

SVM with Multi Class

Note Title

11/25/2006

Task: learn a classifier $h: \mathcal{X} \rightarrow \mathcal{Y}$

$$\text{Hypothesis: } h_{\underline{\omega}}(\underline{x}) = \arg \max_{\underline{y}} \sum_{i=1}^n w_i \cdot f_i(\underline{x}, \underline{y}) = \arg \max_{\underline{y}} \underline{\omega}^T f(\underline{x}, \underline{y})$$

Task: select weights $\{\omega_i\}$ using training data.

$$\left\{ (\underline{x}^{(i)}, \underline{t}^{(i)} = \underline{t}(\underline{x}^{(i)})) : i = 1 \dots N \right\} \quad \underline{t}(\underline{x}) = \text{true label of } \underline{x}.$$

$$\text{Define: } \Delta f_x(y) = \underline{f}(\underline{x}, \underline{t}(\underline{x})) - f(\underline{x}, y)$$

$$\text{Minimize}_{\underline{\omega}} \quad \frac{1}{2} \|\underline{\omega}\|^2 + C \sum_{\underline{x}} \xi_{\underline{x}}$$

$$\text{s.t.} \quad \underline{\omega} \cdot \Delta f_x(y) \geq \Delta t_x(y) - \xi_x \quad \forall \underline{x}, \underline{y}$$

Condition $\Rightarrow \xi_x \geq 0$, by setting $y = \underline{t}(\underline{x})$, $\Delta f_x(t(x)) = 0 = \Delta t_x(t(x))$.

$$\Delta t_x(y) = 1, \text{ if } y \neq \underline{t}(\underline{x}) \\ 0, \text{ otherwise.}$$

Primal Problem

$$L_p(\underline{\omega}, \underline{\xi}: \underline{\alpha}_x)$$

$$= \frac{1}{2} \|\underline{\omega}\|^2 + C \sum_{\underline{x}} \xi_{\underline{x}}$$

$$- \sum_{\underline{x}, y} \underline{\alpha}_x(y) \underline{\omega} \cdot \Delta f_x(y) - \Delta t_x(y) + \xi_x$$

Dual Problem

$$L_d(\underline{\alpha}_x) = \sum_{\underline{x}, y} \underline{\alpha}_x(y) \Delta t_x(y)$$

$$- \frac{1}{2} \left\| \sum_{\underline{x}, y} \underline{\alpha}_x(y) \Delta f_x(y) \right\|^2$$

$$\text{s.t.} \quad \sum_y \underline{\alpha}_x(y) = C, \quad \forall x. \quad \underline{\alpha}_x(y) \geq 0.$$

E.G.

a b | c | d | e | f

(2) Special Case: Markov Network.

$h: \underline{X} \rightarrow \underline{Y}$

Watasee

$\underline{Y} = Y_1 \times \dots \times Y_n$

$S = \{(\underline{x}_i^i, y_i^i) : i = 1 \dots n\}$

Let \underline{y}^i be multivalued.

in form: $y^i = (y_{1i}^i, \dots, y_{ki}^i)$

with $y_{ai}^i \in \{\pm 1\}$, for each a .

basis function:

$f_i: \underline{X} \times \underline{Y} \rightarrow \text{IR}$.

classification function:

$h_{\underline{\omega}}(\underline{x}) = \arg \max_{\underline{y}} \sum_{i=1}^n \omega_i f_i(\underline{x}, \underline{y})$

E.G. OCR example

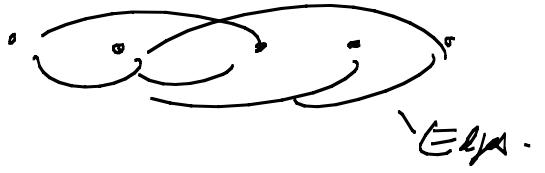
Each y_i is a character

\underline{y} is a full word

webpage classification:

Each y_i is a webpage label.

\underline{y} is a label for an entire website.



(3)

Set of features.

$$f_h(\underline{x}, \underline{y}) = \sum_{(i,j) \in E} f_h(\underline{x}, y_i, y_j)$$

Penalty: $\Delta t_x(y) = \sum_{i=1}^l \Delta t_x(y_i)$ E - edges of graph.

Primal Problem

$$\min \frac{1}{2} \|\underline{\omega}\|^2 + C \sum_{\underline{x}} \underline{\xi}_{\underline{x}}$$

$$\text{s.t. } \underline{\omega}^\top \Delta f_x(\underline{y}) > \Delta t_x(\underline{y}) - \underline{\xi}_{\underline{x}}, \forall \underline{x}, \underline{y}.$$

Dual Problem

$$\max \sum_{\underline{x}, \underline{y}} \underline{\alpha}_x(\underline{y}) \Delta t_x(\underline{y}) - \frac{1}{2} \left\| \sum_{\underline{x}, \underline{y}} \underline{\alpha}_x(\underline{y}) \Delta f_x(\underline{y}) \right\|^2$$

$$\text{s.t. } \sum_{\underline{y}} \underline{\alpha}_x(\underline{y}) = C, \forall \underline{x}, \quad \underline{\alpha}_x(\underline{y}) \geq 0, \forall \underline{x}, \underline{y}.$$

Now we can interpret $\underline{\alpha}_x(\underline{y})$ as a probability distribution \rightarrow with normalization C.

The dual function is a function of the expectation of $\Delta t_x(\underline{y})$ & $\Delta f_x(\underline{y})$.

$$\Delta t_x(\underline{y}) = \sum_i \Delta t_x(y_i), \quad \Delta f_x(\underline{y}) = \sum_{(i,j) \in E} \Delta f_x(y_i, y_j)$$

(4) The dual problem reduces to a formulation on the marginal unary and pairwise distributions

$$\underline{\mu}_x(y_{i,j}) = \sum_{y \sim [y_i, y_j]} \alpha_x(y), \quad \forall (i,j) \in E$$

$$\text{and } y_i, y_j \in \underline{Y}_x$$

$$\underline{\mu}_x(y_i) = \sum_{y \sim [y_i]} \alpha_x(y), \quad \forall i, \forall y_i, \forall x.$$

$y \sim [y_i, y_j]$

full assignment

y consistent with

partial assigned

y_i, y_j .

Must enforce consistency

constraint between the $\underline{\mu}_x(y_i, y_j)$

and $\underline{\mu}_x(y_i)$

$$\sum_{y_j} \underline{\mu}_x(y_i, y_j) = \underline{\mu}_x(y_i),$$

$$\forall y_j$$

$$\forall (i,j) \in E.$$

Because:

$$\sum_y \alpha_x(y) \Delta_x(y) = \sum_y \sum_i \alpha_x(y) \Delta t_x(y_i) = \sum_{i,y_i} \Delta t_x(y_i) \sum_y \alpha_x(y)$$

$$= \sum_{i,y_i} \underline{\mu}_x(y_i) \Delta t_x(y_i)$$

Similarly, the second term can be expressed in form

$$-\frac{1}{2} \sum_{x,x' \in E} \sum_{y_i} \sum_{y_j} \underline{\mu}_x(y_i, y_j) \underline{\mu}_{x'}(y_i, y_j) \Delta f_x(y_i, y_j)^T \Delta f_{x'}(y_i, y_j)$$

$$+ \sum_x \sum_{i,y_i} \underline{\mu}_x(y_i) \Delta t_x(y_i).$$

Algorithm can maximize this, by variant of
Belief Propagation

(6) to 7

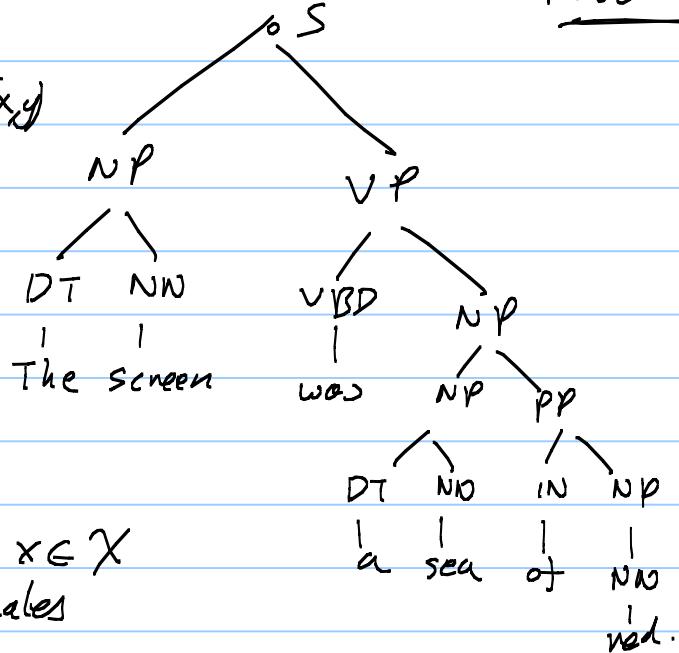
(5)

Max-Margin Parsing.

Production Rules

$$f_{\underline{w}}(x) = \text{argmax}_{y \in G(x)} \underline{\Phi}(x, y)$$

Input x ,
 $G(x)$ is a set of
 pauses.



G maps an input $x \in X$
 to a set of candidates
 parses $\underline{G}(x) \subseteq Y$

$$L : X \times \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$$

$L(x, y, \hat{y})$ is the penalty for making pause y
 on data x when tree parse is \hat{y} .

Want to estimate a linear discriminant form:

$$f_{\underline{w}}(x) = \text{argmax}_{y \in G(x)} \underline{w} \cdot \underline{\Phi}(x, y)$$

The arg max can be performed by
 dynamic programming.

(6)

$$y_c = \arg \max_{y \in G(x)} \underline{w} \cdot \underline{\Phi}(x, y)$$

The margin of the parameters w on example i and parse y

$$\underline{w} \cdot \underline{\Phi}(x_i, y_i) - \underline{w} \cdot \underline{\Phi}(x_i, y) = \underline{w} \cdot (\underline{\Phi}_{i,y_i} - \underline{\Phi}_{i,y})$$

Primal.

$$\min \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

y_i is the true
parse.

$$\text{s.t. } \underline{w} \cdot (\underline{\Phi}_{i,y_i} - \underline{\Phi}_{i,y}) \geq \underline{L}_{i,y} - \xi_i \quad \forall y \in G(x_i)$$
$$\underline{L}_{i,y} = L(x_i, y_i, y)$$

Dual.

$$\max C \sum_{i,y} \alpha_{i,y} \underline{L}_{i,y} - \frac{1}{2} \|C \sum_{i,y} (\underline{L}_{i,y} - \alpha_{i,y}) \underline{\Phi}_{i,y}\|^2$$

$$\text{s.t. } \sum_y \alpha_{i,y} = 1, \quad \forall i, \quad \alpha_{i,y} \geq 0, \quad \forall i, y$$

$\underline{L}_{i,y} = I(x_i, y_i, y)$ indicates
whether y is the true parse y_i .

$$\underline{w}^* = C \sum_{i,y} (\underline{L}_{i,y} - \alpha_{i,y}^+) \underline{\Phi}_{i,y}$$

(At 5)

(7)

PCFG.

Chomsky Normal Form (CNF)

$$A \rightarrow BC$$

A, B, C non-terminal symbols
a terminal symbol

$$A \rightarrow a$$

Two types of parts

First type

$$(A, S, e, i)$$

non-terminal A e.g.

start-part S NP, S, S 'is a sea'

end-part e

sentence i.

e.g.

$$S \rightarrow NP, VP, O, S, T$$

Second type:

$$A \rightarrow BC, S, m, e, i$$

Split part m

Countable set of parts R,

$R(x, y)$ is the set of parts belonging to a particular parse

All rules are in binary-branching form

so $|R(x, y)|$ is constant across different derivations y
for the same input x

$$\Phi(x, y) = \sum_{r \in R(x, y)} \phi(x, r)$$

$$L(x, y, \hat{y}) = \sum_{r \in R(x, y)} l(x, y, r)$$

One possibility - define $\ell(x, y, r) = 0$ only if
the non-terminal A spans words $s, \dots e$
in the derivation y and 1 otherwise

Indicator variable $I(x, y, r) = 1$, if $r \in R(x, y)$
 $= 0$, otherwise

Dual:

$$C \sum_i E_{x_i} [L_{i,y}] - \frac{1}{2} \| C \sum_i \Phi_{i,y_i} - E_x [\Phi_{i,y}] \|^2$$