# Lecture 9.    Machine Learning: Structure & Latent;

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

Recall standard SVM for binary classification.

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

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

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

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

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 variables $y$ is a vector $y = (y_1, \cdots y_w)$. — ie. it could be the state of an MRF, an HMM, or a SCFG.

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

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

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 $y_i = y_i \in \{-1, 1\}$
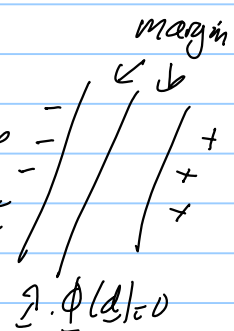(ii) $\underline{\phi}(\underline{d}, y) = y\underline{\phi}(\underline{d})$

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

Hinge loss ___✓ because the function is 0

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

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

$\underline{\lambda} \cdot \underline{\phi}(\underline{d}) = 0$

(2)  This more general formulation is $P(y,d) = \frac{1}{z} e^{\lambda \cdot \phi(y,d)}$

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

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

(*) This requires an <u>inference algorithm</u>, — for binary classification inference only has to compute max $y_i \lambda \cdot \phi(d)$ like in the last few lectures.
$y_i \in \{\pm i\}$ so is trivial.

(#) Also need to be able to minimize $R(\lambda)$ to find $\lambda \to$ hard because the error term $\Delta(\underline{y}_i : \hat{\underline{y}}_i(\lambda))$ is a highly complicated function of $\lambda$.

<u>Modify</u> $R(\lambda)$ to an upper bound $\overline{R}(\lambda)$

$$\overline{R}(\lambda) = \frac{1}{2}|\lambda|^2 + C \sum_{i=1}^m \{ \max_y \{ \Delta(\underline{y}_i : \underline{y}) + \lambda \cdot \phi(\underline{d}_i, \underline{y}) \} - \lambda \cdot \phi(\underline{d}_i : \underline{y}_i) \}$$

which is convex in $\lambda$. hence has a simple minimum.

To get this bounds use two steps:

(step 1)  $\max_{\hat{y}} \{ \Delta(\underline{y}_i : \hat{y}) + \lambda \cdot \phi(\underline{d}_i, \hat{y}) \}$
$\geq \Delta(\underline{y}_i : \hat{\underline{y}}_i(\lambda)) + \lambda \cdot \phi(\underline{d}_i : \hat{\underline{y}}_i(\lambda))$

$\to \hat{\underline{y}}_i(\lambda)$ maximizes $\lambda \cdot \phi$
$\to$ equality if it also maximizes $\Delta + \lambda \phi$

(step 2)  $\lambda \cdot \phi(\underline{d}_i, \hat{\underline{y}}_i(\lambda)) \geq \lambda \cdot \phi(\underline{d}_i, \underline{y}_i)$

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

<u>How</u> to minimize $\overline{R}(\lambda)$?

Several Algorithms  (hot topic)
Some in dual space — like original SVM for binary problems.

<u>Simple</u>:  stochastic gradient descent
pick example $(\underline{d}_i, \underline{y}_i)$
take derivative of $\overline{R}(\lambda)$ w.r.t. $\lambda$.
$$\lambda^{t+1} = \lambda^t - \epsilon C \{ \phi(\underline{d}_i, \hat{y}) - \phi(\underline{d}_i, y_i) \}$$
where $\hat{y} = $ ARG MAX $\{ \Delta(\underline{y}_i, \hat{y}) + \lambda \phi(\underline{x}_i, \hat{y}) \}$

<u>Note</u>: inference algorithm must be adapted to compute this.

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

$\leftarrow$ must be computable by inference algorithm

Want decision rule

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

Training data $\langle (\underline{d}^i, y^i) : i = 1 \text{ to } N \rangle$. the hidden variables are not known.

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

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

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

non-trivial function of $\underline{\lambda}$

replace $R(\underline{\lambda})$ by

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

$$\widetilde{R}(\underline{\lambda}) = f(\underline{\lambda}) + g(\underline{\lambda}), \quad \text{with} \quad g(\underline{\lambda}) = -\max_{h} \underline{\lambda} \cdot \phi(\underline{d}_i, y_i, h)$$

$\underset{\text{convex}}{\uparrow} \quad \underset{\text{concave}}{\uparrow}$

To show Convexity and concavity

Suppose $T(\underline{\lambda}) = \sum_{i=1}^{n} \max_{\tilde{y}_i} \underline{\lambda} \cdot \phi(\underline{d}_i, \hat{y}_i)$

convex if $T(\alpha \underline{\lambda}_1 + (1-\alpha)\underline{\lambda}_2) \leq \alpha T(\underline{\lambda}_1) + (1-\alpha) T(\underline{\lambda}_2)$

$$T(\alpha \underline{\lambda}_1 + (1-\alpha)\underline{\lambda}_2) = \sum_{i=1}^{n} \max_{\tilde{y}_i} \left\langle (\alpha \underline{\lambda}_1 + (1-\alpha)\underline{\lambda}_2) \cdot \phi(\underline{d}_i, \hat{y}_i) \right\rangle$$

$$\alpha T(\underline{\lambda}_1) + (1-\alpha) T(\underline{\lambda}_2) = \alpha \sum_{i=1}^{n} \max_{\tilde{y}_i} \underline{\lambda}_1 \cdot \phi(\underline{d}_i, \tilde{y}_i)$$
$$+ (1-\alpha) \sum_{i=1}^{n} \max_{\tilde{y}_i} \underline{\lambda}_2 \cdot \phi(\underline{d}_i, \tilde{y}_i)$$

but $\max_{\tilde{y}_i} \alpha \underline{\lambda}_1 \cdot \phi(x_i, \tilde{y}_i) + \max_{\tilde{y}_i} (1-\alpha)\underline{\lambda}_2 \cdot \phi(x_i, \tilde{y}_i)$

$$\geq \max_{\tilde{y}_i} \left\langle (\alpha \underline{\lambda}_1 + (1-\alpha)\underline{\lambda}_2) \cdot \phi(\underline{d}_i; \tilde{y}_i) \right\rangle$$

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

(4)

# Apply CCCP algorithm

Two stages: step 1.

$$\frac{\partial}{\partial \lambda} g(\lambda)^t = -\phi(\underline{d_i}, y_i, \bar{h})$$

where $h^* = \arg\max_h \lambda^t \cdot \phi(\underline{d_i}, y_i, h)$

$\lambda^t$ current estimate of $\lambda$.

step 2.  solve

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

This reduces to a modified SVM with known state

$$\min_\lambda \frac{1}{2} |\lambda|^2 + C \sum_{i=1}^{n} \max_{(y, \underline{h})} \{ \lambda \cdot \phi(\underline{d_i}, y, \underline{h}) + \Delta(y_i, y, h) \}$$

$$- C \sum_{i=1}^{n} \lambda \cdot \phi(\underline{d_i}, y_i, \underline{h_i^*})$$

Note: similarities to EM.

→ step 1    involves estimating the hidden state $h_i^*$

→ step 2    estimate $\lambda$

repeat.

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