# Insights into Analogy Completion from the Biomedical Domain

**Denis Newman-Griffis**, Albert Lai, Eric Fosler-Lussier

The Ohio State University    National Institutes of Health, Clinical Center
Washington University in St. Louis

*BioNLP 2017*
August 4, 2017

THE OHIO STATE UNIVERSITY

NIH National Institutes of Health
Clinical Center

Washington University in St. Louis

Increasing work on training embedding-based models for biomedical applications, but not many resources to evaluate on.

# TL;DR

Increasing work on training embedding-based models for biomedical applications, but not many resources to evaluate on.

Analogies have been highly useful in the general domain, so we built an analogy dataset for BioNLP. [1]

[1]https://github.com/OSU-slatelab/BMASS

# TL;DR

Increasing work on training embedding-based models for biomedical applications, but not many resources to evaluate on.

Analogies have been highly useful in the general domain, so we built an analogy dataset for BioNLP. [1]

Findings:

- ▶ Current embeddings are good at direct chemical/biological relationships, not so good at clinical semantics.
- ▶ Changes need to be made to the standard analogy methods to reflect the complexity of real data.
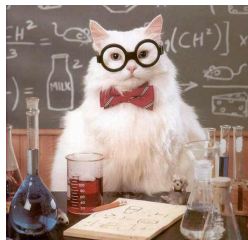
The analogy completion task

BMASS

How can we make analogies more realistic?

Findings and challenges on our dataset

London : England :: Paris : France

London : England :: Paris : _____

London : England :: Paris : _____

$$d^* = \text{argmax}_{d \in V}\big(\cos(d, Eng - Lndn + Paris)\big)$$

London : England :: Paris : _____

$$d^* = \text{argmax}_{d \in V}\big(\cos(d, Eng - Lndn + Paris)\big)$$

England

??

London → Paris

**France**
Switzerland
Italy
urology
swimming
purple
Latvia

✔

London : England :: Paris : _____

**Italy**
Switzerland
France
urology
swimming
purple
Latvia

$$d^* = \text{argmax}_{d \in V}\left(\cos(d, Eng - Lndn + Paris)\right)$$



✗

Unified Medical Language System (UMLS)

# BMASS - BioMedical Analogic Similarity Set

Unified Medical Language System (UMLS)

Normalized concepts

$$C0009443 \longrightarrow \left\{ \begin{array}{c} \text{common cold} \\ \text{cold} \\ \text{acute rhinitis} \end{array} \right\}$$

BMASS - BioMedical Analogic Similarity Set

Unified Medical Language System (UMLS)

Normalized concepts

Relation triples

C0009443 →
{
common cold
cold
acute rhinitis
}

subject
relation
object
⟨
C0450297
RO:has-finding-site
C0341532
⟩

# BMASS - BioMedical Analogic Similarity Set

| ID | Name | Amb |
|----|------|-----|
| *Lab/Rx* | | |
| L1 | form-of | 1.0 |
| L2 | has-lab-number | 1.1 |
| L3 | has-tradename | 1.5 |
| L4 | tradename-of | 1.3 |
| L5 | associated-substance | 1.6 |
| L6 | has-free-acid-or-base-form | 1.0 |
| L7 | has-salt-form | 1.1 |
| L8 | measured-component-of | 1.3 |
| *Hierarchical* | | |
| H1 | refers-to | 1.0 |
| H2 | same-type | 10.4 |
| *Morphological* | | |
| M1 | adjectival-form-of | 1.1 |
| M2 | noun-form-of | 1.0 |

| ID | Name | Amb |
|----|------|-----|
| *Clinical* | | |
| C1 | associated-with-malfunction-of-gene-product | 2.6 |
| C2 | gene-product-malfunction-associated-with-disease | 1.5 |
| C3 | causative-agent-of | 4.6 |
| C4 | has-causative-agent | 2.0 |
| C5 | has-finding-site | 1.9 |
| C6 | associated-with | 1.2 |
| *Anatomy* | | |
| A1 | anatomic-structure-is-part-of | 1.6 |
| A2 | anatomic-structure-has-part | 5.4 |
| A3 | is-located-in | 1.4 |
| *Biology* | | |
| B1 | regulated-by | 1.0 |
| B2 | regulates | 1.0 |
| B3 | gene-encodes-product | 1.1 |
| B4 | gene-product-encoded-by | 2.4 |

Cross-product of 50 samples for each relation:

2,450 analogies for each relation $\Rightarrow$ **61,250 total analogies**

This dataset represents real biomedical relationships...

This dataset represents real biomedical relationships. . .

But it doesn't fit the standard paradigm!

3 key assumptions in evaluation methodology:

- Single Answer
- Same Relationship(s)
- Informativity

Each is violated in recent analogy datasets

- Google[2], BATS[3], Sem-Para[4]

All are problematic in real-world data!

---

[2]Mikolov et al. 2013
[3]Gladkova et al. 2016
[4]Köper et al. 2015

# "Unassuming" the standard assumptions

**Single Answer**     Same Relationship     Informativity

The given analogy has only one correct target.

- ▶ Enforced by argmax over candidates for completing the analogy.
- ▶ If multiple analogies, must get at least one wrong.

—— *Problem cases* ————————————————————————

flu : nausea :: fever :     $\begin{cases} \text{sweating} \\ \text{weakness} \end{cases}$

## Single Answer    Same Relationship    Informativity

The given analogy has only one correct target.

- ▶ Enforced by argmax over candidates for completing the analogy.
- ▶ If multiple analogies, must get at least one wrong.

——— *Problem cases* ————————————————

$$\text{flu : nausea :: fever :} \quad \left\{ \begin{array}{l} \text{sweating} \\ \text{weakness} \end{array} \right\}$$

## Easy fix!

Allow for multiple correct answers; also report on all of them, for fuller picture.

# "Unassuming" the standard assumptions

All information relating exemplars $a$ and $b$ also relates query $c$ to $d$.

- ▶ Enforced by treating full vector difference as relation of interest.
- ▶ Many relations have partial overlap with one another.

——— *Problem cases* ——————————————————————————

# "Unassuming" the standard assumptions

All information relating exemplars $a$ and $b$ also relates query $c$ to $d$.

- ▶ Enforced by treating full vector difference as relation of interest.
- ▶ Many relations have partial overlap with one another.

—— *Problem cases* ——

<div align="center">

brother : sister :: husband : wife

**MaleCounterpart**    **MaleCounterpart**
SiblingOf      MarriedTo

</div>

# "Unassuming" the standard assumptions

All information relating exemplars *a* and *b* also relates query *c* to *d*.

- ▶ Enforced by treating full vector difference as relation of interest.
- ▶ Many relations have partial overlap with one another.

—— *Problem cases* ——————————————

brother : sister :: husband : wife

**MaleCounterpart**    **MaleCounterpart**
SiblingOf        MarriedTo

wrist pain : shoulder pain :: bipolar disorder : trazodone

**CoOccursWith**       **CoOccursWith**
TreatmentFor

# "Unassuming" the standard assumptions

The relationship between exemplars $a$ and $b$ is specific enough to suggest the correct target $d$ for query $c$.

- ▶ Issue with very broad semantic or hierarchical relationships.

——— *Problem cases* ——————————————————————

Single Answer          Same Relationship          **Informativity**

The relationship between exemplars $a$ and $b$ is specific enough to suggest the correct target $d$ for query $c$.

▶ Issue with very broad semantic or hierarchical relationships.

——— *Problem cases* ———

Generally related *(UMLS)*

socks : stockings :: Finns : Finnish language

# "Unassuming" the standard assumptions

The relationship between exemplars $a$ and $b$ is specific enough to suggest the correct target $d$ for query $c$.

▶ Issue with very broad semantic or hierarchical relationships.

———— *Problem cases* ————————————————

Generally related *(UMLS)*

socks : stockings :: Finns : Finnish language

## Fix during dataset generation

Review samples from each relation to ensure they're properly determined.

$$a : b :: c : \underline{\quad} \quad \longrightarrow \quad d^* = \text{argmax}_{d \in V}\big(\cos(d, b - a + c)\big)$$

Evaluating analogies under 3 settings:

$$a : b :: c : \underline{\quad} \quad \longrightarrow \quad d^* = \text{argmax}_{d \in V}\big(\cos(d, b - a + c)\big)$$

Evaluating analogies under 3 settings:

**Single**-**Answer (SA)** Single candidate target for each analogy selected as the only "correct" answer.

$$\text{flu} : \text{nausea} :: \text{fever} : \left\{ \begin{array}{c} \textbf{sweating} \\ \textbf{\sout{weakness}} \end{array} \right\} \quad \longrightarrow \quad \begin{array}{c} \sout{\text{weakness}} \\ \text{temperature} \\ \textbf{sweating} \\ \text{nodule} \end{array} \quad \textcolor{red}{\textbf{✗}}$$

$$a : b :: c : \underline{\quad} \quad \longrightarrow \quad d^* = \text{argmax}_{d \in V}\big(\cos(d, b - a + c)\big)$$

Evaluating analogies under 3 settings:

**Single-Answer (SA)** Single candidate target for each analogy selected as the only "correct" answer.

**Multi-Answer (MA)** All candidate targets for each analogy are considered to be correct.

flu : nausea :: fever : $\left\{ \begin{array}{c} \text{sweating} \\ \text{weakness} \end{array} \right\} \longrightarrow$ 

**weakness**
temperature
**sweating**
nodule

✔

$$a : b :: c : \underline{\quad} \quad \longrightarrow \quad d^* = \text{argmax}_{d \in V} \big( \cos(d, b - a + c) \big)$$

Evaluating analogies under 3 settings:

**Single-Answer (SA)** Single candidate target for each analogy selected as the only "correct" answer.

**Multi-Answer (MA)** All candidate targets for each analogy are considered to be correct.

**All-Info (AI)** Use all possible exemplar objects and all candidate targets.

flu : $\left\{ \begin{array}{c} \text{nausea} \\ \text{cough} \end{array} \right\}$ :: fever : $\left\{ \begin{array}{c} \text{sweating} \\ \text{weakness} \end{array} \right\}$ $\longrightarrow$ $b - a =$ $\frac{1}{2} \left( \begin{array}{c} \text{nausea - flu} \\ \text{+ cough - flu} \end{array} \right)$

Reporting 3 metrics over ranked candidates:

**Acc$_R$** Relaxed accuracy; correct if any valid answer is the top choice
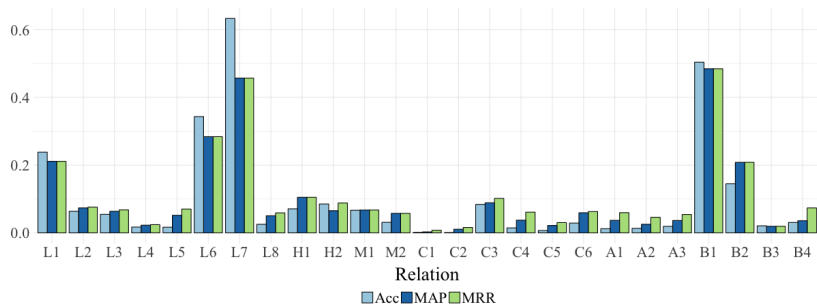
**MAP** Mean average precision

**MRR** Mean reciprocal rank

**weakness**
temperature
sweating
nodule
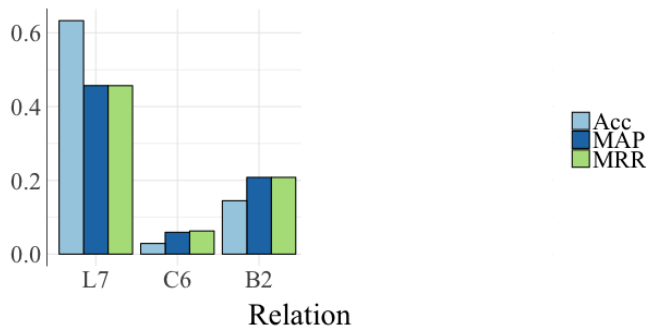
**Acc$_R$** 1.0
**MAP** $\frac{5}{6}$
**MRR** 1.0

NIH⟩

# Overall results



- Results shown for Multi-Answer setting.
- Average performance is around 11% on all metrics, with all embeddings. High variability between relationships!
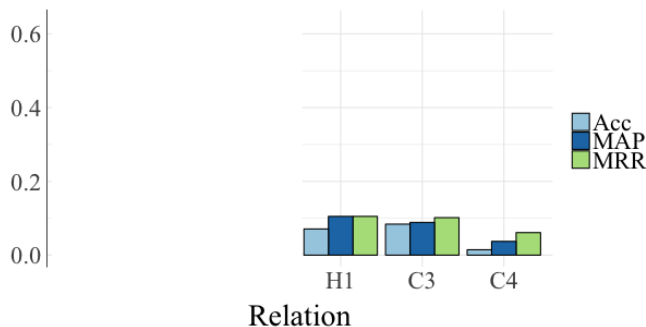- Used 5 different sets of embeddings trained on PubMed.

# MAP/MRR give a better picture



- MAP $<$ Acc$_R$ indicates wider distribution of correct answers on *L7* (has-salt-form)
- MAP $>$ Acc$_R$ shows that even if top answer is wrong, correct answers aren't far down on *C6* (associated-with), *B2* (regulates)
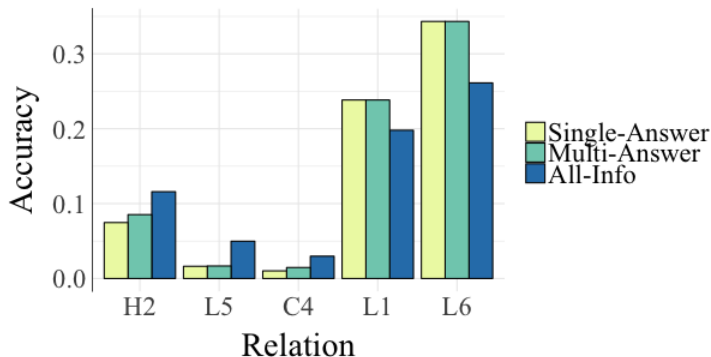
# MAP/MRR give a better picture



- MRR $>$ Acc$_R$ shows that the best correct answer stays near the top on *C4* (has-causative-agent)
- MRR $\approx$ Acc$_R$ reflects more consistent positioning of nearest correct answer on *H1* (refers-to), *C3* (causative-agent-of)

# All-Info benefits vary



- ▶ Extra examples help on *H2* (same-type), *L5* (associated-substance), and *C4* (has-causative-agent).
- ▶ But harm *L1* (form-of) (4% absolute) and *L6* (has-free-acid-or-base-form) (8% absolute)

Single-Answer and Informativity assumptions addressed, but not Same Relationship(s).

Single-Answer and Informativity assumptions addressed, but not Same Relationship(s).

- ▶ Drozd et al (2016) use a parametric logistic regression that can be used to learn affine subspaces.

Single-Answer and Informativity assumptions addressed, but not Same Relationship(s).

- ▶ Drozd et al (2016) use a parametric logistic regression that can be used to learn affine subspaces.

Standard linear offset method does not work for real-world data!

- ▶ Our changes help, but overall performance is still low (as with other recent datasets). Use MAP and MRR!

- ▶ **Analogies are useful!** We need to find better ways to tackle this task.

# Thank you!

Dataset and source code at:
https://www.github.com/OSU-slatelab/BMASS

newman-griffis.1@osu.edu