



Gautam Kunapuli
University of Wisconsin-Madison
kunapuli@wisc.edu

Richard Maclin
University of Minnesota-Duluth
rmaclin@d.umn.edu

Jude W. Shavlik
University of Wisconsin-Madison
shavlik@cs.wisc.edu



This research was supported by the Defense Advanced Research Projects Agency under grant FA8650-06-C-7606 and the National Institute of Health under NLM grant R01-LM008796. Views and conclusions contained here are those of the authors and do not necessarily represent the official opinion or policies, either expressed or implied of the US government or of DARPA.

Abstract. KBSVMs incorporate advice from domain experts to improve generalization. Imperfect advice can lead to significantly poorer models. To learn in this setting, we propose an approach that extends KBSVMs and is able to not only learn from data and advice, but also simultaneously improve the advice. The proposed approach is particularly effective for knowledge discovery in domains with few labeled examples. The additional refinement constraints are bilinear and are solved using two iterative approaches: successive linear programming and a constrained concave-convex approach.

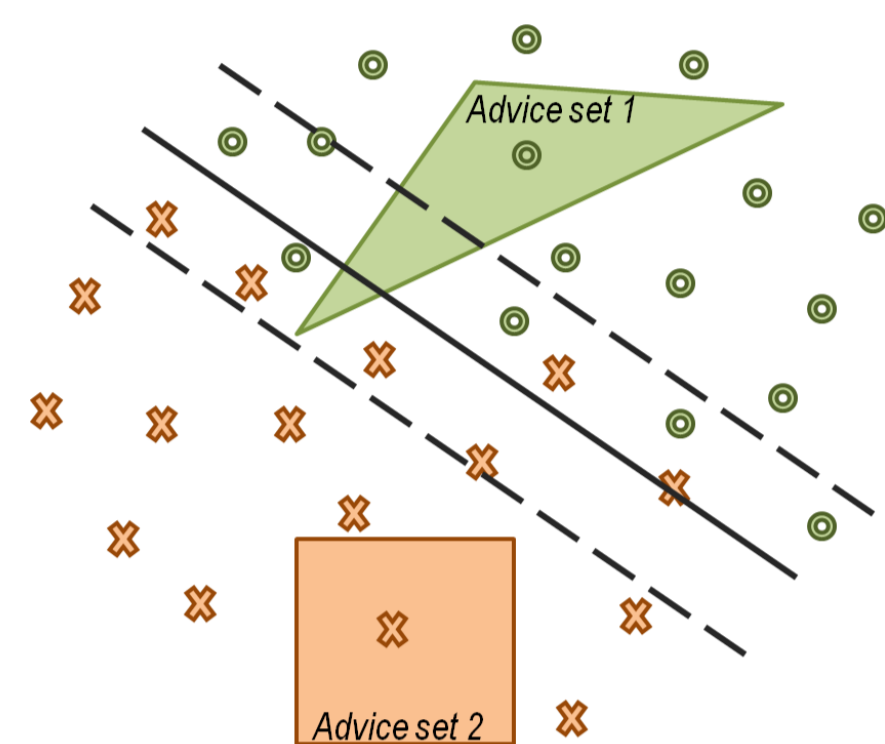
Adding Knowledge to SVMs

Polyhedral knowledge sets in input space can be added to SVMs via *knowledge-based support vector machines* (KBSVMs) [FMS03]. Knowledge sets characterize an area of input space as belonging to one of the two classes and are specified using

$$Dx \leq d \Rightarrow z(w'x - \gamma) \geq 1.$$

The advice specifies that every point $x \in Dx \leq d$ lies above $w'x - \gamma = 1$ (if labeled $z = 1$) or below $w'x - \gamma = -1$ (if labeled $z = -1$).

Example: diabetes diagnosis from two features, blood glucose level and body mass index:



Expert (NIH) Advice for Type-2 Diabetes Diagnosis:

- an obese person ($bmi \geq 30$) with $gluc \geq 126$ is at strong risk
- a person at normal weight ($bmi \leq 25$) with $gluc \leq 100$ is at low risk.

Can be expressed as polyhedral advice in input space:

$$\begin{aligned} (bmi \leq 25) \wedge (gluc \leq 100) &\Rightarrow \text{-diabetes} \\ (bmi \geq 30) \wedge (gluc \geq 126) &\Rightarrow \text{diabetes.} \end{aligned}$$

In general, $Dx \leq d$ allows linear combinations of input space features to specify advice.

With *theorems of the alternative*, this implication can be reformulated as constraints with an **advice vector**, u :

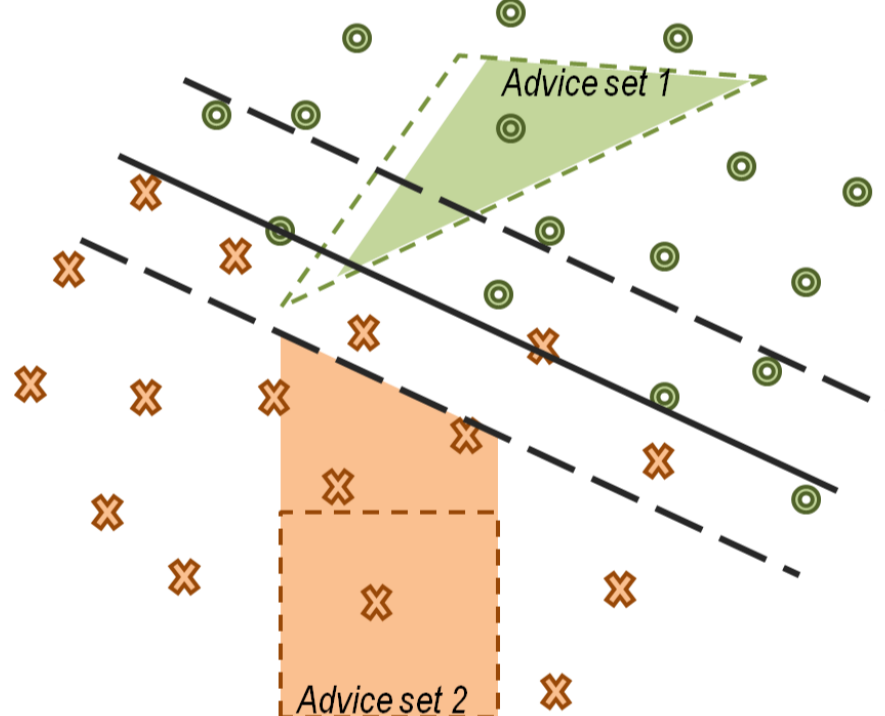
$$D'u + z w = 0, \quad -d'u - z \gamma \geq 1, \quad u \geq 0.$$

The advice vector is analogous to the dual multipliers α in SVMs: constraints of an advice set which have non-zero u 's are called **support constraints** and only they contribute to w . The following is the formulation of the **knowledge-based support vector machine**:

$$\begin{aligned} \min_{w, b, (\xi, u^i, \eta^i, \zeta_i) \geq 0} \quad & \|w\|_1 + \lambda e' \xi + \mu \sum_{i=1}^m (e' \eta^i + \zeta_i) \\ \text{s.t.} \quad & Y(Xw - be) + \xi \geq e, \\ & -\eta^i \leq D_i' u^i + z_i w \leq \eta^i, \\ & -d^i' u^i - z_i b + \zeta_i \geq 1, \quad i = 1, \dots, m. \end{aligned} \quad (1)$$

Learning from Imperfect Advice

Motivation. To learn from small data sets and refine imperfect advice provided by a domain expert.



There are two types of imperfect advice

- **Advice set 1.** Advice that lies on the wrong side of the margin is penalized as advice error through $\sum_i (e' \eta^i + \zeta_i)$ in (1). **Refinement leads to truncation of advice.**

- **Advice set 2.** Advice that lies on the correct side of the margin but far away provides no support vectors, i.e., has all its advice vectors $u^i = 0$ in (1). **Refinement leads to extension of advice.**

Formulation. Introduce terms (F_i, f^i) into the formulation that directly modify boundaries of advice regions:

$$(D_i - F_i)x \leq (d^i - f^i) \Rightarrow z_i(w'x - b) \geq 1, \quad z = \pm 1.$$

Extends **rule-refining support vector machines** (RRSVMs) [MWS⁺07] which performed refinements:

$$D_i x \leq (d^i - f^i) \Rightarrow z_i(w'x - b) \geq 1, \quad z = \pm 1.$$

Advice-Refining Knowledge-Based SVMs (arkSVMs)

Advice refinement is formulated to optimize *model complexity* + λ *training error* + μ *advice error* + ν *refinement error*:

$$\begin{aligned} \min_{w, b, F_i, f^i, (\xi, u^i, \eta^i, \zeta_i) \geq 0} \quad & \|w\|_1 + \lambda e' \xi + \mu \sum_{i=1}^m (e' \eta^i + \zeta_i) + \nu \sum_{i=1}^m (\|F_i\|_1 + \|f^i\|_1) \\ \text{s.t.} \quad & Y(Xw - be) + \xi \geq e, \\ & -\eta^i \leq (D_i - F_i)' u^i + z_i w \leq \eta^i, \\ & -(d^i - f^i)' u^i - z_i b + \zeta_i \geq 1, \quad i = 1, \dots, m. \end{aligned} \quad (2)$$

This is the **Advice-Refining KBSVM** (arkSVMs).

- Objective trades-off the effect of refinement in each of the advice sets via the *refinement parameter* ν ,
- Refining d alone allows only for the *translation* of the boundaries of the polyhedral advice; in arkSVMs the boundaries of the advice can be translated *and rotated*,
- Advice constraints are bilinear. Solved using successive linear programming (also used in [MWS⁺07]), and a concave-convex procedure.

arkSVMs via Successive Linear Programming

Solve a sequence of linear programs while alternately fixing the bilinear variables $(F_i, f^i)_{i=1}^m$ and $\{u^i\}_{i=1}^m$. At the t -th iteration, the algorithm alternates between the following steps:

- **(Estimation Step)** When refinement terms, $(\hat{F}_i^t, \hat{f}^{i,t})_{i=1}^m$, are fixed, resulting LP becomes a standard KB-SVM which attempts to find a data-estimate of the advice vectors $\{u^i\}_{i=1}^m$ using the current refinement of the advice region: $(D_j - \hat{F}_j^t)x \leq (d^j - \hat{f}^{j,t})$.
- **(Refinement Step)** When advice-estimate terms $\{u^i\}_{i=1}^m$ are fixed, the resulting LP solves for $(F_i, f^i)_{i=1}^m$ and attempts to further refine the advice regions based on estimates from data from the previous step.

Algorithm converges to a local solution.

arkSVMs via Successive Quadratic Programming

A general bilinear term $r's$, which is non-convex, can be written as the difference of two convex terms: $\frac{1}{4}\|r+s\|^2 - \frac{1}{4}\|r-s\|^2$. The j -th component of the **bilinear advice constraint (with terms $F_j^i u^i$) in (2) is rewritten:**

$$D_{ij}' u^i + z_j w_j - \eta_j^i + \frac{1}{4} \|F_{ij} - u^i\|^2 \leq \frac{1}{4} \|F_{ij} + u^i\|^2,$$

and both sides of the constraint above are convex and quadratic. We **linearize the right-hand side around an estimate of the bilinear variables** $(\hat{F}_{ij}^t, \hat{u}^{i,t})$, for $(D_i - F_i)' u^i + z_i w + \eta^i \leq 0$ and $-(D_i - F_i)' u^i - z_i w - \eta^i \leq 0$

$$D_{ij}' u^i + z_j w_j - \eta_j^i + \frac{1}{4} \|F_{ij} - u^i\|^2 \leq \frac{1}{4} \|\hat{F}_{ij}^t + \hat{u}^{i,t}\|^2 + \frac{1}{2} (\hat{F}_{ij}^t + \hat{u}^{i,t})' ((F_{ij} - \hat{F}_{ij}^t) + (u^i - \hat{u}^{i,t})),$$

while $d^i' u^i + z_i b + 1 - \zeta_i - f^i' u^i \leq 0$ is replaced by

$$-d^i' u^i + z_i b + 1 - \zeta_i + \frac{1}{4} \|f^i - u^i\|^2 \leq \frac{1}{4} \|\hat{f}^{i,t} + \hat{u}^{i,t}\|^2 + \frac{1}{2} (\hat{f}^{i,t} + \hat{u}^{i,t})' ((f^i - \hat{f}^{i,t}) + (u^i - \hat{u}^{i,t})).$$

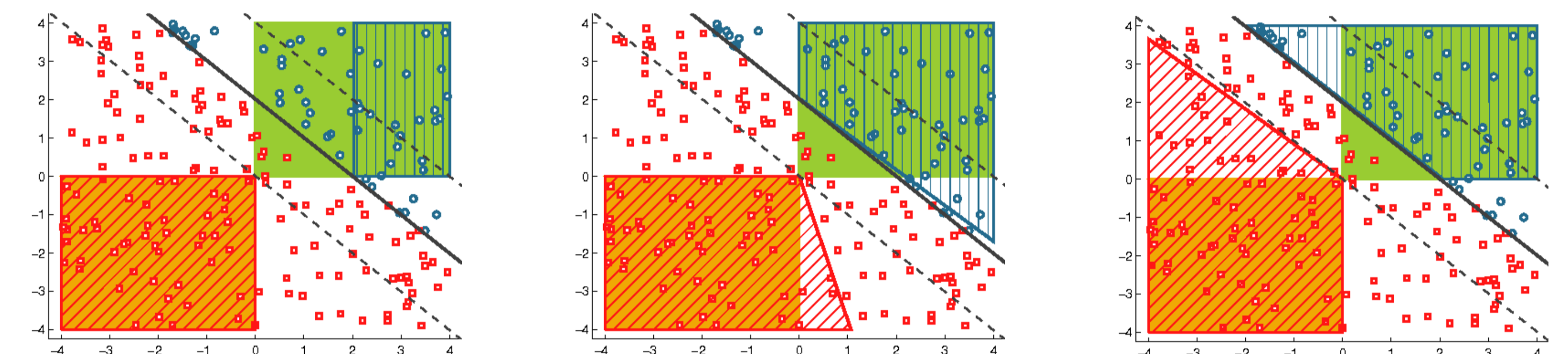
- Convexified relaxations result in a quadratically-constrained linear program (QCLP); rhs are affine and thus, entire set of constraints convex,
- Feasible set of this problem is a subset of the original.

Iteratively solve the QCLP: At the t -th iteration, the current estimate is used to obtain a new feasible point; iterating this procedure produces a sequence of feasible points with decreasing objective values.

Algorithm converges to a local solution.

Synthetic Example

Simple 2-dimensional example consists of 200 points separated by $x_1 + x_2 = 2$. Two advice sets: $\{S_1 : (x_1, x_2) \geq 0 \Rightarrow z = +1\}$, $\{S_2 : (x_1, x_2) \leq 0 \Rightarrow z = -1\}$.



(left) RRSVM (center) arkSVM-slp (right) arkSVM-sqp. Orange and green unhatched regions show the original advice. The dashed lines show the margin, $\|w\|_\infty$. For each method, we show the refined advice: vertically hatched for Class +1, and diagonally hatched for Class -1.

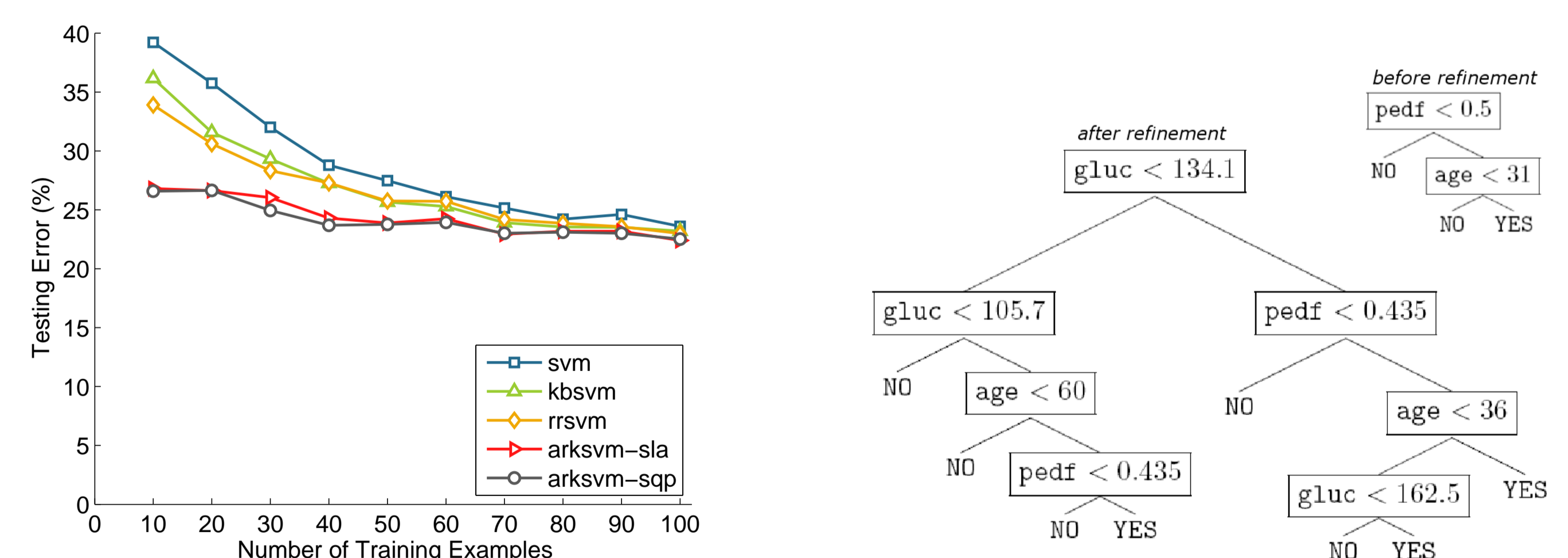
PIMA Indians Diabetes Diagnosis

Predict onset of diabetes in 768 Pima Indian women *within the next 5 years* based on eight features. Studies [HFC⁺98, PHK⁺07] show diabetes incidence among the Pima Indians is significantly higher among subjects with $bmi \geq 30$; diabetes incidence is higher for a person with impaired glucose tolerance.

$$\begin{aligned} (\text{Diabetes Rule 1}) \quad & (gluc \leq 126) && \Rightarrow \text{-diabetes,} \\ (\text{Diabetes Rule 2}) \quad & (gluc \geq 126) \wedge (gluc \leq 140) \wedge (bmi \leq 30) && \Rightarrow \text{-diabetes,} \\ (\text{Diabetes Rule 3}) \quad & (gluc \geq 126) \wedge (gluc \leq 140) \wedge (bmi \geq 30) && \Rightarrow \text{diabetes,} \\ (\text{Diabetes Rule 4}) \quad & (gluc \geq 140) && \Rightarrow \text{diabetes.} \end{aligned}$$

The pedigree function is a feature which provides a measure heredity on the subject's diabetes risk. A subject with high heredity who is at least 31 years old is at a significantly increased risk for diabetes in the next five years [SED⁺88]:

$$\begin{aligned} (\text{Diabetes Rule 5}) \quad & (pedf \leq 0.5) \wedge (age \leq 31) && \Rightarrow \text{-diabetes,} \\ (\text{Diabetes Rule 6}) \quad & (pedf \geq 0.5) \wedge (age \geq 31) && \Rightarrow \text{diabetes.} \end{aligned}$$



(left) Results averaged over 10 runs on a hold-out 412 point test set using all 6 rules; (right) Approximate decision-tree showing Diabetes Rule 6 before and after refinement (if true then left branch). Leaf nodes classify data according to ?diabetes.

arkSVM algorithms produce local solutions from small data sets and imperfect expert advice. Learned models can generalize well and provide refined advice easily interpreted by domain experts.

References

- [FMS03] G. Fung, O. L. Mangasarian, and J. W. Shavlik. Knowledge-based support vector classifiers. In *Advances in Neural Information Processing Systems*, volume 15, pages 521–528, 2003.
- [HFC⁺98] M. I. Harris, K. M. Flegal, C. C. Cowie, M. S. Eberhardt, D. E. Goldstein, R. R. Little, H. M. Wiedmeyer, and D. D. Byrd-Holt. Prevalence of diabetes, impaired fasting glucose, and impaired glucose tolerance in U.S. adults. *Diabetes Care*, 21(4):518–524, 1998.
- [MWS⁺07] R. Maclin, E. W. Wild, J. W. Shavlik, L. Torrey, and T. Walker. Refining rules incorporated into knowledge-based support vector learners via successive linear programming. In *Proc. 22nd AAAI*, pages 584–589, 2007.
- [PHK⁺07] M. E. Pavkov, R. L. Hanson, W. C. Knowler, P. H. Bennett, J. Krakoff, and R. G. Nelson. Changing patterns of Type 2 diabetes incidence among Pima Indians. *Diabetes Care*, 30(7):1758–1763, 2007.
- [SED⁺88] J. W. Smith, J. E. Everhart, W. C. Dickson, W. C. Knowler, and R. S. Johannes. Using the ADAP learning algorithm to forecast the onset of diabetes mellitus. In *Proc. Symp. Comp. Apps. and Medical Care*, pages 261–265, 1988.