

HW3

Mark Vandergon

Summary of modeling results

I trained sixteen different models on the credit dataset, eight estimators chosen from grid search as the best model for accuracy and the same eight models optimized for precision.

The best performing model for accuracy were the regression and forest classifiers, although accuracy varied very little between models. Naïve Bayes was by far the worst performing model. Precision appears to be the most significant differentiator. Gradient Boosting Classifier was the top score, followed by Logistic Regression. Logistic Regression outperformed many other models in terms of accuracy, precision, as well as the two combined.

Because I implemented parallel processing of the grid search, the time difference per model (the time column) didn't impede running these models in reasonable time. However, Gradient Boosting, RandomForest and ExtraTrees Classifiers took significantly longer than Naïve Bayes models.

| i | clf | target | accuracy | precision | recall | time | params |
|----|----------------------------|-----------|----------|-----------|--------|------|--|
| 0 | RandomForestClassifier | accuracy | 0.94 | 0.60 | 0.54 | 0.14 | 'error_score': 'raise', 'cv': None, 'max_depth': 10, 'max_features': 'sqrt', 'min_samples_split': 2, 'n_estimators': 100, 'oob_score': False, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 1 | ExtraTreesClassifier | accuracy | 0.93 | 0.50 | 0.52 | 0.14 | 'error_score': 'raise', 'cv': None, 'max_depth': 10, 'max_features': 'sqrt', 'min_samples_split': 2, 'n_estimators': 100, 'oob_score': False, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 2 | AdaBoostClassifier | accuracy | 0.93 | 0.38 | 0.56 | 0.05 | 'error_score': 'raise', 'cv': None, 'learning_rate': 0.1, 'n_estimators': 100, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 3 | LogisticRegression | accuracy | 0.94 | 0.75 | 0.54 | 0.05 | 'error_score': 'raise', 'cv': None, 'max_iter': 1000, 'multi_class': 'multinomial', 'n_jobs': 1, 'penalty': 'l2', 'random_state': 0, 'solver': 'lbfgs', 'verbose': 0, 'warm_start': False |
| 4 | GradientBoostingClassifier | accuracy | 0.93 | 0.43 | 0.53 | 0.20 | 'error_score': 'raise', 'cv': None, 'learning_rate': 0.1, 'max_depth': 10, 'max_features': 'sqrt', 'min_samples_split': 2, 'n_estimators': 100, 'oob_score': False, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 5 | GaussianNB | accuracy | 0.27 | 0.08 | 0.59 | 0.01 | 'error_score': 'raise', 'cv': None, 'gaussian_prior': False, 'multi_class': 'multinomial', 'priors': [0.1, 0.1, 0.1, 0.1], 'var_smoothing': 1e-05 |
| 6 | DecisionTreeClassifier | accuracy | 0.93 | 0.00 | 0.50 | 0.02 | 'error_score': 'raise', 'cv': None, 'max_depth': 10, 'max_features': 'sqrt', 'min_samples_split': 2, 'n_estimators': 100, 'oob_score': False, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 7 | SGDClassifier | accuracy | 0.93 | 0.00 | 0.50 | 0.01 | 'error_score': 'raise', 'cv': None, 'epsilon': 0.0001, 'learning_rate': 0.01, 'loss': 'log_loss', 'n_jobs': 1, 'n_iter': 1000, 'penalty': 'l2', 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 8 | KNeighborsClassifier | accuracy | 0.93 | 0.00 | 0.50 | 0.01 | 'error_score': 'raise', 'cv': None, 'k': 1, 'leaf_size': 30, 'metric': 'minkowski', 'metric_params': None, 'n_jobs': 1, 'p': 2, 'radius': 1.0, 'weights': 'uniform', 'warm_start': False |
| 9 | RandomForestClassifier | precision | 0.93 | 0.43 | 0.53 | 0.16 | 'error_score': 'raise', 'cv': None, 'max_depth': 10, 'max_features': 'sqrt', 'min_samples_split': 2, 'n_estimators': 100, 'oob_score': False, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 10 | ExtraTreesClassifier | precision | 0.93 | 0.43 | 0.53 | 0.17 | 'error_score': 'raise', 'cv': None, 'max_depth': 10, 'max_features': 'sqrt', 'min_samples_split': 2, 'n_estimators': 100, 'oob_score': False, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 11 | AdaBoostClassifier | precision | 0.93 | 0.38 | 0.56 | 0.05 | 'error_score': 'raise', 'cv': None, 'learning_rate': 0.1, 'n_estimators': 100, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 12 | LogisticRegression | precision | 0.94 | 0.75 | 0.54 | 0.02 | 'error_score': 'raise', 'cv': None, 'max_iter': 1000, 'multi_class': 'multinomial', 'n_jobs': 1, 'penalty': 'l2', 'random_state': 0, 'solver': 'lbfgs', 'verbose': 0, 'warm_start': False |
| 13 | GradientBoostingClassifier | precision | 0.94 | 1.00 | 0.51 | 0.13 | 'error_score': 'raise', 'cv': None, 'learning_rate': 0.1, 'max_depth': 10, 'max_features': 'sqrt', 'min_samples_split': 2, 'n_estimators': 100, 'oob_score': False, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 14 | GaussianNB | precision | 0.27 | 0.08 | 0.59 | 0.01 | 'error_score': 'raise', 'cv': None, 'gaussian_prior': False, 'multi_class': 'multinomial', 'priors': [0.1, 0.1, 0.1, 0.1], 'var_smoothing': 1e-05 |
| 15 | DecisionTreeClassifier | precision | 0.93 | 0.25 | 0.52 | 0.02 | 'error_score': 'raise', 'cv': None, 'max_depth': 10, 'max_features': 'sqrt', 'min_samples_split': 2, 'n_estimators': 100, 'oob_score': False, 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 16 | SGDClassifier | precision | 0.93 | 0.00 | 0.50 | 0.03 | 'error_score': 'raise', 'cv': None, 'epsilon': 0.0001, 'learning_rate': 0.01, 'loss': 'log_loss', 'n_jobs': 1, 'n_iter': 1000, 'penalty': 'l2', 'random_state': 0, 'verbose': 0, 'warm_start': False |
| 17 | KNeighborsClassifier | precision | 0.86 | 0.04 | 0.48 | 0.02 | 'error_score': 'raise', 'cv': None, 'k': 1, 'leaf_size': 30, 'metric': 'minkowski', 'metric_params': None, 'n_jobs': 1, 'p': 2, 'radius': 1.0, 'weights': 'uniform', 'warm_start': False |