

U

O

W

Classifier Evaluation and Ensemble Learning

CSCI316 Big Data Mining Techniques and Implementation



UNIVERSITY
OF WOLLONGONG
AUSTRALIA

Model Evaluation in Classification



UNIVERSITY
OF WOLLONGONG
AUSTRALIA

Model Evaluation and Selection

- Evaluation metrics: How can we measure accuracy? Other metrics to consider?
- Use a separate **test set** of class-labeled tuples instead of training set when assessing accuracy
- Methods for estimating a classifier's accuracy:
 - Holdout method, random subsampling
 - Cross-validation
 - Bootstrap
- Comparing classifiers:
 - Confidence intervals
 - ROC Curves

Model Evaluation

How do we evaluate our classifiers?

- **Positive tuples (P)**: tuples of the main class of interest (e.g., `buys_computers = yes`)
- **Negative tuples (N)**: tuples of all other classes
- **True positives (TP)**: the positive tuples that were correctly labeled by the classifier.
- **True negatives (TN)**: the negative tuples that were correctly labeled by the classifier.
- **False positive/false negative (FP/FN)**: the negative/positive tuples that were incorrectly labeled as positive/negative
 - For convivence, we also use P/N/TP/TN/FP/FN to denote the *number* of P/N/TP/TN/FP/FN, respectively.

Classifier Evaluation Metrics: Confusion Matrix

Confusion Matrix:

Actual class\Predicted class	C_1	$\neg C_1$
C_1	True Positives (TP)	False Negatives (FN)
$\neg C_1$	False Positives (FP)	True Negatives (TN)

Example of Confusion Matrix:

Actual class\Predicted class	buy_computer = yes	buy_computer = no	Total
buy_computer = yes	6954	46	7000
buy_computer = no	412	2588	3000
Total	7366	2634	10000

- More generally, given m classes, an entry, $CM_{i,j}$ in a **confusion matrix** indicates # of tuples in class i that were labeled by the classifier as class j
 - May have extra rows/columns to provide totals

Classifier Evaluation Metrics: Accuracy, Error Rate, Sensitivity and Specificity

A\P	C	¬C	
C	TP	FN	P
¬C	FP	TN	N
	P'	N'	All

- **Classifier Accuracy**, or recognition rate: percentage of test set tuples that are correctly classified

$$\text{Accuracy} = (\text{TP} + \text{TN})/\text{All}$$

- **Error rate**: $1 - \text{accuracy}$, or
 $\text{Error rate} = (\text{FP} + \text{FN})/\text{All}$

- **Class Imbalance Problem:**

- One class may be *rare*, e.g. fraud, or HIV-positive
- Significant *majority of the negative class* and minority of the positive class
- **Sensitivity**: True Positive recognition rate
 - **Sensitivity** = TP/P
- **Specificity**: True Negative recognition rate
 - **Specificity** = TN/N

Classifier Evaluation Metrics:

Precision and Recall, and F-measures

- **Precision:** exactness – what % of tuples that the classifier labeled as positive are actually positive

$$precision = \frac{TP}{TP + FP}$$

- **Recall:** completeness – what % of positive tuples did the classifier label as positive? (equals to sensitivity)

$$recall = \frac{TP}{TP + FN}$$

- Perfect score is 1.0
- In practice, inverse relationship between precision & recall
- **F measure (F_1 or F-score):** harmonic mean of precision and recall,

$$F = \frac{2 \times precision \times recall}{precision + recall}$$

- F_β : weighted measure of precision & recall (common β value: 2 or 0.5)
 - assigns β times as much weight to recall as to precision

$$F_\beta = \frac{(1 + \beta^2) \times precision \times recall}{\beta^2 \times precision + recall}$$



Classifier Evaluation Metrics: Example

Actual Class\Predicted class	cancer = yes	cancer = no	Total	Recognition(%)
cancer = yes	90	210	300	30.00 (<i>sensitivity</i>)
cancer = no	140	9560	9700	98.56 (<i>specificity</i>)
Total	230	9770	10000	96.50 (<i>accuracy</i>)

- *Precision* = $90/230 = 39.13\%$
- *Recall* = $90/300 = 30.00\% = \text{Sensitivity}$

Evaluating Classifier Accuracy: Holdout & Cross-Validation Methods

- **Holdout method**

- Given data is randomly partitioned into two independent sets
 - Training set (e.g., 2/3) for model construction
 - Test set (e.g., 1/3) for accuracy estimation
- Random sampling: a variation of holdout
 - Repeat holdout k times, accuracy = avg. of the accuracies obtained

- **Cross-validation** (k -fold, where $k = 10$ is most popular)

- Randomly partition the data into k *mutually exclusive* subsets, each approximately equal size
- At i -th iteration, use D_i as test set and others as training set
- Leave-one-out: k folds where $k = \#$ of tuples, for small sized data
- ***Stratified cross-validation***: folds are stratified so that class dist. in each fold is approx. the same as that in the initial data

Evaluating Classifier Accuracy: Bootstrap

- **Bootstrap Sampling**

- Works well with small data sets
- Samples the given training tuples uniformly *with replacement*
 - i.e., each time a tuple is selected, it is equally likely to be selected again and re-added to the training set

- Several bootstrap methods, and a common one is **.632 bootstrap**

- A data set with d tuples is sampled d times, with replacement, resulting in a training set of d samples. The data tuples that did not make it into the training set end up forming the test set. About 63.2% of the original data end up in the bootstrap, and the remaining 36.8% form the test set (since $(1 - 1/d)^d \approx e^{-1} = 0.368$)
- Repeat the sampling procedure k times, overall accuracy of the model:

$$Acc(M) = \frac{1}{k} \sum_{i=1}^k (0.632 \times Acc(M_i)_{test_set} + 0.368 \times Acc(M_i)_{train_set})$$



Estimating Confidence Intervals: Classifier Models M_1 vs. M_2

- Suppose we have 2 classifiers, M_1 and M_2 , which one is better?
- Use 10-fold cross-validation to obtain $\overline{err}(M_1)$ and $\overline{err}(M_2)$
- These mean error rates are just *estimates* of error on the true population of *future* data cases
- What if the difference between the 2 error rates is just attributed to *chance*?
 - Use a **test of statistical significance**
 - Obtain **confidence limits** for our error estimates

Estimating Confidence Intervals: Null Hypothesis

- Perform 10-fold cross-validation
- Assume samples follow a **t distribution** with $k-1$ **degrees of freedom** (here, $k=10$)
- Use **t-test** (or **Student's t-test**)
- **Null Hypothesis:** M_1 & M_2 are the same
- If we can **reject** null hypothesis, then
 - we conclude that the difference between M_1 & M_2 is **statistically significant**
 - Chose model with lower error rate

Estimating Confidence Intervals: t-test

- **Pairwise Comparison**

- For i^{th} round of 10-fold cross-validation, the same cross partitioning is used to obtain $err(M_1)_i$ and $err(M_2)_i$
- Average over 10 rounds to get $\overline{err}(M_1)$ and $\overline{err}(M_2)$
- **t-test** computes **t-statistic** with $k-1$ **degrees of freedom**:

$$t = \frac{\overline{err}(M_1) - \overline{err}(M_2)}{\sqrt{var(M_1 - M_2)/k}}$$

where

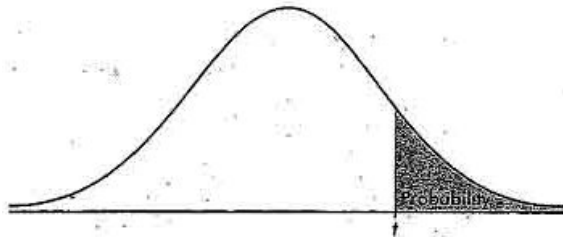
$$var(M_1 - M_2) = \frac{1}{k} \sum_{i=1}^k \left[err(M_1)_i - err(M_2)_i - (\overline{err}(M_1) - \overline{err}(M_2)) \right]^2$$

Estimating Confidence Intervals: Statistical Significance

- Are M_1 & M_2 **significantly different**?
 - Compute t and select *significance level* (e.g. $\text{sig} = 5\%$)
 - Consult the **t-distribution** with $k-1$ (here $k=10$) degree of freedom
 - Look for the value z in the **t-distribution** corresponding to confidence limit $\text{sig}/2$ (here 0.025)
 - **If $t > z$ or $t < -z$** , then t lies in rejection region:
 - **Reject null hypothesis** that mean error rates of M_1 & M_2 are same
 - Conclude: statistically significant difference between M_1 & M_2
 - **Otherwise**, conclude that any difference is **by chance**

Estimating Confidence Intervals:

Table for t-distribution



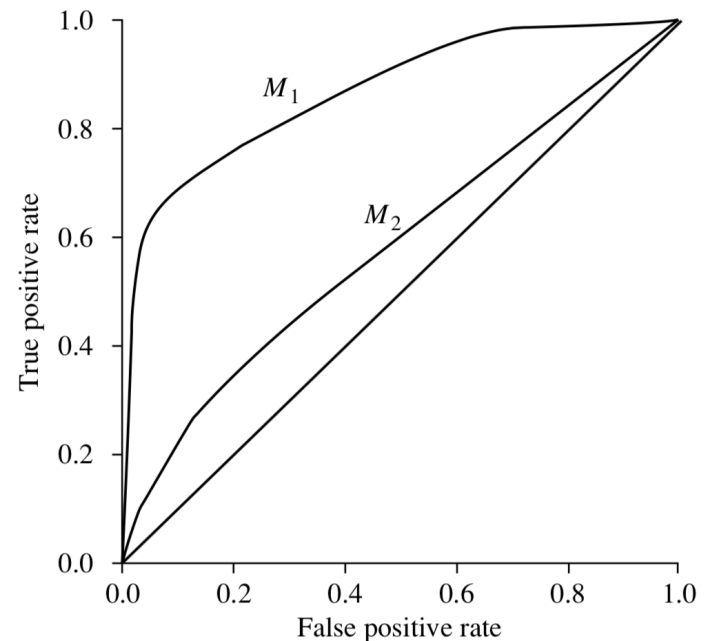
- Symmetric
- **Significance level**, e.g., $\text{sig} = 0.05$ or 5% means M_1 & M_2 are *significantly different* for 95% of population
- **Confidence limit**, $z = \text{sig}/2$

TABLE B: t-DISTRIBUTION CRITICAL VALUES

df	Tail probability p										
	.25	.20	.15	.10	.05	.025	.02	.01	.005	.0025	.001
1	1.000	1.376	1.963	3.078	6.314	12.71	15.89	31.82	63.66	127.3	318.3
2	.816	1.061	1.386	1.886	2.920	4.303	4.849	6.965	9.925	14.09	22.33
3	.765	.978	1.250	1.638	2.353	3.182	3.482	4.541	5.841	7.453	10.21
4	.741	.941	1.190	1.533	2.132	2.776	2.999	3.747	4.604	5.598	7.173
5	.727	.920	1.156	1.476	2.015	2.571	2.757	3.365	4.032	4.773	5.893
6	.718	.906	1.134	1.440	1.943	2.447	2.612	3.143	3.707	4.317	5.208
7	.711	.896	1.119	1.415	1.895	2.365	2.517	2.998	3.499	4.029	4.785
8	.706	.889	1.108	1.397	1.860	2.306	2.449	2.896	3.355	3.833	4.501
9	.703	.883	1.100	1.383	1.833	2.262	2.398	2.821	3.250	3.690	4.297
10	.700	.879	1.093	1.372	1.812	2.228	2.359	2.764	3.169	3.581	4.144
11	.697	.876	1.088	1.363	1.796	2.201	2.328	2.718	3.106	3.497	4.025
12	.695	.873	1.083	1.356	1.782	2.179	2.303	2.681	3.055	3.428	3.930
13	.694	.870	1.079	1.350	1.771	2.160	2.282	2.650	3.012	3.372	3.852
14	.692	.868	1.076	1.345	1.761	2.145	2.264	2.624	2.977	3.326	3.787
15	.691	.866	1.074	1.341	1.753	2.131	2.249	2.602	2.947	3.286	3.733
16	.690	.865	1.071	1.337	1.746	2.120	2.235	2.583	2.921	3.252	3.686
17	.689	.863	1.069	1.333	1.740	2.110	2.224	2.567	2.898	3.222	3.646
18	.688	.862	1.067	1.330	1.734	2.101	2.214	2.552	2.878	3.197	3.611
19	.688	.861	1.066	1.328	1.729	2.093	2.205	2.539	2.861	3.174	3.579
20	.687	.860	1.064	1.325	1.725	2.086	2.197	2.528	2.845	3.153	3.552
21	.686	.859	1.063	1.323	1.721	2.080	2.189	2.518	2.831	3.135	3.527
22	.686	.858	1.061	1.321	1.717	2.074	2.183	2.508	2.819	3.119	3.503
23	.685	.858	1.060	1.319	1.714	2.069	2.177	2.500	2.807	3.104	3.485
24	.685	.857	1.059	1.318	1.711	2.064	2.172	2.492	2.797	3.091	3.467
25	.684	.856	1.058	1.316	1.708	2.060	2.167	2.485	2.787	3.078	3.450
26	.684	.856	1.058	1.315	1.706	2.056	2.162	2.479	2.779	3.067	3.435
27	.684	.855	1.057	1.314	1.703	2.052	2.158	2.473	2.771	3.057	3.421
28	.683	.855	1.056	1.313	1.701	2.048	2.154	2.467	2.763	3.047	3.408
29	.683	.854	1.055	1.311	1.699	2.045	2.150	2.462	2.756	3.038	3.396
30	.683	.854	1.055	1.310	1.697	2.042	2.147	2.457	2.750	3.030	3.385
40	.681	.851	1.050	1.303	1.684	2.021	2.123	2.423	2.704	2.971	3.307
50	.679	.849	1.047	1.299	1.676	2.009	2.109	2.403	2.678	2.937	3.261
60	.679	.848	1.045	1.296	1.671	2.000	2.099	2.390	2.660	2.915	3.232
80	.678	.846	1.043	1.292	1.664	1.990	2.088	2.374	2.639	2.887	3.195
100	.677	.845	1.042	1.290	1.660	1.984	2.081	2.364	2.626	2.871	3.174
1000	.675	.842	1.037	1.282	1.646	1.962	2.056	2.330	2.581	2.813	3.098
∞	.674	.841	1.036	1.282	1.645	1.960	2.054	2.326	2.576	2.807	3.091
	50%	60%	70%	80%	90%	95%	96%	98%	99%	99.5%	99.8%
	Confidence level C										

Model Selection: ROC Curves

- **ROC** (Receiver Operating Characteristics) curves: for visual comparison of classification models
- Originated from signal detection theory
- Shows the trade-off between **true positive rate** ($TP/P = \text{sensitivity}$) and **false positive rate** ($FP/N = 1 - \text{specificity}$)
- The **area under the ROC curve** is a measure of **performance** of the model
- Rank the test tuples in decreasing order: the one that is most likely to belong to the positive class appears at the top of the list
- The closer to the diagonal line (i.e., the closer the area is to 0.5), the less accurate is the model



- Vertical axis represents the true positive rate
- Horizontal axis rep. the false positive rate
- The plot also shows a diagonal line
- A model with perfect accuracy will have an area of 1.0



Issues Affecting Model Selection

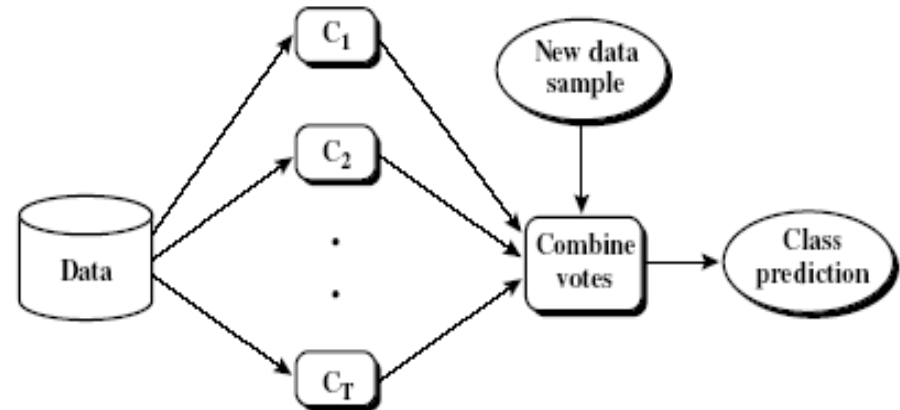
- **Accuracy**
 - classifier accuracy: predicting class label
- **Speed**
 - time to construct the model (training time)
 - time to use the model (classification/prediction time)
- **Robustness**: handling noise and missing values
- **Scalability**: efficiency in disk-resident databases
- **Interpretability**
 - understanding and insight provided by the model
- Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules

Model Enhancement Methods



UNIVERSITY
OF WOLLONGONG
AUSTRALIA

Ensemble Methods:

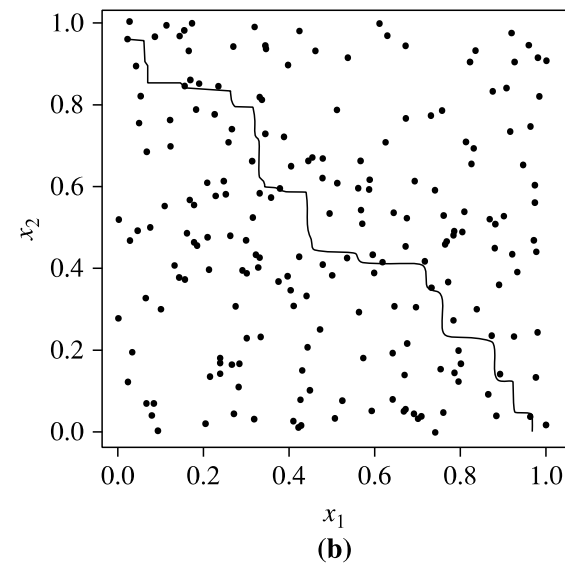
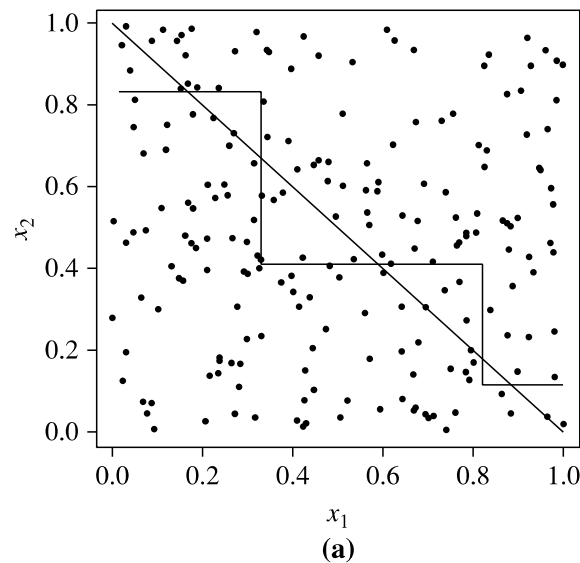


- Ensemble methods
 - Use a combination of models
 - Combine a series of k learned models, M_1, M_2, \dots, M_k , with the aim of creating a combined model M^*
- Popular ensemble methods
 - Bagging: averaging the prediction over a collection of classifiers
 - Boosting: weighted vote with a collection of classifiers
 - Ensemble: combining a set of heterogeneous classifiers

Ensemble Methods:

- Advantages:
 - Increase accuracy: Miss classification occurs only when more than half of base classifiers predict incorrectly (even better if the base classifiers are less correlated).
 - Can deal with data in sheer volume (too many records or attributes)
 - Can run in parallel

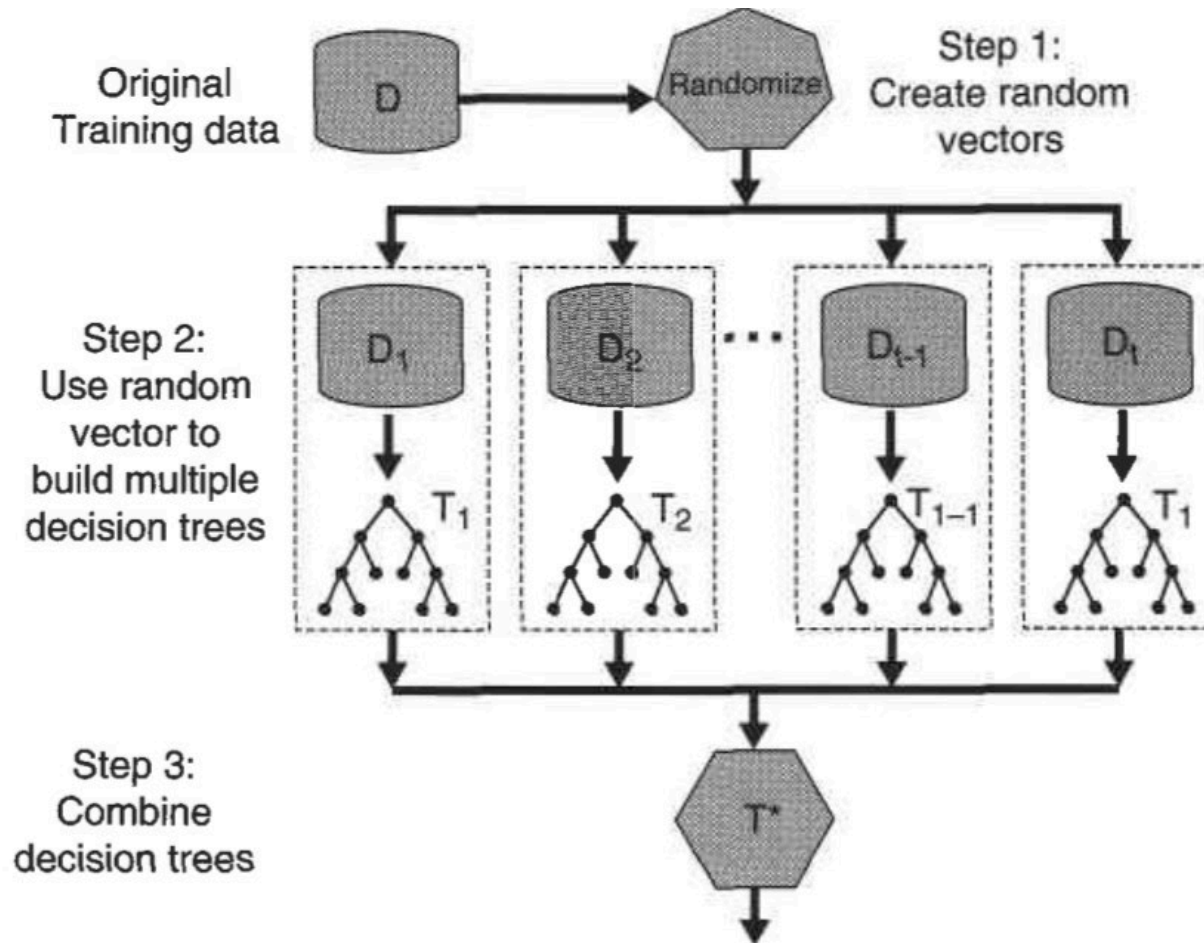
Decision boundary by (a) a single decision tree and (b) a random forest



Random Forest

- Random Forest is a class of ensemble methods specifically designed for decision tree classifiers.
 - It combines the predictions made by multiple decision trees.
 - Each tree is generated randomly based on the training tuples.
 - The final prediction output is produced by a voting function.
- A properly built random forest tends to be more accurate and less biased than individual decision tree classifiers.
 - The accuracy of RF depends on the *strength* of individual classifiers (trees) and a measure of *dependence* between them.
- But the computational cost grows as the number of trees in the forest increases.

Random Forest



Random Forest

- There are 3 common ways to associate randomization with decision trees.
- Bagging: Given a set D of d tuples, bagging works as follows. For iteration i ($i = 1, 2, \dots, k$), a training set D_i of d tuples is sampled *with replacement* from the original set D .
- Note that some of the original tuples of D may not be included in D_i , whereas others may occur more than once.
- A decision tree M_i is learned for each training set, D_i . To classify an unknown tuple X , each classifier M_i returns its class prediction, which counts as one vote.
- The bagged classifier, say, M_* counts the votes and assigns the class with the most votes to X .
- Random Forests can handle datasets that don't fit in memory

Random Forest

- Forest-RI (*random input selection*)
 - When building the tree, randomly select F attributes that are used to determine the split at each node, where F is much smaller than the number of available attributes.
 - The CART methodology is used to grow the trees to maximum size
 - If F is sufficiently small, then the trees tend to become less correlated
 - On the other hand, the strength of the tree classifier tends to improve with a larger F
 - As a trade-off, usually let $F = \log_2 m + 1$, where m is the number of input features
 - Since only a subset of the features is selected at each node, this approach helps to significantly reduce the runtime of the algorithm.

Random Forest

- Forest-RC (*random linear combinations*)
 - creates new attributes (or features) that are a linear combination of the existing attributes.
 - That is, an attribute is generated by specifying L , the number of original attributes to be combined.
 - At a given node, L attributes are randomly selected and added together with coefficients that are uniform random numbers on $[-1,1]$.
 - F linear combinations are generated, and a search is made over these for the best split.
 - This form of random forest is useful when there are only a few attributes available, so as to reduce the correlation between individual classifiers.