

The Solution Path of Multicategory Support Vector Machines

Yoonkyung Lee
Department of Statistics
The Ohio State University
<http://www.stat.ohio-state.edu/~ykleee>

Joint work with Zhenhuan Cui

Method of Regularization

Find $f(\mathbf{x}) \in \mathcal{F}$ minimizing

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(\mathbf{x}_i)) + \lambda J(f).$$

- ▶ \mathcal{F} : a class of candidate functions
- ▶ $J(f)$: complexity of the model f
- ▶ $\lambda > 0$: a regularization parameter
- ▶ Solution path: $\lambda \rightarrow \hat{f}_\lambda$

Classification

- ▶ A training data set $\{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$
- ▶ $\mathbf{x} = (x_1, \dots, x_p) \in \mathbb{R}^p$
- ▶ $y \in \mathcal{Y} = \{1, \dots, k\}$
- ▶ Learn a rule $\phi : \mathbb{R}^p \rightarrow \mathcal{Y}$ from the training data, which can be generalized to novel cases.

Support Vector Machine

Vapnik (1995)

- ▶ $y_i \in \{-1, 1\}$ for $k = 2$
- ▶ Find $f(\mathbf{x}) = b + h(\mathbf{x})$ with $h \in \mathcal{H}_K$ minimizing

$$\frac{1}{n} \sum_{i=1}^n (1 - y_i f(\mathbf{x}_i))_+ + \lambda \|h\|_{\mathcal{H}_K}^2,$$

where K is the reproducing kernel of \mathcal{H}_K , and λ is a regularization parameter.

- ▶ Classification rule: $\phi(\mathbf{x}) = \text{sign} [f(\mathbf{x})]$

Hinge loss

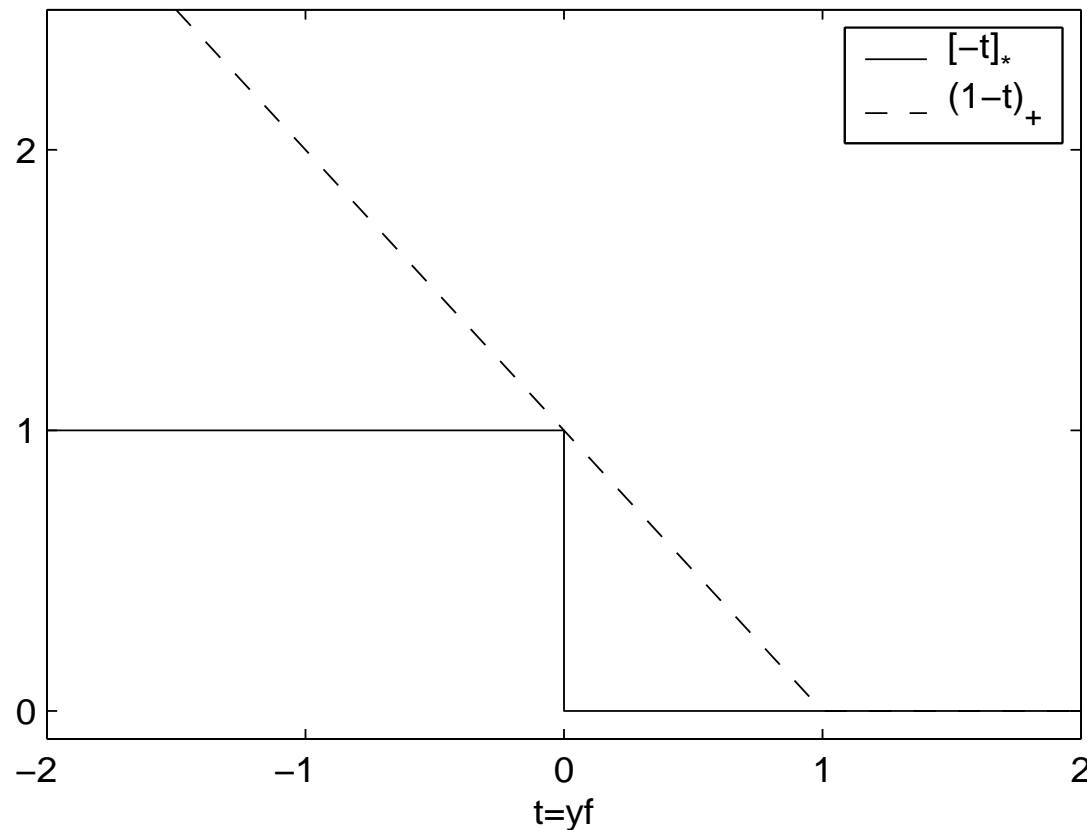


Figure: $(1 - yf(\mathbf{x}))_+$ is an upper bound of the misclassification loss function $l(y \neq \phi(\mathbf{x})) = [-yf(\mathbf{x})]_* \leq (1 - yf(\mathbf{x}))_+$ where $[t]_* = l(t \geq 0)$ and $(t)_+ = \max\{t, 0\}$.

SVM computation

- ▶ Quadratic programming problem
- ▶ Getting solutions for a fixed λ :

e.g. Sequential Minimal Optimization [Platt (1999)], SVM light [Joachims (1999)], LIBSVM [Hsu and Lin (2002)], <http://www.kernel-machines.org/software.html>

- ▶ Getting the entire regularization path ranging from the simplest majority rule to the complete overfit to data as λ decreases.

SVM solution path [Hastie et al. (2004)]

SVM when $k > 2$

Lee, Lin & Wahba, *JASA* (2004)

- ▶ $\mathbf{y} = (y^1, \dots, y^k)$: class code with $y^j = 1$ and $-1/(k-1)$ elsewhere, if $y = j$.
- ▶ Find $\mathbf{f} = (f^1, \dots, f^k) = (b^1 + h^1(\mathbf{x}), \dots, b^k + h^k(\mathbf{x}))$ with $h^j \in \mathcal{H}_K$ and the sum-to-zero constraint minimizing

$$\frac{1}{n} \sum_{i=1}^n \sum_{j \neq y_i} (f^j(\mathbf{x}_i) - y_i^j)_+ + \frac{\lambda}{2} \sum_{j=1}^k \|h^j\|^2.$$

- ▶ Classification rule: $\phi(\mathbf{x}) = \arg \max_j [f^j(\mathbf{x})]$

MSVM hinge loss

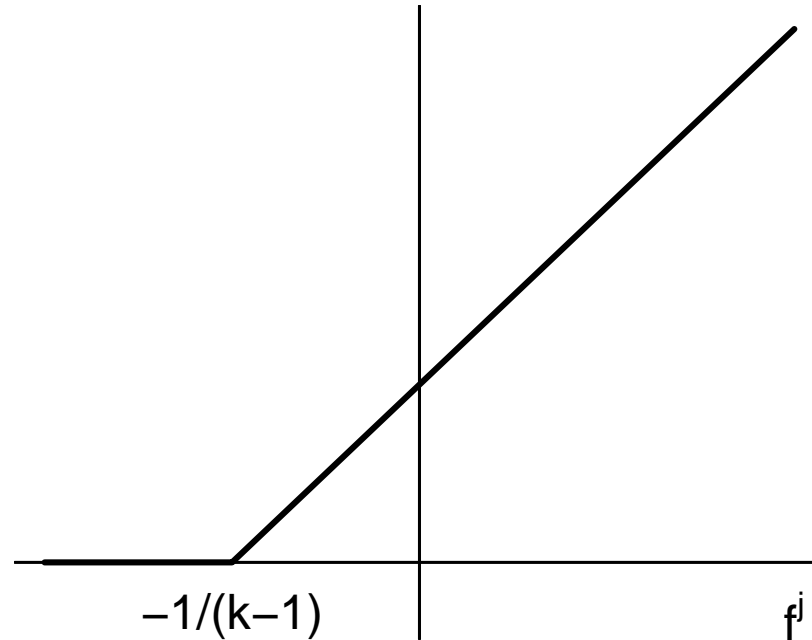


Figure: MSVM component loss $(f^j - y^j)_+$ for $y \neq j$, where $y^j = -1/(k-1)$.

Optimization problem for MSVM

- ▶ By the representer theorem, $\hat{f}^j(\mathbf{x}) = b^j + \sum_{i=1}^n c_i^j K(\mathbf{x}_i, \mathbf{x})$
- ▶ $\mathbf{c}^j = (c_1^j, \dots, c_n^j)^\top$, $\mathbf{K}_n = (K(\mathbf{x}_i, \mathbf{x}_j))$ and $f_i^j = \hat{f}^j(\mathbf{x}_i)$
- ▶ Primal problem: minimize

$$L_P(\mathbf{c}, \mathbf{b}, \boldsymbol{\xi}) = \frac{1}{n} \sum_{i=1}^n \sum_{j \neq y_i} \xi_i^j + \frac{\lambda}{2} \sum_{j=1}^k (\mathbf{c}^j)^\top \mathbf{K}_n \mathbf{c}^j$$

subject to

$$f_i^j - y_i^j \leq \xi_i^j \quad \text{for } i, j$$

$$\xi_i^j \geq 0 \quad \text{for } i, j$$

$$\sum_{j=1}^k (b^j \mathbf{e} + \mathbf{K}_n \mathbf{c}^j) = \mathbf{0}.$$

Optimization problem for MSVM (cont'd)

- ▶ α_i^j : Lagrange multiplier for $f_i^j - y_i^j \leq \xi_i^j$
Set $\alpha_i^j = 0$ for $j = y_i$.
- ▶ $\alpha^j = (\alpha_1^j, \dots, \alpha_n^j)^\top$ and $\mathbf{y}^j = (y_1^j, \dots, y_n^j)^\top$
- ▶ Dual problem: maximize

$$L_D(\alpha) = -\frac{1}{2n} \sum_{j=1}^k (\alpha^j - \bar{\alpha})^\top K_n (\alpha^j - \bar{\alpha}) - \lambda \sum_{j=1}^k (\alpha^j)^\top \mathbf{y}^j$$

subject to

$$0 \leq \alpha_i^j \leq 1 \quad \text{for } i, j$$
$$(\alpha^j - \bar{\alpha})^\top \mathbf{e} = 0 \quad \text{for } j = 1, \dots, k.$$

Optimality conditions

- ▶ By the Karush-Kuhn-Tucker (KKT) complementarity conditions, the MSVM solution at λ satisfies

$$\begin{aligned}\alpha_i^j (f_i^j - y_i^j - \xi_i^j) &= 0, \\ (1 - \alpha_i^j) \xi_i^j &= 0.\end{aligned}$$

Also, $0 \leq \alpha_i^j \leq 1$ and $\xi_i^j \geq 0$.

- ▶ Categorize (i, j) into three sets depending on $f_i^j - y_i^j$.

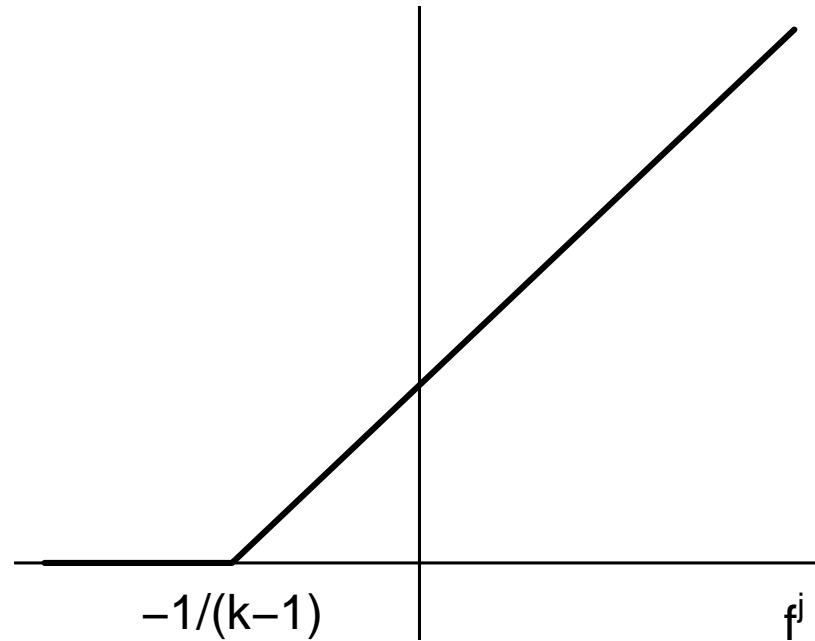


Figure: MSVM component loss $(f^j - y^j)_+$ where $y^j = -1/(k - 1)$.

$$\mathcal{E} = \{(i, j) \mid f_i^j - y_i^j = 0, \xi_i^j = 0, 0 \leq \alpha_i^j \leq 1\}, \text{ Elbow set}$$

$$\mathcal{U} = \{(i, j) \mid f_i^j - y_i^j > 0, \xi_i^j > 0, \alpha_i^j = 1\}, \text{ Upper set}$$

$$\mathcal{L} = \{(i, j) \mid f_i^j - y_i^j < 0, \xi_i^j = 0, \alpha_i^j = 0\}, \text{ Lower set}$$

Characterization of the entire solution path

- ▶ Initialize α_i^j 's for sufficiently large λ so that they correspond to the majority rule.
- ▶ Keep track of the events that change the elbow set.
- ▶ $\lambda_0 > \lambda_1 > \lambda_2 > \dots$, a decreasing sequence of breakpoints of λ at which the elbow set \mathcal{E} changes.
- ▶ Construct the path sequentially by solving a system of linear equations for α_i^j with $(i, j) \in \mathcal{E}_\ell$.

Piecewise linearity of the SVM solution

- ▶ The coefficient c_i^j path of the MSVM is linear in $1/\lambda$ on the interval $(\lambda_{\ell+1}, \lambda_\ell)$.
- ▶ The path of b^j is also linear in $1/\lambda$ on the interval $(\lambda_{\ell+1}, \lambda_\ell)$ if there is at most one empty elbow set \mathcal{E}_ℓ^j at λ_ℓ . Otherwise, the path of b^j for the empty \mathcal{E}_ℓ^j can be arbitrary.

Generating the joints $\{\lambda_\ell\}$

- ▶ Given λ_ℓ , consider the following possible events.
 1. An index (i, j) in \mathcal{E}_ℓ leaves the elbow set, and α_i^j ($0 \leq \alpha_i^j \leq 1$) becomes either 0 or 1.
 2. An index (i, j) in \mathcal{L}_ℓ or \mathcal{U}_ℓ joins the elbow set, and $\hat{f}^j(\mathbf{x}_i)$ is then y_i^j .
- ▶ The two potential events propose candidate values of $\lambda_{\ell+1}$.
- ▶ The next break point $\lambda_{\ell+1}$ is determined by the largest $\lambda < \lambda_\ell$ among the potential values.

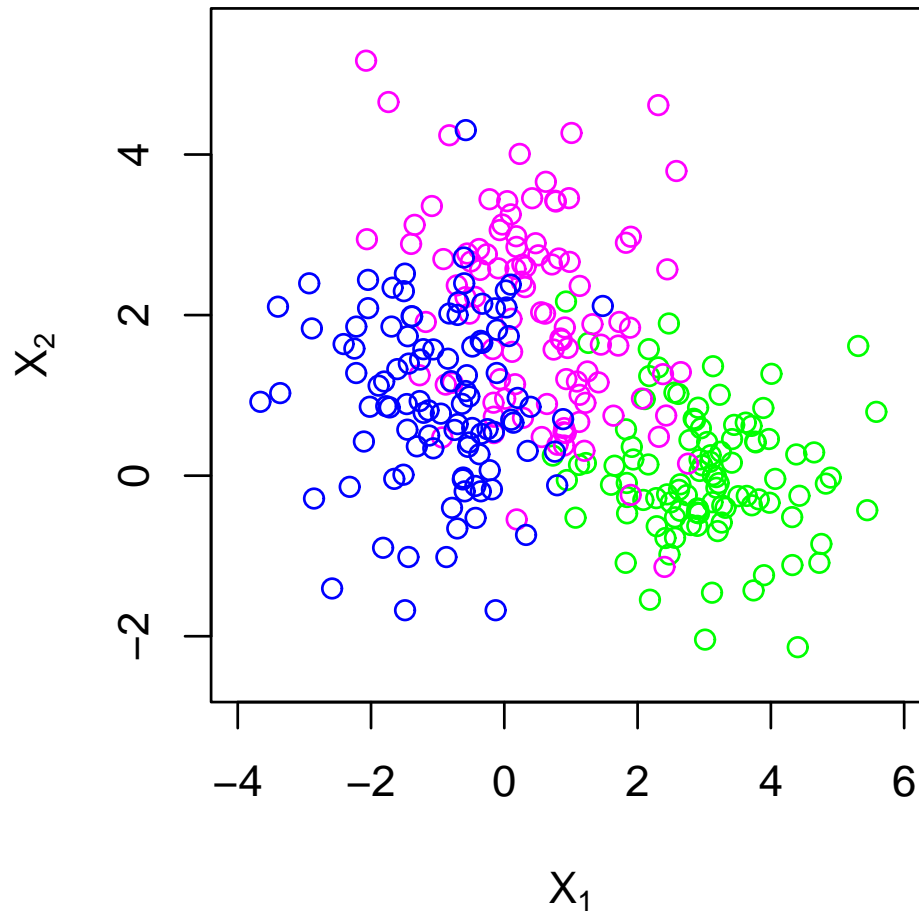


Figure: A scatter plot of the training data. Class 1: green, Class 2: magenta, and Class 3: blue.

Test error rates

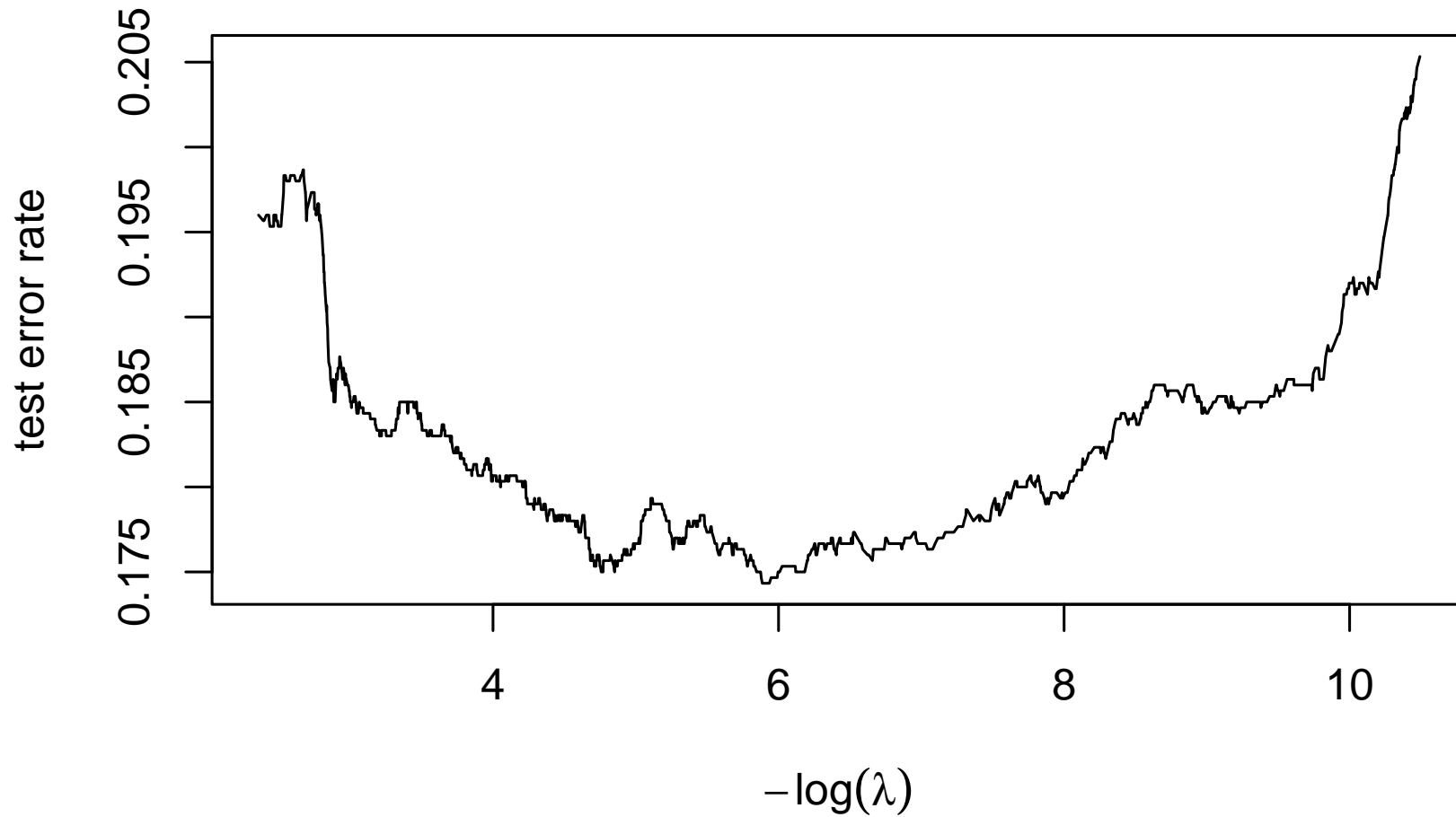


Figure: Test error rate as a function of λ .

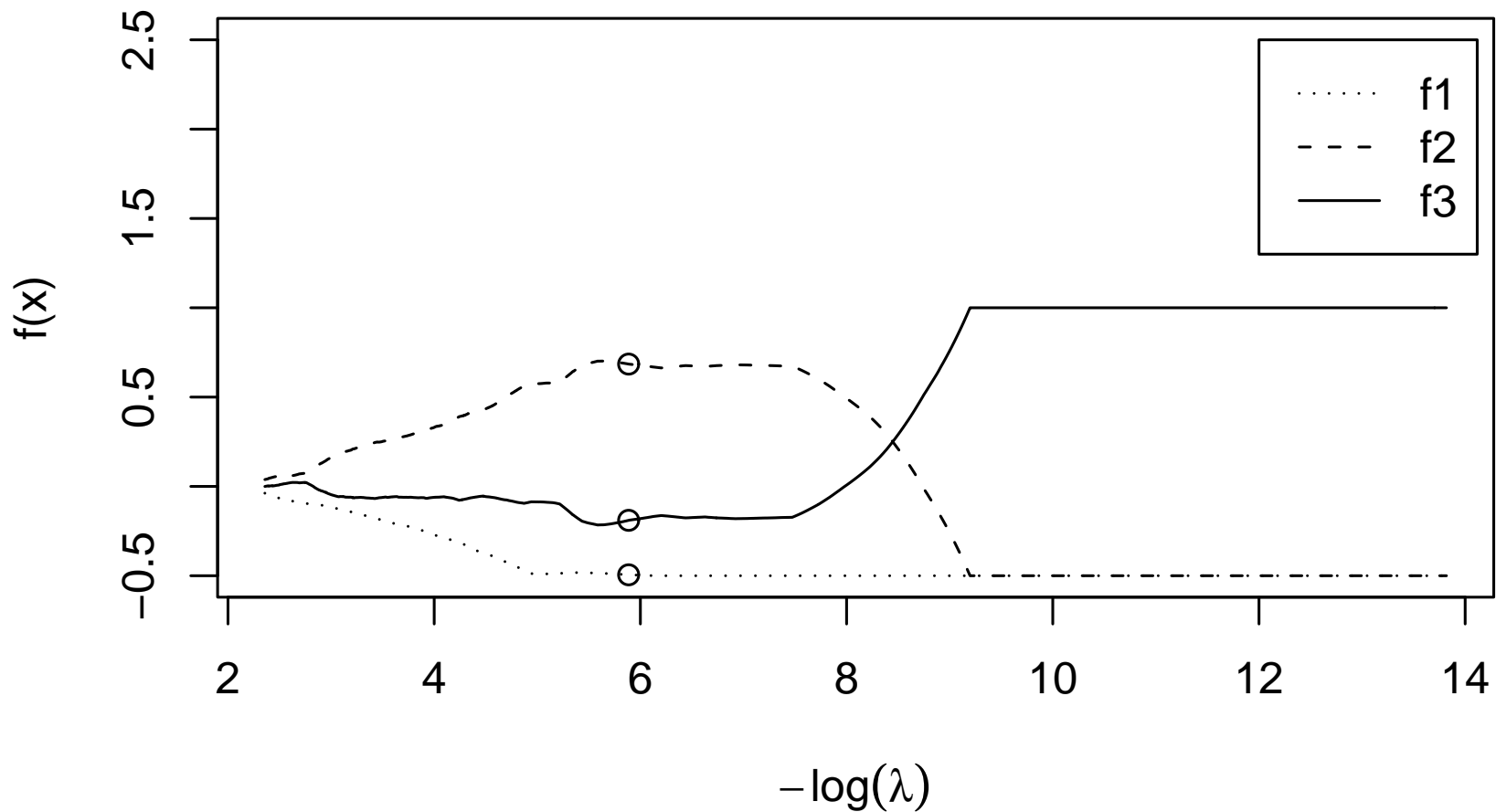


Figure: The entire paths of $\hat{f}_\lambda^1(\mathbf{x}_i)$, $\hat{f}_\lambda^2(\mathbf{x}_i)$, and $\hat{f}_\lambda^3(\mathbf{x}_i)$ for an outlying instance \mathbf{x}_i from class 3. The circles correspond to λ with the minimum test error rate.

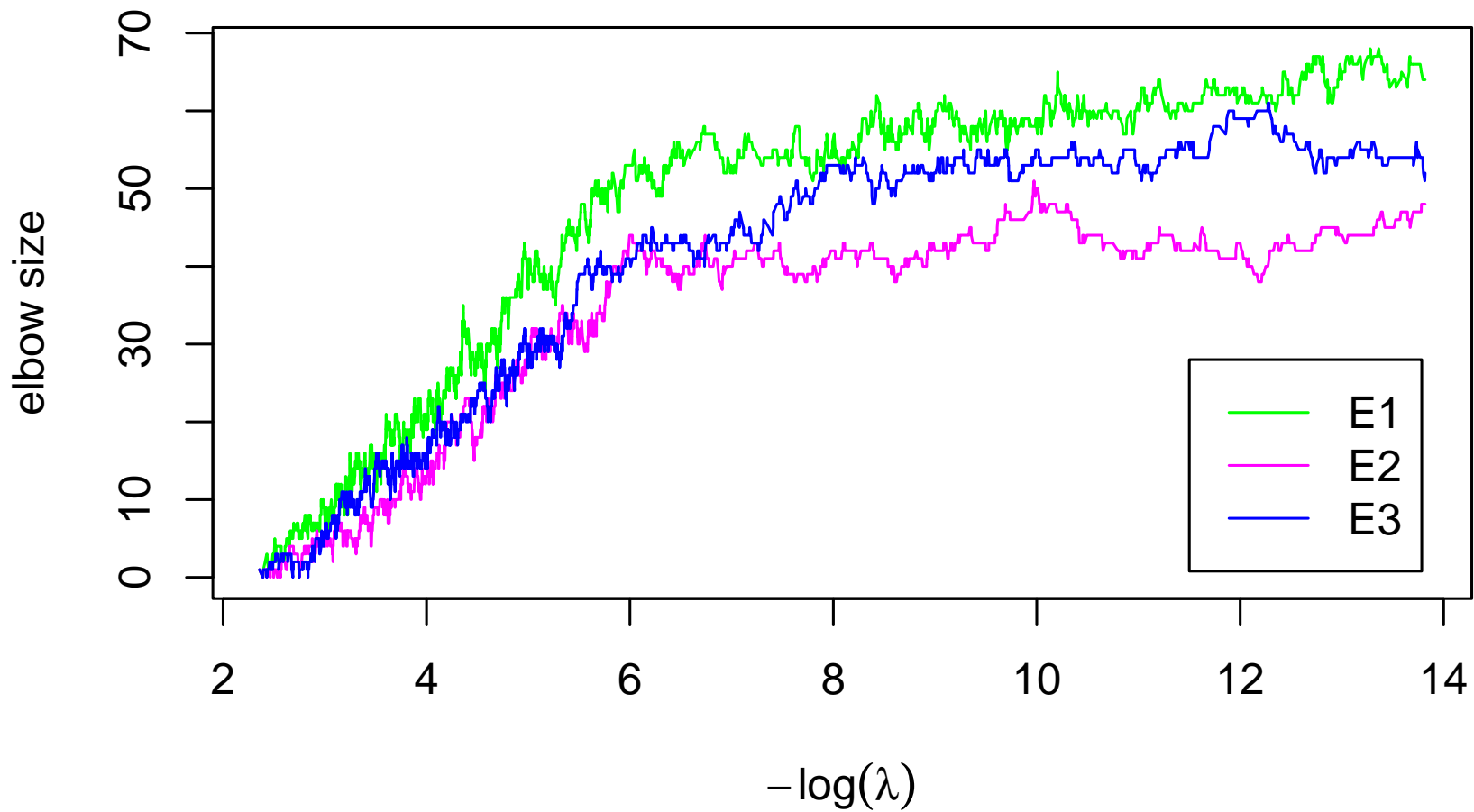


Figure: The size of elbow set \mathcal{E}_ℓ^j for three classes as a function λ .

Further generalization

- ▶ The path finding algorithm can be extended to the case with unequal misclassification costs ($\mathbf{L}(\mathbf{y}) = (L_y^1, \dots, L_y^k)$) or different data weights (w_i).

- ▶ $\mathcal{L}(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) = w_i \sum_{j \neq y_i} L_{y_i}^j (f_i^j - y_i^j)_+$

- ▶ $0 \leq \alpha_i^j \leq w_i L_{y_i}^j$

Complementarity conditions:

$$\alpha_i^j (f_i^j - y_i^j - \xi_i^j) = 0 \text{ and } (w_i L_{y_i}^j - \alpha_i^j) \xi_i^j = 0$$

- ▶ Initialize α_i^j 's for sufficiently large λ so that they correspond to the rule minimizing the total weighted costs

$$\sum_{i=1, y_i \neq j}^n w_i L_{y_i}^j.$$

- ▶ Piecewise linearity of the coefficient path remains true.

Computational complexity

- ▶ Proportional to the number of break points λ_ℓ .
- ▶ At each λ_ℓ , it takes $O(|\mathcal{E}_\ell|^3 + (k - 1)n|\mathcal{E}_\ell^{\mathcal{I}}|)$ operations (a system of linear equations needs to be solved, and the next break point has to be determined.)
- ▶ $|\mathcal{E}_\ell|$, intermediate function evaluations, and the number of break points tend to be proportional to $(k - 1)$.
- ▶ \mathcal{E}_ℓ may grow too big for large data sets and flexible kernel functions.

Scaling down data sets

- ▶ Basis thinning:

Recall $\hat{f}^j(\mathbf{x}) = b^j + \sum_{i=1}^n c_i^j K(\mathbf{x}_i, \mathbf{x})$.

Choose $\mathcal{I}^* \subset \{1, \dots, n\}$ with $|\mathcal{I}^*| \ll n$ such that $\text{span}\{K(\mathbf{x}_{i^*}, \cdot), i^* \in \mathcal{I}^*\} \approx \text{span}\{K(\mathbf{x}_i, \cdot), i = 1, \dots, n\}$.

- ▶ Random sampling
 - ▶ Statistically equivalent blocks [Wilks (1962)]
 - ▶ Cluster analysis: $\{(\tilde{\mathbf{x}}_{i^*}, \tilde{y}_{i^*}) \text{ with weight } w_{i^*}\}$
- ▶ Data squashing [Dumouchel et al. (1999)]:
Binning and generating pseudo data with weights matching the moments of the full data.

Statistically equivalent blocks

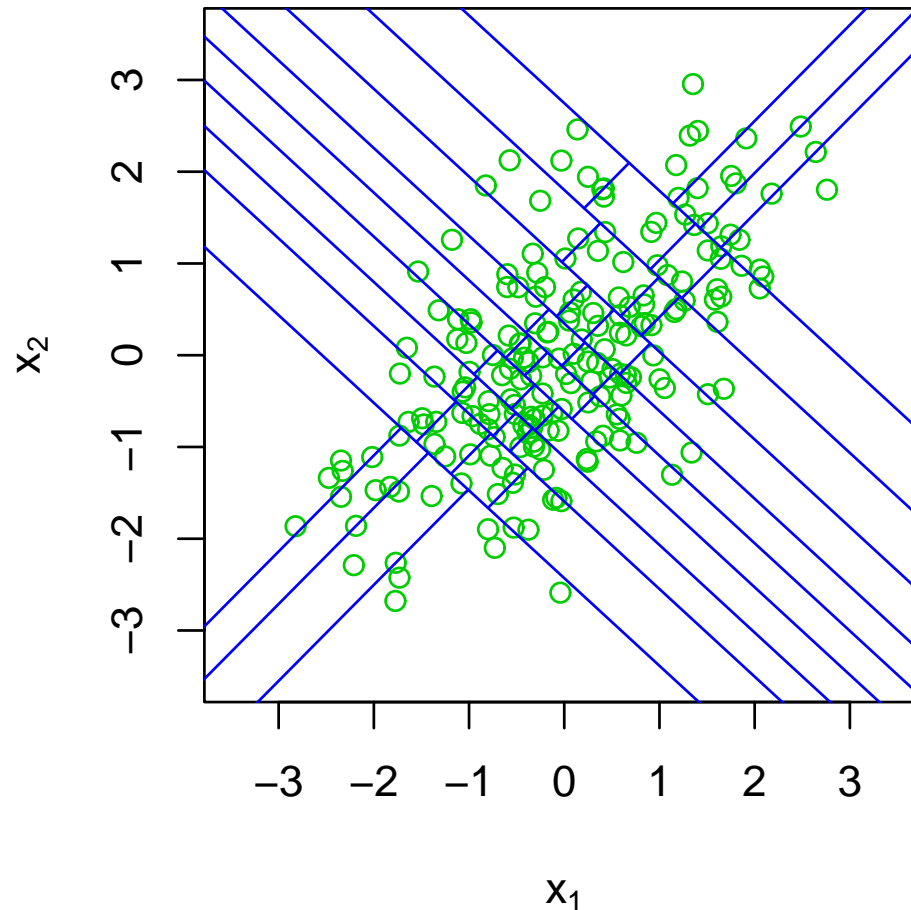


Figure: $n = 200$, 40 (10×4) statistically equivalent blocks obtained by using the principal components as the ordering functions.

Clusters

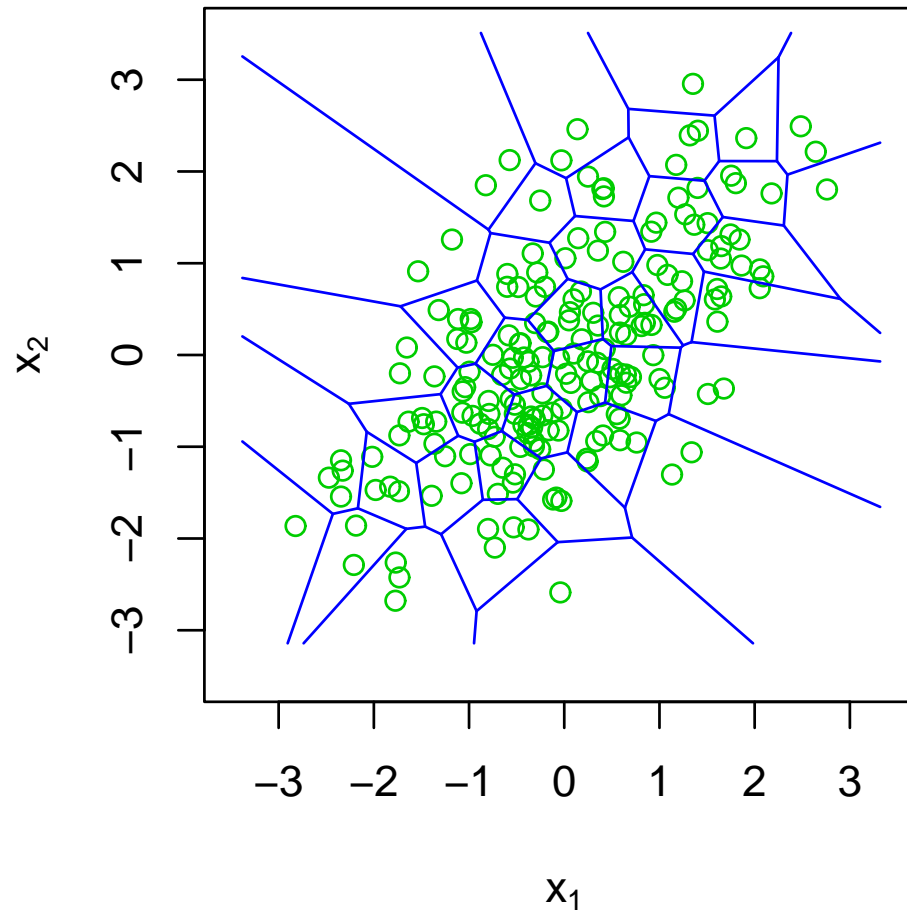


Figure: $n = 200$, 40 clusters formed by the k-means algorithm and the Voronoi tessellation of the cluster centers.

Finding an approximate solution path

- ▶ Scale down the data by choosing a representative subsample for each class.
- ▶ Find the regularized solution path for the reduced data.
- ▶ “Out-of-bag sample” as in bootstrap can be used for tuning.

Comparison of data reduction methods

- ▶ The three-class problem with sample size 120×3 and reduction factor 6 (20×3)
- ▶ Number of replicates: 100
- ▶ Test error rates over 1200×3 cases (the Bayes rate ≈ 0.1773)

Method	Random sampling	SEB	Cluster
Mean	0.19345	0.18307	0.18216
SD	0.01452	0.00768	0.00720

- ▶ The quality of approximate solutions would depend on reduction factor, data partitioning scheme, dimensionality, and the Bayes rate.

Summary

- ▶ Characterize the solution path of MSVM for effective computation and tuning.
- ▶ Direct use of the complementarity conditions and the piecewise linearity of the coefficient path.
- ▶ Find an approximate solution path to cope with large data sets.
- ▶ Data geometry (manifold learning, nonlinear PCA) can guide data reduction procedures.