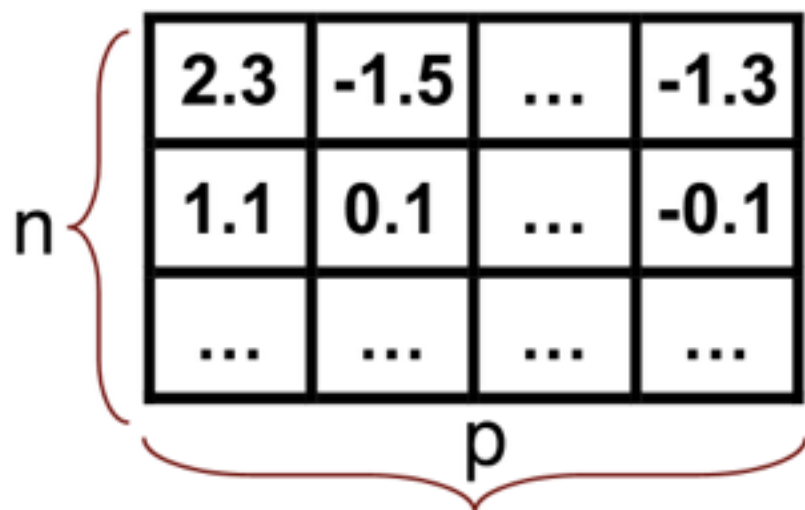


DATA SCIENCE

AN INTRODUCTION

- What are we trying to achieve?
- What is a data scientist?

Types of Data: Flat File Data



A diagram illustrating a data matrix. It consists of a 3x4 grid of cells. The first two rows contain numerical values, and the third row contains ellipses. A red curly brace on the left side of the grid is labeled 'n', indicating the number of rows (objects). A red curly brace at the bottom of the grid is labeled 'p', indicating the number of columns (measurements). The values in the first row are 2.3, -1.5, ..., -1.3. The values in the second row are 1.1, 0.1, ..., -0.1. The values in the third row are ..., ..., ...,

2.3	-1.5	...	-1.3
1.1	0.1	...	-0.1
...

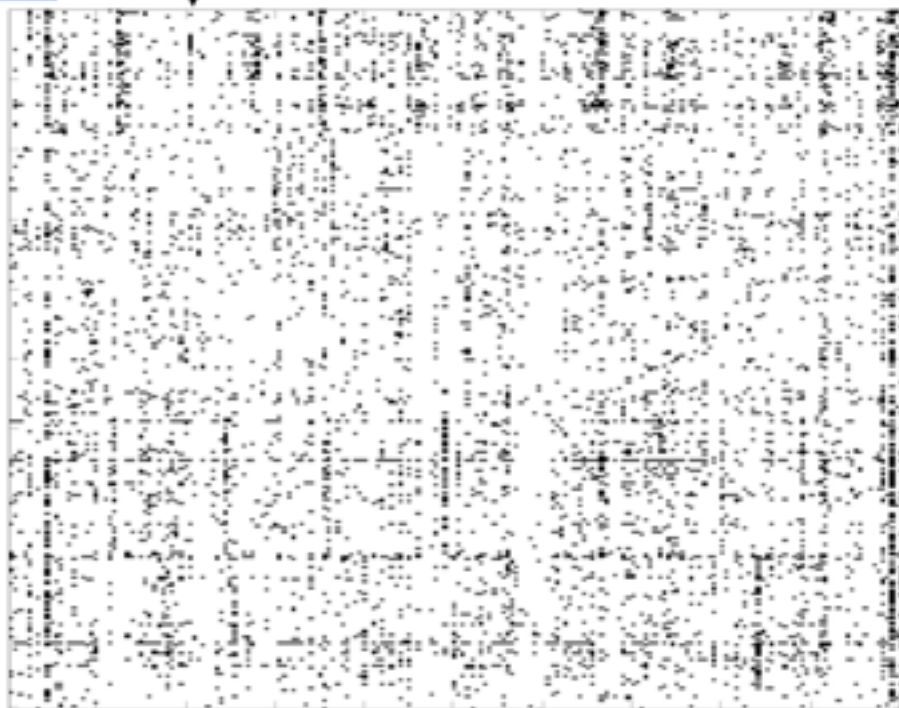
- Rows = objects
- Columns = measurements on objects
- Both n and p can be very large in data mining (also $p \gg n$)
- Matrix can be quite sparse

Types of Data: Text Data

Can be
represented as a
sparse matrix

Obama

Text
Documents



Word ID

Types of Data: Transactional Data

Date stamped events (logs, phone calls):

128.195.36.195, -, 3/22/00, 10:35:11, W3SVC, SRVR1, 128.200.39.181, 781, 363, 875, 200, 0, GET, /top.html, -,
128.195.36.195, -, 3/22/00, 10:35:16, W3SVC, SRVR1, 128.200.39.181, 5288, 524, 414, 200, 0, POST, /spt/main.html, -,
128.195.36.195, -, 3/22/00, 10:35:17, W3SVC, SRVR1, 128.200.39.181, 30, 280, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
128.195.36.101, -, 3/22/00, 16:18:50, W3SVC, SRVR1, 128.200.39.181, 60, 425, 72, 304, 0, GET, /top.html, -,
128.195.36.101, -, 3/22/00, 16:18:58, W3SVC, SRVR1, 128.200.39.181, 8322, 527, 414, 200, 0, POST, /spt/main.html, -,
128.195.36.101, -, 3/22/00, 16:18:59, W3SVC, SRVR1, 128.200.39.181, 0, 280, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
128.200.39.17, -, 3/22/00, 20:54:37, W3SVC, SRVR1, 128.200.39.181, 140, 199, 875, 200, 0, GET, /top.html, -,
128.200.39.17, -, 3/22/00, 20:54:55, W3SVC, SRVR1, 128.200.39.181, 17766, 365, 414, 200, 0, POST, /spt/main.html, -,
128.200.39.17, -, 3/22/00, 20:54:55, W3SVC, SRVR1, 128.200.39.181, 0, 258, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
128.200.39.17, -, 3/22/00, 20:55:07, W3SVC, SRVR1, 128.200.39.181, 0, 258, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
128.200.39.17, -, 3/22/00, 20:55:36, W3SVC, SRVR1, 128.200.39.181, 1061, 382, 414, 200, 0, POST, /spt/main.html, -,
128.200.39.17, -, 3/22/00, 20:55:36, W3SVC, SRVR1, 128.200.39.181, 0, 258, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
128.200.39.17, -, 3/22/00, 20:55:39, W3SVC, SRVR1, 128.200.39.181, 0, 258, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
128.200.39.17, -, 3/22/00, 20:56:03, W3SVC, SRVR1, 128.200.39.181, 1081, 382, 414, 200, 0, POST, /spt/main.html, -,
128.200.39.17, -, 3/22/00, 20:56:04, W3SVC, SRVR1, 128.200.39.181, 0, 258, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
128.200.39.17, -, 3/22/00, 20:56:33, W3SVC, SRVR1, 128.200.39.181, 0, 262, 72, 304, 0, GET, /top.html, -,
128.200.39.17, -, 3/22/00, 20:56:52, W3SVC, SRVR1, 128.200.39.181, 19598, 382, 414, 200, 0, POST, /spt/main.html, -,

Can be represented as a time series:

User 1	2	3	2	2	3	3	3	1	1	1	3	1	3	3	3
User 2	3	3	3	1	1	1									
User 3	7	7	7	7	7	7	7	7							
User 4	1	5	1	1	1	5	1	5	1	1	1	1	1	1	
User 5	5	1	1	5											
...															

Types of Data: Relational Data

128.200.39.17, -, 3/22/00, 20:55:07, W3SVC, SRVR1, 128.200.39.181, 0, 258, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
128.200.39.17, -, 3/22/00, 20:55:36, W3SVC, SRVR1, 128.200.39.181, 1061, 382, 414, 200, 0, POST, /spt/main.html, -,
128.200.39.17, -, 3/22/00, 20:55:36, W3SVC, SRVR1, 128.200.39.181, 0, 258, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
128.195.36.195, -, 3/22/00, 10:35:11, W3SVC, SRVR1, 128.200.39.181, 781, 363, 875, 200, 0, GET, /top.html, -,
128.195.36.195, -, 3/22/00, 10:35:16, W3SVC, SRVR1, 128.200.39.181, 5288, 524, 414, 200, 0, POST, /spt/main.html, -,
128.195.36.195, -, 3/22/00, 10:35:17, W3SVC, SRVR1, 128.200.39.181, 30, 280, 111, 404, 3, GET, /spt/images/bk1.jpg, -,
...

128.195.36.195, Doe, John, 12 Main St, 973-462-3421, Madison, NJ, **07932**
114.12.12.25, Trank, Jill, 11 Elm St, 998-555-5675, Chester, NJ, 07911
...

07911, Chester, NJ, 07954, 34000, , 40.65, -74.12
07932, Madison, NJ, 56000, 40.642, -74.132
...

- Most large data sets are stored in relational data sets
- Special data query language: SQL

Types of Data: Time Series Data



Types of Data: Image Data



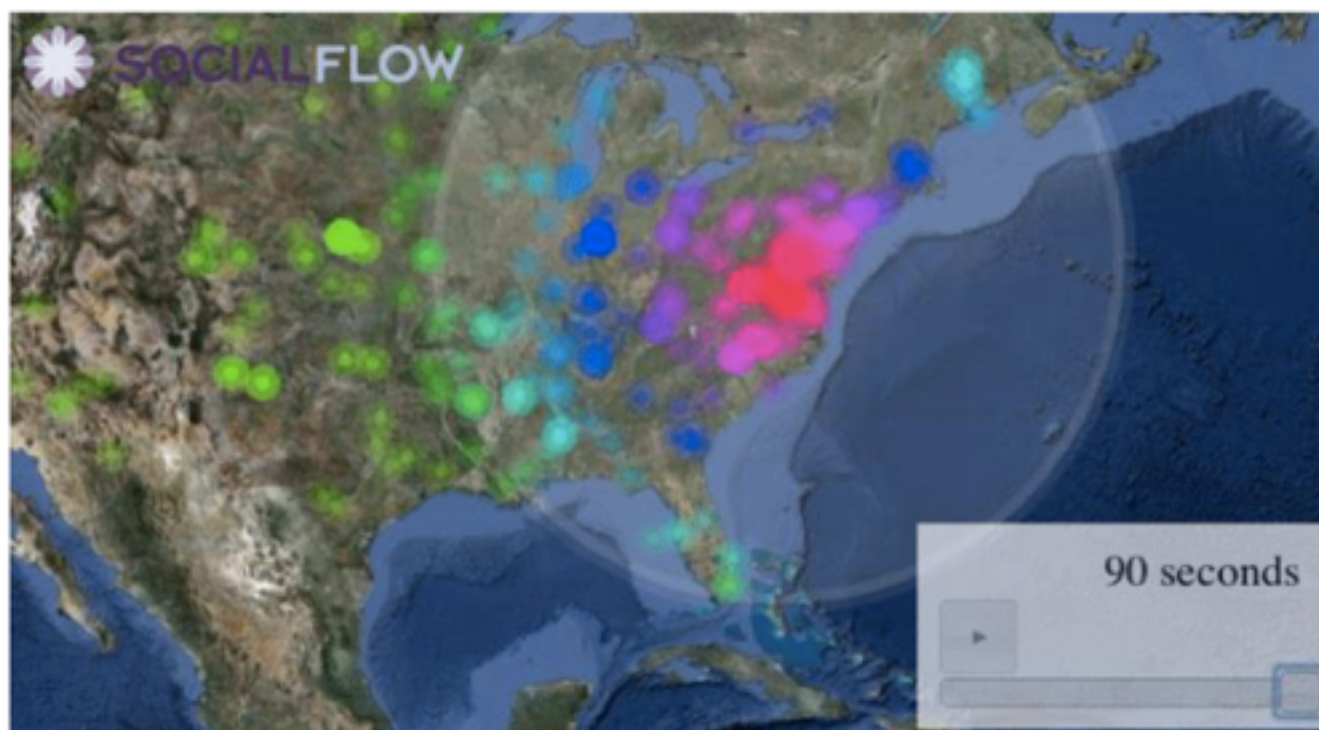
Types of Data: Spatio-Temporal Data



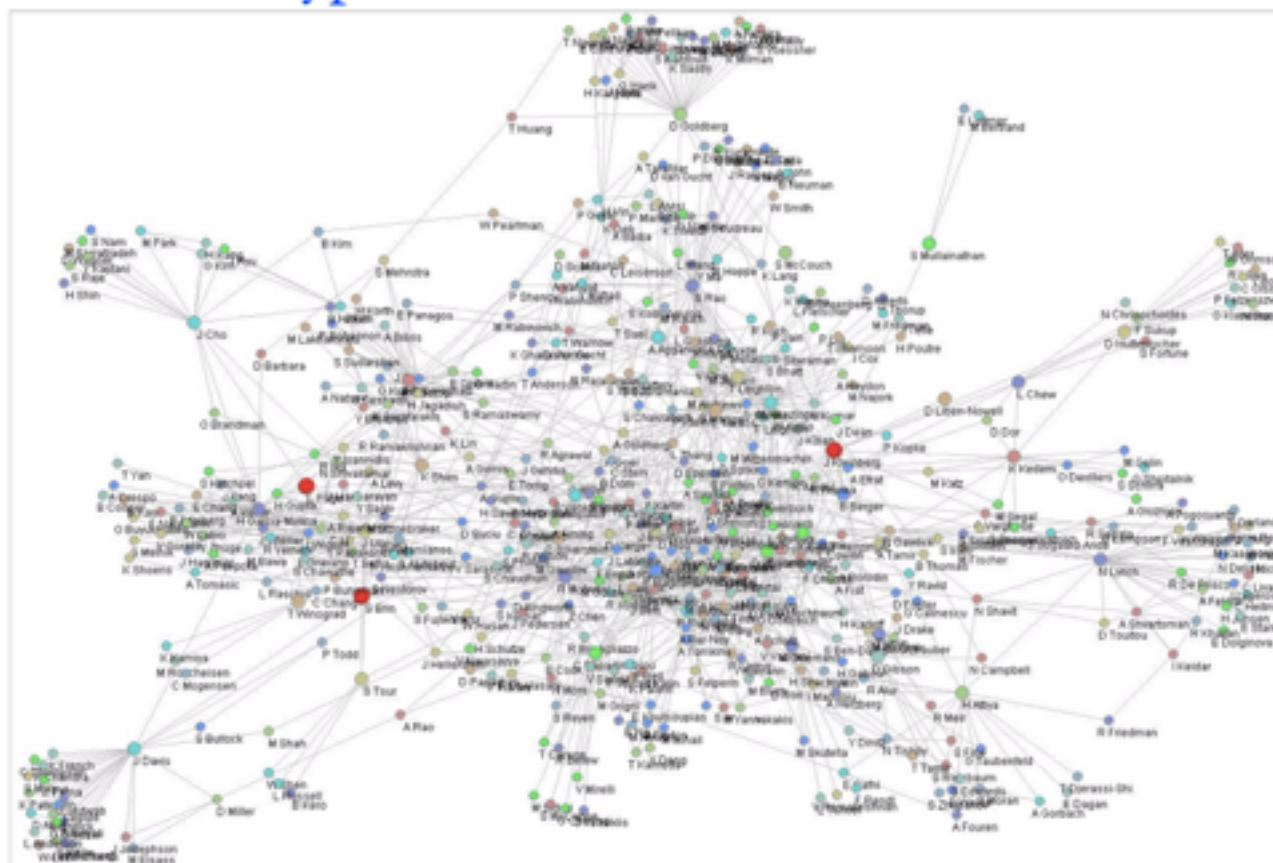
@b_mc817

Glendaaaaa

Omg earthquake!!!



Types of Data: Network Data



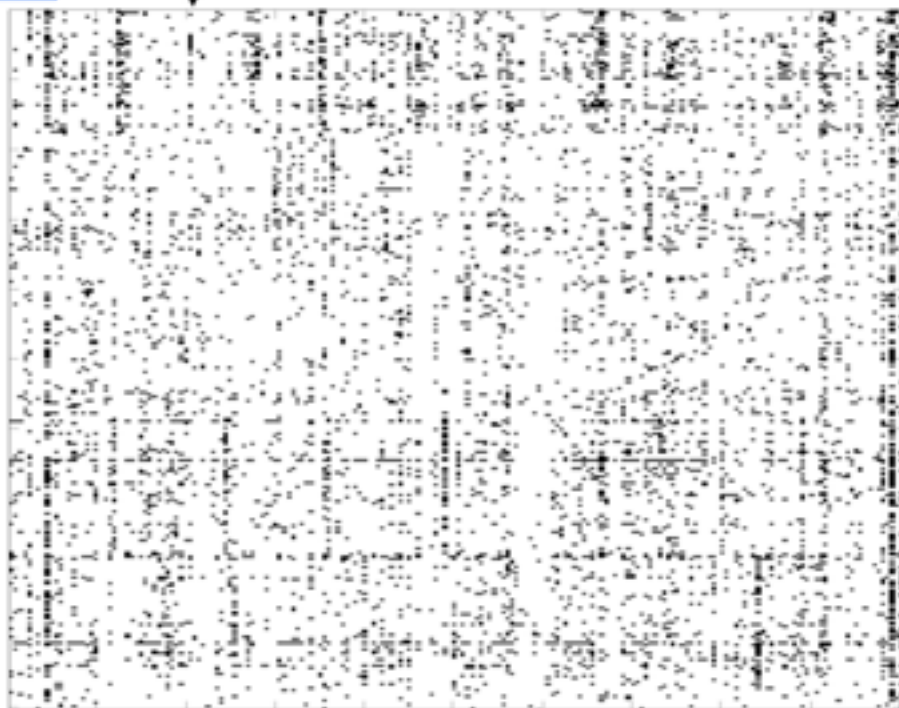
Algorithms for estimating relative importance in networks
S. White and P. Smyth, *ACM SIGKDD*, 2003.

Types of Data: Text Data

Can be
represented as a
sparse matrix

Obama

Text
Documents



Word ID

- What are we trying to achieve?
- **What is a data scientist?**

DATA SCIENCE



Zvi
@nivertech



 Follow

"Data Scientist" is a Data Analyst who lives in California.

 Reply  Retweet  Favorite  More

RETWEETS
140

FAVORITES
40



9:55 PM - 14 Mar 2012

DATA SCIENCE



Javier Nogales
@fjnogales



Follow

Data Scientist (2/2): person who is worse at statistics than any statistician and worse at software engineering than any software engineer



RETWEET

1

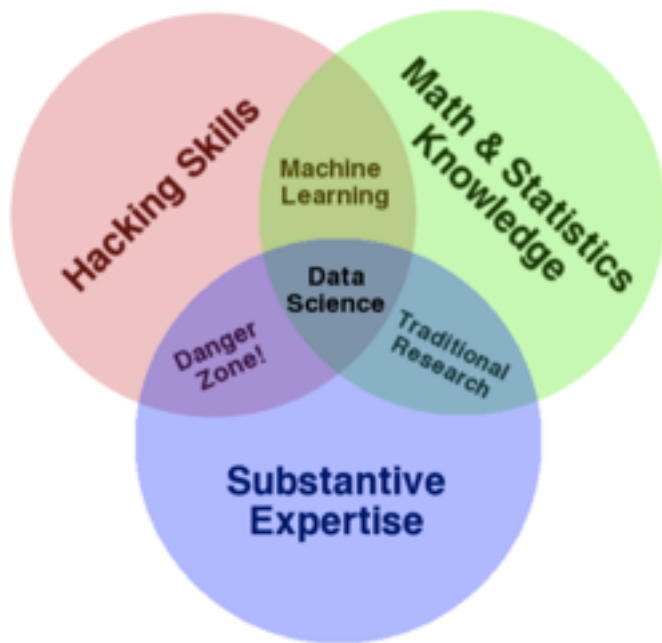
FAVORITES

5



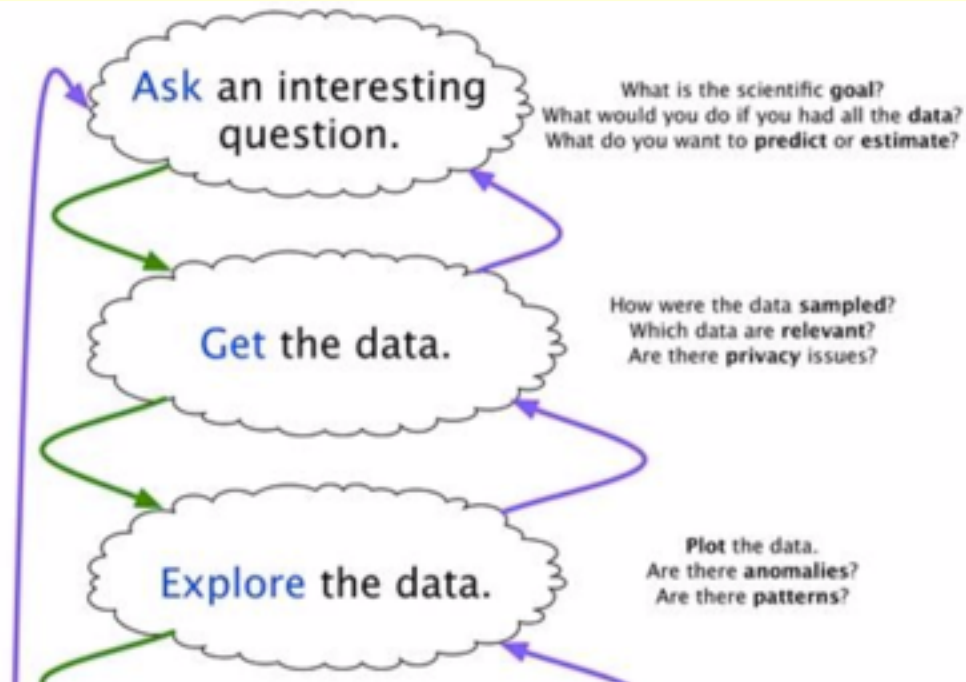
9:08 AM - 27 Jan 2014

DATA SCIENCE

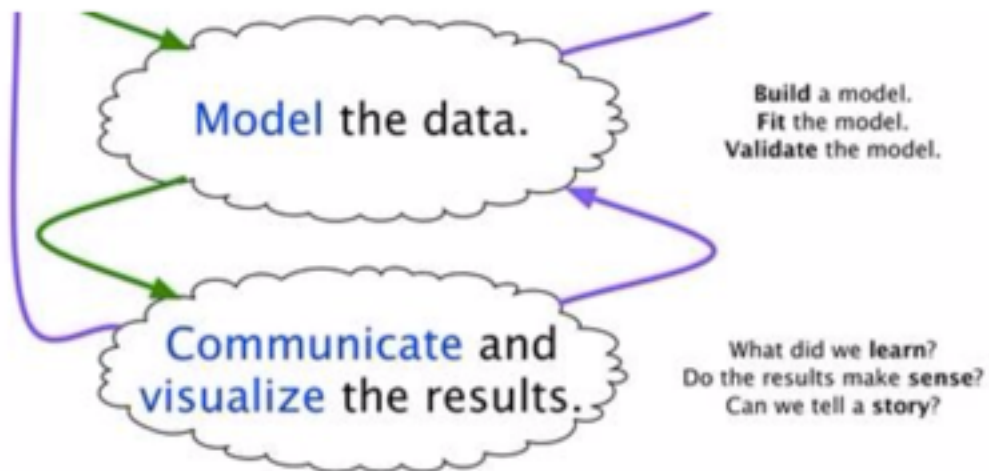


Wide variance in terms of skillsets: many job descriptions are more appropriate for a **team of data scientists!**

DATA SCIENCE



DATA SCIENCE



Problem: Children born prematurely are at high risk of developing infections, many of which are not detected until after the baby is sick

Goal: Detect subtle patterns in the data that predicts infection before it occurs

Data: 16 vital signs such as heart rate, respiration rate, blood pressure, etc...

Impact: Model is able to predict the onset of infection 24 hours before the traditional symptoms of infection appear



Image: <http://www.babycaretips4u.com/wp-content/uploads/2014/03/premature-baby.jpg>

Case Study: <http://www.amazon.com/Big-Data-Revolution-Transform-Think/dp/0544002695>

Problem: Processing disability claims at the Social Security Administration is a time-intensive process, with many claims taking over 2 years to adjudicate

Goal: Automate the approval of a subset of the “simplest” disability claims

Data: Free text in the claims form

Impact: Able to fully automate 20% of the simplest claims. Rating accuracy of the algorithm is higher than the average claims examiner.



QUESTIONS?