

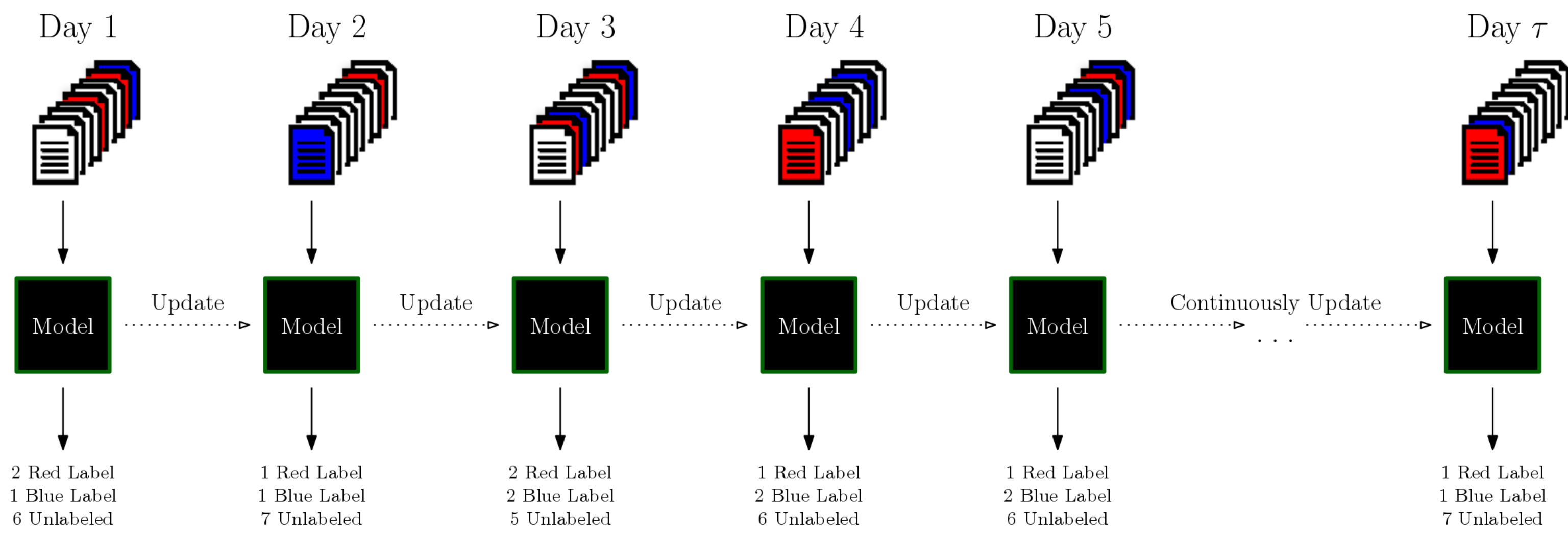


Sequential Maximum Margin Classifiers for Partially Labeled Data

Elizabeth Hou[†], Alfred O. Hero III[†]

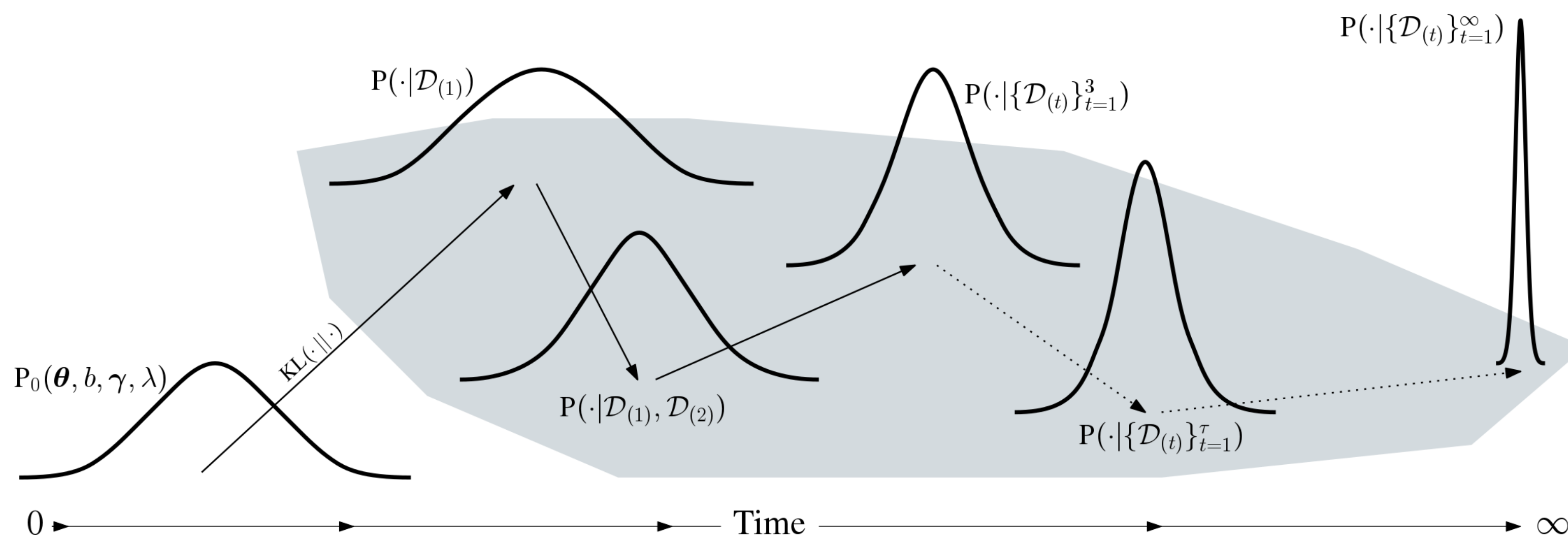
[†] University of Michigan, EECS Department

Motivation



- Data is not collected as a batch, but sequentially over time
 - Speech or streaming text classification
 - Satellite only transmits data daily
 - Agency's quarterly reports
- Often not possible/desirable to wait until complete before analyzing
- Additionally, some of the labels may be missing
- At any time point t , observed data as $\mathcal{D}(t) = \{\mathbf{X}(t), \mathbf{y}(t)\}$
 - $\mathbf{X}(t)$ is a matrix of $n(t)$ samples and p feature variables
 - $l(t) < n(t)$ of the samples have label $y_i = [1, -1]$

Sequential Laplacian MED



- Parameters:
- θ are weights for decision boundary
 - b is a bias term for decision boundary
 - γ_i are margin parameters
 - λ is a regularizer

Objective:

$$\min_{\theta, b, \gamma, \lambda} \text{KL}(\mathcal{P}(\cdot|\{\mathcal{D}(t)\}_{t=1}^{\tau}) || \mathcal{P}_0(\cdot|\{\mathcal{D}(t)\}_{t=1}^{\tau-1}))$$

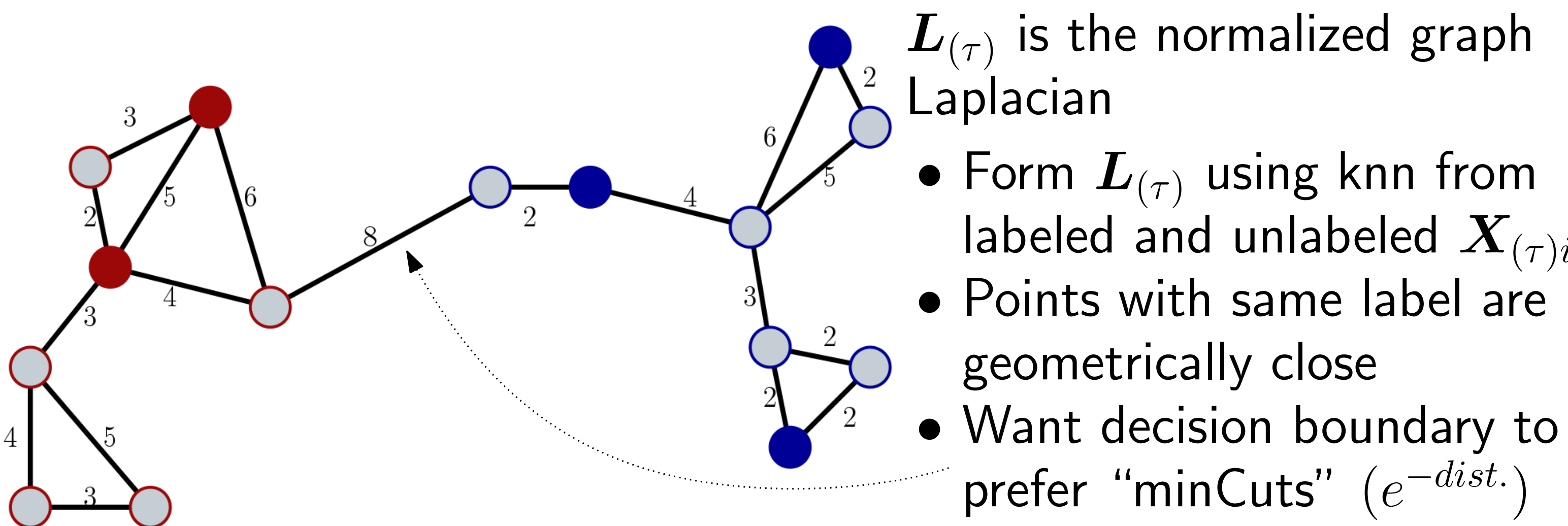
subject to

$$E_{\theta, b, \gamma} (y_{(\tau)i} (\mathbf{X}_{(\tau)i} \theta + b) - \gamma_i | \{\mathcal{D}(t)\}_{t=1}^{\tau}) \geq 0 \quad \forall i$$

Soft Margin Hinge Loss

$$E_{\theta, \lambda} (\theta^T \mathbf{X}_{(\tau)}^T \mathbf{L}_{(\tau)} \mathbf{X}_{(\tau)} \theta - \lambda | \{\mathcal{D}(t)\}_{t=1}^{\tau}) \leq 0$$

Regularize Smoothness w.r.t. distribution of x



$$\mathcal{P}(\theta, b, \gamma, \lambda | \{\mathcal{D}(t)\}_{t=1}^{\tau}) = \frac{\mathcal{P}_0(\cdot | \{\mathcal{D}(t)\}_{t=1}^{\tau-1})}{Z(\alpha_{(\tau)}, \beta_{(\tau)})} \exp \left\{ \sum_{i=1}^{l_{(\tau)}} \alpha_{(\tau)i} H_i + \beta_{(\tau)} S \right\}$$

If prior is exp. family, the minimizing density is of same family [2, 3]

Update Theorem [4]

At time τ , let the MED priors be $\theta \sim N(\mu_{(\tau-1)}, (\mathbf{G}_{(\tau-1)})^{-1})$, $b \sim N(0, \sigma^2)$ and $\lambda \sim \text{Exp}(\nu)$ where $\sigma, \nu \rightarrow \infty$, and $\gamma_i \sim \mathcal{P}_0(C_{(\tau)})$.

Then the posterior also factorizes with $\mathcal{P}(\theta | \{\mathcal{D}(t)\}_{t=1}^{\tau})$ as Gaussian:

mean $\mu_{(\tau)} = \mathbf{G}_{(\tau)}^{-1} (\mathbf{G}_{(\tau-1)} \mu_{(\tau-1)} + \mathbf{X}_{(\tau)}^T \mathbf{J}^T \mathbf{Y}_{(\tau)} \hat{\alpha}_{(\tau)})$ and precision matrix $\mathbf{G}_{(\tau)} = \mathbf{G}_{(\tau-1)} + 2\beta_{(\tau)} \mathbf{X}_{(\tau)}^T \mathbf{L}_{(\tau)} \mathbf{X}_{(\tau)}$.

Decision Rule Corollary [4]

Assume $\beta_{(t)}$ are fixed parameters. Then the decision rule reduces to $\hat{y}_i = \text{sgn}(\mathbf{X}_i \mu_{(\tau)} + \hat{b})$, which is a function of the previous mean and optimal parameters $\hat{\alpha}_{(\tau)} = \arg \max_{\alpha_{(\tau)}} Z(\alpha_{(\tau)}, \beta_{(\tau)}) =$

$$-\frac{1}{2} \alpha_{(\tau)}^T \mathbf{Y}_{(\tau)} \mathbf{J} \mathbf{X}_{(\tau)} \mathbf{G}_{(\tau)}^{-1} \mathbf{X}_{(\tau)}^T \mathbf{J}^T \mathbf{Y}_{(\tau)} \alpha_{(\tau)} + \alpha_{(\tau)}^T (\mathbf{1} - \mathbf{Y}_{(\tau)} \mathbf{J} \mathbf{X}_{(\tau)} \mu_{\tau-1}) + \sum_{i=1}^{l_{(\tau)}} \log \left(1 - \frac{\alpha_{(\tau)i}}{C_{(\tau)}} \right) \quad \text{s.t.} \quad \mathbf{y}_{(\tau)}^T \alpha_{(\tau)} = 0, \alpha_{(\tau)i} \geq 0 \quad \forall i = 1, \dots, l_{(\tau)}$$

and the \hat{b} that satisfy the KKT conditions.

The above can be extended to non-linear decision boundaries with the kernel trick [4].

ACKNOWLEDGMENTS This work was funded partially by the Consortium for Verification Technology under Department of Energy National Nuclear Security Administration award number DE-NA0002534 and partially by the University of Michigan ECE Departmental Fellowship.

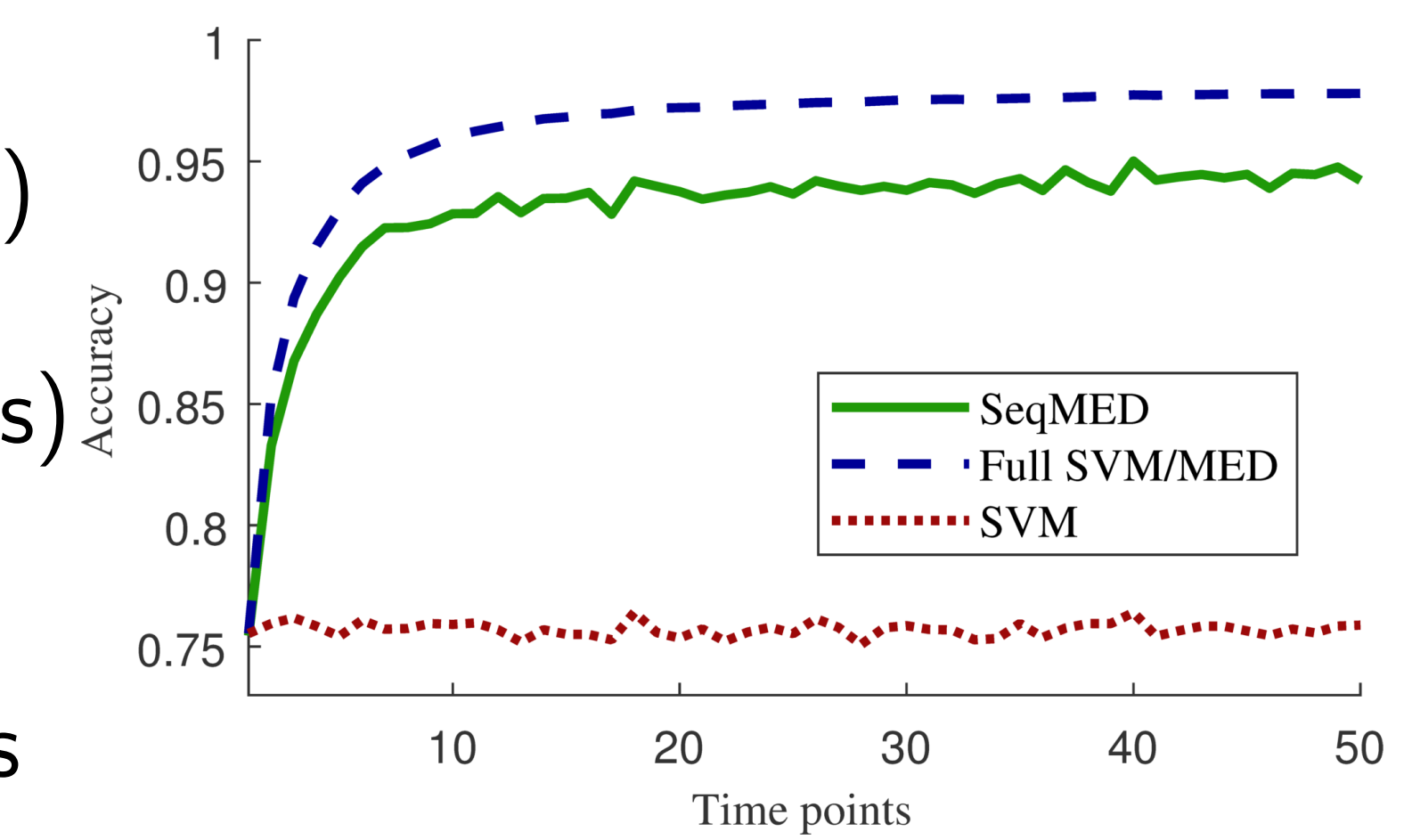
Experiments

Simulation

- Training: ~ 100 per time point ($n(t) = [97, 103]$)
- Test: 1000 test points, Accuracy = $(TP+TN)/1000$

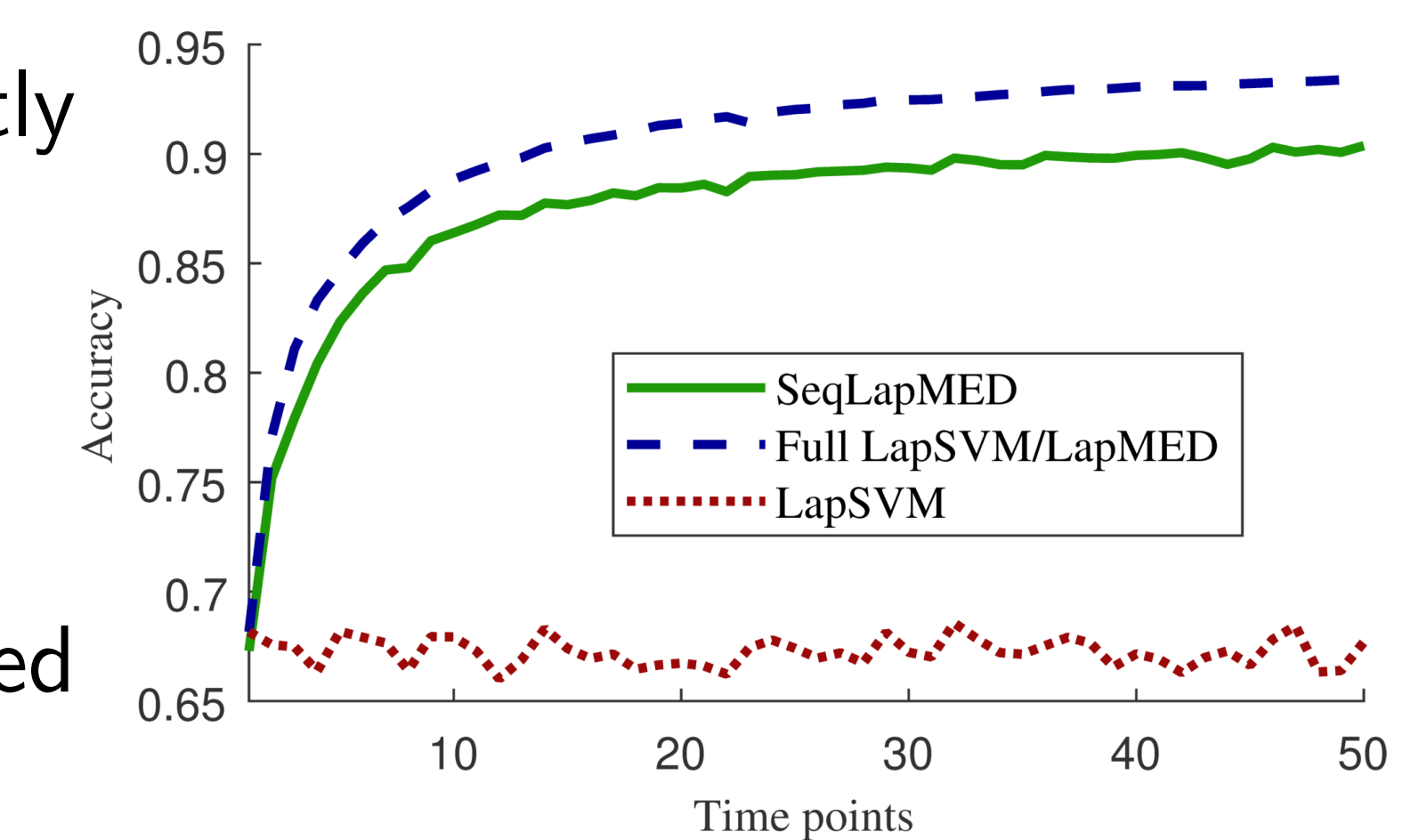
Supervised: Data from 200 categorical distributions

- 100 are sparse (high prob. of 0s)
- 50 are relevant (lower prob. of 0s)
- 50 are used to distinguish between 2 classes



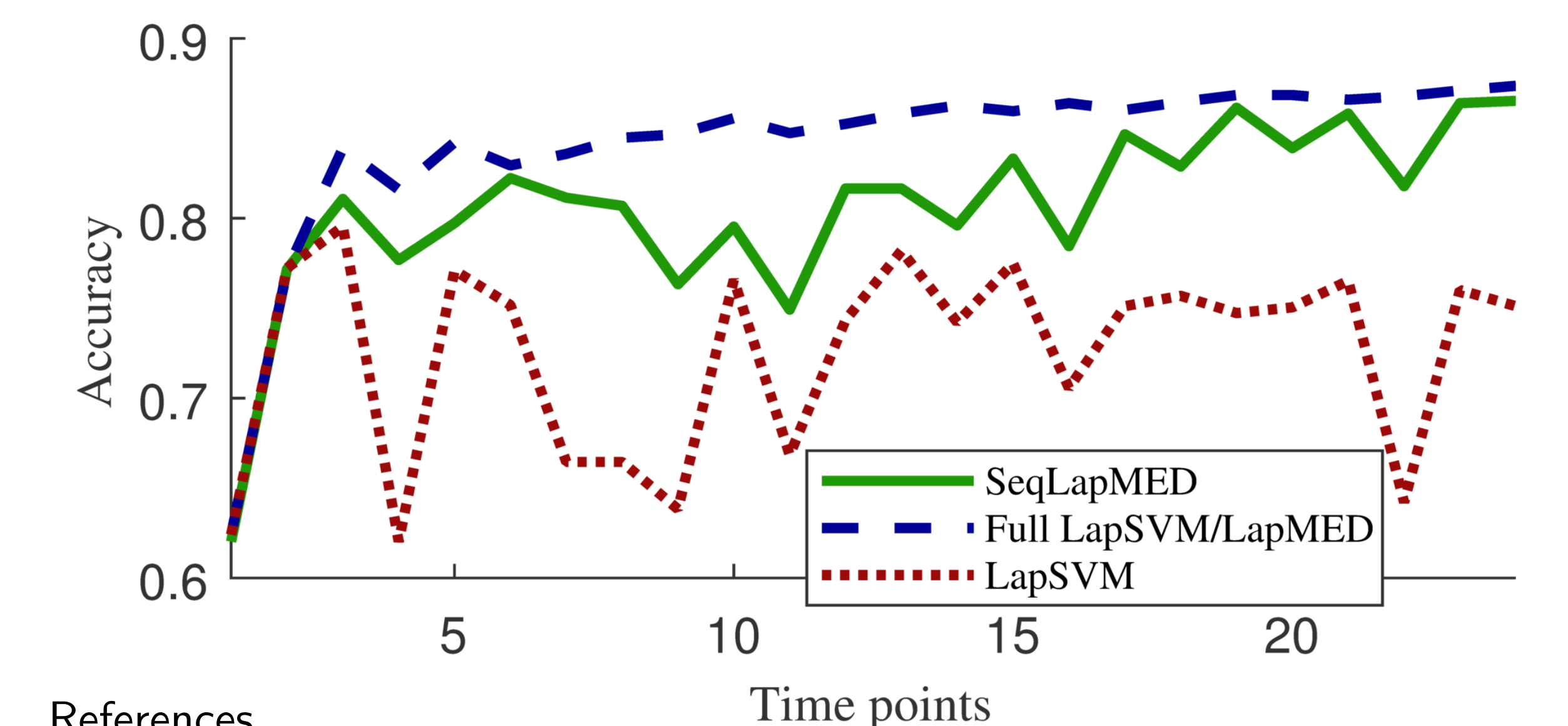
Semi-supervised: Data from interior of a 3-D sphere

- One class is mostly in center/interior
- Other class is on the shell
- Only 10% of the samples are labeled



Isolet Speech Database

- Follows the experimental framework used in [1]
- Training: (isolet1 - isolet4) broken into 24 time points of 5 speakers, only first speaker is labeled
- Test: 1,559 samples from isolet5



References

- M. Belkin, P. Niyogi, and V. Sindhwani. "Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. JMLR (2006).
- O. Koyejo, J. Ghosh, "A representation approach for relative entropy minimization with expectation constraints. In 2013 ICML WDDL workshop.
- E. Hou, K. Sricharan, A. O. Hero. "Latent Laplacian Maximum Entropy Discrimination for Detection of High-Utility Anomalies". IEEE TIFS (2018).
- E. Hou, A. O. Hero. "Sequential Maximum Margin Classifiers for Partially Labeled Data". In 2018 IEEE ICASSP.