A simple formulation where the ~~features~~ all tasks are assumed to ~~low dimensional~~ share the same (low number of) features:

$$\min_{\omega,b,\xi} \quad \frac{1}{2}\left(\sum_{j=1}^{d} \|a_j\|_2\right)^2 + C\sum_{t,i}\xi_{ti}$$

(I) $\quad \underline{s.t.} \quad y_{ti}\left(v_t^T x_{ti}\right) \geq 1-\xi_{ti}, \quad \xi_{ti} \geq 0$

$$
\begin{array}{c}
v_1 \; v_2 \; \cdots \; v_T \\
\downarrow \; \downarrow \quad\quad \downarrow \\
1 \; 2 \; 3 \cdots T
\end{array}
$$

$$
\begin{array}{l}
a_1 \longrightarrow 1 \;\; \omega_{11} \; \omega_{12} \cdots \omega_{1T} \\
a_2 \longrightarrow 2 \;\; \omega_{21} \quad \cdots \quad \omega_{2T} \\
\quad \vdots \\
a_d \longrightarrow d \;\; \omega_{d1} \cdots \omega_{dT}
\end{array}
$$

→ Promotes structured sparsity: if a feature is useful it is employed in all tasks else it is <u>not</u> employed in <u>any</u> task.

Interestingly this formulation can be written as an MKL formulation:

Using "lambda trick" (I) is same as:

$$\min_{\lambda \in \Delta} \; \min_{\omega,b,\xi} \; \frac{1}{2}\sum_{j=1}^{d} \frac{\sum_{t=1}^{T} \omega_{jt}^2}{\lambda_j} + C\sum_{t,i}\xi_{ti} \qquad \text{(II)}$$

$$\underline{s.t.} \quad y_{ti}\left(\sum_{j=1}^{d} \omega_{jt}\,\boxed{x_{tij}}\right) \geq 1-\xi_{ti}, \quad \xi_{ti} \geq 0$$

→ $j^{th}$ feature of $i^{th}$ eg. of $t^{th}$ task.

Let $\bar{\omega}_{jt} = \frac{\omega_{jt}}{\sqrt{\lambda_j}}$, Let $M_t = [0\cdots 0\, I\, 0\cdots 0]_{d\times d}$ ; Let $\Lambda = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_d \end{bmatrix}$

$\quad\quad\quad\quad\quad\quad\quad\quad$ ↓ $d\times d$ zero matrix $\;$ → $d\times d$ identity matrix $\;$ $diag(\lambda_1,\dots,\lambda_d)$

~~Note that~~ Let $\omega = \begin{bmatrix} \omega_{11} \\ \omega_{21} \\ \omega_{d2} \\ \vdots \\ \omega_{d2} \\ \vdots \\ \omega_{1T} \\ \omega_{dT} \end{bmatrix}_{Td\times 1}$ $\quad$ Now $\underline{v_t = M_t\,\omega}$

With this notation it is easy to see the ~~formula~~ (II) is:

$$\min_{\lambda \in \Delta} \quad \min_{\omega, b, \xi} \quad \frac{1}{2}\|\omega\|_2^2 + C\sum_{t,i}\xi_{ti}$$

$$\underline{s.t.} \quad y_{ti}\left(\omega^T M_t^T \Lambda x_{ti} + b\right) \geq 1 - \xi_{ti}, \quad \xi_{ti} \geq 0$$

standard dual (observe there is no 'b' term however)

$$\min_{\lambda \in \Delta} \quad \max_{\alpha} \quad 1^T \alpha - \frac{1}{2}\alpha^T Y K_\lambda Y \alpha$$

$$\underline{s.t.} \quad 0 \leq \alpha \leq C,$$

no term $y^T\alpha = 0$    (III)

$K_\lambda$ is a $Td \times Td$ matrix which is block diagonal. All entries across tasks are zero! Within a task $t$ it is a simple conic combination of ~~kernels~~ bare kernels as linear kernel built ~~are~~ on each feature. Mathematically:

$$K_\lambda(x,y) = \begin{cases} 0 & \text{if } x \text{ \& } y \text{ are not from same task} \\ \lambda_1 K_1 + \dots + \lambda_d K_d & \text{if } x \text{ \& } y \text{ belong to task } a \end{cases}$$

$$K_1(x,y) = x_1 y_1, \dots K_d(x,y) = x_d y_d.$$

simple per feature linear kernels.

*[With this $\alpha^T Y K_\lambda Y \alpha$ decomposes into 'T' quadratics one per each task: $\sum_{t=1}^{T} \alpha_t^T Y_t \left(\sum_{i=1}^{d} \lambda_i K_i^t\right) Y_t \alpha_t$ ]*

Now it is easy to see that (III) can be solved using MD or simple MKL etc.