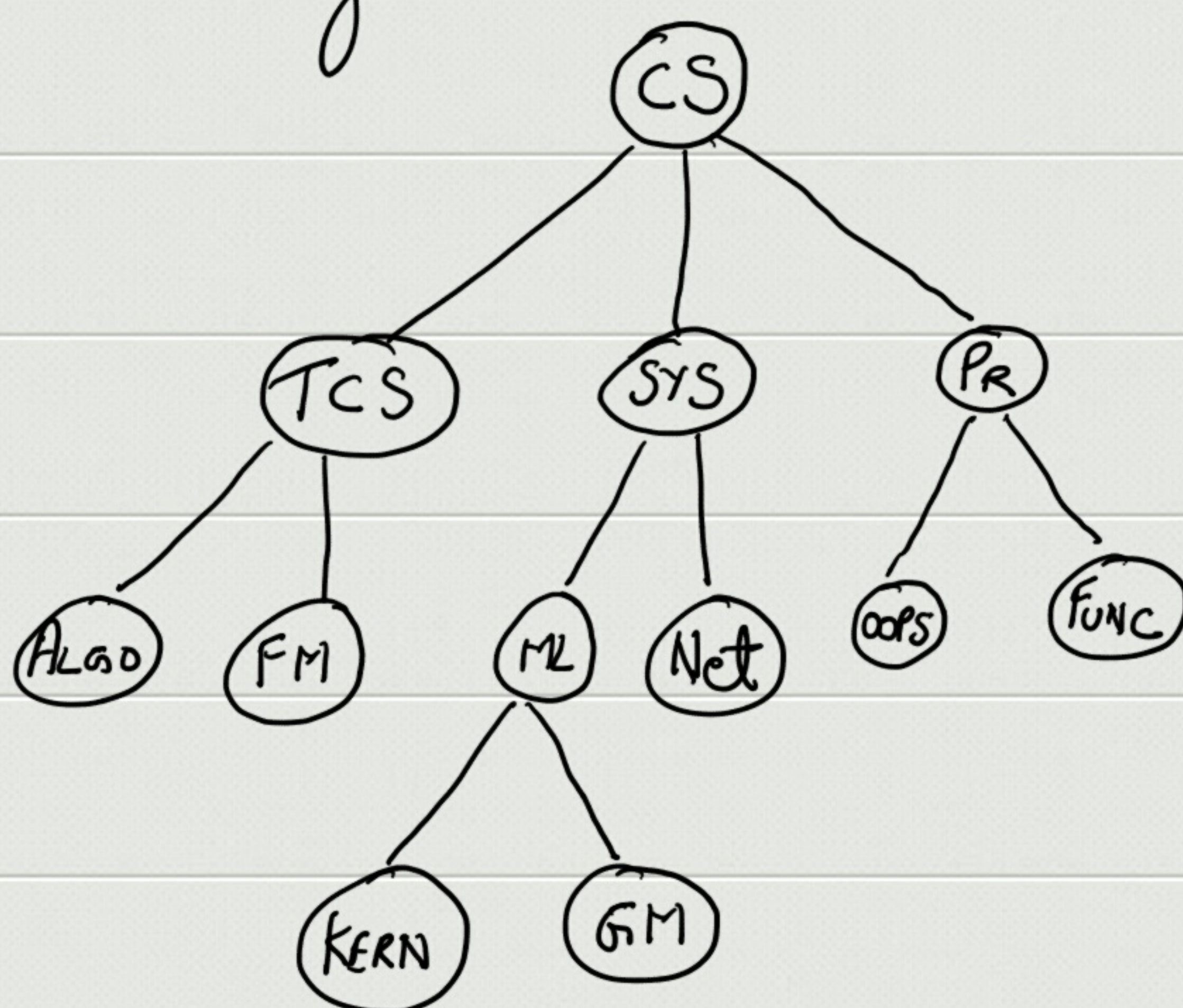


If classes have more 'complex' relations than total ordering, then?

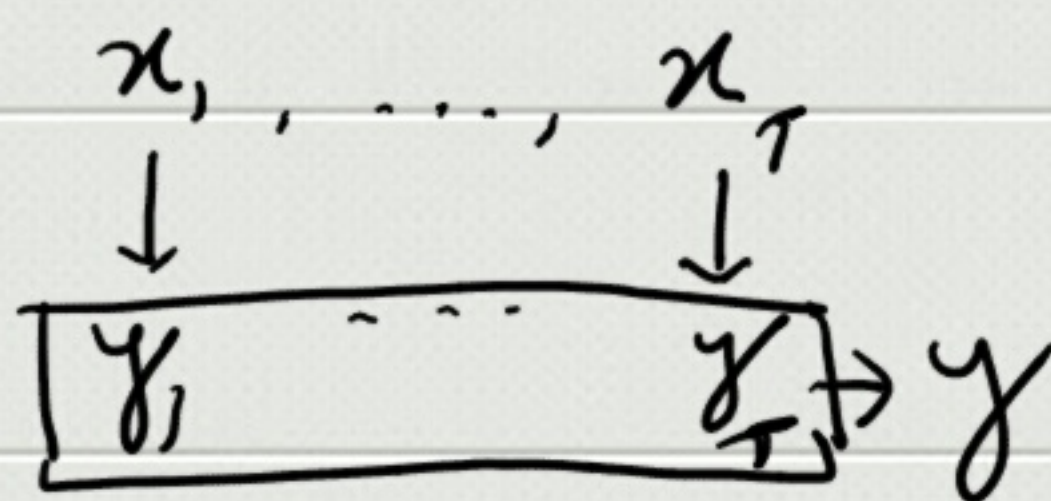
→ eg. classes have a partial order

eg. 1 Taxonomy



→ eg. classes have sub-classes

eg. 2 Sequence labeling



(Exponential no. of classes)

STRUCTURED PREDICTION

\mathcal{Y} has complex struc. other than \mathbb{R} .

→ Cumbersome to (linearly) parameterize $g: \mathcal{X} \rightarrow \mathcal{Y}$.

→ Alternative: pose as regression problem

find $g: \underline{\mathcal{X}} \rightarrow \mathbb{R}$, where $\underline{\mathcal{X}} = \mathcal{X} \times \mathcal{Y}$.

to learn compatibility of pairs (x, y) .

→ Now linear parametrization is known to us.

(i) → How to get training set?

(ii) → How to get kernel over $\mathcal{X} \times \mathcal{Y}$

→ k_1 is kernel over \mathcal{X} , k_2 is kernel over \mathcal{Y} ,

Any analytic function (no coeff.) is a kernel over $\mathcal{X} \times \mathcal{Y}$

(iii) → How to infer label given x using \tilde{g} .

→ find most compatible:

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} g(x, y)$$

Answer (i), consist on,

$$\omega^T \phi(x_i, y_i) - \max_{y \in Y(x_i)} \omega^T \phi(x_i, y) \geq 0$$

Again, it's situation of one free parameter.

→ Maximize margin, i.e. min. $\frac{1}{2} \|\omega\|^2$ while imposing:

$$\omega^T \phi(x_i, y_i) - \max_{y \in Y(x_i)} \omega^T \phi(x_i, y) \geq 1 \quad \forall i$$

→ Allowing for hinge-loss, we have:

$$\min_{\substack{\omega \in \mathbb{R}^n \\ \xi \in \mathbb{R}^m}} \frac{1}{2} \|\omega\|^2 + \frac{C}{m} \sum_{i=1}^m \xi_i$$

$$\text{s.t.} \quad \left. \begin{aligned} \omega^T \phi(x_i, y_i) - \max_{y \in Y(x_i)} \omega^T \phi(x_i, y) &\geq 1 - \xi_i, \\ \xi_i &\geq 0 \end{aligned} \right\} \forall i$$

$$l(y_i, g(x_i)) \equiv \max \left(0, 1 - \left(\omega^T \phi(x_i, y_i) - \max_{y \in Y(x_i)} \omega^T \phi(x_i, y) \right) \right)$$

Equivalently,

$$\min_{\omega \in \mathbb{R}^n, \xi \in \mathbb{R}^m} \frac{1}{2} \|\omega\|^2 + \frac{C}{m} \sum_{i=1}^m \xi_i,$$

$$\text{s.t.} \quad \omega^\top \psi_i(y) \geq 1 - \xi_i, \quad \forall y \in \mathcal{Y} \setminus \{y_i\}, \quad \xi_i \geq 0$$

where, $\psi_i(y) \equiv \varphi(x_i, y_i) - \varphi(x_i, y)$

Let $\alpha_{i,y}$ be the dual variables $[m(|\mathcal{Y}|-1)]$ in number

→ since $|\mathcal{Y}|$ itself may be large, coordinate descent \times

→ alternative is active set method.

Initialize (A)

→ repeat until convergence

→ for every example

→ Solve using (A)

→ Find KKT violator and append to (A)

Dual is:

$$\max_{\alpha \in \mathbb{R}^{m(m-1)}} \sum_{i=1}^m \sum_{y \neq y_i} \alpha_{iy} - \frac{1}{2} \sum_{i=1}^m \sum_{y \neq y_i} \sum_{s=1}^m \sum_{\bar{y} \neq y_s} \alpha_{iy} \alpha_{s\bar{y}} \psi_i(y)^T \psi_s(\bar{y})$$

s.t. $0 \leq \alpha_{iy} \forall y \neq y_i, \forall i; \sum_{y \neq y_i} \alpha_{iy} = \frac{C}{m} \forall i$

KKT conditions:

(i) $\omega^T \psi_i(y) \geq 1 - \xi_i, \alpha_{iy} \geq 0 \forall y \neq y_i, \forall i, \xi_i \geq 0 \forall i$

(ii) $\omega = \sum_{i=1}^m \sum_{y \neq y_i} \alpha_{iy} \psi_i(y), \sum_{y \neq y_i} \alpha_{iy} = \frac{C}{m} \forall i,$

$\alpha_{iy} = 0 \Rightarrow \xi_i \geq 1 - \omega^T \psi_i(y), \xi_i \begin{cases} \geq 0 & \text{if } \sum_{y \neq y_i} \alpha_{iy} = C/m \\ = 0 & \text{otherwise} \end{cases}$

$0 < \alpha_{iy} < \frac{C}{m} \Rightarrow \xi_i = 1 - \omega^T \psi_i(y), \xi_i \begin{cases} \geq 0 & \text{if } \sum_{y \neq y_i} \alpha_{iy} = C/m \\ = 0 & \text{otherwise} \end{cases}$

$\alpha_{iy} = \frac{C}{m} \Rightarrow \xi_i = 1 - \omega^T \psi_i(y) \geq 0$

\therefore No KKT violation if $\mathcal{E}_i \geq \max_{y \neq y_i} 1 - \omega^T \psi_i(y)$,

else the worst violator is given by:

$$\hat{y}_i = \operatorname{argmax}_{y \neq y_i} 1 - \omega^T \psi_i(y) = \operatorname{argmax}_{y \neq y_i} \omega^T \phi(x_i, y)$$

\hat{y}_i can be appended to active set (A).

\rightarrow Training, Inference tractable if \hat{y}_i is tractable.