



# Machine Learning & Non-Parametric Methods for Cost Analysis

Karen Mourikas, Nile Hanov, Joe King, Denise Nelson

ICEAA Workshop, June 2018

# Machine Learning Approach to Cost Analysis

## Machine Learning in General

## ML\* Algorithms for Cost Analysis

## ML Applications related to Cost

- Random Forest Prediction
- Latent Semantic Analysis

## Challenges

\* ML = Machine Learning

# Machine Learning Buzz Words

- Big Data
- Smart Manufacturing
- Deep Learning
- NLP (Natural Language Processing)
- IOT (Internet of Things)
- Predictive Analytics
- Neural Networks
- Autoencoders
- Feature Extraction

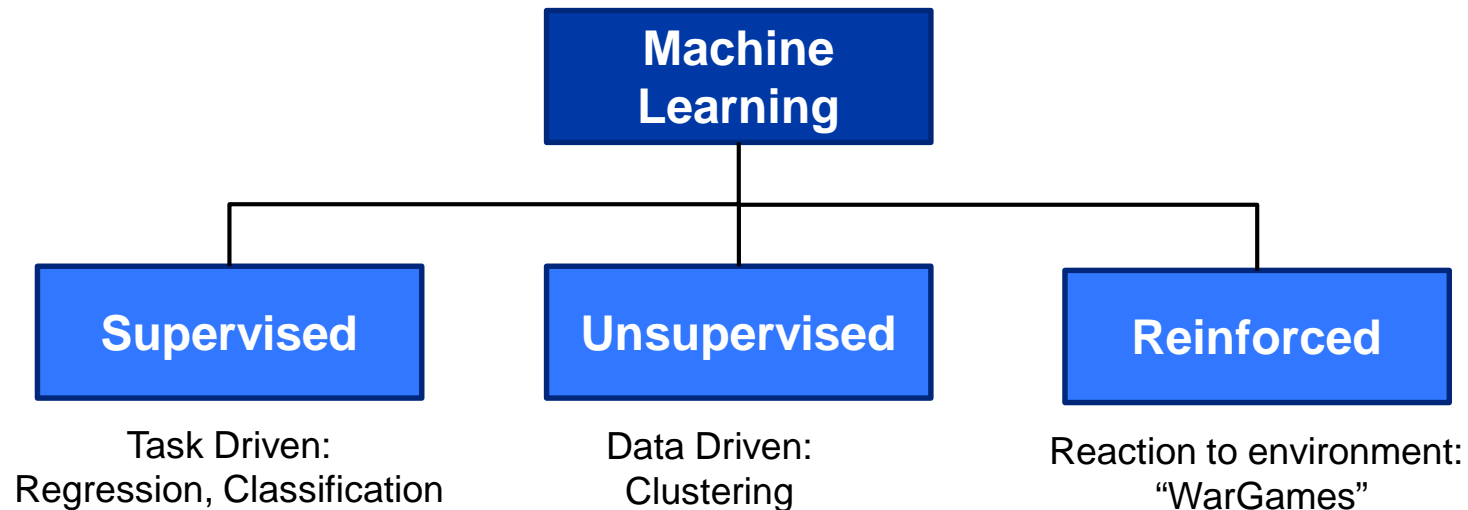
# What is Machine Learning?

Simply,

when a machine mimics "cognitive" functions such as "learning" and "problem solving" \*

Machine Learning (ML) is a method in which algorithms ...

- teach themselves to grow (i.e. learn) from data
- learn without being explicitly programmed



***Machine Learning is a type of Artificial Intelligence***

# What can Machine Learning do?

- Speech recognition
- Autonomous scheduling
- Financial forecasting
- Spam filtering
- Logistics planning
- VLSI layout
- Automatic assembly
- Information extraction
- Market Share Analysis
- Route finding
- Robotics
  - household, surgery, navigation
- Failure prediction
- Fraud detection
- Web search engines
- Autonomous cars
- Energy optimization
- Question answering systems
- Social network analysis
- Medical diagnosis, imaging
- Document summarization

*Many applications for Machine Learning*

# Why is Machine Learning so popular now?

## Machine Learning has been around for a long time

- Has become more popular recently

## Data Explosion

- Much more data available for complex analyses

## Machine Power

- Moore's Law: faster and cheaper computers

## Accuracy of Algorithms

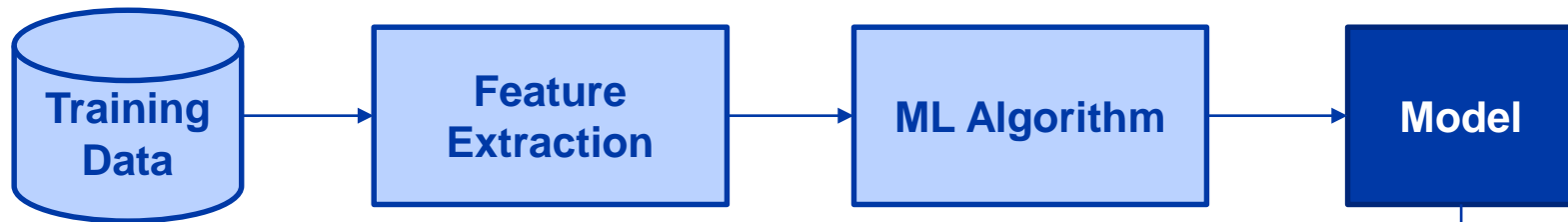
- Reliable enough for usable products

*The Future is Here*

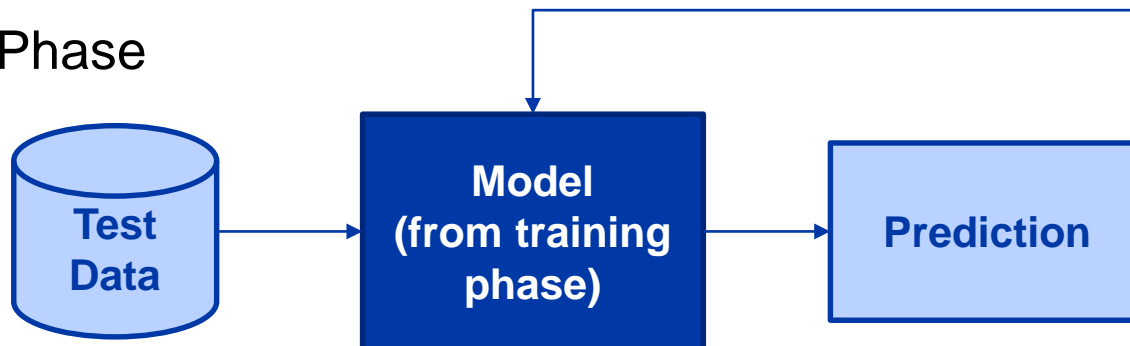
# How does Machine Learning Work?

Typically consists of two stages

- Training phase



- Testing Phase



*General Process*

# Machine Learning Approach to Cost Analysis

## Machine Learning in General

## ML\* Algorithms for Cost Analysis

## ML Applications related to Cost

- Random Forest Prediction
- Latent Semantic Analysis

## Challenges

\* ML = Machine Learning



# Machine Learning for Cost Prediction & Analysis

## Typical Cost Prediction Methods

- Analogies
- Engineering / Bottoms up
- Parametric Equations / Top down

## Machine Learning

- Alternative to traditional cost estimating
- Age of Big Data & Messy Data
- Interactions and non-linear behavior
- Relationship not well understood nor apparent
- Relatively quick & easy to implement

***Could we use Machine Learning techniques for cost prediction?***

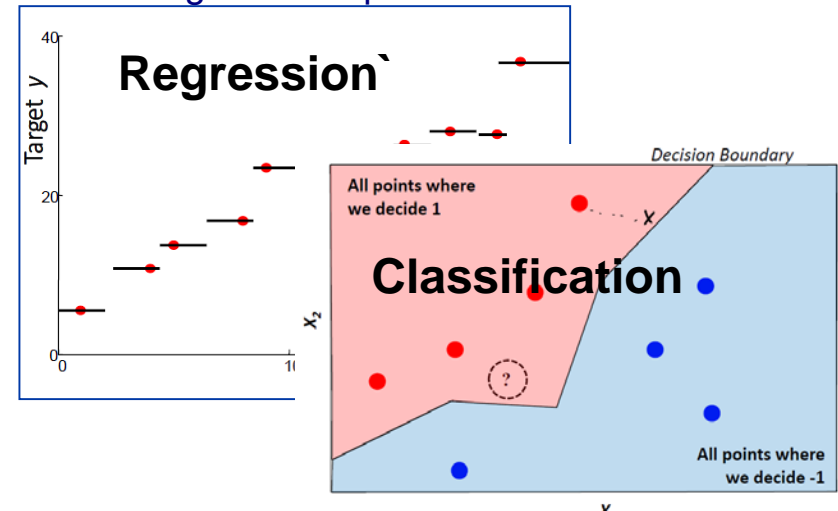
# Supervised Algorithms

## K-Nearest-Neighbors (KNN)

- Clustering approach
- Given new features, finds nearest example and return its value

### Key features

- Regression and Classification



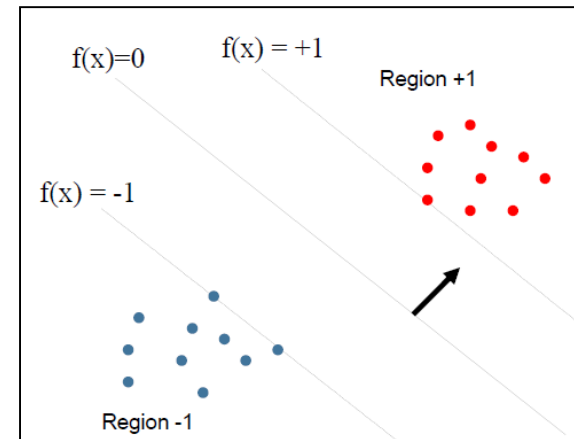
**Fast Classification, Similarity Detection**

## Support Vector Machines (SVM)

- Clustering approach
- Finds the widest margin between classes (boundary decisions)

### Key features

- Able to separate non-linearly- separable regions



**Able to find Optimal Solutions**

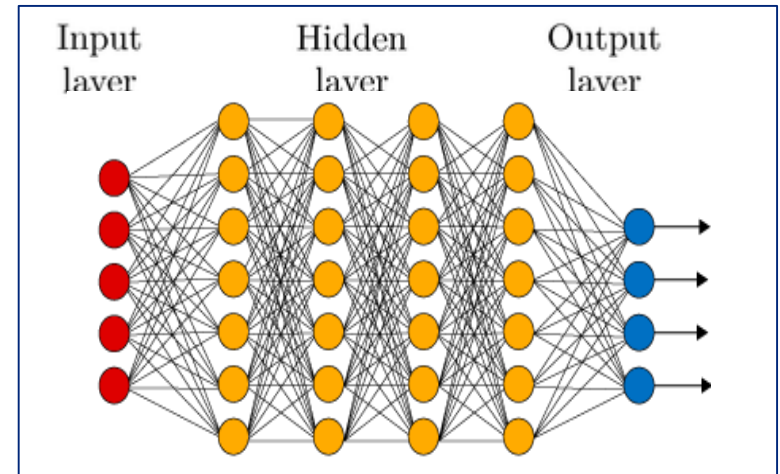
# Supervised Algorithms

## Neural Networks (NN)

- Multi-layer perceptron model
- Finds weights for inputs that optimize the cost function

### Key features

- Very complex shapes/decision boundaries
- Needs a lot of data



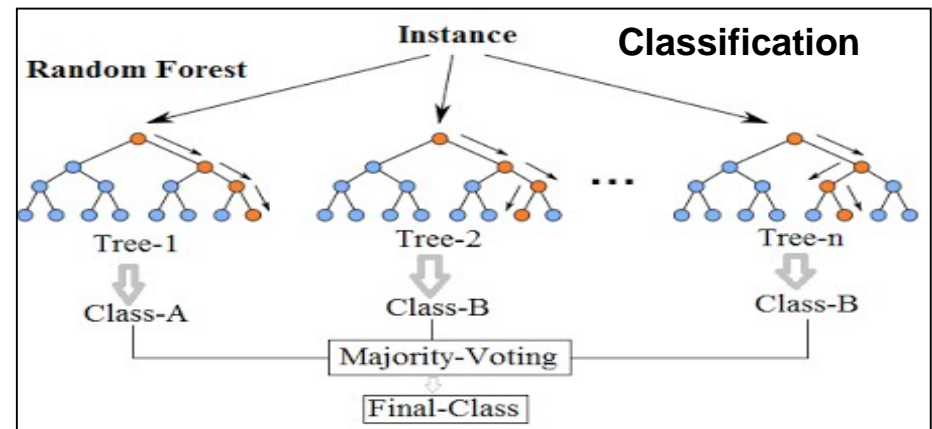
*Finds patterns in large amounts of data*

## Random Forest Prediction

- Decision Tree Ensemble
- Each tree is built from a sample (random) set of features

### Key features

- Training set can be small
- Regression & Classification



*Handles small  $n$ , large  $p$  problems*

# Unsupervised Algorithms

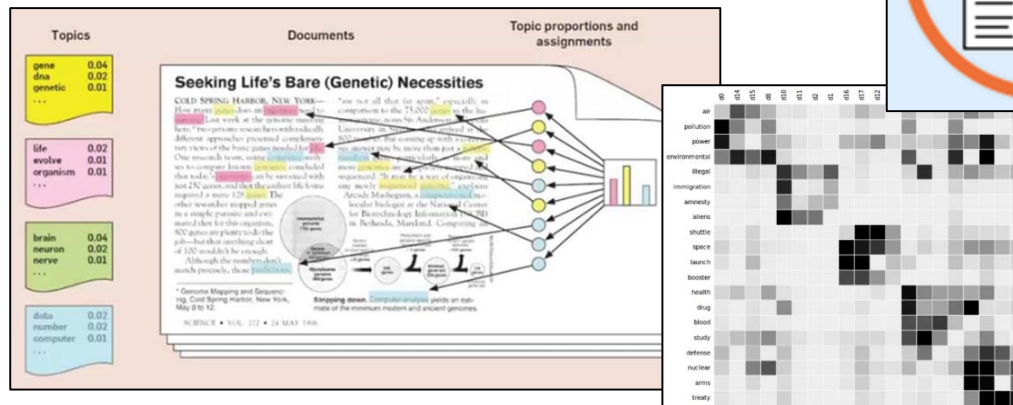
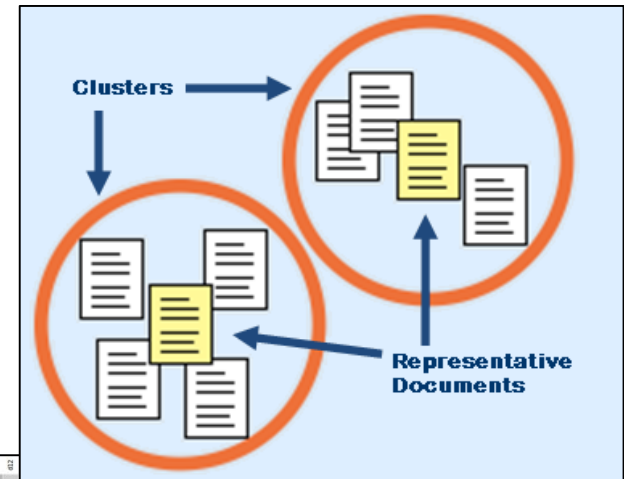
## Natural Language Processing -

### Latent Semantic Analysis (LSA) / Latent Dirichlet Allocation (LDA)

- Document Clustering
- Information retrieval in document groups

### Key features

- Automatic topic detection
- Key term discovery
- Word Clustering



## Automatic Document Grouping

# Machine Learning Approach to Cost Analysis

## Machine Learning in General

## ML Algorithms for Cost Analysis

## ML Applications related to Cost

- **Random Forest Prediction**
- Latent Semantic Analysis

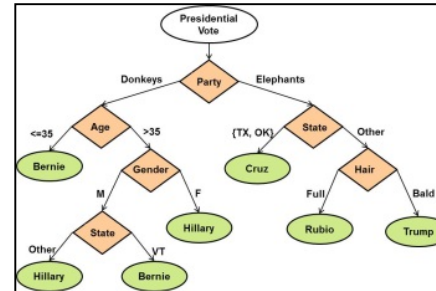
## Challenges

\* ML = Machine Learning

# Trees and Forests

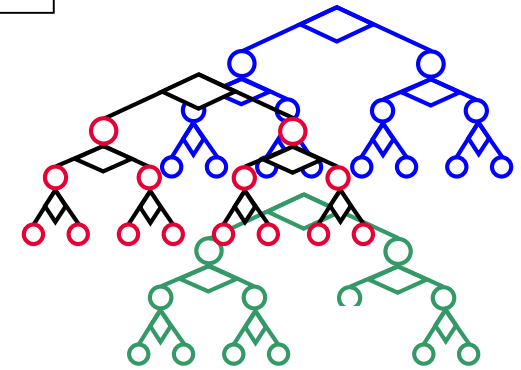
## A Single Decision Tree

- Represents a set of decisions
- Easily interpretable, but ...
- Not a great predictor



## An Ensemble of trees

- Many trees (100s)
- Not as easy to interpret, but ...
- Provides greater prediction accuracy & more stability



## Random Forests

- Ensemble of decision trees “randomly” constructed
- More accurate predictions and reduced error



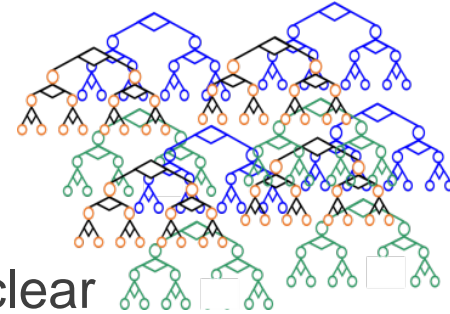
Source: Alexas\_Fotos/Pixabay

***Random Forests Prediction based on Decision Tree Theory***

# Why use Random Forest Prediction?

## Advantages

- Excellent predictors
- Useful if relationship between inputs and outputs is unclear
- Captures non-linear and interaction behavior
- Handles qualitative data as well as missing values
- Relatively stable due to diversity in trees
- Can handle small population size with large number of predictors
- Lower generalization error than other methods
- Runtime very fast, commercial/open source software available



## Disadvantages

- Not so easily interpreted
- Predicts a numeric value (cost) - Not a parametric equation (CER)

***Versatile Black-box Approach***

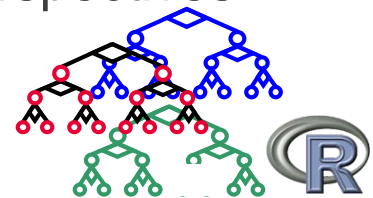
# Application: Logistics Transport Cost Prediction

## Objective

- Predict the shipping cost of products to help determine the best locations to manufacture them

## Analysis Approach

- 1000's of data points, messy, missing values, many potential predictors
- Initial Plan: Multivariate Regression
  - Very cumbersome; required manual partitioning into suitable subsets
- Chosen method: Random Forest Prediction
  - Limited data prep; automatic partitioning / different perspectives
  - Very easy to implement, execute, and analyze



***Random Forest Prediction facilitates logistics transport cost analysis***



# Logistics Transport Cost Prediction Model

## Data Description

- Consists of 150K data points
- Automatically separated into two distinct data sets
  - Domestic with ~ 100K data points
  - International with ~ 50K data points



## Potential Predictors

- Started with 20 potential predictors
- Reduced to 3 key predictors
  - Mode of transportation
  - Origin &/or Destination (country/state)
  - Bill weight



Getty images credits: Mario Gutiérrez – delivery truck; Anucha Sirivisanuwan: barge; hollydc: mailbox; oat autta: cargo truck; JPM: train

***Random Forest Prediction for Big, Messy Data***

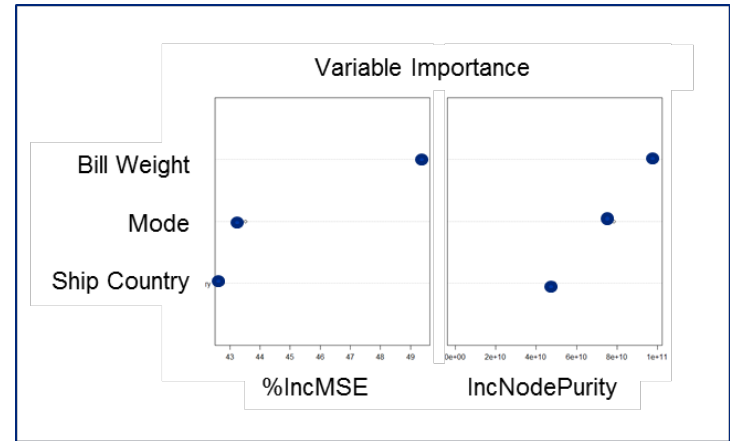
# Analytical Results

## Goodness of fit – Predicted R<sup>2</sup>

- International: 0.83
- Domestic: 0.88

## Graphical Interpretations

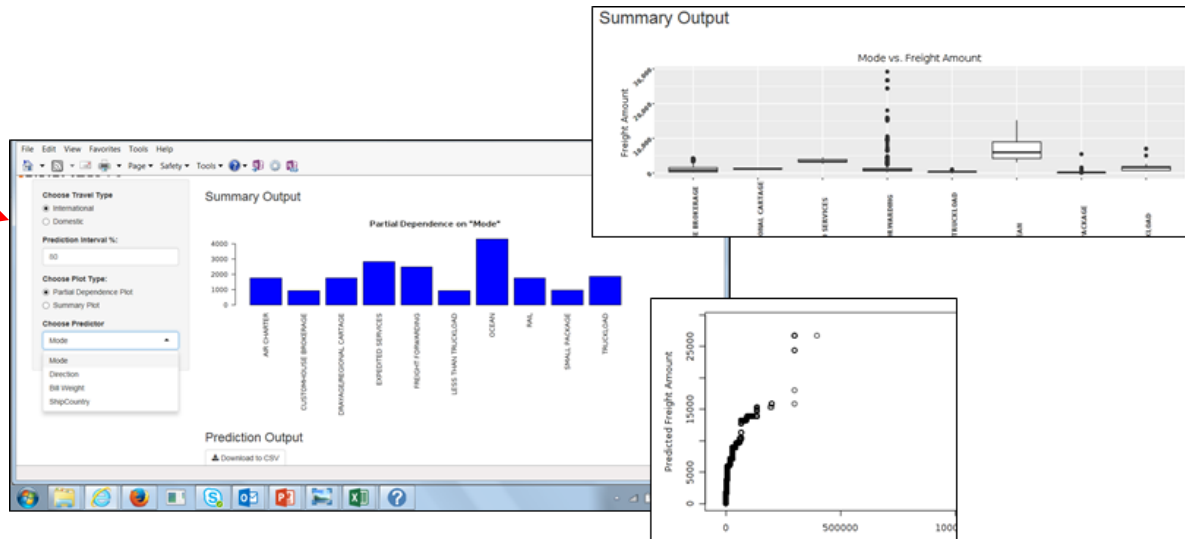
- Quickly produce various charts via R Shiny web-based application



Select Model →

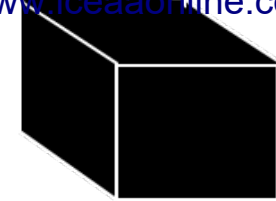
Type of Chart →

Predictor →



**Analysis made easy with R Shiny Package**

# Next Steps: What to do about the ...



## Decision makers want to know what's inside

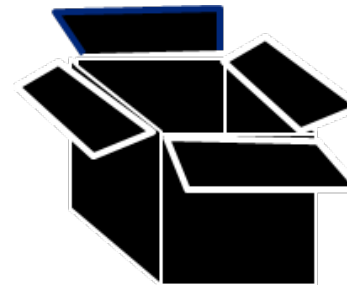
- What can we do?

## Compare results to actuals ...

- Using excel? **Be Careful!**

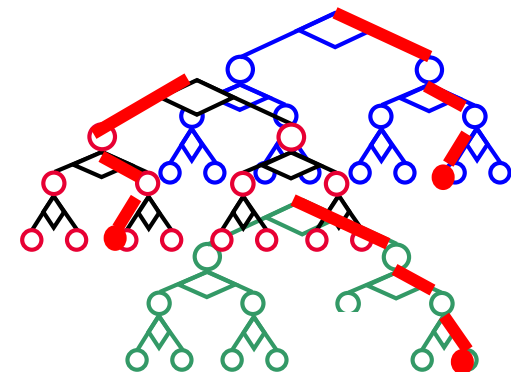
## Develop Interpretation GUI

- R-Shiny to peek inside the black box
- Visualize / Automate standard statistical analyses
- Ability to “play” with the model



## Build algorithm to “create” a CER

- From all the trees, branches, values
- Cost prediction  $\approx f(\text{tree}_i) \quad i = (1..n)$



***Provide ability to “peek” into black box***

# Machine Learning Approach to Cost Analysis

## Machine Learning in General

## ML Algorithms for Cost Analysis

## Applications related to Cost

- Random Forest Prediction
- **Latent Semantic Analysis**

## Challenges

\* ML = Machine Learning

# Application: Analysis of Cost Saving Ideas

## Objective

- Identify best cost savings ideas to apply to other products

## Analysis Approach

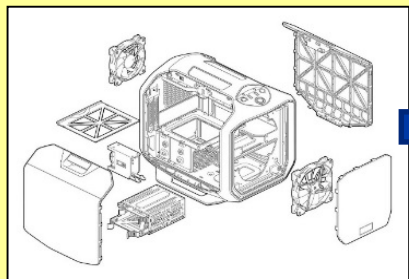
- Collaborative workshops to generate ideas to optimize the product
- 1000's of ideas in free form text from 100's of workshops
  - Could any of these ideas be applicable to other products?
- Natural Language Processing to identify cost-savings ideas for reuse
- Chosen Methods: Latent Semantic Analysis, Latent Dirichlet Allocation
  - Powerful, well-proven, task-invariant algorithms
  - Framework already in place – Open source algorithms

***Natural Language Processing Analyses highlight ideas for reuse***

Presented at the 2018 ICEAA Professional Development & Training Workshop - www.iceaaonline.com

# Generalize Cost Savings Ideas via Text Analytics

## Collaborative Idea Generation



Review Product Detail



Generate Ideas: 10s

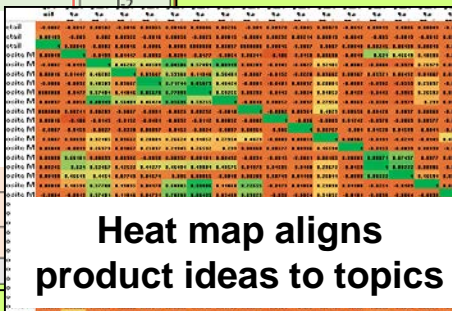
Language	Connectivity	Connectivity Detail	Part Name	Short Description of Idea
Power Distribution	R-25-20A	Electromechanical	Electromechanical	Use LED's connection. Change from copper to stainless connector
Power Distribution	PS30 Panel	Yellow PCB (2005)C	Yellow PCB (2005)C	Reduce space board complexity, use 4-ply boards use low EPOXY material or TCBs.
Power Distribution	PS30 Panel	Servoz in PS30/BA TCBs	Servoz in PS30/BA TCBs	Are the 8-ply boards needed to look PS30/BA TCBs into place? Is cheaper part possible?
Power Distribution	All Panel connectors	Panel Connector replacement to circular	Panel Connector replacement to circular	Use ML, open circular connector instead of Panel. Replace all Panel connectors with circular connectors that are light weight and can be easily disassembled for maintenance. Can this simplify? Check Panel Connector Connector connector lowest
Power Distribution	Connector Plate 10.4	Connector Plate. Plastic Divider between Bus Bars	Connector Plate. Plastic Divider between Bus Bars	Check design for simplified and lower components used accordingly? Specifically the Bus bar dividers

Aggregate ideas from 100s of products  
1000s Unique Ideas

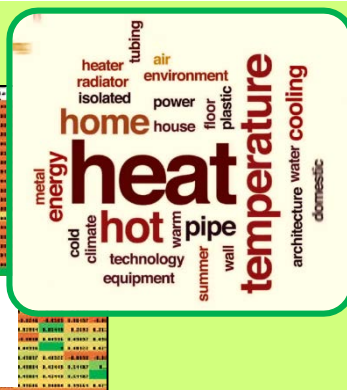
## Machine Learning Analysis

Campaign	Idea	Idea Id
Flight Deck	Alternative Fabric	I-1
	Use off the shelf actuators (rotary actuator)	I-2
Fuel Act	Use of Updated and Lower Cost Transistors	
	Change Material of Solenoid Housing	
Power Distribution		

Group ideas into topics to generalize results



### Identify key terms



Can we identify & apply Ideas from one product to others?

# Similarity Matrices to Align Ideas

## Unstructured Text

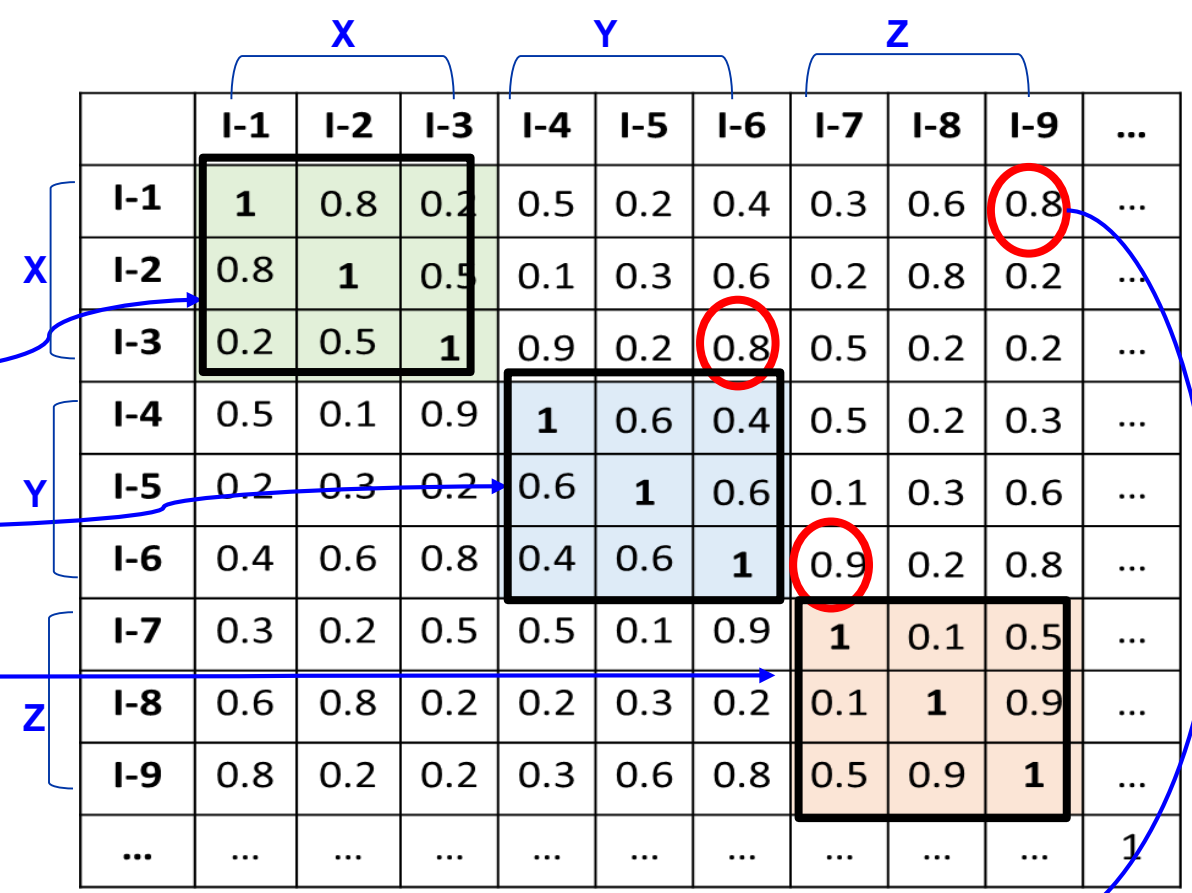
- 100s documents
- 1000s freeform "texts"

Product	Idea	Idea Id
Flight Deck Seats	Alternative Fabric ...	I-1
	Use off the shelf actuators	I-2
	...	I-3
Fuel Valve Actuator	Invest in a design prior to request bids ...	I-4
	Make gears ...	I-5
	...	I-6
Power Distribution	Use of Updated and Lower Cost Transistors	I-7
	Switch ...	I-8
	...	I-9
...	...	...

Product X

Product Y

Product Z

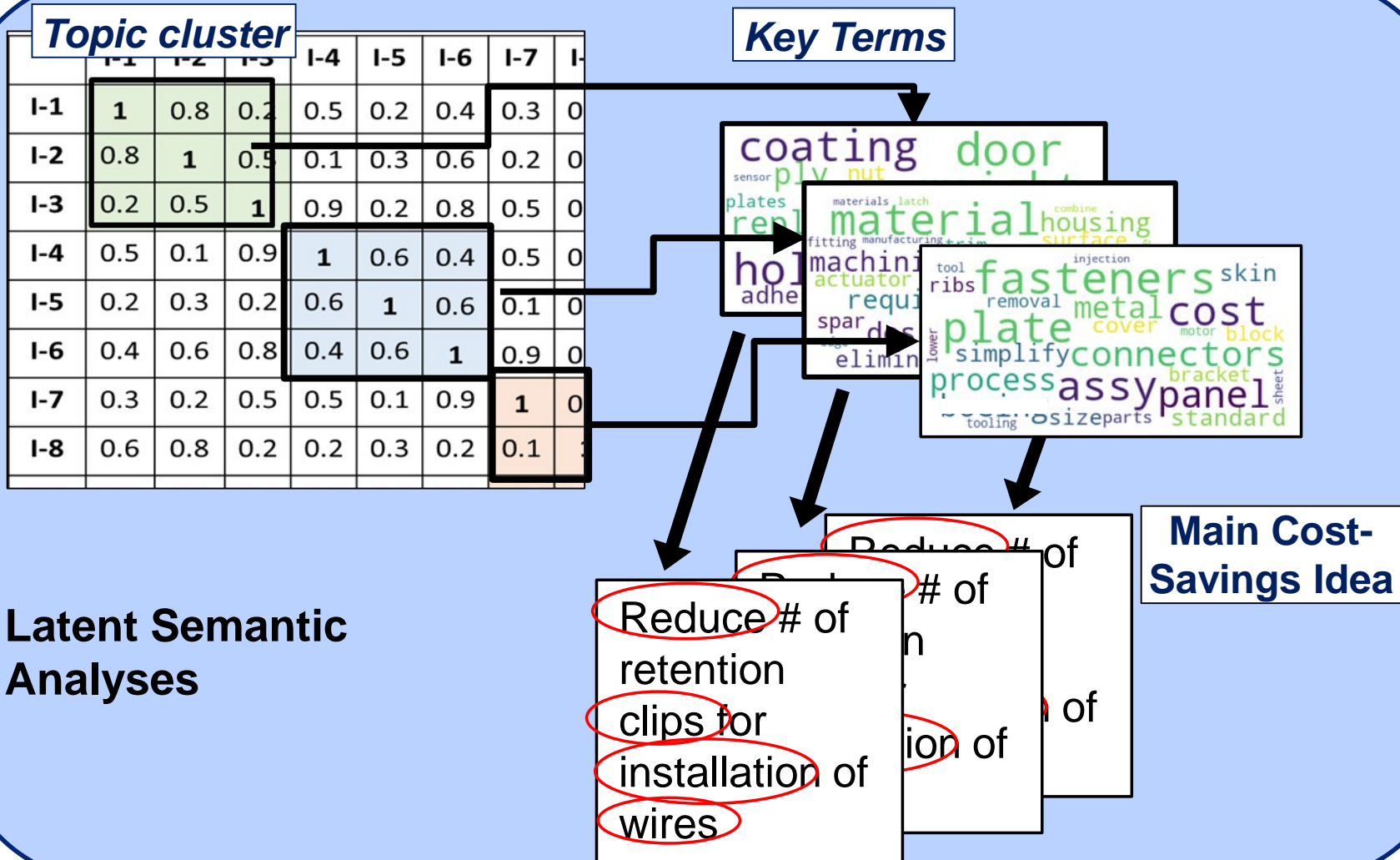


Idea #1 from Product X highly similar to Idea #9 from Product Z

**Cluster similar ideas from unique products via similarity matrices**



# Text Analytics to Identify Reusable Ideas (1 of 2)



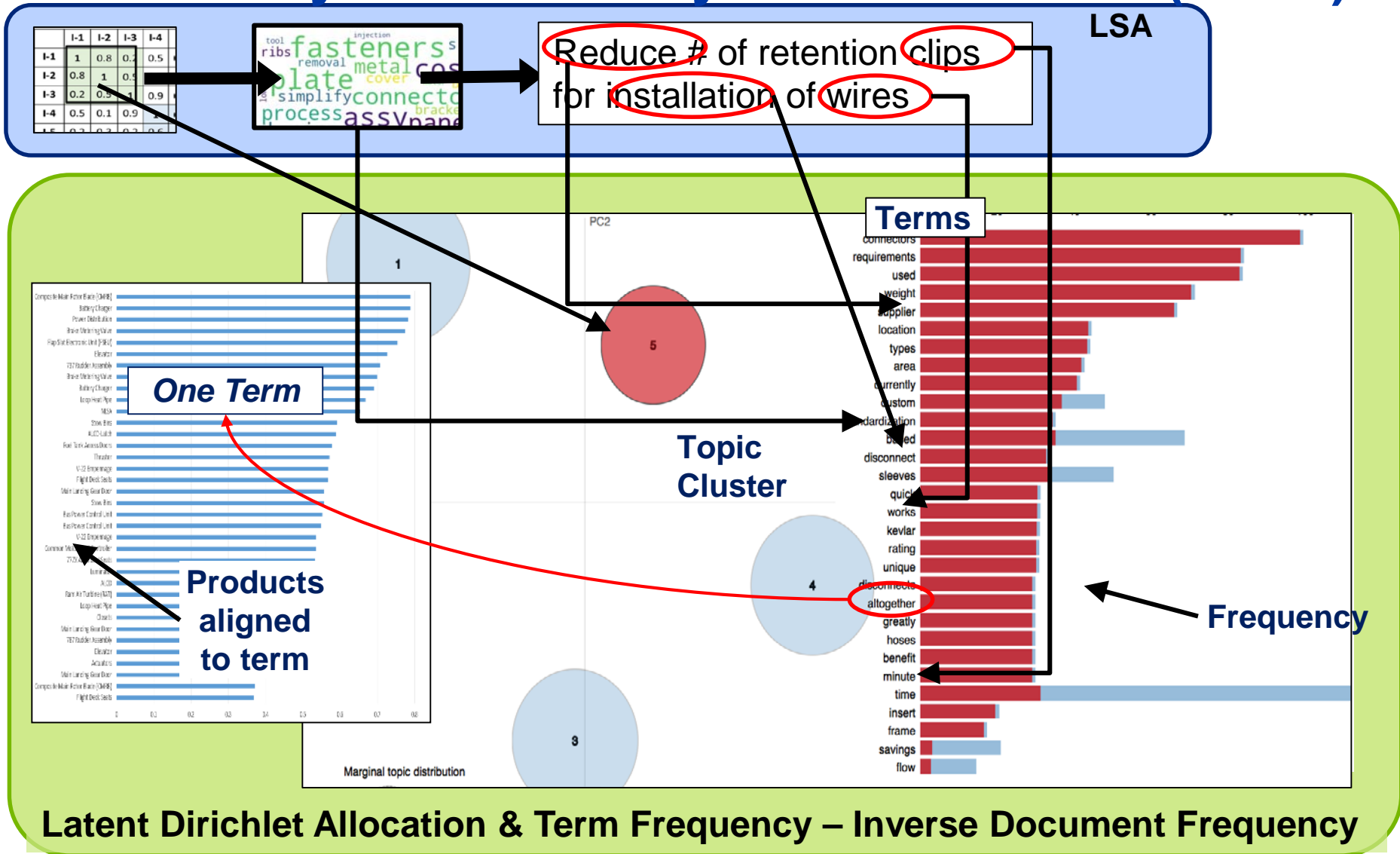
Latent Semantic Analyses

Cluster similar ideas & identify key terms and main concept



Presented at the 2018 ICEAA Professional Development & Training Workshop - www.iceaaonline.com

# Text Analytics to Identify Reusable Ideas (2 of 2)



**Term frequency ~ importance ~ of idea aligned with product**

# Next Steps

## Validate model and verify results

- Modify & Implement existing GUI Framework
- “Evaluate” results – requires thinking!

## Scale to larger population

- Hundreds more workshops & products
- Thousands more ideas

## Capture and incorporate actuals

*Implement cost-saving ideas on other products*

# Machine Learning Approach to Cost Analysis

## Machine Learning in General

## ML Algorithms for Cost Analysis

## Applications related to Cost

- Random Forest Prediction
- Latent Semantic Analysis

## Challenges

\* ML = Machine Learning

# Challenges for Cost Analysis Community

## Machine Learning for cost analysis & estimating

- Different ... from traditional methods
  - Will take time to catch on
- Black box method
  - Not so easy to interpret or follow input-to-output logic
- Regression Algorithms
  - Predict a numeric value (cost) - not a parametric equation (CER)
- ML Algorithms
  - Require pre and post processing for reasonable results

***Do Benefits outweigh Challenges?***

Presented at the 2018 ICEAA Professional Development & Training Workshop - [www.iceaaonline.com](http://www.iceaaonline.com)

## Authors

**Karen Mourikas** is an Associate Technical Fellow at The Boeing Company specializing in Operations Analysis, Affordability, and Systems Optimization. Her current work includes Product Teardown & Should-cost analyses, and Production Systems modeling. Karen has MS degrees in Applied Math and in Operations Research Engineering from the University of Southern California. Karen is a life-time member of ICEAA and has presented at several ICEAA & ISPA/SCEA conferences over the years.

**Nile Hanov** is a Data Scientist at Boeing Research & Technology where he develops novel next gen solutions for commercial and military platforms. In this role, he applies machine learning to event driven data to help organizations better understand and predict failures on board of an aircraft. Nile has four patents under review by the U.S. Patent Office all of which focus on event forecasting and system improvement. He is also currently pursuing a Ph.D. in Computer Science (with a focus on Artificial Intelligence and Machine Learning) at University of California - Irvine.

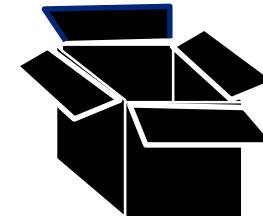
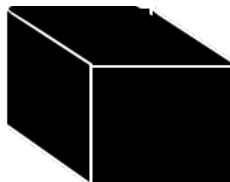
**Joseph King** is a data scientist at The Boeing Company with Boeing Commercial Airplane Analytics, utilizing data to build predictive models and provide analytical solutions. Joseph has contributed to areas such as sensor data analysis, text mining maintenance messages, and customer behavior modeling. Joseph's education background includes a MS in Business Analytics from the University of Tennessee and a background in mathematics and operations research.

**Denise Nelson** is a Systems Analyst at The Boeing Company specializing in software estimating, cost-risk analysis and parametric modeling. Currently, Denise supports Boeing Commercial Airlines Product Development activities. Previous efforts include life-cycle cost analysis; reliability and maintainability analysis; and project management of immersive simulation modeling. Denise graduated from Cal Poly Pomona with an MS in Pure Math and BS in Statistics.

[karen.mourikas@boeing.com](mailto:karen.mourikas@boeing.com) [Nile.Hanov@boeing.com](mailto:Nile.Hanov@boeing.com) [joseph.a.king3@boeing.com](mailto:joseph.a.king3@boeing.com) [Denise.J.Nelson@boeing.com](mailto:Denise.J.Nelson@boeing.com)

# Machine Learning & Non-Parametric Methods for Cost Analysis

The world of big data opens up new opportunities for ICEAA, such as machine learning and non-parametric methods. These methods are more flexible since they do not require explicit assumptions about the structure of the model. However, a large number of observations is needed in order to obtain accurate results. Hence, big data to the rescue! This presentation examines several non-parametric methods, with examples related to our community, and discusses opportunities and limitations going forward.



*Abstract*



**Engineering, Test & Technology**  
Boeing Research & Technology

# Questions?