

The Patient Re-Activation Problem



Sylvia Morrison & Brian Hoar
Database Marketing

3

Patient Reactivation *aka Lowest-Hanging Fruit*

- Patients First
 - 275+ Chronic Conditions
 - Should return – lapsed patients (not “past”)
 - Message = We Care
 - Priority Appointment Line
 - Tested 36, 24, 18 & 12 months

Patient Reactivation



From the desk of Cynthia DeVina, MD

Our records indicate that you have visited Cleveland Clinic in the past. I wanted to remind you that Cleveland Clinic remains an ideal healthcare resource for you and your family.

Whether you opened an article or you live with a chronic illness that requires regular treatment, Cleveland Clinic has several Community Hospitals and Family Health Centers near your home to better serve your needs. We even offer conveniences like on-site pharmacies, evening and weekend hours, same-day appointments and Express Care to better suit your schedule.

Take advantage of world-class care right in your community by making another visit to Cleveland Clinic. Call 866.939.4169 to schedule your appointment at any location.

Cynthia DeVina, MD
Chair, Cleveland Clinic Family Health Centers

P.S. If you recently scheduled your appointment or visited Cleveland Clinic in the last few weeks, thank you for continuing to trust Cleveland Clinic for your healthcare needs.



Complete care. Convenient locations.

Primary care physicians, Medical specialists, Urgent and emergency care. Cleveland Clinic can meet all your healthcare needs—with easy access, close to your home. Call 866.939.4169 to schedule a visit at any location. Ask about same-day appointments.

For a complete list of services at each location, visit clevelandclinic.org/locations.

Strategy

- Monthly mailing
- Geographic versions – essential
- Offer unnecessary (we tested)
- Physician signature
 - Not *their* physician

Empathy Initiative

- All stems from that perspective
- Marketing, after all, is interested in driving volume
- Objective is to shorten interval between encounters – just like retail



Cleveland Clinic Marketing Database

- Patient Count – over 2 million
 - Example: 77% of population of Chicago
 - Average Contribution Margin \$5,780/per patient
- 3,000+ Staff Physicians
- 14,000+ External Referring Physicians



Reactivation Files

- Patient list is created from over 275 chronic condition codes
- No return visit within 6 months (Diabetes/Congestive Heart Failure)
- No return visit within 12 months (Asthma, Hypertension, High Cholesterol, etc.)
- Hierarchical Extract – patient needs to meet one condition
- Once on the patient mail list – suppressed for 3 months

Requirements

- De-coded files for patient privacy (HIPAA)
- Creating patient records with multiple chronic condition indicators
 - CPT procedure
 - ICD-9 diagnosis codes
- Do not mail & deceased patients are suppressed
- Limit response time to 6 months
 - Compensate for cancelled appointments

The Challenge

- 3 consecutive monthly patient mail files
- Supplied complete mail lists
- Supplied complete list of respondents
- Patient Encounter data collected
 - Last visit before mail date - list of chronic conditions, departments seen & demographics
 - First visit after mail date - list of chronic conditions, contribution margin & departments seen

Patient Reactivation Criteria

- Generate return appointments from patients with chronic (or multiple) conditions
- Target Audience: Diabetes (6m), CHF (6m), Arthritis, Back Pain, Hypertension, Gout, Eye Exams, PSA, Children, Women's Health, Injectibles
- Patients from 13-36 months, no treatment code visit in last 12 months

2012 Reactivation Results

- 310,134 patients mailed, 54,426 (18%) were successfully reactivated
- 7% incremental lift over control group
- The program generated
 - 152,411 total encounters (2.8/patient)
 - \$58 million contribution margin
 - **187:1 payback**
 - \$.50/piece in the mail

What we want to learn

- Are patients coming back for the same chronic conditions
- What chronic conditions return the best ROI
- Should we delete any chronic conditions from the reactivation query
- What trends are happening demographically

Innovation Award Judging



NORC, At the University of Chicago

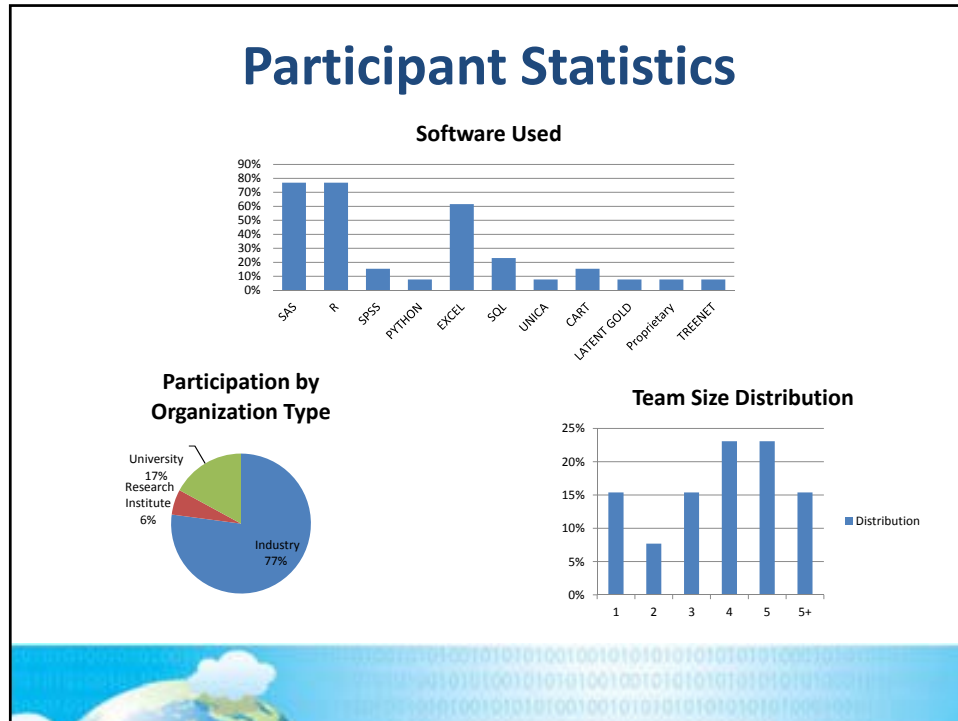
Robert Montgomery

15

Participant Statistics



240 teams from 22 Countries



Observations of Winning Entries

- Everyone used GBM for at least part of the final model.
- As a consequence, everyone used R; two teams R only, others added SAS and some other tools
- All teams but one used an ensemble model of some kind.
- As much or more effort was spent on variable transformation/discovery as on modeling



What Stood Out About Winners?

- Most of the differences were in variable preparation, with a wide variety of techniques being used
- Some variety of techniques in producing ensemble models
- The winner did two things in particular that I think merit the award
 - Using text mining to discover features of the diagnostic codes
 - Modeling the outcome in two stages (response, visit margin) to potentially uncover richer relationships



FINALIST # 1



PricewaterhouseCoopers, Diamond Management &
Technology Consultants

Suman Katragadda
Shreyes Upadhyay

21

Project Team

The team comprises of employees of PwC - DIAC (Analytics) practice from Mumbai and US

- Shreyes Upadhyay
- Rituparna Datta
- Soumya Thakurta
- Suryadip Ghoshal
- Suman Katragadda, PhD

22

Problem Statement

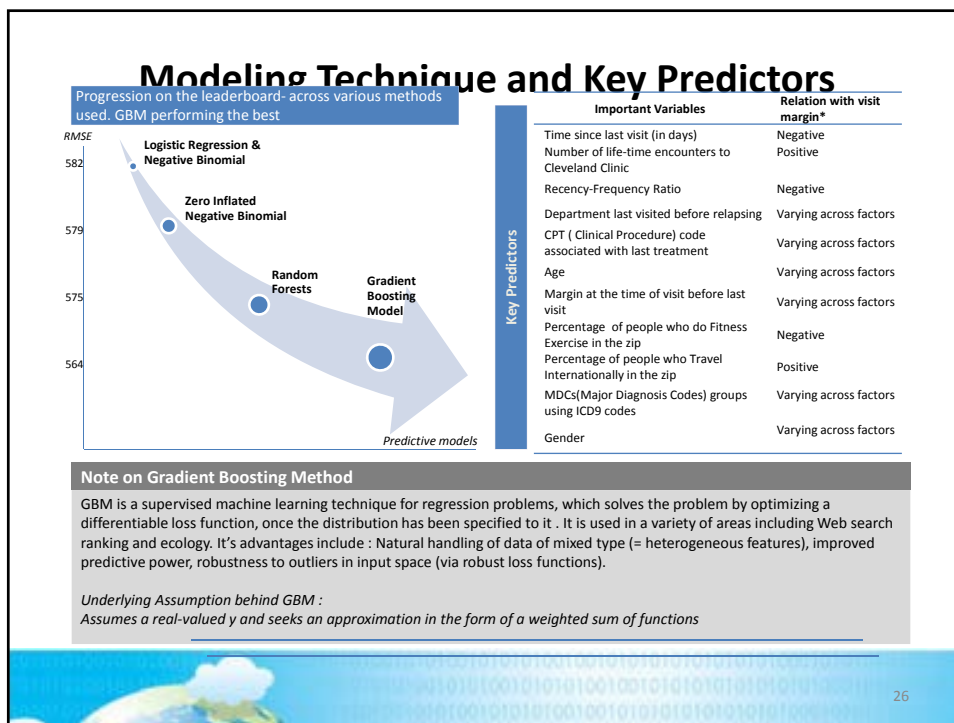
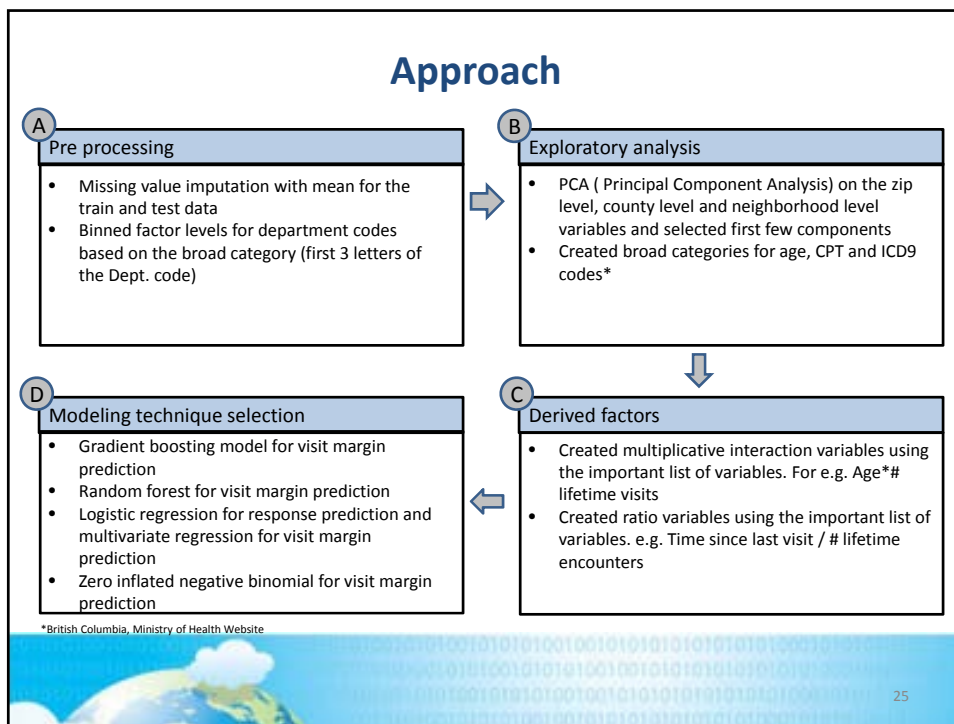
We want to identify the patients who are most likely to respond to a reactivation campaign organized by Cleveland clinic and to predict the margin (revenue – cost) associated with the patients coming back.

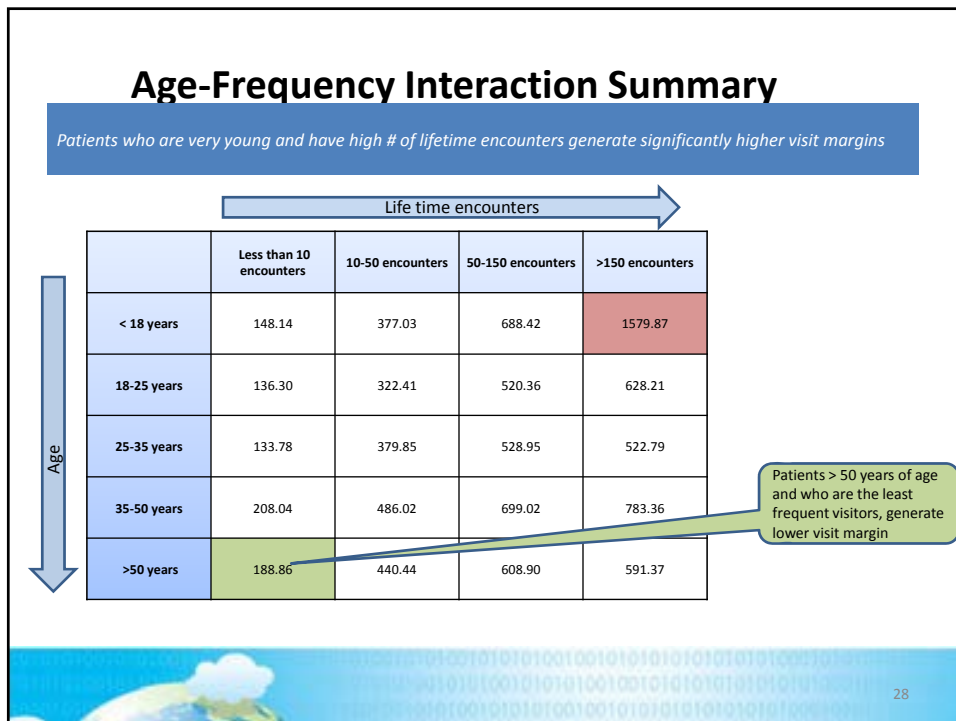
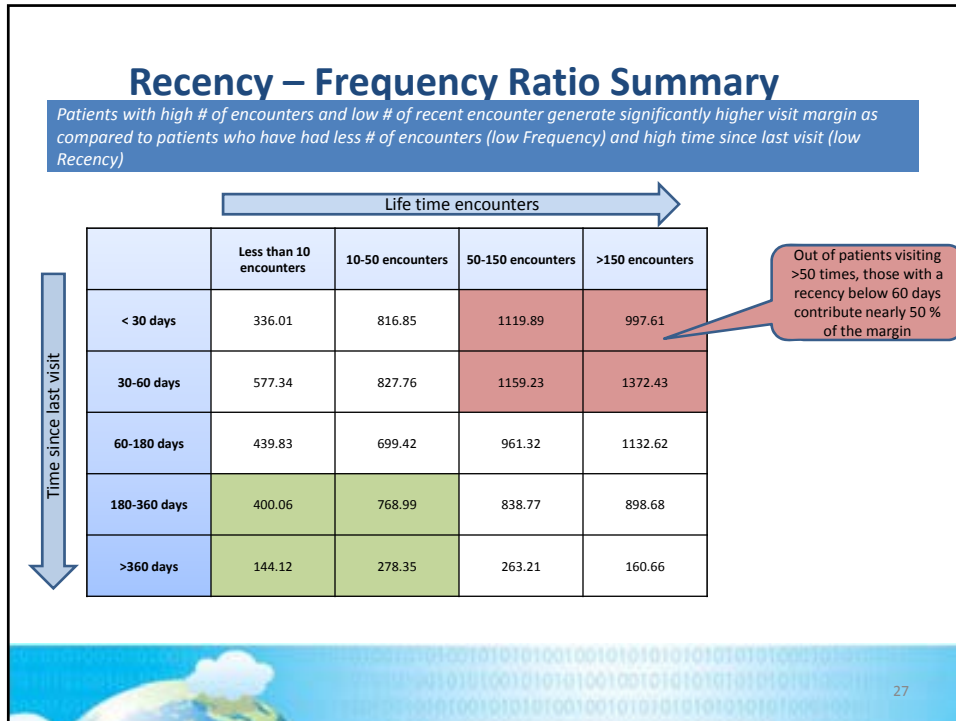
23

Executive Summary

- We tried new methods like machine learning techniques and applied them in traditional set-up ; machine learning techniques worked better in terms of lowering RMSE of the model
- Top 20% of patients who are most likely to respond to reactivation campaign (based on the model) would potentially generate 52% of total visit margin
- Patients with high visit margin are typically
 - Aged men (average age of 58)
 - Recent visitor to the clinic and also frequent visitors (having lower recency/frequency ratio)
 - Belong to regions that are non adherent to fitness regimes (ref slide : #8 , #15)
 - Young patients (<18 yrs) that have greater life time encounters in the past
- Including interaction terms and derived variables (formed out of the key variables) enhanced the predictive power of the model
- Time since last visit (recency) and total # of lifetime visits (frequency) to the clinic are most dominant contributors for predicting visit margin

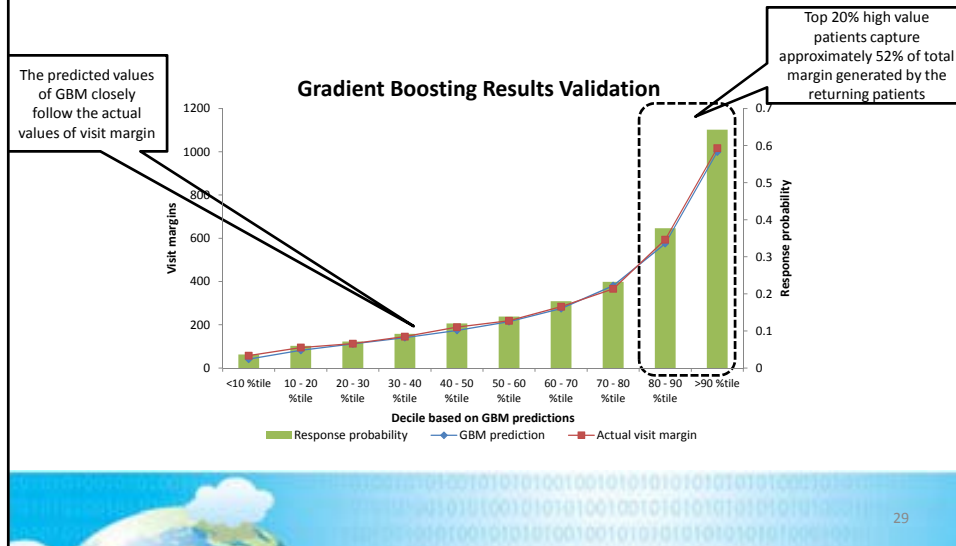
24





Results

Our GBM model accurately identifies the high value patients, in this case the top 20% high value patients, contribute more than 52% of the actual margin generated by the returning patients



Conclusion

What worked for us in terms of techniques and explorations:

- **GBM:** From a technique standpoint, Gradient boosting model which directly predicts “visit margin” gave us the best results in terms of lower RMSE. The model identified the high visit margin patients accurately, which would enable successful execution of a re-activation campaign
- **Principal component analysis:** Reduced dimensionality of datasets : Zip level, county level and Neighborhood level
- **Binning:** Combined factor variables (e.g. CPT codes) based on the mean value of “visit margin” across factor levels
- **Categorization:** Dummy variable creation based on the broad category levels helped to identify category levels highly associated with “visit margin”
- **Interactions:** Creating interaction variables with multiplication and ratio (ratio of “time since last visit” (recency) and “# of life time encounter” (frequency) being the most important interaction variable)

FINALIST # 2

DATA LAB USA

TARGETING BETTER RESULTS

Aaron Davis, VP of Analytics

31

Introduction

- DMA Challenge Philosophies
 - Teamwork!
 - Predictive Modeling Experts With Varying Academic and Professional Backgrounds
 - Iterate
 - Modeling Fosters Business/Data Understanding Leading To Better Models in Future Iterations
 - Ensembling Works
 - Forget the LeaderBoard



Overall Approach

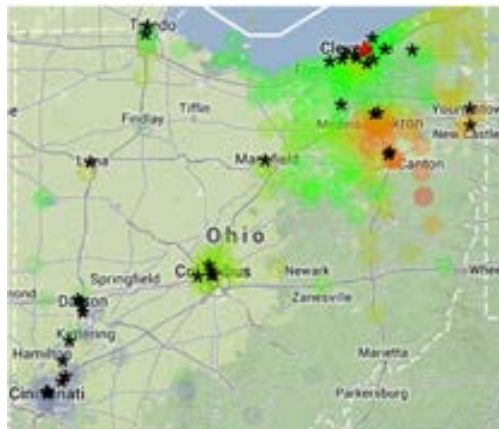
- Phase 1 - Initial Baseline Model
 - Developed in Less than 4 Hours
 - Identifies Areas of Concentration For Future Modeling Iterations
- Phase 2 – Datalab Analytics Challenge 2013
 - 2 weeks
 - Analysts independently developed features, experimented with different approaches and fit models
- Phase 3 – Development of Final Models
 - 1 week
 - Team Members integrated learnings & features from other members
 - Ensembled finalized models

Feature Generation

- Transactional
 - Recency – Time Since Last Visit
 - Frequency - # of Lifetime Encounters
 - Department
 - Procedural Codes From Previous Visit - CPT Codes
 - Disease Codes From Previous Visit - ICD9 Codes
- Demographic –
 - Individual Level – Age/Gender
 - Geo-Demographic
- Geographic

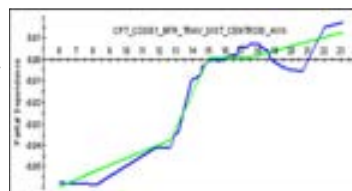
Feature Generation - Geography

- Distance to Treatment Facility
- Distance to Vicinity of Akron General
- Distance to Treatment Facility Vs. Avg Distance of Other Patients w/ same ICD9/CPT Codes
- Lat/Long Coordinates directly into Gradient Boosted Decision Tree Algorithms
- Experimented w/ advanced geosmoothing techniques



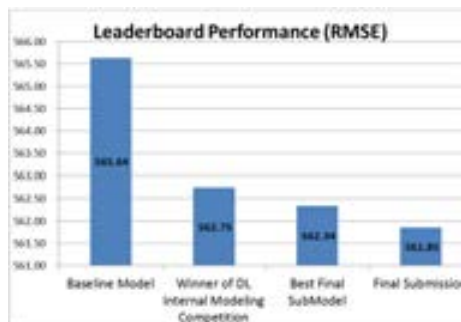
Feature Generation – Disease & Procedural Codes

- Procedural & Disease Codes
 - Very Predictive
 - Very High Order Categorical Variables – 1000's of different values.
 - Impossible to fit relationships directly w/o overfitting.
 - Heavily relied on derived variables that calculated average across ICD9 & CPT Codes:
 - Distance Traveled
 - # of Visits/Time Since Last Visit
 - Age/Gender



Modeling

- Final SubModels
 - Most utilized TreeNet, a gradient boosted decision tree algorithm, for bulk of final predictive models
 - Heavily augmented with internally developed software to identify optimal features & algorithm parameters
 - Final solutions differed heavily in terms of features & parameters
 - Final Models Did Not differ heavily in terms of ultimate performance
- Final MetaModel
 - Straight Average of SubModel Predictions
 - Recalibrated to account for slight skew in Leaderboard data



FINALIST # 3



Fractal Analytics

Rahul Roy Chowdhury

Executive summary

Objective

- Develop predictive model to improve the efficiency of the patient re-activation program for Cleveland Clinic
- Identify and target potential patients who will generate high margin from their return visit to the clinic (Target variable = visit margin)

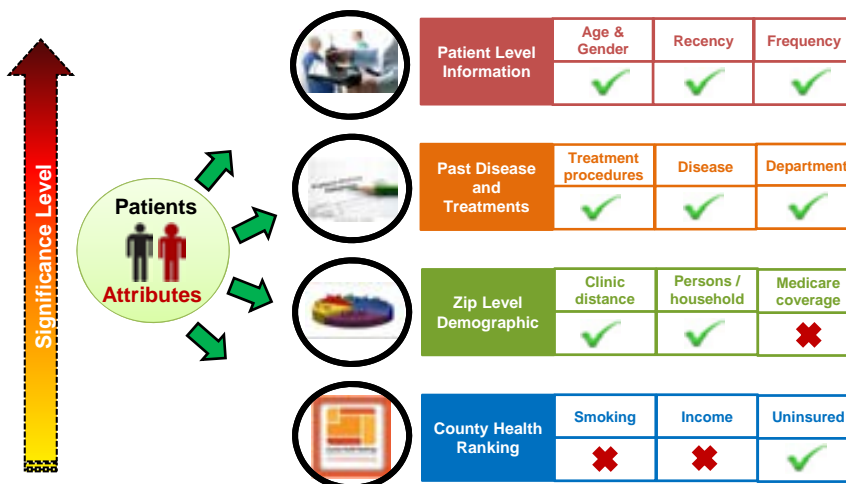
Success criteria

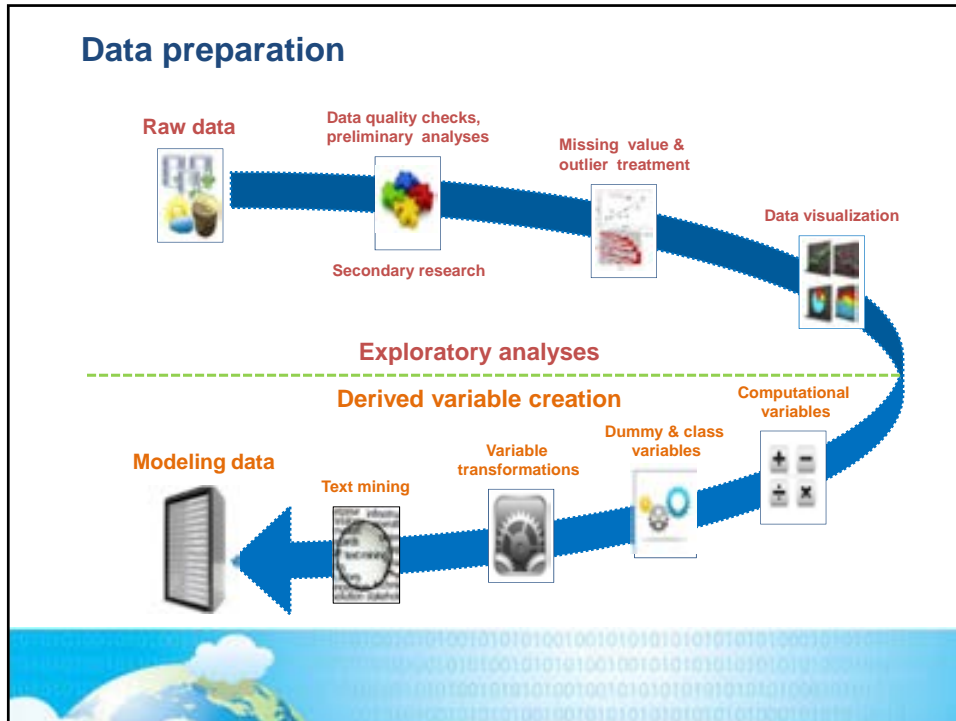
- Minimum error between prediction and actual (Statistic = RMSE)
- Robust model: Evaluation of error on a new set of data

Approach

- Approach 1: Independent model
 - Predict visit margin by a patient post re-activation program
- Approach 2: Two stage model
 - Predict probability that a patient will respond to re-activation program
 - Predict visit margin given a patient responded to re-activation program
 - Expected visit margin is a product of above two predictions
- Final model is an ensemble of approach 1 and 2

Data and insights





Variable creation – key to success

Compute distance from the clinic

Variable transformations

Correlation	Age	1/Age	√Age	Log(Age)	Age^2
Response Rate	-0.35	-0.09	-0.21	-0.07	-0.54

Highest absolute correlation
Inverted u-shape

Text mining

Search for strings in long description of diseases

Nervous system **complications** from surgically implanted device
↪ `complication_ind = 1`

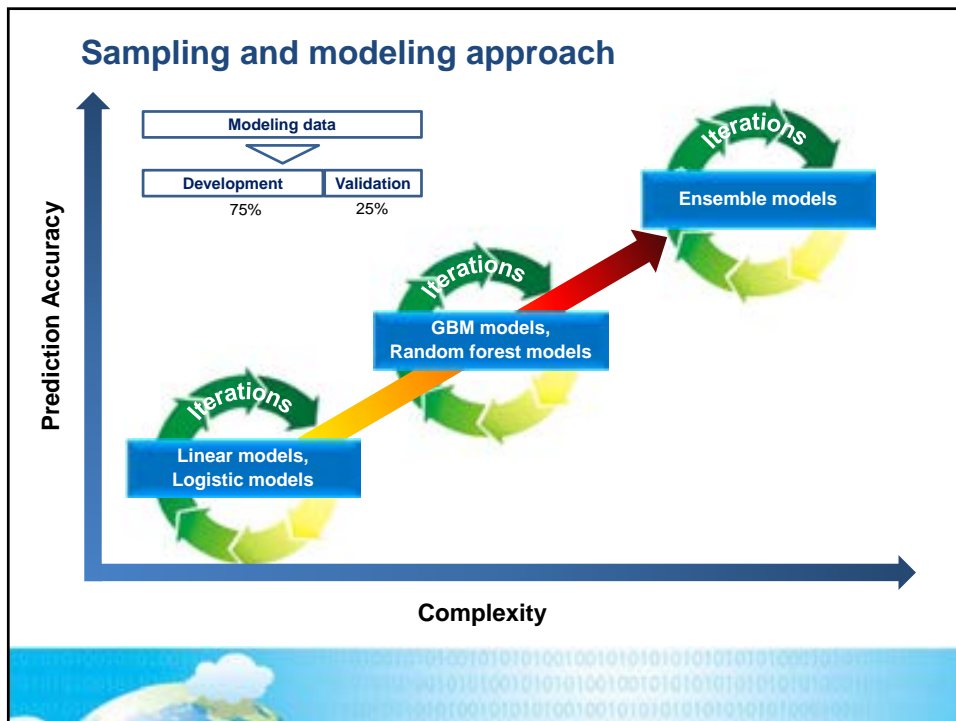
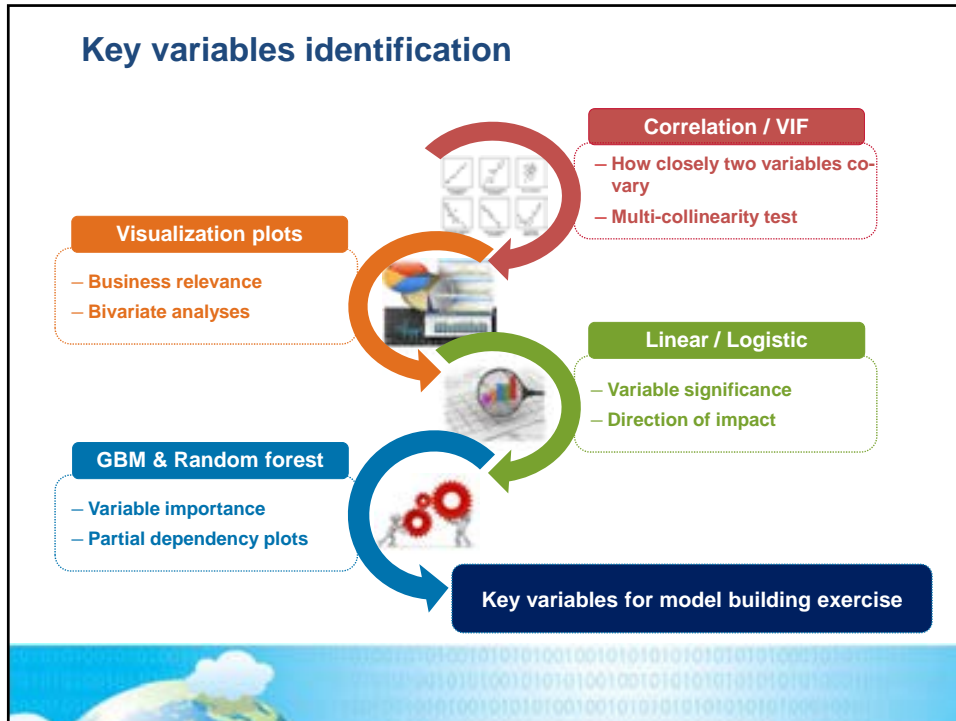
Secondary diabetes mellitus **without** mention of **complication**
↪ `complication_ind = 0`

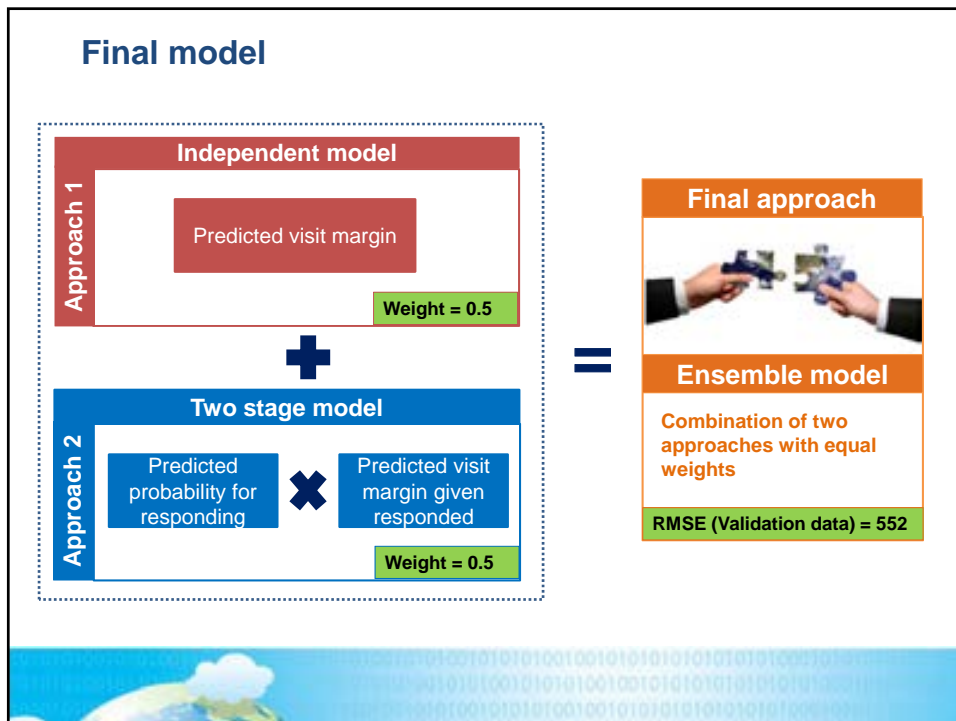
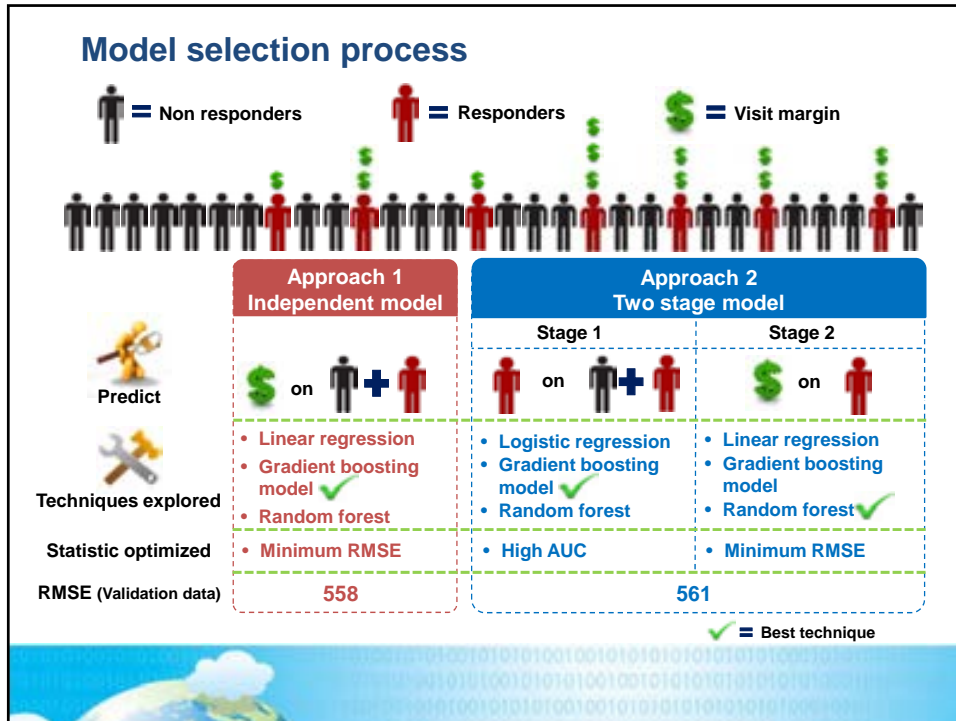
Classification of disease codes



14,568 diseases codes → 19 classifications

- Neoplasms
- Mental disorder
- Injury & poisoning




Source : http://en.wikipedia.org/wiki/List_of_ICD-9_codes







And The Winners Are



47

**You've mastered the theory, now
put it into practice**

Use the Show Guide to match this
session's categories with the
providers in the Hall who can
help execute!
See Page 18 for details

APPENDIX

49

Available Data

We have 75,000 observations for the analysis and 5,000 observations for evaluation

Dependent variables :

- **Responder Indicator:** Indicates whether patient responded to the re-activation campaign. This is not available in the evaluation sample.
- **Visit margin:** Margin associated with visit after re-activation. This is the target variable that needs to be predicted (Values given are transformed to a different scale due to their sensitive nature). This is not available in the evaluation sample.

Other information available (Independent variables):

- Patient Information (treatment history)
- New Neighborhood Life Geo Data (inferred lifestyles provided by Epsilon)
- Zip level demographics
- County Health Ranking Data
- ICD 9 Codes

Note: Our data is unique at patient level, which is captured by Patient ID

50

Variable Transformation

- **ICD grouping** : We mapped the 3,503 unique primary ICD9 codes into 25 MDCs(Major Diagnosis Codes) using our existing mapping list used for other health care projects and created 25 dummies to try in the model
- **CPT grouping**: As we did not have any existing grouping list for CPTs, we grouped it using mean values of visit margin and ind_responder for each CPT code level. For 1,811 unique CPT code, we looked at the range of probabilities for the patients to be reactivated and the average visit margin associated with them and grouped the closest CPTs together into 10 levels
- **Interaction Variables**: We tried different multiplicative interaction variables in the model. For e.g. #of visits*Contribution margin from last visit
- **Ratio Variables**: We tried ratio variables in the model. For e.g. Time since last visit/ Total # of life time Visits

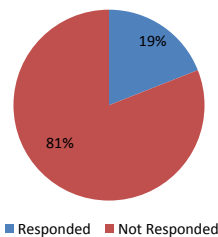
Variables selection:

Used Principal Component Analysis to select important components from the Neighborhood, Zip and County level data and included them in the model

51

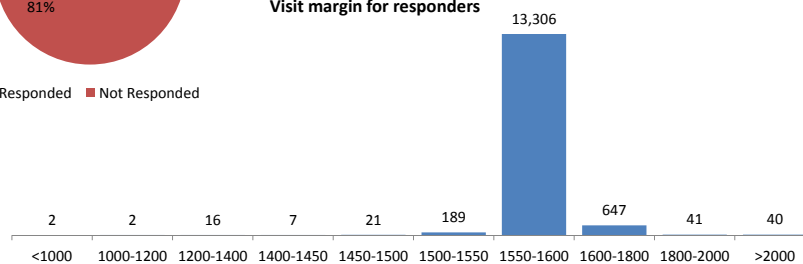
Dependent Variable

Responder Indicator to the re-activation campaign



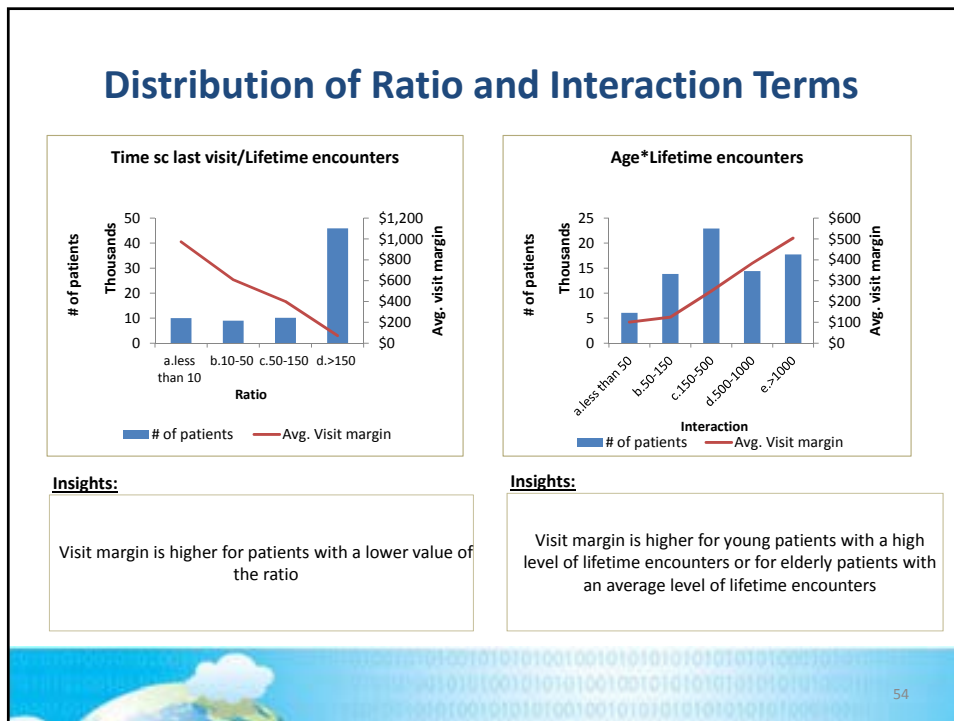
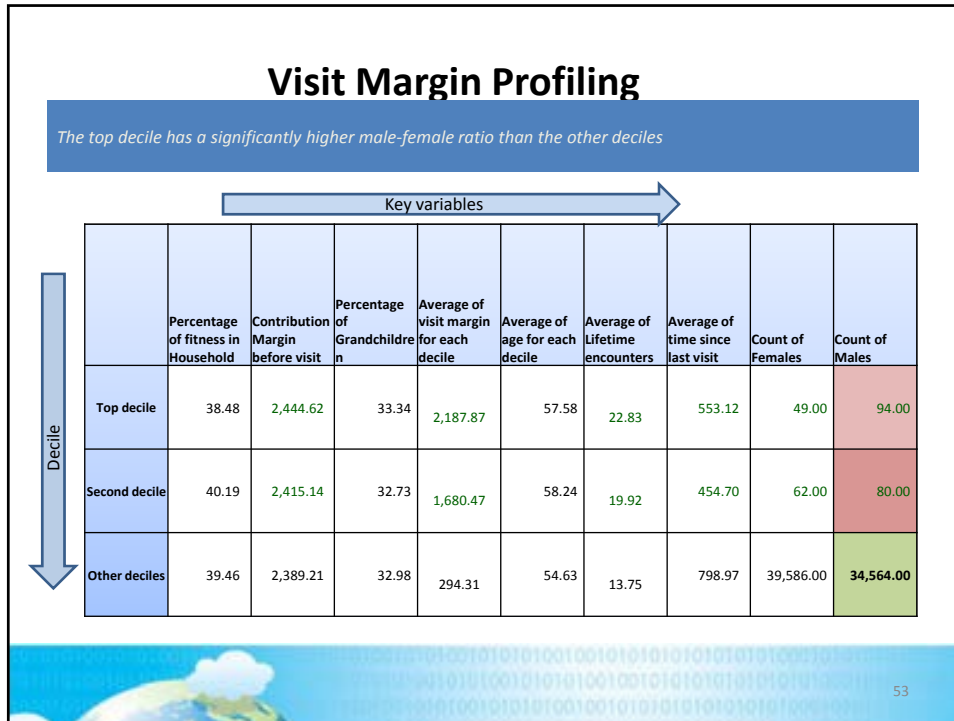
Only 20 % of patients have responded to the reactivation campaign, and most of the responders have a visit margin centered around 1500

Visit margin for responders

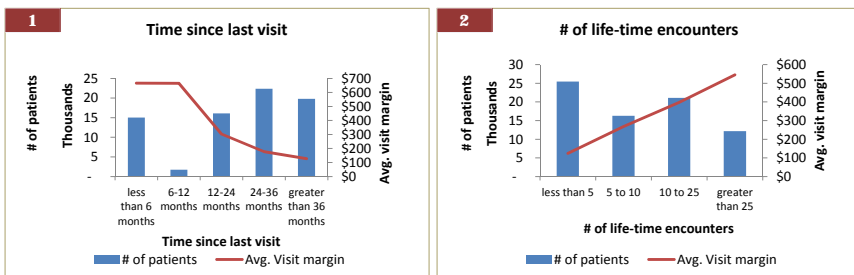


Note: Values of visit margin given are transformed to a different scale due to their sensitive nature

52



Distribution of Important variables(1/5)

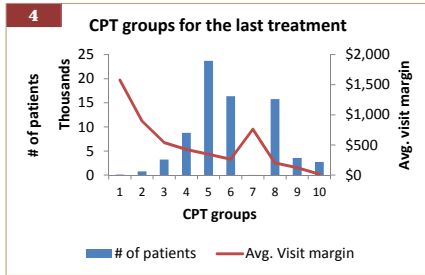
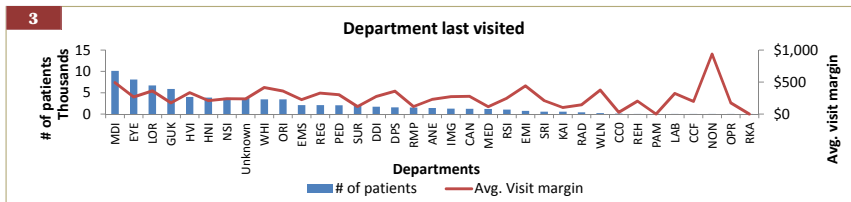


Insights:

Visit margin is higher for the patients who visited last year

Higher the number of lifetime encounters higher is the visit margin

Distribution of Important variables(2/5)

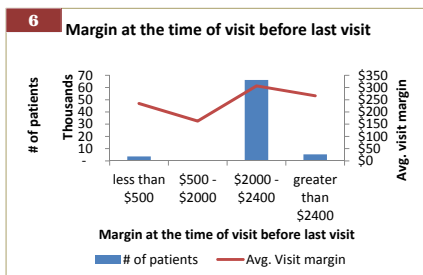
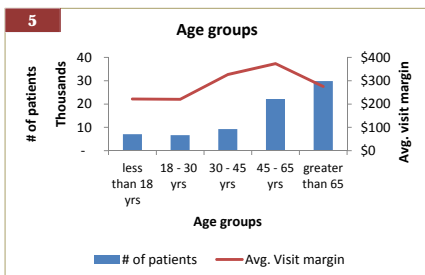


Insights:

of patients and visit margin both are significantly higher for MDI department

Note: Please refer to appendix for CPT grouping

Distribution of Important variables(3/5)

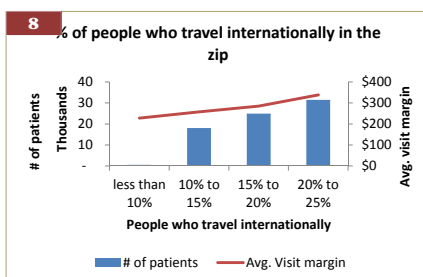
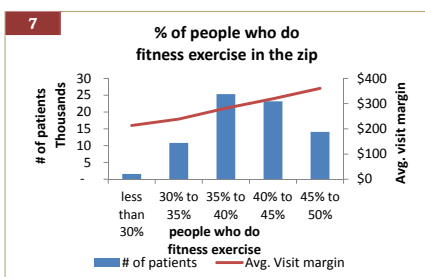


Insights:

Higher the age higher is the visit margin till 65 years of age, beyond that visit margin declines even if the # of visits is more, might be because old people have Medicare

Margin at the time of visit before last visit follows a similar pattern as the visit margin

Distribution of Important variables(4/5)

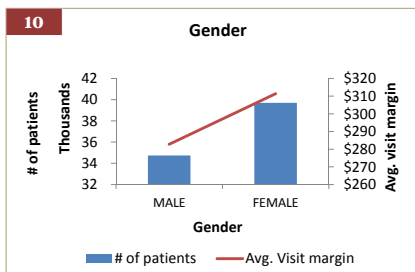
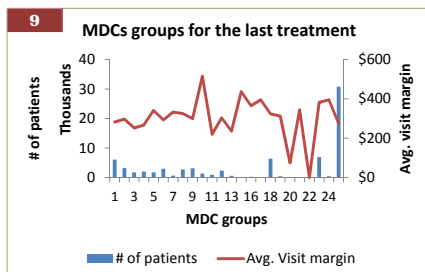


Insights:

More health conscious the patients are, more is their visit margin, may be because they come to the clinic more frequently for regular check-ups

Higher the % of people travelling internationally, higher is the visit margin

Distribution of Important variables(5/5)



Note: Please refer to appendix for MDC grouping

Insights:

Visit margin varies across different MDCs, it being highest for MDC 10, which refers to Endocrine, Nutritional and Metabolic Diseases and Disorders

Visit margin is higher for the females may be because they visit more frequently for regular check-ups

Major Diagnostic Codes (MDCs)

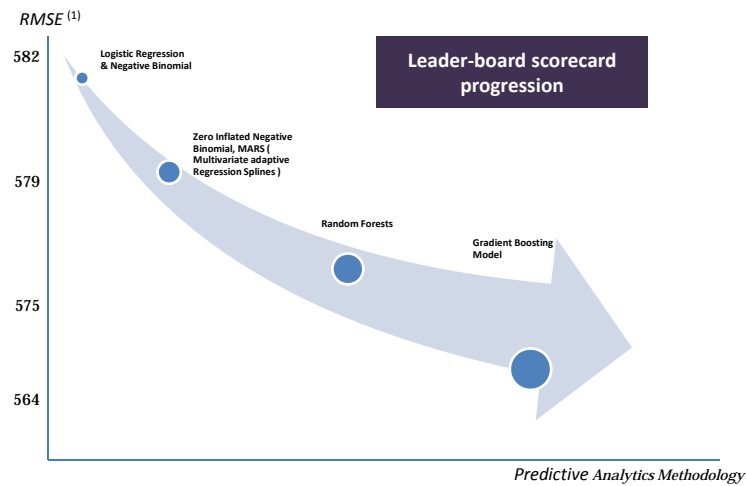
MDC_CAT	MDC_Desc
01	Diseases and Disorders of the Nervous System
02	Diseases and Disorders of the Eye
03	Ear, Nose, Mouth, Throat and Craniofacial Diseases and Disorders
04	Diseases and Disorders of the Respiratory System
05	Diseases and Disorders of the Circulatory System
06	Diseases and Disorders of the Digestive System
07	Diseases and Disorders of the Hepatobiliary System and Pancreas
08	Diseases and Disorders of the Musculoskeletal System and Conn Tissue
09	Diseases and Disorders of the Skin, Subcutaneous Tissue and Breast
10	Endocrine, Nutritional and Metabolic Diseases and Disorders
11	Diseases and Disorders of the Kidney and Urinary Tract
12	Diseases and Disorders of the Male Reproductive System
13	Diseases and Disorders of the Female Reproductive System
14	Pregnancy, Childbirth and the Puerperium
15	Newborns and Other Neonates with Conditions Originating in the PerinatalPeriod
16	Diseases and Disorders of Blood, Blood Forming Organs, and ImmunologicalDisorders
17	Lymphatic, Hematopoietic, Other Malignancies, Chemotherapy and Radiotherapy
18	Infectious and Parasitic Diseases, Systemic or Unspecified Sites
19	Mental Diseases and Disorders
20	Alcohol/Drug Use and Alcohol/Drug Induced Organic Mental Disorders
21	Poisonings, Toxic Effects, Other Injuries and Other Complications of Treatment
22	Burns
23	Rehabilitation, Aftercare, Other Factors Influencing Health Status and OtherHealth Service Contacts
24	Human Immunodeficiency Virus Infections
25	Unknown

ICD and CPT Classification

GROUP	CATEGORY	Ind_responder (MEAN)	Visit_margin (MEAN)
ICD	1	0.99751	1581.447
	2	0.531297	846.3202
	3	0.332029	519.3477
	4	0.222035	350.5661
	5	0.159133	251.4056
	6	0.101378	159.9853
	7	0.049313	78.25415
	8	0	0
CPT	0	0.149191	234.8512
	1	1	1578.65
	2	0.533986	845.7462
	3	0.352233	554.1329
	4	0.263273	420.0865
	5	0.214259	337.2756
	6	0.170037	268.13
	7	0.166667	764.2018
	8	0.123506	196.4531
	9	0.081001	127.6476
10	0.000821	1.31115	

61

Gradient Boosting Method Performs Best



(1) Root Mean Squared Error (RMSE) is the average of squared differences between the predicted and actual values. Lower value is better

62

Summary of Methods Used

	Gradient Boosting Method	Negative Binomial	Zero Inflated Negative Binomial	Random Forests
Description	Gradient boosting is also a machine learning technique for regression problems, which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees.	Used instead of Poisson where data is skewed i.e. Variance exceeds Mean	Used where the event rate is very low (high incidence of zero txns/readmissions)	Its an ensemble learning technique that operates by constructing a set of decision trees at the time of training the model and outputting the class that is the mode of the classes output by individual trees.
Application	Used in a variety of areas including Web search ranking and ecology.	Predicting hospital readmission rates, weather predictions	To model manufacturing defects, genetic defects, rare event modeling (airline failures)	Performance Prediction Challenge and in bio informatics
Advantage	Natural handling of data of mixed type (= heterogeneous features), Predictive power, Robustness to outliers in input space (via robust loss functions)	Effective in managing data with long tail	Isolates cases where event can never occur and models them separately ('True Zeroes' Vs 'Excess Zeroes')	It is one of the most accurate learning algorithms available.
Disadvantage	Scalability, due to the sequential nature of boosting it can hardly be parallelized.	Not good at handling very high variance data, Cannot model negative counts	Does not give good results where all zeroes are driven by one reason only	Random forests have been observed to over fit for some datasets with noisy classification/regression tasks.
Assumption	<i>Assumes a real-valued y and seeks an approximation in the form of a weighted sum of functions</i>	<i>Variance is assumed as quadratic function of mean</i>	<i>High frequency of zeroes in the data. Two separate processes, one always generating zero counts and the other generating both zero and nonzero counts</i>	<i>The only assumption that it relies on is that sampling is representative.</i>

63

Cleveland Clinic – Quick Facts

Quick Facts:

- Founded in 1921 by George Washington Crile
- 5.1 million patient visits per year
- 3,000+ physicians & scientists

Company History:

- Cleveland Clinic is a non-profit academic medical center, provides clinical and hospital care and is a leader in research, education and health information.
- The Cleveland Clinic (formally known as the Cleveland Clinic Foundation) is a multispecialty academic medical center located in Cleveland, Ohio, US.
- The Cleveland Clinic is currently regarded as one of the top 4 hospitals in the US as rated by U.S. News & World Report.
- The Cleveland Clinic was ranked #1 in America for cardiac care from 1994 to 2013.
- Cleveland Clinic has 10 regional hospitals in Northeast Ohio, a hospital and family health center in Florida, and a health center in Toronto, Ontario, Canada, a specialty center in Las Vegas, and a hospital in Abu Dhabi.

64