# Data Jam
## Workshop 2: Weka

Presenter: Mark Voortman
https://datajam.it.pointpark.edu/

(all materials downloadable)

# Introduction – What is this Data Jam About?

## (Big) Data

## HR Attrition
who leaves voluntarily



## Insights 💡

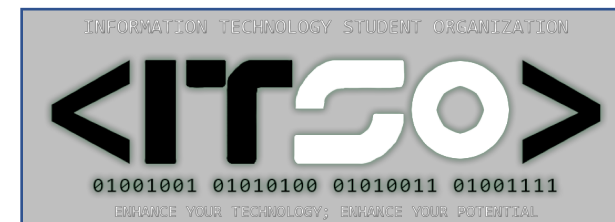# Logistics – Important Dates

- Workshops
    - ~~February 20th (today) – Tableau (data visualization and exploration)~~
    - March 6th – Weka (predictive analytics software)
- Poster Competition
    - April 3rd – Poster Presentations (present your results!)
    - More details ~~next~~ *this* workshop

# Logistics – Random Notes



- Team formation
- You can use any tool you want
  - We teach you Tableau and Weka
  - But feel free to use any other tool (Excel, Python, etc.)
- Judges
  - Industry professionals
  - Very experienced with data and modeling
  - Names, titles, and affiliations to be announced
- The Data Jam is co-organized with ITSO: http://itso.pointpark.edu/
  - Join ITSO if you like this kind of stuff

# Slack – A Tool for Communication

- Slack is a popular **communication** tool used by many **tech** companies
- Go to https://pointparkuniversity.slack.com/ and join
- Use for
  - Reaching out to mentors with questions
  - Team collaboration
- Apps available for iOS, Android, etc.
- See next slides for steps and screen shots

# The Data – How Do I Obtain It?

Download from
[https://datajam.it.pointpark.edu/](https://datajam.it.pointpark.edu/)

*hr-employee-attrition.csv*

**CSV format**
What does that mean?

# The Data – What is the Problem/Goal?

**1. How well can you predict attrition based on other characteristics (e.g., age)?**

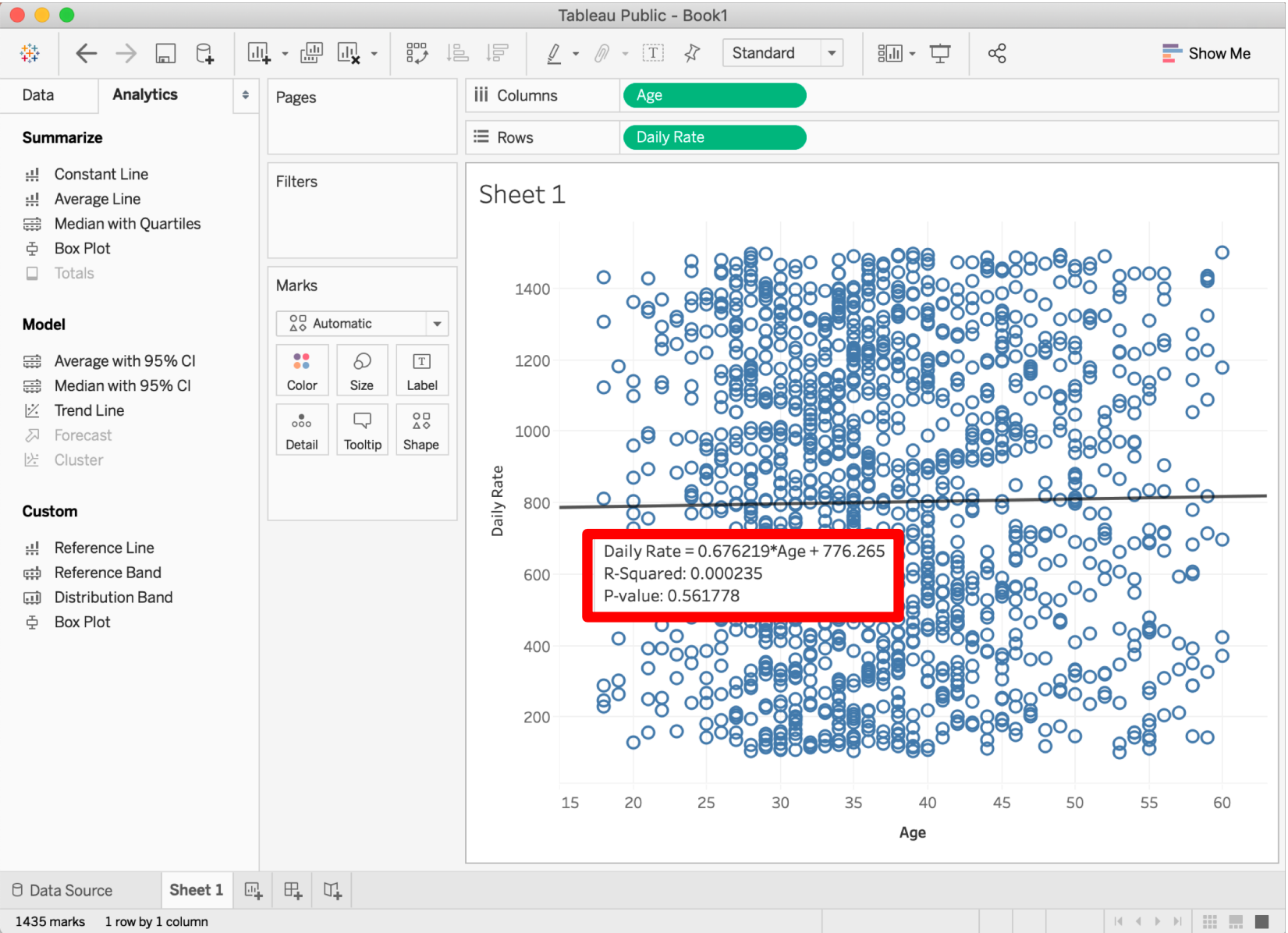Build a model, e.g., if *age >= 65 => attrition=yes*

**2. What drives attrition?**
For example, age

**3. What other general insights can you obtain from the data?**

E.g., what distinguishes high performers?

# Tableau – A Scatter Plot



**Click on Trend Line, drag to the right, and drop on Linear**

**What do you think?**

# Any Questions at This Point?

- Do you need help with anything?

- Were you able to create some visualizations with Tableau?

- Were you able to obtain some insights?

- …

# Posters – Instructions

- Full instructions available at:
https://datajam.it.pointpark.edu/data-jam-poster-guidelines.pdf

- Poster template available at:
https://datajam.it.pointpark.edu/data-jam-poster-template.pptx

- Poster size: 24"x36".  All posters will be printed on foam board and displayed on easels (Data Jam team will coordinate printing).

- Email your PPT poster file to Jaime Ballesteros at
jballesteros@pointpark.edu

- Submit by **Sunday, March 31 at 11:59 p.m**.

# Posters – Guidelines

- **Project title**
- **Full names of all team members**
- **Include information and visuals that address the following:**
  - **Introduction** – State your team's key research question(s)... what are you trying to solve/uncover from the data and why is this relevant?
  - **Method(s) for data analysis** – Tactics to approach... how did you analyze the data?
  - **Results** – Include graphical visualizations of data and key findings. Add legends, captions, or BRIEF explanations if necessary.
  - **Analysis to Insights** – Clearly and concisely explain your findings (what you uncovered through your analysis).
  - **Conclusion** – Link back to your key research question(s) and summarize your impactful findings. Include your team's perspective on the impact of your findings and any recommendations. Also, share problems you encountered.

# Posters – Tips

- **Design your poster as a stand-alone artifact**. Be sure that you "tell the story" of your analysis and findings on the poster. Does it make stand-alone sense without someone there to explain it?

- **Include a brief but descriptive title**. People DO judge a book by its cover… the first thing people will read is your title, so consider your title an invitation to the audience. Your title should let the audience know what your poster is about in a brief sentence or phrase.

- **Emphasize graphics**. Convert information into graphical representations… charts, graphs, and images will capture attention and can effectively communicate data relationships.

- **Keep it clean**. Improve audience engagement and readability… avoid "chart junk" (information not required to understand the graphic), stick to a simple color palette (two to three colors max that don't detract from your content), use dark colors against a light background for better readability when lighting isn't ideal, and leave space between poster elements.

# DATA SCIENCE
## Main Formulas for Machine Learning

**Naïve Bayes**

$$P(a|c) = \frac{P(c|a).P(a)}{P(c)}$$

$$Prob = \Pi P(a|c)$$

**K Nearest Neighbor**

$$D(x_i, x_j) = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

**Support Vector Machines**

$$f(x) = sign[\lambda.y.K(x_i \cdot x_j)]$$

$$K(x_i \cdot x_j) = \sqrt{\frac{\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}}{width}}$$

$$\lambda \rightarrow \nabla L = 0$$

$$y = 1 \wedge y = -1$$

**Perceptron**

$$f(x) = sign\left[\sum_{i=1}^{n} w_j x_{ij}\right]$$

**Neural Networks**

$$f(x) = w_0 + K.\sum_{i=1}^{n} w_i x_i$$

**Backpropagation**

$$\Delta w_{ij}(n) = \eta \delta j x_{ij} + a\Delta w_{ij}(n-1)$$

**Gradient Descent**

$$\theta_{ji} = \theta_j - \alpha \sum_{i=1}^{n}(h(x_i) - y).x_i$$

**Linear Regression**

$$f(x) = \sum_{i=1}^{n} m_i x_i + b$$

**Principal Components Analysis**

$$x_j = x_i - \bar{x}$$

$$Eingenvector = Engeinvalue.[x_i \dots x_n]$$

$$f(x) = Eigenvector^T.[x_{j1} \dots x_{jn}]$$

**Logistic Regression**

$$Odds\ Ratio = log\left(\frac{P(a|c)}{1 - P(a|c)}\right)$$

$$Prob(y = 1) = \frac{1}{1 + e^{-H(\sum_{i=1}^{n} m_i x_i + b)}}$$

**Are you ready for some data mining?**

# Data Mining – Introduction

- Remember this example? *age >= 65 => attrition=yes*

- This is known as a **classification** problem

- Classes: **attrition=yes** and **attrition=no** (binary)
  - Or **yes** and **no** for short

- Prediction is based on **features**
  - Capture the **key** characteristics of the individual
  - In the example above: **age**

**Note: this slide and the next few are all relevant for Weka**

mistake!

separator

attrition=no

attrition=yes

age < 65        65        age > 65

# Data Mining – How Does It Work?



**Which characteristics make someone likely to leave?**

# Data Mining – How Does It Work?



**Rectangular bodies are much more likely to leave than oval bodies!**

# Data Mining – How Does It Work?



**Repeating the process!**

# Data Mining – The Result

This is known as
a **decision tree** or,
more generally,
**(supervised) segmentation**

It can be created **automatically!**

# Data Mining – Overfitting



**Underfitted**                **Sweet spot**                **Overfitted**

**We want a model that is not too general and not too specific!**

# Data Mining – Cross-Validation



Mean and standard deviation of test sample performance

**Cross-validation prevents overfitting by using independent data to assess performance**

**Use 80% to train a model and use 20% to test and repeat 5 times**

**Weka does this for you ☺**

# Data Mining – Assessing Performance

**Accuracy = #CorrectPredictions / #TotalPredictions**

**What should be the baseline to which you compare? 25%? 50%? 75%**
**What would you consider to be good accuracy?**

# Data Mining – Assessing Performance

**Accuracy = #CorrectPredictions / #TotalPredictions**

**What should be the baseline to which you compare? 25%? 50%? 75%**
**What would you consider to be good accuracy?**



**Baseline:**
**1233/(1233+237)**
**= 83.9%**

# Weka – Downloading

**Get the version with the Java VM!**

**Project**    **Software**    **Book**    **Courses**    **Publications**    **People**    **Related**

# Downloading and installing Weka

There are two versions of Weka: Weka 3.8 is the latest stable version and Weka 3.9 is the development version. For the bleeding edge, it is also possible to download nightly snapshots. Stable versions receive only bug fixes, while the development version receives new features.

Weka 3.8 and 3.9 feature a package management system that makes it easy for the Weka community to add new functionality to Weka. The package management system requires an internet connection in order to download and install packages.

# Weka – Install and Run

# Weka – Click Explorer

# Weka – Explorer



**And now we load the data ...**

# Weka – Loading Data



**Make sure to select *.csv**

# Weka – Loading Data



**We now have our data loaded**

# Weka – Change Class to Attrition



**Click on the dropdown and change to Attrition**

# Weka – (Again) Young People Leave!



**Red = stay**

**Blue = leave**

# Weka – Let's Build a Model!



**Click on Classify**

# Weka – Select the Right Target Variable



**Select Attrition as the target variable**

# Weka – Other Things to Note



**ZeroR is the algorithm we run – this happens to be the same as the baseline we discussed**

**Also note that cross-validation is selected by default – no need to worry about it!**

**Now click Start**

# Weka – Baseline Results



**The baseline is 83.9% just like we discussed**

**Many other performance metrics are displayed**

**Try to understand the confusion matrix (yes, it is confusing ☺)**

# Weka – Creating a Decision Tree



**Select J48 under trees**

**Then click Start again**

# Weka – Decision Tree Results



**What do you think?**

# Weka – Visualizing the Tree



**Right click and select Visualize tree**

# Weka – Visualizing the Tree



**A bit messy!**

**But we gain insights: TotalWorkingYears is important!**

# Weka – Tweaking



**Click on the algorithm name**

# Weka – Changing minNumObj



**Change MinNumObj
to 10 and run again**

# Weka – Updated Results



**What do you think?**

**Visualize the tree again: it is less overfitted!**

**(simpler, but better performance)**

# Weka – Logistic Regression (= Classification)



**Click Choose and select Simple Logistic under functions**

**Then run it**

**How about these results?**

# Weka – Logistic Regression (= Classification)



**Scroll up to the first equation**

**The coefficients determine impact (just like with linear regression)**

**For example, traveling frequently is correlated with leaving**

# Next Steps

- Try it out!
- A **lot** of **trial** and **error** (try to be somewhat systematic)
- Use Google
- Try different algorithms and see how well they work
  - Are all of them telling a consistent story?
- Ask questions