

2017 Predictive Analytics Symposium

Session 21, Building Block of Predictive Analytics

Moderator:

Robert Anders Larson, FSA, MAAA

Presenter:

Richard Xu, FSA, Ph.D.

[SOA Antitrust Compliance Guidelines](#)

[SOA Presentation Disclaimer](#)



Building Block for Predictive Analytics

Richard Xu, PhD, FSA

VP & Actuary

Head of Data Science

RGA

September 2017

Agenda

- What actuary already knew
- What actuary may not know
- Basic models beyond OLS
 - Generalized Linear Model
 - Decision tree
 - Clustering
- What is next then?
- Conclusion

What actuary already know

- Are you familiar with the following terms?

- Ordinary Least Square (OLS)
- Time Series

- Linear regression model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \varepsilon$$
$$= \sum_i \beta_i X_i + \varepsilon = \bar{X}\bar{\beta} + \varepsilon$$

- Y target variable, X_i predictor variable, ε error term/noise
- β_i parameters to be estimated
- Underlying Assumptions for a valid LM
 - Normality, $\varepsilon \sim N(0, \sigma^2)$
 - Homogeneity, Y representative of population, Independence between observations
 - Fixed X , error-free
 - (linearity)

Ordinary Least Squares

➤ Ordinary Least Squares(OLS)

$$\hat{\beta} = \arg \min(RSS) = \arg \min(\sum_i (\hat{y}_i - y_i)^2) = \arg \min(\sum_i (\sum_j \beta_j X_{ij} - y_i)^2)$$

- For a simple regression

$$\hat{\beta}_1 = (\sum x_i y_i - \frac{1}{n} \sum x_i \sum y_i) / (\sum x_i^2 - \frac{1}{n} (\sum x_i)^2), \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

➤ Identical to Maximum likelihood estimator

- More robust and consistent approach

$$\hat{\beta} = \arg \max(L(X, Y, \beta)) = \arg \min(-\ln(L(X, Y, \beta))) = \arg \min(\sum_i (y_i - \hat{y}_i(\mu_i))^2)$$

if normal distribution

➤ Use adj R^2 to compare fitness of models

$$1 = \frac{RSS}{TSS} + \frac{ESS}{TSS} \quad \begin{array}{l} \bullet \text{ portion that has been explained by OLS model} \\ \bullet \text{ portion of TSS for the error} \end{array}$$

$$\text{Define } R^2 = \frac{RSS}{TSS} = 1 - \frac{ESS}{TSS} = \frac{\sum_i (Y_i - \hat{Y}_i)^2}{\sum_i (Y_i - \bar{Y})^2}, \text{ but it is biased}$$

$$\text{Adjusted } R^2 = 1 - \frac{ESS}{TSS} * \frac{n-1}{n-k} = 1 - (1-R^2) * \frac{n-1}{n-k}$$

Why actuary did not use OLS

- Processes are inherently linear, or can be well-approximated by LM
- Effectiveness & Completeness
 - OLS makes very efficient use of the data; good results with relatively small data sets
 - Identical to maximum likelihood estimation
- Easy to understand and communicate
 - theory is well-understood; Results are easy to communicate
- Great! But wait ...
- There are several issues with OLS
 - Validation of assumptions - Normal w/ constant σ^2 , independent, homogeneous
 - Unbounded data, non-negative value
- How about insurance application? Distribution of data, variance structure
 - Binomial for rate (mortality/lapse/UW, etc.), $\sigma^2 \sim r(1-r)$
 - Poisson for claim count, \sim mean
- OLS is not applicable in insurance, but you already know lots about modeling

What actuary may not know

Machine Learning & Statistical Techniques

- Random Forest
- XG-boost machine
- Gradient Boosting
- Ada Boosting
- Support vector machine
- Ensemble method
- Survey Data Analysis
- Genetic Algorithms
- Sentiment Analysis
- Markov chain Monte Carlo (MCMC)
- Optimization Methods
- Feature engineering
- Decision Trees (CART/MARS)
- Neural Networks / Deep learning
- Bayesian Analysis
- Classification/Association
- Analysis of Variance
- Mixed Models
- Categorical Data Analysis
- Multivariate Analysis
- Survival Analysis
- Cluster Analysis (e.g. K-Means)
- Non-Parametric Analysis
- Text mining

PM terminology

Supervised vs. Unsupervised Learning

- Supervised: estimate expected value of Y given values of X .
 - GLM, Cox, CART, MARS, Random Forests, SVM, NN, etc.
- Unsupervised: find interesting patterns amongst X ; no target variable Y
 - Clustering, Correlation / Principal Components / Factor Analysis

Classification vs. Regression

- Classification: to segment observations into 2 or more categories
 - fraud vs. legitimate, lapsed vs. retained, UW class
- Regression: to predict a continuous amount.
 - Dollars of loss for a policy, ultimate size of claim

Parametric vs. Non-Parametric

- Parametric Statistics: probabilistic model of data
 - Poisson Regression(claims count), Gamma (claim amount)
- Non-Parametric Statistics: no probability model specified
 - classification trees, NN

Generalized Linear Model

- Generalized Linear Model(GLM)
 - Major focus of PM in insurance industry
 - Include most distributions related to insurance
 - Great flexibility in variance structure
 - OLS model is a special case of GLM
 - (Relatively) Easy to understand and communicate
 - Multiplicative model intuitive & consistent with insurance practice
- 3 components
 - *Random component*
 - *Systematic component*
 - *link function*

Generalized Linear Model

Random component

Observations Y_1, \dots, Y_n are independent w/ density from the exponential family

$$f_i(y_i; \theta_i, \phi) = \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{a_i(\phi)} + c(y_i, \phi) \right\}$$

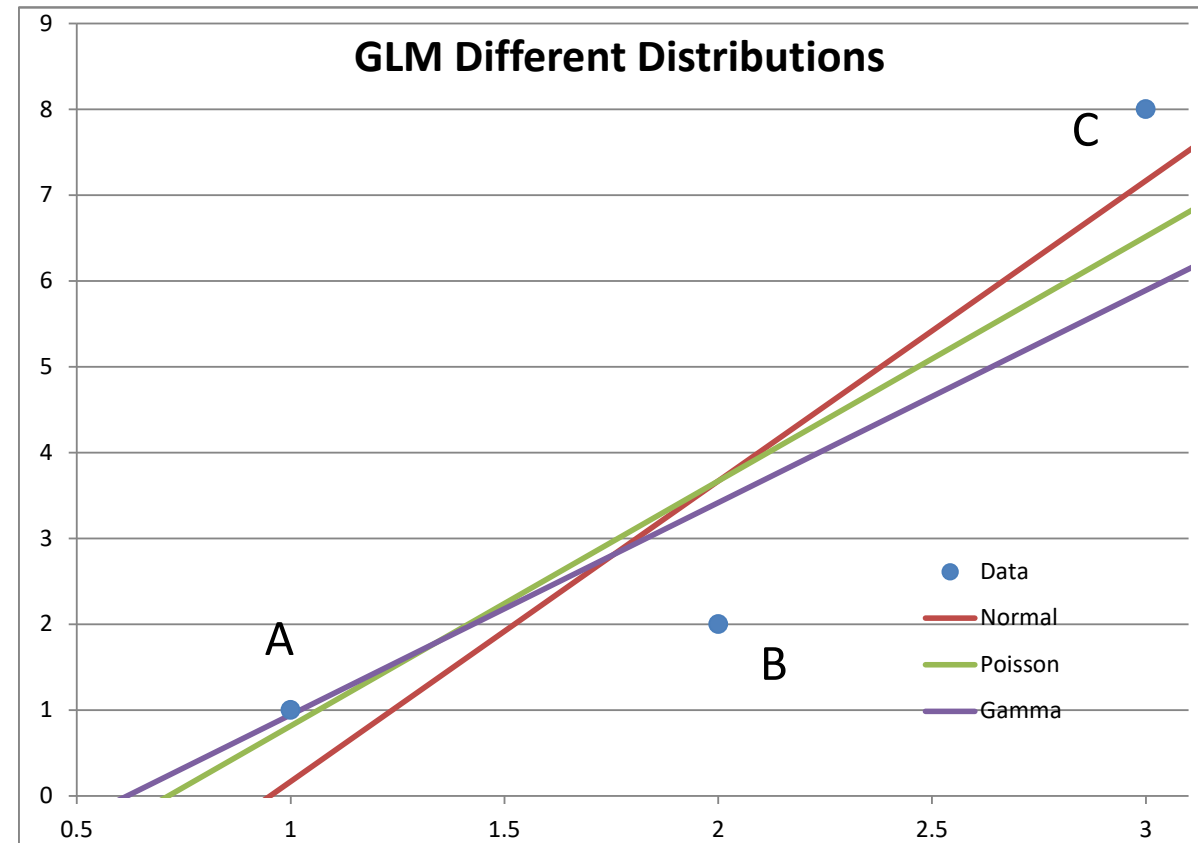
From maximum likelihood theory,

$$E(Y) = \mu = b'(\theta), \quad \text{var}(Y) = b''(\theta)a(\phi) = a(\phi)V(\mu)$$

- Each distribution is specified in terms of mean & variance
- Variance is a function of mean

| | Normal | Poisson | Binomial | Gamma | InverseGaussian |
|---------------|----------------------|----------------|-------------------------|-----------------|---------------------|
| Name | $N(\mu, \sigma^2)$ | $P(\mu)$ | $B(m, \pi)/m$ | $G(\mu, \nu)$ | $IG(\mu, \sigma^2)$ |
| Range | $(-\infty, +\infty)$ | $(0, +\infty)$ | $(0, 1)$ | $(0, +\infty)$ | $(0, +\infty)$ |
| $b(\theta)$ | θ^2 | e^θ | $\ln(1+e^\theta)$ | $-\ln(-\theta)$ | $-(-2\theta)^{1/2}$ |
| $\mu(\theta)$ | θ | e^θ | $e^\theta/(1+e^\theta)$ | $-1/\theta$ | $(-2\theta)^{-1/2}$ |
| $V(\mu)$ | 1 | μ | $\mu(1 - \mu)$ | μ^2 | μ^3 |

Why distribution will affect results



Variance of different distributions

- Gaussian, constant
- Poisson, \sim mean
- Gamma, \sim mean²

Generalized Linear Model

➤ *Systematic* component

A linear predictor $\eta_i = \sum_j x_{ij}\beta_j = X\beta$ for observation i

➤ *link function*

$\eta_i = g(\mu_i)$, random & systematic are connected by a smooth & invertible function

| | Identity | Log | Logit | Reciprocal |
|------------------|----------|----------|---------------------------------|------------|
| $g(\mu_i)$ | x | $\ln(x)$ | $\ln\left(\frac{x}{1-x}\right)$ | $1/x$ |
| $g^{-1}(\eta_i)$ | x | e^x | $\frac{e^x}{1+e^x}$ | $1/x$ |

Log is unique in insurance application s.t. all parameters are multiplicative

- $y = \exp(\sum_j x_{ij}\beta_j) = \prod_j \exp(x_{ij}\beta_j) = \prod_j \exp(\beta_j)^{x_{ij}} = \prod_j f_j^{x_{ij}}$
- Consistent with most insurance practices
- Intuitively easy to understand and communicate

Generalized Linear Model

- Solve for parameters (β) by maximum likelihood
 - Closed form for small data and simple model
 - Iterative numerical techniques for large data set & complex model
 - $\beta_{n+1} = \beta_n - \mathbf{H}^{-1} \cdot s$, similar to Newton's method $x_{n+1} = x_n - f(x_n)/f'(x_n)$
 - Use statistical analysis application, such as *R*
- Compare OSL and GLM

| | Random | Systematic | Link |
|-----|----------------------|---------------------------------|----------------------|
| OLS | Normal only | $\eta_i = \sum_j x_{ij}\beta_j$ | $E(y_i) = \eta_i$ |
| GLM | Various distribution | | $g(E(y_i)) = \eta_i$ |

- Great flexibility
 - Various distribution, variance structure
 - Prior weight and the credibility of data
 - Offset of data

Where we go from here

- More regression models
 - Survival Models (Cox Proportional Hazard)
 - Generalized Additive Models (GAM)
 - Multilevel/Hierarchical Linear Model(HLM)
- Support vector machine
 - Instead of a linear boundary that are affected by all data points to separate classes, an optimal boundary is selected to maximize the gap between classes
- Neural network / Deep Learning
 - Logistic model is the simplest neural network model

Decision Tree Model

- Decision Tree model, or Classification And Regression Tree (CART)
 - ✓ Both classification and regression
 - ✓ Non-parametric approach (no insight in data structure)
- CART tree is generated by repeated partitioning of data set
 - ✓ Data is split into two partitions (binary partition)
 - ✓ Partitions can also be split into sub-partitions (recursive)
 - ✓ Until data in end node(leaf) is homogeneous (more or less)
- Results are very intuitive
 - ✓ Identify specific groups that deviate in target variable
 - ✓ Yet, algorithm is very sophisticated

Recursive Partitioning

- Take all data points
- Consider *all* possible values of *all* variables
- Select the variable/value ($X=t_1$) that produces the greatest “separation” in the target
 - ($X=t_1$) is called a “split”.
- If $X < t_1$ then send the data to the “left”; otherwise, to the “right”
- Repeat same process on these two “nodes”
 - Result is a “tree”; uses *binary* splits
- Stop split data until certain criteria are meet

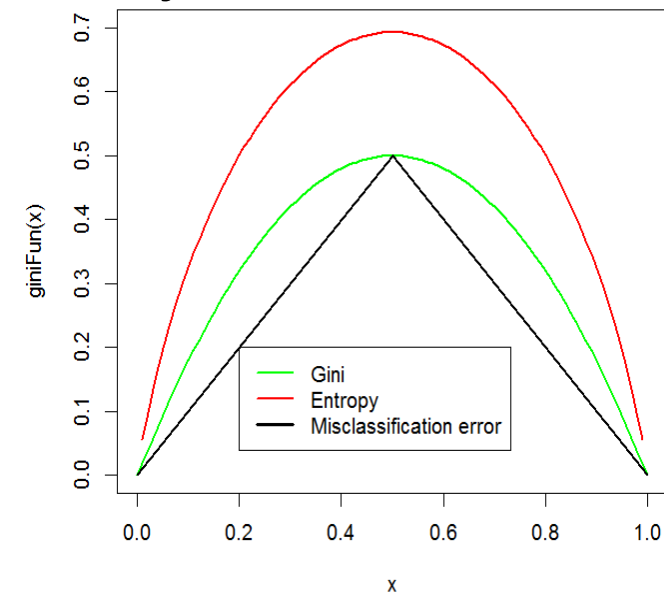
Two Core Questions

- How to find split points
 - Which variable among all, at which value or category, what criterion to use
- When to stop splitting
 - Avoid saturated model

Decision Tree Model

Splitting Point

- “Separation” defined in many ways; different for regression & classification
- Regression Trees: use sum of squared errors
 - $SSE_p = \sum_i (y_i - \mu)^2$
 - $SSE_c = \sum_i (y_i^L - \mu^L)^2 + \sum_i (y_i^R - \mu^R)^2$
 - Select $X=t_1$ such that $\max_{x_i, t} (SSE_p - SSE_c)$
- Classification Trees: use measures of purity/impurity
 - Intuition: an ideal tree model would produce nodes with only either class A or class B - completely pure nodes
 - *Gini Index* - purity of a node $f(p) = p(1 - p)$
 $f(p) = \sum_i p_i (1 - p_i) = 1 - \sum_i p_i^2$, $p_i = \text{freq of class } i$
 - *Entropy* - information index $f(p) = -p \log(p)$
 $f(p) = \sum_i -p_i \log(p_i) = -p \log(p) - (1 - p) \log(1 - p)$



Surrogate Splits

- Problem: for missing data on predictor variable, we don't know how to assign the object
- Solution: we can use a similar split on another variable that is associated (correlated); we use these (surrogate) splits to assign the object to the class
- Missing value can be solved in algorithm level

Greedy Algorithm

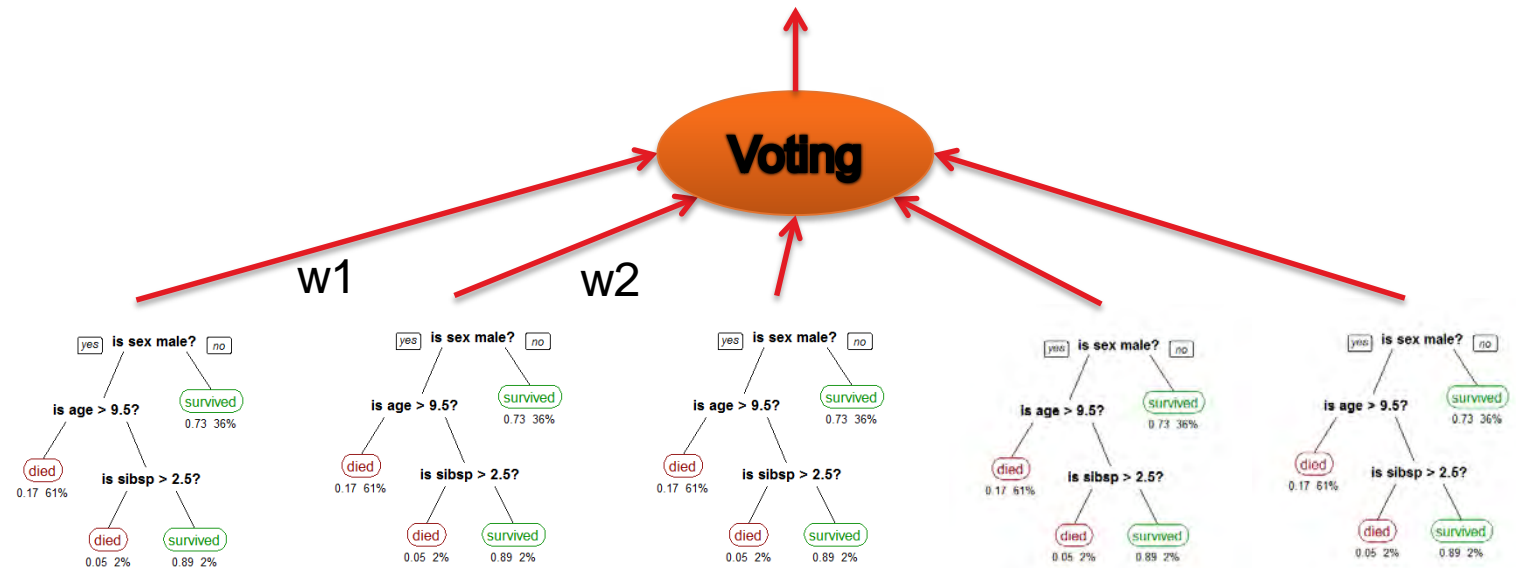
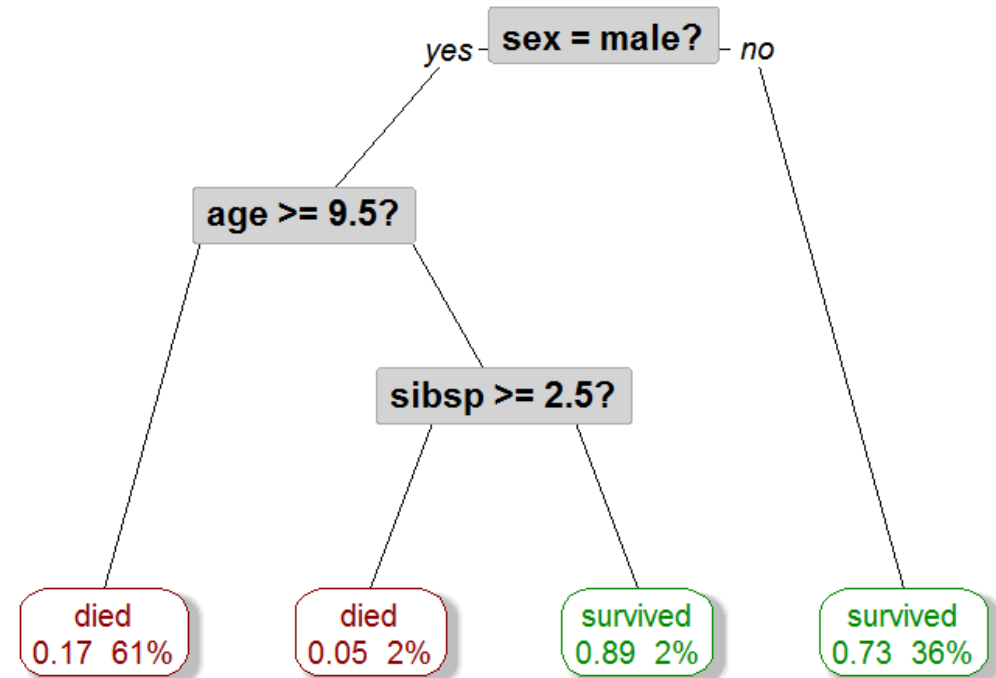
- Define stopping criteria
 - Stop when some minimum node size is reached, or
 - Split only when increase of complexity is larger than the benefits in the overall cost; cross-validation used to balance
- Tree size - balance misclassification error vs. complexity
 - A very large tree may over fit the data
 - A small tree may not structure the important data structure
- Build a very large tree and prune nodes into a small tree

Where we go from here

➤ An example - Titanic survivor model

➤ Decision tree based model

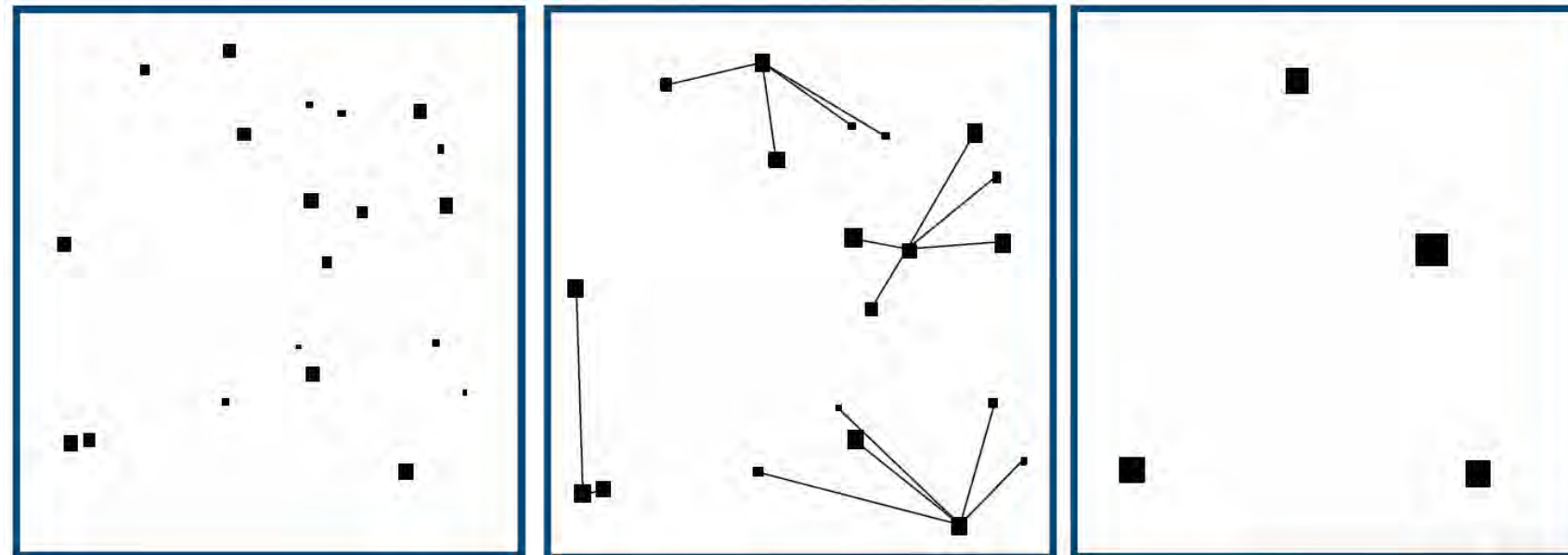
- Random forest
- XG-boost machine
- Gradient Boosting
- Ada Boosting



Data Clustering

Clustering algorithm

- ✓ Find similarities in data according to features in data & group similar objects into clusters
- ✓ Unsurprised (no pre-defined), classification, non-parametric
- ✓ How to measure similarities/dissimilarities, e.g. distance
 - Numeric, categorical, and ordinal variables
- ✓ Partitioning (k-means), Hierarchical, Density-based, etc.



Data Clustering

Algorithm

- Partitioning algorithms - K-means/k-medoids
 - Maintain k clusters with k known; place points into their “nearest” cluster
- Hierarchical (Agglomerative)
 - Objects are more related to nearby objects than to objects farther away; objects are connected by distance; how to define “nearby” object

K-Means Algorithm

1. Select K points as initial centroids, with a given k
2. Repeat
3. Form K clusters by assign each points to its nearest centroid
4. Re-compute the centroids of each cluster
5. Until centroids do not change

Data Clustering

Distance

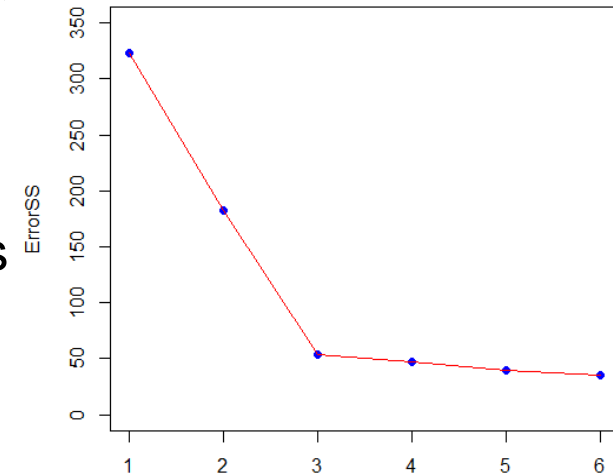
- Euclidean: $d(x_i, x_j) = (\sum_k (x_{ik} - x_{jk})^2)^{1/2}$ ($p=2$), easy to understand, but not scale invariant
- Manhattan: $d(x_i, x_j) = \sum_k |x_{ik} - x_{jk}|$ ($p=1$), city-block distance
- Chebychev: $d(x_i, x_j) = \max_k |x_{ik} - x_{jk}|$ ($p \rightarrow \infty$),
- Minkowski: $d(x_i, x_j) = (\sum_k (x_{ik} - x_{jk})^p)^{1/p}$
- Others like Pearson correlation, Spearman, Canberra, Jaccard, binary, ...

Standardization / Normalization

- Values of variables may have different units
- Variable with high variability/range will dominate metric, & lead to bias

How to determine K

- Business reasons could dictate k
- Try different k, looking at the change in the average distance to centroid, as k increases; error falls rapidly until right k, then changes little



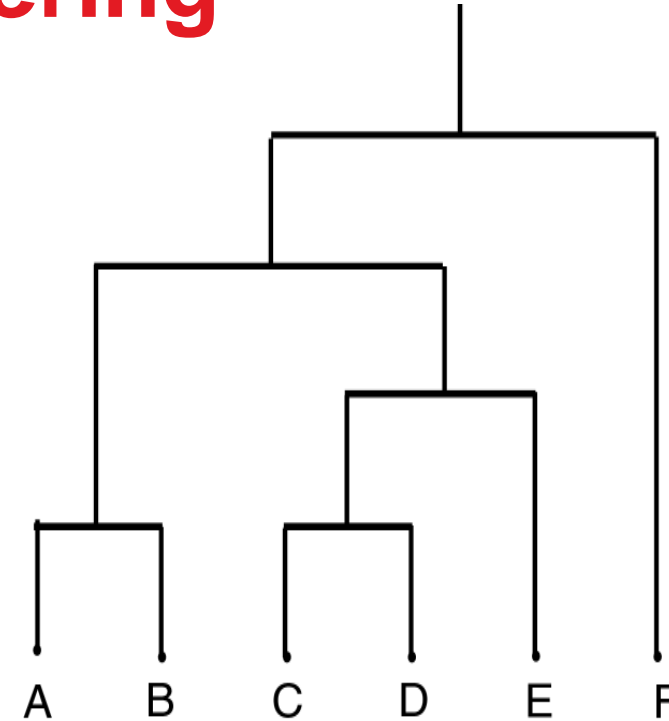
Data Clustering

Comments on K-Means

- Strength: simple, very efficient, & fast
- Weakness
 - Applicable only when *mean* is defined, (categorical?)
 - Need to know *k* in advance
 - Unable to handle noisy data & *outliers*; sensitive to outliers
 - Not suitable for clusters with *non-convex shapes*
 - *Maybe sensitive to initialization*
- There are variants of *k-means*

Hierarchical clustering

- Bottom up (agglomerative) or top down (divisive/deglomerative) produce a dendrogram
- Important questions - how to represent a cluster of more than one point, & how to determine the “nearness” of clusters?
 - **Single Link**: smallest distance between points
 - **Complete Link**: largest distance between points
 - **Average Link**: average distance between points
 - **Centroid**: distance between centroids



Data Clustering

Applications

- Biology – hierarchical classification in structure; genes sequence analysis
- Marketing - distinct groups in customer bases; targeted marketing programs
- Climate- weather patterns in atmosphere & ocean
- Library - book categories in libraries
- Medicine - features & variations of disease; imaging
- WWW - group search results; social network
- Crime analysis: type, location
- City-planning: group houses by type, value, & location
- Earth-quake: cluster epicenters along continent faults
- Data compression: represent a whole cluster by index

What is next then?

- You have solid education background in statistics
- You already have the necessary education components
- Pick up the new skills of data analytics
 - Refresh yourself with the basics of modeling
 - Learn a modeling application / language & practice with examples
 - Attend seminar, conference, training program, etc.
 - Link your new skills with your job & practice if possible

Conclusion

- Advantage of actuary
 - Industry knowledge - domain knowledge is a key in modeling process
 - Expertise in data process - data is always #1 issue in data-driven application
 - Unique position in data analytics
- Opportunity
 - Solid foundation in statistics
 - Education experience in modeling (OLS)
 - Need to pick up new skills & thinking by education, training, and experience
- Actuaries can not miss it
 - Data analytics is here to stay; it is changing insurance industry, and will fundamentally change how we run insurance business
 - Actuaries could and should be on top of it and lead the change



Building Block for Predictive Analytics

Richard Xu, PhD, FSA

VP & Actuary
Head of Data Science
RGA

September 2017