

Quantifying Mental Health from Social Media with Neural User Embeddings

Silvio Amir

*INESC-ID Lisboa, Instituto Superior Técnico, Universidade de Lisboa
Lisboa, Portugal*

SAMIR@INESC-ID.PT

Glen Coppersmith

*Qntfy
Washington DC, United States*

GLEN@QNTFY.COM

Paula Carvalho

*INESC-ID Lisboa, and Universidade Europeia, LIU
Lisboa, Portugal*

PCC@INESC-ID.PT

Mário J. Silva

*INESC-ID Lisboa, Instituto Superior Técnico, Universidade de Lisboa
Lisboa, Portugal*

MJS@INESC-ID.PT

Byron C. Wallace

*Northeastern University
Boston MA, United States*

B.WALLACE@NORTHEASTERN.EDU

Abstract

Mental illnesses adversely affect a significant proportion of the population worldwide. However, the methods traditionally used for estimating and characterizing the prevalence of mental health conditions are time-consuming and expensive. Consequently, best-available estimates concerning the prevalence of mental health conditions are often years out of date. Automated approaches to supplement these survey methods with broad, aggregated information derived from social media content provides a potential means for near real-time estimates at scale. These may, in turn, provide grist for supporting, evaluating and iteratively improving upon public health programs and interventions.

We propose a novel model for automated mental health status quantification that incorporates *user embeddings*. This builds upon recent work exploring *representation learning* methods that induce embeddings by leveraging social media post histories. Such embeddings capture latent characteristics of individuals (e.g., political leanings) and encode a soft notion of homophily. In this paper, we investigate whether user embeddings learned from twitter post histories encode information that correlates with mental health statuses. To this end, we estimated user embeddings for a set of users known to be affected by depression and post-traumatic stress disorder (PTSD), and for a set of demographically matched ‘control’ users. We then evaluated these embeddings with respect to: (i) their ability to capture homophilic relations with respect to mental health status; and (ii) the performance of downstream mental health prediction models based on these features. Our experimental results demonstrate that the user embeddings capture similarities between users with respect to mental conditions, and are predictive of mental health.

1. Introduction

Mental illness is a critically important concern, significantly and adversely affecting a wide swath of the population directly and indirectly. An estimate by the Centers for Disease Control from 2008 (CDC, 2010), suggests that 9% of US adults may meet the criteria for

depression at any given time. While not as prevalent as depression, post traumatic stress disorder (PTSD) issues still cost hundreds of billions of dollars worldwide, according to a conservative estimate from the NIH¹. The collective effect of mental health conditions, as measured by Daily Adjusted Life Years (DALYs), exceeds that of malaria, war, or violence² (?). At the same time, mental health problems are often difficult to identify and thus treat. For example, perhaps half of depressive cases go undetected, in part due to the heterogeneous and complex expression of this condition (Paykel et al., 1997). Another exacerbating factor is that diagnosis generally requires individuals to actively seek out treatment. Yet, the manifestation of this condition and prevailing social stigmas may discline afflicted individuals to seek treatment.

The internet may provide a comfortable medium for people to express their feelings anonymously and connect with health-care professionals (McCaughey et al., 2014) and others affected by similar conditions (De Choudhury et al., 2016). Furthermore, individuals openly discuss mental health challenges on public social network platforms such as Twitter (Coppersmith et al., 2014a, 2015a). Prior work has demonstrated the potential of using social media to investigate mental health issues (Paul and Dredze, 2011), including depression (Schwartz et al., 2014), PTSD (Coppersmith et al., 2014b) and suicidal ideation (Coppersmith et al., 2016; De Choudhury et al., 2016) in individuals. However, models and techniques to identify and quantify mental health related signals from social media are relatively novel. Interest in these applications has motivated the creation of a shared task for the Computational Linguistics and Clinical Psychology workshop (CLPsych)³, which aimed to advance the state-of-the-art in technologies capable of discriminating users affected by mental illness from controls, given their post history (Coppersmith et al., 2015b). A variety of methods have been proposed for this task, but none have achieved consistently superior performance, which implies that improvements may yet be realized by improved models.

Neural representation learning methods have been shown capable of automatically discover good representations (e.g., predictive features) from data, freeing practitioners from the burden of manually designing and encoding task-specific features (Bengio et al., 2015a; Goldberg, 2016). In natural language applications, this has resulted in neural distributed representations becoming the *de-facto* standard representational approach. *Word embeddings* in particular aim to implicitly encode latent word semantics (in the distributional sense), and can be learned in an unsupervised fashion by means of predictive models that exploit word co-occurrence statistics and other regularities in unlabeled corpora (Bengio et al., 2003). These models have been recently extended to infer representations for larger textual units (Le and Mikolov, 2014), and even user representations (Li et al., 2015; Amir et al., 2016). It has been shown that these *user embeddings* also capture latent user aspects and can be used in downstream applications, such as sarcasm detection (Amir et al., 2016) and content recommendation (Yu et al., 2016).

In this paper we investigate whether user representations induced via neural models can inform clinical models operating over social media. In particular we consider whether user embeddings (learned directly from historical data) capture aspects of mental health status. To this end we leverage the dataset created for the CLPsych shared task to address two

1. <https://www.nimh.nih.gov/health/statistics/cost/index.shtml>

2. For a visualization of DALYs see <https://vizhub.healthdata.org/gbd-compare/>

3. <http://clpsych.org>

research questions: (1) To what extent do user embeddings capture information relevant for mental health analysis applications, over social media? (2) Can user embeddings be leveraged to discriminate between users suffering from mental illness and demographically matched controls? We answer the first question by comparing different approaches to induce user representations from a collection of posts. In particular, we investigated whether the induced embeddings capture homophilic relations between users with respect to mental health. To answer the second question, we developed and evaluated predictive models, leveraging user embeddings, to discriminate users affected by depression, PTSD, and age- and gender-match controls.

The main contributions of this paper are as follows: (i) we show that unsupervised user embeddings induced from posting histories capture user similarities, and are predictive of mental health conditions; (ii) we develop a novel neural model that incorporates and refines these embeddings to improve the categorization of users with respect to mental health status; furthermore, we show that the resultant fine-tuned user embeddings better align with mental health conditions.

The remainder of the paper is organized as follows. The next section introduces the aforementioned CLPsych shared task and the corresponding dataset. Section 3 reviews the literature on user modelling for social media analysis and neural embedding learning. In Section 4, we formally describe the user embedding model used in our experiments, and discusses the connections with prior representation learning methods. Section 5 addresses the first research question by evaluating the properties captured by the user embeddings. Section 6 reports on the classification experiments we conducted to answer the second research question. Finally, we present our conclusions in Section 7.

2. Depression and PTSD on Twitter

In 2015, the CLPsych workshop held a shared task to foster progress in NLP technologies with potential for applications related to mental health analysis, over social media streams (Mitchell et al., 2015; Coppersmith et al., 2015b). To that end, a dataset was compiled comprising users that have publicly stated on Twitter that they were diagnosed with depression (327 users) or PTSD (246 users), and an equal number of randomly selected demographically-matched users as *controls*⁴. For each user in this dataset, associated metadata and posting history was also collected — up to the 3000 most recent tweets, per limitations of the Twitter API. For more details on the construction and validation of the data, see (Coppersmith et al., 2015b, 2014a,b).

The participants were then asked to develop models to discriminate between users affected by mental illness from controls, given their posts and metadata. Specifically this entailed three binary sub-tasks: (i) **depression vs control**, (ii) **PTSD vs control** and (iii) **depression vs PTSD**. The proposed systems were based on a wide range of approaches including: rule-based systems leveraging lexical decision lists (Pedersen, 2015), linear classifiers exploiting features based on word clusters and topic models (Preotiuc-Pietro et al., 2015), supervised topic models (Resnik et al., 2015) and systems exploiting character-level language models (Coppersmith et al., 2015b). However, we note that none of the proposed

4. This data was collected according to the ethical protocol of ?, and follows the recommendations spelled out in Mikal et al. (2016).

systems performed consistently better than the others across all the sub-tasks and evaluation metrics, highlighting the difficulty of this problem. Moreover, none of these systems used explicit representations of *users*, which is the innovation we propose here. We did not participate in this shared task, and thus could not obtain the official test data. Therefore, our results are not directly comparable to those of the participating teams. Nevertheless, we compared our proposed approach with the majority of the previously proposed methods.

3. Related Work

Most of the research in social media analysis has been concerned with deriving better models that operate on representations of the texts comprising individual users posts, both via manually crafted features and, more recently, representation learning approaches (Severyn and Moschitti, 2015; Astudillo et al., 2015). Nevertheless, for a variety of problems it is crucial to also capture characteristics of the *users* involved in the communications. These include, information extraction (Yang et al., 2016), opinion mining (Tang et al., 2015), sarcasm detection (Bamman and Smith, 2015) and content recommendation (Yu et al., 2016). The most straightforward approach to induce user representations is by scrapping “profile” information from social websites, to manually extract features based on social ties, demographic attributes, or posting habits (Bamman and Smith, 2015; Rajadesingan et al., 2015). However, these approaches require significant effort for data collection, and for task- and domain-specific feature engineering. Furthermore, the available user profile information depends on specific social websites and may not always be available. It may also be inaccurate or simply outdated.

Neural Embedding Learning

In recent years, models in NLP have moved from *discrete* word representations, based on scalars representing indices into a pre-defined vocabulary, towards *distributed* continuous vector word representations that encode latent semantics — these are usually referred as *word embeddings* (Goldberg, 2016). The general framework to learn unsupervised word embeddings involves associating words with parameter vectors, which are then optimized to be good predictors of other words that occur in the same contexts (Bengio et al., 2003). SKIP-GRAM (Mikolov et al., 2013), one of the most popular word embedding models, operationalizes this approach by sliding a *window* of a pre-specified size across the corpus. At each step, the center word is used to predict the probability of one of the surrounding words, sampled proportionally to the distance to the center word. Le and Mikolov (2014) later expanded this approach with two PARAGRAPH2VEC models that also learn representations for paragraphs (or, more broadly, any sequence of words): (i) PV-DM, tries to predict the center word of the sliding window, given the surrounding words **and** the paragraph (i.e., their respective embeddings); and (ii) PV-DBOW, tries to predict the words of a sliding window within a paragraph, conditioned only on the respective paragraph embedding.

Recently proposed methods to learn user representations use essentially the same approach — associate users with parameter vectors, and optimize these to accurately predict observable attributes or the words used by said user in previous posts (Li et al., 2015; Amir et al., 2016; Yu et al., 2016). As discussed above, leveraging user profile information to collect attributes is not always reliable. Interestingly, however, the embeddings induced

by Amir et al. (2016), using only the previous posts from a user, were shown to capture latent user aspects (e.g. political leanings) and a soft notion of ‘homophily’ — *similar* users were generally represented with similar vectors. Furthermore, these user representations were successfully used to improve a downstream model for sarcasm detection in tweets. Similarly, Yu et al. (2016) used user embeddings to improve a microblog recommendation system. Our work aims to ascertain if user embeddings learnt only from previous posts, can capture useful signals for clinical applications.

4. Learning User Embeddings

To learn user embeddings, we adopted an approach similar to that recently proposed by Amir et al. (2016). The idea is to capture relations between users and the content (i.e., the words) they generate, by optimizing the probability of sentences conditioned on their authors. Formally, let \mathcal{U} be a set of users, \mathcal{C}_j be a collection of posts authored by user $u_j \in \mathcal{U}$, and $S = \{w_1, \dots, w_N\}$ be a post composed of words w_i from a vocabulary \mathcal{V} . The goal is to estimate the parameters of a user vector \mathbf{u}_j , that maximize the conditional probability:

$$P(\mathcal{C}_j|u_j) \propto \sum_{S \in \mathcal{C}_j} \sum_{w_i \in S} \log P(w_i|\mathbf{u}_j) \tag{1}$$

However, directly estimating these quantities (e.g., with a log-linear model) would require calculating a normalizing constant over a potentially large number of words, a computationally expensive operation. Because we are only interested in the user vectors \mathbf{u}_j and not the actual probabilities as such, we can approximate the term $P(w_i|\mathbf{u}_j)$ by minimizing the following Hinge-loss objective:

$$\mathcal{L}(w_i, u_j) = \sum_{\tilde{w}_k \in \mathcal{V}} \max(0, 1 - \mathbf{w}_i \cdot \mathbf{u}_j + \tilde{\mathbf{w}}_k \cdot \mathbf{u}_j) \tag{2}$$

where word \tilde{w}_k (and associated embedding, $\tilde{\mathbf{w}}_k$) is a *negative sample*, i.e. a word not occurring in the sentence under consideration, which was written by user u_j . By learning to discriminate between observed positive examples and *pseudo*-negative examples, the model shifts probability mass to more plausible observations (Smith and Eisner, 2005). Note that we represent both *words* and *users* via d -dimensional embeddings — word embeddings, $\mathbf{w}_i \in \mathbb{R}^d$ which are assumed to have been pre-trained through some neural language model; and user embeddings $\mathbf{u}_j \in \mathbb{R}^d$ to be learned. We will refer to this approach as USER2VEC.⁵

We note that barring some minor operational differences, this model is equivalent to the PV-DBOW variant of PARAGRAPH2VEC — if users are viewed as paragraphs. The key differences are that: (i) USER2VEC predicts **all** the words in a post, whereas PV-DBOW slides a window along the paragraph and only predicts one word per step; and (ii) USER2VEC assumes that the word embeddings are pre-trained, whereas PV-DBOW aims to jointly learn the word and paragraph vectors.

5. This formulation is a simplification of Amir et al. (2016) model. Specifically, we omitted a term in Eq.1, encoding the marginal probability of S ; and we allow the negative samples to be drawn from all the words in \mathcal{V} . These simplifications dramatically reduce training time without significant loss of quality on the resulting embeddings.

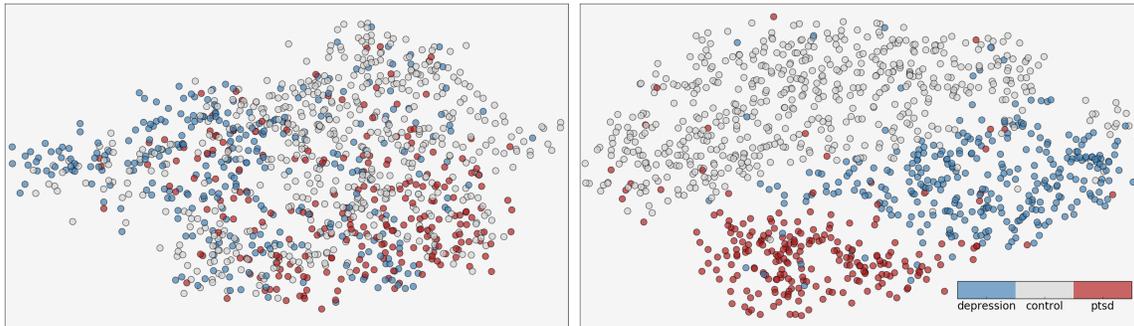


Figure 1: PV-DM embeddings projected into two dimensions, and colored according to the respective cohort. The plot on the left-hand side shows the original unsupervised embeddings — despite being trained without labels, in this space, users tend to be surrounded by others from the same cohort. The plot on the right-hand side, shows the result of an embedding subspace projection induced with the NLSE model — we can see that, in this adapted space, the users are better clustered by cohort.

5. User Embedding Analysis

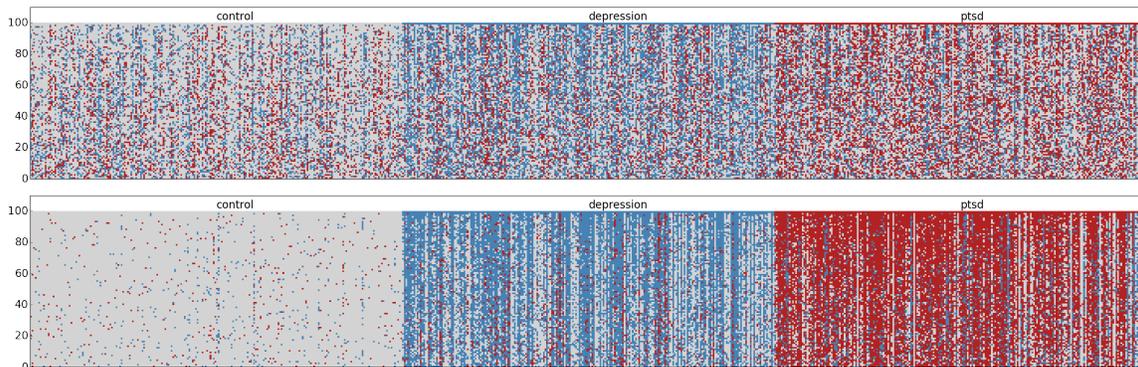
In this section we address our first research question by investigating whether user embeddings learned directly from social media data encode information relevant for public health applications. Previous work has shown that these representations capture latent user aspects and a soft notion of ‘homophily’. If some of these aspects correlate with mental health, then the embeddings could be used to identify risk groups; for example, we might identify and characterize users that ‘look’ like patients affected by depression. The ability to do so would potentially enable scalable, real-time estimates concerning the prevalence of mental health issues in particular populations.

We estimated USER2VEC embeddings from the shared task dataset described in Section 2, as follows.⁶ First, we pre-processed tweets by: lower-casing; reducing character repetitions to at most three repetitions; and replacing usernames and URLs with a canonical form. Users with fewer than 100 tweets were discarded. Next, we pre-trained a set of SKIP-GRAM word vectors from the task data and another large unlabeled Twitter corpus, using the `Gensim`⁷ python package (Řehůřek and Sojka, 2010). Finally, for each user $u_j \in \mathcal{U}$, we sampled a held-out set $\mathcal{H}_j \subset \mathcal{C}_j$ with 10% of the posting history. The rest of the data was used to estimate an embedding \mathbf{u}_j by minimizing Eq. 1 via stochastic gradient descent, using $P(\mathcal{H}_j|\mathbf{u}_j)$ as early stopping criteria.

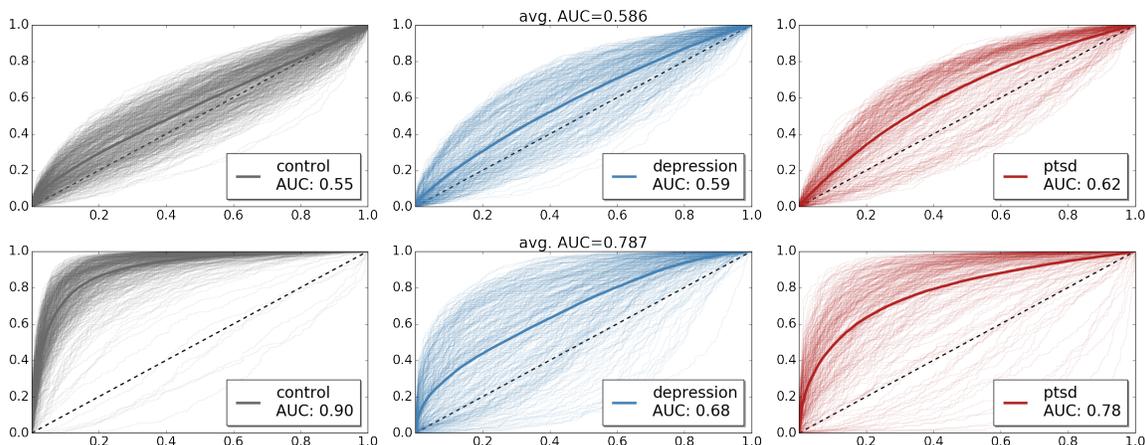
The same dataset was used to derive embeddings with PV-DM and PV-DBOW models (also via `Gensim`). To ensure a fair comparison, we used the same hyper-parameters for all the models, which were set as follows: window size $w = 5$, negative sample size $s = 20$ and vector size $d = 400$.

6. Code will be made available after publication

7. <http://radimrehurek.com/gensim/>



(a) User vector similarity ranking. The first row corresponds to a 'query' user, and the columns show the top 100 most similar users, colored according to their class.



(b) ROC curves and AUC scores of the induced user similarity rankings, per class.

Figure 2: Measuring homophilic relations with respect to mental conditions with vector distances over the user embedding space. The top-most sub-plots refer to rankings induced with the PV-DM model, and the ones at the bottom correspond to rankings obtained with embeddings adapted with the NLSE model.

5.1 Measuring Homophily

To investigate if the induced user vectors capture homophilic relations with respect to mental health status, we first projected the d -dimensional vectors into a 2-dimensional space using t-Stochastic Neighborhood Embedding (TSNE) (Van der Maaten and Hinton, 2008). In Figure 1, we plot the resulting points colored according to the respective class. One can see that, at least to some extent, embeddings do seem to capture some notion of homophily, i.e. users are often surrounded by others of the same cohort.

To better quantify this effect and allow comparison of different user embedding models, we proceeded as follows. For each user in the corpus, we calculated the similarities to all *other* users (i.e., the cosine similarity between their respective embeddings), inducing a ranking of users in terms of similarity to the 'query' user. Intuitively, we would hope

to see that users in the same mental health categories are comparatively similar to one another, i.e., users suffering from depression are most similar to other users also afflicted with depression.

Figure 2 shows the results obtained with PV-DM vectors. The induced similarity ranking is shown in Figure 2a, where the first row correspond to the query users, and each column shows their top $k = 100$ most similar users, colored according to their class. Figure 2b shows the respective Receiver Operating Characteristic (ROC) curves under the induced ranking. In general, we found that all the user embedding models are able to capture user similarities. The embeddings induced with the USER2VEC, PV-DBOW and PV-DM all perform significantly better than chance, with AUC scores of 0.57, 0.57 and 0.59, respectively. Detailed plots can be found in the Appendix. The fact that user vectors are more likely to be close to those of others in the same cohort demonstrates that this approach does indeed capture signals relevant to mental health. This aligns with prior work showing that one’s choice of words can be indicative of psychological states and mental health (Pennebaker et al., 2001).

6. Predicting Mental Health from Twitter Data

To address our second research question, we evaluated the user representations with respect to their predictive performance in downstream mental health analysis applications. Above we showed that representations induced from user posting histories encode relevant signals about mental health. But generic features estimated from unsupervised tasks are sub-optimal for downstream tasks (Astudillo et al., 2015). Neural networks, when trained end-to-end, can refine generic embeddings to specific tasks by modifying their parameters during supervised training (Collobert et al., 2011). However, this strategy requires updating a large number of parameters, which is difficult when only small training datasets are available. The CLPsych dataset comprises 1094 labeled instances (after discarding users with fewer than 100 tweets). Given the modest size of the dataset, we adopted the Astudillo et al. (2015) Non-Linear Subspace Embedding approach, (NLSE), which is able to adapt generic representations to specific tasks with scarce labeled data.

6.1 Proposed Model

The NLSE model adapts generic embeddings to specific applications by learning linear projections into lower-dimensional subspaces, while the keeping the original embeddings fixed. Hence, the resulting *embedding subspaces* capture domain- and task-specific aspects, while preserving the rich information encoded by the original embeddings. More formally, given an user embedding matrix $\mathbf{U} \in \mathbb{R}^{d \times |\mathcal{U}|}$, in which column $\mathbf{U}_{[j]}$ represents user $u_j \in \mathcal{U}$, we induce new representations by factorizing the input as $\mathbf{S} \cdot \mathbf{U}$ where $\mathbf{S} \in \mathbb{R}^{s \times d}$, with $s \ll d$, is a (learned) linear projection matrix. Crucially, this imposes dimensionality reduction on the feature space, simultaneously eliminating noise and reducing the number of free parameters, making the model easier to train with small datasets.

This model is similar to a feed-forward neural network with a word embedding layer and a single hidden layer. The main differences are: (1) the factorization of the embedding layer into two components (the original embedding matrix and a linear projection matrix)

and, (2) the dimensionality reduction induced by the subspace projection, with typical reductions greater than an order of magnitude. Similar to feed-forward networks, we can use the backpropagation algorithm (Rumelhart et al., 1988) to jointly learn task-specific embeddings and the parameters for the classification layer. Using this approach, our proposed mental health prediction model can be formalized as:

$$\begin{aligned} P(\mathcal{Y}|u_j) &\propto \boldsymbol{\beta} \cdot g(u_j) \\ g(u_j) &= \sigma(\mathbf{S} \cdot \mathbf{U}_{[j]}) \end{aligned} \quad (3)$$

where, $\sigma(\cdot)$ denotes an element-wise sigmoid non-linearity, and the matrix $\boldsymbol{\beta} \in \mathbb{R}^{|\mathcal{Y}| \times s}$ maps the embedding subspace to the classification space. Notice that at inference time this model reduces to a linear classifier with embedding subspace features.

6.2 Experimental Setup

We evaluated embeddings induced with the USER2VEC (U2V), and PARAGRAPH2VEC’s PV-DM and PV-DBOW models (Section 5), as features in Logistic Regression (LR) and NLSE classifiers. These were compared against baselines using textual features based on:

1. BOW: bag-of-words vectors with binary weights, $\mathbf{x} \in \{0, 1\}^{|\mathcal{V}|}$;
2. BOE: bag-of-embeddings. We leveraged SKIP-GRAM embeddings to build vectors, $\mathbf{x} = \sum_w \mathbf{E}_{[w]}$, where $\mathbf{E}_{[w]} \in \mathbb{R}^d$ is the embedding of word w ;
3. LDA: bag-of-topics. We induced $t = 100$ topics using Latent Dirichlet Allocation (Blei et al., 2003), to build vectors $\mathbf{x} \in \{0, 1\}^t$ indicating the topics present in user’s posts;
4. BWC: bag-of-word-clusters. We induced $k = 1000$ Brown et al. (1992) word clusters, to build vectors $\mathbf{x} \in \{0, 1\}^k$ mapping words in a user’s posts to their respective clusters;

We also evaluated baselines that combine user vectors with textual features (U2V+BOW and U2V+BOE).

Experiments were conducted with a 10-fold cross-validation protocol; at each iteration, the training partition was split into 80% for model training and the remainder for validation purposes (i.e. hyper-parameter selection and early-stopping). For consistency, we used the same splits for all the models. We performed a grid-search to choose the best ℓ_2 regularization coefficient, over the range $c = \{0.001, 0.01, 0.5, 1, 10, 100\}$, for the LR models; and the optimal subspace size $s = \{10, 15, 20, 25\}$ and learning rate $\alpha = \{0.01, 0.1, 0.5, 1\}$, for the NLSE model.

6.3 Results

The classifiers were mainly evaluated with respect to the macro average F_1 . We also report results in terms of *binary* F_1 , where we only average the scores for the **depression** and **ptsd** classes, to better ascertain the ability of the models to discriminate between mentally afflicted patients, which are less prevalent than the controls, but are the cases that we mostly care about.

The main classification results are shown in Figure 3. The first thing to note is that the BOW is a very strong baseline, essentially outperforming all the other linear classifiers

based on textual features and user embeddings. One reason is that users affected with mental illnesses, often talk about their conditions and the BOW model can easily pick-up on such clues. Regarding the user embeddings, we found that, despite being equivalent, the PV-DBOW performed much worse than the USER2VEC, showing that better embeddings can be obtained by trying to predict **all** the words in users posts, and leveraging pre-trained word vectors. On the other hand, the PV-DM model has a performance comparable to that of the USER2VEC.

As discussed above, generic embeddings are sub-optimal for downstream tasks, and our results are in line with this observation. By inducing task-specific representations via subspace projection, we were able to outperform all the other baselines by a fair margin. Note also that the NLSE approach is particularly better at discriminating the minority classes, i.e. patients suffering from **depression** and **ptsd**, as evidenced by the greater improvements in *binary* F_1 , when compared to the other baselines. To better understand the effects of the embedding adaptation, we repeated the same analysis described in Section 5 over the adapted embeddings (Figures 1 and 2). We can see that the new representations are much better at discriminating the controls from the other users, suggesting that induced embedding subspace captures more fine-grained signal related to mental statuses.

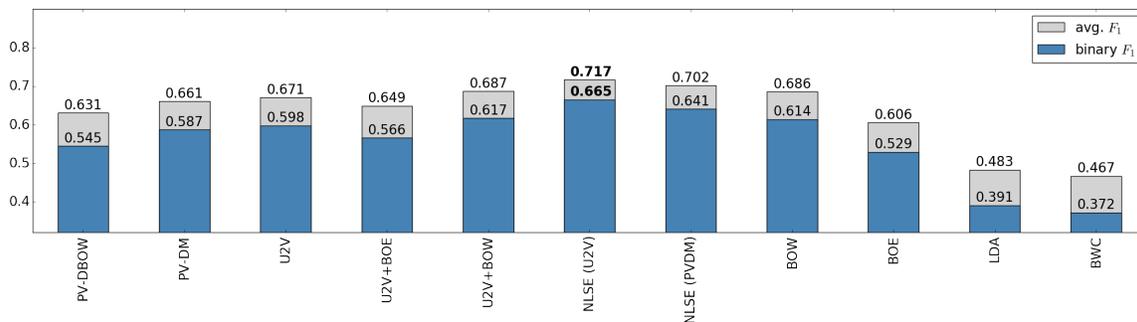


Figure 3: Performance of different models at discriminating users with respect to mental condition, in terms of F_1 and *binary* F_1 .

7. Conclusions

In this paper, we investigated if embeddings induced from user posting histories capture relevant signals for clinical applications. In particular, we compared different user embedding methods, with respect to their ability to capture homophilic relations between users, and their performance as features in downstream mental health prediction models. The evaluation conducted over a dataset comprising of users diagnosed with depression and PTSD, and demographically matched controls, showed that these representations can indeed capture mental health related signals. This is in agreement with prior results from the field of psychology, establishing connections between word usage and mental status (Pennebaker et al., 2001). Interestingly, embeddings induced without knowledge of user labels capture similarities with respect to mental condition. Furthermore, we have shown that

these embeddings can be tailored — with a small amount of task-specific labeled data — to capture more granular information, thus improving the quality of downstream models and applications. Ultimately, this work is a step toward more accurate inference concerning the mental health status of social media users, in turn enabling more accurate epidemiological real-time population-wide monitoring of mental health. Such accurate monitoring — which is currently impossible — may provide empirical support for increased resource allocation to programs dedicated to preventing and alleviating mental health issues.

Moving forward, user embeddings may provide a pivotal piece to allow clinical psychologists to take full advantage of digital phenotyping data. In particular, learned user embeddings may provide a representation at a sweet spot between instantaneous (proximal) state and lifelong (distal) state, which is critical to understanding psychological phenomena and risk of crisis.

References

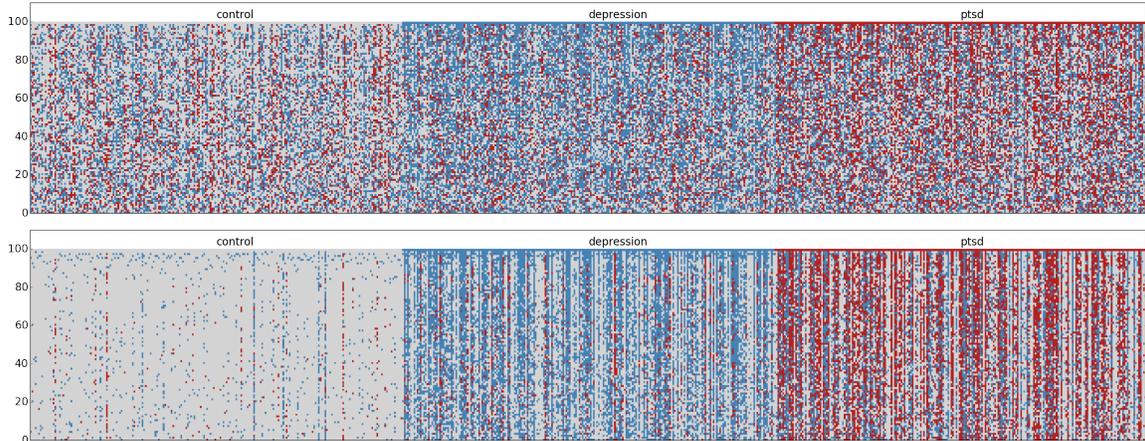
- Silvio Amir, Byron C Wallace, Hao Lyu, and Paula Carvalho Mário J Silva. Modelling context with user embeddings for sarcasm detection in social media. *arXiv preprint arXiv:1607.00976*, 2016.
- Ramón Astudillo, Silvio Amir, Wang Ling, Mario Silva, and Isabel Trancoso. Learning word representations from scarce and noisy data with embedding subspaces. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1074–1084, Beijing, China, July 2015.
- David Bamman and Noah A Smith. Contextualized sarcasm detection on twitter. In *Proceedings of the 9th International Conference on Web and Social Media*, pages 574–77. AAAI Menlo Park, CA, 2015.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic language model. *The Journal of Machine Learning Research*, 3:1137–1155, 2003.
- Yoshua Bengio, Ian J. Goodfellow, and Aaron Courville. Deep learning. Book in preparation for MIT Press, 2015a. URL <http://www.iro.umontreal.ca/~bengioy/dlbook>.
- Yoshua Bengio, Ian J Goodfellow, and Aaron Courville. Deep learning. *Nature*, 521:436–444, 2015b.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022, 2003.
- Peter F. Brown, Peter V. Desouza, Robert L. Mercer, Vincent J. Della Pietra, and Jenifer C. Lai. Class-based n-gram models of natural language. *Computational Linguistics*, 18(4), 1992.
- CDC. Current depression among adults—united states, 2006 and 2008. *MMWR. Morbidity and mortality weekly report*, 59(38):1229, 2010.

- Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *The Journal of Machine Learning Research*, 12:2493–2537, 2011.
- Mike Conway. Ethical issues in using twitter for public health surveillance and research: developing a taxonomy of ethical concepts from the research literature. *Journal of medical Internet research*, 16(12):e290, 2014.
- Glen Coppersmith, Mark Dredze, and Craig Harman. Quantifying mental health signals in Twitter. In *Proceedings of the ACL Workshop on Computational Linguistics and Clinical Psychology*, 2014a.
- Glen Coppersmith, Craig Harman, and Mark Dredze. Measuring post traumatic stress disorder in Twitter. In *Proceedings of the 8th International AAAI Conference on Weblogs and Social Media (ICWSM)*, 2014b.
- Glen Coppersmith, Mark Dredze, Craig Harman, and Kristy Hollingshead. From ADHD to SAD: Analyzing the language of mental health on Twitter through self-reported diagnoses. In *Proceedings of the Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, Denver, Colorado, USA, June 2015a. North American Chapter of the Association for Computational Linguistics.
- Glen Coppersmith, Mark Dredze, Craig Harman, Kristy Hollingshead, and Margaret Mitchell. CLPsych 2015 shared task: Depression and PTSD on Twitter. In *Proceedings of the Shared Task for the NAACL Workshop on Computational Linguistics and Clinical Psychology*, 2015b.
- Glen Coppersmith, Kim Ngo, Ryan Leary, and Tony Wood. Exploratory data analysis of social media prior to a suicide attempt. In *Proceedings of the Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, San Diego, California, USA, June 2016. North American Chapter of the Association for Computational Linguistics.
- Munmun De Choudhury, Emre Kiciman, Mark Dredze, Glen Coppersmith, and Mrinal Kumar. Discovering shifts to suicidal ideation from mental health content in social media. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pages 2098–2110. ACM, 2016.
- Yoav Goldberg. A primer on neural network models for natural language processing. *Journal of Artificial Intelligence Research*, 57:345–420, 2016.
- Quoc V Le and Tomas Mikolov. Distributed representations of sentences and documents. *arXiv preprint arXiv:1405.4053*, 2014.
- Jiwei Li, Alan Ritter, and Dan Jurafsky. Learning multi-faceted representations of individuals from heterogeneous evidence using neural networks. *arXiv preprint arXiv:1510.05198*, 2015.

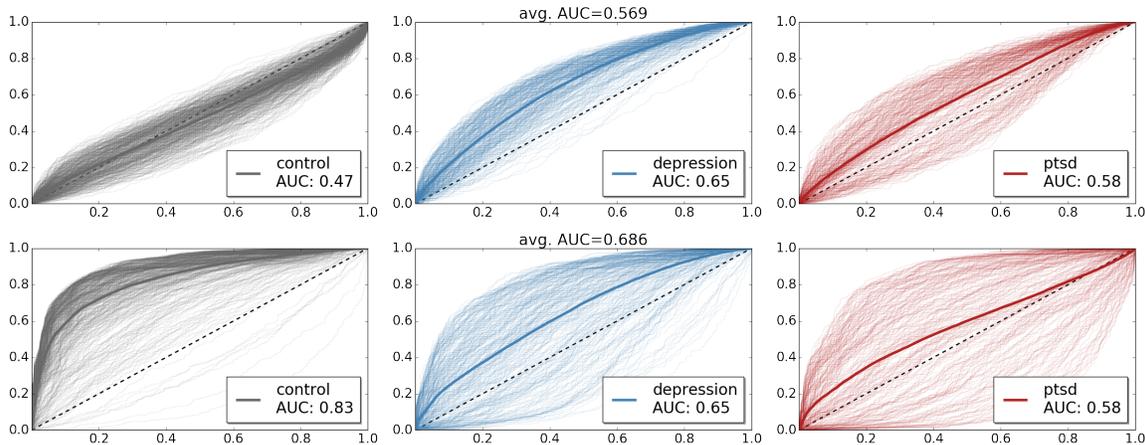
- Deirdre McCaughey, Catherine Baumgardner, Andrew Gaudes, Dominique LaRoche, Kayla Jiaxin Wu, and Tejal Raichura. Best practices in social media: Utilizing a value matrix to assess social media’s impact on health care. *Social Science Computer Review*, 32(5):575–589, 2014.
- Jude Mikal, Samantha Hurst, and Mike Conway. Ethical issues in using twitter for population-level depression monitoring: a qualitative study. *BMC medical ethics*, 17(1):1, 2016.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- Margaret Mitchell, Glen Coppersmith, and Kristy Hollingshead, editors. *Proceedings of the Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*. North American Association for Computational Linguistics, Denver, Colorado, USA, June 2015.
- Michael J Paul and Mark Dredze. You are what you tweet: Analyzing twitter for public health. *Icwsn*, 20:265–272, 2011.
- ES Paykel, A Tylee, A Wright, RG Priest, et al. The defeat depression campaign: psychiatry in the public arena. *The American journal of psychiatry*, 154(6):59, 1997.
- Ted Pedersen. Screening twitter users for depression and ptsd with lexical decision lists. In *Proceedings of the Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 46–53, 2015.
- James W Pennebaker, Martha E Francis, and Roger J Booth. Linguistic inquiry and word count: Liwc 2001. *Mahway: Lawrence Erlbaum Associates*, 71:2001, 2001.
- Daniel Preotiuc-Pietro, Maarten Sap, H Andrew Schwartz, and LH Ungar. Mental illness detection at the world well-being project for the clpsych 2015 shared task. *NAACL HLT 2015*, page 40, 2015.
- Ashwin Rajadesingan, Reza Zafarani, and Huan Liu. Sarcasm detection on twitter: A behavioral modeling approach. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, pages 97–106. ACM, 2015.
- Radim Řehůřek and Petr Sojka. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta, May 2010. ELRA. <http://is.muni.cz/publication/884893/en>.
- Philip Resnik, William Armstrong, Leonardo Claudino, Thang Nguyen, Viet-An Nguyen, and Jordan Boyd-Graber. Beyond lda: exploring supervised topic modeling for depression-related language in twitter. In *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 99–107, 2015.

- David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Cognitive modeling*, 5(3):1, 1988.
- H. Andrew Schwartz, Johannes Eichstaedt, Margaret L. Kern, Gregory Park, Maarten Sap, David Stillwell, Michal Kosinski, and Lyle Ungar. Towards assessing changes in degree of depression through Facebook. In *Proceedings of the ACL Workshop on Computational Linguistics and Clinical Psychology*, 2014.
- Aliaksei Severyn and Alessandro Moschitti. Twitter sentiment analysis with deep convolutional neural networks. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 959–962. ACM, 2015.
- Noah A Smith and Jason Eisner. Contrastive estimation: Training log-linear models on unlabeled data. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 354–362. Association for Computational Linguistics, 2005.
- Duyu Tang, Bing Qin, and Ting Liu. Learning semantic representations of users and products for document level sentiment classification. In *ACL (1)*, pages 1014–1023, 2015.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(2579-2605):85, 2008.
- Yi Yang, Ming-Wei Chang, and Jacob Eisenstein. Toward socially-infused information extraction: Embedding authors, mentions, and entities. *arXiv preprint arXiv:1609.08084*, 2016.
- Yang Yu, Xiaojun Wan, and Xinjie Zhou. User embedding for scholarly microblog recommendation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, volume 2, pages 449–453, 2016.

A. Measuring Homophily (continued)



(a) User vector similarity ranking. The first row corresponds to a 'query' user, and the columns show the top 100 most similar users, colored according to their class.

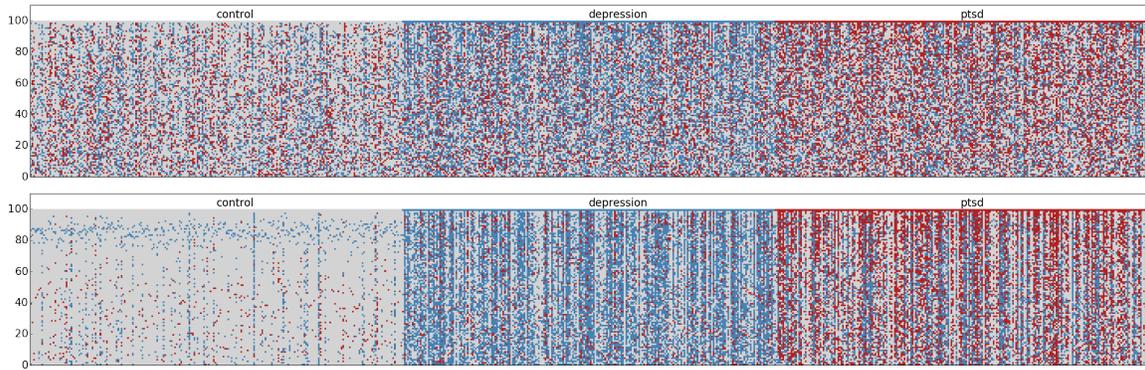


(b) ROC curves and AUC scores of the induced user similarity rankings, per class.

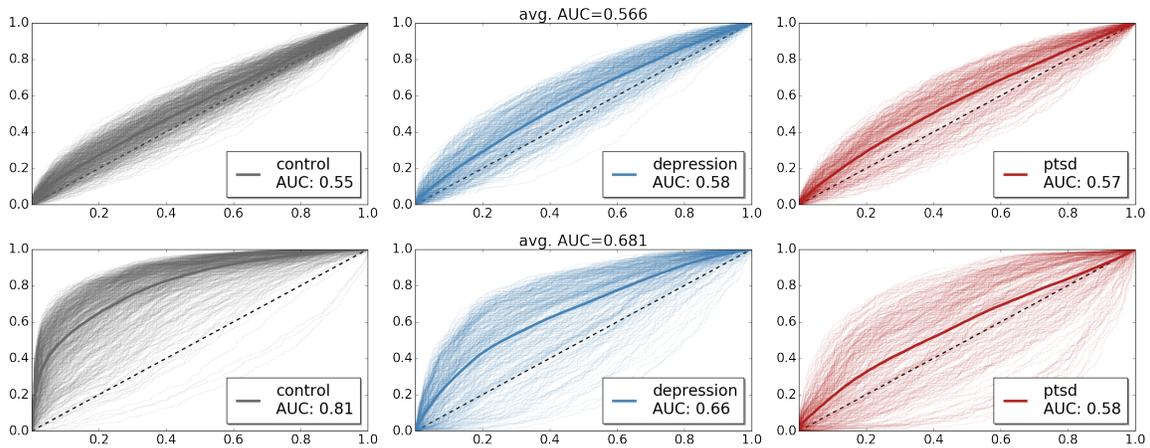
Figure 4: Measuring homophilic relations with respect to mental conditions with vector distances over the user embedding space. The top-most sub-plots refer to rankings induced with the USER2VEC model, and the ones at the bottom correspond to rankings obtained with embeddings adapted with the NLSE model.

B. Visualizing Embedding Subspace Features

The NLSE induces task-specific representations by learning a low dimensional embedding. We exploit the fact that the sigmoid non-linearity (Eq. 3) projects all the feature values into the range $[0; 1]$, to map these values into color intensities in a heatmap. We then used the same sample used to produce the plots in Section 5 (100 users per class), and plotted the representations learnt by model. The resulting plot is shown in Figure 6. We can



(a) User vector similarity ranking. The first row corresponds to a 'query' user, and the columns show the top 100 most similar users, colored according to their class.



(b) ROC curves and AUC scores of the induced user similarity rankings, per class.

Figure 5: Measuring homophilic relations with respect to mental conditions with vector distances over the user embedding space. The top-most sub-plots refer to rankings induced with the PV-DBOW model, and the ones at the bottom correspond to rankings obtained with embeddings adapted with the NLSE model.

observe that specific groups of features are mostly activated in specific classes. From this plot, we can see that the features induced with PV-DM model are sparser than the other models, which might explain why these features can better capture user similarities. On the other hand, the features induced with PV-DM and USER2VEC seem 'noisier' but also seem to capture differences between classes. In Figure 7, we show a similar plot but where we average the feature vectors of *all* the users in each class. Interestingly, it seems to be case that the 'prototypical' class vectors learned with different embeddings are very similar.

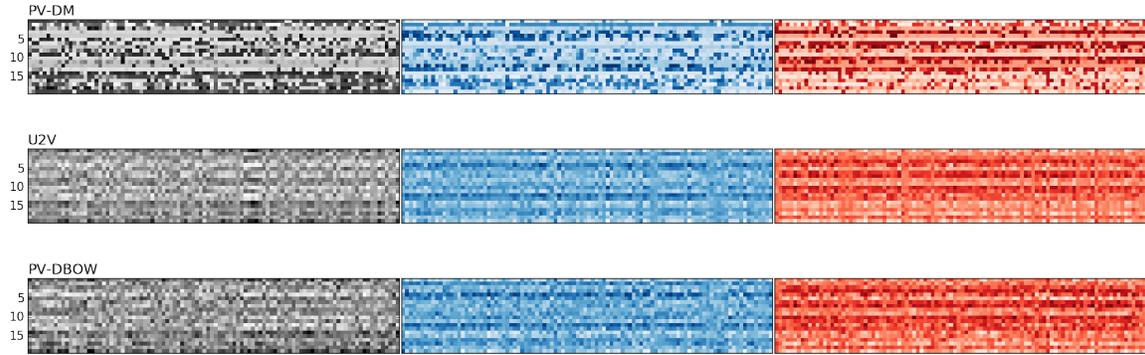


Figure 6: Embedding subspace features by class. Color intensities reflect the magnitude of the feature values

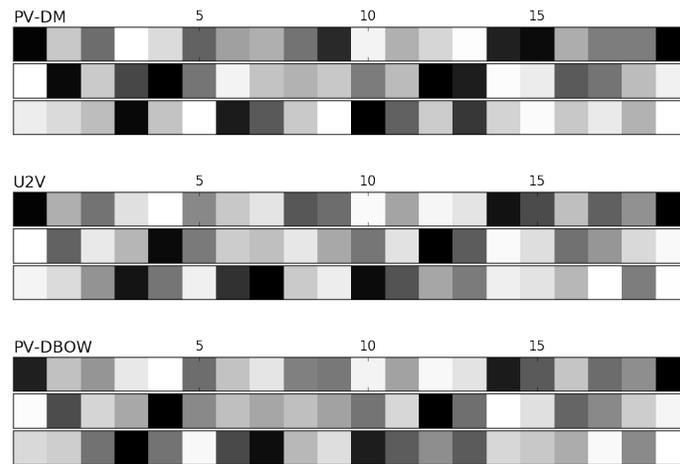


Figure 7: Embedding subspace features for a 'prototypical class vector' obtained by averaging the vectors of all the users in the class.