# Predicting Students' First-year Success at UC Santa Cruz

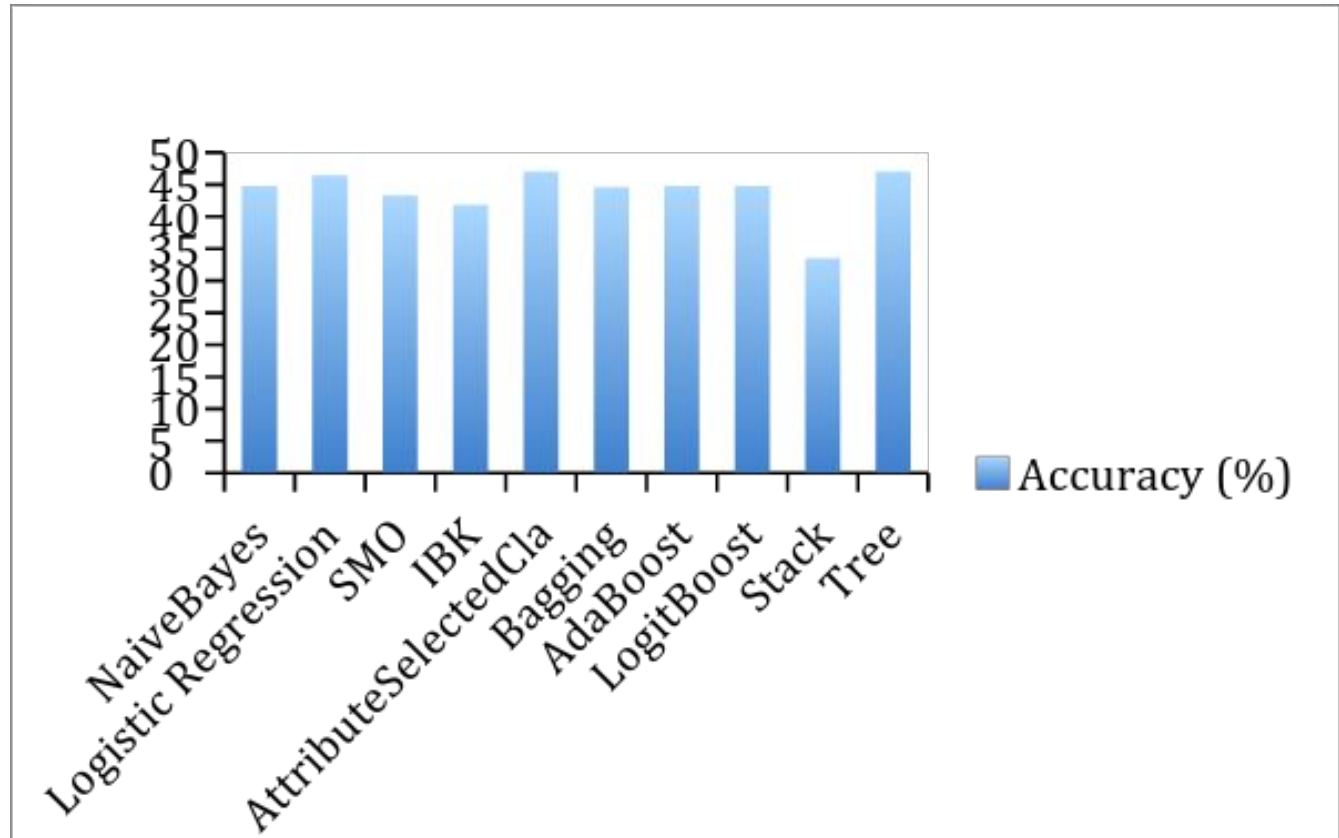**Hairong Wu     Kevin Doyle     Connor McNeill**

## Table1 Best performance of different methods on 10-fold cross-validation

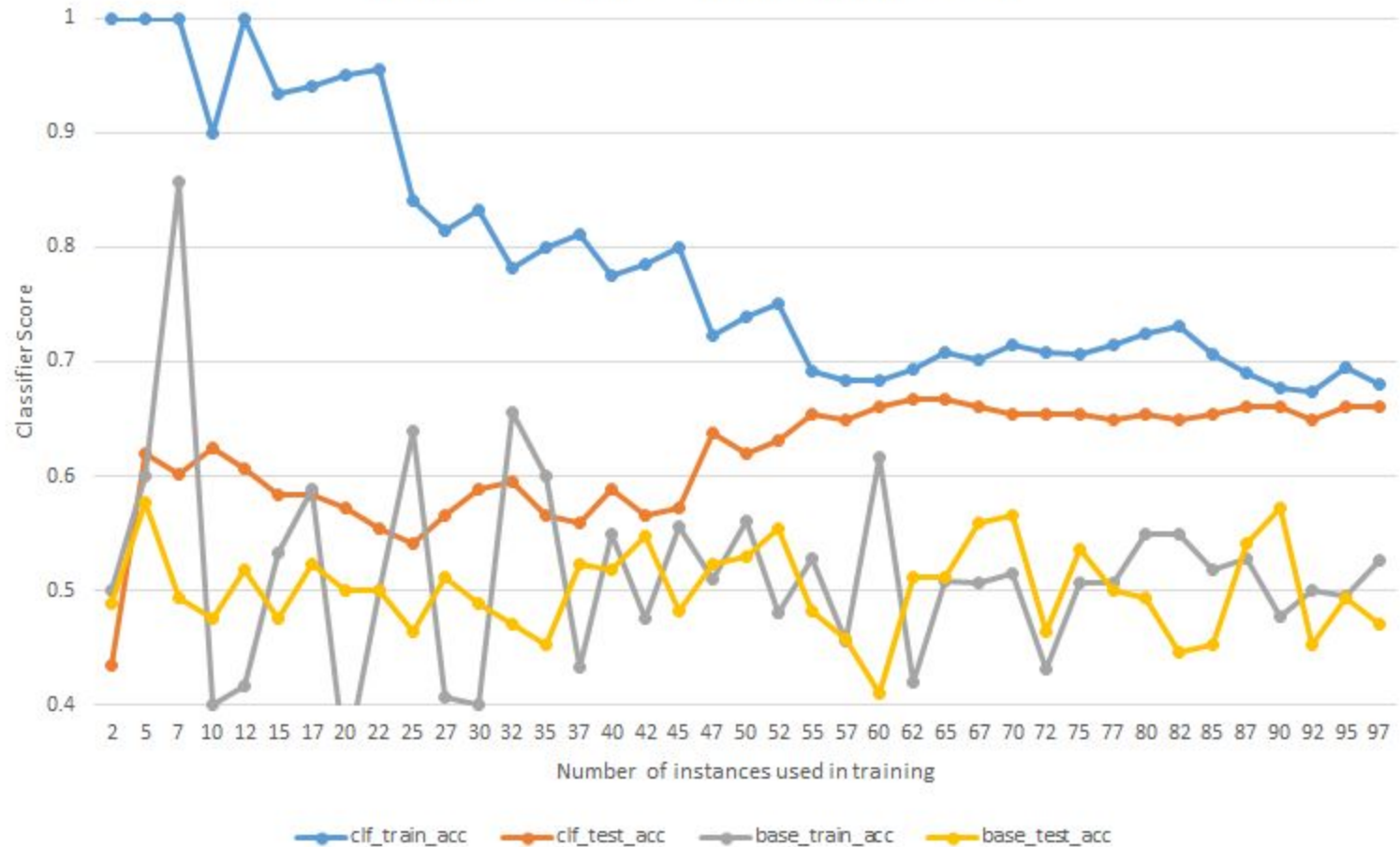| Methods on 10-fold cross-validation | Accuracy (%) |
| --- | --- |
| NaiveBayes | 44.78 |
| Logistic Regression | 46.38 |
| SMO (RBFkernel) | 43.32 |
| IBK (kNN=9, LinearNNSearch) | 41.86 |
| AttributeSelectedClassifier (LMT) | 47.03 |
| Bagging-RandomForest | 44.61 |
| AdaBoostM-HoeffdingTree | 44.78 |
| LogitBoost-DecisionStump | 44.81 |
| Stack-DecisionTable | 33.51 |
| Tree-LMT | 47.03 |

# linear regression

$$Y= 0.2045 * gender$$
$$- 0.0003 * Firgen$$
$$- 0.0004 * SATCRDG$$
$$+ 0.0004 * SATMATH$$
$$+ 0.0011 * SATWRTG$$
$$+ 0.0001 * SATTotal$$
$$+ 0.4487 * HSGPA$$
$$+ 0.0175 * ACTEngWrit$$
$$+ 0.0198 * APIScore$$
$$+ 0.0591 * FirstLang$$
$$- 0.302$$
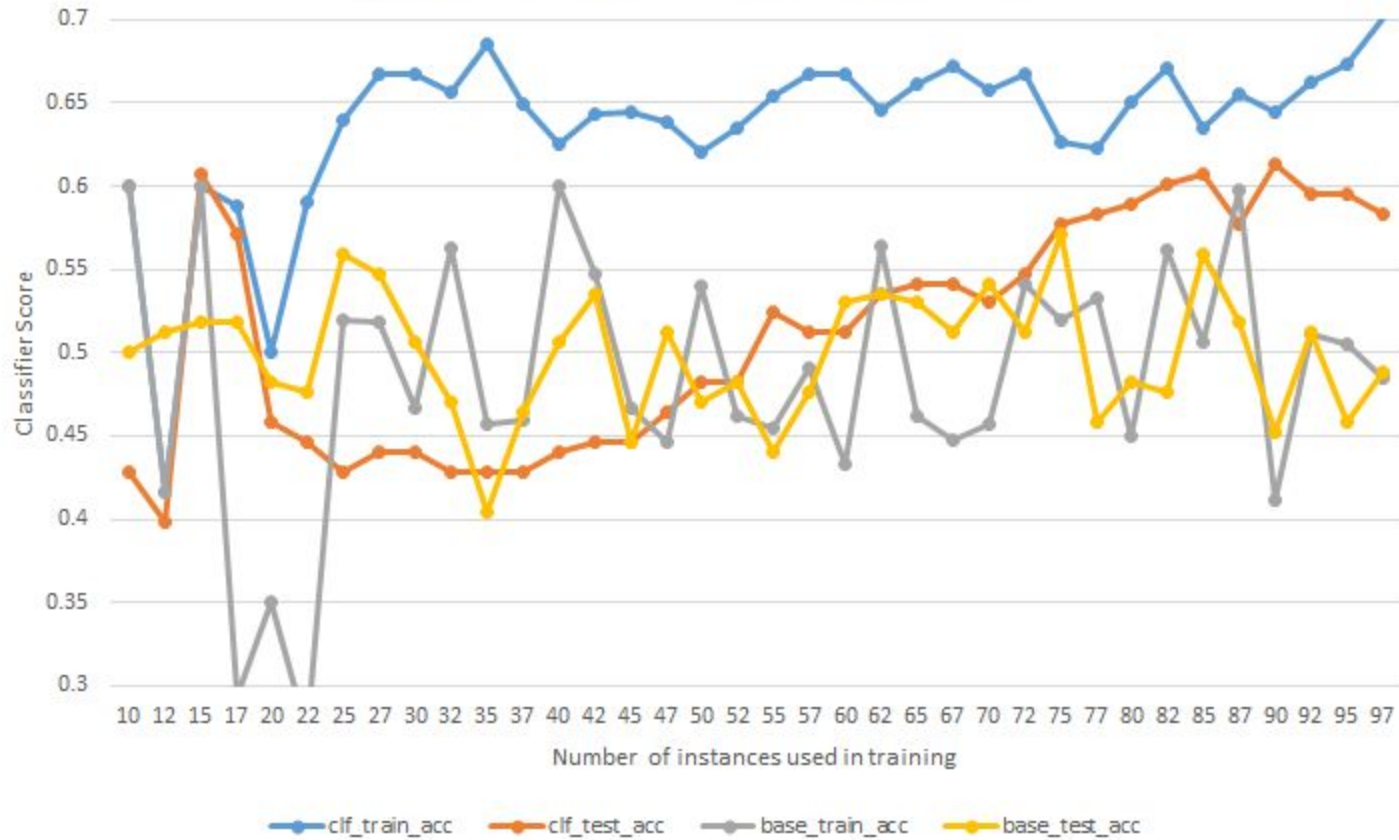
| Properties | Correlation coefficient |
|---|---|
| HSGPA | 0.2074 |
| gender | 0.088 |
| FirstLang | 0.1355 |
| APIScore | 0.1687 |
| ACTEngWrit | 0.1612 |
| SATWRTG | 0.2447 |
| SATCRDG | 0.1915 |
| SATMATH | 0.1764 |
| Firgen | 0.2447 |
| SATTotal | 0.1915 |
| famincome | 0.1764 |
| ACTRead | 0.1590 |

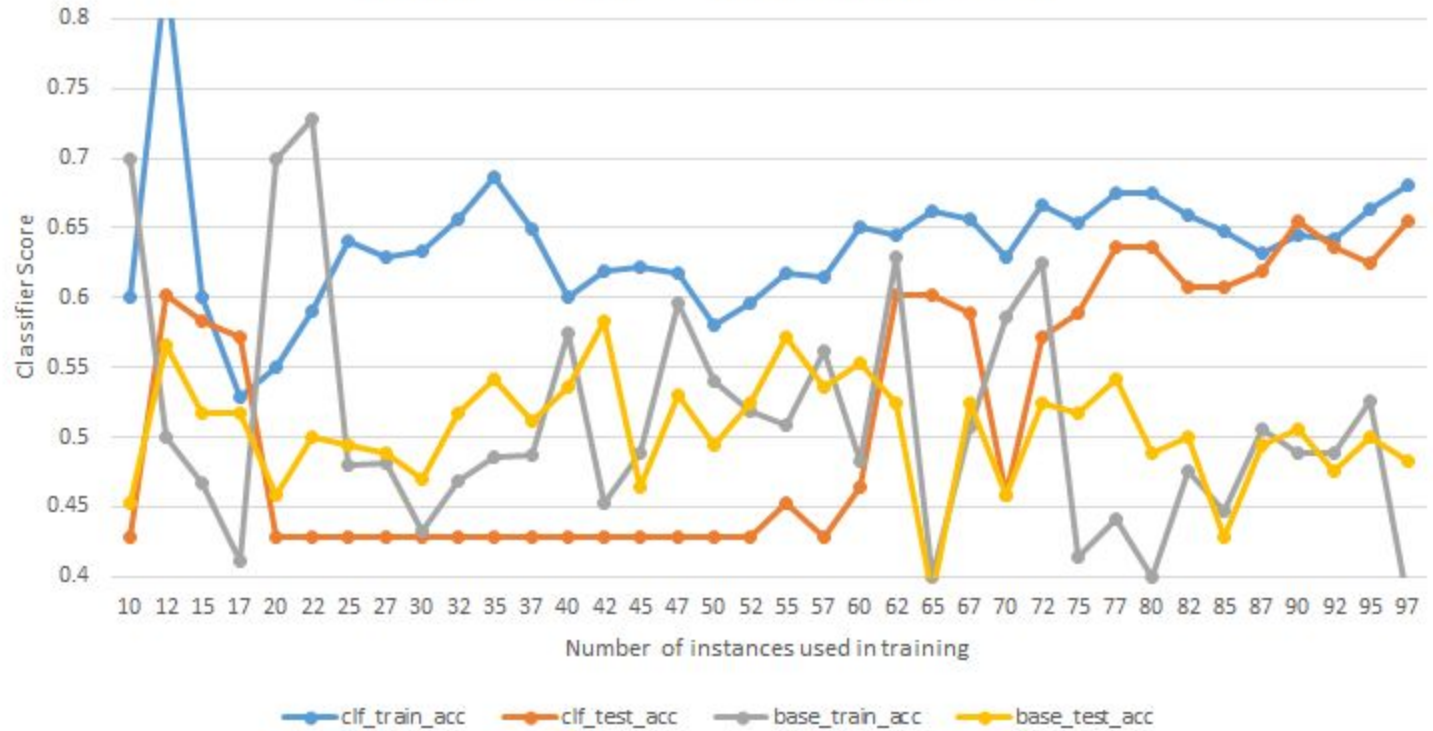| | |
|---|---|
| HSGPA | 0.2074 |
| HSGPA, gender | 0.2199 |
| HSGPA, gender, FirstLang | 0.2482 |
| HSGPA, gender, FirstLang, APIScore | 0.2914 |
| HSGPA, gender, FirstLang, APIScore, ACTEngWrit | 0.3009 |
| HSGPA, gender, FirstLang, APIScore, ACTEngWrit, SATWRTG | 0.3311 |
| HSGPA, gender, FirstLang, APIScore, ACTEngWrit, SATWRTG, SATMATH | 0.3346 |
| HSGPA, gender, FirstLang, APIScore, ACTEngWrit, SATWRTG, SATMATH, SATCRDG | 0.3346 |
| HSGPA, gender, FirstLang, APIScore, ACTEngWrit, SATWRTG, SATMATH, SATCRDG, Firgen | 0.3356 |
| HSGPA, gender, FirstLang, APIScore, ACTEngWrit, SATWRTG, SATMATH, SATCRDG, Firgen, SATTotal | 0.3352 |

NBGaussian: All Features
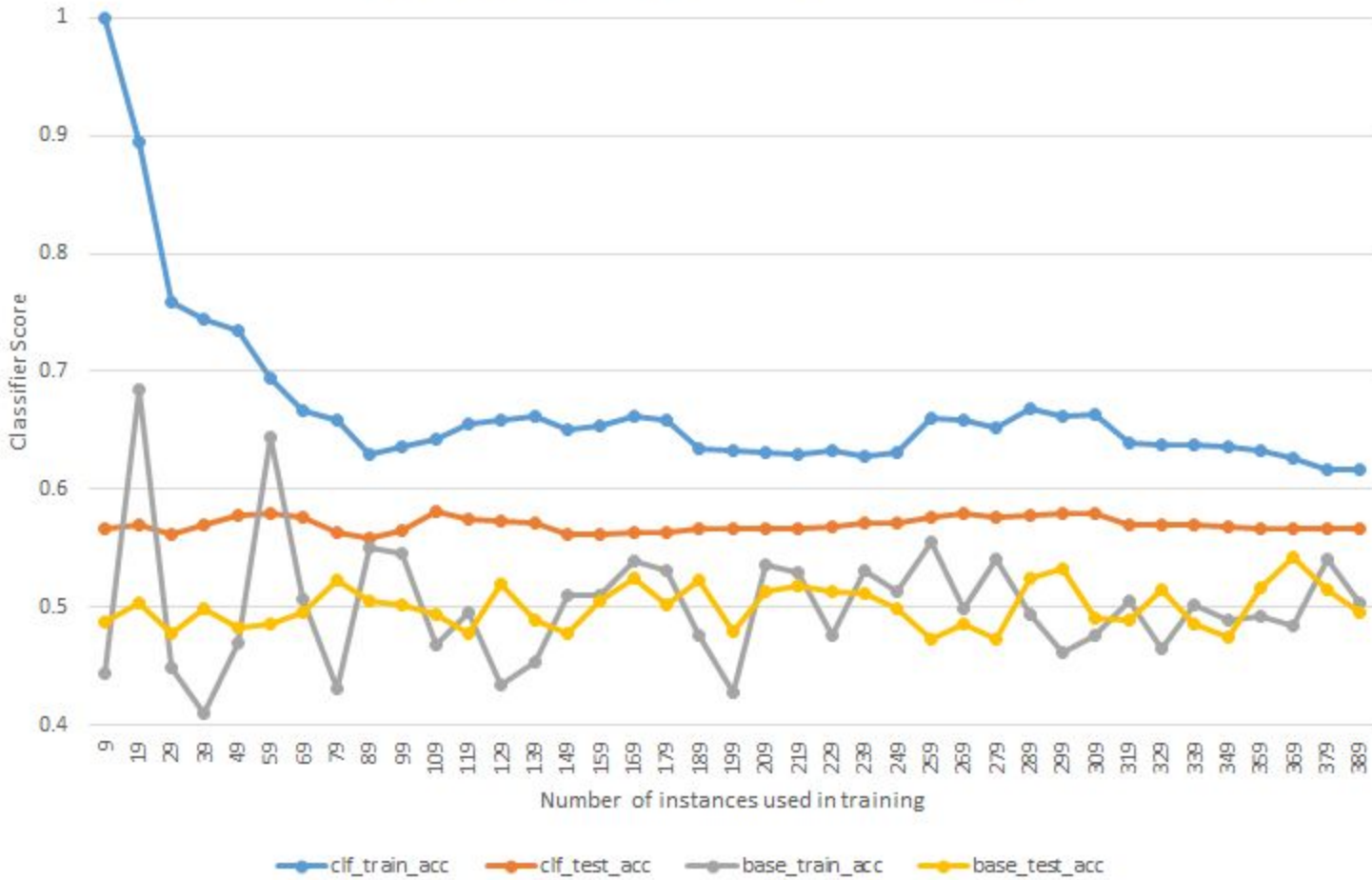Train/Test Accuracy vs Training Instance Count

KNN: All Features
Train/Test Accuracy vs Training Instance Count

clf_train_acc — clf_test_acc — base_train_acc — base_test_acc

SVM: All Features
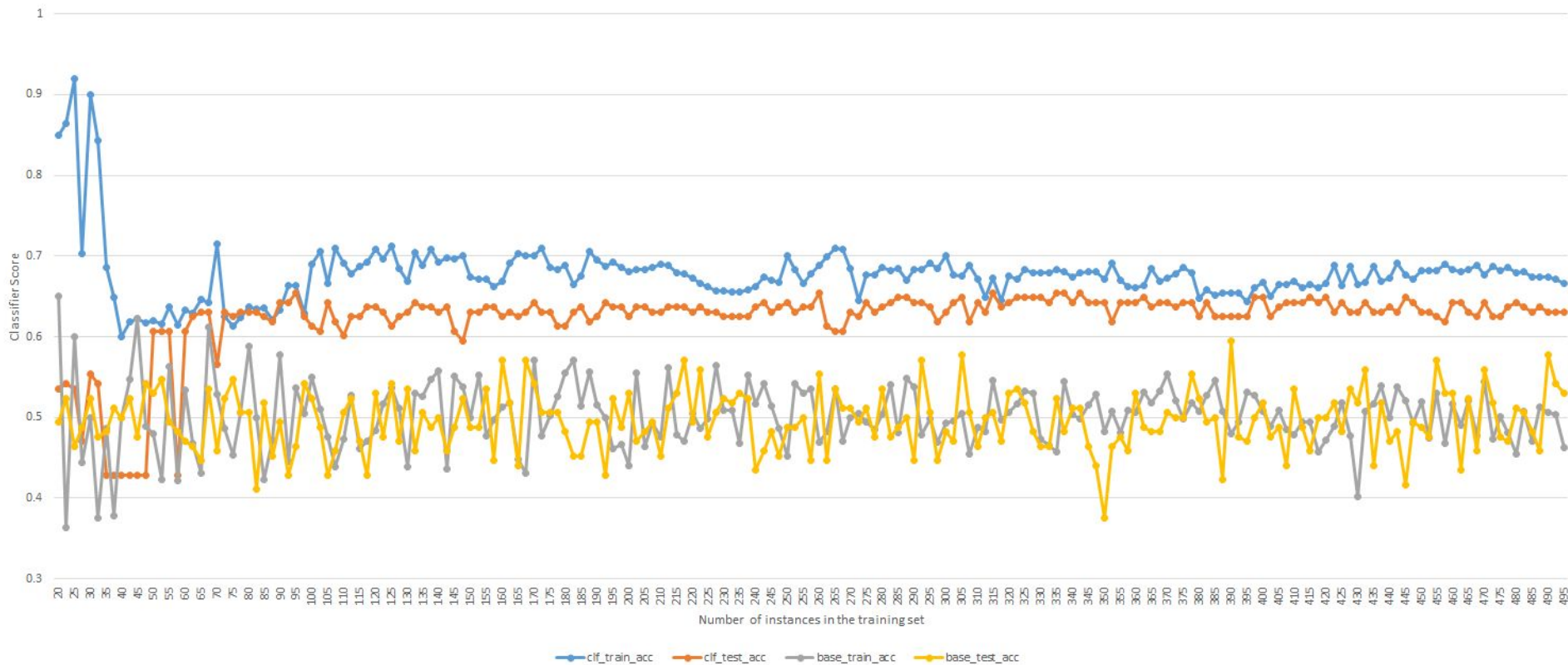Train/Test Accuracy vs Training Instance Count

SVM: Non Academic
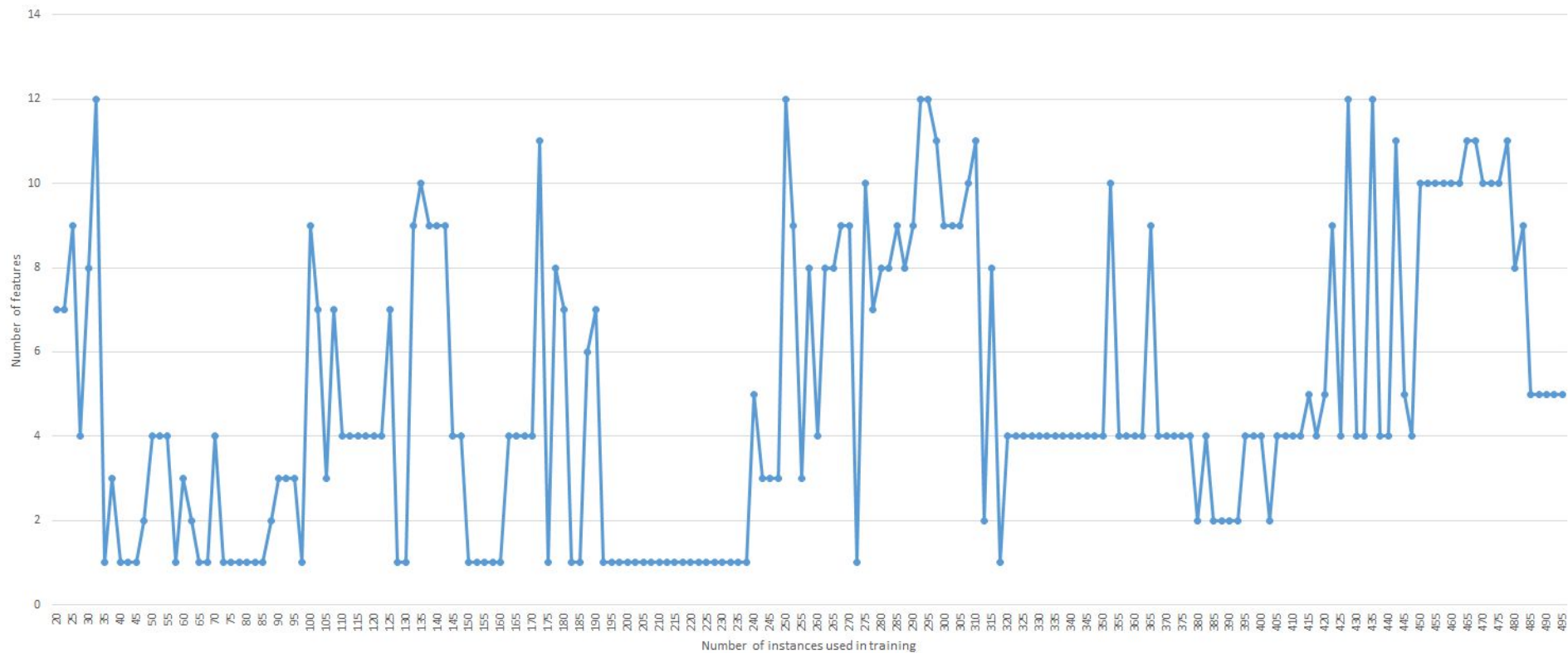Train/Test Accuracy vs Training Instance Count

# Break the bias..? Nope.



Grid Search with PCA + SVM (linear) Pipeline, Across Training Set Size Change

num_feats

c_vals