

IDA: Missing Values :-

→ library(mice), # multivariate imputation by chained equation,

→ . md.pattern()

→ md.pattern(),

* library(VIM) # package for visualization and imputation of missing values.

* aggr()

* marginplot(),

* scattermatrixmiss(),

Hmisc:- * cor()

* prcomp()

complot(package) → package is good.

PLA

- library(scatterplot3d)
- library(ogl).
- library(ggbiplot).
- biplot()
- ggbiplot(),
- scatterplot3d()
- plot3d()

lda and tsne.

LDA in MASS package.

- ~~library~~(MASS)
- lda().
- lda, model %, %, predict(transformed),
- lda needs data scaled and centered.

- `Idahist()`.

- `library("Rtsne")`.

TSNE

- `Rtsne()`.

- Grubbs's test for outliers,

- `library(outliers)`

- `grubbs.test()`, - univariate test.

- `identify()`.

- `outlier()`.

Assignment 4:-

Adaptive predictor modelling Textbook:

* ~~skewness()~~ in ~~e1071 package~~.

Chapter 3:-

Transformations

* skewness() in e1071 package.

* calculates skewness of a predictor.

* Another function, preProcess applies this transformation to a set of predictors.

• spatialSign().

• impute.knn() in impute package

, Preprocess() in lars. package.

(840x_transformation-l.R :-)
* Symbol function in "car" package.

- produces ~~out~~ parallel box plots for different values of lambda.
- Helpful in Transformation.

* boxcox in "EnvStats"

In HW4, My solution

* I used,
* BoxcoxTrans() from caret package
to find the required transformation.

HWS: (Imputations)

Mile:-

• Multivariate Imputation by Chained equations
is used widely for handling missing data.

*.

* How mice works :-

1) Initialization:-

Missing values are initially replaced with simple methods. (e.g. mean imputation) to create a starting point.

2. Chained Equations:-

For each variable with missing data, the missing values are imputed using a regression model that is a fitted based other variables in the data.

3. Iterative process:-

This process is repeated in an iterative manner, cycling through each variable until convergence is achieved.

Convergence occurs when distribution of imputed values stabilizes meaning that further iterations do not significantly change the ~~data~~ results.

H. Multiple Imputations:

This iterative process is repeated m times to create multiple different imputed datasets. Each dataset represent a plausible set of values that could have been observed.

Pooling Analysis:

- After creating multiple datasets, analysis are performed on each dataset separately. (like fitting a regression model on each imputed dataset).
- The results from each dataset are then pooled using Robin's rules to produce a single set of results.
- Pooling involves combining estimates (e.g. means, regression coefficients) and calculating std errors, that account for both within-imputation variability.

(variability within each imputed dataset) and between-imputation variability, (variability between different imputed datasets).

- This provides final parameter estimates and std. errors that incorporate the uncertainty due to missing data.

How to choose m?

Common Recommendations:-

Missing data proportion.	m.
<10%.	5
10-30%.	10-20.
>30%.	30-50.
>50%.	>50
>70-80%.	drop.
>90	obviously drop.

• Rule of thumb from Robin (1987) :-

$$m \approx \frac{100}{F}$$

F is fraction of missing data.

$$F = \frac{\text{Number of missing values}}{\text{Total Number of observations}}$$

Some Important Functions:-

- mice()
- complete.cases().
- with()
- pool()
- md.pattern()
- stripplot()
- density.plot()
- xyplot()
- init()
- complete()
- set.seed()
- md.pairs()

VIM (Visualization and Imputation methods)

- aggr() - Aggregation plot.
- matrixplot()
- marginplot()
- scattMiss()
- histMiss()

Imputation methods:-

- kNN(),
- hotdeck(),
- irnni() - Iterated Robust Model-based Imputation,
- regressionImp() - Regression Imputation,
- robustRegressionImputation(),

Unit 6:- Caret package,

~~(constraint)~~
~~RandomForest~~

(1) train() :-

parameters :- method , tolcontrol ,
tunefreqd , tunelength.

(2) trainControl() :-

parameters :- method , number ,
repeats .

(3) expand.grid() :- parameters changes
according to algorithms.

(4) preProcess() :- Apply transformations like.
scaling , centering , Missing
values , skewness , pca etc.

(5) Predict() :- Uses trained model to generate predictions on new data.

(6) fit\$resample :- Assess the resampling results.

(7) fit\$finalModel :- Assess the underlying model object trained by caret.

Unit 7: Multiple linear Regression

① Caret library:

- `VIF()` - Compute variance inflation factor for a linear model.

② MASS library:

- `lm.ridge` - perform ridge.
- `svd()` - perform singular value decomposition.
- `qplot()` - quick plot function for

Creating simple plots.

③ lars: library :-

- lars() - perform least angle regression or lasso.
- predict.lars()
- cv.lars().

④ glmnet library:-

- cv.glmnet().

⑤ pls library() :-

pls():- perform partial least squares regression.

RMSEP():- compute rmse of prediction for a ~~fit~~ model.

⑥ Stepwise Regression

- `step()`

⑦ earth package:

`earth()`, ~~plot()~~ `plotmo()`.

⑧ rgl package:

- `plot3d()`.
- ~~plot~~ `flares3d()`
- `surface3d()`.

⑨ Robust regression:- (MASS library).

- `rlm()` - fits a robust regression using M estimation.
- `psi.bisquare()` - function specifying the bisquare loss function for robust regression.

Unit 8: Logistic regression :-

- ~~glm()~~ oglm() function used with family specified "binomial".
- glm() is in base R package.

Unit 9:- decision trees:-

~~Decision tree or classifier - CART~~

rpart library :- optional

- rpart() - Builds decision trees using CART algorithm.

- plotcp() - plots the cost complexity pruning path.

- prune() - prunes a tree based on complexity parameter (cp).

partykit library:-

- as.party() - converts a part object to "party" object for visualization.
 - plot() - plots a party tree.
-

caret library:-

- confusionMatrix() - computes confusion matrix & performance metrics
 - varImp() - computes variable importance of models.
 -
-

rattle library :-

- LancyRpartPlot() - plots decision trees with enhanced visual representation.
-

adabag library:

- boosting() - implements boosting algorithms for classification.

ipred library :-

- bagging () - Implements bagging for classification.
- errorest () - computes error estimates using cross validation.

randomForest ~~library~~ library :-

- randomForest () .
- varImpPlot () .
- importance () .

Unit 10 :- Support Vectors

~~delopt1~~ Delopt1 library :-

- svm () - fits SVM model .
- sigest () - Estimates Sigma parameter for radial basis kernels .

Kernlab library :-

- By default caret uses kernlab for SVM.

Unit 11:- Neural networks:-

Unit 12:- Clustering:-

① stats (Base R stats package) :-

- kmeans() - performs k-Means clustering
- princomp() - An alternative function for PCA.

② cluster library :-

- pam() - partitioning around Medoids clustering
- silhouette() - compute silhouette width for clustering evaluation.

- daisy() - computes dissimilarity matrix for mixed data (Numerical and Factor data).

③ Useful library

- fitkMeans() - Fit k-Means clustering with evaluations (eg:-
~~Hartigan's rule~~ Hartigan's rule).

- plotHartigan() - Visualize Hartigan's rule for cluster evaluation.

④ NbClust library

- NbClust() - identifies the optimal number of clusters using multiple indices.

⑤ rgl library

- plot3d() - creates interactive 3d plots.

⑥ Hmisc library:-

- varclus() - performs hierarchical clustering of variables.

⑦ dbSCAN library:-

- dbScan() - perform density-based spatial clustering.
- kNNdist() - Computes the k nearest neighbour distance.
- kNNdistplot() - Visualizes k-nearest neighbour distance to find the "knee" (optimal eps).
- optics() - performs OPTICS clustering.
- extractDBSCAN() - Extracts DBSCAN clustering from OPTICS results.

`hullplot()` - Visualizes clustering results with convex hulls.

Base R:

- `dist()` - computes distance matrix.
- `hclust()` - performs hierarchical clustering
- `rect.hclust()` - draws rectangles around hierarchical clustering.
- `cutree()` - cuts a dendrogram tree into specified clusters.

and the best and the most difficult
and the best and the most difficult
and the best and the most difficult
and the best and the most difficult

and the most difficult

Unit 12:- Clustering

1. k-Means Clustering :-

- kMeans()

- fitkMeans(): Helps identify optimal clusters using Hartigan's rule,

- wssplot(data, n) - custom function to plot elbow method.

- plotHartigan()

2. Partitioning Around Medoids :-

- pam()

- silhouette()

3. Hierarchical Clustering :-

- hclust()

- daisy()

- rect.hclust()