

Table 1: Comparison of different single-machine methods with GREEDY and DDISMI using the SVM classifier.

	MIQ	JMI	mRMR	SpecCMI	QPFS	CMI	Fisher Score	relieff	GREEDY	DDISMI
Colon	78.1 \pm 5.3	83.1 \pm 2.5	79.8 \pm 2.3	72.2 \pm 4.1	78.0 \pm 2.8	69.7 \pm 5.4	81.3 \pm 3.1	71.7 \pm 5.1	84.4 \pm 3.2	83.1 \pm 3.2
Leukemia	92.7 \pm 4.6	88.7 \pm 3.6	99.9 \pm 0.6	89.5 \pm 3.6	82.0 \pm 4.6	82.0 \pm 4.6	99.7 \pm 0.6	98.9 \pm 1.1	96.0 \pm 2.6	96.0 \pm 2.6
Lung-discrete	82.1 \pm 4.2	90.9 \pm 2.2	90.7 \pm 1.7	84.1 \pm 8.7	91.4 \pm 2.6	83.8 \pm 6.7	88.7 \pm 5.1	80.8 \pm 6.6	92.1 \pm 1.6	91.5 \pm 1.7
Lung	86.8 \pm 2.5	90.4 \pm 2.6	96.7 \pm 0.6	95.2 \pm 0.8	96.6 \pm 0.7	96.3 \pm 0.8	89.8 \pm 4.2	92.2 \pm 1.6	96.4 \pm 1.6	95.9 \pm 1.4
Lymphoma	79.6 \pm 5.6	95.8 \pm 1.2	96.7 \pm 0.6	90.9 \pm 2.6	95.6 \pm 2.6	91.9 \pm 4.2	88.9 \pm 4.1	90.6 \pm 5.5	97.8 \pm 1.1	97.8 \pm 1.1
NCIS	90.2 \pm 4.5	74.0 \pm 3.5	74.4 \pm 2.6	51.2 \pm 5.1	25.0 \pm 5.4	32.2 \pm 7.3	64.8 \pm 4.6	50.7 \pm 5.6	83.0 \pm 5.6	82.2 \pm 6.2
Promoter	71.4 \pm 10.3	85.6 \pm 2.5	85.9 \pm 3.6	81.9 \pm 1.7	85.9 \pm 2.5	77.5 \pm 6.6	85.5 \pm 3.1	86.7 \pm 3.2	85.7 \pm 3.2	85.7 \pm 2.8
Spac	82.7 \pm 8.3	80.8 \pm 5.3	99.6 \pm 0.6	98.6 \pm 1.6	99.7 \pm 0.6	99.3 \pm 0.8	98.7 \pm 1.6	94.2 \pm 5.1	99.2 \pm 1.2	98.7 \pm 1.6
TCX-71	61.4 \pm 5.2	57.4 \pm 7.1	83.9 \pm 5.8	77.4 \pm 7.6	72.5 \pm 5.7	88.4 \pm 7.3	75.1 \pm 5.1	70.5 \pm 1.9	99.2 \pm 3.2	98.9 \pm 4.2
Multi-features	85.5 \pm 0.1	82.8 \pm 10.3	96.0 \pm 2.2	95.6 \pm 3.6	96.1 \pm 0.8	95.5 \pm 2.5	95.6 \pm 2.4	96.7 \pm 1.6	96.5 \pm 0.7	96.3 \pm 0.8
Optdigits	81.1 \pm 2.3	96.7 \pm 2.6	96.3 \pm 2.2	96.3 \pm 2.2	95.9 \pm 2.8	95.4 \pm 1.5	96.5 \pm 2.6	95.4 \pm 4.7	96.7 \pm 1.1	96.6 \pm 1.8
USPS	92.6 \pm 3.6	90.9 \pm 3.7	90.9 \pm 2.6	86.9 \pm 2.6	92.5 \pm 4.1	92.3 \pm 3.2	88.5 \pm 8.6	89.7 \pm 5.1	93.9 \pm 2.4	93.9 \pm 2.4
PCMAC	70.2 \pm 0.9	88.0 \pm 1.7	88.5 \pm 1.2	87.9 \pm 1.4	87.9 \pm 1.3	89.2 \pm 1.7	88.5 \pm 2.1	71.5 \pm 1.1	89.8 \pm 2.6	89.9 \pm 1.9
RELATTH	67.1 \pm 2.2	79.8 \pm 3.1	82.5 \pm 2.8	79.9 \pm 3.1	79.2 \pm 3.6	85.4 \pm 3.5	76.2 \pm 6.6	83.9 \pm 3.1	85.0 \pm 4.1	84.7 \pm 4.6
Musk2	90.2 \pm 0.1	90.4 \pm 0.1	90.3 \pm 0.1	90.3 \pm 0.1	90.4 \pm 0.2	90.5 \pm 0.1	90.4 \pm 0.1	90.4 \pm 0.6	90.2 \pm 0.1	90.3 \pm 0.1
WarpAR10P	83.4 \pm 8.1	85.6 \pm 6.5	93.0 \pm 2.8	88.2 \pm 4.8	94.0 \pm 2.8	84.0 \pm 2.2	85.1 \pm 0.4	85.0 \pm 10.3	92.7 \pm 1.4	93.9 \pm 4.0
Pixraw10P	94.1 \pm 2.6	97.2 \pm 1.6	98.0 \pm 0.3	98.0 \pm 2.2	95.9 \pm 4.1	98.5 \pm 0.7	94.9 \pm 5.1	85.0 \pm 13.1	97.9 \pm 2.3	97.6 \pm 2.5
WarpPIE10P	96.9 \pm 4.1	97.2 \pm 1.1	96.1 \pm 0.3	95.0 \pm 1.1	96.9 \pm 1.1	98.1 \pm 1.3	96.2 \pm 1.3	94.0 \pm 2.1	98.5 \pm 1.6	98.8 \pm 0.8
Yale	52.4 \pm 6.1	86.2 \pm 5.3	70.2 \pm 2.8	63.0 \pm 7.1	73.8 \pm 5.3	52.6 \pm 4.5	70.8 \pm 2.1	55.1 \pm 9.2	71.4 \pm 4.1	70.9 \pm 5.1
W17L	19.0 \pm 6.6	14.7 \pm 7.1	11.4 \pm 7.1	17.7 \pm 7.1	12.2 \pm 5.5	11.7 \pm 7.7	14.3 \pm 7.2	15.7 \pm 7.3	27.1 \pm 7.6	27.1 \pm 7.6

present an empirical evaluation to show that the distributed version of our method (DDISMI) is tens of times faster than its centralized variant (GREEDY). Note that GREEDY itself is as fast as state-of-the-art centralized methods due to its efficient greedy strategy. Then we demonstrate the advantages of the greedy approach over the existing local search method for diversity maximization in terms of performance and running time. After that, we investigate the defined objective function by studying the effect of λ value on the results. Finally, to validate the quality of the objective function, we show a high correlation between the objective value and the classification accuracy on two small-sized datasets. Before elaborating on the empirical results, it should be mentioned that unlike the GREEDY algorithm which arbitrarily selects the first feature, in the implementation, we select the one that has the maximum MI with the class labels vector.

Comparison to the state-of-the-art feature selection methods. In this section, we compare the quality of various centralized feature selection methods with the proposed distributed (DDISMI) and centralized (GREEDY) methods. In order to test the sensitivity of our method to the structure of the dataset, we have used several datasets from a variety of domains with various number of features and instances in addition to the classic datasets in the literature of feature selection. We have described these datasets in detail in the supplemental material.

We considered a variety of MI-based filter methods, namely Mutual Information Quotient (MIQ) (Ding and Peng 2005), Minimum Redundancy Maximum Relevance (mRMR) (Ding and Peng 2005), Joint Mutual Information (JMI) (Yang and Moody 1999), Spectral Conditional Mutual Information (SpecCMI) (Nguyen et al. 2014), Quadratic Programming Feature Selection (QPFS) (Rodriguez-Lujan et al. 2010), and Conditional Mutual Information (CMI) (Cheng et al. 2008) as baselines as well as non MI-based methods, fisher score (Duda, Hart, and Stork 2001) and reliefF (Robnik-Šikonja and Kononenko

2003). Note that prior papers have performed extensive studies (Brown et al. 2012) comparing these methods and we have chosen methods that achieve the best results in various domains. To test the quality of the selected features of each method, we feed them into a classifier method M and compare their classification accuracy. All of the experiments are performed with both SVM and 3-NN as the classifiers (M). Note that all of the methods are filter-based and hence are independent from the selection of M . The LIBSVM package (Chang and Lin 2011) is the underlying implementation of the linear SVM with its regularization factor set to 1. We change $|S|$ from 10 to $\min\{100, n\}$, where n is the number of the features in each dataset. Finally, we report the average cross validation (CV) classification accuracy of each method on all of the experiments with different values of $|S|$. A 10-fold CV is being used for the datasets with more than 100 instances and for the others, we employ the leave-one-out CV. In order to compute the probabilities used in MIs and VIs we have discretized the continuous variables using the Minimum Description Length (MDL) method (Irani and Fayyad 1993) with 5 bins. The regularization factor of our algorithm (λ) is set to be 0.8 in all of our experiments. We elaborate on choosing the λ in the experiment about effect of λ value. In order to run DDISMI, we simulate the distributed setting on a single machine and each (simulated) machine only has access to its own feature vectors. Each machine is responsible for processing \sqrt{nk} features when we have n features and want to select k of them. Thus we employ $\sqrt{n/k}$ machines for the first stage. Then, we merge the results of all of the machines together and then process them again, i.e., we select k features from all of these $k\sqrt{n/k} = \sqrt{nk}$ features. For the sake of reproducibility, we have provided all the codes in the supplemental material. Table 1 compares the SVM classification accuracy of the selected features. For each dataset (row), the results with higher accuracy and lower standard deviation are indicated