# Welcome to dotDATA

Wayne Eternika, President
Steve Mandala, Vice President

Data Science



WHAT IS DATA SCIENCE?

**Mission:
Help You Converge**

**Hacking Skills**

**Math & Statistics Knowledge**

Machine Learning

**Data Science**

Danger Zone!

Traditional Research

**Substantive Expertise**
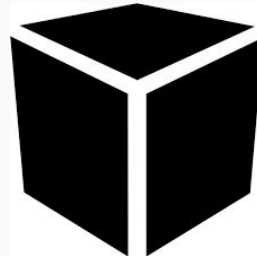
# Motto:  Keep it Parabolic!!!

# Teamwork

# What we are

Connection to resources
Connection to connections
Learning support + framework
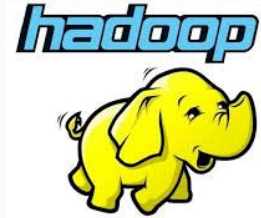Social support

# What we aren't

# Short clips

(First 38 seconds) https://www.youtube.com/watch?v=HPzXlFp4rKE

(First 3 minutes) https://www.youtube.com/watch?v=8pHzROP1D-w

# Technology Behind it:



**Where are you going in the datasphere?**

# Languages in Demand



**Be a polyglot**

# Breakout Groups

Standby for Doodle Polls

Expect exposure to things you haven't thought about

Be a leader, be a member, but just be active

Scheduled times will result from polling as well

Breakout Groups are in addition to general meetings

# Who is he?



A little thing called Machine Learning

# We push MOOCs!



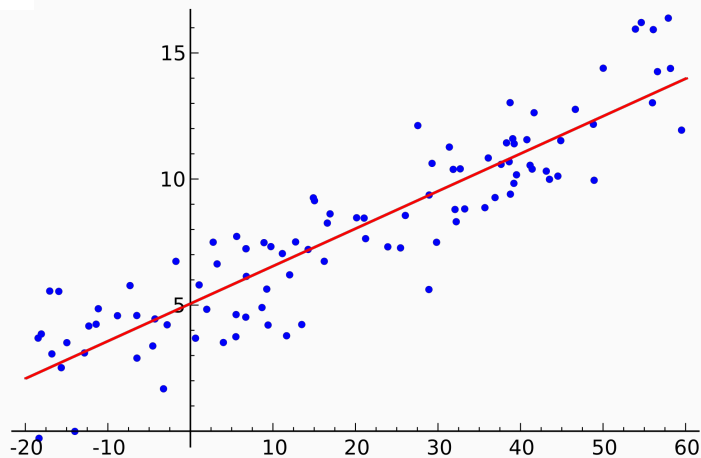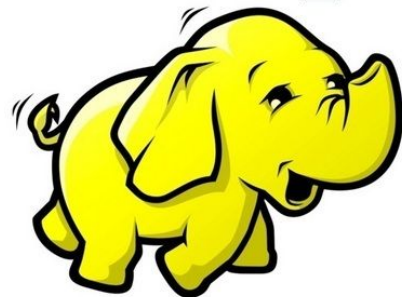Like dealers push drugs.

# Abundance of Information

...just today



Recently ..

# Workshops

# Workshops

# Welcome to Kaggle Competitions

Challenge yourself with real-world machine learning problems

---

## New to Data Science?

Get started with a tutorial on our most popular competition for beginners, Titanic: Machine Learning from Disaster.

## Build a Model

Get the data & use whatever tools or methods you prefer to make predictions.

## Make a Submission

Upload your prediction file for real-time scoring & a spot on the leaderboard.

# Welcome to Kaggle Datasets

The best place to discover and seamlessly analyze publicly-available data

## Dig in

Explore a dataset with our in-browser analytics tool, Kaggle Kernels. You can also download it in an easy-to-read-format.

## Build

Create your data science portfolio. Publish insights and code with Kaggle Kernels and it will be saved to your profile.

## Connect

Engage with other data scientists. Share feedback on other Kagglers' Kernels, or ask a question in a dataset's forum.

# Domain: Stocks

# Example Data Set

Where are our business majors?

# Final Report

https://www.kaggle.com/d/usfundamentals/us-stocks-fundamentals/data-exploring-part-1-indicators/notebook

Apps 　Math 　Actuarial Science 　Statistics 　Discrete Mathematics 　Finance 　General 　AccelaReader

Notebook | Code | Comments (2) | Log | Versions (9) | Forks (1)

**Fork Notebook**

### Data exploring *Part 1*

- How many companies have the same set of indicators?
- How big the set of the common indicators for as many companies as possible?

```
In [1]:  import pandas as pd
         from pandas import Series,DataFrame
         import numpy as np

         # For Visualization
         import matplotlib.pyplot as plt
         import matplotlib
         #%matplotlib inline

         matplotlib.style.use('ggplot')
         df=pd.read_csv('../input/indicators_by_company.csv')
```

```
In [2]:  #number of indicators by company
         df_ind_count = pd.concat([ df[['company_id', 'indicator_id', '2010']].dropna().groupby('com
         pany_id')['indicator_id'].count()
         ,df[['company_id', 'indicator_id', '2011']].dropna().groupby('company_id')
         ['indicator_id'].count()
         ,df[['company_id', 'indicator_id', '2012']].dropna().groupby('company_id')
         ['indicator_id'].count()
         ,df[['company_id', 'indicator_id', '2013']].dropna().groupby('company_id')
         ['indicator_id'].count()
         ,df[['company_id', 'indicator_id', '2014']].dropna().groupby('company_id')
         ['indicator_id'].count()
         ,df[['company_id', 'indicator_id', '2015']].dropna().groupby('company_id')
         ['indicator_id'].count()
         ,df[['company_id', 'indicator_id', '2016']].dropna().groupby('company_id')
         ['indicator_id'].count()
         ], axis=1)
         df_ind_count.columns=['2010','2011','2012','2013','2014','2015','2016']
         df_ind_count.head()
```

Out[2]:  　　2010　2011　2012　2013　2014　2015　2016

---

https://www.kaggle.com/d/usfundamentals/us-stocks-fundamentals/data-exploring-part-1-indicators/notebook

Apps 　Math 　Actuarial Science 　Statistics 　Discrete Mathematics 　Finance 　General 　AccelaReader

Notebook | Code | Comments (2) | Log | Versions (9) | Forks (1)

**Fork Notebook**

Few companies have no more then 1 indicator Some have more then 300 There is a significant number of companies in 180-bin - 250-bin indicators.(except 2010,2011 and 2016) The question is what is the set of these indicators? Are they the same (10-20-30 etc) indocators for the companies or different, not intersectable set?

```
In [12]:  #first 20 indicators which have maximum number of companies
          #each cell contains the num of companies with not empty indicator
          #(one and only one indicator without taking into account any other indicators )
          #in this year
          df_comp_count = pd.concat([
          df[['company_id', 'indicator_id', '2010']].dropna().groupby('indicator_id')['company_id'].c
          ount().sort_values(ascending=False).head(200),
          df[['company_id', 'indicator_id', '2011']].dropna().groupby('indicator_id')['company_id'].c
          ount().sort_values(ascending=False).head(200),
          df[['company_id', 'indicator_id', '2012']].dropna().groupby('indicator_id')['company_id'].c
          ount().sort_values(ascending=False).head(200),
          df[['company_id', 'indicator_id', '2013']].dropna().groupby('indicator_id')['company_id'].c
          ount().sort_values(ascending=False).head(200),
          df[['company_id', 'indicator_id', '2014']].dropna().groupby('indicator_id')['company_id'].c
          ount().sort_values(ascending=False).head(200),
          df[['company_id', 'indicator_id', '2015']].dropna().groupby('indicator_id')['company_id'].c
          ount().sort_values(ascending=False).head(200),
          df[['company_id', 'indicator_id', '2016']].dropna().groupby('indicator_id')['company_id'].c
          ount().sort_values(ascending=False).head(200)
          ], axis=1)

          df_comp_count.columns=['2010','2011','2012','2013','2014','2015','2016']

          df_comp_count.head()
```

Out[12]:

| | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 |
|---|---|---|---|---|---|---|---|
| AccountsPayableAndAccruedLiabilitiesCurrent | NaN | 1029.0 | NaN | NaN | NaN | NaN | NaN |
| AccountsPayableCurrent | 194.0 | 3882.0 | 3959.0 | 3944.0 | 3907.0 | 3541.0 | 7.0 |
| AccountsReceivableNetCurrent | 138.0 | 3351.0 | 3394.0 | 3369.0 | 3368.0 | 3046.0 | 7.0 |
| AccruedIncomeTaxesCurrent | 82.0 | 1055.0 | NaN | NaN | NaN | NaN | NaN |

# Got Your Own Project?

Help get in touch with relevant faculty & staff

2 Project-based courses. CS638 and Stat 679, on data science

Working with professors to facilitate independent projects/directed studies. Similar to Directed Reading Program in Mathematics dept

Stay Tuned!

Local Meetups & Groups

# Reminder



Doodle Polls

Join the org on the WIN page:
https://win.wisc.edu/organization/dotDATA

Like our FB page & join our FB group:
www.facebook.com/dotdatascience/
www.dotdatascience.org

# A few more minutes ..

- Graduate Students
- Undergrads with interest in helping drive cohesive administration
- Questions? Queries? Conundrums? Ideas? Share them with us!