# CFRL: A Python library for counterfactually fair offline reinforcement learning using data preprocessing

**Several Different Contributors**[1]¶

**1** Several Different Departments, University of Michigan ¶ Corresponding author

## Summary

Reinforcement learning (RL) aims to learn a sequential decision-making rule, often referred to as a "policy", that maximizes some pre-specified benefit in an environment across multiple or even infinitely many time steps. It has been widely applied to fields such as healthcare, banking, and autonomous driving. Despite their usefulness, the decisions made by RL algorithms might exhibit systematic bias due to bias in the training data. For example, when using an RL algorithm to assign treatment to patients over time, the algorithm might consistently assign treatment resources to patients of some races while ignoring patients of other races. Concerns have been raised that the deployment of such biased algorithms could exacerbate the discrimination faced by socioeconomically disadvantaged groups.

To address this problem, Wang et al. (2025) extended the concept of single-stage counterfactual fairness (Kusner et al., 2018) to the multi-stage setting and proposed a data preprocessing algorithm that ensures counterfactual fairness in offline reinforcement learning. An RL policy is counterfactually fair if, at every time step, it would assign the same decisions with the same probability for an individual had the individual belong to a different subgroup defined by some sensitive attribute (such as race and gender). At its core, counterfactual fairness views the observed states and rewards as biased proxies of the (unobserved) true underlying states and rewards, where the bias can often be seen as a result of the observed sensitive attribute. In this light, the data preprocessing algorithm ensures counterfactual fairness by removing this bias from the input offline trajectories.

The CFRL library is built upon this definition of RL counterfactual fairness introduced in Wang et al. (2025). It implements the data preprocessing algorithm proposed by Wang et al. (2025) and provides a set of tools to evaluate the value and counterfactual fairness achieved by a given policy. In particular, it takes in an offline RL trajectory and outputs a preprocessed, bias-free trajectory, which could be passed to any off-the-shelf offline RL algorithms to learn a counterfactually fair policy. Additionally, it could also take in an RL policy and return its value and level of counterfactual fairness.

## Statement of Need

Many existing Python libraries implement algorithms that ensure fairness in machine learning. For example, `Fairlearn` (Weerts et al., 2023) and `aif360` (Bellamy et al., 2018) provide tools for mitigating bias in single-stage machine learning predictions under statistical assocoiation-based fairness criterion such as demographic parity and equal opportunity. However, they do not focus on counterfactual fairness, which defines fairness from a causal perspective, and they cannot be easily extended to the reinforcement learning setting in general. Additionally, `ml-fairness-gym` (D'Amour et al., 2020) allows users to simulate unfairness in sequential decision-making, but it neither implement algorithms that reduce unfairness nor address counterfactual fairness. To our current knowledge, Wang et al. (2025) is the first work to

42 study counterfactual fairness in reinforcement learning. Correspondingly, CFRL is also the first
43 code library to address counterfactual fairness in the reinforcement learning setting.

44 The contribution of CFRL is two-fold. First, it implements a data preprocessing algorithm that
45 removes bias from offline RL training data. For each individual (or sample) in the data, the
46 preprocessing algorithm estimates the counterfactual states under different sensitive attribute
47 values and concatenates all of the individual's counterfactual states into a new state variable.
48 The preprocessed data can then be directly used by existing RL algorithms for policy learning,
49 and the learned policy should be approximately counterfactually fair. Second, it provides a
50 platform for assessing RL policies based on counterfactual fairness. After passing in a policy
51 and a trajectory dataset from the environment of interest, users can assess how well the policy
52 performs in the environment of interest in terms of the discounted cumulative reward and a
53 counterfactual fairness metric. This not only allows stakeholders to test their fair RL policies
54 before deployment but also offers RL researchers a hands-on tool to evaluate newly developed
55 counterfactually fair RL algorithms.

## High-level Design

57 The CFRL library is composed of 5 major modules. The functionalities of the modules are
58 summarized in the table below.

| Module | Functionalities |
| --- | --- |
| reader | Implements functions that read tabular trajectory data from either a .csv file or a pandas.Dataframe into a format required by CFRL. Also implements functions that export trajectory data to either a .csv file or a pandas.Dataframe. |
| preprocessor | Implements the data preprocessing algorithm introduced in Wang et al. (2025). |
| agents | Implements a fitted Q-iteration (FQI) algorithm, which learns RL policies and makes decisions based on the learned policy. Users can also pass a preprocessor to the FQI; in this case, the FQI will be able to take in unpreprocessed trajectories, internally preprocess the input trajectories, and directly output counterfactually fair policies. |
| environment | Implements a synthetic environment that produces synthetic data as well as a simulated environment that simulates the transition dynamics of the environment underlying some real-world RL trajectory data. Also implements functions for sampling trajectories from the synthetic and simulated environments. |
| fqe | Implements a fitted Q-evaluation (FQE) algorithm, which can be used to evaluate the value of a policy. |
| evaluation | Implements functions that evaluate the value and counterfactual fairness of a policy. Depending on the user's needs, the evaluation can be done either in a synthetic environment or in a simulated environment. |

59 A general CFRL workflow is as follows: First, simulate a trajectory using environment or
60 read in a trajectory using reader. Then, train a preprocessor using preprocessor to remove
61 the bias in the trajectory data. After that, pass the preprocessed trajectory into the FQI
62 algorithm in agents to learn a counterfactually fair policy. Finally, use functions in evaluation
63 to evaluate the value and counterfactual fairness of the trained policy.

## Data Example

⁶⁵ We provide a data example to demontrate how `CFRL` uses real data to learn a counterfactually
⁶⁶ fair policy and evaluate the value and counterfactual fairness of the learned policy. This is
⁶⁷ only one of the many workflows that `CFRL` can perform. We refer interested readers to the
⁶⁸ "Example Workflows" section of the CFRL documentation for more workflow examples.

⁶⁹ Data Loading

⁷⁰ In this demonstration, we use an offline trajectory generated from a `SyntheticEnvironment`
⁷¹ using some pre-specified transition rules. Although it is actually synthesized, we treat it as if it
⁷² is from some unknown environment for pedagogical convenience.

⁷³ The trajectory contains 500 individuals (i.e. $N = 500$) and 10 transitions (i.e. $T = 10$). The
⁷⁴ sensitive attribute variable and the state variable are both univariate. The sensitive attributes
⁷⁵ are binary (0 or 1). The actions are also binary (0 or 1) and were sampled using a policy that
⁷⁶ selects 0 or 1 randomly with equal probability. The trajectory is stored in a tabular format in a
⁷⁷ `.csv` file. We first load the trajectory from the tabular form into the array format required by
⁷⁸ `CFRL`.

```python
zs, states, actions, rewards, ids = read_trajectory_from_dataframe(
                                    path='../data/sample_data_large_uni.csv',
                                    z_labels=['z1'],
                                    state_labels=['state1'],
                                    action_label='action',
                                    reward_label='reward',
                                    id_label='ID',
                                    T=10)
```

⁷⁹ We then split the trajectory data into a training set (80%) and a testing set (20%). The
⁸⁰ training set is used to train the counterfactually fair policy, while the testing set is used to
⁸¹ evaluate the value and counterfactual fairness metric achieved by the policy.

```python
(
    zs_train, zs_test,
    states_train, states_test,
    actions_train, actions_test,
    rewards_train, rewards_test
) = train_test_split(zs, states, actions, rewards, test_size=0.2)
```

⁸² Preprocessor Training & Trajectory Preprocessing

⁸³ We now train the preprocessor and preprocess the trajectory. We set `cross_folds=5`, which
⁸⁴ reduces overfitting so that we do not need a separate dataset to train the preprocessor. In
⁸⁵ this case, `train_preprocessor()` will internally divide the training data into 5 folds, and
⁸⁶ each fold is preprocessed using a model that is trained on the other 4 folds. We initialize
⁸⁷ the `SequentialPreprocessor`, and `train_preprocessor()` will take care of both preprocessor
⁸⁸ training and trajectory preprocessing.

```python
sp = SequentialPreprocessor(z_space=[[0], [1]], num_actions=2, cross_folds=5,
                            mode='single', reg_model='nn')
states_tilde, rewards_tilde = sp.train_preprocessor(
    zs=zs_train, xs=states_train, actions=actions_train, rewards=rewards_train)
```

⁸⁹ Policy Learning

⁹⁰ Now we train a policy using `FQI` and the preprocessed data with `sp` as its internal preprocessor.
⁹¹ Note that the training data `state_tilde` and `rewards_tilde` are already preprocessed. Thus,

<sup>92</sup> we set `preprocess=False` during training so that the input trajectory will not be preprocessed
<sup>93</sup> again by the internal preprocessor (i.e. `sp`).

```
agent = FQI(num_actions=2, model_type='nn', preprocessor=sp)
agent.train(zs=zs_train, xs=states_tilde, actions=actions_train,
            rewards=rewards_tilde, max_iter=100, preprocess=False)
```

<sup>94</sup> `SimulatedEnvironment` Training

<sup>95</sup> Before moving on to the evaluation stage, there is one more thing to do: We need to train a
<sup>96</sup> `SimulatedEnvironment` that mimics the transition rules of the true environment that generated
<sup>97</sup> the training trajectory, which will be used by the evaluation functions to simulate the true
<sup>98</sup> data-generating environment. To do so, we initialize a `SimulatedEnvironment` and train it on
<sup>99</sup> the whole trajectory data (i.e. training set and testing set combined).

```
env = SimulatedEnvironment(num_actions=2,
                           state_model_type='nn',
                           reward_model_type='nn')
env.fit(zs=zs, states=states, actions=actions, rewards=rewards)
```

<sup>100</sup> Value and Counterfactual Fairness Evaluation

<sup>101</sup> We now estimate the value and counterfactual fairness achieved by the trained policy
<sup>102</sup> when interacting with the environment of interest using `evaluate_value_through_fqe()`
<sup>103</sup> and `evaluate_fairness_through_model()`, respectively. The counterfactual fairness is repre-
<sup>104</sup> sented by a metric from 0 to 1, with 0 representing perfect fairness and 1 indicating complete
<sup>105</sup> unfairness. We use the testing set for evaluation.

```
value = evaluate_reward_through_fqe(zs=zs_test, states=states_test,
                                    actions=actions_test, rewards=rewards_test,
                                    policy=agent, model_type='nn')
cf_metric = evaluate_fairness_through_model(env=env, zs=zs_test, states=states_test,
                                            actions=actions_test, policy=agent)
```

<sup>106</sup> The output value is `7.3576775` and cf_metric is `0.041818181818181824`, which indicates
<sup>107</sup> our policy is close to being perfectly counterfactually fair. Indeed, the CF metric should be
<sup>108</sup> exactly 0 if we know the true dynamics of the environment of interest; the reason why it is not
<sup>109</sup> exactly 0 here is because we need to estimate the dynamics of the environment of interest
<sup>110</sup> during preprocessing, which can introduce errors.

<sup>111</sup> Bonus: Assessing a Fairness-through-unawareness Policy

<sup>112</sup> Fairness-through-unawareness proposes to ensure fairness by excluding the sensitive attribute
<sup>113</sup> from the agent's decision-making. However, it might still be unfair because of the indirect bias
<sup>114</sup> in the states and rewards. In this section, we use the same trajectory data to train a policy
<sup>115</sup> following fairness-through-unawareness and estimate its value and counterfactual fairness.

```
agent_unaware = FQI(num_actions=2, model_type='nn', preprocessor=None)
agent_unaware.train(zs=zs_train, xs=states_train, actions=actions_train,
                    rewards=rewards_train, max_iter=100, preprocess=False)
value_unaware = evaluate_reward_through_fqe(
            zs=zs_test, states=states_test, actions=actions_test,
            rewards=rewards_test, policy=agent_unaware, model_type='nn')
cf_metric_unaware = evaluate_fairness_through_model(
            env=env, zs=zs_test, states=states_test,
            actions=actions_test, policy=agent_unaware)
```

<sup>116</sup> The output value is `8.588442` and cf_metric is `0.44636363636363635`. The fairness-through-
<sup>117</sup> unawareness policy is much less fair than the policy learned using the preprocessed trajectory.

---

<sub>118</sub> This suggests that the preprocessing method likely reduced the bias in the training trajectory
<sub>119</sub> effectively.

## Conclusions

<sub>121</sub> CFRL is a Python library that empowers counterfactually fair reinforcement learning through
<sub>122</sub> data preprocessing. It also provides tools to evaluate the value and counterfactual fairness
<sub>123</sub> of a given policy. As far as we know, it is the first library to address counterfactual fairness
<sub>124</sub> problems in the context of reinforcement learning. Nevertheless, despite this, CFRL also admits
<sub>125</sub> a few limitations. For example, the current CFRL implementation requires every individual
<sub>126</sub> in the offline dataset to have the same number of time steps. Extending the library to
<sub>127</sub> accommodate variable-length episodes can improve its flexibility and usefulness. Besides, CFRL
<sub>128</sub> could also be made more well-rounded by integrating the preprocessor with popular offline RL
<sub>129</sub> algorithm libraries such as d3rlpy (Seno & Imai, 2022), or connecting the evaluation functions
<sub>130</sub> with established RL environment libraries such as gym (Towers et al., 2024). We leave these
<sub>131</sub> extensions to future updates.

## Acknowledgements

<sub>133</sub> This is the acknowledgements.

## References

<sub>135</sub> Bellamy, R. K. E., Dey, K., Hind, M., Hoffman, S. C., Houde, S., Kannan, K., Lohia, P.,
<sub>136</sub>    Martino, J., Mehta, S., Mojsilovic, A., Nagar, S., Ramamurthy, K. N., Richards, J., Saha,
<sub>137</sub>    D., Sattigeri, P., Singh, M., Varshney, K. R., & Zhang, Y. (2018). *AI Fairness 360: An*
<sub>138</sub>    *extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias*.
<sub>139</sub>    https://arxiv.org/abs/1810.01943

<sub>140</sub> D'Amour, A., Srinivasan, H., Atwood, J., Baljekar, P., Sculley, D., & Halpern, Y. (2020).
<sub>141</sub>    Fairness is not static: Deeper understanding of long term fairness via simulation studies.
<sub>142</sub>    *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 525–534.
<sub>143</sub>    https://doi.org/10.1145/3351095.3372878

<sub>144</sub> Kusner, M. J., Loftus, J. R., Russell, C., & Silva, R. (2018). *Counterfactual fairness*. https:
<sub>145</sub>    //arxiv.org/abs/1703.06856

<sub>146</sub> Seno, T., & Imai, M. (2022). d3rlpy: An offline deep reinforcement learning library. *Journal of*
<sub>147</sub>    *Machine Learning Research*, *23*(315), 1–20. http://jmlr.org/papers/v23/22-0017.html

<sub>148</sub> Towers, M., Kwiatkowski, A., Terry, J., Balis, J. U., De Cola, G., Deleu, T., Goulão, M.,
<sub>149</sub>    Kallinteris, A., Krimmel, M., KG, A., & others. (2024). Gymnasium: A standard interface
<sub>150</sub>    for reinforcement learning environments. *arXiv Preprint arXiv:2407.17032*.

<sub>151</sub> Wang, J., Shi, C., Piette, J. D., Loftus, J. R., Zeng, D., & Wu, Z. (2025). *Counterfactually fair*
<sub>152</sub>    *reinforcement learning via sequential data preprocessing*. https://arxiv.org/abs/2501.06366

<sub>153</sub> Weerts, H., Dudík, M., Edgar, R., Jalali, A., Lutz, R., & Madaio, M. (2023). Fairlearn:
<sub>154</sub>    Assessing and improving fairness of AI systems. In *Journal of Machine Learning Research*
<sub>155</sub>    (No. 257; Vol. 24, pp. 1–8). http://jmlr.org/papers/v24/23-0389.html