

## Introduction to Graphical Models

### lecture 9 - Structural Support Vector Machines

Ulf Brefeld  
TU Berlin

- support vector machines
- multi-class classification

1/??

## Recall:

- relation between
  - $P(y|x)$
  - model  $f(x, y)$
  - log-Viterbi algorithm
- intuition
  - $f(x, y)$  = how good does  $y$  fits to  $x$
  - log-Viterbi: find top-scoring  $y$

2/??

## Towards Structured Support Vector Machines

- Add confidence to decision
- Incorporate arbitrary (structured) loss functions
- Impact of ordering resolved by quadratic programming

3/??

## Confidence Term

- Perceptron:

$$\langle \mathbf{w}, \Phi(\mathbf{x}, [N, V, A, N]) \rangle > \langle \mathbf{w}, \Phi(\mathbf{x}, [A, A, A, A]) \rangle$$

$$\langle \mathbf{w}, \Phi(\mathbf{x}, [N, V, A, N]) \rangle > \langle \mathbf{w}, \Phi(\mathbf{x}, [A, A, A, N]) \rangle$$

$$\langle \mathbf{w}, \Phi(\mathbf{x}, [N, V, A, N]) \rangle > \langle \mathbf{w}, \Phi(\mathbf{x}, [A, A, N, A]) \rangle$$

$$\vdots > \vdots$$

- Now, add a confidence  $\bar{\gamma}$ :

$$\langle \mathbf{w}, \Phi(\mathbf{x}, [N, V, A, N]) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}, [A, A, A, A]) \rangle \geq \bar{\gamma}$$

$$\langle \mathbf{w}, \Phi(\mathbf{x}, [N, V, A, N]) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}, [A, A, A, N]) \rangle \geq \bar{\gamma}$$

$$\langle \mathbf{w}, \Phi(\mathbf{x}, [N, V, A, N]) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}, [A, A, N, A]) \rangle \geq \bar{\gamma}$$

$$\vdots \qquad \qquad \qquad \vdots$$

4/??

## Optimization Problem

$$\begin{aligned} \max_{\bar{\gamma}, \mathbf{w}} \quad & \frac{\bar{\gamma}}{\|\mathbf{w}\|} \\ \text{s.t.} \quad & \forall_{i=1}^n, \forall_{\bar{y} \neq y_i} : \langle \mathbf{w}, \Phi(\mathbf{x}_i, y_i) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}_i, \bar{y}) \rangle \geq \bar{\gamma} \end{aligned}$$

- We call
  - $\mathbf{w}$  the weight vector
  - $\bar{\gamma}$  the functional margin
  - $\gamma = \frac{\bar{\gamma}}{\|\mathbf{w}\|}$  the geometrical margin
- Problem:  $\bar{\gamma}$  and  $\mathbf{w}$  interdepend!
  - Remedy: fix one, solve for the other
  - Common approach:  $\bar{\gamma} = 1$ .

5/??

## Structural Hard-margin SVM

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \forall_{i=1}^n, \forall_{\bar{y} \neq y_i} : \langle \mathbf{w}, \Phi(\mathbf{x}_i, y_i) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}_i, \bar{y}) \rangle \geq 1 \end{aligned}$$

- Converges only when data is linear separable
- Remedy: allow for pointwise relaxations of the margin constraint
  - introduce slack variables  $\xi_i$  for input examples

6/??

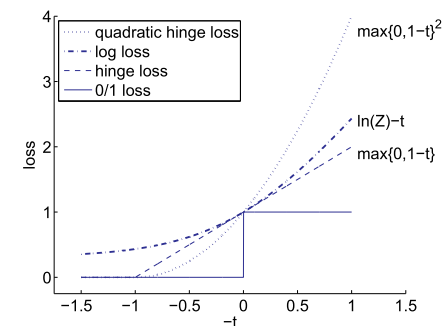
## Structural Soft-margin SVM

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \forall_{i=1}^n, \forall_{\bar{y} \neq y_i} : \langle \mathbf{w}, \Phi(\mathbf{x}_i, y_i) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}_i, \bar{y}) \rangle \geq 1 - \xi_i \\ & \forall_{i=1}^n : \xi_i \geq 0 \end{aligned}$$

- Maximize margin between true  $y_i$  and best runner-up  $\bar{y}$ 
  - Sum of slacks upper bounds 0/1 loss
  - Trade-off parameter  $C > 0$
- Alternative formulation:
  - slack  $\xi_{i\bar{y}}$  are bound to constraint  $\langle \mathbf{w}, \Phi(\mathbf{x}_i, y_i) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}_i, \bar{y}) \rangle$
  - computationally demanding

7/??

## Hinge-loss



- SVM implicitly implements a hinge loss (solve for slacks)
- Hinge loss can be rescaled to incorporate arbitrary loss functions
  - Let  $\Delta(y_i, \hat{y})$  denote a structural loss.
  - $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_0^+$ .
  - $\Delta(y_i, y_i) = 0$

8/??

## Exemplary Loss Functions

- 0/1 loss:  $\Delta(\mathbf{y}, \bar{\mathbf{y}}) = [[\mathbf{y} == \bar{\mathbf{y}}]]$

- Hamming loss for sequences

$$\Delta(\mathbf{y}, \bar{\mathbf{y}}) = T - \sum_{t=1}^T [[y_t == \bar{y}_t]]$$

- Property: decomposes across the cliques!

9/??

## Risk Minimization

- We want to minimize the theoretical risk (the generalization error)

$$R(f) = \int_{\mathbf{x} \times \mathcal{Y}} \Delta(\mathbf{y}, \operatorname{argmax}_{\bar{\mathbf{y}}} f(\mathbf{x}, \bar{\mathbf{y}})) dP(\mathbf{x}, \mathbf{y})$$

- In general, we don't know  $P(\mathbf{x}, \mathbf{y})$ 
  - Remedy: Use training sample instead!

- Minimize the empirical risk

$$\hat{R}(f) = \sum_{i=1}^n \Delta(\mathbf{y}_i, \operatorname{argmax}_{\bar{\mathbf{y}}} f(\mathbf{x}_i, \bar{\mathbf{y}}))$$

10/??

## Idea

- SVMs minimize the (regularized) empirical risk:

$$\hat{R}(f) = \sum_{i=1}^n \Delta(\mathbf{y}_i, \operatorname{argmax}_{\bar{\mathbf{y}}} f(\mathbf{x}_i, \bar{\mathbf{y}}))$$

- Sum of slacks  $\sum_i \xi_i$  upper bounds empirical risk
- Slack variable  $\xi_i$  denotes the error for input  $\mathbf{x}_i$ 
  - Now: Find maximal error wrt  $\Delta$

11/??

## Margin-rescaling

- Taskar et al. (2004)
- Rescale the (functional) margin by actual loss

$$\min_{\mathbf{w}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad \forall_{i=1}^n, \forall_{\bar{\mathbf{y}} \neq \mathbf{y}_i} : \langle \mathbf{w}, \Phi(\mathbf{x}_i, \mathbf{y}_i) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}_i, \bar{\mathbf{y}}) \rangle \geq \Delta(\mathbf{y}_i, \bar{\mathbf{y}}) - \xi_i$$

$$\forall_{i=1}^n : \xi_i \geq 0$$

- Implicit hinge loss upper bounds  $\Delta$
- Most strongly violated constraint:

$$\operatorname{argmax}_{\bar{\mathbf{y}} \neq \mathbf{y}_i} \left( \underbrace{\Delta(\mathbf{y}_i, \bar{\mathbf{y}}) - (\langle \mathbf{w}, \Phi(\mathbf{x}_i, \mathbf{y}_i) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}_i, \bar{\mathbf{y}}) \rangle)}_{\xi_i} \right)$$

12/??

## Slack-rescaling

- Tsochantaridis et al. (2005)
- Rescale slack variables by actual loss

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \forall_{i=1}^n, \forall_{\bar{y} \neq y_i} : \langle \mathbf{w}, \Phi(\mathbf{x}_i, y_i) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}_i, \bar{y}) \rangle \geq 1 - \frac{\xi_i}{\Delta(\mathbf{y}_i, \bar{y})} \\ & \forall_{i=1}^n : \xi_i \geq 0 \end{aligned}$$

- Implicit hinge loss upper bounds  $\Delta$
- Most strongly violated constraint:

$$\operatorname{argmax}_{\bar{y} \neq y_i} \left( \underbrace{1 - \Delta(\mathbf{y}_i, \bar{y}) \times (\langle \mathbf{w}, \Phi(\mathbf{x}_i, y_i) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}_i, \bar{y}) \rangle)}_{\xi_i} \right)$$

13/??

## Implications

- Loss  $\Delta$  decomposes across the cliques of the graph
  - Margin-rescaling is easily integrated into inference
  - Slack-rescaling difficult
- Loss not decomposable
  - Both difficult!
- In practice, slack-rescaling often better than margin-rescaling
  - rarely applicable (needs good approximation or enumerable sets)

14/??

## Example

- Margin rescaling for sequences / Viterbi algorithm
- Remainder: 0/1 loss for simplicity

15/??

## Towards Dual SVMs

- Integrate constraints into objective
  - Apply Lagrange's Theorem
  - Lagrange multipliers  $\alpha_i(\bar{y})$  and  $\mu_i$
- Build Lagrangian:

$$\begin{aligned} L = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ & - \sum_{i=1}^n \sum_{\bar{y} \neq y_i} \alpha_i(\bar{y}) \langle \mathbf{w}, \Phi(\mathbf{x}_i, y_i) \rangle - \langle \mathbf{w}, \Phi(\mathbf{x}_i, \bar{y}) \rangle - 1 + \xi_i \\ & - \sum_{i=1}^n \beta_i \xi_i \end{aligned}$$

- Minimum of Lagrangian is a saddle-point
  - max wrt  $\alpha, \mu$ , min wrt  $\mathbf{w}, \xi$

16/??

## Partial Derivatives: $\xi_i$

- Compute partial derivatives wrt  $\xi$ :

$$\frac{\partial L}{\partial \xi_i} = C - \sum_{\bar{y} \neq y_i} \alpha_i(\bar{y}) - \beta_i \stackrel{!}{=} 0 \quad (1)$$

- Using the non-negativity of  $\alpha$  and  $\beta$  yields

$$\forall_{i=1}^n : 0 \leq \sum_{\bar{y} \neq y_i} \alpha_i(\bar{y}) \leq C$$

17/??

## Partial Derivatives: $w$

- Compute partial derivatives wrt  $w$ :

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^n \sum_{\bar{y} \neq y_i} \alpha_i(\bar{y}) (\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{y})) \stackrel{!}{=} 0.$$

- We obtain:

$$w = \sum_{i=1}^n \sum_{\bar{y} \neq y_i} \alpha_i(\bar{y}) (\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{y}))$$

- Recall: dual perceptron!
  - Definition of  $w$  is equivalent
  - $\alpha$ 's act like counters

18/??

## Putting all together...

- Dual SVM for structured output variables:

$$\max_{\alpha} \sum_{i=1}^n \sum_{\bar{y} \neq y_i} \alpha_i(\bar{y}) - \frac{1}{2} \sum_{i,j=1}^n \sum_{\bar{y} \neq y_i} \sum_{\bar{y} \neq y_j} (\alpha_i(\bar{y}) \alpha_j(\bar{y})) \langle \Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{y}), \Phi(\mathbf{x}_j, \mathbf{y}_j) - \Phi(\mathbf{x}_j, \bar{y}) \rangle$$

– subject to the constraints

$$\forall_{i=1}^n \sum_{\bar{y} \neq y_i} \alpha_i(\bar{y}) \leq C$$

$$\forall_{i=1}^n \forall i : \alpha_i(\bar{y}) \geq 0.$$

19/??

## Optimization

- Dual SVM is a QP (quadratic program)
- Sketch of algorithm:
  - 1 Loop
  - 2 Loop  $i = 1, \dots, n$
  - 3 Compute most strongly violated constraint  $\bar{y}$
  - 4 If constraint causes more slack than before ( $\xi_i$ )
  - 5 Add constraint  $\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{y})$  to working set,
  - 6 Initialize  $\alpha_i(\bar{y})$
  - 8 Solve QP
  - 9 Until convergence

(For details see Tsochantaridis et al. JMLR,2005)

20/??

## Comments

- Working set contains active constraints
  - Only a few of the constraints will be active
  - remove constraints if  $\alpha_i(\bar{y}) = 0$  after optimization
- Convergence in polynomial time
  - if computation of argmax takes at most polynomial time
  - proof: see Tsochantaridis et al. (2005)
- Implementation on the web:
  - SVM<sup>struct</sup> by Thorsten Joachims (extension of SVM<sup>light</sup>)
  - [http://svmlight.joachims.org/svm\\_struct.html](http://svmlight.joachims.org/svm_struct.html)

21/??

## Named Entity Recognition

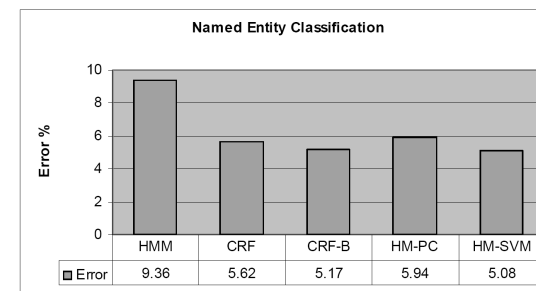


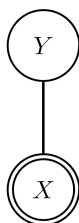
Figure 1. Test error of NER task over a window of size 3 using 5-fold cross validation.

Example: Como (O) contrapartida (O) Deutsche (C-B) Telekom (C-I) vender (O) al (O) consorcio (O) francs (O) su (O) participacion (O) del (O) por (O) ciento (O) en (O) el (O) empresa (O) mixta (O) britnica (O) MetroHoldings (C-B).  
(see Altun et al. (2003))

22/??

## CRF: Special Case

- Sequential approach =
  - a chain of multi-class classification problems
  - linked via label-label transitions



23/??

## Multi-class Classification

- Let  $\psi(\mathbf{x})$  the feature vector of observation  $\mathbf{x}$
- Let  $\Sigma = \{\sigma_1, \dots, \sigma_k\}$  be the set of classes
- Then, the joint feature map is given by

$$\Phi(\mathbf{x}, y) = ([y = \sigma_1]\psi(\mathbf{x})^\top, \dots, [y = \sigma_k]\psi(\mathbf{x})^\top)^\top.$$

- dimension of  $\Phi$  is  $\dim(\Phi) = |\Sigma|\dim(\psi)$
- see also Crammer & Singer (2001), Weston & Watkins (1998)

24/??

## Outlook

- What happens if only observations but no labels are observed?
  - i.e., training set consists of  $n$  observations  $\mathbf{x}_1, \dots, \mathbf{x}_n$
- Probabilistic planning...