# Algorithm for Dimensional Sampling in HoloClean

Step 1 identifies the dimensional attributes which aide to our sampling methodology. It helps us to eliminate training points which generate similar features during featurization and fail to add any new information to the existing model during training.

---

**Algorithm 1** Identify Dimensional Columns

---

1: **procedure** $get\_dim\_columns()$
2:      **if** $domain\_knowledge$ **is** $true$ **then**
3:          $dim\_attr \leftarrow import(user\_input)$
4:      **else**
5:          **for** a **in** attributes **do**
6:              **if** a **in** (list of columns in CD or FD) **then**
7:                  $dim\_attr[\,].append(a)$
8:              **end if**
9:          **end for**
10:      **end if**
11:      **return** $dim\_attr[\,]$
12: **end procedure**

---

Step 2 shows the modification made to the domain engine, which iterates through each cell in the dataset to generate the domain. At this stage, we identify which cells should be sampled based on our dimensional attributes. We do not remove any cells here to ensure the initial dataset is entirely available to make the train-test split.

---

**Algorithm 2** Modify Domain Engine to mark sampling cells

---

     **procedure** $generate\_domain()$
2:      $dim\_duplicate[\,] \leftarrow duplicated(dataset, [group\_by = dimensional\_attr])$
     $samp\_ind[\,] \leftarrow empty\_list()$
4:      **for** t **in** tuples **do**
         **for** a **in** attributes **do**
6:              **if** (a **in** dim_attr) AND (t **in** dim_duplicate) **then**
                 $samp\_ind.append(0)$
8:              **else**
                 $samp\_ind.append(1)$
10:              **end if**
             $cell\_domain \leftarrow$ `<generate domain on each cell with existing code>`
12:              $cell\_domain.append\_attribute(samp\_ind[\,])$
         **end for**
14:      **end for**
     **return** $cell\_domain[\,]$
16: **end procedure**

---

Step 3 shows the modification made to the data featurizer. Here the training set is generated from the initial input dataset. Existing code filters out cells which are erroneous and for which the weak label is NULL. We add one extra layer to also filter out cells which are not going to aide to the training model, based on our dimensional attributes.

---

**Algorithm 3** Modify Dataset Featurizer to filter training points

---

     **procedure** $generate\_weak\_labels()$
         **for** cd **in** cell_domain **do**
3:          cd.filter(weak_labels is True)
         cd.filter(samp_ind is False)
     **end for**
6:      $labels \leftarrow$ `<Generate training points with cd as input with existing code>`
     **return** $labels$
     **end procedure**

---