

COM S 778
**A Support Vector Method
for
Multivariate Performance Measures**

by Thorsten Joachims
International Conference on Machine Learning 2005

Presented by: Thorsten Joachims

Thanks to
Rich Caruana, Alexandru Niculescu-Mizil, Pierre Dupont,
Jérôme Callut

Supervised Learning

- Find function from input space X to output space Y

$$h : X \rightarrow \{+1, -1\}$$

such that the ~~prediction error~~ is low.

- Text Classification:**
- F1-Score
 - Precision/Recall Break-Even (PRBEP)
- Medical Diagnosis:**
- ROC Area
- Information Retrieval:**
- Precision at 10

Related Work

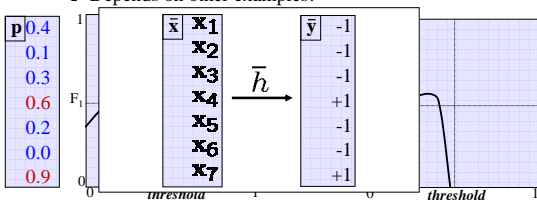
- Approach "Estimate Probabilities"**
 - E.g. [Platt, 2000] [Langford & Zadrozny, 2005] [Niculescu-Mizil & Caruana, 2005]
 - Potentially solve harder problem than required
- Approach "Optimize Substitute Loss, then Post-Process"**
 - E.g. [Lewis, 2001] [Yang, 2001] [Abe et al. 2004] [Caruana & Niculescu-Mizil, 2004]
 - Typically multi-step approach, cross-validation
- Approach "Directly Optimize Desired Loss"**
 - Linear cost models: e.g. [Morik et al., 1999] [Lin et al., 2002]
 - ROC-Area: e.g. [Herbrich et al. 2000] [Rakotomamonjy, 2004] [Cortes & Mohri, 2003] [Freund et al., 1998] [Yan et al., 2003] [Ferri et al., 2002]
 - F1-Score: difficult [Musicant et al. 2003]

Overview

- Formulation of Support Vector Machine for**
 - any loss function that can be computed from the contingency table.
 - F1-score, Error Rate, Linear Cost Models, etc.
 - any loss function that can be computed from contingency tables with cardinality constraints.
 - PRBEP, Prec@k, Rec@k, etc.
 - ROC-Area
- Polynomial Time Algorithm**
- Conventional classification SVM is special case**
 - New optimization problem
 - New representation and (extremely sparse) support vectors

Optimizing F_1 -Score

- F_1 -score is non-linear function of example set
 - F_1 -score: harmonic average of precision and recall
- $$F_1 = \frac{2 \text{Prec} \text{Rec}}{\text{Prec} + \text{Rec}}$$
- For example vector x_1 . Predict $y_j=1$, if $P(y_j=1|x_1)=0.4$?
→ Depends on other examples!



Approach: Multivariate Prediction

- Training Data:** $S = ((x_1, y_1), \dots, (x_n, y_n)) \sim_{i.i.d} \text{Pr}(x, y)$
- Conventional Setting:** learn $h : X \rightarrow \{-1, +1\}$

$$R^\delta(h) = \int \delta(h(x'), y') d\text{Pr}(x', y')$$

$$\hat{R}_S^\delta(h) = \frac{1}{n} \sum_{i=1}^n \delta(h(x_i), y_i)$$

- Multivariate Setting:** learn $\bar{h} : X^n \rightarrow \{-1, +1\}^n$

$$R^\Delta(\bar{h}) = \int \Delta(\bar{h}(x'_1, \dots, x'_n), (y'_1, \dots, y'_n)) d\text{Pr}(S')$$

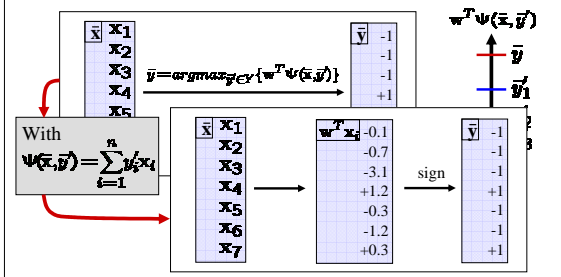
$$\hat{R}_S^\Delta(\bar{h}) = \Delta(\bar{h}(x_1, \dots, x_n), (y_1, \dots, y_n))$$

Note:
If $\Delta(\bar{h}(x_1, \dots, x_n), (y_1, \dots, y_n)) = \frac{1}{n} \sum_{i=1}^n \delta(h(x_i), y_i)$
then both settings are equivalent.

Multivariate Support Vector Machine

Approach: Linear Discriminant [Collins 2002] [Lafferty et al. 2002] [Taskar et al. 2004] [Tsochantaridis et al. 2004] etc.

– “Learn weights \mathbf{w} so that $\mathbf{w}^T \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ is max for correct $\bar{\mathbf{y}}$ ”



Multivariate SVM Optimization Problem

Approach: Structural SVM [Taskar et al. 04] [Tsochantaridis et al. 04]

Hard-margin optimization problem:

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

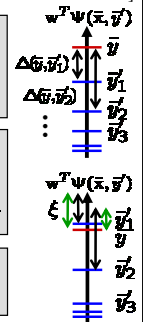
s.t. $\forall \bar{\mathbf{y}} \in \mathcal{Y} \setminus \{\bar{\mathbf{y}}\} : \mathbf{w}^T \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \geq \mathbf{w}^T \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) + 1$

Soft-margin optimization problem:

$$\min_{\mathbf{w}, \xi \geq 0} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \xi$$

s.t. $\forall \bar{\mathbf{y}} \in \mathcal{Y} : \mathbf{w}^T \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \geq \mathbf{w}^T \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) + \Delta(\bar{\mathbf{y}}, \bar{\mathbf{y}}) - \xi$

Theorem: At the solution, the training loss is upper bounded by $\Delta(\bar{\mathbf{y}}, \hat{\mathbf{h}}(\bar{\mathbf{x}})) \leq \xi$.



Multivariate SVM Generalizes Classification SVM

Theorem: The solutions of the multivariate SVM with number of errors as the loss function and an (unbiased) classification SVM are equal.

Multivariate SVM optimizing Error Rate:

$$\min_{\mathbf{w}, \xi \geq 0} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \xi$$

s.t. $\forall \bar{\mathbf{y}} \in \mathcal{Y} : \mathbf{w}^T \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \geq \mathbf{w}^T \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) + 2E_{\text{err}}(\bar{\mathbf{y}}) - \xi$

Classification SVM (unbiased):

$$\min_{\mathbf{w}, \xi \geq 0} \frac{1}{2} \mathbf{w}^T \mathbf{w} + 2C \sum_{i=1}^n \xi_i$$

s.t. $y_1(\mathbf{w}^T \mathbf{x}_1) \geq 1 - \xi_1, \dots, y_n(\mathbf{w}^T \mathbf{x}_n) \geq 1 - \xi_n$

Sparse Approximation Algorithm for Multivariate SVM

Approach: Sparse Approx. Structural SVM [Tsochantaridis et al. 04]

• Input: $\bar{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_n), \bar{\mathbf{y}} = (y_1, \dots, y_n), C, \epsilon$

• $\mathcal{S} \leftarrow \emptyset, \mathbf{w} \leftarrow \mathbf{0}, \xi \leftarrow 0$

• REPEAT

– compute $\bar{\mathbf{y}} = \text{argmax}_{\bar{\mathbf{y}} \in \mathcal{Y}} \{\Delta(\bar{\mathbf{y}}, \bar{\mathbf{y}}) + \mathbf{w}^T \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}})\}$

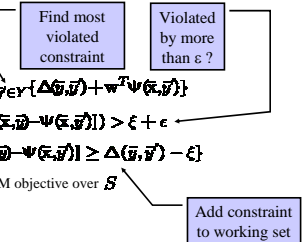
– IF $(\Delta(\bar{\mathbf{y}}, \bar{\mathbf{y}}) - \mathbf{w}^T [\Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) - \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}})]) > \xi + \epsilon$

• $\mathcal{S} \leftarrow \mathcal{S} \cup \{\Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) - \Psi(\bar{\mathbf{x}}, \bar{\mathbf{y}})\} \geq \Delta(\bar{\mathbf{y}}, \bar{\mathbf{y}}) - \xi$

• $[\mathbf{w}, \xi] \leftarrow \text{optimize SVM objective over } \mathcal{S}$

– ENDIF

• UNTIL \mathcal{S} has not changed during iteration



Polynomial Convergence Bound

• **Theorem [Tsochantaridis et al., 2004]:** The sparse-approximation algorithm finds a solution to the soft-margin optimization problem after adding at most

$$\max \left\{ \frac{2L}{\epsilon}, \frac{8Cn^2R^2L}{\epsilon^2} \right\}$$

constraints to the working set \mathcal{S} , so that the Kuhn-Tucker conditions are fulfilled up to a precision ϵ . The loss has to be bounded $0 \leq \Delta(\bar{\mathbf{y}}, \bar{\mathbf{y}}) \leq L$, and $R = \max_i \|\mathbf{x}_i\|$

ARGMAX for Contingency Table

• **Problem:**

$$\text{argmax}_{\bar{\mathbf{y}} \in \{1, -1\}^n} \{\Delta(\bar{\mathbf{y}}, \bar{\mathbf{y}}) + \mathbf{w}^T \sum_{i=1}^n y_i \mathbf{x}_i\}$$

• **Key Insight:**

– Only n^2 different contingency tables exist.
– ARGMAX for each table easy to compute via sorting.
– Time $O(n^2)$

• **Applies to:**

– Errorrate, F1, Prec@k, Rec@k, PRBEP, etc.

```

1: Input:  $\bar{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_n), \bar{\mathbf{y}} = (y_1, \dots, y_n), \mathcal{Y}$ 
2:  $\{i_1, \dots, i_{\text{pos}}\} \leftarrow \text{sort}\{i : y_i = 1\}$  by  $\mathbf{w}^T \mathbf{x}_i$ 
3:  $\{j_1, \dots, j_{\text{neg}}\} \leftarrow \text{sort}\{i : y_i = -1\}$  by  $\mathbf{w}^T \mathbf{x}_i$ 
4: for  $a \in \{0, \dots, \#i_{\text{pos}}\}$  do
5:    $c \leftarrow \#i_{\text{pos}} - a$ 
6:   set  $y_{i_1}, \dots, y_{i_a}$  to 1
7:   set  $y_{j_1}, \dots, y_{j_c}$  to -1
8:   for  $d \in \{0, \dots, \#i_{\text{neg}}\}$  do
9:      $b \leftarrow \#i_{\text{neg}} - d$ 
10:    set  $y_{i_1}, \dots, y_{i_a}$  to 1
11:    set  $y_{j_1}, \dots, y_{j_c}$  to -1
12:     $v \leftarrow \Delta(a, b, c, d) + \mathbf{w}^T \sum_{i=1}^n y_i \mathbf{x}_i$ 
13:    if  $v$  is the largest so far then
14:       $\bar{\mathbf{y}} \leftarrow (y_1, \dots, y_n)$ 
15:    end if
16:  end for
17: end for
18: return( $\bar{\mathbf{y}}$ )
    
```

ARGMAX for ROC-Area

Problem:

$$\arg\max_{\vec{y} \in \{1,-1\}^n} (\Delta(\vec{x}, \vec{y}) + w^T \sum_{i=1}^n y_i x_i)$$

Key Insight:

- ROC Area is proportional to "swapped pairs"
- Loss decomposes linearly over pairs
- Find argmax via sort in time $O(n \log n)$
- Represent n^2 pairs as

$$\Psi(\vec{x}, \vec{y}) = \sum_{i=1}^n c_i x_i \text{ with } c_i = \begin{cases} \sum_{j=1}^{\#neg} y_j, & \text{if } (y_i = 1) \\ -\sum_{j=1}^{\#pos} y_j, & \text{if } (y_i = -1) \end{cases}$$

```

1: Input:  $\vec{x} = (x_1, \dots, x_n)$ ,  $\vec{y} = (y_1, \dots, y_n)$ 
2: For  $i \in \{i : y_i = 1\}$  do  $a_i \leftarrow -0.25 + w^T x_i$ 
3: For  $i \in \{i : y_i = -1\}$  do  $a_i \leftarrow 0.25 + w^T x_i$ 
4:  $(r_1, \dots, r_n) \leftarrow \text{sort } \{1, \dots, n\}$  by  $a_i$ 
5:  $a_p = \#pos$ ,  $a_n = 0$ 
6: For  $i \in \{1, \dots, n\}$  do
7:   If  $y_{r_i} > 0$  then
8:      $a_i \leftarrow (\#neg - 2 a_n)$ 
9:      $a_p \leftarrow a_p - 1$ 
10:  else
11:     $a_i \leftarrow (-\#pos + 2 a_p)$ 
12:     $a_n \leftarrow a_n + 1$ 
13:  end if
14: end for
15: return  $(c_1, \dots, c_n)$ 
    
```

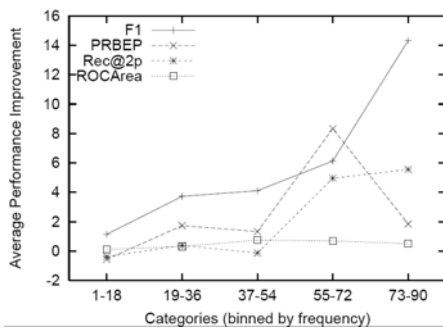
Experiment: Generalization Performance

Experiment Setup

- Macro-average over all classes in dataset
- Baseline: classification SVM with linear cost model
- Select C and cost ratio j via 2/3 - 1/3 holdout test
- Two-tailed Wilcoxon (**=95%, *=90%)

Dataset	Method	F_1	PRBEP	Rec@2p	ROCArea
Reuters (90 classes) Examples: 9603/3299 Features: 27658	SVM _{multi} $^\Delta$	62.0	68.2	78.3	99.1
	SVM _{org}	56.1	65.7	77.2	98.6
	win/lose	(51/20)**	(16/8)**	(14/8)	(43/33)*
Arxiv (14 classes) Examples: 1168/32487 Features: 13525	SVM _{multi} $^\Delta$	58.8	58.4	73.3	92.8
	SVM _{org}	49.6	57.9	74.4	92.7
	win/lose	(9/5)*	(9/4)	(1/13)**	(8/6)
Optdigits (10 classes) Examples: 3823/1797 Features: 64	SVM _{multi} $^\Delta$	92.5	92.7	98.4	99.4
	SVM _{org}	91.5	91.5	98.7	99.4
	win/lose	(8/2)*	(5/1)*	(1/5)	(6/4)
Covertype (7 classes) Examples: 1000/2000 Features: 54	SVM _{multi} $^\Delta$	73.8	72.1	93.1	94.6
	SVM _{org}	73.9	71.0	94.7	94.1
	win/lose	(3/4)	(5/2)	(2/5)	(4/3)

Experiment: Unbalanced Classes in Reuters



Experiment: Number of SV

Corollary: For error rate as the loss function, the hard-margin solution after the first iteration is equal to Rocchio Algorithm.

$$\mathbf{w}_1 \sim \sum_{\{pos: y_{pos}=1\}} \mathbf{x}_{pos} - \sum_{\{neg: y_{neg}=-1\}} \mathbf{x}_{neg}$$

Dataset	Method	F_1	PRBEP	Rec@2p	ROCArea	Err
Reuters (90 classes) Examples: 9603/3299 Features: 27658	SVM _{multi} $^\Delta$	62.0	45.0	46.3	5.1	86.3
	SVM _{org}					371.4
Arxiv (14 classes) Examples: 1168/32487 Features: 13525	SVM _{multi} $^\Delta$	129.5	43.4	45.3	26.8	177.7
	SVM _{org}					645.3
Optdigits (10 classes) Examples: 3823/1797 Features: 64	SVM _{multi} $^\Delta$	19.6	14.6	14.0	3.9	25.0
	SVM _{org}					556.9
Covertype (7 classes) Examples: 1000/2000 Features: 54	SVM _{multi} $^\Delta$	12.5	12.0	9.4	5.0	17.1
	SVM _{org}					372.8

Conclusions

- Generalization of SVM to multivariate loss functions**
 - Classification SVMs are special case
- Polynomial time training algorithms for**
 - any loss function based on contingency table.
 - ROC-Area.
- New representation of SVM optimization problem**
 - Support Vectors represent vector of classifications
 - Can be extremely sparse
- Future work**
 - Other performance measures, other methods (e.g. boosting)
 - Faster training algorithm exploiting special structure