

DIAGNOSIS ENGINE: A FEATURE SELECTION ALGORITHM FOR PERSONALIZED DIAGNOSIS

Jinsung Yoon, Mihaela van der Schaar

Department of Electrical Engineering, University of California, Los Angeles

ABSTRACT

Diagnosis decision support systems (DDSS) play an increasingly important role helping clinicians make informed diagnosis decisions. DDSS use the high-dimensional clinical data stored in electronic health records and learn from it to recommend diagnosis for new patients. However, discovering which data is relevant to consider when diagnosing a patient can be very challenging. We present a novel DDSS, the diagnosis engine (DE), that discovers which patient data/features are most relevant (informative) to determine a specific diagnosis and subsequently use this information to make diagnosis recommendation for new patients. While DE is general and can be applied to diagnosing various diseases, we evaluate its performance in the context of breast cancer and show that DE significantly outperforms state-of-the-art solutions in terms of prediction error rates (by 7.47%) and false positive rates (by 13.25%) when the false negative rates is fixed to be below a threshold set by medical practice (usually 2%).

Index Terms— Relevance Learning, Feature Selection, Healthcare Informatics, Diagnosis Decision Support System

1. INTRODUCTION

Clinicians are routinely faced with the practical challenge of integrating high-dimensional data in order to select proper tests and diagnosis for a given patient. As the understanding of complex diseases (such as cancer) is progressing and, together with it, the range of available tests grows as well, the difficulty of determining the appropriate diagnosis for a particular patient increases as well. Furthermore, statistics show that diagnostic errors yield around 10% of patient deaths; moreover, they represent the primary type of medical malpractice claims in United States [1]. This underscores the urgent need for building diagnosis decision support systems (DDSS) that can assist clinicians in determining the correct diagnosis [2]. DDSS can capitalize on the wealth of information that is being routinely collected in the electronic health records (EHR) of patients. This provides an unprecedented opportunity to 1) correctly diagnose a patient by appropriately considering the diversity of available information about him/her and 2) use historical information about similar patients and their diseases to learn the correct diagnosis for the current patient [3]. However, capitalizing on this information is difficult precisely because there is too much of it; thus, DDSS needs to extract the

information that is actually relevant for diagnosis and diagnose among the wealth of available information [4].

In this paper, we present a novel DDSS approach – which we refer to as the Diagnosis Engine (DE) – that is able to discover out of the vast available EHR data the patient features (i.e. the intrinsic characteristics of a patient and/or his/her medical test results) that are relevant to establish a specific diagnosis and then use this information to issue personalized diagnosis recommendations for the current patient to the attending clinician. The feature discovery component of DE – which we refer to as Diagnosis-Relevant Feature Selection (DiReFS) – is capable of learning which features are most informative to consider in order to make an accurate diagnosis for a patient.

While in this paper we apply DE to breast cancer diagnosis, its approach is general and can also be used for diagnosing other diseases. Moreover, while here we only show the applicability of DE to personalized diagnosis recommendations, the proposed method can be relatively easily extended to personalized treatment recommendations in addition to personalized diagnosis [5-6].

The primary contributions of this paper are:

- We develop a new method (DiReFS) for discovering what features are most relevant to consider when making a diagnosis.
- Using the discovered relevant features, we developed a diagnosis recommendation system (DE) which can be used by the clinicians when attending to a patient. (Alternatively, DE can also be used by the clinicians and/or patients, to get a second, independent opinion.)
- We apply DE to the diagnosis of breast cancer from images of cellular samples obtained from fine needle aspiration (FNA) of breast mass. DiReFS is used to discover which features are relevant to make a correct diagnosis and then use this knowledge to build a diagnosis recommendation system. We test DE on a well-known dataset and show that it consistently and significantly outperforms diagnostic systems based on state-of-the-art machine learning and feature extraction methods.

2. RELATION TO PRIOR WORK

Current medical practice relies on manually curated systematic reviews and clinical guidelines that provide diagnosis recommendations for large groups of patients

rather than personalized diagnosis that are tailored to individual patients. DDSS have been proposed before to help clinicians make more informed decision, but many of them do not consider the specific characteristics (features) of patients and do not provide personalized diagnosis recommendations; hence, they are not very accurate and have only limited applicability in practice [7-8]. Some DDSSs issue accurate diagnosis recommendations for certain diseases, but based only on a small number of manually selected features [9-11]. Whenever the number of features is large (as it is the case for breast cancer diagnosis), these methods are not applicable [12]. Instead, our DE can robustly issue accurate diagnosis for patients even when the number of features is large by identifying the features that are most relevant to consider when diagnosing a patient.

Another strand of literature related to this work is that on machine learning techniques such as Support Vector Machines (SVMs), AdaBoost, logistic regression etc. However, as shown in the experiments section, these methods are not able to issue accurate recommendations. The reason is that they cannot accurately capture the nuanced relationships between patient characteristics and various diagnosis decisions.

Finally, feature selection algorithms such as correlation feature selection (CFS) and mutual information feature selection (MIFS) [13-15] are also related to DiReFS. However, DiReFS is very different from existing feature selection algorithms which focus on the patients' characteristics and not on how these characteristics differently impact on different diagnosis: our approach is capable of discovering *different* features that are relevant to *different* diagnosis. This makes DiReFS similar to our prior work [16-17] – the RELEAF algorithm. However, unlike RELEAF, which is very slow because it must compare all combinations of features, DiReFS is able to discover the relevant features in a very fast and efficient manner because it adopts a new sequential feature selection approach.

3. PROBLEM FORMULATION

In this section, we introduce the proposed Diagnosis Engine (DE) which consists of the diagnosis-relevant feature selection and the diagnosis recommendation algorithm. Figure 1 depicts the proposed system as applied to breast cancer diagnosis: it issues a diagnosis recommendation (tumor is benign or malignant) for a patient based on the relevant features extracted from images of cellular samples obtained from FNA of breast mass.

Let $\mathbf{x} = (x_1, x_2, \dots, x_D)$ denote the patient's feature information where D is the total number of features extracted from the imaging of a patient such as tumor radius, texture, perimeter, etc.; $a \in A \triangleq \{a_1, a_2, \dots, a_K\}$ denotes the action (i.e., diagnosis recommendation) for the patient. For the breast cancer diagnosis used for illustration in this paper, the action/diagnosis recommendation is simply whether the tumor of the patient is benign or malignant

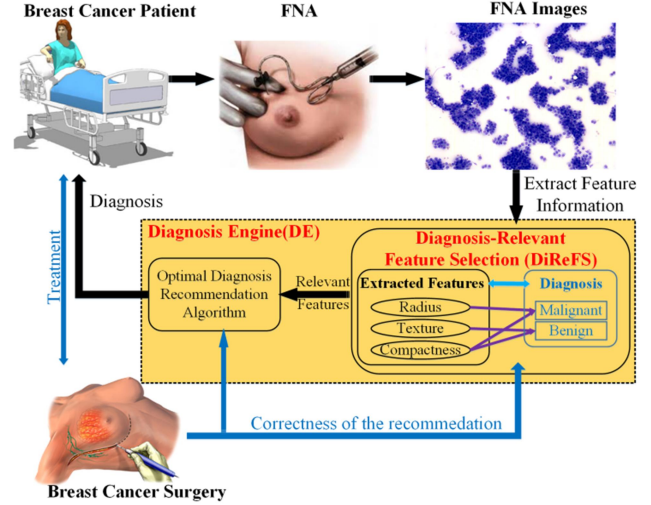


Figure 1: Diagnosis decision support system using DE

(binary action). However, DE is applicable to any discrete number of actions. Each feature is denoted as $f \in F \triangleq \{f_1, f_2, \dots, f_D\}$. y is defined as the prediction accuracy of the diagnosis: if the diagnosis is correct, y is 1, otherwise, y is 0. Let $\mathbf{x}(n), a(n), y(n)$ be the n -th patient information which includes her imaging, action and prediction accuracy and $\mathcal{H}_N = (\mathbf{x}(n), a(n), y(n))_{n=1}^N$ be the information available for the N previously seen patients as stored in EHR. This information becomes the basis for making diagnostic decisions for the $(N + 1)$ -th patient.

The diagnosis of breast cancer does not depend on all the features that can be extracted from the FNA images. We assume that the diagnosis a depends only on a subset of features $\mathcal{R}(a) \subseteq F$ which we refer to as the relevant features for diagnosis a . A key challenge is that the features that are relevant for recommending the breast cancer diagnosis are not known *a priori*; they need to be *discovered/learned*. Hence, we should discover the relevant features of each diagnosis a (this may be different for each diagnosis).

The recommended diagnosis based on the relevant features extracted from the FNA images is determined as:

$$a^*(\mathbf{x}) \triangleq \arg \max_a \mathbb{E}_{y|a, \mathbf{x}_{\mathcal{R}(a)}}(y|a, \mathbf{x}_{\mathcal{R}(a)})$$

where $a^*(\mathbf{x})$ is the diagnosis that yields the highest prediction accuracy for a patient whose imaging is characterized by the information vector \mathbf{x} .

4. ALGORITHMS

4.1. Diagnosis-Relevant Feature Selection (DiReFS)

The proposed DiReFS algorithm sequentially discovers the relevant features which yield maximum relevance to the specific diagnosis with minimum redundancy (compared with the previously discovered relevant features).

To describe DiReFS, we start by introducing a few notations. Let $\hat{y}_a^S(\mathbf{x}_S)$ and $N_a^S(\mathbf{x}_S)$ be the sample mean prediction accuracy estimator and the number of patients (whose feature information contains \mathbf{x}_S and was provided

diagnosis a). Let \hat{y}_a and N_a be the sample mean prediction accuracy estimator and the number of patients who received the diagnosis a .

First, we define a relevance metric $h_f^r(a)$ which measures how the expected diagnosis accuracy for patients having the feature x_f differs from that obtained for the entire set of patients in \mathcal{H}_N (previously defined in section 3) when diagnosis a is chosen. We formalize this as:

$$h_f^r(a) \triangleq \sum_{x_f} \frac{N_a^f(x_f)}{N_a} |\hat{y}_a^f(x_f) - \hat{y}_a|$$

Second, we define a redundancy metric $h_{f,s}^d(a)$ which measures how the expected diagnosis accuracy made for a patient is affected by considering an additional feature x_s when diagnosis a is chosen. We formalize this as:

$$h_{f,s}^d(a) = - \sum_{x_f, x_s} \frac{N_a^{f,s}(x_f, x_s)}{N_a(x_s)} |\hat{y}_a^{f,s}(x_f, x_s) - \hat{y}_a^s(x_s)|$$

Next, we use the minimum-redundancy-maximum-relevance (mRMR) criterion [15] to combine the above metrics to select diagnosis-relevant features. Before we describe DiReFS, let us define $\mathcal{U}_f(a)$ as the utility obtained if feature x_f is selected as a relevant feature for diagnosis a . If $\hat{\mathcal{R}}(a)$ is the relevant features set discovered by DiReFS for diagnosis a , the utility $\mathcal{U}_f(a)$ is determined as:

$$\mathcal{U}_f(a) = h_f^r(a) - \frac{1}{|\hat{\mathcal{R}}(a)|} \sum_{s \in \hat{\mathcal{R}}(a)} h_{f,s}^d(a),$$

where $1/|\hat{\mathcal{R}}(a)|$ is used as a normalization factor.

The main steps of the DiReFS are outlined below:

Step 1: Define $\hat{\mathcal{R}}(a)$ as the relevant feature set discovered by DiReFS for diagnosis a and $\hat{\mathcal{R}}^c(a)$ as the complementary set of $\hat{\mathcal{R}}(a)$. G and H are defined as selected relevant features in step 2 and step 3, respectively. For each diagnosis a , initialize $\hat{\mathcal{R}}(a)$ as the empty set (i.e. \emptyset) and $\hat{\mathcal{R}}^c(a)$ as the set of all features.

Step 2: The algorithm selects the first relevant feature that maximizes the relevance metric ($h_f^r(a)$), i.e.,

$$G = \arg \max_{f \in \hat{\mathcal{R}}^c(a)} h_f^r(a)$$

$$\hat{\mathcal{R}}(a) = \hat{\mathcal{R}}(a) \cup G$$

Step 3: The algorithm finds the subsequent relevant feature that maximizes utility function ($\mathcal{U}_f(a)$), i.e.,

$$H = \arg \max_{f \in \hat{\mathcal{R}}^c(a)} \mathcal{U}_f(a)$$

$$\hat{\mathcal{R}}(a) = \hat{\mathcal{R}}(a) \cup H$$

Step 4: The algorithm iteratively runs Step 3 until m -th relevant feature is selected, where m is an input parameter for the algorithm.

4.2. Diagnosis recommendation algorithm

The proposed diagnosis recommendation algorithm is a modified contextual multi-armed bandit algorithm [16-19] which uses the contexts (features) selected by DiReFS to recommend the optimal diagnosis for each patient. The main steps of the recommendation engine are outlined below:

Step 1: Find the set of underexplored actions for the patient with information vector $\mathbf{x}_{\hat{\mathcal{R}}(a)}$.

$$U = \{a \in A \mid N_a^{\hat{\mathcal{R}}(a)}(\mathbf{x}_{\hat{\mathcal{R}}(a)}) < C \cdot \log(n)\}$$

where $C \cdot \log(n)$ is a control function. If there are underexplored actions, DE abstains from making diagnosis recommendation and only updates $N_a^{\hat{\mathcal{R}}(a)}(\mathbf{x}_{\hat{\mathcal{R}}(a)})$ and $\hat{y}_a^{\hat{\mathcal{R}}(a)}(\mathbf{x}_{\hat{\mathcal{R}}(a)})$ based on the actual label (benign/malignant) obtained from post-examination. Hence, note that DE only issues recommendations when it is sufficiently confident about its predictions and it abstains otherwise.

Step 2: If there is no underexplored actions for the patient with information vector $\mathbf{x}_{\hat{\mathcal{R}}(a)}$, the optimal diagnosis with respect to the relevant feature set $\hat{\mathcal{R}}(a)$ is determined as

$$\hat{a}(\mathbf{x}) = \arg \max_a \hat{y}_a^{\hat{\mathcal{R}}(a)}(\mathbf{x}_{\hat{\mathcal{R}}(a)})$$

This optimization selects the action with the highest estimated prediction accuracy for the patient with information vector $\mathbf{x}_{\hat{\mathcal{R}}(a)}$. The pseudo-code of DE is given in Algorithm 1.

Algorithm 1 Diagnosis Engine (DE)

Input: m, C

Initialize: $\hat{\mathcal{R}}(a) = \emptyset$, $\hat{\mathcal{R}}^c(a) = \{f_1, f_2, \dots, f_D\}$ for each a
for each diagnosis a

$$G = \arg \max_{f \in \hat{\mathcal{R}}^c(a)} h_f^r(a)$$

$$\hat{\mathcal{R}}(a) = \hat{\mathcal{R}}(a) \cup G$$

do

$$H = \arg \max_{f \in \hat{\mathcal{R}}^c(a)} \mathcal{U}_f(a)$$

$$\hat{\mathcal{R}}(a) = \hat{\mathcal{R}}(a) \cup H$$

while ($|\hat{\mathcal{R}}(a)| < m$)

end for

$$U = \{a \in A \mid N_a^{\hat{\mathcal{R}}(a)}(\mathbf{x}_{\hat{\mathcal{R}}(a)}) < C \cdot \log(n)\}$$

if ($U = \emptyset$)

$$\hat{a}(\mathbf{x}) = \arg \max_a \hat{y}_a^{\hat{\mathcal{R}}(a)}(\mathbf{x}_{\hat{\mathcal{R}}(a)})$$

end if

Update $N_a^{\hat{\mathcal{R}}(a)}(\mathbf{x}_{\hat{\mathcal{R}}(a)})$, $\hat{y}_a^{\hat{\mathcal{R}}(a)}(\mathbf{x}_{\hat{\mathcal{R}}(a)})$ based on the actual label (benign/malignant)

5. EXPERIMENTS

In this section we evaluate the performance of DE for breast cancer diagnosis using the well-known UCI dataset [20]. The dataset contains 30 patient features extracted from FNA images. The diagnosis (label) for each patient is either malignant or benign.

We compare the performance of DE algorithms with four existing machine learning algorithms and three existing feature selection algorithms:

- Logistic Regression (LogitR);
- Linear Regression (LinearR);
- Support Vector Machines (SVMs); we use a radial basis function (RBF) kernel SVM;

- Adaptive Boosting (AdaBoost);
- Correlation Feature Selection (CFS): a well-known feature selection algorithm based on correlation [14];
- Mutual Information Feature Selection (MIFS): a well-known feature selection algorithm based on mutual information [15];
- Relevance Learning with Feedback (RELEAF): an action dependent relevance learning algorithm based on the expected rewards [16, 17];

5.1. Simulation Setup

First, we compare our DE algorithm against state-of-the-art machine learning algorithms: LogitR, LinearR, SVM and AdaBoost. The training set contains 10% of the patients in the dataset and standard 50-fold stratified cross-validation was applied in the simulation.

Second, to highlight the importance of DiReFS, we performed two additional sets of simulations. In the first set, we compare the performance of our DE system using DiReFS with the performance of the DE system where DiReFS was replaced with one of the three different feature selection algorithms: CFS [14], MIFS [15], and RELEAF [16-17]. This comparison shows the impact of DiReFS on the overall performance of the DE.

In the second set of simulations, we use the features selected by DiReFS in conjunction with the diagnosis recommendation made by the benchmark algorithms - linear regression, logistic regression, SVM - to highlight the specific impact of our feature selection algorithm.

5.2. Measuring Success

Given a patient, DE as well as the other benchmark algorithms classify the tumors as malignant or benign. To quantify the performance, we apply three performance metrics: the prediction error rate (*PER*), the false positive rate (*FPR*), and the false negative rate (*FNR*). *PER* is defined as the fraction of times the classification of our algorithm is different from the actual label. *FPR* and *FNR* are defined as the diagnosis error rate for benign tumors and the diagnosis error rate for malignant tumors, respectively. The goal of DDSS is to minimize the FPR given an allowable threshold for FNR as selected by the clinicians. (In practice, this is often set to be below 2% [21].)

Comparison with machine learning algorithms: As the table 1 shows, our DE algorithm has 2.23% prediction error rates and 2.62% false positive rates which is 7.47% and 13.25% better than the second best algorithm (LogitR) when the tolerable threshold of FNR is set to below 2%. There are two reasons for the outstanding performance of the DE algorithm. First, our diagnosis recommendation algorithm yields high accuracy for classification, because it is able to provide personalized diagnosis, while other comparable

Table 1: Comparison with typical machine learning algorithms

%	DE	LogitR	LinearR	SVMs	AdaBoost
PER	2.23	9.70	35.44	10.15	11.94
FPR	2.62	15.87	45.19	16.22	18.35
FNR	1.92	1.94	1.96	1.98	1.99

algorithms apply the same model for all patients. Second, DE can discover different relevant features for different diagnosis based on DiReFS, while the other algorithms base their decisions on all the features.

Comparison with feature selection algorithms: In this subsection, we demonstrate the impact of DiReFS algorithm on the DE system. We compare the performance of the DE using DiReFS with the performance of DE using different feature selection algorithms. As it can be seen in table 2, DiReFS significantly outperforms all the other feature selection algorithms when the tolerable threshold of FNA is set to below 2%. This is because DiReFS is capable of discovering diagnosis relevant features based on their impact on the expected diagnosis accuracies. Although RELEAF also considers the dependence between diagnosis and feature selection, it is extremely slow and not able to exploit the redundancy existing among features.

Table 2: Performance of DE with other feature selection methods

%	DiReFS	RELEAF	CFS	MIFS
PER	2.23	18.37	5.69	9.97
FPR	2.62	24.11	9.81	16.4
FNR	1.92	1.96	1.98	1.89

Next, we replace the recommendation part of DE with conventional machine learning algorithms and demonstrate the importance of DiReFS when used for diagnosis decisions in conjunction with such alternative decision methods. As it can be seen in table 3, DiReFS is capable of improving the performance of all the benchmark algorithms because it is able to discover and selecting different relevant features for different diagnosis.

Table 3: Impact of the DiReFS in conjunction with alternative machine learning algorithms

	PER (%)		FPR (%)	
	DiReFS	w/o DiReFS	DiReFS	w/o DiReFS
Linear R	21.30	35.44	26.31	45.19
Logit R	6.32	9.70	10.28	15.87
SVMs	6.51	10.15	10.76	16.22

6. CONCLUSION

We describe a Diagnosis Engine (DE) which uses past patients' information (medical tests and diagnosis) to discover the relevant features extracted from images of cellular samples obtained from FNA of breast mass and uses these features to provide personalized diagnosis for the current patient. When applied to a well-known breast cancer dataset, our results demonstrate that DE is capable of significantly outperforming (by 13.25%) existing techniques in terms of false positive rates. *This improvement is extremely important because it saves numerous patients unnecessary distress and saves spending on unnecessary treatments.* We also show that our diagnosis-relevant feature selection, DiReFS, can be applied in conjunction with other machine learning algorithms to significantly improve their performance by discovering *different* features that are relevant to *different* diagnosis.

7. REFERENCES

- [1] M. Frellick, *Landmark Report Urges Reform to Avert Diagnostic Errors*, Medscape, 2015.
- [2] E.S. Berner, *Clinical Decision Support Systems*, Springer , New York, 2007.
- [3] E. Çomak, A. Arslan, and I. Türkoğlu, “A decision support system based on support vector machines for diagnosis of the heart valve diseases,” *Computers in Biology and Medicine*, 37(1), pp. 21-27, 2007.
- [4] D. L. Hunt, R. B. Haynes, S. E. Hanna, and K. Smith, “Effects of computer-based clinical decision support systems on physician performance and patient outcomes: a systematic review,” *JAMA*, 280(15), pp. 1339-1346, 1998.
- [5] R. Bellazzi, and B. Zupan, “Predictive data mining in clinical medicine: current issues and guidelines,” *International journal of medical informatics*, 77(2), pp. 81-97, 2008.
- [6] S. Balakrishnan, R. Narayanaswamy, N. Savarimuthu, and R. Samikannu, “SVM ranking with backward search for feature selection in type II diabetes databases,” *In Systems, Man and Cybernetics, IEEE International Conference* pp. 2628-2633, 2008.
- [7] A. X. Garg, N. K. Adhikari, H. McDonald, M. P. Rosas-Arellano, P. J. Devereaux, J. Beyene, and R. B. Haynes, “Effects of computerized clinical decision support systems on practitioner performance and patient outcomes: a systematic review,” *JAMA*, 293(10), pp. 1223-1238, 2005.
- [8] R. Ahmed, A. Temko, W. Marnane, G. Boylan, and G. Lightbody, “Grading brain injury in neonatal EEG using SVM and supervector kernel,” *In Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference* pp. 5894-5898, 2014.
- [9] K. Polat, and S. Güneş, “Breast cancer diagnosis using least square support vector machine,” *Digital Signal Processing*, 17(4), pp. 694-701, 2007
- [10] L. Song, W. Hsu, J. Xu, and M. van der Schaar, M, “Using Contextual Learning to Improve Diagnostic Accuracy: Application in Breast Cancer Screening”, *IEEE J. Biomedical and Health Informatics*, 2015
- [11] J. Xu, D. Sow, D. Turaga, and M. van der Schaar, “Online Transfer Learning for Differential Diagnosis Determination,” *AAAI Workshop on the World Wide Web and Public Health Intelligence*, 2015
- [12] P. Hall, J. S. Marron, and A. Neeman, “Geometric representation of high dimension, low sample size data,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(3), pp. 427-444. 2005.
- [13] A. L. Blum, and P. Langley, “Selection of relevant features and examples in machine learning,” *Artificial intelligence*, 97(1), pp. 245-271, 1997.
- [14] M. A. Hall, “Correlation-based feature selection for machine learning,” *Doctoral dissertation*, The University of Waikato, 1999
- [15] H. Peng, F. Long, and C. Ding, “Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy,” *Pattern Analysis and Machine Intelligence, IEEE Transactions*, 27(8), pp. 1226-1238, 2005.
- [16] C. Tekin, and M. van der Schaar, “Discovering, learning and exploiting relevance,” *In Advances in Neural Information Processing Systems*, pp. 1233-1241, 2014.
- [17] C. Tekin, and M. van der Schaar, “RELEAF: An algorithm for learning and exploiting relevance,” *IEEE Journal of Selected Topics in Signal Processing, Special Issue on Signal Processing for Big Data*, vol. 9, no. 4, pp. 716-727, 2015.
- [18] A. Slivkins, “Contextual bandits with similarity information,” *The Journal of Machine Learning Research*, 15(1), 2533-2568, 2014.
- [19] A. Badanidiyuru, J. Langford, and A. Slivkins, “Resourceful contextual bandits,” *The 27th Conference on Learning Theory*, pp. 1109-1134, 2014.
- [20] K. Bache and M. Lichman, “UCI machine learning repository,” 2013. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [21] C. I. Flowers, C. O’Donoghue, D. Moore, A. Goss, D. Kim, and J.H. Kim, “Reducing false-positive biopsies: a pilot study to reduce benign biopsy rates for BI-RADS 4A/B assessments through testing risk stratification and new thresholds for intervention,” *Breast Cancer Research and Treatment*, vol. 139, no. 3, pp. 769-777, 2013.
- [22] A. Lanata, A. Greco, G. Valenza, and E. P. Scilingo, “A pattern recognition approach based on electrodermal response for pathological mood identification in bipolar disorders,” *In Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference*, pp. 3601-3605, 2014.