# LEADS 3-Day Camp
# Day 3 Session 2:
# Automated Data Science Tools

Il-Yeol Song, Ph.D.

Professor

College of Computing & Informatics

Drexel University

Philadelphia, PA 19104

Song@drexel.edu

http://www.cci.drexel.edu/faculty/song/

Il-Yeol Song, Ph.D.

1

---

# Quick Survey

- How many of you are still not comfortable with R coding?

- Do you feel it will be good if there is an automated tool that uses only clicks and drops, without coding?

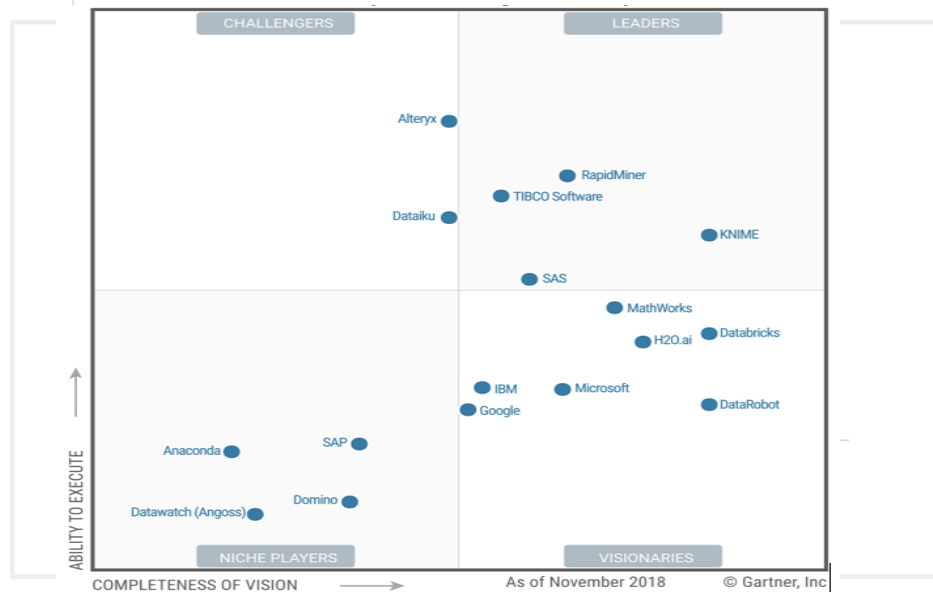- Do you believe automated data science tools are possible?

Il-Yeol Song, Ph.D.

2

# Magic Quadrant for Data Science and Machine-Learning Platforms (2019)



https://www.gartner.com

# Magic Quadrant for Analytics and BI Platforms (2019)



Source: Gartner (February 2019)

https://www.gartner.com

# 19 Data Science and Machine Learning Tools for people who Don't Know Programming

**ML Tools**
- **RapidMiner**
- **DataRobot**
- **BigML**
- **Google Cloud AutoML**
- **MLBase (Berkeley)**
- **WEKA**
- **Driverless AI**
- **Microsoft Azure ML Studio**
- **MLJar**
- **Amazon Lex (Chatbot builder)**
- **IBM Watson Studio**

- **KNIME**
- **Pure Predictive**
- **Logical Glue**
- **FeatureLab**
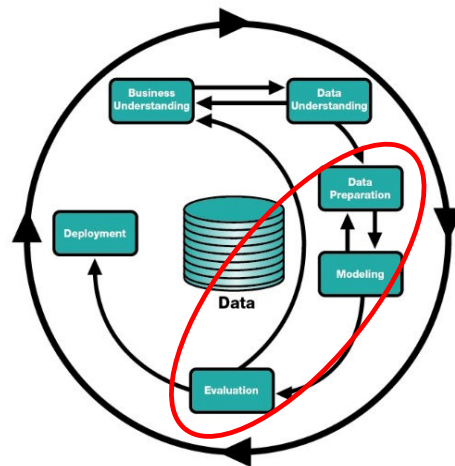- **MarketSwitch**
- **Automatic Statistician**

**Data Preparation Tools**
- **Paxata**
- **Trifacta**

https://www.analyticsvidhya.com/blog/2018/05/19-data-science-tools-for-people-dont-understand-coding/

Il-Yeol Song, Ph.D.

5

---

# Automated DS Tools

- DS Lifecycle

**ML Tasks**

Sub-tasks of Data Science including
- Preparation (feature extraction, etc.)
- Modeling (algorithm selection, hyper-parameters tuning, etc.)
- Evaluation (Validation)



Il-Yeol Song, Ph.D.

6

3

# Future of Automated ML Platform

- Industry pushes automated DS/ML platforms
- *Automated tools support data scientists, but do not replace data scientists*
  - Frees data scientists from the burden of repetitive and time-consuming tasks of
    - developing ideas on where to start,
    - preprocessing,
    - selecting models and parameters,
    - evaluating models and output, and
    - optimizing models.

Il-Yeol Song, Ph.D.

7

---

# Automated DS Tools

- *Automated ML is not the same as Automated DS*

- "*In industry, data scientists will be evaluated on the value added to the business, rather than algorithm accuracy. A project with 99% classification accuracy, but that isn't deployed in production, is bringing no value to the company."*
  **-- Sandro Saitta**

  Source: **Data Science Automation: Debunking Misconceptions**
  https://www.kdnuggets.com/2016/08/data-science-automation-debunking-misconceptions.html

Il-Yeol Song, Ph.D.

8

# Automated DS Tools

- Steps that are difficult to fully automate:
  - Defining problems to solve
  - Getting data and integrating data
  - Exploring data
  - Deploying the project
  - Debugging and monitoring solutions

Il-Yeol Song, Ph.D.

9

---

**Data Science is Changing and Data Scientists will Need to Change Too – Here's Why and How**

Posted by William Vorhies on January 16, 2018 at 8:14am ✉ Send Message 🔭 View Blog

*Summary: Deep changes are underway in how data science is practiced and successfully deployed to solve business problems and create strategic advantage. These same changes point to major changes in how data scientists will do their work. Here's why and how.*

- Gartner says that by 2020 more than 40% of data science tasks will be automated.
- Advanced analytic platforms are becoming one-stop automated system with full support of analytics lifecycle

- Algorithm Selection and Tuning Will No Longer Matter
- Data Prep will be Mostly Automated
- The ability to correctly understand the business problem and translate that into a data science problem will be key.
- Feature Engineering and Model Validation Become a Focus
- Data Science will Increasingly be a Team Sport

Source: https://www.datasciencecentral.com/profiles/blogs/data-science-is-changing-and-data-scientists-will-need-to-change-
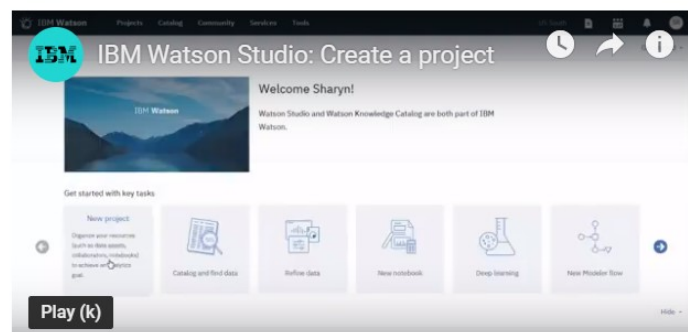
Il-Yeol Song, Ph.D.

10
10

# A Summary of Automated DS Tools

- *Automated ML is not the same as Automated DS*
- *DS tools will become increasingly automated*
- *Automated tools support data scientists, but do not replace data scientists*
  - Frees data scientists from the burden of repetitive and time-consuming tasks of developing ideas on where to start, preprocessing, selecting models and parameters, evaluating models and output, and optimizing models.
- *Citizen Data Scientists will use automated tools*
- *Some tasks difficult to automate include:*
  - Ability to translate business questions into data science problems are critical
  - Feature engineering and model validation are becoming more important

Il-Yeol Song, Ph.D.

11

---

# Watson Analytics



**IBM Watson Studio**



Il-Yeol Song, Ph.D.

12

## Hands-on Experience with Watson Analytics

- Access to Watson Analytics

http://www.ibm.com/analytics/watson-analytics/

Step 1: register a Watson Analytics account

Upload your IRIS data set.



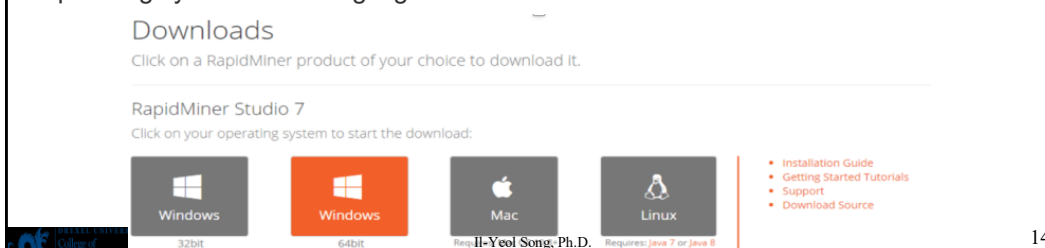Il-Yeol Song, Ph.D.

13

## Hands-on Experience with RapidMiner

Follow these instructions to download RapidMiner Studio:
1. To download the application, go to the RapidMiner website
(https://rapidminer.com/products/studio/)

2. Click the **Download** button in the upper right corner.



3. Log in if you haven't already. See next slide if you need to create an account. If you just want to download RapidMiner Studio without logging in, click on **Downloads**.



4.Click on your preferred operating system to begin the download. Your current operating system will be highlighted.



Il-Yeol Song, Ph.D.

14

7

# RapidMiner Process



Il-Yeol Song, Ph.D.                                                                    15
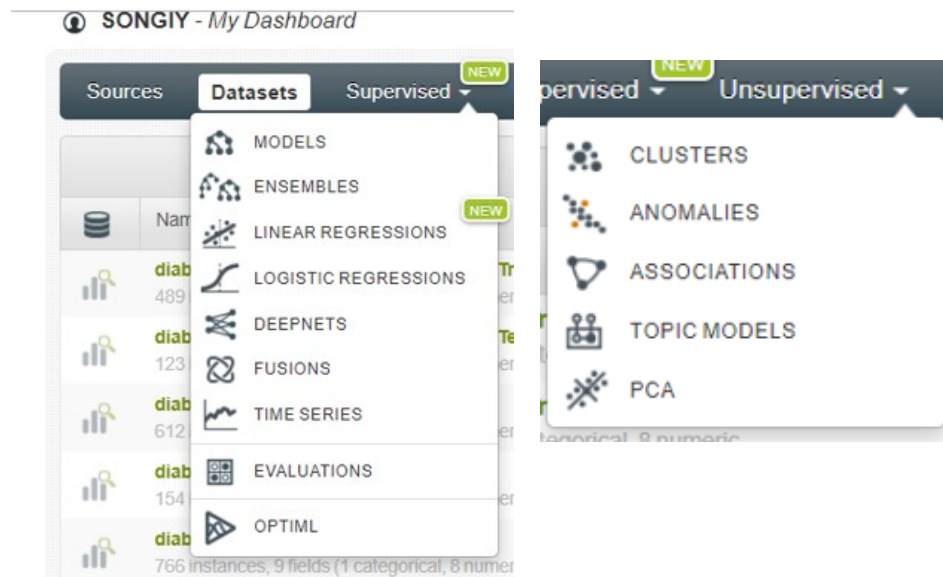
# Hands-on Experience with BigML

• Access to BigML:  https://bigml.com/

Step 1: Create your own ID and PW. In order to get a free ID for one year, you must use an email ID with **.edu** domain, not gmail or other private IDs.
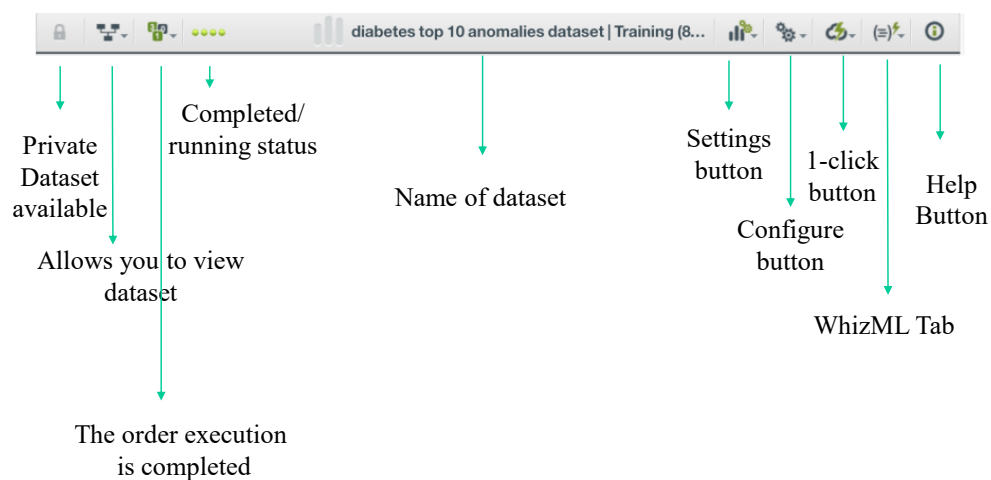


Il-Yeol Song, Ph.D.                                                                    16

# ML Methods supported in BigML

17

# BigML Dashboard



Private Dataset available

Completed/ running status

Allows you to view dataset

Name of dataset

Settings button

Configure button

1-click button

WhizML Tab

Help Button

The order execution is completed

18

# Diabetes Data File Download and Loading

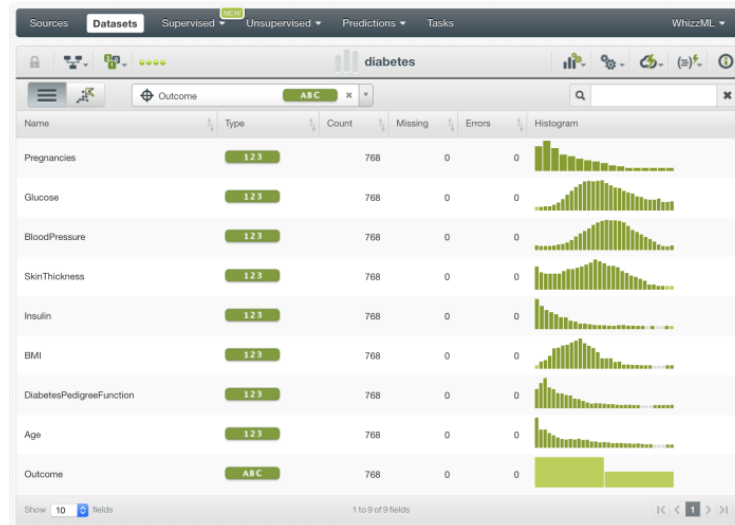- Downloaded a diabetes data from the UCI dataset:

  https://www.kaggle.com/uciml/pima-indians-diabetes-database

- Or Google "pima indian diabetes data set"
- You will find the file description and can download the file
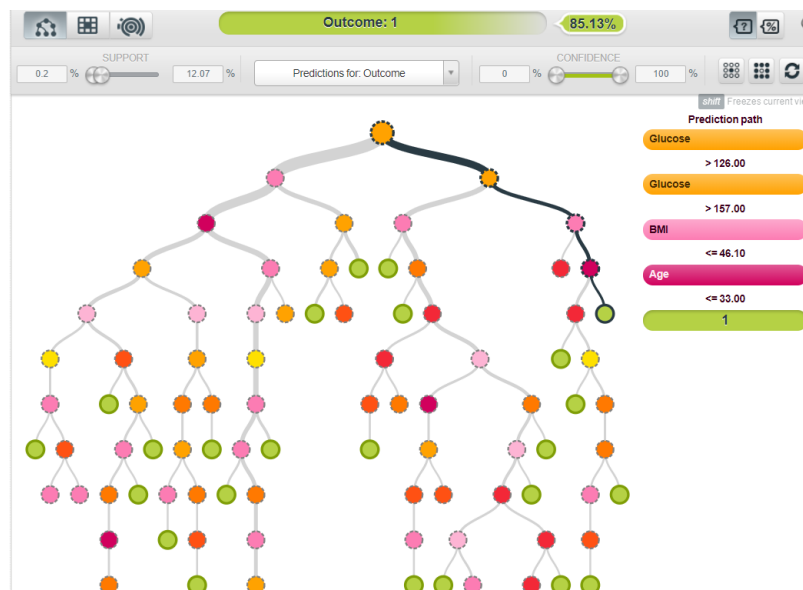
Il-Yeol Song, Ph.D.

19

# Data Description

✓ The Diabetics dataset consists of **9** variables
✓ There are a total of **768 rows**.

- **Pregnancies**: Number of times pregnant
- **Glucose**: Plasma glucose concentration a 2 hours in an oral glucose tolerance test
- **BloodPressure**: Diastolic blood pressure (mm Hg)
- **SkinThickness**: Triceps skin fold thickness (mm)
- **Insulin**: 2-Hour serum insulin (mu U/ml)
- **BMI**: Body mass index (weight in kg/(height in m)^2)
- **DiabetesPedigreeFunction**: Diabetes pedigree function
- **Age**: Age (years)
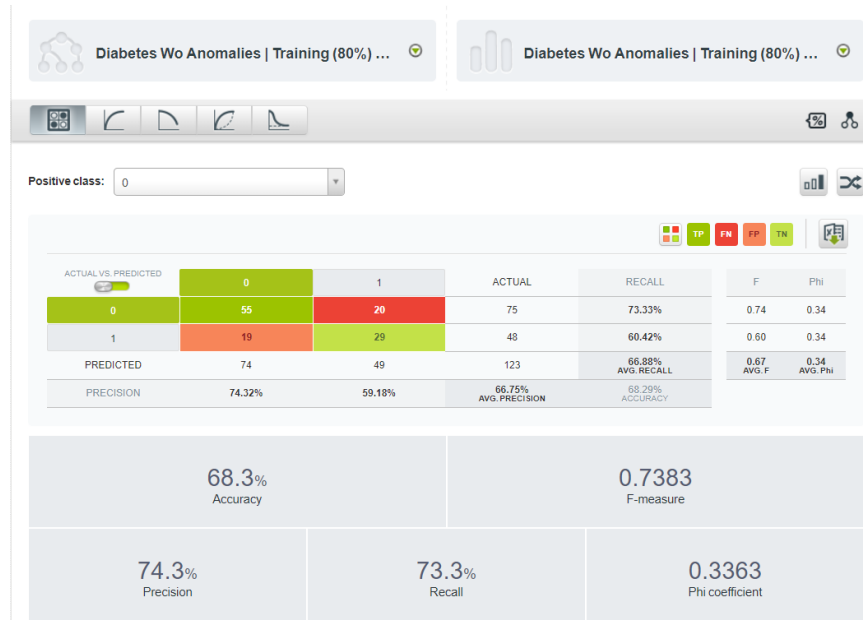- **Outcome**: Class variable (0 or 1) 268 of 768 are 1, the others are 0

Il-Yeol Song, Ph.D.

20

# Diabetes Data Description

✓ The Diabetics dataset consists of **9** variables

✓ There are a total of **768 rows**.



Il-Yeol Song, Ph.D.

21

# Decision Tree Model



Il-Yeol Song, Ph.D.

22

# Confusion Matrix

---

# YouTube Tutorial Links on BigML

- **BigML Interface**
  https://www.youtube.com/watch?v=6xbNpILmQYo
- **Anomaly Detection:**
  https://www.youtube.com/watch?v=a5Q7b4e7lqg&list=PL1bKyu9GtNYHAk0PUojkLYZzASoYVcsTQ&index=13
- **Model Link (Decision Tree):**
  https://www.youtube.com/watch?v=hnt7z24wvxs&list=PL1bKyu9GtNYHAk0PUojkLYZzASoYVcsTQ&index=5
- **Ensemble Link:**
  https://www.youtube.com/watch?v=zqFj6l2WZCU&list=PL1bKyu9GtNYHAk0PUojkLYZzASoYVcsTQ&index=7
- **Evaluation:**
  https://www.youtube.com/watch?v=cPErxYP9CmQ&list=PL1bKyu9GtNYHAk0PUojkLYZzASoYVcsTQ&index=11