

Data Science in Action

Ji Li

Data Scientist

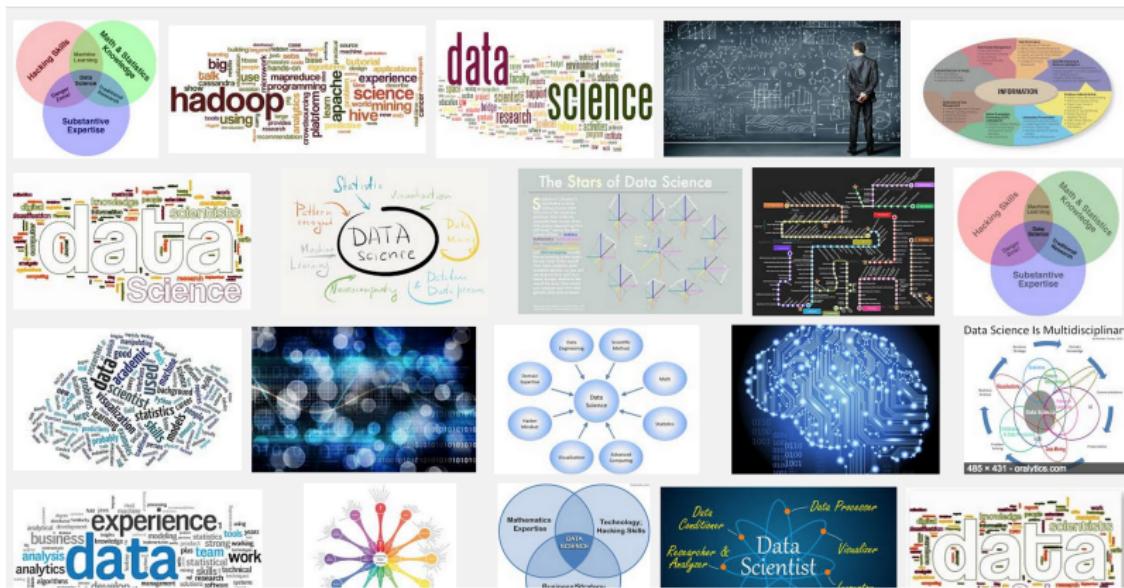
March 25, 2015

Table of Contents

- 1 What is Data Science
 - Overview
 - Data Science Workflow
 - 2 Churn Model
 - Business Understanding
 - Demo in R
 - Top Features
 - Prescriptive Analysis
 - 3 Yesware Email Analysis
 - Email Analysis
 - 4 Data Science in the Industry
 - Data Science Toolbox
 - Big Data

Overview

Data science on Google search



Data science from my point of view

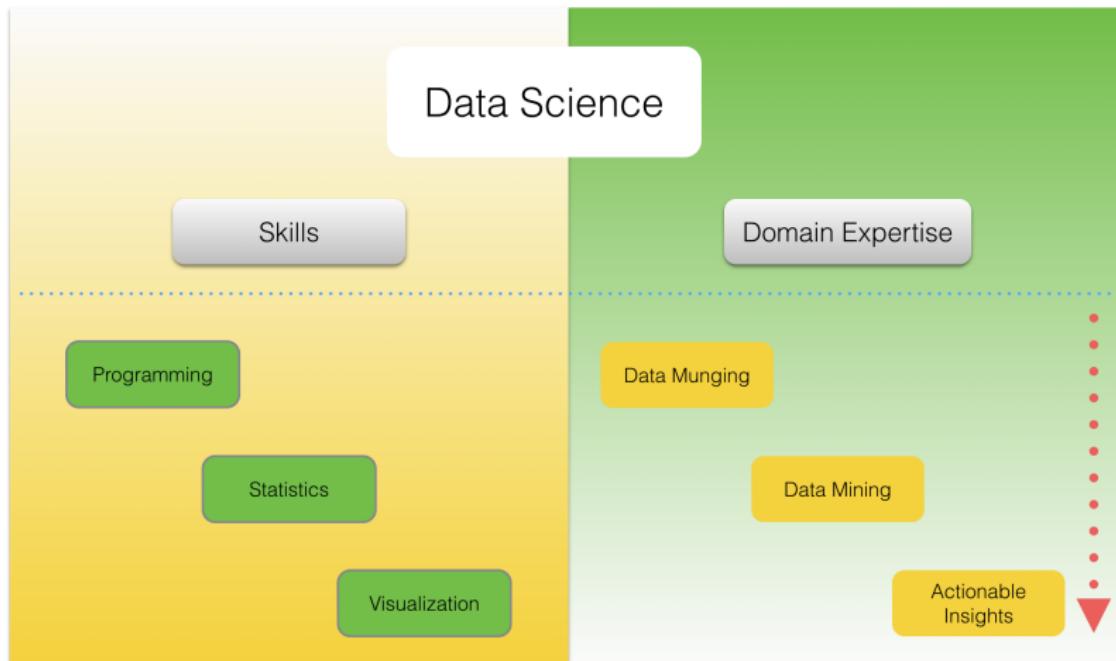


Table of Contents

- 1 What is Data Science
 - Overview
 - Data Science Workflow

 - 2 Churn Model
 - Business Understanding
 - Demo in R
 - Top Features
 - Prescriptive Analysis

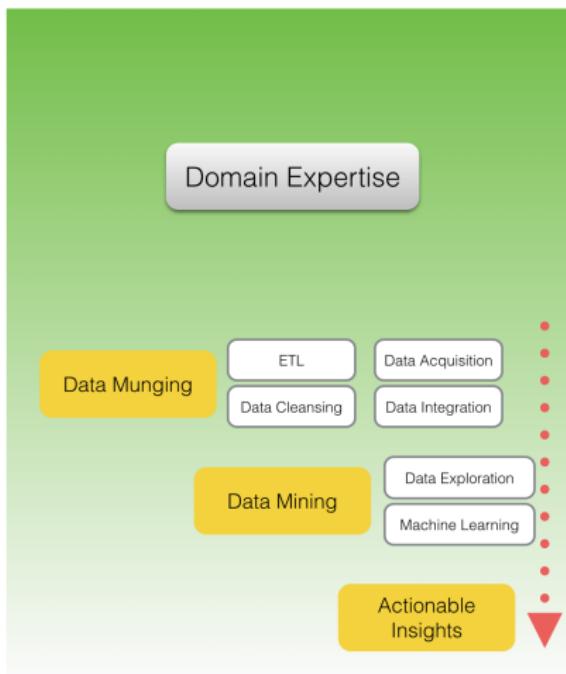
 - 3 Yesware Email Analysis
 - Email Analysis

 - 4 Data Science in the Industry
 - Data Science Toolbox
 - Big Data

Data Science Workflow

Data Science Workflow

- Data Munging
 - Data Mining
 - Delivery of actionable Insights



Data Munging

Data Munging

Data Munging means some or all of the following tasks:

- ETL
- Data Integration
- Data Cleansing

Data Munging

Data Munging

Data Munging means some or all of the following tasks:

- ETL
 - Data Integration
 - Data Cleansing

ETL

The process of extract,
transform, and load data.

- To acquire data from external sources.
 - To migrate multiple data sources internally.

Data Munging

Data Munging

Data Munging means some or all of the following tasks:

- ETL
 - Data Integration
 - Data Cleansing

Data Integration

To combine data from disparate sources into meaningful and valuable information.

Data Munging

Data Munging

Data Munging means some or all of the following tasks:

- ETL
 - Data Integration
 - Data Cleansing

Data Cleansing

Data cleansing, also called data scrubbing, is the process of amending or removing data in a database that is incorrect, incomplete, improperly formatted, or duplicated.

Data Mining

Data Mining

Data Mining is the key step to turn data into insights:

- Data Exploration
 - Machine Learning
 - Model Evaluation

Data Mining

Data Mining

Data Mining is the key step to turn data into insights:

- Data Exploration
 - Machine Learning
 - Model Evaluation

Data Exploration

The process of visually examining and exploring the data.

- To gain basic understanding of the data.
 - To identify relationships between different attributes.
 - To answer basic questions using data.

Data Mining

Data Mining

Data Mining is the key step to turn data into insights:

- Data Exploration
 - Machine Learning
 - Model Evaluation

Machine Learning

To obtain statistical models, we usually need to go through multiple steps like the following:

- ① Construct new features.
 - ② Remove redundant features.
 - ③ Choose one or more suitable machine learning algorithm.

Data Mining

Data Mining

Data Mining is the key step to turn data into insights:

- Data Exploration
 - Machine Learning
 - Model Evaluation

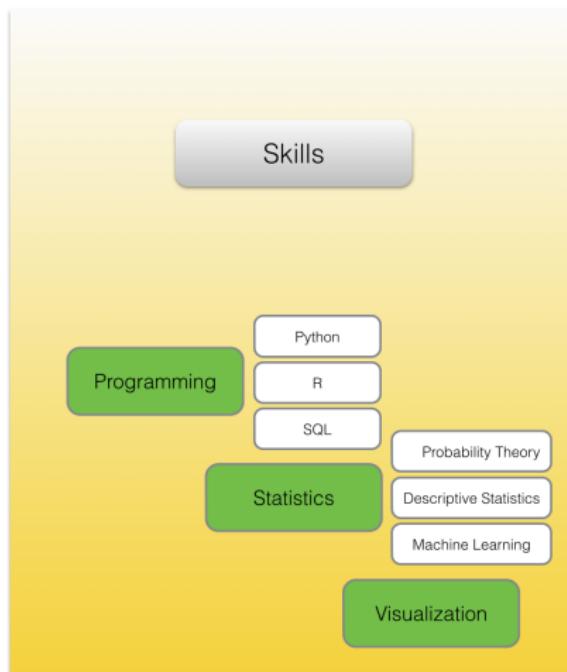
Model Evaluation

Model evaluation is often used not only to select the best model from the set of models, but also to get ready for producing actionable insights.

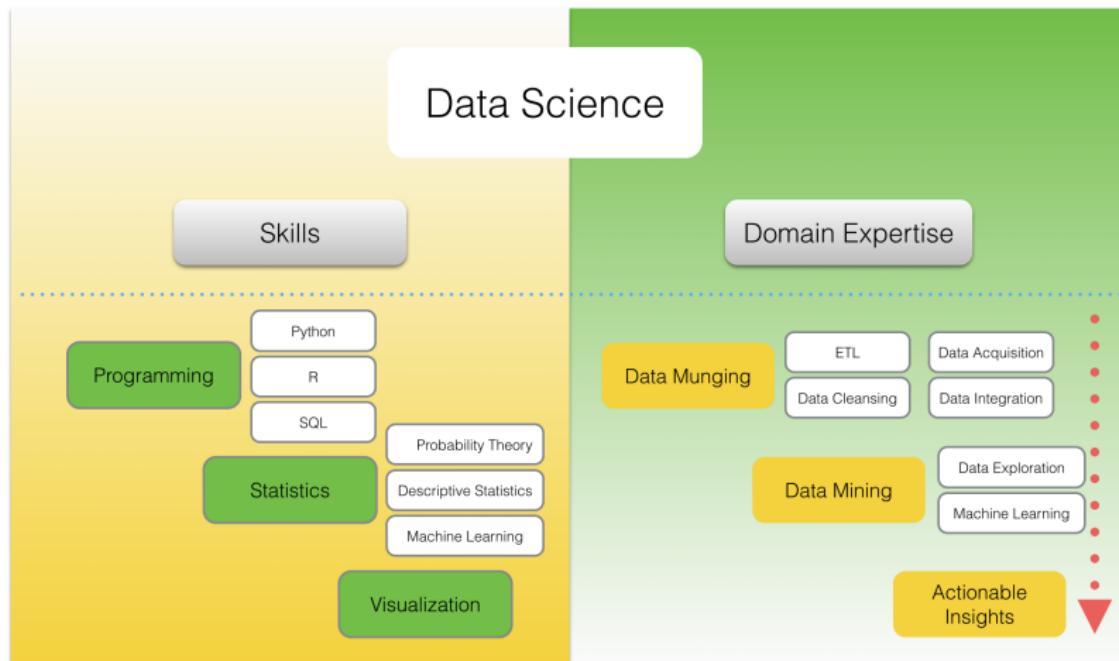
Skills of a data scientist

Data Science skills

- Programming
 - Statistics
 - Visualization



Doing data science



Business Understanding

Table of Contents

1 What is Data Science

- Overview
- Data Science Workflow

2 Churn Model

- Business Understanding
- Demo in R
- Top Features
- Prescriptive Analysis

3 Yesware Email Analysis

- Email Analysis

4 Data Science in the Industry

- Data Science Toolbox
- Big Data

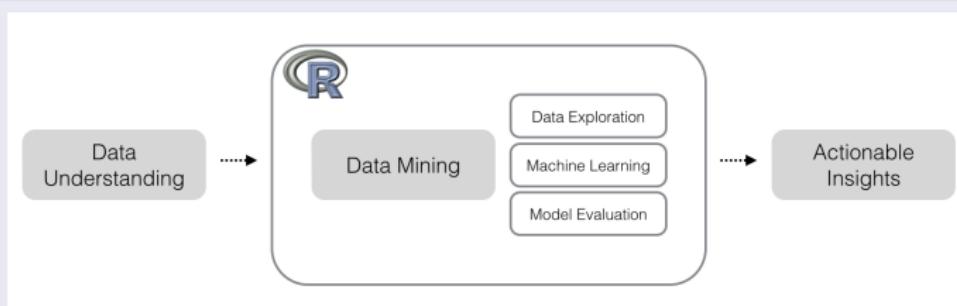
Business Understanding

The Goal

Goal

- To predict which customers will churn.

Workflow

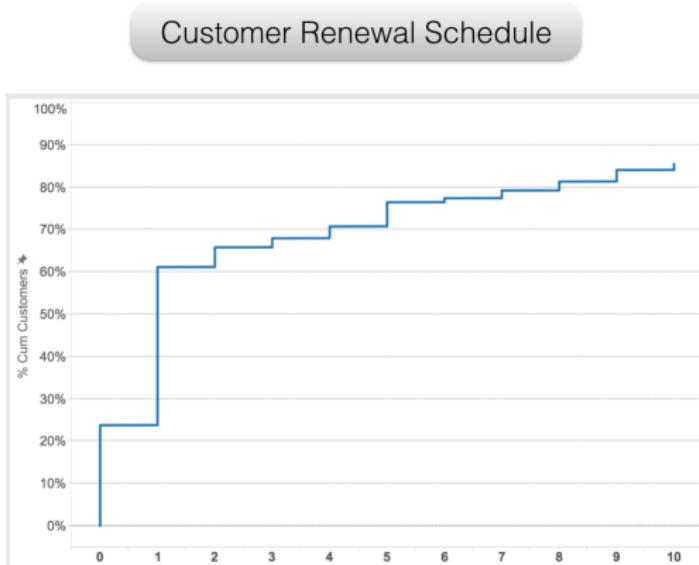


Business Understanding

Customer renewal data

| customer_id | is_churn | days_renew |
|-------------|----------|------------|
| 213 | FALSE | 5 |
| 102 | TRUE | NULL |
| 31 | TRUE | NULL |
| 921 | FALSE | 5 |
| ... | ... | ... |

- **Days_Renew** is the number of days after subscription expiration before the customer renewed.
 - If the customer has churned, then this value will be NULL.



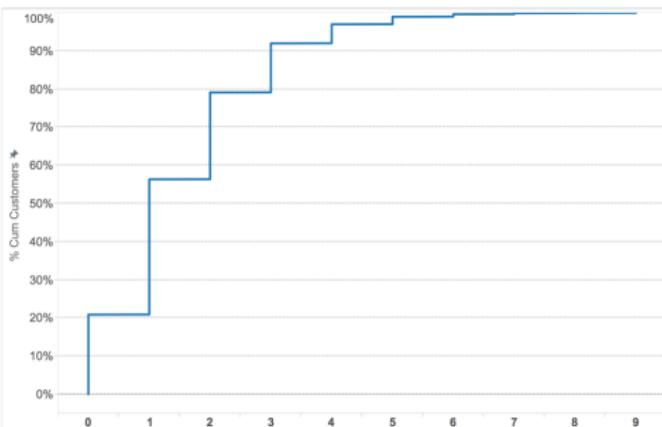
Business Understanding

Customer service call data

| customer_id | cust_surv_calls |
|-------------|-----------------|
| 213 | 0 |
| 102 | 3 |
| 31 | 1 |
| 921 | 2 |
| ... | ... |

- 80% of customers made at most 2 phone calls to the customer service.
 - On average, customers made 1.6 customer service calls.
 - Some customers made as many as 9 calls to customer service.

Customer Call
Cumulative Distribution



Business Understanding

Customer plan type and demographic

| customer_id | state | account_len | area_code | is_intl_plan | is_vmail_plan |
|-------------|-------|-------------|-----------|--------------|---------------|
| 1 | KS | 128 | 415 | no | yes |
| 2 | OH | 107 | 415 | no | yes |
| 3 | NJ | 137 | 415 | no | no |
| 4 | OH | 84 | 408 | yes | no |
| 5 | OK | 75 | 415 | yes | no |
| 6 | AL | 118 | 510 | yes | no |
| 7 | MA | 121 | 510 | no | yes |
| 8 | MO | 147 | 415 | yes | no |
| 9 | LA | 117 | 408 | no | no |
| 10 | WV | 141 | 415 | yes | yes |
| ... | ... | ... | ... | ... | ... |

Business Understanding

Customer usage data

Business Understanding

Churn data structure

```
$ customer_id      : int  1 2 3 4 5 6 7 8 9 10 ...
$ state            : Factor w/ 51 levels "AK","AL","AR",...: 17 36 32 36 37 2 20 25 19 50 ...
$ account_len     : int  128 107 137 84 75 118 121 147 117 141 ...
$ area_code        : int  415 415 415 408 415 510 510 415 408 415 ...
$ phone            : Factor w/ 3333 levels "327-1058","327-1319",...: 1927 1576 1118 1708 111 2254 1048 81
$ is_intl_plan    : Factor w/ 2 levels "no","yes": 1 1 1 2 2 2 1 2 1 2 ...
$ is_vmail_plan   : Factor w/ 2 levels "no","yes": 2 2 1 1 1 2 1 1 2 ...
$ vmail_messages  : int  25 26 0 0 0 24 0 0 37 ...
$ day_mins         : num  265 162 243 299 167 ...
$ day_calls        : int  110 123 114 71 113 98 88 79 97 84 ...
$ day_charge       : num  45.1 27.5 41.4 50.9 28.3 ...
$ eve_mins         : num  197.4 195.5 121.2 61.9 148.3 ...
$ eve_calls        : int  99 103 110 88 122 101 108 94 80 111 ...
$ eve_charge       : num  16.78 16.62 10.3 5.26 12.61 ...
$ night_mins       : num  245 254 163 197 187 ...
$ night_calls      : int  91 103 104 89 121 118 118 96 90 97 ...
$ night_charge     : num  11.01 11.45 7.32 8.86 8.41 ...
$ intl_mins        : num  10 13.7 12.2 6.6 10.1 6.3 7.5 7.1 8.7 11.2 ...
$ intl_calls       : int  3 3 5 7 3 6 7 6 4 5 ...
$ intl_charge      : num  2.7 3.7 3.29 1.78 2.73 1.7 2.03 1.92 2.35 3.02 ...
$ cust_surv_calls: int  1 1 0 2 3 0 3 0 1 0 ...
$ is_churn          : Factor w/ 2 levels "False.","True.": 1 1 1 1 1 1 1 1 1 1 ...
$ days_renew       : int  0 0 0 0 0 0 0 0 0 0 ...
```

Table of Contents

1 What is Data Science

- Overview
- Data Science Workflow

2 Churn Model

- Business Understanding
- Demo in R
- Top Features
- Prescriptive Analysis

3 Yesware Email Analysis

- Email Analysis

4 Data Science in the Industry

- Data Science Toolbox
- Big Data

Demo in R

R Script

<http://goo.gl/IV3HDs>

Top Features

Table of Contents

1 What is Data Science

- Overview
- Data Science Workflow

2 Churn Model

- Business Understanding
- Demo in R
- Top Features
- Prescriptive Analysis

3 Yesware Email Analysis

- Email Analysis

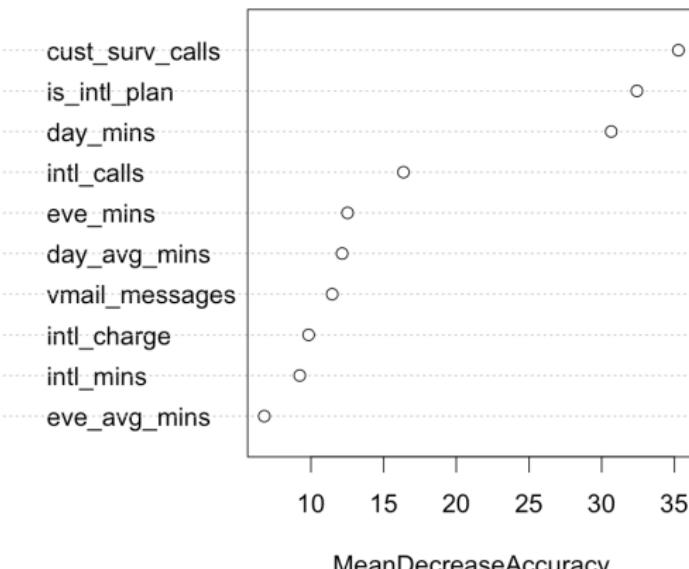
4 Data Science in the Industry

- Data Science Toolbox
- Big Data

Top Features

Variable importance from random forest

Top Features from Random Forest

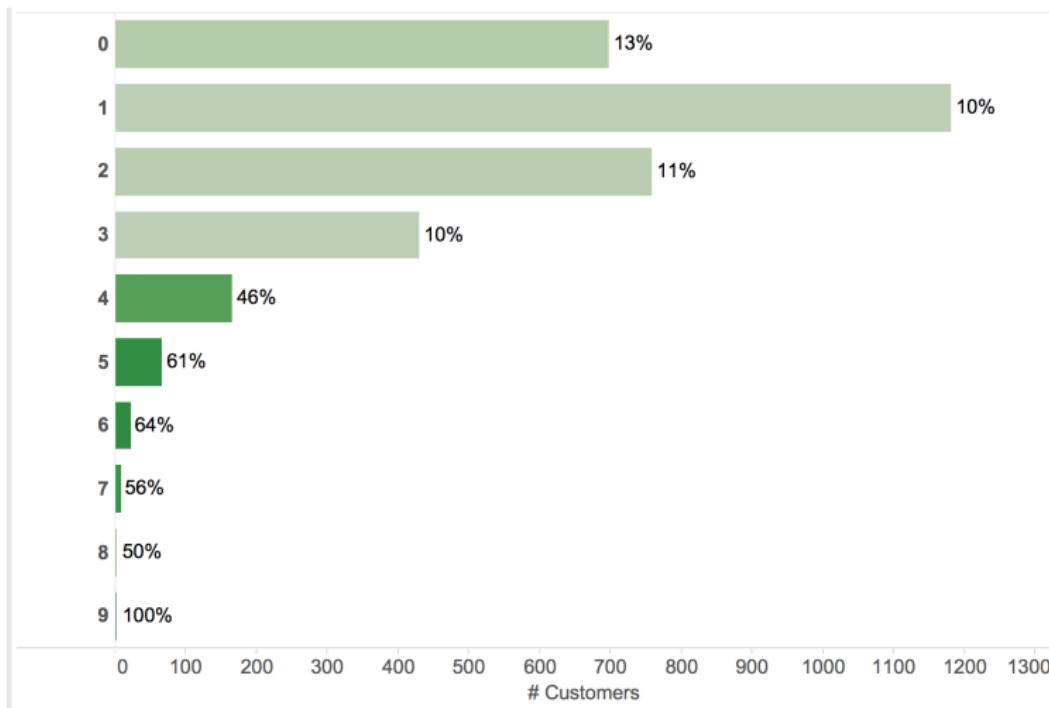


Top Features

- cust surv calls
- is intl plan
- day mins

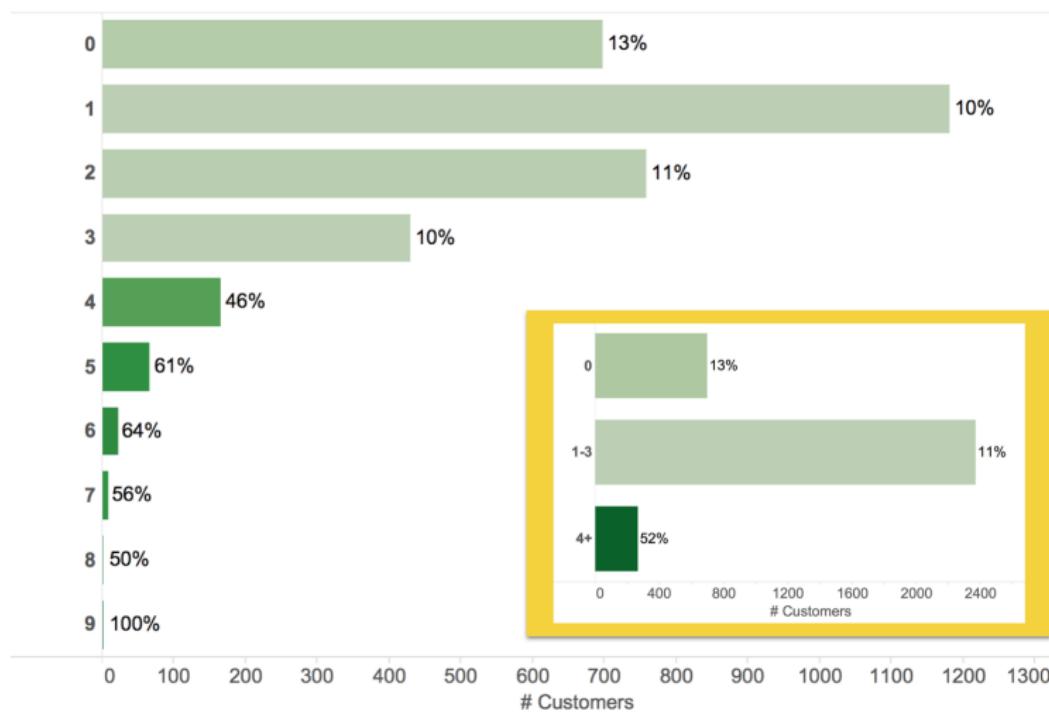
Top Features

Top feature - customer service calls



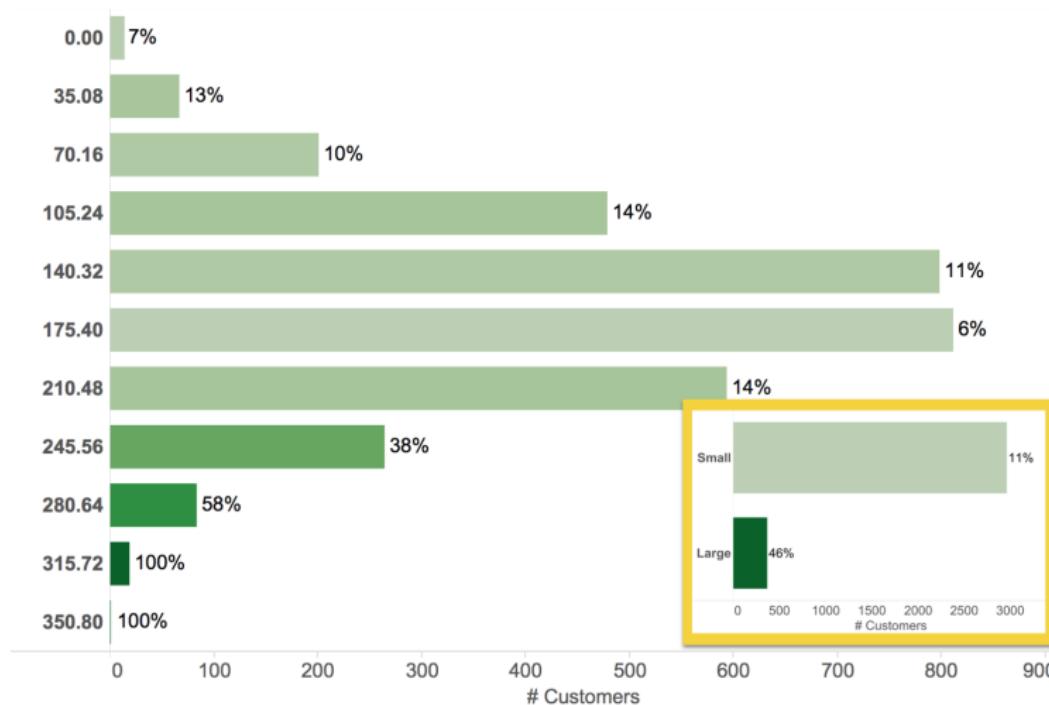
Top Features

Top feature - customer service calls



Top Features

Top feature - customer day time minutes



Top Features

Top feature - international plan and international calls

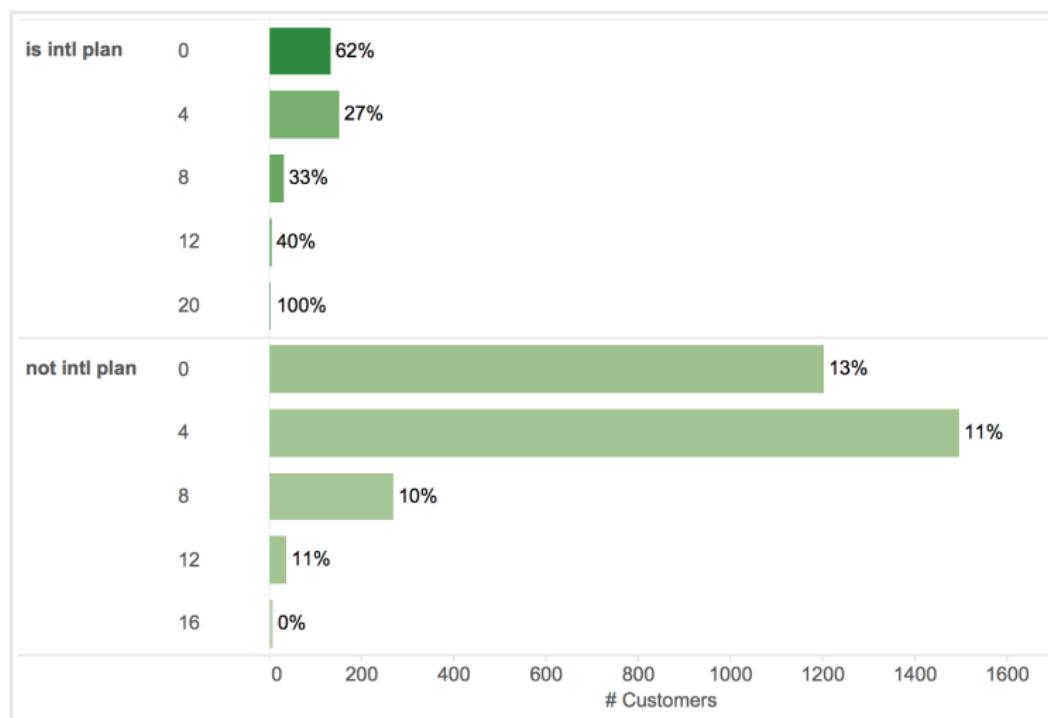


Table of Contents

1 What is Data Science

- Overview
- Data Science Workflow

2 Churn Model

- Business Understanding
- Demo in R
- Top Features
- Prescriptive Analysis

3 Yesware Email Analysis

- Email Analysis

4 Data Science in the Industry

- Data Science Toolbox
- Big Data

Prescriptive Analysis

Churn analysis

| Is Intl Plan | Cust Surv Calls (group) | Day Mins Group | % Churn | Number of Records | Summary |
|--------------------|-------------------------|----------------|------------|-------------------|-------------|
| is intl plan | 0 | Large | 60% | 15 | Small Group |
| | | Small | 44% | 68 | Small Group |
| | 1-3 | Large | 56% | 36 | Small Group |
| | | Small | 34% | 176 | |
| | 4+ | Large | 100% | 1 | Outlier |
| | | Small | 67% | 27 | Small Group |
| | not intl plan | 0 | Large | 48% | 60 |
| | | | 4% | 554 | Low Churn |
| | 1-3 | Large | 43% | 229 | |
| | | Small | 4% | 1,928 | Low Churn |
| | 4+ | Large | 44% | 27 | Small Group |
| | | Small | 50% | 212 | |
| Grand Total | | | 14% | 3,333 | |

Prescriptive Analysis

Churn analysis

| Is Intl Plan | Cust Surv Calls (group) | Day Mins Group | % Churn | # Customers | Summary | Avg Charge | Revenue if we got 10% of the churned customers back |
|--------------------|-------------------------|----------------|------------|--------------|-------------|-------------|---|
| is intl plan | 0 | Large | 60% | 15 | Small Group | \$ 72 | \$ 65 |
| | | Small | 44% | 68 | Small Group | \$ 56 | \$ 168 |
| | 1-3 | Large | 56% | 36 | Small Group | \$ 72 | \$ 144 |
| | | Small | 34% | 176 | | \$ 55 | \$ 325 |
| | 4+ | Large | 100% | 1 | Outlier | \$ 78 | \$ 8 |
| | | Small | 67% | 27 | Small Group | \$ 56 | \$ 101 |
| | not intl plan | 0 | Large | 48% | 60 | Small Group | \$ 73 |
| | | | Small | 4% | 554 | Low Churn | \$ 55 |
| | | 1-3 | Large | 43% | 229 | | \$ 72 |
| | | | Small | 4% | 1,928 | Low Churn | \$ 55 |
| | | 4+ | Large | 44% | 27 | Small Group | \$ 74 |
| | | | Small | 50% | 212 | | \$ 54 |
| Grand Total | | | 14% | 3,333 | | \$ 57 | \$ 2,660 |

Table of Contents

1 What is Data Science

- Overview
- Data Science Workflow

2 Churn Model

- Business Understanding
- Demo in R
- Top Features
- Prescriptive Analysis

3 Yesware Email Analysis

- Email Analysis

4 Data Science in the Industry

- Data Science Toolbox
- Big Data

When to send an email

To Optimize Email Reply

- Yesware users want to know how to get more replies.

When to send an email

To Optimize Email Reply

- Yesware users want to know how to get more replies.

An Email Reply Model

- ➊ Construct features from the email data.
- ➋ Create a model to predict the reply on each email.
- ➌ Identify some top features that contributed most to the reply:
 - Sent Hour
 - Sent Weekday

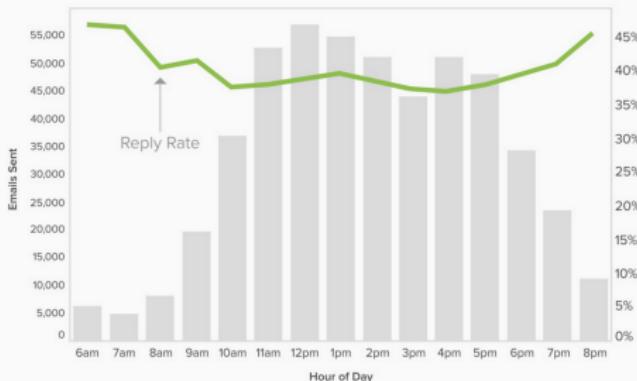
Email Analysis

When to send an email

To Optimize Email Reply

- Yesware users want to know how to get more replies.
- Emails Sent and Reply Rate by Sent Hour

Send Emails in the Early Morning or Evening



When to send an email

To Optimize Email Reply

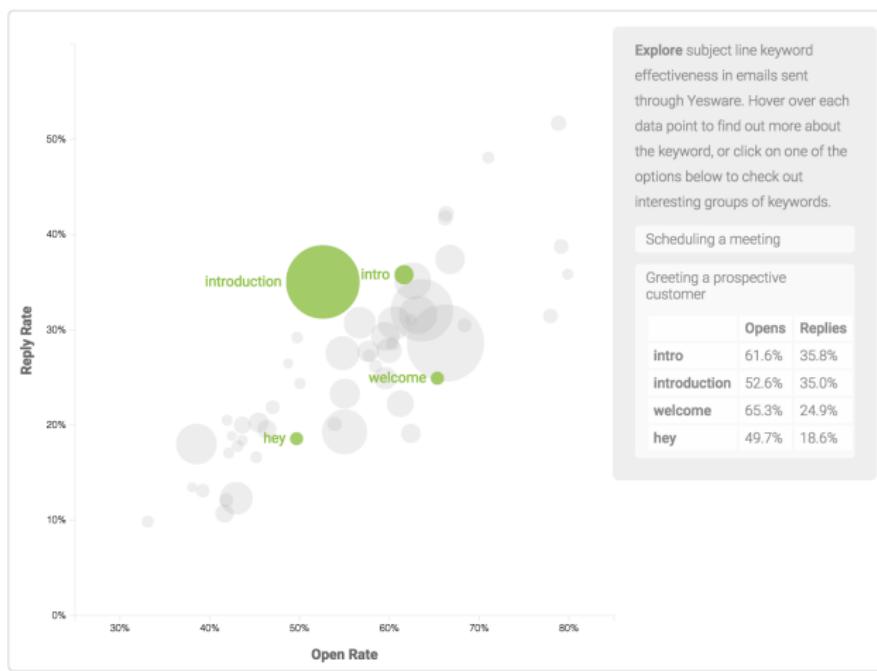
- Yesware users want to know how to get more replies.
- Emails Sent and Reply Rate by Sent Hour
- Emails Sent and Reply Rate by Sent Weekday

Email Reply Rates are Highest
on the Weekends

| | Emails Sent | % Open | % Reply | % Reply Same Day |
|----------|-------------|--------|---------|------------------|
| Week Day | 525,742 | 66.3% | 39.1% | 33.1% |
| Weekend | 5,278 | 73.6% | 45.8% | 32.6% |

Email Analysis

D3 Visualization: Subject Line Keywords



Subject line keywords: <http://goo.gl/PK9xh0>

Table of Contents

1 What is Data Science

- Overview
- Data Science Workflow

2 Churn Model

- Business Understanding
- Demo in R
- Top Features
- Prescriptive Analysis

3 Yesware Email Analysis

- Email Analysis

4 Data Science in the Industry

- Data Science Toolbox
- Big Data

Why do we need a toolbox?

| | Academia | Industry |
|---------------------|--|--|
| Goal | Improve human knowledge | Make money |
| Success Criteria | Publish papers | Create and deliver business value |
| Approach | Finding a better way to do a new thing | Finding the fastest way to do lots of things |
| Importance of Speed | Not the most important | Very important |

Tools that helped me do data science fast

| | Python | R | Unix | SQL | Scala |
|-----------------------------|--------|-------|------|------|-------|
| Powerful Packages / Library | ***** | ***** | ** | * | ***** |
| Community Support | ***** | ***** | **** | *** | ***** |
| Data Munging | **** | *** | **** | **** | ***** |
| Data Exploration | **** | ***** | ** | *** | *** |
| Machine Learning | **** | ***** | * | ** | ***** |

Data visualization tools

| | Excel | R | Tableau | D3 |
|------------------------------------|-------|------|---------|-------|
| Ease of Learning | ***** | *** | ***** | * |
| Is Free | No | Yes | No | Yes |
| Good for Data Exploration | ** | **** | ***** | *** |
| Flexibility in Data Representation | ** | **** | **** | ***** |
| Good for Reporting and Sharing | **** | **** | *** | **** |

D3: <http://d3js.org/>

Table of Contents

1 What is Data Science

- Overview
- Data Science Workflow

2 Churn Model

- Business Understanding
- Demo in R
- Top Features
- Prescriptive Analysis

3 Yesware Email Analysis

- Email Analysis

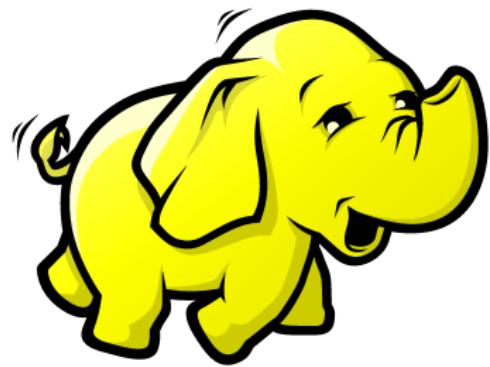
4 Data Science in the Industry

- Data Science Toolbox
- Big Data

Big Data

Big Data Ecosystem

- Hadoop - file system



Big Data

Big Data

Big Data Ecosystem

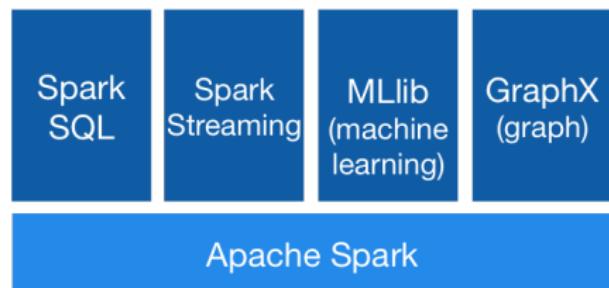
- Hadoop - file system
- Spark - computing system



Big Data

Big Data Ecosystem

- Hadoop - file system
- Spark - computing system
- Spark Stack



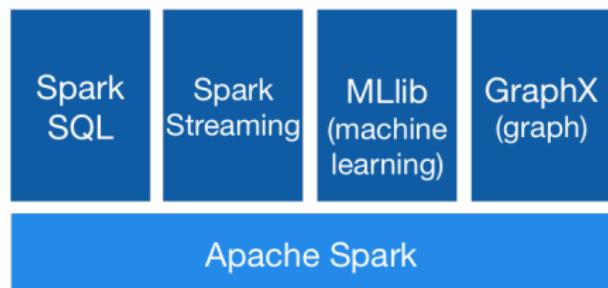
<https://spark.apache.org/>

Big Data

Big Data

Big Data Ecosystem

- Hadoop - file system
- Spark - computing system
- Spark Stack
 - Spark SQL - Data Munging
 - Spark Streaming - Real Time Processing
 - MLlib - Machine Learning
 - GraphX - Visualization



<https://spark.apache.org/>

Big Data

Thank You

Please send your questions and feedback to me!