

Probabilistic Graphical Model (Bayesian Network)

Advanced Distributed Machine Learning Labs (ADML)

Essex Lake Group LLC

Knowledge Sharing Session

July, 2015

Essex Lake Group LLC

www.essexlg.com

Objective: Embracing Bayesian approaches for business problem solving

Technical Objectives:

- How can we identify the probabilistic relationships among the variables in the given data
- For any given target, how can we identify the most influencing variables ?
- Can Bayesian networks help us in finding better insights than the existing Pivot approach ?

Business Objective :

- For the MET004 Account level data, how do we evaluate the performance of “*Participation Rate*” using PGM

Agenda

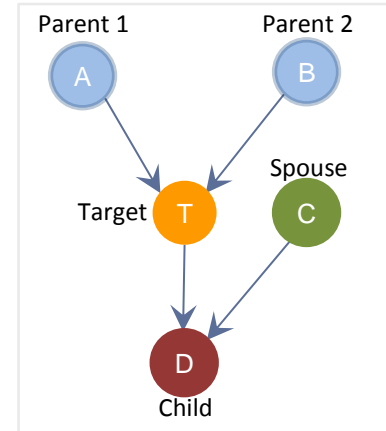
- **How do we construct Probabilistic Graphical Models (PGMs) ?**
- PGMs on Sales Ratio Data
- PGM with respect to CART

Network Terminology

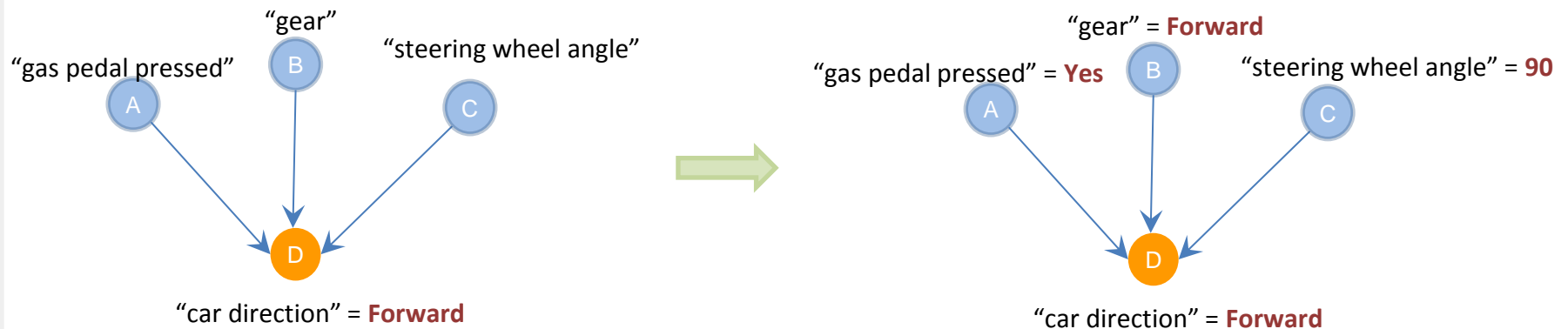
Markov Blanket: For a given target variable **T** What are the most influencing variables ?

The **Markov Blanket** of a (Target) node consists of all nodes that make this Target conditionally independent of all the other nodes in the model:

- The **Parents** (blue nodes in the example) are used for cutting the information coming from their ascendants
- The **Children** (red node) are used for cutting the information coming from their descendants
- The **Spouses** (or co-parents, dark green) are used for cutting the information coming from the ascendants of the Children

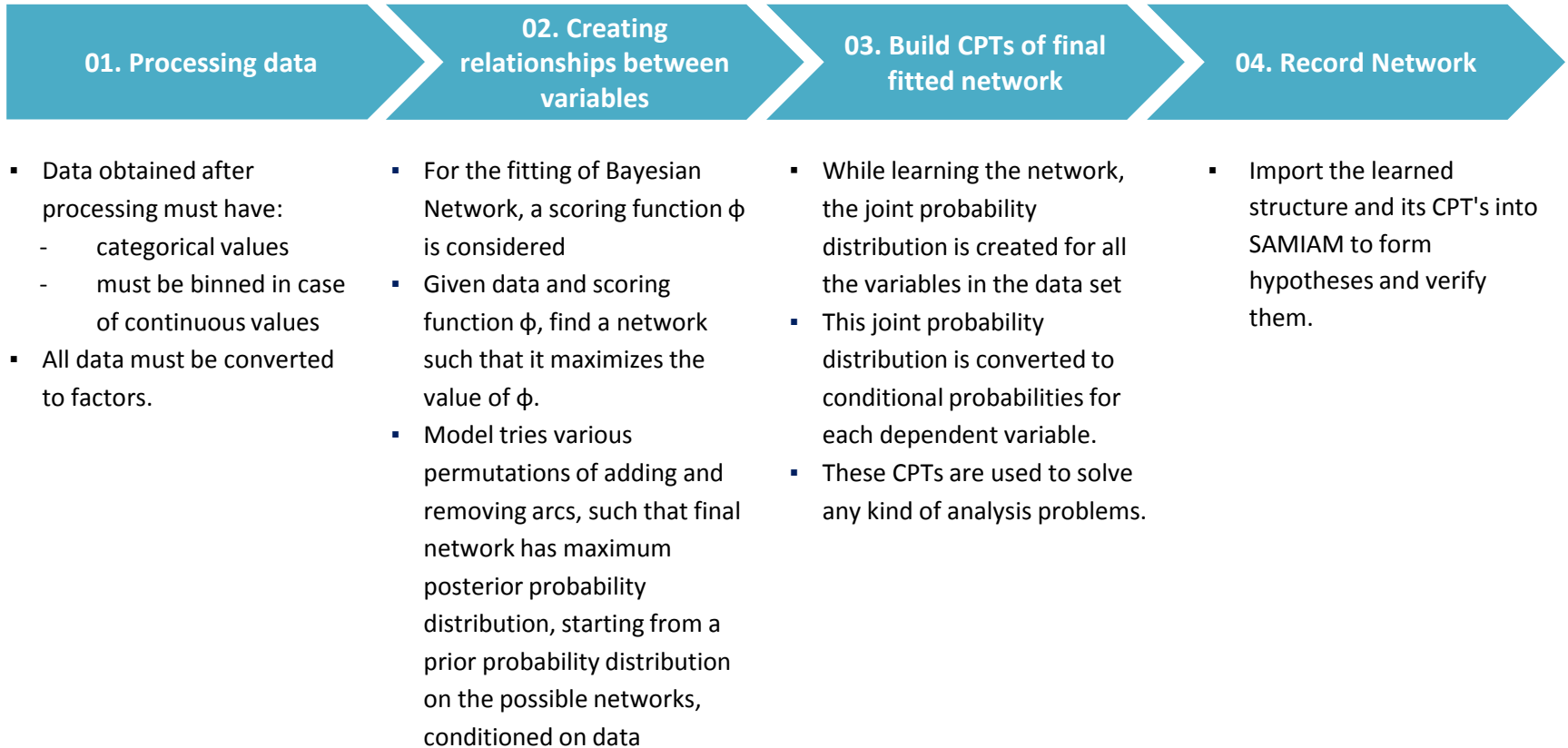


Dependency among the variables of a system (**Car**)



If one Variable takes on some value then other variables are going to take that value

Constructing Bayesian Network - Process Flow



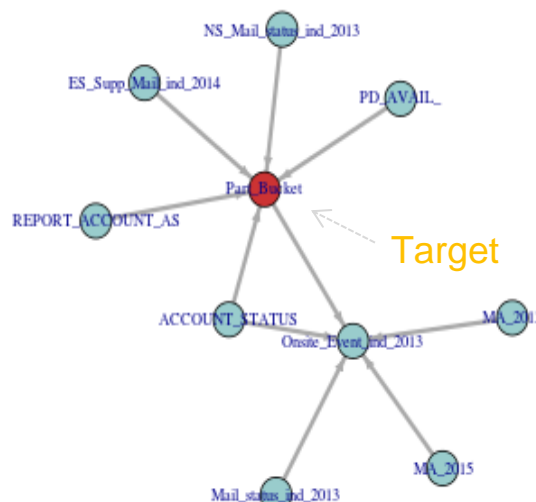
How do we interpret the Bayesian Network for problem solving ?

Example: 1

A. Forward Analysis :

Q1. I have not done any mailing activity this year. How did this affect my participation percentage?

Q2. To check for conditional dependencies like mailing activity not being helpful if the account is already active.

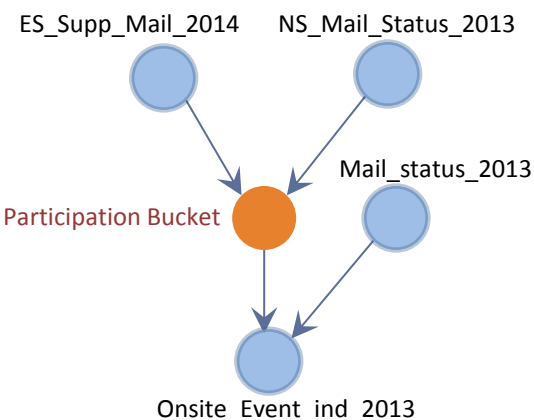


B. Backward Analysis

Q3. If we know that our participation is high, then find responsible variables and their values to replicate the result in the future.

Q4. For my participation bucket to lie in the higher bucket, should I offer payroll deduct or not?

Fig: Markov Blanket for a target variable



Variable Name	Evidence Set	Participation Bucket value before evidence	Participation Bucket value after evidence	Change in Values
Sponsored Mail 2014	Y	15.15%	19.46%	+4.31%
Non Sponsored Mail 2013	Y	15.15%	18.85%	+3.70%
Mail Status 2013	Y	15.15%	18.68%	+3.53%
Sponsored Mail 2014	N	15.15%	13.8%	-1.35%
Non Sponsored Mail 2013	N	15.15%	10.30%	-4.85%
Mail Status 2013	N	15.15%	10.32%	-4.83%

All Probability Values for Participation being in the 2.5%-5% (high) bucket.

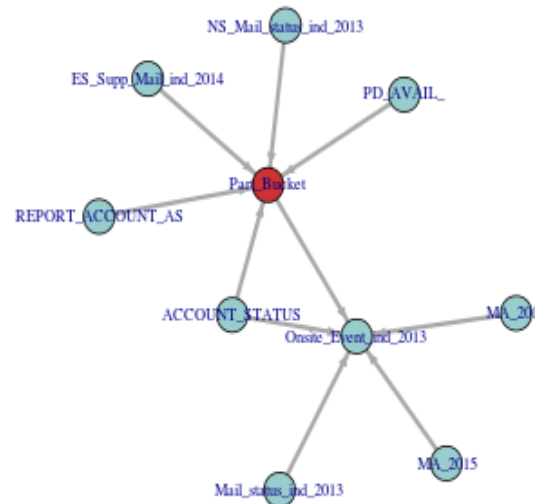
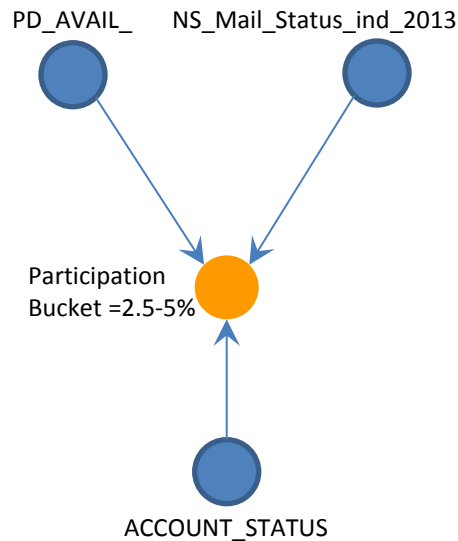
One more example → How do we interpret the Bayesian Network for problem solving ?

Example: 2

A. Forward Analysis :

Q1. I have not done any mailing activity this year. How did this affect my participation percentage?

Q2. To check for conditional dependencies like mailing activity not being helpful if the account is already active.



B. Backward Analysis

Q3. If we know that our participation is high, then find responsible variables and their values to replicate the result in the future.

Q4. For my participation bucket to lie in the higher bucket, should I offer payroll deduct or not?

Evidence Value	Non-Sponsored Mail 2013	Payroll Deduct	Account Status
Participation Bucket = 2.5-5%	Y	Y	ACTIVE

Participation Bucket	Variable	Evidence	Participation Bucket Value before Evidence	Participation Bucket Value before Evidence	Change in Values
2.5% - 5%	Non-Sponsored Mail 2013 Payroll Deduct Account Status	Y Y ACTIVE	15.15%	28.55%	+13.40%

All Probability Values for Participation being in the 2.5%-5% (high) bucket.

Most Probable Explanation for given variable (MAP): *Process Flow*

01. Selection of Target Variable

- The network consists of all the variables in the data set
- To perform any kind of analysis, the target variable must be selected.
- The network is built without keeping a target variable so in case analysis needs to be done on any other variable, it can also be selected in the same network.

02. Find Target's Markov Blanket

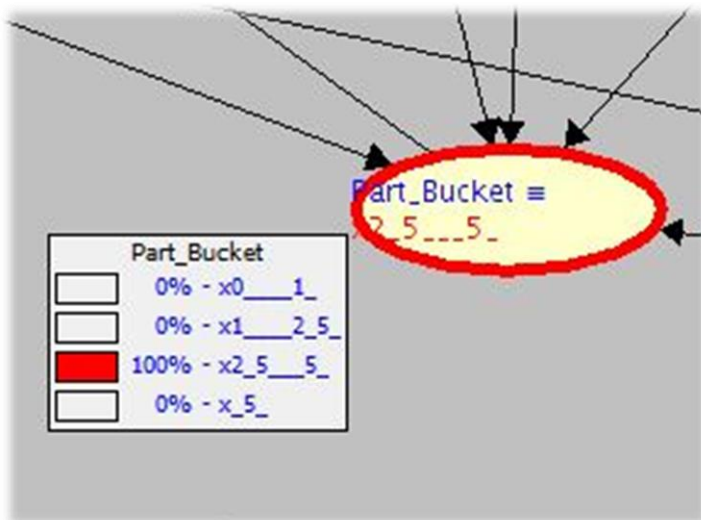
- Target's parent nodes, child nodes and all other parents of target's child nodes form it's Markov Blanket.
- These variables make the target variable independent from the rest of the variables in the network.

03. Set Evidence on Target

- In Samiam, for the given Bayesian Network, set evidence to the target variable, which is the variable of our interest.
- We need to find what values should other variables have so as to observe the evidenced value of target variable.

04. Run MAP Query on Markov Blanket

- Select the target variable's Markov Blanket variables for which MAP values are to be computed.
- Run MAP query.
- This query returns the value of the selected variables such that the posterior probability of the evidenced target is maximized.



MAP Computation

☐ Approximate ☒ Exact ☐ Sloppy? Slop: 0.5

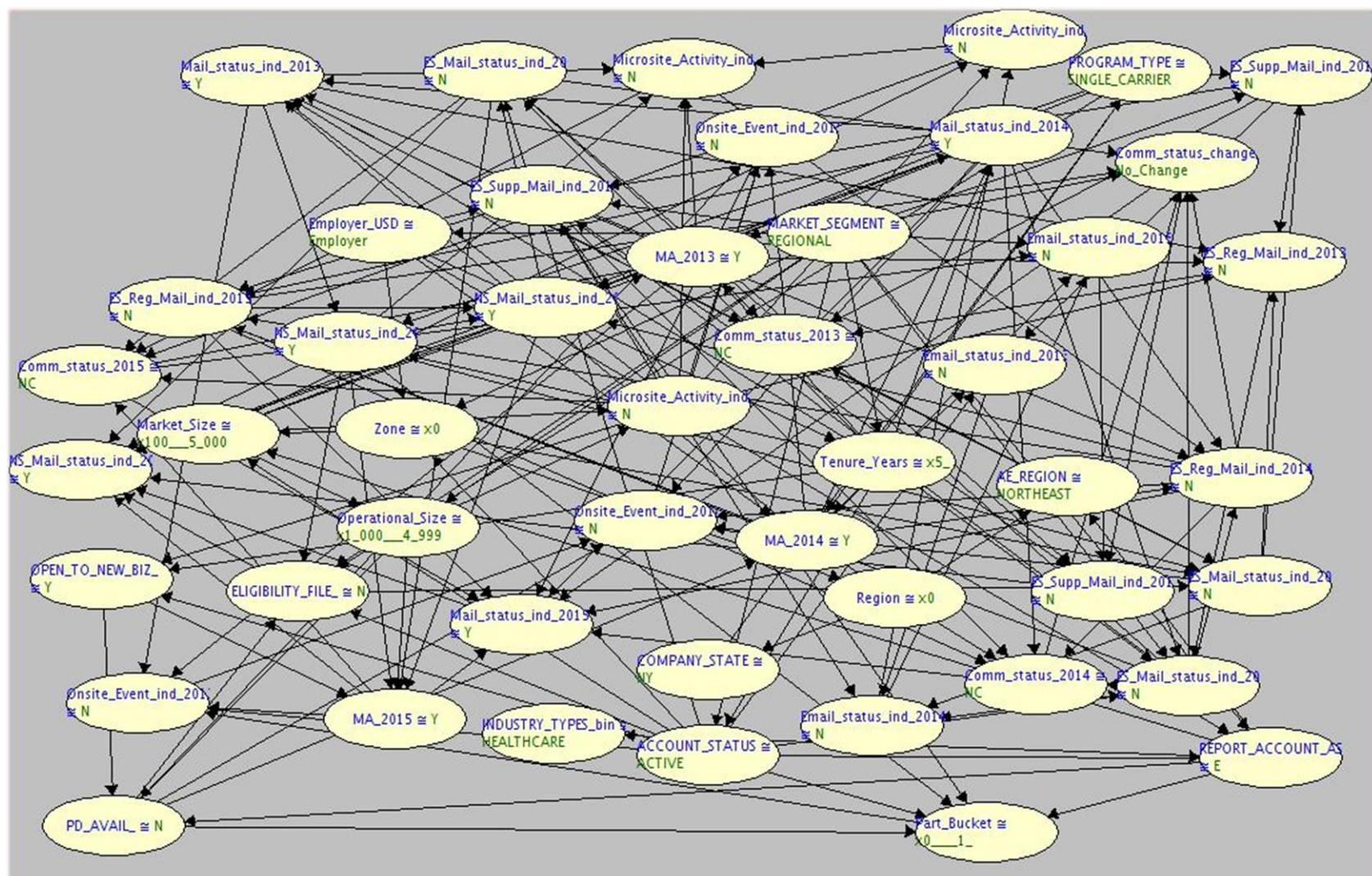
Time out (secs): Width barrier (0=none):

9 MAP Variables: [Variable Selection Tool](#)

Variable	Value
ACCOUNT_STATUS	ACTIVE
ES_Supp_Mail_ind_2014	N
MA_2013	Y
MA_2015	Y
Mail_status_ind_2013	Y
NS_Mail_status_ind_2013	Y
Onsite_Event_ind_2013	N
PD_AVAIL	Y
REPORT_ACCOUNT_AS	E

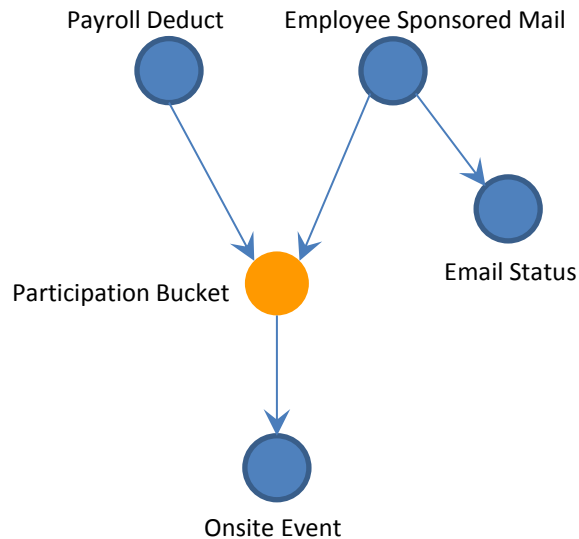
Find:

P(MAP,e)=0.03747422054239237
P(MAP|e)=0.3035436139245486
Result is exact.



SAMIAM and R-SHINY Demo

More details on Network Structure:



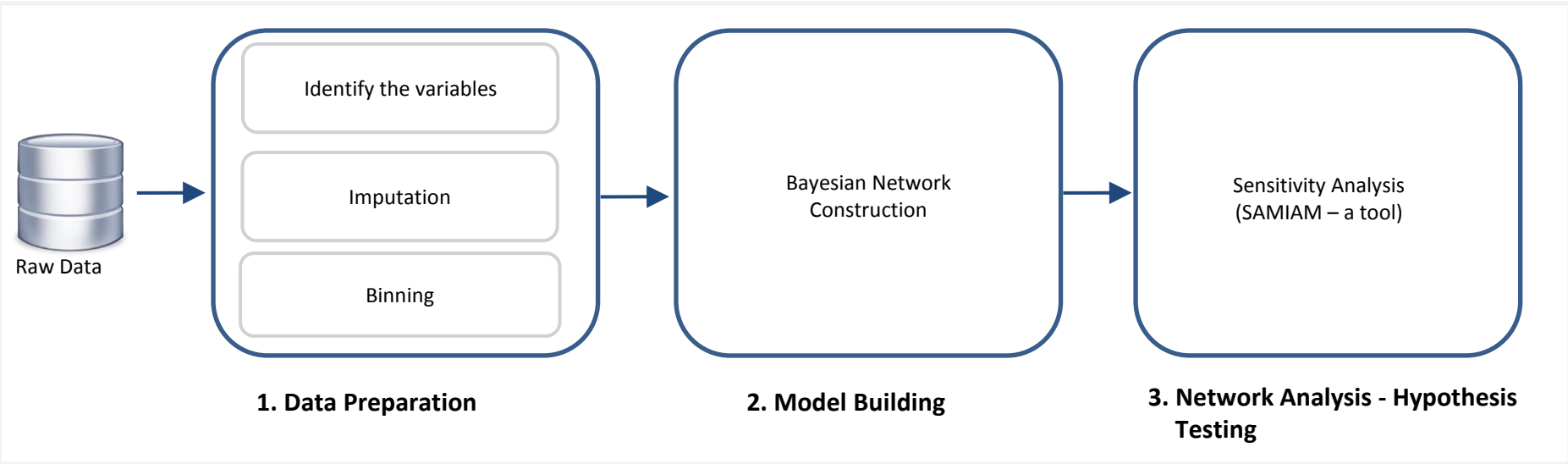
- Payroll Deduct, Employee Sponsored Mail and Participation Bucket for a V-Structure.
- This V-Structure means that Payroll Deduct and Employee Sponsored Mail are **independent** of each other.
- This means, any change in Payroll Deduct **won't influence** Employee Sponsored Mail and vice versa.
- Also, any change in Email Status which is the sibling of Participation Bucket **will result** in change in Participation Bucket as well due to the presence of an active trail between the two.
- Payroll deduct can however, **affect** Onsite Event through Participation Bucket but not Email Status due to the presence of a V-Structure between Payroll Deduct and Employee Sponsored Mail which is the parent of Email Status.
- Similarly, Employee Sponsored Mail **can affect** Onsite Event due to the presence of the active trail between them.
- Also, if Participation Bucket is evidenced, then, The V-Structure will be activated and **influence** can flow between Payroll Deduct and Employee Sponsored Mail.

Agenda

- What are Probabilistic Graphical Models (PGM)
- **PGMs on Sales Ratio Data**
- PGM with respect to CART

Our work in 3 modules:

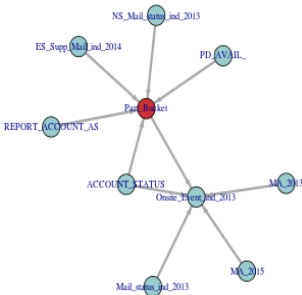
Module 1:



Module 2: Pivot Approach Vs. Bayesian Network Approach

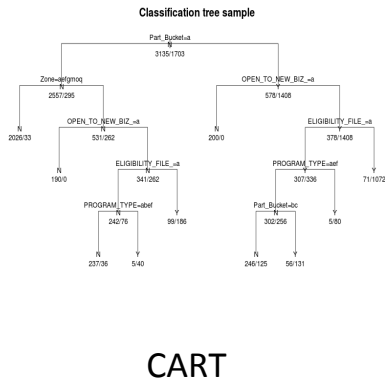
Sno	No. of V	Variable Name	levels	MAP Query Results for Markov
1	7	par_conv	low medium high extra_high	participant_band part_perc_bin MA_2014 Tenure_Years
2	13	participant_band	High Low	ACCOUNT_STATUS OPEN_TO_NEW_BIZ PD_AVAIL REPORT_ACCOUNT_AS
3	16	part_perc_bin	1	PROGRAM_TYPE ACCOUNT_STATUS OPEN_TO_NEW_BIZ PD_AVAIL

Pivots

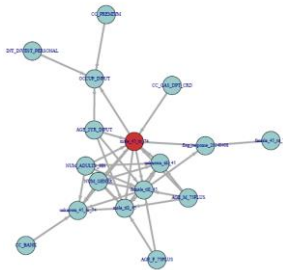


Bayesian network

Module 3: Observations on CART and Bayesian Networks



CART

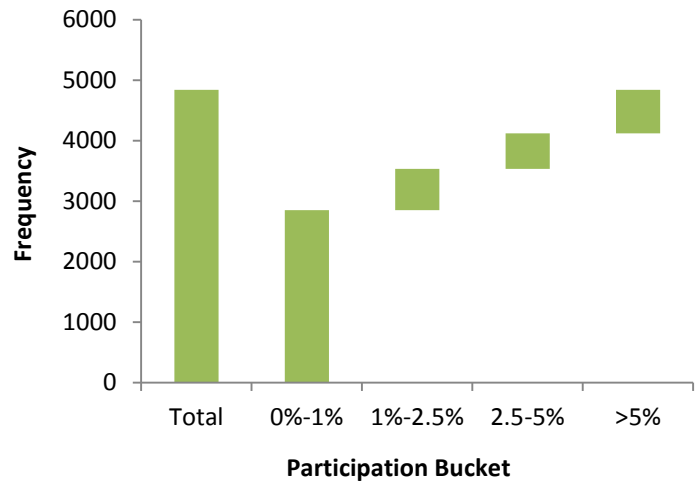


Bayesian Network

Data Details

1. Sales Account Data (2013 to 2015)
2. Demographic details
3. MetLife's mailing
4. Nature and tenure policy
5. Agent and employer type details
6. No. of participants and eligible accounts

1. Target Variable : Participation Bucket
(Ratio of no. of participants and no. of eligible)
2. No. of observations : 5007



Account Variables

- Company Name
- Market Segment
- Company State

Mail Variables

- Marketing Activity
- Non-Sponsored Mail
- Employee Sponsored Mail

Flag Variables

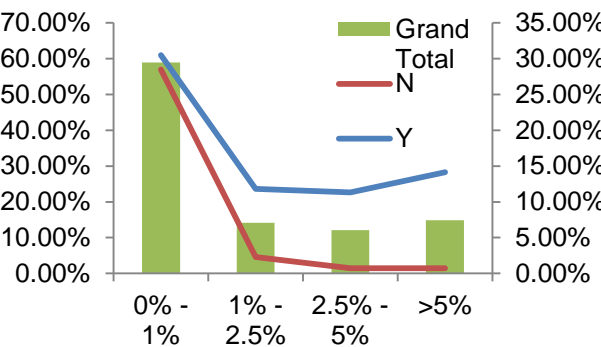
- New or Existing account flag

Policy Variables

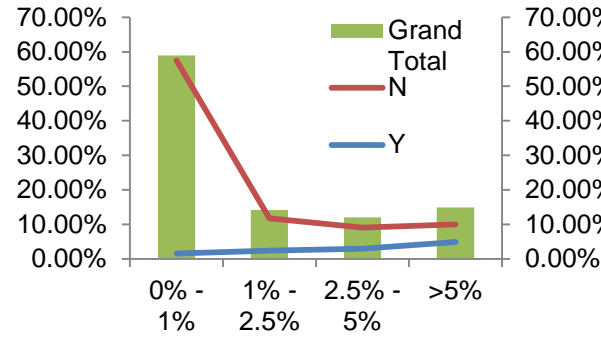
- Program Type
- Payroll Deduct
- Employer Type

1. Company Name and Company State depict the account details. Market segment could be local, national, regional, USD or other. Majority of the accounts are regional
2. Mailing variable(Non Sponsored Mail) is crucial for participation bucket
3. Policy Variable (Payroll Deduct) is essential to improve participation bucket
4. Most of the accounts are of Single Carrier type and don't have Payroll Deduct

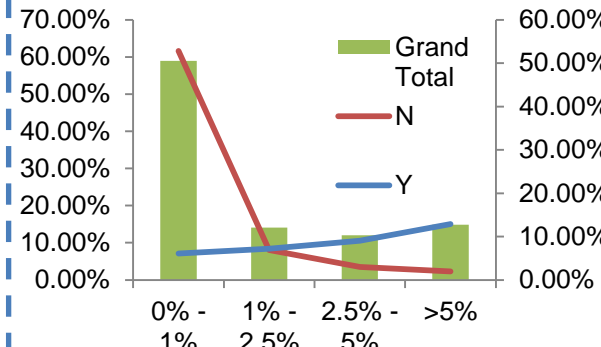
Bivariate Analysis on Target Variable "Participation Bucket"



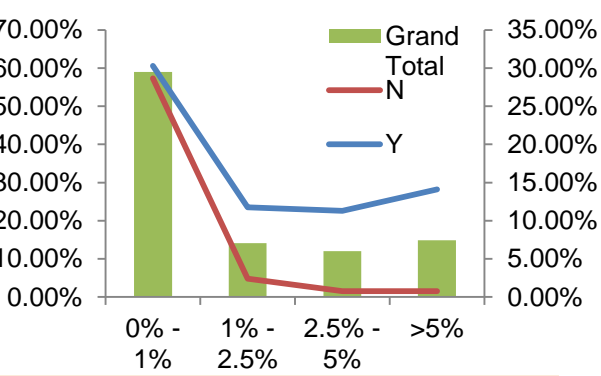
Marketing Activity 2013 - For Participation Bucket in range 2.5%-5%, 11% accounts have marketing activity in 2013



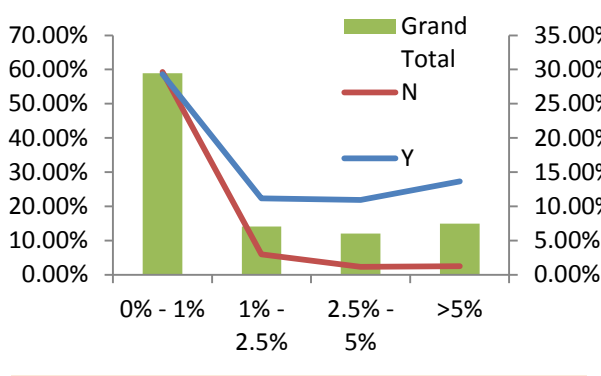
Onsite Event 2013 - For Participation Bucket in range 2.5%-5%, 3% accounts have Onsite Events in 2013



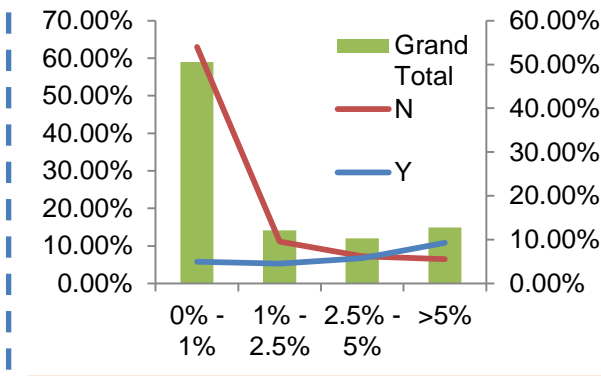
Payroll Deduct - For Participation Bucket in range 2.5%-5%, 9% accounts have Payroll Deduct



Mail Status 2013 - For Participation Bucket in range 2.5%-5%, 11% accounts have Mail Status 2013

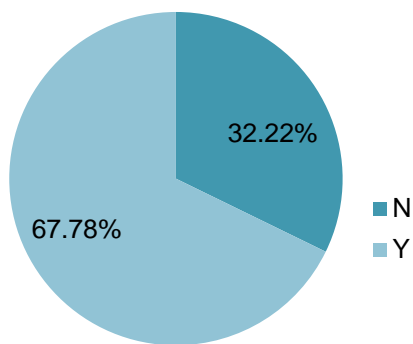


Non Sponsored Mail 2013 - For Participation Bucket in range 2.5%-5%, 11% accounts have Non Sponsored Mail 2013



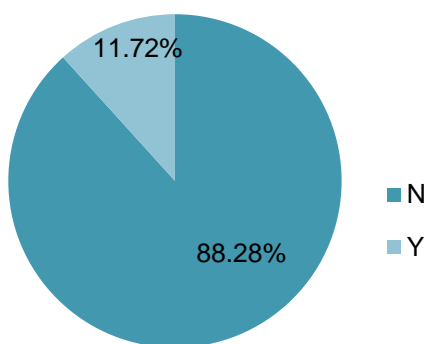
Employee Sponsored Mail 2013 - For Participation Bucket in range 2.5%-5%, 5.81% accounts have Employee Sponsored Mail 2013

Profiling of Important Variables against Target Variable – “Participation Bucket”



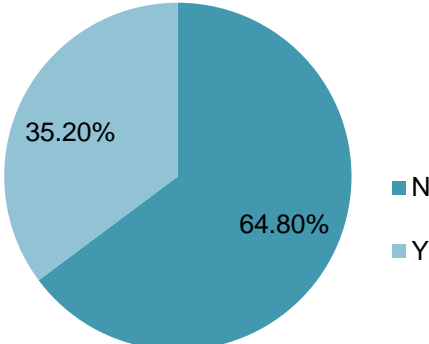
Marketing Activity 2013

68% of the accounts have Marketing Activity 2013



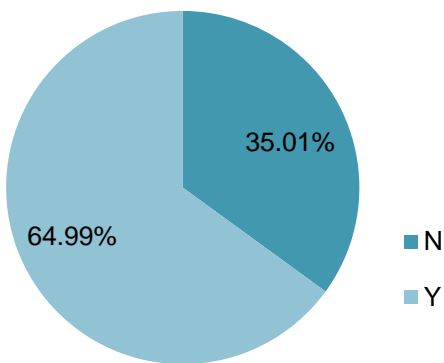
Onsite Event 2013

88% of the accounts don't have Onsite Event 2013



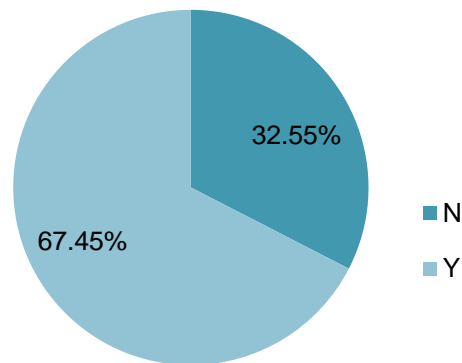
Payroll Deduct

65% of the accounts don't have Payroll Deduct



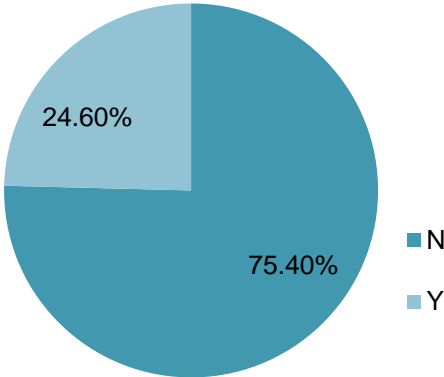
Non Sponsored Mail 2013

65% of the accounts have Non Sponsored Mail 2013



Mail Status 2013

65% of the accounts have Non Sponsored Mail 2013



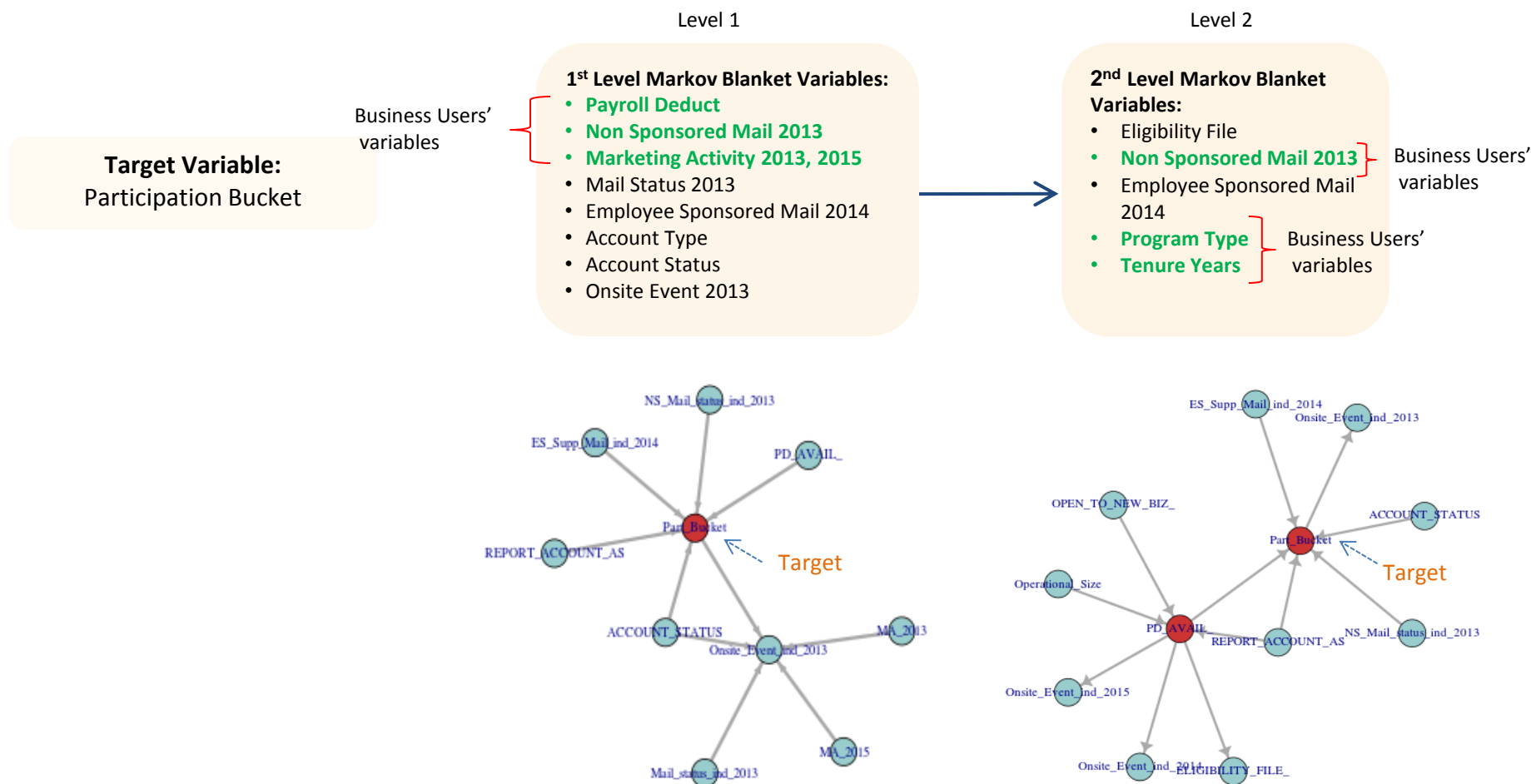
Employee Sponsored Mail 2013

25% of the accounts received employee sponsored mails.

Bayesian Network discovers NOT ONLY the variables identified by the business users AND ALSO discovers other most influencing variables

Description: To show how to analyze the variables in a Bayesian Network, we tried to find the values of the variables that influence the target variable the most, so as to maximize/ minimize the desired target variable class.

The following are the **variables found using the Markov Blanket** of the Target Variable. When considered together, these make the target Variable independent of the rest of the network. Thus, it can be said that **these hold the maximum influence over the target**.

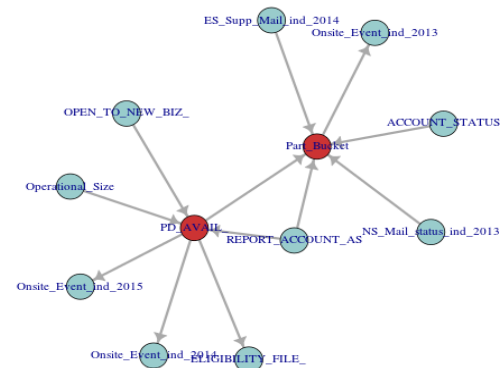


Module 1

Better Results with Bayesian Network (All Variables)

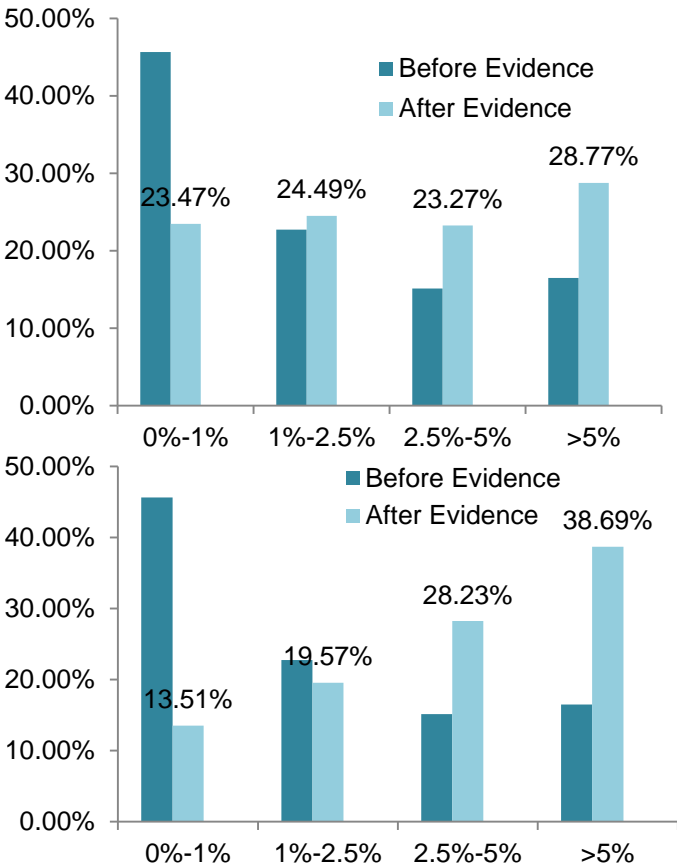
Description: To show how to analyze the variables in a Bayesian Network, we tried to find the values of the variables that influence the target variable the most, so as to maximize/ minimize the desired target variable class.

The following are the variables found using the Markov Blanket of the Target Variable. When considered together, these make the target Variable independent of the rest of the network. Thus, it can be said that these hold the maximum influence over the target.



Variables Selected: Evidence
Payroll Deduct -> Yes

Variables Selected: Evidence
Payroll Deduct -> Yes
Non- Sponsored Mail Status 2013 -> Yes
Mail Status 2013 -> Yes



Inference:

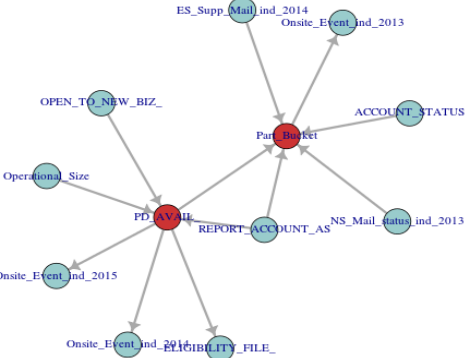
- Payroll Deduct is Critical to account performance
- Payroll Deduct without the support of any Mailing Activity cannot drive account performance.

Module 1

Full Model Analysis on Bayesian Network

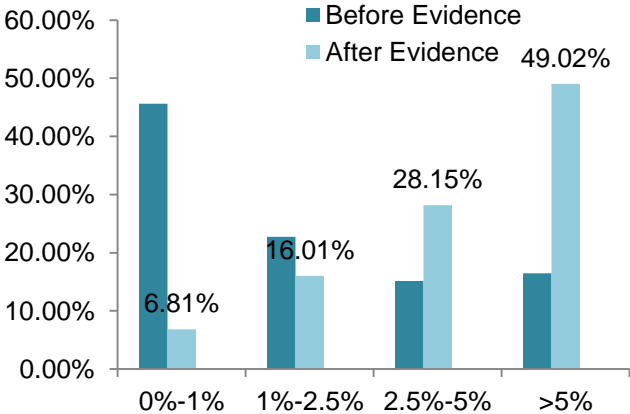
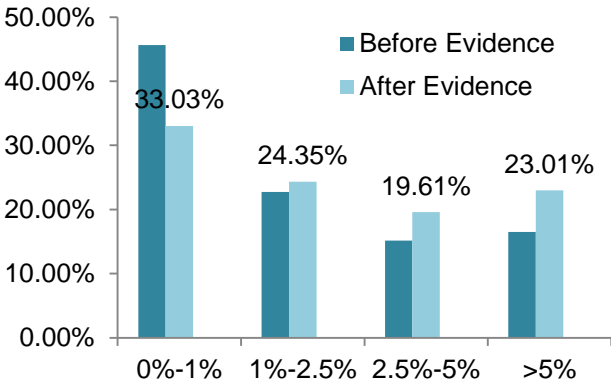
Description: To show how to analyze the variables in a Bayesian Network, we tried **to find the values of the variables that influence the target variable the most**, so as to maximize/ minimize the desired target variable class.

The following are the **variables found using the Markov Blanket** of the Target Variable. When considered together, these make the target Variable independent of the rest of the network. Thus, it can be said that **these hold the maximum influence over the target**.



Variables Selected: Evidence
Eligibility File -> Yes

Variables Selected: Evidence
Eligibility File -> Yes
Non- Sponsored Mail -> Yes
Employee Sponsored Mail ->Yes



Inference:

- Payroll Deduct causes the maximum change in the Target Variable. Apart from that, Payroll Deduct’s Markov blanket variables like Eligibility File, Employee Sponsored Mail should be considered as well.

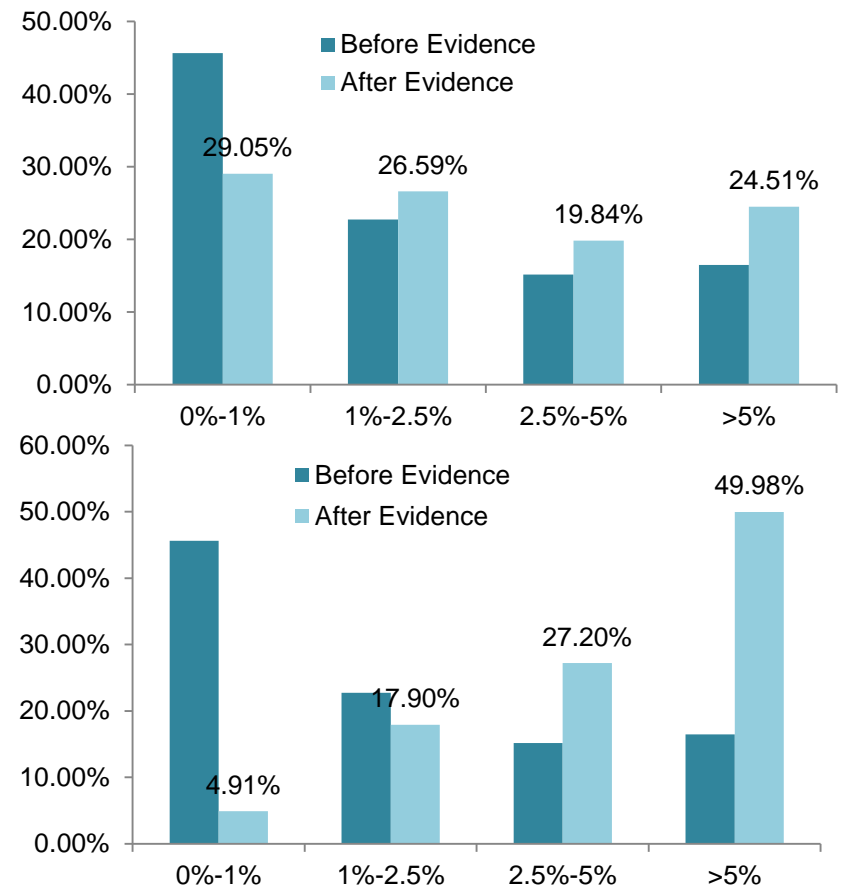
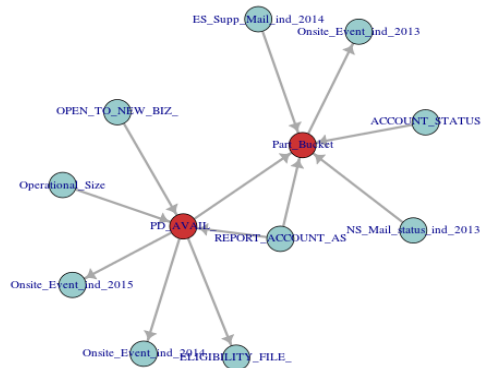
Module 1

Full Model Analysis on Bayesian Network

Hypothesis: Mailing activity significantly improves participation bucket with mailing activity in 2013 having more influence than 2014

Variables Selected: Evidence
Employee Sponsored Mail 2014 -> Yes
Non-Sponsored Mail 2014 -> Yes
Onsite Event 2014 -> Yes

Variables Selected: Evidence
Employee Sponsored Mail 2013 -> Yes
Non-Sponsored Mail 2013 -> Yes
Onsite Event 2013 -> Yes



Inference :

- Improving Mailing Activity for the year 2013 has more influence than in the year 2014
- Certain combinations of marketing activities tend to more effective, in particular - a combination of Non-sponsored, employee sponsored and On-site Events for the year 2013 is yielding higher Participation Bucket (12.05% higher)

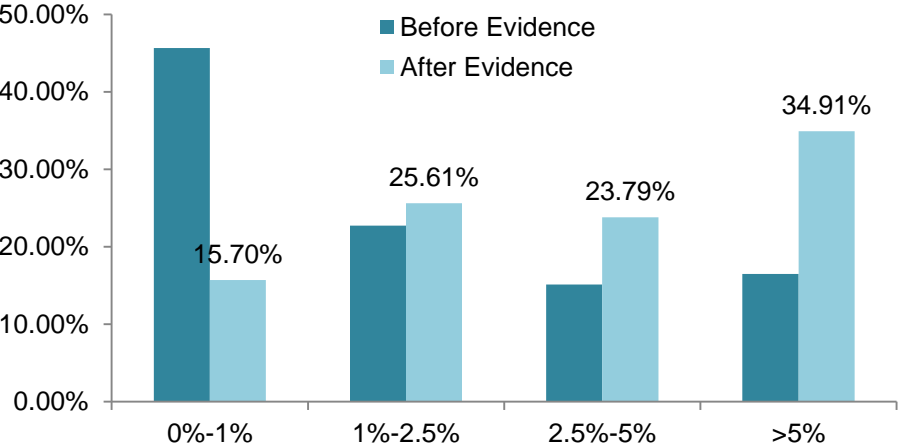
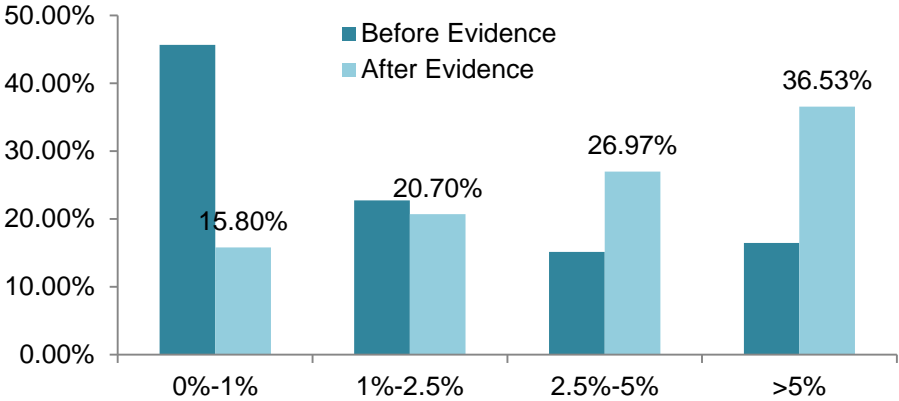
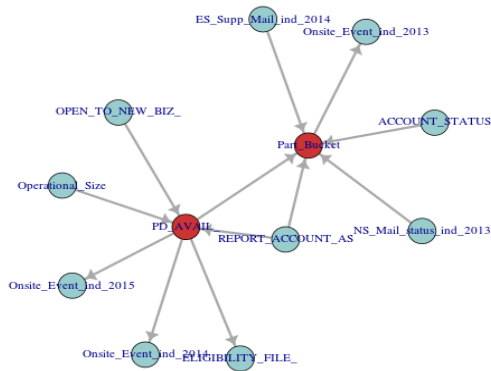
Module 1

Full Model Analysis on Bayesian Network

Hypothesis: Participation Bucket **depends on type of Marketing Channel** – Eligibility File, Employee Sponsored Mail, Onsite Events, Payroll Deduct.

Variables Selected: Evidence
Marketing Activity 2013 -> **Yes**
Payroll Deduct -> **Yes**

Variables Selected: Evidence
Marketing Activity 2013 -> **Yes**
Onsite Event 2013 -> **Yes**



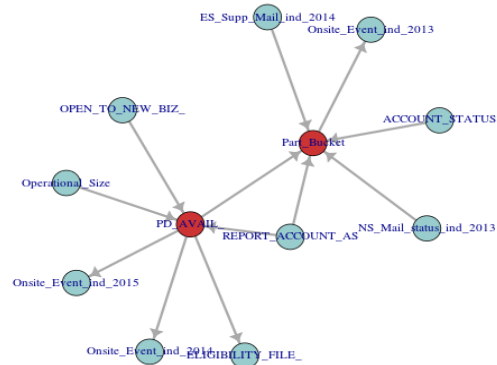
Inference :

- Along with Marketing Activity, performance also depends on the type of channels – Payroll Deduct, and Onsite Event Activity 2013 being the most important

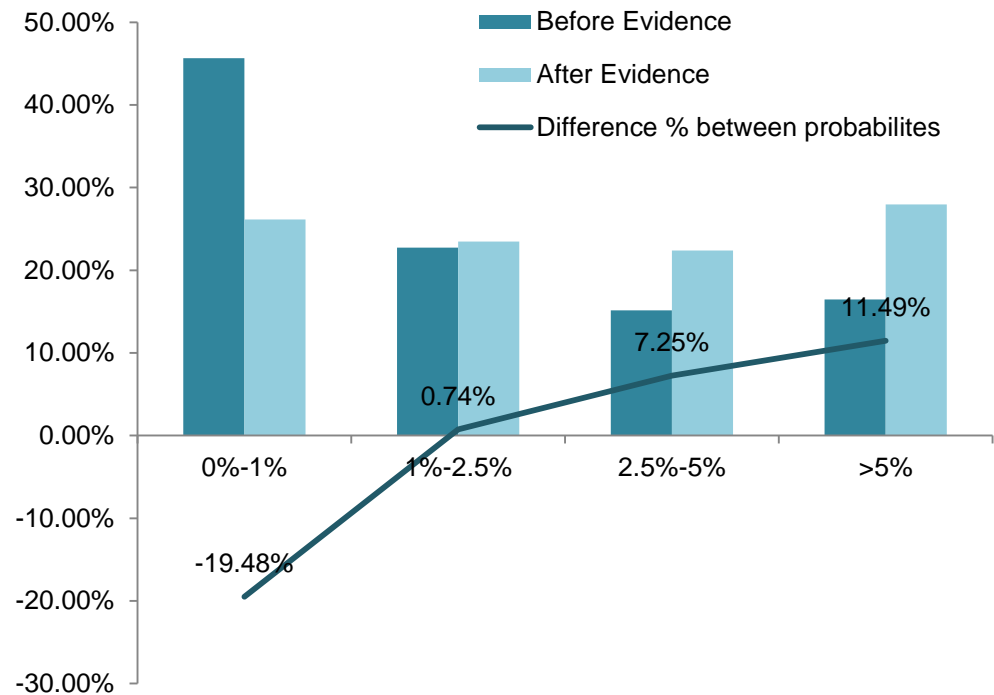
Module 1

Full Model Analysis on Bayesian Network

Hypothesis: Newer accounts can improve their Participation Bucket by Mailing Activity



Variables Selected: Evidence
Tenure Years -> 0
Non-Sponsored Mail 2013 -> Y
Employee Sponsored Mail 2013 -> Yes
Onsite Event 2013 -> Yes



Inference :

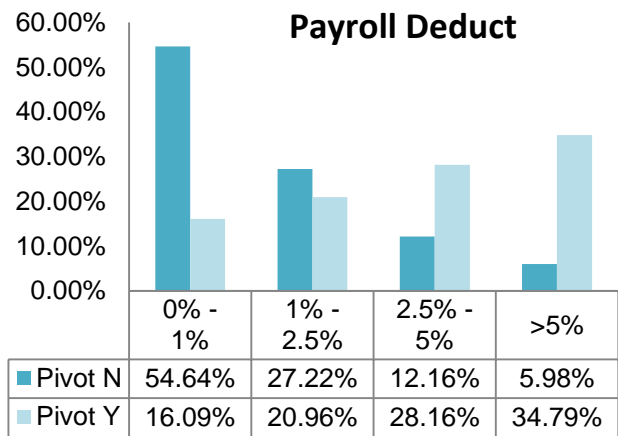
- Newer accounts can improve their Participation Bucket by Mailing Activity - specifically, sending non sponsored and employee sponsored mails and doing some onsite events.

Module 2

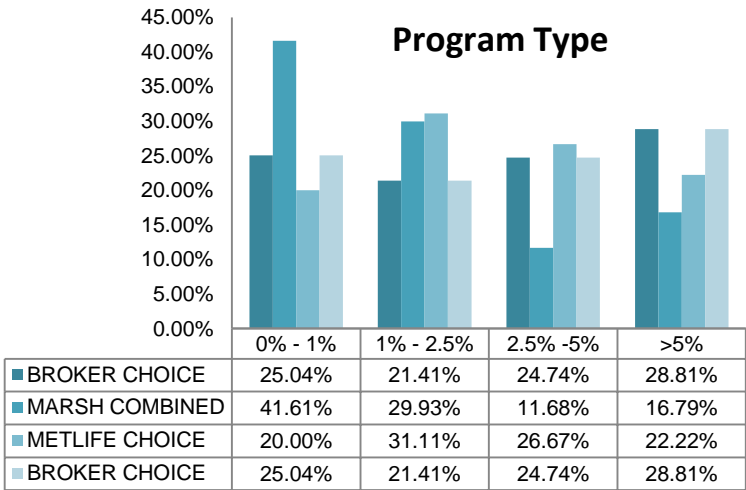
Same insights from the both Pivot and Bayesian Analysis (3 variables)

Description: To test the accuracy of the Bayesian network, we built models with only those variables that were considered in pivot analysis. We then **observed the target variable** i.e. the Participation Bucket **with respect to selected variables**. The following results were obtained:

Pivot Analysis

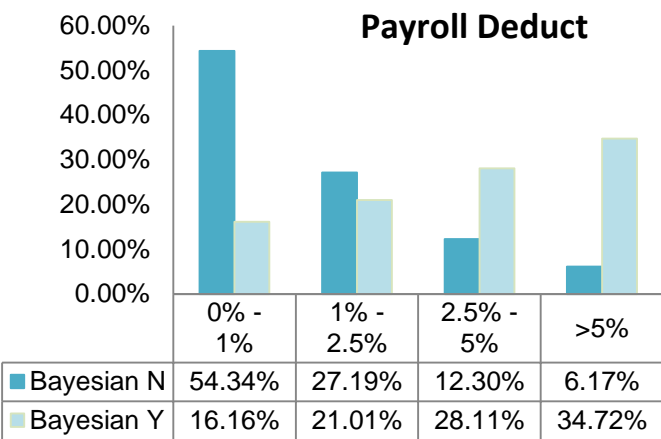


Participation Bucket

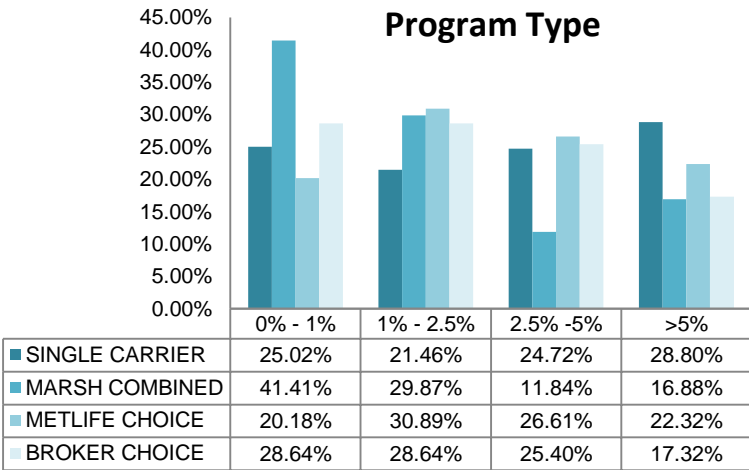


Participation Bucket

Bayesian Analysis



Participation Bucket

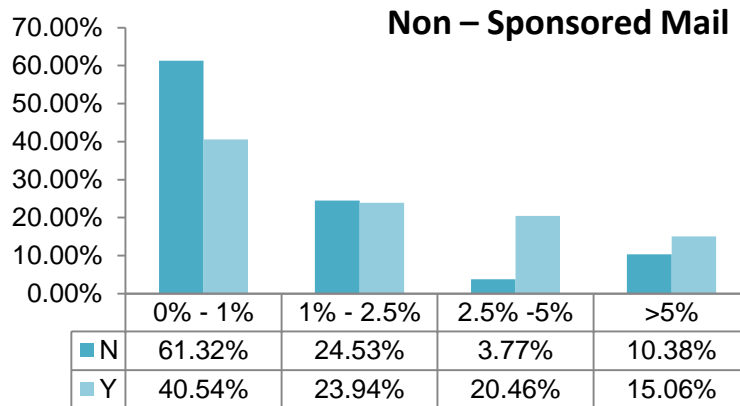


Participation Bucket

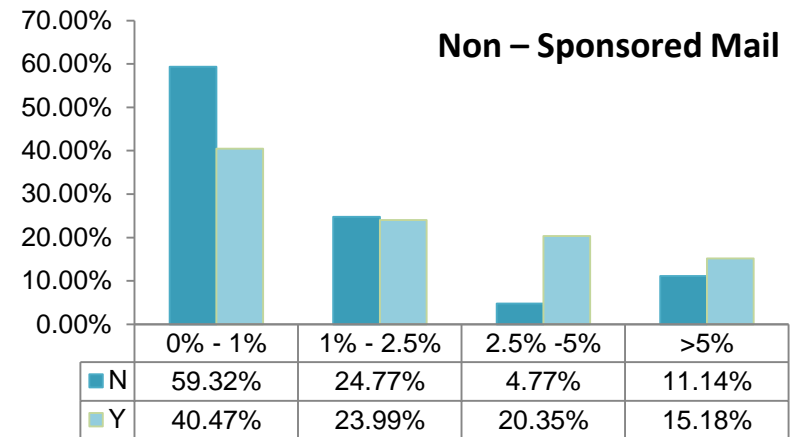
More insights with more interactions: Pivot and Bayesian (6 variables)

Description: To test the accuracy of the Bayesian network, we built models with only those variables that were considered in pivot analysis. We then **observed the target variable** i.e. the Participation Bucket **with respect to selected variables**. The following results were obtained:

Pivot Analysis

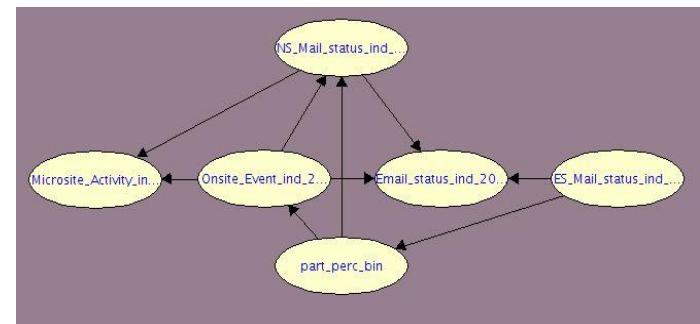


Bayesian Analysis



Variables Selected:

- Non Sponsored Mail 2014
- Employee Sponsored Mail 2014
- Onsite Event 2014
- Microsite Activity 2014
- Participation Bucket 2014
- Email Status 2013



Bayesian Advantage:

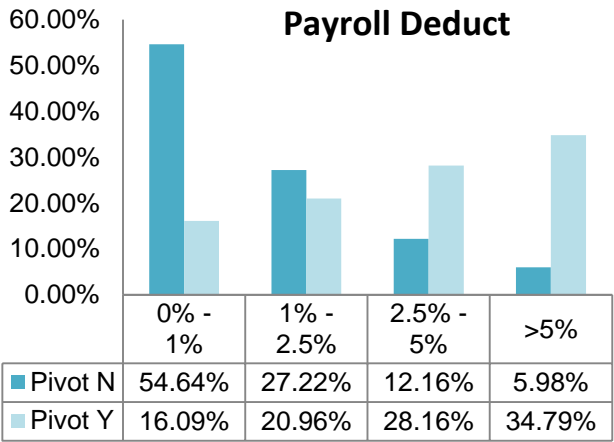
- Pivot Tables often return zero values in case it cannot return values according to the set filter when multiple filters are applied.
- Bayesian Networks do not work on values but on conditional probabilities. Thus, even though data rows may be unavailable, probability of the occurrence of each bucket according to the evidence given to filter variables in the network will still be returned.

Module 2

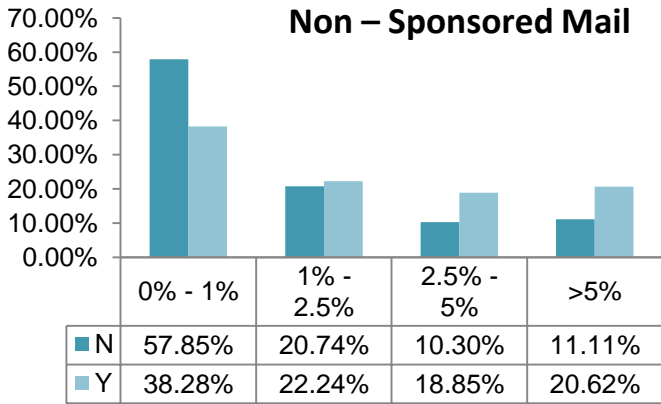
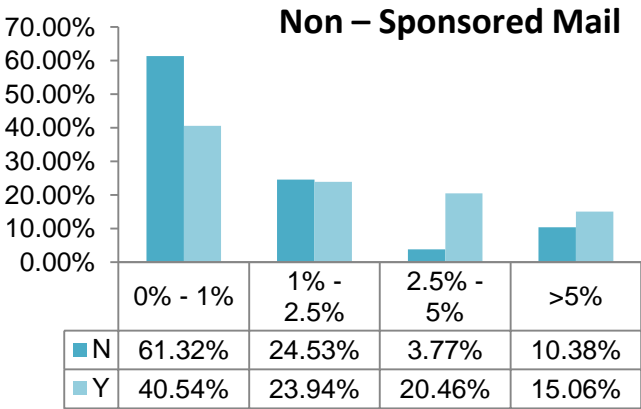
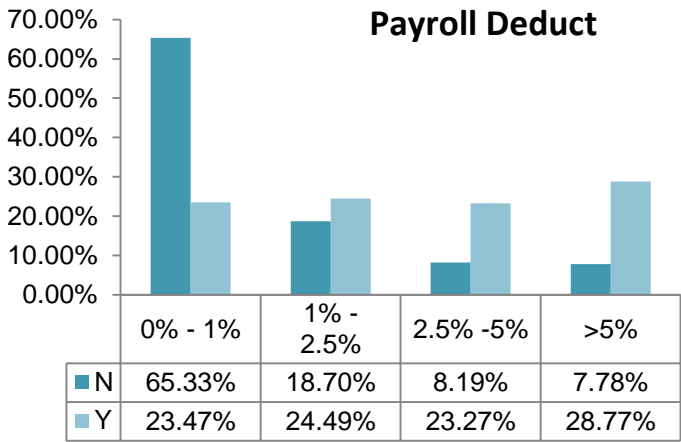
Better Insights with Bayesian (Interaction among all variables)

Description: To test the accuracy of the Bayesian network, we built models with only those variables that were considered in pivot analysis. We then **observed the target variable** i.e. the Participation Bucket **with respect to selected variables**. The following results were obtained:

Pivot Analysis



Bayesian Analysis



Agenda

- What are Probabilistic Graphical Models (PGM)
- PGMs on Sales Ratio Data
- **PGM with respect to CART**

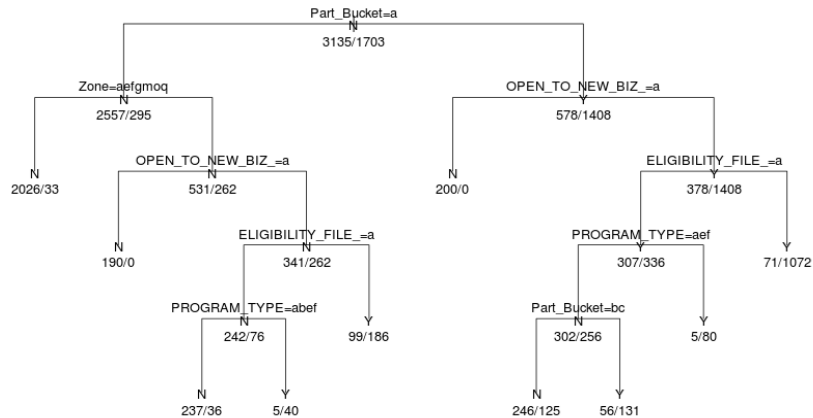
Module 5

Observations On Bayesian Networks and CART

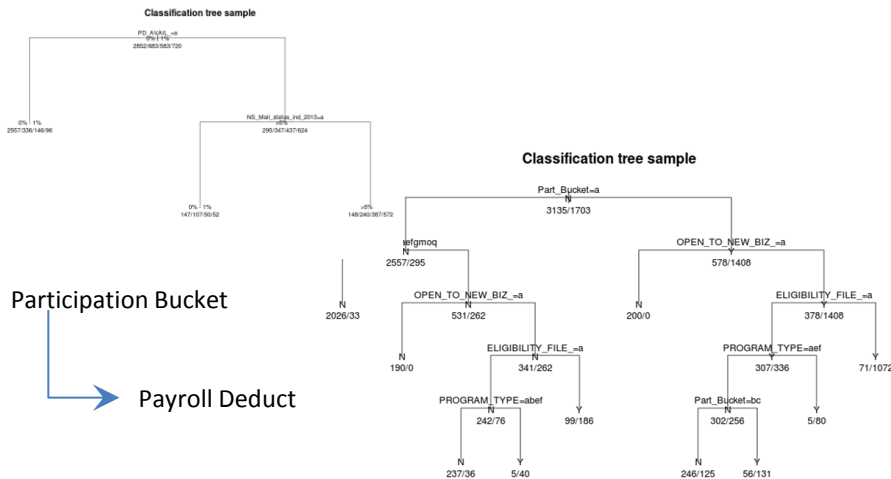
CART is target variable dependent. Problem specific and rigid in terms of modification.

Bayesian **Network** doesn't involve fixed target variable. It is easier to build and modify.

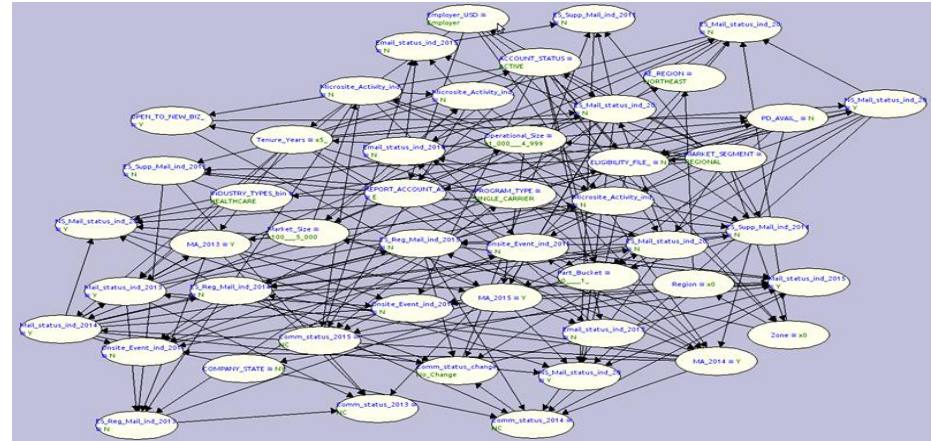
Classification tree sample



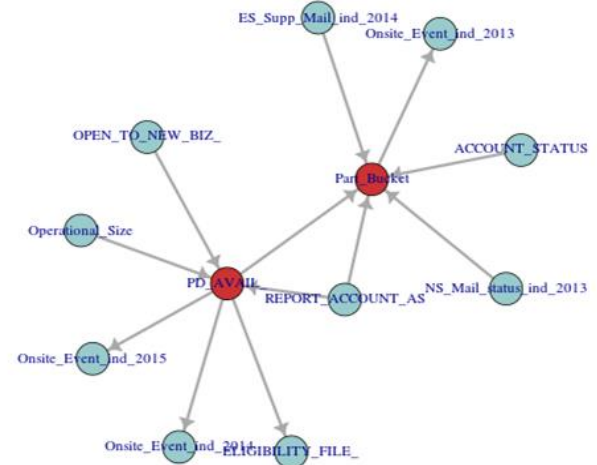
To perform second level analysis we need to construct a new classification tree taking important variables as a target variable.



No variable can be added to the tree forcibly or otherwise if needed.



Second Level of analysis can be done by considering Markov Blanket of the important variable in target's Markov Blanket.

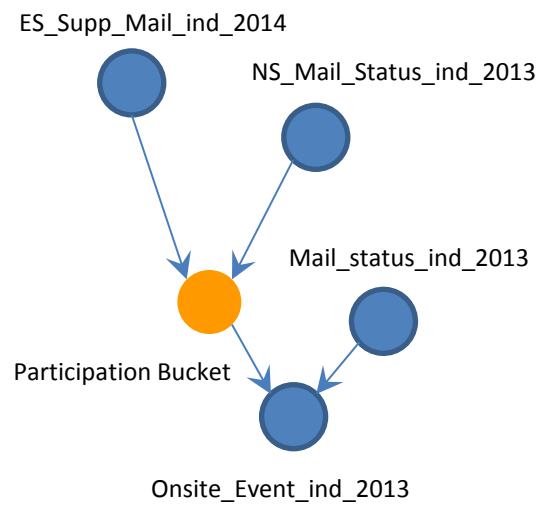
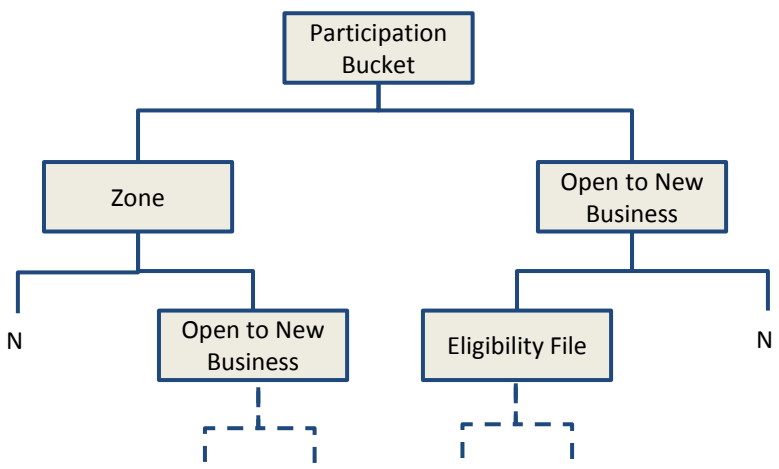


As all variables of data set are already in the model, such need would never arise.

Module 3

Observations On Bayesian Networks and CART

CART (Classification And Regression Tree): For the same data set, we try to see how many business questions can the Bayesian Network and CART models answer.



Questions	CART	Bayesian Network
If we know that our participation is high, then find responsible variables and their values to replicate the result in the future.	Yes	Yes
For my participation bucket to lie in the higher bucket, should I offer payroll deduct or not?	Yes	Yes
I have not done any mailing activity this year. How did this affect my participation percentage?	No	Yes
To check for conditional dependencies like mailing activity not being helpful if the account is already active.	No	Yes
Is the model easily interpretable?	Yes	Yes

Further Steps

- Optimization of Bayesian Network
- CART comparison