# Startup Database and Recommendation Engine

K Iwasaki
kaiwasaki@berkeley.edu

# Table of Contents

**Introduction**

- Problem

- Project Overview

**Data collection and preprocessing**

- Data Collection and Preprocessing Scheme

- Data Collection and Preprocessing Approach

- Data Scraping

- Data Extraction: Company Name from Article

- Data Extraction: Company Address

- Data Extraction: Industry Attributes

**Data Exploration and Visualization**

- Data Exploration to Recommendation Generation

- Summary Statistics, Variable Category, and NAs

- Univariate analysis

- Bivariate analysis

**Recommendation Engine**

- Recommendation Engine - KNN

- Feature Transformation and Engineering

- Recommendation Output

**Closing**

- Future Development

# Problem: Finding Good Startup is Hard

For job seekers, finding a startup that matches their interests is hard because:

Overwhelming information available online

So many information sources to check

Need to synthesize information

Need to check information frequently

So many texts to read

Job seekers are highly biased

Frustration, Time waste, Not finding company that matches your interests

# Project Overview: Finding the Best Startup For You

Create end-to-end solution from data collection, to database generation, to generation of recommendation for startups that matches your interests.

**Data Collection/Preprocessing**

10 articles x 100 pages

**TechCrunch**

~300 searches

**Google**

~ 300 company profiles

**Linked in**

**Bloomberg**

**Startup Database**

Company name
Company size
Money raised
Industry
Description
HQ Location

...

**Recommendation Engine**

**Data Exploration and Visualization**

# Data Collection and Preprocessing Scheme

| | | |
|---|---|---|
| Article link | Article title | Article excerpt |

| | | | |
|---|---|---|---|
| Published at | Company Name | Funding Round | Money Raised |

Link to LinkedIn Company Profile

| | | | |
|---|---|---|---|
| Company Name | Founded | Industry | Specialties |
| Website Link | Description | Also-viewed | Location |

| | |
|---|---|
| Industry attributes | Location |

**Scrape the articles about Series C fundraising from TechCrunch (article.csv)**

Preprocessing1. Extract
Preprocesisng2. Extract company names
Preprocessing3. Extract funding_round and money_raised

**Scrape the website links to LinkedIn Company Profiles from Google Search (linkedin_link_list.csv)**

Preprocessing4. Merge the two CSV files
Preprocessing5. Validate company names

**Scrape the company profile for each company from LinkedIn (linkedin_profiles.csv)**

Preprocessing6. Merge the two CSV files
Preprocessing7. Extract locations
Preprocessing8. Assign industry attributes

5

# Data Collection and Preprocessing Approach

*Think what information we want for the database and for the recommendation engine*

*Write codes and extract information from the target source*

*Error analysis: Confirm if we get what we want and verify missing data and why*

Company name
Company size
Money raised
Industry
Description
HQ Location
…

# Data Scraping

**Tasks:** Use Selenium and BeautifulSoup to scrape information the target websites.

*get_seriesC_news_from_techcrunch.py*
Input: key words "raises Series C"
Output: articles in csv file

*get_link_to_linkedin_from_google.py*
Input: company name
Output: a link to company profile at LinkedIn

*get_profile_from_linkedin.py*
Input: a link to company key profile
Output: company profiles in csv file

# Data Extraction: Company Name from Article Title

Once we collect the articles, next step is extract company names from the article titles.
**Challenge:** a company name is irregular: it can be one word, two words, or more. It often is a mix of verb, noun, or others. Below are typical patterns that a company name shows up in an article title.

"*Stash raises $40 million Series C to make investing more approachable*"

"*Data Storage Company Scale Computing Raises $17 Million Series C*"

"*Pivotal confirms Series C round is actually over $650 million*"

"*After bump in the road, Movinga raises $17M Series C*"

"*Carwow, a UK startup that helps you buy a new car, raises $39M Series* "

"*Confirmed: London fintech Curve raises $10M Series A*"

*Company names, Key verbs, decorative words*

# Algorithm for Company Name Extraction

**Solution:** algorithm to extract a company name, leveraging sentence structures of the articles that are scraped from TechCrunch. Also double-check the company name when googling it later to look for a link for a company profile page at LinkedIn. Check ***company_from_title.py for the codes***

Step1:

- Split the sentence by a key verb and keep the head

- Remove ", word word ... ,"
  - If one or two words remained=> done    Else: => Step2

Step2:

- Split the sentence by a key noun and keep the tail
  - If one or two words remained=> done    Else: => Step2

Step3:

- Split the sentence by "$" and keep the head

- Split the sentence by "Series" and keep the head

# Data Extraction: Company Address

Company address is import input for the recommendation engine because many of us care where we work at.

**Challenge:** Some companies don't input their company address at LinkedIn. Some companies are based outside of US and thus their addresses have different formatting.



It only says United States for Headquarters.



It doesn't have address at all.

10

# Solution to extract/revise Company Address

**Solution:** Multi-step approach: first focus on label countries and then focus on US companies to extract zip code.

- Step1: Complete labeling by countries
  - Country list scraped from Wikipedia --- *get_country_from_wiki.py*
  - Extract country information from features collected so far

- Step2: Focus on the US companies and get zip code for them
  For missing or insufficient information
  - Google Search ---  *get_comnay_address.py*
  - Company Website --- *get_location_from_company_website.py*
  - Bloomberg  --- *get_company_address_from_Bloomberg.py*

- Step3: Gain state, city, geo location from the zip code for US companies
  Use two python modules to capture city and state because both of them have some missing data

# Data Extraction: Industry Attributes from Description

**Challenges:** Industries have been arbitrarily assigned to companies. As a result, there are 49 unique industries for about 300 companies. There are three problems in order for recommendation engine to work:

1) **Some industries are quite similar** thus should be merged.
2) **Some industries have lots of companies** such as Computer Software. They should be split into more smaller segment.
3) **one industry is not sufficient to describe a nature of a company** because its business is often a combination of different elements. For example, the company below is internet x financial service, instead of internet alone

```
data.Industry.unique()
```

```
array(['Financial Services', 'Information Technology and Services',
       'Human Resources', 'Computer Software',
       'Logistics and Supply Chain', 'Internet',
       'Computer & Network Security', 'Food & Beverages',
       'Marketing and Advertising', 'Medical Devices', 'E-Learning',
       'Consumer Services', 'Sports', 'Consumer Electronics',
       'Computer Hardware', 'Education Management', 'Apparel & Fashion',
       'Entertainment', 'Consumer Goods', 'Biotechnology',
       'Management Consulting', 'Real Estate', 'Fund-Raising',
       'Commercial Real Estate', 'Food Production', 'Online Media',
       'Mechanical or Industrial Engineering', 'Renewables & Environment',
       'Farming', 'Electrical/Electronic Manufacturing',
       'Leisure, Travel & Tourism', 'Sporting Goods', 'Retail',
       'Semiconductors', 'Cosmetics', 'Insurance', 'Telecommunications',
       'Health, Wellness and Fitness', 'Textiles',
       'Staffing and Recruiting', 'Nanotechnology',
       'Luxury Goods & Jewelry'], dtype=object)
```

**LendingClub**
Internet
1001-5000 employees

Home    Careers

**LendingClub**

LendingClub is America's largest online marketplace connecting borrowers and investors, facilitating personal loans, business loans, and financing for elective medical procedures and K-12 education and tutoring. Borrowers access lower interest rate loans through a fast and easy online or mobile interface. Investors provide the capital to enable many of the loans in exchange for earning interest. We operate fully online with no branch infrastructure, and use technology to lower cost and deliver an amazing experience. We pass the cost savings to borrowers in the form of lower rates and investors in the form of attractive returns. We're transforming the banking system into a frictionless, transparent and highly efficient online marketplace, helping people achieve their financial goals every day.

Since launching in 2007 we've built a trusted brand with a track record of delivering exceptional value and satisfaction to both borrowers and investors. LendingClub's awards include being named to the Inc. 500 in 2014 and a CNBC Disruptor 50 for the second year in a row, one of Forbes' America's Most Promising Companies three years in a row, one of The World's 10 Most Innovative Companies in Finance by Fast Company in 2013 and a 2012 World Economic Forum Technology Pioneer.

(Notes by Prospectus - https://www.lendingclub.com/info/prospectus.action)

**Specialties**
Personal Loans, Investing, Peer-to-Peer Lending, Patient Financing, Marketplace Lending, Business Loans, Education Financing

| Website | Industry | Type |
|---|---|---|
| http://www.lendingclub.com | Internet | Public Company |

| Headquarters | Company Size | Founded |
|---|---|---|
| 71 Stevenson Street Suite 300 San Francisco, CA 94105 United States | 1001-5000 employees | 2006 |

# Algorithm to Assign Industry Attributes to Each Company

**Solution Part1:** algorithm to simplify the industry classification by merging some industries so that minor industry labels are eliminated

| Industry (Original) | Industry_consolidated (New) |
| --- | --- |
| ["Apparel & Fashion", "Consumer Goods", "Consumer Services", "Cosmetics", "Luxury Goods & Jewelry", "Retail", "Leisure, Travel & Tourism", "Sporting Goods", "Textiles"] | Consumers Goods  & Services |
| ["Computer Software"] | Computer Software |
| ["Computer & Network Security", "Computer Hardware"] | Computer & Network Security & Hardware |
| ['E-Learning', 'Education Management'] | Education |
| ["Entertainment"] | Entertainment |
| ["Marketing and Advertising"] | Marketing and Advertising |
| ["Farming", "Food & Beverages", "Food Production", "Restaurants"] | Food Business |
| ["Insurance", "Fund-Raising", "Financial Services"] | Financial Services |
| ["Information Technology and Services"] | Information Technology and Services |
| ["Internet", "Online Media"] | Internet |
| ["Commercial Real Estate", "Real Estate"] | Real Estate |
| ['Health, Wellness and Fitness', 'Medical Devices', "Sports"] | Healthcare_health |
| ["Human Resources", "Staffing and Recruiting"] | Human Resources |
| ["Telecommunications", "Renewables & Environment", "Logistics and Supply Chain"] | Infrastructure |
| ["Semiconductors", "Nanotechnology", "Biotechnology", "Management Consulting", "Electrical/Electronic Manufacturing" "Mechanical or Industrial Engineering" ] | Niche |

13

# Algorithm to Assign Industry Attributes to Each Company

**Solution Part2:** Algorithm to add new features to represent company businesses better based on the key words in appeared in company profiles

```python
key_words_dict = {
    "Food Business": ["restaurant", "farm", "greenhouse", "Gastronomie"],
    "Education": ["Online Learning", "Education", "Tutor"],
    "Financial Services": ["payment", "loan", "financ", "fundraising",
                "investing", "lending"],
    "Healthcare_health": ["healthcare", "medical", "genetic", "therapy", "disease",
                "fitness", "wellness", "welfare","wearable", "gym"],
    "Human Resources": ["recruit", "workforce", "Human Resource"],
    "Logistics and Supply Chain": ["delivery", "drone",
                                    "transportation", "supply chain"],
    "Entertainment": ["entertainment", "game"],
    "Computer & Network Security & Hardware": ["storage","backup", "recovery",
                                    "privacy"],
    "Real Estate": ["Real Estate"],
    "Marketing and Advertising": ["marketing", "advertising", "advertisement"],

    "commerce": ["eCommerce", "Commerce", "Retail"],
    "mobile" : ["mobile"],
    "app": ["mobile app", "app\s"],
    "analysis": ["analytics", "analysis"],
    "developer": ["developer"],
    "security" : ["fraud", "detection", "protection"],
    "social": ["Social Media"],
    "ds": ["artificial intelligence", "machine learning",
        "deep learning", "big data"],
    "travel": ["Travel"],
    "booking_ticketing": ["booking", "ticket"],
    "Apparel": ["fashion", "clothing", "shoes", "Sporting Goods"],
    "cloud": ["cloud"],
    "API": ["API"],
    "device": ["device"],
    "design": ["design"],
    "enterprise": ["enterprise", "productivity", "collaboration"],
    "robotics_manufacturing": ["Manufact", "robotics", "3d"]
}
```

14

# Now Database is Set! --- 232 rows by 47 columns

# Data Exploration to Recommendation Generation

**Data Exploration/ Visualization**

Data exploration is to gain insights for algorithm selection, feature selection, feature transformation through following steps:

- Summary statistics, variable category, NA value detection

- Univariate analysis

- Bivariate analysis

*Documentation and codes:*

**Recommendation Engine**

Based on the inputs from the data exploration, we create the recommendation and generate recommendations in the following steps:

- Algorithm Selection

- Feature Transformation and Engineering

- Recommendation Output

*Documentation and codes:*

# Summary Statistics, Variable Category, and NA values

Refer the document for the details: explatory_data_analysis/explatory_data_analysis

```
summary(data[, 1:13 ])
```

```
##      published_at   funding_round money_raised_float
## 1/14/2016 :  3              : 12   Min.   : 10.00
## 7/29/2015 :  3    Series B: 48     1st Qu.: 15.00
## 11/15/2016:  2    series C:  1     Median : 25.00
## 11/3/2015 :  2    Series C:152     Mean   : 41.17
## 12/11/2013:  2    Series D: 14     3rd Qu.: 45.00
## 2/11/2008 :  2    Series E:  5     Max.   :793.50
## (Other)   :218
##                                        CompanyName      CompanySize
## 2U                                           :  1   51-200   :101
## 3D Robotics                                  :  1   201-500  : 64
## aCommerce - Ecommerce Solutions for Southeast Asia:  1   Nov-50   : 27
## Affle                                        :  1   1001-5000: 21
## App Annie                                    :  1   501-1000 : 14
## Appear Here                                  :  1   10,001+  :  2
## (Other)                                      :226   (Other)  :  3
##     Founded               City    address_check         Country
## Min.   :1939   San Francisco:49   False: 59     United States :178
## 1st Qu.:2007                :46   True :173     United Kingdom: 17
## Median :2010   New York     :27                 Germany       :  8
## Mean   :2009   Mountain View:11                 Canada        :  4
## 3rd Qu.:2012   San Mateo    : 8                 India         :  4
## Max.   :2017   Boston       : 6                 Singapore     :  4
##                (Other)      :85                 (Other)       : 17
```

```
# check columns 1:13. Columns 13: have same format.
str(data[, 1:13 ])
```

```
## 'data.frame':    232 obs. of  13 variables:
##  $ published_at              : Factor w/ 209 levels "1/11/2010","1/14/2016",..: 189 176 160 15
##  $ funding_round             : Factor w/ 6 levels "","Series B",..: 4 4 4 4 4 4 4 4 4 4 ...
##  $ money_raised_float        : num  45 39 40 48 90 20.2 29 32 36 20 ...
##  $ CompanyName               : Factor w/ 232 levels "2U","3D Robotics",..: 25 29 185 126 39 12
##  $ CompanySize               : Factor w/ 8 levels "10-Jan","10,001+",..: 4 7 7 7 4 4 7 7 7 .
##  $ Founded                   : num  2013 2013 2015 2011 2013 ...
##  $ City                      : Factor w/ 68 levels "","Arlington",..: 1 1 36 1 55 1 52 36 52 2
##  $ address_check             : Factor w/ 2 levels "False","True": 1 1 2 1 2 1 2 2 2 2 ...
##  $ Country                   : Factor w/ 23 levels "Belgium","Brazil",..: 8 22 23 22 23 7 23 2
##  $ latitude                  : num  NA NA 40.7 NA 37.4 ...
##  $ longitude                 : num  NA NA -74 NA -122 ...
##  $ Industry_consolidated     : Factor w/ 16 levels "Computer & Network Security & Hardware",..
##  $ spc_Logistics.and.Supply.Chain: int  0 0 0 1 0 0 0 0 0 0 ...
```
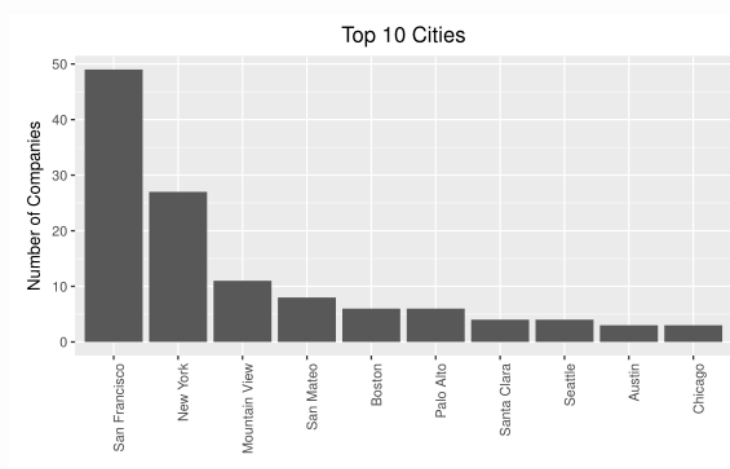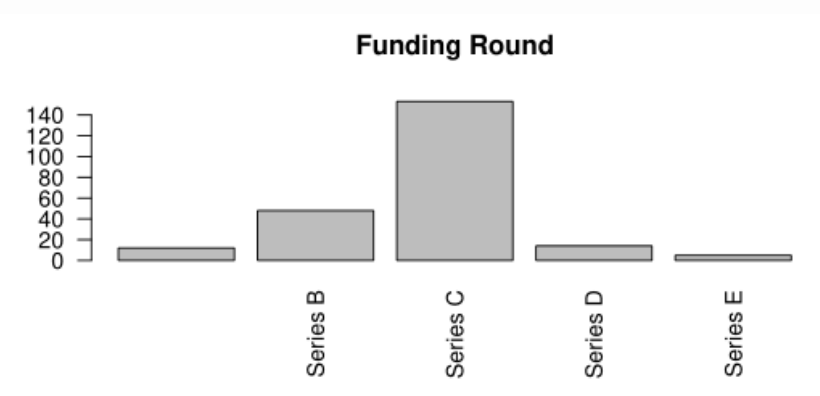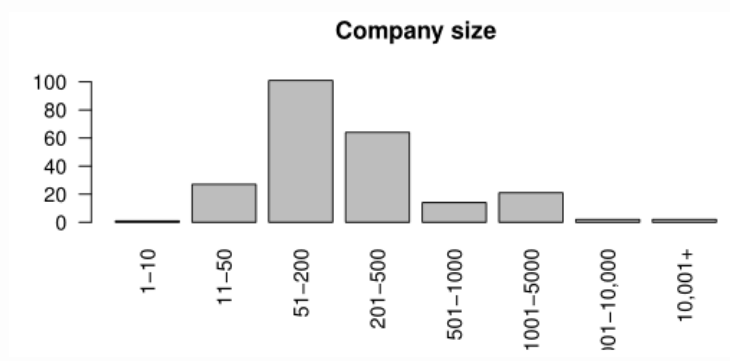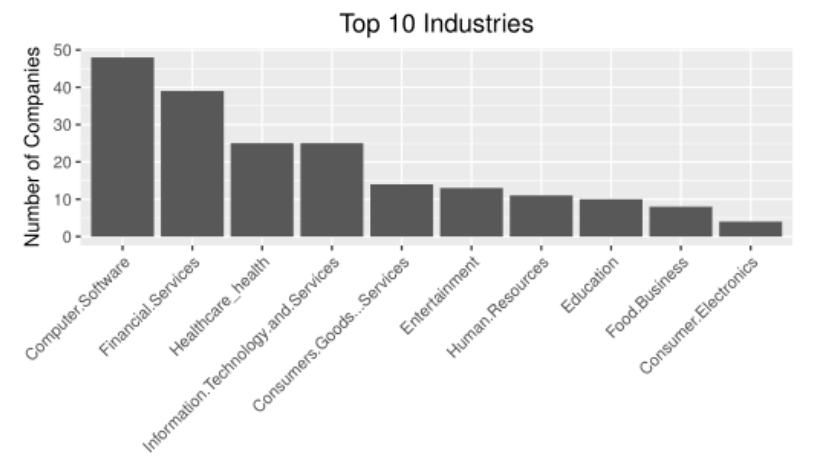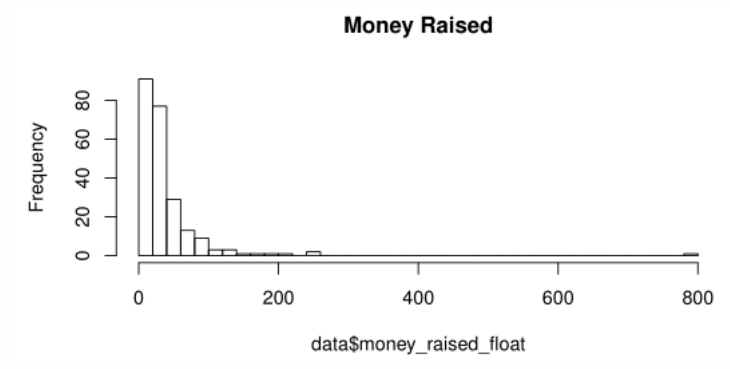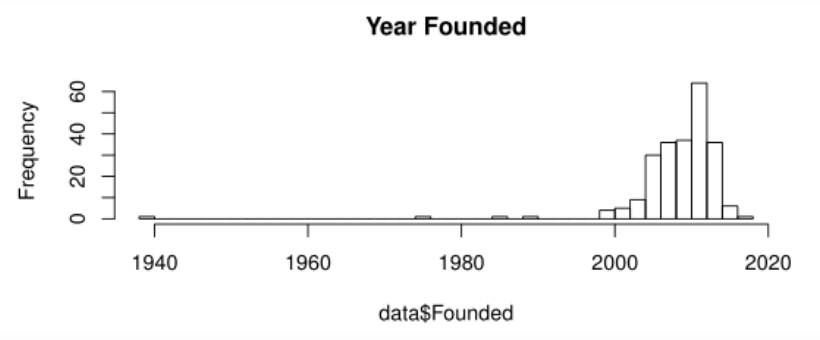
```
# show columns with na
na = lapply(data, function(x) sum(ifelse(is.na(x) | x == "" | x == "not found", TRUE, FALSE)))
na[na > 0]
```

```
## $funding_round
## [1] 12
##
## $City
## [1] 46
##
## $latitude
## [1] 46
##
```
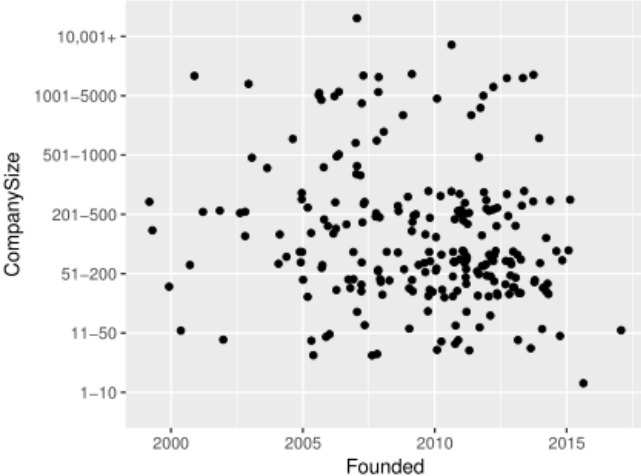
# Univariate analysis
Refer the document for the details: explatory_data_analysis/explatory_data_analysis
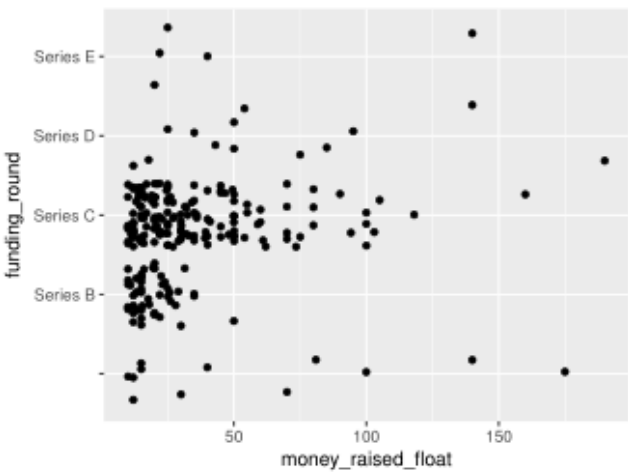
# Bivariate analysis

Refer the document for the details: explatory_data_analysis/explatory_data_analysis

# Bivariate analysis
Refer the document for the details: explatory_data_analysis/explatory_data_analysis

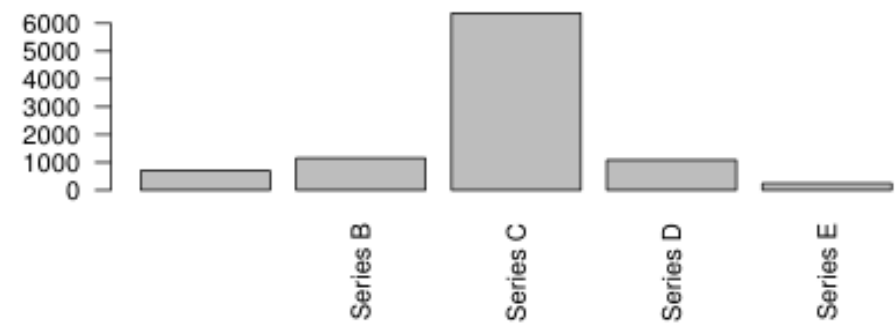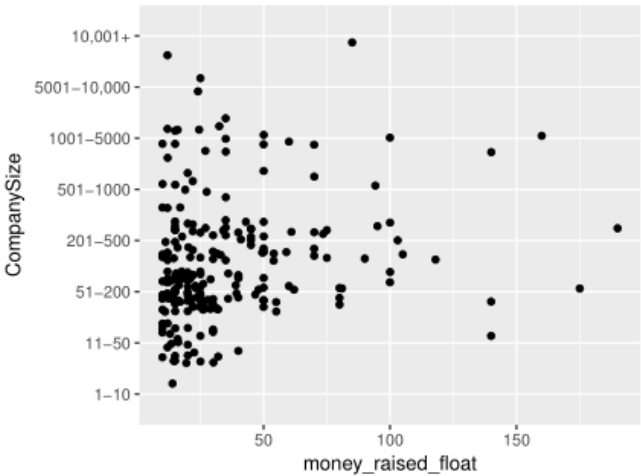# Recommendation Engine – K-Nearest Neighbor



**How the recommendation engine should work:**

Given user inputs such as industry, company size, and year founded, it provides a few companies that matches the inputs.

**Algorithm choice:** K-nearest neighbors (KNN)

**Justification for the choice:** KNN works well for multi-class problems like this problem where we want to assign the user input to a label (company) as outputs out of all the different labels.  It also produces several neighbors which we can use as a secondary recommendations for the user.

**How KNN works:** A green dot as the user input and other dots are startups in the database. KNN calculates the distance between the green dot and other dots and come up with K dots that are closest to the green dot. The shorter the distance is is, the better matches between the input and the neighbors are. These neighbors become the recommendation.

# Feature Transformation and Engineering

KNN requires features to be scaled properly because KNN is distance-based algorithm and calculates a selected distance metric between the user inputs and each example of the training data.
This implies that KNN only takes a numeric variable and a dummy variable. Thus, I made transformations as followings for the features.

| Features | Money Raised | Company Size | Location | Year Founded | Industry |
|---|---|---|---|---|---|
| Transformations | Log transformation to handle outliers | Binning | Convert city to geo coordinate | Removal of outliers | Key word matching |
| | Min max transformation | Min max transformation | Min max transformation | Min max transformation | Dummy variables |

Refer the document for the details:  recommendation/recommendation

22

# Recommendation Output

```
In [51]:  generate_recommendation(train, test.ix[0:0,], x_test.ix[0:0,])

          ----------------------------------------------------------------------
          Thank you for providing your interests! Below are the summary of your interests

          Headquarters:        San Francisco
          Year founded:        2015
          Company size:        11-50
          Industry:            Education & Internet
          ----------------------------------------------------------------------
          We recommend to check 'Edmodo' that matches your interests!

          About the start up

          Our mission is to connect all learners to the people and resources needed to achieve their full potential. We are the world's l
          eading global education network that provides communication, collaboration, and coaching tools for all members of the school co
          mmunity. We were founded in 2008 and currently have over 70 million members across 350,000+ schools in 150 countries.

          The investors backing Edmodo are some of the best-recognized firms in the world, including Benchmark Capital, Greylock Venture
          s, Index Ventures, Union Square Ventures, Learn Capital and our Chairman is Reid Hoffman, founder of LinkedIn.

          So join the team that is changing how teachers and students learn - change lives, build your career and rack up the karma.

          Company details

          Website:             http://www.edmodo.com
          Headquarters:        San Mateo, CA
          Year founded:        2008
          Company size:        51-200
          Techcrunch article:  https://techcrunch.com/2012/07/19/nea-leads-educational-network-edmodos-25-million-series-c/


          ----------------------------------------------------------------------
          We also suggest checking following startups

                       Company  Money_raised  Founded Company Size         City
          99       Varsity Tutors            50    2,007    201-500    Saint Louis
          20             Coursera            50    2,012    201-500  Mountain View
          158  Boomerang Commerce            12    2,012     51-200  Mountain View
          42          Engine Yard            19    2,006     51-200  San Francisco
          135             Base CRM            15    2,009    201-500  Mountain View
          ----------------------------------------------------------------------
```

Summery of user inputs

Summary of the top recommendation

Secondary recommendations

Refer the document for the details:  recommendation/recommendation

# Future Development

| | |
|---|---|
| Data collection | Incorporate more data sources such as Glassdoor.<br>Create data pipeline that is based on once a day batch processing from multiple data sources.<br>Improve algorithms for various data extraction works by utilizing existing NLP packages. |
| Data storage | Store the data in database such as PostgreSQL for better data management and data retrieving capability. |
| Data preprocessing | Clean up codes and streamline the process.<br>Incorporate better handlings. |
| Recommendation Engine | Store companies also-viewed for each company profile at linked in Graph DB such as Neo4j and generate startup recommendations based on the DB. |
| Interface | Create a Web application using Flask and develop GUI to enable users to input their preferences and to view recommendation outputs. |