# Big Data Analytics in
# Large Scale Socio-economic Development

Ihab F. Ilyas

University of Waterloo
@ihabilyas

# Data Science: Many Definitions and One Goal

- **Extract Value from Data**



Statistics + Machine Learning + Data Management + Systems ..

# Case 1
## The Effectiveness of the World Bank Funded Projects

# The World Bank

- Increase the transparency and accountability of international development projects

- Visualize the location of Bank-financed projects to better monitor development impact

- Integrate location data with procurement and spending on local projects

- Measure spending against local economic and development indicators

# Extract, Integrate, Map

*Project Reports*



*Procurement Data*

# Read and Process Project Reports

- Geo-tagging

- Link location-specific procurement data

- Link to local Socio-economic indicators

# Trace back Origins of Fact

# Case 2
## Media Promotion Analytics

# Internews - Information changes lives

- An international NGO

- Ensure access to trusted, quality information that empowers people to have a voice in their future

# Linking Internews Projects and Aid Data

# Case 3
## Behavioral Insights and Policy

# How *Data Science* affects Policy

- How people react to different messages

- Behavioural insights combines statistics, economics, and psychology

- Use data science to answer the question:

  *"What works for whom?"*

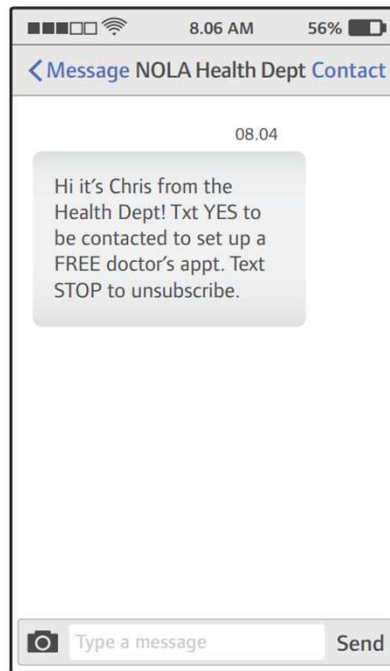- Tailor policy based on demographics/populations

# Predictive Analytics - Randomized Controlled Trials

1. Send different messages to different populations

2. Analyze how they respond

3. Build predictive models

# Behavioural Insights - Healthcare Appointments

People respond differently to different messages

- Simple and generic

- Feeling special

- Promoting social life



Simple — Unique — Prosocial

# Behavioural Insights - Healthcare Appointments

People respond
differently to diffe...
messages

- Simple and g...

- Feeling spec...

- Promoting social life



Hi, It's Chris from the Health Dept! Txt YES to be connected to set up a FREE doctor's appt. Text STOP to unsubscribe
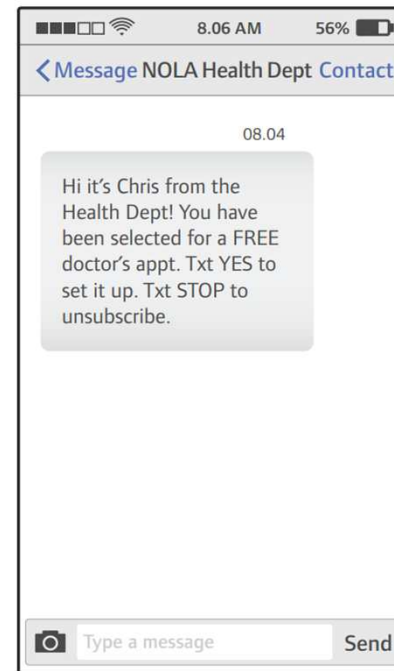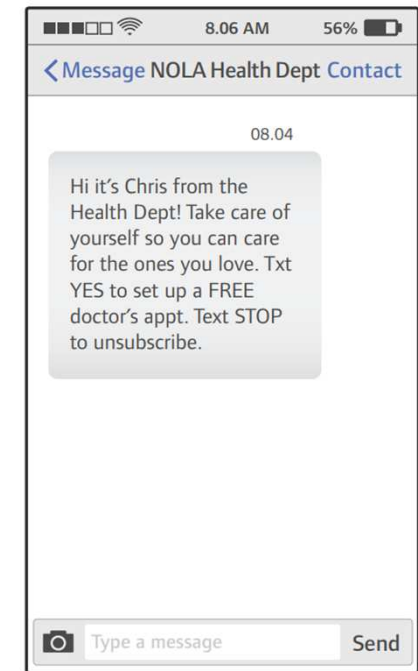
Simple

Unique

Prosocial

# Behavioural Insights - Healthcare Appointments

People respond differently to different messages

- Simple and generic

- Feeling special

- Promoting social life



Hi, It's Chris from the Health Dept! You have been selected for a FREE doctor's appt. Txt YES set it up. Text STOP to unsubscribe

Simple

Unique

Prosocial

# Behavioural Insights - Healthcare Appointments

People respond differently to different messages

- Simple and generic
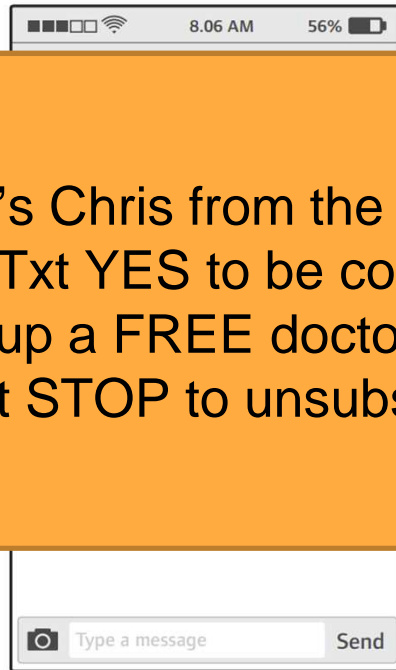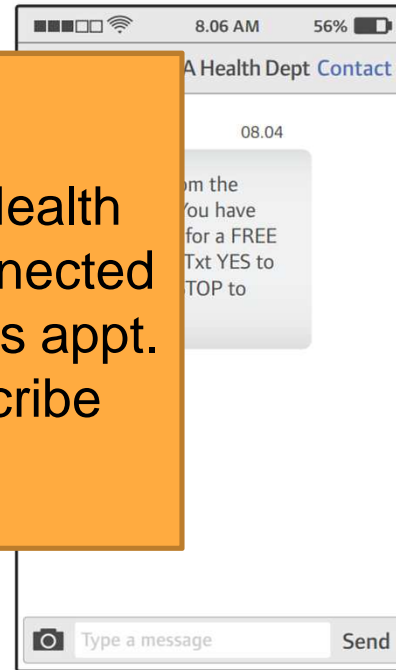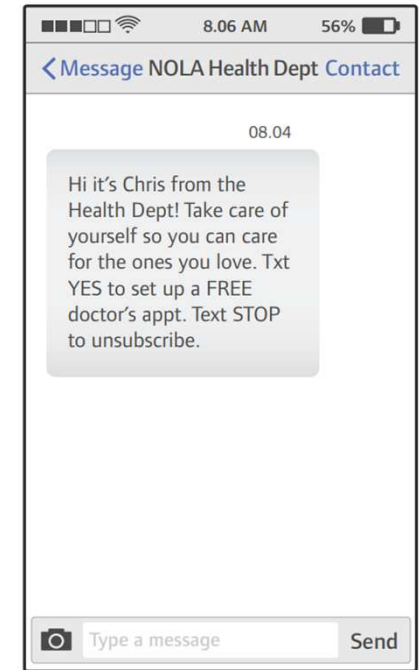
- Feeling special

- Promoting social life



**Simple:**

08.04

Hi it's Chris from the Health Dept! Txt YES to be contacted to set up a FREE doctor's appt. Text STOP to unsubscribe.

**Prosocial:**

Hi, It's Chris from the Health Dept! Take care of yourself so you can care for the ones you love. Txt YES to set up a FREE doctor's appt. Text STOP to unsubscribe

Simple        Unique        Prosocial

# Behavioural Insights – Healthcare Appointments

People respond
differently to diff
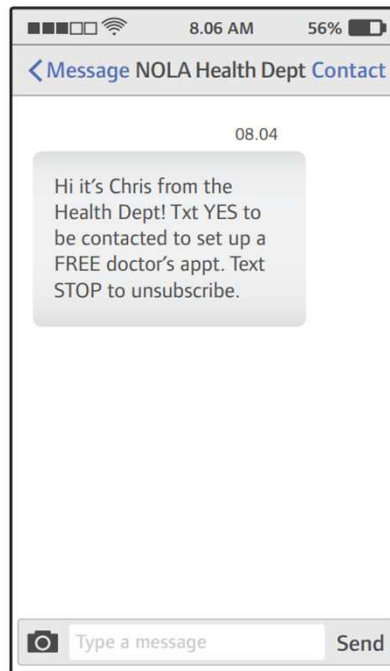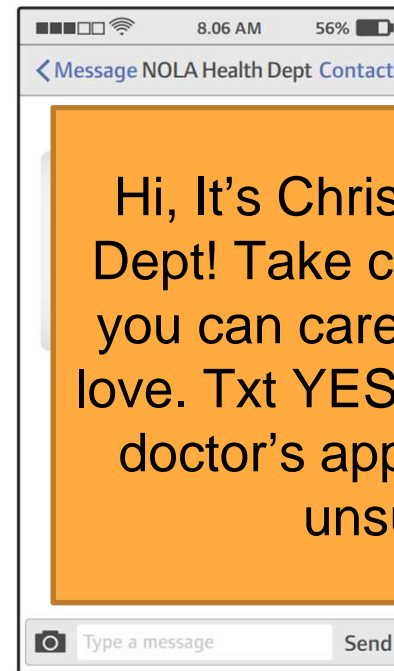messages

● Simple and

● Feeling spe

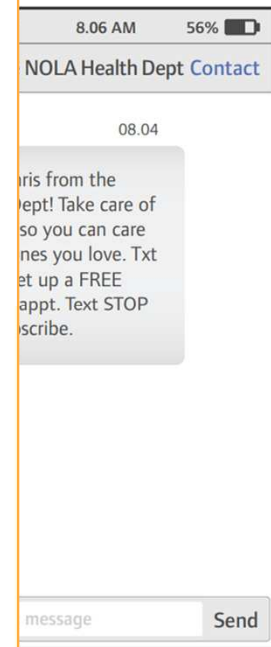● Promoting s



Effect of text messages on preventative healthcare take up

# Data Analytics to Fight Crime

- Predict Criminal Profiles
- Detect Human Trafficking

# Pre-trial Criminal Justice

- Pre-trial has been estimated at over $9 billion per year in the U.S

- Data science can improve how decisions are made during the earliest part of the criminal justice process

# Pre-trial Criminal Justice - Public Safety

- The District Attorney of New Jersey used data analysis
    - Camden, New Jersey, reduced murder by 41% and crime in the city general by 26%

- Use rigorous statistical analysis to <u>classify</u> defendants into
  **low**, **moderate**, and **high**-risk
    - Avoid releasing dangerous people
    - Avoid the cost of keeping low-risk offenders in jail waiting for trials
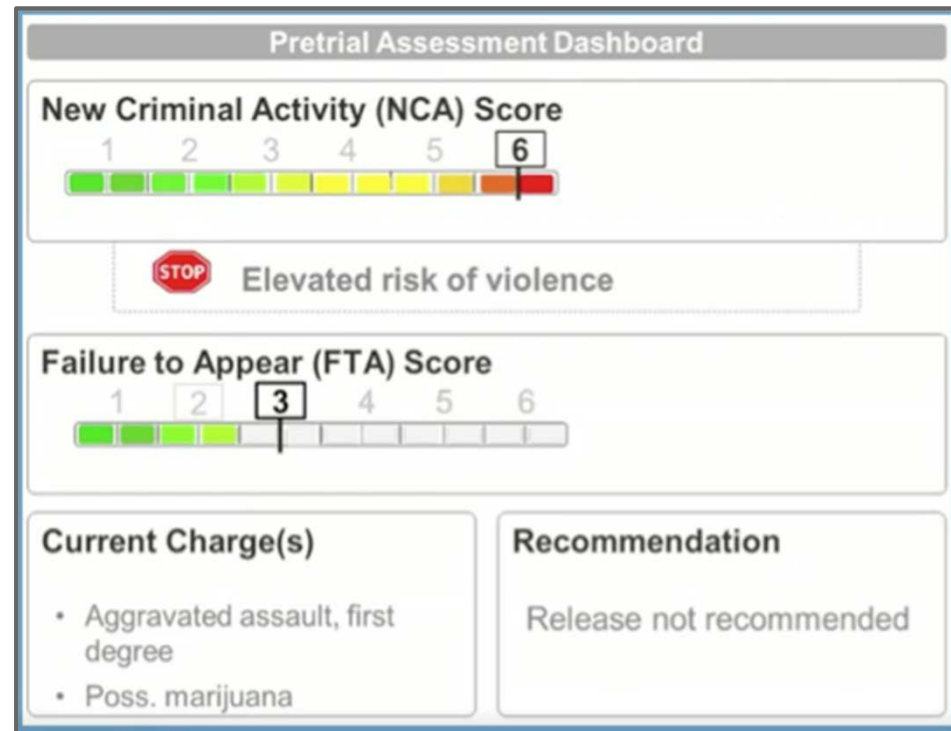
# Predictive Analytics for Criminal Justice

A team of Data Scientist examined

**1.5 Million cases**

Build a **risk-assessment tool**

Predict whether or not,
if someone is released, they will:

1. Commit a new crime.
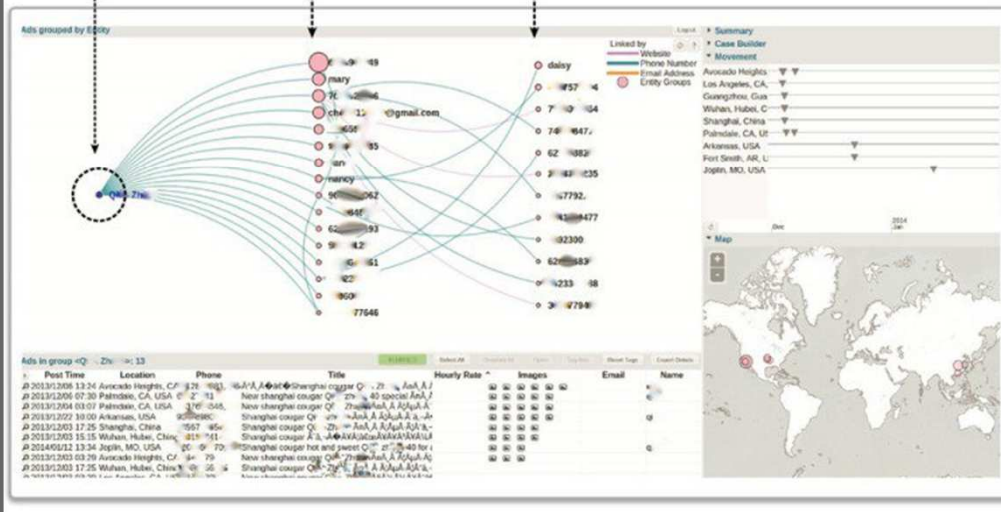2. Commit an act of violence.
3. Come back to court

# Human Trafficking in the Deep Web



*Information Extraction*

*Machine Learning*

*Statistical Inference*

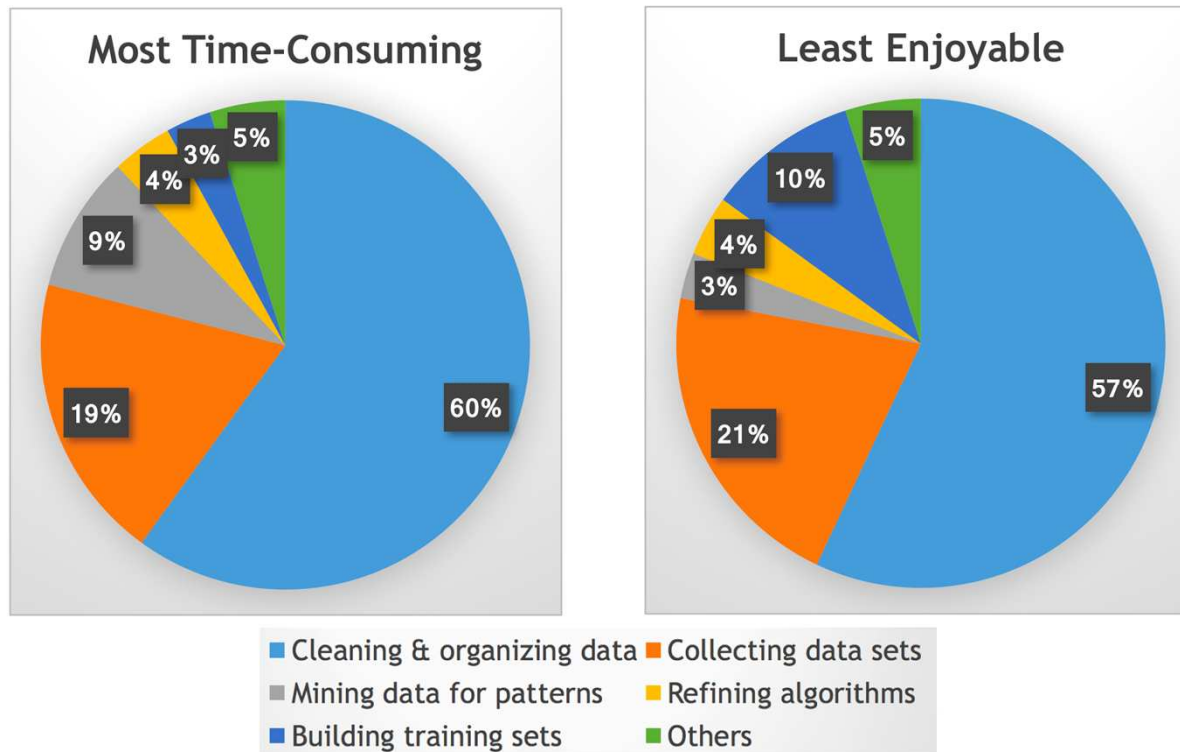# Human Trafficking - MEMEX

# **Major Challenges**
## Data Management and Quality

# Data Curation: Most Time-Consuming, Least Enjoyable



**Most Time-Consuming**
- 60%
- 19%
- 9%
- 4%
- 3%
- 5%

**Least Enjoyable**
- 57%
- 21%
- 3%
- 4%
- 10%
- 5%

Legend:
- Cleaning & organizing data
- Collecting data sets
- Mining data for patterns
- Refining algorithms
- Building training sets
- Others

# Unification: Dedup

**.Record Linkage and Deduplication**



| Part Description | Part Number |
|---|---|
| O-Ring Gasket | 27-00091 |

| Desc. | PN |
|---|---|
| Gasket, O-Rin | 27-00091 |

| Item | Part |
|---|---|
| O-Ring Gasket | 27-00091 |

| Item Descrip | Part # |
|---|---|
| Gasket | 27-00091 |

Record 1

Record 2

Record 3

Record 4

Unified Record

# Cleaning: Missing Values

**.Missing Values (the curse of `Nulls`)**



*Real data is full of of N/A or nulls, special values (99999) etc.*

# Cleaning: Rule Violations

## .Integrity Constraints

| ID | FN | LN | ROLE | CITY | ST | SAL |
|----|------|-------|------|------|----|------|
| 105 | Anne | Nash | M | NYC | NY | 110 |
| 211 | Mark | White | E | SJ | CA | 80 |
| 386 | Mark | Lee | E | NYC | AZ | 75 |
| 235 | John | Smith | M | NYC | NY | 1200 |

**Business Rule**

*Two employees of the same role, the one who lives in NYC cannot earn less than the one who does not live in NYC*

$$\forall t_\alpha, t_\beta \in Emp, \neg(t_\alpha.ROLE = t_\beta.ROLE \wedge t_\alpha.CITY = \text{``}NYC\text{''}$$
$$\wedge t_\beta.CITY \neq \text{``}NYC\text{''} \wedge t_\alpha.SAL < t_\beta.SAL)$$

*Rarely expressed in practice. Most curation tools are rule-based implemented in imperative language*

# And They Don't Come Piece-meal

| ID | Name | ZIP | City | State | Income |
|----|------|-----|------|-------|--------|
| 1 | Green | 60610 | Chicago | IL | 30k |
| 2 | Green | 60611 | Chicago | IL | 32k |
| 3 | Peter | | New Yrk | NY | 40k |
| 4 | John | 11507 | New York | NY | 40k |
| 5 | Gree | 90057 | Los Angeles | CA | 55k |
| 6 | Chuck | 90057 | San Francisco | CA | 30k |

Missing Value

Duplicates

Integrity Constraint Violation

Value/Syntactic Error

# Multiple Efforts for Automation

*8443322821*
*vs 844-332-2821*

*"Jeff Bezos"*
*vs "J. Bezos"*

*VLDB17, Chicago*
*Vs VLDB17, Munich*

Pattern enforcement

Entity Resolution
Schema Integration

Error repairs & imputation







Program synthesis

ML+ Expert sourcing

Scalable Statistical Inference

# Machine Learning is at the Heart of Data Science

- **Engineering and plumbing is ~ 80% of the exercise**

  - In-situ data preparation and signals computation: **Feature engineering**

  - Expert registration and engagement: **Training data management**

  - Blocking and pruning the candidate space: **Scale**

  - Provenance and lineage maintenance: **Explain and rollback**

  - Continuous model monitoring and validation: **Model management**

**Data Repair is a statistical leaning and inference problem**

Thank You!

@ihabilyas