

Experiment 9

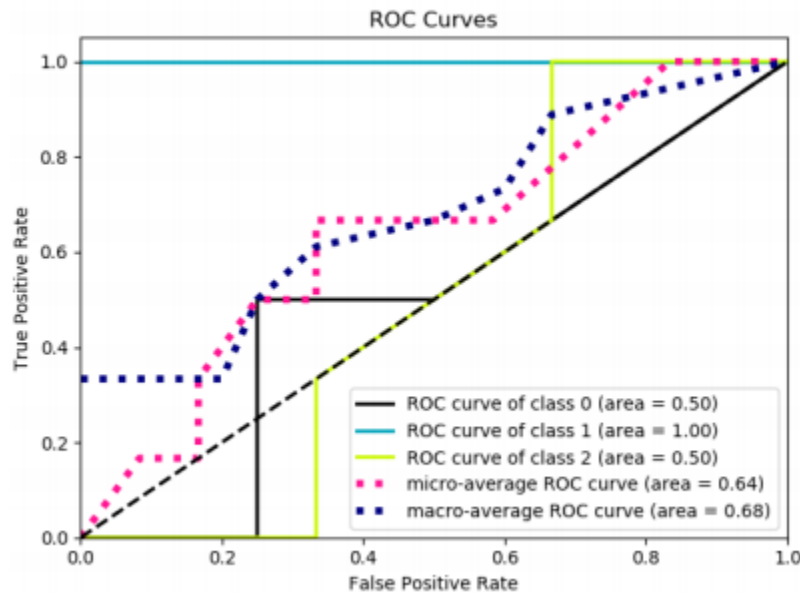
Aim: To perform ROC analysis, Forecast analysis and Survival analysis on any sample dataset

Theory:

Receiver Operator Curve(ROC) Analysis

A receiver operating characteristic curve, or ROC curve, is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied. The ROC curve is created by plotting the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings. The true-positive rate is also known as sensitivity, recall or probability of detection in machine learning. The false-positive rate is also known as the fall-out or probability of false alarm and can be calculated as $(1 - \text{specificity})$. It can also be thought of as a plot of the power as a function of the Type I Error of the decision rule (when the performance is calculated from just a sample of the population, it can be thought of as estimators of these quantities). The ROC curve is thus the sensitivity as a function of fall-out. In general, if the probability distributions for both detection and false alarm are known, the ROC curve can be generated by plotting the cumulative distribution function (area under the probability distribution from negative infinity to the discrimination threshold) of the detection probability in the y-axis versus the cumulative distribution function of the false-alarm probability on the x-axis.

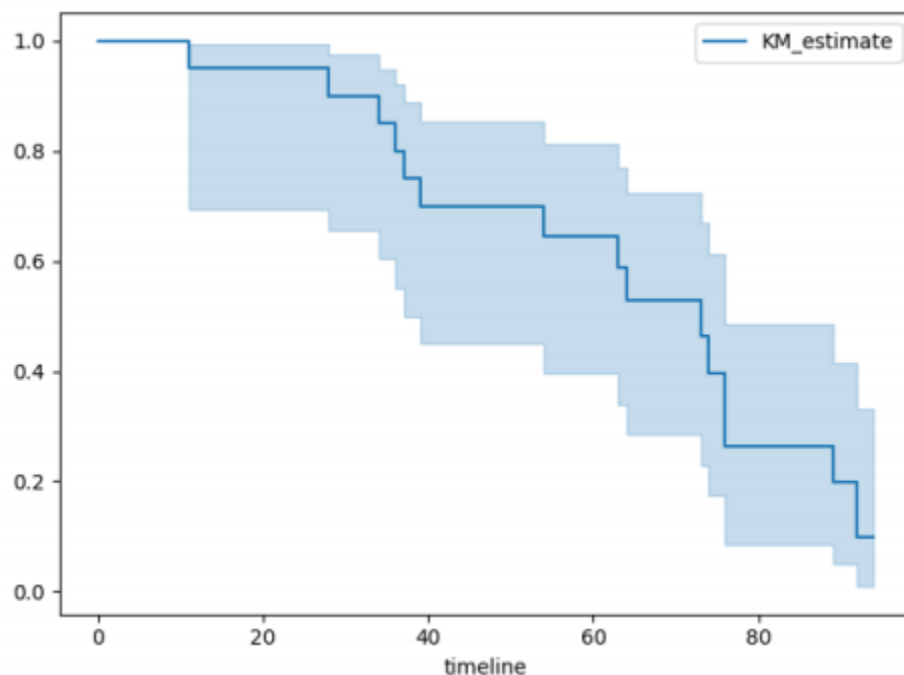
ROC analysis provides tools to select possibly optimal models and to discard suboptimal ones independently from the cost context or the class distribution. ROC analysis is related in a direct and natural way to cost/benefit analysis of diagnostic decision making.



Survival Analysis

Survival analysis is a branch of statistics for analyzing the expected duration of time until one or more events happen, such as death in biological organisms and failure in mechanical systems. This topic is called reliability theory or reliability analysis in engineering, duration analysis or duration modelling in economics, and event history analysis in sociology. Survival analysis attempts to answer questions such as: what is the proportion of a population which will survive past a certain time? Of those that survive, at what rate will they die or fail? Can multiple causes of death or failure be taken into account? How do particular circumstances or characteristics increase or decrease the probability of survival? To answer such questions, it is necessary to define "lifetime". In the case of biological survival, death is unambiguous, but for mechanical reliability, failure may not be well-defined, for there may well be mechanical systems in which failure is partial, a matter of degree, or not otherwise localized in time. Even in biological problems, some events (for example, heart attack or other organ failure) may have the same ambiguity. The theory outlined below assumes well-defined events at specific times; other cases may be better treated by models which explicitly account for ambiguous events.

More generally, survival analysis involves the modelling of time to event data; in this context, death or failure is considered an "event" in the survival analysis literature – traditionally only a single event occurs for each subject, after which the organism or mechanism is dead or broken. Recurring event or repeated event models relax that assumption. The study of recurring events is relevant in systems reliability, and in many areas of social sciences and medical research.



Forecast Analysis

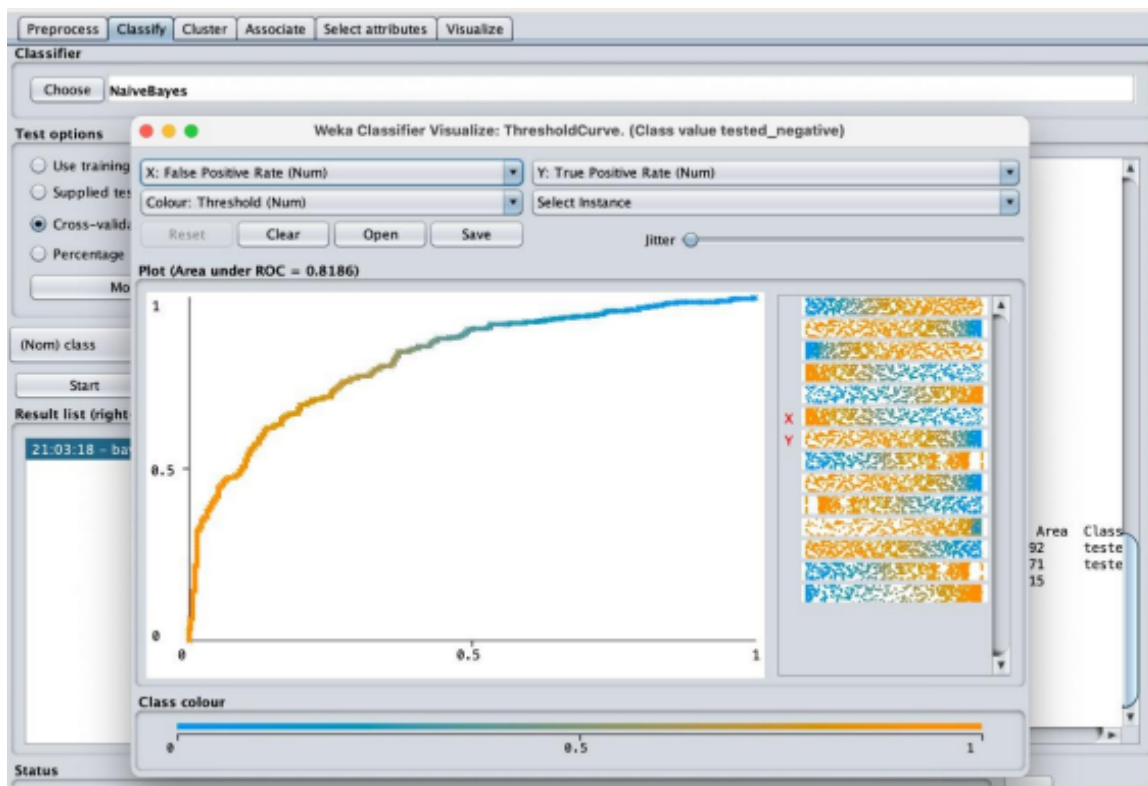
Forecasting is the process of making predictions of the future based on past and present data and most commonly by analysis of trends. A commonplace example might be estimation of some variable of interest at some specified future date. Prediction is a similar, but more general term. Both might refer to formal statistical methods employing time series, cross-sectional or longitudinal data, or alternatively to less formal judgmental methods. Usage can differ between areas of application: for example, in hydrology the terms "forecast" and "forecasting" are sometimes reserved for estimates of values at certain specific future times, while the term "prediction" is used for more general estimates, such as the number of times floods will occur over a long period.

Risk and uncertainty are central to forecasting and prediction; it is generally considered good practice to indicate the degree of uncertainty attaching to forecasts. In any case, the data must be up to date in order for the forecast to be as accurate as possible. In some cases the data used to predict the variable of interest is itself forecasted.

Output:

- **ROC Analysis**

ROC Curve for class tested_negative



- **Survival Analysis using evolutionary fuzzy classifier**

Preprocess Classify Cluster Associate Select attributes Visualize Forecast

Classifier

Choose **MultiObjectiveEvolutionaryFuzzyClassifier** -generations 20 -populationSize 100 -seed 1 -maxSimilarity 0.4 -minV 30.0 -maxV 2.0 -maxRules -1 -maxI

Test options

☐ Use training set
☐ Supplied test set Set...
☒ Cross-validation Folds
☐ Percentage split %
 More options...

(Nom) play

Start Stop

Result list (right-click for options)

21:34:37 - rules.MultiObjectiveEvolutionaryFuzzyClassifier

Classifier output

```

=== Run information ===

Scheme:      weka.classifiers.rules.MultiObjectiveEvolutionaryFuzzyClassifier -gene
Relation:    weather
Instances:    14
Attributes:   5
              outlook
              temperature
              humidity
              windy
              play
Test mode:    10-fold cross-validation

=== Classifier model (full training set) ===

Number of generations: 20      Population size: 100
Seed: 1
Maximum similarity: 0.4
Minimum variance parameter: 30.0
Maximum variance parameter: 2.0
Minimum number of rules: 2
Maximum number of rules: 12
Maximum number of labels: 5
Evaluation criteria: acc
Algorithm: ENORA
Report frequency: 20
Log file: /Users/manintirkey

Objective 1: acc, fmin = 0.0, fmax = 1.0, maximize
Objective 2: number of rules, fmin = 2.0, fmax = 12.0, minimize

Initial population
Time: 0.0 ns
  
```

Status

OK Log

Preprocess Classify Cluster Associate Select attributes Visualize Forecast

Classifier

Choose **MultiObjectiveEvolutionaryFuzzyClassifier** -generations 20 -populationSize 100 -seed 1 -maxSimilarity 0.4 -minV 30.0 -maxV 2.0 -maxRules -1 -maxLabels 5 -evaluationMeasure 0 -algorithm 0 -reportFrequency 20 -logFile /Users/manintirkey -run-decimal-pl

Test options

☐ Use training set
☐ Supplied test set Set...
☒ Cross-validation Folds
☐ Percentage split %
 More options...

(Nom) play

Start Stop

Result list (right-click for options)

21:34:37 - rules.MultiObjectiveEvolutionaryFuzzyClassifier

Classifier output

```

0.5714      11.0      0.4      d u f
0.3571      9.0      0.4      d u f
0.5      8.0      0.4      d u f
0.5      12.0      0.4      d u f
0.4286      11.0      0.4      d u f
0.5      10.0      0.4      d u f
0.3571      9.0      0.4      d u f
0.3571      4.0      0.4      d u f
0.4286      8.0      0.4      d u f
0.5714      9.0      0.4      d u f
0.5714      5.0      0.4      d u f

Generation number: 20
Time: 143.0 ms
Current mean acc: 0.6886
Normalized non-dominated space ratio: 0.15
Cross operators:
No cross: 0.31
Rule crossover: 0.33
Rule incremental crossover: 0.2
Fuzzy set crossover: 0.35
Mutation operators:
No mutate: 0.09
Gaussian set center mutation: 0.19
Gaussian set variance mutation: 0.11
Fuzzy set mutation: 0.19
Rule incremental mutation: 0.15
Integer mutation: 0.27
Amplitude for gaussian set center mutation: 0.4047
Amplitude for gaussian set variance mutation: 0.4962

u f -> Unfeasible Individuals
d -> Dominated Individuals
nd -> Non-dominated Individuals
s -> Selected Individual

acc      number of rules      similarity
0.3857      2.0      0.4      nd
0.8571      3.0      0.4      nd s
0.8571      3.0      0.4      nd
0.7143      2.0      0.4      d
0.7143      2.0      0.4      d
0.3857      3.0      0.4      d
0.3857      3.0      0.4      d
0.8571      4.0      0.4      d
0.8571      4.0      0.4      d
0.4286      2.0      0.4      d
0.4286      2.0      0.4      d
0.7143      3.0      0.4      d
0.7143      3.0      0.4      d
0.8571      5.0      0.4      d
0.8571      5.0      0.4      d
0.5      2.0      0.4      d
0.5      2.0      0.4      d
0.4286      1.0      0.4      d
  
```

Status

OK Log

Preprocess Classify Cluster Associate Select attributes Visualize Forecast

Classifier

Choose MultiObjectiveEvolutionaryFuzzyClassifier -generations 20 -populationSize 100 -seed 1 -maxSimilarity 0.4 -minV 30.0 -maxV 2.0 -maxRules -1 -maxLabels 5 -evaluationMeasure 0 -algorithm 0 -reportFrequency 20 -logFile /Users/manimkey -num-decimal-pla

Test options

☐ Use training set

☐ Supplied test set

☒ Cross-validation Folds 10

☐ Percentage split % 60

More options...

(Print) play

Start Stop

Result list (right-click for options)

21:34:37 - rules.MultiObjectiveEvolutionaryFuzzyClassifier

Classifier output

Objectives:
acc = 0.8571
number of rules = 3.0
Constraint:
similarity = 6.4
Classifier:
RULE 1:
IF
outlook IS sunny
AND temperature IS null
AND humidity IS Low
AND windy IS FALSE
THEN play IS yes
RULE 2:
IF
outlook IS rainy
AND temperature IS High
AND humidity IS High
AND windy IS TRUE
THEN play IS yes
RULE 3:
IF
outlook IS sunny
AND temperature IS High
AND humidity IS High
AND windy IS TRUE
THEN play IS no

Time taken to build model: 0.2 seconds

=== Stratified cross-validation ===
=== Summary ===

Metric	Value	Percentage
Correctly Classified Instances	5	35.7143 %
Incorrectly Classified Instances	9	64.2857 %
Kappa statistic	-0.4651	
Mean absolute error	0.4429	
Root mean squared error	0.6658	
Relative absolute error	135 %	
Root relative squared error	163.5137 %	
Total Number of Instances	14	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC	AUC Area	PRC Area	Class
Weighted Avg.	0.357	0.643	0.321	0.357	0.338	-0.471	0.278	0.498	no

=== Confusion Matrix ===

a b == classified as

5 4 | a = yes

5 0 | b = no

Status

OK Log x0

● Forecast Analysis

Preprocess Classify Cluster Associate Select attributes Visualize Forecast

Basic configuration Advanced configuration

Target Selection

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> passenger_numbers

Parameters

Number of time units to forecast 1

Time stamp Date

Periodicity <Detect automatically>

Skip list

Confidence intervals Level (%) 95

Perform evaluation

Start Stop Help

Result list

21:11:34 - LinearRegression

Output/Visualization

Output Train future pred.

=== Run information ===

Scheme:
LinearRegression -S 0 -R 1.0E-8 -num-decimal-places 4

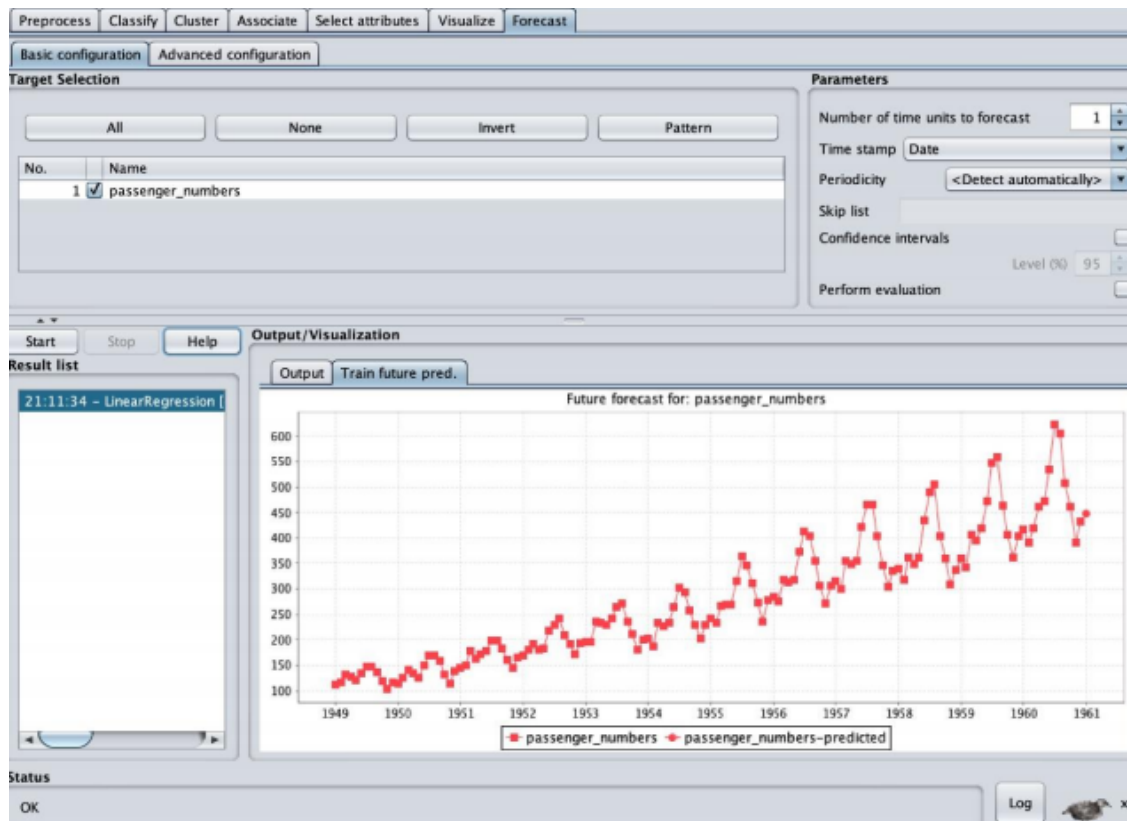
Lagged and derived variable options:
-F passenger_numbers -L 1 -M 12 -G Date -month -quarter

Relation: airline_passengers
Instances: 144
Attributes: 2
passenger_numbers
Date

Transformed training data:

passenger_numbers
Month
Quarter
Date-remapped
lag_passenger_numbers-1

Status



Findings and Learnings :

1. What is ROC Curve and how ROC analysis is useful.
2. What is survival analysis and where it is used.
3. Various ways to perform forecast analysis and its importance.