

## scripts/ml\_dm\_visualization.py

```
42 # %% Load data into dataframes, set roles and create views
43 import getml
44
45 import cora.helpers as helpers
46
47 getml.engine.launch()
48 getml.engine.set_project("cora_visualization")
49
50 conn = getml.database.connect_mysql(
51     host="db.relational-data.org",
52     dbname="CORA",
53     port=3306,
54     user="guest",
55     password="relational",
56 )
57
58 df_paper = getml.data.DataFrame.from_db(name="df_paper", table_name="paper", conn=conn)
59 df_cites = getml.data.DataFrame.from_db(name="df_cites", table_name="cites", conn=conn)
60 df_content = getml.data.DataFrame.from_db(name="df_content", table_name="content", conn=conn)
61
62 df_paper.set_role("paper_id", getml.data.roles.join_key)
63 df_paper.set_role("class_label", getml.data.roles.categorical)
64
65 df_cites.set_role(["cited_paper_id", "citing_paper_id"], getml.data.roles.join_key)
66
67 df_content.set_role("paper_id", getml.data.roles.join_key)
68 df_content.set_role("word_cited_id", getml.data.roles.categorical)
69
70 v_paper_multi_varget = getml.data.make_target_columns(df_paper, "class_label")
71
72 # %% Define the data model
73 import getml.data.relationship as rel
74 from getml.data.placeholder import Placeholder
75
76 dm = getml.data.DataModel(population=Placeholder(name="ph_population", roles=
v_paper_multi_varget.roles))
77
78 dm.add(Placeholder(name="ph_cites", roles=df_cites.roles))
79 dm.add(Placeholder(name="ph_cites", roles=df_cites.roles))
80
81 dm.add(Placeholder(name="ph_content", roles=df_content.roles))
82 dm.add(Placeholder(name="ph_paper", roles=df_paper.roles))
83
84 dm["ph_population"].join(dm["ph_cites"][0], on=("paper_id", "cited_paper_id"))
85 dm["ph_cites"][0].join(dm["ph_content"], on=("citing_paper_id", "paper_id"))
86 dm["ph_cites"][0].join(dm["ph_paper"], on=("citing_paper_id", "paper_id"), relationship=
rel.many_to_one)
87
88 dm["ph_population"].join(dm["ph_cites"][1], on=("paper_id", "citing_paper_id"))
89 dm["ph_cites"][1].join(dm["ph_content"], on=("cited_paper_id", "paper_id"))
90 dm["ph_cites"][1].join(dm["ph_paper"], on=("cited_paper_id", "paper_id"), relationship=
rel.many_to_one)
91
92 dm["ph_population"].join(dm["ph_content"], on="paper_id")
```

```

93
94 # %% Define a split and create a container connecting the dataframes with the data model
95 split = getml.data.split.random(train=0.7, test=0.3, seed=0)
96 container = getml.data.Container(population=v_paper_multi_varget, split=split)
97 container.add(ph_cites=df_cites, ph_content=df_content, ph_paper=df_paper)
98 container.freeze()
99
100 # %% Define a pipeline and fit it to the data
101 fast_prop = getml.feature_learning.FastProp(
102     loss_function=getml.feature_learning.loss_functions.CrossEntropyLoss,
103     num_threads=10,
104     aggregation=getml.feature_learning.aggregations.fastprop.Minimal,
105     num_features=600,
106 )
107
108 pipe = getml.pipeline.Pipeline(
109     tags=["fast_prop_mapping"],
110     preprocessors=[getml.preprocessors.Mapping()],
111     data_model=dm,
112     feature_learners=[fast_prop],
113     predictors=[getml.predictors.XGBoostClassifier(objective="binary:logistic")],
114 )
115
116 pipe.fit(container.train)
117 prob_vecs_1_vs_all_test = pipe.predict(container.test)
118
119 predicted_labels_test = helpers.probs_1_vs_all_to_label_via_argmax(
120     prob_vecs_1_vs_all_test, class_labels=df_paper.class_label.unique()
121 )
122 gt_labels_test = df_paper[split == "test"].class_label.to_numpy()
123 accuracy = (gt_labels_test == predicted_labels_test).sum() / len(gt_labels_test)

```