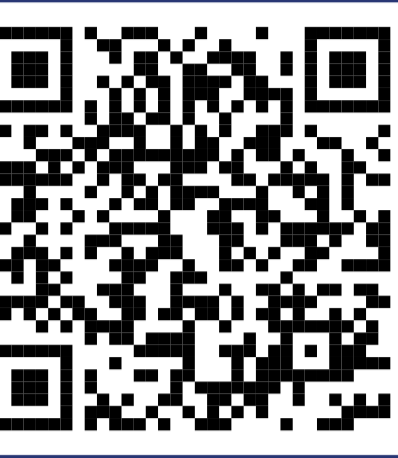




Fast Training of Support Vector Machines for Survival Analysis

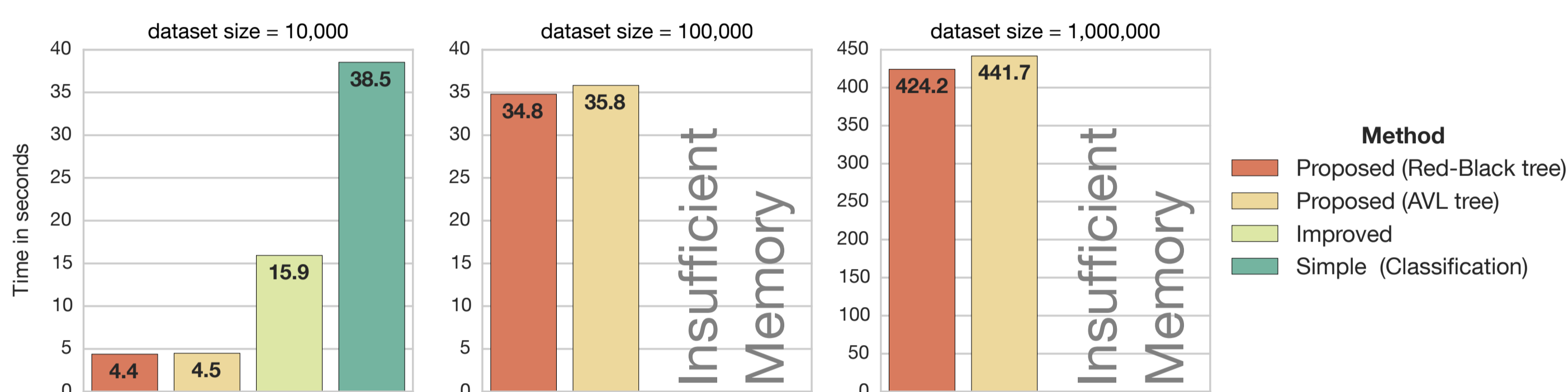
Sebastian Pölsterl¹, Nassir Navab^{1,2}, and Amin Katouzian¹

1: Computer Aided Medical Procedures, Technische Universität München, Munich, Germany
2: Johns Hopkins University, Baltimore MD, USA



Overview

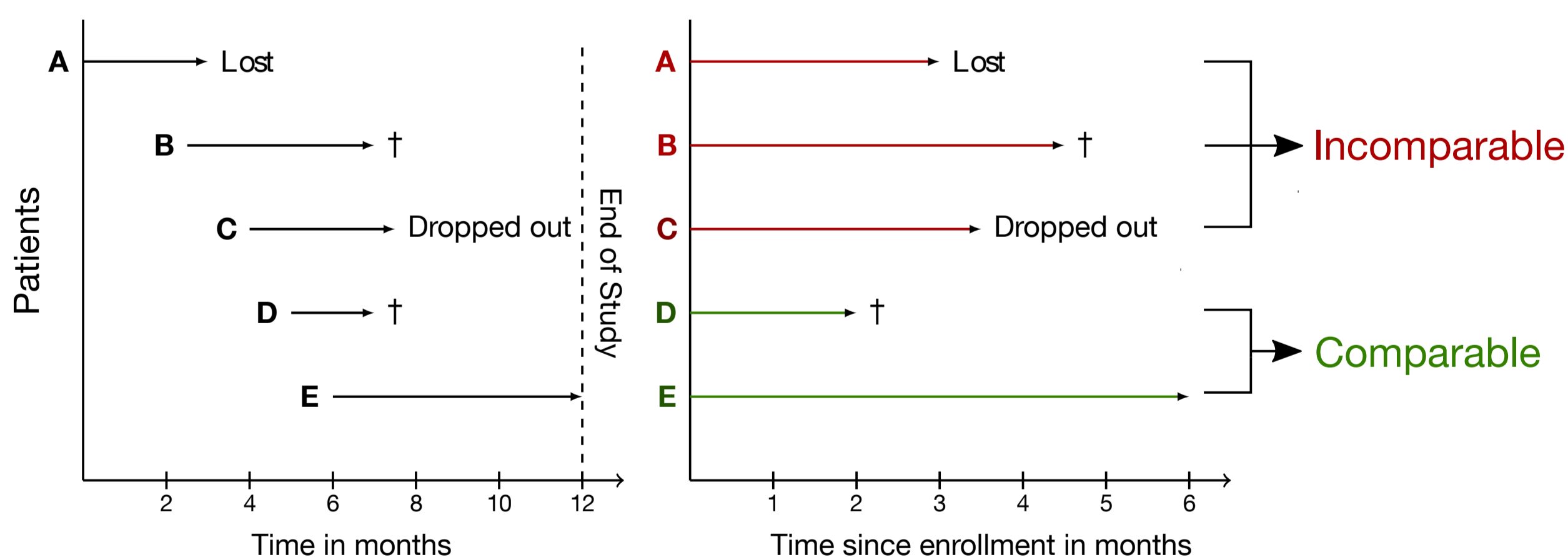
- Problem:** Training a Survival SVM by converting the ranking problem into a classification problem and using a standard dual SVM solver requires $O(dn^4)$ time and $O(n^2)$ space.¹
- Current RankSVM solvers² require $O(nd + n \log n + n \log r + d)$.
- Solution:**
 - Optimization in the primal using truncated Newton optimization.
 - Use order statistics trees to lower complexities.
- Conclusion:** Requires $O(n)$ space and has time complexity $[O(n \log n) + O(nd + d + n \log n)] \cdot \bar{N}_{CG} \cdot N_{Newton}$.



- Same optimization scheme can be applied to:
 - Non-linear Survival SVM: Transform data with Kernel PCA before training in primal.³
 - Hybrid ranking-regression to predict exact time of event.

Survival Analysis

- Objective:** to establish a connection between a set of features and the time of an event (*survival time*).
- Patients that remain event-free during the study-period are **right censored**, because it is unknown whether an event has or has not occurred after the study ended. Only **partial information** about their survival is available.



Survival Support Vector Machine

- Survival analysis can be formulated as a ranking problem, with the objective to predict the order of patients according to their survival time.
- Survival SVM¹ is closely related to RankSVM⁴, but in addition accounts for right censoring: subject with *smaller survival time* must always be *uncensored*.
- Number of relevance levels (here: *survival times*) are of the order of number of patients.
- Survival data consists of a feature vector ($x_i \in \mathbb{R}^d$), the time of an event/censoring ($y_i > 0$), and an event indicator ($\delta_i \in \{0, 1\}$).

Truncated Newton Optimization

- Objective function:**

$$\mathcal{P} = \{(i, j) \mid y_i > y_j \wedge \delta_j = 1\}_{i,j=1,\dots,n}$$

$$f(w) = \frac{1}{2} w^T w + \frac{\gamma}{2} \sum_{(i,j) \in \mathcal{P}} \max(0, 1 - (w^T x_i - w^T x_j))^2$$

$$= \frac{1}{2} w^T w + \frac{\gamma}{2} (p_w + w^T X^T (A_w^T A_w X w - 2A_w^T))$$

$$(A_w)_{k,q} = \begin{cases} 1 & \text{if } y_q < y_s \wedge \delta_q = 1 \wedge w^T x_s < w^T x_q + 1, \\ -1 & \text{if } y_q > y_s \wedge \delta_s = 1 \wedge w^T x_s > w^T x_q - 1, \\ 0 & \text{else,} \end{cases} \quad (1)$$

$$(A_w)_{k,q} = \begin{cases} -1 & \text{if } y_q > y_s \wedge \delta_s = 1 \wedge w^T x_s > w^T x_q - 1, \\ 0 & \text{else,} \end{cases} \quad (2)$$

where $s, q \in \{1, \dots, n\}, k \in \{1, \dots, p_w\}$

- Example:** $\mathcal{P} = \{(1, 2); (1, 3); (2, 3); (4, 2); (4, 3)\}$

i	1	2	3	4
y_i	9	6	5	8
δ_i	0	1	1	0
$w^T x_i$	-0.7	-0.1	0.15	1.6

$$A_w = \begin{pmatrix} -1 & 1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & -1 & 1 & 0 \end{pmatrix}$$

- Do not store Hessian, but only compute Hessian-vector product $Hv = v + \gamma X^T A_w^T A_w X v$.

Using Order Statistics Trees

- Matrix $A_w^T A_w$ has size $O(n^2)$ and has to be recomputed in each iteration, which would be impractical. Instead we compactly express one entry as

$$(A_w^T A_w)_{i,j} = \begin{cases} l_i^+ + l_i^- & \text{if } i = j, \\ -1 & \text{if } i \neq j, \text{ and } j \in SV_i^+ \text{ or } i \in SV_j^-, \\ 0 & \text{else.} \end{cases}$$

$$SV_i^+ = \{s \mid y_s > y_i \wedge \delta_i = 1 \wedge w^T x_s < w^T x_i + 1\} \quad l_i^+ = |SV_i^+|$$

$$SV_i^- = \{s \mid y_s < y_i \wedge \delta_s = 1 \wedge w^T x_s > w^T x_i - 1\} \quad l_i^- = |SV_i^-|$$

- Sets can be updated incrementally by adding y_i and $w^T x_i$ to an order statistics tree (balanced binary search tree).

- Example:**

i	1	2	3	4	5	6	7	8	9
$w^T x_i$	-0.7	-0.1	0.15	0.2	0.3	0.8	1.6	1.7	2.3
y_i	1	9	6	5	8	2	7	3	4
δ_i	0	0	1	0	1	1	1	0	0

