



Bottom-up Standardization for Data Preparation



Eugenie Lai, Yuze Lou, Brit Youngmann, Michael Cafarella

Background & Motivation

Data preparation has been seen as "janitor work" yet essential in data-to-insight pipelines. Now we have **more access to more data**, but **challenges** including:

- People have limited technical and domain expertise
- People who have the skills keep writing one-off programs

Core observation:

- Scripts **embed thought process** and knowledge
- There is a **latent space** of data preprocessing steps in script corpus

Recently, automating data preparation scripts improves many aspects of the pipeline, including **data quality**, **research reproducibility**, and user productivity. We propose an approach to **automatically standardize scripts**.

Problem Definition

Informally, our goal is to **standardize** user input scripts s^* using a **sequence of transformations** informed by the **script corpus** S .

Specifically, the output script \hat{s} should (1) leverage the domain knowledge embedded in the script collection; (2) be easily consumed by the user

We use relative entropy to define a score for our standardization goal.

$$RE(P||Q) = \sum_{x \in X} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

Our objective is to minimize relative entropy, where

- P represents the distribution of data preprocessing steps in the **input** s^*
- Q represents the distribution of data preprocessing steps in the **collective knowledge** S

Intuition: (1) The latent space of dataset-specific data preprocessing steps in the script corpus manifests itself in the distribution of its data preprocessing steps; (2) with the optimal sequence of transformations, s^* becomes closer to S with control over the semantics.

We use table jaccard to measure how much script semantics have changed.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}, \text{ where } A, B \text{ are sets.}$$

Related Work

Roughly three topics:

- **Predictive data transformation:** step/pipeline recommendations (tldr: by-example approach, program synthesis)
- **Program synthesis** for data preparation: automatically generates code for given user intent (e.g., Copilot, AutoPandas)
- Exploratory data analysis (**EDA**) recommendation

Contributions

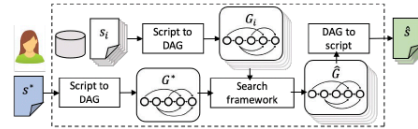
New things we do:

- We define a **new problem of data preparation program standardization**
- We propose a novel approach to transform programs while allowing **control over script semantics**
- We empirically show that our approach mostly outperforms learned methods (e.g., Copilot, Codex)

Framework & Data Model

At high level, our approach

- (1) uses collection $S = \{s_0, s_1, s_2, \dots, s_n\}$, transform each into its DAG representation, and obtain a latent space of data preprocessing steps $Q(x)$
- (2) takes an input s^* and get its DAG representation G^*
- (3) uses the search-based framework with $G = \{G_0, G_1, G_2, \dots, G_n\}$ and G^* as input to find a sequence of transformations \hat{F} that transforms G^* to \hat{G}
- (4) returns \hat{s} as the script representation of \hat{G}



Choice of script-to-AST and AST-to-DAG:

Improving script quality requires a way to take it apart and transform it. The challenge is that the improved script should be free from execution errors. AST already has many desired properties.

- We construct DAG $G = (V, E)$ from AST for implicit flows
- V : PL construct (e.g., function call); E : data flow

Algorithm 1: A meta-level greedy search framework

```

Input:  $S$  (Script collection),  $s^*$  (input script),  $K$  (budget),  $\theta$  (data correctness threshold),  $\alpha$  (early checking switch)
Output:  $\hat{s}$  ( $\hat{s}$  a standardized version of  $s^*$ )
1  $C = \{s^*\}$ ; /* Set of all candidates. */
2 while not StoppingCriteria( $C$ ) do
3    $C' = \emptyset$ ; /* Set of new candidates. */
4   foreach  $s \in C$  do
5      $\mathcal{F} = \text{GetSteps}(s, S)$ ; /*  $\mathcal{F}$ : A ranked set of next steps based on RE score. */
6      $C' = \text{GetTopKBears}(C', s', K, \alpha, \mathcal{F})$ 
7    $C = C'$ 
8  $C = \text{VerifyConstraints}(C, \theta)$ ;
9  $\hat{s} = \text{GetTopCandidate}(C)$ ;
10 return  $\hat{s}$ ;

```

Methods & Optimization

Challenges:

- Huge **search space**, due to (1) possible sequences grow exponentially as the steps increase; (2) # of unique V is large
- Execution **errors** and data **semantics** are difficult to estimate in the process. Since we only care about them at the end, they are allowed to be off the given thresholds during the process, which poses challenges on pruning the candidate sequences

We propose two **GetTopKBears()**.

- **Difference:** diversity measure
- Define heuristics to capture operator functional similarity

Algorithm 2: GetTopKBears

```

Input:  $C$  (Current candidates),  $s$  (a script),  $K$  (budget),  $\alpha$  (early checking switch),  $\mathcal{F}$  (a ranked set of next steps)
Output:  $C'$  (Updated candidates)
1 valid = True;
2  $C' = C$ ;
3 foreach  $f \in \mathcal{F}$  do
4    $s' = \text{ApplyStep}(s, f)$ ;
5   if  $\alpha$  then valid = CheckIfExecutes( $s'$ );
6   if Score( $s'$ ) < min $s \in C$  Score( $s$ ) or  $|C| \leq K$  then
7     if valid then  $C' = \text{UpdateKBears}(C', s', K)$ ;
8 return  $C'$ ;

```

Algorithm 3: GetDiverseTopKBears

```

Input:  $C$  (Current candidates),  $s$  (a script),  $K$  (budget),  $\alpha$  (early checking switch),  $\mathcal{F}$  (a ranked set of next steps)
Output:  $C'$  (Updated candidates)
1  $C' = C$ ;
2  $\mathcal{F}_{diff} = \text{ClusterSteps}(\mathcal{F})$ ;
3 foreach  $\mathcal{F}_{diff}$  do
4    $C' = \text{GetTopKBears}(C', s', \frac{K}{|\mathcal{F}_{diff}|}, \alpha, \mathcal{F})$ ;
5 return  $C'$ ;

```

Experiments & Discussion

Setup: Two datasets from Kaggle; for each of ~50 scripts, use one as s^* , the rest as S ; repeat for transformation seq length in [2, 4, 8, 16]; beam in [1, 3]; data threshold in [1.0, .95, .8, .7].

% improvement = (start_score - end_score) / start_score

