

# Machine Learning with the Experts: School Budgets

## 2). Creating a simple first model:

### a). Setting up a train test split

```
# Create the new DataFrame: numeric_data_only
```

```
numeric_data_only = df[NUMERIC_COLUMNS].fillna(-1000)
```

```
# Get labels and convert to dummy variables: label_dummies
```

```
label_dummies = pd.get_dummies(df[LABELS])
```

```
# Create training and test sets
```

```
X_train, X_test, y_train, y_test = multilabel_train_test_split(numeric_data_only, label_dummies, size=0.2, seed=123)
```

```
# Print the info
```

```
print("X_train info:")
```

```
print(X_train.info())
```

```
print("\nX_test info:")
```

```
print(X_test.info())
```

```
print("\ny_train info:")
```

```
print(y_train.info())
```

```
print("\ny_test info:")
```

```
print(y_test.info())
```

```
<script.py> output:
```

```
X_train info:
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
Int64Index: 1040 entries, 198 to 101861
```

```
Data columns (total 2 columns):
```

```
FTE      1040 non-null float64
```

```
Total   1040 non-null float64
```

**b). Training a Model**

```
# Import classifiers
```

```
from sklearn.linear_model import LogisticRegression
```

```
from sklearn.multiclass import OneVsRestClassifier
```

```
# Create the DataFrame: numeric_data_only
```

```
numeric_data_only = df[NUMERIC_COLUMNS].fillna(-1000)
```

```
# Get labels and convert to dummy variables: label_dummies
```

```
label_dummies = pd.get_dummies(df[LABELS])
```

```
# Create training and test sets
```

```
X_train, X_test, y_train, y_test = multilabel_train_test_split(numeric_data_only, label_dummies, size=0.2, seed=123)
```

```
# Instantiate the classifier: clf
```

```
clf = OneVsRestClassifier(LogisticRegression())
```

```
# Fit the classifier to the training data
```

```
clf.fit(X_train,y_train)
```

```
# Print the accuracy
```

```
print("Accuracy: {}".format(clf.score(X_test, y_test)))
```

<script.py> output:

Accuracy: 0.0

**c). Use your model to predict values on holdout data**

**# Instantiate the classifier: clf**

```
clf = OneVsRestClassifier(LogisticRegression())
```

**# Fit it to the training data**

```
clf.fit(X_train, y_train)
```

**# Load the holdout data: holdout**

```
holdout = pd.read_csv('HoldoutData.csv', index_col=0)
```

**# Generate predictions: predictions**

```
predictions = clf.predict_proba(holdout[NUMERIC_COLUMNS].fillna(-1000))
```

**d). Writing out your results to a csv for submission****# Generate predictions: predictions****predictions = clf.predict\_proba(holdout[NUMERIC\_COLUMNS].fillna(-1000))****# Format predictions in DataFrame: prediction\_df****prediction\_df = pd.DataFrame(columns=pd.get\_dummies(df[LABELS].columns),  
 index=holdout.index,  
 data=predictions)****# Save prediction\_df to csv****prediction\_df.to\_csv('predictions.csv')****# Submit the predictions for scoring: score****score = score\_submission(pred\_path='predictions.csv')****# Print score****print('Your model, trained with numeric data only, yields logloss score: {}'.format(score))**

**e). Writing out your result to a csv submission****# Generate predictions: predictions****predictions = clf.predict\_proba(holdout[NUMERIC\_COLUMNS].fillna(-1000))****# Format predictions in DataFrame: prediction\_df****prediction\_df = pd.DataFrame(columns=pd.get\_dummies(df[LABELS]).columns, index=holdout.index, data=predictions)****# Save prediction\_df to csv****prediction\_df.to\_csv('predictions.csv')****# Submit the predictions for scoring: score****score = score\_submission(pred\_path='predictions.csv')****# Print score****print('Your model, trained with numeric data only, yields logloss score: {}'.format(score))**

&lt;script.py&gt; output:

Your model, trained with numeric data only, yields logloss score: 1.9067227623381413

**f). Creating a bag-of-words in scikit-learn**

```
# Import CountVectorizer
```

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
# Create the token pattern: TOKENS_ALPHANUMERIC
```

```
TOKENS_ALPHANUMERIC = '[A-Za-z0-9]+(?=\s+)'
```

```
# Fill missing values in df.Position_Extra
```

```
df.Position_Extra.fillna('',inplace=True)
```

```
# Instantiate the CountVectorizer: vec_alphanumeric
```

```
vec_alphanumeric = CountVectorizer(token_pattern=TOKENS_ALPHANUMERIC)
```

```
# Fit to the data
```

```
vec_alphanumeric.fit(df.Position_Extra)
```

```
# Print the number of tokens and first 15 tokens
```

```
msg = "There are {} tokens in Position_Extra if we split on non-alpha numeric"
```

```
print(msg.format(len(vec_alphanumeric.get_feature_names())))
```

```
print(vec_alphanumeric.get_feature_names()[:15])
```

<script.py> output:

There are 123 tokens in Position\_Extra if we split on non-alpha numeric

['1st', '2nd', '3rd', 'a', 'ab', 'additional', 'adm', 'administrative', 'and', 'any', 'art', 'assessment', 'assistant', 'asst', 'athletic']

**g). Combining text column for tokenization**

```
# Define combine_text_columns()
```

```
def combine_text_columns(data_frame, to_drop=NUMERIC_COLUMNS + LABELS):
```

```
    """ converts all text in each row of data_frame to single vector """
```

```
    # Drop non-text columns that are in the df
```

```
    to_drop = set(to_drop) & set(data_frame.columns.tolist())
```

```
    text_data = data_frame.drop(to_drop, axis=1)
```

```
    # Replace nans with blanks
```

```
    text_data.fillna("", inplace=True)
```

```
    # Join all text items in a row that have a space in between
```

```
    return text_data.apply(lambda x: " ".join(x), axis=1)
```

**h).What is in a token?****# Import the CountVectorizer****from sklearn.feature\_extraction.text import CountVectorizer****# Create the basic token pattern****TOKENS\_BASIC = '\\S+(?=\\s+)'****# Create the alphanumeric token pattern****TOKENS\_ALPHANUMERIC = '[A-Za-z0-9]+(?=\\s+)'****# Instantiate basic CountVectorizer: vec\_basic****vec\_basic = CountVectorizer(token\_pattern=TOKENS\_BASIC)****# Instantiate alphanumeric CountVectorizer: vec\_alphanumeric****vec\_alphanumeric = CountVectorizer(token\_pattern=TOKENS\_ALPHANUMERIC)****# Create the text vector****text\_vector = combine\_text\_columns(df)****# Fit and transform vec\_basic****vec\_basic.fit\_transform(text\_vector)****# Print number of tokens of vec\_basic****print("There are {} tokens in the dataset".format(len(vec\_basic.get\_feature\_names())))****# Fit and transform vec\_alphanumeric****vec\_alphanumeric.fit\_transform(text\_vector)****# Print number of tokens of vec\_alphanumeric****print("There are {} alpha-numeric tokens in the  
dataset".format(len(vec\_alphanumeric.get\_feature\_names())))**

&lt;script.py&gt; output:

There are 1405 tokens in the dataset

There are 1117 alpha-numeric tokens in the dataset