

Model Selection by Evaluation Metrics and Learning Curves

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

%matplotlib inline
```

Import cleaned emails from the original emails dataset

Since fold validation need more dataset for a better evaluation so I decide to use the whole dataset which has 5728 instances

```
In [2]: nlp_emails = pd.read_csv('cleaned_emails.csv')
```

```
In [3]: nlp_emails.shape
```

```
Out[3]: (5728, 4)
```

```
In [4]: nlp_emails[0:3]
```

```
Out[4]:
```

	Unnamed: 0		text	spam	nlp_X
0	0	Subject: naturally irresistible your corporate...		1	subject naturally irresistible corporate ident...
1	1	Subject: the stock trading gunslinger fanny i...		1	subject stock trading gunslinger fanny merrill...
2	2	Subject: unbelievable new homes made easy im ...		1	subject unbelievable new homes easy im wanting...

```
In [5]: nlp_X = np.array(nlp_emails['nlp_X'])
```

```
In [6]: y = np.array(nlp_emails['spam'])
```

Text Vectorization

```
In [7]: from sklearn.feature_extraction.text import TfidfVectorizer

tv = TfidfVectorizer(max_features= 2500)
tv_nlp_X = tv.fit_transform(nlp_X)
tv_nlp_X = tv_nlp_X.toarray()
```

Splitting the dataset into the Training set and Test set

```
In [8]: from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(tv_nlp_X, y, test_
```

Will focus on accuracy and precision

Unlike something like cancer detection, recall is the main matrix features need to be pay more attention; in spam detection machine learning, we need to pay more attention on precision since we don't want the machine mistakenly predict an important email as a spam ending up the receiver could not receive the email on time.

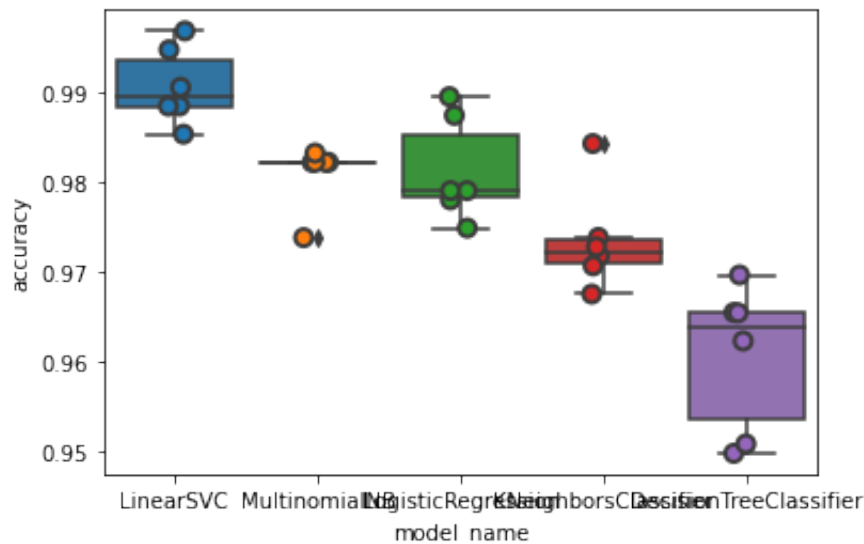
I generated several models to compare their accuracy and precision rate. RandomForest and LDA are among the lowest so I remove them for saving the time.

```
In [9]: import warnings
warnings.filterwarnings("ignore")
from sklearn.linear_model import LogisticRegression
#from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import LinearSVC
from sklearn.naive_bayes import MultinomialNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
#from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.model_selection import cross_val_score
models = [
    #RandomForestClassifier(n_estimators=50, max_depth=3, random_state
    LinearSVC(),
    MultinomialNB(),
    LogisticRegression(random_state=0),
    KNeighborsClassifier(),
    #LinearDiscriminantAnalysis(),
    DecisionTreeClassifier().
```

```

]
CV = 6
cv_df = pd.DataFrame(index=range(CV * len(models)))
entries = []
for model in models:
    model_name = model.__class__.__name__
    accuracies = cross_val_score(model, tv_nlp_X, y, scoring='accuracy',
    for fold_idx, accuracy in enumerate(accuracies):
        entries.append((model_name, fold_idx, accuracy))
cv_df = pd.DataFrame(entries, columns=['model_name', 'fold_idx', 'accuracy'])
import seaborn as sns
sns.boxplot(x='model_name', y='accuracy', data=cv_df)
sns.stripplot(x='model_name', y='accuracy', data=cv_df,
              size=8, jitter=True, edgecolor="gray", linewidth=2)
plt.show()

```



```
In [10]: cv_df.groupby('model_name').accuracy.mean()
```

```

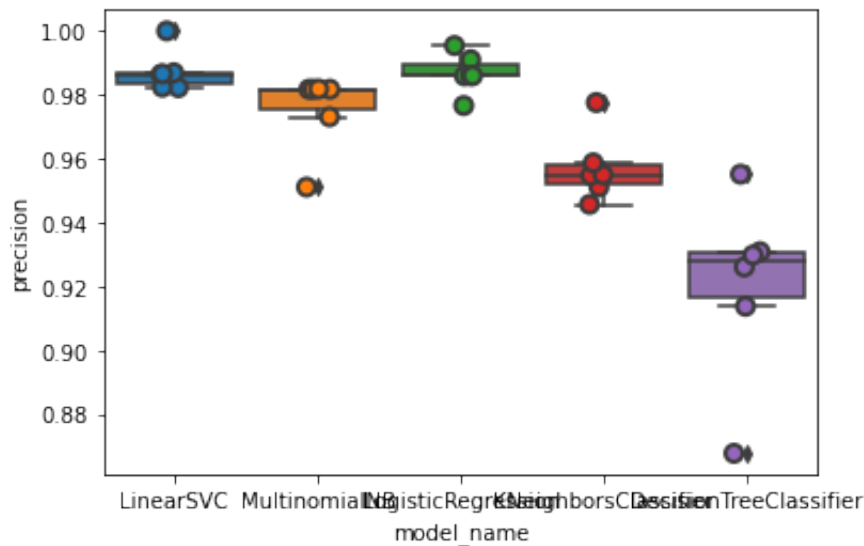
Out[10]: model_name
DecisionTreeClassifier    0.960546
KNeighborsClassifier      0.973462
LinearSVC                 0.990747
LogisticRegression       0.981319
MultinomialNB            0.980970
Name: accuracy, dtype: float64

```

```

In [11]: models = [
    #RandomForestClassifier(n_estimators=50, max_depth=3, random_state=0),
    LinearSVC(),
    MultinomialNB(),
    LogisticRegression(random_state=0),
    KNeighborsClassifier(),
    #LinearDiscriminantAnalysis(),
    DecisionTreeClassifier(),
]
CV = 6
cv_df = pd.DataFrame(index=range(CV * len(models)))
entries = []
for model in models:
    model_name = model.__class__.__name__
    recall = cross_val_score(model, tv_nlp_X, y, scoring='precision', cv=CV)
    for fold_idx, recall in enumerate(recall):
        entries.append((model_name, fold_idx, recall))
cv_df = pd.DataFrame(entries, columns=['model_name', 'fold_idx', 'precision'])
import seaborn as sns
sns.boxplot(x='model_name', y='precision', data=cv_df)
sns.stripplot(x='model_name', y='precision', data=cv_df,
              size=8, jitter=True, edgecolor="gray", linewidth=2)
plt.show()

```



Which metrics should you focus on?

Choice of metric depends on your business objective

- Spam filter* (positive class is "spam"): Optimize for precision or specificity because false negatives (spam goes to the inbox) are more acceptable than false positives (non-spam is caught by the spam filter)
- Fraudulent transaction* detector (positive class is "fraud"): Optimize for sensitivity because false positives (normal transactions that are flagged as possible fraud) are more acceptable than false negatives (fraudulent transactions that are not detected)

Precision

false positives (FP): We predicted yes, but they don't actually have the disease. (Also known as a "Type I error."). In this project we should reduce the Type I errors.

precision, also known as the true positive rate (TPR), is the proportion of the total amount of relevant instances that were actually retrieved. It answers the question "What proportion of actual positives was identified correctly?"

```
In [12]: cv_df.groupby('model_name').precision.mean()
```

```
Out[12]: model_name
DecisionTreeClassifier    0.920584
KNeighborsClassifier      0.957135
LinearSVC                 0.987402
LogisticRegression       0.986820
MultinomialNB            0.975205
Name: precision, dtype: float64
```

Model Evaluation & Selection

Train/test/split sampling

Select the right size of test dataset to avoid Overfitting or Variable Bias(underfitting)

```
In [13]: from sklearn.model_selection import train_test_split, StratifiedShuffleSplit
from sklearn.metrics import accuracy_score
from sklearn.metrics import recall_score
import seaborn as sns
from scipy import stats
from sklearn.model_selection import StratifiedKFold
```

Since LinearSVC reaches the highest scores both in Accuracy and Recall. I will use LinearSVC to do the model evaluation by different test splitting sampling.

```
In [14]: from sklearn.svm import LinearSVC
LSVCclf = LinearSVC()
LSVCclf.fit(X_train,y_train)
```

Out[14]: LinearSVC()

```
In [15]: y_pred = LSVCclf.predict(X_test)
```

```
In [16]: from sklearn.metrics import accuracy_score,classification_report,confusion_matrix
print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.99	0.99	0.99	1398
1	0.99	0.96	0.97	493
accuracy			0.99	1891
macro avg	0.99	0.98	0.98	1891
weighted avg	0.99	0.99	0.99	1891

```
In [17]: # Importing the metrics package from sklearn library
from sklearn import metrics
# Creating the confusion matrix
cm = metrics.confusion_matrix(y_test, y_pred)
# Assigning columns names
cm_df = pd.DataFrame(cm,
                      columns = ['Predicted Negative', 'Predicted Positive'],
                      index = ['Actual Negative', 'Actual Positive'])
# Showing the confusion matrix
cm_df
```

Out[17]:

	Predicted Negative	Predicted Positive
Actual Negative	1391	7
Actual Positive	19	474

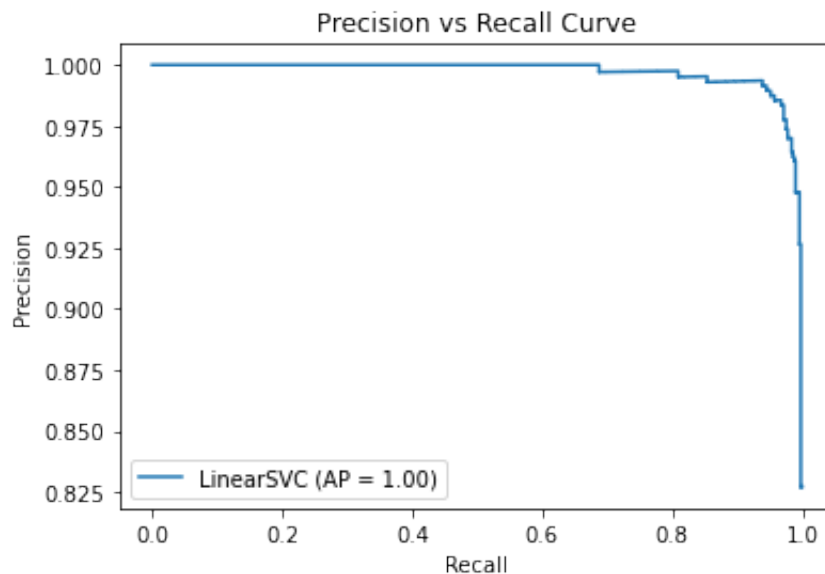
```
In [18]: print('True:', y_test[0:25])
print('Pred:', y_pred[0:25])
```

```
True: [0 0 0 0 0 1 0 0 0 1 1 0 1 0 0 1 0 0 0 1 1 0 0 0 0]
Pred: [0 0 0 0 0 1 0 0 0 1 1 0 0 0 0 1 0 0 0 1 1 0 0 0 0]
```

```
In [19]: from sklearn.metrics import precision_recall_curve
from sklearn.metrics import plot_precision_recall_curve
import matplotlib.pyplot as plt
from sklearn.preprocessing import label_binarize

disp = plot_precision_recall_curve(LSVCclf, X_test, y_test)
disp.ax_.set_title('Precision vs Recall Curve')
```

Out[19]: Text(0.5, 1.0, 'Precision vs Recall Curve')



Sensitivity: When the actual value is positive, how often is the prediction correct? How "sensitive" is the classifier to detecting positive instances? Also known as "True Positive Rate" or "Recall"

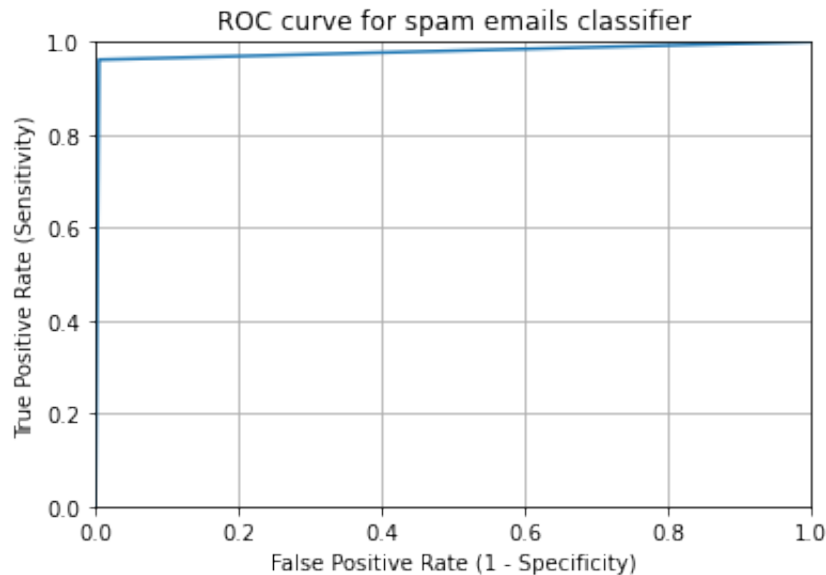
Specificity: When the actual value is negative, how often is the prediction correct? How "specific" (or "selective") is the classifier in predicting positive instances?

ROC Curves and Area Under the Curve (AUC)

Question: Wouldn't it be nice if we could see how sensitivity and specificity are affected by various thresholds, without actually changing the threshold?

Answer: Plot the ROC curve


```
In [20]: fpr, tpr, thresholds = metrics.roc_curve(y_test, y_pred)
plt.plot(fpr, tpr)
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.0])
plt.title('ROC curve for spam emails classifier')
plt.xlabel('False Positive Rate (1 - Specificity)')
plt.ylabel('True Positive Rate (Sensitivity)')
plt.grid(True)
```



```
In [21]: # define a function that accepts a threshold and prints sensitivity and
def evaluate_threshold(threshold):
    print('Sensitivity:', tpr[thresholds > threshold][-1])
    print('Specificity:', 1 - fpr[thresholds > threshold][-1])
evaluate_threshold(0.5)
```

```
Sensitivity: 0.9614604462474645
Specificity: 0.9949928469241774
```

Calculate the AUC

AUC is useful as a single number summary of classifier performance. If you randomly chose one positive and one negative observation, AUC represents the likelihood that your classifier will assign a higher predicted probability to the positive observation. AUC is useful even when there is high class imbalance (unlike classification accuracy).

```
In [22]: print(metrics.roc_auc_score(y_test, y_pred))

0.978226646585821
```

Stratified sampling

Split train and test set will result a better testing accuracy, but it provides a high variance estimate since changing which observation happen to be in the testing set can signifivantly change testing accuracy

Stratified sampling While splitting, we need to ensure that the distribution of features as well as target remains the same in the training and test sets. For ex: Consider a problem where we're trying to classify an observation as fraudulent or not. While splitting, if the majority of fraud cases went to the test set, the model won't be able to learn the fraudulent patterns, as it doesn't have access to many fraud cases in the training data. In such cases, stratified sampling should be done, as it maintains the proportion of different classes in the train and test set. In this project: spam emails are the minority cases so that using Stratified sampling will be better

Avoid High Bias or High Variance

A good choice of hyperparameters ensures that parameters learnt are corresponding to a good loss minima (more generalizable model). Generalizable models are less prone to overfitting. They have consistently good performance on train and test data.

High bias implies our estimate based on the observed data is not close to the true parameter. (aka underfitting). **High variance** implies our estimates are sensitive to sampling. They'll vary a lot if we compute them with a different sample of data (aka overfitting).

Validation strategies can be broadly divided into 2 categories: **Holdout validation** and **cross validation**.

a)Single holdout: Varying test Size by testing the recall scores

Implementation

The basic idea is to split our data into a training set and a holdout test set. Train the model on the training set and then evaluate model performance on the test set. We take only a single holdout—hence the name

step1: split target data into 2 subsets

step2: Choose LinearSVC as the learning algorithm

step3: Predict on the test data using the trained model. Choose an appropriate metric for performance estimation(I choose recall to the classification task). Assess predictive performance by comparing predictions and ground truth.

Step 4: If the performance estimate computed in the previous step is satisfactory, combine the train and test subset to train the model on the full data with the same hyperparameters.

Stratified sampling

While splitting, we need to ensure that the distribution of features as well as target remains the same in the training and test sets. For ex: Consider a problem where we're trying to classify an observation as fraudulent or not. While splitting, if the majority of fraud cases went to the test set, the model won't be able to learn the fraudulent patterns, as it doesn't have access to many fraud cases in the training data. In such cases, stratified sampling should be done, as it maintains the proportion of different classes in the train and test set.

In this project: spam emails are the minority cases so that using Stratified sampling will be better

```
In [23]: x= tv_nlp_X  
Y=y
```

```

In [31]: from sklearn.metrics import precision_score
# varying hold out size
test_size = np.arange(0.05, 0.55, 0.05)

trn_precision = []
tst_precision = []

for sz in test_size:
    #stratified sampling
    sss = StratifiedShuffleSplit(n_splits=1, test_size=sz, random_state=0)

    #train-test split
    for trn_idx, tst_idx in sss.split(x, y):
        x_trn, y_trn, x_tst, y_tst = x[trn_idx], y[trn_idx], x[tst_idx], y[tst_idx]

        #model fitting
        clf = LinearSVC()
        clf.fit(x_trn, y_trn)

        #model prediction
        pred_tst = clf.predict(x_tst)
        pred_trn = clf.predict(x_trn)

        #performance evaluation
        tst_precision.append(precision_score(y_tst, pred_tst))
        trn_precision.append(precision_score(y_trn, pred_trn))

```

```

In [35]: # 95% CI calculation using normal approximation method
ui = []
li = []
for i, n in enumerate(test_size):
    p = tst_precision[i]
    sigma = np.sqrt(p * (1 - p) / (n * 10000))
    ui.append(p + 1.96 * sigma)
    li.append(p - 1.96 * sigma)

```

```

In [36]: #CI plot
def lineplotCI(x_data, y_data, low_CI, upper_CI, x_label, y_label, title):
    # Create the plot object
    _, ax = plt.subplots(figsize=(20, 10))
    # Plot the data, set the linewidth, color and transparency of the
    # line, provide a label for the legend
    ax.plot(
        x_data,
        y_data,
        lw=1,
        color='#FF335B',

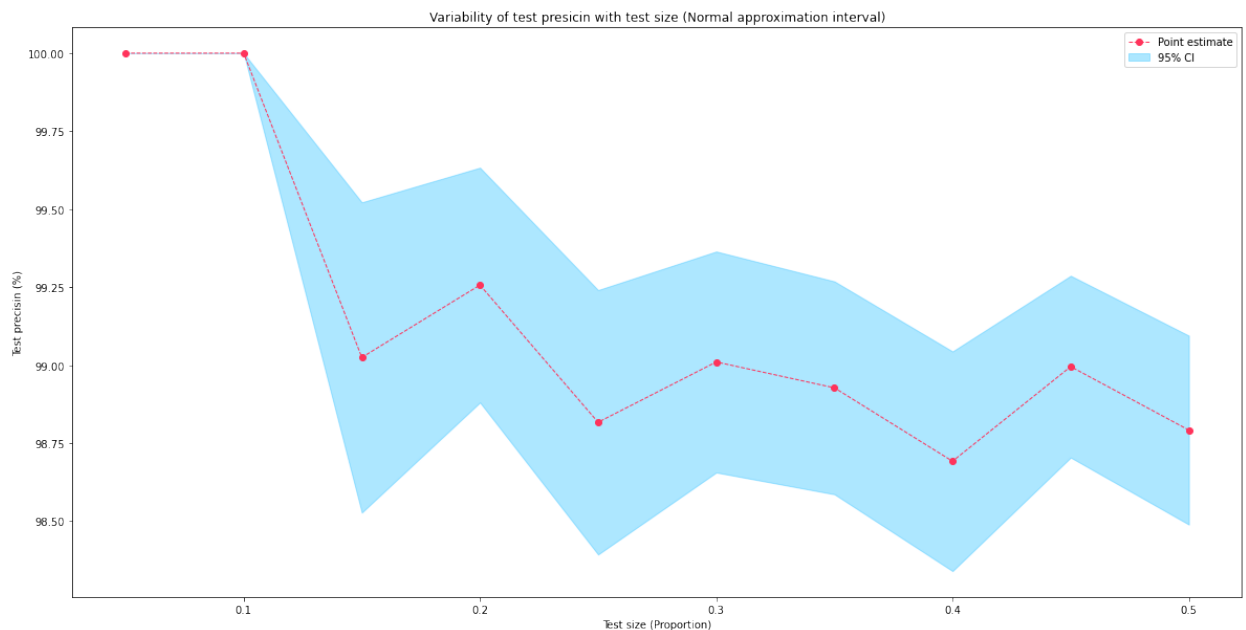
```

```

        alpha=1,
        label='Point estimate',
        linestyle='--',
        marker='o')
# Shade the confidence interval
ax.fill_between(
    x_data, low_CI, upper_CI, color='#33C4FF', alpha=0.4, label='95% CI')
# Label the axes and provide a title
ax.set_title(title)
ax.set_xlabel(x_label)
ax.set_ylabel(y_label)
# Display legend
ax.legend(loc='best')

# Call the function to create plot
lineplotCI(
    x_data=test_size,
    y_data=100 * np.array(tst_precision),
    low_CI=100 * np.array(li),
    upper_CI=100 * np.array(ui),
    x_label='Test size (Proportion)',
    y_label='Test precision (%)',
    title=
        'Variability of test precision with test size (Normal approximation interval)')

```



Choice of test size

Keeping aside a large amount of data for the test can result in an underestimation of predictive power (high bias). **But the estimate will be more stable (low variance)**, as shown in the figure above. This consideration is more relevant for smaller datasets. When test size are in the range of 10%, it reaches the highest precision rate.

Repeated HoldOut for the precision scores

```
In [59]: test_size = 0.15
# repeating for 50 different seeds
seed = np.random.randint(0, 1000, 50)
trn_precision_1 = []
tst_precision_1 = []
for state in seed:
    sss = StratifiedShuffleSplit(
        n_splits=1, test_size=test_size, random_state=state)

    for trn_idx, tst_idx in sss.split(x, y):
        x_trn, y_trn, x_tst, y_tst = x[trn_idx], y[trn_idx], x[tst_idx], y[tst_idx]

    clf = LinearSVC()

    clf.fit(x_trn, y_trn)

    pred_tst = clf.predict(x_tst)
    pred_trn = clf.predict(x_trn)

    tst_precision_1.append(precision_score(y_tst, pred_tst))
    trn_precision_1.append(precision_score(y_trn, pred_trn))
```

```
In [45]: test_size = 0.5
seed = np.random.randint(0, 1000, 50)
trn_precision_2 = []
tst_precision_2 = []
for state in seed:
    sss = StratifiedShuffleSplit(
        n_splits=1, test_size=test_size, random_state=state)

    for trn_idx, tst_idx in sss.split(x, y):
        x_trn, y_trn, x_tst, y_tst = x[trn_idx], y[trn_idx], x[tst_idx], y[tst_idx]

        clf = LinearSVC()

        clf.fit(x_trn, y_trn)

        pred_tst = clf.predict(x_tst)
        pred_trn = clf.predict(x_trn)

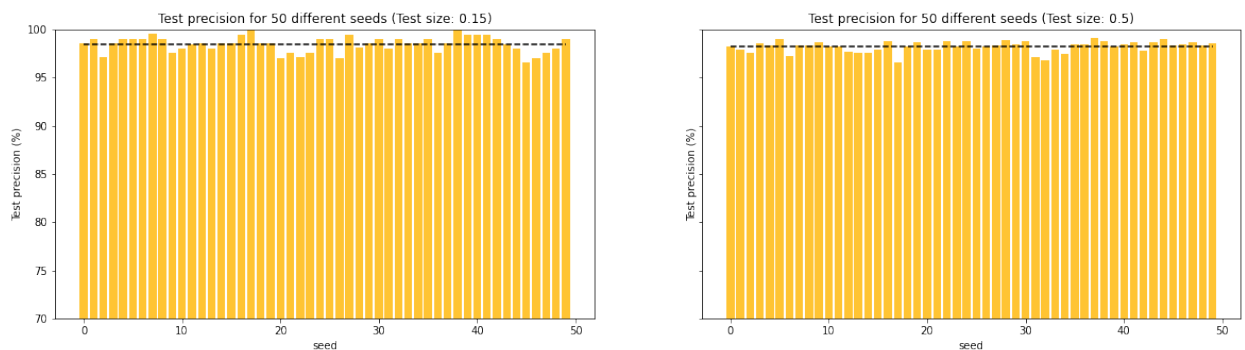
        tst_precision_2.append(precision_score(y_tst, pred_tst))
        trn_precision_2.append(precision_score(y_trn, pred_trn))
```

```
In [61]: df = pd.DataFrame()

df['0.15'] = tst_precision_1
df['0.5'] = tst_precision_2
df['seed'] = df.index
```

```
In [62]: # plotting accuracy score for 50 different iterations
_, axes = plt.subplots(1, 2, figsize=(20, 5), sharey=True)
axes[0].bar(df.seed, 100 * df['0.15'], color='#FFC433')
axes[0].plot(df.seed, [100 * np.mean(df['0.15'])] * df.shape[0], "k--")
axes[0].set_title("Test precision for 50 different seeds (Test size: 0.15)")
axes[1].bar(df.seed, 100 * df['0.5'], color='#FFC433')
axes[1].plot(df.seed, [100 * np.mean(df['0.5'])] * df.shape[0], "k--")
axes[1].set_title("Test precision for 50 different seeds (Test size: 0.5)")
plt.ylim([70, 100])
axes[0].set_xlabel("seed")
axes[0].set_ylabel("Test precision (%)")
axes[1].set_xlabel("seed")
axes[1].set_ylabel("Test precision (%)")
```

Out[62]: Text(0, 0.5, 'Test precision (%)')



```
In [64]: print(f"Mean precision for test size (0.15): {100*np.mean(df['0.15'])}")
```

Mean precision for test size (0.15): 98.48023669821966

```
In [65]: print(f"Mean precision for test size (0.5): {100*np.mean(df['0.5'])}")
```

Mean precision for test size (0.5): 98.23198601765745

How confident are we in our estimates? From the above steps, we'll get a point estimate of the true predictive power of our model. But this single number doesn't mean anything unless we know how confident we are in this estimate. Defining the confidence interval around this point estimate would tell us how much this estimate can vary for a different set of model inputs. Let's discuss a way of estimating this interval.

Normal approximation interval

Suppose we're choosing accuracy as the proxy for predictive power of the model. Let's look at the calculation for the confidence interval (CI) in this case: Normal approximation interval

$$SE_{repeated} = \sqrt{\left(\sum_{i=1}^k (ACC_i - ACC_{avg})^2\right) / (k - 1)}$$

$$CI = ACC_{avg} \pm t * SE_{repeated}$$

In the above formula, SE is the standard error and t is the value coming from t-distribution with degree of freedom as k-1. We're using t-distribution because we're calculating SE from the sample.

```
In [66]: stats.t.ppf([.025, .975], 49) #t-value for 95% confidence
```

```
Out[66]: array([-2.00957523,  2.00957523])
```

```
In [67]: ui, li = np.mean(
            df['0.15']) + np.sqrt(np.var(df['0.15']) * 50 / 49) * 2.00957523,
            df['0.15']) - np.sqrt(np.var(df['0.15']) * 50 / 49) * 2.00957523

            li, ui
```

```
Out[67]: (0.9683439322421953, 1.0012608017221978)
```

```
In [68]: ui, li = np.mean(
            df['0.5']) + np.sqrt(np.var(df['0.5']) * 50 / 49) * 2.00957523, np
            df['0.5']) - np.sqrt(np.var(df['0.5']) * 50 / 49) * 2.00957523

            li, ui
```

```
Out[68]: (0.971205140908077, 0.993434579445072)
```

Empirical interval

Empirical interval is suggested when our samples don't follow a normal distribution and the value of k is high

```
In [69]: np.percentile(df['0.15'],97.5)
```

```
Out[69]: 0.9989024390243902
```

```
In [70]: np.percentile(df['0.15'],2.5)
```

```
Out[70]: 0.9703948321661434
```

```
In [71]: np.percentile(df['0.5'],97.5)
```

```
Out[71]: 0.989545200643489
```

```
In [72]: np.percentile(df['0.5'],2.5)
```

```
Out[72]: 0.9687432031119291
```

K-fold cross validation

The most common approach for model evaluation is cross validation. Let's quickly go through the steps:

- Choose a value of k and divide the data into k equal subsets
- Combine k-1 subsets and consider it as a training fold and the remaining one as a test fold
- Conduct the holdout method to get test performance (let's choose recall for now)
- Repeat 2nd and 3rd steps, k times with a different subset as test fold
- Point estimate of predictive power is the average of the k different test accuracies

Choice of k

Small k: High bias (less data for training in each fold) but low variance (more data in test)

High k: Low bias but high variance Below is the comparison of variance and bias for 5-fold CV, 10-fold CV and 10-fold CV repeated multiple times from this project. It get a little bit higher recall in the fold of 10, but fold of 5 get less erros of standard deviation.

```
In [73]: from sklearn.model_selection import cross_val_score
skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=0)
folds_score = []
for trn_idx, tst_idx in skf.split(x, y):
    x_trn, y_trn, x_tst, y_tst = x[trn_idx], y[trn_idx], x[tst_idx], y[tst_idx]
    clf = LinearSVC()

    clf.fit(x_trn, y_trn)

    pred_tst = clf.predict(x_tst)

    folds_score.append(precision_score(y_tst, pred_tst))

mean_precision = np.mean(folds_score)

std_precision = np.std(folds_score)

mean_precision, std_precision
```

Out[73]: (0.9845816688101117, 0.00802927388349656)

```
In [74]: skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=0)
folds_score = []
for trn_idx, tst_idx in skf.split(x, y):
    x_trn, y_trn, x_tst, y_tst = x[trn_idx], y[trn_idx], x[tst_idx], y[tst_idx]
    clf = LinearSVC()

    clf.fit(x_trn, y_trn)

    pred_tst = clf.predict(x_tst)

    folds_score.append(precision_score(y_tst, pred_tst))

mean_precision = np.mean(folds_score)

std_precision = np.std(folds_score)

mean_precision, std_precision
```

Out[74]: (0.984661987785689, 0.011763937960986202)

K_fold Cross Validation for multiple algorithms

```
In [75]: from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
```

```

from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.naive_bayes import GaussianNB
from sklearn.svm import SVC
from sklearn.svm import LinearSVC

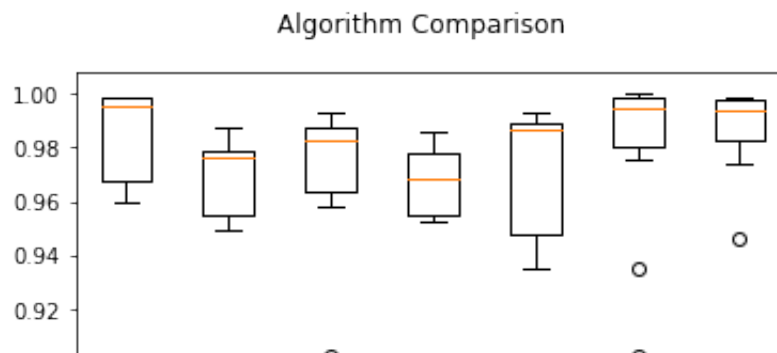
seed = 7
# prepare models
models = []
models.append(('LR', LogisticRegression()))
models.append(('LDA', LinearDiscriminantAnalysis()))
models.append(('KNN', KNeighborsClassifier()))
models.append(('CART', DecisionTreeClassifier()))
models.append(('NB', GaussianNB()))
models.append(('SVM', SVC()))
models.append(('LSVC', LinearSVC()))
# evaluate each model in turn
results = []
names = []
scoring = 'accuracy'
for name, model in models:
    kfold = model_selection.KFold(n_splits=10, random_state=seed)
    cv_results = model_selection.cross_val_score(model, tv_nlp_X, y, cv=kfold)
    results.append(cv_results)
    names.append(name)
    msg = "%s: %f (%f)" % (name, cv_results.mean(), cv_results.std())
    print(msg)
# boxplot algorithm comparison
fig = plt.figure()
fig.suptitle('Algorithm Comparison')
ax = fig.add_subplot(111)
plt.boxplot(results)
ax.set_xticklabels(names)
plt.show()

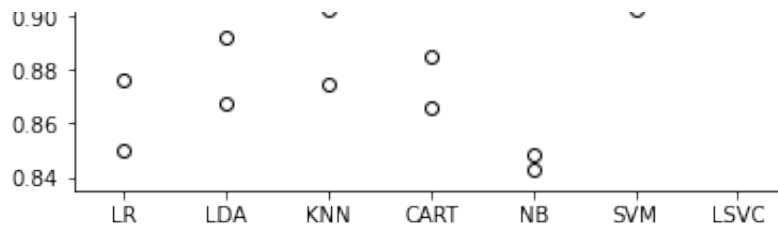
```

```

NB: 0.954447 (0.030702)
SVM: 0.979231 (0.031774)
LSVC: 0.986909 (0.015730)

```





Evaluate the learning curves

Learning curves constitute a great tool to diagnose bias and variance in any supervised learning algorithm. For classification training learning score is base on the accuracy, the higher the better.

In [76]:

```
from sklearn.model_selection import learning_curve
from sklearn.model_selection import ShuffleSplit

def plot_learning_curve(estimator, title, x, y, axes=None, ylim=None,
                        n_jobs=None, train_sizes=np.linspace(.1, 1.0,

if axes is None:
    _, axes = plt.subplots(1, 3, figsize=(20, 5))

axes[0].set_title(title)
if ylim is not None:
    axes[0].set_ylim(*ylim)
axes[0].set_xlabel("Training examples")
axes[0].set_ylabel("Score")

train_sizes, train_scores, test_scores, fit_times, _ = \
    learning_curve(estimator, x, Y, cv=cv, n_jobs=n_jobs,
                  train_sizes=train_sizes,
                  return_times=True)
train_scores_mean = np.mean(train_scores, axis=1)
train_scores_std = np.std(train_scores, axis=1)
test_scores_mean = np.mean(test_scores, axis=1)
test_scores_std = np.std(test_scores, axis=1)
fit_times_mean = np.mean(fit_times, axis=1)
fit_times_std = np.std(fit_times, axis=1)

# Plot learning curve
axes[0].grid()
axes[0].fill_between(train_sizes, train_scores_mean - train_scores_std,
                    train_scores_mean + train_scores_std, alpha=0.1,
                    color="r")
axes[0].fill_between(train_sizes, test_scores_mean - test_scores_std,
                    test_scores_mean + test_scores_std, alpha=0.1,
```

```

        color="g")
axes[0].plot(train_sizes, train_scores_mean, 'o-', color="r",
             label="Training score")
axes[0].plot(train_sizes, test_scores_mean, 'o-', color="g",
             label="Cross-validation score")
axes[0].legend(loc="best")

# Plot n_samples vs fit_times
axes[1].grid()
axes[1].plot(train_sizes, fit_times_mean, 'o-')
axes[1].fill_between(train_sizes, fit_times_mean - fit_times_std,
                    fit_times_mean + fit_times_std, alpha=0.1)
axes[1].set_xlabel("Training examples")
axes[1].set_ylabel("fit_times")
axes[1].set_title("Scalability of the model")

# Plot fit_time vs score
axes[2].grid()
axes[2].plot(fit_times_mean, test_scores_mean, 'o-')
axes[2].fill_between(fit_times_mean, test_scores_mean - test_scores_std,
                    test_scores_mean + test_scores_std, alpha=0.1)
axes[2].set_xlabel("fit_times")
axes[2].set_ylabel("Score")
axes[2].set_title("Performance of the model")

return plt
fig, axes = plt.subplots(3, 2, figsize=(10, 15))

title = "Learning Curves (Naive Bayes)"
# Cross validation with 100 iterations to get smoother mean test and train
# score curves, each time with 20% data randomly selected as a validation set
cv = ShuffleSplit(n_splits=100, test_size=0.2, random_state=0)

estimator = GaussianNB()
plot_learning_curve(estimator, title, x, Y, axes=axes[:, 0], ylim=(0.7, 1.0),
                    cv=cv, n_jobs=4)

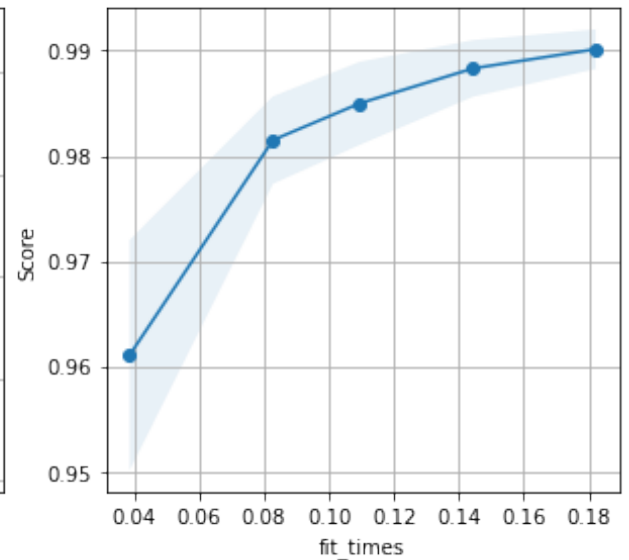
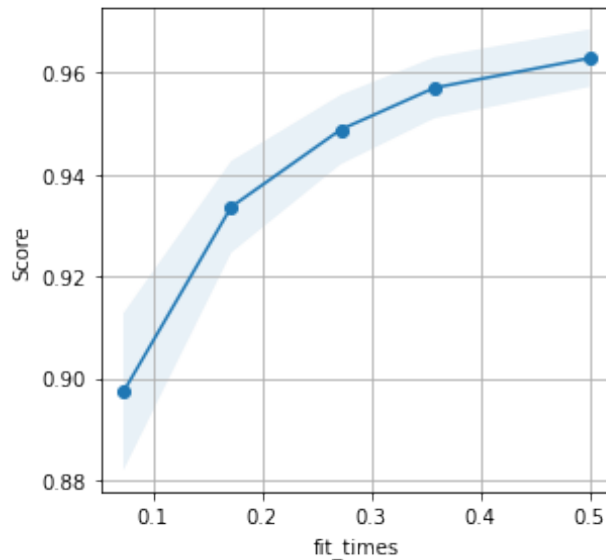
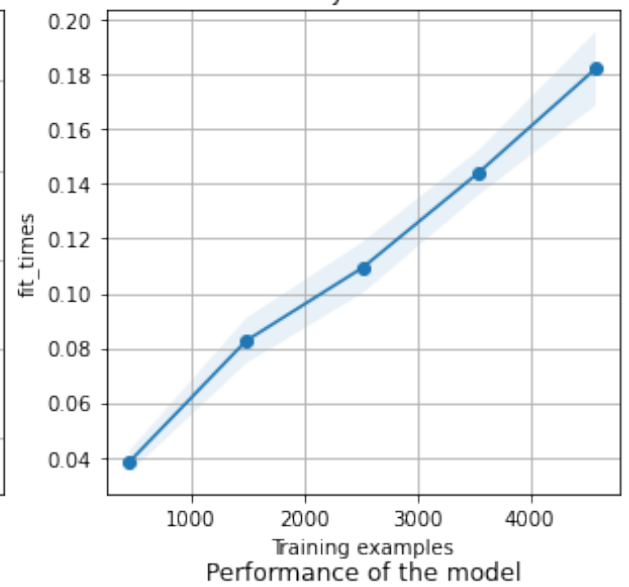
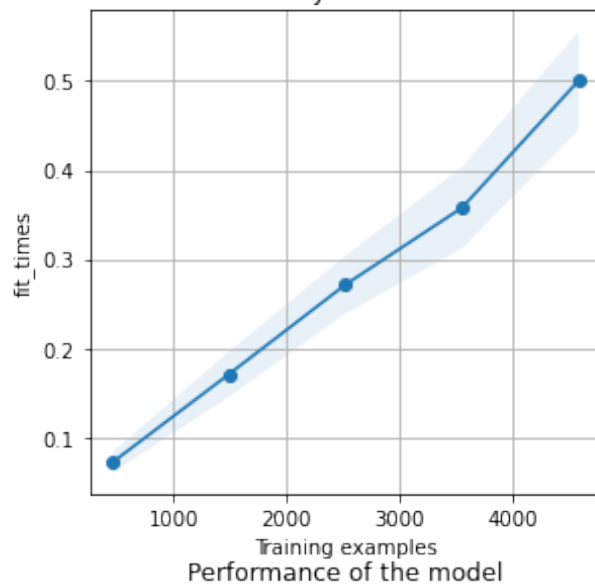
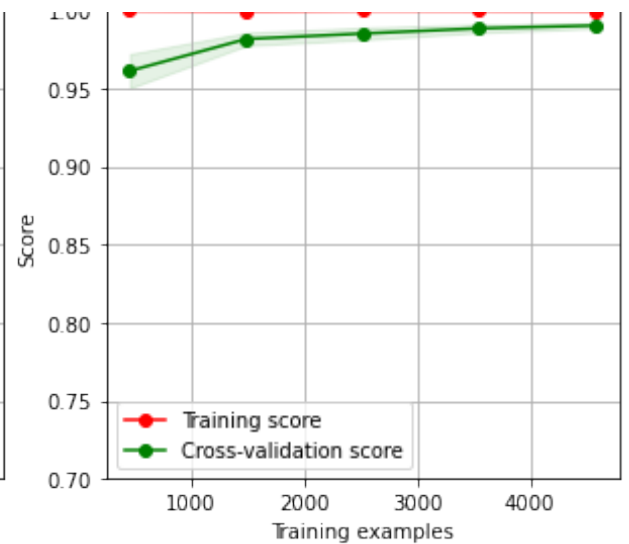
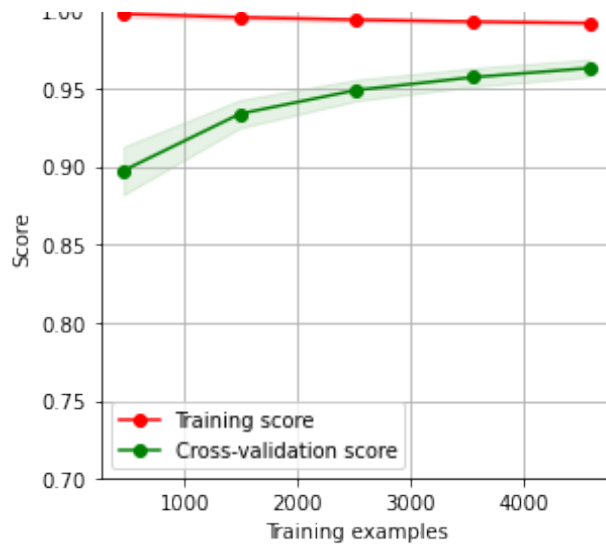
title = r"Learning Curves (LinearSVC)"

cv = ShuffleSplit(n_splits=10, test_size=0.2, random_state=0)
estimator = LinearSVC()
plot_learning_curve(estimator, title, x, Y, axes=axes[:, 1], ylim=(0.7, 1.0),
                    cv=cv, n_jobs=4)

plt.show()

```





Ensemble Methods

Bagging

```
In [77]: import itertools

import seaborn as sns
import matplotlib.gridspec as gridspec

from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.svm import LinearSVC
from sklearn.naive_bayes import MultinomialNB

from sklearn.ensemble import BaggingClassifier
from sklearn.model_selection import cross_val_score, train_test_split

from mlxtend.plotting import plot_learning_curves
from mlxtend.plotting import plot_decision_regions

np.random.seed(0)
```

Use PCA for dimension reduction

```
In [78]: from sklearn.decomposition import TruncatedSVD

svd = TruncatedSVD(n_components=2)
svd.fit(tv_nlp_X)
X_new=svd.fit_transform(tv_nlp_X)
```

```
In [79]: X_new
```

```
Out[79]: array([[ 0.12490899, -0.32372055],
 [ 0.08279241, -0.02436395],
 [ 0.11413207, -0.14975481],
 ...,
 [ 0.26748046,  0.04119078],
 [ 0.27993117,  0.14929532],
 [ 0.1609578 , -0.17484529]])
```



```
In [84]: (criterion='entropy', max_depth=1)
         _neighbors=1)

base_estimator=clf1, n_estimators=10, max_samples=0.8, max_features=0.8
base_estimator=clf2, n_estimators=10, max_samples=0.8, max_features=0.8
```

```
In [85]: label = ['Decision Tree', 'K-NN', 'Bagging Tree', 'Bagging K-NN']

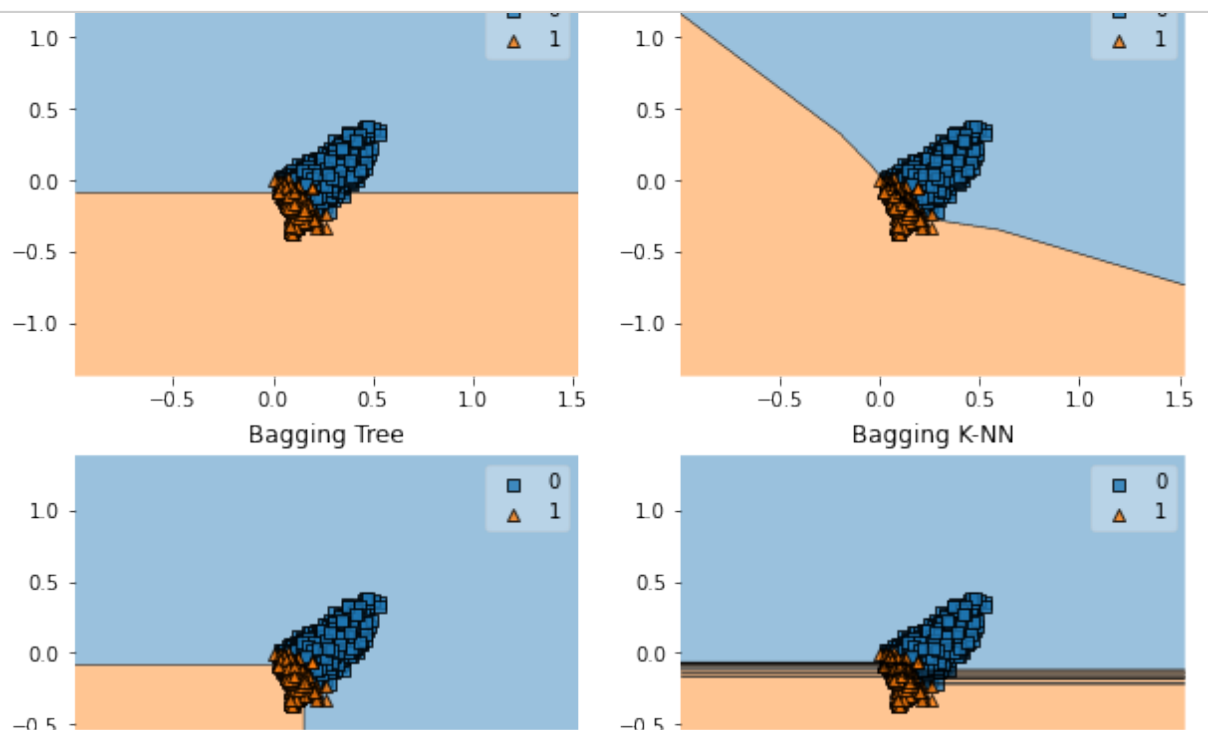
clf_list = [clf1, clf2, bagging1, bagging2]

fig = plt.figure(figsize=(10, 8))
gs = gridspec.GridSpec(2, 2)
grid = itertools.product([0,1],repeat=2)

for clf, label, grd in zip(clf_list, label, grid):
    scores = cross_val_score(clf, X_new, y, cv=3, scoring='accuracy')
    print("Accuracy: %.2f (+/- %.2f) [%s]" % (scores.mean(), scores.std(), label))

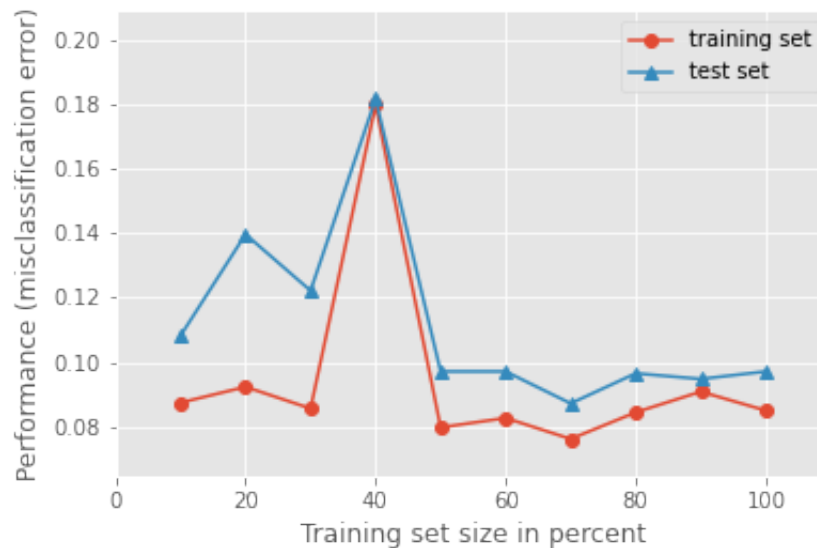
    clf.fit(X_new, y)
    ax = plt.subplot(gs[grd[0], grd[1]])
    fig = plot_decision_regions(X=X_new, y=y, clf=clf)
    plt.title(label)

plt.show()
```



```
In [86]: #plot learning curves
X_train, X_test, y_train, y_test = train_test_split(X_new, y, test_size=0.4)

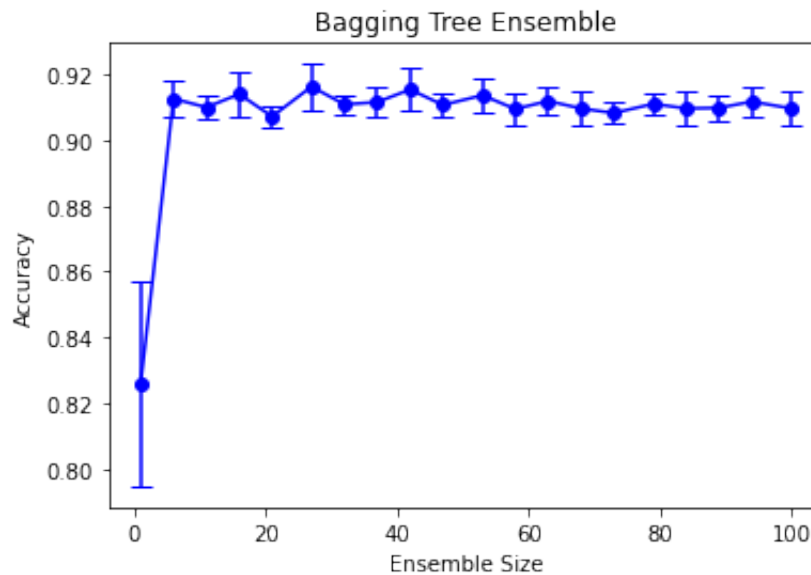
plt.figure()
plot_learning_curves(X_train, y_train, X_test, y_test, bagging1, print_results=True)
plt.show()
```



The figure above shows learning curves for the bagging tree ensemble. We can see an average error of 0.9 on the training data and a triangle-shaped error curve for the testing data. The smallest gap between training and test errors occurs at around 40% of the training set size.

```
In [87]: #Ensemble Size
num_est = np.linspace(1,100,20).astype(int)
bg_clf_cv_mean = []
bg_clf_cv_std = []
for n_est in num_est:
    bg_clf = BaggingClassifier(base_estimator=clf1, n_estimators=n_est)
    scores = cross_val_score(bg_clf, X_new, y, cv=3, scoring='accuracy')
    bg_clf_cv_mean.append(scores.mean())
    bg_clf_cv_std.append(scores.std())
```

```
In [88]: plt.figure()
( _, caps, _ ) = plt.errorbar(num_est, bg_clf_cv_mean, yerr=bg_clf_cv_st
for cap in caps:
    cap.set_maheredgewidth(1)
plt.ylabel('Accuracy'); plt.xlabel('Ensemble Size'); plt.title('Bagging
plt.show()
```



The figure above shows how the test accuracy improves with the size of the ensemble. Based on cross-validation results, we can see the accuracy increases until approximately 10 base estimators and then plateaus afterwards. Thus, adding base estimators beyond 10 only increases computational complexity without accuracy gains for the Iris dataset. A commonly used class of ensemble algorithms are forests of randomized trees. In **random forests**, each tree in the ensemble is built from a sample drawn with replacement (i.e. a bootstrap sample) from the training set. In addition, instead of using all the features, a random subset of features is selected further randomizing the tree. As a result, the bias of the forest increases slightly but due to averaging of less correlated trees, its variance decreases resulting in an overall better model. In **extremely randomized trees** algorithm randomness goes one step further: the splitting thresholds are randomized. Instead of looking for the most discriminative threshold, thresholds are drawn at random for each candidate feature and the best of these randomly-generated thresholds is picked as the splitting rule. This usually allows to reduce the variance of the model a bit more, at the expense of a slightly greater increase in bias.

Boosting

Boosting refers to a family of algorithms that are able to convert weak learners to strong learners. The main principle of boosting is to fit a sequence of weak learners (models that are only slightly better than random guessing, such as small decision trees) to weighted versions of the data, where more weight is given to examples that were mis-classified by earlier rounds. The predictions are then combined through a weighted majority vote (classification) or a weighted sum (regression) to produce the final prediction. The principal difference between boosting and the committee methods such as bagging is that base learners are trained in sequence on a weighted version of the data.

```
In [89]: import itertools
import numpy as np

import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.gridspec as gridspec

from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression

from sklearn.ensemble import AdaBoostClassifier
from sklearn.model_selection import cross_val_score, train_test_split

from mlxtend.plotting import plot_learning_curves
from mlxtend.plotting import plot_decision_regions
```

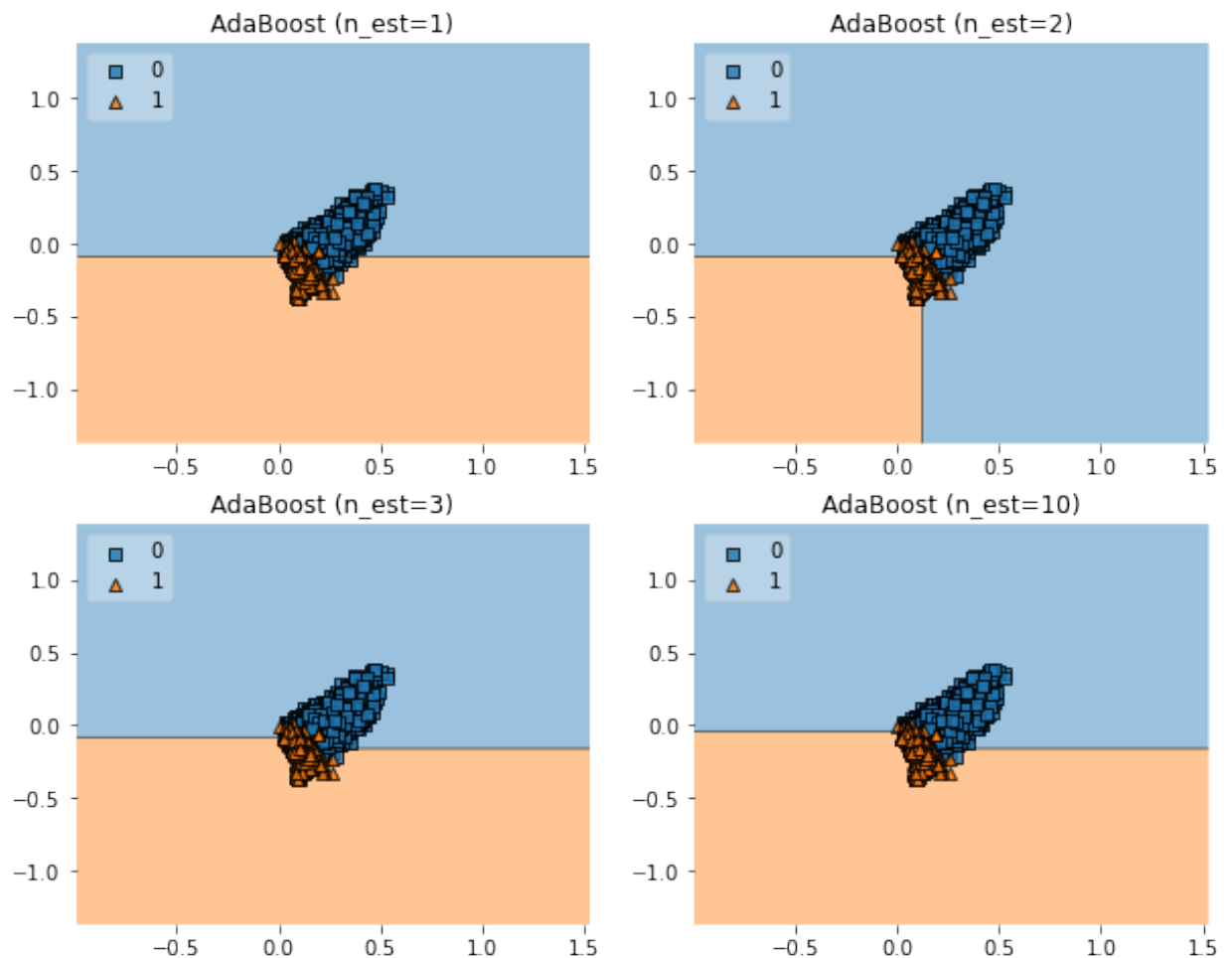
```
In [90]: clf = DecisionTreeClassifier(criterion='entropy', max_depth=1)

num_est = [1, 2, 3, 10]
label = ['AdaBoost (n_est=1)', 'AdaBoost (n_est=2)', 'AdaBoost (n_est=
```

```
In [91]: fig = plt.figure(figsize=(10, 8))
gs = gridspec.GridSpec(2, 2)
grid = itertools.product([0,1],repeat=2)

for n_est, label, grd in zip(num_est, label, grid):
    boosting = AdaBoostClassifier(base_estimator=clf, n_estimators=n_e
    boosting.fit(X_new, y)
    ax = plt.subplot(gs[grd[0], grd[1]])
    fig = plot_decision_regions(X=X_new, y=y, clf=boosting, legend=2)
    plt.title(label)

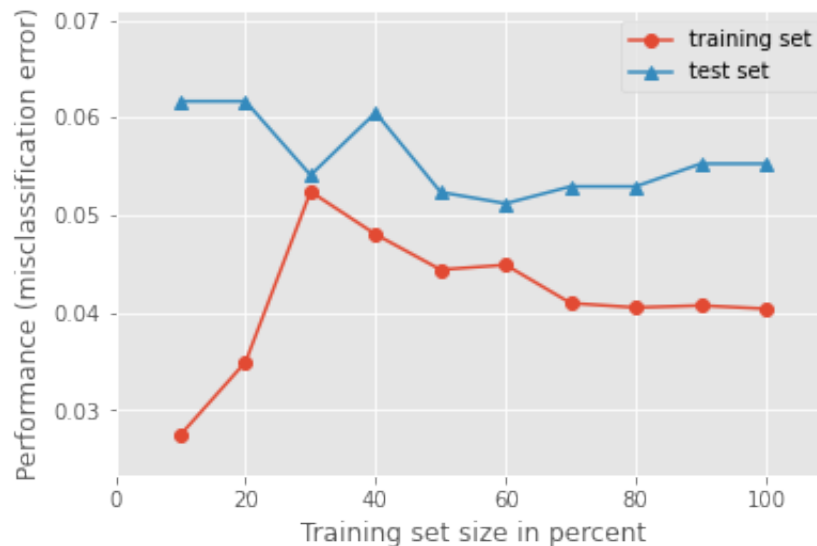
plt.show()
```



```
In [92]: #plot learning curves
X_train, X_test, y_train, y_test = train_test_split(X_new, y, test_size=0.2)

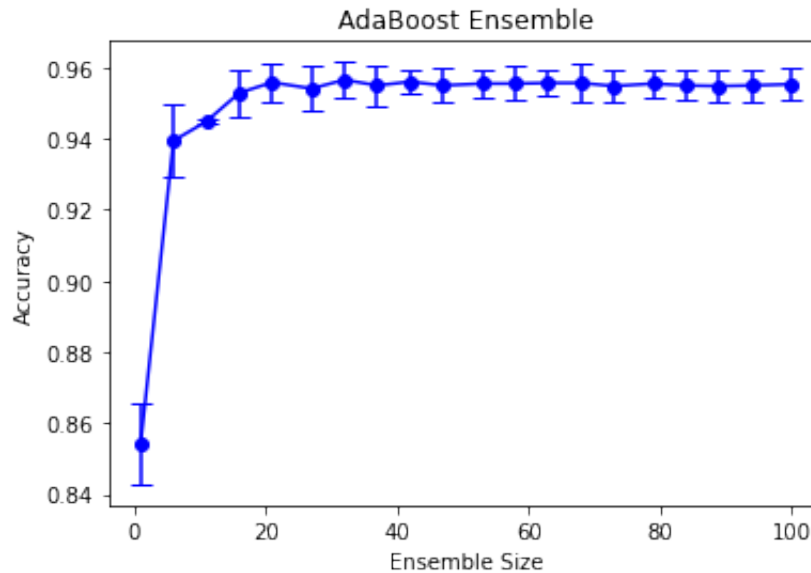
boosting = AdaBoostClassifier(base_estimator=clf, n_estimators=10)

plt.figure()
plot_learning_curves(X_train, y_train, X_test, y_test, boosting, print_results=True)
plt.show()
```



```
In [93]: #Ensemble Size
num_est = np.linspace(1,100,20).astype(int)
bg_clf_cv_mean = []
bg_clf_cv_std = []
for n_est in num_est:
    ada_clf = AdaBoostClassifier(base_estimator=clf, n_estimators=n_est)
    scores = cross_val_score(ada_clf, X_new, y, cv=3, scoring='accuracy')
    bg_clf_cv_mean.append(scores.mean())
    bg_clf_cv_std.append(scores.std())
```

```
In [171]: plt.figure()
          (_, caps, _) = plt.errorbar(num_est, bg_clf_cv_mean, yerr=bg_clf_cv_st
          for cap in caps:
              cap.set_movedgewidth(1)
          plt.ylabel('Accuracy'); plt.xlabel('Ensemble Size'); plt.title('AdaBoc
          plt.show()
```



The figure above shows how the test accuracy improves with the size of the ensemble. (size increased to 20 to 22 reached highest accuracy.

Gradient Tree Boosting is a generalization of boosting to arbitrary differentiable loss functions. It can be used for both regression and classification problems.

Stacking

Stacking is an ensemble learning technique that combines multiple classification or regression models via a meta-classifier or a meta-regressor. The base level models are trained based on complete training set then the meta-model is trained on the outputs of base level model as features. The base level often consists of different learning algorithms and therefore stacking ensembles are often heterogeneous.

```
In [94]: from mlxtend.classifier import StackingClassifier
          from sklearn.ensemble import RandomForestClassifier
```

```
In [95]: clf1 = LinearSVC()
         clf2 = RandomForestClassifier(random_state=1)
         clf3 = GaussianNB()
         lr = LogisticRegression()
         sclf = StackingClassifier(classifiers=[clf1, clf2, clf3],
                                   meta_classifier=lr)
```

```
In [96]: label = ['LSVC', 'Random Forest', 'Naive Bayes', 'Stacking Classifier']
         clf_list = [clf1, clf2, clf3, sclf]

         fig = plt.figure(figsize=(10,8))
         gs = gridspec.GridSpec(2, 2)
         grid = itertools.product([0,1],repeat=2)

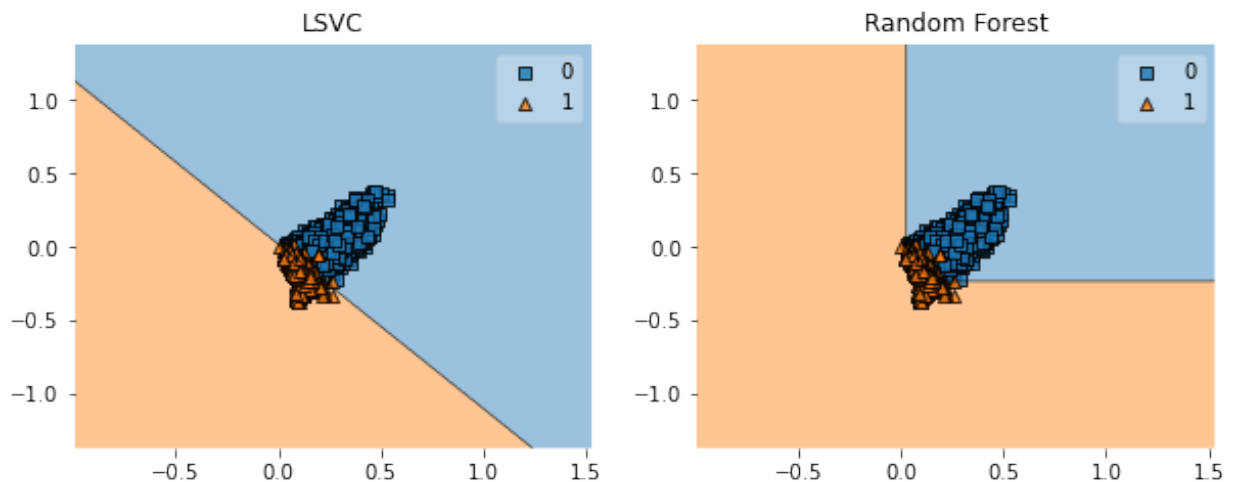
         clf_cv_mean = []
         clf_cv_std = []
         for clf, label, grd in zip(clf_list, label, grid):

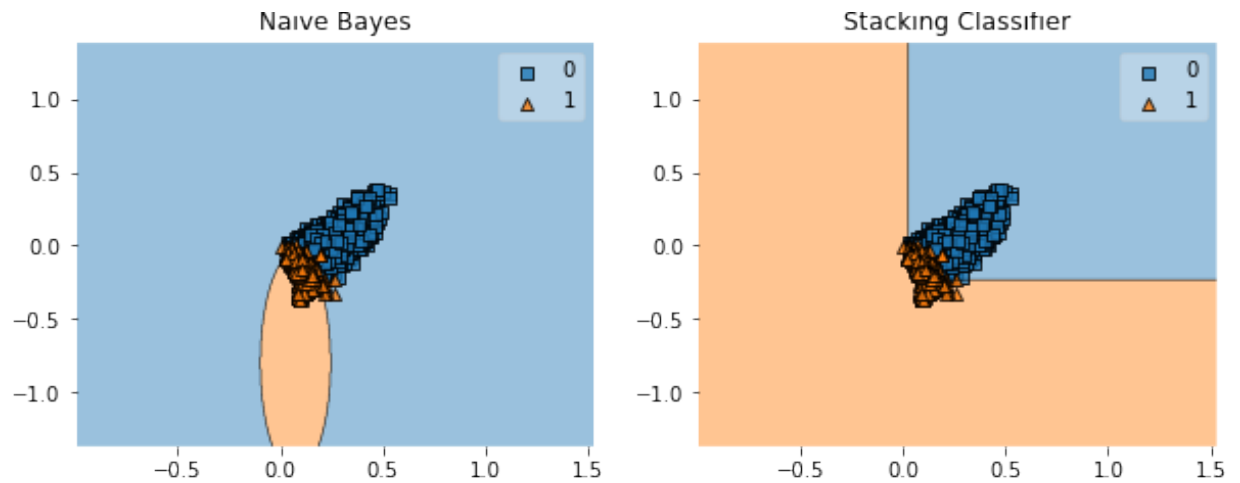
             scores = cross_val_score(clf, X_new, y, cv=3, scoring='accuracy')
             print("Accuracy: %.2f (+/- %.2f) [%s]" %(scores.mean(), scores.std(), label))
             clf_cv_mean.append(scores.mean())
             clf_cv_std.append(scores.std())

             clf.fit(X_new, y)
             ax = plt.subplot(gs[grd[0], grd[1]])
             fig = plot_decision_regions(X=X_new, y=y, clf=clf)
             plt.title(label)

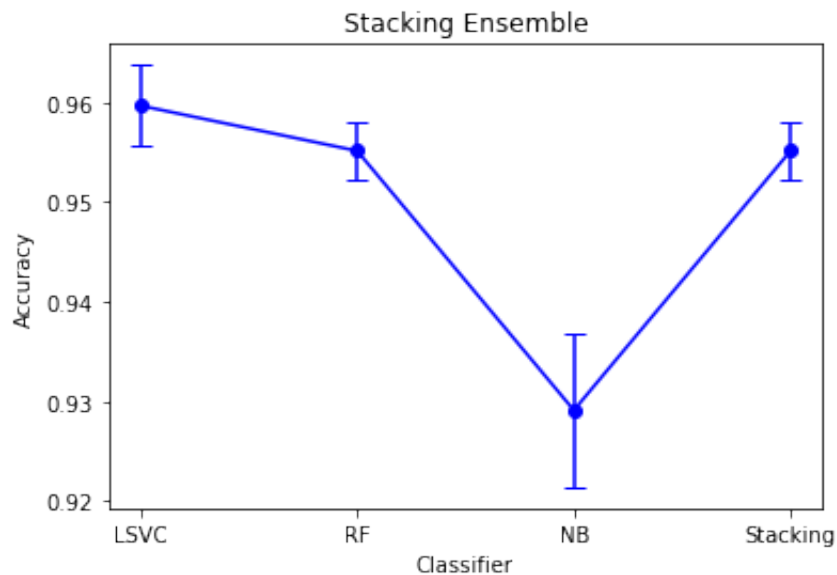
         plt.show()
```

```
Accuracy: 0.96 (+/- 0.00) [LSVC]
Accuracy: 0.96 (+/- 0.00) [Random Forest]
Accuracy: 0.93 (+/- 0.01) [Naive Bayes]
Accuracy: 0.96 (+/- 0.00) [Stacking Classifier]
```





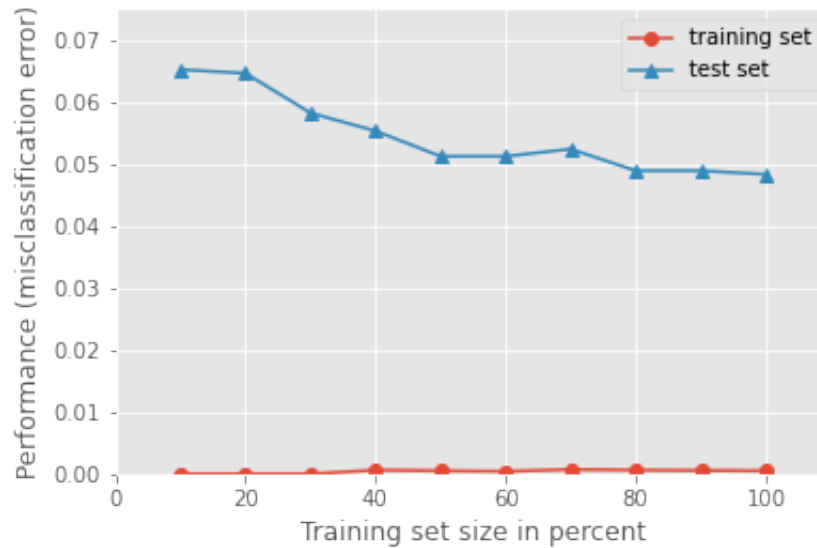
```
In [97]: #plot classifier accuracy
plt.figure()
(_, caps, _) = plt.errorbar(range(4), clf_cv_mean, yerr=clf_cv_std, c=
for cap in caps:
    cap.set_maheredgewidth(1)
plt.xticks(range(4), ['LSVC', 'RF', 'NB', 'Stacking'])
plt.ylabel('Accuracy'); plt.xlabel('Classifier'); plt.title('Stacking
plt.show()
```



In [98]: *#plot learning curves*

```
X_train, X_test, y_train, y_test = train_test_split(X_new, y, test_size=0.2)

plt.figure()
plot_learning_curves(X_train, y_train, X_test, y_test, scrf, print_model=False)
plt.show()
```



In the training set, we can see that stacking achieves higher accuracy than individual classifiers and based on learning curves, it shows no signs of overfitting. However the test set has more errors than that of training set