



**Nightmare at Test Time:
Robust Learning by Feature Deletion**

Präsentieren: Yang yang

17.11.07

Outline

- ❖ Introduction
- ❖ FDRQP: tractable quadratic program for training robust classifiers
- ❖ Experiments
 - ⌘ Handwritten Digit Classification
 - ⌘ Spam Filtering
- ❖ Summary

Introduction

- ❖ Testing comes *after* training
 - ⌘ too much weight to any single input feature
 - ❖ with nonstationary feature distribution
 - ❖ with input sensor failure
- ❖ A common approach
 - ⌘ Regularization which spreads the weight
 - ⌘ Very generic and cannot induce robustness

Introduction

❖ Solution

⌘ New algorithm

- ❖ avoiding single feature over-weighting
- ❖ Using quadratic programming

❖ The application of our methodes on

- ⌘ Handwritten digit recongnition
- ⌘ Spam filtering

Worst Case Deletion

Input:

❖ Labeled Sample (x_i, y_i) ($i = 1, \dots, n$),

❖ Feature vektor $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{\pm 1\}$

❖ Number of features deleted from each sample point X : K

output:

❖ a linear classifier: $y(\mathbf{x}) = \text{Sign}(\mathbf{w} \cdot \mathbf{x})$

Worst Case Deletion

Output:

❖ Performance Measure: Regularized hinge

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i h^{wc}(\mathbf{w}, y_i \mathbf{x}_i)$$

❖ Hinge loss:

$$h^{wc}(\mathbf{w}, y_i \mathbf{x}_i) = \begin{array}{ll} \max & [1 - y_i \mathbf{w} \cdot (\mathbf{x}_i \circ (1 - \boldsymbol{\alpha}_i))]_+ \\ \text{s.t.} & \boldsymbol{\alpha}_i \in \{0, 1\} \\ & \sum_j \alpha_{ij} = K \end{array}$$

Hinge loss

- ❖ a convex upper bound on the zero

$$l_{zo}(\mathbf{w}, y, \mathbf{x}) \leq \sum_i [1 - y_i \mathbf{w} \cdot \mathbf{x}_i]_+$$

- ❖ Find \mathbf{w} which minimizes the worst case hinge loss

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_i h^{wc}(\mathbf{w}, y_i \mathbf{x}_i)$$

- ❖ Minimizing hinge loss \rightarrow minimizing on the training error

The FDROP

How can people Solving the minimax

$$h^{wc}(\mathbf{w}, y_i \mathbf{x}_i) = \left[1 - y_i \mathbf{w}^T \mathbf{x}_i + s_i \right]_+$$
$$s_i = \max_{\substack{\alpha_i \in \{0, 1\} \\ \sum_j \alpha_{ij} = K}} y_i \mathbf{w} \cdot (\mathbf{x}_i \circ \boldsymbol{\alpha}_i)$$

s_i is the maximum contribution of K features to the margin of sample x_i

$$s_i = \max_{\substack{\alpha_i \in \{0, 1\} \\ \sum_j \alpha_{ij} = K}} y_i (\mathbf{w} \circ \mathbf{x}_i) \cdot \boldsymbol{\alpha}_i$$
$$\text{s.t. } 0 \leq \alpha_i \leq 1$$
$$\sum_j \alpha_{ij} = K$$

The FDRPOP

- ❖ The maximization problem for s_i has an LP

$$s_i = \min \quad K z_i + \sum_j v_{ij}$$
$$s.t. \quad z_i + \mathbf{v}_i \geq (y_i \mathbf{x}_i \circ \mathbf{w}), \mathbf{v}_i \geq 0$$

- ❖ Linear in all variables

FDRPOP:

$$\min \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i [1 - y_i \mathbf{w}^T \mathbf{x}_i + t_i]_+$$
$$s.t. \quad t_i \geq K z_i + \sum_j v_{ij}$$
$$z_i + \mathbf{v}_i \geq (y_i \mathbf{x}_i \circ \mathbf{w}) \quad , \quad \mathbf{v}_i \geq 0$$

FDRP VS SVM

- ❖ FDRP is variant of SVM

- ☞ linear classifier

- ☞ the training objektiv is measured using a

- ☞ regularized hinge loss

FDRP VS SVM

❖ FDRP is variant of SVM

∞ differently error term compare to FDRP

	<u>DSVM:</u>		<u>DFDRP:</u>
min	$\frac{1}{2} \ \mathbf{w}\ ^2 - \sum_i \alpha_i$	min	$\frac{1}{2} \ \mathbf{w}\ ^2 - \sum_i \alpha_i$
s.t.	$\mathbf{w} = \sum_i y_i \alpha_i \mathbf{x}_i$	s.t.	$\mathbf{w} = \sum_i y_i \alpha_i \mathbf{x}_i \circ (\mathbf{1} - \lambda_i)$
	$0 \leq \alpha \leq C$		$0 \leq \alpha \leq C$
			$0 \leq \lambda_i \leq 1$
			$\sum_j \lambda_{ij} = K$

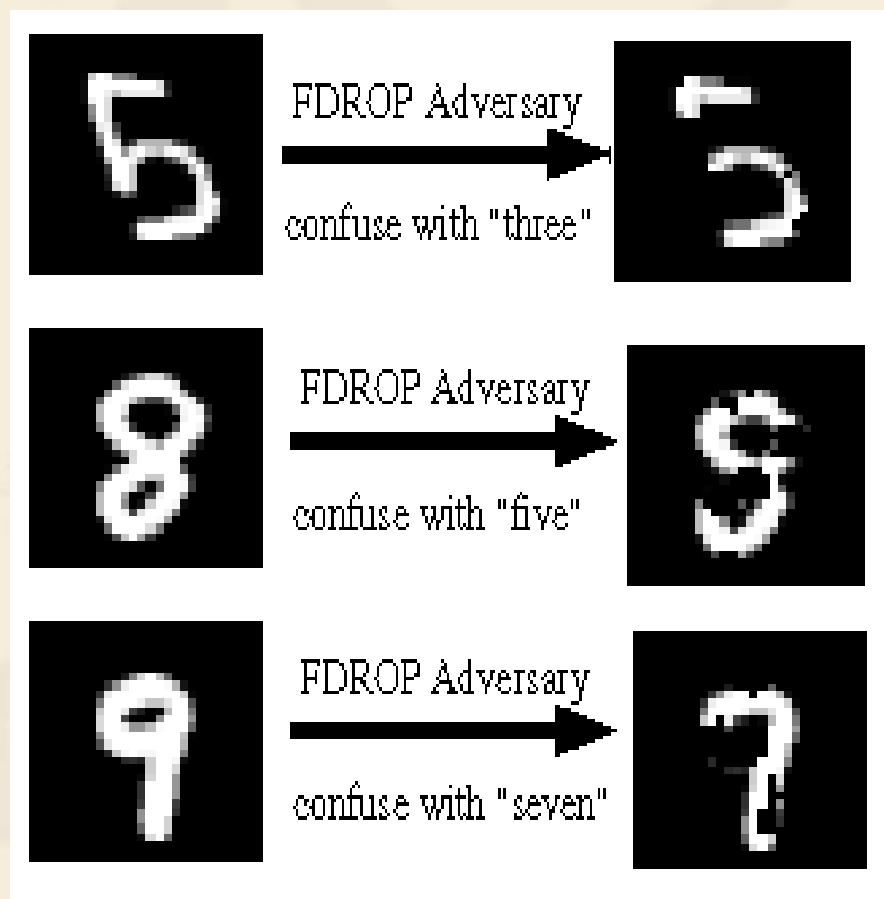
Handwritten Digit Classification

- ❖ investigated the application of FDROP to classifying handwritten digits
- ❖ robustness to pixel deletion in these images
- ❖ Binary problems
- ❖ Small training sets of 50 samples per digit
- ❖ Chosen pairs which hard to distinguish

∞ (5,3),(8,5),(9,7)....

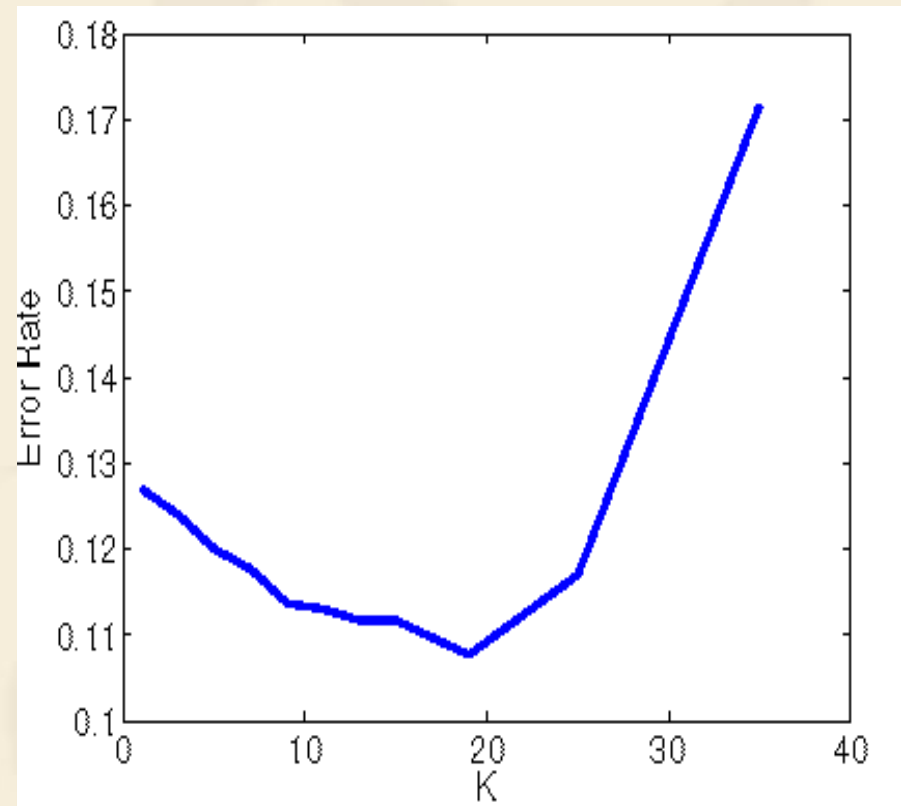
Handwritten Digit classification

- ❖ visual representation of the feature deletion process
- ❖ K destructive feature deleted (K=50)
- ❖ maximize the resemblance between the given digit and the digit in the other class



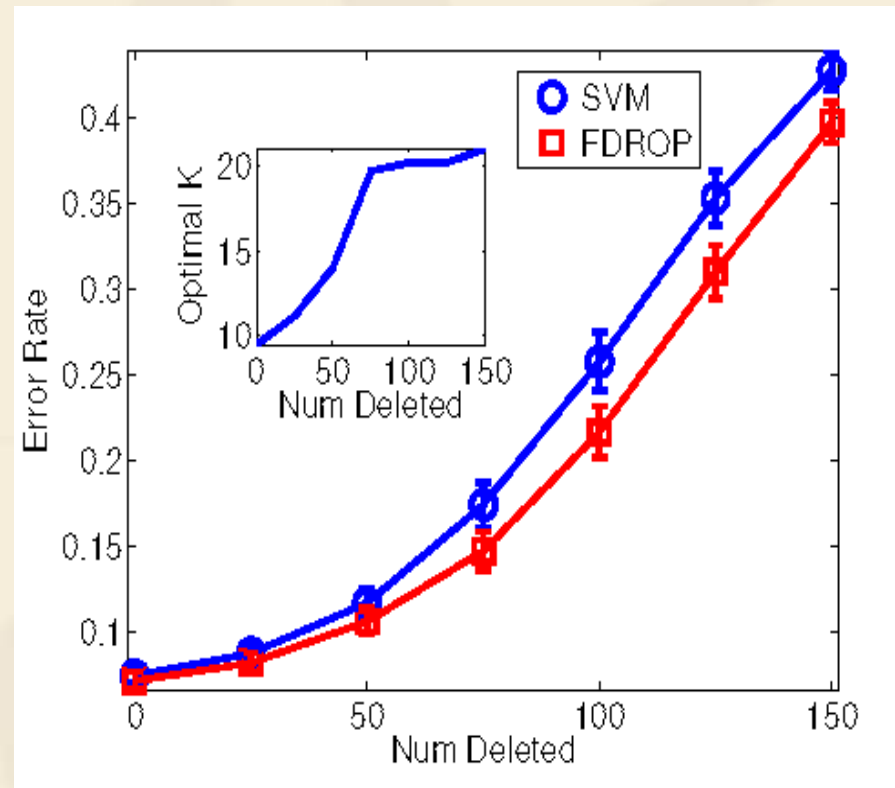
Handwritten Digit Classification

- ❖ Classification error for the digit pair (4; 7)
- ❖ $K=50$
- ❖ dependence on K
- ❖ e.g.
 - ⌘ Book and exam
 - ⌘ howmuch book read in order to better point



Handwritten Digit Classification

- ❖ the dependence of classification error on the number of deleted features
- ❖ FDROP suffers less degradation in error when compared to SVM
- ❖ optimal K grows monotonously
- ❖ features dropped randomly



Summary

- ❖ Presented a new classification algorithm that is robust to worst case feature deletion
 - ∞ FDRP
 - ∞ Hinge loss
- ❖ Handwritten Digit classification

The image features a traditional Chinese ink wash painting of a plum blossom branch. The branch is dark and gnarled, with small, delicate blossoms and buds. The background is a light, warm tone with large, faint calligraphic characters. The text is centered and reads:

Discussion
Thanks for you
attention