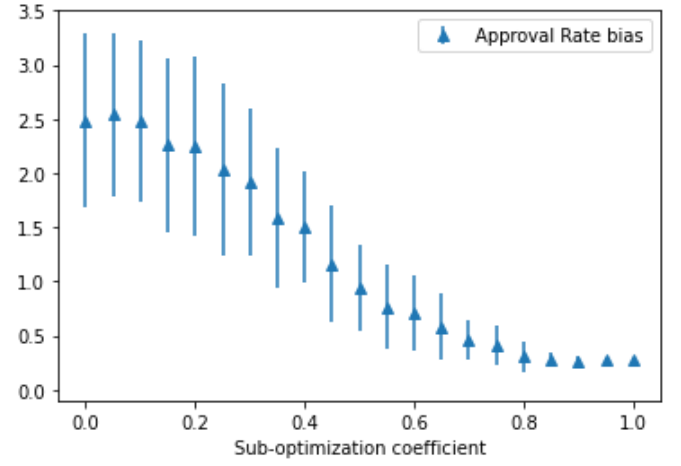


Figure A.10: Final Approval Rate bias spread depending on θ when $bias_metric = bias_rates$

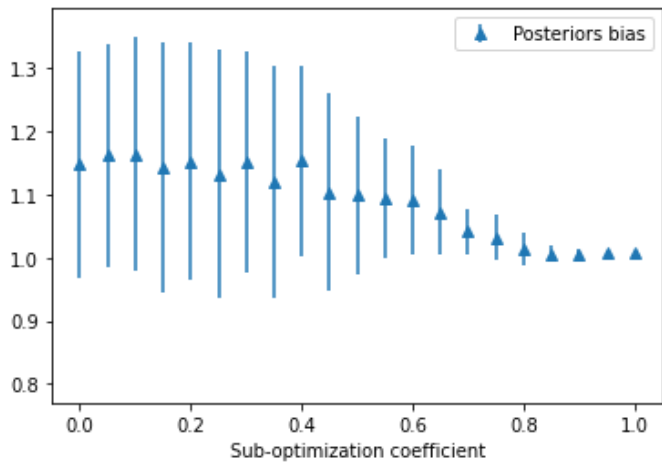


Figure A.11: Final Posteriors bias spread depending on θ when $bias_metric = bias_rates$

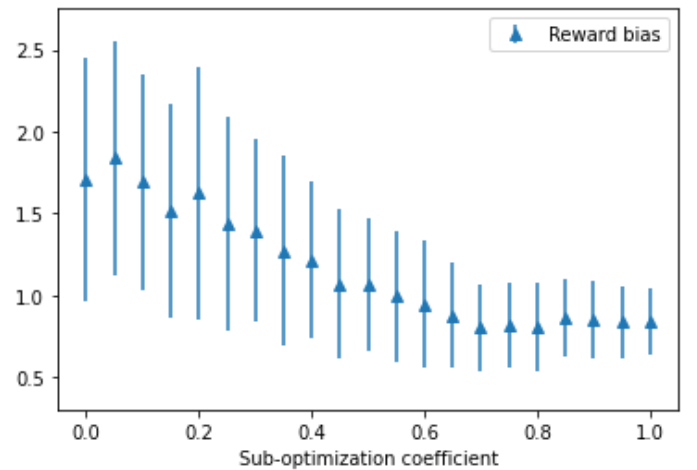


Figure A.12: Final Rewards bias spread depending on θ when $bias_metric = bias_rates$

c) Posteriors Bias

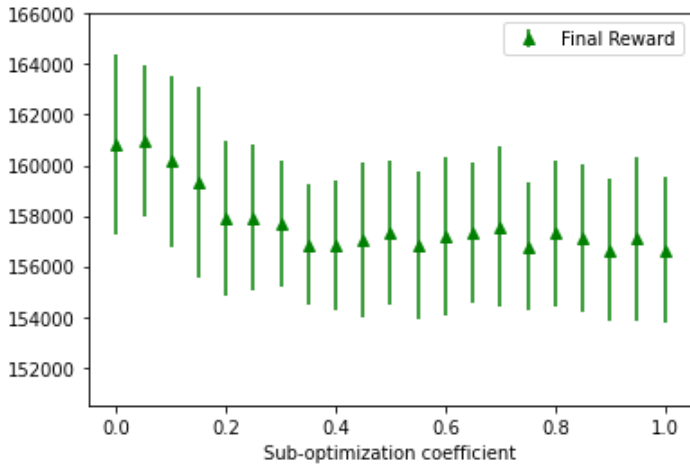


Figure A.13: Final Reward spread depending on θ when $bias_metric = bias_posteriors$

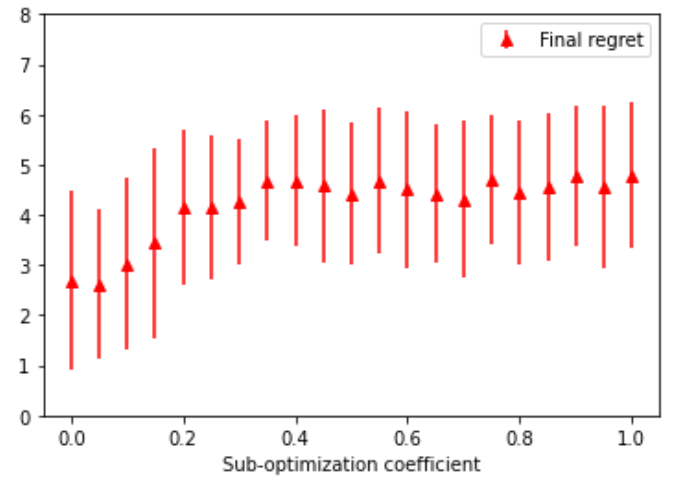


Figure A.14: Final Regret spread depending on θ when $bias_metric = bias_posteriors$

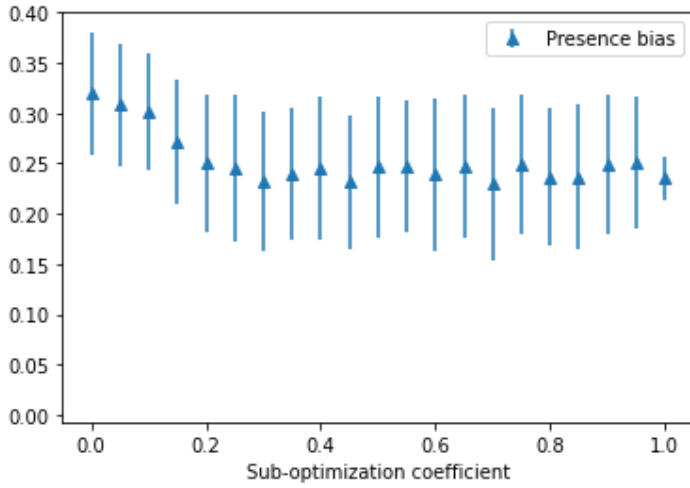


Figure A.15: Final Presence bias spread depending on θ when $bias_metric = bias_posteriors$

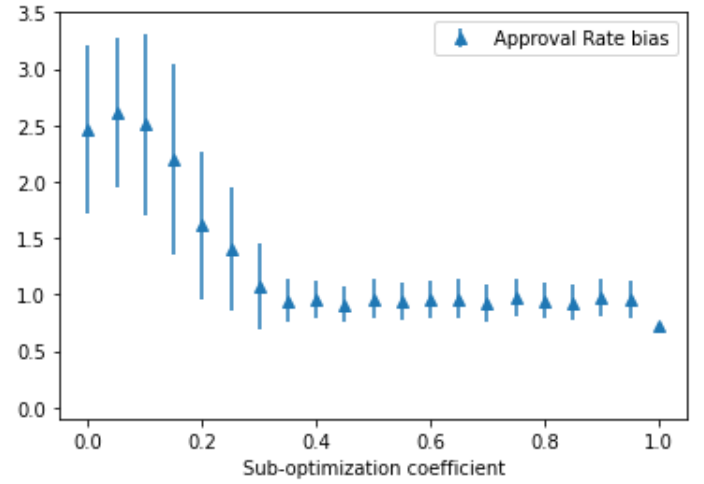


Figure A.16: Final Approval Rate bias spread depending on θ when $bias_metric = bias_posteriors$

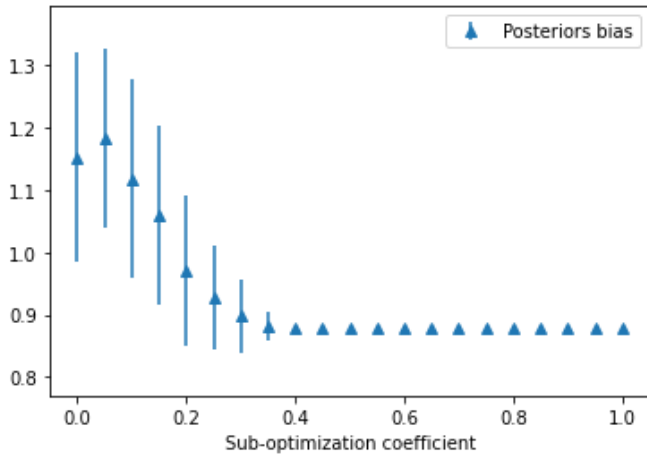


Figure A.17: Final Posteriors bias spread depending on θ when $bias_metric = bias_posteriors$

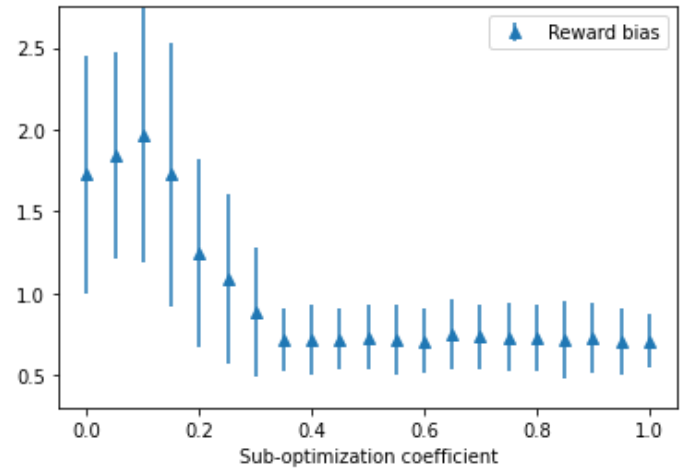


Figure A.18: Final Rewards bias spread depending on θ when $bias_metric = bias_posteriors$

d) Rewards Bias

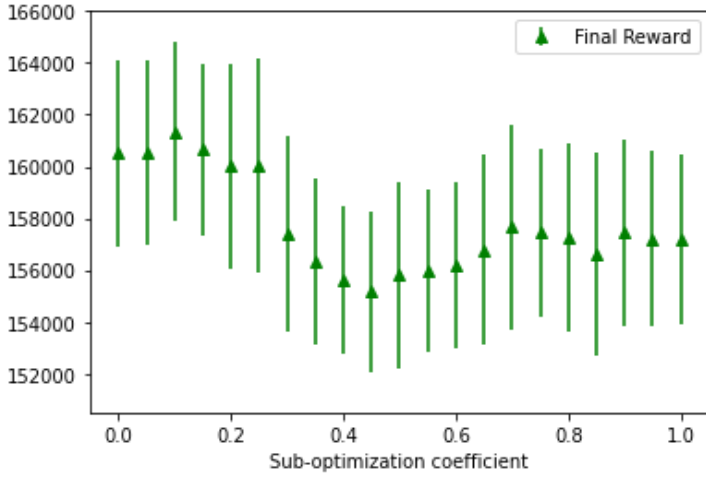


Figure A.19: Final Reward spread depending on θ when $bias_metric = bias_rewards$

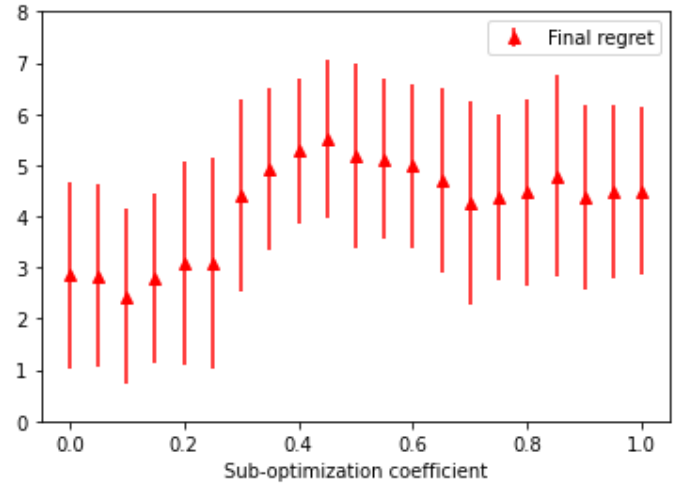


Figure A.20: Final Regret spread depending on θ when $bias_metric = bias_rewards$

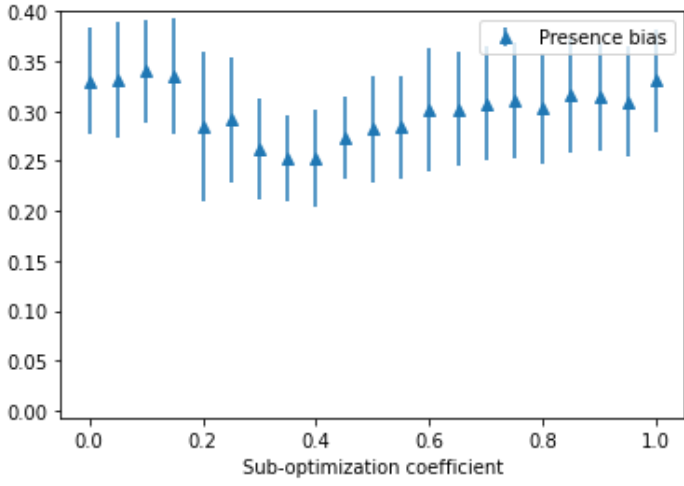


Figure A.21: Final Presence bias spread depending on θ when $bias_metric = bias_rewards$

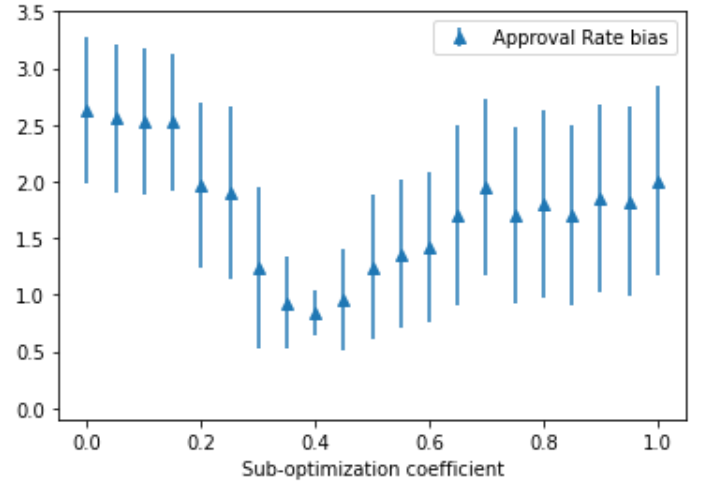


Figure A.22: Final Approval Rate bias spread depending on θ when $bias_metric = bias_rewards$

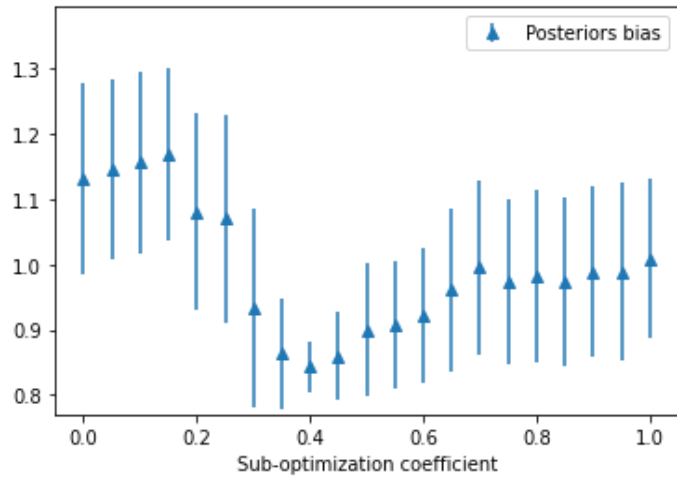


Figure A.23: Final Posteriors bias spread depending on θ when $bias_metric = bias_rewards$

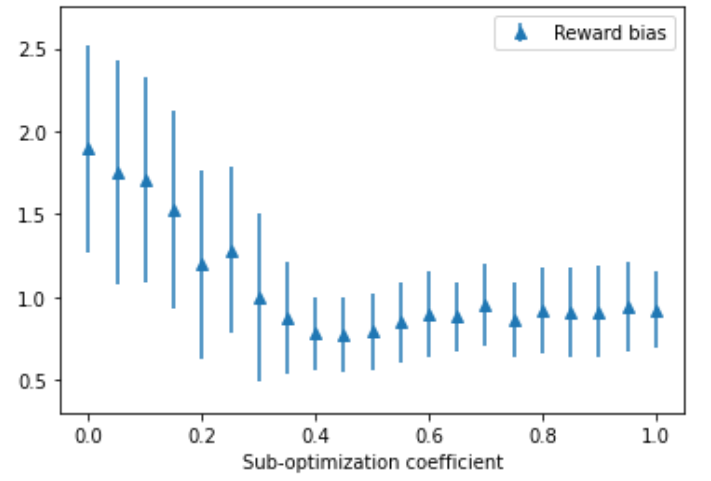


Figure A.24: Final Rewards bias spread depending on θ when $bias_metric = bias_rewards$

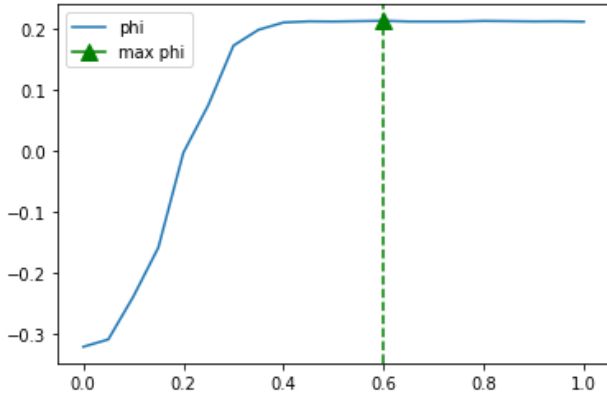


Figure 5.1: ϕ when $\text{bias_metric} = \text{bias_presence}$

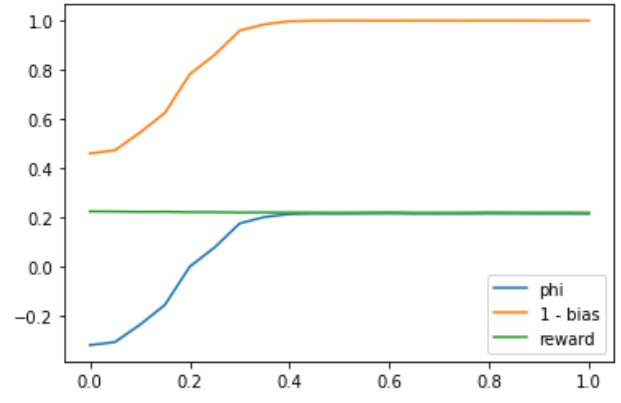


Figure 5.2: Normalized values of A.1 and A.3

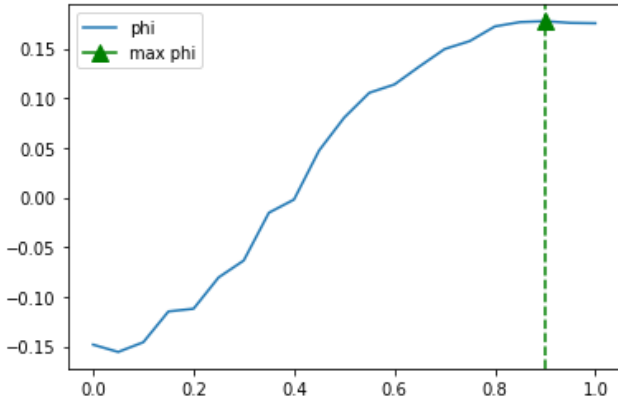


Figure 5.3: ϕ when $\text{bias_metric} = \text{bias_rate}$

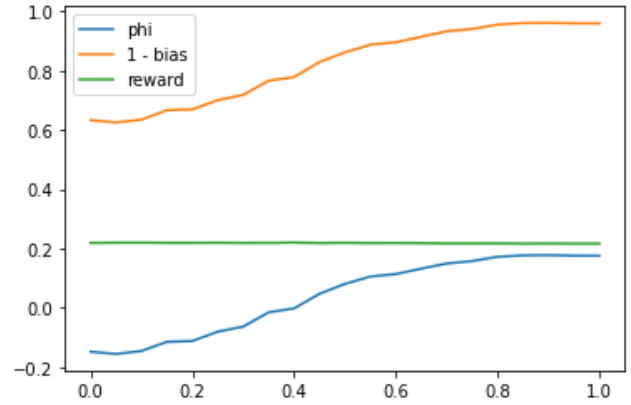


Figure 5.4: Normalized values of A.7 and A.10

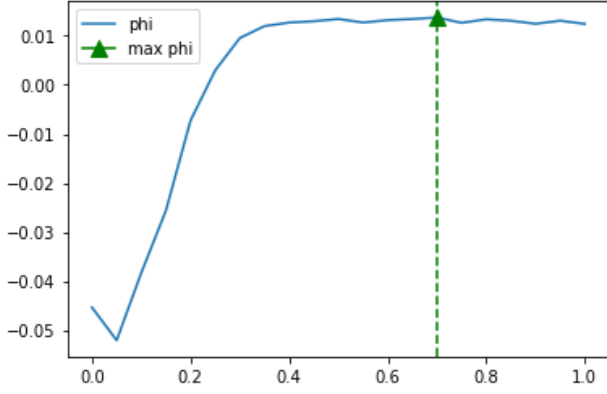


Figure 5.5: ϕ when $\text{bias_metric} = \text{bias_posteriors}$

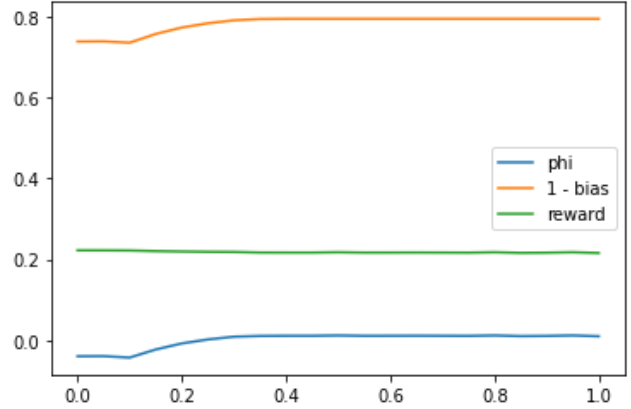


Figure 5.6: Normalized values of A.13 and A.17

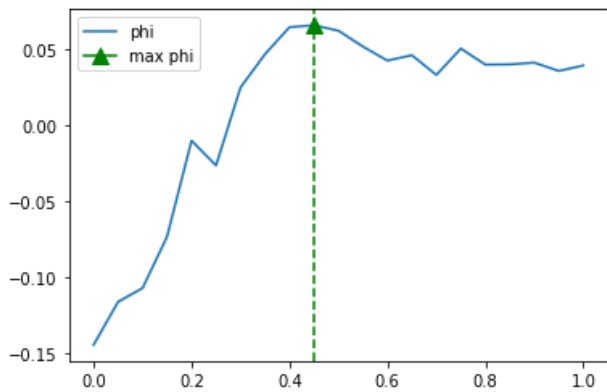


Figure 5.7: ϕ when $\text{bias_metric} = \text{bias_rewards}$

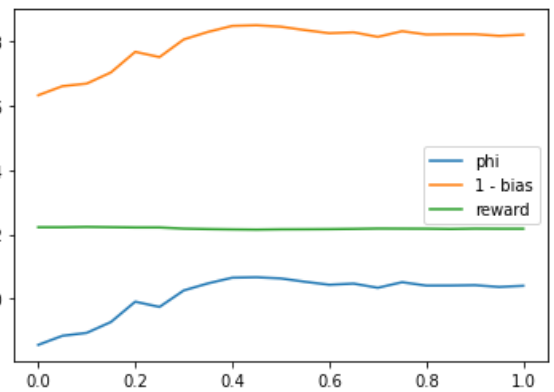


Figure 5.8: Normalized values of A.19 and A.24



Results Evaluation

Presence bias:

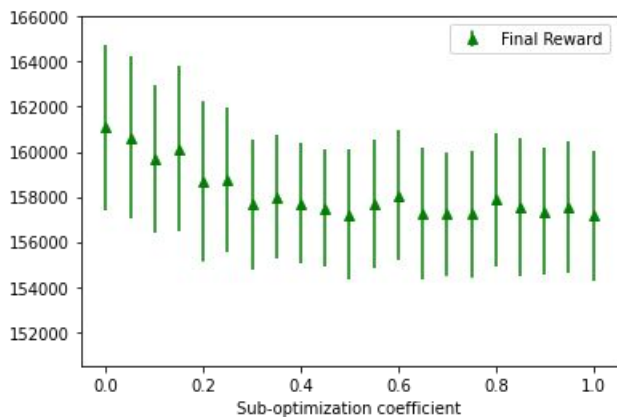


Figure A.1: Final Reward spread depending on θ when $bias_metric = bias_presence$

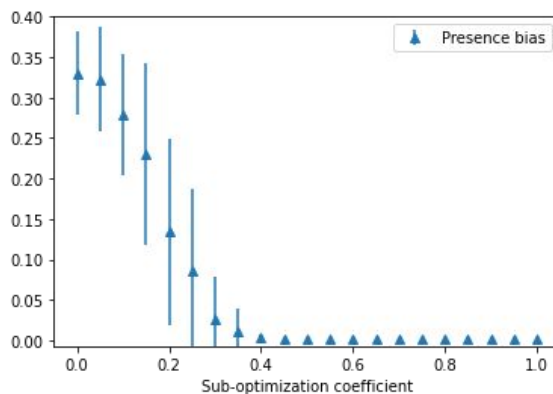


Figure A.3: Final Presence bias spread depending on θ when $bias_metric = bias_presence$

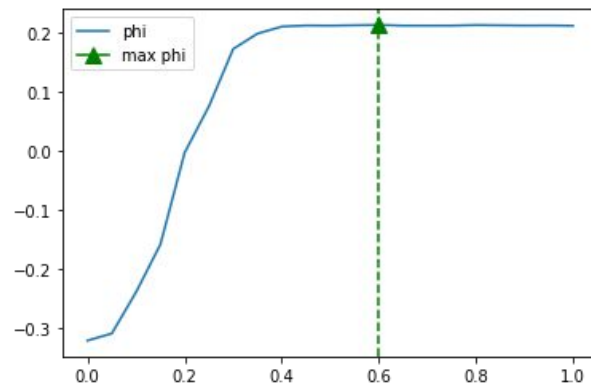


Figure 5.1: ϕ when $bias_metric = bias_presence$



Results Evaluation

Approval Rate bias:

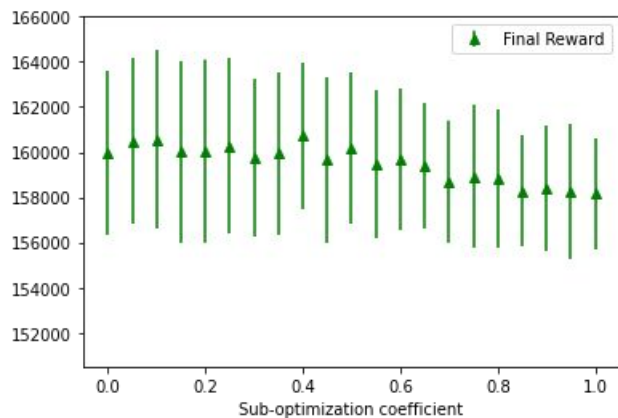


Figure A.7: Final Reward spread depending on θ when $bias_metric = bias_rates$

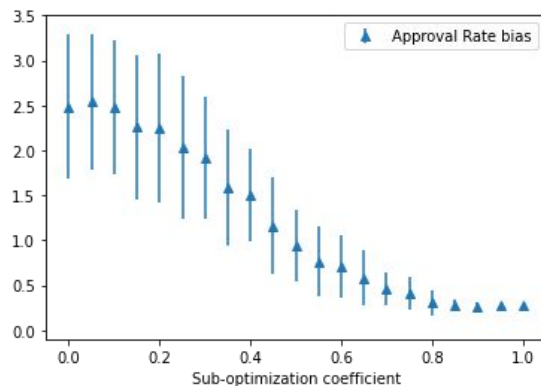


Figure A.10: Final Approval Rate bias spread depending on θ when $bias_metric = bias_rates$

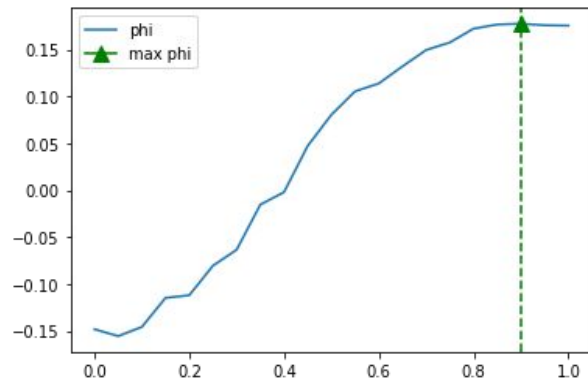


Figure 5.3: ϕ when $bias_metric = bias_rate$

Results Evaluation

⦿ Posteriors bias:

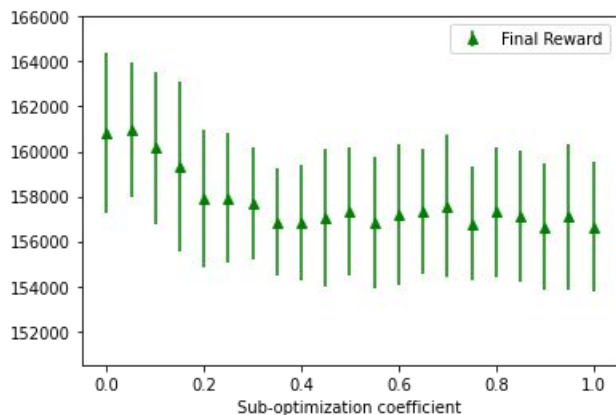


Figure A.13: Final Reward spread depending on θ when $bias_metric = bias_posteriors$

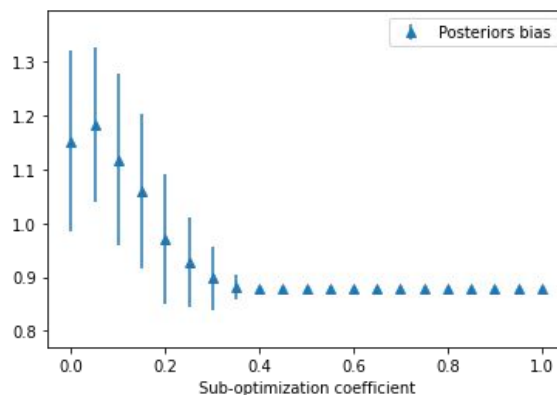


Figure A.17: Final Posteriors bias spread depending on θ when $bias_metric = bias_posteriors$

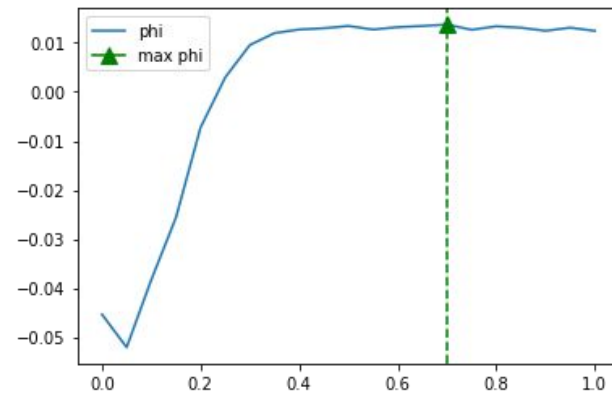


Figure 5.5: ϕ when $bias_metric = bias_posteriors$

Results Evaluation

⊙ Rewards bias:

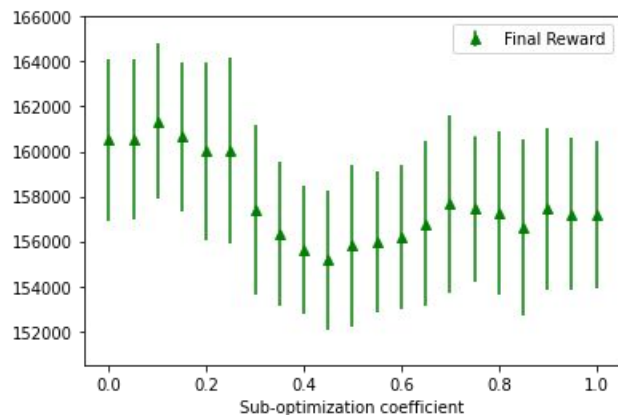


Figure A.19: Final Reward spread depending on θ when $bias_metric = bias_rewards$

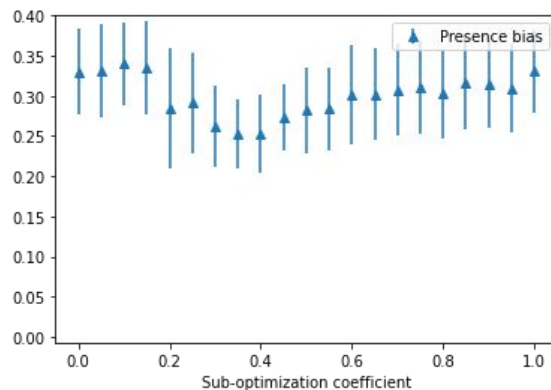


Figure A.21: Final Rewards bias spread depending on θ when $bias_metric = bias_rewards$

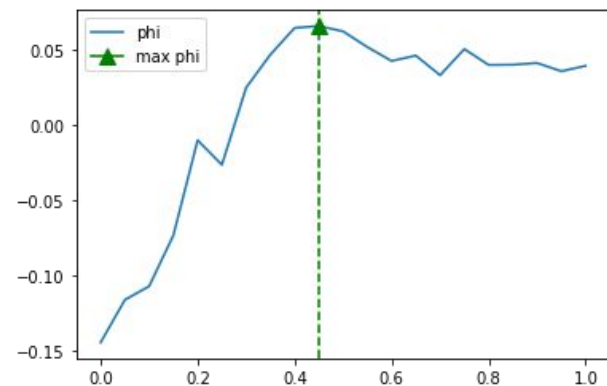


Figure 5.7: ϕ when $bias_metric = bias_rewards$

In [70]:

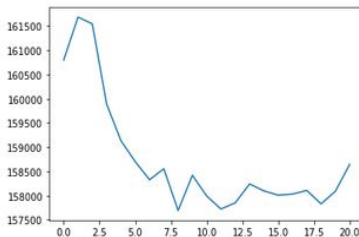
```
print(yr)
print([round(i,3) for i in yr/np.sum(yr)])
print([round(i,3) for i in yr/np.linalg.norm(yr)])
print([round(i,3) for i in minmax(yr)])
```

```
[160800.0, 161686.0, 161546.0, 159890.0, 159134.0, 158702.0, 158326.0, 158554.0, 157692.0, 158418.0, 157990.0, 157722.0, 157850.0, 15824
0.0, 158096.0, 158008.0, 158030.0, 158106.0, 157828.0, 158088.0, 158644.0]
[0.048, 0.049, 0.048, 0.048, 0.048, 0.048, 0.047, 0.048, 0.047, 0.048, 0.047, 0.047, 0.047, 0.047, 0.047, 0.047, 0.047, 0.0
47, 0.048]
[0.221, 0.222, 0.222, 0.22, 0.219, 0.218, 0.218, 0.218, 0.217, 0.218, 0.217, 0.217, 0.217, 0.218, 0.217, 0.217, 0.217, 0.217, 0.21
7, 0.218]
[0.778, 1.0, 0.965, 0.55, 0.361, 0.253, 0.159, 0.216, 0.0, 0.182, 0.075, 0.008, 0.04, 0.137, 0.101, 0.079, 0.085, 0.104, 0.034, 0.099, 0.
238]
```

In [71]:

```
plt.plot(yr)
```

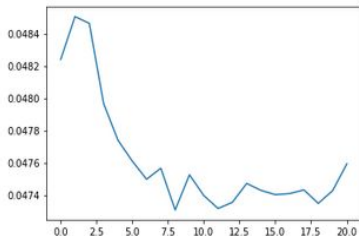
Out[71]:



In [72]:

```
plt.plot(yr/np.sum(yr))
```

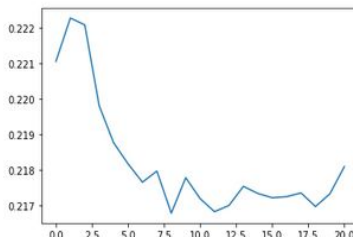
Out[72]:



In [73]:

```
plt.plot(yr/np.linalg.norm(yr))
```

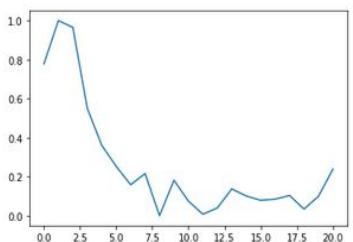
Out[73]:



In [74]:

```
plt.plot(minmax(yr))
```

Out[74]:



Non-Stationary
Environment

Appendix B: Application in a Non-Stationary Environment

a) Presence Bias ($\theta=0.6$)

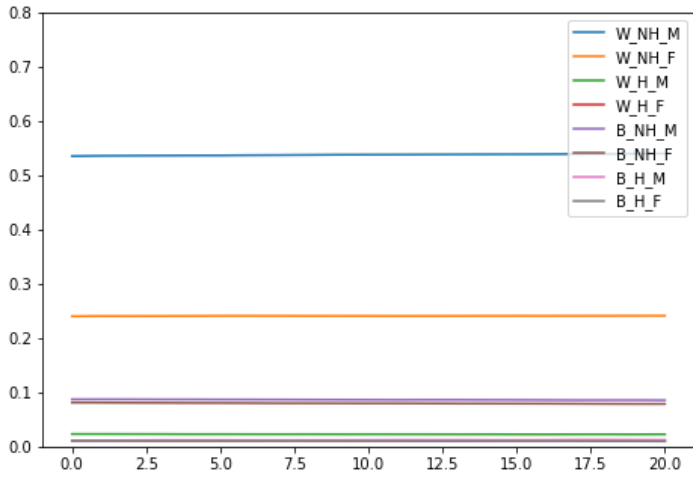


Figure B.1: Evolution of population weights every 5 epochs in a Presence-constrained environment

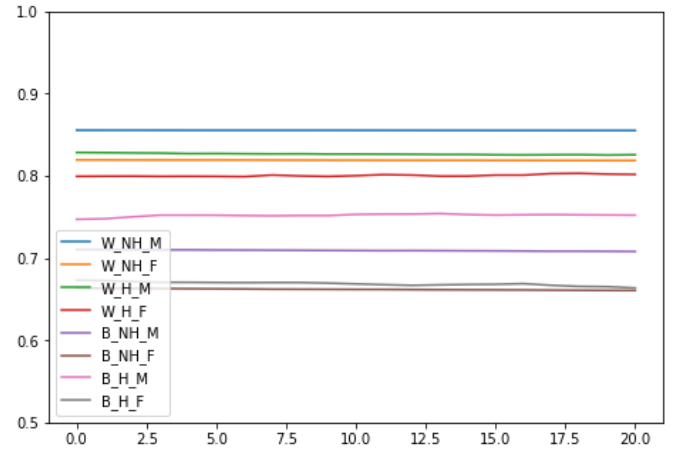


Figure B.2: Evolution of reward estimates every 5 epochs in a Presence-constrained environment

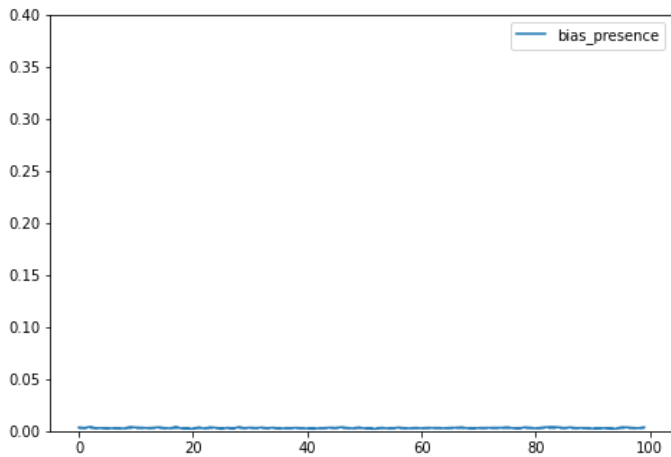


Figure B.3: Evolution of final Presence bias for every epoch in a Presence-constrained environment

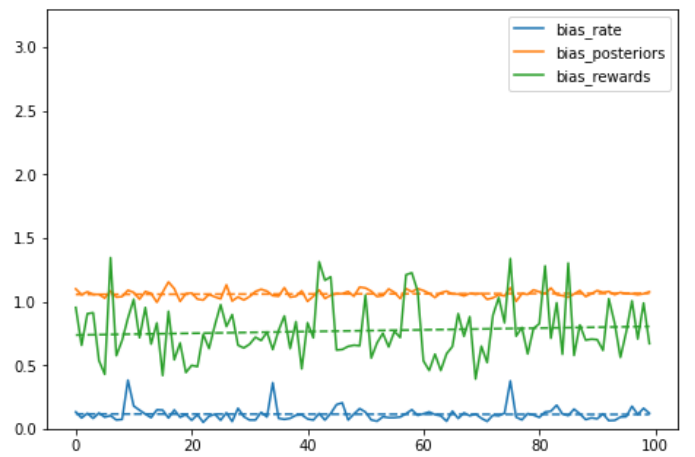


Figure B.4: Evolution of final Rate, Posteriors and Rewards biases in a Presence-constrained environment

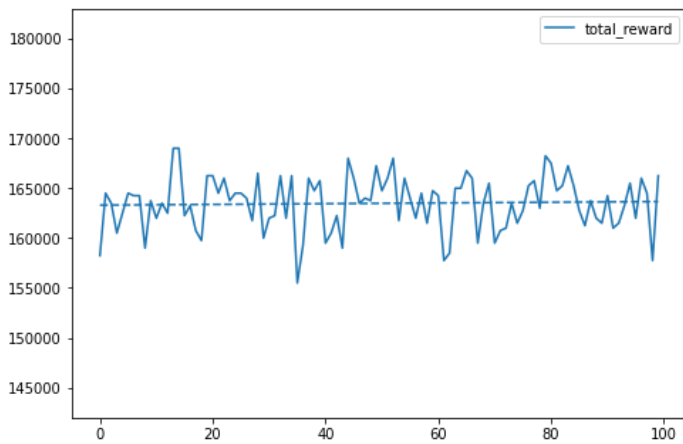


Figure B.5: Evolution of final Total Reward for every epoch in a Presence-constrained environment

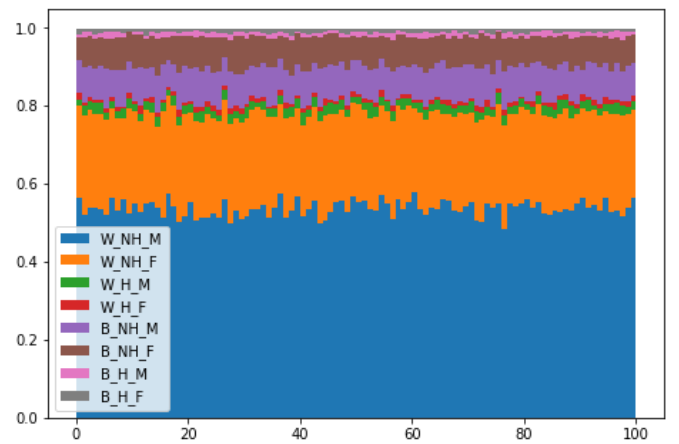


Figure B.6: Evolution of contribution to total approvals for every epoch in a Presence-constrained environment

b) Approval Rate Bias ($\theta=0.9$)

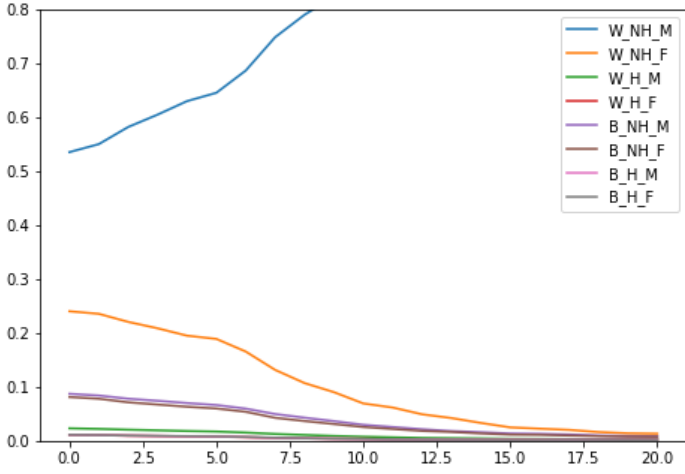


Figure B.7: Evolution of population weights every 5 epochs in an Approval Rate-constrained environment

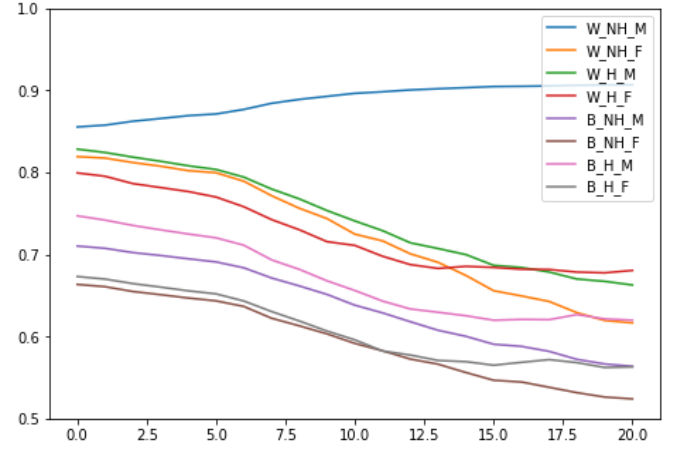


Figure B.8: Evolution of reward estimates every 5 epochs in a Approval Rate-constrained environment

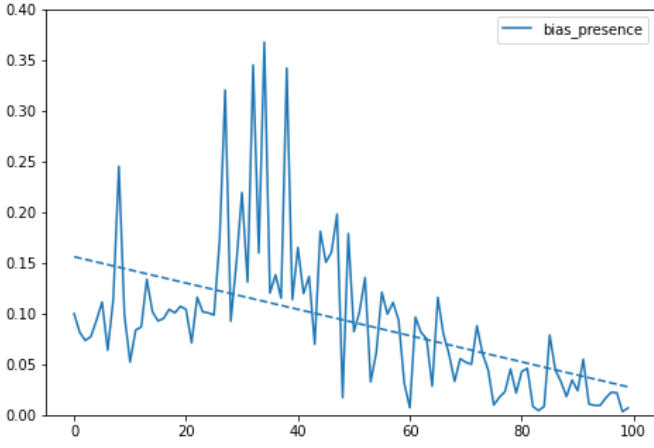


Figure B.9: Evolution of final Presence bias for every epoch in a Approval Rate-constrained environment

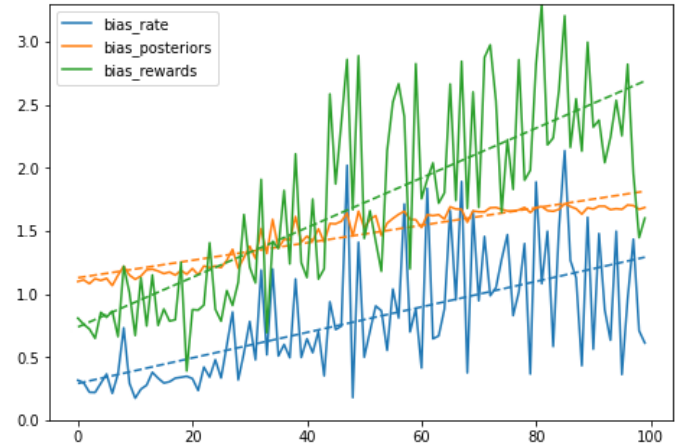


Figure B.10: Evolution of final Rate, Posteriors and Rewards biases in a Approval Rate-constrained environment

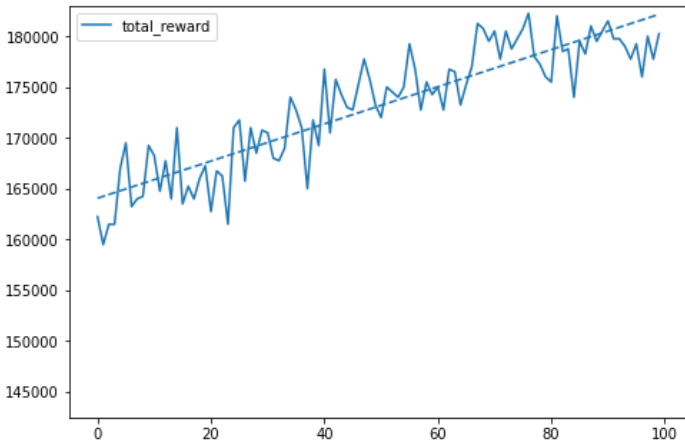


Figure B.11: Evolution of final Total Reward for every epoch in a Approval Rate-constrained environment

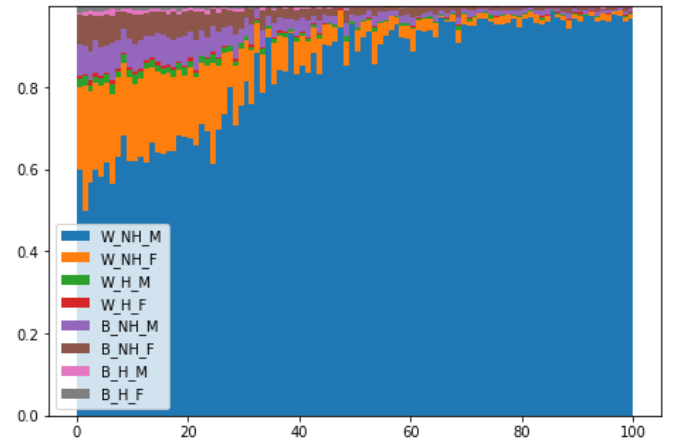


Figure B.12: Evolution of contribution to total approvals for every epoch in a Approval Rate-constrained environment

c) Posteriors Bias ($\theta=0.7$)

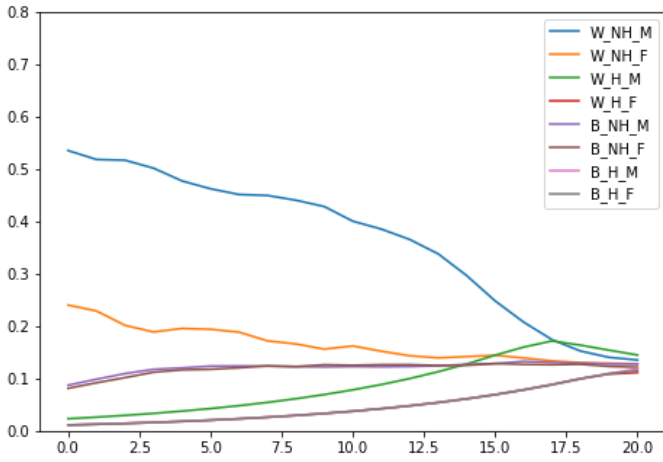


Figure B.13: Evolution of population weights every 5 epochs in a Posteriors-constrained environment

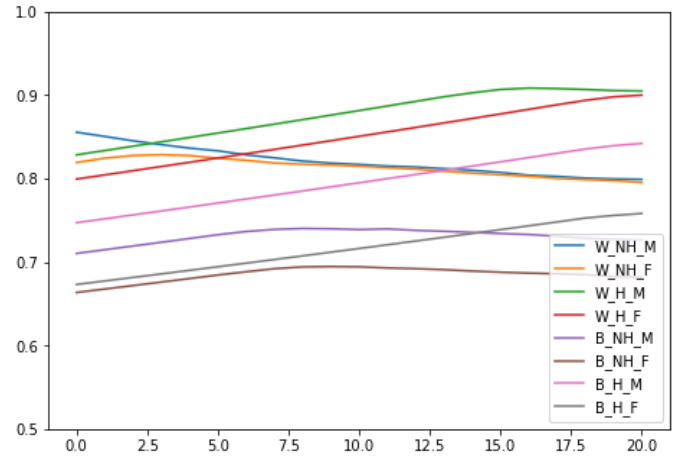


Figure B.14: Evolution of reward estimates every 5 epochs in a Posteriors-constrained environment

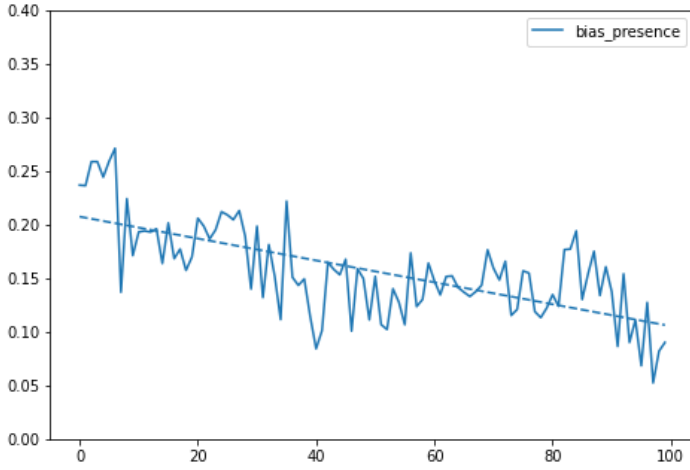


Figure B.15: Evolution of final Presence bias for every epoch in a Posteriors-constrained environment

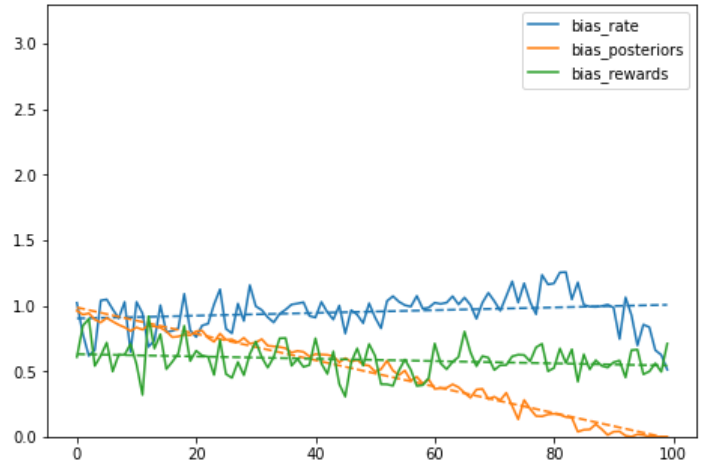


Figure B.16: Evolution of final Rate, Posteriors and Rewards biases in a Posteriors-constrained environment

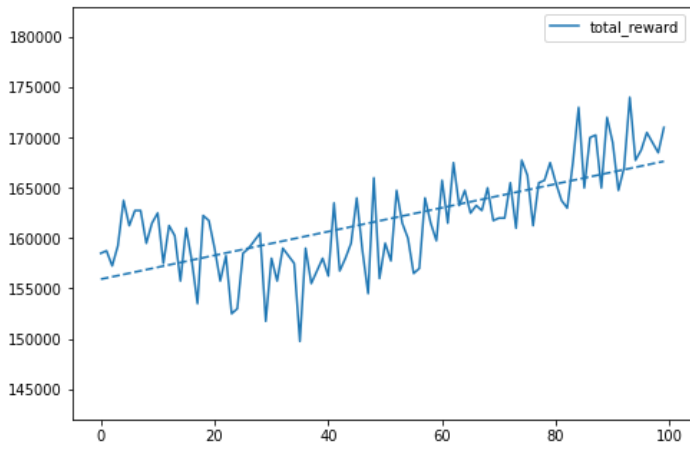


Figure B.17: Evolution of final Total Reward for every epoch in a Posteriors-constrained environment

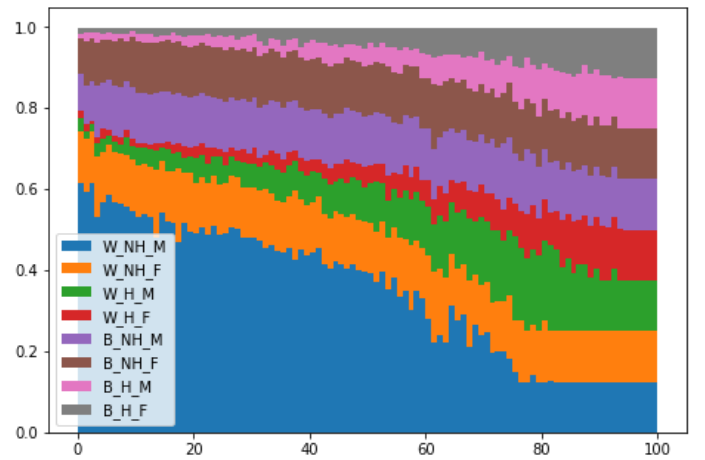


Figure B.18: Evolution of contribution to total approvals for every epoch in a Posteriors-constrained environment

d) Rewards Bias ($\theta=0.45$)



Figure B.19: Evolution of population weights every 5 epochs in a Rewards-constrained environment

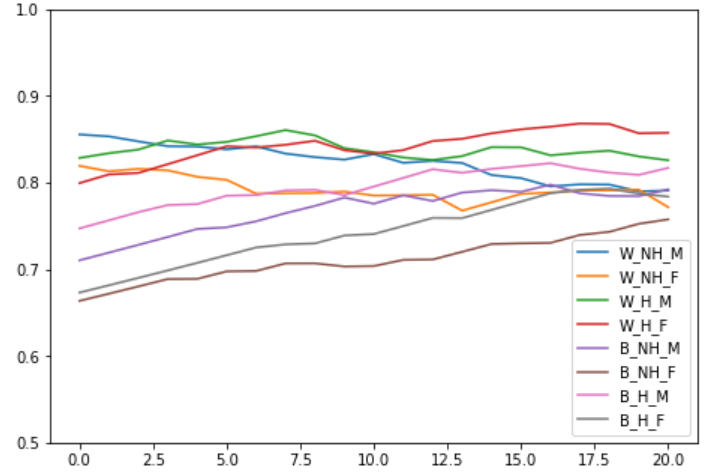


Figure B.20: Evolution of reward estimates every 5 epochs in a Rewards-constrained environment

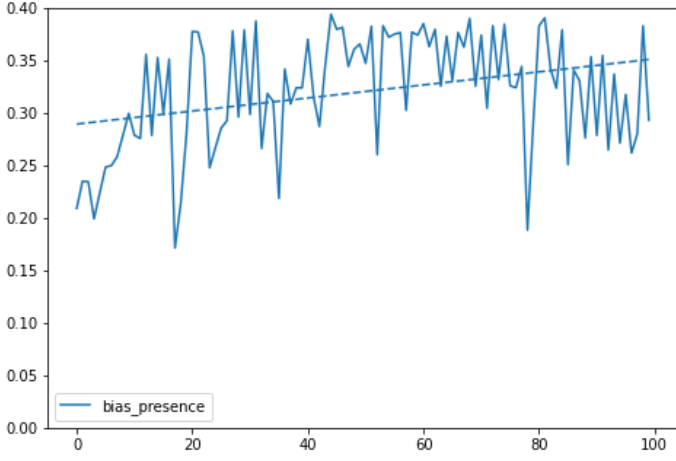


Figure B.21: Evolution of final Presence bias for every epoch in a Rewards-constrained environment

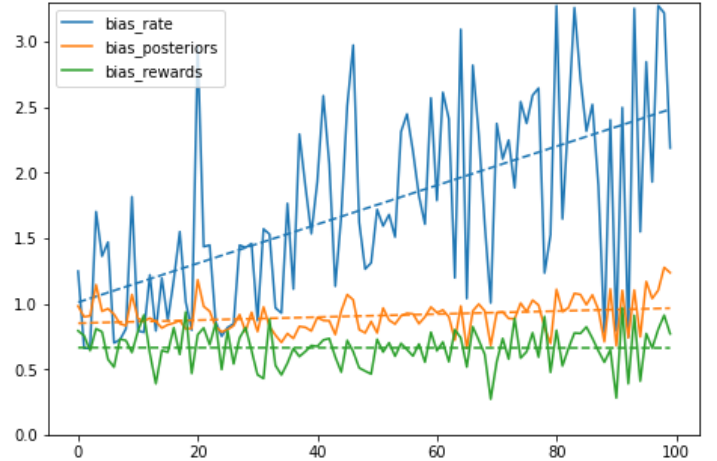


Figure B.22: Evolution of final Rate, Posteriors and Rewards biases in a Rewards-constrained environment

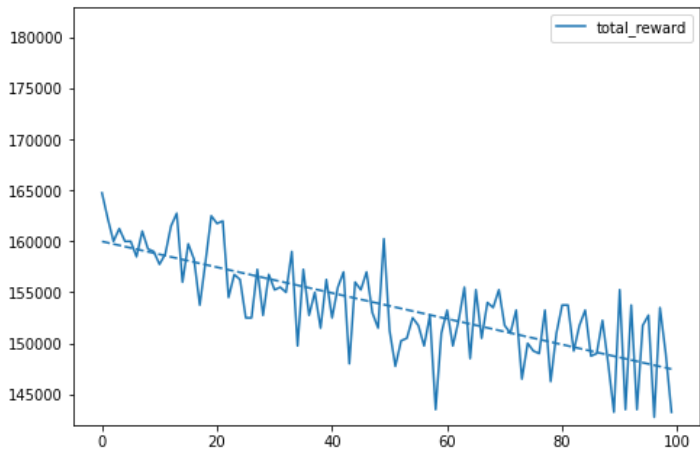


Figure B.23: Evolution of final Total Reward for every epoch in a Rewards-constrained environment

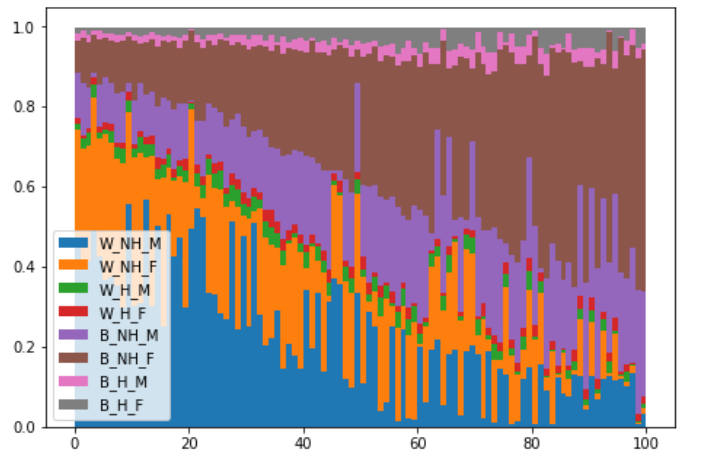


Figure B.24: Evolution of contribution to total approvals for every epoch in a Rewards-constrained environment

$$approval_rates = (|C_a \cup C'| / |C_a| : a \in A)$$

$$bias_rates = \sum_{a=1}^k (|\overline{approval_rates} - approval_rates_a|) \quad bias_rates = \sum_{a=1}^k (|\max(approval_rates) - approval_rates_a|)$$

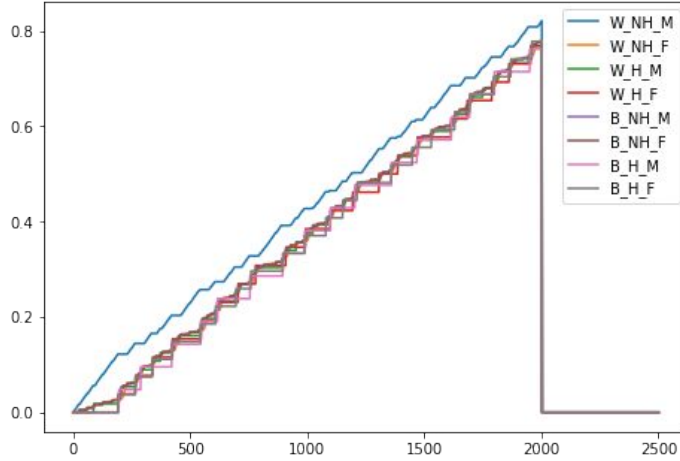


Figure 5.9: Evolution of Approval Rates within a single execution using the faulty Approval Rate constraint as-is

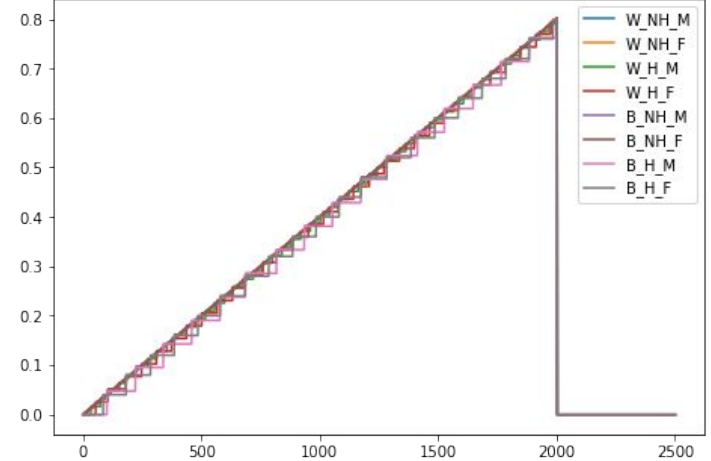


Figure 5.10: Evolution of Approval Rates within a single execution using the corrected Approval Rate constraint

$$approval_rates = (|C_a \cup C'| / |C_a| : a \in A)$$

$$bias_rates = \sum_{a=1}^k (|\overline{approval_rates} - approval_rates_a|) \quad bias_rates = \sum_{a=1}^k (|\max(approval_rates) - approval_rates_a|)$$

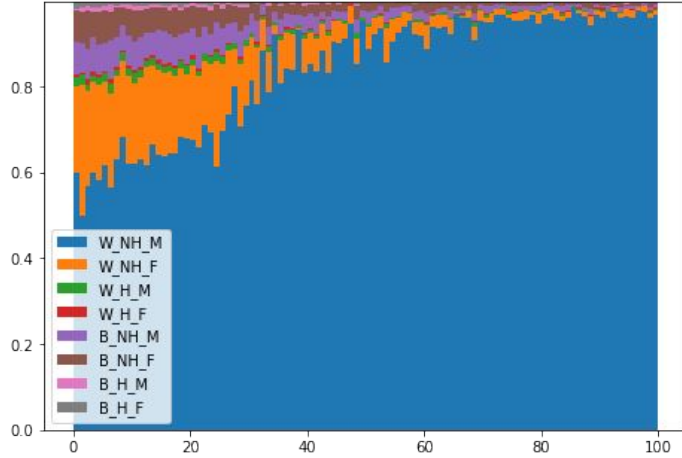


Figure B.12: Evolution of contribution to total approvals for every epoch in a Approval Rate-constrained environment

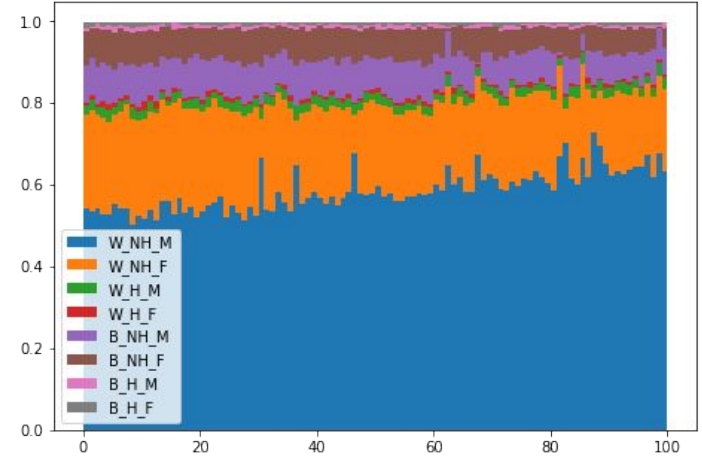


Figure 5.12: Evolution of contribution to total approvals for every epoch using the corrected Approval Rate constraint



Conclusions and Future Work

- ◉ We can summarize what was achieved in this work as the following:
 - Studied and **identified** indications of **bias** in the **original HMDA** dataset.
 - Showcased the **harmful effect** of those underlying biases on a **live environment** when a MAB model computed with a **UCB** algorithm is applied to it sequentially over a **sustained period of time**.
 - Developed **4 bias metrics** adequate to the scope of our simulation.
 - **Examined the effect** of using these metrics in a new version of the UCB algorithm, **BC-UCB**, both in an **isolated** environment as well as in the same **live** environment that changed according to its performance.



Conclusions and Future Work

- The present work could be expanded in the following directions:
 - Adaptation to domains of work other than mortgage lending
 - Expansion of the bias metrics used, either in definition or formalization
 - Development of more adequate ϕ formula, and deliberate usage of β
 - More thorough transformation processes from real data to simulated one
 - Refinement of the environment's modelization
 - Refinement of the (MAB) problem formulation



IMAGE REFERENCES

1. Ledford, Heidi. “Millions of Black People Affected by Racial Bias in Health-Care Algorithms.” *Nature*, vol. 574, no. 7780, 2019, pp. 608–609, doi:10.1038/d41586-019-03228-6.
2. Dastin, Jeffrey. “Amazon Scraps Secret AI Recruiting Tool That Showed Bias against Women.” *Reuters*, 2018, <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>.
3. Angwin, Julia, et al. “Machine Bias.” *ProPublica*, 2016, <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.

Thank you for your attention

Quim De Las Heras Molins (u160402)

quim.delasheras01@estudiant.upf.edu

Juliol 2022

Treball de Fi de Grau

Enginyeria Informàtica

Universitat Pompeu Fabra