

# Efficient learning of email similarities for customer support

Jelle Bakker and Kerstin Bunte

University of Groningen - Faculty of Science and Engineering  
P.O. Box 407, 9700 AK Groningen - The Netherlands

**Abstract.** One way to increase customer satisfaction is efficient and consistent customer email support. In this contribution we investigate the use of dimensionality reduction, metric learning and classification methods to predict answer templates that can be used by an employee or retrieve historic conversations with potential suitable answers given an email query. The strategies are tested on email data and the publicly available Reuters data. We conclude that prototype-based metric learning is fast to train and the parameters provide a compressed representation of the database enabling efficient content based retrieval. Furthermore, learning customer email embeddings based on the similarity of employee answers is a promising direction for computer aided customer support.

## 1 Introduction

One way to increase customer satisfaction is to provide very good customer support: answering emails quick and helpful for the specific problem. However, answering emails efficiently can be difficult dependent on the experience of employees and difficulty of the customers question. We collaborate with a Dutch e-commerce company which spends approximately 40 FTEs on answering emails from customers about order information, product details, warranty and prices. Historic emails and example responses which have been used in similar context could serve as default response that can be sent immediately or with little adaptation. An efficient computer aided system supporting the staff member to answer emails more quickly and consistently, based on similar historic correspondences, is therefore highly desirable. For frequently asked questions (FAQs) the company creates templates, such that a response can be sent quickly only needing small changes. Therefore, one possibility is to formulate the problem as text classification, aiming in training a system with email questions which were answered using a set of templates. A variety of text and email classification methods have been proposed over the years, ranging from kernel methods, such as Support Vector Machines (SVMs) [1], which have been shown to perform better than Naive Bayes, Decision Trees and Rocchio classification [2, 3] and more recently deep Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) [4]. The deep learning models outperform the older methods both in text representation and classification accuracy, but are highly complex, difficult to train and require large amounts of labelled training data. A general introduction to machine learning for text can be found in [5].

In this contribution we propose to use a prototype based method called Learning Vector Quantization (LVQ) and its metric learning extensions [6] for automatic email classification and content based retrieval of email texts. This choice offers several advantages compared to alternative techniques: it is fast to train, intuitive to understand and the model parameters are interpretable providing a compressed representation of the labelled data. The learned prototypes

and dissimilarity measures have been proven useful in versatile applications, such as dimensionality reduction, visualization and content based retrieval [7, 8].

Furthermore, we propose using t-distributed Stochastic Neighbor Embedding (t-SNE) [9] training an explicit mapping [10] based on the similarity of answers.

## 2 Data and Preprocessing

We compare the metric learning methods on the well-known publicly available Reuters-21578 [11] text dataset to show its applicability. The ModApte version from the Reuters newswire corpus contains 8,293 preprocessed documents in 65 categories split into 5,946 training and 2,347 testing documents. We compare our results to the results obtained by [12], SVM and Rocchio. Furthermore, we investigate a dataset of 433,170 emails spanning approximately a year and a half from mid 2016 to the begin of 2018 of customer support of a Dutch e-commerce company. The current system enables us to search through large sets of emails quickly and identify threads of customer and support staff correspondences. We extract features using the Term Frequency-Inverse Document Frequency (TF-IDF) [13], which is a well-known and established statistic to determine the relevance of words. In order to capture relations and similarities between words and reduce the dimensionality we use Latent Semantic Analysis (LSA) [14, 15], based on the term-document co-occurrence matrix.

We refer to the last email sent by an employee as *answer* and the emails sent by the client within the same thread as *question*. A reasonable aim is to predict the topic of the thread such that incoming emails are classified by the response that should follow and provide answers of similar threads from the database dealing with the same topic. Questions can be from any customer and are therefore extremely diverse, even when concerning the same topic. However, the number of employees is much smaller and a topic is answered more consistently, especially for FAQs for which templates are used. Therefore, we can infer a reasonable topic class label, by grouping questions which were most likely answered using the same template, and use it to train a classifier to identify questions dealing with similar topics. Only if the probability based on co-occurring words of answers and a template  $T_j$  is above 0.75 and the next largest probability to belong to another template  $T_k$  is at least 0.15 points smaller we consider  $T_j$  a match. We prefer an accurate labeling of the data since we will use them as “ground truth” for evaluation. The total number of answers is 116,128 and 2,989 answers could be labeled by using these criteria. This suggests that only a small portion is currently based on templates and not all employees are using them. A computer aided system is likely to improve consistency in answering.

## 3 Methods

### 3.1 Generalized Matrix Learning Vector Quantization

LVQ is a prototype based classification method, which is easy to implement, intuitive to understand and has model parameters that are interpretable. It exhibits a runtime complexity of  $\mathcal{O}(nk)$ , with  $n$  the number of training samples and  $k$  the number of prototypes. Metric learning extensions, such as Generalized Matrix LVQ (GMLVQ) and Localized GMLVQ (LGMLVQ) [6], train adaptive distances to discriminate the classes. These methods exhibit excellent classification performance, often comparable to SVMs while faster to train. Assume a set of  $n$  training vectors, accompanied by a label  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and a set of  $k$  labeled

prototypes  $\{(\mathbf{w}_j, c(\mathbf{w}_j))\}_{j=1}^k$  with  $\{\mathbf{x}_i, \mathbf{w}_j\} \in \mathbb{R}^N$  and  $\{y_i, c(\mathbf{w}_j)\} \in \{1, \dots, C\}$  classes. The parameterized distance measure  $d^\Lambda(\mathbf{x}, \mathbf{w}) = (\mathbf{x} - \mathbf{w})^\top \Lambda (\mathbf{x} - \mathbf{w})$  and prototypes  $\mathbf{w}$  are trained by minimizing the following cost function:

$$E_{\text{GMLVQ}} = \sum_{i=1}^n \Phi \left( \frac{d_J^\Lambda - d_K^\Lambda}{d_J^\Lambda + d_K^\Lambda} \right),$$

where  $\Phi$  is a monotonically increasing function. The quantities  $d_L^\Lambda = d^\Lambda(\mathbf{x}_i, \mathbf{w}_L)$  with  $L \in \{J, K\}$  denote the parameterized distance of the  $i^{\text{th}}$  sample to the closest prototype  $\mathbf{w}_J$  with the same class label  $c(\mathbf{w}_J) = y_i$  and the closest prototype with a different label  $c(\mathbf{w}_K) \neq y_i$  respectively. The positive semi-definite metric tensor  $\Lambda \in \mathbb{R}^{N \times N}$  can be substituted by  $\Lambda = \Omega^\top \Omega$  with  $\Omega \in \mathbb{R}^{M \times N}$ . If  $M < N$  the rank of  $\Lambda$  is reduced limiting the number of free parameters and the resulting distance measure corresponds to the squared Euclidean distance in a linearly transformed space of lower dimension  $M$ :  $d^\Lambda = [\Omega(\mathbf{x} - \mathbf{w})]^2$ . The set of prototypes  $\mathbf{w}$  and matrix  $\Omega$  can be optimized using for example gradient methods. The runtime complexity for GMLVQ is therefore  $\mathcal{O}(nkNM)$ . While GMLVQ produces piecewise linear decision boundaries, the localized variant LGMLVQ assumes localized dissimilarities  $d^\Lambda(\mathbf{x}, \mathbf{w}_j) = (\mathbf{x} - \mathbf{w}_j)^\top \Lambda_j (\mathbf{x} - \mathbf{w}_j)$ , resulting in more complex non-linear boundaries. The learned prototypes and similarities can not only be used for classification, but also for supervised dimensionality reduction and content based retrieval [7, 8].

### 3.2 Linear t-SNE mapping

Since many emails in the data base do not follow a template it would be beneficial if we could use the similarities between employee answers to map the corresponding customer questions closer together. This could improve the classification and/or be used for retrieval. Embedding techniques such as t-SNE became very popular for visualization and non-linear dimensionality reduction. It uses the similarity relationship between high-dimensional data vectors to embed low-dimensional representatives, aiming to preserve the original data neighborhoods as much as possible. The original formulation of t-SNE embeds the data implicitly, so it does not provide an explicit mapping function, which makes out-of-sample extension tedious. In [10] a general framework for parameterized dimensionality reduction was proposed. This includes extensions of t-SNE allowing to learn linear and non-linear explicit embedding functions, which can readily be applied to embed new data. With respect to our email-answer pair data we would like to embed the email questions such that their neighborhood relation resembles the neighborhood relation of the corresponding answers as depicted in Fig. 1. We minimize the cost function:

$$E_{t^*} = \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad \text{with} \quad \begin{cases} p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n} \quad \text{and} \quad p_{j|i} = \frac{\exp\left(\frac{-d^v(\mathbf{v}_i, \mathbf{v}_j)}{2\sigma_i}\right)}{\sum_{k \neq i} \exp\left(\frac{-d^v(\mathbf{v}_i, \mathbf{v}_k)}{2\sigma_i}\right)} \\ q_{ij} = \frac{(1 + \|\boldsymbol{\xi}_i - \boldsymbol{\xi}_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|\boldsymbol{\xi}_k - \boldsymbol{\xi}_i\|^2)^{-1}} \end{cases}$$

whereas  $d^v(\mathbf{v}_i, \mathbf{v}_j)$  denote the dissimilarities of answer vectors,  $\sigma_i$  corresponds to the effective number of neighbors found by line search using a hyper-parameter called perplexity [9] and  $\boldsymbol{\xi}$  represents the mapped question emails. Since we want

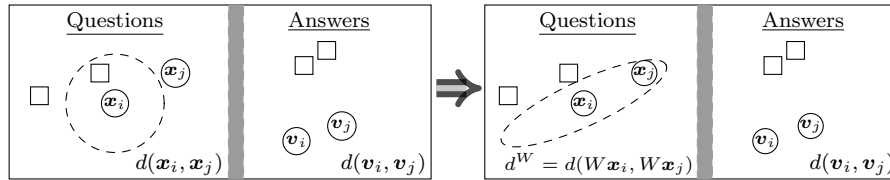


Fig. 1: Mapping email questions  $\mathbf{x}_i$  and  $\mathbf{x}_j$  with similar answers  $\mathbf{v}_i$  and  $\mathbf{v}_j$  closer to each other using a parameterized mapping function.

to map new emails for which no answer yet exists, we define a mapping  $f_W$ , which in the simplest case can be chosen as  $f_W : \mathbf{x}_i \rightarrow \xi_i = f_W(\mathbf{x}_i) = W \cdot \mathbf{x}_i$  and found using gradient methods. An example showing the effect of the trained  $W \in \mathbb{R}^{2 \times N}$  embedding on an LSA preprocessed email question-answer set is shown in Fig. 2. Note, the samples shown were not used to train this mapping.

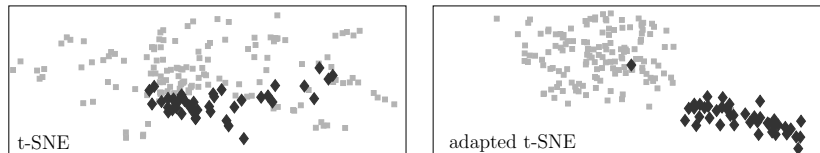


Fig. 2: Visualization of test email questions of two classes using t-SNE, based on the distances between question vectors (left) and based on the adapted mapping trained using a questions and answers training set (right).

## 4 Experiments

*Email data:* Experimenting with preprocessing and Rocchio classification on the labeled email question data we conclude that TF-IDF scores with normalized binary counts is promising. To measure the quality of the projections we use the Local Continuity Meta-Criterion (LCMC) [16] extracted from the neighborhood co-ranking matrix [17]. The criterion indicates that the LSA exhibits better preservation scores than PCA while aiming to preserve  $K$ -ary neighborhoods with  $K = 20$  and reducing the dimensions between 100 and 800. Also in successive Rocchio classification we observe consistently better accuracies using LSA, as measured in 10-fold-cross-validation compared on the test set. Therefore we project the data to 200 dimensions using LSA to compare the different classification methods on the email data. Since we aim for a computer aided system that retrieves answers to similar questions from the database, we evaluate the percentage of emails for which the right class label is among the  $m$  closest prototypes. Since each class has a single prototype, these prototypes are of  $m$  distinct classes. The results are shown in Table 1. Especially the localized LGMLVQ shows performance either superior or comparable to SVM while exhibiting linear training complexity with respect to the number of points.

*Reuters data:* For the Reuters dataset we again compute the TF-IDF scores and project using LSA, resulting in a 200-dimensional vector representation. The classification results among the first  $m$  matches are summarized in Table 1. Since the documents are derived from news stories they use more formal language and contain less errors compared to the real world email data, which explains

Table 1: Percentage of correct label retrieved among the first  $m$  matches.

$m$	Rocchio		SVM		GMLVQ		LGMLVQ		
	train	test	train	test	train	test	train	test	
Emails	1	71.76	69.77 ± 2.43	<b>93.94</b>	77.45 ± 2.13	81.20	75.27 ± 7.67	92.29	<b>78.52 ± 3.53</b>
	3	96.21	94.16 ± 1.23	<b>100.0</b>	<b>95.87 ± 1.18</b>	95.12	94.06 ± 2.59	99.26	94.66 ± 1.33
	5	98.47	97.42 ± 0.82	<b>100.0</b>	<b>98.09 ± 0.76</b>	97.13	96.44 ± 1.31	99.26	97.35 ± 0.78
Reuters	1	79.10	82.62	<b>99.60</b>	87.81	90.26	87.18	97.81	<b>88.45</b>
	3	92.92	90.29	<b>100.0</b>	93.61	95.44	92.37	99.06	<b>94.03</b>
	5	95.63	92.07	<b>100.0</b>	95.23	96.55	94.12	99.13	<b>95.57</b>

the increase in performance. Note that LGMLVQ shows again comparable or superior performance on the test set, making it a good choice for text data.

*Email-answer embedding:* In a final experiment we compare the LSA and both t-SNE embeddings using perplexity 15: first based on the email questions only and second the adapted version learning a linear mapping using the answers. Instead of measuring the mismatch between high-dimensional and low-dimensional neighborhoods using the LCMC criterion, we compute the mismatch between question neighborhoods and their respective answers in LCMC (see Fig 3 left). This way we quantify how much closer similarly answered email questions are using the trained linear t-SNE embedding. We compare the original 200-dimensional embedding of LSA used as preprocessing in the classification experiments before, as well as LSA, t-SNE and adapted t-SNE embeddings to 10 dimensions. The original high-dimensional data shows the lowest question-answer neighborhood overlap, followed by the t-SNE question embedding, improved using LSA to 10 dimensions and the best overlap is found using the adapted t-SNE. The average KNN classification error on the embedded data for different values of  $K$  is shown on the right side of Fig 3. We conclude that the email answer pair trained embedding is beneficial for the classification of email data.

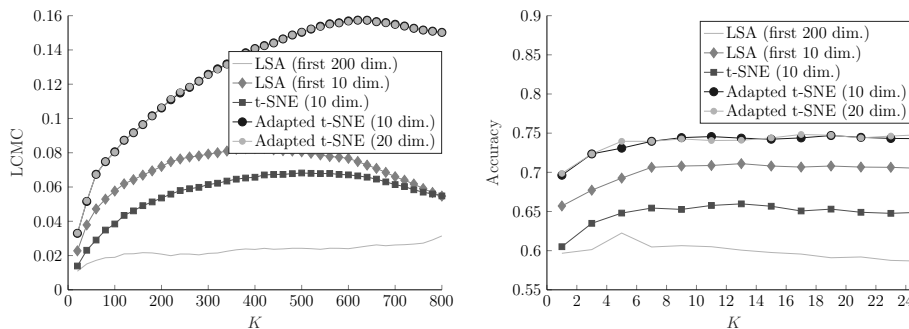


Fig. 3: Left panel: LCMC co-ranking evaluation curves of answer data and embeddings of the question data for varying neighborhood size  $K$ . Right panel: average KNN classification accuracy for several embeddings and values of  $K$ .

## 5 Conclusion and Outlook

In this contribution we demonstrate the usefulness of metric learning for text and email classification. On the publicly available Reuters dataset we demonstrate comparable or superior performance to SVM using LGMLVQ, a prototype based machine learning technique using localized adaptive distances. This method has several advantages compared to alternative techniques: it is of linear complexity

with the number of training samples, learns non-linear decision boundaries and provides a compressed representation of the dataset. We furthermore investigate a real world dataset of customer support emails from a Dutch e-commerce company. In this study a subset of the data most likely answered using FAQ templates serves as ground truth for evaluating two strategies: 1) extracting template classes from the database and training the classifier to predict possible template classes for new incoming customer emails and 2) an adaptive embedding technique based on t-SNE trained to map similarly answered emails closer to each other. Both strategies learn similarities of email threads based on how employees have answered in the past. The adapted similarities can readily be used to retrieve similar emails from the database given a query for computer aided customer support. Future work includes using more powerful feature extraction techniques, the combination of the embedding and classification strategy for non-linear email mappings and life-long/continuous learning adaptations, detecting if new email topics differ significantly from historic ones, computer aided template creation, and comparison with CNNs.

**Acknowledgement:** We thank support by the European Commission's and University of Groningen's COFUND Rosalind Franklin Fellowship program. Furthermore we thank the company Belsimpel for their support for this project.

## References

- [1] J. D. Brutlag and C. Meek. Challenges of the email domain for text classification. In *ICML*, pages 103–110, 2000.
- [2] Y. Yang. An evaluation of statistical approaches to text categorization. *Information Retrieval*, 1:69–90, 1999.
- [3] S. Dumais, J. Platt, D. Heckerman, and M. Sahami. Inductive learning algorithms and representations for text categorization. In *Proc. of the 7th CIKM*, pages 148–155, 1998.
- [4] S. Lai, L. Xu, K. Liu, and J. Zhao. Recurrent convolutional neural networks for text classification. In *AAAI*, volume 333, pages 2267–2273, 2015.
- [5] C. C. Aggarwal. Machine learning for text: An introduction. In *Machine Learning for Text*, chapter 1, pages 1–16. Springer, Cham, 2018.
- [6] P. Schneider, M. Biehl, and B. Hammer. Adaptive relevance matrices in learning vector quantization. *Neural Computation*, 21(12):3532–3561, December 2009.
- [7] K. Bunte, B. Hammer, A. Wismüller, and M. Biehl. Adaptive Local Dissimilarity Measures for Discriminative Dimension Reduction of Labeled Data. *Neurocomputing*, 73(7-9):1074–1092, March 2010.
- [8] K. Bunte, M. Biehl, M. F. Jonkman, and N. Petkov. Learning effective color features for content based image retrieval in dermatology. 44(9):1892–1902, 2011.
- [9] L. van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.
- [10] K. Bunte, M. Biehl, and B. Hammer. A general framework for dimensionality-reducing data visualization mapping. *Neural Computation*, 24(3):771–804, 2012.
- [11] D. Cai. Popular text data sets in matlab format, 2018.
- [12] M. T. Martín-Valdivia, L. A. Ureña-López, and M. García-Vega. The learning vector quantization algorithm applied to automatic text classification tasks. *Neural Networks*, 20(6):748–756, 2007.
- [13] G. Salton and C. Buckley. Term-weighting approaches in automatic text retrieval. *Information processing & management*, 24(5):513–523, 1988.
- [14] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman. Indexing by latent semantic analysis. *J. Assoc. Inf. Sci. Technol.*, 41(6):391–407, 1990.
- [15] M. W. Berry, S. T. Dumais, and G. W. O'Brien. Using linear algebra for intelligent information retrieval. *SIAM review*, 37(4):573–595, 1995.
- [16] L. Chen and A. Buja. Local multidimensional scaling for nonlinear dimension reduction, graph drawing, and proximity analysis. *J AM STAT ASSOC*, 104(485):209–219, 2009.
- [17] J. Lee and M. Verleysen. Quality assessment of dimensionality reduction: Rank-based criteria. *Neurocomputing*, 72(7-9):1431–1443, 2009.