

Minimizing Row Removal to Achieve Desired Average Treatment Effect (ATE) in Datasets

Amit Grosman

Technion - Israel Institute of Technology

Under the guidance of assistant professors Brit Youngmann and Shaul Almagor

October 2024

Contents

1	Introduction	3
2	Problem Definition	4
2.1	General Description and Formal Definition	4
2.2	Calculation of Average Treatment Effect (ATE)	4
3	Proposed Algorithms	5
3.1	Brute-Force Approach	5
3.1.1	Description of Brute-Force Method	5
3.1.2	Pseudocode for Brute-Force Method	5
3.1.3	Why Brute-Force Was Not Utilized	5
3.2	Naive Removal Algorithm	5
3.2.1	Description of Naive Method	5
3.3	Binary Search Removal Algorithm	6
3.3.1	Description of Binary Search Method	6
3.4	Greedy Combined Algorithm	6
3.4.1	Description of Greedy Combined Algorithm	6
3.4.2	Pseudocode	6
3.4.3	Algorithm Explanation	7
3.5	Per Group Removal Algorithm with Confounders	8
3.5.1	Description of Per Group Removal Algorithm	8
3.5.2	Pseudocode	8
3.5.3	Algorithm Explanation	9
3.5.4	Implementation Details	10
4	Experiments	11
4.1	Explanation of Epsilon Factor	11
4.2	Experiment Configuration	11
4.2.1	First Experiment: Without Confounders	11
4.2.2	Second Experiment: With Binary Confounder W	11
4.2.3	Third Experiment: Real-World Data	11
4.3	Experimental Procedure	12
5	Results	13
5.1	First Experiment: Without Confounders	13
5.1.1	Epsilon Factor = 1,000	13
5.1.2	Summary	14
5.2	Second Experiment: With Binary Confounder W	15
5.2.1	Epsilon Factor = 100	15
5.2.2	Summary	16
5.3	Third Experiment: Real-World Data	16
5.3.1	Description of Third Experiment	16
5.3.2	Experiment Setup	16
5.3.3	Results of Third Experiment	17
6	Discussion	19
6.1	Analysis of Results	19
6.2	Scalability	19
6.3	Impact of Confounders	19
7	Conclusion	20

1 Introduction

The Average Treatment Effect (ATE) is a fundamental metric in causal inference, measuring the mean difference in outcomes between treated and control groups. Accurate estimation of ATE is crucial across various disciplines, including medicine, economics, and public health. However, observational datasets often contain confounders—variables that influence both the treatment assignment and the outcome—resulting in biased ATE estimates.

To mitigate this bias, selectively removing data rows can adjust the ATE to a desired range. The primary challenge is to minimize the number of rows removed to preserve the dataset's integrity while achieving precise causal inferences. This project develops and implements algorithms aimed at efficiently identifying and removing the minimal set of rows necessary to attain a specified ATE range. The methodology progresses from handling datasets with only treatment and outcome variables to more complex scenarios incorporating confounders. The proposed algorithms are evaluated on both synthetic and real-world datasets, demonstrating their effectiveness in balancing data preservation with statistical objectives.

2 Problem Definition

2.1 General Description and Formal Definition

The primary objective is to adjust the Average Treatment Effect (ATE) of a dataset to fall within a specified target range by selectively removing data rows. Formally, given a dataset D comprising n rows, each with a binary treatment indicator T , an outcome variable O , and potentially additional confounders W , the goal is to find a subset $D' \subseteq D$ such that:

$$v - \epsilon \leq \text{ATE}(D') \leq v + \epsilon$$

where:

- v is the target ATE.
- ϵ is the acceptable deviation from the target.

2.2 Calculation of Average Treatment Effect (ATE)

The Average Treatment Effect (ATE) is calculated as the difference in the average outcomes between the treated and control groups. Mathematically, it is defined as:

$$\text{ATE} = \mathbb{E}[O|T = 1] - \mathbb{E}[O|T = 0]$$

where:

- $\mathbb{E}[O|T = 1]$ is the expected outcome for the treated group.
- $\mathbb{E}[O|T = 0]$ is the expected outcome for the control group.

When confounders W are present, they are variables that influence both the treatment assignment T and the outcome O . To account for confounders, the ATE can be adjusted using causal inference techniques to isolate the effect of the treatment from the confounding variables. This ensures that the estimated ATE reflects the true causal impact of the treatment.

The adjusted ATE considering confounders can be expressed as:

$$\text{ATE} = \sum_w (\mathbb{E}[O|T = 1, W = w] - \mathbb{E}[O|T = 0, W = w]) P(W = w)$$

where:

- $\mathbb{E}[O|T = 1, W = w]$ and $\mathbb{E}[O|T = 0, W = w]$ are the expected outcomes for treated and control groups within each confounder value $W = w$.
- $P(W = w)$ is the probability of confounder W having the value w .

3 Proposed Algorithms

3.1 Brute-Force Approach

3.1.1 Description of Brute-Force Method

A brute-force approach to this problem involves evaluating all possible subsets of the dataset D and selecting the one that satisfies the ATE condition with the minimal number of removals. Specifically, for each possible number of rows to remove k (starting from 0 up to n), the algorithm checks every combination of k rows to determine if removing them results in D' with the desired ATE range.

3.1.2 Pseudocode for Brute-Force Method

Algorithm 1 Brute-Force ATE Adjustment

Require: Dataset D with binary attributes T , outcome O , target ATE v , threshold ϵ

Ensure: Minimal number of tuples to remove to achieve $v - \epsilon \leq \text{ATE} \leq v + \epsilon$

```
1: for  $k = 0$  to  $n$  do
2:   for each subset  $S \subseteq D$  with  $|S| = k$  do
3:      $D' = D \setminus S$ 
4:     Compute  $\text{ATE}(D')$ 
5:     if  $v - \epsilon \leq \text{ATE}(D') \leq v + \epsilon$  then
6:
7:       return  $k$ 
8:     end if
9:   end for
10: end for
11: return No solution found
```

3.1.3 Why Brute-Force Was Not Utilized

Although the brute-force approach provides an exhaustive search for the minimal number of tuples to remove, its exponential time complexity $O(2^n)$ makes it computationally infeasible for large datasets that are commonly encountered in real-world applications. As the size of the dataset increases, the number of possible subsets grows exponentially, resulting in impractical computation times. Therefore, the focus shifted to more efficient algorithms like the **Greedy Combined Algorithm** and the **Per Group Removal Algorithm**, which offer significant reductions in computational complexity while still effectively achieving the desired ATE adjustments.

3.2 Naive Removal Algorithm

3.2.1 Description of Naive Method

The **Naive Removal Algorithm** iteratively removes the tuple with the highest or lowest outcome O from the treated or control group, respectively, until the ATE falls within the desired range. This straightforward approach focuses on getting to the desired range as fast as possible, but in the price of not being able to fine tune the removal process when being close to the desired range. While simple and easy to implement, the naive method may not always achieve the most optimal removal of tuples.

3.3 Binary Search Removal Algorithm

3.3.1 Description of Binary Search Method

The **Binary Search Removal Algorithm** enhances the naive approach by incorporating a binary search strategy to identify the optimal row for removal every iteration. Instead of solely removing the highest or lowest O values, this method simulates the removal of candidates to determine which row most effectively brings the ATE closer to the target. By evaluating the impact of potential removals, the binary search method achieves greater precision in adjusting the ATE, but at the cost of increased computational effort compared to the naive method.

3.4 Greedy Combined Algorithm

3.4.1 Description of Greedy Combined Algorithm

The **Greedy Combined Algorithm** derives its name from the integration of the **Naive Removal Algorithm** and the **Binary Search Removal Algorithm**. Initially, it employs the naive strategy of removing the highest or lowest O values to adjust the ATE. Upon detecting a sign change in the difference between the current ATE and the target (indicating that the removal has overshoot the desired range), the algorithm transitions to a binary search approach. This hybrid strategy leverages the speed of the naive method and the precision of the binary search method, aiming to achieve optimal ATE adjustments efficiently.

3.4.2 Pseudocode

Below is the pseudocode for the Greedy Combined Algorithm:

Algorithm 2 Greedy Combined Algorithm

Require: Dataset D with a binary attribute T and outcome O , a target difference v and a threshold ϵ

Ensure: The number of tuples deleted from D

```
1: Initialize the groups and their sums, sizes, and current difference:
2:    $D0, D1, v_c \leftarrow \text{InitializeGroups}(D, T, O)$ 
3:  $\text{num\_removals} \leftarrow 0$ 
4:  $\text{prev\_diff} \leftarrow v_c - v$ 
5:  $\text{switched} \leftarrow \text{False}$ 
   /* Greedy naive approach loop */
6: while  $v_c \notin [v - \epsilon, v + \epsilon]$  and  $D0, D1$  not empty do
7:   if  $v_c > v + \epsilon$  then
8:     Remove tuple with max  $O$  from  $D1$ 
9:   else if  $v_c < v - \epsilon$  then
10:    Remove tuple with max  $O$  from  $D0$ 
11:   end if
12:   Update  $v_c \leftarrow \text{UpdateDifference}(D0, D1)$ 
13:    $\text{num\_removals} \leftarrow \text{num\_removals} + 1$ 
14:    $\text{curr\_diff} \leftarrow v_c - v$ 
15:   if  $\text{prev\_diff} \times \text{curr\_diff} < 0$  then
16:      $\text{switched} \leftarrow \text{True}$ 
17:     break
18:   end if
19:    $\text{prev\_diff} \leftarrow \text{curr\_diff}$ 
20: end while
   /* Switch to binary search for fine-tuning */
21: if  $\text{switched}$  then
22:   while  $v_c \notin [v - \epsilon, v + \epsilon]$  and  $D0, D1$  not empty do
23:      $\text{best\_candidate} \leftarrow \text{BinarySearchBestCandidate}(\text{candidates}, v)$ 
24:     Remove  $\text{best\_candidate}$ 
25:     Update  $v_c \leftarrow \text{UpdateDifference}(D0, D1)$ 
26:      $\text{num\_removals} \leftarrow \text{num\_removals} + 1$ 
27:   end while
28: end if
29: return  $\text{num\_removals}$ 
```

3.4.3 Algorithm Explanation

1. Initialization:

- **Groups Setup:** Divide the dataset D into treated $D1$ and control $D0$ groups based on the treatment indicator T .
- **Aggregates Calculation:** Compute initial sums and sizes for each group and the current ATE difference v_c .
- **Counters Initialization:** Initialize a counter to track the number of removals.
- **Difference Tracking:** Initialize variables to monitor the difference relative to the target v .

2. Greedy Removal Loop:

- **Row Removal Decision:**
 - If $v_c > v + \epsilon$: Remove the tuple with the maximum O from the treated group $D1$.

- If $v_c < v - \epsilon$: Remove the tuple with the maximum O from the control group $D0$.
- **Aggregate Update:** After removal, update the sums and counts for the respective groups.
- **Difference Recalculation:** Recompute the ATE difference v_c .
- **Counter Increment:** Increment the removal counter.
- **Switching Check:** If the product of the previous difference and the current difference is negative ($\text{prev_diff} \times \text{curr_diff} < 0$), a sign change has occurred, indicating that the algorithm has overshoot the desired range. This triggers a switch to the binary search phase for fine-tuning.

3. Switch to Binary Search for Fine-Tuning:

- **Binary Search Loop:**
 - Continue removing the best candidate identified through binary search until the ATE falls within the desired range or no further removals are possible.
- **Row Removal:** Identify and remove the best candidate row that brings the ATE closest to the target.
- **ATE Update:** Recalculate the ATE using the updated dataset.
- **Counter Increment:** Increment the removal counter.

4. Termination:

- The algorithm terminates when the ATE difference v_c is within the desired range or when no further removals can be made without violating group sizes.

3.5 Per Group Removal Algorithm with Confounders

3.5.1 Description of Per Group Removal Algorithm

To handle datasets with confounders, we extend the Greedy Combined Algorithm to the **Per Group Removal Algorithm**. This method ensures that row removals are balanced across different confounder values. Confounder groups refer to the unique combinations of treatment and confounder values, enabling targeted and balanced adjustments.

3.5.2 Pseudocode

Below is the pseudocode for the Per Group Removal Algorithm:

Algorithm 3 Per Group Removal Algorithm

Require: Dataset D with binary attributes T , Confounders W_1, W_2, \dots , and outcome O ; target difference v ; threshold ϵ

Ensure: Number of tuples deleted from D

```
1: Initialize:
2:   num_removals  $\leftarrow 0$ 
3:   sorted_indices  $\leftarrow$  Sort rows by  $O$  for each  $(T, \text{Confounder Combination})$  group
4:   current_ATE  $\leftarrow$  CalculateATE( $D$ )
5: while  $\text{current\_ATE} \notin [v - \epsilon, v + \epsilon]$  and rows exist in groups do
6:   if  $\text{current\_ATE} > v + \epsilon$  then
7:     candidates  $\leftarrow$  Select highest  $O$  from each group
8:   else
9:     candidates  $\leftarrow$  Select lowest  $O$  from each group
10:  end if
11:  best_diff  $\leftarrow \infty$ 
12:  row_to_remove  $\leftarrow$  None
13:  for each candidate in candidates do
14:    new_ATE  $\leftarrow$  SimulateRemove( $D$ , candidate)
15:    diff  $\leftarrow |\text{new\_ATE} - v|$ 
16:    if diff < best_diff then
17:      best_diff  $\leftarrow$  diff
18:      row_to_remove  $\leftarrow$  candidate
19:    end if
20:  end for
21:  if row_to_remove is None then
22:    BREAK
23:  end if
24:  RemoveTuple( $D$ , row_to_remove)
25:  num_removals  $\leftarrow$  num_removals + 1
26:  current_ATE  $\leftarrow$  CalculateATE( $D$ )
27: end while
28: return num_removals
```

3.5.3 Algorithm Explanation

1. Initialization:

- **Counters Initialization:** Initialize a counter to track the total number of removals.
- **Sorted Indices:** Sort the dataset rows by the outcome variable O within each treatment T and confounder W group. Binning may be employed to reduce the number of confounder combinations, thereby lowering computational costs.
- **Initial ATE Calculation:** Compute the initial ATE using the full dataset.

2. Main Removal Loop:

- **Row Removal Decision:**
 - If $\text{current_ATE} > v + \epsilon$: Select the row with the highest O from each confounder group.
 - Else: Select the row with the lowest O from each confounder group.
- **Best Candidate Identification:** Among all candidates, determine the row whose removal results in the smallest difference from the target ATE.

- **Row Removal:** Remove the identified row from the dataset and increment the removal counter.
- **ATE Update:** Recalculate the ATE using the updated dataset.
- **Termination Check:** If the ATE falls within the desired range, terminate the loop.

3.5.4 Implementation Details

The extended algorithm incorporates confounders by dividing the dataset based on the treatment T and confounders W values. Confounder groups refer to the unique combinations of these variables, enabling targeted and balanced row removals. Binning may be applied to confounders with multiple values to reduce the number of confounder groups, thereby decreasing computational complexity and improving the algorithm's scalability.

4 Experiments

This section outlines the experimental setup used to evaluate the effectiveness of the proposed algorithms in minimizing row removals to achieve the desired ATE range.

4.1 Explanation of Epsilon Factor

The **Epsilon Factor** (ϵ) is a parameter used in the experiments to define the acceptable deviation from the target ATE v . It determines the precision of the ATE adjustment by specifying the range $[v - \epsilon, v + \epsilon]$ within which the adjusted ATE must fall. Mathematically, it is defined as:

$$\epsilon = \frac{|\text{Original ATE} - \text{New ATE}|}{\text{Epsilon Factor}}$$

New ATE is decided by randomly removing 10% of the dataset and calculating the partial dataset's ATE. A higher epsilon factor results in a smaller epsilon, demanding more precise adjustments to the ATE. This precision can lead to more computational resources and time being required to achieve the desired ATE within the specified range.

4.2 Experiment Configuration

4.2.1 First Experiment: Without Confounders

- **Database Sizes:** {1,000,000, 2,000,000, 3,000,000, 4,000,000, 5,000,000, 6,000,000, 7,000,000, 8,000,000, 9,000,000, 10,000,000}
- **Number of Trials:** 3 per configuration to ensure statistical significance.
- **Epsilon Adjustment Factors:** {10, 100, 1,000, 10,000, 100,000}

4.2.2 Second Experiment: With Binary Confounder W

- **Database Sizes:** {1,000,000, 2,000,000, 3,000,000, 4,000,000, 5,000,000, 6,000,000, 7,000,000, 8,000,000, 9,000,000, 10,000,000}
- **Number of Trials:** 3 per configuration.
- **Epsilon Adjustment Factors:** {10, 100, 1,000, 10,000}

4.2.3 Third Experiment: Real-World Data

- **Database:** StackOverflow Survey Dataset
- **Size:** 38,090 records
- **Required ATE Modification:** Increase by \$500
- **Epsilon Factor:** 1
- **Row Removal Strategy:** The **Per Group Removal Algorithm** was applied to selectively remove the lowest-paid individuals from appropriate confounder groups to achieve the desired ATE increase with minimal row removal.

4.3 Experimental Procedure

For each combination of dataset size and epsilon factor in the first and second experiments, the following steps were performed:

1. **Dataset Generation:** Synthetic datasets were generated with predefined treatment assignments and outcomes. In the second experiment, a binary confounder W was added.
2. **Row Removal Simulation:** Initially, 10% of rows were removed at random (from $W = 0$ in the second experiment) in order to set the desired ATE.
3. **Algorithm Execution:** The respective algorithms—**Greedy Combined Algorithm**, **Naive Removal Algorithm**, and **Binary Search Removal Algorithm** for the first experiment, and **Per Group Removal Algorithm** for the second and third experiments—were executed to adjust the ATE within the specified range.
4. **Measurement:** The number of tuples removed and the time taken for each algorithm were recorded.
5. **Averaging:** Results were averaged over the three trials for each configuration to account for variability.

5 Results

This section presents the results of the experiments conducted to evaluate the performance of the algorithms across different dataset sizes and epsilon factors.

5.1 First Experiment: Without Confounders

5.1.1 Epsilon Factor = 1,000

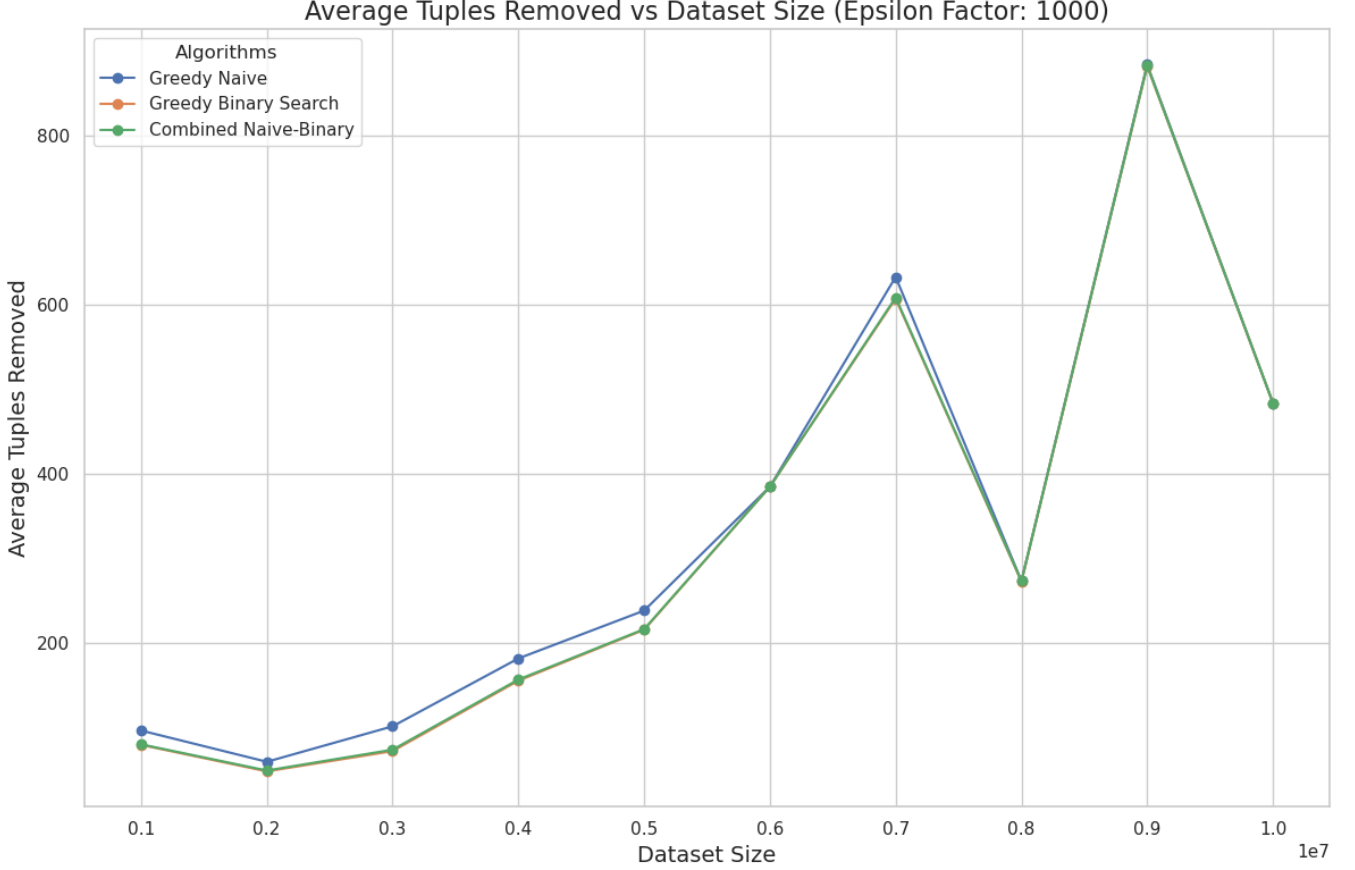


Figure 1: Average Tuples Removed vs. Dataset Size for Epsilon Factor = 1,000 (No Confounder)

Figure 1: Average Tuples Removed vs. Dataset Size for Epsilon Factor = 1,000 (No Confounder) This plot illustrates the average number of tuples removed by the **Greedy Combined Algorithm**, **Naive Removal Algorithm**, and **Binary Search Removal Algorithm** across different dataset sizes when the epsilon factor is set to 1,000. The **Greedy Combined Algorithm** and the **Binary Search Removal Algorithm** consistently removed a similar number of tuples, closely followed by the **Naive Removal Algorithm**, which performed slightly behind. However, the differences among the three algorithms were negligible, indicating that the combined approach effectively matches the precision of the binary search method while maintaining the efficiency of the naive method.

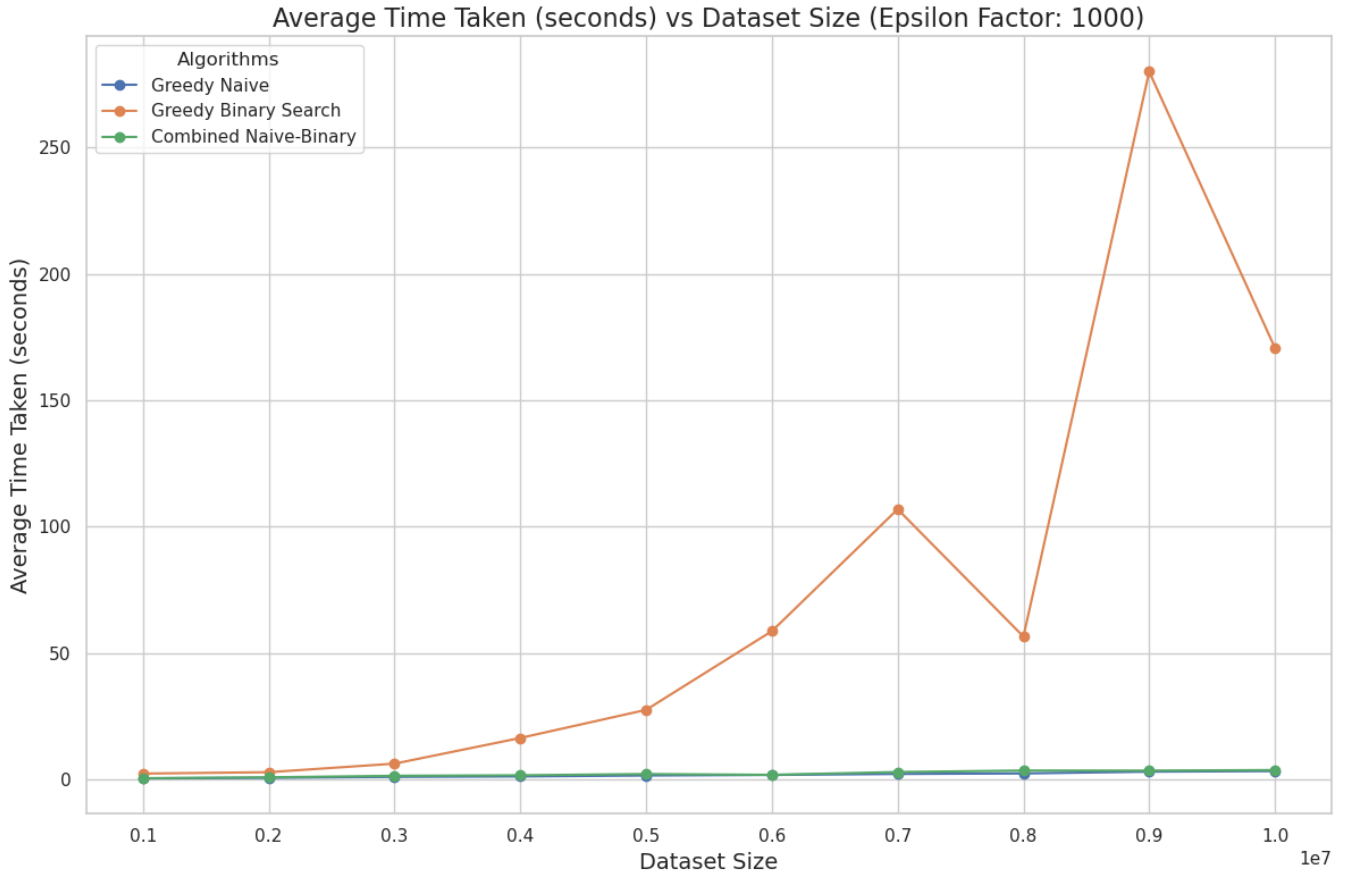


Figure 2: Average Time Taken vs. Dataset Size for Epsilon Factor = 1,000 (No Confounder)

Figure 2: Average Time Taken vs. Dataset Size for Epsilon Factor = 1,000 (No Confounder) This plot depicts the average time taken by the **Greedy Combined Algorithm**, **Naive Removal Algorithm**, and **Binary Search Removal Algorithm** to perform the ATE adjustment across various dataset sizes with an epsilon factor of 1,000. The **Naive Removal Algorithm** and the **Greedy Combined Algorithm** demonstrated significantly faster execution times compared to the **Binary Search Removal Algorithm**, which required more computational time due to its intensive search process.

5.1.2 Summary

Across all dataset sizes ranging from 1,000,000 to 10,000,000, the **Greedy Combined Algorithm** consistently removed a similar number of tuples as the **Binary Search Removal Algorithm**, while matching the **Naive Removal Algorithm** in terms of computational speed. This indicates that the combined approach effectively balances the precision of causal adjustments with operational efficiency, making it a superior choice for large-scale applications.

5.2 Second Experiment: With Binary Confounder W

5.2.1 Epsilon Factor = 100

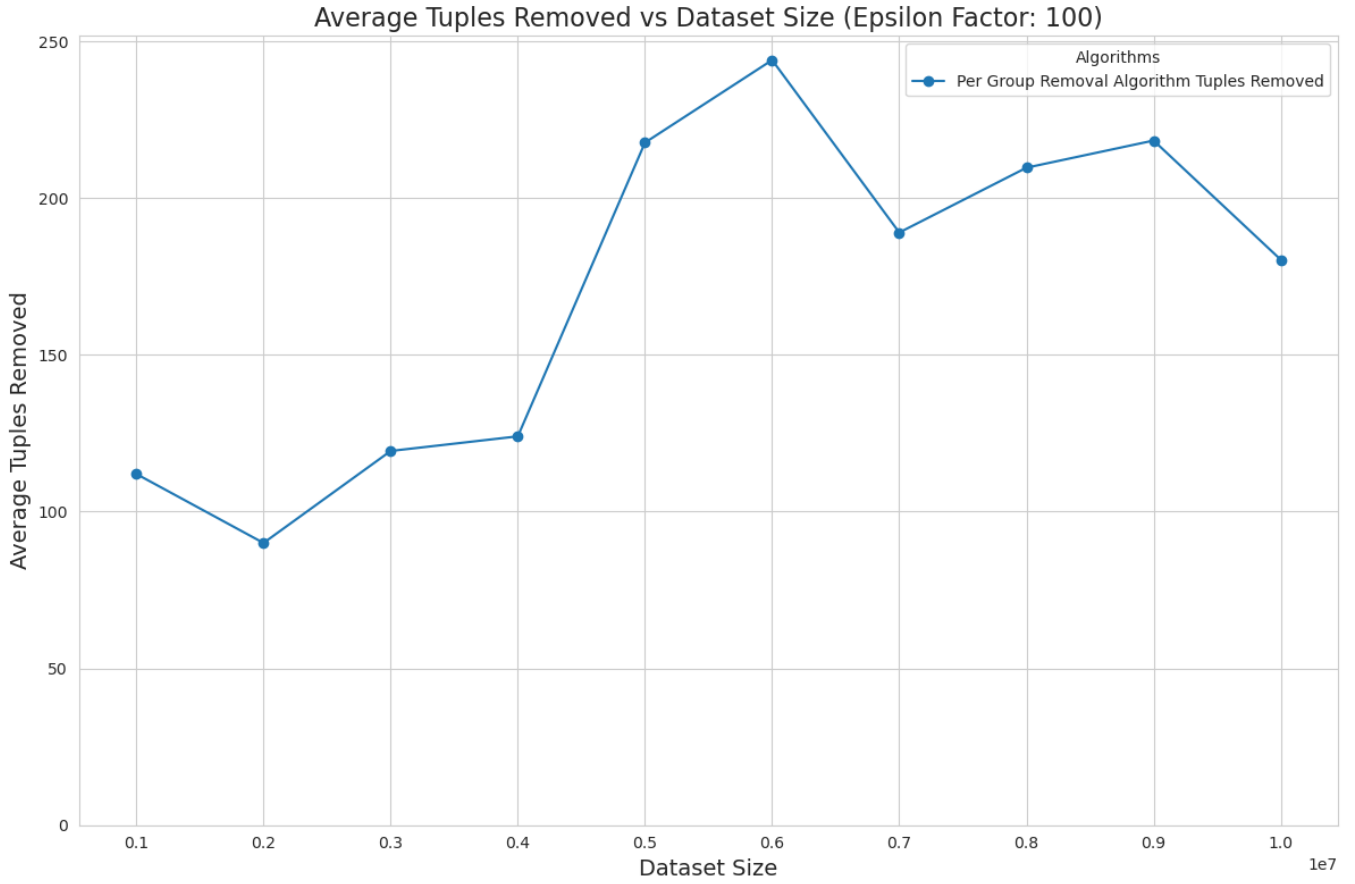


Figure 3: Average Tuples Removed vs. Dataset Size for Epsilon Factor = 100 (Per Group Removal)

Figure 3: Average Tuples Removed vs. Dataset Size for Epsilon Factor = 100 (Per Group Removal) This plot shows the average number of tuples removed by the **Per Group Removal Algorithm** across different dataset sizes when the epsilon factor is 100. The algorithm maintained a consistent number of tuples removed, ranging from 100 to 250. This consistency indicates the algorithm's effectiveness in managing confounders without being significantly affected by the scale of the dataset.

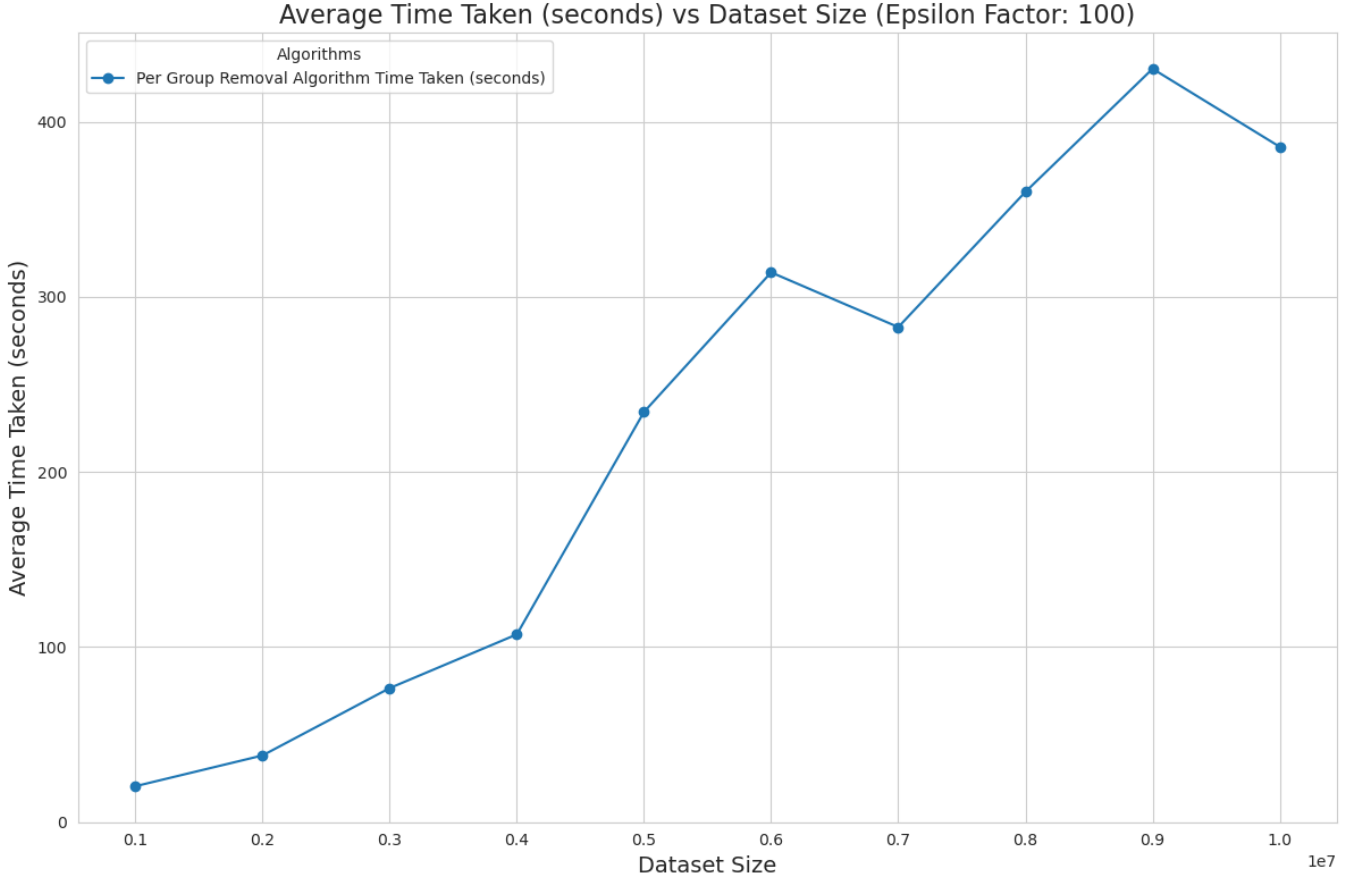


Figure 4: Average Time Taken vs. Dataset Size for Epsilon Factor = 100 (Per Group Removal)

Figure 4: Average Time Taken vs. Dataset Size for Epsilon Factor = 100 (Per Group Removal) This plot presents the average time taken by the **Per Group Removal Algorithm** to adjust the ATE across various dataset sizes with an epsilon factor of 100. The results show an almost linear growth with increasing dataset sizes. This behavior is expected, as larger datasets inherently require more computational effort to process, especially when handling multiple confounder groups.

5.2.2 Summary

In the presence of the binary confounder W , the **Per Group Removal Algorithm** effectively removed rows from both $W = 0$ and $W = 1$ groups to achieve the desired ATE. Specifically, for an epsilon factor of 100, the algorithm removed a consistent number of tuples across all dataset sizes, maintaining dataset integrity while aligning the ATE within the specified range. The time taken by the algorithm scaled almost linearly with dataset size, demonstrating its scalability and efficiency in managing confounders.

5.3 Third Experiment: Real-World Data

5.3.1 Description of Third Experiment

In the third experiment, the **Per Group Removal Algorithm** was applied to a real-world dataset derived from a survey conducted by StackOverflow among its users. The dataset comprises 38,090 records.

5.3.2 Experiment Setup

- **Treatment Variable T :** FormalEducation (Yes/No).

- **Outcome Variable O :** ConvertedSalary.
- **Confounders:** UndergradMajor, Continent, RaceEthnicity.
- **Desired ATE Range:** The ATE was set to increase by \$500 from its original value, with $\epsilon = 1$ to allow a minimal margin of error.
- **Row Removal Strategy:** To increase the ATE as required, the algorithm selectively removed the lowest-paid individuals from the treated groups within each confounder combination. This targeted removal ensures that the ATE increases while minimizing the total number of rows removed.

5.3.3 Results of Third Experiment

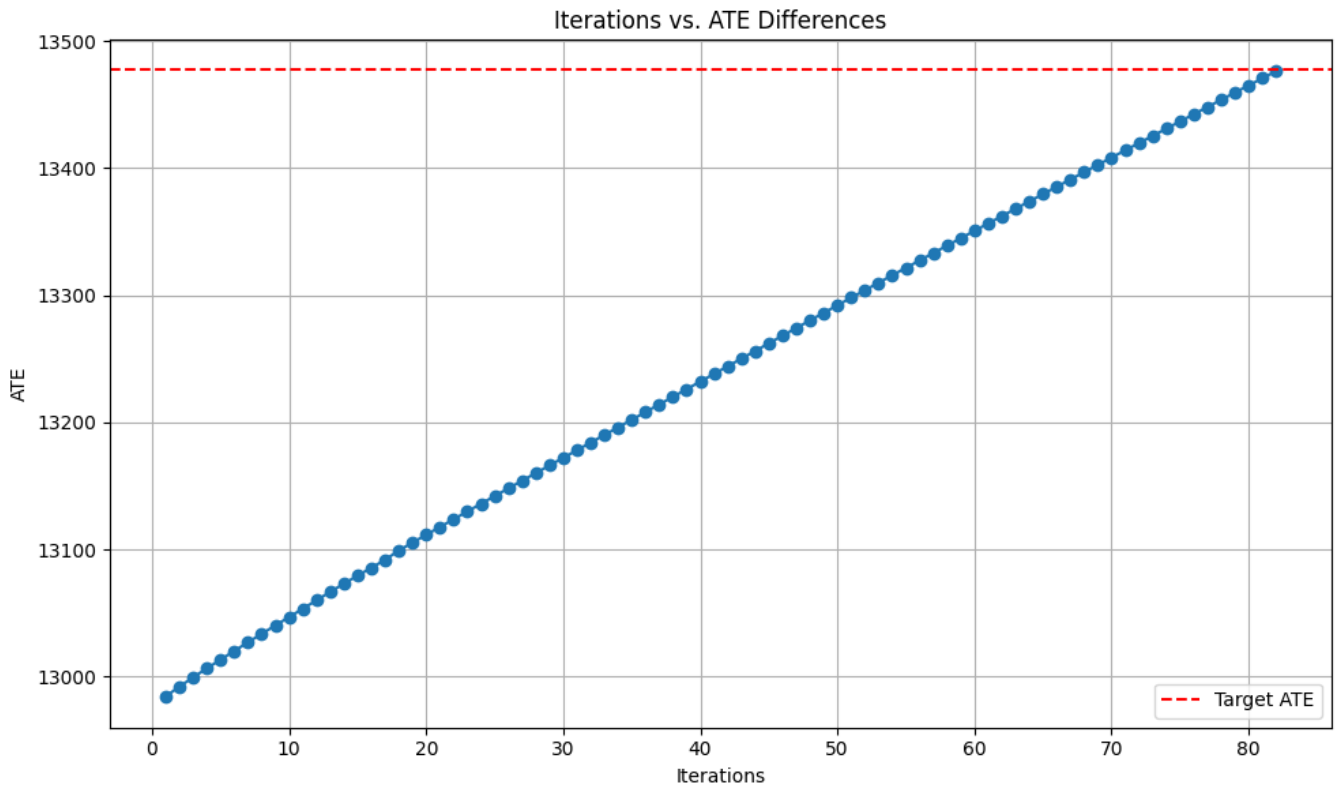


Figure 5: ATE Progression Over Iterations for Real-World Dataset

Plot: Iterations vs. ATE

Figure 5: ATE Progression Over Iterations for Real-World Dataset This plot illustrates the progression of the ATE as the **Per Group Removal Algorithm** iteratively removed rows from the StackOverflow survey dataset. Each iteration represents a row removal, and the corresponding ATE is plotted to show how it approaches the desired range. The plot demonstrates a steady convergence towards the target ATE, highlighting the algorithm’s effectiveness in adjusting the ATE with minimal data loss.

Summary of Third Experiment In the third experiment, the **Per Group Removal Algorithm** was successfully applied to a real-world StackOverflow survey dataset containing 38,090 records. The algorithm aimed to increase the ATE of ConvertedSalary by \$500 within a tight epsilon margin of 1. To achieve this, the algorithm strategically removed the lowest-paid individuals from specific confounder groups, ensuring that the ATE increased while minimizing the total number of rows removed.

Key findings from the experiment include:

- **Total Removals:** 82 tuples were removed to achieve the desired ATE increase.
- **Average Salary of Removed Rows:** \$21,332.10 per year, indicating that the algorithm targeted lower-income individuals to effectively raise the ATE.
- **Groups Affected:** The majority of removals occurred in groups with Formal Education as "Yes," UndergradMajors in "Engineering & Math" and "Other Disciplines," Continent as "North America" and "OC," and RaceEthnicity as "Asian Descent" and "European Descent."
- **Execution Time:** The algorithm completed the ATE adjustment in approximately 145 seconds, demonstrating reasonable computational efficiency given the dataset size.
- **ATE Progression:** The ATE steadily increased towards the target range with each row removal, as depicted in Figure 5, showcasing the algorithm's effectiveness in real-world scenarios.

The experiment confirms that the **Per Group Removal Algorithm** can effectively adjust the ATE in real-world datasets by removing the least impactful rows, thereby preserving data integrity and achieving the desired causal inference with minimal data loss.

6 Discussion

6.1 Analysis of Results

The experimental results indicate that the **Greedy Combined Algorithm** consistently achieves the desired ATE adjustment with minimal data removal and superior time efficiency compared to the **Binary Search Removal Algorithm**. Specifically, across all dataset sizes and epsilon factors in the first experiment, the combined approach balanced the removal of influential rows with computational speed, making it effective for large-scale applications.

In the second experiment, incorporating the binary confounder W and utilizing the **Per Group Removal Algorithm**, the results demonstrate the algorithm’s capability to handle confounders effectively.

The third experiment further validates the algorithms’ practical applicability by applying them to real-world data. The successful adjustment of ATE in the StackOverflow survey dataset shows the methodologies’ efficiency in handling complex data scenarios.

6.2 Scalability

Both algorithms exhibit strong scalability, effectively handling datasets ranging from 1,000,000 to 10,000,000 rows. The time taken increases proportionally with dataset size, however, the **Greedy Combined Algorithm** and **Per Group Removal Algorithm** remain practical even for the largest datasets tested.

6.3 Impact of Confounders

The inclusion of confounders introduces additional complexity in ATE adjustment. The **Per Group Removal Algorithm** handles this by dividing the dataset based on treatment and confounders combinations values, ensuring that row removals are balanced across confounder groups. This approach mitigates potential biases and maintains the dataset’s representativeness, improving the accuracy of ATE estimations.

7 Conclusion

This project successfully developed and evaluated methodologies to minimize row removals in datasets to achieve a desired Average Treatment Effect (ATE). The **Greedy Combined Algorithm** proved to be highly effective in the first experiment, consistently removing a similar number of tuples as the **Binary Search Removal Algorithm** while matching the **Naive Removal Algorithm** in terms of computational speed. This indicates that the combined approach effectively balances the precision of causal adjustments with operational efficiency, making it a superior choice for large-scale applications. In the second experiment, the **Per Group Removal Algorithm** effectively managed confounders, ensuring balanced data removals and maintaining dataset integrity. The successful application of these algorithms to real-world data in the third experiment further validates their practical utility and robustness in complex data scenarios.