

γ SVMs + SMO

P. Agius – L8, Spring 2008

Contents

- The γ parameter
- γ SVC – primal and dual
- γ SVR – primal and dual
- Properties of γ
- Choosing γ
- Results
- SMO – another efficient method

P. Agius – L8, Spring 2008

The U parameter

In our earlier models , we had a 'c' parameter which we used to regulate between function complexity and training error.

What happens if we trade 'c' for something more intuitive?

Here comes U ...

This U parameter will be used to weight the slacks and to encourage a wider margin.

P. Agius – L8, Spring 2008

U SVC - primal

Using m samples ...

$$\begin{array}{ll} \text{SVC} & \min_{w,\xi,b} \quad \frac{1}{2}w^T w + C \sum_i \xi_i \\ & \text{s.t.} \quad y_i(w^T x_i + b) \geq 1 - \xi_i \\ & \quad \quad \xi_i \geq 0 \end{array}$$

$$\begin{array}{ll} \text{U SVC} & \min_{w,\xi,b,\rho} \quad \frac{1}{2}w^T w - \nu\rho + \frac{1}{m} \sum_i \xi_i \\ & \text{s.t.} \quad y_i(w^T x_i + b) \geq \rho - \xi_i \\ & \quad \quad \xi_i \geq 0, \rho \geq 0 \end{array}$$

We have a new variable ρ to optimize and a new parameter U

So the two classes are now separated by a margin of width 2ρ

P. Agius – L8, Spring 2008

U SVC - dual

Using m samples ...

$$\begin{aligned} \text{SVM} \quad & \max_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i x_j \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned}$$

$$\begin{aligned} \text{U SVC} \quad & \max_{\alpha} \quad - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i x_j \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq \frac{1}{m} \\ & \sum_i \alpha_i y_i = 0 \\ & \sum_i \alpha_i \geq \nu \end{aligned}$$

P. Agius – L8, Spring 2008

SVR

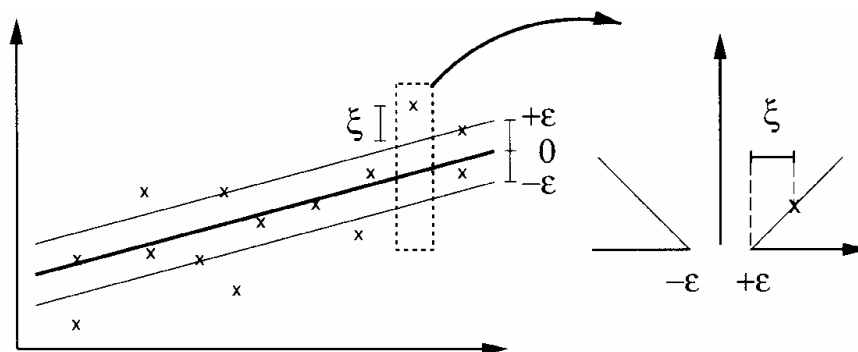


Figure 1: In SV regression, a desired accuracy ε is specified a priori. It is then attempted to fit a tube with radius ε to the data. The trade-off between model complexity and points lying outside the tube (with positive slack variables ξ) is determined by minimizing the expression 1.5.

P. Agius – L8, Spring 2008

U SVR - primal

Using m samples ...

$$\begin{array}{ll} \text{SVR} & \min_{w,b,\xi,\hat{\xi}} \quad \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i^2 + \hat{\xi}_i^2) \\ & \text{s.t.} \quad (\langle w, \phi(x_i) \rangle + b) - y_i \leq \epsilon + \xi_i \\ & \quad y_i - (\langle w, \phi(x_i) \rangle + b) \leq \epsilon + \hat{\xi}_i \\ & \quad \xi_i, \hat{\xi}_i \geq 0 \end{array}$$

$$\begin{array}{ll} \text{U SVR} & \min_{w,b,\xi,\hat{\xi},\epsilon} \quad \frac{1}{2} \|w\|^2 + C[\nu\epsilon + \frac{1}{m} \sum_i (\xi_i^2 + \hat{\xi}_i^2)] \\ & \text{s.t.} \quad (\langle w, \phi(x_i) \rangle + b) - y_i \leq \epsilon + \xi_i \\ & \quad y_i - (\langle w, \phi(x_i) \rangle + b) \leq \epsilon + \hat{\xi}_i \\ & \quad \xi_i, \hat{\xi}_i \geq 0 \end{array}$$

P. Agius – L8, Spring 2008

U SVR - dual

Using m samples ...

$$\begin{array}{ll} \text{SVR} & \max \quad \sum_i (\hat{\alpha}_i - \alpha_i) y_i - \epsilon \sum_i (\hat{\alpha}_i + \alpha_i) \\ & \quad - 0.5 \sum_{i,j} (\hat{\alpha}_i - \alpha_i) (\hat{\alpha}_j - \alpha_j) \kappa(x_i, x_j) \\ & \text{s.t.} \quad \sum_i (\hat{\alpha}_i - \alpha_i) = 0 \\ & \quad 0 \leq \alpha_i, \hat{\alpha}_i \leq C \end{array}$$

$$\begin{array}{ll} \text{U SVR} & \max \quad \sum_i (\hat{\alpha}_i - \alpha_i) y_i - 0.5 \sum_{i,j} (\hat{\alpha}_i - \alpha_i) (\hat{\alpha}_j - \alpha_j) \kappa(x_i, x_j) \\ & \text{s.t.} \quad \sum_i (\hat{\alpha}_i - \alpha_i) = 0 \\ & \quad 0 \leq \alpha_i, \hat{\alpha}_i \leq \frac{C}{m} \\ & \quad \sum_i (\hat{\alpha}_i + \alpha_i) \leq C\nu \end{array}$$

P. Agius – L8, Spring 2008

Properties of U

Proposition 1 Assume $\varepsilon > 0$. The following statements hold:

- (i) ν is an upper bound on the fraction of errors.
- (ii) ν is a lower bound on the fraction of SVs.
- (iii) Suppose the data (2) were generated iid from a distribution $P(\mathbf{x}, y) = P(\mathbf{x})P(y|\mathbf{x})$ with $P(y|\mathbf{x})$ continuous. With probability 1, asymptotically, ν equals both the fraction of SVs and the fraction of errors.

P. Agius – L8, Spring 2008

Intuition behind U properties

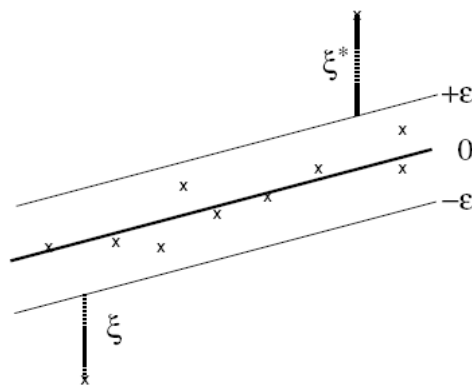


Figure 1: Graphical depiction of the ν -trick. Imagine increasing ε , starting from 0. The first term in $\nu\varepsilon + \frac{1}{\ell} \sum_{i=1}^{\ell} (\xi_i + \xi_i^*)$ (cf. (4)) will increase proportionally to ν , while the second term will decrease proportionally to the fraction of points outside of the tube. Hence, ε will grow as long as the latter fraction is larger than ν . At the optimum, it therefore must be $\leq \nu$ (Proposition 1, (i)). Next, imagine decreasing ε , starting from some large value. Again, the change in the first term is proportional to ν , but this time, the change in the second term is proportional to the fraction of SVs (even points on the edge of the tube will contribute). Hence, ε will shrink as long as the fraction of SVs is smaller than ν , eventually leading to Proposition 1, (ii).

P. Agius – L8, Spring 2008

Intuition behind U properties

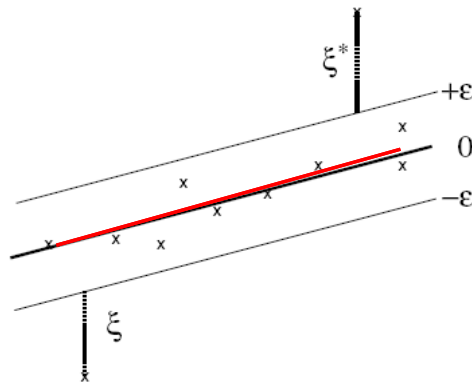


Figure 1: Graphical depiction of the ν -trick. Imagine increasing ϵ , starting from 0. The first term in $\nu\epsilon + \frac{1}{T} \sum_{i=1}^T (\xi_i + \xi_i^*)$ (cf. (4)) will increase proportionally to ν , while the second term will decrease proportionally to the fraction of points outside of the tube. Hence, ϵ will grow as long as the latter fraction is larger than ν . At the optimum, it therefore must be $\leq \nu$ (Proposition 1, (i)). Next, imagine decreasing ϵ , starting from some large value. Again, the change in the first term is proportional to ν , but this time, the change in the second term is proportional to the fraction of SVs (even points on the edge of the tube will contribute). Hence, ϵ will shrink as long as the fraction of SVs is smaller than ν , eventually leading to Proposition 1, (ii).

P. Agius – L8, Spring 2008

Varying the U parameter

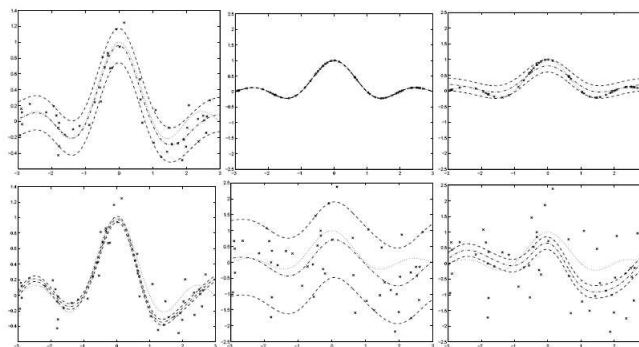


Figure 2: *Left*: ν -SV regression with $\nu = 0.2$ (top) and $\nu = 0.8$ (bottom). The larger ν allows more points to lie outside the tube (see Sec. 2). The algorithm automatically adjusts ϵ to 0.22 (top) and 0.04 (bottom). Shown are the sinc function (dotted), the regression f and the tube $f \pm \epsilon$. *Middle*: ν -SV regression on data with noise $\sigma = 0$ (top) and $\sigma = 1$ (bottom). In both cases, $\nu = 0.2$. The tube width automatically adjusts to the noise (top: $\epsilon = 0$, bottom: $\epsilon = 1.19$). *Right*: ϵ -SV regression (Vapnik, 1995) on data with noise $\sigma = 0$ (top) and $\sigma = 1$ (bottom). In both cases, $\epsilon = 0.2$ — this choice, which has to be specified a priori, is ideal for neither case: in the top figure, the regression estimate is biased; in the bottom figure, ϵ does not match the external noise (cf. Smola et al., 1998).

P. Agius – L8, Spring 2008

Results using U models

Table 1: Results for the Boston housing benchmark; *top*: ν -SVR, *bottom*: ϵ -SVR. **MSE**: Mean squared errors, **STD**: standard deviations thereof (100 trials), **Errors**: fraction of training points outside the tube, **SVs**: fraction of training points which are SVs.

ν	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
automatic ϵ	2.6	1.7	1.2	0.8	0.6	0.3	0.0	0.0	0.0	0.0
MSE	9.4	8.7	9.3	9.5	10.0	10.6	11.3	11.3	11.3	11.3
STD	6.4	6.8	7.6	7.9	8.4	9.0	9.6	9.5	9.5	9.5
Errors	0.0	0.1	0.2	0.2	0.3	0.4	0.5	0.5	0.5	0.5
SVs	0.3	0.4	0.6	0.7	0.8	0.9	1.0	1.0	1.0	1.0

ϵ	0	1	2	3	4	5	6	7	8	9	10
MSE	11.3	9.5	8.8	9.7	11.2	13.1	15.6	18.2	22.1	27.0	34.3
STD	9.5	7.7	6.8	6.2	6.3	6.0	6.1	6.2	6.6	7.3	8.4
Errors	0.5	0.2	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SVs	1.0	0.6	0.4	0.3	0.2	0.1	0.1	0.1	0.1	0.1	0.1

P. Agius – L8, Spring 2008

Faster SVMs – related to SVMlight (some of Kostja's talk)

- Chunking and decomposition
- SMO – Sequential Minimization Optimization

**Using Analytic QP and Sparseness to Speed
Training of Support Vector Machines**

John C. Platt
Microsoft Research
1 Microsoft Way
Redmond, WA 98052
jplatt@microsoft.com

P. Agius – L8, Spring 2008

Chunking

Large datasets → large kernels ...
storage ☹ computational costs ☹

But solution can be derived using only the support vectors!

Chunking

- Pick a chunk of data and train a SVM on it.
- Keep the SVs and toss the rest
- Identify M points that most violate the KKT conditions
- Update 'working set' with SVs and M 'bad boys'
- Terminate at some stopping criterion

P. Agius – L8, Spring 2008

Decomposition

Goal: Optimize the problem by acting on small subsets at a time
(does not find ALL active constraints)

Every time you add a new point to the 'working set', a point is removed.
Keep working set fixed at a certain size (SVMlight)

Convergence of chunking and decomposition has not been proved.
But in practice they have been found to work really well.

What happens if we take decomposition to the extreme?
i.e. solve only two points at a time?!!!

P. Agius – L8, Spring 2008

SMO – Sequential Minimization Optimization

Break down the problem by solving it two points at a time!
No QP required!!!

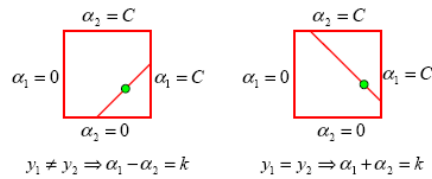


Figure 1: The Lagrange multipliers α_1 and α_2 must fulfill all of the constraints of the full problem. The inequality constraints cause the Lagrange multipliers to lie in the box. The linear equality constraint causes them to lie on a diagonal line.

P. Agius – L8, Spring 2008

SMO Performance

Experiment	SMO Time (sec)	SVM ^{lin} Time (sec)	Chunking Time (sec)	SMO Scaling Exponent	SVM ^{lin} Scaling Exponent	Chunking Scaling Exponent
AdultLin	13.7	217.9	20711.3	1.8	2.1	3.1
AdultLinD	21.9	n/a	21141.1	1.0	n/a	3.0
WebLin	339.9	3980.8	17164.7	1.6	2.2	2.5
WebLinD	4589.1	n/a	17332.8	1.5	n/a	2.5
AdultGaussK	442.4	284.7	11910.6	2.0	2.0	2.9
AdultGauss	523.3	737.5	n/a	2.0	2.0	n/a
AdultGaussKD	1433.0	n/a	14740.4	2.5	n/a	2.8
AdultGaussD	1810.2	n/a	n/a	2.0	n/a	n/a
WebGaussK	2477.9	2949.5	23877.6	1.6	2.0	2.0
WebGauss	2538.0	6923.5	n/a	1.6	1.8	n/a
WebGaussKD	23365.3	n/a	50371.9	2.6	n/a	2.0
WebGaussD	24758.0	n/a	n/a	1.6	n/a	n/a
MNIST	19387.9	38452.3	33109.0	n/a	n/a	n/a

Table 2: Timings of algorithms on various data sets.

P. Agius – L8, Spring 2008