

FEAT2VEC: SEMANTIC MODULARITY FOR CONTINUOUS REPRESENTATIONS OF DATA WITH FEATURES

Anonymous authors

Paper under double-blind review

ABSTRACT

Methods that calculate dense embeddings for features in unstructured data—such as words in a document—have proven to be very successful for knowledge representation. Surprisingly, very little work has focused on methods for structured datasets where there is more than one feature type—this is, datasets that have arbitrary features beyond words.

We study how to estimate continuous representations for multiple feature types within a dataset, where each feature type exists in a different higher-dimensional space. Feat2Vec is a novel method that calculates embeddings for data with multiple feature types enforcing *semantic modularity* across features—all different feature types exist in a common space. We demonstrate our work on two datasets. The first one is collected from a leading educational technology firm, and thus we are able to discover a common continuous space for entities such as universities, courses, textbooks and students. The second dataset is the public IMDB dataset, and we discover embeddings for entities that include actors, movies and producers. Our experiments suggest that Feat2Vec significantly outperforms existing algorithms that do not leverage the structure of the data.

1 PRELIMINARIES

Informally, in machine learning an *embedding* of a vector $\vec{x} \in \mathbb{R}^n$ is another vector $\vec{y} \in \mathbb{R}^r$ that has much lower dimensionality ($r \ll n$) than the original representation, and can be used to replace the original vector in downstream prediction tasks. Embeddings have multiple advantages, as they enable more efficient training (Mikolov et al., 2013b), and unsupervised learning (Schnabel et al., 2015). For example, when applied to text, semantically similar words are mapped to nearby points.

Word2Vec (Mikolov et al., 2013a;b) is an extremely successful software package that contains two embedding functions with the same domain and codomain:

$$\text{Word2Vec} : \{ \vec{x} \in \mathbb{R}^n \mapsto \vec{y} \in \mathbb{R}^r \} \quad (1)$$

Word2Vec is suited for calculating embeddings for datasets that consist of documents of words with a vocabulary size of n . Here, \vec{x} is sparse because words (and categorical variables, in general) are modeled using one-hot encoding. In this paper we study how to generalize embedding functions for arbitrary datasets with multiple feature types by leveraging their internal structure. We call this approach Feat2Vec:

$$\text{Feat2Vec} : \{ \vec{x} \in \vec{\mathcal{F}} \mapsto \vec{y} \in \mathbb{R}^r \} \quad (2)$$

$$\vec{\mathcal{F}} = [\mathbb{R}^{d_1}, \mathbb{R}^{d_2}, \dots, \mathbb{R}^{d_n}] \quad (3)$$

Here, $\vec{\mathcal{F}}$ is a structured feature set—where each dimension represents a different feature type. Notice that it is possible to use Word2Vec on structured data by simply flattening the input. However, Feat2Vec leverage the features’ structure to enforce *semantic modularity* across features. That is, the embedding of an observation is additive in respect to the embeddings of its feature types; thus, each feature type \mathcal{F}_i must be projected into the same embedding space. Semantic modularity is useful because it enables making principled comparisons between different feature types. For example, on an educational dataset that contains many feature types, including students and courses, semantic modularity enables inferring course preferences for students.

2 FEAT2VEC

The rest of this section describes the three parts of a Feat2Vec implementation. § 2.1 describes a model that infers an embedding from an observation. To learn this model we need positive and negative examples: a positive example is “semantic” and it is observed during training, but negative examples are not. This is similar how a positive example in Word2Vec are grammatical co-occurrences, and negative examples are generated. § 2.2 describes our novel sampling strategy. Finally, in § 2.3 we describe how to learn a Feat2Vec model from data.

2.1 STRUCTURED DEEP-IN FACTORIZATION MACHINE

Levy & Goldberg (2014) showed that a Word2Vec model can be formalized as a Shallow Factorization Machine (Rendle, 2010) with two features types: a word and its context. This factorization model is a binary classifier that scores the likelihood of an observation $\vec{x} \in \mathbb{R}^n$ being labeled $y \in \{0, 1\}$, as proportional to the sum of the factorized pairwise interactions:

$$p(Y = 1|\vec{x}) = \sigma \left(\sum_{i=1}^n \sum_{j=i}^n \langle \vec{\beta}_i x_i, \vec{\beta}_j x_j \rangle \right) \quad (4)$$

Here, Y is the random variable that defines whether the observation is “semantic” (which in practice means whether it occurs in the training data), $\vec{\beta}_i$ is a rank- r vector of factors, and σ is a sigmoid:

$$\sigma(x) = \log \left(\frac{\exp(x)}{\exp(x) + 1} \right) \quad (5)$$

However, features may have some structure that is known beforehand. For example, consider multiple discrete entity types, like universities, students, and books. To model this data, a Shallow Factorization Machine would ignore the structure and simply concatenate the one-hot encoding representation of each entity into a single feature vector. Feat2Vec relies on an novel extension to factorization called Structured Deep-In Factorization Machine, which we describe in detail in a [companion paper](#) (Anonymized, Under review). While Shallow Factorization Machine learns an embedding per feature, the Structured Deep-In model allows greater flexibility. For example, consider using images or text on these factorization models. The shallow model learns an embedding for each word, or an embedding per pixel. The structured model enables higher-level of abstraction and flexibility, and it can learn an embedding per passage of text, or an embedding per image.

Structured Deep-In Factorization Machine inputs $\vec{\kappa}$ as a vector of vectors that define the structure of the groups of features. Each entry i is a vector of size d_i and it is used to represent the feature types of \mathcal{F}_i . In a shallow model, a feature interacts with all other features, but in the structured model they only interact with features from different groups. For example, consider a model with four features. If we define $\vec{\kappa} = \begin{bmatrix} [1, 2], [3, 4] \end{bmatrix}$, the structure of the model would not allow x_1 to interact with x_2 , or x_3 to interact with x_4 . Additionally, for each group of features, the model allows applying a $d_i \times r$ feature extraction function ϕ_i . More formally, this is:

$$p(Y = 1|\vec{x}; \vec{\phi}, \vec{\kappa}) = \sigma \left(s(\vec{x}, \vec{\phi}) \right) = \sigma \left(\sum_{i=1}^{|\vec{\kappa}|} \sum_{j=i}^{|\vec{\kappa}|} \langle \phi_i(\vec{x}_{\vec{\kappa}_i}), \phi_j(\vec{x}_{\vec{\kappa}_j}) \rangle \right) \quad (6)$$

In this notation, $\vec{x}_{\vec{\kappa}_i}$ is a subvector that contains all of the features that belong to the group κ_i , and $s(\vec{x}, \vec{\phi})$ is the score function. Thus, $\vec{x}_{\vec{\kappa}_i} = [x_{i_a} | a \in \{1, 2, \dots, |\kappa_i|\}]$. The simplest implementation for ϕ_i is a linear fully-connected layer, where the output of the r -th entry is:

$$\phi_i(\vec{x}_i; \vec{\beta})_r = \sum_{a=1}^{d_i} \beta_{r_a} x_{i_a}$$

We build the Feat2Vec embedding for an observation \vec{x} as:

$$\text{Feat2Vec}(x) = [\phi_1(\vec{x}_{\vec{\kappa}_1}), \phi_2(\vec{x}_{\vec{\kappa}_2}), \dots, \phi_n(\vec{x}_{\vec{\kappa}_n})] \quad (7)$$

2.2 SAMPLING

The training dataset for a Feat2Vec model consists of only semantic observations. In natural language, these would be documents written by humans. Since Structured Deep-In Factorization Machine (Equation 6) requires positive and negative examples, we also need to supply observations that are not semantic. Consider a feature type $\mathcal{F}_i \in \mathbb{R}^{d_i}$, that exists in very high dimensional space (d_i is large). For example, this could happen because we are modeling with one-hot encoding a categorical variable with large number of possible values. In such scenario, it is overwhelmingly costly to feed the model all negative labels, particularly if the model is fairly sparse.

A shortcut around this is a concept known as *implicit sampling*, where instead of using all of the possible negative labels, one simply samples a fixed number (k) from the set of possible negative labels for each positively labelled record. Word2Vec makes use of an algorithm called Negative Sampling (Dyer, 2014). In short, their approach samples a negative observation w from a noise distribution q :

$$\text{Draw } w \sim p(W|Y = 0, X) \quad (8)$$

$$= q(\{w_1, w_2, \dots, w_n\}, \alpha) \quad (9)$$

Here q is the empirical frequency of w in the dataset, exponentiated by α , a flattening hyper-parameter:

$$q(W, \alpha) = \frac{c_X(w)^\alpha}{\sum_{w' \in W} c_X(w')^\alpha} \quad (10)$$

Here, $c_X(w)$ is the number of times a feature w appeared in the training dataset X . With $\alpha = 1$, the noise distribution corresponds to the empirical frequency distribution, and with $\alpha = 0$ to the uniform distribution. Intermediate values are intended to be continuous progressions to these two extremes.

Our main contribution is introducing a new implicit sampling method that enables learning embedding for structured feature sets. We can learn the correlation of features within a dataset by imputing negative labels, simply by generating “unobserved” records as our negative sample. Unlike Word2Vec, we do not constraint features types to be words. Features types can be individual numeric columns, but they need not to be. By defining feature types as groups of subfeatures (using κ parameter in Equation 6), the model can reason on more complex entities. For example, in our experiments on a movie dataset, we use a “genre” feature type, where we group non-mutually exclusive indicators for comedy, action, and drama films under one umbrella. For more involved feature types (e.g., an image), one would need to define a function ϕ that builds intermediate layers to map the entity to an embedding.

We start with a dataset S^+ of records with p feature types. We then mark all observed records in the training set as positive examples. For each positive record, we generate k negative labels using the following 2-step algorithm:

Algorithm 1 Implicit sampling algorithm for Feat2Vec

```

1: function FEAT2VEC_SAMPLE( $S^+, k, \alpha_1, \alpha_2$ )
2:    $S^- \leftarrow \emptyset$ 
3:   for  $s \in S^+$  do
4:     Draw a random feature type  $\kappa_i \sim q_1(\{d_i\}_{i=1}^p, \alpha_1)$ 
5:     for  $j \in \{1, \dots, k\}$  do
6:        $\vec{e} \leftarrow \vec{x}^{(s)}$  ▷ set initially to be equal to a positive sample
7:       Draw a random (sub)feature  $\vec{r} \sim q_2(X_{\kappa_i}, \alpha_2)$ 
8:        $\vec{e}_{\kappa_i} \leftarrow \vec{r}$  ▷ substitute the  $i$ -th feature type with a sampled one
9:        $S^- \leftarrow S^- + \{\vec{e}\}$ 
10:    end for
11:  end for
12:  return  $S^-$ 
13: end function

```

Here $q_2(X_{\kappa_i}, \alpha_2)$ is the “flattened” empirical multinomial distribution over the frequency of the feature values of feature κ_i , as in the negative sampling procedure used in Word2Vec :

$$x_j \sim q_2(X_{\kappa_i}, \alpha_2), \text{ where } P_{q_2}(x_j = x | X_{\kappa_i}, \alpha_2) = \frac{c_X(x)^{\alpha_2}}{\sum_{x' \in X_{\kappa_i}} c_X(x')^{\alpha_2}}, \alpha \in [0, 1]$$

Where $P_q(x)$ denotes the probability of x under a distribution q . Analogously, $q_1(\{d_i\}_{i=1}^p, \alpha_1)$ draws a feature to negatively sample from a multinomial distribution with probabilities proportional to a feature’s complexity, d_i . Here, by complexity, we mean the number of parameters the algorithm learns associated with a feature group:

$$P_{q_1}(\kappa_i | \{d_i\}_{i=1}^p, \alpha_1) = \frac{d_i^{\alpha_1}}{\sum_{k=1}^p d_k^{\alpha_1}}$$

d_i may be the dimensionality of a feature, but might also capture other parameters used for learning embeddings, where applicable. For example, if we have multiple intermediate layers to process a set of pixels embeddings into an “image embedding” (such as deep convolutional layers), these parameters would count towards d_i .

Explained in words, our negative sampling method is to randomly select one of the features κ from the noise distribution $q_1(\cdot)$, and, conditional on the feature to be resampled, choose another value x_j for feature κ from a noise distribution $q_2(\cdot)$. In our application, we use the same class of noise distributions (flattened multinomial) for both levels of sampling, but this need not necessarily be the case.

This method will sometimes by chance generate negatively labeled samples that *do* exist in our sample of observed records. For this reason, among others, we employ a technique known as noise contrastive estimation (Gutmann & Hyvärinen, 2010) which uses basic probability laws to adjust the structural statistical model $p(Y = 1 | \vec{\phi}, \vec{x})$ to account for the possibility of random negative labels that appear identical to positively labeled data. An additional burden of this method, however, is we need to learn additional nuisance parameters $Z_{\vec{x}}$ for each unique record type \vec{x} that transform the score function $s(\cdot)$ into a well-behaved probability distribution that integrates to 1. This introduces an astronomical amount of new parameters and greatly increase the complexity of the model. Instead of estimating these, we appeal to the work of (Mnih & Teh, 2012), who show in the context of language models that setting the $Z_{\vec{x}} = 1$ in advance effectively does not change the performance of the model.¹ Written explicitly, the new structural probability model is:

$$\hat{p}(Y = 1 | \vec{\phi}, \vec{x}_i) = \frac{\exp(s(\vec{x}_i, \vec{\phi}))}{\exp(s(\vec{x}_i, \vec{\phi})) + P_q(\vec{x}_i | \alpha_1, \alpha_2)}$$

where $P_q(\cdot)$ denotes the unconditional probability of a record \vec{x}_i being drawn from our negative sampling algorithm, and $s(x, \vec{\phi})$ is the score function for Factorization Machines, in the syntax of (Gutmann & Hyvärinen, 2010). We can show, with relatively little effort, that our 2-step sampler has some interesting theoretical properties. Namely, the embeddings learned under this model will be equivalent to a convex combination of the embeddings learned from p individual Factorization Machines.

2.3 LEARNING FROM DATA

We optimize Noise Contrastive Estimation (Gutmann & Hyvärinen, 2010) loss, which uses basic probability laws to adjust the structural statistical model $p(Y = 1 | \vec{\phi}, \vec{x})$ to account for the possibility of randomly generated negative labels that appear identical to positively labeled data. Written explicitly, the new structural probability model is just:

$$\hat{p}(Y = 1 | \vec{\phi}, \vec{x}_i) = \frac{\exp(s(\vec{x}_i, \vec{\phi}))}{\exp(s(\vec{x}_i, \vec{\phi})) + P_q(\vec{x}_i | \alpha_1, \alpha_2)} \times Z_{\vec{x}} \quad (11)$$

¹(Mnih & Teh, 2012) suggests one can set the normalizing constant to 1 because the neural net model has enough free parameters that it will effectively just learn the probabilities itself so that it does not over/under predict the probabilities on average (since that will result in penalties on the loss function). Note that while we follow the same procedure, one could set $Z_{\vec{x}}$ to any positive value in advance and the same logic would follow, if one was worried about astronomically low probabilities.

where $P_q(\cdot)$ denotes the unconditional probability of a record \vec{x}_i being drawn from our implicit sampling algorithm, $s(x, \vec{\phi})$ is the score function for Factorization Machines, in the syntax of (Gutmann & Hyvärinen, 2010), and $Z_{\vec{x}_i}$ is a positive learned parameter to guarantee that the probability distribution integrates to 1:

$$\sum_{\vec{x}_i \in X} Z_{\vec{x}_i} \times \frac{\exp(s(\vec{x}_i, \vec{\phi}))}{\exp(s(\vec{x}_i, \vec{\phi})) + P_q(\vec{x}_i | \alpha_1, \alpha_2)} = 1$$

This additional parameter for each unique record, $Z_{\vec{x}_i}$, introduces an astronomical amount of new parameters and greatly increase the complexity of the model. Instead of estimating these, we appeal to the work of (Mnih & Teh, 2012), who show in the context of language models that setting the $Z_{\vec{x}} = 1$ in advance effectively does not change the performance of the model.²

We can show, with relatively little effort, that our 2-step sampler has some interesting theoretical properties. Namely, the embeddings learned under this model will be equivalent to a convex combination of the embeddings learned from p individual Factorization Machines that also use an NCE loss:

Theorem 1. *The embeddings learned with Feat2Vec are a convex combination of the embeddings learned from p targeted Factorization Machines for each feature in the data.*

See the proof in the Appendix.

3 LIMITATIONS

Continuous features require a sensible feature function, and future work could focus on them.

Feature engineering (grouping) is important.

4 EMPIRICAL RESULTS

5 RELATION TO PRIOR WORK

6 CONCLUSION

A APPENDIXES

A.1 PROOF TO THEOREM 1

Theorem 1. *The embeddings learned with Feat2Vec are a convex combination of the embeddings learned from p targeted Factorization Machines for each feature in the data.*

Proof. Let S_f^+ denote the positively labeled records whose corresponding negative samples resample feature f . We can express the loss function $L(\cdot)$, the binary cross-entropy of the data given the Feat2Vec model, as follows:

²(Mnih & Teh, 2012) suggests one can set the normalizing constant to 1 because the neural net model has enough free parameters that it will effectively just learn the probabilities itself so that it does not over/under predict the probabilities on average (since that will result in penalties on the loss function). Note that while we follow the same procedure, one could set $Z_{\vec{x}}$ to any positive value in advance and the same logic would follow, if one was worried about astronomically low probabilities.

$$\begin{aligned}
L(\vec{x}|\vec{\phi}) &= \frac{1}{|S^+|} \sum_{i \in S^+} \left(\log(\tilde{p}(Y=1|\vec{\phi}, \vec{x}_i)) + \sum_{\vec{x}_j \sim q(\cdot|\vec{x}_i)}^k \log(\tilde{p}(Y=0|\vec{\phi}, \vec{x}_j)) \right) \\
&= \frac{1}{|S^+|} \sum_{i \in S^+} \left(\log(\tilde{p}(Y=1|\vec{\phi}, \vec{x}_i, i \in S_f)p(i \in S_f^+)) + \sum_{\vec{x}_j \sim q(\cdot|\vec{x}_i)}^k \log(\tilde{p}(Y=0|\vec{\phi}, \vec{x}_j, i \in S_f)p(i \in S_f^+)) \right) \\
&= \frac{1}{|S^+|} \sum_{f=1}^p \sum_{i \in S_f^+} \left(\log\left(\frac{e^{s(\vec{x}_i, \vec{\phi})}p(i \in S_f^+)}{e^{s(\vec{x}_i, \vec{\phi})} + P_q(\vec{x}_i|\vec{x}_i, i \in S_f^+)}\right) + \sum_{\vec{x}_j \sim q(\cdot|\vec{x}_i, i \in S_f^+)}^k \log\left(\frac{P_q(\vec{x}_j|\vec{x}_i, i \in S_f)p(i \in S_f)}{e^{s(\vec{x}_j, \vec{\phi})} + P_q(\vec{x}_j|\vec{x}_i, i \in S_f^+)}\right) \right)
\end{aligned}$$

Note now that $P_q(\vec{x}_k|\vec{x}_i, i \in S_f^+)$ is simply the probability of the record's feature value $\vec{x}_{k,f}$ under the second step noise distribution $q_2(\mathbf{X}_f, \alpha_2)$: $P_q(\vec{x}_k|\vec{x}_i, i \in S_f^+) = P_{q_2}(\vec{x}_{k,f})$, where k refers to record i or a record j negatively sampled from i .

$$\begin{aligned}
&= \frac{1}{|S^+|} \sum_{f=1}^p \sum_{i \in S_f^+} \left(\log\left(\frac{e^{s(\vec{x}_i, \vec{\phi})}p(i \in S_f^+)}{e^{s(\vec{x}_i, \vec{\phi})} + P_{q_2}(\vec{x}_{i,f})}\right) + \sum_{\vec{x}_j \sim q(\cdot|\vec{x}_i, i \in S_f)}^k \log\left(\frac{P_{q_2}(\vec{x}_{j,f})p(i \in S_f)}{e^{s(\vec{x}_j, \vec{\phi})} + P_{q_2}(\vec{x}_{j,f})}\right) \right) \\
&= \frac{1}{|S^+|} \sum_{f=1}^p \sum_{i \in S_f^+} \left(\log\left(\frac{e^{s(\vec{x}_i, \vec{\phi})}}{e^{s(\vec{x}_i, \vec{\phi})} + P_{q_2}(\vec{x}_{i,f})}\right) + \log(p(i \in S_f^+)^{k+1}) \right. \\
&\quad \left. + \sum_{\vec{x}_j \sim q(\cdot|\vec{x}_i, i \in S_f^+)}^k \log\left(\frac{P_{q_2}(\vec{x}_{j,f})}{e^{s(\vec{x}_j, \vec{\phi})} + P_{q_2}(\vec{x}_{j,f})}\right) \right)
\end{aligned}$$

We now drop the term containing the probability of assignment to a feature type $p(i \in S_f^+)$ since it is outside of the learned model parameters $\vec{\phi}$ and fixed in advance:

$$\begin{aligned}
&\propto \frac{1}{|S^+|} \sum_{f=1}^p \sum_{i \in S_f^+} \left(\log\left(\frac{e^{s(\vec{x}_i, \vec{\phi})}}{e^{s(\vec{x}_i, \vec{\phi})} + P_{q_2}(\vec{x}_{i,f})}\right) + \sum_{\vec{x}_j \sim q(\cdot|\vec{x}_i, i \in S_f^+)}^k \log\left(\frac{P_{q_2}(\vec{x}_{j,f})}{e^{s(\vec{x}_j, \vec{\phi})} + P_{q_2}(\vec{x}_{j,f})}\right) \right) \\
&\xrightarrow{|S^+| \rightarrow \infty} \sum_{f=1}^p p(i \in S_f^+) E \left[\log\left(\frac{e^{s(\vec{x}_i, \vec{\phi})}}{e^{s(\vec{x}_i, \vec{\phi})} + P_{q_2}(\vec{x}_{i,f})}\right) + \sum_{\vec{x}_j \sim q(\cdot|\vec{x}_i, i \in S_f^+)}^k \log\left(\frac{P_{q_2}(\vec{x}_{j,f})}{e^{s(\vec{x}_j, \vec{\phi})} + P_{q_2}(\vec{x}_{j,f})}\right) \right] \\
&= \sum_{f=1}^p p(i \in S_f^+) E \left[L(\vec{x}|\vec{\phi}, \text{target} = f) \right]
\end{aligned}$$

Thus, the loss function is just a convex combination of the loss functions of the targeted classifiers for each of the p features, and by extension so is the gradient since:

$$\frac{\partial}{\partial \phi} \sum_{f=1}^p p(i \in S_f^+) E \left[L(\vec{x}|\vec{\phi}, \text{target} = f) \right] = \sum_{f=1}^p p(i \in S_f^+) \frac{\partial}{\partial \phi} E \left[L(\vec{x}|\vec{\phi}, \text{target} = f) \right]$$

Thus the algorithm will, at each step, learn a convex combination of the gradient for a targeted classifier on feature f , with weights proportional to the feature type sampling probabilities in step 1 of the sampling algorithm. \square

A.2 PROBLEMS WITH FLATTENING MULTINOMIAL OF Word2Vec

Write

REFERENCES

- Anonymized. Structured deep factorization machine: Towards general-purpose architectures. In *6th International Conference on Learning Representations*, Under review.
- Chris Dyer. Notes on noise contrastive estimation and negative sampling. *arXiv preprint arXiv:1410.8251*, 2014.
- Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 297–304, 2010.
- Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. In *Advances in neural information processing systems*, pp. 2177–2185, 2014.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013a.
- 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*, pp. 3111–3119, 2013b.
- Andriy Mnih and Yee Whye Teh. A fast and simple algorithm for training neural probabilistic language models. In *In Proceedings of the International Conference on Machine Learning*, 2012.
- Steffen Rendle. Factorization machines. In *Data Mining (ICDM), 2010 IEEE 10th International Conference on*, pp. 995–1000. IEEE, 2010.
- Tobias Schnabel, Igor Labutov, David M Mimno, and Thorsten Joachims. Evaluation methods for unsupervised word embeddings. In *EMNLP*, pp. 298–307, 2015.