

```

def create_features(df, df_test=None):
    X = df.copy()
    y = X.pop("SalePrice")
    mi_scores = make_mi_scores(X, y)

    # Combine splits if test data is given
    #
    # If we're creating features for test set
    # predictions, we should
    # use all the data we have available. After
    # creating our features,
    # we'll recreate the splits.
    if df_test is not None:
        X_test = df_test.copy()
        X_test.pop("SalePrice")
        X = pd.concat([X, X_test])

    # Lesson 2 - Mutual Information
    X = drop_uninformative(X, mi_scores)

    # Lesson 3 - Transformations
    X = X.join(mathematical_transforms(X))
    X = X.join(interactions(X))
    X = X.join(counts(X))
    # X = X.join(break_down(X))
    X = X.join(group_transforms(X))

```

```
# Lesson 4 - Clustering
# X = X.join(cluster_labels(X, cluster_features, n_clusters=20))
# X = X.join(cluster_distance(X, cluster_features, n_clusters=20))
```

```
# Lesson 5 - PCA
X = X.join(pca_inspired(X))
# X = X.join(pca_components(X, pca_features))
# X = X.join(indicate_outliers(X))
```

```
X = label_encode(X)
```

```
# Reform splits
```

```
if df_test is not None:
    X_test = X.loc[df_test.index, :]
    X.drop(df_test.index, inplace=True)
```

```
# Lesson 6 - Target Encoder
```

```
encoder = CrossFoldEncoder(MEstimateEncoder, m=1)
```

```
X = X.join(encoder.fit_transform(X, y, cols=["MSSubClass"]))
```

```
if df_test is not None:
    X_test = X_test.join(encoder.transform(X_test))
```



```

    encoder = CrossFoldEncoder(MEstimateEnc
oder, m=1)
    X = X.join(encoder.fit_transform(X, y,
cols=["MSSubClass"]))
    if df_test is not None:
        X_test = X_test.join(encoder.transf
orm(X_test))

    if df_test is not None:
        return X, X_test
    else:
        return X

df_train, df_test = load_data()
X_train = create_features(df_train)
y_train = df_train.loc[:, "SalePrice"]

score_dataset(X_train, y_train)

```

Out[23]:

0.13863986787521657

```
X_train = create_features(df_train)
y_train = df_train.loc[:, "SalePrice"]

xgb_params = dict(
    max_depth=6,          # maximum depth of
    # each tree - try 2 to 10
    learning_rate=0.01,   # effect of each
    # tree - try 0.0001 to 0.1
    n_estimators=1000,    # number of trees
    # (that is, boosting rounds) - try 1000 to 800
    # 0
    min_child_weight=1,   # minimum number
    # of houses in a leaf - try 1 to 10
    colsample_bytree=0.7, # fraction of fea
    # tures (columns) per tree - try 0.2 to 1.0
    subsample=0.7,        # fraction of ins
    # tances (rows) per tree - try 0.2 to 1.0
    reg_alpha=0.5,        # L1 regularizati
    # on (like LASSO) - try 0.0 to 10.0
    reg_lambda=1.0,       # L2 regularizati
    # on (like Ridge) - try 0.0 to 10.0
    num_parallel_tree=1,  # set > 1 for boo
    # sted random forests
)
```

```
xgb = XGBRegressor(**xgb_params)
score_dataset(X_train, y_train, xgb)
```