

# A

---

## Appendix: SIAM Text Mining Competition 2007

### Overview

The 2007 Text Mining Workshop held in conjunction with the Seventh SIAM International Conference on Data Mining was the first to feature a text mining competition. Members of the Intelligent Data Understanding group at NASA Ames Research Center in Moffett Field, California, organized and judged the competition. Being the first such competition held as a part of the workshop, we did not expect the large number of contestants that more established competitions such as KDDCUP<sup>1</sup> have, but we did receive five submissions, though one person later withdrew.

Matthew E. Otey, Ashok N. Srivastava  
Santanu Das, and Pat Castle

### A.1 Classification Task

The contest focused on developing text mining algorithms for document classification. The documents making up the corpus used in the competition are aviation safety reports documenting one or more problems that occurred on certain flights. These documents come from the Aviation Safety Reporting System (ASRS), and they are publicly available in their raw form at <http://asrs.arc.nasa.gov>. Since each of these documents can describe one or more anomalies that occurred during a given flight, they can be tagged with one or more labels describing a class of anomalies. The goal is to label the documents according to the classes to which they belong, while maximizing both precision and recall, as well as the classifier's confidence in its labeling. This second criterion concerning confidence is useful for presenting results to end-users, as an end-user may be more forgiving of a misclassification if he or she knew that the classifier had little confidence in its labeling.

---

<sup>1</sup> See <http://www.kdnuggets.com/datasets/kddcup.html>.

The competition entries were scored using the following cost function that accounts for both the accuracy and confidence of the classifier. When describing the cost function, we will use the following notation. Let  $D$  be the number of documents in the corpus, and let  $L$  be the number of labels. Let  $F_j$  be the fraction of documents having label  $j$ . Let  $t_{i,j} \in \{-1, +1\}$  be the target (true) value for label  $j$  of document  $i$ . Let  $p_{i,j} \in \{-1, +1\}$  be the predicted value for label  $j$  of document  $i$ , and let  $q_{i,j} \in [0, 1]$  be the classifier's confidence of the value  $p_{i,j}$ . Finally, let  $A_j$  be the area under the ROC curve for the classifier's predictions for label  $j$ . We define the intermediate cost function  $Q_j$  for label  $j$  as

$$Q_j = (2A_j - 1) + \frac{1}{D} \sum_{i=1}^D q_{i,j} t_{i,j} p_{i,j}. \quad (\text{A.1})$$

This function has a maximum value of 2. The final cost function  $Q$  is the average  $Q_j$  for all labels:

$$Q = \frac{1}{C} \sum_{j=1}^C Q_j. \quad (\text{A.2})$$

In case of ties, we also formulated a figure of merit (FOM), defined as

$$FOM = \frac{1}{C} \sum_{j=1}^C \frac{(F - F_j)}{F} Q_j \quad (\text{A.3})$$

where

$$F = \sum_{j=1}^C F_j. \quad (\text{A.4})$$

This figure of merit measures the quality of predictions across all anomaly categories, giving a lower weighting to those categories making up a larger fraction of the dataset. Hence, better  $FOM$  are achieved by obtaining higher accuracies and confidences on rarer classes.

We wrote a small Java program implementing these functions to use during our judging of the submissions. Before the contest deadline, we made this program and its source code available to the contestants so that they could validate its correctness, and so that they could use our implementation to tune their algorithms.

## A.2 Judging the Submissions

A training dataset was provided over a month in advance of the deadline, giving the contestants time to develop their approaches. Two days before the deadline we released the test dataset. Both of these datasets had been run through PLADS, a system that performs basic text processing operations such as stemming and acronym expansion. The contestants submitted their labeling of the test dataset, their confidences in the labeling, and the source code implementing their approach. The scores of the

submissions were calculated using the score function described above. In addition to scoring the submissions, we ran each contestant's code to ensure that it worked and produced the same output that was submitted, and we inspected the source code to ensure that the contestants properly followed the rules of the contest.

### **A.3 Contest Results**

The submissions of the contestants all successfully ran and passed our inspection, and we announced our three winners. First place was awarded to Cyril Goutte at the NRC Institute for Information Technology, Canada, with a score of 1.69. A team consisting of Edward G. Allan, Michael R. Horvath, Christopher V. Kopek, Brian T. Lamb, and Thomas S. Whaples of Wake Forest University, and their advisor, Michael W. Berry of the University of Tennessee, Knoxville, came in second with a score of 1.27. The third place team of Mostafa Keikha and Narjes Sharif Razavian of the University of Tehran in Iran, and their advisor, Farhad Oroumchian of the University of Wollongong, Dubai, United Arab Emirates, scored a 0.97. At NASA, we evaluated Schapire's and Singer's BoosTexter approach, and achieved a maximum score of 0.82 on the test data, showing that the contestants made some significant improvements over standard approaches.

---

# Index

- N*-gram, 220, 223
- k*-means, 45, 58, 76, 96, 103
  - batch, 75
  - bisecting, 49, 58
  - incremental, 75
  - kernel, 52, 56, 60
  - smoothed (*smoka*), 76, 77
- Allan, Edward, 202, 235
- AlSumait, Loulwah, 85
- Alternating Least Squares (ALS), 151
- annealing
  - deterministic, 65, 77, 190
  - simulated, 77
- anti-spam methods, 167
- Aono, Masaki, 109
- ArcRank, 36–41
- attribute selection, 172
  - $\chi^2$ , 172
- authority, 34, 35
- Aviation Safety Reporting System (ASRS),  
203, 205, 219, 233
- Bader, Brett, 147
- bag-of-words, 188
- balanced iterative reducing and clustering  
(BIRCH), 65, 70
- Berry, Michael W., 147, 202, 235
- Blondel, Vincent, 23
- BoosTexter, 235
- break-even point, 197
- Browne, Murray, 147
- calibration, 192
- CANDECAMP, 148
- Castle, Pat, 233
- categorization
  - multiclass, multilabel, 191
- category description, 193
- centroid, 6, 66, 221
  - method, 17, 221
- cluster
  - indicator vector problem, 49
  - measure, 27, 28
  - quality, 7
- clustering, 28, 29, 32, 65, 87
  - algorithm, 45, 65
  - hierarchical, 45
  - partitional, 45
  - spherical *k*-means, 205
- co-occurrence model, 188, 189
- confidence estimation, 187, 192
- corpus, 25–29, 31, 32, 42
  - bilingual, 32
- cosine, 26–29, 42
- COV with selective scaling (COV-SS), 119
- covariance matrix analysis (COV), 114
- cross-validation, 194
- curse of dimensionality, 88, 90
- Das, Santanu, 233
- descriptive keywords, 193
- dictionary
  - bilingual, 32
  - monolingual, 25, 26, 32, 33, 42
  - synonym, 25, 26, 33, 41
- dictionary graph, 33, 34, 36–38, 41

- dimensionality reduction, 5, 171, 178
- discriminant analysis, 6
- distance, 32, 36, 39–42
  - semantic, 88, 93, 94
- Distributed National ASAP Archive (DNAA), 210
- divergence
  - Bregman, 67
  - Csiszár, 67, 68, 204
  - Kullback–Leibler, 67
- document
  - frequency thresholding, 171
  - representation, 220
- document type definition (DTD), 130
- document vector space, 26–28
- Domeniconi, Carlotta, 85
  
- eigenvalue, 53, 54, 73
- eigenvector, 71, 73
  - principal, 35–37
- email classification, 167
- Enron, 147
- Enron email corpus, 65, 78, 147, 152
- ensemble, 26, 32
- entropy, 172
  - relative, 32
- expectation maximization (EM), 187, 190
  
- F1 measure, 98, 100
- factor analysis, 17
- false-positive rate (FPR), 177, 229
- feature
  - reduction, 178
  - selection, 95, 172
- figure of merit (FOM), 207, 234
- Fisher linear discriminant analysis, 204
- folding in, 190
- frequent itemset mining, 95
- function
  - closed, proper, convex, 67
  - distance-like, 69
  
- Gallopoulos, Efstratios, 44
- Gansterer, Wilfried, 163
- General Text Parser (GTP), 154, 205
- generalized singular value decomposition (GSVD), 8
- generalized vector space model (GVSM), 90, 93
  
- Goutte, Cyril, 187, 235
- Gram matrix, 53
- grammatical context, *see* syntactical context
- graph mining, 26, 34, 42
- Green measure, 42
  
- Hadamard product, 149
- hierarchical clustering
  - agglomerative, 46
  - divisive, 46
- HITS, 33, 34, 42
- Horvath, Michael, 202, 235
- Howland, Peg, 3
- hub, 34, 35
- human plausible reasoning, 224
  
- indifferent discriminator, 28, 29
- inductive inference, 188, 191
- information gain, 172
- information retrieval, 25, 26, 28, 29, 42, 75, 167
- inverse document frequency, 221
  
- Jaccard measure, 30
- Janecek, Andreas, 163
  
- Keikha, Mostafa, 217, 235
- kernel
  - clustering, 54
  - local, 88, 93
  - method, 88, 93, 130, 132, 137
- kernel PCA (KPCA), 52
  - Gaussian kernels, 53, 61
  - kernel trick, 52
  - Mercer’s Theorem, 52
  - polynomial kernels, 53, 60
- keywords
  - highest probability, 198
  - largest divergence, 194, 198
- Khatri-Rao product, 149, 151
- Kobayashi, Mei, 109
- Kogan, Jacob, 64
- Kopek, Christopher, 202, 235
- Kronecker product, 149, 152
- Kullback–Leibler divergence, 204
  
- Lamb, Brian, 202, 235
- latent semantic indexing (LSI), 13, 89, 112, 165, 166, 168, 176

- truncation, 169, 176
- latent semantic kernels (LSK), 89, 101
- LDA/GSVD algorithm, 11
- lexical analysis, 29
- local
  - feature relevance, 88
  - feature selection, 90
- locally adaptive clustering (LAC) algorithm, 88, 90
- logical
  - statement, 224
  - terms, 224
- MATLAB Tensor Toolbox, 154
- matricizing, 149
- matrix
  - covariance, 53
  - feature-document, 168
  - proximity, 89, 93
  - scatter, 6
  - semantic, 89
  - semantic dissimilarity, 93
  - term-document, 45, 59, 89
- maximum  $F$ -score threshold, 197
- maximum likelihood (ML), 189, 190
- mean
  - arithmetic, 68
  - geometric, 68
- MEDLINE, 18
- Metropolis algorithm, 77
- minimum error rate threshold, 197
- mixture of multinomial, 188
- multidimensional scaling (MDS), 130
  
- naïve Bayes, 189, 191
- natural language processing (NLP), 220, 228
- neighborhood graph, 33–36, 38, 39, 41
- Neumayer, Robert, 163
- Newton-type method, 130, 137
- Nicholas, Charles, 64
- nonconvexity, 204
- nonnegative matrix factorization (NMF), 152, 190, 203
  - constrained, 191
  - multiplicative update, 152, 205
- nonnegative tensor factorization (NNTF), 151, 155, 158
- normalization factor, 221
  
- Ohsumed collection, 57
- optimization
  - bound-constrained, 204
  - quasi-Newton, 204
- Oroumchian, Farhad, 217, 235
- orthogonal projection, 73
- Otey, Matthew, 233
  
- PageRank, 36, 37, 42
- PARAFAC, 148
  - model, 150
  - nonnegative factors, 151
  - uniqueness, 150
- Park, Haesun, 3
- part of speech, 29, 30, 38
  - tags, 224
- partition
  - $\Pi_A, \Pi_B$ , 71
  - quality of, 66
- phrase, 220, 224
- PLIR, 220, 222, 225
- precision, 26, 28, 29, 98
- principal component analysis (PCA), 13, 46, 52, 109, 114
- principal direction divisive partitioning (PDDP), 45, 46, 58, 59, 66, 71
  - $k$ -means steering, 49
  - geometric interpretation, 48
  - KPDDP-2MEANS, 56
  - KPDDP-OCPC, 56
  - KPDDP: Kernel PDDP, 47, 54, 56, 60
  - PDDP( $l$ ), 48, 58
  - PDDP-2MEANS, 49, 58
  - PDDP-OC, 50, 58
  - PDDP-OC-2MEANS, 50, 58
  - PDDP-OCPC, 52, 58, 59
- principal directions divisive partitioning (PDsDP), 66, 69, 71
- probabilistic
  - categorizer, 187
  - model, 188
  - profiles, 187, 193
- probabilistic latent semantic analysis (PLSA), 188, 190
- PROPACK, 56
  
- Qi, Houduo, 127
  
- recall, 26, 28, 29, 98

- receiver operating characteristic (ROC), 210, 228
- related words, *see* similar words
- Reuters-21578 collection, 57, 225
  - Modapte split, 57
- rich document representation (RDR), 219, 220, 224, 227
- Salton, Gerard, 3, 109
- SDP, *see* semidefinite programming
- search engine, 32, 33
- semantic LAC algorithm, 93
- semantic network, 89
  - WordNet, 89, 102, 103
- semidefinite programming (SDP), 130
- Senellart, Pierre, 23
- SEXTANT, 26, 29–31, 42
- Seyed Razi, Hassan, 217
- Sharif Razavian, Narjes, 217, 235
- similar words, 25–43
- similarity, 25–30, 32–35, 42
- singular value decomposition (SVD), 3, 46, 165, 166, 168, 176, 178
- spam filtering, 166, 169
- SpamAssassin, 165, 166, 169
- sparsity, 79
- Srivastava, Ashok, 233
- stemming, 37, 42, 220, 224
- stop word, 29, 38, 225
  - stoplist, 154, 206
- stop-word, 42
- subspace clustering, 88, 93, 101, 103
- supervised learning, 188
- support vector machine (SVM), 52, 130
- synonym, 25–27, 32–34, 36, 38, 39, 41, 43
- syntactical analysis, 26, 29–32
- syntactical context, 26, 29, 31, 32
- tensor
  - decomposition, 148
  - norm, 149
- outer product, 149
- term
  - discrimination value, 26, 28
  - frequency, 221
  - vector space model, 26, 28
- term vector space, 27, 28
- term weighting
  - log-entropy, 154
- term-weighting, 89, 93
- inverse document frequency (*idf*), 89
- log-entropy, 206
- text
  - categorization, 188
  - classification, 219
- text to matrix generator (TMG), 57
- tf-idf, 28, 42
- thesaurus, 25, 26, 28, 29, 32, 36, 42, 43
- trace optimization, 7
- transductive inference, 188, 191
- tree edit distance, 131
- true-positive rate (TPR), 177, 229
- unfolding, 149
- unsolicited bulk email (UBE), 165
- unsolicited commercial email (UCE), 165
- vector space model (VSM), 5, 45, 88, 111, 165, 166, 168, 176, 220
- Wedderburn rank reduction, 5
- weighted Euclidean distance, 91
- weighting schema, 221
- Whaples, Thomas, 202, 235
- Wiacek, Mike, 64
- WordNet, 38–41
- World Wide Web, 25–27, 32–34, 41, 42
- Xia, Zhonghang, 127
- Xing, Guangming, 127
- Zeimpekis, Dimitrios, 44