

A SYNTHETIC GENERATION

In this section, we introduce a few additional findings on coverage redundant Simpson's paradoxes from a data generative perspective. We first discuss the generation of separate instances of Simpson's paradox in Section A.1. We then introduce the process of generating coverage redundant Simpson's paradoxes in Section A.2. Finally, we summarize the overall data synthesization procedure in Section A.3.

A.1 Generating Simpson's Paradoxes

Recall from Definition 2.2 that an association configuration (AC) $p = (s_1, s_2, X, Y)$ is a Simpson's paradox if:

- (1) $P(Y|s_1) \leq P(Y|s_2)$; and
- (2) $P(Y|s_1 \langle X \rightarrow v \rangle) \geq P(Y|s_2 \langle X \rightarrow v \rangle), \forall v \in \text{Dom}(X)$.

In particular, the frequency statistics $P(Y|s_j)$ for $j = 1, 2$ in condition (1), is obtained by weighted averaging their sub-population frequency statistics in condition (2), where the weights are determined by the relative coverage sizes of each sub-population. Specifically, we have that:

$$\begin{aligned} P(Y|s_j) &= \sum_{v \in \text{Dom}(X)} \frac{|\text{cov}(s_j \langle X \rightarrow v \rangle)|}{|\text{cov}(s_j)|} \cdot P(Y|s_j \langle X \rightarrow v \rangle) \\ &= Q(s_j|X) \cdot P(s_j|Y, X)^\top, \end{aligned}$$

where

$$Q(s_j|X) = \left[\frac{|\text{cov}(s_j \langle X \rightarrow v_1 \rangle)|}{|\text{cov}(s_j)|}, \frac{|\text{cov}(s_j \langle X \rightarrow v_2 \rangle)|}{|\text{cov}(s_j)|}, \dots, \frac{|\text{cov}(s_j \langle X \rightarrow v_{|\text{Dom}(X)|} \rangle)|}{|\text{cov}(s_j)|} \right]$$

is the sample distribution of s_j partitioned under X , and

$$\begin{aligned} P(s_j|Y, X) &= [P(Y|s_j \langle X \rightarrow v_1 \rangle), P(Y|s_j \langle X \rightarrow v_2 \rangle), \\ &\quad \dots, P(Y|s_j \langle X \rightarrow v_{|\text{Dom}(X)|} \rangle)] \end{aligned}$$

is the frequency statistics of s_j 's sub-populations partitioned by X . With this, we can rephrase Definition 2.2 by substituting the terms, that (s_1, s_2, X, Y) is a Simpson's paradox if:

$$Q(s_1|X) \cdot P(s_1|Y, X)^\top < Q(s_2|X) \cdot P(s_2|Y, X)^\top; \text{ and} \quad (1)$$

$$P(s_1|Y, X)[j] > P(s_2|Y, X)[j], 1 \leq j \leq |\text{Dom}(X)|. \quad (2)$$

Therefore, the essence of generating an instance of Simpson's paradox is to (a) find the set of sub-population frequency statistics $P(s_1|Y, X)$ and $P(s_2|Y, X)$ that satisfy inequality (2); and (b) solve for the sample distributions $Q(s_1|X)$ and $Q(s_2|X)$ that satisfy inequality (1). We discuss each in the following paragraphs.

Sub-population Frequency Statistics. The first step of generating an instance of Simpson's paradox is to ensure that each sub-population in s_2 , partitioned by X , has an frequency statistics value

smaller than its sibling sub-population in s_1 (inequality (2)). A simple pattern that achieves this is:

$$\begin{aligned} P(s_2|Y, X)[1] &< P(s_1|Y, X)[1] < \\ P(s_2|Y, X)[2] &< P(s_1|Y, X)[2] < \\ &\dots \\ P(s_2|Y, X)[|\text{Dom}(X)|] &< \\ P(s_1|Y, X)[|\text{Dom}(X)|]. \end{aligned} \quad (3)$$

This pattern ensures that for any value v_j ($1 \leq j \leq |\text{Dom}(X)|$), we have $P(s_1|Y, X)[j] > P(s_2|Y, X)[j]$, satisfying inequality (2). With this pattern, we can now focus on finding sample distributions that satisfy inequality (1).

Sample Distributions. Given the sub-population frequency statistics pattern established above, we now need to solve for sample distributions $Q(s_1|X)$ and $Q(s_2|X)$ that satisfy inequality (1). To achieve this, we formulate the problem as a quadratic program:

$$\begin{aligned} &\text{minimize} \quad \sum_{j=1}^2 \left\| Q(s_j|X) - \frac{1}{|\text{Dom}(X)|} \mathbf{1} \right\|_2^2 \\ &\text{subject to} \quad (i) \quad \sum_{k=1}^{|\text{Dom}(X)|} Q(s_j|X)[k] = 1, \quad j = 1, 2; \\ &\quad (ii) \quad Q(s_j|X)[k] > 0, \quad j = 1, 2 \\ &\quad \quad \text{and } 1 \leq k \leq |\text{Dom}(X)|; \\ &\quad (iii) \quad Q(s_2|X) \cdot P(s_2|Y, X)^\top \\ &\quad \quad > Q(s_1|X) \cdot P(s_1|Y, X)^\top. \end{aligned} \quad (4)$$

We incorporate inequality (1) as a linear constraint specified in condition (iii). Moreover, the objective function of the QP aims to minimize the squared distance between each sample distribution and the uniform distribution $\frac{1}{|\text{Dom}(X)|} \mathbf{1}$ to promote uniformity. Our synthetic generator also supports optimizing towards other distribution patterns such as normal or Zipfian distributions to better mimic the statistics of real-world data.

Synthesis of Simpson's Paradox. To summarize, generating an instance of Simpson's paradox requires establishing sub-population frequency statistics satisfying inequality (2) and solving for sample distributions that produce the reversal effect outlined in inequality (1). Algorithm 7 formalizes this generation procedure, taking as input an AC (s_1, s_2, X, Y) and crucially, a size parameter U that controls the number of records covered by the generated Simpson's paradox. The algorithm proceeds in three main steps: first, it generates the sub-population frequency statistics following the pattern in Equation (3); second, it solves the quadratic program in (4) to obtain optimal sample distributions; and third, it populates the output table with $2 \cdot U$ records for the Simpson's paradox (s_1, s_2, X, Y) , distributing records across sub-populations according to the sample distributions and assigning label values according to probabilities given by the sub-population frequency statistics.

A.2 Realizing Coverage Redundancies

Having established a method for generating individual instances of Simpson's paradox, we now discuss creating (coverage) redundant Simpson's paradoxes. Building on the discussions in Section 3.1 and

Algorithm 7 Generate non-redundant Simpson's paradox.

Input: An AC (s_1, s_2, X, Y) , paradox size U

Output: Data records T for the Simpson's paradox (s_1, s_2, X, Y)

- 1: Obtain the sub-population frequency statistics $P(s_1|Y, X)$, $P(s_2|Y, X)$ following Equation (3);
 - 2: Obtain the sample distributions $Q(s_1|X)$, $Q(s_2|X)$ by solving the quadratic program (QP) in (4);
 - 3: // Populate data records following the obtained sample distributions and sub-population aggregate statistics
 - 4: **for** each $1 \leq k \leq |\text{Dom}(X)|$ and each $j \in [1, 2]$ **do**
 - 5: // Find the number of records for each sub-population
 - 6: Let $U_{j,k} \leftarrow U \cdot Q(s_j|X)[k]$;
 - 7: // Assign labels according to sub-population aggr. stats.
 - 8: Add $U_{j,k}$ copies of $s_j \langle X \rightarrow v_k \rangle$ as records to T and assign $U_{j,k} \cdot P(s_j|Y, X)[k]$ of them with $(Y = 1)$;
 - 9: **return** T .
-

Definition 3.8, (coverage) redundancies fundamentally arise when distinct populations have identical coverage within the data. To systematically realize (coverage) redundancies, we must understand the conditions under which populations would share the same coverage, as this property serves as the foundation for realizing both sibling child and separator equivalences. To this end, we remark on the following proposition.

PROPOSITION A.1 (IMPOSSIBILITY OF COVERAGE IDENTICALITY). *Let T be a base table with n categorical attributes $\{X_1, \dots, X_n\}$. If the set of unique records in T corresponds exactly to all possible combinations of attribute values (i.e., the complete Cartesian product $\prod_{i=1}^n \text{Dom}(X_i)$), then no two distinct populations of T share the same coverage.*

PROOF. Assume, by contradiction, that there exist two distinct populations s and s' with $\text{cov}(s) = \text{cov}(s')$. Since $s \neq s'$, let X_{k_0} be an attribute for which $s[k_0] \neq s'[k_0]$. Without loss of generality, assume $s[k_0] = v$ for a fixed value $v \in \text{Dom}(X_{k_0})$ and $s'[k_0] = *$.

Let us consider two records, r_1 and r_2 , which are identical on all attributes except X_{k_0} . For r_1 , let $r_1.X_{k_0} = v$, and for r_2 , let $r_2.X_{k_0} = v'$ where $v' \neq v$ and $v' \in \text{Dom}(X_{k_0})$. For all other attributes X_j where $j \neq k_0$, if $s[j] \neq *$, then $r_1.X_j = r_2.X_j = s[j]$. Observe that both records r_1 and r_2 must exist in T because T contains the complete Cartesian product of all attribute domains.

Now, r_1 is covered by both s and s' . However, record r_2 is covered by s' (because $s'[k_0] = *$) but not by s (because $s[k_0] = v$ but $r_2.X_{k_0} = v', v' \neq v$).

This means $r_2 \in \text{cov}(s')$ but $r_2 \notin \text{cov}(s)$, which contradicts our assumption that $\text{cov}(s) = \text{cov}(s')$. Therefore, no two distinct populations can share the same coverage when T contains the complete Cartesian product of all attribute domains. \square

The contrapositive of Proposition A.1 implies that populations sharing identical coverage can only exist when the dataset contains a proper subset of the complete Cartesian product of attribute domains. Therefore, to facilitate the generation of coverage redundant Simpson's paradoxes, we impose a size threshold t that is significantly smaller than $\prod_{i=1}^n |\text{Dom}(X_i)|$ to constrain the number of unique records in the generated dataset.

Having established the condition for realizing populations with identical coverage, we now proceed to develop methods for realizing each of the three types of coverage equivalences: sibling child equivalence, separator equivalence, and statistic equivalence.

Sibling Child Equivalence. Suppose we have a set of data records T that produces a Simpson's paradox $p_1 = (s_1, s_2, X, Y)$, where $s_1 = s \langle X_0 \rightarrow u_1 \rangle$ and $s_2 = s \langle X_0 \rightarrow u_2 \rangle$ are siblings from a common parent s . To realize sibling equivalence, our goal is update the records in T such that they also produce another Simpson's paradox $p_2 = (s'_1, s'_2, X, Y)$, where $s'_1 = s' \langle X'_0 \rightarrow v_1 \rangle$ and $s'_2 = s' \langle X'_0 \rightarrow v_2 \rangle$ are siblings from a common parent s' , and that p_2 is sibling child equivalent to p_1 . According to Definition 3.8, sibling child equivalence requires $\text{cov}(s_1) = \text{cov}(s'_1)$ and $\text{cov}(s_2) = \text{cov}(s'_2)$. Based on the relationship between (s, s_1, s_2) and (s', s'_1, s'_2) , we have three scenarios:

- (1) **Scenario 1:** $s \neq s'$, $X_0 = X'_0$, and $\{u_1, u_2\} = \{v_1, v_2\}$. In this case, sibling child equivalence is achieved by ensuring $\text{cov}(s) = \text{cov}(s')$. To this, for each categorical attribute X_k where $s[k] \neq *$ or $s'[k] \neq *$, we update every record r in $\text{cov}(s)$ (within T), such that $r.X_k = s[k]$ or $r.X_k = s'[k]$. If $s[k] = s'[k] = *$, then no update is needed for $r.X_k$.
- (2) **Scenario 2:** $s = s'$, $X_0 \neq X'_0$, and $\{u_1, u_2\} \neq \{v_1, v_2\}$. In this case, to achieve sibling child equivalence, we establish a one-to-one mapping $f : \{u_1, u_2\} \mapsto \{v_1, v_2\}$ such that $f(u_1) = v_1$ and $f(u_2) = v_2$. For each record r in T , we set $r.X'_0 = f(r.X_0)$ when $r.X_0 \in \{u_1, u_2\}$. This ensures $\text{cov}(s \langle X_0 \rightarrow u_k \rangle) = \text{cov}(s \langle X'_0 \rightarrow f(u_k) \rangle)$ for $k = 1, 2$.
- (3) **Scenario 3:** $s \neq s'$, $X_0 \neq X'_0$, and $\{u_1, u_2\} \neq \{v_1, v_2\}$. This combines the previous scenarios. To achieve sibling child equivalence, we first ensure $\text{cov}(s) = \text{cov}(s')$ as in Scenario 1, then establish the one-to-one mapping as in Scenario 2.

Example A.2. Consider the data records in a slightly perturbed version of Table 2 where attribute values in D are randomized. Supposed the perturbed Table 2 is populated as a result of generating the Simpson's paradox $p_1 = ((*, b_1, *, *), (*, b_2, *, *), A, Y_1)$. To create a sibling equivalent Simpson's paradox $p_2 = ((*, *, *, d_1), (*, *, *, d_2), A, Y_1)$, we apply both Scenarios 1 and 2.

For Scenario 1, the parent population is $(*, *, *, *)$ for both p_1 and p_2 , which are identical, so no adjustment to records in the perturbed table is needed.

For Scenario 2, we define a one-to-one mapping $f : \{b_1, b_2\} \rightarrow \{d_1, d_2\}$ where $f(b_1) = d_1$ and $f(b_2) = d_2$. We then update each record in the perturbed table with $B = b_1$ gets $D = d_1$ and each record with $B = b_2$ gets $D = d_2$. This establishes that $\text{cov}((*, b_1, *, *)) = \text{cov}((*, *, *, d_1))$ and $\text{cov}((*, b_2, *, *)) = \text{cov}((*, *, *, d_2))$.

In this way, we make p_1 and p_2 to be sibling equivalent, as verified in Example 3.2. \square

Algorithm 8 formalizes this process of generating sibling-child-equivalent Simpson's paradoxes.

Separator Equivalence. Recall from Proposition 3.3 and Definition 3.8, Simpson's paradoxes $p_1 = (s_1, s_2, X_1, Y)$ and $p_2 = (s_1, s_2, X'_1, Y)$ are separator equivalent if there exists a one-to-one mapping f between $\text{Dom}(X_1)$ and $\text{Dom}(X'_1)$ such that for every

Algorithm 8 Realizing sibling child equivalence.

Input: Data records T producing the Simpson's paradox $p_1 = (s_1, s_2, X, Y)$ where $s_1 = s\langle X_0 \rightarrow u_1 \rangle$ and $s_2 = s\langle X_0 \rightarrow u_2 \rangle$, sibling populations (s_1, s_2) where $s'_1 = s\langle X'_0 \rightarrow v_1 \rangle$ and $s'_2 = s\langle X'_0 \rightarrow v_2 \rangle$

Output: Updated data records T producing a sibling-child-equivalent Simpson's paradox $p_2 = (s'_1, s'_2, X, Y)$

```

1: // Scenario 1: Ensure  $\text{cov}(s) = \text{cov}(s')$ 
2: for each record  $r \in T$  s.t.  $r \in \text{cov}(s)$  do
3:   for each attribute  $X_k$  s.t.  $s[k] \neq *$  or  $s'[k] \neq *$  do
4:      $\text{Set } r.X_k \leftarrow s[k]$  if  $s[k] \neq *$ , else  $r.X_k \leftarrow s'[k]$ ;
5: // Scenario 2: Establish the one-to-one mapping
6: Establish the mapping  $f$  where  $f(u_j) = v_j$  for  $j = 1, 2$ ;
7: for each record  $r \in T$  do
8:    $\text{Set } r.X'_0 \leftarrow f(r.X_0)$  if  $r.X_0 \in \{u_1, u_2\}$ .
```

$v \in \text{Dom}(X_1)$ and $j = 1, 2$,

$$\text{cov}(s_j\langle X_1 \rightarrow v \rangle) = \text{cov}(s_j\langle X'_1 \rightarrow f(v) \rangle).$$

To achieve this, for every record r in T , we set $r.X'_1 = f(r.X_1)$, which we formalize the process in Algorithm 9.

Example A.3. Consider a perturbed version of Table 2 where attribute values in C are initially randomized. Suppose the perturbed Table 2 is populated as a result of generating the Simpson's paradox $p_1 = ((*, b_1, *, *), (*, b_2, *, *), A, Y_1)$. To create a separator equivalent Simpson's paradox $p_2 = ((*, b_1, *, *), (*, b_2, *, *), C, Y_1)$, we establish a one-to-one mapping f between $\text{Dom}(A) = a_1, a_2$ and $\text{Dom}(C) = c_1, c_2$ where $f(a_1) = c_1$ and $f(a_2) = c_2$. We then update each record in the perturbed table so that whenever $A = a_1$, we set $C = c_1$, and whenever $A = a_2$, we set $C = c_2$. This ensures that $\text{cov}((*, b_1, *, *)\langle A \rightarrow a_k \rangle) = \text{cov}((*, b_1, *, *)\langle C \rightarrow c_k \rangle)$ and $\text{cov}((*, b_2, *, *)\langle A \rightarrow a_k \rangle) = \text{cov}((*, b_2, *, *)\langle C \rightarrow c_k \rangle)$ for $k \in \{1, 2\}$. As verified in Example 3.4, p_1 and p_2 are separator equivalent. \square

Algorithm 9 Realizing separator equivalence.

Input: Set of data records T producing the Simpson's paradox (s_1, s_2, X_1, Y) , a separator attribute X'_1

Output: Updated set of data records T producing a separator equivalent Simpson's paradox $p_2 = (s_1, s_2, X'_1, Y)$

```

1: Let  $f : \text{Dom}(X_1) \mapsto \text{Dom}(X'_1)$  be the one-to-one map;
2: for each record  $r \in T$  do
3:    $\text{Set } r.X'_1 \leftarrow f(r.X_1)$ .
```

Statistic Equivalence. Recall from Proposition 3.5 and Definition 3.8, Simpson's paradoxes $p_1 = (s_1, s_2, X, Y_2)$ and $p_2 = (s_1, s_2, X, Y'_2)$ are statistic equivalent if for each s_j ($j = 1, 2$) $P(Y_2|s_j) = P(Y'_2|s_j)$, and for every value $v \in \text{Dom}(X)$, $P(Y_2|s_j\langle X \rightarrow v \rangle) = P(Y'_2|s_j\langle X \rightarrow v \rangle)$. To achieve this, we simply ensure that each record has identical values for both label attributes Y_2 and Y'_2 . We formalize this process in Algorithm 10.

Example A.4. Consider a perturbed version of Table 2 where attribute values in Y_2 are initially randomized. Suppose the perturbed Table 2 is populated as a result of generating the Simpson's paradox $p_1 = ((*, b_1, *, *), (*, b_2, *, *), A, Y_1)$. To create a statistic equivalent

Simpson's paradox $p_2 = ((*, b_1, *, *), (*, b_2, *, *), A, Y_2)$, we update each record in the perturbed table so that $Y_2 = Y_1$ for all records. This ensures that $P(Y_1|s_j) = P(Y_2|s_j)$ and $P(Y_1|s_j\langle A \rightarrow a_k \rangle) = P(Y_2|s_j\langle A \rightarrow a_k \rangle)$ for $j \in \{1, 2\}$ and $k \in \{1, 2\}$. As verified in Example 3.6, p_1 and p_2 are statistic equivalent. \square

Algorithm 10 Realizing statistic equivalence.

Input: Data records T producing the Simpson's paradox $p_1 = (s_1, s_2, X, Y_2)$, a label attribute Y'_2

Output: Updated data records T producing a statistic equivalent Simpson's paradox $p_2 = (s_1, s_2, X, Y'_2)$

```

1: for each record  $r \in T$  do
2:    $\text{Set } r.Y'_2 \leftarrow r.Y_2$ .
```

A.3 Data Generation Workflow**Algorithm 11** Generate redundant Simpson's paradoxes.

Input: Categorical attributes $\{X_i\}_{i=1}^n$, label attributes $\{Y_j\}_{j=1}^m$, size threshold $t \ll \prod_{i=1}^n |\text{Dom}(X_i)|$

Output: Data table T (initially empty)

```

1: Let  $R \leftarrow \emptyset$  to collect the set of unique data records;
2: Let  $P \leftarrow \emptyset$  to collect the set of generated Simpson's paradoxes;
3: while  $|R| < t$  do
4:   Let  $p_1 = (s_1, s_2, X_1, Y_3)$  be an AC not in  $P$ ;
5:   // Step 1: Generate distinct Simpson's paradox
6:   Populate  $T'$  for the Simpson's paradox  $p_1$  using Alg. 7;
7:   // Step 2: Introduce coverage redundancies
8:   // Sibling child equivalence
9:   Apply Alg. 8 to  $T'$  to create a sibling-child-equivalent Simpson's paradox  $p_2 = (s'_1, s'_2, X_1, Y_2)$ ;
10:  // Separator equivalence
11:  Apply Alg. 9 to  $T'$  to create a separator equivalent Simpson's paradox  $p_3 = (s_1, s_2, X'_1, Y_3)$ ;
12:  // Statistic equivalence
13:  Apply Alg. 10 to  $T'$  to create a statistic equivalent Simpson's paradox  $p_4 = (s_1, s_2, X_1, Y'_2)$ ;
14:  Add  $p_1, p_2, p_3$ , and  $p_4$  to  $P$ ;
15:  Add  $T'$  to  $T$  and unique records of  $T'$  to  $R$ ;
16: return  $T$ .
```

Building upon the techniques established in Sections A.1 and A.2, we formulate a systematic approach for synthetic data generation that integrates both individual Simpson's paradox generation and coverage redundancy realization. The process employs a two-phase strategy: first generating distinct instances of Simpson's paradoxes (Section A.1), then systematically introducing coverage redundancies through sibling child, separator, and statistics equivalences (Section A.2). These phases are iterated until reaching a specified threshold $t \ll \prod_{i=1}^n |\text{Dom}(X_i)|$ of unique records populated, which per Proposition A.1 ensures the dataset contains populations with identical coverage necessary for redundancy.

Algorithm 11 formalizes this process, taking categorical attributes $\{X_i\}_{i=1}^n$, label attributes $\{Y_j\}_{j=1}^m$, and the size threshold t as input, and producing a synthetic data table containing groups of (coverage) redundant Simpson's paradoxes.