

1 **A Assigning Infill**

2 Each transformation category has some default infill convention to serve as precursor to ML infill.
3 For cases where a user wishes to override defaults assignments can be passed in the `assigninfill`
4 parameter. ML infill is applied for columns not otherwise assigned, which default can be deactivated
5 with the `MLinfill` parameter. When `MLinfill` is deactivated, columns not explicitly assigned
6 will have infill per the default initialization associated with a transformation category. Here we
7 demonstrate deactivating the ML infill default and assigning infill types zero infill to column1 and
8 ML infill just to column2.

9 Not shown, if the data includes a label set or other features such as an index column appropriate for ex-
10 clusion from ML infill basis, they should be designated with `labels_column` or `trainID_column`.

```
11 assigninfill = {'zeroinfill' : ['column1'],  
12                'MLinfill'   : ['column2'] }  
13  
14 train, train_ID, labels, \  
15 val, val_ID, val_labels, \  
16 test, test_ID, test_labels, \  
17 postprocess_dict = \  
18 am.automunge(df_train,  
19              MLinfill = False,  
20              assigninfill = assigninfill)
```

21 Note that the column headers can be assigned in `assigninfill` using the received column headers
22 to apply consistent infill to all sets derived from an input column, or may alternatively be assigned
23 using the returned column headers with transformation suffix appenders to assign infill to distinct
24 returned columns, which take precedence.

25 **B ML Infill Parameters**

26 The default ML infill architecture is a Scikit-Learn random forest with default parameters. Alternate
27 auto ML options are currently available as CatBoost, FLAML, and AutoGluon. Parameters can be
28 passed to the models with `ML_cmdnd`.

29 First we'll demonstrate applying ML infill with the CatBoost library. Note that we can either defer to
30 the library default parameters or also pass parameters to the model initializations or fit operations.
31 Here we also demonstrate assigning a particular GPU device number.

```
32 ML_cmdnd = {'autoML_type' : 'catboost'}  
33  
34 #GPU device assignment takes place in model initialization  
35 ML_cmdnd.update({'MLinfill_cmdnd' :  
36                {'catboost_classifier_model' :  
37                  {'task_type' : 'GPU', 'devices' : 0 },  
38                  'catboost_regressor_model' :  
39                  {'task_type' : 'GPU', 'devices' : 0 }}})  
40  
41 train, train_ID, labels, \  
42 val, val_ID, val_labels, \  
43 test, test_ID, test_labels, \  
44 postprocess_dict = \  
45 am.automunge(df_train,  
46              MLinfill = True,  
47              ML_cmdnd = ML_cmdnd)
```

48 The FLAML library is also available. Here we demonstrate setting a training time budget (in seconds)
49 for each imputation model.

```
50 ML_cmdnd = {'autoML_type' : 'flaml'}
```

```

51
52 ML_cmdnd.update({'MLinfill_cmdnd' :
53                 {'flaml_classifier_fit' : {'time_budget' : 60 },
54                 'flaml_regressor_fit' : {'time_budget' : 60 }}})
55
56 train, train_ID, labels, \
57 val, val_ID, val_labels, \
58 test, test_ID, test_labels, \
59 postprocess_dict = \
60 am.automunge(df_train,
61              MLinfill = True,
62              ML_cmdnd = ML_cmdnd)

```

63 As another demonstration, here is an example of applying the AutoGluon library for ML infill and also
64 applying the `best_quality` option which causes AutoGluon to train extra models for the aggregated
65 ensembles. (Note this will likely result in large disk space usage, especially when applying to every
66 column, so recommend saving this for final production if at all.)

```

67 ML_cmdnd = {'autoML_type' : 'autogluon'}
68
69 ML_cmdnd.update({'MLinfill_cmdnd' :
70                 {'AutoGluon' :
71                 {'presets' : 'best_quality' }}})
72
73 train, train_ID, labels, \
74 val, val_ID, val_labels, \
75 test, test_ID, test_labels, \
76 postprocess_dict = \
77 am.automunge(df_train,
78              MLinfill = True,
79              ML_cmdnd = ML_cmdnd)

```

80 To be complete, here we'll demonstrate passing parameters to the Scikit-Learn random forest models.
81 Note that for random forest there are built in methods to perform grid search or random search
82 hyperparameter tuning when parameters are passed as lists or distributions instead of static figures.
83 Here we'll demonstrate performing tuning of the `n_estimators` parameter (which otherwise would
84 default to 100).

```

85 ML_cmdnd = {'autoML_type' : 'randomforest'}
86
87 #by passing random forest parameters as a list
88 #each imputation model will perform a grid search
89
90 ML_cmdnd.update({'MLinfill_cmdnd' :
91                 {'RandomForestClassifier' :
92                 {'n_estimators' : [100, 222, 444] },
93                 'RandomForestRegressor' :
94                 {'n_estimators' : [100, 222, 444] }}})
95
96 train, train_ID, labels, \
97 val, val_ID, val_labels, \
98 test, test_ID, test_labels, \
99 postprocess_dict = \
100 am.automunge(df_train,
101              MLinfill = True,
102              ML_cmdnd = ML_cmdnd)

```

103 **C Broader Impacts**

104 The following discussions are somewhat speculative in nature. At the time of this writing Automunge
105 has yet to establish what we would consider a substantial user base and there may be a bias towards
106 optimism at play in how we have been proceeding, which we believe is our sole leverage of bias.

107 From an ethical standpoint, we believe the potential benefits of our platform far outweigh any
108 negative aspects. We have sought to optimize the `postmunge(.)` function for speed, used as a proxy
109 for computational efficiency and carbon intensity. As a rule of thumb, processing times for equivalent
110 data in the `postmunge(.)` function, such as could be applied to streams of data in inference, have
111 shown to operate on the order of twice the speed of initial preparations in the `automunge(.)` function,
112 although for some specific transforms like those implementing string parsing that advantage may
113 be considerably higher. While the overhead may prevent achieving the speed of directly applying
114 manually specified transformations to a dataframe, the `postmunge(.)` speed gets close to manual
115 transformations with increasing data size.

116 We believe too that the impact to the machine learning community of a formalized open source
117 standard to tabular data preprocessing could have material benefits to ensuring reproducibility of
118 results. There for some time has been a gap between the wide range of open source frameworks
119 for training neural networks in comparison to options for prerequisites of data pipelines. I found
120 some validation for this point from the tone of the audience Q&A at a certain 2019 NeurIPS keynote
121 presentation by the founder of a commercial data wrangling package. In fact it may be considered
122 a potential negative impact of this research in the risk to commercial models of such vendors, as
123 Automunge's GNU GPL 3.0 license coupled with patent pending status on the various inventions
124 behind our library (including parsed categoric encodings, family tree primitives, ML infill, and etc.)
125 will preclude commercial platforms offering comparable functionality. We expect that the benefits to
126 the machine learning community in aggregate will far outweigh the potential commercial impacts
127 to this narrow segment. Further, benefits of automating machine learning derived infill to missing
128 data may result in a material impact to the mainstream data science workflow. That old rule of thumb
129 often thrown around about how 80% of a machine learning project is cleaning the data may need to
130 be revised to a lower figure.

131 Regarding consequence of system failure, it should be noted that Automunge is an industry agnostic
132 toolset, with intention to establish users across a wide array of tabular data domains, potentially
133 ranging from the trivial to mission critical. We recognize that with this exposure comes additional
134 scrutiny and responsibility. Our development has been performed by a professional engineer and
135 we have sought to approach validations, which has been an ongoing process, with a commensurate
136 degree of rigor.

137 Our development has followed an incremental and one might say evolutionary approach to systems
138 engineering, with frequent and sequential updates as we iteratively added functionality and transforms
139 to the library within defined boundaries of the data science workflow. The intent has always been to
140 transition to a more measured pace at such time as we may establish a more substantial user base.

141 **D Intellectual Property Disclaimer**

142 Automunge is released under GNU General Public License v3.0. Full license details available on
143 GitHub. Contact available via (anonymized). Copyright (C) 2021 - All Rights Reserved. Patent
144 Pending, applications (anonymized)