

Online Discriminative Kernel Density Estimation

Matej Kristan, Aleš Leonardis
 Faculty of Computer and Information Science
 University of Ljubljana
 Ljubljana, Slovenia
 {matej.kristan},{ales.leonardis}@fri.uni-lj.si

Abstract—We propose a new method for online estimation of probabilistic discriminative models. The method is based on the recently proposed online Kernel Density Estimation (oKDE) framework which produces Gaussian mixture models and allows adaptation using only a single data point at a time. The oKDE builds reconstructive models from the data, and we extend it to take into account the interclass discrimination through a new distance function between the classifiers. We arrive at an online discriminative Kernel Density Estimator (odKDE). We compare the odKDE to oKDE, batch state-of-the-art KDEs and support vector machine (SVM) on a standard database. The odKDE achieves comparable classification performance to that of best batch KDEs and SVM, while allowing online adaptation, and produces models of lower complexity than the oKDE.

Index Terms—Online Estimation; Discriminative Models; Kernel Density Estimation

I. INTRODUCTION

Building discriminative models of some process from the observed data is a central task of many applications in machine learning. A popular approach to generating models is to estimate the probability density function (pdf) associated with the observed data. In this respect, *reconstructive* models such as the Gaussian mixture models, (GMM), (e.g., [1]) have been successfully applied in *batch* operation, i.e., in situations in which all the data is observed in advance. In contrast to the reconstructive models, the discriminative models capitalize on the discriminative information, however, this may lead to decreased robustness [2]. A significant drawback of the purely batch methods is that their estimation becomes increasingly difficult when processing extremely large amounts of data. Furthermore, in real-world environments, all the data may not be available in advance, or we even want to observe some process for an indefinite duration, while continually providing the best estimate of the model from the data observed so far. This generates the need for models that can be constructed in an *online* operation.

Adapting the existing reconstructive GMM methods to work with online cases, in which as little as a single data-point may be observed at a time, is a nontrivial task. In contrast to the batch incremental models (e.g., [3]) who store and revisit all the data in multiple passes, the online models have to adapt from a (single) new data-point and then discard that data-point. The main difficulty is therefore that the online models have to maintain sufficient information to generalize well to the yet unobserved data and have to adjust their complexity without having access to all the observations (future as well

as past). There have been various attempts to extend the *reconstructive* GMMs to online operation, however, these either imply strong spatio-temporal constraints on the data [4], [5], assume constraints on the shape of the target distribution [6] or require tuning of parameters to a specific application [7]. Recently, we have proposed a non-parametric approach called the *online Kernel Density Estimation* (oKDE) [8]. In contrast to the other approaches, the oKDE does not impose any of the above constraints but assumes only that the target pdf is sufficiently smooth and produces models with a high reconstructive performance. In [9] we have also considered a variant of the oKDE that allows adaptation from positive as well as negative examples.

While the purely reconstructive models may contain redundant information required for discrimination, the discriminative models disregard the reconstructive information required for online adaptation. Indeed, Fidler et al. [2] have shown that even in batch methods accounting for the reconstructive information leads to improved robustness of the discriminative models. Following their results, we adapt the oKDE framework to account for the discriminative power of the models along with the reconstructive, thus arriving at an *online discriminative Kernel Density Estimator*, which is the main contribution of the paper. The proposed method allows online adaptation of the discriminative models, by maintaining enough reconstructive power to efficiently adapt to new observations, and can be used to develop online classifiers. The remainder of the paper is structured as follows. In Section II we briefly review the oKDE framework, in Section III we extend the oKDE to discriminative models, in Section IV we evaluate the approach and the Section V concludes the paper.

II. THE ONLINE KERNEL DENSITY ESTIMATION

The online Kernel Density Estimation (oKDE) produces a generative model from the d -dimensional streaming data as an N -component Gaussian mixture model

$$p(\mathbf{x}) = \sum_{i=1}^N w_i \phi_{\Sigma_i}(\mathbf{x} - \mathbf{x}_i), \quad (1)$$

where $\phi_{\Sigma}(\mathbf{x} - \mu) = (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$ is a Gaussian kernel centered at μ with covariance matrix Σ . We give here only a brief overview of the oKDE framework and refer the reader to [8] for more details.

Broadly speaking, the oKDE proceeds in two steps:

Update: Starting with the GMM from the previous time-step $p_{t-1}(\mathbf{x})$ and the new observation \mathbf{x}_t , the oKDE augments the $p_{t-1}(\mathbf{x})$ by a Gaussian kernel $\phi_{\Sigma_t}(\mathbf{x} - \mathbf{x}_t)$ centered at \mathbf{x}_t , it automatically calculates the optimal covariance Σ_t for that kernel and accordingly readjusts the covariances of the existing kernels using the multivariate online plug-in rule proposed in [8]. If required, it also refines the GMM by splitting up some of the components.

Compress: To maintain a low complexity (i.e., low number of components) the GMM is simplified from time to time. The oKDE generates a binary tree among the components and iteratively merges pairs of components until some threshold of a cost function is exceeded. In [8], the cost function represents the distance between the distribution before and after the compression, which penalizes the reconstruction errors.

III. DISCRIMINATIVE KDE

We pose the online discriminative learning as a task to estimate from a stream of data a set of K discriminative classes, each class c_i described by a Gaussian mixture model $p(\mathbf{x}|c_i)$ and a prior probability $p(c_i)$. In principle, we could use the oKDE to construct each of these classes, but, due to its reconstructive nature, the produced models will be likely redundant for classification. Indeed, we require the models to contain just so much of the information to prevent a degraded classification. Recall that the oKDE simplifies the models under a certain cost function which measures the reconstruction error induced by compression [8]. This means that by redefining this cost function to rather take into account the classification error, the compression step in the oKDE will lead to models that *reduce their complexity while retaining their discriminative power*.

Assume that we want to compress the c_i -th class mixture model $p(\mathbf{x}|c_i)$ into $p_{\text{cmp}}(c_i|\mathbf{x})$, while minimizing the induced classification errors. First we have to rewrite this model into a *classification model*. We consider the class c_i as a *positive example* class C^+ , described by a mixture model $p(C^+|\mathbf{x}) = p(c_i|\mathbf{x})$. Then we collect *all the other classes* to form a single *negative example* class C^- , $p(C^-|\mathbf{x}) = \sum_{j \neq i} p(c_j|\mathbf{x})$. The posterior over the resulting two-class model is then defined as

$$p(C|\mathbf{x}) = \delta_{C^+}(C)p(C^+|\mathbf{x}) + \delta_{C^-}(C)p(C^-|\mathbf{x}), \quad (2)$$

where $\delta_{C^*}(C)$ is a Dirac function centered at C^* . The compressed counterpart of the posterior (2), is obtained by setting $p_{\text{cmp}}(C^+|\mathbf{x}) = p_{\text{cmp}}(c_i|\mathbf{x})$:

$$p_{\text{cmp}}(C|\mathbf{x}) = \delta_{C^+}(C)p_{\text{cmp}}(C^+|\mathbf{x}) + \delta_{C^-}(C)p(C^-|\mathbf{x}). \quad (3)$$

From the classification point of view we can say that $p(\mathbf{x}|c_i)$ can be compressed into $p_{\text{cmp}}(\mathbf{x}|c_i)$ as long as the distance between the corresponding posteriors $p(C|\mathbf{x})$ and $p_{\text{cmp}}(C|\mathbf{x})$, does not change significantly. We therefore require a distance measure between the posterior before and after compression.

A. Distance between two classifiers

We define the distance between the posterior $p(C|\mathbf{x})$ and its compression $p_{\text{cmp}}(C|\mathbf{x})$, given some value of \mathbf{x} , using the Hellinger distance [10],

$$\begin{aligned} D^2(p, p_{\text{cmp}}|\mathbf{x}) &\triangleq \frac{1}{2} \int_C (p(C|\mathbf{x})^{\frac{1}{2}} - p_{\text{cmp}}(C|\mathbf{x})^{\frac{1}{2}})^2 \\ &= \frac{1}{2} \sum_{C \in [C^+, C^-]} (p(C|\mathbf{x})^{\frac{1}{2}} - p_{\text{cmp}}(C|\mathbf{x})^{\frac{1}{2}})^2. \end{aligned} \quad (4)$$

Integrating (4) over the relevant feature space \mathbf{x} gives the *expected* Hellinger distance

$$\hat{D}^2(p, p_{\text{cmp}}) = \int D^2(p, p_{\text{cmp}}|\mathbf{x})p_0(\mathbf{x})d\mathbf{x}, \quad (5)$$

where the expectation is calculated over the distribution $p_0(\mathbf{x}) = p(\mathbf{x}|C^+)p(C^+) + p(\mathbf{x}|C^-)p(C^-)$, with the priors $P(C^+) = p(c_i)$ and $P(C^-) = \sum_{j \neq i} p(c_j)$. In our case, $p_0(\mathbf{x})$ can be written in the form of a Gaussian mixture model $p_0(\mathbf{x}) = \sum_{i=1}^M w_i \phi_{\Sigma_i}(\mathbf{x} - \mathbf{x}_i)$ and (5) becomes

$$\hat{D}^2(p, p_{\text{cmp}}) = \sum_{i=1}^M w_i \int D^2(p, p_{\text{cmp}}|\mathbf{x}) \phi_{\Sigma_i}(\mathbf{x} - \mathbf{x}_i) d\mathbf{x}. \quad (6)$$

Note that while $D(\cdot, \cdot|\mathbf{x})$ is a metric, constrained to the interval $[0, 1]$, it is a nonlinear function of \mathbf{x} , and the integrals in (6) cannot be evaluated analytically. However, they can be numerically approximated using the *unscented transform*, which has been proposed by [11] for calculating nonlinear transformations of Gaussian variables. Similarly to a Monte Carlo integration, the *unscented transform* relies on evaluating integrals using carefully placed points, called *the sigma points*, over the support of the integral. Therefore, (6) is approximated as

$$\hat{D}^2(p, p_{\text{cmp}}) \approx \sum_{i=1}^M w_i \sum_{j=0}^{2d+1} D^2(p, p_{\text{cmp}}|^{(j)}\mathcal{X}_i)^{(j)}\mathcal{W}_i, \quad (7)$$

where $\{^{(j)}\mathcal{X}_i, ^{(j)}\mathcal{W}_i\}_{j=0:d}$ are the weighted sets of sigma points corresponding to the i -th Gaussian $\phi_{\Sigma_i}(\mathbf{x} - \mathbf{x}_i)$, and are defined as

$$\begin{aligned} ^{(0)}\mathcal{X}_i &= \mathbf{x}_i; \quad ^{(0)}\mathcal{W}_i = \frac{\kappa}{1 + \kappa} \\ ^{(j)}\mathcal{X}_i &= \mathbf{x}_i + s_j \sqrt{1 + \kappa} (\sqrt{d\Sigma_i})_j; \\ ^{(j)}\mathcal{W}_i &= \frac{\kappa}{2(1 + \kappa)}; \quad s_j = \begin{cases} 1 & ; \quad j \leq d \\ -1 & ; \quad \text{otherwise} \end{cases} \end{aligned} \quad (8)$$

with $\kappa = \max([0, m - d])$, and $(\sqrt{d\Sigma_i})_j$ is the j -th column of the matrix square root of Σ_i . Specifically, let $\mathbf{U}\mathbf{D}\mathbf{U}^T$ be a singular value decomposition of covariance matrix Σ , such that $\mathbf{U} = \{U_1, \dots, U_d\}$ and $\mathbf{D} = \text{diag}\{\lambda_1, \dots, \lambda_d\}$, then $(\sqrt{\Sigma})_k = \sqrt{\lambda_k}U_k$. In line with the discussion on the properties of the unscented transform in [11], we set the parameter m to $m = 3$.

B. The online discriminative KDE

The distance function $\hat{D}(p, p_{\text{cmp}})$ from Section III-A penalizes any change of the classification posterior which is induced through the compression of the mixture models. Small values, i.e., $\hat{D}(p, p_{\text{cmp}}) \approx 0$, mean that the classification does not change, while $\hat{D}(p, p_{\text{cmp}}) = 1$ implies a complete change. In an online operation, we can therefore simply adapt the oKDE (Section II) to build and compress a mixture model for each class separately from the observed data. To penalize the loss of discrimination during the compression, we can use the distance measure $\hat{D}(p, p_{\text{cmp}})$ as the compression cost function. We arrive at an online discriminative KDE (odKDE). In practice, we let the odKDE compress the model until the cost function exceeds some (small) threshold \hat{D}_{th} .

During learning, at time-step $t - 1$, we have a set of K models $\{p_{t-1}(\mathbf{x}|c_i), p_{t-1}(c_i)\}_{i=1:K}$. For simplicity assume that at time-step t , K new observations $\{\mathbf{z}_i\}_{i=1:K}$, one per each class¹, arrive and the models are updated into $\{p_t(\mathbf{x}|c_i), p_t(c_i)\}_{i=1:K}$. A single time-step iteration of the approach is outlined in Algorithm 1. In the classification phase a new observation \mathbf{z} is classified into a class \hat{c} by applying the Bayesian rule

$$\hat{c} = \arg \max_{c_i} p(\mathbf{z}|c_i)p(c_i). \quad (9)$$

Algorithm 1 : The online discriminative KDE

Input:

$\{p_{t-1}(\mathbf{x}|c_i), p_{t-1}(c_i)\}_{i=1:K}$... the input models.
 $\{\mathbf{z}_i\}_{i=1:K}$... observations (one per each class)

Output:

$\{p_t(\mathbf{x}|c_i), p_t(c_i)\}_{i=1:K}$... the output models.

Procedure:

- 1: **for** $i = 1 : K$ **do**
 - 2: Update $p_{t-1}(\mathbf{x}|c_i)$ with \mathbf{z}_i into $\tilde{p}_t(\mathbf{x}|c_i)$ using the original update step of oKDE (Section II).
 - 3: Update the prior $p_t(c_i)$.
 - 4: **end for**
 - 5: **for** $i = 1 : K$ **do**
 - 6: Construct the two-class classification model (2) by treating $\tilde{p}_t(\mathbf{x}|c_i)$ as a positive example and the rest $K-1$ models as the negative example.
 - 7: Compress $\tilde{p}_t(\mathbf{x}|c_i)$ into $p_t(\mathbf{x}|c_i)$ by hierarchical merging components such that $\hat{D}(\tilde{p}_t, p_t) \leq \hat{D}_{\text{th}}$.
 - 8: **end for**
-

IV. EXPERIMENTS

We have compared the classification performance of the odKDE with the online reconstructive KDE, oKDE [8], and three state-of-the-art batch KDEs: the cross-validation (CV) KDE [12], the reduced-set density estimator [13] (RSDE) initialized by the CV, and the Hall's KDE [14] (Hall). For the

¹This restriction serves only for clarity of the presentation. Note that, in general, our approach can also handle cases in which observations come only from a subset of classes at a time.

baseline classification, we have applied a multiclass support vector machine (SVM) with an RBF kernel [15]. The methods were compared on a set of public classification problems [16] (Table I). In all experiments, the distance parameter in the odKDE was set to $\hat{D}_{\text{th}} = 0.005$.

The odKDE and oKDE were initialized for each class using the first 10 samples and the rest were added one at a time; this experiment was repeated via four-fold cross validation and for three random data orderings. This amounted to twelve repetitions per dataset. The parameter for the SVM kernel was determined separately in each experiment via cross validation on the training dataset. Table I shows the classification score and the number of components in the models after observing all the samples, while Fig. 1 shows the evolution of the results with respect to (w.r.t.) the number of observations for the oKDE and the doKDE. For reference, the graphs also show results for the batch methods after observing all the samples. From the results in Table I we see that the odKDE generated models with comparable classification performance as the oKDE, but generally with a significantly lower complexity. We can verify that this was also true during the online estimation from Fig. 1. Although the odKDE was learnt only by observing a *single* example at a time, the resulting models exhibit classification performance similar to the best batch KDE approaches and the SVM, who optimized their structure having access to *all* the data. By further inspection of the results we can also observe a general trend that the number of components initially significantly increases in doKDE (as well as in oKDE), but then stabilizes for larger number of samples, while the recognition score further improves. For example, in the case of the *Letter* dataset, the bound on the complexity is reached relatively early on after observing 200 samples per class (i.e., after 5200th sample), while the classification performance further increased through the model refinement. This makes the odKDE a very appropriate tool for online operation since it produces compressed models with good classification performance, while at the same time maintains sufficient reconstructive information to allow online refinements of the models from new observations.

V. CONCLUSION

We have proposed an approach for online estimation of discriminative models by adapting the framework of online Kernel Density Estimation. We have defined a distance measure which measures the discrimination of models and used this measure in the oKDE as a cost function for compression. Results demonstrate that the proposed odKDE produces comparable classification performance to the state-of-the-art, and produces models of significantly lower complexity while allowing online adaptation. This makes the approach ideal for online estimation of classifiers from streaming data. In our future work we will study how different distance measures and data orderings influence the performance of the proposed method as well as test how the method handles noise in labels.

TABLE I
AVERAGE CLASSIFICATION RESULTS ALONG WITH \pm ONE STANDARD DEVIATION. THE NUMBER OF SAMPLES IN EACH DATASET, THE DIMENSIONALITY AND THE NUMBER OF CLASSES ARE DENOTED BY N_S , N_D AND N_C , RESPECTIVELY.

dataset	N_S	N_D	N_C	Recognition accuracy [%] (Number of components per class)					
				odKDE	oKDE	CV	RSDE	Hall	SVM
Iris	150	4	3	97 \pm 4%(6.3 \pm 1)	97 \pm 3%(31 \pm 0)	96 \pm 3%(38 \pm 0)	96 \pm 2%(10 \pm 5)	97 \pm 4%(38 \pm 0)	96 \pm 3%(16 \pm 1)
Pima	768	8	2	71 \pm 3%(108 \pm 6)	70 \pm 1%(162 \pm 3)	72 \pm 2%(288 \pm 0)	65 \pm 3%(48 \pm 10)	67 \pm 2%(288 \pm 0)	78 \pm 3%(160 \pm 4)
Wine	178	13	3	97 \pm 2%(2.5 \pm 1)	99 \pm 2%(45 \pm 0)	92 \pm 4%(45 \pm 0)	94 \pm 4%(44 \pm 0)	99 \pm 2%(45 \pm 0)	98 \pm 2%(22 \pm 4)
Letter	20000	16	26	94 \pm 0%(16 \pm 1)	95 \pm 0%(222 \pm 2)	96 \pm 0%(613 \pm 0)	55 \pm 0%(25 \pm 0)	95 \pm 0%(613 \pm 0)	96 \pm 0%(322 \pm 0)

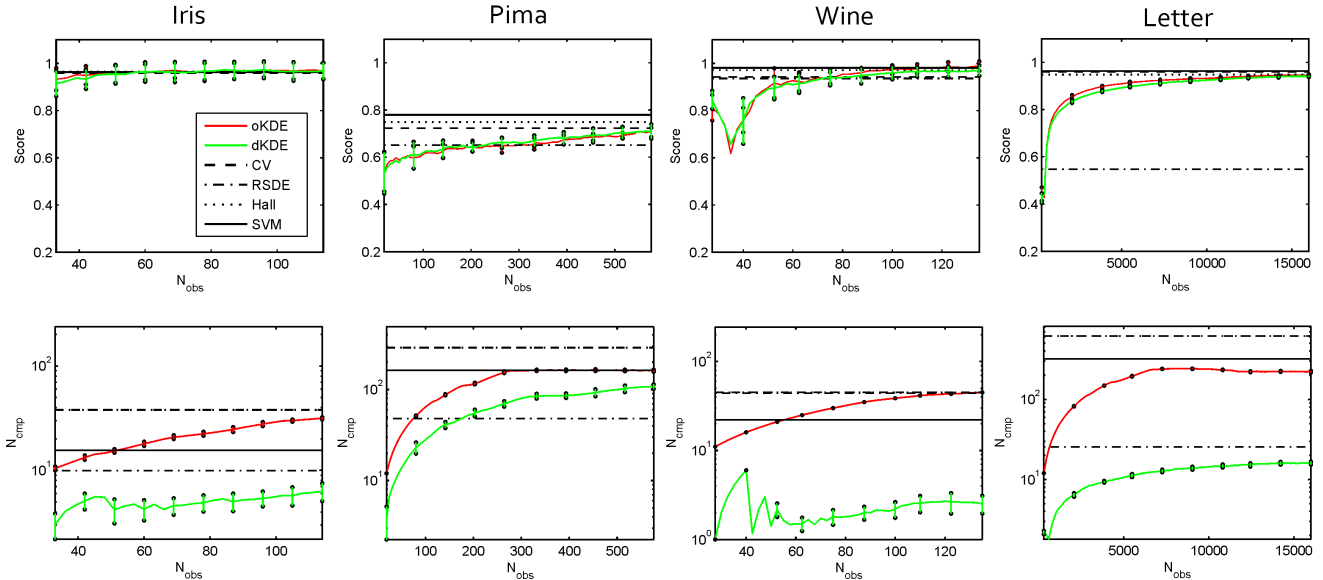


Fig. 1. Upper row shows the classification results (Score) and the lower shows the number of components per class (N_{cmp}) w.r.t. the number of samples (N_{obs}). The results for the oKDE and odKDE are depicted by darker (red) line and bright (green) line, respectively along with one standard deviation bars. For reference, we also show results for the batch methods after observing all samples.

ACKNOWLEDGMENT

This research was supported by: RP P2-0214 and P2-0094 (RS), ARRS project "Learning a large number of visual object categories for content-based retrieval in image and video databases", and EU FP7-ICT215181-IP project CogX.

REFERENCES

- [1] M. A. F. Figueiredo and A. K. Jain, "Unsupervised learning of finite mixture models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 3, pp. 381–396, 2002.
- [2] S. Fidler, D. Skočaj, and A. Leonardis, "Combining reconstructive and discriminative subspace methods for robust classification and regression by subsampling," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 3, pp. 337–350, March 2006. [Online]. Available: <http://vicos.fri.uni-lj.si/data/publications/fidlerPAMI06.pdf>
- [3] D. Mansjurić and B. Juang, "Incremental learning of mixture models for simultaneous estimation of class distribution and inter-class decision boundaries," in *International Conference on Pattern Recognition*, 2008, pp. 1–4.
- [4] O. Arandjelovic and R. Cipolla, "Incremental learning of temporally-coherent gaussian mixture models," in *British Machine Vision Conference*, 2005, pp. 759–768.
- [5] M. Song and H. Wang, "Highly efficient incremental estimation of Gaussian mixture models for online data stream clustering," in *SPIE: Intelligent Computing: Theory and Applications*, 2005, pp. 174–183.
- [6] B. Han, D. Comaniciu, Y. Zhu, and L. S. Davis, "Sequential kernel density approximation and its application to real-time visual tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 7, pp. 1186–1197, 2008.
- [7] W. F. Szewczyk, "Time-evolving adaptive mixtures," National Security Agency, Tech. Rep., 2005.
- [8] M. Kristan and A. Leonardis, "Multivariate online kernel density estimation," in *Computer Vision Winter Workshop*, 2010, pp. 77–84. [Online]. Available: <http://vicos.fri.uni-lj.si/data/publications/KristanCVWW2010.pdf>
- [9] M. Kristan, D. Skočaj, and A. Leonardis, "Online kernel density estimation for interactive learning," *Image and Vision Computing*, vol. 28, no. 7, pp. 1106–1116, 2010.
- [10] D. E. Pollard, *A user's guide to measure theoretic probability*. Cambridge University Press, 2002.
- [11] S. Julier and J. Uhlmann, "A general method for approximating nonlinear transformations of probability distributions," Department of Engineering Science, University of Oxford, Tech. Rep., 1996.
- [12] J. M. L. Murillo and A. A. Rodriguez, "Algorithms for gaussian bandwidth selection in kernel density estimators," in *Neural Inf. Proc. Systems*, 2008.
- [13] M. Girolami and C. He, "Probability density estimation from optimally condensed data samples," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 25, no. 10, pp. 1253–1264, 2003. [Online]. Available: <http://dblp.uni-trier.de/db/journals/IEEEPAMI/IEEEPAMI25.html>
- [14] P. Hall, S. J. Sheater, M. C. Jones, and J. S. Marron, "On optimal data-based bandwidth selection in kernel density estimation," *Biometrika*, vol. 78, no. 2, pp. 263–269, 1991.
- [15] C. C. Chang and C. J. Lin, *LIBSVM: A library for support vector machines*, 2001. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
- [16] A. Asuncion and D. Newman, "UCI machine learning repository," 2007. [Online]. Available: <http://www.ics.uci.edu/~mllearn/MLRepository.html>