Prototype Support Vector Machines: Supervised Classification in Complex Datasets

April Shen

Williams College Computer Science Department Advisor: Andrea Danyluk

May 13, 2013









- Finding the right model without knowing anything about your data is tricky
- Requires making assumptions about data distributions

- Freund and Schapire, 1995 [3]
- Iteratively train classifiers on different samples of training data
- Focus on hard-to-capture data
- Requires choice of base classifier

Support Vector Machines (SVMs)



- Cortes and Vapnik, 1995 [1]
- Simplest case: separating hyperplane with minimum error and maximum margin

Support Vector Machines (SVMs)



- Cortes and Vapnik, 1995 [1]
- Simplest case: separating hyperplane with minimum error and maximum margin
- Unseparable case

Support Vector Machines (SVMs)



- Can nonlinearly map to high-dimensional spaces using kernels
- Requires choice of kernel function

Exemplar SVMs (ESVMs)

- Malisiewicz et al., 2011 [2]
- Train classifier with single positive instance and many negative instances
- Ensemble of these can "vote" on new instances



- Strengths
 - Ensemble of classifiers allows high accuracy
 - Placement of models is specific to data distribution
- Weaknesses
 - Fine-tuned to object recognition domain
 - Data does sometime occur in clusters

From Exemplars to Prototypes



Initialization

- Seeding the positive sets
- Seeding the negative sets

















- Choose best iteration's ensemble based on accuracy on validation set
- If a model classifies a new instance positively, makes a weighted vote for its class
- Class with maximum number of votes is the predicted class

Strengths

- Inherits benefits of ESVMs
- Balances flexibility with performance
- No need for model selection
- Weaknesses
 - Potentially large ensemble sizes
 - Slower than other algorithms
 - Too many models for simple distributions
 - Too few models for complex distributions

- Sample data, so don't overdo simple regions
- Resample using AdaBoost-style weighting, so focus on difficult regions

Weighted Sampling



Weighted Sampling



Weighted Sampling



- Strengths
 - Greater sensitivity to complex structure
 - Smaller ensemble sizes
- Weaknesses
 - Underperforms in simpler/smaller datasets
 - More parameters must be optimized

Remember: there is no free lunch!

Algorithms

- C4.5 (decision tree)
- AdaBoosted C4.5
- Linear SVM (LibSVM)
- AdaBoosted Linear SVM
- Polynomial SVM
- Multilayer Perceptron

Datasets

- Synthetic
- UCI benchmarks
- Twitter

(Results averaged over 10-fold cross validation, using Weka's implementations and default parameters)

Synthetic Datasets

isolated



Synthetic Datasets

striated



Synthetic Datasets

spirals



Dataset	PSVM	C4.5	Boosted C4.5	Linear SVM
isolated	90.9 ± 4.07	98.1±1.79 •	98.5±1.46 •	$62.5\pm0.0~\circ$
striated	$\textbf{97.3} \pm \textbf{1.66}$	69.4±13.4 °	$90.1{\pm}7.38{}\circ$	$48.5 \!\pm\! 3.94 \circ$
spirals	21.8 ± 18.9	$0.0\pm0.0\circ$	$0.0\pm0.0\circ$	$0.0\pm0.0\circ$
iris	96.0 ± 4.42	94.0 ± 6.29	94.0 ± 5.54	98.7 ± 2.67
glass	52.8 ± 8.72	69.1±6.40 •	72.9±7.85•	63.5±8.08 •
vehicle	$\textbf{79.4} \pm \textbf{4.49}$	$73.8{\pm}4.48\circ$	75.7 ± 3.56	80.4 ± 4.50
segment	95.2 ± 1.18	97.1±0.93 ●	98.1±0.85 •	96.3±0.93 •
twitter	55.8 ± 5.12	54.5 ± 2.89	60.5 ± 7.99	62.2±3.66 •

Dataset	PSVM	Boosted Linear SVM	Polynomial SVM	Multilayer Perceptron
isolated	90.9 ± 4.07	$62.5\pm0.0\circ$	$81.9 {\pm} 4.38 \circ$	$80.6{\pm}2.81\circ$
striated	$\textbf{97.3} \pm \textbf{1.66}$	$53.4{\pm}16.3\circ$	$72.8 {\pm} 4.18 \circ$	$75.0{\pm}25.0\circ$
spirals	21.8 ± 18.9	5.63±8.97 o	$0.0\pm0.0\circ$	20.6 ± 28.9
iris	96.0 ± 4.42	$\textbf{97.3} \pm \textbf{3.27}$	96.7 ± 4.47	$\textbf{97.3} \pm \textbf{4.42}$
glass	52.8 ± 8.72	63.1±7.73•	$69.7{\pm}6.55 \bullet$	70.6±8.82 •
vehicle	$\textbf{79.4} \pm \textbf{4.49}$	80.4 ± 4.23	80.4 ± 4.53	79.7 ± 4.61
segment	95.2 ± 1.18	96.1 ± 0.89	95.8 ± 1.38	96.2 ± 1.30
twitter	55.8 ± 5.12	62.2±4.15 •	52.3 ± 5.54	52.3 ± 5.54

Results

Dataset	Sampling	Regular
isolated	84.0 ± 7.0	$90.9 \pm 4.07 ~\bullet~$
striated	96.9 ± 3.84	$\textbf{97.3} \pm \textbf{1.66}$
spirals	60.7 ± 34.9	$21.8\pm18.9\circ$
iris	86.7 ± 13.0	96.0 ± 4.42
glass	47.6 ± 12.9	52.8 ± 8.72
vehicle	75.4 ± 3.83	$\textbf{79.4} \pm \textbf{4.49}$
segment	94.4 ± 1.86	95.2 ± 1.18
twitter	54.0 ± 8.86	55.8 ± 5.12

- Synthetic and benchmarks are unrealistically clean
- Inject 10% noise in class labels

Dataset	PSVM	C4.5	Boosted C4.5	Linear SVM
isolated	71.9 ± 7.71	84.8±4.36 •	83.6±4.65 •	59.3±3.67 o
striated	$\textbf{76.8} \pm \textbf{7.42}$	60.3±7.24 o	$64.3 {\pm} 6.55 \circ$	$52.0{\pm}2.86\circ$
spirals	29.0 ± 15.4	$9.79\!\pm\!6.28\circ$	$9.79\!\pm\!6.28\circ$	$11.9\!\pm\!10.2\circ$
iris	88.7 ± 5.21	84.7 ± 7.33	83.3 ± 7.45	90.0 ± 6.83
glass	50.0 ± 5.50	62.2±7.46 •	66.3±8.87•	53.7 ± 8.91
vehicle	68.1 ± 4.43	$63.1{\pm}1.89\circ$	67.6 ± 2.54	71.5 ± 3.57
segment	84.3 ± 2.18	85.5 ± 2.30	85.4 ± 2.23	84.1 ± 1.82

Dataset	PSVM	Boosted Linear SVM	Polynomial SVM	Multilayer Perceptron
isolated	71.9 ± 7.71	$57.9{\pm}3.71\circ$	75.5 ± 4.91	$\textbf{73.9} \pm \textbf{3.51}$
striated	$\textbf{76.8} \pm \textbf{7.42}$	$56.8 {\pm} 7.36 \circ$	$66.8 {\pm} 3.12 \circ$	67.4 ± 18.5
spirals	29.0 ± 15.4	16.9 ± 17.1	$10.3 {\pm} 6.95 \circ$	17.1 ± 13.4
iris	88.7 ± 5.21	90.0 ± 6.15	90.0 ± 5.37	91.3 ± 6.70
glass	50.0 ± 5.50	57.5±7.90 ●	61.6±8.95 •	61.7±4.14 ●
vehicle	68.1 ± 4.43	69.1 ± 4.24	$57.4 {\pm} 5.12 \circ$	71.0 ± 2.61
segment	84.3 ± 2.18	84.0 ± 2.12	$67.8 {\pm} 3.95 \circ$	85.4 ± 2.51

Dataset	Sampling	Regular
isolated	68.8 ± 5.81	71.9 ± 7.71
striated	75.1 ± 10.1	$\textbf{76.8} \pm \textbf{7.42}$
spirals	57.8 ± 28.6	$29.0\pm15.4\circ$
iris	76.7 ± 10.4	88.7 ± 5.21 •
glass	43.4 ± 9.40	$50.0\pm5.50~\bullet$
vehicle	59.1 ± 5.69	68.1 ± 4.43
segment	80.6 ± 3.15	$84.3\pm2.18~\bullet$

		Noiseless		Noisy	
	ESVM	Regular	Sampling	Regular	Sampling
isolated	720	323	124	122	70
striated	720	149	88	170	93
spirals	175	47	29	56	20
iris	135	54	13	24	11
glass	193	35	15	37	20
vehicle	761	110	61	113	64
segment	2079	453	172	323	164
twitter	540	42	32	Х	Х

Win-Lose-Tie counts:

	Noiseless	Noisy
Regular	16 - 13 - 19	14 - 7 - 21
Sampling	14 - 23 - 11	13 - 23 - 6

- Performs very well in most difficult datasets
- Performs as well as other algorithms in all datasets
- Degrades gracefully with noise
- Standard PSVM outperforms version with sampling

- Further empirical studies
 - More datasets, especially real-world data
 - More thorough comparison in noisy datasets
 - How dataset complexity affects ensemble size
- Algorithm improvements
 - Feature selection for each model (e.g. 1-norm SVMs [4])
 - Different metrics besides Euclidean distance
- Prove can approximate functions to within certain error
 - Given a minimum size of training set?
 - Within a certain number of shifting iterations?
- Explicitly incorporate clustering

References

Corinna Cortes and Vladimir Vapnik. Support-vector networks. Machine Learning, 20(3), 1995.

- Tomasz Malisiewicz, Abhinav Gupta, and Alexei A. Efros. Ensemble of exemplar-SVMs for object detection and beyond. In Proceedings of the 2011 International Conference on Computer Vision, ICCV '11, 2011.
 - Robert E. Schapire.
 - A brief introduction to boosting.

In Proceedings of the 1999 International Joint Conference on Artificial Intelligence, IJCAI '99, 1999.

J. Zhu, S. Rosset, T. Hastie, and R. Tibshirani.
1-norm support vector machines.
In Advances in Neural Information Processing Systems 16, 2004.

BONUS MATERIAL

Final Ensembles



Final Ensembles



TRAIN

Input: set of data, T **Parameters:** number of iterations, s 1: Split data into training and validation sets, D and V. 2: $P \leftarrow [[d_i] \text{ for } d_i \in D]$ 3: $N \leftarrow [CHOOSENEGATIVES(D, d_i) \text{ for } d_i \in D]$ 4: for j = 0, ..., s do 5: $E_i \leftarrow []$ 6: for $P_i \in P$ and $N_i \in N$ do 7: Train a linear SVM, using P_i and N_i . 8: Add this SVM to E_i . 9: $a_i \leftarrow \text{TEST}(V, E_i)$ 10: $P, N \leftarrow \text{SHIFT}(D, E_i, P, N)$ 11: **return** ensemble E_i with highest accuracy on V

Input: set of training data, D, and training instance, d_i **Parameters:** number of negatives to return, k

- 1: $N_i \leftarrow [$]
- 2: $D_i \leftarrow$ all instances in D of a different class label from d_i
- 3: Compute Euclidean distance from d_i to each element of D_i .
- 4: Sort D_i in ascending order by distance.
- 5: $x \leftarrow \text{closest negative in } D_i$

6:
$$\vec{n} \leftarrow (x - d_i)/||x - d_i||$$

7: for
$$d_j \in D_i$$
 do

8: if
$$\vec{n} \cdot (d_j - d_i) > 0$$
 then

9: Add
$$d_j$$
 to N_i

- 10: **if** $|N_i| = k$ then
- 11: **return** *N_i*

Shift

Input: set of training data, D; ensemble of models, E; positive and negative sets for each model, P and N**Parameters:** probability to add to negative set, p 1: $C \leftarrow [[] \text{ for } d_i \in D]$ (candidate models) 2: for $m_i \in E$ and $d_i \in D$ do if m_i classifies d_i positively then 3: if d_i 's class matches m_i 's class then 4: 5: Compute the distance of d_i to m_i 's exemplar. Add m_i and its distance to the list of candidates C_i . 6: 7. else Add d_i to m_i 's negative set N_i with some probability p. 8: 9: for $d_i \in D$ do Add d_i to the positive set P_k of closest candidate model m_k 10: 11: for $m_i \in E$ do if m_i did not classify anything positively then 12: 13: Remove m_i from the ensemble 14: **return** *P*, *N*

Input: set of testing data, D, and ensemble of models, E

- 1: for $m_i \in E$ and $d_j \in D$ do
- 2: **if** m_i classifies d_j positively **then**
- 3: Record weighted vote for m_i 's class.
- 4: for $d_j \in D$ do
- 5: Output class with max votes as prediction for d_j .
- 6: $a \leftarrow$ classification accuracy over all data in D.
- 7: return a

Input: set of training data, D

Parameters: number of iterations, r

- 1: Initialize all weights to 1/n.
- 2: $S \leftarrow \text{SAMPLE}(D)$
- 3: for r iterations do
- 4: $E \leftarrow \text{TRAIN}(S)$
- 5: $a \leftarrow \text{Test}(D, E)$
- 6: REWEIGHT(D, a)
- 7: $S \leftarrow S \cup \text{SAMPLE}(D)$
- 8: return E

Input: set of training data, D, and accuracy of the ensemble, a

1: $\varepsilon = 1 - a$ 2: $\alpha = \frac{1}{2} \ln \left(\frac{1 - \varepsilon}{\varepsilon} \right)$ 3: for $d_i \in D$ do 4: if d_i was correctly classified then 5: weight $(d_i) = 0$ 6: else

7:
$$weight(d_i) = \frac{weight(d_i) \times e^{\alpha}}{z}$$
 $(z = normalization factor)$