

Anomaly Detection

Hyperplane-based methods

Paul Irofti
Cristian Rusu
Andrei Pătrașcu

Computer Science Department
University of Bucharest

- One class classification: OC-SVM, SVDD. Algorithms.
- Robust versions. Algorithms.



Consider the training data is provided:

$$x_1, x_2, \dots, x_m \in X.$$

where m is the number of the observations and X some space (i.e. compact subset of \mathbb{R}^n).

Question

What is a "good" binary function f that captures the "region" of the most of datapoints where returns $+1$, and -1 elsewhere ?



Consider the training data is provided:

$$x_1, x_2, \dots, x_m \subset X.$$

where m is the number of the observations and X some space (i.e. compact subset of \mathbb{R}^n).

A simple answer: convex bodies such as

- hyperplanes
- spheres
- ellipsoids

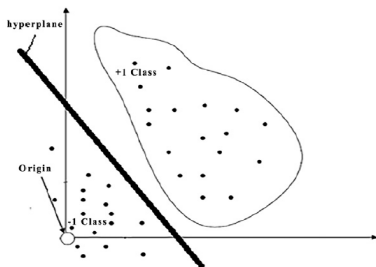


One class classification - hyperplane

Consider the training data is provided:

$$x_1, x_2, \dots, x_m \subset X.$$

where m is the number of the observations and X some space (e.g. compact subset of \mathbb{R}^n).



4. Hyperplane in one-class support vector machine.

Idea: *The inliers are grouped in a region far from the origin.*



One class classification - hyperplane

Consider the training data is provided:

$$x_1, x_2, \dots, x_m \in X.$$

where m is the number of the observations and X some space (e.g. compact subset of \mathbb{R}^n).

Prior assumptions:

- the inliers are distributed in a region separable from the origin by a hyperplane
- the outliers lies near the origin

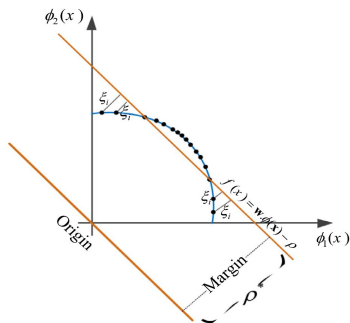
Take the decision function as $f(x) = \text{sgn}(w^T x - \rho)$, then find the optimal parameters (w, ρ) of the hyperplane such that:

$$w^T x_{inlier} \geq \rho$$

$$w^T x_{outlier} < \rho.$$



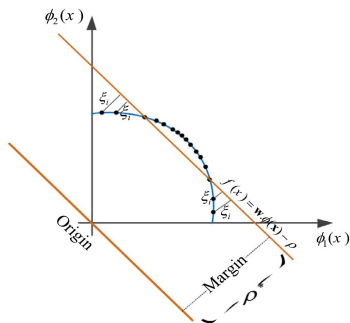
One class classification - hyperplane



- In the separable case, there is an infinite number of hyperplanes of choice



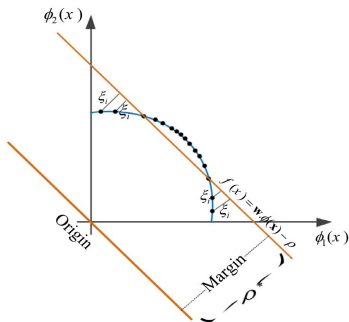
One class classification - hyperplane



- In the separable case, there is an infinite number of hyperplanes of choice
- Let the distance from the origin (of a hyperplane) be named as *margin*, then we desire to obtain the hyperplane with the maximal margin.



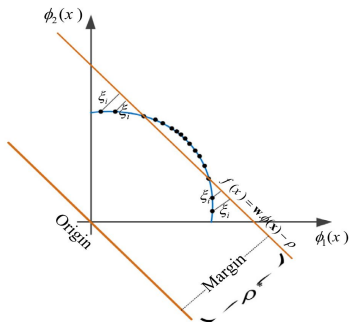
One class classification - hyperplane



- In the separable case, there is an infinite number of hyperplanes of choice
- Let the distance from the origin (of a hyperplane) be named as *margin*, then we desire to obtain the hyperplane with the maximal margin.
- The distance from x to hyperplane $\{x : w^T x = \rho\}$ is $\frac{|w^T x - \rho|}{\|w\|}$.



One class classification - hyperplane



- In the separable case, there is an infinite number of hyperplanes of choice
- Let the distance from the origin (of a hyperplane) be named as *margin*, then we desire to obtain the hyperplane with the maximal margin.
- The distance from x to hyperplane $\{x : w^T x = \rho\}$ is $\frac{|w^T x - \rho|}{\|w\|}$.
- Thus we maximize $\frac{\rho}{\|w\|}$.



- In order to maximize $\frac{\rho}{\|w\|}$, we minimize $\frac{1}{2}\|w\|^2 - \rho$
- In the nonseparable case we allow slack variables ξ to encode the outlyingness of nonseparable data:

$$\begin{aligned}w^T x_{inlier} &\geq \rho \\ \xi = \rho - w^T x_{outlier} &\geq 0.\end{aligned}$$

$$\begin{aligned}\min_{w, \xi \in \mathbb{R}^m, \rho \in \mathbb{R}} \quad & \frac{1}{2}\|w\|^2 + \frac{1}{m\nu} \sum_{i=1}^m \xi_i - \rho \\ \text{s.t.} \quad & \langle w, x_i \rangle \geq \rho - \xi_i, \xi_i \geq 0 \quad \forall i \in \{1, \dots, m\}.\end{aligned}$$



$$\min_{\mathbf{w}, \xi \in \mathbb{R}^m, \rho \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m \xi_i - \rho$$

s.t. $\langle \mathbf{w}, \mathbf{x}_i \rangle \geq \rho - \xi_i, \xi_i \geq 0 \quad \forall i \in \{1, \dots, m\}.$

- convex QP with m linear inequalities constraints
- regularization: $\|\mathbf{w}\|^2$ (justify the minimum margin hyperplane)
- error slack variables: ξ
- penalty parameter: $1/\nu \in [1, \infty)$



The primal problem is convex (linear constraints):

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m \xi_i - \rho \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{x}_i \geq \rho - \xi_i \quad \forall i = 1, \dots, m \\ & \xi \geq 0. \end{aligned}$$

The large number of constraints makes the primal hard to handle. Therefore, we take steps toward the dual: let the Lagrange multipliers $\lambda, \gamma \geq 0$

$$\mathcal{L}(\mathbf{w}, \rho, \xi, \lambda, \gamma) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m \xi_i - \rho - \sum_{i=1}^m \lambda_i [\mathbf{w}^T \mathbf{x}_i - \rho + \xi_i] - \sum_{i=1}^m \gamma_i \xi_i.$$



Kuhn-Tucker optimality conditions:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \rho, \xi, \lambda, \gamma) = \mathbf{w} - \sum_{i=1}^m \lambda_i \mathbf{x}_i = \mathbf{0}$$

$$\nabla_{\rho} \mathcal{L}(\mathbf{w}, \rho, \xi, \lambda, \gamma) = \sum_{i=1}^m \lambda_i - 1 = 0$$

$$\nabla_{\xi_i} \mathcal{L}(\mathbf{w}, \rho, \xi, \lambda, \gamma) = \frac{1}{m\nu} - \lambda_i - \gamma_i.$$

$$\gamma \geq \mathbf{0}, \mathbf{0} \leq \lambda \leq \frac{1}{m\nu}.$$

We observe that at optimality we have:

$$\mathbf{w}^* = \sum_{i=1}^m \lambda_i^* \mathbf{x}_i.$$

The datapoints \mathbf{x}^i such that $\lambda_i^* = 0$ do not contribute to the problem solution, those for which $\lambda_i^* > 0$ are **support vectors**: $\langle \mathbf{w}^*, \mathbf{x}^i \rangle - \rho^* = 0$.



The dual problem:

$$\begin{aligned} \max_{\lambda} \quad & -\frac{1}{2} \lambda^T X X^T \lambda \\ \text{s.t.} \quad & e^T \lambda = 1, \quad 0 \leq \lambda_i \leq \frac{1}{m\nu}. \end{aligned}$$

where $X = [x_1 \quad x_2 \quad \cdots \quad x_m]$.

- A convex quadratic objective with a high-dimensional Hessian ($m \times m$).
- Strict inliers $\lambda_i^* = 0$, SVs $\lambda_i^* > 0$.
- Once optimal λ^* is computed then $w^* = \sum_{i=1}^m \lambda_i^* x_i$.



Back to our data:

$$x_1, x_2, \dots, x_m \subset X.$$

In the non-separable case, we map the training data in high-dimensional spaces by choosing a function $\phi : X \rightarrow \mathcal{F}$ such that the inner product between the images of ϕ can be evaluated some simple kernel:

$$k(x, y) := \langle \phi(x), \phi(y) \rangle.$$

Example: Gaussian kernel

$$k(x, y) := e^{-\frac{\|x-y\|^2}{\sigma}}$$



$$\phi(x_1), \phi(x_2), \dots, \phi(x_m) \in \mathcal{F}.$$

Primal nonlinear problem:

$$\begin{aligned} \min_{w, \xi \in \mathbb{R}^m, \rho \in \mathbb{R}} \quad & \frac{1}{2} \|w\|^2 + \frac{1}{m\nu} \sum_{i=1}^m \xi_i - \rho \\ \text{s.t.} \quad & \langle w, \phi(x_i) \rangle \geq \rho - \xi_i, \xi_i \geq 0 \quad \forall i \in \{1, \dots, m\}. \end{aligned}$$

Now the dimension of w is the dimension of feature space. The problem is still convex and has a similar dual as in the linear case.



The dual problem:

$$\begin{aligned} \max_{\lambda} \quad & -\frac{1}{2} \lambda^T K \lambda \\ \text{s.t.} \quad & \mathbf{e}^T \lambda = 1, \quad 0 \leq \lambda_i \leq \frac{1}{m\nu}. \end{aligned}$$

where $K = [\phi(x_1) \quad \phi(x_2) \quad \cdots \quad \phi(x_m)]^T [\phi(x_1) \quad \phi(x_2) \quad \cdots \quad \phi(x_m)]$.

- Hessian $K_{ij} = \langle \phi(x_i), \phi(x_j) \rangle$ and $K_{ii} = 1$.
- Strict inliers $\lambda_i^* = 0$, SVs $\lambda_i^* > 0$.
- Once optimal λ^* is computed then $w^* = \sum_{i=1}^m \lambda_i^* \phi(x_i)$.



How to solve the Dual problem ?

Possible algorithms:

- off-the-shelf QP solvers: cvxpy, quadprog, MOSEK etc. $O(m^3)$



How to solve the Dual problem ?

Possible algorithms:

- off-the-shelf QP solvers: cvxpy, quadprog, MOSEK etc. $O(m^3)$
- Simplex-type feasible set ($m \log(m)$ to project on)



Possible algorithms:

- off-the-shelf QP solvers: cvxpy, quadprog, MOSEK etc. $O(m^3)$
- Simplex-type feasible set ($m \log(m)$ to project on)
- The cost of a first-order iteration: gradient evaluation $O(m^2)$ + projection onto the simple $O(m \log(m))$.



How to solve the Dual problem ?

Possible algorithms:

- off-the-shelf QP solvers: cvxpy, quadprog, MOSEK etc. $O(m^3)$
- Simplex-type feasible set ($m \log(m)$ to project on)
- The cost of a first-order iteration: gradient evaluation $O(m^2)$ + projection onto the simple $O(m \log(m))$.
- Best suggestion: approach large-scale instances by coordinate descent algorithms (scikit-learn uses libsvm for training)



$$\begin{aligned} \min_{\lambda} \quad & \frac{1}{2} \lambda^T H \lambda + b^T \lambda \\ \text{s.t.} \quad & e^T \lambda = 1, \quad 0 \leq \lambda_i \leq C. \end{aligned}$$

Idea of CD

Instead of approximating the whole λ at each iteration, update only a small block of variables.



$$\begin{aligned} \min_{\lambda} \quad & \frac{1}{2} \lambda^T H \lambda + \mathbf{b}^T \lambda \\ \text{s.t.} \quad & \mathbf{e}^T \lambda = 1, \quad 0 \leq \lambda_i \leq C. \end{aligned}$$

Idea of CD

Instead of approximating the whole λ at each iteration, update only a small block of variables.

Example CD

Exact 2-coordinate descent: choose $(i, j) \in \{1, \dots, m\}$; let $\lambda_{ij} = [\lambda_i \ \lambda_j]^T$

$$\begin{aligned} \lambda_{ij}^+ := \arg \min_{\lambda_{ij}} \quad & \frac{1}{2} (\lambda_i^2 + \lambda_j^2) + K_{ij} \lambda_i \lambda_j + \mathbf{b}_{ij}^T \lambda_{ij} \\ \text{s.t.} \quad & \lambda_i + \lambda_j = \Delta, \quad 0 \leq \lambda_i, \lambda_j \leq C. \end{aligned}$$

$$\begin{aligned} \min_{\lambda} \quad & \frac{1}{2} \lambda^T H \lambda + \mathbf{b}^T \lambda \\ \text{s.t.} \quad & \mathbf{e}^T \lambda = 1, \quad 0 \leq \lambda_i \leq C. \end{aligned}$$

Exact 2-coordinate descent main loop:

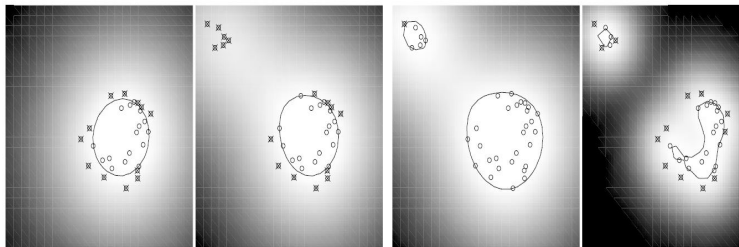
- 1 choose $(i, j) \in \{1, \dots, m\}$ (cyclic, random etc.)
- 2 update:

$$\begin{aligned} \lambda_{ij}^* := \arg \min_{\lambda_{ij}} \quad & \frac{1}{2} (\lambda_i^2 + \lambda_j^2) + K_{ij} \lambda_i \lambda_j + \mathbf{b}_{ij}^T \lambda_{ij} \\ \text{s.t.} \quad & \lambda_i + \lambda_j = \Delta, \quad 0 \leq \lambda_i, \lambda_j \leq C. \end{aligned}$$

where $\Delta = 1 - \sum_{t \neq i, j} \lambda_t$.

- 3 $\lambda_{ij}^+ := \lambda_{ij}^*$ and $\lambda_{ij}^- := \lambda_{ij}$





ν , width c

0.5, 0.5

0.5, 0.5

0.1, 0.5

0.5, 0.1

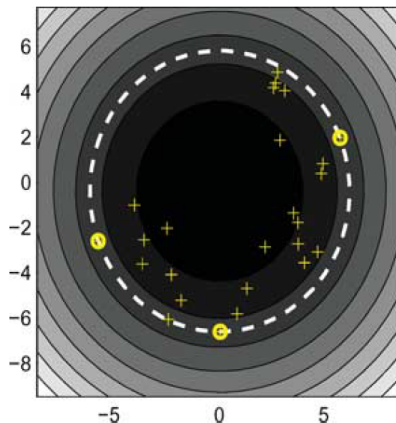


One class classification - hypersphere

Consider the training data is provided:

$$x_1, x_2, \dots, x_m \subset X.$$

where m is the number of the observations and X some space (i.e. compact subset of \mathbb{R}^n).



Take the decision function as $f(x) = \text{sgn}(\|c - x_i\| - R)$, then finding the optimal parameters of the hypersphere reduces to solving:

$$\begin{aligned} \min_{c, R \geq 0, \xi \in \mathbb{R}^m} \quad & R^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & \|c - x_i\| \leq R^2 + \xi_i, \xi_i \geq 0 \quad \forall i \in \{1, \dots, m\}. \end{aligned}$$



The dual problem:

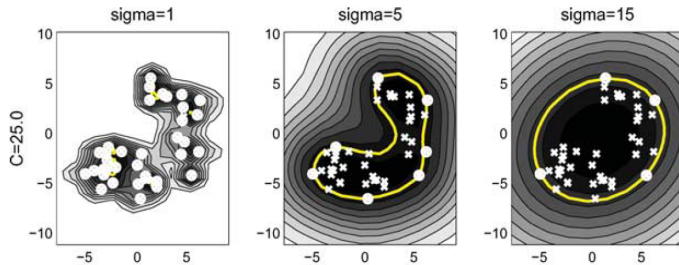
$$\begin{aligned} \max_{\lambda} \quad & -\frac{1}{2} \lambda^T X^T X \lambda + \lambda^T \text{diag}(X^T X) \\ \text{s.t.} \quad & \mathbf{e}^T \lambda = 1, \quad 0 \leq \lambda_i \leq C. \end{aligned}$$

where $X = [x_1 \quad x_2 \quad \cdots \quad x_m]$.

- 1 $\|x - c\|^2 < R^2 \rightarrow \lambda_i = 0, \gamma_i = 0$
- 2 $\|x - c\|^2 = R^2 \rightarrow 0 < \lambda_i < C, \gamma_i = 0$
- 3 $\|x - c\|^2 > R^2 \rightarrow \lambda_i = C, \gamma_i > 0$

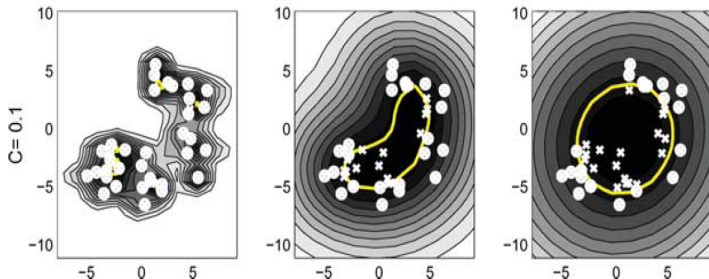
When the data is normalized then SVDD is equivalent with OCSVM.





- large C means high penalty of outlyingness (thus large coverage of data)





- large C means high penalty of outlyingness (thus large coverage of data)
- low C means high margin of the hyperplane (thus large robustness)
- however, any datapoint has a certain influence on the decision boundary



- One class classification: OC-SVM, SVDD. Algorithms.
- **Robust versions. Algorithms.**



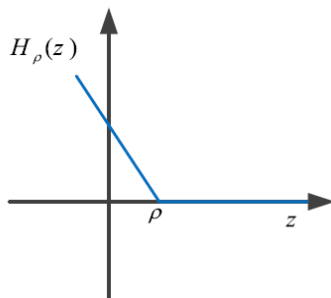
$$\begin{aligned} \min_{\mathbf{w}, \xi \in \mathbb{R}^m, \rho \in \mathbb{R}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m \xi_i - \rho \\ \text{s.t.} \quad & \langle \mathbf{w}, \mathbf{x}_i \rangle \geq \rho - \xi_i, \xi_i \geq 0 \quad \forall i \in \{1, \dots, m\}. \end{aligned}$$

By elimination of ξ we obtain:

$$\min_{\mathbf{w}, \rho \in \mathbb{R}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m \underbrace{\max\{0, \rho - \langle \mathbf{w}, \mathbf{x}_i \rangle\}}_{H_\rho(\langle \mathbf{w}, \mathbf{x}_i \rangle)} - \rho.$$

We denote hinge penalty (convex) function: $H_\rho(z) := \max\{0, \rho - z\}$.



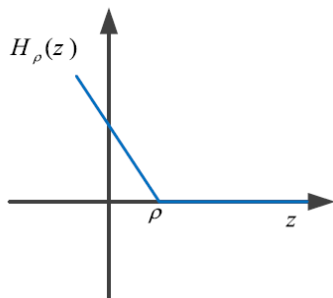


a) The Hinge Loss
for OC-SVM

- if a datapoint falls above the hyperplane $w^T z \geq \rho$, then no penalty $H_\rho(z) = 0$



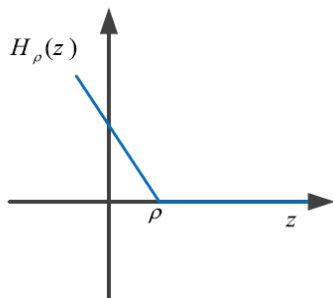
One class classification - reformulation



a) The Hinge Loss
for OC-SVM

- if a datapoint falls above the hyperplane $w^T z \geq \rho$, then no penalty $H_\rho(z) = 0$
- otherwise, if $w^T z < \rho$, then a penalty corresponding to the distance of this point to the hyperplane will be applied

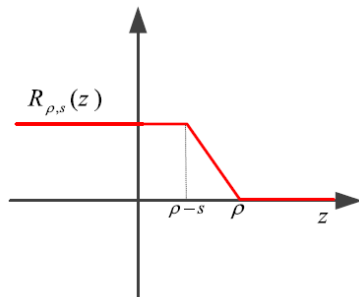




a) The Hinge Loss
for OC-SVM

- if a datapoint falls above the hyperplane $w^T z \geq \rho$, then no penalty $H_\rho(z) = 0$
- otherwise, if $w^T z < \rho$, then a penalty corresponding to the distance of this point to the hyperplane will be applied
- one can "robustify" $H_\rho(z)$ by limiting the penalty to a given threshold





b) The Ramp Loss
for OC-SVM

- Ramp function $R_{\rho,s}(z) = \begin{cases} 0, & \text{if } z \geq \rho \\ \rho - z, & \text{if } \rho - s < z < \rho \\ s, & \text{if } z < \rho - s. \end{cases}$



Given parameters s, ν then

$$\min_{w, \rho \in \mathbb{R}} \frac{1}{2} \|w\|^2 + \frac{1}{m\nu} \sum_{i=1}^m R_{\rho, s}(\langle w, x_i \rangle) - \rho.$$

- This new problem is nonconvex nondifferentiable



Given parameters s, ν then

$$\min_{\mathbf{w}, \rho \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m R_{\rho, s}(\langle \mathbf{w}, \mathbf{x}_i \rangle) - \rho.$$

- This new problem is nonconvex nondifferentiable
- Notice that $R_{\rho, s}(z) = H_{\rho}(z) - H_{\rho-s}(z)$ (difference of convex function)



Given parameters s, ν then

$$\min_{\mathbf{w}, \rho \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m R_{\rho, s}(\langle \mathbf{w}, \mathbf{x}_i \rangle) - \rho.$$

- This new problem is nonconvex nondifferentiable
- Notice that $R_{\rho, s}(z) = H_{\rho}(z) - H_{\rho-s}(z)$ (difference of convex function)
- Based on this observation we can derive a simple iterative first-order algorithm.



Given parameters s, ν then

$$\begin{aligned} \min_{\mathbf{w}, \rho \in \mathbb{R}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m R_{\rho, s}(\langle \mathbf{w}, \mathbf{x}_i \rangle) - \rho \\ = \quad & \underbrace{\frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m H_{\rho}(\langle \mathbf{w}, \mathbf{x}_i \rangle) - \rho}_{\text{convex}} - \underbrace{\frac{1}{m\nu} \sum_{i=1}^m H_{\rho-s}(\langle \mathbf{w}, \mathbf{x}_i \rangle)}_{\text{convex}}. \end{aligned}$$



DC algorithm

- 1 Initialize (\mathbf{w}_1, ρ_1) and $k := 0$
- 2 $(\mathbf{w}_{k+1}, \rho_{k+1}) =$

$$\arg \min \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{m\nu} \sum_{i=1}^m H_\rho(\langle \mathbf{w}, \mathbf{x}_i \rangle) - \rho - \frac{1}{m\nu} \sum_{i=1}^m \left\langle \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} H'_{\rho_k - s}(\langle \mathbf{w}_k, \mathbf{x}_i \rangle), (\mathbf{w}, \rho) \right\rangle$$
- 3 If $(\mathbf{w}_{k+1}, \rho_{k+1})$ satisfies the convergence criterion, then STOP; otherwise, $k := k + 1$ and reiterate.



Dual DC algorithm

- 1 Compute $\delta_i = \begin{cases} -\frac{1}{m\nu} & \rho - (w^k)^T \phi(x_i) > s \\ 0, & \text{otherwise} \end{cases}$.
 - 2 $\lambda^{k+1} := \max_{\lambda} -\frac{1}{2} \lambda^T K \lambda$ s.t. $e^T \lambda = 1$, $-\nu m \delta_i \leq \lambda_i \leq \frac{1}{m\nu} - \nu m \delta_i$.
 - 3 If λ^{k+1} satisfies the convergence criterion, then STOP; otherwise, $k := k + 1$ and reiterate.
- If the number of iterations is T then Dual DC solves T QP dual problems.



Dual DC algorithm

- 1 Compute $\delta_i = \begin{cases} -\frac{1}{m\nu} & \rho - (w^k)^T \phi(x_i) > s \\ 0, & \text{otherwise} \end{cases}$.
 - 2 $\lambda^{k+1} := \max_{\lambda} -\frac{1}{2} \lambda^T K \lambda$ s.t. $e^T \lambda = 1$, $-\nu m \delta_i \leq \lambda_i \leq \frac{1}{m\nu} - \nu m \delta_i$.
 - 3 If λ^{k+1} satisfies the convergence criterion, then STOP; otherwise, $k := k + 1$ and reiterate.
- If the number of iterations is T then Dual DC solves T QP dual problems.
 - DC provides the pair λ^*, ρ^* .

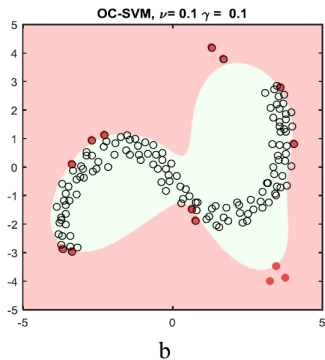
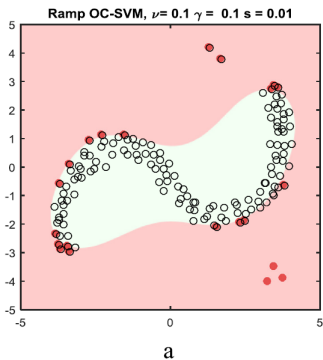


Dual DC algorithm

- 1 Compute $\delta_i = \begin{cases} -\frac{1}{m\nu} & \rho - (w^k)^T \phi(x_i) > s \\ 0, & \text{otherwise} \end{cases}$.
 - 2 $\lambda^{k+1} := \max_{\lambda} -\frac{1}{2} \lambda^T K \lambda$ s.t. $e^T \lambda = 1$, $-\nu m \delta_i \leq \lambda_i \leq \frac{1}{m\nu} - \nu m \delta_i$.
 - 3 If λ^{k+1} satisfies the convergence criterion, then STOP; otherwise, $k := k + 1$ and reiterate.
- If the number of iterations is T then Dual DC solves T QP dual problems.
 - DC provides the pair λ^*, ρ^* .
 - Test on new sample x : evaluate $\text{sgn}(\sum_i \lambda_i^* k(x_i, x) - \rho^*)$



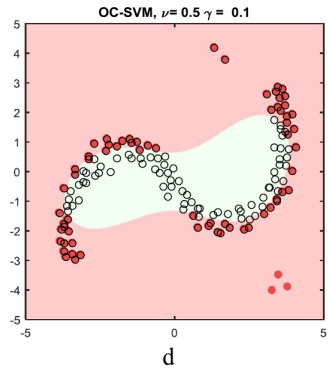
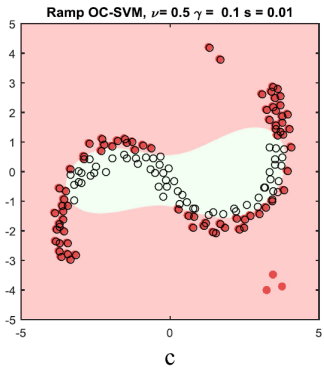
Experiments (synthetic in 2D)



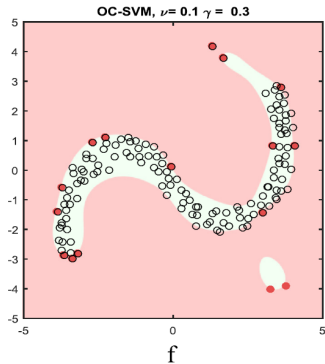
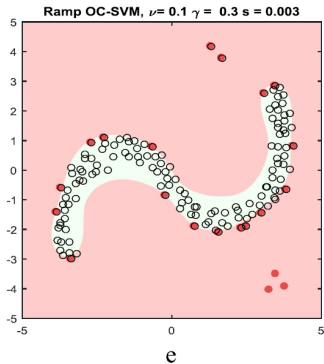
- ν estimate the ratio of outliers
- outliers have a lower impact over Ramp-OCSVM



Experiments (synthetic in 2D)



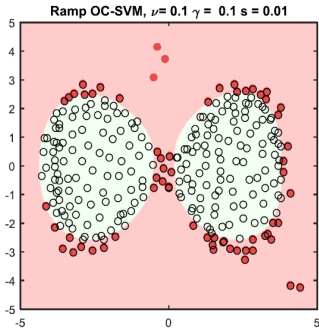
Experiments (synthetic in 2D)



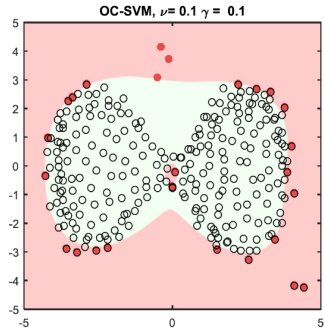
- for small ν OCSVM shift towards outliers
- Ramp-OCSVM controls this shifting through parameter s



Experiments (synthetic in 2D)



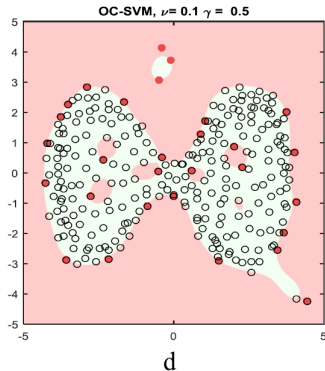
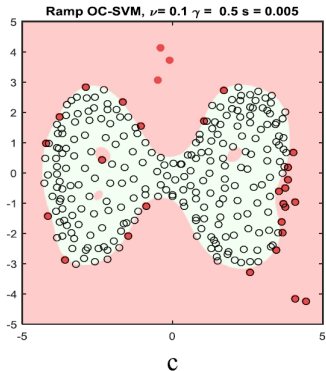
a



b



Experiments (synthetic in 2D)



- behaviour comparison against changing the kernel parameter γ



- 1 Tian, Yingjie, et al. "Ramp loss one-class support vector machine; a robust and effective approach to anomaly detection problems." *Neurocomputing* 310 (2018): 223-235.
- 2 Alam, Shamshe, et al. "One-class support vector classifiers: A survey." *Knowledge-Based Systems* 196 (2020): 105754.
- 3 Tax, David MJ, and Robert PW Duin. "Support vector data description." *Machine learning* 54 (2004): 45-66.
- 4 Kaiser, Ferdinand. *Robust Support Vector Machines For Implicit Outlier Removal*. MS thesis. 2013.
- 5 Schölkopf, Bernhard, et al. "Estimating the support of a high-dimensional distribution." *Neural computation* 13.7 (2001): 1443-1471.

