

- Can we predict whether an officer

- Data Structures: Numpy and

- Data Extraction: DataGrip (Pos

- ```
from sklearn.model_selecti
```

[illegible]

The Invisible Institute classified police officers into different groups or cohorts based on a topic modeling analysis, as a result, we have a baseline set of target data for every officer in the dataset.

Based on the data\_officercrew table, we selected and aggregated data about officers such as the average number of coaccusals, complaint percentile, and allegation severity to produce a table of approximately 20,000 officers. Given the classification of each officer into a group (1 for crew, 2 for community, and 3 for unaffiliated), we present the potential to identify what features of an officer best predict their membership.

```
# change the datatype in the dataframe
df['leakent_id'] = df['leakent_id'].astype('category')
df['leakent_id'].cat.rename('leakent_id', inplace=True)
```

```
df['office_id'] = df['office_id'].astype('category', copy=False)
df.dtypes
```

|                |         |
|----------------|---------|
| avg_coaccusals | float64 |
| avg_coaccusals | float64 |

```
In [ ]: # all minmaxscaler + avg_allegation_severity to scale the data
mapper_features = DataFeaturizer([
    ('avg_coaccusals', sklearn.preprocessing.MinMaxScaler()),
    ('avg_years_on_force_at_incident', sklearn.preprocessing.MinMaxScaler()),
    ('avg_age_at_incident', sklearn.preprocessing.MinMaxScaler()),
    ('gender', sklearn.preprocessing.LabelBinarizer()),
    ('avg_complaint_percentile', sklearn.preprocessing.MinMaxScaler()),
    ('avg_disciplined_count', sklearn.preprocessing.MinMaxScaler()),
```

- average years on force at incident

- average complaint percentile
- average allegation severity
- maximum allegation severity

After dropping missing data, we mapped feature and target data with SKLearn MinMax Scaler and fit\_transform. The results of the

```
Out[ ]: array([[0.006, 0.436, 0.538, ..., 0.612, 0. , 0.085],
               [0.017, 0.427, 0.432, ..., 0.72 , 0. , 0.269],
               [0.004, 0.691, 0.7 , ..., 0.49 , 0. , 0.213],
               ...,
               [0.009, 0.528, 0.591, ..., 0.452, 0.1 , 0.256],
               [0.05 , 0.453, 0.478, ..., 0.348, 0.857, 0.142],
               [0.005, 0.523, 0.63 , ..., 0. , 0. , 0.003]])
```

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

```
Out[ ]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
In [ ]: # make predictions of cohorts
prediction = linreg.predict(X_val)
y_valnp = y_val.to_numpy()
predictionnp = np.around(prediction)
```

```
In [ ]: # coefficient of the linear model
linreg.coef_

Out[ ]: array([ -0.4657824 ,  0.24479679,  0.55500357, -0.09472194, -1.21751588,
                0.4029604 , -0.72255634])

In [ ]: # train a logistic regression on the data to check difference
logreg = linear_model.LogisticRegression(multi_class='ovr')
```

```
In [ ]: predictionlog = logreg.predict(X_val)
        metrics.accuracy_score(y_valnp, predictionlog)

Out[ ]: 0.7152180361910412
```

On the other hand, the number of years on force and age at incident appeared to have less effect on the prediction. This runs counter to an initial assumption that maturity might have an impact on misconduct.

**The current model has an accuracy of about 70%.**

```
In [ ]: # read the officer and cohort data into pandas dataframe
url = 'https://raw.githubusercontent.com/NorthwesternData-Sci-Seminar/Invisible-Institute-Chicago-Rep
ort-Collaboration-Public/master/The%20Spectacular%20Sailors/Checkpoint_4/src/officers_crews_ml_3.csv'
```

```
original_headers = list(df1.columns.values)
```

| df1 |           |              |                |               |               |              |                |                               |
|-----|-----------|--------------|----------------|---------------|---------------|--------------|----------------|-------------------------------|
|     | cohort_id | member_count | years_on_force | percent_black | percent_white | percent_male | percent_female | internal_complaints_per_persu |
| 0   | 824       | 7            | 19.428571      | 14.285714     | 28.571429     | 71.428571    | 28.571429      | 4.2857                        |

|      |      |     |           |           |           |            |           |
|------|------|-----|-----------|-----------|-----------|------------|-----------|
| ...  | ...  | ... | ...       | ...       | ...       | ...        | ...       |
| 2329 | 115  | 9   | 22.000000 | 0.000000  | 77.777778 | 44.444444  | 33.333333 |
| 2330 | 883  | 13  | 9.500000  | 7.692308  | 30.769231 | 84.615385  | 7.692308  |
| 2331 | 558  | 12  | 24.181818 | 83.333333 | 0.000000  | 75.000000  | 16.666667 |
| 2332 | 830  | 7   | 30.000000 | 0.000000  | 85.714286 | 100.000000 | 0.000000  |
| 2333 | 1114 | 7   | 25.571429 | 0.000000  | 14.285714 | 100.000000 | 0.000000  |

```
member_count          int64
years_on_force        float64
percent_black          float64
percent_white          float64
percent_male           float64
percent_female         float64
internal_complaints_per_person  float64
detected_crew          category
dtype: object
```

```

        ('percent_maire', sklearn.preprocessing.MinMaxScaler()),
        (['internal_complaints_per_person', sklearn.preprocessing.MinMaxScaler()])
    ])

In [ ]: # use the mapper to scale the data
features1 = np.round(mapper_features1.fit_transform(df1.copy()), 3)
features1

```

```
In [ ]: # train test split
X_train1, X_test1, y_train1, y_test1 = train_test_split(features1, df1['detected_crew'], test_size=0.3,
shuffle=True)

In [ ]: # train the linear regression modell using train data
linreg1 = linear_model.LinearRegression()
linreg1.fit(X_train1, y_train1)

Out[ ]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

|      |      |     |      |      |      |      |      |      |      |      |      |      |
|------|------|-----|------|------|------|------|------|------|------|------|------|------|
| -0.1 | 0.1  | 0.1 | -0.1 | 0.1  | 0.1  | 0.1  | -0.1 | 0.1  | 1.1  | 0.1  | 0.1  | 0.1  |
| -0.1 | 0.1  | 0.1 | 0.1  | 0.1  | -0.1 | 0.1  | 0.1  | -0.1 | 0.1  | -0.1 | 0.1  | 0.1  |
| 0.1  | 0.1  | 0.1 | 0.1  | 0.1  | 0.1  | -0.1 | -0.1 | 0.1  | 0.1  | -0.1 | -0.1 | 0.1  |
| 0.1  | 0.1  | 0.1 | -0.1 | -0.1 | 0.1  | 0.1  | -0.1 | 0.1  | -0.1 | 0.1  | -0.1 | -0.1 |
| -0.1 | -0.1 | 1.1 | -0.1 | -0.1 | -0.1 | -0.1 | 0.1  | 0.1  | 0.1  | -0.1 | 0.1  | 0.1  |
| 0.1  | 0.1  | 0.1 | 0.1  | -0.1 | 0.1  | -0.1 | 0.1  | 0.1  | 0.1  | 1.1  | -0.1 | 0.1  |
| 1.1  | 0.1  | 0.1 | 0.1  | 0.1  | 0.1  | -0.1 | 0.1  | 0.1  | 0.1  | 0.1  | -0.1 | 0.1  |
| -0.1 | 0.1  | 0.1 | 0.1  | 0.1  | 0.1  | 0.1  | 0.1  | 0.1  | 0.1  | 0.1  | -0.1 | 0.1  |
| 0.1  | 0.1  | 0.1 | 0.1  | 0.1  | -0.1 | -0.1 | -0.1 | -0.1 | -0.1 | 0.1  | 0.1  | -0.1 |

[illegible]

|      |      |      |      |      |      |      |      |      |      |      |      |      |
|------|------|------|------|------|------|------|------|------|------|------|------|------|
| -0.0 | 0.0  | 0.0  | 0.0  | -0.0 | -0.0 | -0.0 | -0.0 | 0.0  | -0.0 | -0.0 | -0.0 | 0.0  |
| 0.0  | -0.0 | -0.0 | 0.0  | 0.0  | -0.0 | 0.0  | 0.0  | 0.0  | 0.0  | 0.0  | 0.0  | -0.0 |
| 0.0  | -0.0 | -0.0 | -0.0 | 0.0  | 0.0  | 0.0  | -0.0 | 0.0  | 0.0  | -0.0 | -0.0 | 0.0  |
| -0.0 | 0.0  | 0.0  | 0.0  | 0.0  | 0.0  | -0.0 | 0.0  | 0.0  | -0.0 | 0.0  | 0.0  | -0.0 |
| 0.0  | 0.0  | -0.0 | 1.0  | -0.0 | -0.0 | 0.0  | -0.0 | -0.0 | 0.0  | 0.0  | -0.0 | 0.0  |
| 0.0  | 0.0  | 0.0  | 0.0  | -0.0 | 0.0  | -0.0 | -0.0 | 0.0  | 0.0  | 0.0  | 0.0  | 0.0  |
| 0.0  | 0.0  | 0.0  | 0.0  | -0.0 | -0.0 | 0.0  | 0.0  | 0.0  | 0.0  | -0.0 | -0.0 | 0.0  |
| -0.0 | 0.0  | 0.0  | 0.0  | 0.0  | -0.0 | -0.0 | 0.0  | 0.0  | 0.0  | -0.0 | -0.0 | -0.0 |

```
0., -0., -0., -0., 0., 1., 0., -0., 0., 0., 0., 0.))

In [ ]: #show accuracy of model 1
        metrics.accuracy_score(y_test1, predictionspl)

Out[ ]: 0.9258202567760342

In [ ]: # coefficient of the linear model
```

```
Out[ ]: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                             intercept_scaling=1, l1_ratio=None, max_iter=100,
                             multi_class='auto', n_jobs=None, penalty='l2',
                             random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
                             warm_start=False)

In [ ]: predictionlog1 = logreg1.predict(X_test1)
         metrics.accuracy_score(y_test1, predictionlog1)
```

Similar to the earlier model, we want to know which feature of a group impacts the prediction of crew and whether such a model is accurate. This time, based on the `data_crew` table, we selected features that average information about the officers in their ranks such as member count, race, and complaints per person.

According to our current model, the internal complaints per person and member count are most heavily weighted when predicting whether a group is classified as a crew or not.

It is possible to predict and model officer misconduct.

At an individual level, officers with high complaint percentiles are most heavily weighted when predicting their cohort.

If these metrics are tracked, it may be possible to detect whether individual officers are at higher risk of continued misconduct.