**UCLA** Computational Medicine

# KNN/Patient Like Me Update
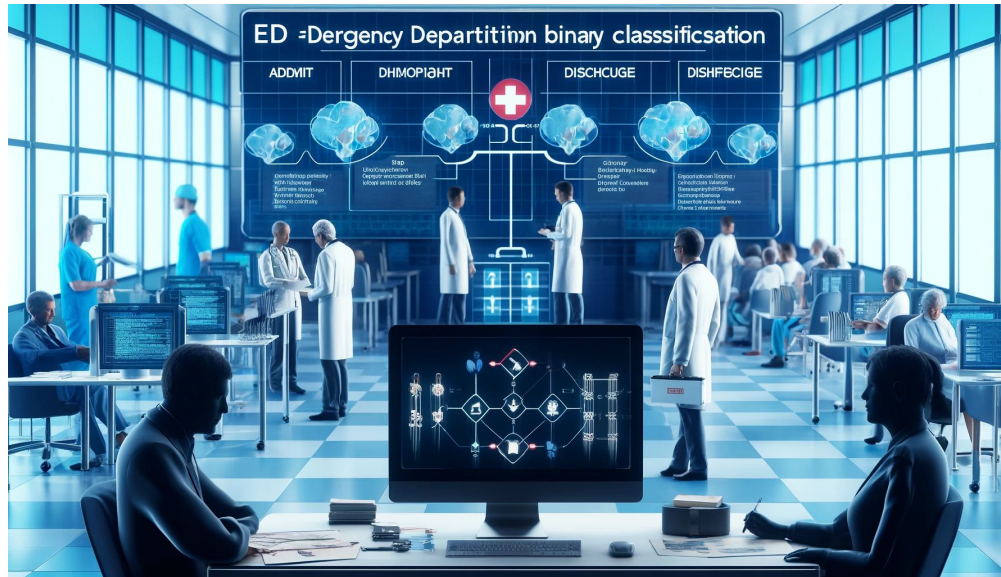
Simon Lee
July 23, 2024

# Motivation

We developed Pseudo-notes which is an interpretable readout of tabular EHR datasets. It also creates an interface to interact with pre-trained Language Models and Large Language Models

Pseudo-notes also addresses issues regarding the data covering several biological and time scales which is hard to account for and synergize.

**UCLA** Computational Medicine

# Prediction Task (ED Disposition)

Binary Classification Task to determine whether a patient should be admitted to the Emergency Department or not. (1- ADMIT 0 - HOME)



UCLA Computational Medicine

# My Update

**Designed 2 KNN Algorithms**

1. **Bootstrapping KNN algorithm:** This provides a memory efficient solution and potentially robust model for generating prediction if we cannot fit all embeddings in the KNN

---

**Algorithm 1** Bootstrapping KNN

1: **function** KNN_BOOTSTRAP($query\_text, data\_source, model\_name, n\_neighbors = 5, n\_bootstrap = 1000, sample\_size = 1000$)
2:    $data\_frame \leftarrow$ DataFrame($data\_source$)
3:    $all\_embeddings \leftarrow$ Encode all patient information texts to vectors using model_name
4:    $query\_embedding \leftarrow$ Encode the query text to a vector using model_name
5:    $bootstrap\_predictions \leftarrow$ initialize an empty list for storing predictions
6:    **for** $i \leftarrow 1$ **to** $n\_bootstrap$ **do**
7:       $sample\_indices \leftarrow$ randomly select indices from data_frame with replacement
8:       $sample\_data\_frame \leftarrow$ create a data frame from selected indices
9:       $sample\_embeddings \leftarrow$ get embeddings for the sampled data frame
10:      $knn \leftarrow$ initialize a KNN model with cosine distance metric
11:      $knn.fit(sample\_embeddings)$
12:      $distances, indices \leftarrow knn.kneighbors(query\_embedding)$
13:      $results \leftarrow$ get the rows from sample_data_frame corresponding to indices
14:      $prediction \leftarrow$ apply majority voting on the 'discharge' field of results
15:      $bootstrap\_predictions.append(prediction)$
16:   **end for**
17:   $final\_prediction \leftarrow$ determine the most common prediction from bootstrap_predictions
18:   $confidence \leftarrow$ calculate confidence of the final_prediction based on frequency
19:   **return** $final\_prediction, confidence$
20: **end function**

# My Update

**Designed 2 KNN Algorithms**

1. **Bootstrapping KNN algorithm:** This provides a memory efficient solution and potentially robust model for generating prediction if we cannot fit all embeddings in the KNN

```
Algorithm 1 Bootstrapping KNN
 1: function KNN_BOOTSTRAP(query_text, data_source, model_name, n_neighbors =
    5, n_bootstrap = 1000, sample_size = 1000)
 2:     data_frame ← DataFrame(data_source)
 3:     all_embeddings ← Encode all patient information texts to vectors using model_name
 4:     query_embedding ← Encode the query text to a vector using model_name
 5:     bootstrap_predictions ← initialize an empty list for storing predictions
 6:     for i ← 1 to n_bootstrap do
 7:         sample_indices ← randomly select indices from data_frame with replacement
 8:         sample_data_frame ← create a data frame from selected indices
 9:         sample_embeddings ← get embeddings for the sampled data frame
10:         knn ← initialize a KNN model with cosine distance metric
11:         knn.fit(sample_embeddings)
12:         distances, indices ← knn.kneighbors(query_embedding)
13:         results ← get the rows from sample_data_frame corresponding to indices
14:         prediction ← apply majority voting on the 'discharge' field of results
15:         bootstrap_predictions.append(prediction)
16:     end for
17:     final_prediction ← determine the most common prediction from bootstrap_predictions
18:     confidence ← calculate confidence of the final_prediction based on frequency
19:     return final_prediction, confidence
20: end function
```

UCLA Computational Medicine

# My Update

**Designed 2 KNN Algorithms**

**2. Evidence Based KNN:** Returns a list of patient ID's, "decisions", distance, where clinicians can further inspect the nearest neighbors. Very similar design minus bootstrapping. **Assumption is that we have enough compute to store these embeddings.**

**Algorithm 2** Returning the Results KNN

1: **function** RETURNRESULTSKNN($query\_text, data\_source, model\_name, n\_neighbors = 5$)
2:     $data\_frame \leftarrow$ DataFrame($data\_source$)
3:     $text\_embeddings \leftarrow$ Encode all texts in data_frame using model_name
4:     $knn \leftarrow$ initialize K-Nearest Neighbors model with cosine metric
5:     $knn.fit(text\_embeddings)$
6:     $query\_embedding \leftarrow$ Encode the query text to a vector using model_name
7:     $distances, indices \leftarrow knn.kneighbors(query\_embedding)$
8:     $results \leftarrow$ retrieve rows from data_frame corresponding to indices
9:     $results['distance'] \leftarrow distances$
10:     $prediction \leftarrow$ apply majority vote on the 'eddischarge' field of results
11:     $most\_common\_label \leftarrow$ find the most frequent label from prediction
12:     **return** $results, most\_common\_label$
13: **end function**

# Results 1: KNN Bootstrapping method

Just as a sanity check. I tried a base bert model versus a clinical bert model

```
In [46]:   query_text = df["patient_info"].iloc[0]
           data_source = temp
           model_name = "bert-base-uncased"
           prediction, confidence = knn(query_text, data_source, model_name)
           print(f"Predicted eddischarge: {prediction}")
           print(f"Ground Truth: {df['eddischarge'].iloc[0]}")
           print(f"Confidence: {confidence:.2f}")
```

10k

```
100%|████████| 10000/10000 [02:05<00:00, 79.70it/s]
embeddings are generated
100%|████████| 1/1 [00:00<00:00, 77.52it/s]
embeddings are generated
Bootstrapping: 100%|████████| 1000/1000 [00:03<00:00, 271.03it/s]
Predicted eddischarge: 0
Ground Truth: 1
Confidence: 0.78
```

1k

UCLA Computational Medicine

# Results

Just as a sanity check. I tried a base bert model versus a clinical bert model

In [47]:

```python
query_text = df["patient_info"].iloc[0]
data_source = temp
model_name = "medicalai/ClinicalBERT"
prediction, confidence = knn(query_text, data_source, model_name)
print(f"Predicted eddischarge: {prediction}")
print(f"Ground Truth: {df['eddischarge'].iloc[0]}")
print(f"Confidence: {confidence:.2f}")
```

```
100%|████████| 10000/10000 [01:11<00:00, 140.37it/s]
embeddings are generated
100%|████████| 1/1 [00:00<00:00, 132.01it/s]
embeddings are generated
Bootstrapping: 100%|████████| 1000/1000 [00:03<00:00, 277.14it/s]
Predicted eddischarge: 1
Ground Truth: 1
Confidence: 1.00
```

UCLA Computational Medicine

# Results 2: "Evidence Based" KNN

In [38]:
```python
# Example
query_text = df["patient_info"].iloc[0]
data_source = temp
model_name = "bert-base-uncased" # or any other model name
results, prediction = query_knn_embeddings(query_text, data_source, model_name)
print(results)
print("--------")
print(f"Predicted eddischarge: {prediction}")
```

```
100%|████████| 1000/1000 [00:12<00:00, 80.19it/s]
embeddings are generated
100%|████████| 1/1 [00:00<00:00, 77.90it/s]
embeddings are generated
                   ID                          patient_info  \
32394613    18125751    Patient 18125751, a 57 year old white female, ...
34022538    14187451    Patient 14187451, a 53 year old black/african ...
35671330    18242530    Patient 18242530, a 76 year old white female, ...
37661549    14062869    Patient 14062869, a 41 year old white female, ...
32129835    12019283    Patient 12019283, a 25 year old black/african ...

            eddischarge   distance
32394613              0   0.013725
34022538              1   0.016040
35671330              0   0.016175
37661549              0   0.016219
32129835              0   0.017117
--------
Predicted eddischarge: 0
Ground Truth: 1
```

# Results 2: "Evidence Based" KNN

```
In [39]:  query_text = df["patient_info"].iloc[0]
          data_source = temp
          model_name = "medicalai/ClinicalBERT"
          results, prediction = query_knn_embeddings(query_text, data_source, model_name)
          print(results)
          print("-------")
          print(f"Predicted eddischarge: {prediction}")
          print(f"Ground Truth: {df['eddischarge'].iloc[0]}")
```

```
100%|██████████| 1000/1000 [00:07<00:00, 141.10it/s]
embeddings are generated
100%|██████████| 1/1 [00:00<00:00, 134.61it/s]
embeddings are generated
                   ID                                    patient_info  \
36599058  16315929    Patient 16315929, a 62 year old white female, ...
32576195  13374041    Patient 13374041, a 58 year old white female, ...
33509281  10018862    Patient 10018862, a 56 year old white female, ...
31725842  18307993    Patient 18307993, a 45 year old black/african ...
32638903  13471464    Patient 13471464, a 73 year old white female, ...

          eddischarge  distance
36599058            1  0.000650
32576195            1  0.000689
33509281            1  0.000711
31725842            1  0.000743
32638903            1  0.000770
-------
Predicted eddischarge: 1
Ground Truth: 1
```
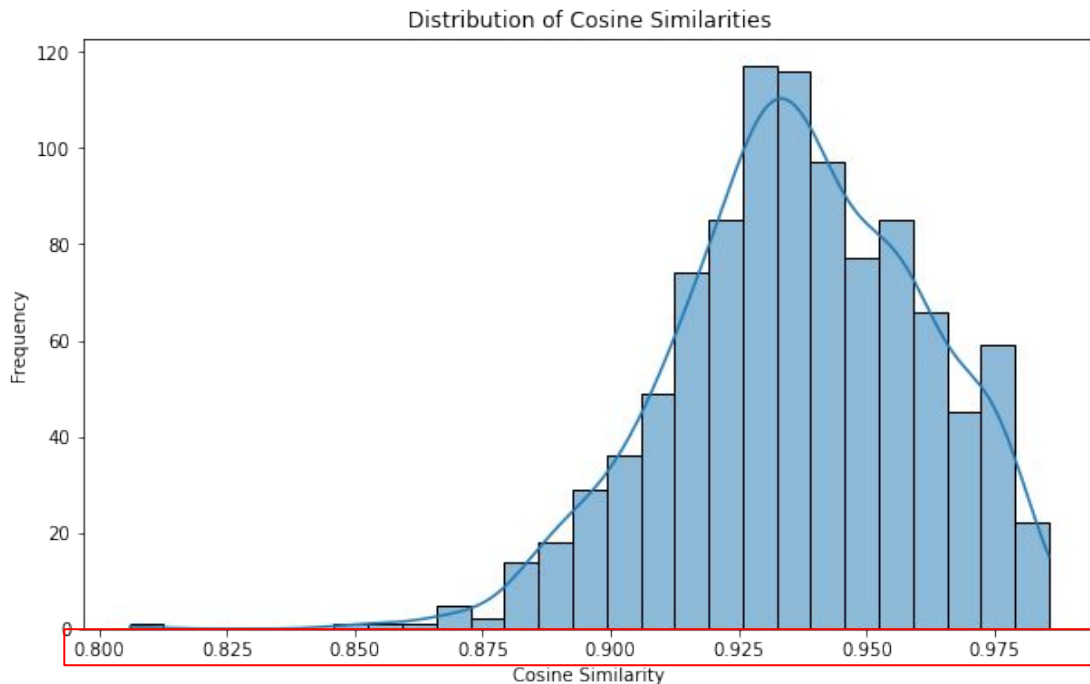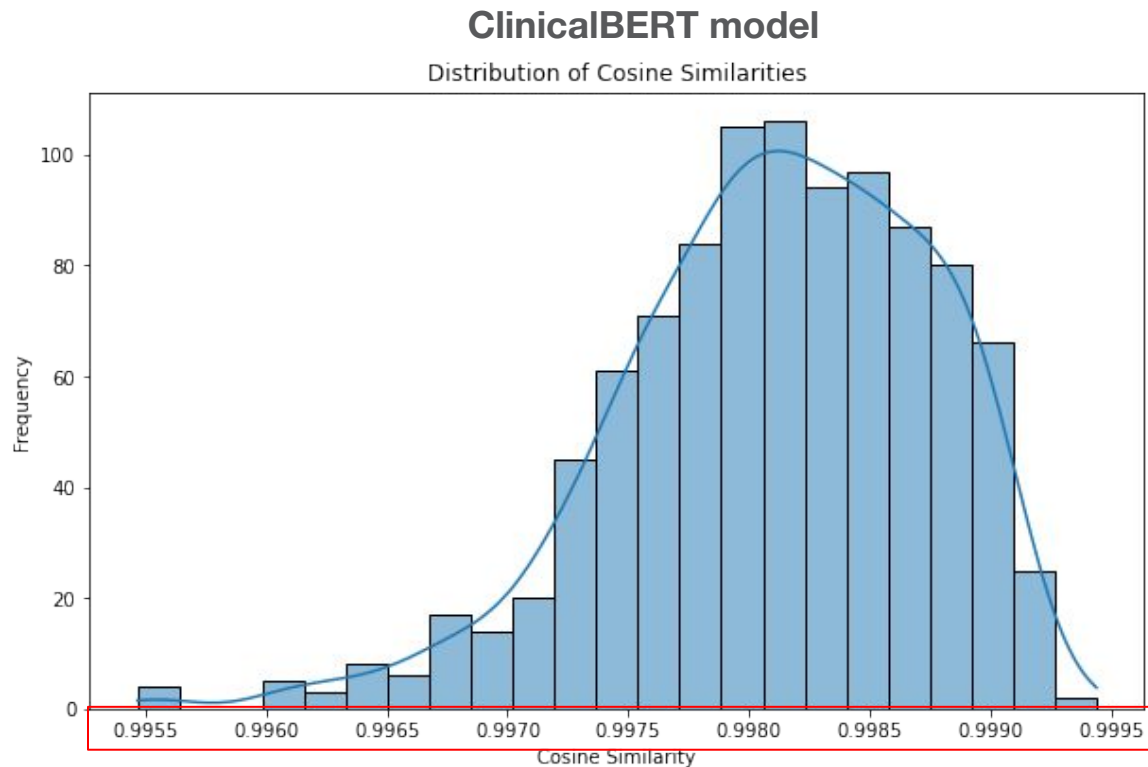
UCLA Computational Medicine

# Similarities

# Similarity Distribution for patient 0

**DistilBERT model**
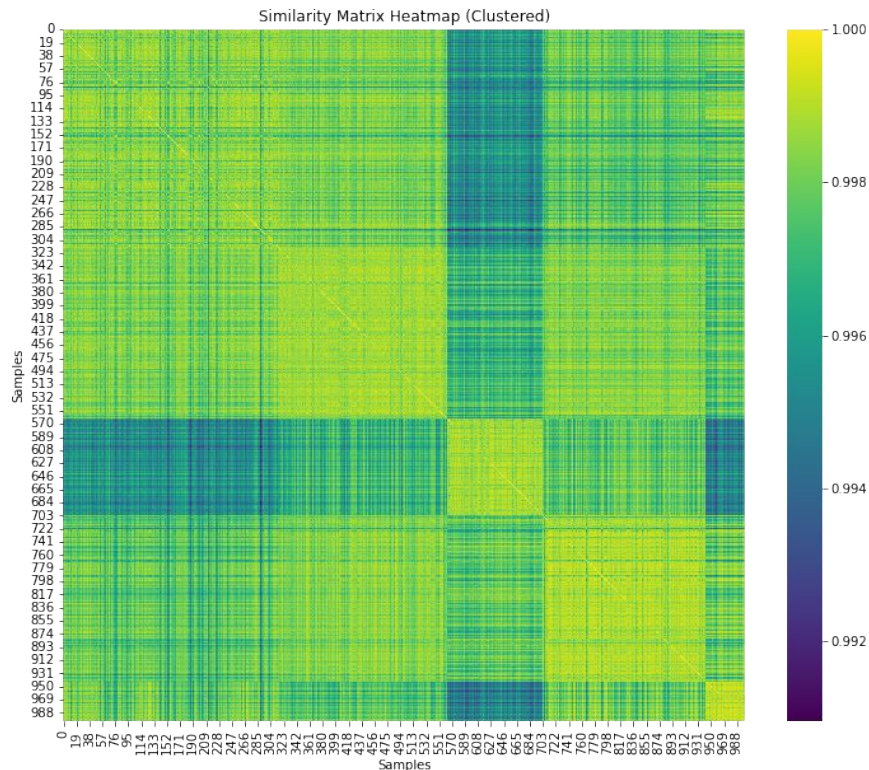
# Similarity Distribution for patient 0

# Similarity Matrices to see Groups

## DistilBERT model using Agglomerative Clustering

# Similarity Matrices to see Groups

## ClinicalBERT model using Agglomerative Clustering



Similarity Matrix Heatmap (Clustered)

# Similar to BERTopic

# Differences Between Bertopic and Cosine Similarity of Embedding

**BERTopic (Topic Modeling)**

**Purpose:** Extracting topics from a large collection of text. It helps in identifying and clustering similar documents based on the content

**Cosine Similarity**

**Purpose:** Measure the cosine of the angle between two embeddings in a multidimensional space. It qualifies how similar two documents are based on their content

**Clustering methodology?**

I wonder how different clustering methodologies affect model (HDBSCAN vs Agglomerative)?

- **Resource:** FAISS (Facebook AI Simiarlity Search) Package

UCLA | Computational Medicine

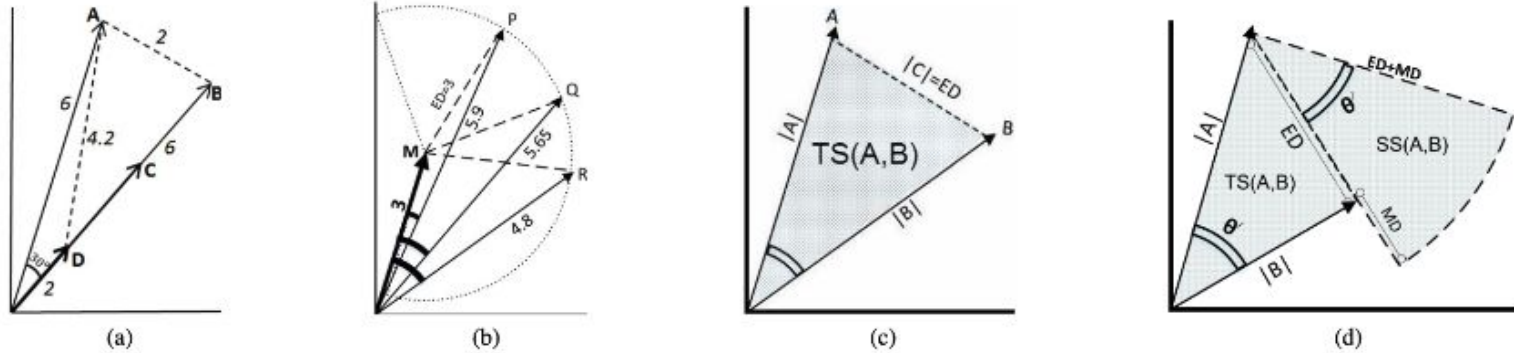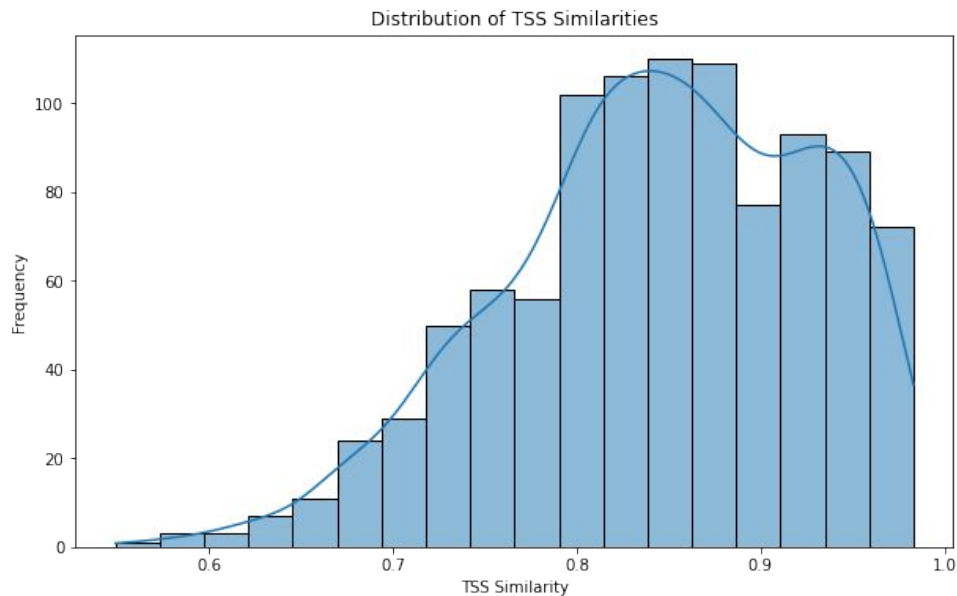# Triangle similarity section similarity (TS-SS)



Figure 1: (a)Example of Cosine drawback. (b)Example of ED drawback. (c)Triangle Similarity (TS). (d)Triangle Similarity-Section Similarity (TS-SS)

**Distance Metric**

The triangle similarity section similarity (TS-SS) considers both angle and magnitude difference, whereas something like cosine only considers angle
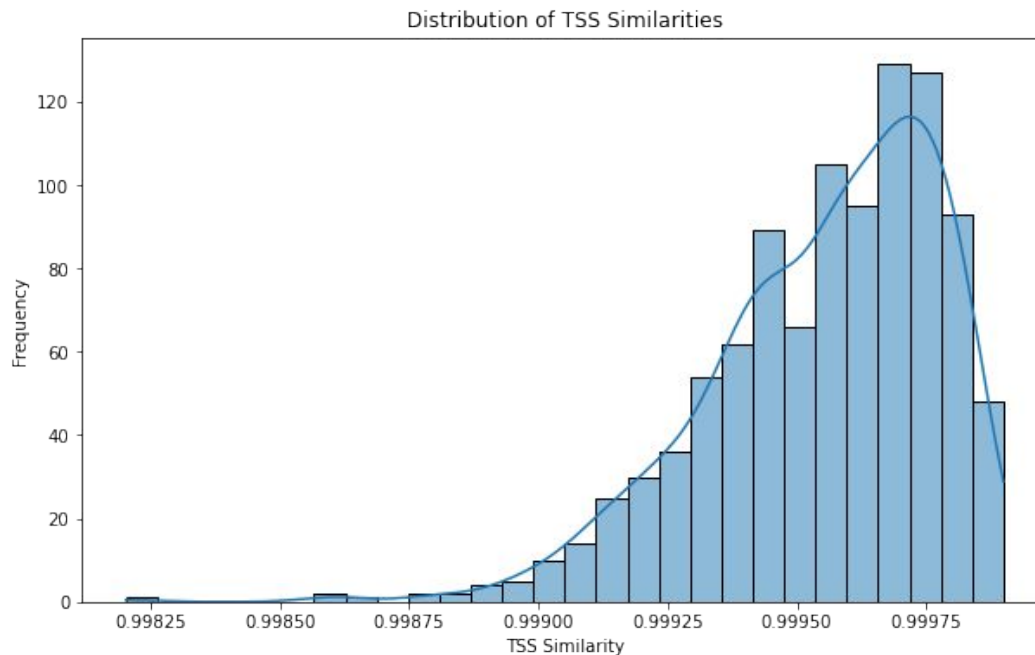
# Now what if we change the distance metric

**DistilBERT model**



Distribution of TSS Similarities

# Now what if we change the distance metric

## ClinicalBERT model

# Scaling up the Experiments

UCLA Computational Medicine

# AUROC and AUPRC Curves (test set = 2k)

**Both models appear to be working pretty well.**

# AUROC and AUPRC Curves (test set = 2k)

**DistilBERT**

**ClinicalBERT**

**Cosine Sim.**

```
Predictions completed for test data.
AUROC: 0.88 (95% CI: 0.86-0.89)
AUPRC: 0.76 (95% CI: 0.73-0.79)
F1 Score: 0.81 (95% CI: 0.79-0.82)
```

```
Predictions completed for test data.
AUROC: 0.95 (95% CI: 0.94-0.96)
AUPRC: 0.88 (95% CI: 0.86-0.90)
F1 Score: 0.90 (95% CI: 0.89-0.92)
```

**TS-SS**

```
Results:
AUROC: 0.867 (95% CI: 0.850-0.882)
AUPRC: 0.727 (95% CI: 0.692-0.761)
F1: 0.789 (95% CI: 0.771-0.808)
```

```
Results:
AUROC: 0.946 (95% CI: 0.936-0.956)
AUPRC: 0.866 (95% CI: 0.838-0.890)
F1: 0.895 (95% CI: 0.881-0.909)
```

**UCLA** Computational Medicine

# XAI...

# XAI… ??

We compute every token of the query and take the dot product with each token in the top k most similar exemplars. However as expected, the subword tokenizer does not give very explainable results…
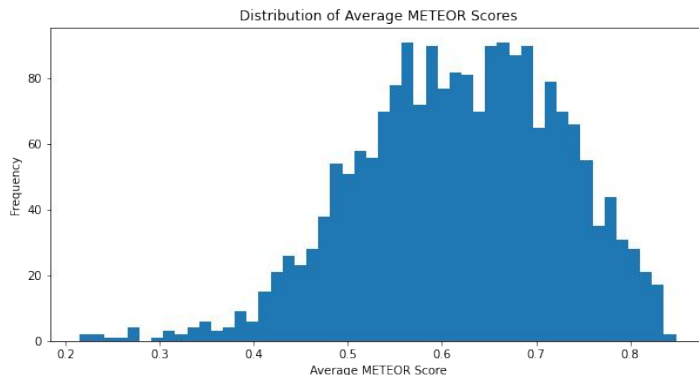
**One Example:**

```
Important tokens for neighbor 1:
    Query: pressure, Sample: 0, Importance: 0.0851
    Query: ##ic, Sample: ##ic, Importance: 0.0669
    Query: sp, Sample: ##tives, Importance: 0.0297
    Query: acu, Sample: ab, Importance: 0.0285
    Query: 0, Sample: ##ic, Importance: 0.0265
    Query: satu, Sample: ., Importance: 0.0242
    Query: ##ic, Sample: 77, Importance: 0.0238
    Query: ##d, Sample: :, Importance: 0.0229
    Query: dis, Sample: med, Importance: 0.0220
    Query: ##tero, Sample: gi, Importance: 0.0214
```
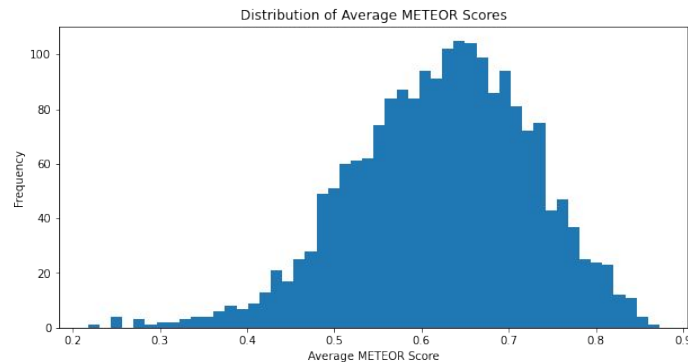
**UCLA** Computational Medicine

# Current Work

# METEOR Scores (not done…)
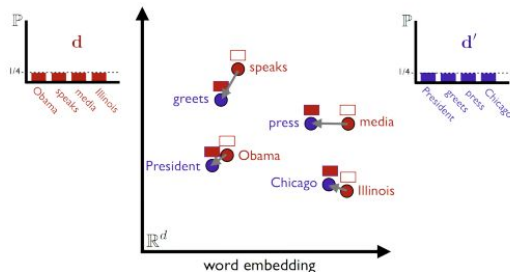
**DistilBERT**



**ClinicalBERT**



**METEOR Score:**

Aligns/matches words in the hypothesis to reference by sequentially applying exact match, stemmed match and wordnet based synonym match. In case there are multiple matches the match which has the least number of crossing is chosen.

# Mover's Distance

**Should we use this pre-trained model for even our text?**

Did limited exploration but I'm not sure how much sense it makes given our clinical text



```
import gensim.downloader as api
model = api.load('word2vec-google-news-300')
```

```
sentence_orange = preprocess('Oranges are my favorite fruit')
distance = model.wmdistance(sentence_obama, sentence_orange)
print('distance = %.4f' % distance)
```

**Earth Mover's Distance**
Word mover's distance is based on the Earth Mover's Distance. Our goal is to calculate the distance that a word embedding needs to travel to reach the word embedding of the other word embedding. Now, when we talk about sentences, we need to calculate the distance of each word in one sentence to every other word in another sentence.

**UCLA** Computational Medicine