# openXBOW – Introducing the Passau Open-Source Crossmodal Bag-of-Words Toolkit

**Maximilian Schmitt**                    MAXIMILIAN.SCHMITT@UNI-PASSAU.DE
**Björn W. Schuller** *                    BJOERN.SCHULLER@UNI-PASSAU.DE
*Chair of Complex and Intelligent Systems*
*University of Passau*
*Passau, Germany*

## Abstract

We introduce OPENXBOW, an open-source toolkit for the generation of bag-of-words (BoW) representations from multimodal input. In the BoW principle, word histograms were first used as features in document classification, but the idea was and can easily be adapted to, e.g., acoustic or visual low-level descriptors, introducing a prior step of vector quantisation. The OPENXBOW toolkit supports arbitrary numeric input features and text input and concatenates computed subbags to a final bag. It provides a variety of extensions and options. To our knowledge, OPENXBOW is the first publicly available toolkit for the generation of crossmodal bags-of-words. The capabilities of the tool are exemplified in two sample scenarios: time-continuous speech-based emotion recognition and sentiment analysis in tweets where improved results over other feature representation forms were observed.

**Keywords:** bag-of-words, multimodal signal processing, histogram feature representations

## 1. Introduction

The bag-of-words (BoW) principle is a common practice in natural language processing (NLP) (Wöllmer et al., 2012; Weninger et al., 2013). In this method, *word histograms* are generated, i.e., within a text document, the frequencies of words from a dictionary are counted. The resulting word-frequency vector is then input to a classifier, such as *naïve Bayes* or a *support vector machine* (SVM), i.e., machine learning schemes which are known to cope well with possibly irrelevant features and large, yet sparse feature vectors.

One major drawback of the BoW approach is that the order of the words in a document, which often implies important information, is not taken into account. In order to overcome this problem, *n-grams* have been employed, where sequences of *n* words or characters (*n-character-grams*) are counted instead of single words (Wallach, 2006; Schuller et al., 2009).

BoW has been adopted by the visual community, where it is known under the name bag-of-visual-words (BoVW) (Fei-Fei and Perona, 2005; Sivic et al., 2005). Instead of lexical words, local image features are extracted from an image and then their general distribution is modelled by a histogram according to a learnt codebook, which substitutes the dictionary from BoW.

In recent years, the principle has also been employed successfully in the field of audio classification, where it is known under the term bag-of-audio-words (BoAW). Acoustic low-level descriptors (LLDs), such as mel-frequency cepstral coefficients (MFCC), are extracted

---

*. B. W. Schuller is also with the Department of Computing, Imperial College London, UK.

from the audio signal, then, the LLD vectors from single frames are quantised according to a codebook (Pancoast and Akbacak, 2012). This codebook can be the result of either k-means clustering (Pokorny et al., 2015) or random sampling of LLD vectors (Rawat et al., 2013). Other approaches employ expectation maximisation (EM) clustering, which leads to a soft vector quantisation step (Grzeszick et al., 2015). A histogram finally describes the distribution of the codebook vectors over the whole audio signal or one segment. Major applications of BoAW are acoustic event detection and multimedia event detection (Liu et al., 2010; Pancoast and Akbacak, 2012; Rawat et al., 2013; Plinge et al., 2014; Lim et al., 2015) but they have also been successfully used in music information retrieval (Riley et al., 2008) and emotion recognition from speech (Pokorny et al., 2015).

In this contribution, we introduce the first open-source toolkit for the generation of BoW representations across modalities, thus named 'OPENXBOW' ("open cross-BoW"). The motivation behind OPENXBOW is to ease the generation of a fused BoW-based representation from different modalities. These modalities can be the audio or the visual domain, providing numeric LLDs, and written documents or transcriptions of speech, providing text. However, in principle arbitrary other types of modalities "$x$" can be thought off, such as stemming from physiological measurement or feature streams as used in brain computing, etc. As a placeholder for arbitrary modalities, we thus introduce the notion of bag-of-x-words or BoXW for short. In case of multimodal or 'crossmodal' usage, the output of the toolkit is a concatenated feature vector consisting of histogram representations per modality or combinations of these. Multimodal BoW have already been employed, e.g., for depression monitoring in (Joshi et al., 2013), exploiting both the audio and the video domain.

OPENXBOW provides a multitude of options, e.g., different modes of vector quantisation, codebook generation, term frequency weighting, and methods known from natural language processing to process the textual features. To the knowledge of the authors, such a toolkit has not been published, so far, whereas there are already some libraries implementing BoVW, such as *DBoW2* (Gálvez-López and Tardos, 2012).

In the next section, we give an overview of the OPENXBOW tool, its structure and its options. In section 3, we give results from two exemplary applications of the tool. We conclude and give an outlook on future developments in OPENXBOW in the final section 4.

## 2. Overview

OPENXBOW is implemented in Java and can thus be used on any platform. It has been published on GitHub as a public repository[1], including both the source code and a compiled jar file for those who do not have a Java Development Kit installed. The software and the source code can be freely used by the research community for non-commercial purposes.

OPENXBOW supports three different file formats for input of LLDs and text and the output of the BoW feature vector:

- ARFF (Attribute-Relation File Format), used in the machine learning software *Weka* (Hall et al., 2009)

- CSV (Comma separated values), with separator semicolon (;)

---
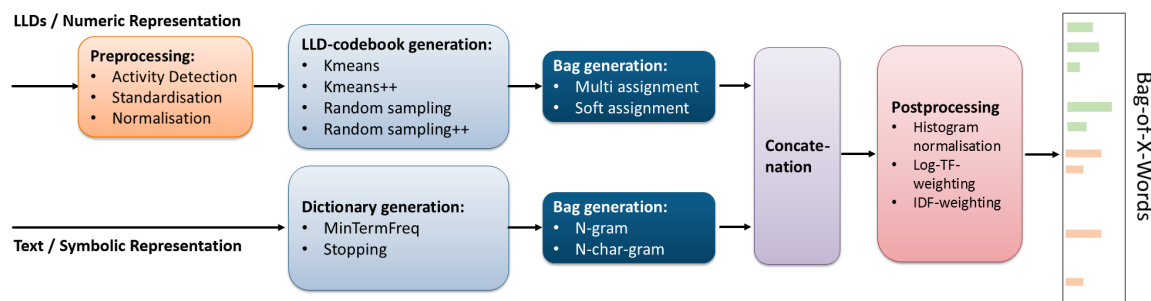
1. https://github.com/openXBOW/openXBOW

Figure 1: Overview of the workflow of openXBOW.

- LIBSVM file format, used in LIBSVM (Chang and Lin, 2011) and LIBLINEAR (Fan et al., 2008) (only output)

The input is processed by openXBOW in the way shown in figure 1. First, there is an optional preprocessing stage, where input LLDs with a low activity can be excluded from further processing. This is especially relevant in the field of speech recognition, where LLDs at instants of time without voice activity should not be considered as they describe only the background noise.

In case feature types with (significantly) different ranges of values are combined in one BoW, normalisation or standardisation of the LLDs is essential. Both options are available, the corresponding parameters are stored in the codebook file, so that they can be applied in the same way to a given test file (online approach).

All options with a corresponding help text are displayed in the command line window if openXBOW is started without any arguments or the argument (-h), e. g.,

```
java -jar openXBOW.jar -h
```

All configurations are made in the command line call.

For codebook generation, there are four different methods available at the time; *random sampling* means that the codebook vectors are picked randomly from the input LLDs, whereas *random sampling++* favours far-off vectors as proposed in (Arthur and Vassilvitskii, 2007) as an improved initialisation step for *kmeans*, called in that case *kmeans++*. For the latter two methods, up to 500 updates of cluster centroids and cluster assignments are executed after initialisation.

In case of nominal class labels, the codebook generation can also be done in a supervised manner, learning a codebook from all LLDs in one class separately, first, and then concatenating these codebooks to form a *super-codebook* (Grzeszick et al., 2015).

In case of large LLD vectors, *split vector quantisation* (SVQ) can be suitable, as proposed in (Pokorny et al., 2015), where subvectors are first quantised and then the indexes of the quantised subvectors are processed in the usual scheme. However, it is also possible to split the input LLDs manually into well-defined subvectors in order to generate different codebooks for different feature types.

A way of doing *soft* vector quantisation without EM-clustering, is applying *Gaussian encoding* to the term frequencies (TF), i. e., the number of occurrences of each word from the

3

codebook (Pancoast and Akbacak, 2014). The TF is then weighted in each word assignment by the distance to the word in the codebook. This can be especially useful in combination with *multiple assignments*, where not only the closest word from the codebook is considered but a certain number $N_a$ of closest words.

For the text domain, standard processing techniques are available, such as *stopping*, *n-grams*, and *n-character-grams*.

Finally, the BoW representations of different domains are fused into one feature vector. For post-processing, logarithmic TF-weighting ($\text{TF}_{\log} = \lg(\text{TF}+1)$), and inverse document-frequency (IDF) weighting ($\text{TF}_{\text{IDF}} = \text{TF} \cdot \lg \frac{\text{N}}{\text{DF}}$, N: Number of instances, DF: Number of instances where the word is present) can be applied (Riley et al., 2008).

The resulting histograms can then be normalised, a step which is essential if the input instances to be classified have different sizes or if (voice or any other modality's) activity detection is used and so different segments of the input signal have different numbers of assigned words.

## 3. Experiments

In this section, the usage of OPENXBOW is exemplified in two scenarios: time-continuous processing of speech input for emotion recognition, and BoW processing of tweets.

In fact, there are comparably few suitable multimodal databases with acoustic, visual, and text input available. In this introductory paper, we exemplify the principle on separate tasks for audio and text words to showcase the principle. However, results of OPENXBOW with multimodal input will be shown in our future efforts.

### 3.1 Time- and value-continuous emotion recognition from speech

Emotion recognition from speech has been conducted on the RECOLA (Remote Collaborative and Affective Interactions) corpus (Ringeval et al., 2013). In this database, 46 French participants have been recorded in dyadic remote collaboration for 5 minutes, each. During their collaboration, their face, speech and some physiological measures have been recorded. In our present example, we focus on the audio domain.

A gold standard in terms of *arousal* and *valence*, generated from the annotations by six different persons, was used as a target in our experiments. The corpus was split into three partitions, a training partition (16 subjects) in order to learn the codebook and train the classifier, a validation partition to optimise the parameters of OPENXBOW and support vector regression (SVR) with a linear kernel and a test partition to prove the universality of the learnt regressor. The input LLDs (MFCCs 1–12 and log-energy) were extracted with our toolkit openSMILE (Eyben et al., 2013).

To learn a codebook and generate a BoAW representation (`BoAW_arousal_train.arff`) from the training file `LLD_train.csv`, the following options have been applied:

```
-i LLD_train.csv -o BoAW_arousal_train.arff -l arousal_train.csv -t 8.0 0.8
-standardizeInput -size 1000 -c random++ -B codebook.txt -a 20 -log
```

The option `-t 8.0 0.8` means that the sequence of input LLDs is segmented into windows of 8.0 seconds width and the hop size between two successive windows is 0.8 seconds. Labels

for arousal in the file `arousal_train.csv` must be available for exactly those instants in time (see option `-h` for information on the labels file format).

After standardisation, a codebook of 1 000 words is generated by *random sampling++* and stored in the file `codebook.txt` together with information on the parameters of standardisation (mean, and standard deviation per LLD). Each LLD is then assigned to the 20 closest (`-a 20`) words in the codebook. Finally, the number TFs are compressed applying logarithmic TF-weighting. Also, this information is stored in the codebook files, as it must be used in the same way with the validation and test instances. The order of the options does not have any effect.

After the optimum configuration of BoAW has been found (let us assume that it is the one described above), the BoAW can be generated from the validation (and also test) files (`valid`) in the following way:

```
-i LLD_valid.csv -o BoAW_arousal_valid.arff -l arousal_valid.csv -t 8.0 0.04
-b codebook.txt -a 20
```

The codebook is loaded with the command `-b`. Please note that, the file `codebook.txt` implies also the standardisation and the log-TF-weighting. The hop size is chosen differently for the validation partition; it is only 0.04 seconds (`-t 8.0 0.04`) which is not an issue if the corresponding labels are included in the labels file (`arousal_valid.csv`).

Table 1 shows the results of BoAW for emotion recognition compared to the baseline of the AVEC 2016 challenge at ACM Multimedia 2016 (Valstar et al., 2016), which is also carried out on RECOLA under the same conditions. However, the number of recordings in the respective training, validation and test set in the AVEC 2016 challenge were smaller than the ones we used.

The optimum window size for the prediction of valence is 10.0 seconds, compared to 8.0 seconds for arousal. All labels have been shifted to the front in time in order to compensate the natural delay between the shown emotion of the subject and the reaction of the annotator. A shift of 4.0 seconds prove to be an optimum. The complexity of SVR has been optimised for each configuration. As a metric for evaluation, the concordance correlation coefficient (CCC) is used, which takes also scaling of the outputs into account, compared to the linear correlation coefficient. Note that, the CCC is also used as competition measure in the AVEC 2016 competition.

Table 1: *Performance of speech-based emotion recognition on RECOLA using BoAW compared to the baseline of AVEC 2016. Shown are results on the official Valid(ation) and Test sets. Results printed in **bold** are statistically significantly better than the baseline (level of significance: 0.01)*

| Model | CCC | | | |
|---|---|---|---|---|
| | Arousal | | Valence | |
| | Valid | Test | Valid | Test |
| Baseline AVEC 2016 (audio only) | .796 | .648 | .455 | .375 |
| openXBOW (BoAW) | .793 | **.753** | **.550** | **.430** |

It can be clearly seen that, our approach significantly outperforms the baseline on the validation and test sets, except for the performance for arousal on the validation partition, where a similar result is achieved.

### 3.2 Twitter sentiment analysis

Our second experiment shows obtainable performances on written text on another well-defined task. A large corpus of short messages from the social network *Twitter*[2], known as 'tweets', in English language is provided by *Thinknook*[3]. This dataset has been collected by the *University of Michigan* and *Niek Sanders*. It includes $1\,578\,627$ tweets with a 2-class annotation of either positive or negative sentiment. An accuracy of $75\,\%$ is reported as state-of-the-art.

In order to train a dictionary and create a BoW representation from the training set `senti-train.csv`, the following command line arguments can be used:

```
-i senti-train.csv -attributes ncr0 -o BoW.arff -B dictionary.txt
-minTermFreq 1000 -maxTermFreq 30000 -nGram 2 -log -idf
```

The option `-attributes` is needed if the data structure of the input is not the standard defined in the OPENXBOW help text. It specifies each input feature in the input file, i. e., each column in an input CSV file. `n` means that, the corresponding column (the first column here) specifies the *name* (or an index) as a unique ID this line belongs to. All lines in the input with the same name are put into one bag later on. The second column in the Twitter dataset (`c`) specifies the class label. The third column (`r`, remove) can be discarded as it does not contain any relevant information, the last column of the input contains the text to be classified. The digit `0` always specifies text input, the digits `1` to `9` specify numeric input. A separate codebook is generated for every digit if at least one feature with a digit is present. More information about the input format can be found in the help text.

`-minTermFreq` and `-maxTermFreq` implement stopping as known in NLP. Thereby, either very rare words or very common words (such as 'and', or 'or') can be excluded from the dictionary as/if they are not likely to carry meaningful information. 2-grams (`-nGram 2`) and log-TF-IDF weighting (`-log -idf`) are also used in this sample call.

To apply the same model to the test set, the following line is used:

```
-i senti-test.csv -attributes ncr0 -o BoW.arff -b dictionary.txt -nGram 2
```

Note that, the TF-weighting is stored in the dictionary and is thus not needed; the option `-nGram 2` must, however, be repeated.

In our experiments, the best results have been achieved without using (higher order) n-grams, i. e., simply employing uni-grams or each word by itself, respectively. We split the whole corpus into two partitions, the first $1\,000\,000$ instances form the training set, the remaining $578\,627$ instances form the test set. Support vector machine with linear kernel was again found to be suitable to handle the high-dimensional BoW vector.

A *minimum term frequency* of 500 and a *maximum term frequency* of $100\,000$, leading to a dictionary size of $1\,875$ terms was found suitable. With an SVM complexity of 0.1,

---

2. http://twitter.com
3. http://thinknook.com/twitter-sentiment-analysis-training-corpus-dataset-2012-09-22/

a weighted accuracy (WA) of 77.28 % and an unweighted accuracy (UA) of 77.29 % have been achieved. These results appear to match or outperform a little bit the state-of-the-art. Even better performance might be achieved with more sophisticated approaches, such as *long short-term memory recurrent neural networks* (Schuller et al., 2015).

## 4. Conclusions and outlook

We introduced our novel OPENXBOW toolkit – a first of its kind – for the generation of BoW representations from multimodal symbolic (including text), but also numeric (such as audio or video feature streams) information representations. In two (monomodal) examples we showed the potential of the toolkit and the underlying BoXW principle by outperforming the state-of-the-art on a well-defined modern speech-based emotion recognition competition task. The full potential is, however, likely to be revealed once targeting actual crossmodal tasks. Likewise, it stands to reason that in the future, the performance can be improved using input from the visual domain, the physiological domain, and the transcribed speech on the emotion recognition task considered.

We have also shown that our toolkit already provides state-of-the-art results in text classification.

Future work on OPENXBOW will include further soft vector quantisation techniques such as using EM clustering or non-negative matrix factorisation-based soft clustering and methods taking the order of the crossmodal words into account, such as temporal augmentation (Grzeszick et al., 2015) and n-grams for numeric features (Pancoast and Akbacak, 2013). Another goal is to add a graphical user interface for configuration in order to make the work more convenient.

The public repository[4] will be regularly updated. The authors kindly request for feedback.

## Acknowledgments

## References

D. Arthur and S. Vassilvitskii. k-means++: The advantages of careful seeding. In *Proc. of the 18th annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035. Society for Industrial and Applied Mathematics, 2007.

C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at `http://www.csie.ntu.edu.tw/~cjlin/libsvm`.

---

4. https://github.com/openXBOW/openXBOW

F. Eyben, F. Weninger, F. Groß, and B. Schuller. Recent developments in openSMILE, the munich open-source multimedia feature extractor. In *Proc. of the 21st ACM International Conference on Multimedia*, pages 835–838, Barcelona, Spain, October 2013. ACM.

R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. Liblinear: A library for large linear classification. *The Journal of Machine Learning Research*, 9:1871–1874, 2008.

L. Fei-Fei and P. Perona. A bayesian hierarchical model for learning natural scene categories. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 2, pages 524–531. IEEE, 2005.

D. Gálvez-López and J. D. Tardos. Bags of binary words for fast place recognition in image sequences. *Robotics, IEEE Transactions on*, 28(5):1188–1197, 2012.

R. Grzeszick, A. Plinge, and G. A. Fink. Temporal acoustic words for online acoustic event detection. In *Proc. 37th German Conf. Pattern Recognition*, Aachen, Germany, 2015.

M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. The weka data mining software: an update. *ACM SIGKDD explorations newsletter*, 11(1):10–18, 2009.

J. Joshi, R. Goecke, S. Alghowinem, A. Dhall, M. Wagner, J. Epps, G. Parker, and M. Breakspear. Multimodal assistive technologies for depression diagnosis and monitoring. *Journal on MultiModal User Interfaces*, 7(3):217–228, 2013.

H. Lim, M. J. Kim, and H. Kim. Robust sound event classification using lbp-hog based bag-of-audio-words feature representation. In *Proc. INTERSPEECH*, pages 3325–3329, Dresden, Germany, September 2015.

Y. Liu, W.-L. Zhao, C.-W. Ngo, C.-S. Xu, and H.-Q. Lu. Coherent bag-of audio words model for efficient large-scale video copy detection. In *Proc. of the ACM International Conference on Image and Video Retrieval*, pages 89–96. ACM, 2010.

S. Pancoast and M. Akbacak. Bag-of-audio-words approach for multimedia event classification. In *Proc. INTERSPEECH*, pages 2105–2108, Portland, USA, September 2012.

S. Pancoast and M. Akbacak. N-gram extension for bag-of-audio-words. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 778–782. IEEE, 2013.

S. Pancoast and M. Akbacak. Softening quantization in bag-of-audio-words. In *Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on*, pages 1370–1374. IEEE, 2014.

A. Plinge, R. Grzeszick, and G. A. Fink. A bag-of-features approach to acoustic event detection. In *IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3732–3736. IEEE, 2014.

F. Pokorny, F. Graf, F. Pernkopf, and B. Schuller. Detection of negative emotions in speech signals using bags-of-audio-words. In *Proc. 1st International Workshop on Automatic Sentiment Analysis in the Wild (WASA 2015) held in conjunction with the 6th biannual Conference on Affective Computing and Intelligent Interaction (ACII 2015)*, pages 879–884, Xi'an, P. R. China, September 2015. AAAC, IEEE.

S. Rawat, P. F. Schulam, S. Burger, D. Ding, Y. Wang, and F. Metze. Robust audio-codebooks for large-scale event detection in consumer videos. In *Proc. INTERSPEECH*, pages 2929–2933, Lyon, France, August 2013.

M. Riley, E. Heinen, and J. Ghosh. A text retrieval approach to content-based audio hashing. In *ISMIR 2008, 9th International Conference on Music Information Retrieval*, pages 295–300, Philadelphia, PA, USA, September 2008.

F. Ringeval, A. Sonderegger, J. Sauer, and D. Lalanne. Introducing the recola multimodal corpus of remote collaborative and affective interactions. In *Face and Gestures 2013, Workshop on Emotion Representation, Analysis and Synthesis in Continuous Time and Space (EmoSPACE)*, 2013.

B. Schuller, J. Schenk, G. Rigoll, and T. Knaup. "The Godfather" vs. "Chaos": Comparing Linguistic Analysis based on Online Knowledge Sources and Bags-of-N-Grams for Movie Review Valence Estimation. In *Proceedings 10th International Conference on Document Analysis and Recognition, ICDAR 2009*, pages 858–862, Barcelona, Spain, July 2009. IAPR, IEEE.

B. Schuller, A. E.-D. Mousa, and V. Vasileios. Sentiment Analysis and Opinion Mining: On Optimal Parameters and Performances. *WIREs Data Mining and Knowledge Discovery*, 5:255–263, September/October 2015.

J. Sivic, B. C. Russell, A. A. Efros, A. Zisserman, and W. T. Freeman. Discovering object categories in image collections. 2005.

M. Valstar, J. Gratch, B. Schuller, F. Ringeval, D. Lalanne, M. T. Torres, S. Scherer, G. Stratou, R. Cowie, and M. Pantic. Avec 2016-depression, mood, and emotion recognition workshop and challenge. *arXiv preprint arXiv:1605.01600*, 2016.

H. M. Wallach. Topic modeling: Beyond bag-of-words. In *Proceedings of the 23rd international conference on Machine learning*, pages 977–984. ACM, 2006.

F. Weninger, P. Staudt, and B. Schuller. Words that Fascinate the Listener: Predicting Affective Ratings of On-Line Lectures. *International Journal of Distance Education Technologies, Special Issue on Emotional Intelligence for Online Learning*, 11(2):110–123, April–June 2013.

M. Wöllmer, M. Kaiser, F. Eyben, F. Weninger, B. Schuller, and G. Rigoll. Fully Automatic Audiovisual Emotion Recognition – Voice, Words, and the Face. In T. Fingscheidt and W. Kellermann, editors, *Proceedings of Speech Communication; 10. ITG Symposium*, pages 1–4, Braunschweig, Germany, September 2012. ITG, IEEE.