

```
@tracer(cat_col = ['personal_status_and_sex'], numerical_col = ['age'])
def german_pipeline_normal(f_path_1='../data/german_titled_split_1.csv', f_path_2='../data/german_titled_split_2.csv'):
    # load data
    dataSplit1 = pd.read_csv(f_path_1, index_col = 0)
    dataSplit2 = pd.read_csv(f_path_2, index_col = 0)

    # join
    data = dataSplit1.merge(dataSplit2, on='identifier')

    # drop first col
    data.drop(data.columns[0], axis=1, inplace = True)

    # projection
    data = data[['duration_in_month', 'credit_his', 'credit_amt', 'preset_emp', 'personal_status_and_sex', 'guarantors', 'present_residence', 'property', 'age', 'label']]

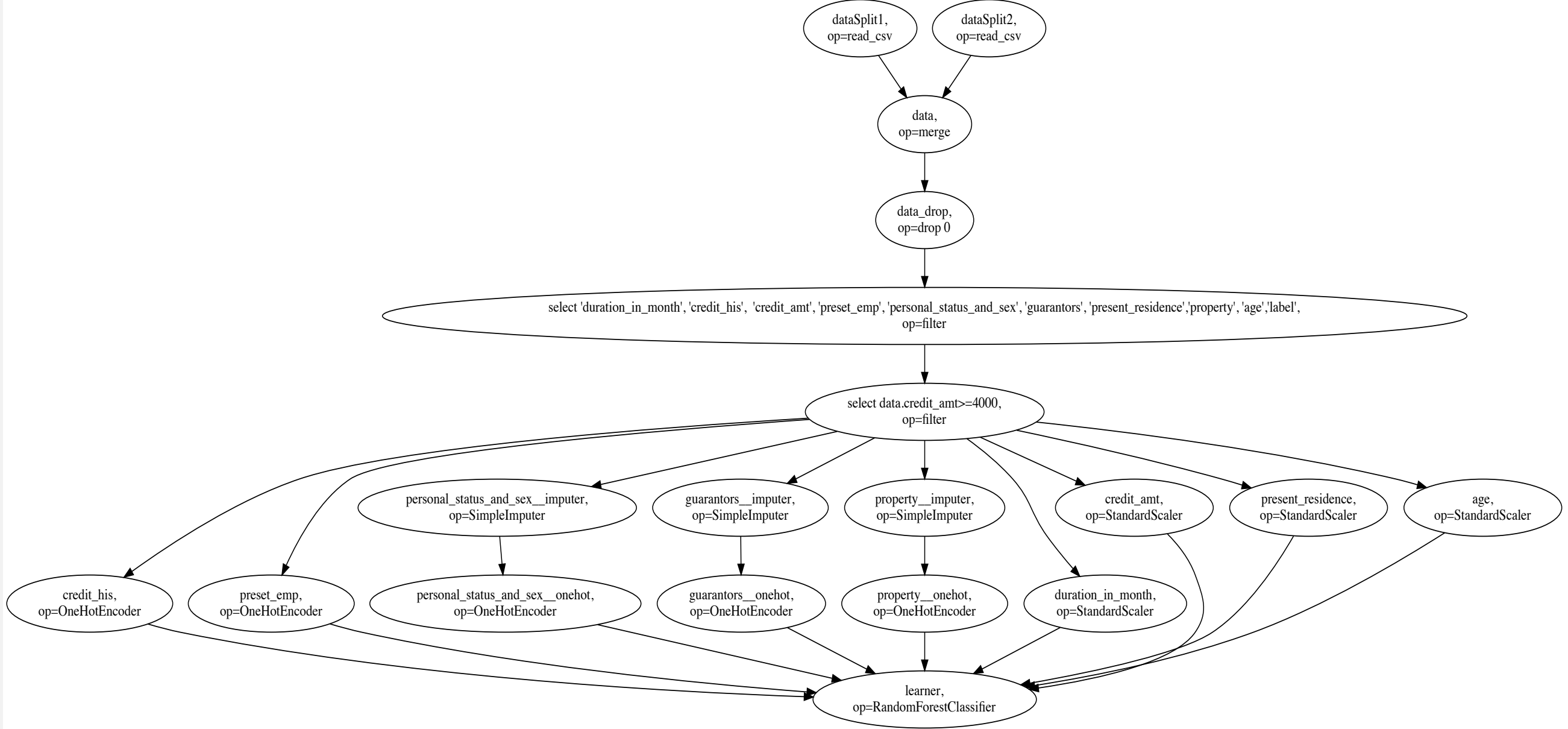
    # filtering
    data = data.loc[(data.credit_amt>=4000)]

    #start sklearn pipeline
    one_hot_and_impute = Pipeline([
        ('imputer', SimpleImputer(strategy='most_frequent')),
        ('onehot', OneHotEncoder())
    ])

    featurizer = ColumnTransformer(transformers=[
        ('onehot', OneHotEncoder(), ['credit_his', 'preset_emp']),
        ('impute_onehot', one_hot_and_impute, ['personal_status_and_sex', 'guarantors', 'property']),
        ('std_scaler', StandardScaler(), ['duration_in_month', 'credit_amt', 'present_residence', 'age'])
    ])

    pipeline = Pipeline([
        ('features', featurizer),
        ('learner', RandomForestClassifier())
    ])

    return pipeline
```



Start Pandas Opeation

Inpected dataSplit1 = pd.read_csv(f_path_1, index_col = 0)

Changes in numerical features!

	count	missing_count	median	mad	range
age	-inf	-inf	-inf	-inf	-inf

Inpected dataSplit2 = pd.read_csv(f_path_2, index_col = 0)

Inpected data = dataSplit1.merge(dataSplit2, on='identifier')

Inpected data.drop(data.columns[0], axis=1, inplace = True)

Inpected data = data[['duration_in_month', 'credit_his', 'credit_amt', 'preset_emp', 'personal_status_and_s
ex', 'guarantors', 'present_residence', 'property', 'age', 'label']]

Changes in numerical features!

	count	missing_count	median	mad	range
age	-754.0	0.0	0.5	0.7413	-1.0

Changes in categorical features!

	missing_count	num_class	class_count	class_percent
personal_status_and_sex	0.0	0.0	{'A93': -384, 'A92': -251, 'A91': -37, 'A94': -82}	{'A93': 0.1187, 'A92': -0.0702, 'A91': 0.0028, 'A94': -0.0513}

Inpected data = data.loc[(data.credit_amt>=4000)]

Start Sklearn Pipeline

Operations SimpleImputer on personal_status_and_sex

Operations OneHotEncoder on personal_status_and_sex

Changes in categorical features!

	personal_status_and_sex
missing_count	0
num_class	-2
class_count	{0.0: 233, 1.0: 13}
class_percent	{0.0: 0.9472, 1.0: 0.0528}

Operations StandardScaler on age

Changes in numerical features!

	age
count	0.0000
missing_count	0.0000
median	-33.7344
mad	-10.1331
range	-50.1208