

(1)

Multi-Class Multi-State SVM

note title

11/21/2011

$$L(\underline{w}, \underline{z}; \underline{\alpha}) = \frac{1}{2} \|\underline{w}\|^2 + C \sum_i z_i - \sum_i \sum_y \alpha_y^i \{ z_i - \ell(y, y_i) - \underline{w} \cdot \underline{\phi}(x_i, y) + \underline{w} \cdot \underline{\phi}(x_i, y_i) \}$$

solution $\hat{y}(x; \underline{w}) = \underset{y}{\operatorname{argmax}} \underline{w} \cdot \underline{\phi}(x; y)$

Constraint $z_i - \ell(y, y_i) - \underline{w} \cdot \underline{\phi}(x_i, y) + \underline{w} \cdot \underline{\phi}(x_i, y_i) \geq 0$
 $z_i \geq \max_y \{ \ell(y, y_i) + \underline{w} \cdot \underline{\phi}(x_i, y) - \underline{w} \cdot \underline{\phi}(x_i, y_i) \}$

$$\frac{1}{2} \|\underline{w}\|^2 + C \sum_i \max_y \{ \ell(y, y_i) + \underline{w} \cdot \underline{\phi}(x_i, y) - \underline{w} \cdot \underline{\phi}(x_i, y_i) \}$$

Note: no need to separately impose $z_i \geq 0$ because if we set $y = y_i$, we see $z_i \geq \ell(y_i, y_i) + \underline{w} \cdot \underline{\phi}(x_i, y_i) - \underline{w} \cdot \underline{\phi}(x_i, y_i) = 0$

Solve the primal problem

$$\frac{\partial}{\partial \underline{w}} L_p(\underline{w}, \underline{z}; \underline{\alpha}) = 0$$

$$\Rightarrow \underline{w} = \sum_i \sum_y \alpha_y^i \{ \underline{\phi}(x_i, y_i) - \underline{\phi}(x_i, y) \}$$

like support vectors

note: if $\underline{w} \cdot \underline{\phi}(x_i, y_i) > \underline{w} \cdot \underline{\phi}(x_i, y) + \ell(y, y_i)$ then $\alpha_y^i = 0$

$$\frac{\partial}{\partial z_i} L_p(\underline{w}, \underline{z}; \underline{\alpha}) = 0$$

$$\Rightarrow C = \sum_y \alpha_y^i$$

so $\frac{\alpha_y^i}{C}$ is a probability distribution.

This gives solution $\underline{w}(\underline{\alpha}), \underline{z}(\underline{\alpha})$

substituting back gives the dual energy:

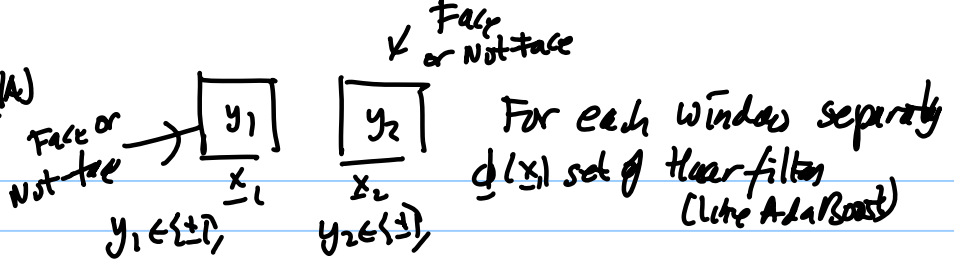
$$L_d(\underline{\alpha}) = L_p(\underline{w}(\underline{\alpha}), \underline{z}(\underline{\alpha}); \underline{\alpha}) = \sum_i \sum_y \alpha_y^i \ell(y, y_i) - \frac{1}{2} \sum_{i,j} \sum_{y_1, y_2} \alpha_{y_1}^i \alpha_{y_2}^j \{ \underline{\phi}(x_i, y_1) - \underline{\phi}(x_i, y_2) \} \cdot \{ \underline{\phi}(x_j, y_2) - \underline{\phi}(x_j, y_1) \}$$

The case of binary classification can be recovered by setting $\ell(y, y_i) = 1$, if $y \neq y_i$
 $= 0$, if $y = y_i$

$$\underline{\phi}(x, y) = y \underline{\phi}(x)$$

(2)

Example (A)



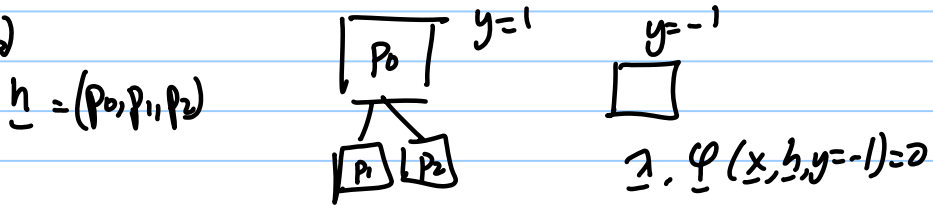
Together

$$\underline{y} = (y_1, y_2) \quad \underline{x} = (x_1, x_2)$$

$$\underline{\psi}(y, x) = (\delta_{y_{1,1}}, \delta_{y_{2,1}}, \delta_{y_{1,0}}, \delta_{y_{2,0}}, \delta_{y_{1,1}}, \delta_{y_{1,0}}, \delta_{y_{2,1}}, \delta_{y_{2,0}})$$

$$\omega \cdot \underline{\phi}(y, x) = y^1 \omega^1 \cdot \underline{\phi}(x^1) + y^2 \omega^2 \cdot \underline{\phi}(x^2) + \omega^3 \cdot \underline{\psi}(y_1, y_2)$$

Example (B)



$$\lambda \cdot \underline{\phi}(x, h, y = 1) \quad h = p_0 p_1 p_2$$

$$= \lambda \cdot \underline{\phi}(x, p_0 p_1 p_2, y = 1)$$

$$= \lambda_1 \cdot \langle (p_0 - p_1), (p_0 - p_1) \rangle$$

$$+ \lambda_2 \cdot \langle (p_0 - p_1), (p_0 - p_2) \rangle$$

$$+ \mu^0 \cdot \underline{\phi}(x, p_0)$$

$$+ \mu^1 \cdot \underline{\phi}(x, p_1)$$

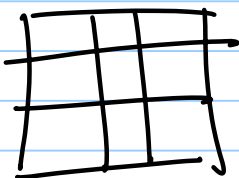
$$+ \mu^2 \cdot \underline{\phi}(x, p_2)$$

During learning the position p_0, p_1, p_2 of the object parts are unspecified

$$\underline{x} = \{x_{ij} : i, j \in \mathbb{D}\}$$

$$\underline{x}^P = \{x_{ij} : i, j \in \mathbb{D}^P\}$$

\mathbb{D} is the image lattice
 x_{ij} is the intensity at lattice site ij
 \mathbb{D}^P is a subregion of \mathbb{D} .



filters

$$\mathbb{D}^P = \bigcup_{\alpha=1}^n \mathbb{D}_\alpha^P$$

$$\underline{f} = (f_1, \dots, f_n)$$

$$\mathbb{D}^P$$

$$f_i = f(x_{ij} \in \mathbb{D}_i^P)$$

$$= f(x_{ij} : (ij) \in \mathbb{D}_i^P)$$

Quantize the response values

$$f(x_{ij} \in \mathbb{D}^P) \in \{\alpha : \alpha \in \Lambda\}$$

Total response: histogram $n_\alpha = \sum_{ij} \delta_{f(x_{ij}), \alpha}$

Quantize

cluster the set of values $\{f_{ij}\}$

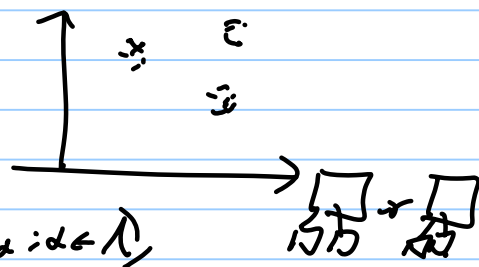
use k-means

each datapoint (feature values) is

associated to the closest mean $\{m_\alpha : \alpha \in \Lambda\}$

Extension, to multiple models (to allow for different viewpoints)

$$\lambda \cdot \underline{\phi}(x, m^1, p^1, p^2, p^3, y = 1) \quad \mu \text{ label models}$$



Lecture 9. Machine Learning: Structure & Latents

Note Title

2/4/2010

Structure Max-Margin extends binary-classification methods so they can be applied to learn the parameters of an MRF, HMM, SFG or other model.

Recall standard SVM for binary classification.

$$R(\underline{\lambda}) = \frac{1}{2} \|\underline{\lambda}\|^2 + C \sum_{i=1}^m \max\{0, 1 - y_i \underline{\lambda} \cdot \underline{\phi}(d_i)\}$$

Training Data $\{(y_i, d_i)\}$ $y_i \in \{\pm 1\}$

e.g. to get a plane, set $\underline{\phi}(d) = d$.

Decision rule: $\hat{y}_i(\underline{\lambda}) = \arg \max_y y \underline{\lambda} \cdot \underline{\phi}(d_i) = \text{sgn} \underline{\lambda} \cdot \underline{\phi}(d_i)$

The task is to minimize $R(\underline{\lambda})$ w.r.t. $\underline{\lambda}$ which maximizes the 'margin' $1/\|\underline{\lambda}\|$.

Here is a more general formulation that can be used if the output variable y is a vector $y = (y_1, \dots, y_n)$ - i.e. it could be the state of an MRF, an HMM, or a SFG.

$$R(\underline{\lambda}) = \frac{1}{2} \|\underline{\lambda}\|^2 + C \sum_{i=1}^m \Delta(y_i; \hat{y}_i(\underline{\lambda}))$$

decision rule: $\hat{y}_i(\underline{\lambda}) = \arg \max_y \underline{\lambda} \cdot \underline{\Phi}(d, y)$

the error function $\Delta(y_i; \hat{y}_i(\underline{\lambda}))$ is any measure of distance between the true solution y_i and the estimate $\hat{y}_i(\underline{\lambda})$

→ to obtain binary-value. → i) set $y_i = y_i \in \{-1, 1\}$

(ii) $\underline{\Phi}(d, y) = y \underline{\Phi}(d)$

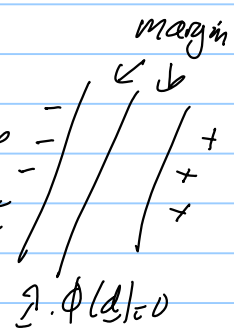
(iii) $\Delta(y_i; \hat{y}_i(\underline{\lambda})) = \max\{0, 1 - y_i \underline{\lambda} \cdot \underline{\phi}(d_i)\}$

Hinge loss ✓ because the function is 0

i) $y_i \underline{\lambda} \cdot \underline{\phi}(d_i) > 1$ (ie point is on the right side of the margin)

and the function increases linearly with $\underline{\lambda} \cdot \underline{\phi}(d_i)$

(iv) $\hat{y}_i(\underline{\lambda}) = \arg \max_y y \underline{\lambda} \cdot \underline{\phi}(d)$



(2) This more general formulation is $P(y, d) = \frac{e^{\lambda \cdot \phi(y, d)}}{Z}$
 $R(\lambda) = \frac{1}{2} \lambda^2 + C \sum_{i=1}^m \Delta(y_i; \hat{y}_i(\lambda))$

$$\hat{y}_i = \arg \max_y \lambda \cdot \phi(y, d_i)$$

(*) This requires an inference algorithm, like in the last few lectures.
 for binary classification inference only has to compute $\max_{y_i \in \{0, 1\}} y_i \lambda \cdot \phi(d_i)$ so is trivial.

(*) Also need to be able to minimize $R(\lambda)$ to find $\lambda \rightarrow$ hard because the error term $\Delta(y_i; \hat{y}_i(\lambda))$ is a highly complicated function of λ .

Modify $R(\lambda)$ to an upper bound $\bar{R}(\lambda)$

$$\bar{R}(\lambda) = \frac{1}{2} \lambda^2 + C \sum_{i=1}^m \left(\max_y \{ \Delta(y_i; \hat{y}) + \lambda \cdot \phi(d_i; \hat{y}) - \lambda \cdot \phi(d_i; y_i) \} \right)$$

which is convex in λ . hence has a single minimum.

To get this bounds use two steps:

(step 1) $\max_{\hat{y}} \{ \Delta(y_i; \hat{y}) + \lambda \cdot \phi(d_i; \hat{y}) \} \geq \Delta(y_i; \hat{y}_i(\lambda)) + \lambda \cdot \phi(d_i; \hat{y}_i(\lambda))$
 $\rightarrow \hat{y}_i(\lambda)$ maximizes $\lambda \cdot \phi$
 \rightarrow equality if it also maximizes $\Delta + \lambda \cdot \phi$

(step 2) $\lambda \cdot \phi(d_i; \hat{y}_i(\lambda)) \geq \lambda \cdot \phi(d_i; y_i)$

Note: bounds are 'tight' because if we can find a good solution then $y_i \approx \hat{y}_i(\lambda)$.

How to minimize $\bar{R}(\lambda)$?

Several Algorithms (hot topic)

Some in dual space - like original SVM for binary problem.

Simple: stochastic gradient descent

pick example (d_i, y_i)

take derivative of $\bar{R}(\lambda)$ w.r.t. λ .

$$\lambda^{t+1} = \lambda^t - \epsilon C \{ \phi(d_i; \hat{y}) - \phi(d_i; y_i) \}$$

where $\hat{y} = \arg \max_y \{ \Delta(y_i; \hat{y}) + \lambda \cdot \phi(d_i; \hat{y}) \}$

Note: inference algorithm must be adapted to compute this.

(3)

How to extend to models with latent (hidden) variables? Denote these variables by \underline{h} .

Want decision rule

$$(\hat{y}, \hat{h}) = \arg \max_{(y, h) \in Y \times H} \lambda \cdot \phi(\underline{d}, \underline{y}, \underline{h})$$

← must be computable by inference algorithm

Training data $\{(\underline{d}^i, \underline{y}^i) : i = 1, \dots, n\}$. the hidden variables are not known.

Loss function $\Delta(\underline{y}_i; \hat{y}_i(\lambda), \hat{h}_i(\lambda))$

depends on the truth \underline{y}_i
the estimate of $\underline{y}_i(\lambda), \hat{h}_i(\lambda)$ from the model

$$R(\lambda) = \frac{1}{2} |\lambda|^2 + C \sum_{i=1}^m \Delta(\underline{y}_i; \hat{y}_i(\lambda), \hat{h}_i(\lambda))$$

non-trivial function of λ

replace $R(\lambda)$ by m

$$\tilde{R}(\lambda) = \frac{1}{2} |\lambda|^2 + C \sum_{i=1}^m \left\{ \max_{(\underline{y}, \underline{h})} \left\{ \Delta(\underline{y}_i; \underline{y}, \underline{h}) + \lambda \cdot \phi(\underline{d}_i; \underline{y}, \underline{h}) \right\} - \max_{\underline{h}} \lambda \cdot \phi(\underline{d}_i; \underline{y}_i, \underline{h}) \right\}$$

$$\tilde{R}(\lambda) = \underbrace{f(\lambda)}_{\text{convex}} + \underbrace{g(\lambda)}_{\text{concave}}, \quad \text{with } g(\lambda) = -\max_{\underline{h}} \lambda \cdot \phi(\underline{d}_i; \underline{y}_i, \underline{h})$$

To show convexity and concavity

suppose $\tau(\lambda) = \sum_{i=1}^n \max_{\underline{y}_i} \lambda \cdot \phi(\underline{d}_i; \underline{y}_i)$

convex if $\tau(\alpha \lambda_1 + (1-\alpha) \lambda_2) \leq \alpha \tau(\lambda_1) + (1-\alpha) \tau(\lambda_2)$

$$\tau(\alpha \lambda_1 + (1-\alpha) \lambda_2) = \sum_{i=1}^n \max_{\underline{y}_i} \left\{ (\alpha \lambda_1 + (1-\alpha) \lambda_2) \cdot \phi(\underline{d}_i; \underline{y}_i) \right\}$$

$$\alpha \tau(\lambda_1) + (1-\alpha) \tau(\lambda_2) = \alpha \sum_{i=1}^n \max_{\underline{y}_i} \lambda_1 \cdot \phi(\underline{d}_i; \underline{y}_i) + (1-\alpha) \sum_{i=1}^n \max_{\underline{y}_i} \lambda_2 \cdot \phi(\underline{d}_i; \underline{y}_i)$$

but $\max_{\underline{y}_i} \alpha \lambda_1 \cdot \phi(\underline{d}_i; \underline{y}_i) + \max_{\underline{y}_i} (1-\alpha) \lambda_2 \cdot \phi(\underline{d}_i; \underline{y}_i) \geq \max_{\underline{y}_i} \left\{ (\alpha \lambda_1 + (1-\alpha) \lambda_2) \cdot \phi(\underline{d}_i; \underline{y}_i) \right\}$

hence $f(\cdot)$ is convex
and $g(\cdot)$ is concave.

(4)

Apply CCCP algorithm.

Two stages: step 1.

$$\frac{\partial g(\underline{\lambda})}{\partial \underline{\lambda}} = - \underline{\phi}(\underline{d}_i, \underline{y}_i, \underline{h}^*)$$

where $\underline{h}^* = \arg \max_{\underline{h}} \underline{\lambda}^t \cdot \underline{\phi}(\underline{d}_i, \underline{y}_i, \underline{h})$

$\underline{\lambda}^t$ / current estimate of $\underline{\lambda}$.

step 2. solve

$$\underline{\lambda}^{t+1} = \arg \min_{\underline{\lambda}} \left(f(\underline{\lambda}) + \underline{\lambda} \cdot \frac{\partial g(\underline{\lambda})}{\partial \underline{\lambda}} \right)$$

This reduces to a modified SVM with known ~~data~~

$$- \min_{\underline{\lambda}} \frac{1}{2} \|\underline{\lambda}\|^2 + C \sum_{i=1}^n \max \left\{ \underline{\lambda} \cdot \underline{\phi}(\underline{d}_i, \underline{y}_i, \underline{h}) + \Delta(\underline{y}_i, \underline{y}, \underline{h}) \right. \\ \left. - C \sum_{i=1}^n \underline{\lambda} \cdot \underline{\phi}(\underline{d}_i, \underline{y}_i, \underline{h}_i^*) \right\}$$

Note: Similarities to EM.

→ step 1 involves estimating the hidden state \underline{h}_i^*

→ step 2 estimates $\underline{\lambda}$

repeat.

Note: like EM there is no guarantee that this will converge to the global optimum.