Ranc	scription of the algorithm
tree It use pred	om forest, like its name implies, consists of a large number of individual decision trees that operate as an ensemble. Each individual name the random forest spits out a class prediction and the class with the most votes becomes our model's prediction.  Es bagging and feature randomness when building each individual tree to try to create an uncorrelated forest of trees whose ction by committee is more accurate than that of any individual tree.  Es some trees may be wrong, many other trees will be right, so as a group the trees are able to move in the correct direction.  Training Dataset
	Bagging $ \begin{bmatrix} x^{\frac{1}{2}} & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ x^{\frac{1}{2}} & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ x^{\frac{1}{2}} & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ x^{\frac{1}{2}} & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ x^{\frac{1}{2}} & 0 & 0 & 0 & \frac{1}{2} \\ x^{\frac{1}{2}} & 0 & 0 & 0 & 0 & \frac{1}{2} \\ x^{\frac{1}{2}} & 0 & 0 & 0 & 0 & 0 & \frac{1}{2} \\ x^{\frac{1}{2}} & 0 & 0 & 0 & 0 & 0 & \frac{1}{2} \\ x^{\frac{1}{2}} & 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{2} \\ x^{\frac{1}{2}} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ x^{\frac{1}{2}} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ x^{\frac{1}{2}} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ x^{\frac{1}{2}} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 &$
-	1st Tree  2nd Tree  Nth Tree  Ist weak classifier  2nd weak classifier  Nth weak classifier
Key •	refits and Challenges  Benefits  Reduced risk of overfitting: Decision trees run the risk of overfitting as they tend to tightly fit all the samples within training dath dowever, when there's a robust number of decision trees in a random forest, the classifier won't overfit the model since the average uncorrelated trees lowers the overall variance and prediction error.
•	Provides flexibility: Since random forest can handle both regression and classification tasks with a high degree of accuracy, it is copular method among data scientists. Feature bagging also makes the random forest classifier an effective tool for estimating missing values as it maintains accuracy when a portion of the data is missing.  Easy to determine feature importance: Random forest makes it easy to evaluate variable importance, or contribution, to the material feature importance. Gini importance and mean decrease in impurity (MDI) are usually used to measure how much the model's accuracy decreases when a given variable is excluded. However, permutation importance, also known as mean decrease accuracy (MDA), is another importance measure. MDA identifies the average decrease in accuracy by randomly permutating the feature values in oob samples.
•	Challenges  Time-consuming process: Since random forest algorithms can handle large data sets, they can be provide more accurate prediction to the slow to process data as they are computing data for each individual decision tree.  Requires more resources: Since random forests process larger data sets, they'll require more resources to store that data.  More complex: The prediction of a single decision tree is easier to interpret when compared to a forest of them.
In or Usefi	der to generate my dataset i used the make_classification function from sklearn datasets  ul parameter to know:  n_samples refers to the number of samples.  n_features refers to the number of columns/features of dataset. These comprise n_informative and n_redundant.  n_informative refers to the number of features that will be useful in helping to classify your test dataset. In other words, if you
•	perform PCA or another dimensionality reduction algorithm, you should be able to explain nearly 100 % of the variance in the provided in the number of redundant features. These features are generated as random linear combinations of the informative features.  In the number of redundant features. These features are generated as random linear combinations of the informative features.  In the number of reduction algorithm, you should be able to explain nearly 100 % of the variance in the provided in the prov
X, pli pli	<pre>Define dataset y = make_classification(n_samples=300, n_features=20, n_informative=15, n_redundant=5, random_state=3 c.figure() c.title("Original Dataset") c.scatter(X[:,0], X[:,1], s=70, c=y, alpha=0.8) tplotlib.collections.PathCollection at 0x1afc5275f70&gt;</pre>
4 - 2 - 0 2 4 6 -	
The the set we have the set we	dation Set  raining error does not provide good estimates of the test error.  The ple way, which is much better than the training error, to estimate the test error is to use a hold-out set: a small part of the train thich is held-out from training.  It is case i use the train_test_split function to split the starting dataset in 30% test set and 70% training set.  The ple way, which is much better than the training error, to estimate the test error is to use a hold-out set: a small part of the training thick is held-out from training.  The ple way, which is much better than the training error, to estimate the test error is to use a hold-out set: a small part of the training thick is held-out from training.  The ple way, which is much better than the training error, to estimate the test error is to use a hold-out set: a small part of the training thick is held-out from training.
To bu	Instruction and Training of the model  wild and train the model i used the RandomForestClassifier from the sklearn module.  Interpretation in order to get a detailed report of the model.  It print_score (model, X_train, y_train, X_test, y_test, train=False):  if train:
	<pre>pred = model.predict(X_train) # Prediction  model_report = pd.DataFrame(classification_report(y_train, pred, output_dict=True)) print("\033[94m Train Result:\n====================================</pre>
	<pre>model_report = pd.DataFrame(classification_report(y_test, pred, output_dict=True))     print("\033[94m Test Result:\n====================================</pre>
# 1 Y_1 Y_1 Pr:	fit the model on the whole dataset  del.fit(X_train, y_train)  make a single prediction  train_pred = print_score(model, X_train, y_train, X_test, y_test, train=True)  test_pred = print_score(model, X_train, y_train, X_test, y_test, train=False)  int("\033[1m Labels of the Test set: \033[0m\n", y_test)  int("\n\033[1m Predicted Labels of Test set: \033[0m\n", y_test_pred)  ain Result:
	reuracy Score: 100.00%  Infusion Matrix:  Ill 0  0 99]]  St Result:  Curacy Score: 85.56%  Infusion Matrix:  37 4]
Lai [1 0 0 Pr [1 0 0 0	Dels of the Test set:  1 0 0 1 1 1 1 1 1 1 0 1 0 1 1 1 1 0 0 0 0 1 1 0 0 1 1 0 1 1 0 1 1 0 0 0 1  1 1 1 0 1 1 1 1
fic	s part i plot the first 5 trees of the random forest, as we can see they are all uncorrelated, that's a good thing.  g, axes = plt.subplots(nrows = 1,ncols = 5,figsize = (10,2), dpi=900)  sindex in range(0, 5): tree.plot_tree(model.estimators_[index],
	Estimator: 0 Estimator: 2 Estimator: 3 Estimator  A A A A A A A A A A A A A A A A A A A
<b>Eval</b> The	perimental Assessment  uation of the prediction  calcError function compute the error in the predictions of a certain model.  calcError(Ypred, Ytrue): return (np.count_nonzero(Ypred != Ytrue)) / len(Ytrue)
In th Class no: tre fo:	mparison between the Decision Tree and the Random Forest  s part i tried different noises during the generation of the datasets in order to compare the predictions using the Random Forest ifier and the Decision Tree Classifier.  se_array = [0.1, 0.2, 0.3, 0.4, 0.5] se_error = [] rest_error = [] st i in noise_array:
	<pre># data generation with noise X, y = make_classification(n_samples=300, n_features=20, n_informative=15, n_redundant=5, random_state, y = mixGauss([[0,0],[1,1]], [i, i], 300) X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3) # dtc = decision tree classifier dtc = tree.DecisionTreeClassifier() dtc.fit(X_train, y_train) dtc_pred = dtc.predict(X_test) # rfc = random forest classifier</pre>
	<pre>rfc = RandomForestClassifier() rfc.fit(X_train, y_train) rfc_pred = rfc.predict(X_test)  # error dtc_err = calcError(dtc_pred, y_test) rfc_err = calcError(rfc_pred, y_test)  tree_error.append(dtc_err) forest_error.append(rfc_err)</pre>
fic ax ax ax ax ax	Plot the validation and test errors  g, ax = plt.subplots() set_title("Error on the Test Set") plot(noise_array, tree_error, label="Tree") plot(noise_array, forest_error, label="Forest") set_xlabel("Noise") set_ylabel("Missclassified %") legend();  Error on the Test Set  Forest
Missclassified %	40 - 35 - 30 - 25 - 20 - 25 - 0.30 - 0.35 - 0.40 - 0.45 - 0.50 - Noise
K-F In th	F KFoldCV(Xtr, Ytr, num_folds, k_list, name):
	<pre>"""Run K-Fold CV for the Random Forest model  Parameters: Xtr : np.array     the full training set data - Ytr : np.array     the full training set labels - num_folds : int     the number of folds - k_list : List[int]     the values of k to try.</pre>
	Returns: best_k : int     The value of k (in k_list) which obtains the best average validation error - best_k_idx : int     The index of the best_k element in k_list - tr_err_mean : np.array     A 1D array of the same length as k_list, with the average training error for each tested k tr_err_std : np.array     A 1D array of the same length as k_list, with the standard deviation     of the training error for each tested k.
	<pre>- val_err_mean : np.array     A 1D array of the same length as k_list, with the average validation error for each tested k val_err_std : np.array     A 1D array of the same length as k_list, with the standard deviation     of the validation error for each tested k. """  rnd_state = np.random.RandomState() # Ensures that k_list is a numpy array k_list = np.array(k_list) num_k = len(k_list)  n_tot = Xtr.shape[0]</pre>
	<pre>n_val = n_tot // num_folds  # We want to compute 1 error for each `k` and each fold  tr_errors = np.zeros((num_k, num_folds))  val_errors = np.zeros((num_k, num_folds))  # `split_idx`: a list of arrays, each containing the validation indices for 1 fold  rand_idx = rnd_state.choice(n_tot, size=n_tot, replace=False)  split_idx = np.array_split(rand_idx, num_folds)  for fold_idx in range(num_folds):     # Set the indices in boolean mask for all validation samples to `True`</pre>
	<pre>val_mask = np.zeros(n_tot, dtype=bool) val_mask[split_idx[fold_idx]] = True # Split training set in training part and validation part # Hint: you can use boolean mask as index vector to split Xtr and Ytr  x_train = Xtr[~val_mask] y_train = Ytr[~val_mask] x_val = Xtr[val_mask] y_val = Ytr[val_mask]  for k_idx, current_k in enumerate(k_list):  if name=='estimators':</pre>
	<pre>rfc = RandomForestClassifier(n_estimators=current_k) elif name=='features':     rfc = RandomForestClassifier(max_features=current_k) elif name=='samples':     rfc = RandomForestClassifier(max_samples=current_k) elif name=='depth':     rfc = RandomForestClassifier(max_depth=current_k) else:     rfc = RandomForestClassifier()  rfc.fit(x_train, y_train) rfc_pred_train = rfc.predict(x_train)</pre>
	<pre>rfc_pred_val = rfc.predict(x_val)  # Compute the training error of the kNN classifier for the given value of k tr_errors[k_idx, fold_idx] = calcError(rfc_pred_train, y_train) # Compute the validation error of the kNN classifier for the given value of k val_errors[k_idx, fold_idx] = calcError(rfc_pred_val, y_val)  # Calculate error statistics along the repetitions: # 1) mean training error, training error standard deviation tr_mean = np.mean(tr_errors, axis=1) tr_var = np.var(tr_errors, axis=1)</pre>
Pa	<pre># 2) mean validation error, validation error standard deviation val_mean = np.mean(val_errors, axis=1) val_var = np.var(val_errors, axis=1)  # 3) best k (k which minimize mean validation error) and index of best k in k_list best_k_idx = np.argmin(val_mean) best_k = k_list[best_k_idx]  return best_k, best_k_idx, tr_mean, tr_var, val_mean, val_var</pre> ndom Forest Hyperparameters
In the performance of the perfor	s section, i will take a closer look at some of the hyperparameters tuning for the random forest ensemble and their effect on modernance.  i use the make_classification function to create a synthetic binary classification problem. Then, through the in_test_split function, i split the dataset into a 70% training set and 30% validation set.  y = make_classification(n_samples=300, n_features=20, n_informative=15, n_redundant=5, random_state=3 train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3)
	figure(figsize=(15,7))
pli pli pli pli	<pre>c.subplot(1,2,1) c.title("Training set") c.scatter(X_train[:,0], X_train[:,1], s=70, c=y_train, alpha=0.8) c.subplot(1,2,2) c.title("Validation set") c.scatter(X_test[:,0], X_test[:,1], s=70, c=y_test, alpha=0.8) c.plotlib.collections.PathCollection at 0x1afc5568c70&gt; Training set  Validation set</pre>
pla pla pla pla	<pre>c.subplot(1,2,1) c.title("Training set") c.scatter(X_train[:,0], X_train[:,1], s=70, c=y_train, alpha=0.8) c.subplot(1,2,2) c.title("Validation set") c.scatter(X_test[:,0], X_test[:,1], s=70, c=y_test, alpha=0.8) c.scatter(X_test[:,0], X_test[:,1], s=70, c=y_test, alpha=0.8)</pre>
pl: pl: pl: cma 6	<pre>c.subplot(1,2,1)title("Training set")scatter(X_train(:,0), X_train(:,1), s=70, c=y_train, alpha=0.8)subplot(1,2,2)scatter(X_train(:,0), X_test(:,1), s=70, c=y_test, alpha=0.8) .tplotlib.collections.PathCollection at 0xlafc5568c70&gt; Training set  Validation set  6 4 2 0 4 6 4 6 4 6 4 7 7 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0</pre>
plipliplication over train  plipliplication over train	<pre>c.subplot(1,2,1) c.title("Training set") c.scatter(X_train[:,0], X_train[:,1], s=70, c=y_train, alpha=0.8) c.subplot(1,2,2) c.title("Validation set") c.scatter(X_test[:,0], X_test[:,1], s=70, c=y_test, alpha=0.8) c.plotlib.collections.PathCollection at Oxlafc5568c70&gt;  Training set  Validation set  6- 4- 4- 4- 4- 4- 6- 4- 6-</pre>
plipliplication over the content of	Little "Training set") Listle ("Training set") Listle ("Training set") Listle ("Validation set")
plipliplication of the control of th	Little (Teaming set") Little (Teaming set ") Litt
The is for train.  # in ax	Little (Trislation) set (") Li
# fic ax ax ax Tex On	subject (1,2,1) statical carriers age ** stati
# fice ax	instances are producted to the number of trees changes of the section of the sect
The strain	Label (1997)   Secret (1997)
The standard	The proposal process of the process
# fix ax	consideration (and consideration) and consideration (and consideration
The state of the s	contribution (%):    Contribution (%):   Contr
The start of the s	Consideration of the control of the
The series of th	The number of trees  The numbe
The service of the se	The complete of trees  The number of trees  The num
The series of th	And the company of th
Using the second of the second	and the control of th
The series of th	The control of the co
In the series of	The surface of the control of the co
In the second of	The second control of trees  The number of trees  T
1. The second of	The strategies of the state of
In the second of	The control of the co
In the second of	The control of the co
The series of th	The second of the control of the con
In the State of th	The resulting of these processes and a process of the second of the seco
The second of th	See Committee of these  For Co
# In the second of the second	The company of the co
In the Start of th	Section 1997 1997 1997 1997 1997 1997 1997 199
In the second of	
In the series of	The content of the
Using the second of the second	The control of the co
1. The size of the second of t	

In this case, we can see that larger depth results in better model performance, with the default of no maximum depth we achieving the best performance on this dataset.	0.36 - 0.35 - 0.34 - 0.33 - 0.32 - 0.31 -	10-Fold CV  Validation  depth  Validation				
		e that larger depth results in be	tter model performance, w	ith the default of no maximum	n depth we achieving the	