



UNIVERSITY OF PISA

DEPARTMENT OF COMPUTER SCIENCE

COMPUTATIONAL MATHEMATICS  
WILDCARD PROJECT NR. 5 WITH MACHINE LEARNING  
GROUP 35

---

# Support Vector Machines

---

*Author:*

**Donato Meoli**  
d.meoli@studenti.unipi.it

March, 2021

# Contents

<b>List of Figures</b>	<b>2</b>
<b>List of Tables</b>	<b>3</b>
<b>1 Track</b>	<b>4</b>
<b>2 Abstract</b>	<b>4</b>
<b>3 Linear Support Vector Classifier</b>	<b>5</b>
3.1 Hinge loss . . . . .	5
3.1.1 Primal formulation . . . . .	5
3.1.2 Wolfe Dual formulation . . . . .	7
3.1.3 Lagrangian Dual formulation . . . . .	9
3.2 Squared Hinge loss . . . . .	10
3.2.1 Primal formulation . . . . .	10
<b>4 Linear Support Vector Regression</b>	<b>11</b>
4.1 Epsilon-insensitive loss . . . . .	12
4.1.1 Primal formulation . . . . .	12
4.1.2 Wolfe Dual formulation . . . . .	12
4.1.3 Lagrangian Dual formulation . . . . .	14
4.2 Squared Epsilon-insensitive loss . . . . .	16
4.2.1 Primal formulation . . . . .	16
<b>5 Nonlinear Support Vector Machines</b>	<b>17</b>
5.1 Polynomial kernel . . . . .	17
5.2 Gaussian RBF kernel . . . . .	17
<b>6 Gradient Descent</b>	<b>19</b>
6.1 Momentum . . . . .	19
6.1.1 Standard . . . . .	19
6.1.2 Nesterov . . . . .	19
<b>7 AdaGrad</b>	<b>20</b>
<b>8 Sequential Minimal Optimization</b>	<b>21</b>
8.1 Classification . . . . .	21
8.2 Regression . . . . .	22
<b>9 Experiments</b>	<b>23</b>
9.1 Support Vector Classifier . . . . .	23
9.1.1 Hinge loss . . . . .	23
9.1.2 Squared Hinge loss . . . . .	25
9.2 Support Vector Regression . . . . .	25
9.2.1 Epsilon-insensitive loss . . . . .	25
9.2.2 Squared Epsilon-insensitive loss . . . . .	26
<b>10 Conclusions</b>	<b>27</b>
<b>References</b>	<b>29</b>

List of Figures

1	Linear SVC hyperplane . . . . .	6
2	SVC Hinge loss with optimization steps . . . . .	7
3	SVC Squared Hinge loss with optimization steps . . . . .	11
4	Linear SVR hyperplane . . . . .	12
5	SVR Epsilon-insensitive loss with optimization steps . . . . .	13
6	SVC Squared Epsilon-insensitive loss with optimization steps . . . . .	16
7	Polynomial SVM hyperplanes . . . . .	17
8	Gaussian SVM hyperplanes . . . . .	18

## List of Tables

1	SVC Primal formulation results with Hinge loss . . . . .	23
2	Linear SVC Wolfe Dual formulation results with Hinge loss . . . . .	23
3	Linear SVC Lagrangian Dual formulation results with Hinge loss . . . . .	23
4	Nonlinear SVC Wolfe Dual formulation results with Hinge loss . . . . .	24
5	Nonlinear SVC Lagrangian Dual formulation results with Hinge loss . . . . .	24
6	SVC Primal formulation results with Squared Hinge loss . . . . .	25
7	SVR Primal formulation results with Epsilon-insensitive loss . . . . .	25
8	Linear SVR Wolfe Dual formulation results with Epsilon-insensitive loss . . . . .	26
9	Linear SVR Lagrangian Dual formulation results with Epsilon-insensitive loss . . . . .	27
10	Nonlinear SVR Lagrangian Dual formulation results with Epsilon-insensitive loss . . . . .	28
11	SVR Primal formulation results with Squared Epsilon-insensitive loss . . . . .	29

## 1 Track

(M1.1) is a *Support Vector Classifier (SVC)* with the *hinge* loss.

(A1.1.1) is the *AdaGrad* algorithm [1], a *deflected subgradient* method for solving the SVC in its *primal* formulation.

(A1.1.2) is the *Sequential Minimal Optimization (SMO)* algorithm [2] (see [3] for improvements), an ad hoc *active set* method for training a SVC in its *Wolfe dual* formulation with *linear*, *polynomial* and *gaussian* kernels.

(A1.1.3) is the *AdaGrad* algorithm [1], a *deflected subgradient* method for solving the SVC in its *Lagrangian dual* formulation with *linear*, *polynomial* and *gaussian* kernels.

(M1.2) is a *Support Vector Classifier (SVC)* with the *squared hinge* loss.

(A1.2.1) is a *momentum descent* approach, an *accelerated gradient* method for solving the SVC in its *primal* formulation.

(M2.1) is a *Support Vector Regression (SVR)* with the *epsilon-insensitive* loss.

(A2.1.1) is the *AdaGrad* algorithm [1], a *deflected subgradient* method for solving the SVR in its *primal* formulation.

(A2.1.2) is the *Sequential Minimal Optimization (SMO)* algorithm [4] (see [5] for improvements), an ad hoc *active set* method for training a SVR in its *Wolfe dual* formulation with *linear*, *polynomial* and *gaussian* kernels.

(A2.1.3) is the *AdaGrad* algorithm [1], a *deflected subgradient* method for solving the SVR in its *Lagrangian dual* formulation with *linear*, *polynomial* and *gaussian* kernels.

(M2.2) is a *Support Vector Regression (SVR)* with the *squared epsilon-insensitive* loss.

(A2.2.1) is a *momentum descent* approach, an *accelerated gradient* method for solving the SVR in its *primal* formulation.

## 2 Abstract

A *Support Vector Machine* is a learning model used both for *classification* and *regression* tasks whose goal is to construct a *maximum margin separator*, i.e., a decision boundary with the largest distance from the nearest training data points.

The aim of this report is to compare the *primal*, the *Wolfe dual* and the *Lagrangian dual* formulations of this model in terms of *numerical precision*, *accuracy* and *complexity*.

Firstly, I will provide a detailed mathematical derivation of the model for all these formulations, then I will propose two algorithms to solve the optimization problem in case of *constrained* or *unconstrained* formulation of the problem, explaining their theoretical properties, i.e., *convergence* and *complexity*.

Finally, I will show some experiments for *linearly* and *nonlinearly* separable generated datasets to compare the performance of different *kernels*, also by comparing the *custom* results with *sklearn* SVM implementations, i.e., *liblinear* and *libsvm* implementations, and *cvxopt* QP solver.

### 3 Linear Support Vector Classifier

Given  $n$  training points, where each input  $x_i$  has  $m$  attributes, i.e., is of dimensionality  $m$ , and is in one of two classes  $y_i = \pm 1$ , i.e., our training data is of the form:

$$\{(x_i, y_i), x_i \in \mathbb{R}^m, y_i = \pm 1, i = 1, \dots, n\} \quad (1)$$

For simplicity we first assume that data are (not fully) linearly separable in the input space  $x$ , meaning that we can draw a line separating the two classes when  $m = 2$ , a plane for  $m = 3$  and, more in general, a hyperplane for an arbitrary  $m$ .

Support vectors are the examples closest to the separating hyperplane and the aim of support vector machines is to orientate this hyperplane in such a way as to be as far as possible from the closest members of both classes, i.e., we need to maximize this margin.

This hyperplane is represented by the equation  $w^T x + b = 0$ . So, we need to find  $w$  and  $b$  so that our training data can be described by:

$$\begin{aligned} w^T x_i + b &\geq +1 - \xi_i, \forall y_i = +1 \\ w^T x_i + b &\leq -1 + \xi_i, \forall y_i = -1 \\ \xi_i &\geq 0 \quad \forall_i \end{aligned} \quad (2)$$

where the positive slack variables  $\xi_i$  are introduced to allow misclassified points. In this way data points on the incorrect side of the margin boundary will have a penalty that increases with the distance from it.

These two equations can be combined into:

$$\begin{aligned} y_i(w^T x_i + b) &\geq 1 - \xi_i \quad \forall_i \\ \xi_i &\geq 0 \quad \forall_i \end{aligned} \quad (3)$$

The margin is equal to  $\frac{1}{\|w\|}$  and maximizing it subject to the constraint in 3 while as we are trying to reduce the number of misclassifications is equivalent to finding:

$$\begin{aligned} \min_{w, b, \xi} \quad & \|w\| + C \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall_i \\ & \xi_i \geq 0 \quad \forall_i \end{aligned} \quad (4)$$

Minimizing  $\|w\|$  is equivalent to minimizing  $\frac{1}{2}\|w\|^2$ , but in this form we will deal with a convex optimization problem that has more desirable convergence properties. So we need to find:

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2}\|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall_i \\ & \xi_i \geq 0 \quad \forall_i \end{aligned} \quad (5)$$

where the parameter  $C$  controls the trade-off between the slack variable penalty and the size of the margin.

#### 3.1 Hinge loss

##### 3.1.1 Primal formulation

The general primal unconstrained formulation takes the form:

$$\min_{w, b} \mathcal{R}(w, b) + C \sum_{i=1}^n \mathcal{L}(w, b; x_i, y_i) \quad (6)$$

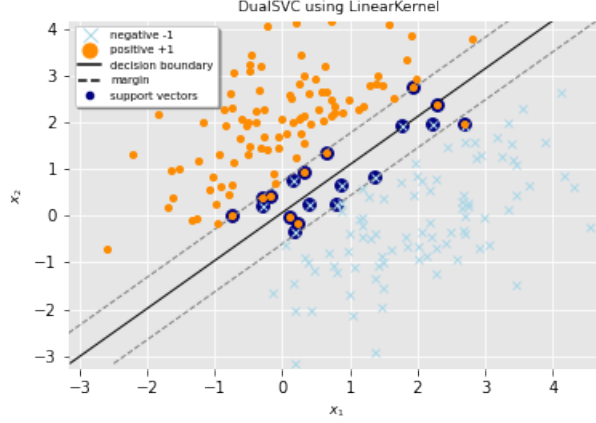


Figure 1: Linear SVC hyperplane

where  $\mathcal{R}(w, b)$  is the *regularization term* and  $\mathcal{L}(w, b; x_i, y_i)$  is the *loss function* associated with the observation  $(x_i, y_i)$ .

The quadratic optimization problem 5 can be equivalently formulated as:

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b)) \quad (7)$$

where we make use of the *hinge* loss defined as:

$$\mathcal{L}_1 = \begin{cases} 0 & \text{if } y(w^T x + b) \geq 1 \\ 1 - y(w^T x + b) & \text{otherwise} \end{cases} \quad (8)$$

or, equivalently:

$$\mathcal{L}_1 = \max(0, 1 - y(w^T x + b)) \quad (9)$$

The above formulation penalizes slacks  $\xi$  linearly and is called  $\mathcal{L}_1$ -SVC.

The *hinge* loss is a convex function and it is nondifferentiable due to its nonsmoothness in 1, but has a subgradient wrt  $w$  that is given by:

$$\frac{\partial \mathcal{L}_1}{\partial w} = \begin{cases} -yx & \text{if } y(w^T x + b) < 1 \\ 0 & \text{otherwise} \end{cases} \quad (10)$$

To simplify the notation and so also the design of the algorithms, the simplest approach to learn the bias term  $b$  is that of including that into the *regularization term*; so we can rewrite 7 and 38 as follows:

$$\min_{w, b} \frac{1}{2} (\|w\|^2 + b^2) + C \sum_{i=1}^n \mathcal{L}(w; x_i, y_i) \quad (11)$$

or, equivalently, by augmenting the weight vector  $w$  with the bias term  $b$  and each instance  $x_i$  with an additional dimension, i.e., with constant value equal to 1:

$$\begin{aligned} \min_w \quad & \frac{1}{2} \|\bar{w}\|^2 + C \sum_{i=1}^n \mathcal{L}(w; \bar{x}_i, y_i) \\ \text{where} \quad & \bar{w}^T = [w^T, b] \\ & \bar{x}_i^T = [x_i^T, 1] \end{aligned} \quad (12)$$

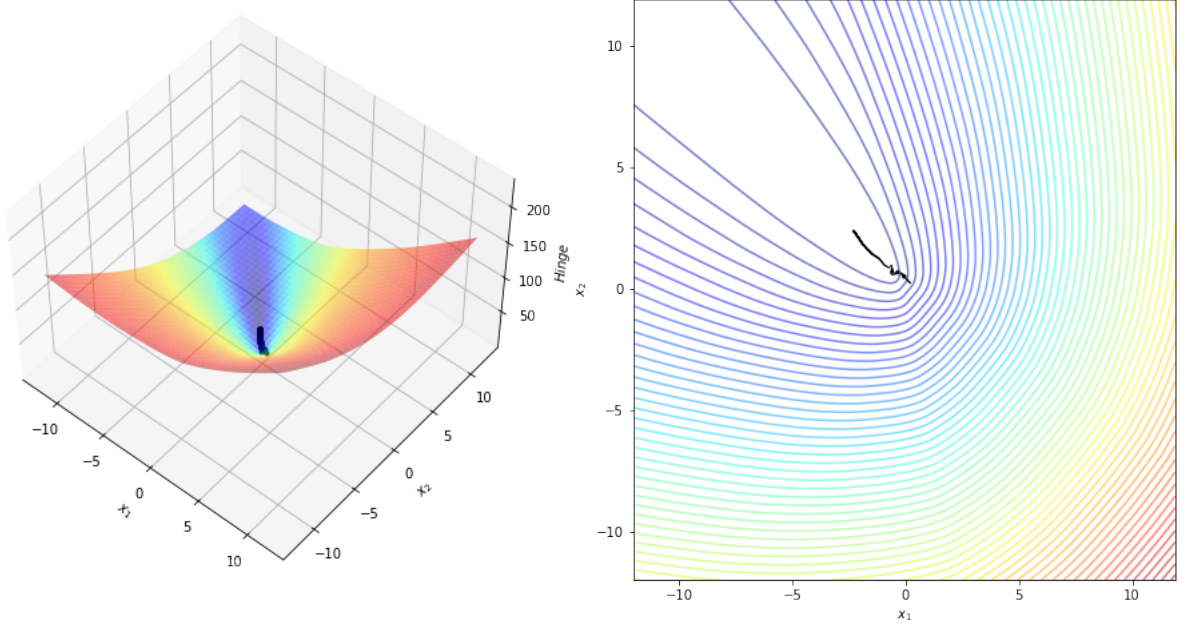


Figure 2: SVC Hinge loss with optimization steps

with the advantages of having convex properties of the objective function useful for convergence analysis and the possibility to directly apply algorithms designed for models without the bias term.

Notice that in terms of numerical optimization the formulations 7 and 38 are not equivalent to 11 or 12 since in the first one the bias term  $b$  does not contribute to the *regularization term*, so the SVM formulation is based on an unregularized bias term  $b$ , as highlighted by the *statistical learning theory*. But, in machine learning sense, numerical experiments in [6] show that the accuracy does not vary much when the bias term  $b$  is embedded into the weight vector  $w$ .

### 3.1.2 Wolfe Dual formulation

To reformulate the 5 as a *Wolfe dual*, we need to allocate the Lagrange multipliers  $\alpha_i \geq 0, \mu_i \geq 0 \forall_i$ :

$$\max_{\alpha, \mu} \min_{w, b, \xi} \mathcal{W}(w, b, \xi, \alpha, \mu) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i (w^T x_i + b) - 1 + \xi_i) - \sum_{i=1}^n \mu_i \xi_i \quad (13)$$

We wish to find the  $w$ ,  $b$  and  $\xi_i$  which minimizes, and the  $\alpha$  and  $\mu$  which maximizes  $\mathcal{W}$ , provided  $\alpha_i \geq 0, \mu_i \geq 0 \forall_i$ . We can do this by differentiating  $\mathcal{W}$  wrt  $w$  and  $b$  and setting the derivatives to 0:

$$\frac{\partial \mathcal{W}}{\partial w} = w - \sum_{i=1}^n \alpha_i y_i x_i \Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i \quad (14)$$

$$\frac{\partial \mathcal{W}}{\partial b} = - \sum_{i=1}^n \alpha_i y_i \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0 \quad (15)$$

$$\frac{\partial \mathcal{W}}{\partial \xi_i} = 0 \Rightarrow C = \alpha_i + \mu_i \quad (16)$$



Substituting 14 and 15 into 13 together with  $\mu_i \geq 0 \forall_i$ , which implies that  $\alpha \leq C$ , gives a new formulation being dependent on  $\alpha$ . We therefore need to find:

$$\begin{aligned} \max_{\alpha} \mathcal{W}(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i Q_{ij} \alpha_j \text{ where } Q_{ij} = y_i y_j \langle x_i, x_j \rangle \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \alpha^T Q \alpha \text{ subject to } 0 \leq \alpha_i \leq C \forall_i, \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned} \quad (17)$$

or, equivalently:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha + q^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C \forall_i \\ & y^T \alpha = 0 \end{aligned} \quad (18)$$

where  $q^T = [1, \dots, 1]$ .

By solving 18 we will know  $\alpha$  and, from 14, we will get  $w$ , so we need to calculate  $b$ .

We know that any data point satisfying 15 which is a support vector  $x_s$  will have the form:

$$y_s(w^T x_s + b) = 1 \quad (19)$$

and, by substituting in 14, we get:

$$y_s \left( \sum_{m \in S} \alpha_m y_m \langle x_m, x_s \rangle + b \right) = 1 \quad (20)$$

where  $s$  denotes the set of indices of the support vectors and is determined by finding the indices  $i$  where  $\alpha_i > 0$ , i.e., nonzero Lagrange multipliers.

Multiplying through by  $y_s$  and then using  $y_s^2 = 1$  from 2:

$$y_s^2 \left( \sum_{m \in S} \alpha_m y_m \langle x_m, x_s \rangle + b \right) = y_s \quad (21)$$

$$b = y_s - \sum_{m \in S} \alpha_m y_m \langle x_m, x_s \rangle \quad (22)$$

Instead of using an arbitrary support vector  $x_s$ , it is better to take an average over all of the support vectors in  $S$ :

$$b = \frac{1}{N_s} \sum_{s \in S} y_s - \sum_{m \in S} \alpha_m y_m \langle x_m, x_s \rangle \quad (23)$$

We now have the variables  $w$  and  $b$  that define our separating hyperplane's optimal orientation and hence our support vector machine. Each new point  $x'$  is classified by evaluating:

$$y' = \text{sgn} \left( \sum_{i=1}^n \alpha_i y_i \langle x_i, x' \rangle + b \right) \quad (24)$$

From 18 we can notice that the equality constraint  $y^T \alpha = 0$  arises from the stationarity condition  $\partial_b \mathcal{W} = 0$ . So, again, for simplicity, we can again consider the bias term  $b$  embedded into the weight vector. We report

below the box-constrained dual formulation [6] that arises from the primal 11 or 12 where the bias term  $b$  is embedded into the weight vector  $w$ :

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T (Q + yy^T) \alpha + q^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C \quad \forall_i \end{aligned} \quad (25)$$

### 3.1.3 Lagrangian Dual formulation

In order to relax the constraints in the *Wolfe dual* formulation 18 we define the problem as a *Lagrangian dual* relaxation by embedding them into objective function, so we need to allocate the Lagrangian multipliers  $\mu \geq 0, \lambda_+ \geq 0, \lambda_- \geq 0$ :

$$\begin{aligned} \max_{\mu, \lambda_+, \lambda_-} \min_{\alpha} \mathcal{L}(\alpha, \mu, \lambda_+, \lambda_-) &= \frac{1}{2} \alpha^T Q \alpha + q^T \alpha - \mu^T (y^T \alpha) - \lambda_+^T (u - \alpha) - \lambda_-^T \alpha \\ &= \frac{1}{2} \alpha^T Q \alpha + (q - \mu y + \lambda_+ - \lambda_-)^T \alpha - \lambda_+^T u \end{aligned} \quad (26)$$

where the upper bound  $u^T = [C, \dots, C]$ .

Taking the derivative of the Lagrangian  $\mathcal{L}$  wrt  $\alpha$  and settings it to 0 gives:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = 0 \Rightarrow Q\alpha + (q - \mu y + \lambda_+ - \lambda_-) = 0 \quad (27)$$

With  $\alpha$  optimal solution of the linear system:

$$Q\alpha = -(q - \mu y + \lambda_+ - \lambda_-) \quad (28)$$

the gradient wrt  $\mu, \lambda_+$  and  $\lambda_-$  are:

$$\frac{\partial \mathcal{L}}{\partial \mu} = -y\alpha \quad (29)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_+} = \alpha - u \quad (30)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_-} = -\alpha \quad (31)$$

If the Hessian matrix  $Q$  is indefinite, i.e., the Lagrangian function is not strictly convex since it will be linear along the eigenvectors correspondent to the null eigenvalues, the Lagrangian dual relaxation will be nondifferentiable, so it will have infinite solutions and for each of them it will have a different subgradient. In order to compute the gradient, we will choose  $\alpha$  in such a way as the one that minimizes the residue, i.e. the least-squares solution:

$$\begin{aligned} \min_{\alpha \in K_n(Q, b)} \quad & \|Q\alpha - b\| \\ \text{where} \quad & b = -(q - \mu y + \lambda_+ - \lambda_-) \end{aligned} \quad (32)$$

Since we are dealing with a symmetric but indefinite linear system we will choose a well-known Krylov method that performs the Lanczos iterate, i.e., symmetric Arnoldi iterate, called *minres*, i.e., symmetric *gmres*, which computes the vector  $\alpha$  that minimizes  $\|Q\alpha - b\|$  among all vectors in  $K_n(Q, b) = \text{span}(b, Qb, Q^2b, \dots, Q^{n-1}b)$ .

From 18 we can notice that the equality constraint  $y^T \alpha = 0$  arises from the stationarity condition  $\partial_b \mathcal{W} = 0$ . So, again, for simplicity, we can again consider the bias term  $b$  embedded into the weight vector. In this way the

dimensionality of 26 is reduced of 1/3 by removing the multipliers  $\mu$  which was allocated to control the equality constraint  $y^T \alpha = 0$ , so we will end up solving exactly the problem 25.

$$\begin{aligned} \max_{\lambda_+, \lambda_-} \min_{\alpha} \mathcal{L}(\alpha, \lambda_+, \lambda_-) &= \frac{1}{2} \alpha^T (Q + yy^T) \alpha + q^T \alpha - \lambda_+^T (u - \alpha) - \lambda_-^T \alpha \\ &= \frac{1}{2} \alpha^T (Q + yy^T) \alpha + (q + \lambda_+ - \lambda_-)^T \alpha - \lambda_+^T u \end{aligned} \quad (33)$$

where, again, the upper bound  $u^T = [C, \dots, C]$ .

Now, taking the derivative of the Lagrangian  $\mathcal{L}$  wrt  $\alpha$  and settings it to 0 gives:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = 0 \Rightarrow (Q + yy^T) \alpha + (q + \lambda_+ - \lambda_-) = 0 \quad (34)$$

With  $\alpha$  optimal solution of the linear system:

$$(Q + yy^T) \alpha = -(q + \lambda_+ - \lambda_-) \quad (35)$$

the gradient wrt  $\lambda_+$  and  $\lambda_-$  are:

$$\frac{\partial \mathcal{L}}{\partial \lambda_+} = \alpha - u \quad (36)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_-} = -\alpha \quad (37)$$

## 3.2 Squared Hinge loss

### 3.2.1 Primal formulation

Since smoothed versions of objective functions may be preferred for optimization, we can reformulate 7 as:

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b))^2 \quad (38)$$

where we make use of the *squared hinge* loss that quadratically penalized slacks  $\xi$  and is called  $\mathcal{L}_2$ -SVC.

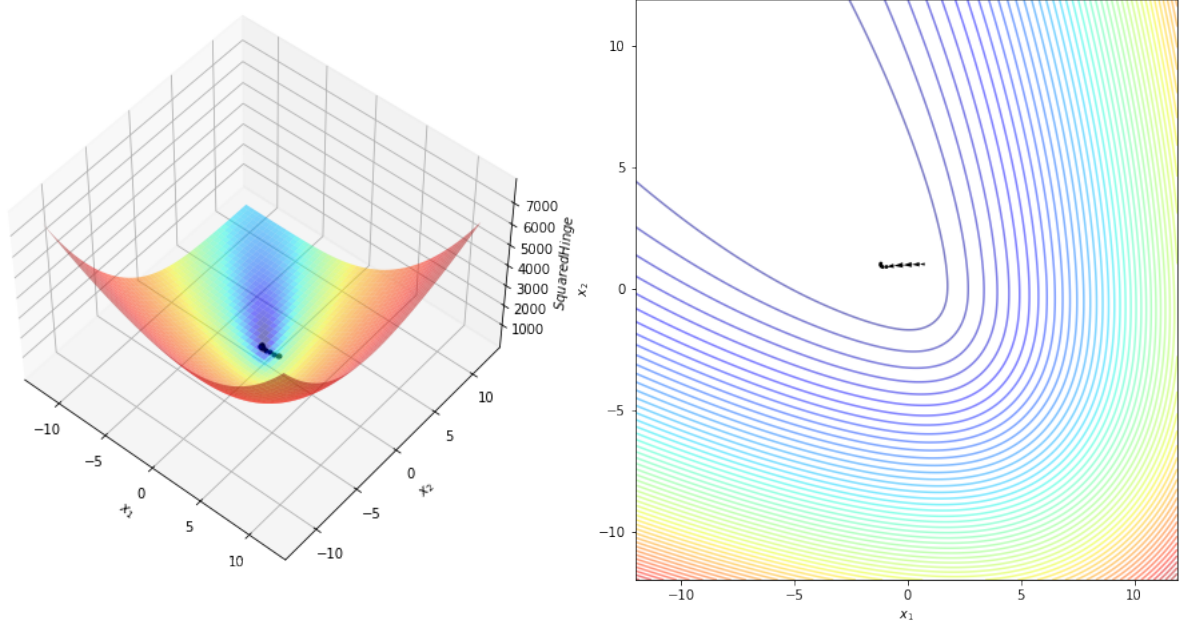


Figure 3: SVC Squared Hinge loss with optimization steps

## 4 Linear Support Vector Regression

In the case of regression the goal is to predict a real-valued output for  $y'$  so that our training data is of the form:

$$\{(x_i, y_i), x \in \mathbb{R}^m, y_i \in \mathbb{R}, i = 1, \dots, n\} \quad (39)$$

The regression SVM use a loss function that not allocating a penalty if the predicted value  $y'_i$  is less than a distance  $\epsilon$  away from the actual value  $y_i$ , i.e., if  $|y_i - y'_i| \leq \epsilon$ , where  $y'_i = w^T x_i + b$ . The region bound by  $y'_i \pm \epsilon \forall_i$  is called an  $\epsilon$ -insensitive tube. The output variables which are outside the tube are given one of two slack variable penalties depending on whether they lie above,  $\xi^+$ , or below,  $\xi^-$ , the tube, provided  $\xi^+ \geq 0$  and  $\xi^- \geq 0 \forall_i$ :

$$\begin{aligned} y_i &\leq y'_i + \epsilon + \xi^+ \forall_i \\ y_i &\geq y'_i - \epsilon - \xi^- \forall_i \\ \xi_i^+, \xi_i^- &\geq 0 \forall_i \end{aligned} \quad (40)$$

The objective function for SVR can then be written as:

$$\begin{aligned} \min_{w, b, \xi^+, \xi^-} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i^+ + \xi_i^-) \\ \text{subject to} \quad & y_i - w^T x_i - b \leq \epsilon + \xi_i^+ \forall_i \\ & w^T x_i + b - y_i \leq \epsilon + \xi_i^- \forall_i \\ & \xi_i^+, \xi_i^- \geq 0 \forall_i \end{aligned} \quad (41)$$



Figure 4: Linear SVR hyperplane

## 4.1 Epsilon-insensitive loss

### 4.1.1 Primal formulation

The general primal unconstrained formulation takes the same form of 6.

The quadratic optimization problem 41 can be equivalently formulated as:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, |y_i - (w^T x_i + b)| - \epsilon) \quad (42)$$

where we make use of the *epsilon-insensitive* loss defined as:

$$\mathcal{L}_\epsilon = \begin{cases} 0 & \text{if } |y - (w^T x + b)| \leq \epsilon \\ |y - (w^T x + b)| - \epsilon & \text{otherwise} \end{cases} \quad (43)$$

or, equivalently:

$$\mathcal{L}_\epsilon = \max(0, |y - (w^T x + b)| - \epsilon) \quad (44)$$

The above formulation penalizes slacks  $\xi$  linearly and is called  $\mathcal{L}_1$ -SVR.

As the *hinge* loss, also the *epsilon insensitive* loss is a convex function and it is nondifferentiable due to its nonsmoothness in  $\pm\epsilon$ , but has a subgradient wrt  $w$  that is given by:

$$\frac{\partial \mathcal{L}_\epsilon}{\partial w} = \begin{cases} (y - (w^T x + b))x & \text{if } |y - (w^T x + b)| > \epsilon \\ 0 & \text{otherwise} \end{cases} \quad (45)$$

### 4.1.2 Wolfe Dual formulation

To reformulate the 41 as a *Wolfe dual*, we introduce the Lagrange multipliers  $\alpha_i^+ \geq 0, \alpha_i^- \geq 0, \mu_i^+ \geq 0, \mu_i^- \geq 0 \forall_i$ :

$$\begin{aligned} \max_{\alpha^+, \alpha^-, \mu^+, \mu^-} \min_{w, b, \xi^+, \xi^-} \mathcal{W}(w, b, \xi^+, \xi^-, \alpha^+, \alpha^-, \mu^+, \mu^-) &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i^+ + \xi_i^-) - \sum_{i=1}^n (\mu_i^+ \xi_i^+ + \mu_i^- \xi_i^-) \\ &\quad - \sum_{i=1}^n \alpha_i^+ (\epsilon + \xi_i^+ + y'_i - y_i) - \sum_{i=1}^n \alpha_i^- (\epsilon + \xi_i^- - y'_i + y_i) \end{aligned} \quad (46)$$

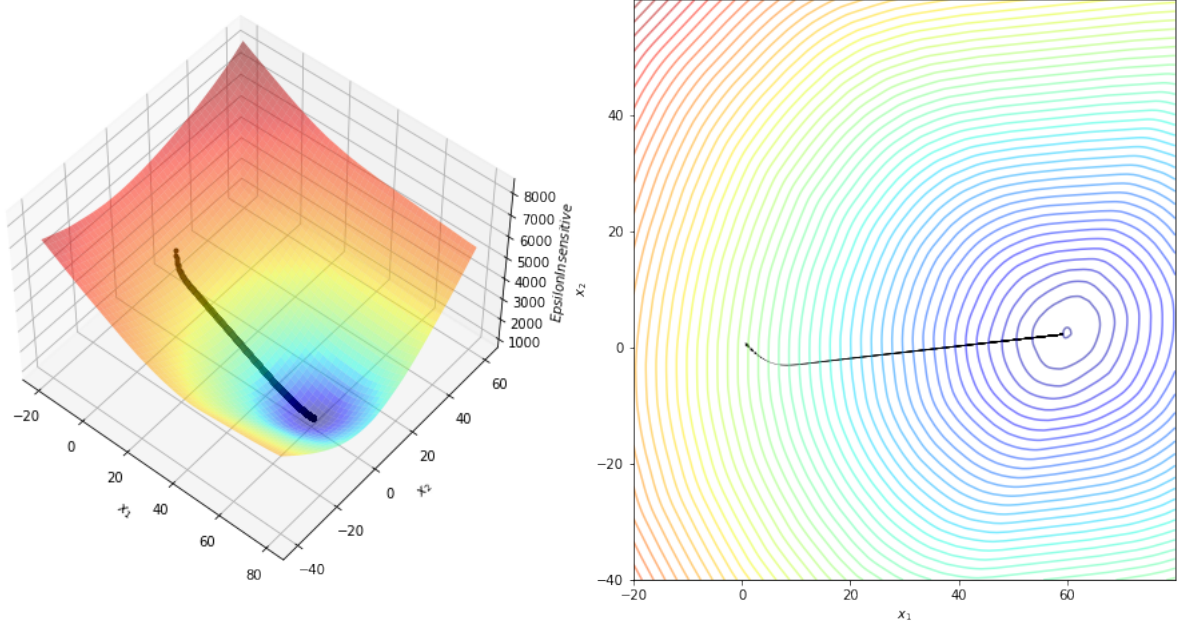


Figure 5: SVR Epsilon-insensitive loss with optimization steps

Substituting for  $y_i$ , differentiating wrt  $w, b, \xi^+, \xi^-$  and setting the derivatives to 0 gives:

$$\frac{\partial \mathcal{W}}{\partial w} = w - \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) x_i \Rightarrow w = \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) x_i \quad (47)$$

$$\frac{\partial \mathcal{W}}{\partial b} = - \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) \Rightarrow \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) = 0 \quad (48)$$

$$\frac{\partial \mathcal{W}}{\partial \xi_i^+} = 0 \Rightarrow C = \alpha_i^+ + \mu_i^+ \quad (49)$$

$$\frac{\partial \mathcal{W}}{\partial \xi_i^-} = 0 \Rightarrow C = \alpha_i^- + \mu_i^- \quad (50)$$

Substituting 47 and 48 in, we now need to maximize  $\mathcal{W}$  wrt  $\alpha_i^+$  and  $\alpha_i^-$ , where  $\alpha_i^+ \geq 0$ ,  $\alpha_i^- \geq 0 \forall i$ :

$$\max_{\alpha^+, \alpha^-} \mathcal{W}(\alpha^+, \alpha^-) = \sum_{i=1}^n y_i (\alpha_i^+ - \alpha_i^-) - \epsilon \sum_{i=1}^n (\alpha_i^+ + \alpha_i^-) - \frac{1}{2} \sum_{i,j} (\alpha_i^+ - \alpha_i^-) \langle x_i, x_j \rangle (\alpha_j^+ - \alpha_j^-) \quad (51)$$

Using  $\mu_i^+ \geq 0$  and  $\mu_i^- \geq 0$  together with 47 and 48 means that  $\alpha_i^+ \leq C$  and  $\alpha_i^- \leq C$ . We therefore need to find:

$$\begin{aligned} \min_{\alpha^+, \alpha^-} & \quad \frac{1}{2} (\alpha^+ - \alpha^-)^T K (\alpha^+ - \alpha^-) + \epsilon q^T (\alpha^+ + \alpha^-) - y^T (\alpha^+ - \alpha^-) \\ \text{subject to} & \quad 0 \leq \alpha_i^+, \alpha_i^- \leq C \forall i \\ & \quad q^T (\alpha^+ - \alpha^-) = 0 \end{aligned} \quad (52)$$

where  $q^T = [1, \dots, 1]$ .

We can write the 52 in a standard quadratic form as:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - q^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C \quad \forall_i \\ & e^T \alpha = 0 \end{aligned} \quad (53)$$

where the Hessian matrix  $Q$  is  $\begin{bmatrix} K & -K \\ -K & K \end{bmatrix}$ ,  $q$  is  $\begin{bmatrix} -y \\ y \end{bmatrix} + \epsilon$ , and  $e$  is  $\begin{bmatrix} 1 \\ -1 \end{bmatrix}$ .

Each new predictions  $y'$  can be found using:

$$y' = \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) \langle x_i, x' \rangle + b \quad (54)$$

A set  $S$  of support vectors  $x_s$  can be created by finding the indices  $i$  where  $0 \leq \alpha \leq C$  and  $\xi_i^+ = 0$  or  $\xi_i^- = 0$ . This gives us:

$$b = y_s - \epsilon - \sum_{m \in S} (\alpha_m^+ - \alpha_m^-) \langle x_m, x_s \rangle \quad (55)$$

As before it is better to average over all the indices  $i$  in  $S$ :

$$b = \frac{1}{N_s} \sum_{s \in S} y_s - \epsilon - \sum_{m \in S} (\alpha_m^+ - \alpha_m^-) \langle x_m, x_s \rangle \quad (56)$$

From 53 we can notice that the equality constraint  $e^T \alpha = 0$  arises from the stationarity condition  $\partial_b \mathcal{W} = 0$ . So, again, for simplicity, we can again consider the bias term  $b$  embedded into the weight vector. We report below the box-constrained dual formulation [6] that arises from the primal 11 or 12 where the bias term  $b$  is embedded into the weight vector  $w$ :

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T (Q + ee^T) \alpha + q^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C \quad \forall_i \end{aligned} \quad (57)$$

#### 4.1.3 Lagrangian Dual formulation

In order to relax the constraints in the *Wolfe dual* formulation 52 we define the problem as a *Lagrangian dual* relaxation by embedding them into objective function, so we need to allocate the Lagrangian multipliers  $\mu \geq 0, \lambda_+ \geq 0, \lambda_- \geq 0$ :

$$\begin{aligned} \max_{\mu, \lambda_+, \lambda_-} \min_{\alpha} \mathcal{L}(\alpha, \mu, \lambda_+, \lambda_-) &= \frac{1}{2} \alpha^T Q \alpha + q^T \alpha - \mu^T (e^T \alpha) - \lambda_+^T (u - \alpha) - \lambda_-^T \alpha \\ &= \frac{1}{2} \alpha^T Q \alpha + (q - \mu e + \lambda_+ - \lambda_-)^T \alpha - \lambda_+^T u \end{aligned} \quad (58)$$

where the upper bound  $u^T = [C, \dots, C]$ .

Taking the derivative of the Lagrangian  $\mathcal{L}$  wrt  $\alpha$  and settings it to 0 gives:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = 0 \Rightarrow Q \alpha + (q - \mu e + \lambda_+ - \lambda_-) = 0 \quad (59)$$

With  $\alpha$  optimal solution of the linear system:

$$Q \alpha = -(q - \mu e + \lambda_+ - \lambda_-) \quad (60)$$

the gradient wrt  $\mu$ ,  $\lambda_+$  and  $\lambda_-$  are:

$$\frac{\partial \mathcal{L}}{\partial \mu} = -e\alpha \quad (61)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_+} = \alpha - u \quad (62)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_-} = -\alpha \quad (63)$$

If the Hessian matrix  $Q$  is indefinite, i.e., the Lagrangian function is not strictly convex since it will be linear along the eigenvectors correspondent to the null eigenvalues, the Lagrangian dual relaxation will be nondifferentiable, so it will have infinite solutions and for each of them it will have a different subgradient. In order to compute the gradient, we will choose  $\alpha$  in such a way as the one that minimizes the residue, i.e. the least-squares solution:

$$\min_{\alpha \in K_n(Q, b)} \|Q\alpha - b\| \quad (64)$$

where  $b = -(q - \mu e + \lambda_+ - \lambda_-)$

Since we are dealing with a symmetric but indefinite linear system we will choose a well-known Krylov method that performs the Lanczos iterate, i.e., symmetric Arnoldi iterate, called *minres*, i.e., symmetric *gmres*, which computes the vector  $\alpha$  that minimizes  $\|Q\alpha - b\|$  among all vectors in  $K_n(Q, b) = \text{span}(b, Qb, Q^2b, \dots, Q^{n-1}b)$ .

From 53 we can notice that the equality constraint  $e^T \alpha = 0$  arises from the stationarity condition  $\partial_b \mathcal{W} = 0$ . So, again, for simplicity, we can again consider the bias term  $b$  embedded into the weight vector. In this way the dimensionality of 58 is reduced of 1/3 by removing the multipliers  $\mu$  which was allocated to control the equality constraint  $e^T \alpha = 0$ , so we will end up solving exactly the problem 57.

$$\begin{aligned} \max_{\lambda_+, \lambda_-} \min_{\alpha} \mathcal{L}(\alpha, \lambda_+, \lambda_-) &= \frac{1}{2} \alpha^T (Q + ee^T) \alpha + q^T \alpha - \lambda_+^T (u - \alpha) - \lambda_-^T \alpha \\ &= \frac{1}{2} \alpha^T (Q + ee^T) \alpha + (q + \lambda_+ - \lambda_-)^T \alpha - \lambda_+^T u \end{aligned} \quad (65)$$

where, again, the upper bound  $u^T = [C, \dots, C]$ .

Now, taking the derivative of the Lagrangian  $\mathcal{L}$  wrt  $\alpha$  and settings it to 0 gives:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = 0 \Rightarrow (Q + ee^T) \alpha + (q + \lambda_+ - \lambda_-) = 0 \quad (66)$$

With  $\alpha$  optimal solution of the linear system:

$$(Q + ee^T) \alpha = -(q + \lambda_+ - \lambda_-) \quad (67)$$

the gradient wrt  $\lambda_+$  and  $\lambda_-$  are:

$$\frac{\partial \mathcal{L}}{\partial \lambda_+} = \alpha - u \quad (68)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_-} = -\alpha \quad (69)$$



## 4.2 Squared Epsilon-insensitive loss

### 4.2.1 Primal formulation

To provide a continuously differentiable function the optimization problem 42 can be formulated as:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, |y_i - (w^T x_i + b)| - \epsilon)^2 \quad (70)$$

where we make use of the *squared epsilon-insensitive* loss that quadratically penalized slacks  $\xi$  and is called  $\mathcal{L}_2$ -SVR.

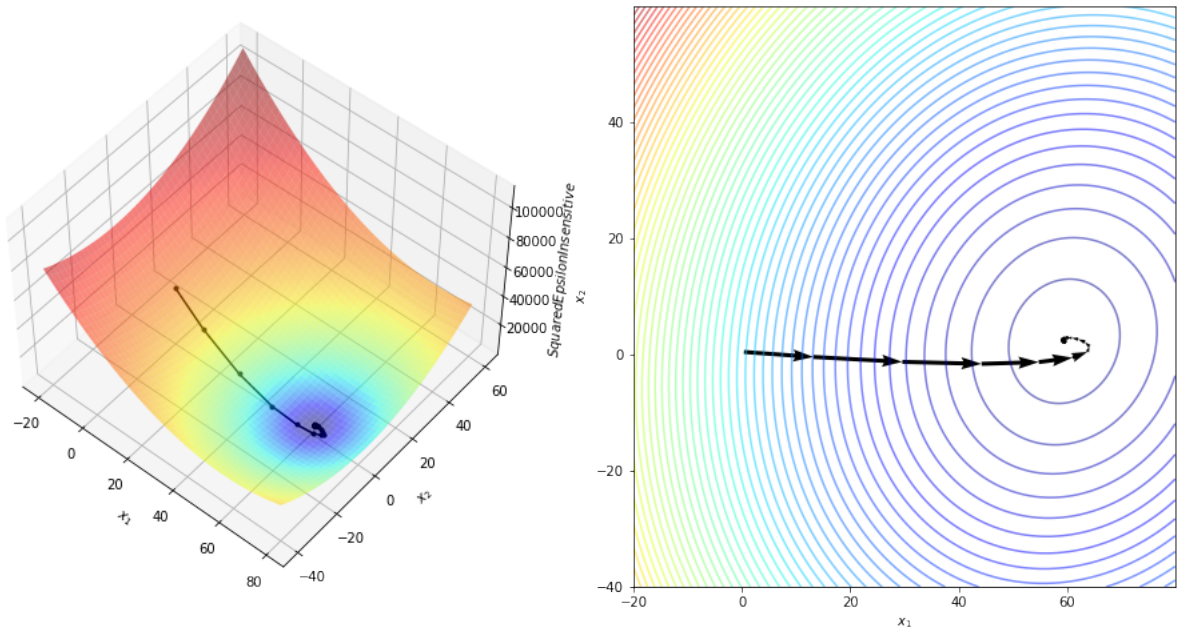


Figure 6: SVC Squared Epsilon-insensitive loss with optimization steps

## 5 Nonlinear Support Vector Machines

When applying our SVC to linearly separable data we have started by creating a matrix  $Q$  from the dot product of our input variables:

$$Q_{ij} = y_i y_j k(x_i, x_j) \quad (71)$$

or, a matrix  $K$  from in the SVR case:

$$K_{ij} = k(x_i, x_j) \quad (72)$$

where  $k(x_i, x_j)$  is an example of a family of functions called *kernel functions* and:

$$k(x_i, x_j) = \langle x_i, x_j \rangle = x_i^T x_j \quad (73)$$

is known as *linear* kernel.

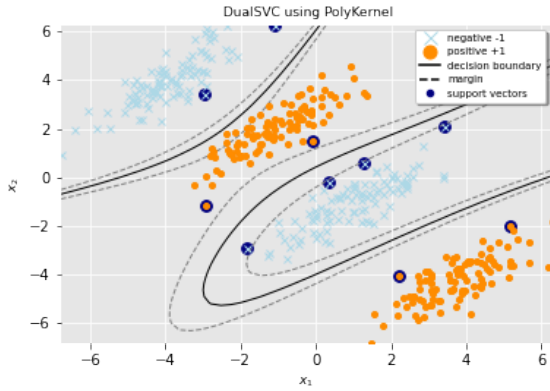
The reason that this *kernel trick* is useful is that there are many classification/regression problems that are nonlinearly separable/regressable in the *input space*, which might be in a higher dimensionality *feature space* given a suitable mapping  $x \rightarrow \phi(x)$ .

### 5.1 Polynomial kernel

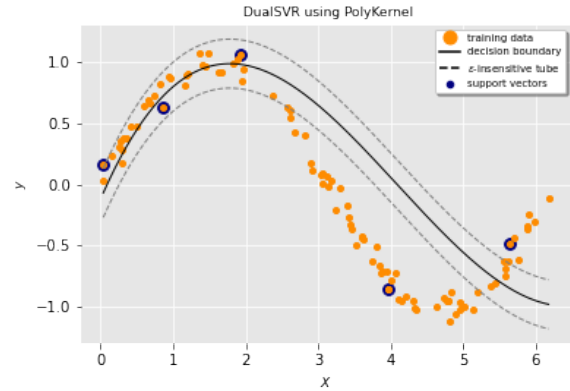
The *polynomial* kernel is defined as:

$$k(x_i, x_j) = (\gamma \langle x_i, x_j \rangle + r)^d \quad (74)$$

where  $\gamma$  define how far the influence of a single training example reaches (low values meaning ‘far’ and high values meaning ‘close’).



(a) Polynomial SVC hyperplane



(b) Polynomial SVR hyperplane

Figure 7: Polynomial SVM hyperplanes

### 5.2 Gaussian RBF kernel

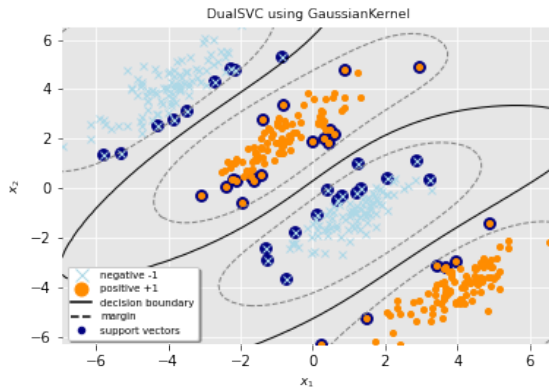
The *gaussian* kernel is defined as:

$$k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right) \quad (75)$$

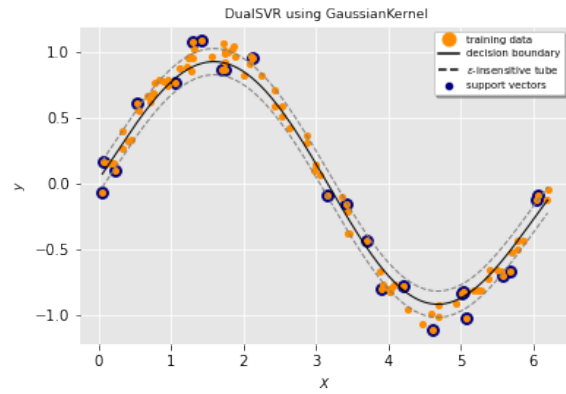
or, equivalently:

$$k(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2) \quad (76)$$

where  $\gamma = \frac{1}{2\sigma^2}$  define how far the influence of a single training example reaches (low values meaning ‘far’ and high values meaning ‘close’).



(a) Gaussian SVC hyperplane



(b) Gaussian SVR hyperplane

Figure 8: Gaussian SVM hyperplanes

## 6 Gradient Descent

### 6.1 Momentum

#### 6.1.1 Standard

---

**Algorithm 1** Standard Momentum Accelerated Gradient Descent. The learning rate  $\eta$ , the  $\alpha$  term and the maximum number of iterations are given.

---

**Require:** Learning rate  $\eta$  and momentum parameter  $\alpha$

**Require:** Maximum number of iteration and error threshold

```

1: procedure MOMENTUM DESCENT
2:   Initialize  $\mathbf{w}$  and  $\mathbf{v}$ 
3:    $k \leftarrow 0$ 
4:   while  $k < \text{max\_iterations} \ \&\& \ \text{error\_th} < e$  do
5:     if Nesterov Momentum then
6:        $\tilde{\mathbf{W}} \leftarrow \mathbf{w} + \alpha \mathbf{v}$ 
7:     end if
8:     Compute gradient estimate:  $\mathbf{g} \leftarrow \frac{1}{n} \nabla \sum_i L(\tilde{\mathbf{W}})$ 
9:     Compute velocity update:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \eta \mathbf{g}$ 
10:    Apply update:  $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v}$ 
11:  end while
12: end procedure
```

---

#### 6.1.2 Nesterov

---

**Algorithm 2** Nesterov Momentum Accelerated Gradient Descent. The learning rate  $\eta$ , the  $\alpha$  term and the maximum number of iterations are given.

---

**Require:** Learning rate  $\eta$  and momentum parameter  $\alpha$

**Require:** Maximum number of iteration and error threshold

```

1: procedure MOMENTUM DESCENT
2:   Initialize  $\mathbf{w}$  and  $\mathbf{v}$ 
3:    $k \leftarrow 0$ 
4:   while  $k < \text{max\_iterations} \ \&\& \ \text{error\_th} < e$  do
5:     if Nesterov Momentum then
6:        $\tilde{\mathbf{W}} \leftarrow \mathbf{w} + \alpha \mathbf{v}$ 
7:     end if
8:     Compute gradient estimate:  $\mathbf{g} \leftarrow \frac{1}{n} \nabla \sum_i L(\tilde{\mathbf{W}})$ 
9:     Compute velocity update:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \eta \mathbf{g}$ 
10:    Apply update:  $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v}$ 
11:  end while
12: end procedure
```

---

## 7 AdaGrad

Due to the nondifferentiability of the *hinge* loss, we might end up in a situation where some components of the gradient are very small and others large. So, given a learning rate, a standard gradient descent approach might end up in a situation where it decreases too quickly the small weights or too slowly the large ones.

*AdaGrad* [1] addresses this problem by introducing the aggregate of the squares of previously observed gradients to adjust the learning rate. This has two benefits: first, we no longer need to decide just when a gradient is large enough. Second, it scales automatically with the magnitude of the gradients. Coordinates that routinely correspond to large gradients are scaled down significantly, whereas others with small gradients receive a much more gentle treatment.

We use the variable  $s_t$  to accumulate past gradient variance as follows:

$$\begin{aligned} g_t &= \partial_{w_t} \mathcal{L}(y_t, f(x_t, w)) \\ s_t &= s_{t-1} + g_t^2 \\ w_{t+1} &= w_t - \frac{\eta}{\sqrt{s_t + \epsilon}} \cdot g_t \end{aligned} \tag{77}$$

where  $\epsilon$  is an additive constant that ensures that we do not divide by 0.

---

**Algorithm 3** AdaGrad Algorithm. The learning rate  $\eta$ , the  $\alpha$  term and the maximum number of iterations are given.

---

**Require:** Learning rate  $\eta$  and momentum parameter  $\alpha$

**Require:** Maximum number of iteration and error threshold

```

1: procedure MOMENTUM DESCENT
2:   Initialize  $\mathbf{w}$  and  $\mathbf{v}$ 
3:    $k \leftarrow 0$ 
4:   while  $k < \text{max\_iterations} \ \&\& \ \text{error\_th} < e$  do
5:     if Nesterov Momentum then
6:        $\tilde{\mathbf{W}} \leftarrow \mathbf{w} + \alpha \mathbf{v}$ 
7:     end if
8:     Compute gradient estimate:  $\mathbf{g} \leftarrow \frac{1}{n} \nabla \sum_i L(\tilde{\mathbf{W}})$ 
9:     Compute velocity update:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \eta \mathbf{g}$ 
10:    Apply update:  $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v}$ 
11:  end while
12: end procedure
```

---

## 8 Sequential Minimal Optimization

The *Sequential Minimal Optimization (SMO)* [2] method is the most popular approach for solving the SVM QP problem without any extra  $Q$  matrix storage required by common QP methods. The advantage of SMO lies in the fact that it performs a series of two-point optimizations since we deal with just one equality constraint, i.e.,  $y^T \alpha = 0$ , so the Lagrange multipliers can be solved analytically.

At each iteration, SMO chooses two  $\alpha_i$  to jointly optimize, let  $\alpha_1$  and  $\alpha_2$ , finds the optimal values for these multipliers and update the SVM to reflect these new values. In order to solve for two Lagrange multipliers, SMO first computes the constraints over these and then solves for the constrained minimum. Since there are only two multipliers, the bound constraints cause the Lagrange multipliers to lie within a box, while the linear equality constraint causes the Lagrange multipliers to lie on a diagonal line inside the box. So, the constrained minimum must lie there.

### 8.1 Classification

The ends of the diagonal line segment in terms of  $\alpha_2$  can be expressed as follow if the target  $y_1 \neq y_2$ :

$$\begin{aligned} L &= \max(0, \alpha_2 - \alpha_1) \\ H &= \min(C, C + \alpha_2 - \alpha_1) \end{aligned} \quad (78)$$

or, alternatively, if the target  $y_1 = y_2$ :

$$\begin{aligned} L &= \max(0, \alpha_2 + \alpha_1 - C) \\ H &= \min(C, \alpha_2 + \alpha_1) \end{aligned} \quad (79)$$

The second derivative of the objective quadratic function along the diagonal line can be expressed as:

$$\eta = K(x_1, x_1) + K(x_2, x_2) - 2K(x_1, x_2) \quad (80)$$

that will be grather than zero if the kernel matrix will be positive definite, so there will be a minimum along the linear equality constraints that will be:

$$\alpha_2^{new} = \alpha_2 + \frac{y_2(E_1 - E_2)}{\eta} \quad (81)$$

where  $E_i = u_i - y_i$  is the error on the  $i$ -th training example and  $u_i$  is the output of the SVM for the same.

Then, the box-constrained minimum is found by clipping the unconstrained minimum to the ends of the line segment:

$$\alpha_2^{new, clipped} = \begin{cases} H & \text{if } \alpha_2^{new} \geq H \\ \alpha_2^{new} & \text{if } L < \alpha_2^{new} < H \\ L & \text{if } \alpha_2^{new} \leq L \end{cases} \quad (82)$$

Finally, the value of  $\alpha_1$  is computed from the new clipped  $\alpha_2$  as:

$$\alpha_1^{new} = \alpha_1 + s(\alpha_2 - \alpha_2^{new, clipped}) \quad (83)$$

where  $s = y_1 y_2$ .

Since the *Karush-Kuhn-Tucker* conditions are necessary and sufficient conditions for optimality of a positive definite QP problem and the KKT conditions for the problem 18 are:

$$\begin{aligned} \alpha_i &= 0 \Leftrightarrow y_i u_i \geq 1 \\ 0 < \alpha_i < C &\Leftrightarrow y_i u_i = 1 \\ \alpha_i &= C \Leftrightarrow y_i u_i \leq 1 \end{aligned} \quad (84)$$

the steps described above will be iterate as long as there will be an example that violates these KKT conditions.

## 8.2 Regression

## 9 Experiments

The following experiments refer to 3-fold cross-validation over *linearly* and *nonlinearly* separable generated datasets of size 100, so the reported results are to considered as a mean over the 3 folds.

### 9.1 Support Vector Classifier

Below experiments are about the SVC for which I tested different values for the  $C$  complexity hyperparameter, i.e., from *soft* to *hard margin*, and in case of nonlinearly separable data also different *kernel functions* mentioned above.

#### 9.1.1 Hinge loss

Table 1: SVC Primal formulation results with Hinge loss

		fit_time	train_accuracy	val_accuracy	n_iter	train_n_sv	val_n_sv
solver	C						
adagrad	1	0.045228	0.967475	0.969848	66	20	10
	10	0.247653	0.972468	0.974898	398	12	7
	100	0.301816	0.972468	0.974898	593	10	7
liblinear	1	0.001328	0.982494	0.980024	153	12	5
	10	0.001203	0.987487	0.989974	633	7	4
	100	0.001864	0.985000	0.984999	831	7	3

Table 2: Linear SVC Wolfe Dual formulation results with Hinge loss

		fit_time	train_accuracy	val_accuracy	n_iter	train_n_sv	val_n_sv
solver	C						
cvxopt	1	0.015791	0.972487	0.959897	1	15	15
	10	0.025156	0.972487	0.959897	1	11	11
	100	0.028108	0.972487	0.959897	1	38	38
smo	1	0.042172	0.972487	0.959897	23	15	15
	10	0.114690	0.972487	0.959897	76	11	11
	100	0.281406	0.972487	0.959897	411	10	10
libsvm	1	0.002600	0.987506	0.985075	209	12	12
	10	0.004442	0.990012	0.980100	213	8	8
	100	0.002549	0.990012	0.980100	757	7	7

Table 3: Linear SVC Lagrangian Dual formulation results with Hinge loss

		fit_time	train_accuracy	val_accuracy	n_iter	train_n_sv	val_n_sv
dual	C						
bcqp	1	0.006401	0.992519	0.99005	1	130	130
	10	0.006864	0.992519	0.99005	1	130	130
	100	0.005661	0.992519	0.99005	1	130	130
qp	1	0.005798	0.995006	0.99005	1	132	132
	10	0.006755	0.995006	0.99005	1	132	132
	100	0.006650	0.995006	0.99005	1	132	132



Table 4: Nonlinear SVC Wolfe Dual formulation results with Hinge loss

solver	kernel	C	fit_time	train_accuracy	val_accuracy	n_iter	train_n_sv	val_n_sv
cvxopt	poly	1	0.128965	0.852492	0.728519	1	30	30
		10	0.117962	0.923686	0.760857	1	10	10
		100	0.090554	0.956197	0.835615	1	8	8
	rbf	1	0.132866	1.000000	0.997494	1	43	43
		10	0.123202	1.000000	0.997494	1	15	15
		100	0.094838	1.000000	0.994987	1	12	12
libsvm	poly	1	0.009136	1.000000	0.992481	287	30	30
		10	0.005259	1.000000	0.984962	504	12	12
		100	0.003439	1.000000	0.974937	460	8	8
	rbf	1	0.005949	1.000000	0.997512	121	46	46
		10	0.003889	1.000000	1.000000	214	15	15
		100	0.002700	1.000000	1.000000	105	11	11
smo	poly	1	0.452768	0.851243	0.730988	114	29	29
		10	0.383016	0.922437	0.763364	71	10	10
		100	0.343767	0.947429	0.803296	132	9	9
	rbf	1	0.287587	1.000000	0.997494	41	42	42
		10	0.259508	1.000000	0.994987	43	13	13
		100	0.319339	1.000000	0.994987	101	12	12

Table 5: Nonlinear SVC Lagrangian Dual formulation results with Hinge loss

dual	kernel	C	fit_time	train_accuracy	val_accuracy	n_iter	train_n_sv	val_n_sv
bcqp	poly	1	0.070312	0.806196	0.553791	4	211	211
		10	0.051106	0.806196	0.553791	4	211	211
		100	0.033508	0.806196	0.553791	4	211	211
	rbf	1	0.057051	1.000000	0.994987	1	246	246
		10	0.024727	1.000000	0.994987	1	246	246
		100	0.023880	1.000000	0.994987	1	246	246
qp	poly	1	0.601361	0.804934	0.551378	102	178	178
		10	0.598967	0.804934	0.551378	102	178	178
		100	0.455699	0.804934	0.551378	102	178	178
	rbf	1	1.002897	0.776332	0.600363	126	151	151
		10	0.731749	0.891207	0.745670	155	197	197
		100	0.297122	0.860001	0.728089	49	160	160

Table 6: SVC Primal formulation results with Squared Hinge loss

solver	C	momentum	fit_time	train_accuracy	val_accuracy	n_iter	train_n_sv	val_n_sv
gd	1	nesterov	0.071938	0.985000	0.985075	87	27	14
		none	0.140007	0.982512	0.975049	158	26	14
		standard	0.081723	0.982512	0.980100	91	26	15
	10	nesterov	0.026135	0.982512	0.985075	24	15	11
		none	0.049573	0.985000	0.985075	48	16	11
		standard	0.034831	0.982512	0.990050	29	15	9
	100	nesterov	0.008988	0.982531	0.975124	10	4	2
		none	0.017715	0.987525	0.980100	17	8	4
		standard	0.009709	0.982531	0.975124	10	4	1
liblinear	1	-	0.001800	0.982512	0.974974	314	17	9
	10	-	0.002610	0.982512	0.969998	1001	12	5
	100	-	0.002244	0.977500	0.969998	1001	10	3

Table 7: SVR Primal formulation results with Epsilon-insensitive loss

solver	C	epsilon	fit_time	train_r2	val_r2	n_iter	train_n_sv	val_n_sv
adagrad	1	0.1	0.692084	0.919185	0.915674	876	66	33
		0.2	0.701444	0.920036	0.916541	895	66	33
		0.3	0.684244	0.920119	0.916694	890	65	33
	10	0.1	2.629582	0.977826	0.972929	3560	65	32
		0.2	2.797290	0.977797	0.972832	3526	65	32
		0.3	2.594063	0.977778	0.972872	3446	65	32
	100	0.1	2.818639	0.978125	0.974200	4000	66	33
		0.2	2.674827	0.978118	0.974194	4000	66	32
		0.3	2.300388	0.978119	0.974199	4000	66	32
liblinear	1	0.1	0.000605	0.918867	0.916885	13	65	33
		0.2	0.000638	0.918866	0.916725	13	65	32
		0.3	0.000583	0.919243	0.916993	13	65	32
	10	0.1	0.000654	0.977853	0.972064	156	65	33
		0.2	0.000704	0.977852	0.972041	164	65	33
		0.3	0.000654	0.977870	0.972142	107	64	33
	100	0.1	0.000954	0.977723	0.974270	681	65	33
		0.2	0.001165	0.977642	0.973861	710	65	33
		0.3	0.001070	0.977654	0.974052	809	65	33

### 9.1.2 Squared Hinge loss

## 9.2 Support Vector Regression

### 9.2.1 Epsilon-insensitive loss

solver	kernel	C	epsilon	fit_time	train_r2	val_r2	n_iter	train_n_sv	val_n_sv
cvxopt	poly	1	0.1	0.012785	0.414638	-26.864749	-	18	18
libsvm	poly	1	0.1	0.046200	0.979295	-8.978090	134487	21	21
smo	poly	1	0.1	52.024504	0.309358	-30.592107	101185	18	18
cvxopt	rbf	1	0.1	0.011190	0.971967	-1.895743	-	12	12
libsvm	rbf	1	0.1	0.001463	0.985722	-0.927310	136	17	17
smo	rbf	1	0.1	0.030710	0.975638	-2.110523	27	11	11
cvxopt	poly	10	0.1	0.010302	0.756741	-37.096432	-	19	19
libsvm	poly	10	0.1	0.415621	0.979466	-9.594846	1475764	20	20
smo	poly	10	0.1	732.300189	0.506437	-42.271049	1596205	18	18
cvxopt	rbf	10	0.1	0.008236	0.978627	-1.546229	-	9	9
libsvm	rbf	10	0.1	0.003731	0.986915	-0.380431	631	14	14
smo	rbf	10	0.1	0.084678	0.975736	-1.668655	127	9	9

Table 8: Linear SVR Wolfe Dual formulation results with Epsilon-insensitive loss

			fit_time	train_r2	val_r2	n_iter	train_n_sv	val_n_sv
solver	C	epsilon						
cvxopt	1	0.1	0.046771	0.917772	0.914479	1	67	67
		0.2	0.043559	0.918341	0.915058	1	67	67
		0.3	0.043129	0.918942	0.915614	1	66	66
	10	0.1	0.031975	0.977920	0.972466	1	67	67
		0.2	0.020005	0.977926	0.972474	1	67	67
		0.3	0.015315	0.977954	0.972562	1	66	66
	100	0.1	0.013129	0.977788	0.974150	1	67	67
		0.2	0.013004	0.977742	0.974033	1	67	67
		0.3	0.012722	0.977737	0.973956	1	67	67
smo	1	0.1	0.053333	0.917773	0.914442	15	66	66
		0.2	0.061615	0.918341	0.915019	13	66	66
		0.3	0.085741	0.918942	0.915576	60	66	66
	10	0.1	0.114405	0.977920	0.972445	56	66	66
		0.2	0.194250	0.977926	0.972457	219	65	65
		0.3	0.072651	0.977953	0.972544	38	65	65
	100	0.1	0.610673	0.977788	0.974139	1508	66	66
		0.2	0.297791	0.977742	0.974022	394	66	66
		0.3	0.487787	0.977737	0.973939	900	66	66
libsvm	1	0.1	0.005714	0.917627	0.915448	64	66	66
		0.2	0.008318	0.918194	0.915985	103	66	66
		0.3	0.005886	0.918786	0.916554	55	66	66
	10	0.1	0.005883	0.977852	0.972051	283	66	66
		0.2	0.005896	0.977851	0.972025	194	65	65
		0.3	0.006089	0.977870	0.972135	594	65	65
	100	0.1	0.010685	0.977723	0.974270	2622	66	66
		0.2	0.014737	0.977673	0.974122	2710	66	66
		0.3	0.010669	0.977655	0.974045	4142	66	66

### 9.2.2 Squared Epsilon-insensitive loss

Table 9: Linear SVR Lagrangian Dual formulation results with Epsilon-insensitive loss

			fit_time	train_r2	val_r2	n_iter	train_n_sv	val_n_sv
dual	C	epsilon						
bcqp	1	0.1	0.925850	0.731073	0.721200	522	67	67
		0.2	0.779710	0.731073	0.721199	524	67	67
		0.3	0.745185	0.731073	0.721199	526	67	67
	10	0.1	0.812003	0.733638	0.723925	539	67	67
		0.2	0.742510	0.733638	0.723924	541	67	67
		0.3	0.780790	0.733638	0.723924	543	67	67
	100	0.1	0.766943	0.733638	0.723925	539	67	67
		0.2	0.695929	0.733638	0.723924	541	67	67
		0.3	0.506877	0.733638	0.723924	543	67	67
qp	1	0.1	1.054467	0.876534	0.870926	653	67	67
		0.2	1.103491	0.876534	0.870927	653	67	67
		0.3	0.868793	0.876534	0.870927	653	67	67
	10	0.1	0.667056	0.731825	0.722021	519	67	67
		0.2	0.649676	0.731825	0.722021	524	67	67
		0.3	0.629018	0.731825	0.722020	530	67	67
	100	0.1	0.634145	0.731825	0.722021	519	67	67
		0.2	0.614166	0.731825	0.722021	524	67	67
		0.3	0.673949	0.731825	0.722020	530	67	67

## 10 Conclusions

For what about the SVM formulations, it is known, in general, that the *primal* formulation, is suitable for large linear training since the complexity of the model grows with the number of features or, more in general, when the number of examples  $n$  is much larger than the number of features  $m$ ,  $n \gg m$ ; meanwhile the *dual* formulation, is more suitable in case the number of examples  $n$  is less than the number of features  $m$ ,  $n \ll m$ , since the complexity of the model is dominated by the number of examples.

From all these experiments we can see as, for what about the *primal* formulations, the results provided from the *custom* implementations are strongly similar to those of *sklearn* implementations, i.e., *liblinear* implementations, with a slight exception about the time gap obviously due to the different core implementation languages, Python and C respectively.

Meanwhile, for what about the *dual* formulations we can notice as *cvxopt* underperforms the *sklearn* implementations, i.e., *libsvm* implementations, in terms of time since it is a general-purpose QP solver and it does not exploit the structure of the problem, as SMO does. Despite this, the *custom* implementations does not overperform the *cvxopt* probably due to the gap generated from the different core implementation languages, again Python and C respectively. For these reasons, *sklearn* provides better results in terms of time wrt the other implementations since it is designed to work in a large-scale context and its core is implemented in C. Furthermore, in the SVC example with the polynomial kernel of degree 5, we can see that the time gap is significantly, properly two different orders of magnitude ( $\simeq 29\text{min}$  vs.  $\simeq 19\text{ms}$ ), and this could not depend just only by the different implementation languages; it's probable that *liblinear* adopts some heuristics, i.e., low rank approximations of the kernel matrix, to deal with the polynomial kernel in case of high degree.

Important consideration involves the number of support vector machines: the *Lagrangian dual* formulation tends to select all the data points as support vectors, so it makes the model complex and it tends to give low scores wrt the equivalent *Wolfe dual* formulation. In particular, the *Lagrangian relaxation* resulting from the *Wolfe dual* always gives rise to a nonsmooth optimization with an exception for the SVC with a Gaussian kernel where the two formulations solve exactly the same problem. In all the other cases the goodness of the solution depends on the residue in the solution of the *Lagrangian dual* at each step; one of the worst results certainly concerns the SVC with the polynomial kernel of degree 3, where the residue is in the order of  $+02/03$  and so

Table 10: Nonlinear SVR Lagrangian Dual formulation results with Epsilon-insensitive loss

dual	kernel	C	epsilon	fit_time	train_r2	val_r2	n_iter	train_n_sv	val_n_sv
bcqp	poly	1	0.1	0.081454	5.144076e-01	-9.076995e+00	35	67	67
			0.2	0.060330	5.079650e-01	-5.349260e+00	40	67	67
			0.3	0.088997	4.488241e-01	-4.556107e+00	64	67	67
		10	0.1	0.056643	5.144076e-01	-9.076995e+00	35	67	67
			0.2	0.076490	5.079650e-01	-5.349260e+00	40	67	67
			0.3	0.097213	4.488241e-01	-4.556107e+00	64	67	67
		100	0.1	0.060653	5.144076e-01	-9.076995e+00	35	67	67
			0.2	0.063831	5.079650e-01	-5.349260e+00	40	67	67
			0.3	0.105566	4.488241e-01	-4.556107e+00	64	67	67
	rbf	1	0.1	0.077817	7.396010e-01	-1.390908e+00	33	67	67
			0.2	0.185861	7.395830e-01	-1.392041e+00	90	67	67
			0.3	0.466749	5.929126e-01	-2.757236e+00	196	67	67
		10	0.1	0.108501	7.396010e-01	-1.390908e+00	33	67	67
			0.2	0.193361	7.395830e-01	-1.392041e+00	90	67	67
			0.3	0.457909	5.929126e-01	-2.757236e+00	196	67	67
		100	0.1	0.076684	7.396010e-01	-1.390908e+00	33	67	67
			0.2	0.211798	7.395830e-01	-1.392041e+00	90	67	67
			0.3	0.254640	5.929126e-01	-2.757236e+00	196	67	67
qp	poly	1	0.1	0.135663	4.113756e-01	-1.020445e+01	63	66	66
			0.2	0.528343	3.492985e-01	-7.300395e+00	347	64	64
			0.3	0.549913	-1.552200e+15	-1.787602e+12	380	46	46
		10	0.1	0.102671	4.113756e-01	-1.020445e+01	63	66	66
			0.2	0.550925	3.492985e-01	-7.300395e+00	347	64	64
			0.3	0.569163	-1.552200e+15	-1.787602e+12	380	46	46
		100	0.1	0.098125	4.113756e-01	-1.020445e+01	63	66	66
			0.2	0.407858	3.492985e-01	-7.300395e+00	347	64	64
			0.3	0.401141	-1.552200e+15	-1.787602e+12	380	46	46
	rbf	1	0.1	0.338488	6.913982e-01	-1.635557e+00	132	67	67
			0.2	0.442192	6.866212e-01	-1.660805e+00	184	67	67
			0.3	0.636889	6.115361e-01	-2.297081e+00	257	67	67
		10	0.1	0.128231	7.148036e-01	-1.440764e+00	54	67	67
			0.2	0.276627	7.066835e-01	-1.470430e+00	106	67	67
			0.3	0.336929	6.456839e-01	-1.981865e+00	137	67	67
		100	0.1	0.177903	7.148036e-01	-1.440764e+00	54	67	67
			0.2	0.272468	7.066835e-01	-1.470430e+00	106	67	67
			0.3	0.305527	6.456839e-01	-1.981865e+00	137	67	67

the approximation is horrible. Finally, we can see as fitting the intercept in an explicit way, i.e., by adding Lagrange multipliers to control the equality constraint, always get lower scores wrt the *Lagrangian relaxation* of the same problem with the bias term embedded into the weight matrix.

Table 11: SVR Primal formulation results with Squared Epsilon-insensitive loss

solver	C	momentum	epsilon	fit_time	train_r2	val_r2	n_iter	train_n_sv	val_n_sv
gd	1	nesterov	0.1	0.258087	0.978130	0.973982	183	66	32
			0.2	0.175323	0.978129	0.973978	181	66	32
			0.3	0.181325	0.978129	0.973978	179	66	32
		none	0.1	0.429835	0.978126	0.973975	352	66	32
			0.2	0.279171	0.978125	0.973974	349	66	32
			0.3	0.326429	0.978125	0.973972	346	66	32
		standard	0.1	0.312475	0.978130	0.973982	180	66	32
			0.2	0.127938	0.978130	0.973979	178	66	32
			0.3	0.184472	0.978129	0.973977	175	66	32
	10	nesterov	0.1	0.019044	0.978184	0.973958	26	66	33
			0.2	0.016165	0.978184	0.973958	25	66	33
			0.3	0.015675	0.978184	0.973958	25	66	33
		none	0.1	0.033184	0.978184	0.973958	48	66	33
			0.2	0.032136	0.978184	0.973957	46	66	33
			0.3	0.026559	0.978183	0.973955	45	66	32
		standard	0.1	0.017513	0.977871	0.975106	25	65	33
			0.2	0.016943	0.977874	0.975095	25	65	33
			0.3	0.017088	0.977872	0.975098	25	65	33
	100	nesterov	0.1	0.004475	-1629.737655	-1601.158587	6	67	33
			0.2	0.005785	-1637.506672	-1609.161144	6	67	33
			0.3	0.004273	-1638.224513	-1609.008403	6	67	33
		none	0.1	0.004214	-18.594624	-18.441771	6	67	33
			0.2	0.004430	-18.653692	-18.498210	6	67	33
			0.3	0.004730	-18.969060	-18.800697	6	67	33
		standard	0.1	0.021219	0.978184	0.973963	29	66	33
			0.2	0.019697	0.978184	0.973963	29	66	33
			0.3	0.018692	0.978184	0.973968	29	66	33
liblinear	1	-	0.1	0.002499	0.978134	0.974000	90	67	32
			0.2	0.002433	0.978132	0.974007	89	66	32
			0.3	0.002493	0.978130	0.974013	87	66	32
	10	-	0.1	0.009630	0.978183	0.973959	784	66	33
			0.2	0.009278	0.978183	0.973968	780	66	33
			0.3	0.008759	0.978183	0.973983	759	66	32
	100	-	0.1	0.011892	0.977919	0.972826	1001	67	33
			0.2	0.012556	0.978147	0.973841	1001	66	32
			0.3	0.011528	0.978086	0.975133	1001	66	32

## References

- [1] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- [2] John Platt. Sequential minimal optimization: A fast algorithm for training support vector machines. 1998.
- [3] S. Sathya Keerthi, Shirish Krishnaji Shevade, Chiranjib Bhattacharyya, and Karuturi Radha Krishna Murthy. Improvements to Platt’s SMO algorithm for SVM classifier design. *Neural computation*, 13(3):637–649, 2001.

- [4] Gary William Flake and Steve Lawrence. Efficient SVM regression training with SMO. *Machine Learning*, 46(1):271–290, 2002.
- [5] SK Shevade, SS Keerthi, C Bhattacharyya, and KRK Murthy. Improvements to SMO algorithm for SVM regression (Tech. Rep. No. CD-99-16). *Singapore: Control Division Department of Mechanical and Production Engineering*, 1999.
- [6] Chih-Wei Hsu and Chih-Jen Lin. A simple decomposition method for support vector machines. *Machine Learning*, 46(1):291–314, 2002.