# SDTM Mapping Based on TF-IDF and Neural Network Probabilistic Models
## Sam Tomioka

NJ CDISC Users Group Meeting

March 12, 2019

# Brief Introduction

- Sam Tomioka

- Director, Clinical Data Programming – Data Science

- Sunovion Pharmaceuticals

- Current ML Projects
  - SDTM Mapping
  - Protocol Optimization
  - Digital Endpoints (Seizure)

- Past ML Projects
  - Adverse Events
  - Dose Titration
  - Digital Endpoints (Stroke, Depression)
  - SDTM Mapping

# AGENDA

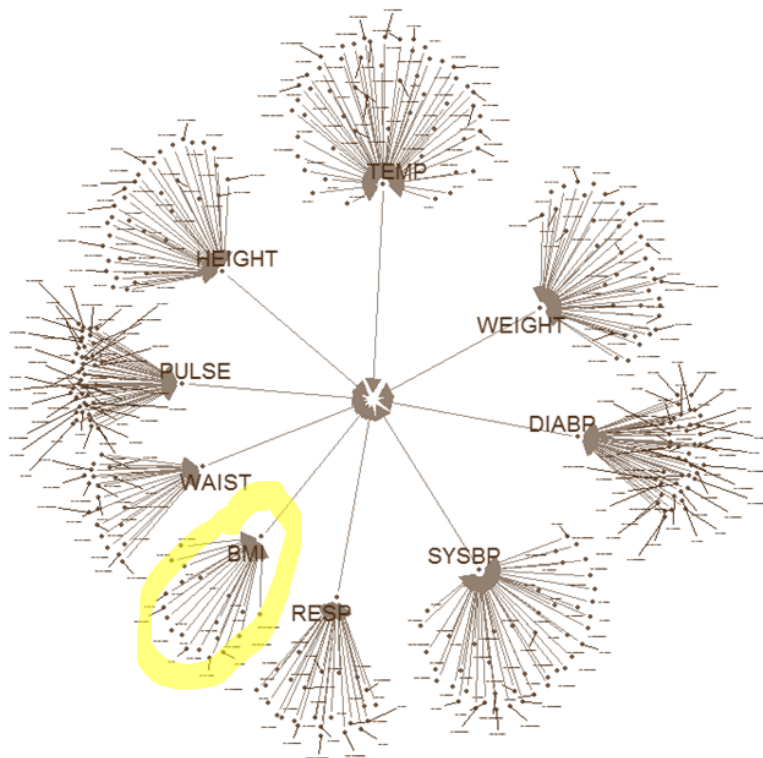Problem to Solve

SDTM Mapping with "Machine Learning"

Tools Used

Thought

# ENDLESS MAPPING…



SDTM.VS.VSORRES mapping for
12 Sunovion studies delivered by one CRO
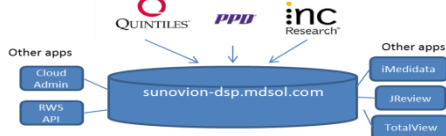
TEMP
HEIGHT
WEIGHT
PULSE
DIABP
WAIST
BMI
SYSBP
RESP

20 sources for BMI

"VS.BMI"
"VS.BMI_RAW"
"VS.BMI_Z"
"VS.BMI_Z_RAW"
"VS.BMIS"
"VS.BMIS_RAW"
"VS1.BMI"
"VS1.BMI_RAW"
"VS1.BMI_Z"
"VS1.BMI_Z_RAW"
"VS1.VS1BMI"
"VS2.BMI"
"VS2.BMI_RAW"
"VS2.D_BMI"
"VS2.D_BMI_RAW"
"VS2.VS2BMI"
"VSMSTR.BMI"
"VSMSTR.BMI_RAW"
"VSMSTR.D_BMI"
"VSMSTR.D_BMI_RAW"

Data Source: ~/My Documents/r/programming_challenges_201711xx/inventory20171114_t.xlsx

# Embrace Metadata Standards

# WHAT NEXT?

Can I use **natural language** model and **machine learning** algorithms to map raw data variables to SDTM variables?

CDISC

# AGENDA

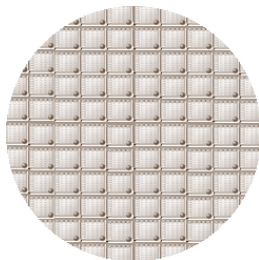Problem to Solve

SDTM Mapping with "Machine Learning"

Tools Used

Thought

# ML based SDTM mapping for
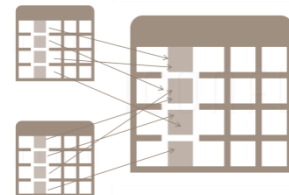# <u>fast</u>, <u>accurate</u>, <u>consistent</u> SDTM generation

**Input**

**NLP & ML**

**Output**

Raw data



Mapping Specification

Wordcloud for AE raw metadata in raw







CDASHIG, SDTM IG



CDISC

# Steps

# Create mapping specifications (with Human Intelligence)

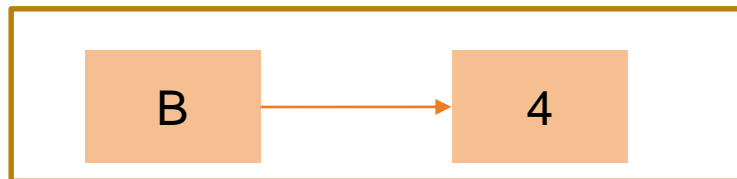| Raw Variable | SDTM Variable |
|---|---|
| PT | AEDECOD |
| SOC | AEBODSYS |
| PTNAME | AEDECOD |
| SOCNAME | AEBODSYS |
| … | … |
| | |

*illustration purpose only

CDISC

# Natural Language Models

# Natural Language Model A
## TF-IDF algorithm: Weighing terms

- Words occur nearby frequently are important than words that only appear once or twice

  Frequency (TF) $\quad \text{tf}_{t,d} = \begin{cases} 1 + \log_{10} \text{count}(t,d) & \text{if } \text{count}(t,d) > 0 \\ 0 & \text{otherwise} \end{cases}$

- Words that are too frequent are not important

  Inverse Document Frequency (DF) $\quad \text{idf}_t = \log_{10} \left( \dfrac{N}{\text{df}_t} \right)$

- Weight $\quad w_{t,d} = \text{tf}_{t,d} \times \text{idf}_t$

**Logistic regression, accuracy=0.64058**

AETERM definition in CDASH IG

The reported or pre-specified name of the adverse event.

0.6331503    0.6105753                    0.6247964    0.6611213

$\longrightarrow$    LR    $\longrightarrow$    $\hat{y}$

Document term matrix of 17 studies data and IG
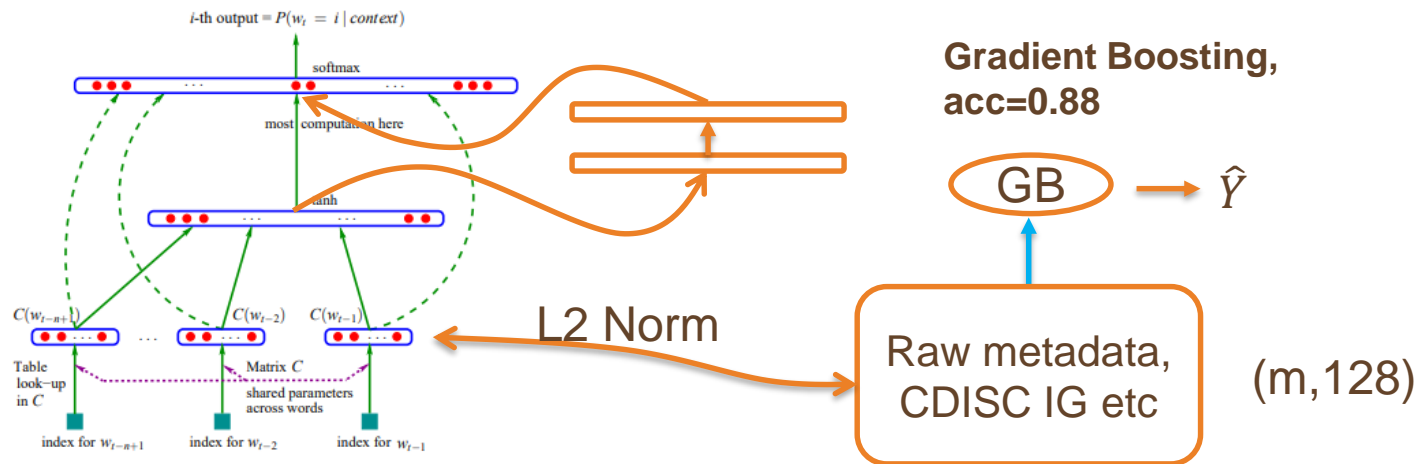
# Natural Language Model B



Figure 1: Neural architecture: $f(i, w_{t-1}, \cdots, w_{t-n+1}) = g(i, C(w_{t-1}), \cdots, C(w_{t-n+1}))$ where $g$ is the neural network and $C(i)$ is the $i$-th word feature vector.

**Gradient Boosting, acc=0.88**

L2 Norm

Raw metadata, CDISC IG etc

(m,128)

200,000,000,000 English Google News corpus

Figure from Yoshua Bengio, Réjean Ducharme, Pascal Vincent, Christian Jauvin. A Neural Probabilistic Language Model. Journal of Machine Learning Research, 3:1137-1155, 2003.

# Natural Language Models

**Simple**

TFIDF

**1**
**0.769**

**High dimension**

**2**
**0.827**

CDISC
key words + Simple

**3**
**0.965**
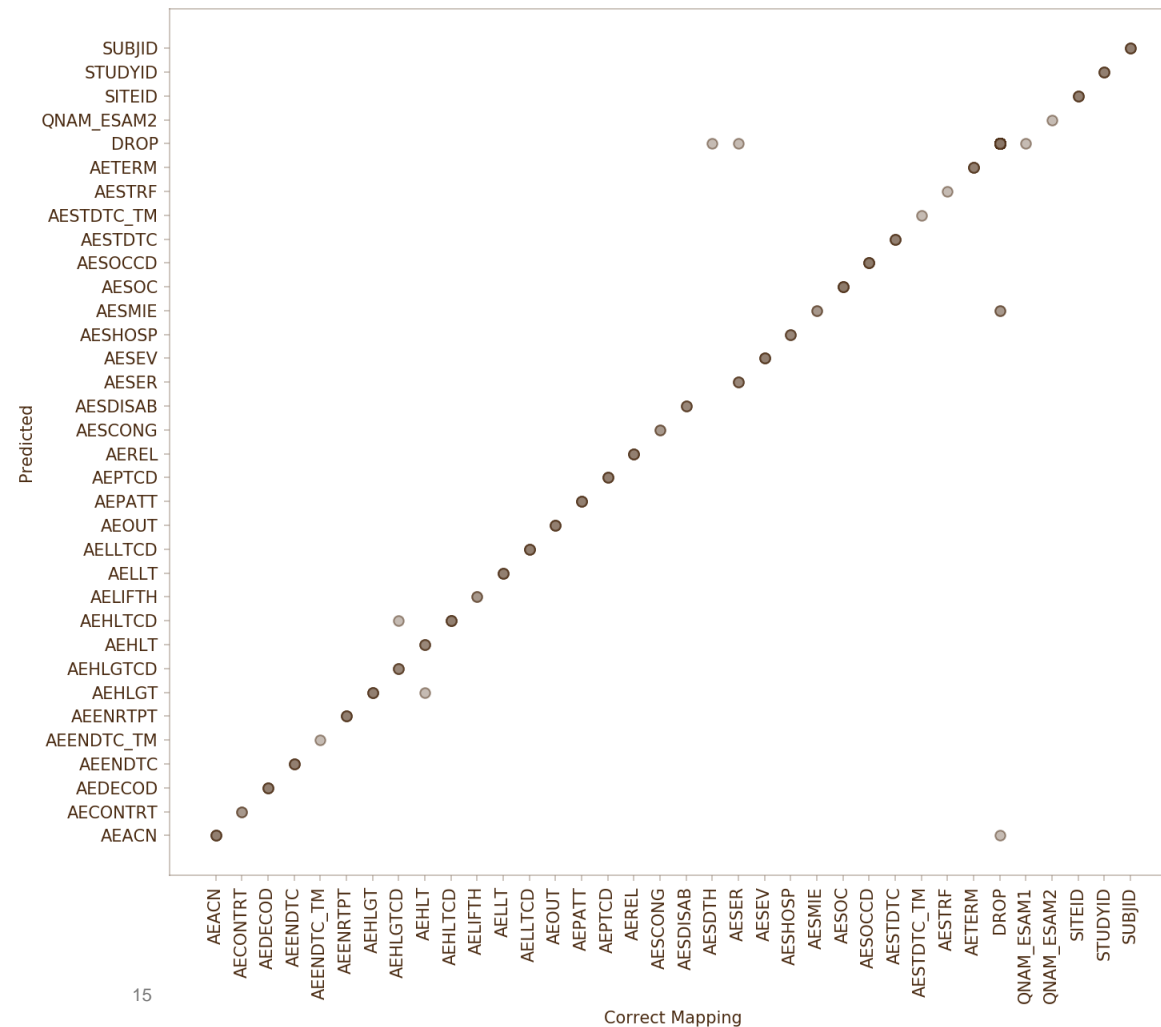
NNLM

**4**
**0.88**

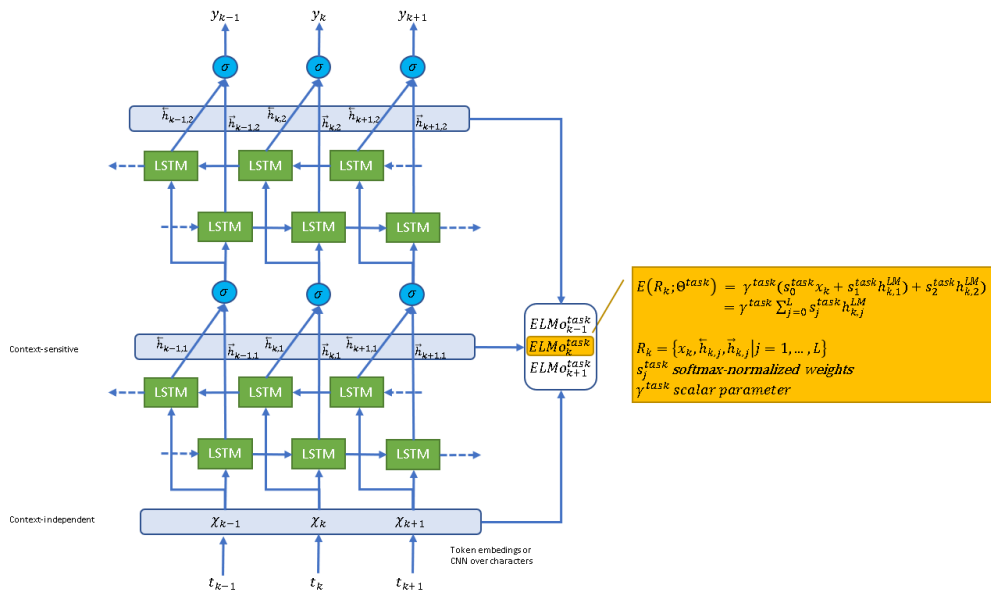## Mapping Results



Figure 6. Results of SDTM mapping (Accuracy with 95% CI and Kappa)

Final Ensemble Model

Mapping Accuracy on data from 3 new studies using Final Ensemble Model

0.97

# More robust approach

# AGENDA

Problem to Solve

SDTM Mapping with "Machine Learning"

Tools Used

Thought

# Tools used for POC

| Programming | IDE | ML Framework | Purpose |
|---|---|---|---|
|  |  | NA | Metadata extraction from sas7bdat |
|  |  | caret | NLP ML Visualizations |
|  |  |  | Transfer learning ML Visualizations |

# AGENDA

Problem to Solve

SDTM Mapping with "Machine Learning"

Tools Used

Thought

## Summary

- This proof of concept demonstrated that **machine learning** along with a **natural language model** can produce a pretty accurate SDTM mapping specification document.

- As in any ML models, as you feed more mapping specs, the model will learn them and become more robust.

CDISC

# Thank You

www.linkedin.com/in/STomioka

Sam.Tomioka @ Sunovion.com

CDISC