

An Empirical Comparison of Supervised Learning Algorithms for Classification

Dillon Ford, PID: A16092047

Abstract—This paper examines multiple supervised machine learning algorithms used to solve classification problems. Here we perform an empirical analysis and comparison between the supervised learning methods: logistic regression, k-nearest neighbors, random forests, and decision trees. Performance is measured across multiple trials and datasets for each classifier.

I. INTRODUCTION

For a classification problem, finding the appropriate classifier with optimal parameters is crucial for success. Many of the ideas for this paper stem from prior work conducted to introduce methodologies and practices in classification [1] that optimized and empirically analyzed the performance of classifiers over many datasets. This paper will examine four different classifiers' performance in three trials across three datasets to solve supervised machine learning algorithms' problems. The Mushroom, Heart Disease, and Drug Consumption datasets are obtainable from the UCI Machine Learning Repository [2].

A note on datasets: The instructor provided approval for using datasets in this paper.

The Mushroom dataset (MU) contains records drawn from The Audubon Society Field Guide to North American Mushrooms [3]. The goal is to distinguish edible mushrooms from poisonous mushrooms based on 22 features describing the mushroom characteristics such as odor, ring type, and gill size and color. MU classification problem was attractive to me because I lived in the Pacific Northwest for some time and became fascinated with hunting and identifying mushrooms in the Redwood Forest.

The Heart Disease dataset (HD) contains records from medical centers in Budapest, Switzerland, and California [4]. The classification operation is to separate patients with heart disease from patients without heart disease utilizing 14 different attributes including, chest pain type, resting blood pressure, and age, to name a few. Being able to classify whether or not a patient has heart disease is crucial as heart disease is the number one leading cause of death in the United States [5] which is why I elected to run this classification problem.

The Drug Consumption dataset (DRG) contains records from institutes in the UK, including Rampton Hospital, University of Nottingham, and University of Leicester [6]. Respondents were questioned about their use of 18 legal and illegal drugs and asked about how frequently they used the drug. For this project, a transformation of the dataset into a binary

classification problem took place to examine whether or not the participant is a "user" or "non-user" of a particular drug based on 32 attributes such as the big five personality traits. The drug examined in this project is marijuana since it is the most balanced in terms of use and non-use in this dataset. Favorite topics of mine to study are pharmacology and the psychology of drug use. I was interested in seeing whether or not personality and behavior could determine whether or not a person uses a particular drug.

II. METHODS

Four learning algorithms are tested across the datasets consisting of three trials in each. Learning algorithms include Logistic Regression, K-Nearest Neighbors, Random Forests, and Decision Trees. Modeling of the algorithm parameters come from the Caruna and Niculescu-Mizil paper [1] as described below. Implementation of Scikit-learn's GridSearchCV [7] is used to find the optimal parameters associated with each algorithm across each trial. The performance of each learning algorithm is evaluated using the performance metrics: accuracy and precision. Scores between all metrics are averaged across the three trials for each dataset and averaged between the three datasets.

The parameter spaces searched for each of the classifiers in each of the datasets are shown below.

Logistic Regression (LOGREG): in this learning algorithm, regularized and unregularized models are implemented. The regularization parameter C varies by factors of ten from 10^{-8} to 10^4 .

K-Nearest Neighbors (KNN): in this learning algorithm, KNN with Euclidean distance is implemented and weighted uniformly and by distance. The number of neighbors is 20, varying by 15, ranging from 1 to 300 for MU and DRG. Neighbor values adjusted to 20 varying by 5, ranging from 1 to 60 for HD due to it being a significantly smaller dataset.

Random Forests (RF): in this learning algorithm, the forests have 1024 trees. The size of the feature set considered at each split is 1,2,4,6,8,12,16, or 20.

Decision Trees (DT): in this learning algorithm, the decision tree classifier CART with the gini criterion is used with a max depth of each ranging from 1 to 10. The best splitting method is used with 2 samples to split.

III. EXPERIMENT

For each dataset and learning algorithm, 3 trials are performed. For each trial in the datasets, randomly select 80 percent of the training cases and the remaining 20 percent of cases for testing. A Pipeline GridSearchCV for all learning algorithms is performed with 5-fold cross-validation on the training cases to return the best parameters for the mean over 5-folds. Tuning of GridSearchCV for scoring metrics occurred to refit for accuracy and precision. In each case, precision and accuracy, the model is trained a final time on all the training cases, selecting optimal parameters, and measuring performance on the test cases.¹

For each trial and learning algorithm, visualization of training performance is presented through a learning curve, plotting training score against cross-validation score across training examples, a plot for the scalability of the model, showing fit times across training examples, and a plot of model performance, showing score across fit times. Please see the Appendix for an example of such graphics.

IV. RESULTS

Table 1 . Algorithm Mean Test Set Performance By Dataset

		KNN	RF	LR	DT
HD	Accuracy	0.809	0.831	0.847	0.781
	Precision	0.789	0.855	0.839	0.794
DRG	Accuracy	0.801	0.791	0.811	0.774
	Precision	0.857	0.845	0.850	0.851
MU	Accuracy	1.000	1.000	0.966	1.000
	Precision	1.000	1.000	0.964	1.000

In the HD dataset, LOGREG performed best in accuracy, and RF showed the highest precision scores. Similarly, in DRG, LOGREG also showed the highest accuracy results, but KNN scored the highest in precision. In contrast, MU, LOGREG did not perform as well as other learning algorithms. Table 1 summarizes the results of each dataset and algorithm accuracy and precision performance. A paired t-test was used between each algorithm across the three datasets to find the p-value, and a calculation of Cohen's d-value to measure effect size; these values are in Appendix 1 along with dataset raw algorithm test scores for the trial. With $p = 0.05$, significant results are shown in DRG, where a comparison between KNN and DT for accuracy resulted in $p = 0.035$ and $d = 2.610$ as well as LOGREG and DT, with $p = 0.007$ and $d = 4.583$. In MUSH significant results are shown between LR and all other classifiers, with $p = 0.01$ and $d = 8.057$ for accuracy, and $p = 0.026$ and $d = 4.978$ for precision.

¹A .py file containing functions used to perform these calculations and a notebook outlining execution and analysis are in the Appendix; for complete documentation of code and notebooks, please refer to my GitHub, https://github.com/dillon4d/COGS118A_FinalProject

Table 2 . Algorithm Mean Test Set Performance

	KNN	RF	LR	DT
Accuracy	0.870	0.874	0.875	0.852
Precision	0.882	0.900	0.885	0.881

An analysis of the overall performance of each classification algorithm ensued. In terms of accuracy, LOGREG and RF performed almost equally, followed by KNN and DT. For precision, RF outperformed all other learning algorithms. Table 2 summarizes the mean test set performance for each algorithm. As in Appendix 1, respective p-values and d-values are in Appendix 2. With $p = 0.05$, no significant results were shown between classifiers when averaging test performances across the three datasets, and there was no significance in effect size.

Table 3 . Algorithm Mean Optimal Train Set Performance

		KNN	RF	LR	DT
HD	Accuracy	0.814	0.821	0.817	0.749
	Precision	0.856	1.000	0.815	0.864
DRG	Accuracy	0.802	0.807	0.813	0.723
	Precision	0.903	1.000	0.857	0.904
MU	Accuracy	0.999	1.000	0.965	0.998
	Precision	1.000	1.000	0.964	1.000

In comparison to algorithm mean test set performance, as in Table 1, scores in Table 3 reflect the mean train set performance of algorithms adjusted with optimal hyperparameters. In HD, algorithm accuracy performance was similar to test set performance; KNN accuracy slightly increased while the other algorithms slightly decreased. Precision scores improved significantly for HD classifiers when adjusted, notably RF, while LR saw a slight decrease. In DRG, accuracy performance increased somewhat in all classifiers with adjusted parameters except DT, and precision scores improved in all classifiers. In MU, performance with adjusted parameters in the train set remained nearly the same for all classifiers as in their test set performance. Raw scores for dataset algorithm trainset with optimal parameters are in Appendix 3.

V. CONCLUSION

Overall, RF showed the most precise results among the four classifiers tested across datasets and had near equal accuracy to that of LR, which showed the highest mean test accuracy performance. As shown in Appendix Table 2, the t-test on average performances between algorithms shows no significant results nor significance in effect size. In retrospect, there is room for improvement; the testing of more classifiers, trials, and metrics would provide insight into the results achieved here. Further assessment of the datasets utilizing a bootstrapping analysis would also prove to be beneficial. Given more time, a more in-depth analysis of more drugs, not just marijuana, in DRG could provide greater insight into the personality and the attractiveness of drugs.

VI. BONUS POINTS

This project deserves bonus points for performing tests on an extra learning algorithm and evaluating the algorithms on an accuracy metric and a precision metric. Furthermore, during each trial, for each algorithm, learning curves are plotted to show performance and other visualizations for each trial algorithm, such as a confusion matrix and ROC curve. Each notebook contains an extensive exploratory data analysis, including comprehensive data visualizations. In HD and MU, perform a preemptive random forest classifier to retrieve the dataset's top features. The top 5 explored extensively before model training.

REFERENCES

- [1] R. Caruana and A. Niculescu-Mizil, "An empirical comparison of supervised learning algorithms," in *Proceedings of the 23rd International Conference on Machine Learning*, ser. ICML '06. New York, NY, USA: Association for Computing Machinery, 2006, p. 161–168. [Online]. Available: <https://doi.org/10.1145/1143844.1143865>
- [2] D. Dua and C. Graff, "UCI machine learning repository," 2017. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [3] G. H. Lincoff, *The Audubon Society Field Guide to North American Mushrooms*. New York: Alfred A. Knopf, 1981.
- [4] R. Detrano, A. Janosi, W. Steinbrunn, M. Pfisterer, J.-J. Schmid, S. Sandhu, K. H. Guppy, S. Lee, and V. Froelicher, "International application of a new probability algorithm for the diagnosis of coronary artery disease," *The American Journal of Cardiology*, vol. 64, no. 5, pp. 304 – 310, 1989. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/0002914989905249>
- [5] N. S. Shah, R. Molsberry, J. S. Rana, S. Sidney, S. Capewell, M. O'Flaherty, M. Carnethon, D. M. Lloyd-Jones, and S. S. Khan, "Heterogeneous trends in burden of heart disease mortality by subtypes in the united states, 1999-2018: observational analysis of vital statistics," *BMJ*, vol. 370, 2020. [Online]. Available: <https://www.bmjjournals.org/content/370/bmj.m2688>
- [6] E. Fehrman, A. K. Muhammad, E. M. Mirkes, V. Egan, and A. N. Gorban, "The five factor model of personality and evaluation of drug consumption risk," in *Data Science*, F. Palumbo, A. Montanari, and M. Vichi, Eds. Cham: Springer International Publishing, 2017, pp. 231–242.
- [7] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

Appendix

Below the appendix tables is a .py file that outlines much of the learning algorithms' calculations (varies slightly by metric and dataset). Also attached is a trial from one of the notebooks. This process is repeated across all three trials in every dataset. The full documentation of notebooks for the datasets and their .py files are found on my GitHub.

https://github.com/dillon4d/COGS118A_FinalProject

Appendix 1 . p-values and d-values for algorithm comparisons by dataset

Dataset	Metric	KNN & RF	KNN & LR	KNN & DT	RF & DT	RF & LR	LR & DT
HD	Accuracy	p = 0.746	p = 0.463	p = 0.607	p = 0.502	p = 0.814	p = 0.395
		d = 0.289	d = 0.664	d = 0.456	d = 0.608	d = 0.099	d = 0.780
	Precision	p = 0.363	p = 0.297	p = 0.922	p = 0.438	p = 0.827	p = 0.459
		d = 0.907	d = 1.019	d = 0.088	d = 0.708	d = 0.193	d = 0.671
DRG	Accuracy	p = 0.380	p = .291	p = 0.035	p = 0.127	p = 0.097	p = 0.007
		d = 0.807	d = 1.030	d = 2.61	d = 1.601	d = 1.937	d = 4.583
	Precision	p = 0.609	p = 0.679	p = 0.757	p = 0.821	p = .809	p = 0.982
		d = 0.456	d = 0.368	d = 0.272	d = 0.198	d = 0.216	d = 0.020
MU	Accuracy	p = nan	p = 0.010	p = nan	p = nan	p = 0.010	p = 0.010
		d = nan	d = 8.057	d = nan	d = nan	d = 8.057	d = 8.057
	Precision	p = nan	p = 0.026	p = nan	p = nan	p = 0.026	p = 0.026
		d = nan	d = 4.978	d = nan	d = nan	d = 4.978	d = 4.978

Appendix 1.1. Raw Test Set Scores By Dataset

		KNN	RF	LR	DT
HD	Accuracy	Trial 1: 0.787	Trial 1: 0.803	Trial 1: 0.803	Trial 1: 0.771
		Trial 2: 0.771	Trial 2: 0.754	Trial 2: 0.820	Trial 2: 0.721
	Precision	Trial 3: 0.869	Trial 3: 0.934	Trial 3: 0.918	Trial 3: 0.853
		Trial 1: 0.750	Trial 1: 0.815	Trial 1: 0.793	Trial 1: 0.759
DRG	Accuracy	Trial 2: 0.800	Trial 2: 0.784	Trial 2: 0.816	Trial 2: 0.744
		Trial 3: 0.816	Trial 3: 0.967	Trial 3: 0.909	Trial 3: 0.879
	Precision	Trial 1: 0.796	Trial 1: 0.801	Trial 1: 0.804	Trial 1: 0.775
		Trial 2: 0.814	Trial 2: 0.796	Trial 2: 0.817	Trial 2: 0.783
MU	Accuracy	Trial 3: 0.793	Trial 3: 0.777	Trial 3: 0.812	Trial 3: 0.764
		Trial 1: 0.875	Trial 1: 0.878	Trial 1: 0.859	Trial 1: 0.868
	Precision	Trial 2: 0.834	Trial 2: 0.827	Trial 2: 0.835	Trial 2: 0.823
		Trial 3: 0.861	Trial 3: 0.832	Trial 3: 0.857	Trial 3: 0.861

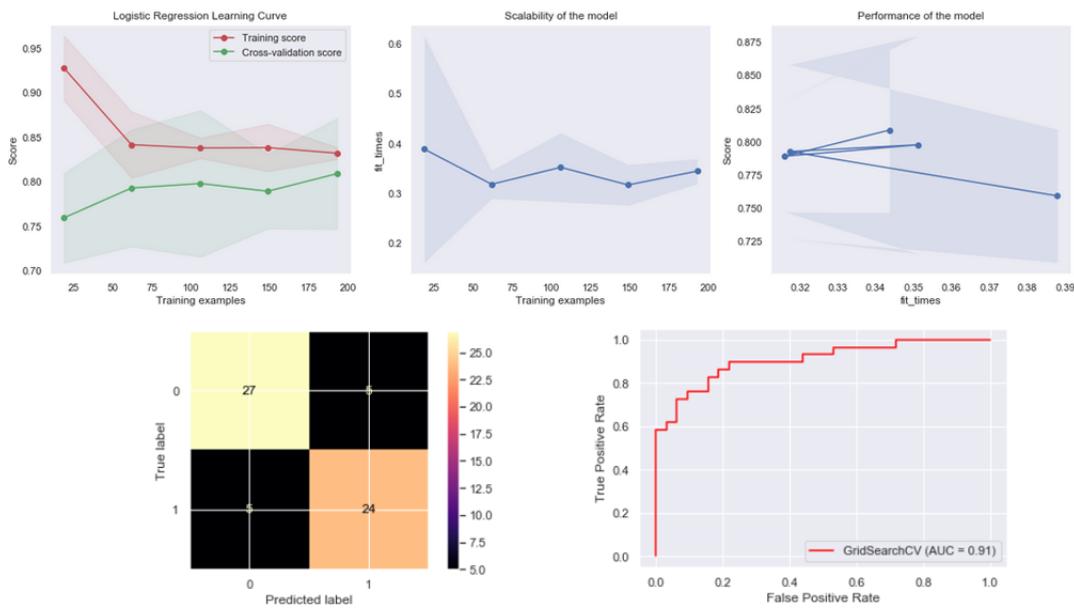
Appendix Table 2. p-values and d-values for algorithm comparisons

	KNN & RF	KNN & LR	KNN & DT	RF & DT	RF & LR	LR & DT
Accuracy	p = 0.935	p = 0.914	p = 0.727	p = 0.678	p = 0.989	p = 0.630
	d = 0.038	d = 0.052	d = 0.167	d = 0.200	d = 0.006	d = 0.232
Precision	p = 0.680	p = 0.942	p = 0.994	p = 0.682	p = 0.687	p = 0.937
	d = 0.198	d = 0.035	d = 0.004	d = 0.197	d = 0.194	d = 0.038

Appendix 3. Raw Optimal Train Set Scores By Dataset

		KNN	RF	LR	DT
HD	Accuracy	Trial 1: 0.827	Trial 1: 0.818	Trial 1: .831	Trial 1: 0.736
		Trial 2: 0.826	Trial 2: 0.831	Trial 2: 0.806	Trial 2: 0.814
	Precision	Trial 3: 0.789	Trial 3: 0.814	Trial 3: 0.814	Trial 3: 0.699
		Trial 1: 0.825	Trial 1: 1.000	Trial 1: 0.813	Trial 1: 0.744
		Trial 2: 0.872	Trial 2: 1.000	Trial 2: 0.840	Trial 2: 0.970
		Trial 3: 0.871	Trial 3: 1.000	Trial 3: 0.791	Trial 3: 0.879
DRG	Accuracy	Trial 1: 0.801	Trial 1: 0.810	Trial 1: 0.816	Trial 1: 0.787
		Trial 2: 0.801	Trial 2: 0.803	Trial 2: 0.813	Trial 2: 0.694
		Trial 3: 0.804	Trial 3: 0.808	Trial 3: 0.812	Trial 3: 0.688
	Precision	Trial 1: 0.858	Trial 1: 1.000	Trial 1: 0.855	Trial 1: 0.903
		Trial 2: 0.852	Trial 2: 1.000	Trial 2: 0.860	Trial 2: 0.909
		Trial 3: 1.000	Trial 3: 1.000	Trial 3: 0.855	Trial 3: 0.898
MU	Accuracy	Trial 1: 0.999	Trial 1: 1.000	Trial 1: 0.962	Trial 1: 0.998
		Trial 2: 0.999	Trial 2: 1.000	Trial 2: 0.968	Trial 2: 1.000
		Trial 3: 0.999	Trial 3: 1.000	Trial 3: 0.964	Trial 3: 0.996
	Precision	Trial 1: 1.000	Trial 1: 1.000	Trial 1: 0.961	Trial 1: 1.000
		Trial 2: 1.000	Trial 2: 1.000	Trial 2: 0.968	Trial 2: 1.000
		Trial 3: 1.000	Trial 3: 1.000	Trial 3: 0.963	Trial 3: 1.000

Example of visualization for an accuracy trial of LR in HD



```
1 import numpy as np
2 from sklearn.metrics import classification_report, plot_confusion_matrix, plot_roc_curve
3 from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
4 from sklearn.model_selection import cross_val_score
5 from sklearn.model_selection import GridSearchCV
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.neighbors import KNeighborsClassifier
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn import tree
10 from sklearn.svm import SVC
11 from sklearn.pipeline import Pipeline
12 from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score, roc_auc_score
13 from sklearn.model_selection import StratifiedKFold
14 from sklearn.preprocessing import StandardScaler, MinMaxScaler
15
16 #KNN, RandomForest, SVM
17
18 # Initializing Classifiers
19
20 clf1 = KNeighborsClassifier()
21 clf2 = RandomForestClassifier(n_estimators = 1024)
22 clf3 = LogisticRegression()
23 clf4 = tree.DecisionTreeClassifier()
24
25 # Building the pipelines
26
27 pipe1 = Pipeline([('std', StandardScaler()),
28                   ('classifier', clf1)])
29
30 pipe2 = Pipeline([('std', StandardScaler()),
31                   ('classifier', clf2)])
32
33 pipe3 = Pipeline([('std', StandardScaler()),
34                   ('classifier', clf3)])
35
36 pipe4 = Pipeline([('std', StandardScaler()),
37                   ('classifier', clf4)])
38
39 # Declaring some parameter values
40 C_list = np.power(10., np.arange(-8, 4))
41 D_list = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
42 F_list = [1, 2, 4, 6, 8, 12, 16, 20]
43 G_list = [0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 2]
44 K_list = [n*15 for n in range(1,21)]
45 penalty_list = ['l1', 'l2']
46 weight_list = ['uniform', 'distance']
47
```

```

48 # list of the param lists
49 all_param_lists = [C_list, D_list, F_list, G_list, G_list, K_list, penalty_list, weight_list]
50 knn_param_list = [weight_list, K_list]
51 lr_param_list =[penalty_list, C_list]
52
53 # Setting up the parameter grids
54
55 param_grid1 = [{classifier__weights: ['uniform', 'distance'],
56                 'classifier__n_neighbors': K_list}]
57
58 param_grid2= [{classifier__max_features: F_list}]
59
60 param_grid3 = [{classifier__C: C_list,
61                 'classifier__penalty': ['l1','l2']}]
62
63 param_grid4 = [{classifier__max_depth:D_list}]
64
65 # Setting up multiple GridSearchCV objects, 1 for each algorithm
66
67 scoring_metrics = ['accuracy','precision', 'recall', 'f1']
68 gridcvs = {}
69
70 for pgrid, est, name in zip((param_grid1, param_grid2, param_grid3, param_grid4),
71                             (pipe1, pipe2, pipe3, pipe4),
72                             ('KNN','RandomForest','Logistic', 'DecisionTree')):
73     gcv = GridSearchCV(estimator=est,
74                         param_grid=pgrid,
75                         scoring=scoring_metrics,
76                         n_jobs=3,
77                         cv=5, # 5-fold inner
78                         verbose=0,
79                         return_train_score=True,
80                         refit='accuracy')
81     gridcvs[name] = gcv
82
83 cv_scores = {name: [] for name, gs_est in gridcvs.items()}
84 skfold = StratifiedKFold(n_splits=5, shuffle=True, random_state=1)
85
86 def train_algo(gridcvs, X_train, y_train, folds):
87     """
88         Train the algorithm on the classifiers and parameters defined in the gridsearch and returns
89         parameters
90
91         gridcvs = the setup gridsearchCV to train the algorithms on
92         X_train, y_train = data split (random state of data varies by trial)
93         folds = number of k folds to perform on the data
94     """
95

```

```

96     cv_scores = {name: [] for name, gs_est in gridcvs.items()}
97     skfold = StratifiedKFold(n_splits=folds, shuffle=True, random_state=1)
98     # The outer loop for algorithm selection
99     c = 1
100    for outer_train_idx, outer_valid_idx in skfold.split(X_train,y_train):
101        for name, gs_est in sorted(gridcvs.items()):
102            print('outer fold %d/5 | tuning %-8s' % (c, name), end=' ')
103            # The inner loop for hyperparameter tuning
104            gs_est.fit(X_train[outer_train_idx], y_train[outer_train_idx])
105            y_pred = gs_est.predict(X_train[outer_valid_idx])
106            acc = accuracy_score(y_true=y_train[outer_valid_idx], y_pred=y_pred)
107            print(' | inner ACC %.2f%% | outer ACC %.2f%%' %
108                  (gs_est.best_score_ * 100, acc * 100))
109            cv_scores[name].append(acc)
110        c += 1
111
112    for name in cv_scores:
113        print('%-8s | outer CV acc. %.2f%% +\/- %.3f' % (
114            name, 100 * np.mean(cv_scores[name]), 100 * np.std(cv_scores[name])))
115    print()
116    for name in cv_scores:
117        print('{} best parameters'.format(name), gridcvs[name].best_params_)
118
119 def t1_algo(algo):
120
121     if algo =='KNN':
122         t1_algo_knn = gridcvs[algo]
123         return t1_algo_knn
124
125     elif algo == 'RandomForest':
126         t1_algo_rf = gridcvs[algo]
127         return t1_algo_rf
128
129     elif algo == 'Logistic':
130         t1_algo_lr = gridcvs[algo]
131         return t1_algo_lr
132
133     elif algo == 'DecisionTree':
134         t1_algo_dt = gridcvs[algo]
135         return t1_algo_dt
136
137 def t2_algo(algo):
138
139     if algo =='KNN':
140         t2_algo_knn = gridcvs[algo]
141         return t2_algo_knn
142
143     elif algo == 'RandomForest':

```

```

144     t2_algo_rf = gridcvs[algo]
145     return t2_algo_rf
146
147     elif algo == 'Logistic':
148         t2_algo_lr = gridcvs[algo]
149         return t2_algo_lr
150
151     elif algo == 'DecisionTree':
152         t2_algo_dt = gridcvs[algo]
153         return t2_algo_dt
154
155 def t3_algo(algo):
156
157     if algo =='KNN':
158         t3_algo_knn = gridcvs[algo]
159         return t3_algo_knn
160
161     elif algo == 'RandomForest':
162         t3_algo_rf = gridcvs[algo]
163         return t3_algo_rf
164
165     elif algo == 'Logistic':
166         t3_algo_lr = gridcvs[algo]
167         return t3_algo_lr
168
169     elif algo == 'DecisionTree':
170         t3_algo_dt = gridcvs[algo]
171         return t3_algo_dt
172
173 def fit_algo(trial_algo, algo, X_train, y_train, X_test, y_test, accu_dict, trial_accu, train_d
174 """
175     trial_algo = variable name associated with trial number and algorithm to store metrics, e.g.
176     algo = string of algo to get results of e.g. "SVM"
177     X_train, X_test, y_train, y_test = data split (random state of data varies by trial)
178     accu_dict = name of dictionary to store average accuracy across 5 folds
179     trial_accu = label for value stored in accu_dict
180     train_dict, test_dict = name of dictionaries to store train and test results
181 """
182
183     # Fitting model to the whole training set
184
185     #using the "best" algorithm
186     #trial_algo = gridcvs[algo]
187
188     # fitting the algorithm to training data
189     trial_algo.fit(X_train, y_train)
190     #gathering train/test accuracy scores
191     train_acc = accuracy_score(y_true=y_train, y_pred=trial_algo.predict(X_train))

```

```

192 test_acc = accuracy_score(y_true=y_test, y_pred=trial_algo.predict(X_test))
193
194 train_score_list = [train_acc]
195
196 test_score_list = [test_acc]
197
198 # print out results
199 print('Accuracy %.2f%% (average over CV test folds)' % (100 * trial_algo.best_score_))
200 print('Best Parameters: %s' % gridcvs[algo].best_params_)
201 print('Training Accuracy: %.2f%%' % (100 * train_acc))
202 print('Test Accuracy: %.2f%%' % (100 * test_acc))
203
204
205 # store accuracy result average over CV test folds
206 accu_dict[trial_accu] = trial_algo.best_score_
207
208 # store results in dictionary using the train_test_results function
209 train_test_results(train_dict, algo, train_metric_list, train_score_list)
210 train_test_results(test_dict, algo, test_metric_list, test_score_list)
211
212 #optimal parameters knn
213 def optimized(weight, K, F, C, D, penalty, X_train, y_train, X_test, y_test, train_dict, test_dict):
214 """
215 weight = optimal weight parameter for KNN
216 K = optimal number of n_neighbors for KNN
217 F = optimal max_features for Random Forest
218 C = optimal C value for Logistic Regression
219 D = optimal max_depth of Decision Tree
220 penalty = optimal 'l1' or 'l2' for Logistic Regression
221 """
222
223 clf_knn = KNeighborsClassifier(weights=weight, n_neighbors=K)
224 clf_knn.fit(X_train, y_train)
225 knn_accuracy = cross_val_score(clf_knn, X_train, y_train)
226 pred_knn = clf_knn.predict(X_test)
227 print('KNN Train Accuracy', np.mean(knn_accuracy))
228 print('KNN Test Accuracy', clf_knn.score(X_test, y_test))
229 print(classification_report(y_test, pred_knn))
230 train_dict['KNN Train Accuracy'] = np.mean(knn_accuracy)
231 test_dict['KNN Test Accuracy'] = clf_knn.score(X_test, y_test)
232
233 clf_rf = RandomForestClassifier(n_estimators = 1024, max_features=F)
234 clf_rf.fit(X_train, y_train)
235 rf_accuracy = cross_val_score(clf_rf, X_train, y_train)
236 pred_rf = clf_rf.predict(X_test)
237 print('Random Forest Train Accuracy', np.mean(rf_accuracy))
238 print('Random Forest Test Accuracy', clf_rf.score(X_test, y_test))
239 print(classification_report(y_test, pred_rf))

```

```

240     train_dict['Random Forest Train Accuracy'] = np.mean(rf_accuracy)
241     test_dict['Random Forest Test Accuracy'] = clf_rf.score(X_test, y_test)
242
243     clf_lr = LogisticRegression(penalty=penalty, C = C)
244     clf_lr.fit(X_train,y_train)
245     lr_accuracy = cross_val_score(clf_lr, X_train,y_train)
246     pred_lr = clf_lr.predict(X_test)
247     print('Logistic Train Accuracy',np.mean(lr_accuracy))
248     print('Logistic Test Accuracy',clf_lr.score(X_test, y_test))
249     print(classification_report(y_test, pred_lr))
250     train_dict['Logistic Train Accuracy'] = np.mean(lr_accuracy)
251     test_dict['Logistic Test Accuracy'] = clf_lr.score(X_test, y_test)
252
253     clf_dt = tree.DecisionTreeClassifier(max_depth = D)
254     clf_dt.fit(X_train,y_train)
255     dt_accuracy = cross_val_score(clf_dt, X_train,y_train)
256     pred_dt = clf_dt.predict(X_test)
257     print('Decision Tree Train Accuracy',np.mean(dt_accuracy))
258     print('Decision Tree Test Accuracy',clf_dt.score(X_test, y_test))
259     print(classification_report(y_test, pred_dt))
260     train_dict['Decision Tree Train Accuracy'] = np.mean(dt_accuracy)
261     test_dict['Decision Tree Test Accuracy'] = clf_dt.score(X_test, y_test)
262
263     train_metric_list = ['Training Accuracy','Training Precision','Training Recall','Training F1',''
264     test_metric_list = ['Test Accuracy','Test Precision','Test Recall','Test F1','Test ROC AUC']
265     #train_score_list = [train_acc, train_prec, train_recall, train_f1, train_roc]
266     #test_score_list = [test_acc, train_prec, test_prec, test_recall, test_f1, test_roc]
267
268     def train_test_results(dict_name,algo, metric_list, score_list):
269         """
270             dict_name = name of dictionary to put results
271             algo = defined grid string of algo name to be used e.g. 'KNN','RandomForest', 'SVM', 'Logis
272             metric_list = train/test metric list
273             score_list = train/score lists defined in model fitting
274         """
275
276         for score, metric in zip(score_list, metric_list):
277             dict_name[algo+' '+metric] = score
278
279     #train_string_list = ['mean_train_accuracy','mean_train_precision','mean_train_recall','mean_tr
280     #                      'mean_train_roc_auc',]
281
282     #test_string_list = ['mean_test_accuracy','mean_test_precision','mean_test_recall','mean_test_f
283     #                      'mean_test_roc_auc']
284
285     #train_metric_list = ['Training Accuracy','Training Precision','Training Recall','Training F1',
286     #test_metric_list = ['Test Accuracy','Test Precision','Test Recall','Test F1','Test ROC AUC']
287
288     def optimal_results(trial_algo, param_list, algo, param, dict_name, string_list, metric_list):

```

```

288
289     trial_algo = attained algorithm results for trial, e.g. t1_algo_svm => gridcvs[SVM]
290     param_list = input optimal parameter's list for algo e.g C_list, K_list, D_list ect..
291     algo = defined grid string of algo name to be used e.g. 'KNN', 'RandomForest', 'SVM', 'Logis
292     param = corresponding parameter name to param_list e.g. C_list => 'classifier__C'
293     dict_name = name of defined dictionary to store results in
294     string_list = train_string_list for training results and test_string_list for testing resul
295     metric_list = train_metric list for training and test_metric_list for testing
296     """
297
298     for i,j in enumerate(param_list):
299         if j == gridcvs[algo].best_params_[param]:
300             for item, metric in zip(string_list, metric_list):
301                 dict_name[algo+' '+metric] = trial_algo.cv_results_[item][i]
302
303     # Example Function Input:
304     #optimal_results(t1_algo_svm, C_list,'SVM', 'classifier__C',t1_optimal_train_dict, train_string
305
306     # Example Function Input:
307     #train_test_results(t1_train_dict, 'SVM', train_metric_list, train_score_list)
308     #train_test_results(t1_test_dict, 'SVM', test_metric_list, test_score_list)
309
310     # function to draw heat map of hyperparameters against performance
311     def draw_heatmap(acc, acc_desc, C_list, character):
312
313         plt.figure(figsize = (2,4))
314         ax = sns.heatmap(acc, annot=True, fmt='%.3f', yticklabels=C_list, xticklabels[])
315         ax.collections[0].colorbar.set_label("accuracy")
316         ax.set(ylabel='$' + character + '$')
317         plt.title(acc_desc + ' w.r.t $' + character + '$')
318         #sns.set_style("whitegrid", {'axes.grid' : False})
319         b, t = plt.ylim()
320         b += 0.5
321         t -= 0.5
322         custom_ylim = (b, t)
323         plt.setp(ax, ylim=custom_ylim)
324         plt.show()

```

Trial 1

```
In [29]: # To store results of models in trial 1
t1_accu_dict = {}

t1_train_dict = {}
t1_test_dict = {}

t1_optimal_train_dict = {}
t1_optimal_test_dict = {}
```

```
In [30]: from sklearn.model_selection import train_test_split
# split the data into testing and training sets
X1_train, X1_test, y1_train, y1_test = train_test_split(X, y, train_size = 0.8, test_size=0.20, random_state=42)

print("Shape of input data X_train: {} and shape of target variable y_train: {}".format(X1_train.shape, y1_train.shape))
print("Shape of input data X_test: {} and shape of target variable y_test: {}".format(X1_test.shape, y1_test.shape))
```

Accuracy Performance

Model 'refit' set to accuracy

```
In [31]: %%time
# run the train_algo function and train the data on the models
train_algo(gridcvs, X1_train, y1_train, 5)
```

```
outer fold 1/5 | tuning DecisionTree | inner ACC 90.64% | outer ACC 7  
2.19%  
outer fold 1/5 | tuning KNN       | inner ACC 85.13% | outer ACC 82.4  
5%  
outer fold 1/5 | tuning Logistic | inner ACC 85.23% | outer ACC 82.7  
8%  
outer fold 1/5 | tuning RandomForest | inner ACC 83.60% | outer ACC 8  
1.79%  
outer fold 2/5 | tuning DecisionTree | inner ACC 90.55% | outer ACC 6  
8.54%  
outer fold 2/5 | tuning KNN       | inner ACC 85.12% | outer ACC 81.4  
6%  
outer fold 2/5 | tuning Logistic | inner ACC 85.12% | outer ACC 84.1  
1%  
outer fold 2/5 | tuning RandomForest | inner ACC 83.38% | outer ACC 8  
2.45%  
outer fold 3/5 | tuning DecisionTree | inner ACC 90.41% | outer ACC 6  
6.56%  
outer fold 3/5 | tuning KNN       | inner ACC 85.47% | outer ACC 78.4  
8%  
outer fold 3/5 | tuning Logistic | inner ACC 85.72% | outer ACC 78.8  
1%  
outer fold 3/5 | tuning RandomForest | inner ACC 84.66% | outer ACC 7  
7.81%  
outer fold 4/5 | tuning DecisionTree | inner ACC 89.71% | outer ACC 7  
0.10%  
outer fold 4/5 | tuning KNN       | inner ACC 85.30% | outer ACC 79.0  
7%  
outer fold 4/5 | tuning Logistic | inner ACC 85.13% | outer ACC 80.0  
7%  
outer fold 4/5 | tuning RandomForest | inner ACC 84.12% | outer ACC 7  
8.41%  
outer fold 5/5 | tuning DecisionTree | inner ACC 90.45% | outer ACC 7  
0.43%  
outer fold 5/5 | tuning KNN       | inner ACC 85.42% | outer ACC 79.4  
0%  
outer fold 5/5 | tuning Logistic | inner ACC 85.42% | outer ACC 80.7  
3%  
outer fold 5/5 | tuning RandomForest | inner ACC 84.32% | outer ACC 8  
0.73%  
KNN       | outer CV acc. 80.17% +\-\ 1.518  
RandomForest | outer CV acc. 80.24% +\-\ 1.831  
Logistic | outer CV acc. 81.30% +\-\ 1.904  
DecisionTree | outer CV acc. 69.56% +\-\ 1.897
```

```
KNN best parameters {'classifier__n_neighbors': 75, 'classifier__weig  
hts': 'distance'}  
RandomForest best parameters {'classifier__max_features': 12}  
Logistic best parameters {'classifier__C': 1.0, 'classifier__penalty  
': 'l2'}  
DecisionTree best parameters {'classifier__max_depth': 1}  
Wall time: 2min 39s
```

In [32]:

```
%time
# plot learning curves for training
plot_learning_curve(gridcvs['KNN'], 'KNN Learning Curve', X1_train, y1_train);
# run the fit algo function and fit all the training data to KNN
fit_algo(t1_algo('KNN'), 'KNN', X1_train, y1_train, X1_test, y1_test, t1_accu_dict, 'KNN Accuracy', t1_train_dict, t1_test_dict)
# plot confusion matrix on the test samples
plot_confusion_matrix(gridcvs['KNN'], X1_test, y1_test, cmap = 'magma')
# plot ROC curve on the test samples
plot_roc_curve(gridcvs['KNN'], X1_test, y1_test, c ='red');
```

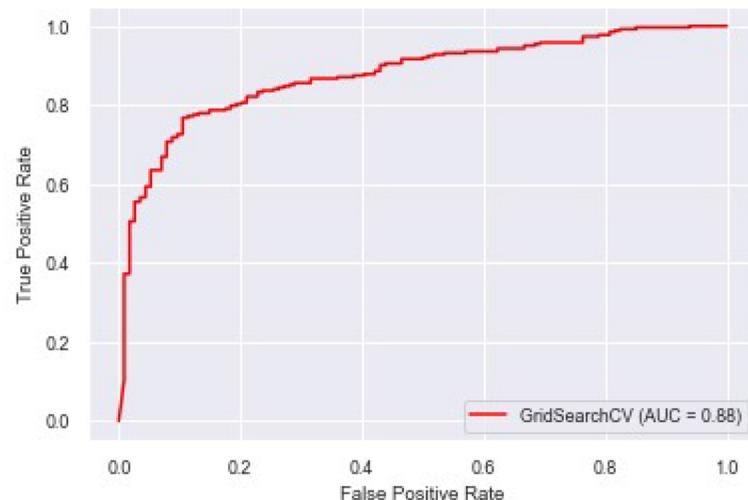
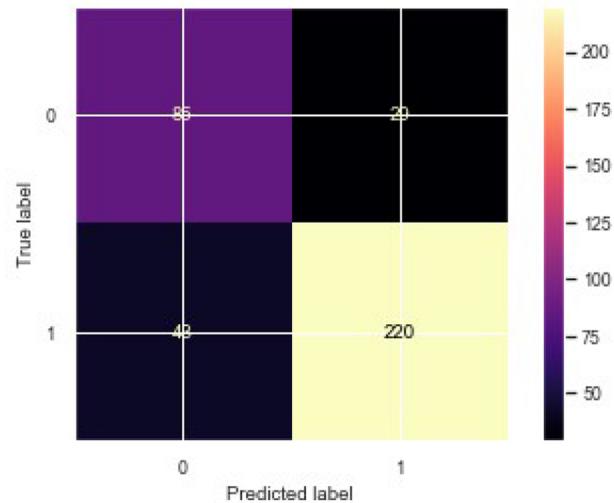
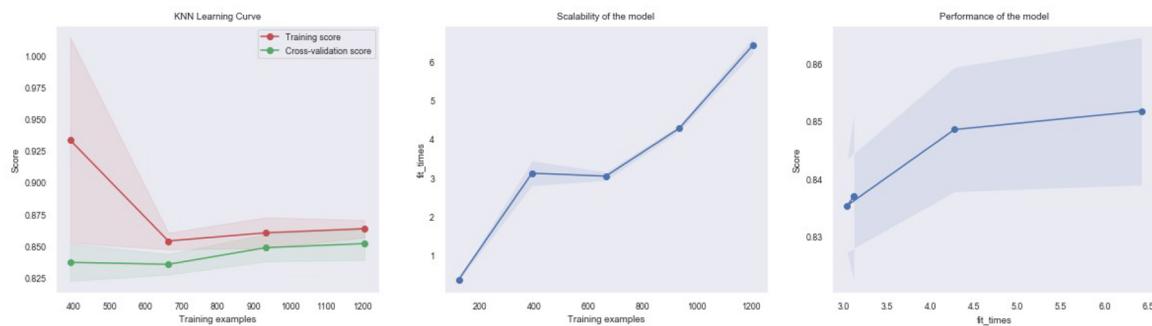
Accuracy 81.30% (average over CV test folds)

Training Accuracy: 100.00%

Test Accuracy: 79.58%

Wall time: 1min 35s

Out[32]: <sklearn.metrics._plot.roc_curve.RocCurveDisplay at 0x1a60f56a188>



In [33]:

```
%%time
# plot learning curves for training
plot_learning_curve(gridcvs['RandomForest'], 'Random Forest Learning Curve', X1_train, y1_train);
# run the fit algo function and fit all the training data to RF
fit_algo(t1_algo('RandomForest'), 'RandomForest', X1_train, y1_train, X1_test, y1_test, t1_accu_dict, 'RandomForest Accuracy', t1_train_dict, t1_test_dict)
# plot confusion matrix on the test samples
plot_confusion_matrix(gridcvs['RandomForest'], X1_test, y1_test, cmap = 'magma')
# plot ROC curve on the test samples
plot_roc_curve(gridcvs['RandomForest'], X1_test, y1_test, c ='red');
```

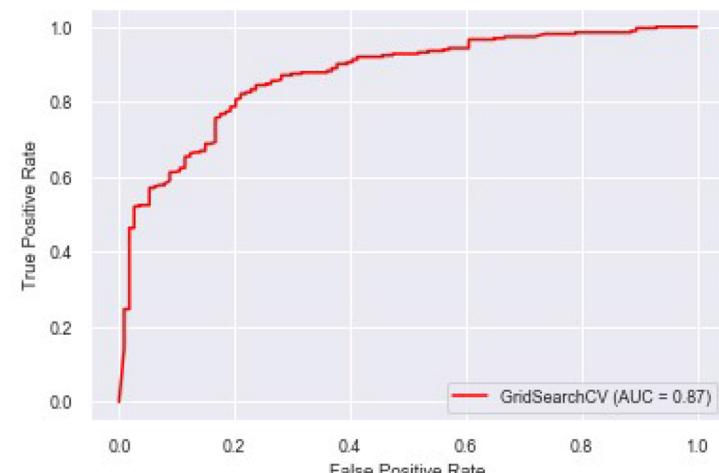
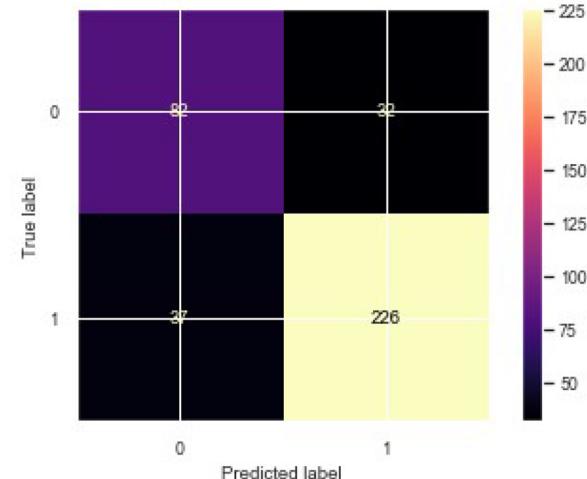
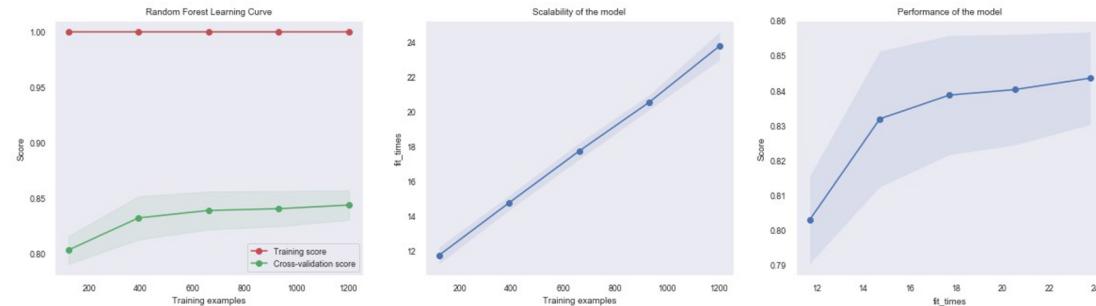
Accuracy 81.56% (average over CV test folds)

Training Accuracy: 100.00%

Test Accuracy: 80.11%

Wall time: 7min 54s

Out[33]: <sklearn.metrics._plot.roc_curve.RocCurveDisplay at 0x1a60f682f88>



In [34]:

```
%time
# plot learning curves for training
plot_learning_curve(gridcvs['Logistic'], 'Logistic Regression Learning Curve', X1_train, y1_train);
# run the fit algo function and fit all the training data to Logistic Regression
fit_algo(t1_algo('Logistic'), 'Logistic', X1_train, y1_train, X1_test,
y1_test, t1_accu_dict,'Logistic Accuracy', t1_train_dict, t1_test_dict)
# plot confusion matrix on the test samples
plot_confusion_matrix(gridcvs['Logistic'], X1_test, y1_test, cmap = 'magma')
# plot ROC curve on the test samples
plot_roc_curve(gridcvs['Logistic'], X1_test, y1_test, c ='red');
```

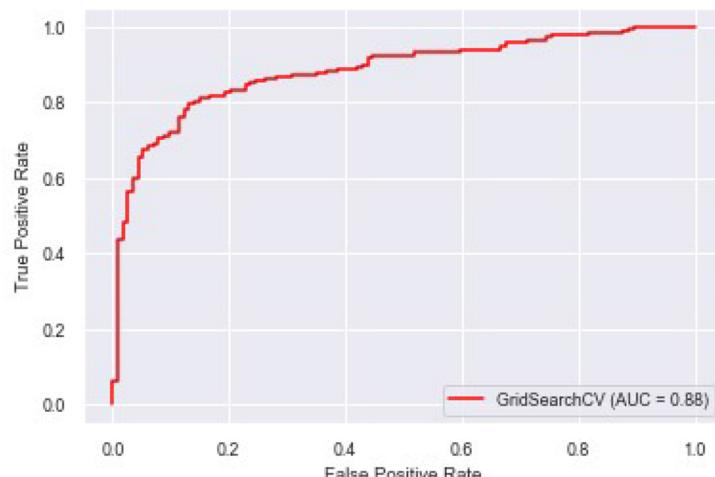
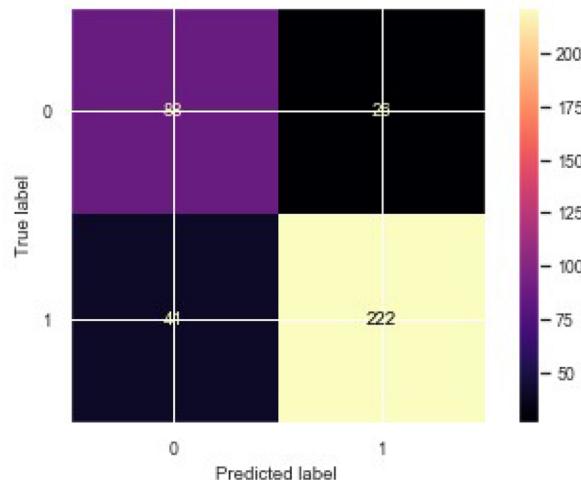
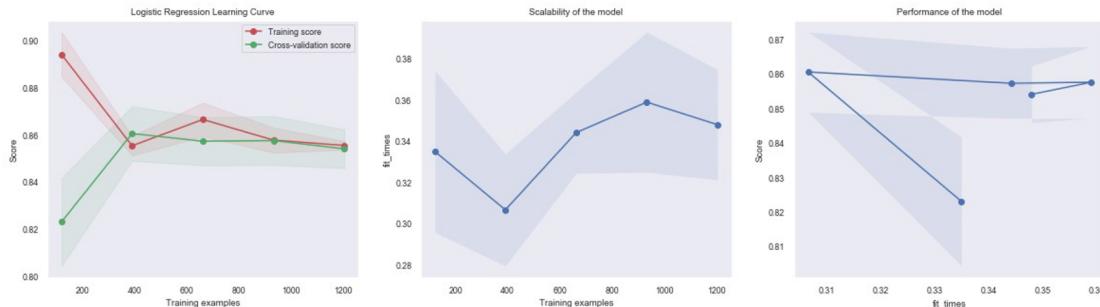
Accuracy 81.63% (average over CV test folds)

Training Accuracy: 81.83%

Test Accuracy: 80.37%

Wall time: 8.93 s

Out[34]: <sklearn.metrics._plot.roc_curve.RocCurveDisplay at 0x1a60f7d588>



In [35]:

```
%%time
# plot learning curves for training
plot_learning_curve(gridcvs['DecisionTree'], 'Decision Tree Learning Curve',
                     X1_train, y1_train);
# run the fit algo function and fit all the training data to Decision Tree
fit_algo(t1_algo('DecisionTree'), 'DecisionTree', X1_train, y1_train, X1_test,
          y1_test, t1_accu_dict, 'DecisionTree Accuracy', t1_train_dict, t1_test_dict)
# plot confusion matrix on the test samples
plot_confusion_matrix(gridcvs['DecisionTree'], X1_test, y1_test, cmap =
'magma')
# plot ROC curve on the test samples
plot_roc_curve(gridcvs['DecisionTree'], X1_test, y1_test, c ='red');
```

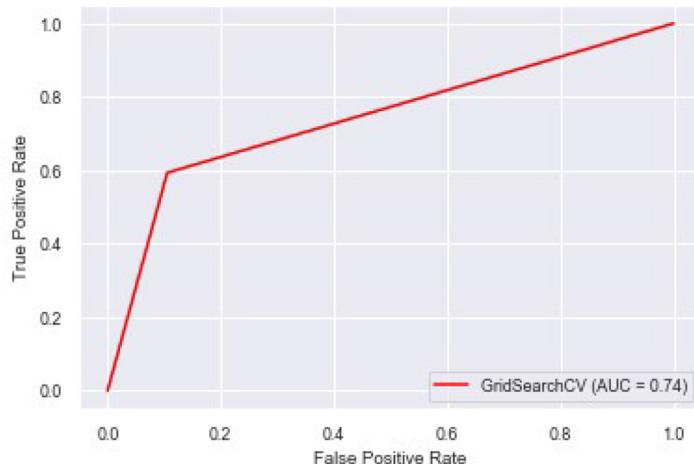
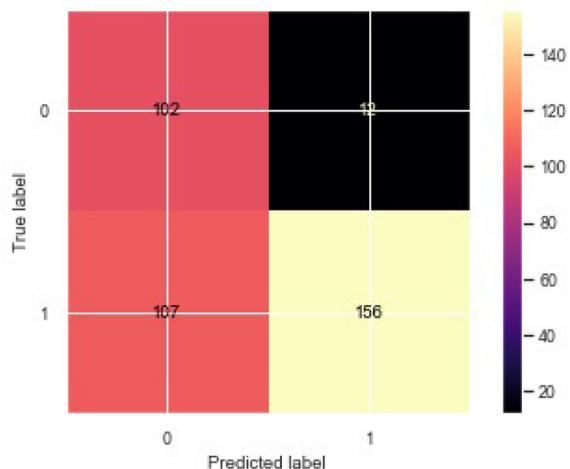
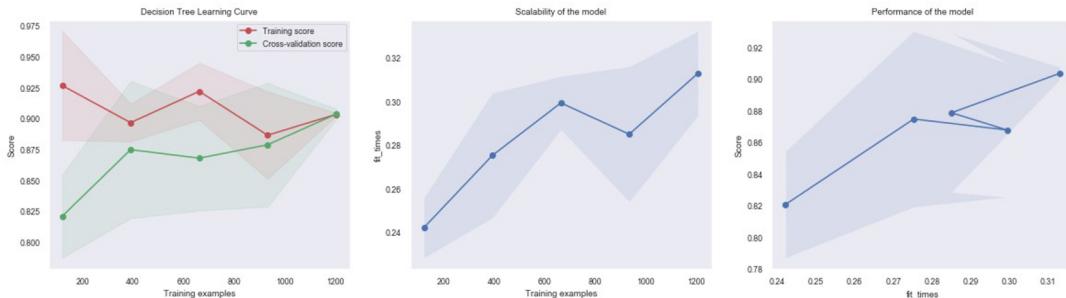
Accuracy 78.71% (average over CV test folds)

Training Accuracy: 81.43%

Test Accuracy: 77.45%

Wall time: 7.52 s

Out[35]: <sklearn.metrics._plot.roc_curve.RocCurveDisplay at 0x1a6102cbf88>



In [37]:

```
%%time
# run the optimized function and apply optimal parameters to algorithms
# and look at training/testing scores
optimized('distance', 75, 12, 1, 1, '12', X1_train, y1_train, X1_test,
y1_test,t1_optimal_train_dict,t1_optimal_test_dict)
```

KNN Train Accuracy 0.8010450815163583

KNN Test Accuracy 0.7771883289124668

	precision	recall	f1-score	support
0	0.62	0.69	0.65	114
1	0.86	0.81	0.84	263
accuracy			0.78	377
macro avg	0.74	0.75	0.74	377
weighted avg	0.79	0.78	0.78	377

Random Forest Train Accuracy 0.8103298057248465

Random Forest Test Accuracy 0.8090185676392573

	precision	recall	f1-score	support
0	0.67	0.73	0.70	114
1	0.88	0.84	0.86	263
accuracy			0.81	377
macro avg	0.77	0.79	0.78	377
weighted avg	0.81	0.81	0.81	377

Logistic Train Accuracy 0.8156212184550394

Logistic Test Accuracy 0.8010610079575596

	precision	recall	f1-score	support
0	0.68	0.66	0.67	114
1	0.85	0.86	0.86	263
accuracy			0.80	377
macro avg	0.76	0.76	0.76	377
weighted avg	0.80	0.80	0.80	377

Decision Tree Train Accuracy 0.78645134320477

Decision Tree Test Accuracy 0.7745358090185677

	precision	recall	f1-score	support
0	0.61	0.72	0.66	114
1	0.87	0.80	0.83	263
accuracy			0.77	377
macro avg	0.74	0.76	0.75	377
weighted avg	0.79	0.77	0.78	377

Wall time: 11.2 s

Precision Performance

Model 'refit' set to precision

```
[38]: train_algo_prec(gridcvs, X1_train, y1_train, 5)
```

```
outer fold 1/5 | tuning DecisionTree | inner prec 90.64% | outer prec  
89.80%  
outer fold 1/5 | tuning KNN       | inner prec 85.13% | outer prec 87.  
00%  
outer fold 1/5 | tuning Logistic | inner prec 85.23% | outer prec 85.  
31%  
outer fold 1/5 | tuning RandomForest | inner prec 84.08% | outer prec  
83.80%  
outer fold 2/5 | tuning DecisionTree | inner prec 90.55% | outer prec  
89.55%  
outer fold 2/5 | tuning KNN       | inner prec 85.12% | outer prec 84.  
04%  
outer fold 2/5 | tuning Logistic | inner prec 85.12% | outer prec 84.  
93%  
outer fold 2/5 | tuning RandomForest | inner prec 83.32% | outer prec  
83.78%  
outer fold 3/5 | tuning DecisionTree | inner prec 90.41% | outer prec  
89.60%  
outer fold 3/5 | tuning KNN       | inner prec 85.47% | outer prec 86.  
89%  
outer fold 3/5 | tuning Logistic | inner prec 85.72% | outer prec 86.  
17%  
outer fold 3/5 | tuning RandomForest | inner prec 84.54% | outer prec  
83.51%  
outer fold 4/5 | tuning DecisionTree | inner prec 89.71% | outer prec  
92.97%  
outer fold 4/5 | tuning KNN       | inner prec 85.30% | outer prec 86.  
63%  
outer fold 4/5 | tuning Logistic | inner prec 85.13% | outer prec 87.  
23%  
outer fold 4/5 | tuning RandomForest | inner prec 84.08% | outer prec  
85.20%  
outer fold 5/5 | tuning DecisionTree | inner prec 90.45% | outer prec  
89.93%  
outer fold 5/5 | tuning KNN       | inner prec 85.42% | outer prec 83.  
17%  
outer fold 5/5 | tuning Logistic | inner prec 85.42% | outer prec 84.  
47%  
outer fold 5/5 | tuning RandomForest | inner prec 84.33% | outer prec  
81.90%  
KNN      | outer CV prec. 85.55% +/- 1.612  
RandomForest | outer CV prec. 83.64% +/- 1.052  
Logistic | outer CV prec. 85.62% +/- 0.981  
DecisionTree | outer CV prec. 90.37% +/- 1.307  
  
KNN best parameters {'classifier__n_neighbors': 75, 'classifier__weig  
hts': 'distance'}  
RandomForest best parameters {'classifier__max_features': 12}  
Logistic best parameters {'classifier__C': 1.0, 'classifier__penalty  
': 'l2'}  
DecisionTree best parameters {'classifier__max_depth': 1}
```

In [39]:

```
%time
# plot learning curves for training
plot_learning_curve(gridcvs['KNN'], 'KNN Learning Curve', X1_train, y1_train);
# run the fit algo function and fit all the training data to KNN
fit_algo_prec(tl_algo('KNN'), 'KNN', X1_train, y1_train, X1_test, y1_test,
t1_accu_dict, 'KNN Precision', t1_train_dict, t1_test_dict)
# plot confusion matrix on the test samples
plot_confusion_matrix(gridcvs['KNN'], X1_test, y1_test, cmap = 'magma')
# plot ROC curve on the test samples
plot_roc_curve(gridcvs['KNN'], X1_test, y1_test, c ='red');
```

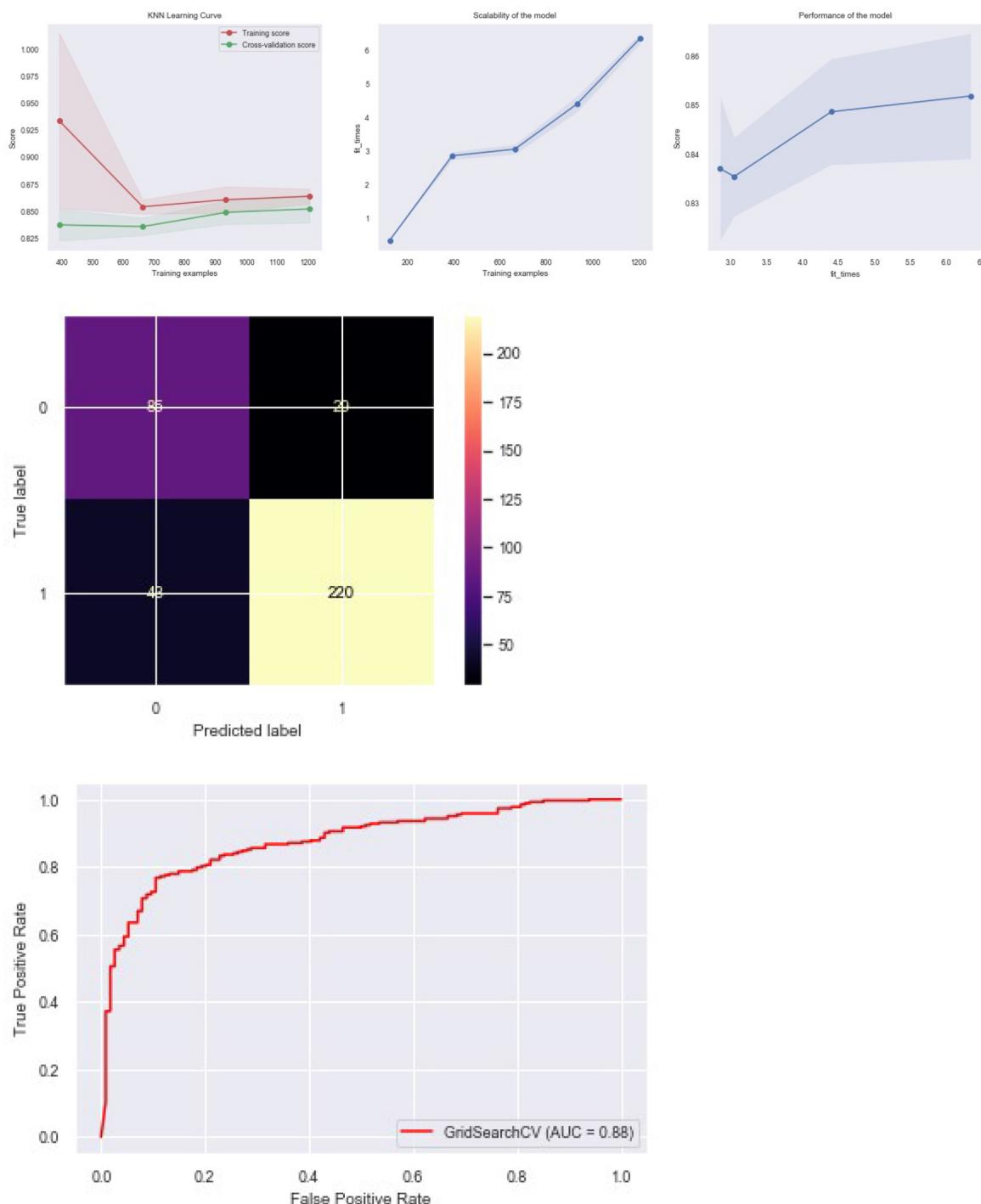
Precision 81.30% (average over CV test folds)

Training Precision: 100.00%

Test Precision: 87.50%

Wall time: 1min 35s

Out[39]: <sklearn.metrics._plot.roc_curve.RocCurveDisplay at 0x1a610226748>



In [40]:

```
%time
# plot learning curves for training
plot_learning_curve(gridcvs['RandomForest'], 'Random Forest Learning Curve',
                     X1_train, y1_train);
# run the fit algo function and fit all the training data to KNN
fit_algo_prec(t1_algo('RandomForest'), 'RandomForest', X1_train, y1_train,
               X1_test, y1_test, t1_accu_dict, 'Random Forest Precision', t1_train_
               dict, t1_test_dict)
# plot confusion matrix on the test samples
plot_confusion_matrix(gridcvs['RandomForest'], X1_test, y1_test, cmap =
'magma')
# plot ROC curve on the test samples
plot roc curve(gridcvs['RandomForest'], X1 test, y1 test, c ='red');
```

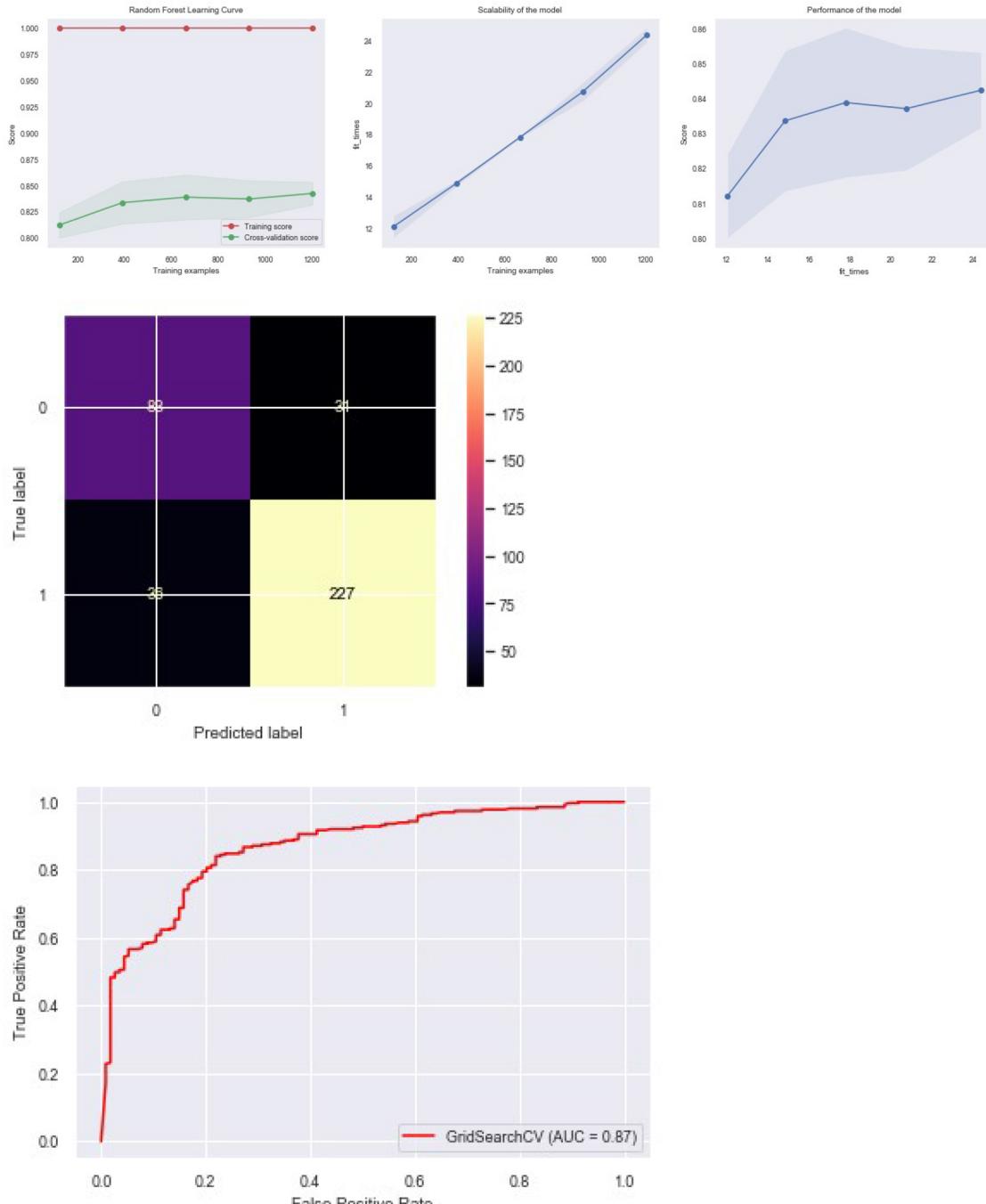
Precision 81.56% (average over CV test folds)

Training Precision: 100.00%

Test Precision: 87.75%

Wall time: 8min 3s

Out[40]: <sklearn.metrics._plot.roc_curve.RocCurveDisplay at 0x1a6102aa648>



In [41]:

```
%%time
# plot learning curves for training
plot_learning_curve(gridcvs['Logistic'], 'Logistic Regression Learning Curve', X1_train, y1_train);
# run the fit algo function and fit all the training data to KNN
fit_algo_prec(t1_algo('Logistic'), 'Logistic', X1_train, y1_train, X1_test, y1_test, t1_accu_dict, 'Logistic Precision', t1_train_dict, t1_test_dict)
# plot confusion matrix on the test samples
plot_confusion_matrix(gridcvs['Logistic'], X1_test, y1_test, cmap = 'magma')
# plot ROC curve on the test samples
plot_roc_curve(gridcvs['Logistic'], X1_test, y1_test, c ='red');
```

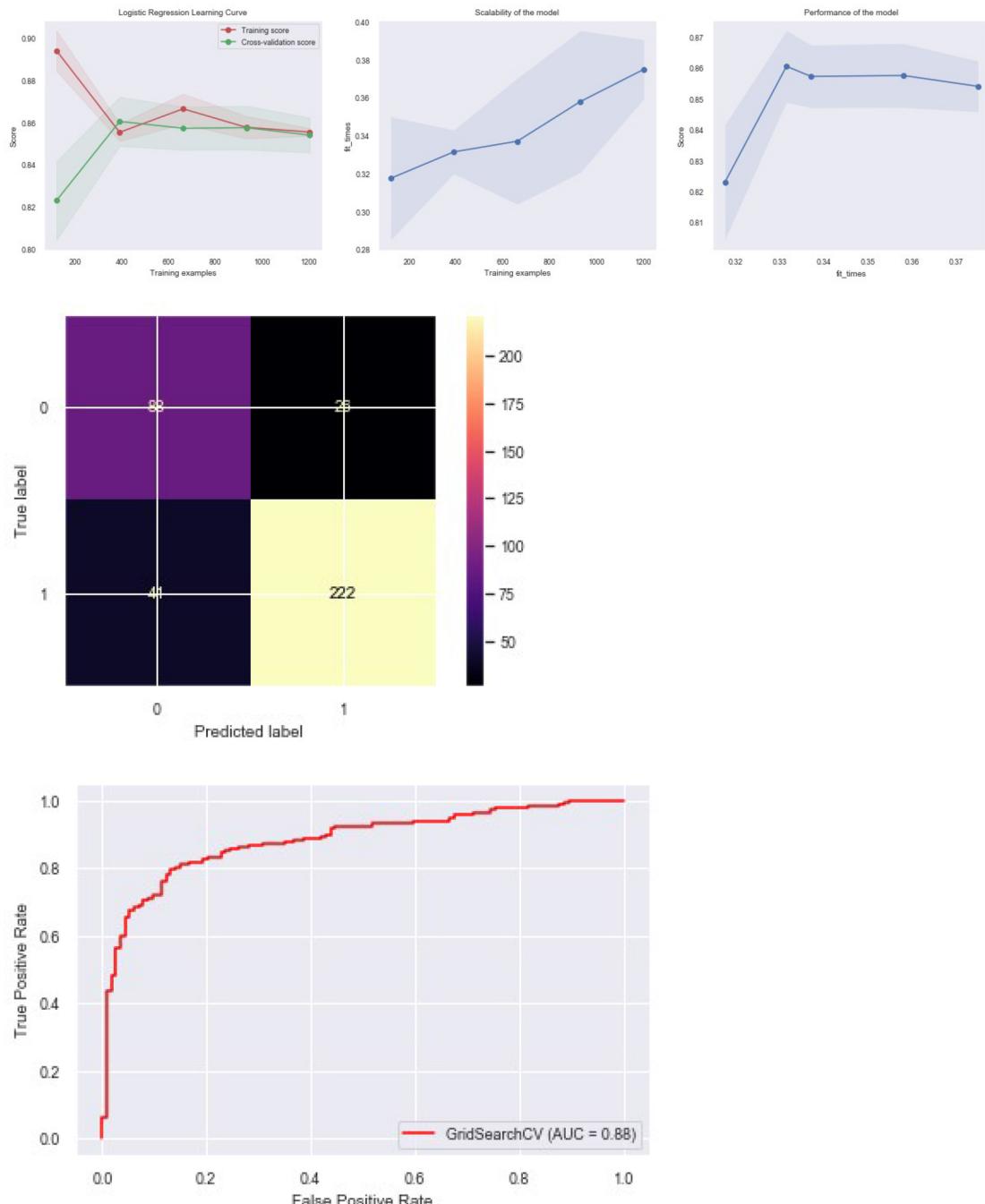
Precision 81.63% (average over CV test folds)

Training Precision: 84.60%

Test Precision: 85.93%

Wall time: 9.09 s

Out[41]: <sklearn.metrics._plot.roc_curve.RocCurveDisplay at 0x1a60f7bdbcc8>



In [42]:

```
%time
# plot learning curves for training
plot_learning_curve(gridcvs['DecisionTree'], 'Decision Tree Learning Curve',
                     X1_train, y1_train);
# run the fit algo function and fit all the training data to KNN
fit_algo_prec(t1_algo('DecisionTree'), 'DecisionTree', X1_train, y1_train,
               X1_test, y1_test, t1_accu_dict, 'Decision Tree Precision', t1_train_
               dict, t1_test_dict)
# plot confusion matrix on the test samples
plot_confusion_matrix(gridcvs['DecisionTree'], X1_test, y1_test, cmap =
'magma')
# plot ROC curve on the test samples
plot_roc_curve(gridcvs['DecisionTree'], X1_test, y1_test, c ='red');
```

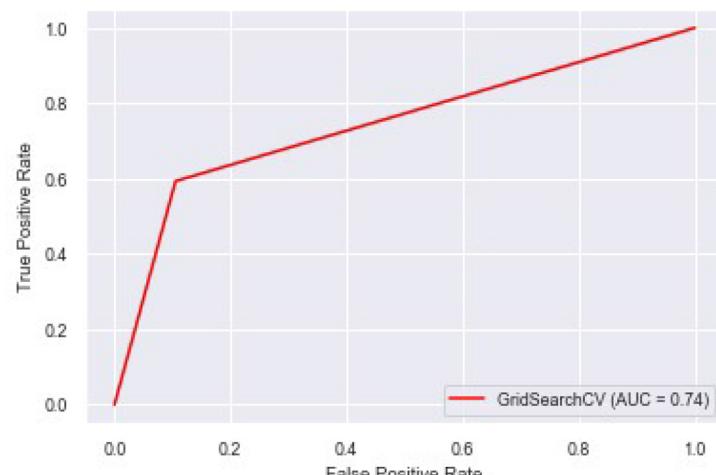
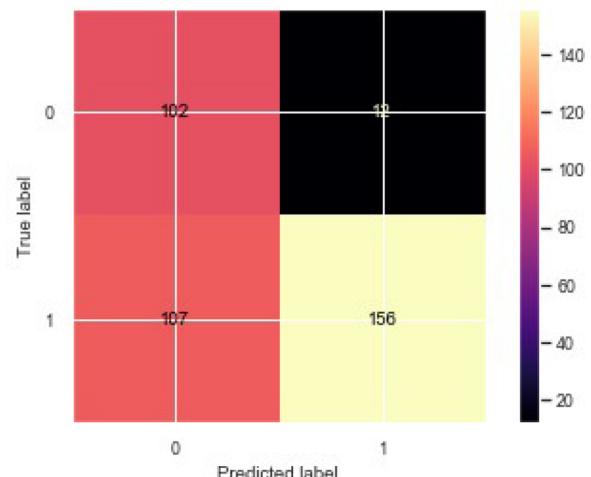
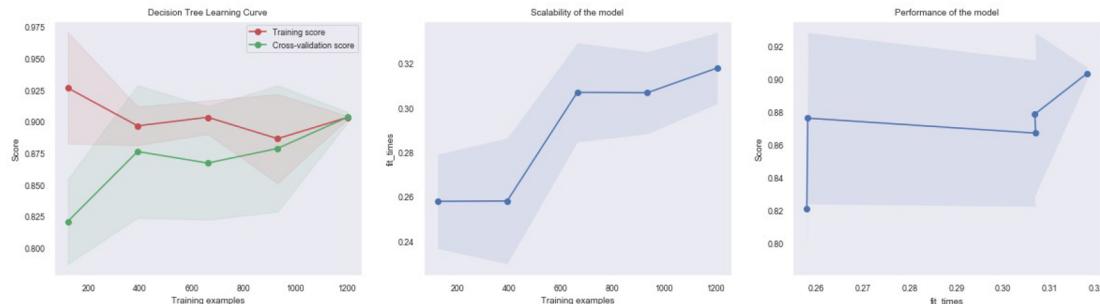
Precision 78.58% (average over CV test folds)

Training Precision: 86.10%

Test Precision: 86.78%

Wall time: 7.7 s

Out[42]: <sklearn.metrics._plot.roc_curve.RocCurveDisplay at 0x1a6189b2ec8>



```
In [43]: optimized_prec('uniform', 75, 12, 1, 1, '12', X1_train, y1_train, X1_test, y1_test,t1_optimal_train_dict,t1_optimal_test_dict)
```

KNN Train Precision 0.857566765578635

KNN Test Precision 0.8870967741935484

	precision	recall	f1-score	support
0	0.67	0.75	0.71	114
1	0.89	0.84	0.86	263
accuracy			0.81	377
macro avg	0.78	0.80	0.78	377
weighted avg	0.82	0.81	0.81	377

Random Forest Train Precision 1.0

Random Forest Test Precision 0.8844621513944223

	precision	recall	f1-score	support
0	0.67	0.75	0.71	114
1	0.88	0.84	0.86	263
accuracy			0.81	377
macro avg	0.78	0.79	0.79	377
weighted avg	0.82	0.81	0.82	377

Logistic Train Precision 0.8553149606299213

Logistic Test Precision 0.8906882591093117

	precision	recall	f1-score	support
0	0.67	0.76	0.71	114
1	0.89	0.84	0.86	263
accuracy			0.81	377
macro avg	0.78	0.80	0.79	377
weighted avg	0.82	0.81	0.82	377

Decision Tree Train Precision 0.9034175334323923

Decision Tree Test Precision 0.9285714285714286

	precision	recall	f1-score	support
0	0.49	0.89	0.63	114
1	0.93	0.59	0.72	263
accuracy			0.68	377
macro avg	0.71	0.74	0.68	377
weighted avg	0.80	0.68	0.70	377

Results: Trial 1

```
In [44]: t1_df_accu_dict = pd.DataFrame.from_dict(t1_accu_dict,orient = "index", columns=["Score"]).sort_values(by=["Score"], ascending=False)
t1_df_train_dict = pd.DataFrame.from_dict(t1_train_dict,orient = "index",columns=["Score"]).sort_values(by=["Score"], ascending=False)
t1_df_test_dict = pd.DataFrame.from_dict(t1_test_dict,orient = "index",columns=["Score"]).sort_values(by=["Score"], ascending=False)
t1_df_optimal_train_dict = pd.DataFrame.from_dict(t1_optimal_train_dict,orient = "index",columns=["Score"]).sort_values(by=["Score"], ascending=False)
t1_df_optimal_test_dict = pd.DataFrame.from_dict(t1_optimal_test_dict,orient = "index",columns=["Score"]).sort_values(by=["Score"], ascending=False)
```

```
In [45]: print("Average Accuracy Across Folds")
print(t1_df_accu_dict)
print("")
print("Training Scores")
print(t1_df_train_dict)
print("")
print("Testing Scores")
print(t1_df_test_dict)
print("")
print("Training Scores For Optimal Parameter")
print(t1_df_optimal_train_dict)
print("")
print("Testing Scores For Optimal Parameter")
print(t1_df_optimal_test_dict)
```

Average Accuracy Across Folds

	Score
Logistic Accuracy	0.816283
Logistic Precision	0.816283
Random Forest Precision	0.815630
RandomForest Accuracy	0.815628
KNN Accuracy	0.812979
KNN Precision	0.812979
DecisionTree Accuracy	0.787116
Decision Tree Precision	0.785787

Training Scores

	Score
KNN Training Accuracy	1.000000
RandomForest Training Accuracy	1.000000
KNN Train Precision	1.000000
RandomForest Train Precision	1.000000
DecisionTree Train Precision	0.861000
Logistic Train Precision	0.846008
Logistic Training Accuracy	0.818302
DecisionTree Training Accuracy	0.814324

Testing Scores

	Score
RandomForest Test Precision	0.877470
KNN Test Precision	0.875000
DecisionTree Test Precision	0.867769
Logistic Test Precision	0.859316
Logistic Test Accuracy	0.803714
RandomForest Test Accuracy	0.801061
KNN Test Accuracy	0.795756
DecisionTree Test Accuracy	0.774536

Training Scores For Optimal Parameter

	Score
Random Forest Train Precision	1.000000
Decision Tree Train Precision	0.903418
KNN Train Precision	0.857567
Logistic Train Precision	0.855315
Logistic Train Accuracy	0.815621
Random Forest Train Accuracy	0.810330
KNN Train Accuracy	0.801045
Decision Tree Train Accuracy	0.786451

Testing Scores For Optimal Parameter

	Score
Decision Tree Test Precision	0.928571
Logistic Test Precision	0.890688
KNN Test Precision	0.887097
Random Forest Test Precision	0.884462
Random Forest Test Accuracy	0.809019
Logistic Test Accuracy	0.801061
KNN Test Accuracy	0.777188
Decision Tree Test Accuracy	0.774536

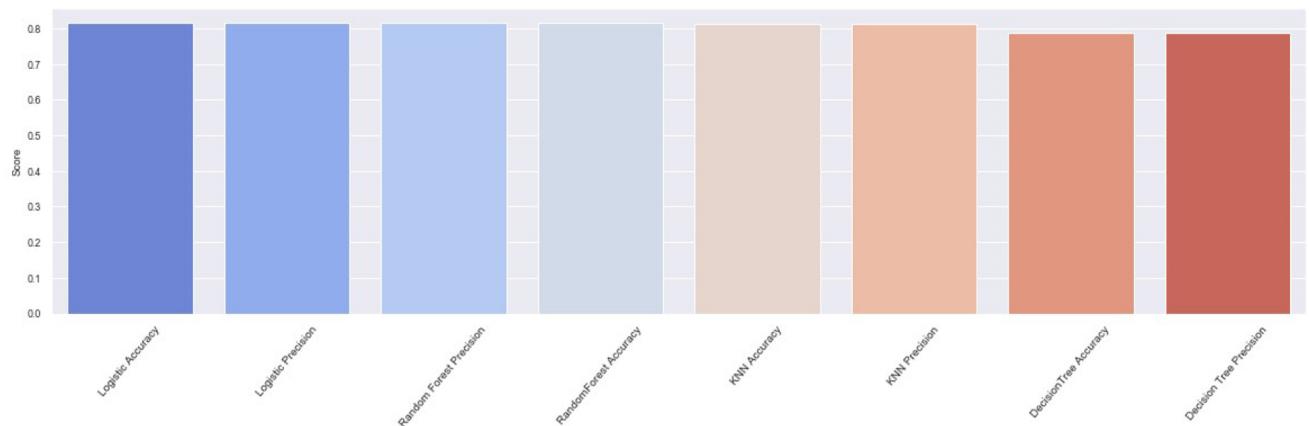
Visualizations

Average Accuracies Across Folds

```
In [46]: fig = plt.figure(figsize = (20,5))

sns.barplot(x = t1_df_accu_dict.index,y = t1_df_accu_dict.Score,
            palette = 'coolwarm').set_xticklabels(t1_df_accu_dict.index,rotation = 50, fontsize = 10)

plt.show()
```

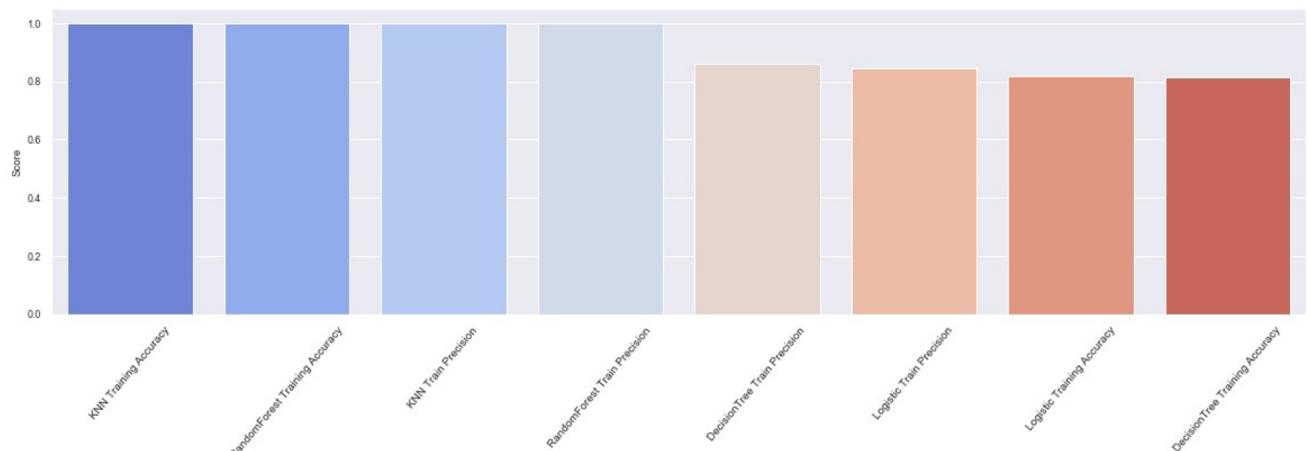


Algorithm Training Scores

```
In [47]: fig = plt.figure(figsize = (20,5))

sns.barplot(x = t1_df_train_dict.index,y = t1_df_train_dict.Score,
            palette = 'coolwarm').set_xticklabels(t1_df_train_dict.index,rotation = 50, fontsize = 10)

plt.show()
```

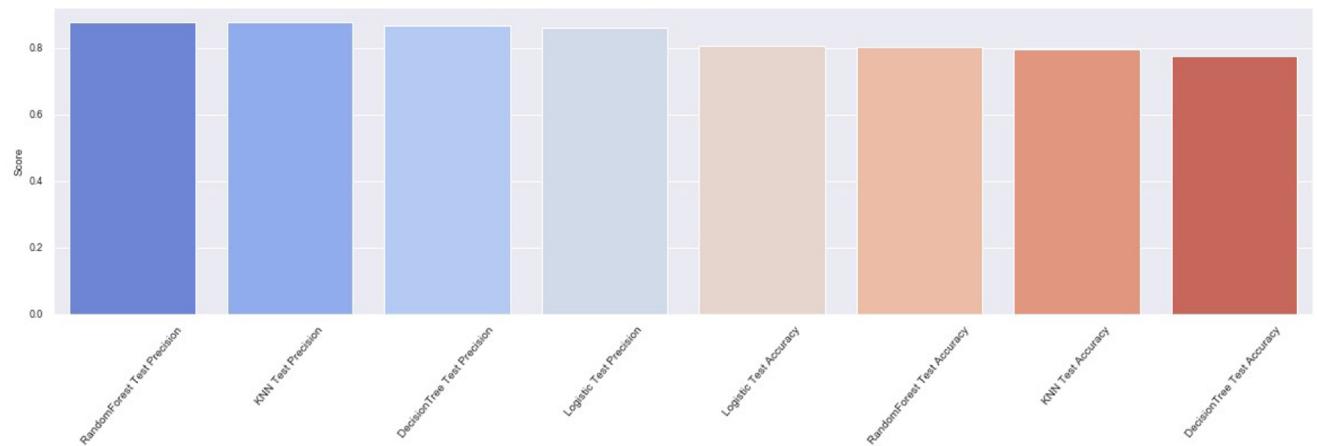


Algorithm Testing Scores

```
In [48]: fig = plt.figure(figsize = (20,5))

sns.barplot(x = t1_df_test_dict.index,y = t1_df_test_dict.Score,
            palette = 'coolwarm').set_xticklabels(t1_df_test_dict.index,rotation = 50, fontsize = 10)

plt.show()
```

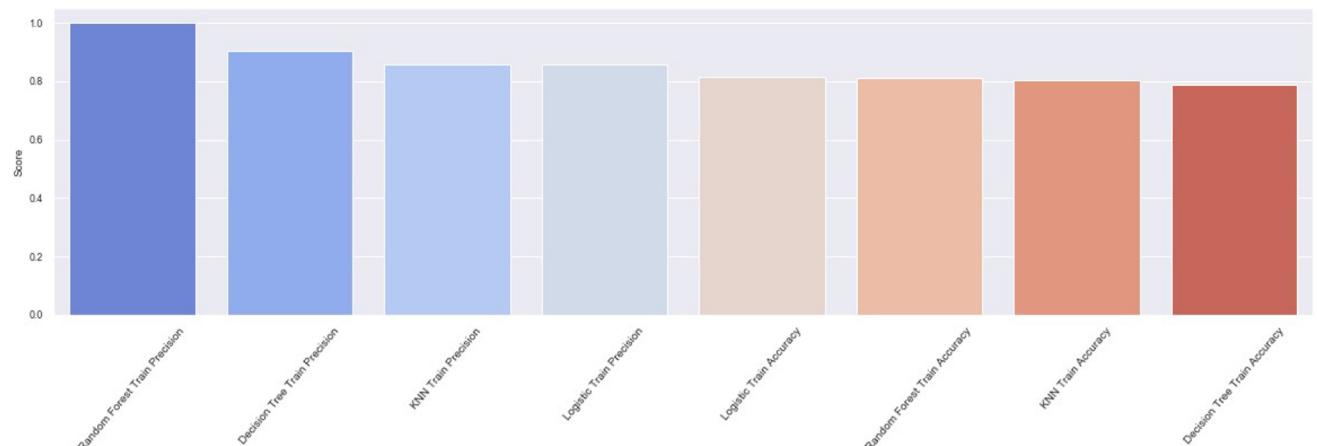


Algorithm Training Scores With Optimal Hyperparameter

```
In [49]: fig = plt.figure(figsize = (20,5))

sns.barplot(x = t1_df_optimal_train_dict.index,y = t1_df_optimal_train_dict.Score,
            palette = 'coolwarm').set_xticklabels(t1_df_optimal_train_dict.index,rotation = 50, fontsize = 10)

plt.show()
```



Algorithm Testing Scores With Optimal Hyperparameter

```
In [50]: fig = plt.figure(figsize = (20,5))
```

```
    sns.barplot(x = t1_df_optimal_test_dict.index,y = t1_df_optimal_test_dict.Score,  
                palette = 'coolwarm').set_xticklabels(t1_df_optimal_test_dict.index,rotation = 50, fontsize = 10)
```

```
    plt.show()
```

