



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

MULTIPLE LINEAR REGRESSION-PART V

EXHAUSTIVE SEARCH

LECTURE 26

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



MULTIPLE LINEAR REGRESSION

- Steps to reduce the no. of predictors
 - Summary statistics and graphs
 - Statistical methods using computational power
 - Exhaustive search: all possible combinations
 - Partial-iterative search: algorithm based
- Exhaustive Search
 - Large no. of subsets
 - Criteria to compare models
 - Adjusted R²



MULTIPLE LINEAR REGRESSION

- Adjusted R²

$$R^2_{adj} = 1 - \frac{n-1}{n-p-1} (1 - R^2)$$

Where R² is proportion of explained variability in the model

$$R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

- R² is called coefficient of determination



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE

MULTIPLE LINEAR REGRESSION

- R^2 would be equal to squared correlation in a single predictor model, that is how R^2 gets its name
- Adjusted R^2 introduces a penalty on the no. of predictors to trade-off between artificial increase vs. amount of information
- High adjusted R^2 values \rightarrow low $\hat{\sigma}^2$



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE

MULTIPLE LINEAR REGRESSION

- Exhaustive Search
 - Criteria to compare models
 - Mallow's C_p
- Mallow's C_p

$$C_p = \frac{SSR}{\hat{\sigma}_f^2} + 2(p + 1) - n$$

Where $\hat{\sigma}_f^2$ is estimated value of σ^2 in the full model

$$\text{and } SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})$$



MULTIPLE LINEAR REGRESSION

- Mallow's C_p
 - Assumption: full model with all predictors is unbiased
 - Predictors elimination would reduce the variability
 - Best subset model would have $C_p \sim p+1$ and p would be a small value
 - Requires high n value for the training partition relative to p
- Open RStudio



MULTIPLE LINEAR REGRESSION

- Partial-iterative search
 - Computationally cheaper
 - Best subset is not guaranteed
 - Potential of missing “good” sets of predictors
 - Produce close-to-best subsets
 - Preferred approach for large no. of predictors
 - For moderate no. of predictors, exhaustive search is better
- Trade-off between computation cost vs. potential of finding best subset



MULTIPLE LINEAR REGRESSION

- Partial-iterative search algorithms
 - Forward selection
 - Add predictors one by one
 - Strength as a single predictor is used
 - Backward elimination
 - Drop predictors one by one
 - Stepwise regression
 - Add predictors one by one and consider dropping insignificant ones
- Open RStudio



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

MULTIPLE LINEAR REGRESSION-PART VI

Partial Iterative Search

LECTURE 27

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



MULTIPLE LINEAR REGRESSION

- Partial-iterative search
 - Computationally cheaper
 - Best subset is not guaranteed
 - Potential of missing “good” sets of predictors
 - Produce close-to-best subsets
 - Preferred approach for large no. of predictors
 - For moderate no. of predictors, exhaustive search is better
- Trade-off between computation cost vs. potential of finding best subset



MULTIPLE LINEAR REGRESSION

- Partial-iterative search algorithms
 - Forward selection
 - Add predictors one by one
 - Strength as a single predictor is used
 - Backward elimination
 - Drop predictors one by one
 - Stepwise regression
 - Add predictors one by one and consider dropping insignificant ones
- Open RStudio



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

MACHINE LEARNING TECHNIQUE k-NEAREST NEIGHBORS (k-NN)

LECTURE 28

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



k-NEAREST NEIGHBORS (k-NN)

- k-NN
 - No assumptions about the form of relationship between outcome variable and the set of predictors
 - Non-parametric method
 - No parameters from the assumed functional form to estimate
 - Useful information for modeling is extracted using the similarities between the records based on predictors' values
 - Typically, distance based similarity measures are used



k-NEAREST NEIGHBORS (k-NN)

- k-NN: distance metrics
 - Most popular metric is Euclidean distance
For two records having values of the predictors denoted by (x_1, x_2, \dots, x_p) and (w_1, w_2, \dots, w_p)
$$D_{Eu} = \sqrt{(x_1 - w_1)^2 + (x_2 - w_2)^2 + \dots + (x_p - w_p)^2}$$
 - Low computation costs
 - Other distance metrics: statistical distance or Mahalanobis distance and Manhattan distance
 - Euclidean distance is preferred in k-NN due to many distance computations



k-NEAREST NEIGHBORS (k-NN)

- k-NN
 - Scaling of predictors: standardized values of predictors
- k-NN for Classification task
 - Main idea is to find k records in the training partition which are neighboring the new observation to be classified
 - These k neighbors are used to classify the new observation into a class
 - Predominant class among the neighbors



k-NEAREST NEIGHBORS (k-NN)

- k-NN: Finding neighbors and Classification
 - Compute the distance between the new observation and training partition records
 - Determine k nearest or closest records to the new observation
 - Find most prevalent class among k neighbors and it would be the predicted class of new observation
- Open RStudio



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

MACHINE LEARNING TECHNIQUE k-NEAREST NEIGHBORS (k-NN)- Part 2

LECTURE 29

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



k-NEAREST NEIGHBORS (k-NN)

- k-NN
 - Choosing appropriate value of k
 - k=1: powerful for large no. of records in training partition
 - k>1: smoothing effects (control overfitting issues)
 - Low value of k -> more likely to fit the noise
 - High value of k -> more likely to ignore the local patterns in the data
 - Trade-off between benefits from local pattern vs global effects
 - k=n: naïve rule



k-NEAREST NEIGHBORS (k-NN)

- k-NN
 - Value of k: depends on nature of the data as well
 - Low value of k for data with complex and irregular structures
 - Typical value of k: between ‘1-20’
 - Odd value of k is preferred to avoid ties in majority class decisions
- Best value of k
 - Classification performance on validation partition
- Open RStudio



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

MACHINE LEARNING TECHNIQUE k-NEAREST NEIGHBORS (k-NN) PART 3

LECTURE 30

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



k-NEAREST NEIGHBORS (k-NN)

- k-NN: Finding neighbors and Classification
 - Compute the distance between the new observation and training partition records
 - Determine k nearest or closest records to the new observation
 - Find most prevalent class among k neighbors and it would be the predicted class of new observation
- Open RStudio



k-NEAREST NEIGHBORS (k-NN)

- k-NN
 - Choosing appropriate value of k
 - k=1: powerful for large no. of records in training partition
 - k>1: smoothing effects (control overfitting issues)
 - Low value of k -> more likely to fit the noise
 - High value of k -> more likely to ignore the local patterns in the data
 - Trade-off between benefits from local pattern vs global effects
 - k=n: naïve rule



k-NEAREST NEIGHBORS (k-NN)

- k-NN
 - Value of k: depends on nature of the data as well
 - Low value of k for data with complex and irregular structures
 - Typical value of k: between ‘1-20’
 - Odd value of k is preferred to avoid ties in majority class decisions
- Best value of k
 - Classification performance on validation partition
- Open RStudio



k-NEAREST NEIGHBORS (k-NN)

- Majority decision rule vs. cutoff probability
 - Two class scenario: majority rule \equiv cutoff value of 0.5
- k-NN for multi-class scenario
- Class of interest
 - Instead of the majority rule, compare proportion of k neighbors belonging to class of interest to a user-specified cut off value



k-NEAREST NEIGHBORS (k-NN)

- k-NN for Prediction task
 - Main idea is to find k records in the training partition which are neighboring the new observation to be predicted
 - These k neighbors are used to predict the value of new observation
 - Average value of the outcome variable among the neighbors
 - Weighted average wherein weight for a neighbor decreases as its distance from new observation increases
 - Performance metric: RMSE or some other prediction error metric



k-NEAREST NEIGHBORS (k-NN)

- k-NN: Finding neighbors and Prediction
 - Compute the distance between the new observation and training partition records
 - Determine k nearest or closest records to the new observation
 - Compute the average or weighted average of outcome variable values among k neighbors and it would be the predicted value of new observation



k-NEAREST NEIGHBORS (k-NN)

- Further Comments on k-NN algorithm
 - Computation time to find nearest neighbors for large training partition
 - Dimension reduction techniques
 - Steps to find neighbors can be optimized using efficient data structures for search operations like trees
 - Identification and pruning of redundant records from training partition which will not be included in neighbor search steps
 - Curse of dimensionality
 - Sample size requirement depends on no. of predictors
 - Leads to more computations for neighbors



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

NAÏVE BAYES

LECTURE 31

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



NAÏVE BAYES

- Complete or Exact Bayes for classification
 - Search for records in training partition having same predictors' values as the new observation to be classified
 - Find the most prevalent class of the outcome variable among the records
 - Assign this class to the new observation
- Class of interest
 - User specified cut off value for the class of interest



NAÏVE BAYES

- Class of interest
 - Search for records in training partition having same predictors' values as the new observation to be classified
 - Find the probability of a record belonging to the class of interest among the records
 - If computed probability value > cut off value, assign the new observation to the class of interest



NAÏVE BAYES

- Concept of conditional probability
 - For an outcome variable with m classes $\{C_1, C_2, \dots, C_m\}$ and p predictors $\{x_1, x_2, \dots, x_p\}$, we are interested in the following probability value:

$$P(C_i | x_1, x_2, \dots, x_p) = \frac{P(x_1, x_2, \dots, x_p | C_i)P(C_i)}{P(x_1, x_2, \dots, x_p | C_1)P(C_1) + \dots + P(x_1, x_2, \dots, x_p | C_m)P(C_m)}$$

- Assign the new observation to the class with highest probability value
- Or, if the probability value for the class of interest > cut off value for the same, assign the new observation to the class of interest



NAÏVE BAYES

- Bayes Model for classification
 - Predictors should also be categorical
 - Numerical variables will have to be converted into categorical variables through binning
- Open Excel
- Complete or Exact Bayes Limitations
 - For a model even with small no. of predictors, many new observations to be classified might not get exact matches
 - Probability of a match might reduce significantly on adding just one variable to the set of predictors



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

NAÏVE BAYES PART-2

LECTURE 32

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



NAÏVE BAYES

- Bayes Model for classification
 - Predictors should also be categorical
 - Numerical variables will have to be converted into categorical variables through binning
- Open Excel
- Complete or Exact Bayes Limitations
 - For a model even with small no. of predictors, many new observations to be classified might not get exact matches
 - Probability of a match might reduce significantly on adding just one variable to the set of predictors



NAÏVE BAYES

- Instead of Complete or Exact Bayes, switch to Naïve Bayes
 - In Naïve Bayes, all the records are used instead of relying on just the matching records
- Naïve Bayes Modification
 - For class i of outcome variable, compute the probabilities (P_1, P_2, \dots, P_p) of belonging to class i for each predictor's value (x_1, x_2, \dots, x_p) taken by the new observation to be classified
 - Compute $P_1 \times P_2 \times \dots \times P_p \times P(C_i)$
 - Execute previous two steps for all the classes



NAÏVE BAYES

- Naïve Bayes Modification
 - To compute the probability of the new observation belonging to class i, divide the value computed in step 2 by the summation of values computed in step 2 for all the classes
 - Execute previous step for all the classes
 - Classify the new observation to the class with the highest probability value



NAÏVE BAYES

- Naïve Bayes formula

$$P(C_i | x_1, x_2, \dots, x_p) = \frac{[P(x_1 | C_i)P(x_2 | C_i) \dots P(x_p | C_i)]P(C_i)}{[P(x_1 | C_1)P(x_2 | C_1) \dots P(x_p | C_1)]P(C_1) + \dots + [P(x_1 | C_m)P(x_2 | C_m) \dots P(x_p | C_m)]P(C_m)}$$

- Naïve Bayes formula is directly derived from the exact Bayes formula after making following assumption:
- Predictors' values $\{x_1, x_2, \dots, x_p\}$ occur independent of each other for a given class

$$P(x_1, x_2, \dots, x_p | C_i) \equiv P(x_1 | C_i)P(x_2 | C_i) \dots P(x_p | C_i)$$



NAÏVE BAYES

- Naïve Bayes formula
 - For classification, naïve Bayes formula works quite well
 - Since we don't require probabilities values to be accurate in absolute term, rather just a reasonably accurate rank ordering of these values
 - For the same reason, we should use the numerator only and drop the denominator which is common for all the classes
- Steps when we have a class of interest
 - User specified cut off value for the class of interest



NAÏVE BAYES

- Steps when we have a class of interest
 - Compute the probabilities (P_1, P_2, \dots, P_p) of belonging to class of interest for each predictor's value (x_1, x_2, \dots, x_p) taken by the new observation to be classified
 - Compute $P_1 \times P_2 \times \dots \times P_p \times P(\text{Class of interest})$
 - Execute previous two steps for all the classes
 - To compute the probability of the new observation belonging to class of interest, divide the value computed in step 2 by the summation of values computed in step 2 for all the classes



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

NAÏVE BAYES PART-3

LECTURE 33

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



NAÏVE BAYES

- Bayes Model for classification
 - Predictors should also be categorical
 - Numerical variables will have to be converted into categorical variables through binning
- Open Excel
- Complete or Exact Bayes Limitations
 - For a model even with small no. of predictors, many new observations to be classified might not get exact matches
 - Probability of a match might reduce significantly on adding just one variable to the set of predictors



NAÏVE BAYES

- Instead of Complete or Exact Bayes, switch to Naïve Bayes
 - In Naïve Bayes, all the records are used instead of relying on just the matching records
- Naïve Bayes Modification
 - For class i of outcome variable, compute the probabilities (P_1, P_2, \dots, P_p) of belonging to class i for each predictor's value (x_1, x_2, \dots, x_p) taken by the new observation to be classified
 - Compute $P_1 \times P_2 \times \dots \times P_p \times P(C_i)$
 - Execute previous two steps for all the classes



NAÏVE BAYES

- Naïve Bayes Modification
 - To compute the probability of the new observation belonging to class i, divide the value computed in step 2 by the summation of values computed in step 2 for all the classes
 - Execute previous step for all the classes
 - Classify the new observation to the class with the highest probability value



NAÏVE BAYES

- Naïve Bayes formula

$$P(C_i | x_1, x_2, \dots, x_p) = \frac{[P(x_1 | C_i)P(x_2 | C_i) \dots P(x_p | C_i)]P(C_i)}{[P(x_1 | C_1)P(x_2 | C_1) \dots P(x_p | C_1)]P(C_1) + \dots + [P(x_1 | C_m)P(x_2 | C_m) \dots P(x_p | C_m)]P(C_m)}$$

- Naïve Bayes formula is directly derived from the exact Bayes formula after making following assumption:
- Predictors' values $\{x_1, x_2, \dots, x_p\}$ occur independent of each other for a given class

$$P(x_1, x_2, \dots, x_p | C_i) \equiv P(x_1 | C_i)P(x_2 | C_i) \dots P(x_p | C_i)$$



NAÏVE BAYES

- Naïve Bayes formula
 - For classification, naïve Bayes formula works quite well
 - Since we don't require probabilities values to be accurate in absolute term, rather just a reasonably accurate rank ordering of these values
 - For the same reason, we should use the numerator only and drop the denominator which is common for all the classes
- Steps when we have a class of interest
 - User specified cut off value for the class of interest



NAÏVE BAYES

- Steps when we have a class of interest
 - Compute the probabilities (P_1, P_2, \dots, P_p) of belonging to class of interest for each predictor's value (x_1, x_2, \dots, x_p) taken by the new observation to be classified
 - Compute $P_1 \times P_2 \times \dots \times P_p \times P(\text{Class of interest})$
 - Execute previous two steps for all the classes
 - To compute the probability of the new observation belonging to class of interest, divide the value computed in step 2 by the summation of values computed in step 2 for all the classes



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

NAÏVE BAYES PART-4

LECTURE 34

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



NAÏVE BAYES

- Naïve Bayes formula
 - For classification, naïve Bayes formula works quite well
 - Since we don't require probabilities values to be accurate in absolute term, rather just a reasonably accurate rank ordering of these values
 - For the same reason, we should use the numerator only and drop the denominator which is common for all the classes
- Steps when we have a class of interest
 - User specified cut off value for the class of interest



NAÏVE BAYES

- Steps when we have a class of interest
 - Compute the probabilities (P_1, P_2, \dots, P_p) of belonging to class of interest for each predictor's value (x_1, x_2, \dots, x_p) taken by the new observation to be classified
 - Compute $P_1 \times P_2 \times \dots \times P_p \times P(\text{Class of interest})$
 - Execute previous two steps for all the classes
 - To compute the probability of the new observation belonging to class of interest, divide the value computed in step 2 by the summation of values computed in step 2 for all the classes



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

NAÏVE BAYES PART-5

LECTURE 35

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



NAÏVE BAYES

- Further Comments on Naïve Bayes
 - Good performance despite assumption of independent predictors' values being far from true
 - Requires large no. of records
 - Few classes of predictors might not be represented in the training partition records
 - Zero probability is assumed
 - Good for classification but not for estimating probabilities of class membership



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

CLASSIFICATION & REGRESSION TREES

LECTURE 36

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- CART
 - A data-driven method
 - Based on separating observations into homogeneous subgroups by creating splits on predictors
 - Used for both prediction and classification tasks
 - Model is represented by a tree diagram
 - Easy to interpret logical rules
 - CART algorithm grows binary trees
 - Adoption across domains



CLASSIFICATION & REGRESSION TREES

- Classification Trees
 - Recursive partitioning
 - About partitioning p -dimensional space of predictors using training partition, where p is no. of predictors
 - Pruning
 - About pruning the built tree using validation data



CLASSIFICATION & REGRESSION TREES

- Recursive Partitioning
 - Partitioning p-dimensional space of predictors into non-overlapping multi-dimensional rectangles
 - The partitioning process is recursive in nature
 - Applied on the results of previous partitions
- Steps for Recursive Partitioning
 - An optimal combination of one of the predictors, x_i and its value v_i is selected to create first split of p-dimensional space into two parts
 - Part I: $x_i \leq v_i$
 - Part II: $x_i > v_i$



CLASSIFICATION & REGRESSION TREES

- Steps for Recursive Partitioning
 - Step 1 is applied again on the two parts and process continues to create more rectangular parts
 - The partitioning process continues till we reach pure homogeneous parts
 - All the observations in the part belong to just one of the classes
- Open RStudio



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

CLASSIFICATION & REGRESSION TREES PART-2

LECTURE 37

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- Impurity Measures

- Gini index and Entropy measure

- Gini Index

For an outcome variable with m classes, Gini impurity index for a rectangular part is defined as

$$gini = 1 - \sum_{k=1}^m P_k^2$$

Where P_k is the proportion of rectangular part observations belonging to class k



CLASSIFICATION & REGRESSION TREES

- Gini Index
 - Gini values lie in the range $\{0, (m-1)/m\}$ for m-class scenario and $\{0, 0.5\}$ for two-class scenario
- Entropy Measure

For an outcome variable with m classes, entropy for a rectangular part is defined as

$$\text{Entropy} = - \sum_{k=1}^m P_k \log_2(P_k)$$



CLASSIFICATION & REGRESSION TREES

- Entropy Measure
 - Entropy values lie in the range $\{0, \log_2(m)\}$ for m-class scenario and $\{0, 1\}$ for two-class scenario
- Open RStudio
- Tree diagram or tree structure
 - Each split of p-dimensional space into two parts can be depicted as a split of a node in a decision tree into two child nodes
 - First split creates branches of root node



CLASSIFICATION & REGRESSION TREES

- Two types of nodes in tree structure
 - Decision node: Depicted with a circle
 - Terminal or leaf node: Depicted with a rectangle
 - Correspond to Final rectangular parts
- Steps to classify a new observation using tree based models
 - New observation to be classified is dropped down the tree starting from root node
 - At each decision node, the appropriate branch is taken until we reach a leaf node



CLASSIFICATION & REGRESSION TREES

- Steps to classify a new observation using tree based models
 - At leaf node, majority class is assigned to the new observation
 - For class of interest scenario, proportion of records belonging to the class of interest is compared with the user specified cut off value for the same
- Open RStudio



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

CLASSIFICATION & REGRESSION TREES PART-3

LECTURE 38

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- Two types of nodes in tree structure
 - Decision node: Depicted with a circle
 - Terminal or leaf node: Depicted with a rectangle
 - Correspond to Final rectangular parts
- Steps to classify a new observation using tree based models
 - New observation to be classified is dropped down the tree starting from root node
 - At each decision node, the appropriate branch is taken until we reach a leaf node



CLASSIFICATION & REGRESSION TREES

- Steps to classify a new observation using tree based models
 - At leaf node, majority class is assigned to the new observation
 - For class of interest scenario, proportion of records belonging to the class of interest is compared with the user specified cut off value for the same
- Open RStudio



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE

Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

CLASSIFICATION & REGRESSION TREES PART-4

LECTURE 39

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- CART
 - A data-driven method
 - Based on separating observations into homogeneous subgroups by creating splits on predictors
 - Used for both prediction and classification tasks
 - Model is represented by a tree diagram
 - Easy to interpret logical rules
 - CART algorithm grows binary trees
 - Adoption across domains



CLASSIFICATION & REGRESSION TREES

- Classification Trees
 - Recursive partitioning
 - About partitioning p -dimensional space of predictors using training partition, where p is no. of predictors
 - Pruning
 - About pruning the built tree using validation data



CLASSIFICATION & REGRESSION TREES

- CART example has been discussed in the lecture video



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE

CLASSIFICATION & REGRESSION TREES

- Further Comments on CART
 - Can be used as a variable selection approach
 - No variable transformation is required
 - Robust to outliers
 - Non-linear and non-parametric technique
 - Handle missing values
 - Sensitive to sample data changes
 - Predictor's strength as a single variable is modeled and not as part of a group of predictors



CLASSIFICATION & REGRESSION TREES

- Further Comments on CART
 - Might not fit linear structures or relationships between predictors
 - New predictors based on hypothesized relationships can be used
 - Require a large dataset
 - High computation time



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

CLASSIFICATION & REGRESSION TREES PART-5

LECTURE 40

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- CART example has been discussed in the lecture video



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE

CLASSIFICATION & REGRESSION TREES

- Further Comments on CART
 - Can be used as a variable selection approach
 - No variable transformation is required
 - Robust to outliers
 - Non-linear and non-parametric technique
 - Handle missing values
 - Sensitive to sample data changes
 - Predictor's strength as a single variable is modeled and not as part of a group of predictors



CLASSIFICATION & REGRESSION TREES

- Further Comments on CART
 - Might not fit linear structures or relationships between predictors
 - New predictors based on hypothesized relationships can be used
 - Require a large dataset
 - High computation time



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

CLASSIFICATION & REGRESSION TREES PART-6

LECTURE 41

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- Steps to classify a new observation using tree based models
 - At leaf node, majority class is assigned to the new observation
 - For class of interest scenario, proportion of records belonging to the class of interest is compared with the user specified cut off value for the same
- Open RStudio



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE

CLASSIFICATION & REGRESSION TREES

- Pruning
 - Avoid overfitting
 - Full grown tree leads to complete overfitting of data
 - Poor performance on new data
 - Overall error of tree models
 - Expected to decrease until the point where relationships between outcome variable and predictors are fitted
 - Then tree models start fitting to the noise and overall error starts increasing
 - Due to splits involving small number of observations



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

Pruning Process

LECTURE 42

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Avoid overfitting
 - Full grown tree leads to complete overfitting of data
 - Poor performance on new data
 - Overall error of tree models
 - Expected to decrease until the point where relationships between outcome variable and predictors are fitted
 - Then tree models start fitting to the noise and overall error starts increasing
 - Due to splits involving small number of observations



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Stop tree growth before it starts overfitting data or fitting noise
 - No. of splits or tree depth level
 - No. of observations in a node to attempt the split
 - Accepted level of reduction in impurity
 - Difficulties in determining the stopping point for such rules
 - Prune the full grown tree back to a level where it doesn't overfit data or fit noise
 - Use validation partition to prune the tree modeled with training partition
 - Idea is to remove the tree branches which don't reduce the error rate further



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Prune the full grown tree back to a level where it doesn't overfit data or fit noise
 - Find the point where error rate on validation partition starts to increase
 - Cost complexity parameter or complexity parameter (CP) in CART algorithm
$$CP = Err + PF * TL$$
Where Err is misclassification error, PF is penalty factor for tree length (TL)
 - Minimum error tree
 - Tree with minimum misclassification error on validation partition



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Best pruned tree
 - Adjustment for sampling error on minimum error tree
 - Smallest tree in the pruning sequence which lies within one std. err. (of error rate) of minimum error tree
- Open RStudio
- Classification Rules
 - Each terminal node in a tree model is equivalent to a classification rule
 - Simplify and remove redundant rules



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

Pruning Process Part-2

LECTURE-43

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Avoid overfitting
 - Full grown tree leads to complete overfitting of data
 - Poor performance on new data
 - Overall error of tree models
 - Expected to decrease until the point where relationships between outcome variable and predictors are fitted
 - Then tree models start fitting to the noise and overall error starts increasing
 - Due to splits involving small number of observations



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Stop tree growth before it starts overfitting data or fitting noise
 - No. of splits or tree depth level
 - No. of observations in a node to attempt the split
 - Accepted level of reduction in impurity
 - Difficulties in determining the stopping point for such rules
 - Prune the full grown tree back to a level where it doesn't overfit data or fit noise
 - Use validation partition to prune the tree modeled with training partition
 - Idea is to remove the tree branches which don't reduce the error rate further



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Prune the full grown tree back to a level where it doesn't overfit data or fit noise
 - Find the point where error rate on validation partition starts to increase
 - Cost complexity parameter or complexity parameter (CP) in CART algorithm
$$CP = Err + PF * TL$$
Where Err is misclassification error, PF is penalty factor for tree length (TL)
 - Minimum error tree
 - Tree with minimum misclassification error on validation partition



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Best pruned tree
 - Adjustment for sampling error on minimum error tree
 - Smallest tree in the pruning sequence which lies within one std. err. (of error rate) of minimum error tree
- Open RStudio
- Classification Rules
 - Each terminal node in a tree model is equivalent to a classification rule
 - Simplify and remove redundant rules



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

Pruning Process Part-3

LECTURE-44

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Avoid overfitting
 - Full grown tree leads to complete overfitting of data
 - Poor performance on new data
 - Overall error of tree models
 - Expected to decrease until the point where relationships between outcome variable and predictors are fitted
 - Then tree models start fitting to the noise and overall error starts increasing
 - Due to splits involving small number of observations



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Stop tree growth before it starts overfitting data or fitting noise
 - No. of splits or tree depth level
 - No. of observations in a node to attempt the split
 - Accepted level of reduction in impurity
 - Difficulties in determining the stopping point for such rules
 - Prune the full grown tree back to a level where it doesn't overfit data or fit noise
 - Use validation partition to prune the tree modeled with training partition
 - Idea is to remove the tree branches which don't reduce the error rate further



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Prune the full grown tree back to a level where it doesn't overfit data or fit noise
 - Find the point where error rate on validation partition starts to increase
 - Cost complexity parameter or complexity parameter (CP) in CART algorithm
$$CP = Err + PF * TL$$
Where Err is misclassification error, PF is penalty factor for tree length (TL)
 - Minimum error tree
 - Tree with minimum misclassification error on validation partition



CLASSIFICATION & REGRESSION TREES

- Pruning
 - Best pruned tree
 - Adjustment for sampling error on minimum error tree
 - Smallest tree in the pruning sequence which lies within one std. err. (of error rate) of minimum error tree
- Open RStudio
- Classification Rules
 - Each terminal node in a tree model is equivalent to a classification rule
 - Simplify and remove redundant rules



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

REGRESSION TREES

LECTURE 45

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



CLASSIFICATION & REGRESSION TREES

- Regression Trees
 - Outcome variable should be numerical
 - Steps to build tree model are similar to that of classification trees
 - Prediction step, impurity measures and performance metrics are different
- Prediction step
 - Value of a leaf node is predicted value for a new observation that fell in that leaf node
 - Value of a leaf node is computed by taking average of training partition records constituting that leaf node



CLASSIFICATION & REGRESSION TREES

- Impurity Measures
 - Sum of squared deviations from mean of leaf node
 - Equivalent to squared errors since mean value of leaf node is predicted value
 - Lowest impurity is zero when all the observations that fell in a leaf node have same actual value of outcome variable



CLASSIFICATION & REGRESSION TREES

- Further Comments on CART
 - Can be used as a variable selection approach
 - No variable transformation is required
 - Robust to outliers
 - Non-linear and non-parametric technique
 - Handle missing values
 - Sensitive to sample data changes
 - Predictor's strength as a single variable is modeled and not as part of a group of predictors



CLASSIFICATION & REGRESSION TREES

- Further Comments on CART
 - Might not fit linear structures or relationships between predictors
 - New predictors based on hypothesized relationships can be used
 - Require a large dataset
 - High computation time



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

LOGISTIC REGRESSION

LECTURE 46

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



LOGISTIC REGRESSION

- Equivalent of linear regression for categorical outcome variable
 - Predictors can be categorical or continuous
- Applied in following tasks
 - Classification task
 - Predicting the class of a new observation
 - Profiling
 - Understanding similarities and differences among groups



LOGISTIC REGRESSION

- Steps for logistic regression
 - Estimate probabilities of class memberships
 - Classify observations using probabilities values
 - Most probable class method: assign the observation to the class with highest probability value
 - Equivalently, for a two-class case, cutoff value of 0.5 can be used
 - Class of interest: user specified cutoff value
 - For a two-class case, typically a value greater than average probability value for class of interest, but less than 0.5 can be used



LOGISTIC REGRESSION

- Logistic Regression Model
 - Used typically in cases when structured model is preferred over data-driven models for classification tasks
 - Categorical outcome variable cannot be directly modeled as a linear function of predictors
 - Inability to apply various mathematical operators
 - Variable type mismatches
 - Range reasonability issues
 - LHS range={0, ..., m-1}
 - RHS range=(-∞, ∞)



LOGISTIC REGRESSION

- Logistic Regression Model
 - Instead of using outcome variable (Y) in the model, a function of Y , called *logit* is used
 - Logit
 - Think about modeling probability value as a linear function of predictors, specifically in a two-class case
- If P is the probability of class 1 membership

$$P = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p$$

Where p is the no. of predictors



LOGISTIC REGRESSION

- Logit
 - LHS range improves from {0, 1} to [0, 1], however still cannot match RHS
 - Can we bring RHS range to [0,1]?
 - Nonlinear approach
 - Typically, a nonlinear function of the following form is used to perform the required transformation

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p)}}$$

This function is called *logistic response function*



LOGISTIC REGRESSION

- Logit
 - Rearrange the previous equation as below:

$$\frac{P}{1 - P} = e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}$$

LHS is expression for *odds*, another measure of class membership

$$odds = \frac{P}{1 - P}$$

- Odds of belonging to a class is defined as ratio of probability of class 1 membership to probability of class 0 membership
 - This metric is popular in sports, horse racing, gambling, and many other areas



LOGISTIC REGRESSION

- Logit

- Previous equation can be rewritten as

$$odds = e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}$$

- Range is now $(0, \infty)$
 - Take log on both sides of previous equation

$$\log(odds) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

- Standard logistic model
 - Now, LHS and RHS both have same range $(-\infty, \infty)$
- Log(odds) is called logit
 - Logit is used as the outcome variable in the model instead of categorical Y



LOGISTIC REGRESSION

- Odds and logit can be written as a function of probability of class 1 membership
 - Open RStudio
- In logistic regression model, we predict the logit values and therefore corresponding probability of a categorical outcome
 - Predicted probabilities values become the basis for classification
 - A prediction model for classification task



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

LOGISTIC REGRESSION- PART 2

LECTURE 47

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



LOGISTIC REGRESSION

- Logit
 - LHS range improves from {0, 1} to [0, 1], however still cannot match RHS
 - Can we bring RHS range to [0,1]?
 - Nonlinear approach
 - Typically, a nonlinear function of the following form is used to perform the required transformation

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p)}}$$

This function is called *logistic response function*



LOGISTIC REGRESSION

- Logit
 - Rearrange the previous equation as below:

$$\frac{P}{1 - P} = e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}$$

LHS is expression for *odds*, another measure of class membership

$$odds = \frac{P}{1 - P}$$

- Odds of belonging to a class is defined as ratio of probability of class 1 membership to probability of class 0 membership
 - This metric is popular in sports, horse racing, gambling, and many other areas



LOGISTIC REGRESSION

- Logit

- Previous equation can be rewritten as

$$odds = e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}$$

- Range is now $(0, \infty)$
 - Take log on both sides of previous equation

$$\log(odds) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

- Standard logistic model
 - Now, LHS and RHS both have same range $(-\infty, \infty)$
- Log(odds) is called logit
 - Logit is used as the outcome variable in the model instead of categorical Y



LOGISTIC REGRESSION

- Odds and logit can be written as a function of probability of class 1 membership
 - Open RStudio
- In logistic regression model, we predict the logit values and therefore corresponding probability of a categorical outcome
 - Predicted probabilities values become the basis for classification
 - A prediction model for classification task



LOGISTIC REGRESSION

- Estimation Technique
 - Least squares method used in multiple linear regression cannot be used
 - Non-linear formulation of logistic regression
 - Maximum likelihood method is used
 - Estimates are optimized in order to maximize the likelihood of obtaining the observations used in training the model
 - Less robust than estimation techniques used in linear regression
 - Reliability of estimates
 - Outcome variable categories should have adequate proportion
 - Adequate sample size w.r.t no. of estimates



LOGISTIC REGRESSION

- Estimation Technique
 - Maximum likelihood method is used
 - Collinearity issues similar to linear regression
- Interpretation of Results
 - Logit model
 - Additive factor (β)
 - If $\beta < 0$, increase in $x \Rightarrow$ decrease in logit values
 - If $\beta > 0$, increase in $x \Rightarrow$ increase in logit values
 - For any value of x , interpretative statements of results are same



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

LOGISTIC REGRESSION PART-3

LECTURE 48

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



LOGISTIC REGRESSION

- Odds and odds ratios
 - Odds is a ratio of two probability values (prob. of class 1/prob. Of Class 0)
 - Odds ratio is ratio of two odds (odds of class m1/odds of class m2)
 - Odds ratio $> 1 \Rightarrow$ odds of class m1 are higher than class m2
- Open RStudio



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

LOGISTIC REGRESSION PART-4

LECTURE 49

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



LOGISTIC REGRESSION

- Odds and logit can be written as a function of probability of class 1 membership
 - Open RStudio
- In logistic regression model, we predict the logit values and therefore corresponding probability of a categorical outcome
 - Predicted probabilities values become the basis for classification
 - A prediction model for classification task

LOGISTIC REGRESSION

- Estimation Technique
 - Least squares method used in multiple linear regression cannot be used
 - Non-linear formulation of logistic regression
 - Maximum likelihood method is used
 - Estimates are optimized in order to maximize the likelihood of obtaining the observations used in training the model
 - Less robust than estimation techniques used in linear regression
 - Reliability of estimates
 - Outcome variable categories should have adequate proportion
 - Adequate sample size w.r.t no. of estimates



LOGISTIC REGRESSION

- Estimation Technique
 - Maximum likelihood method is used
 - Collinearity issues similar to linear regression
- Interpretation of Results
 - Logit model
 - Additive factor (β)
 - If $\beta < 0$, increase in $x \Rightarrow$ decrease in logit values
 - If $\beta > 0$, increase in $x \Rightarrow$ increase in logit values
 - For any value of x , interpretative statements of results are same



LOGISTIC REGRESSION

- Interpretation of Results
 - Odds model
 - Multiplicative factor (e^{β})
 - If $\beta < 0$, increase in $x \Rightarrow$ decrease in odds
 - If $\beta > 0$, increase in $x \Rightarrow$ increase in odds
 - For any value of x , interpretative statements of results are same
 - Probability model
 - For a unit increase in a particular predictor, corresponding change in the probability value is not a constant, while holding all other predictors constant
 - Depends on the specific values of the predictor
 - Interpretative statements of results depend on specific values of x



LOGISTIC REGRESSION

- Odds and odds ratios
 - Odds is a ratio of two probability values (prob. of class 1/prob. Of Class 0)
 - Odds ratio is ratio of two odds (odds of class m1/odds of class m2)
 - Odds ratio $> 1 \Rightarrow$ odds of class m1 are higher than class m2
- Open RStudio



LOGISTIC REGRESSION

- Linear Regression for a categorical outcome variable?
 - Can be done by treating the outcome variable as continuous and coding it numerically
 - However, anomalies will lead to spurious modeling
 - Predictions can take any value, not just dummy values {0,1}
 - Outcome variable or residuals don't follow normal distribution
 - binomial distribution
 - Variance of outcome variable is not constant across all records (violation of homoscedasticity)
 - $np(1-p)$



LOGISTIC REGRESSION

- Logistic Regression for Profiling Task
 - Apart from model performance on validation partition
 - Model's fit to data is assessed on training partition
 - However, still avoid overfitting
 - Usefulness of predictors is examined
 - Goodness of fit metrics
 - Overall fit of the model
 - Deviance (equivalent to SSE in linear regression)
 - $1 - \text{Deviance}/\text{Null Deviance}$ (equivalent to multiple R^2 in linear regression)
 - Single predictors



LOGISTIC REGRESSION

- Outcome variable with m classes ($m > 2$)
 - Multinomial logistic regression
 - Separate binary logistic regression model for $m-1$ classes (one class is treated as reference class)
 - Ordinal logistic regression
 - Large no. of ordinal classes: treat ordinal variable as continuous variable and apply multiple linear regression



LOGISTIC REGRESSION

- Outcome variable with m classes (m>2)
 - Ordinal logistic regression
 - Small no. of ordinal classes: Proportional odds or cumulative logit method
 - Separate binary logistic regression model for m-1 cumulative probabilities
- For a three class case: C1, C2, and C3 and a single predictor x1
$$\text{logit}(C1) = \alpha_0 + \beta_1 x_1$$
$$\text{logit}(C1 \text{or } C2) = \beta_0 + \beta_1 x_1$$
- RStudio

Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE



IIT ROORKEE



NPTEL ONLINE
CERTIFICATION COURSE

LOGISTIC REGRESSION PART-5

LECTURE 50

DR. GAURAV DIXIT
DEPARTMENT OF MANAGEMENT STUDIES



LOGISTIC REGRESSION

- Interpretation of Results
 - Odds model
 - Multiplicative factor (e^{β})
 - If $\beta < 0$, increase in $x \Rightarrow$ decrease in odds
 - If $\beta > 0$, increase in $x \Rightarrow$ increase in odds
 - For any value of x , interpretative statements of results are same
 - Probability model
 - For a unit increase in a particular predictor, corresponding change in the probability value is not a constant, while holding all other predictors constant
 - Depends on the specific values of the predictor
 - Interpretative statements of results depend on specific values of x



LOGISTIC REGRESSION

- Odds and odds ratios
 - Odds is a ratio of two probability values (prob. of class 1/prob. Of Class 0)
 - Odds ratio is ratio of two odds (odds of class m1/odds of class m2)
 - Odds ratio $> 1 \Rightarrow$ odds of class m1 are higher than class m2
- Open RStudio



LOGISTIC REGRESSION

- Linear Regression for a categorical outcome variable?
 - Can be done by treating the outcome variable as continuous and coding it numerically
 - However, anomalies will lead to spurious modeling
 - Predictions can take any value, not just dummy values {0,1}
 - Outcome variable or residuals don't follow normal distribution
 - binomial distribution
 - Variance of outcome variable is not constant across all records (violation of homoscedasticity)
 - $np(1-p)$



LOGISTIC REGRESSION

- Logistic Regression for Profiling Task
 - Apart from model performance on validation partition
 - Model's fit to data is assessed on training partition
 - However, still avoid overfitting
 - Usefulness of predictors is examined
 - Goodness of fit metrics
 - Overall fit of the model
 - Deviance (equivalent to SSE in linear regression)
 - $1 - \text{Deviance}/\text{Null Deviance}$ (equivalent to multiple R^2 in linear regression)
 - Single predictors



LOGISTIC REGRESSION

- Outcome variable with m classes ($m > 2$)
 - Multinomial logistic regression
 - Separate binary logistic regression model for $m-1$ classes (one class is treated as reference class)
 - Ordinal logistic regression
 - Large no. of ordinal classes: treat ordinal variable as continuous variable and apply multiple linear regression



LOGISTIC REGRESSION

- Outcome variable with m classes (m>2)
 - Ordinal logistic regression
 - Small no. of ordinal classes: Proportional odds or cumulative logit method
 - Separate binary logistic regression model for m-1 cumulative probabilities
- For a three class case: C1, C2, and C3 and a single predictor x1
$$\text{logit}(C1) = \alpha_0 + \beta_1 x_1$$
$$\text{logit}(C1 \text{or } C2) = \beta_0 + \beta_1 x_1$$
- RStudio



Key References

- Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data by EMC Education Services (2015)
- Data Mining for Business Intelligence: Concepts, Techniques, and Applications in Microsoft Office Excel with XLMiner by Shmueli, G., Patel, N. R., & Bruce, P. C. (2010)



Thanks...



IIT ROORKEE



NPTEL
ONLINE
CERTIFICATION COURSE