Support Vector Machines Meet Goldfarb Cubes

Bernd Gärtner, IPAM Conference "Quo Vadis Hirsch Conjecture?"

January 19, 2011

Wednesday, January 19, 2011

Disproving the Hirsch Conjecture of Machine Learning

Bernd Gärtner, IPAM Conference "Quo Vadis Hirsch Conjecture?"

January 19, 2011

Wednesday, January 19, 2011

Outline

- Introduction to Support Vector Machines, and their interpretation as parametric *quadratic* programs
- Statement of the main result: disproof of a conjecture by Hastie et al. concerning the solution path complexity of support vector machines
- Parametric *linear* programs, the shadow vertex pivot rule, and the Goldfarb cubes
- Putting it together: constructing exponentially long solution paths for support vector machines, using Goldfarb cubes (sketch)

Support Vector Machines (SVM)

 Suppose an optical character recognition needs to distinguish between the letters 'A' and 'B'



 Suppose an optical character recognition needs to distinguish between the letters 'A' and 'B'



 Training phase: the system gets to see letters 'A' and 'B' plus the information whether it's an 'A' or a 'B'

 Suppose an optical character recognition needs to distinguish between the letters 'A' and 'B'



- Training phase: the system gets to see letters 'A' and 'B' plus the information whether it's an 'A' or a 'B'
- After the training phase, the system is supposed to decide on its own which letter it sees.

Solution: map training letters to points in some high-dimensional space, with label 'A' (blue) or 'B' (red), based e.g. on a pixel-based representation



Solution: map training letters to points in some high-dimensional space, with label 'A' (blue) or 'B' (red), based e.g. on a pixel-based representation

 Separate the 'A's from the 'B's by a simple shape, for example the maximum-margin hyperplane (separating hyperplane of maximum distance to any point).

R

A

B

Solution: map training letters to points in some high-dimensional space, with label 'A' (blue) or 'B' (red), based e.g. on a pixel-based representation

- Separate the 'A's from the 'B's by a simple shape, for example the maximum-margin hyperplane (separating hyperplane of maximum distance to any point).
- * Classify an unknown letter according to this hyperplane ($\bigcirc = 'B'$)

R

A

B

Support Vector Machines The primal program

Computation of Maximum-Margin Hyperplane is a convex quadratic program:

$$\begin{array}{lll} \min_{\mathbf{w},b} & \frac{1}{2}\mathbf{w}^T\mathbf{w} \\ \text{subject to} & \mathbf{w}^T\mathbf{p}_i - b & \geq & 1, \quad \mathbf{p}_i \in A \\ & \mathbf{w}^T\mathbf{p}_i - b & \leq & 1, \quad \mathbf{p}_i \in B \end{array}$$

$$\mathbf{w}^T \mathbf{x} = b$$

Support Vector Machines The primal program, class labels

Computation of Maximum-Margin Hyperplane is a convex quadratic program:

$$\begin{array}{ll} \min_{\mathbf{w},b} & \frac{1}{2}\mathbf{w}^T\mathbf{w} \\ \text{subject to} & y_i(\mathbf{w}^T\mathbf{p}_i-b) & \geq 1, \quad i=1,2,\ldots,n \end{array}$$

Support Vector Machines Dual program

 \min_{λ_i} subject to

$$\|\mathbf{p} - \mathbf{q}\|^{2}$$

$$\sum_{\mathbf{p}_{i} \in A} \lambda_{i} \mathbf{p}_{i} = \mathbf{p}$$

$$\sum_{\mathbf{p}_{i} \in B} \lambda_{i} \mathbf{p}_{i} = \mathbf{q}$$

$$\sum_{\mathbf{p}_{i} \in A} \lambda_{i} = 1$$

$$\sum_{\mathbf{p}_{i} \in B} \lambda_{i} = 1$$

$$\lambda_{i} \geq 0 \quad i = 1, 2, \dots, n$$

竹

Support Vector Machines Dual program

 \min_{λ_i} subject to

$$\begin{aligned} \|\mathbf{p} - \mathbf{q}\|^2 \\ \sum_{\mathbf{p}_i \in A} \lambda_i \mathbf{p}_i &= \mathbf{p} \\ \sum_{\mathbf{p}_i \in B} \lambda_i \mathbf{p}_i &= \mathbf{q} \\ \sum_{\mathbf{p}_i \in A} \lambda_i &= 1 \\ \sum_{\mathbf{p}_i \in B} \lambda_i &= 1 \\ \lambda_i &\geq 0 \quad i = 1, 2 \end{aligned}$$

Geometrically, this is the problem of finding the distance between the two polytopes *conv*(*A*) and *conv*(*B*)



 \dots, n

Support Vector Machines Dual program

 \min_{λ_i} subject to

$$\begin{aligned} \|\mathbf{p} - \mathbf{q}\|^2 \\ \sum_{\mathbf{p}_i \in A} \lambda_i \mathbf{p}_i &= \mathbf{p} \\ \sum_{\mathbf{p}_i \in B} \lambda_i \mathbf{p}_i &= \mathbf{q} \\ \sum_{\mathbf{p}_i \in A} \lambda_i &= 1 \\ \sum_{\mathbf{p}_i \in B} \lambda_i &= 1 \\ \lambda_i &\geq 0 \quad i = 1, \end{aligned}$$

The points with positive λ_i are the support vectors

 $2,\ldots,n$

Geometrically, this is the problem of finding the distance between the two polytopes conv(A) and conv(B)



Support Vector Machines Dual program, class labels

 \min_{λ_i}

 $\sum_{i,j}\lambda_i\lambda_jy_iy_j\mathbf{p}_i^T\mathbf{p}_j$ subject to $\sum_{\mathbf{p}_i \in A} \lambda_i$ $= 1 \\ \geq 0 \quad i = 1, 2, \dots, n$ $\sum_{\mathbf{p}_i \in B} \lambda_i$ λ_i



Support Vector Machines Dual program, class labels

 $\begin{array}{l} \min_{\lambda_i} \\ \text{subject to} \end{array}$

$$\sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{p}_i^T \mathbf{p}_j$$

$$\sum_{\mathbf{p}_i \in A} \lambda_i = 1$$

$$\sum_{\mathbf{p}_i \in B} \lambda_i = 1$$

$$\lambda_i \geq 0 \quad i = 1, 2, \dots, n$$

The dual program does not need to know the points, just their pairwise scalar products!



Support Vector Machines Nonlinear separators

 Other separating shapes (e.g. spheres) may be the "right" shapes for the given problem



Support Vector Machines Nonlinear separators

 Other separating shapes (e.g. spheres) may be the "right" shapes for the given problem



* Apply *lifting* to reduce to separating hyperplane in higher dimension!





* Use preimage of maximum-margin hyperplane in lifted space

* Apply lifting map $\ell: (x, y) \mapsto (x, y, x^2 + y^2)$

Dual quadratic program becomes

 $\min_{\lambda_i} \sum_{\substack{i,j \ \lambda_i \lambda_j y_i y_j k(\mathbf{p}_i, \mathbf{p}_j) \\ \text{subject to}} \sum_{\substack{\mathbf{p}_i \in A \ \lambda_i \\ \sum_{\mathbf{p}_i \in B} \lambda_i \\ \lambda_i} = 1 \\ \geq 0 \quad i = 1, 2, \dots, n }$

$$k(\mathbf{p},\mathbf{q}) = \mathbf{p}_x \mathbf{q}_x + \mathbf{p}_y \mathbf{q}_y + (\mathbf{p}_x^2 + \mathbf{p}_y^2)(\mathbf{q}_x^2 + \mathbf{q}_y)^2$$

* Apply lifting map $\ell: (x, y) \mapsto (x, y, x^2 + y^2)$

Dual quadratic program becomes

$$\min_{\lambda_i} \sum_{\substack{i,j \ \lambda_i \lambda_j y_i y_j k(\mathbf{p}_i, \mathbf{p}_j) \\ \text{subject to}} \sum_{\substack{p_i \in A \ \lambda_i \\ \sum_{\mathbf{p}_i \in B} \lambda_i \\ \lambda_i} = 1 \\ \geq 0 \quad i = 1, 2, \dots, n$$

$$k(\mathbf{p}, \mathbf{q}) = \mathbf{p}_x \mathbf{q}_x + \mathbf{p}_y \mathbf{q}_y + (\mathbf{p}_x^2 + \mathbf{p}_y^2)(\mathbf{q}_x^2 + \mathbf{q}_y)^2 \longleftarrow function$$

* Let $k : \mathbf{R}^d \times \mathbf{R}^d \to \mathbf{R}$ be symmetric and positive semidefinite, meaning that

$$\sum_{i,j} \lambda_i \lambda_j k(\mathbf{p}_i, \mathbf{p}_j) \ge 0, \quad \forall \mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n, \lambda_1, \lambda_2, \dots, \lambda_n$$

* Let $k : \mathbf{R}^d \times \mathbf{R}^d \to \mathbf{R}$ be symmetric and positive semidefinite, meaning that

$$\sum_{i,j} \lambda_i \lambda_j k(\mathbf{p}_i, \mathbf{p}_j) \ge 0, \quad \forall \mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n, \lambda_1, \lambda_2, \dots, \lambda_n$$

* Then the kernel *k* is a scalar product in some (possibly infinite dimensional) space *H* (*Mercer's Theorem*)

* Let $k : \mathbf{R}^d \times \mathbf{R}^d \to \mathbf{R}$ be symmetric and positive semidefinite, meaning that

$$\sum_{i,j} \lambda_i \lambda_j k(\mathbf{p}_i, \mathbf{p}_j) \ge 0, \quad \forall \mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n, \lambda_1, \lambda_2, \dots, \lambda_n$$

* Then the kernel *k* is a scalar product in some (possibly infinite dimensional) space *H* (*Mercer's Theorem*)

The quadratic program implicitly computes the maximum-margin hyperplane in *H*

$$\min_{\lambda_i} \sum_{\substack{i,j \\ \lambda_i \lambda_j y_i y_j k(\mathbf{p}_i, \mathbf{p}_j) \\ \sum_{\substack{p_i \in A \\ \sum_{\mathbf{p}_i \in B \\ \lambda_i}} \lambda_i = 1 \\ \lambda_i \geq 0 \quad i = 1, 2, \dots, n }$$

Then the kernel k is a scalar product in some (possibly infinite dimensional) space H (Mercer's Theorem)

The quadratic program implicitly computes the maximum-margin hyperplane in *H*

 $\min_{\lambda_i} \sum_{\substack{i,j \ \lambda_i \lambda_j y_i y_j k(\mathbf{p}_i, \mathbf{p}_j) \\ \text{subject to}} \sum_{\substack{\sum_{\mathbf{p}_i \in A} \lambda_i \\ \sum_{\mathbf{p}_i \in B} \lambda_i \\ \lambda_i} = 1 \\ \lambda_i \geq 0 \quad i = 1, 2, \dots, n }$

 Separating shape in the original space may be "weird" and hard to compute, but we can still classify unknown points based on it!

* Then the kernel k is a scalar product in some (possibly infinite dimensional) space *H* (*Mercer's Theorem*)

 \min_{λ_i}

The quadratic program implicitly computes the maximum-margin hyperplane in H

$$\min_{\lambda_i} \sum_{\substack{i,j \ \lambda_i \lambda_j y_i y_j k(\mathbf{p}_i, \mathbf{p}_j) \\ \text{subject to}} \sum_{\substack{\mathbf{p}_i \in A \ \lambda_i \\ \sum_{\mathbf{p}_i \in B} \lambda_i \\ \lambda_i} = 1 \\ \lambda_i \geq 0 \quad i = 1, 2, \dots, n$$

* Separating shape in the original space may be "weird" and hard to compute, but we can still classify unknown points based on it!

* Compute
$$\mathbf{w}^T \mathbf{x} - b$$
, where $\mathbf{w}^T \mathbf{x} = \sum_i \lambda_i y_i k(\mathbf{p}_i, \mathbf{x})$ support
 $b = \sum_i^i \lambda_i y_i k(\mathbf{p}_i, \mathbf{p}_j) - y_j$

Visualization of nonlinear separator in original space

$$k(\mathbf{p}, \mathbf{q}) = e^{-\mu \|\mathbf{p} - \mathbf{q}\|}$$

"Radial Basis Function" (RBF)



 In many cases, a proper separation cannot be expected, due to noise and outliers



- In many cases, a proper separation cannot be expected, due to noise and outliers
- * Allow misclassifications of training points, but penalize them!



- In many cases, a proper separation cannot be expected, due to noise and outliers
- * Allow misclassifications of training points, but penalize them!
- * v-SVM, where v is a fixed parameter

$$\begin{array}{ll} \min_{\mathbf{w},b,\rho,\xi_{i}} & \frac{1}{2}\mathbf{w}^{T}\mathbf{w} - \nu\rho & + & \frac{1}{n}\sum_{i=1}^{n}\xi_{i} \\ \text{subject to} & y_{i}(\mathbf{w}^{T}\mathbf{p}_{i} - b) & \geq & \rho - \xi_{i}, \quad i = 1, 2, \dots, n \\ & \xi_{i} & \geq & 0, \quad i = 1, 2, \dots, n \\ & \rho & \geq & 0 \\ & & & & \text{previously,} \\ & & & & \text{we had 1 here} \end{array}$$

- In many cases, a proper separation cannot be expected, due to noise and outliers
- * Allow misclassifications of training points, but penalize them!
- * dual v-SVM (with $\mu := \frac{2}{n\nu}$):

 $\min_{\lambda_{i}} \sum_{i,j} \lambda_{i} \lambda_{j} y_{i} y_{j} \mathbf{p}_{i}^{T} \mathbf{p}_{j}$ subject to $\sum_{\mathbf{p}_{i} \in A} \lambda_{i} = 1$ $\sum_{\mathbf{p}_{i} \in B} \lambda_{i} = 1$ $\lambda_{i} \geq 0 \quad i = 1, 2, \dots, n$ $\max \longrightarrow \lambda_{i} \leq \mu \quad i = 1, 2, \dots, n$

Support Vector Machines Soft Margin: The Geometric View

 The dual v-SVM is the problem of finding the distance between two reduced convex hulls



Support Vector Machines Soft Margin: The Geometric View

 The dual v-SVM is the problem of finding the distance between two reduced convex hulls



* $\mu = 1$: no reduction

 $\mu = 1/2$: edges start to disappear

Support Vector Machines Soft Margin: Parameter Choice

* How does v (or μ in the dual) influence the separator?



Support Vector Machines Soft Margin: Parameter Choice

* How does v (or μ in the dual) influence the separator?



Small v: small margin (overfitting)

Large v: large margin (underfitting)


* How to choose the parameter in a given application?

- * How to choose the parameter in a given application?
- Standard method: grid search (discretize the relevant parameter range and try all values in the discretized range; one QP per value)

- * How to choose the parameter in a given application?
- Standard method: grid search (discretize the relevant parameter range and try all values in the discretized range; one QP per value)
- Solution path method: compute the solution of the dual SVM as a piecewise linear function in μ, then grid search (one lookup per value)

$$\mu = 1 \qquad \mu = 0.84 \qquad \min_{\lambda_i} \qquad \sum_{\substack{i,j \\ j \\ i \\ j \\ i \\ j \\ i \\ k_i \\ \lambda_i \\ \lambda_i \\ k_i \\$$

- * How to choose the parameter in a given application?
- Standard method: grid search (discretize the relevant parameter range and try all values in the discretized range; one QP per value)
- Solution path method: compute the solution of the dual SVM as a piecewise linear function in μ, then grid search (one lookup per value)



* Performance depends on the number of bends in the solution path

Support Vector Machines The Hastie et al. Conjecture

* Hastie et al. 2004: The entire regularization path for the support vector machine

Support Vector Machines The Hastie et al. Conjecture

- Hastie et al. 2004: The entire regularization path for the support vector machine
- * **Conjecture:** The number of bends is O (min (|A|, |B|) = O(n)

Support Vector Machines The Hastie et al. Conjecture

- Hastie et al. 2004: The entire regularization path for the support vector machine
- * **Conjecture:** The number of bends is O (min (|A|, |B|) = O(n)
- Conjecture was repeatedly stated in other articles, and it was experimentally "confirmed" in a large number of applications

Support Vector Machines The Conjecture is False

- * Hastie et al. 2004: The entire regularization path for the support vector machine
- * **Conjecture:** The number of bends is O (min (|A|, |B|) = O(n)
- Conjecture was repeatedly stated in other articles, and it was experimentally "confirmed" in a large number of applications
- * **Theorem (G., Jaggi, Maria 2010) :** There are point sets *A* and *B* in ddimensional space with |A| = 2d, |B| = 2, for which the solution path of the dual v-SVM has an *exponential* number of at least $2^d / 4$ many bends.

Parametric Linear Programs

 Recall: primal v-SVM is a parametric quadratic program (parameter in the objective function):

$$\min_{\mathbf{w},b,\rho,\xi_{i}} \quad \frac{1}{2}\mathbf{w}^{T}\mathbf{w} - \nu\rho \quad + \quad \frac{1}{n}\sum_{i=1}^{n}\xi_{i} \\ \text{subject to} \quad y_{i}(\mathbf{w}^{T}\mathbf{p}_{i} - b) \geq \rho - \xi_{i}, \quad i = 1, 2, \dots, n \\ \xi_{i} \quad \geq \quad 0, \quad i = 1, 2, \dots, n \\ \rho \quad \geq \quad 0$$

Parametric Linear Programs

 Recall: primal v-SVM is a parametric quadratic program (parameter in the objective function):

$$\min_{\mathbf{w},b,\rho,\xi_{i}} \quad \frac{1}{2} \mathbf{w}^{T} \mathbf{w} - \nu \rho \quad + \quad \frac{1}{n} \sum_{i=1}^{n} \xi_{i}$$
subject to
$$y_{i}(\mathbf{w}^{T} \mathbf{p}_{i} - b) \geq \rho - \xi_{i}, \quad i = 1, 2, \dots, n$$

$$\xi_{i} \qquad \geq \quad 0, \quad i = 1, 2, \dots, n$$

$$\rho \qquad \geq \quad 0$$

* Parametric linear program:

$$\begin{array}{ll} \min & (\mathbf{c}^T + \nu \mathbf{d})^T \mathbf{x} \\ \text{subject to} & A\mathbf{x} = b & 0 \le \nu \le \infty \\ & \mathbf{x} \ge 0 \end{array}$$

* Start with $\nu = \infty$ (**d** rules)

С

* Start with $\nu = \infty$ (**d** rules)

Decrease v...



С

* Start with $\nu = \infty$ (**d** rules)

Decrease v...



C









Same picture in higher d, with blue polygon the shadow of the LP

The Shadow Vertex / Parametric Objective Simplex Algorithm $_{\min c^T x}$



The Goldfarb Cubes

- Worst case inputs for the shadow vertex simplex algorithm
- The *d*-dimensional Goldfarb cube is a deformed unit cube with all its 2^d vertices appearing on a 2-dimensional shadow

The Goldfarb Cubes

- Worst case inputs for the shadow vertex simplex algorithm
- The *d*-dimensional Goldfarb cube is a deformed unit cube with all its 2^d vertices appearing on a 2-dimensional shadow

Solution Paths of Parametric Linear Programs...

- * can be exponentially long in the worst case
- It is not implausible that a similar result holds for the solution path of SVM (specific parametric *quadratic* programs)

Solution Paths of Parametric Linear Programs...

- * can be exponentially long in the worst case
- It is not implausible that a similar result holds for the solution path of SVM (specific parametric *quadratic* programs)
- * Obstacles:
 - SVM are very special parametric quadratic programs
 - * Goldfarb cube is not an instance of an SVM
 - * we want a *nondegenerate* SVM with a long solution path

* Outline:

* Work in the dual (μ-SVM: distance between reduced convex hulls)

* Outline:

- * Work in the dual (μ-SVM: distance between reduced convex hulls)
- First convex hull: dual Goldfarb cube (2d points); there exists a 2d plane S such that the dual Goldfarb cube intersects S in a polygon with 2^d edges



* Outline:

 $w_{(3,-)}$

3d

- * Work in the dual (μ-SVM: distance between reduced convex hulls)
- * First convex hull: **dual Goldfarb cube** (2*d* points); there exists a 2*d* plane S such that the dual Goldfarb cube intersects S in a polygon with 2^d edges

 $w_{(3,1)}$

8d

* Outline:

- * Work in the dual (μ-SVM: distance between reduced convex hulls)
- * First convex hull: **dual Goldfarb cube** (2*d* points); there exists a 2*d* plane S such that the dual Goldfarb cube intersects S in a polygon with 2^d edges

 $w_{(3,1)}$

Second convex hull: line segment in S

 $\mathbf{w}_{(3)}$

* Outline:



* Outline:



* Outline:



* Outline:



* Outline:



* Outline:



* Outline:



- Technical difficulties:
 - Point of optimal distance in the dual Goldfarb cube may be far away from S and does not "walk along" S



- Technical difficulties:
 - Point of optimal distance in the dual Goldfarb cube may be far away from S and does not "walk along" S
 - Solution: stretching!

S

stretched dual Goldfarb cube almost walks on intersection with S, same number of bends

- Technical difficulties:
 - Reduction of stretched dual Goldfarb cube may affect the intersection with S



- * Technical difficulties:
 - Reduction of stretched dual Goldfarb cube may affect the intersection with S
 - Solution: make line segment so long that everything happens for μ very close to 1


 The solution path of a support vector machine may be exponentially long in the number of training points; the linear-length conjecture of Hastie et al. fails dramatically

Summary

- The solution path of a support vector machine may be exponentially long in the number of training points; the linear-length conjecture of Hastie et al. fails dramatically
- For "us" polytope people, this result is what we would expect, knowing the exponential lower bounds for solution paths of parametric *linear* programs (Goldfarb cubes)

Summary

- The solution path of a support vector machine may be exponentially long in the number of training points; the linear-length conjecture of Hastie et al. fails dramatically
- For "us" polytope people, this result is what we would expect, knowing the exponential lower bounds for solution paths of parametric *linear* programs (Goldfarb cubes)
- * *Approximate* solution paths are always short (Giesen, Jaggi, Laue)

Summary

- The solution path of a support vector machine may be exponentially long in the number of training points; the linear-length conjecture of Hastie et al. fails dramatically
- For "us" polytope people, this result is what we would expect, knowing the exponential lower bounds for solution paths of parametric *linear* programs (Goldfarb cubes)
- * *Approximate* solution paths are always short (Giesen, Jaggi, Laue)
- The machine learning community is proficient in statistics, but not in geometry; it can be beneficial to look at some of their problems from a geometric viewpoint