

Robust and Effective Factorization Machines

Anonymous

Optimization Algorithm

The original objective for classification task takes the form:

$$\min_{\mathbf{w} \in \mathbf{R}^d, \mathbf{Z} \in \mathbf{S}_+^{d \times d}} \sum_{i=1}^n e_i(\max\{\max(y_i(\mathbf{w}^\top \mathbf{x}_i + \langle \mathbf{Z}, \mathbf{x}_i \mathbf{x}_i^\top \rangle), 0) - \epsilon_1, 0\})^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2 + \sum_s \min\{\lambda_s^2, \epsilon_3\}, \quad (1)$$

$$e_i = \begin{cases} \frac{1}{2error}, & 0 < error \leq \epsilon_2; \\ 0, & otherwise \end{cases}$$

where $error = \max(y_i(\mathbf{w}^\top \mathbf{x}_i + \langle \mathbf{Z}, \mathbf{x}_i \mathbf{x}_i^\top \rangle), 0) - \epsilon_1$

The subgradient with respect to \mathbf{Z} is

$$\nabla_{\mathbf{Z}, I} = \sum_{i=1}^b \mathbf{x}_i \mathbf{x}_i^\top + \beta \mathbf{P}_M \mathbf{P}_M^\top \mathbf{Z} \quad (2)$$

To incrementally calculate the SVD of $\mathbf{Z} - \eta \nabla_{\mathbf{Z}, I}$. Let the symmetric and low rank matrix \mathbf{Z} has rank k and its economy SVD is $\mathbf{Z} = \mathbf{P}_k \Sigma_k \mathbf{P}_k^\top$. As matrix $\nabla_{\mathbf{Z}, I}$ is symmetric and low rank, we can represent it as $\mathbf{A} \mathbf{A}^\top$.

$$\nabla_{\mathbf{Z}, I} = \mathbf{X} \mathbf{X}^\top + \beta \mathbf{P}_M \mathbf{P}_M^\top \mathbf{Z} \quad (3)$$

Experiments

In this section, we empirically investigate whether our proposed RobFM method can achieve better and robust performance compared to original factorization machine model with fixed rank on benchmark datasets.

Experimental Testbeds and Setup

We conduct our experiments on four public datasets. Table 1 gives a brief summary of these datasets. All the datasets are normalized to have zero mean and unit variance in each dimension. To make fair comparison, all the algorithms are conducted over 5 experimental runs of different random permutations. We apply hinge loss for training and evaluate the performance of our proposed methods for classification task by measuring accuracy and hinge loss. For parameter settings, we perform grid search to choose the best parameters for each algorithm on the training set.

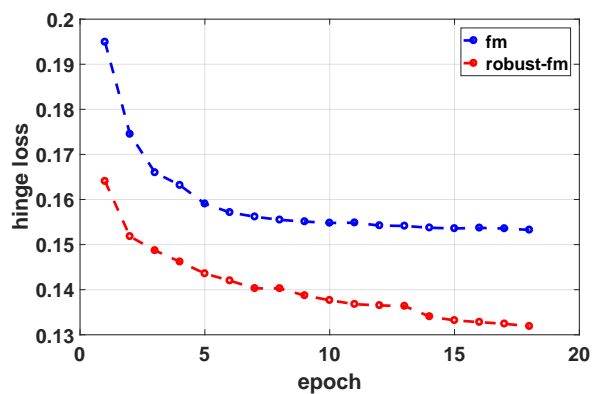
Dataset	#Training	#Test	#Feature	#class
phishing	7370	3685	68	2
w8a	49749	14951	300	2
protein	17766	6621	357	3
IJCNN	49990	91701	22	2
Covtype	387342	193670	54	2
MNIST	60000	10000	784	10

Table 1: Summary of datasets used in our experiments.

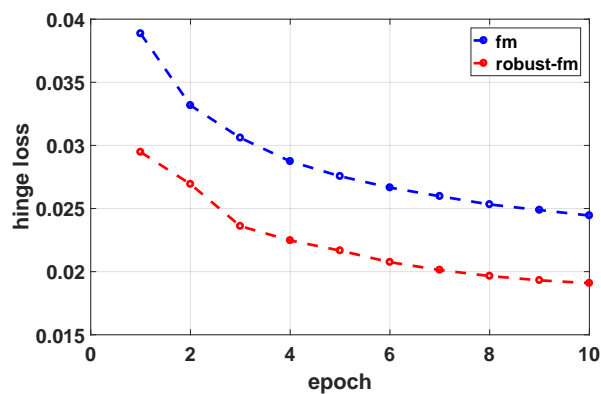
Experimental Results

phishing	Train loss	Test loss	Acc(%)
Ridge	0.3037 \pm 0.0002	0.3172 \pm 0.0015	92.65 \pm 0.11
SVM	0.1451 \pm 0.0002	0.1590 \pm 0.0058	93.26 \pm 0.30
FM	0.1397 \pm 0.0002	0.1534 \pm 0.0006	93.35 \pm 0.10
RobFM-IncSVD	\pm	\pm	\pm
RobFM	0.1145\pm0.0004	0.1357\pm0.0051	94.59 \pm 0.33
w8a	Train loss	Test loss	Acc(%)
Ridge	0.0706 \pm 0.0004	0.0730 \pm 0.0096	98.34 \pm 0.02
SVM	0.0305 \pm 0.0006	0.0316 \pm 0.0003	98.65 \pm 0.02
FM	0.0234 \pm 0.0010	0.0245 \pm 0.0002	98.86 \pm 0.07
RobFM-IncSVD	\pm	\pm	\pm
RobFM	0.0178\pm0.0007	0.0190\pm0.0003	99.10 \pm 0.06
protein	Train loss	Test loss	Acc(%)
Ridge	1.8807 \pm 0.0002	1.8828 \pm 0.0035	68.95 \pm 0.06
SVM	1.5369 \pm 0.0002	1.5198 \pm 0.0014	68.61 \pm 0.07
FM	1.3757 \pm 0.0002	1.4668 \pm 0.0016	69.21 \pm 0.10
RobFM-IncSVD	\pm	\pm	\pm
RobFM	\pm	\pm	\pm
IJCNN	Train loss	Test loss	Acc(%)
Ridge	0.3147 \pm 0.0007	0.3173 \pm 0.0016	90.50 \pm 0.00
SVM	0.1770 \pm 0.0004	0.1843 \pm 0.0003	91.35 \pm 0.02
FM	0.0930 \pm 0.0005	0.0955 \pm 0.0004	96.66 \pm 0.09
RobFM-IncSVD	\pm	\pm	\pm
RobFM	0.0712\pm0.0001	0.0744\pm0.0013	97.83 \pm 0.15
Covtype	Train loss	Test loss	Acc(%)
Ridge	0.7779 \pm 0.0002	0.7780 \pm 0.0001	70.07 \pm 0.04
SVM	0.6080 \pm 0.0003	0.6102 \pm 0.0002	68.76 \pm 0.06
FM	0.5111 \pm 0.0001	0.5088 \pm 0.0030	77.31 \pm 0.10
RobFM-IncSVD	\pm	\pm	\pm
RobFM	0.4894\pm0.0001	0.4844\pm0.0066	79.65 \pm 0.25
MNIST	Train loss	Test loss	Acc(%)
Ridge	0.8434 \pm 0.0002	0.8257 \pm 0.0069	91.56 \pm 0.12
SVM	0.5957 \pm 0.0008	0.5776 \pm 0.0007	91.64 \pm 0.13
FM	0.2224 \pm 0.0005	0.2949 \pm 0.0034	96.62 \pm 0.08
RobFM-IncSVD	\pm	\pm	\pm
RobFM	\pm	\pm	\pm

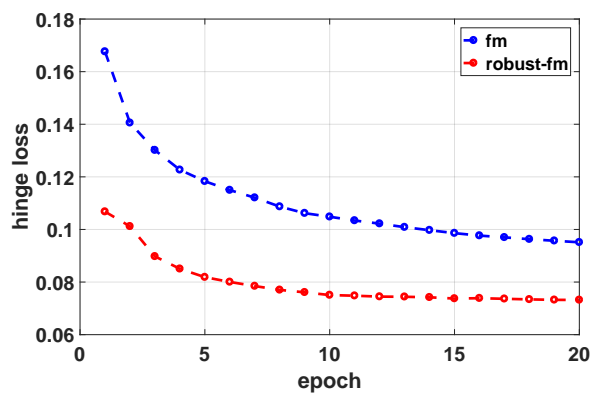
Table 2: Comparison of different algorithms in terms of train loss, test loss, classification accuracy



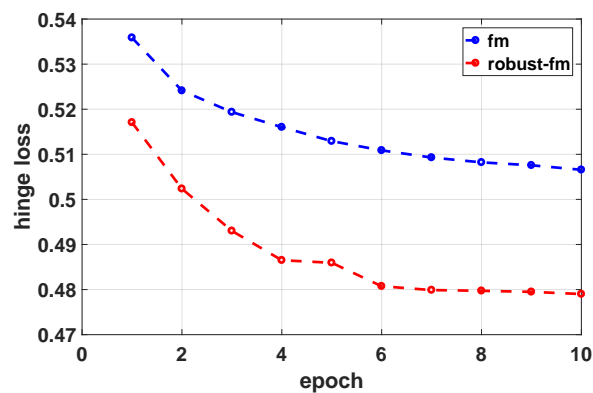
(a) phishing



(b) w8a



(c) IJCNN



(d) Covtype

Figure 1: Epoch-wise demonstration of different algorithms with hinge loss on test data