# Statistical Inference for NLP algorithms

*Kishore Basu & TJ Ayoub*

# What are the authors trying to solve?

Computer Science/Linguistics

Statistical Approaches

◎ In NLP, statistical techniques used infrequently

◎ Uncertainty quantification is difficult

◎ word2vec has strong statistical underpinnings

"History ⟨Diabetes⟩ high blood pressure"

| $C = c$ (Center Word) | $C' = c'$ (Neighbouring Word) | $D = 0/1$ (Boolean neighbour) |
|---|---|---|
| diabetes | history | $D = 1$ |
| diabetes | diabetes | $D = 1$ |
| diabetes | high | $D = 1$ |
| diabetes | blood | $D = 1$ |
| diabetes | pressure | $D = 0$ |
| pressure | history | $D = 0$ |

$PMI > 0$
$- C, C'$ associated

$PMI < 0$
$-C, C'$ not associated

$$PMI = log \frac{P(C = c \,|C' = c', D = 1)}{P(C = c \,|D = 1)}$$

3

*Can statistical techniques improve our understanding of diabetes classification?*

# Yes, but we need to do some work first.

◎ PMI model

◎ MLE estimation

◎ Multivariate delta method

$$PMI = log \frac{P_1(C = c \mid C' = c', D = 1)}{P_2(C = c \mid D = 1)}$$

MLE       MVDM

$$P_1(C = c \mid C' = c', D = 1) = \frac{\exp(\beta^T X)}{1 + \exp(\beta^T X)} \quad \rightarrow \quad \hat{\beta} \quad \rightarrow \quad \widehat{\sigma_1}^2 = Var(\widehat{P_1})$$

DM

$$\widehat{PMI}$$

$$P_2(C = c \mid D = 1) = \frac{\exp(\alpha)}{1 + \exp(\alpha)} \quad \rightarrow \quad \hat{\alpha} \quad \rightarrow \quad \widehat{\sigma_2}^2 = Var(\widehat{P_2})$$

# Building a predictive model from data

◎ Goal: Identify type-2 diabetes in EHR

◎ Two summary statistics for two groups

◎ Word association pairs, words with the 'diabet' stem

$$(diabet-, c')$$

$$m_k PMI_d^d$$

$$m_k PMI_d^{xd}$$

$$(c, c')$$

$$m_k PMI^d$$

$$m_k PMI^{xd}$$

# RESULTS

# 1000 patients
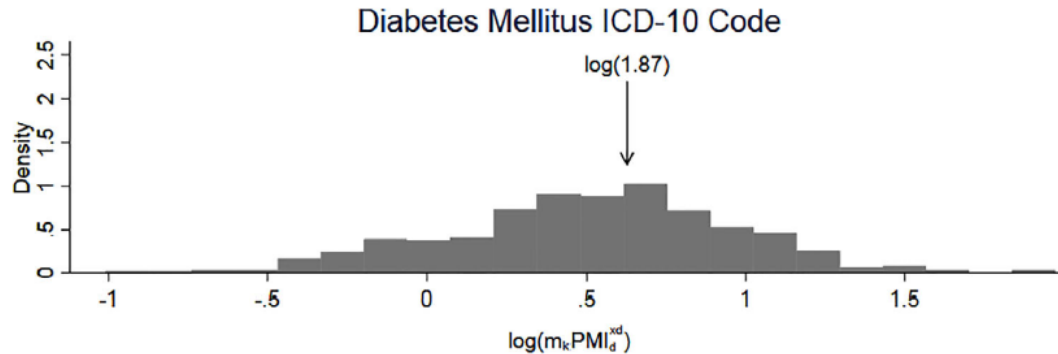
500 Diabetes, 500 non-diabetes

# 61,489 total words

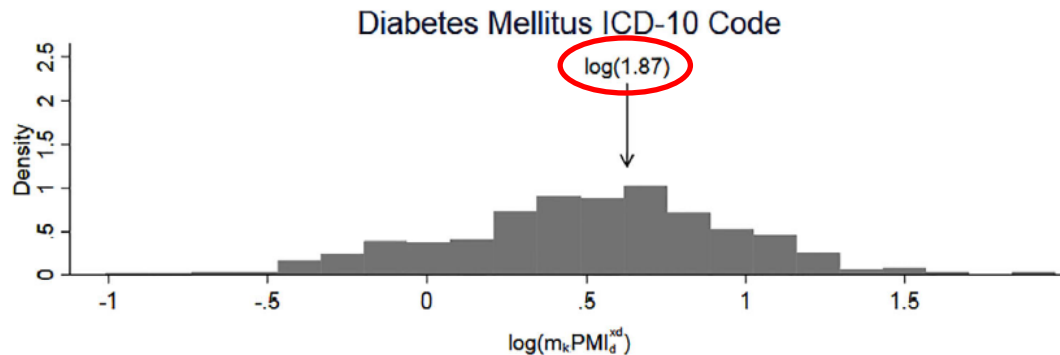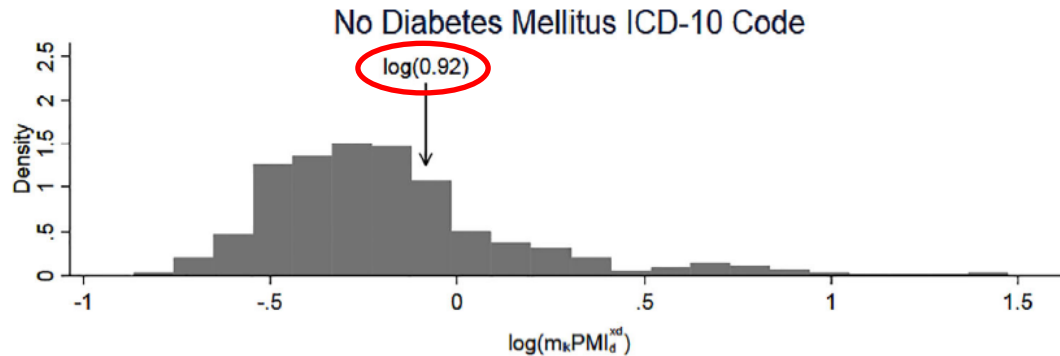Each patient has many records (notes)

# 1500 words kept

Reduces dimensionality of dataset

# PREDICTIVE MODEL



**Mean PMI over non-diabetics where $C' = \text{`}diabet'$**

# PREDICTIVE MODEL



No Diabetes Mellitus ICD-10 Code

log(0.92)

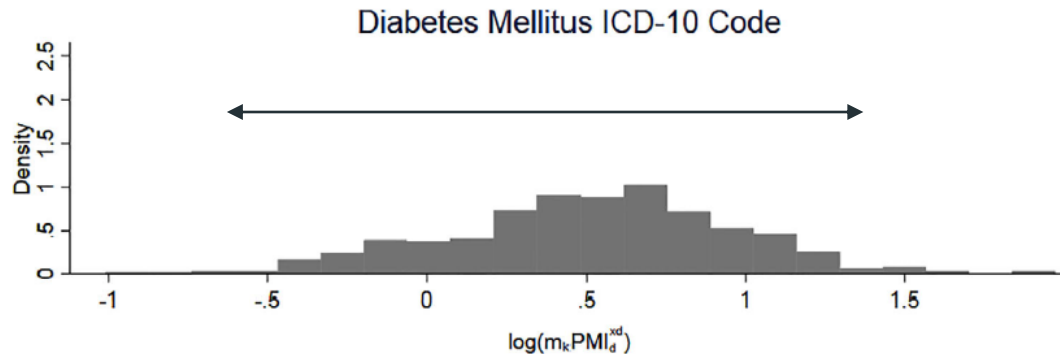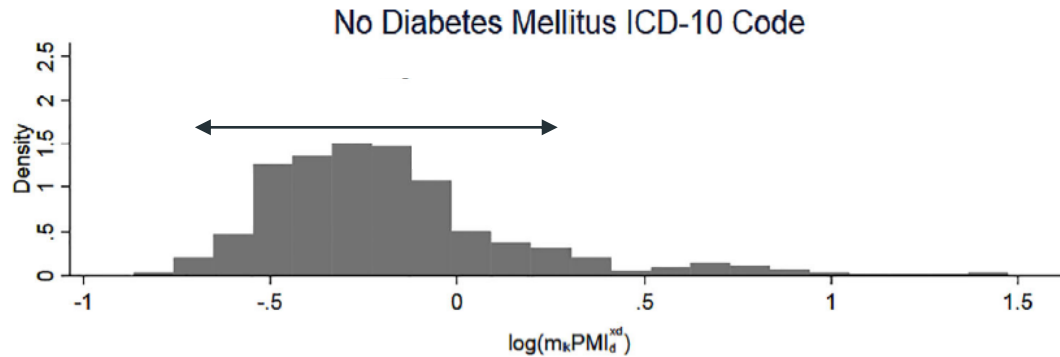Diabetes Mellitus ICD-10 Code

log(1.87)

**Mean PMI over non-diabetics where $C' = `diabet'$**

◎ **Key Takeaway:**

Less words associated with 'diabet' in non-diabetics

More words associated with 'diabet' in diabetics

# PREDICTIVE MODEL



No Diabetes Mellitus ICD-10 Code

Diabetes Mellitus ICD-10 Code

**Mean PMI over non-diabetics where $C' =$ `$diabetes'$**

◎ **Key Takeaway:**

Non-diabetics have tighter distribution

PMI's in non-diabetes group are better at predicting who doesn't have diabetes

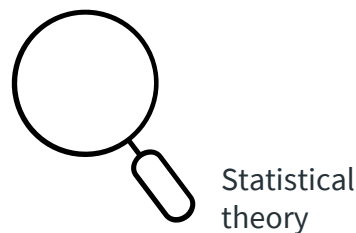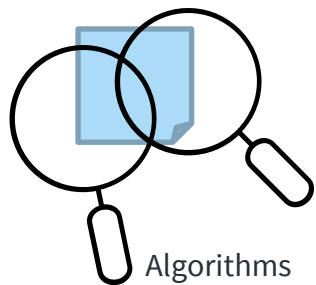So understanding the distribution of the PMI can tell us a lot!

# In summary,

◎ Novel framework to calculate SEs for PMI

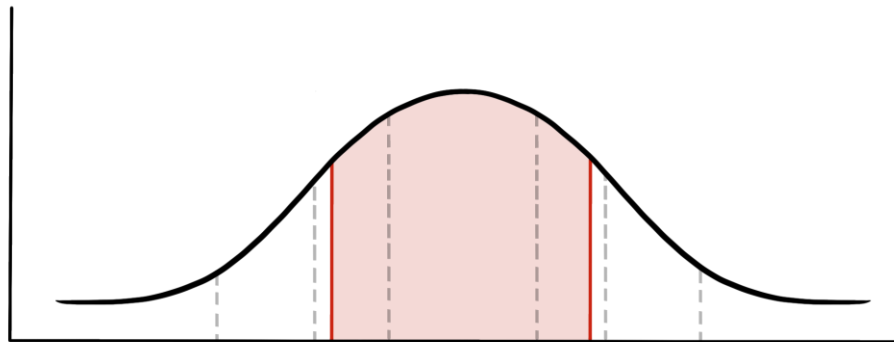$$PMI = log\frac{P(C = c \mid C' = c', D = 1)}{P(C = c \mid D = 1)} \qquad \rightarrow \qquad \begin{array}{c} \widehat{PMI} \\ Var(\widehat{PMI}) \end{array}$$

◎ Importance of statistical analysis in NLP and data science



Algorithms

Computer Science

Statistical theory

# In summary,

◎ High relevance contribution in a sparse literature space

# Thank you for listening