

# Fair Feature Selection with a Lexicographic Multi-Objective Genetic Algorithm

James Brookhouse<sup>[0000-0002-9802-7070]</sup> and Alex Freitas<sup>[0000-0001-9825-4700]</sup>

School of Computing, University of Kent, Canterbury, UK  
james@brookhou.se, A.A.Freitas@kent.ac.uk

**Abstract.** There is growing interest in learning from data classifiers whose predictions are both accurate and fair, avoiding discrimination against sub-groups of people based e.g. on gender or race. This paper proposes a new Lexicographic multi-objective Genetic Algorithm for Fair Feature Selection (LGAFFS). LGAFFS selects a subset of relevant features which is optimised for a given classification algorithm, by simultaneously optimising one measure of accuracy and four measures of fairness. This is achieved by using a lexicographic multi-objective optimisation approach where the objective of optimising accuracy has higher priority over the objective of optimising the four fairness measures. LGAFFS was used to select features in a pre-processing phase for a random forest algorithm. The experiments compared LGAFFS' performance against two feature selection approaches: (a) the baseline approach of letting the random forest algorithm use all features, i.e. no feature selection in a pre-processing phase; and (b) a Sequential Forward Selection method. The results showed that LGAFFS significantly improved fairness measures in several cases, with no significant difference regarding predictive accuracy, across all experiments.

## 1 Introduction

Recently, there has been an increased focus on the fairness of the decisions made by automated processes [17,1]; since algorithms that learn from biased data often produce biased predictive models. We address fairness in the classification task of machine learning, where a predictive feature (e.g. gender or race) is set as a sensitive feature. The values of a sensitive feature are used to split individuals (instances in a dataset) into protected and unprotected groups. The protected group contains individuals likely to be victims of discrimination, who are more likely to obtain a negative outcome (class label) than the unprotected group.

A large number of fairness measures have been proposed to capture some notion of fairness in a model learned from data [14,21]. These fairness measures can be categorised into group-level and individual-level fairness measures.

An example of a group-level fairness metrics is the discrimination score [2], which measures the difference between the predicted positive-class probabilities of the protected and unprotected groups. Some group-level metrics of fairness measure the difference between the false positive error rate and/or the false negative error rate between the protected and unprotected groups [3]. Group-level

fairness measures have the limitation of not considering fairness at the individual level; i.e., they do not penalise models where two very similar individuals within the same group unfairly receive different outcomes (class labels).

An individual-level fairness metric avoids this limitation, by measuring similarities among individuals. Consistency is an individual fairness metric which compares an individual to its  $k$ -nearest neighbours; if all of an individual’s neighbours have the same class as the current individual, this test is considered maximally satisfied for that individual, this is then repeated for each individual and an average taken [22]. However, as the number of features grows, the notion of “nearest neighbours” become increasingly meaningless, as the distances between individuals tend to increase, leading to comparisons being made between increasingly different individuals.

In practice, no single fairness measure can be deemed the best in general, and it has also been proved that there is a clear trade-off among some fairness measures, which cannot be simultaneously optimised [3,10].

Hence, intuitively it makes sense to use multiple fairness measures, with different pros and cons, and try to optimise those multiple measures at the same time, in order to achieve more robust fairness results. This is precisely the focus of this paper, where we propose a new multi-objective Genetic Algorithm (GA) for fair feature selection, The GA uses the lexicographic approach to optimise two objectives in decreasing priority order: predictive accuracy and fairness. The accuracy objective involves one measure, but the fairness objective is more complex and involves four measures. Hence, we propose a new procedure for aggregating four fairness measures into a single fairness objective by systematically considering all permutations of lexicographic ordering of those four measures, as described in detail later.

The GA selects a subset of relevant features for a given classification algorithm in a data pre-processing phase [13]. This is a difficult task for two reasons. First, the search space’s size is exponential in the number of features, with  $2^n - 1$  candidate solutions (feature subsets), where  $n$  is the number of features in the dataset (the “ $- 1$ ” discounts the empty feature subset). Second, intuitively the search space is rugged (highly non-convex) with many local optima, even in a single-objective scenario, with the problem being aggravated in the multi-objective scenario.

We focus on GAs for two main reasons. First, they are robust global search methods, being less likely to get trapped into local optima in the search space, by comparison with conventional local search methods [18,7], and so they tend to cope better with feature interaction (a key issue in feature selection). Second, the fact that they evolve a population of candidate solutions facilitates multi-objective optimisation [5,19], as proposed in this work.

This paper is organised as follows. Section 2 describes the proposed multi-objective genetic algorithm for fair feature selection. Section 3 describes the datasets used in the experiments and the experimental setup. Section 4 reports experimental results and Section 5 presents the conclusions and future work.

---

**Algorithm 1:** Ramped Population Initialisation

---

**Data:** `population_size`, `MIN_P`, `MAX_P`  
**Result:** Population of Individuals

```

1 Function initialise_population():
2   step_size = (MAX_P - MIN_P) / population_size
3   for i to population_size do
4     p = MIN_P + (i * step_size)
5     population += Individual.initialise(p)
6   return population

```

---

## 2 A Lexicographic-Optimisation Genetic Algorithm for Fair Feature Selection

This section describes our new Lexicographic-optimisation Genetic Algorithm for Fair Feature Selection (LGAFFS), which selects a subset of relevant features for a classification algorithm in a data pre-processing phase. LGAFFS selects individuals for reproduction based on the principle of lexicographic optimisation to combine predictive accuracy and fairness measures, as described later.

In LGAFFS, each individual of the population represents a candidate feature subset. More precisely, each individual consists of a string of  $N$  bits (genes), where  $N$  is the number of features in the dataset, and the  $i$ -th gene takes the value 1 or 0 to indicate whether or not (respectively) the  $i$ -th feature is selected.

LGAFFS follows a wrapper approach to feature selection [13], where a base classification algorithm is used to learn a classification model based on the feature subset selected by an individual, and that model’s quality (in terms of accuracy and fairness) is used to compute that individual’s fitness. Hence, the GA aims at finding the best subset of features for the base classification algorithm. Fitness computation is performed by using a well-known internal cross-validation procedure, which uses only the training set (i.e. not using the test set).

LGAFFS uses uniform crossover and bit-flip mutation as genetic operators to generate new individuals in each generation. However, the population initialisation, tournament selection and elitism selection are non-standard procedures, and hence these are described in detail in the next subsections.

### 2.1 Population Initialisation

When creating the initial population, each individual has a different probability that each gene (feature) will be selected or not. This ramping initialisation is described in Algorithm 1. As shown in line 4, each individual has the probability ( $p$ ) that a gene (feature) will be switched on increased by `step_size` compared to the previous individual, where `step_size` is defined in line 2 as a function of the maximum and minimum probabilities for a feature to be selected – denoted `MAX_P` and `MIN_P`, which are input arguments for Algorithm 1.

The motivation for this ramped population initialisation procedure is to promote diversity in the population. If all individuals had the same probability  $p$

---

**Algorithm 2:** Pseudo-code of Lexicographic Tournament selection.

---

```

Data: Instances, Population,  $\epsilon$ , fair_win_ $\epsilon$ 
1 Function tournament_selection():
2   | i1, i2 = select_random_individuals()
3   | if not  $|i1.accuracy - i2.accuracy| > \epsilon$  then
4     |   | i1_win, i2_win = fairness_aggregation(i1,i2) // See Algorithm 3
5     |   |   | if  $|i1\_win - i2\_win| > fair\_win\_epsilon$  then
6     |   |   |   | return fairest_individual
7   | return best_accuracy_individual

```

---

that a single gene is switched on or off, then, as the number of genes (features) in an individual increases, the number of features selected (switched on genes) in each individual would tend to converge to  $p \times N_{genes}$ , the mean of a binomial distribution, where  $p$  is the probability of each gene being switched on and  $N_{genes}$  is the number of genes. Hence, all individuals would tend to have a similar number of selected features, and individuals with low or high numbers of selected features in the initial population would be rare, limiting the search of these areas. The ramped population initialisation avoids this problem, giving each individual a different probability  $p$  of switching a gene, sweeping from MIN\_P to MAX\_P.

## 2.2 Lexicographic Tournament Selection

Lexicographic tournament selection, with tournament size of two, is used to select individuals for reproduction. In the lexicographic-optimisation approach [8], we compare the two individuals in a tournament considering the objectives in decreasing order of priority. Let  $V_1$  and  $V_2$  denote the values of the current objective for individuals 1 and 2. When those two individuals are compared based on the first objective, if  $|V_1 - V_2| > \epsilon$  (where  $\epsilon$  is a very small value), then the best individual is the tournament winner. Otherwise, the two individuals' objective values are deemed equivalent (negligible difference) based on that objective; then the next objective is considered in the same way, and so on. This is repeated until a significant (greater than  $\epsilon$ ) difference is observed and a best individual selected. If there is no significant difference between two individuals for all objectives, then the individual with the best value of the highest priority objective is selected.

The pseudo-code for this is shown in Algorithm 2. When comparing two individuals, the lexicographic approach requires the objectives to be ordered. LGAFFS considers accuracy as the highest-priority objective to be optimised, followed by a lower-priority set of fairness measures which are aggregated into a single fairness objective to be optimised as described in section 2.4.

Note that the lexicographic method avoids the specification of ad-hoc weights to each objective, which would be the case if using a weighted sum of objectives [8]. The lexicographic approach simply requires that an order of priority for the objectives be defined; and intuitively it is easier for users to specify a priority order of objectives than ad-hoc numerical weights. The lexicographic approach

requires a small threshold parameter ( $\epsilon$ ); but again, it is intuitively easier for users to specify this parameter than to specify ad-hoc weights for each objective.

Note that an alternative to the lexicographic approach would be the well-known Pareto dominance approach [5]. However, the Pareto approach is not suitable for our feature selection task where the objective of accuracy has higher priority than the objective of fairness, since the Pareto approach ignores this objective prioritisation. In particular, if we used the Pareto approach, once the fairest model is found it would tend to be preserved by the selection operator and remain in the Pareto front along the GA run even if its accuracy was very low; but that model would be a bad solution, given the objectives' priority order. In this case, the GA would waste computational resources searching on areas of the Pareto front around bad solutions (like areas with maximal fairness but low accuracy). In contrast, a lexicographic approach would never select that fairest model due to its very low accuracy (as the highest-priority objective).

### 2.3 The Four Fairness Measures and the Accuracy Measure

No single fairness measure captures all nuances of a fair model, so LGAFFS optimises four fairness measures to get more robust fairness results. To define these measures we use the following nomenclature:

- $S$ : Protected/sensitive feature: 0  $\rightarrow$  unprotected group, 1  $\rightarrow$  protected group
- $\hat{Y}$  : the predicted class;  $Y$ : the actual class; taking class labels 1 or 0
- $TP, FP, TN, FN$ : Number of True Positives, False Positives, True Negatives and False Negatives, respectively

The first measure is the discrimination score (DS) [2], which is defined as:

$$DS = 1 - \left| P(\hat{Y} = 1|S = 0) - P(\hat{Y} = 1|S = 1) \right| \quad (1)$$

DS is a group-level fairness measure that takes the optimal value of 1 if both protected and unprotected groups have an equal probability of being assigned to the positive class by the classifier. If DS is used on unbalanced datasets, those where the data shows a large difference between the probability of a positive outcome for both groups, to satisfy DS will require a reduction in accuracy. In this case the lexicographic approach is robust to such selective pressures as the ordering of the objectives prioritises accuracy over the fairness measures.

The second measure used is consistency[22], defined as:

$$C = 1 - \frac{1}{Nk} \sum_i \sum_{j \in kNN(x_n)} |\hat{y}_i - \hat{y}_j| \quad (2)$$

Consistency is an individual-level similarity metric that compares the class predicted by a classifier to each instance in the dataset to the class predicted by the classifier to that instance's  $k$  nearest instances (neighbours) in the dataset. If all these neighbours have the same predicted class as the current instance, then that instance is considered consistent. The measure computes the average

degree of consistency over all instances in the dataset. A fully consistent model has a consistency of 1 and an inconsistent model has a value of 0.

Thirdly, the False Positive Error Rate Balance Score (FPERBS) [3,4] is:

$$FPERBS = 1 - \left| \frac{FP_{S=0}}{FP_{S=0} + TN_{S=0}} - \frac{FP_{S=1}}{FP_{S=1} + TN_{S=1}} \right| \quad (3)$$

FPERBS measures the difference in the probability that a truly negative instance is incorrectly predicted as positive between protected and unprotected groups. Fourthly, the False Negative Error Rate Balance Score (FNERBS) [3,9,11] measures the difference in the probability that a truly positive instance is incorrectly predicted as negative between protected and unprotected groups:

$$FNERBS = 1 - \left| \frac{FN_{S=0}}{FN_{S=0} + TP_{S=0}} - \frac{FN_{S=1}}{FN_{S=1} + TP_{S=1}} \right| \quad (4)$$

A score of 1 indicates an optimally fair result for both FPERBS and FNERBS.

As the accuracy measure to be optimised, LGAFFS uses the geometric mean of Sensitivity and Specificity (Equation 5). This measure was chosen because it incentivises the correct classification of both positive-class and negative-class instances, to counteract pressure from the fairness measures to produce maximally fair models that trivially predict the same class for all instances.

$$\begin{aligned} Sensitivity &= \frac{TP}{TP + TN}, & Specificity &= \frac{TN}{TN + FP}, \\ GM_{Sen \times Spec} &= \sqrt{Sensitivity \cdot Specificity} \end{aligned} \quad (5)$$

## 2.4 Aggregating fairness measures

As discussed earlier, LGAFFS optimises one accuracy measure and four fairness measures. We consider accuracy as the highest-priority objective (as usual in machine learning), and the four fairness measures as lower-priority objectives. Among those fairness measures, there is no consensus in the literature about what is the best one, and so it would be “unfair” to prioritise one fairness measure over the others. Hence, we aggregate the four fairness measures into a single objective to be optimised by the GA (in addition to the accuracy objective), by computing all possible 24 (4!) permutations of the four fairness measures.

Algorithm 3 shows how two individuals are compared regarding fairness. Each permutation defines a lexicographic order for the fairness measures which can be evaluated to find the first significant difference between the individuals, at which point the best individual is given a win. A significant difference is one greater than the very small  $\epsilon$ . After all permutations have been evaluated, the individual with the higher number of wins is declared the best individual overall.

---

**Algorithm 3:** Aggregating fairness measures.

---

**Data:**  $Ind_1, Ind_2, \epsilon$   
**Result:** Number of wins for each individual

```

1 Function fairness_aggregation():
2    $i1\_win = i2\_win = 0$ 
3    $permutations = generate\_permutations(measures)$ 
4   forall  $permutations$  do
5     forall  $permutation.measures$  do
6        $i1, i2 = compute\_fairness\_measure(measure, Ind_1, Ind_2)$ 
7       if  $|i1 - i2| > \epsilon$  then
8         if  $i1 > i2$  then
9            $i1\_win++$ 
10          break // Exit inner forall
11        else
12           $i2\_win++$ 
13          break // Exit inner forall
14  return  $i1\_win, i2\_win$ 

```

---

## 2.5 Lexicographic elitism

Recall that the lexicographic approach requires the ranking of objectives, and our GA prioritises the accuracy objective (the geometric mean of Sensitivity and Specificity) over the four fairness measures. The fairness measures are aggregated into a single objective (see Section 2.4). To find the **best individual** the procedure in Algorithm 4 is used for implementing elitism.

First the population is sorted by accuracy (line 2 of Algorithm 4), where any individuals with accuracy within  $\epsilon$  of the most accurate individual are shortlisted for fairness comparison (line 4). These shortlisted individuals have their average rank of fairness computed across all 24 permutations generated as described in Section 2.4. For each permutation of the four fairness measures, the set of

---

**Algorithm 4:** Lexicographic Selection of the Best Individual

---

**Data:**  $population, \epsilon, fair\_rank\_e$   
**Result:** Best individual

```

1 Function get_best_individual():
2    $accuracy\_rank = sort\_population\_by\_accuracy(population)$ 
3    $best\_accur\_indiv = accuracy\_rank.head()$ 
4    $individuals = select\_all\_individuals\_within\_accuracy\_e(accuracy\_rank)$ 
5    $avg\_fair\_rank = average\_rank\_of\_fairness\_permutations(individuals)$ 
6   if  $(avg\_fair\_rank.head().average\_rank - best\_accur\_indiv.average\_rank) >$ 
7      $fair\_rank\_e$  then
8     | return  $avg\_fair\_rank.head()$ 
9   else
10  | return  $best\_accur\_indiv$ 

```

---

**Table 1.** Datasets used in all experiments, detailing the number of instances, features and the sensitive features for each dataset.

Data set	Instances	Features	Sensitive Features
Adult Income (US Census)	48842	14	Race, Gender, Age
German Credit	1000	20	Age, Gender
Credit Card Default	30000	24	Gender
Communities and Crime	1994	128	Race
Student Performance (Portuguese)	650	30	Age, Gender, Relationship
Student Performance (Maths)	396	30	Age, Gender, Relationship
ProPublica recidivism	6167	52	Race, Gender

shortlisted individuals is arranged by its lexicographic order, where the first (last) measure in the permutation is considered the most (least) important. Fairness values within the threshold  $\epsilon$  are considered equivalent and the less important metrics are considered until a significant difference is found.

If the fairest shortlisted individual (with the lowest average rank) has a significantly better rank than the most accurate shortlisted individual (i.e. the difference between their average ranks is greater than  $\text{fair\_rank\_}\epsilon$ ), the former is selected by elitism as the best individual in preference over the most accurate individual – since the difference in accuracy between those two shortlisted individuals is considered non-significant, i.e., within  $\epsilon$ . Otherwise, the most accurate shortlisted individual is selected by elitism.

LGAFFS’ Python code is available at <https://github.com/bunu/LGAFFS>.

## 2.6 Related Work

Quadrianto et al. [16] and Valvidia et al. [20] proposed a GA for fair classification. Both GAs were designed for optimising (hyper)-parameters of a classification algorithm, rather than feature selection; and both GAs use Pareto dominance rather than the lexicographic approach used here. The Pareto approach is sound in general, but as noted earlier, it is not suitable for our feature selection task prioritising accuracy over fairness. La Cava and Moore [12] proposed genetic programming (GP) for feature construction, which can implicitly perform feature selection, but feature construction has a much larger search space than feature selection. Their GP uses lexicase selection, a broadly lexicographic approach. However, instead of ordering the objectives based on user-defined priorities like in LGAFFS; their GP uses *randomised* lexicographic orderings of different subgroups of instances (with different sensitive feature values). The GP evaluates multiple fairness-violation events, each for a different subgroup of instances; but each event is evaluated by the same fairness formula: the difference of error rates (either FP or FN error rates) between all instances and a sub-group of instances. In addition, unlike those three algorithms, LGAFFS combines group-level and individual-level fairness measures, increasing fairness robustness.



**Table 2.** Results for Research Question 1: Comparing the performance of random forest trained with the features selected by LGAFFS in a pre-processing phase against the performance of random forest trained with all features. Showing the values for all five measures being optimised by LGAFFS.

Dataset	Sensitive Feature	$GM_{Sen \times Spec}$		Discrimination Score		Consistency		FPERBS		FNERBS	
		LGAFFS	All feats	LGAFFS	All feats	LGAFFS	All feats	LGAFFS	All feats	LGAFFS	All feats
Adult	Age	0.6475	<b>0.7602</b>	<b>0.8485</b>	0.7557	<b>0.8656</b>	0.7887	<b>0.9522</b>	0.9084	<b>0.9189</b>	0.6416
Adult	Race	0.7409	<b>0.7632</b>	<b>0.9356</b>	0.8955	<b>0.8201</b>	0.7889	<b>0.9862</b>	0.9589	<b>0.9902</b>	0.9004
Adult	Sex	0.7420	<b>0.7623</b>	<b>0.8498</b>	0.8152	<b>0.8163</b>	0.7877	<b>0.9490</b>	0.9217	<b>0.9416</b>	0.9110
German Credit	Age	<b>0.5901</b>	0.5774	<b>0.9361</b>	0.8394	0.7642	<b>0.8032</b>	<b>0.8633</b>	0.7626	0.8911	<b>0.8968</b>
German Credit	Gender	<b>0.6036</b>	0.5637	<b>0.9399</b>	0.9087	0.7510	<b>0.8114</b>	<b>0.8993</b>	0.8489	0.9230	<b>0.9349</b>
Student Maths	Age	<b>0.9208</b>	0.9046	0.7964	<b>0.8169</b>	0.8444	<b>0.8483</b>	<b>0.9022</b>	0.8781	0.8951	<b>0.8974</b>
Student Maths	Dalc	0.8951	<b>0.9072</b>	<b>0.7563</b>	0.7214	0.8377	<b>0.8447</b>	0.8051	<b>0.8312</b>	0.8157	<b>0.8160</b>
Student Maths	Famrel	<b>0.9052</b>	0.8914	0.7039	<b>0.7148</b>	0.8361	<b>0.8412</b>	<b>0.8760</b>	0.8480	<b>0.9222</b>	0.9189
Student Maths	Romantic	0.8984	<b>0.9008</b>	<b>0.9076</b>	0.9055	<b>0.8468</b>	0.8443	<b>0.9151</b>	0.9090	<b>0.9106</b>	0.8989
Student Maths	Sex	<b>0.9027</b>	0.8977	0.8397	<b>0.8652</b>	0.8387	<b>0.8488</b>	0.7646	<b>0.8233</b>	0.9012	<b>0.9251</b>
Student Maths	Walc	0.9000	<b>0.9012</b>	<b>0.8364</b>	0.8023	0.8376	<b>0.8503</b>	0.8206	<b>0.8362</b>	<b>0.9196</b>	0.8741
Student Portuguese	Age	<b>0.8196</b>	0.7864	<b>0.8638</b>	0.8536	0.9106	<b>0.9192</b>	<b>0.7038</b>	0.6686	<b>0.9578</b>	0.9428
Student Portuguese	Dalc	<b>0.7867</b>	0.7834	0.8470	<b>0.8758</b>	0.9128	<b>0.9251</b>	<b>0.6230</b>	0.5511	0.9088	<b>0.9376</b>
Student Portuguese	Famrel	0.8031	<b>0.8035</b>	0.8019	<b>0.8305</b>	0.9097	<b>0.9177</b>	<b>0.6450</b>	0.6286	0.9205	<b>0.9583</b>
Student Portuguese	Romantic	<b>0.7846</b>	0.7825	0.9370	<b>0.9452</b>	0.9211	0.9189	0.7608	<b>0.7775</b>	0.9686	<b>0.9796</b>
Student Portuguese	Sex	<b>0.8110</b>	0.7994	0.9200	<b>0.9313</b>	0.9134	<b>0.9239</b>	<b>0.7646</b>	0.7029	<b>0.9708</b>	0.9650
Student Portuguese	Walc	<b>0.8035</b>	0.7831	0.9271	<b>0.9380</b>	0.9186	<b>0.9205</b>	0.7214	<b>0.7248</b>	<b>0.9679</b>	0.9671
Communities & Crime	Race	0.8303	<b>0.8419</b>	<b>0.6623</b>	0.5957	0.6056	<b>0.6126</b>	<b>0.9182</b>	0.8553	<b>0.8313</b>	0.7762
Default of Credit	Sex	0.5896	<b>0.5910</b>	<b>0.9755</b>	0.9724	0.8354	<b>0.8395</b>	0.9739	<b>0.9776</b>	<b>0.9852</b>	0.9832
Propublica Recidivism	Race	0.7306	<b>0.7515</b>	<b>0.8324</b>	0.7744	<b>0.6849</b>	0.6836	<b>0.9022</b>	0.8169	<b>0.9105</b>	0.8834
Propublica Recidivism	Sex	0.7113	<b>0.7518</b>	<b>0.9408</b>	0.9189	0.6756	<b>0.6876</b>	<b>0.9254</b>	0.9125	<b>0.9451</b>	0.9400
Number of Wins		10	11	13	8	6	15	15	6	13	8
Wilcoxon Signed-Rank Test		0.9442		0.06288		0.07346		<b>▲ 0.0088</b>		0.18024	

### 3 Datasets and Experimental Setup

Table 1 describes the 7 binary classification datasets used. When a dataset has multiple sensitive features – a sensitive feature is one which represents a protected characteristic and/or a group that is unfairly treated – the algorithm is ran multiple times using a different sensitive feature each time. 6 datasets are from the UCI Machine Learning repository [6]. The 7th dataset is from ProPublica, investigating biases in predicting if criminals would re-offend [1].

For all datasets, except Adult Income, the experiments use a well-known 10-fold cross-validation procedure. Adult Income is already partitioned into a training and test set, so this partition is used instead of cross validation. LGAFFS’ parameters were not optimised and were set as follows:  $\epsilon$ : 0.01 (threshold for significant differences in Algorithms 2, 3 and 4), `fair_rank_` $\epsilon$ : 1 (threshold for significant fairness-rank differences in Algorithm 4), `fair_win_` $\epsilon$ : 1 (threshold for significant difference in the number of wins among 24 permutations of fairness measures in Algorithm 2); `population_size`: 100, `MAX_P`: 0.5 and `MIN_P`: 0.1 (for population initialisation in Algorithm 1), `internal cross validation folds`: 3, `max_iterations`: 50, `tournament_size`: 2, `crossover_probability`: 0.9, `mutation_probability`: 0.05.

## 4 Experimental Results

We addresses two research questions. First, we compare the use of LGAFFS to select features in a pre-processing phase against the baseline of no feature selection in that phase. Second, we compare LGAFFS to the popular Sequential Forward Selection (SFS) method. Both LGAFFS and SFS use the wrapper approach to feature selection; i.e., they repeatedly use a base classification algorithm to evaluate feature subsets. The base algorithm was Random Forest from scikit-learn [15], with default parameter settings; which was chosen because it is a very popular and powerful classification algorithm. Note that the Random Forest algorithm performs embedded feature selection (during its run), but that feature selection considers only accuracy; whilst using LGAFFS to perform feature selection in a pre-processing phase we optimise both accuracy and fairness.

We also calculated the Pearson’s linear correlation coefficient for each of the 6 pairs of fairness measures for LGAFFS, the coefficients were: 0.71 for (DS,FNERBS), -0.59 for (C,FPERBS), 0.42 for (C,FNERBS), 0.24 for (DS,C), 0.08 for (DS,FPERBS) and 0.02 for (FPERBS,FNERBS). So, 4 pairs of fairness measures have an absolute value of correlation smaller than 0.5.

#### 4.1 RQ1: Does LGAFFS select a better subset than the full set?

The first research question asks whether using LGAFFS to select features in a pre-processing phase leads to better results than the baseline approach of not performing any feature selection. That is, does the random forest algorithm perform better (regarding accuracy and fairness) when it is trained with the features selected by LGAFFS or when it is trained with the full feature set?

**Table 3.** Results for Research Question 2: Comparing the performance of random forest trained with the features selected by LGAFFS vs. Sequential Forward Selection (both selecting features in a pre-processing phase). Showing the values for all five measures being optimised by LGAFFS.

Dataset	Feature	Sensitive		$GM_{Sen \times Spec}$		Discrimination Score		Consistency		FPERBS		FNERBS	
		LGAFFS	SFS	LGAFFS	SFS	LGAFFS	SFS	LGAFFS	SFS	LGAFFS	SFS	LGAFFS	SFS
Adult	Age	0.6475	<b>0.7482</b>	<b>0.8485</b>	0.7968	<b>0.8656</b>	0.8142	<b>0.9522</b>	0.9437	<b>0.9189</b>	0.8306		
Adult	Race	0.7409	<b>0.7465</b>	0.9356	<b>0.9358</b>	<b>0.8201</b>	0.8149	0.9862	<b>0.9887</b>	<b>0.9902</b>	0.9444		
Adult	Sex	0.7420	<b>0.7480</b>	<b>0.8498</b>	0.8331	<b>0.8163</b>	0.8145	<b>0.9490</b>	0.9396	<b>0.9416</b>	0.9057		
German Credit	Age	<b>0.5901</b>	0.4653	<b>0.9361</b>	0.8872	0.7642	<b>0.8110</b>	<b>0.8633</b>	0.8095	<b>0.8911</b>	0.8890		
German Credit	Gender	<b>0.6036</b>	0.4851	<b>0.9399</b>	0.9237	0.7510	<b>0.8212</b>	<b>0.8993</b>	0.8457	<b>0.9230</b>	0.9202		
Student Maths	Age	<b>0.9208</b>	0.9041	<b>0.7964</b>	0.7954	<b>0.8444</b>	0.8407	<b>0.9022</b>	0.8468	<b>0.8951</b>	0.8619		
Student Maths	Dalc	<b>0.8951</b>	0.8921	<b>0.7563</b>	0.6893	0.8377	<b>0.8392</b>	<b>0.8051</b>	0.7940	0.8157	<b>0.8177</b>		
Student Maths	Famrel	0.9052	<b>0.9069</b>	<b>0.7039</b>	0.6707	0.8361	<b>0.8367</b>	<b>0.8760</b>	0.7993	<b>0.9222</b>	0.9096		
Student Maths	Romantic	0.8984	<b>0.9033</b>	<b>0.9076</b>	0.9008	<b>0.8468</b>	0.8311	<b>0.9151</b>	0.8906	0.9106	<b>0.9170</b>		
Student Maths	Sex	0.9027	<b>0.9265</b>	0.8397	<b>0.8468</b>	<b>0.8387</b>	0.8361	0.7646	<b>0.8602</b>	0.9012	<b>0.9521</b>		
Student Maths	Walc	0.9000	<b>0.9084</b>	<b>0.8364</b>	0.8151	<b>0.8376</b>	0.8347	0.8206	<b>0.8796</b>	<b>0.9196</b>	0.9026		
Student Portuguese	Age	<b>0.8196</b>	0.7812	0.8638	<b>0.8909</b>	<b>0.9106</b>	0.9050	0.7038	<b>0.7445</b>	<b>0.9578</b>	0.9576		
Student Portuguese	Dalc	0.7867	<b>0.7917</b>	0.8470	<b>0.8580</b>	<b>0.9128</b>	0.9090	<b>0.6230</b>	0.6138	0.9088	<b>0.9460</b>		
Student Portuguese	Famrel	<b>0.8031</b>	0.7946	<b>0.8019</b>	0.7392	<b>0.9097</b>	0.9050	<b>0.6450</b>	0.5911	<b>0.9205</b>	0.8538		
Student Portuguese	Romantic	<b>0.7846</b>	0.7833	0.9370	<b>0.9378</b>	<b>0.9211</b>	0.9062	<b>0.7608</b>	0.7273	0.9686	<b>0.9716</b>		
Student Portuguese	Sex	<b>0.8110</b>	0.8017	<b>0.9200</b>	0.9148	<b>0.9134</b>	0.9032	<b>0.7646</b>	0.7504	<b>0.9708</b>	0.9554		
Student Portuguese	Walc	<b>0.8035</b>	0.7845	<b>0.9271</b>	0.9179	<b>0.9186</b>	0.9087	0.7214	<b>0.7393</b>	<b>0.9679</b>	0.9541		
Communities and Crime	Race	<b>0.8303</b>	0.7923	<b>0.6623</b>	0.5740	<b>0.6056</b>	0.6049	<b>0.9182</b>	0.7935	<b>0.8313</b>	0.6992		
Default of Credit	Sex	<b>0.5896</b>	0.5639	0.9755	<b>0.9816</b>	0.8354	<b>0.8580</b>	0.9739	<b>0.9807</b>	0.9852	<b>0.9889</b>		
Propublica Recidivism	Race	<b>0.7306</b>	0.7200	<b>0.8324</b>	0.7763	<b>0.6849</b>	0.6617	<b>0.9022</b>	0.8026	<b>0.9105</b>	0.8769		
Propublica Recidivism	Sex	0.7113	<b>0.7199</b>	<b>0.9408</b>	0.8660	<b>0.6756</b>	0.6661	<b>0.9254</b>	0.8492	<b>0.9451</b>	0.9058		
Number of Wins		12	9	15	6	16	5	15	6	15	6		
Wilcoxon Signed-Rank Test		0.17068		<b>▲ 0.00634</b>		0.05876		<b>▲ 0.04236</b>		<b>▲ 0.030</b>			

Table 2 show the experimental results for this research question. In this table, the first two columns show the dataset and the sensitive feature. The following ten columns show the accuracy and fairness results of training the Random Forest algorithm with features selected by LGAFFS or with the full feature set.

In each row of this table (i.e. for each pair of a dataset and a sensitive feature), for each pair of columns comparing the accuracy or fairness of LGAFFS against the full feature set, the best result is shown in boldface. The last but one row of the table shows the number of wins for each approach for each of the five measures of performance, whilst the last row shows the p-value obtained by the Wilcoxon signed-rank statistical significance test. Statistically significant results, at the conventional significance level of  $\alpha = 0.05$ , are marked with a red triangle. In Table 2, there was no substantial difference in the number of wins regarding accuracy. LGAFFS achieved substantially more wins in three of the four fairness measures, with statistical significance in one measure: FPERBS.

#### 4.2 RQ2: Does LGAFFS perform better than SFS?

The second question involves the comparison of LGAFFS to a popular local search-based feature selection method, viz. Sequential Forward Selection (SFS). The SFS method is not aware of the 4 fairness measures; it is just optimising the accuracy measure, i.e., the geometric mean of sensitivity and specificity.

Table 3 presents the results for the Random Forest algorithm when using LGAFFS or SFS to select features. LGAFFS achieved more wins in all 5 measures, with statistical significance shown in 3 of the 4 fairness measures.

## 5 Conclusions

We have proposed a new lexicographic-optimisation Genetic Algorithm for fair feature selection, which selects a feature subset optimised for a classification algorithm based on both predictive accuracy and 4 fairness measures capturing different aspects of fairness, including both group-level and individual-level fairness. No single fairness measure reflects all the nuances of fairness; LGAFFS optimises multiple fairness measures to obtain more robust fairness results.

LGAFFS was compared with 2 other feature selection approaches (no feature selection and Sequential Forward Selection) using Random Forest as the classification algorithm. There was no significant difference in the predictive accuracies of models learned when using LGAFFS versus the 2 other approaches. Regarding fairness, when comparing LGAFFS against the 2 other approaches across the 4 fairness measures, LGAFFS achieved significantly better results in 4 of the 8 comparisons, and there was no significant differences between LGAFFS and the 2 other approaches in the other 4 comparisons.

Future work could include extending SFS to make it a fairness-aware method.

**Acknowledgements:** This work was funded by a research grant from The Leverhulme Trust, UK, reference number RPG-2020-145.

## References

1. Angwin, J., Larson, J., Mattu, S., Kirchner, L.: Machine bias: There's software used across the country to predict future criminals. and it's biased against blacks. (2016), <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
2. Calders, T., Verwer, S.: Three naive bayes approaches for discrimination-free classification. *Data Mining and Knowledge Discovery* **21**(2), 277–292 (2010)
3. Chouldechova, A.: Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data* **5**(2), 153–163 (2017)
4. Corbett-Davies, S., Goel, S.: The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023* (2018)
5. Deb, K.: *Multi-objective optimization using evolutionary algorithms*. John Wiley & Sons (2002)
6. Dua, D., Graff, C.: *UCI machine learning repository* (2017), <http://archive.ics.uci.edu/ml>
7. Freitas, A.: *Data Mining and Knowledge Discovery with Evolutionary Algorithms*. Springer (2002)
8. Freitas, A.A.: A critical review of multi-objective optimization in data mining: a position paper. *ACM SIGKDD Explorations Newsletter* **6**(2), 77–86 (2004)
9. Hardt, M., Price, E., Srebro, N.: Equality of opportunity in supervised learning. In: *Advances in neural information processing systems*. pp. 3315–3323 (2016)
10. Kleinberg, J., Mullainathan, S., Raghavan, M.: Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807* (2016)
11. Kusner, M.J., Loftus, J., Russell, C., Silva, R.: Counterfactual fairness. In: *Advances in neural information processing systems*. pp. 4066–4076 (2017)
12. La Cava, W., Moore, J.: Genetic programming approaches to learning fair classifiers. In: *Proc. Genetic and Evolutionary Computation Conference (GECCO-2020)*. pp. 967–975 (2020)
13. Li, J., Cheng, K., Wang, S., Morstatter, F., Trevino, R., Tang, J., Liu, H.: Feature selection: a data perspective. *ACM Computing Surveys* **50**(6), 94:1–94:45 (2017)
14. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., Galstyan, A.: A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635* (2019)
15. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: *Scikit-learn: Machine learning in Python*. *Journal of Machine Learning Research* **12**, 2825–2830 (2011)
16. Quadrianto, N., Sharmanska, V.: Recycling privileged learning and distributed matching for fairness. In: *Proceedings of the 31st Conference on Neural Information Processing Systems (NIPS 2017)*. pp. 677–688 (2017)
17. Skeem, J.L., Lowenkamp, C.T.: Risk, race, & recidivism: Predictive bias and disparate impact.(2016). *Criminology* **54**, 680 (2016)
18. Telikani, A., Tahmassebi, A., Banzhaf, W., Gandomi, A.: Evolutionary machine learning: a survey. *ACM Computing Surveys* **54**(8), 161:1–161:35 (2021)
19. Tian, Y., Si, L., Zhang, X., Cheng, R., He, C., Tan, K.C., Jin, Y.: Evolutionary large-scale multi-objective optimization: a survey. *ACM Computing Surveys* **54**(8), 174:1–174:34 (2021)
20. Valdivia, A., Sanchez-Monedero, J., Casillas, J.: How fair can we go in machine learning? assessing the boundaries of accuracy and fairness. *International Journal of Intelligent Systems* **36**(4), 1619–1643 (2021)

21. Verma, S., Rubin, J.: Fairness definitions explained. In: 2018 IEEE/ACM International Workshop on Software Fairness (FairWare). pp. 1–7. IEEE (2018)
22. Zemel, R., Wu, Y., Swersky, K., Pitassi, T., Dwork, C.: Learning fair representations. In: International Conference on Machine Learning. pp. 325–333 (2013)