# Toward Automated Early Sepsis Alerting: Identifying Infection Patients from Nursing Notes

**Emilia Apostolova**
Language.ai
Chicago, IL, USA
emilia@language.ai

**Tom Velez**
Vivace Health Solutions
Cardiff, CA, USA
tom.velez@cta.com

## Abstract

Severe sepsis and septic shock are conditions that affect millions of patients and have close to 50% mortality rate. Early identification of at-risk patients significantly improves outcomes. Electronic surveillance tools have been developed to monitor structured Electronic Medical Records and automatically recognize early signs of sepsis. However, many sepsis risk factors (e.g. symptoms and signs of infection) are often captured only in free text clinical notes. In this study, we developed a method for automatic monitoring of nursing notes for signs and symptoms of infection. We utilized a creative approach to automatically generate an annotated dataset. The dataset was used to create a Machine Learning model that achieved an F1-score ranging from 79 to 96%.

## 1 Introduction

Severe sepsis and septic shock are rapidly progressive, life-threatening conditions caused by complications from an infection. They are major healthcare problems that affect millions of patients globally each year (Kim and Hong, 2016). The mortality rate for severe sepsis and septic shock is approaching 50% (Nguyen et al., 2006).

A key goal in critical care medicine is the early identification and treatment of infected patients with early stage sepsis. The most recent guidelines for the management of severe sepsis and septic shock include early recognition and management of these conditions as medical emergencies, immediate administration of resuscitative fluids, frequent reassessment, and empiric antibiotics as soon as possible following recognition (Dellinger et al., 2008).

Early recognition of infections that can lead to sepsis, severe sepsis and/or septic shock can be challenging for several reasons: 1) these conditions can quickly develop from any form of common infections (bacterial, viral or fungal) and can be localized or generalized; 2) culture-dependent diagnosis of infection is commonly slow and prior use of antibiotics may make cultures falsely negative (Vincent, 2016); 3) systemic inflammatory response syndrome, traditionally associated with sepsis, may be the result of other noninfectious disease processes (Bone et al., 1992). Consequently, clinicians frequently rely on a myriad of non-specific symptoms of infections and physiological signs of systemic inflammatory response for rapid diagnosis. Each hour of delay in the administration of recommended therapy is associated with a linear increase in the risk of mortality rate (Kumar et al., 2006; Han et al., 2003), driving the need for automation of early sepsis recognition.

In response to this need, electronic surveillance tools have been developed to monitor for the arrival of new patient electronic medical record (EMR) data, automatically recognize early signs of sepsis risk in specific patients, and trigger alerts to clinicians to help guide timely, effective interventions (Herasevich et al., 2011; Azzam et al., 2009; Koenig et al., 2011). The automated decision logic used in many existing sepsis screening tools, for example (Nguyen et al., 2014; Hooper et al., 2012; Nelson et al., 2011), relies on consensus criteria-based rules.

Structured EMR data, such as diagnostic codes, vital signs and orders for tests, imaging and medications, can be a reliable source of sepsis criteria. However, many sepsis risk factors (e.g. symptoms and signs of infection) are routinely captured solely in free text clinical notes. The aim of this

study is to develop a system for the detection of signs and symptoms of infection from free-text nursing notes. The output of the system is later used, in conjunction with available structured data, as an input to an electronic surveillance tool for an early detection of sepsis.

## 2  Task Definition and Dataset

Depending on the infection source and the specifics of the patient history, signs and symptoms of infection can vary widely. In addition, similar symptoms can be stated using a number of synonymous expressions, complicated by the presence of abbreviations, variant spellings and misspellings. Table 1 lists nursing note snippets indicating a possible presence of infection with various degrees of certainty. For example, line items 1, 8, and 12 indicate increased temperature; line items 1, 3, 6, 7, and 13 indicate the use of infection-treating antibiotics; line items 4, 7, 9, 11, and 13 mention specific infectious diseases. A number of examples mention additional infection symptoms or infection detecting tests.

| 1 | Afebrile on antibiotics. |
|---|---|
| 2 | Very large copious amount of secretions,sputum |
| 3 | ... medicated with iv cefazolin dose 2 of 3 |
| 4 | UA positive for UTI. |
| 5 | Blood culture pos for gram neg organisms. |
| 6 | Elevated WBC count, on clindamycin IV. |
| 7 | ... continues on clindamycin(D6), pen-G(D5) and doxycycline(D4) for LLL pneumonia |
| 8 | 84 year old male with h/o throat cancer who presented on [DATE] to [LOCATION] with fever, diffuse rash, renal failure and altered mental status. |
| 9 | Pt had a positive sputum specs for GBS and GPC, and he has HSV on lips. |
| 10 | ... blood, urine and sputum culture sent today ... |
| 11 | PEG tube to gravity, episodes of vomitting on day shift,NPO. Contact precautions MRSA. |
| 12 | Pt had temp spike of 102.4 |
| 13 | Levaquin started for pnuemonia. |
| 14 | Had infected knee prosthesis which led to wash out of joint yesterday. |

Table 1: MIMIC-III nursing note snippets indicating a presence of infection with various degrees of certainty. Abbreviations and misspellings are preserved to demonstrate the task challenges.

In the literature of identifying patient phenotype cohorts using electronic health records, most studies map textual elements to standard vocabularies, such as the Unified Medical Language System (UMLS) (Shivade et al., 2014). The standard vocabulary concepts are later used in rule-based, and, in some cases, Machine Learning (ML) approaches to identify patient cohorts.

In the context of identifying infection from clin-ical notes, however, such an approach poses a number of challenges. Symptoms can vary widely depending on the source of infection, for example, *redness, sputum, swelling, pus, phlegm, vomiting, increased white blood cell count*, etc. The same symptom can also be expressed in a large number of ways, for example, *afebrile, temp spike of 102.4, fever*, etc. There is a large number of conditions indicating infections, for example *UTI, strep throat, hepatitis, HIV/AIDS*, etc. In addition, abbreviations and misspellings are quite common in the context of ICU care, for example, *pneumonia, PNA, pnuemonia, pneu*, etc.

Due to their nature, the dataset and task are better suited for ML approaches that are not relying on standard vocabularies or a structured set of features. As with most ML tasks in the clinical domain, the challenge in this approach is obtaining a sufficient amount of training data (Chapman et al., 2011).

To address these challenges, we utilized the MIMIC-III dataset (Johnson et al., 2016) and developed a creative solution to automatically generate training data as described in section 3. MIMIC (Medical Information Mart for Intensive Care) is a large, freely-available database comprising deidentified health-related data associated with over 40,000 patients who stayed in critical care units of the Beth Israel Deaconess Medical Center between 2001 and 2012. The dataset contains over 2 million free-text clinical notes. We focused only on nursing notes for adult patients, and our dataset consists of a total of 634,369 nursing notes.

## 3  Rule-based Training Dataset Creation

To obtain a sizable training dataset we explored the use of available MIMIC-III structured data, such as test orders and results, prescribed medications, and diagnosis codes. However, this approach did not translate to accurately identifying nursing notes suggesting infection for a number of reasons[1]. Instead, we utilized a simple heuristic. We observed that whenever there is an existing infection or a suspicion of infection, the nursing notes describe the fact that the patient is taking or is prescribed infection-treating antibiotics. Thus, identifying nursing notes describing the use of antibiotics will, in most cases, also identify nursing notes describing signs and symptoms of infection.

---

[1]Challenges include missing or incorrect data, discontinuous or disordered EMR data entry timestamps, etc.

To identify positive mentions of administered antibiotics, we used a list of the 60 most commonly administered infection-treating antibiotics in the MIMIC dataset (Misquitta, 2013). This initial list was then extended to include additional antibiotic names, brands, abbreviations, spelling variations, and common misspellings. We semi-automated this laborious task by utilizing *word embeddings* (Mikolov et al., 2013). Word embeddings were generated utilizing all available MIMIC-III nursing notes[2]. The initial set of antibiotics was then extended using the closest *word embeddings* in terms of cosine distance. For example, the closest words to the antibiotic *amoxicillin* are *amox, amoxacillin, amoxycillin, cefixime, suprax, amoxcillin, amoxicilin*. As shown, this includes misspellings, abbreviations, similar drugs and brand names. The extended list was then manually reviewed. The final infection-treating antibiotic list consists of 402 unambiguous expressions indicating antibiotics.

Antibiotics, however, are sometimes negated and are often mentioned in the context of allergies (e.g. *allergic to penicillin*). To distinguish between affirmed, negated, and speculated mentions of administered antibiotics, we also developed a set of rules in the form of keyword triggers. Similarly to the NegEx algorithm (Chapman et al., 2001), we identified phrases that indicate uncertain or negated mentions of antibiotics that precede or follow a list of hand-crafted expression at the sentence and clause levels. Word embeddings were again used to extend the list of triggers with synonyms, spelling variations, abbreviations, and misspellings. For example, the words *allergic, anaphylaxis, anaphalaxis, allerg,* and *anaphylaxsis* are all used as triggers indicating the negation of an antibiotic use. The full list of keywords indicating antibiotics, negation/speculation triggers and conjunctions is available online[3].

The described approach identified 186,158 nursing notes suggesting the unambiguous presence of infection (29%) and 3,262 notes suggesting possible infection. The remaining 448,211 notes (70%) were considered to comprise our negative dataset, i.e. not suggesting infection.

# 4 Machine Learning Results

We modeled the task as a binary classification of free-form clinical notes. It has been shown that Support Vector Machines (Cortes and Vapnik, 1995) achieve superior results in most text classifications tasks and were selected as a sensible first choice. The individual nursing notes were represented as a bag-of-words (1-grams). The tokens were all converted to lower case and non-alphanumeric characters were discarded. Tokens that are present in more than 60% of all samples or less than 6 times were also discarded. The tokens were weighted using the tf-idf scheme (Salton and McGill, 1986). We trained the model using linear kernel SVMs[4] (Chang and Lin, 2011). We set the positive class weight to 2 to address the unbalanced dataset. 70% of the automatically generated dataset was used for training and the remaining 30% for testing. This resulted in a precision of 93.12 and a recall of 99.04 as shown in Table 2.

|  | Precision | Recall | F1-score |
|---|---|---|---|
| SVM[auto] | 93.12 | 99.04 | 95.99 |
| SVM[gold] | 92.10 | 68.46 | 78.53 |

Table 2: Classification Results. SVM[auto]=Results from applying the SVM model on an automatically generated test set of 190,000 nursing notes; SVM[gold]=Results from applying the SVM model on a manually reviewed dataset of 200 nursing notes.

As the training dataset was automatically created, the above results do not truly reflect the model performance. To evaluate the model on the ground truth, a qualified professional manually reviewed 200 randomly selected nursing notes. These results are also shown in Table 2. While the model precision remained high (92.10), the recall dropped significantly to 68.46.

The drop in recall can be partially attributed to the manner in which the testing data was created. Nursing notes describing signs of infection but failing to mention the use of antibiotics were considered (incorrectly) negative examples. However, an error analysis revealed that the majority of the false negatives (contributing to the low recall) were actually all indicating low level of suspicion of infection. For example, the human annotator considered the following snippets sufficient to indicate a possible infection *afebrile, bld cx's sent; monitor temp, wbc's, await stool cx results; lungs coarse, thick yellow secretions suctioned from ett;*

---

[2]We used vector size 200, window size 7, and continuous bag-of-words model.

[3]https://github.com/ema-/antibiotic-dictionary

[4]We used the LibSVM library with both the cost and gamma parameters set to 2 (obtained via grid-search parameter estimation).

*awaiting results of CT, malignancy vs pneumonia.* In all cases, the note expresses only a suspicion for infection, pending further tests.

We further attempted to improve the system performance by utilizing Paragraph Vectors (Le and Mikolov, 2014). Unsupervised algorithms have been used to represent variable pieces of texts such as paragraphs and documents as fixed-length feature representations (Paragraph Vectors). Studies have shown that Paragraph Vectors outperform bag-of-words models on some text classification tasks. We used the text from all nursing notes to create Paragraph Vectors. We generated document embeddings using a distributed memory model and distributed bag-of-words model, each of size 300 with a window size of 7. Combining the vectors of the distributed memory model and the distributed bag-of-words model, we represented each document as a vector of size 600. The paragraph vectors of the training instances were then fed to a logistic regression, K-nearest neighbors, and an SVM classifier. Results significantly under-performed the SVM bag-of-words model and we were able to achieve a maximum precision and recall of 63% and 77% respectively.

## 5   Related Work

A review of approaches to identifying patient phenotype cohorts using EMR data (Shivade et al., 2014) describes a number of studies using clinical notes, most often in combination with additional structured information, such as diagnosis codes. The study asserts that clinical notes are often the only source of information from which to infer important phenotypic characteristics.

Demner-Fushman et al. (2009) note that clinical events monitoring is one of the most common and essential tasks of Clinical Decision Support systems. The task is in many respects similar to the task of identifying patient phenotype cohorts and it has been observed that free text clinical notes are again the best source of information. For example, Murff et al. (2003) found the electronic discharge summaries to be an excellent source for detecting adverse events. They also note that simple keywords and triggers are not sufficient to detect such events.

In the context of identifying infection from clinical text, most studies map textual elements to standard vocabularies, such as UMLS. For example, Matheny et al. (2012) develop a system for detecting infectious symptoms from emergency de-

partment and primary care clinical documentation, utilizing keywords and SNOMED-CT concepts. Bejan et al. (2012) describe a system for pneumonia identification from narrative reports using n-grams and UMLS concepts. Similarly, Elkin et al. (2008) encoded radiology reports using SNOMED-CT concepts and developed a set of rules to identify pneumonia cases.

Horng et al. (2017) develop an automated trigger for sepsis clinical decision support at emergency department triage. They utilize machine learning and establish that free text drastically improves the discriminatory ability of identifying infection (increase in AUC from 0.67 to 0.86). Arnold et al. (2014) develop an EHR screening tool to identify sepsis patients. They utilize NLP applied to clinical documentation, providing greater clinical context than laboratory and vital sign screening alone. DeLisle et al. (2010) used a combination of structured EMR parameters and text analysis to detect acute respiratory infections. Murff et al. (2011) develop a natural language processing search approach to identify postoperative surgical complications within a comprehensive electronic medical record.

Halpern et al. (2014) describe a system for learning to estimate and predict clinical state variables without labeled data. Similar to our approach, they use a combination of domain expertise and vast amounts of unlabeled data, without requiring labor-intensive manual labeling. In their system, an expert encodes a certain amount of domain knowledge (identifying anchor variables) which is later used to train classifiers. Elkan and Noto (2008) show that a classifier trained on positive and unlabeled examples predicts probabilities that differ by only a constant factor from the true conditional probabilities of being positive.

## 6   Discussion

We presented an approach to identifying nursing notes describing the suspicion or presence of an infection. We utilized the MIMIC-III dataset and a creative approach to obtain an ample amount of annotated data. We then applied ML methods to the task and achieved performance sufficient for practical applications. The ultimate goal of this study is to utilize free-text notes, in combination with structured EMR data, to build an automated surveillance system for early detection of patients at risk of sepsis.

# References

R Arnold, J Isserman, S Smola, and E Jackson. 2014. Comprehensive assessment of the true sepsis burden using electronic health record screening augmented by natural language processing. *Critical Care* 18(1):P244.

Helen C Azzam, Satjeet S Khalsa, Richard Urbani, Chirag V Shah, Jason D Christie, Paul N Lanken, and Barry D Fuchs. 2009. Validation study of an automated electronic acute lung injury screening tool. *Journal of the American Medical Informatics Association* 16(4):503–508.

Cosmin Adrian Bejan, Fei Xia, Lucy Vanderwende, Mark M Wurfel, and Meliha Yetisgen-Yildiz. 2012. Pneumonia identification using statistical feature selection. *Journal of the American Medical Informatics Association* 19(5):817–823.

Roger C Bone, William J Sibbald, and Charles L Sprung. 1992. The accp-sccm consensus conference on sepsis and organ failure. *CHEST journal* 101(6):1481–1483.

Chih-Chung Chang and Chih-Jen Lin. 2011. LIB-SVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology* 2:27:1–27:27. Software available at http://www.csie.ntu.edu.tw/~cjlin/libsvm.

Wendy W Chapman, Will Bridewell, Paul Hanbury, Gregory F Cooper, and Bruce G Buchanan. 2001. A simple algorithm for identifying negated findings and diseases in discharge summaries. *Journal of biomedical informatics* 34(5):301–310.

Wendy W Chapman, Prakash M Nadkarni, Lynette Hirschman, Leonard W D'avolio, Guergana K Savova, and Ozlem Uzuner. 2011. Overcoming barriers to nlp for clinical text: the role of shared tasks and the need for additional creative solutions.

Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning* 20(3):273–297.

Sylvain DeLisle, Brett South, Jill A Anthony, Ericka Kalp, Adi Gundlapallli, Frank C Curriero, Greg E Glass, Matthew Samore, and Trish M Perl. 2010. Combining free text and structured electronic medical record entries to detect acute respiratory infections. *PloS one* 5(10):e13377.

R Phillip Dellinger, Mitchell M Levy, Jean M Carlet, Julian Bion, Margaret M Parker, Roman Jaeschke, Konrad Reinhart, Derek C Angus, Christian Brun-Buisson, Richard Beale, et al. 2008. Surviving sepsis campaign: international guidelines for management of severe sepsis and septic shock: 2008. *Intensive care medicine* 34(1):17–60.

Dina Demner-Fushman, Wendy W Chapman, and Clement J McDonald. 2009. What can natural language processing do for clinical decision support? *Journal of biomedical informatics* 42(5):760–772.

Charles Elkan and Keith Noto. 2008. Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, pages 213–220.

Peter L Elkin, David Froehling, Dietlind Wahner-Roedler, Brett E Trusko, Gail Welsh, Haobo Ma, Armen X Asatryan, Jerome I Tokars, S Trent Rosenbloom, and Steven H Brown. 2008. Nlp-based identification of pneumonia cases from free-text radiological reports. In *AMIA*.

Yoni Halpern, Youngduck Choi, Steven Horng, and David Sontag. 2014. Using anchors to estimate clinical state without labeled data. In *AMIA Annual Symposium Proceedings*. American Medical Informatics Association, volume 2014, page 606.

Yong Y Han, Joseph A Carcillo, Michelle A Dragotta, Debra M Bills, R Scott Watson, Mark E Westerman, and Richard A Orr. 2003. Early reversal of pediatric-neonatal septic shock by community physicians is associated with improved outcome. *Pediatrics* 112(4):793–799.

Vitaly Herasevich, Mykola Tsapenko, Marija Kojicic, Adil Ahmed, Rachul Kashyap, Chakradhar Venkata, Khurram Shahjehan, Sweta J Thakur, Brian W Pickering, Jiajie Zhang, et al. 2011. Limiting ventilator-induced lung injury through individual electronic medical record surveillance. *Critical care medicine* 39(1):34–39.

Michael H Hooper, Lisa Weavind, Arthur P Wheeler, Jason B Martin, Supriya Srinivasa Gowda, Matthew W Semler, Rachel M Hayes, Daniel W Albert, Norment B Deane, Hui Nian, et al. 2012. Randomized trial of automated, electronic monitoring to facilitate early detection of sepsis in the intensive care unit. *Critical care medicine* 40(7):2096.

Steven Horng, David A Sontag, Yoni Halpern, Yacine Jernite, Nathan I Shapiro, and Larry A Nathanson. 2017. Creating an automated trigger for sepsis clinical decision support at emergency department triage using machine learning. *PloS one* 12(4):e0174708.

Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific data* 3.

Won-Young Kim and Sang-Bum Hong. 2016. Sepsis and acute respiratory distress syndrome: recent update. *Tuberculosis and respiratory diseases* 79(2):53–57.

Helen C Koenig, Barbara B Finkel, Satjeet S Khalsa, Paul N Lanken, Meeta Prasad, Richard Urbani, and Barry D Fuchs. 2011. Performance of an automated electronic acute lung injury screening system in intensive care unit patients. *Critical care medicine* 39(1):98–104.

Anand Kumar, Daniel Roberts, Kenneth E Wood, Bruce Light, Joseph E Parrillo, Satendra Sharma, Robert Suppes, Daniel Feinstein, Sergio Zanotti, Leo Taiberg, et al. 2006. Duration of hypotension before initiation of effective antimicrobial therapy is the critical determinant of survival in human septic shock. *Critical care medicine* 34(6):1589–1596.

Quoc V Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *ICML*. volume 14, pages 1188–1196.

Michael E Matheny, Fern FitzHenry, Theodore Speroff, Jennifer K Green, Michelle L Griffith, Eduard E Vasilevskis, Elliot M Fielstein, Peter L Elkin, and Steven H Brown. 2012. Detection of infectious symptoms from va emergency department and primary care clinical documentation. *International journal of medical informatics* 81(3):143–156.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. pages 3111–3119.

Donald Misquitta. 2013. *Early Prediction of Antibiotics in Intensive Care Unit Patients*. Ph.D. thesis, The Center for Biomedical Informatics at the Harvard Medical School.

Harvey J Murff, Fern FitzHenry, Michael E Matheny, Nancy Gentry, Kristen L Kotter, Kimberly Crimin, Robert S Dittus, Amy K Rosen, Peter L Elkin, Steven H Brown, et al. 2011. Automated identification of postoperative complications within an electronic medical record using natural language processing. *Jama* 306(8):848–855.

Harvey J Murff, Vimla L Patel, George Hripcsak, and David W Bates. 2003. Detecting adverse events for patient safety research: a review of current methodologies. *Journal of biomedical informatics* 36(1):131–143.

Jessica L Nelson, Barbara L Smith, Jeremy D Jared, and John G Younger. 2011. Prospective trial of real-time electronic surveillance to expedite early care of severe sepsis. *Annals of emergency medicine* 57(5):500–504.

H Bryant Nguyen, Emanuel P Rivers, Fredrick M Abrahamian, Gregory J Moran, Edward Abraham, Stephen Trzeciak, David T Huang, Tiffany Osborn, Dennis Stevens, David A Talan, et al. 2006. Severe sepsis and septic shock: review of the literature and emergency department management guidelines. *Annals of emergency medicine* 48(1):54–e1.

Su Q Nguyen, Edwin Mwakalindile, James S Booth, Vicki Hogan, Jordan Morgan, Charles T Prickett, John P Donnelly, and Henry E Wang. 2014. Automated electronic medical record sepsis detection in the emergency department. *PeerJ* 2:e343.

Gerard Salton and Michael J McGill. 1986. Introduction to modern information retrieval .

Chaitanya Shivade, Preethi Raghavan, Eric Fosler-Lussier, Peter J Embi, Noemie Elhadad, Stephen B Johnson, and Albert M Lai. 2014. A review of approaches to identifying patient phenotype cohorts using electronic health records. *Journal of the American Medical Informatics Association* 21(2):221–230.

Jean-Louis Vincent. 2016. The clinical challenge of sepsis identification and monitoring. *PLoS Med* 13(5):e1002022.