

MACHINE LEARNING

Support Vector Machine For Classification

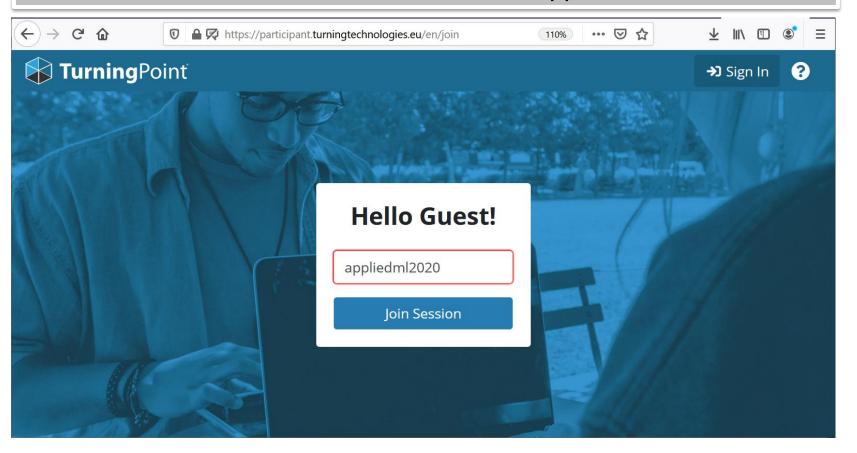
Interactive Lecture



Launch polling system

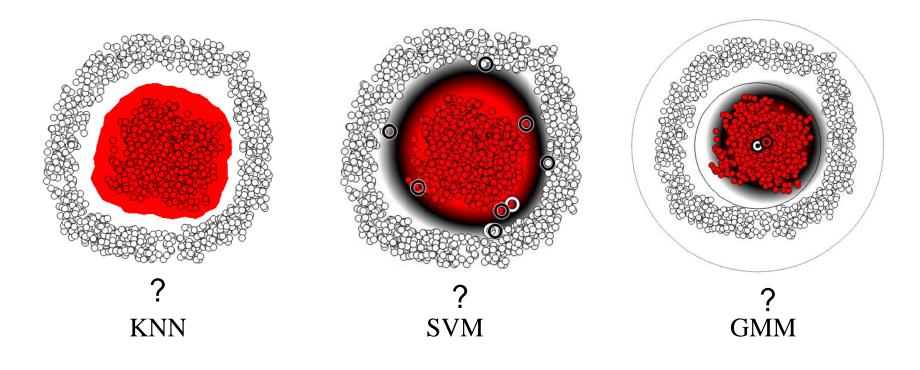
https://participant.turningtechnologies.eu/en/join

Acces as GUEST and enter the session id: appliedml2020



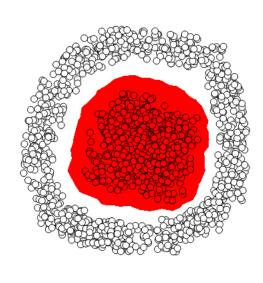


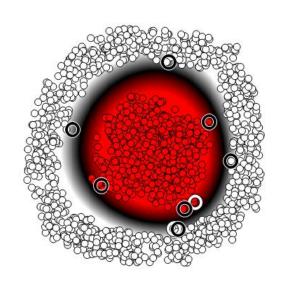
Can you recognize the classifiers?



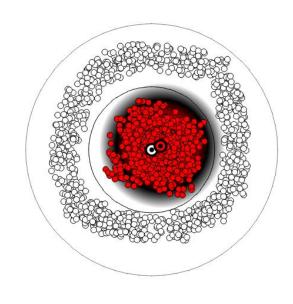


Which of the three solutions requires the least number of parameters for prediction at testing time?



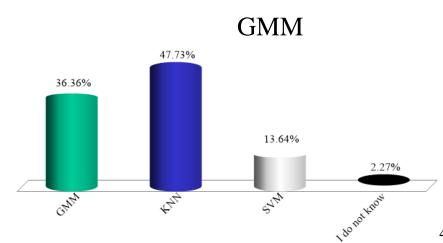


SVM



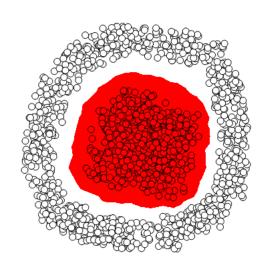
KNN

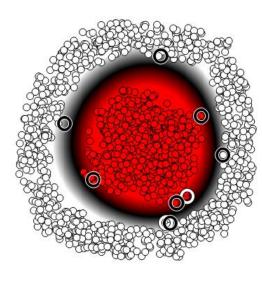
- **GMM**
- **KNN** В.
- **SVM**
- I do not know

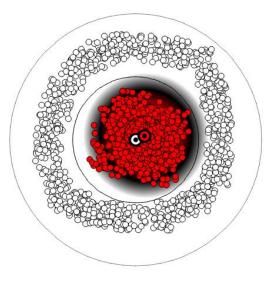




Which of the three solutions requires the least number of parameters for prediction at testing time?







KNN

KNN is the most costly as it requires to keep all Datapoints at testing.

SVM

SVM requires 9 SVs, each 2D

- \rightarrow 9*2=18 parameters
- \rightarrow + 8 parameters for the α
- \rightarrow +1 for b

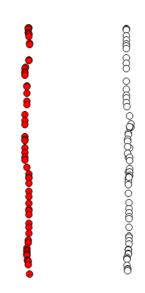
$$f(x) = sgn\left(\sum_{i=1}^{M} \alpha_i y^i k(x, x^i) + b\right)$$

GMM

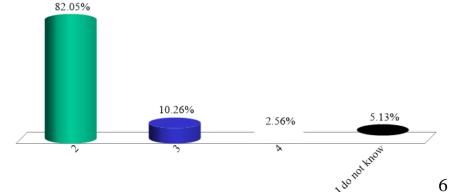
GMM is the least costly as it requires 2 spherical Gauss fcts. Each fct is represented by its 2-dim. mean and its variance → (2+1)*2=6 parameters + 2 priors on Gauss fcts



How many Support Vectors (SV) do you need at minimum when using linear SVM?

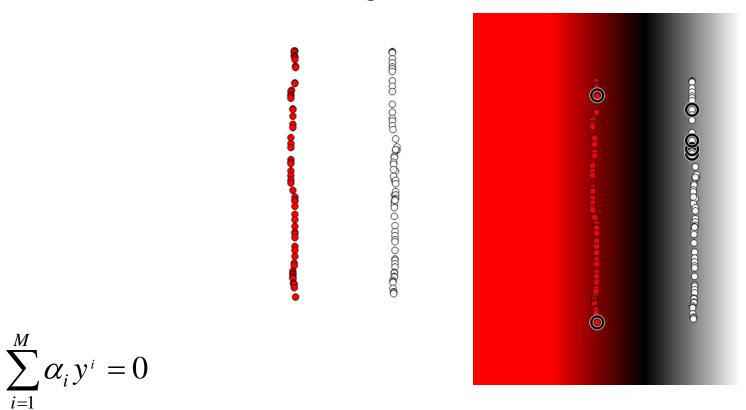


- 3
- I do not know





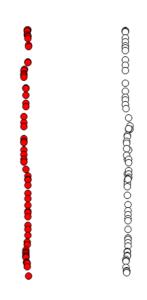
How many Support Vectors (SV) do you need at minimum when using linear SVM?



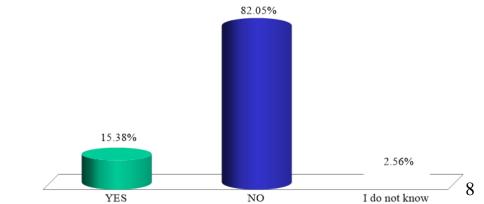
At minimum one needs only 2 SV-s, one in each class.

SVM usually finds more than required number of SV as the standard optimization problem is not constrained to find the minimum.

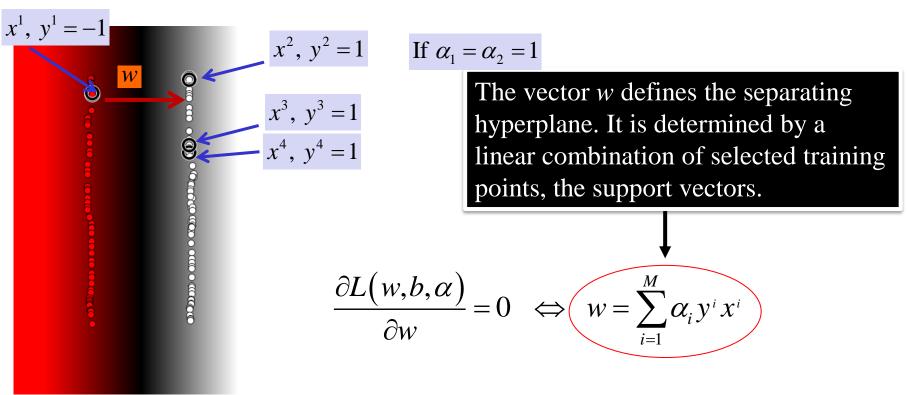




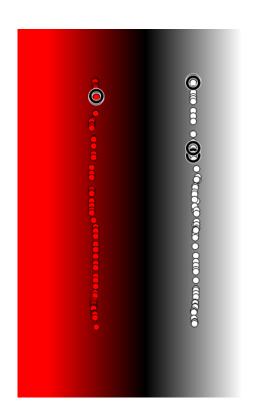
- A. YES
- B. NO
- C. I do not know

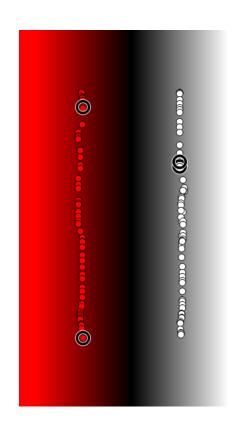


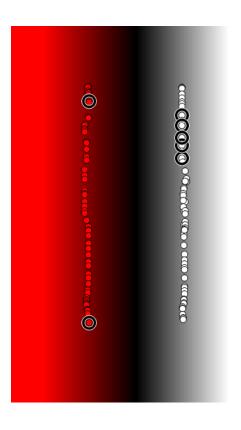






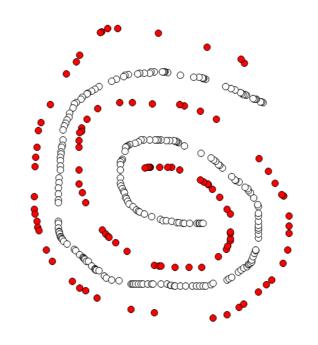




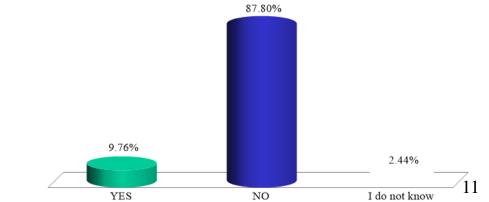


Convex optimization has a global optimum but several solutions may have the same optimum on the objective function. Here the same separating plane can be found when using different combinations of datapoints.

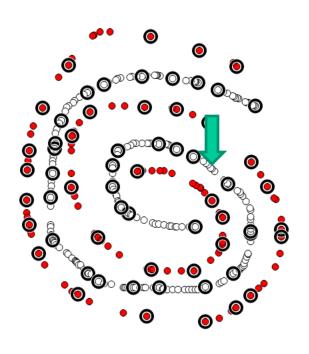


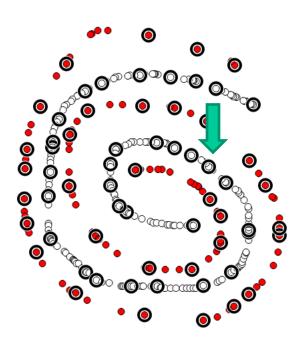


- A. YES
- B. NO
- C. I do not know



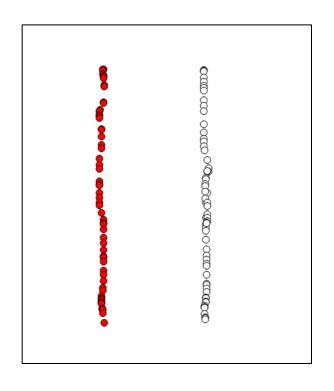




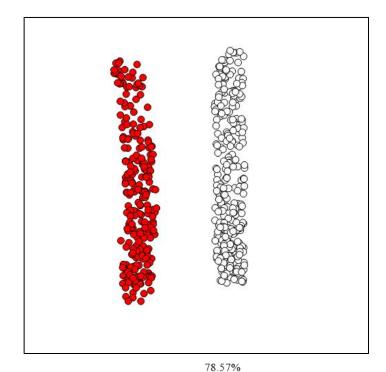


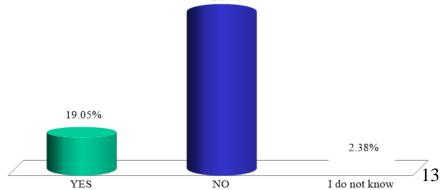
The solution is not unique either but as data become more complex, less combination will lead to global optimum. Only redundant points may be discarded from one run to the next.



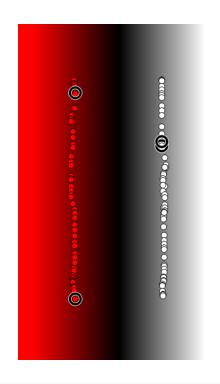


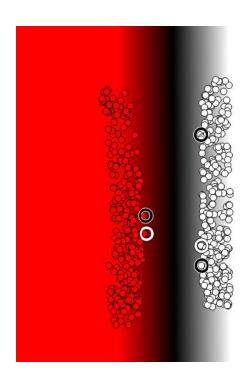
- A. YES
- B. NO
- C. I do not know





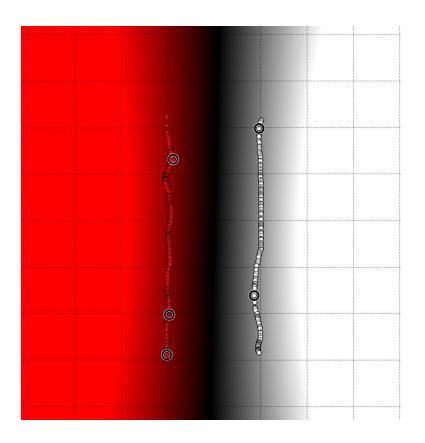


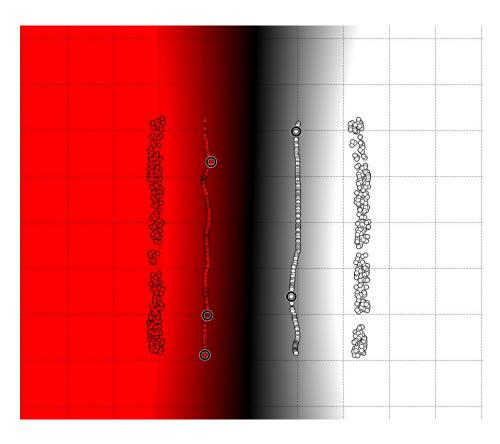




In general, the SVs may be different, but if we have the same number of datapoints close to the border of the margin, the SV-s may be the same, as only the datapoints at the border of the margin matter.

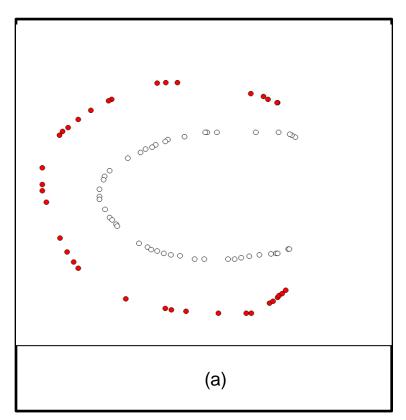




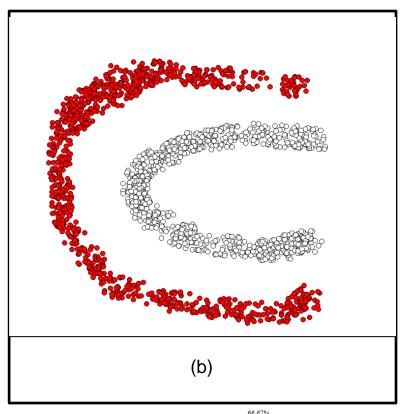


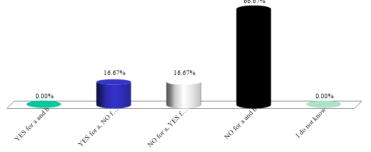
Adding datapoints away from the boundary will not affect classification.



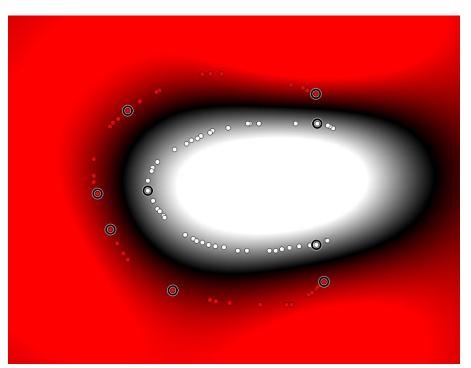


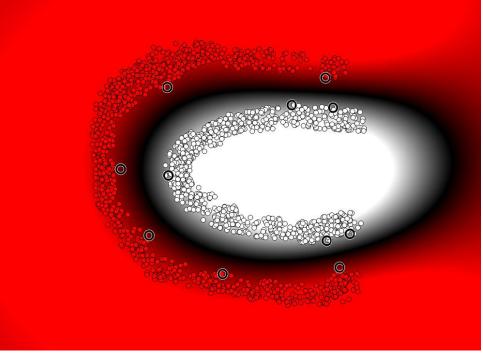
- A. YES for a and b
- B. YES for a, NO for b
- C. NO for a, YES for b
- D. NO for a and b
- E. I do not know





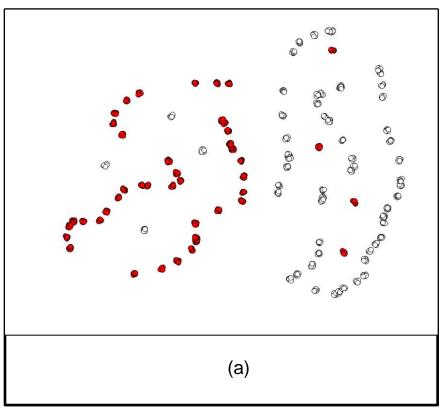




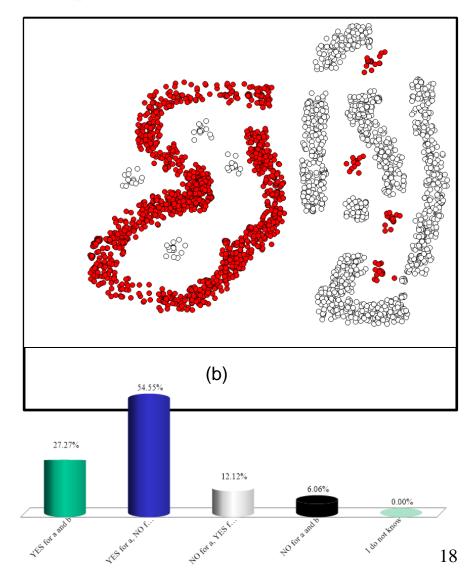


Barely as the boundary is not affected

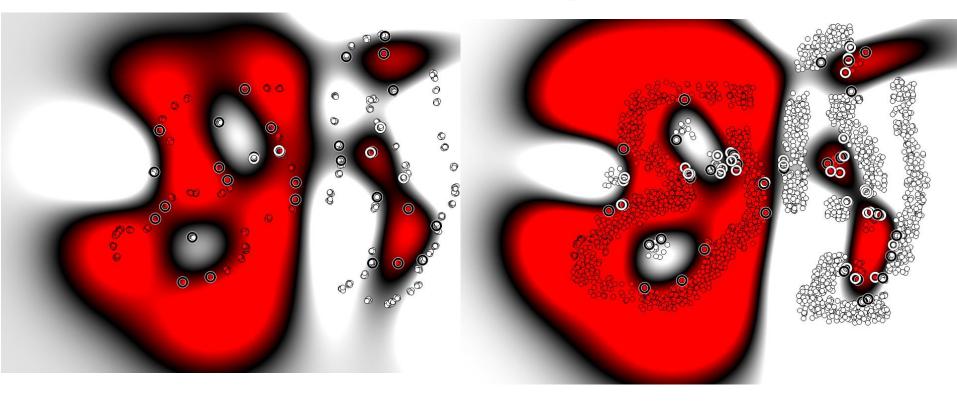




- A. YES for a and b
- B. YES for a, NO for b
- C. NO for a, YES for b
- D. NO for a and b
- E. I do not know



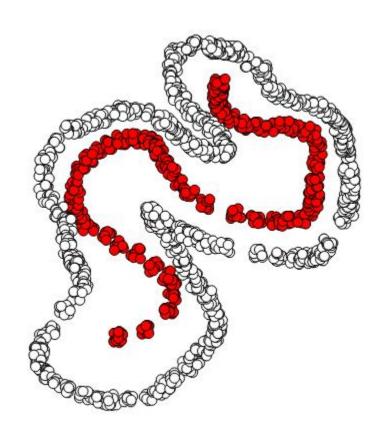




Yes, many new SV-s are added. To maintain the boundary around the regions with few datapoints requires a tight fit This implies to use a small kernel width. As a result, each SV has only a local influence.



Non-Linear Classification



The decision function is given by:

$$f(x) = sgn\left(\sum_{i=1}^{M} \alpha_i y_i k(x, x^i) + b\right)$$

How can we build the decision boundary?



The decision function is given by:

$$f(x) = sgn\left(\sum_{i=1}^{M} \alpha_i y_i k(x, x^i) + b\right)$$

RBF (Gaussian) kernel:
$$k(x, x^i) = e^{-\frac{\|x - x^i\|}{2\sigma^2}}$$
, $\sigma \in \mathbb{R}$.

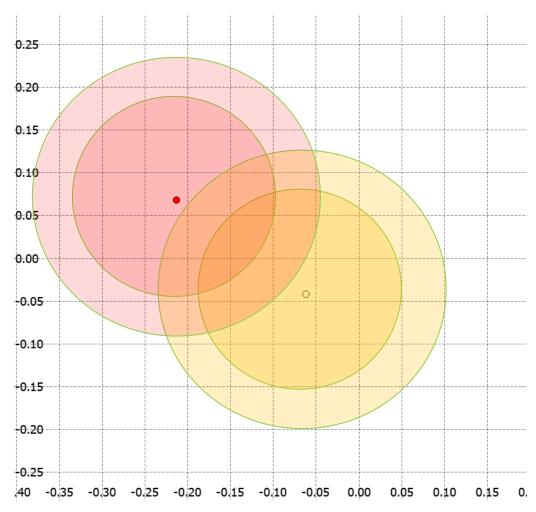


The decision function is given by:

$$f(x) = sgn(\alpha_1 y_1 k(x, x^1) + \alpha_2 y_2 k(x, x^2) + b)$$

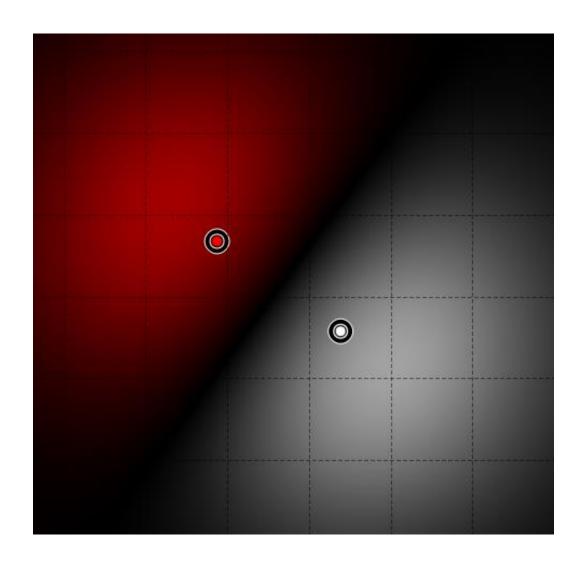
What are α_1, α_2 and b? -0.30 -0.25 -0.20 -0.15 -0.10 -0.05 0.00 0.05 0.10 0.15



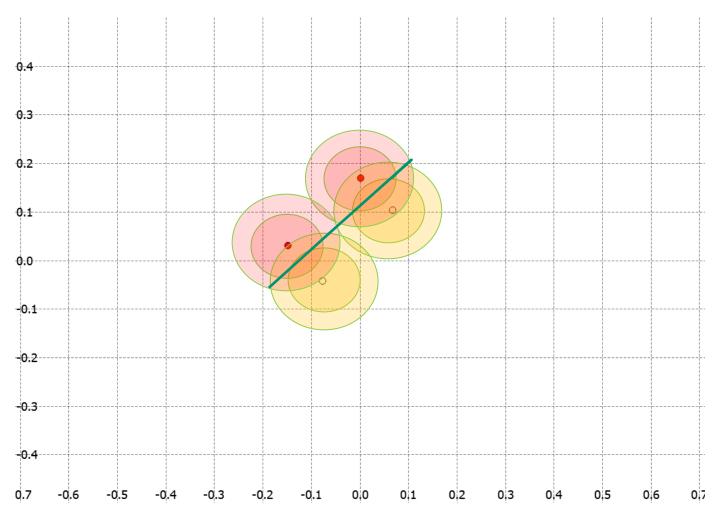


Draw the decision boundary assuming equal effect of both points





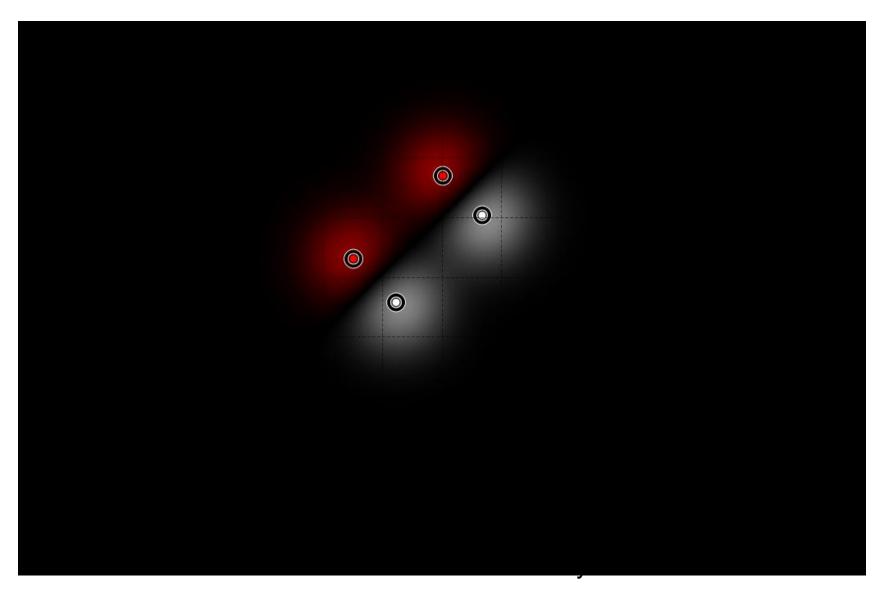




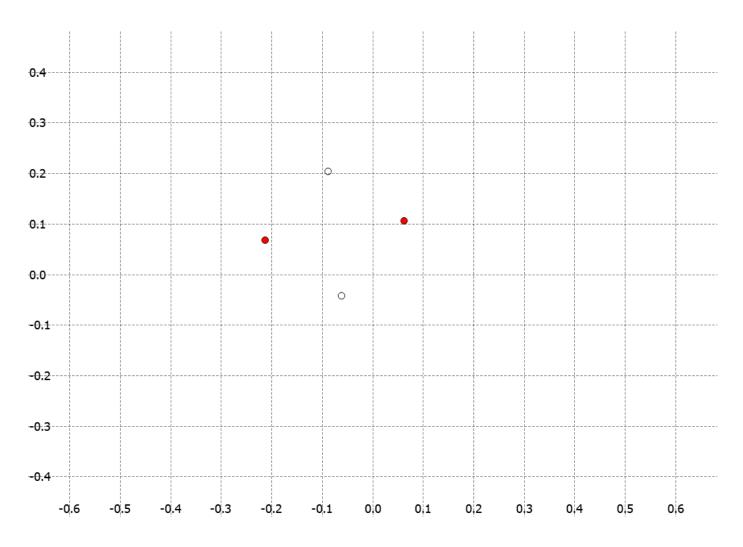
Draw the decision boundary?



Solution



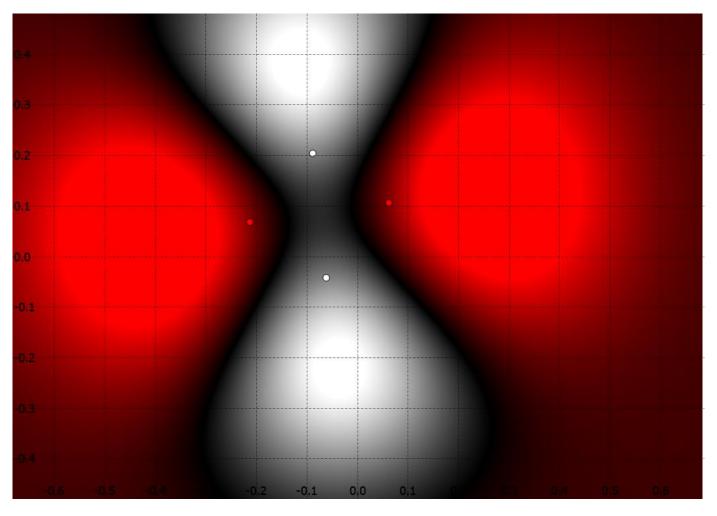




Draw the decision boundary?



Solution

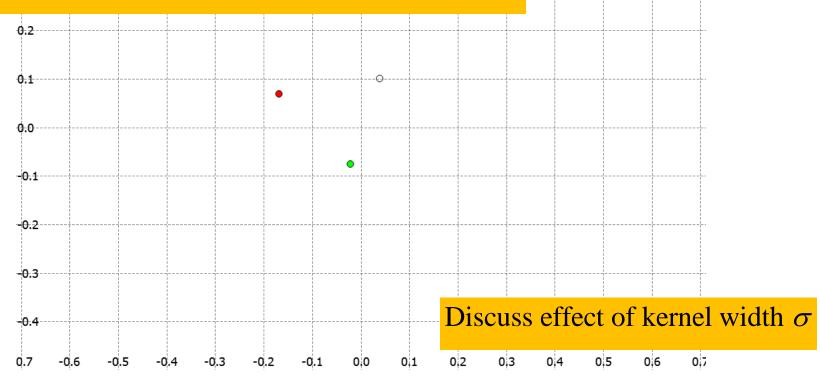


Draw the decision boundary?



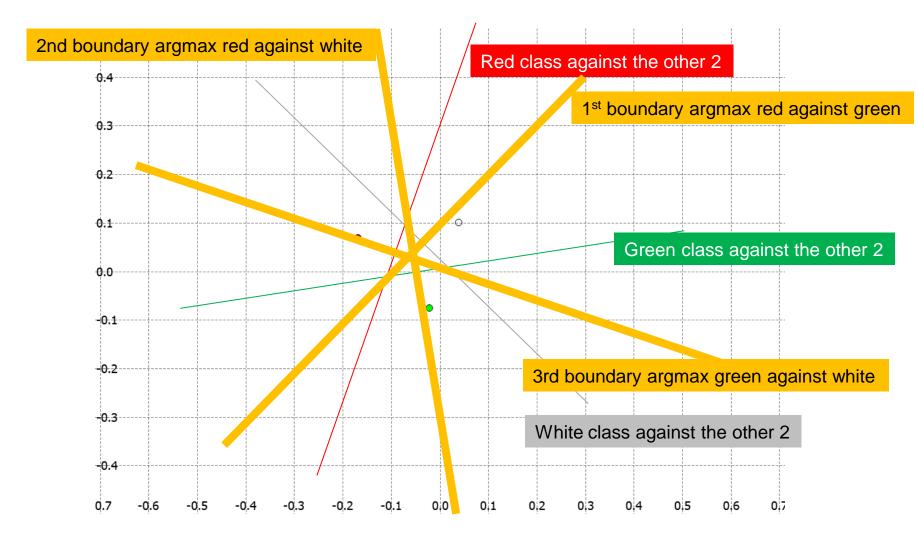
Compute the class label in a winner-take-all approach:

$$j=\underset{j=1,\dots K}{\operatorname{arg\,max}}\left(\sum_{i=1}^{M}y^{i}\alpha_{i}^{j}k\left(x,x^{i}\right)+b^{j}\right)$$

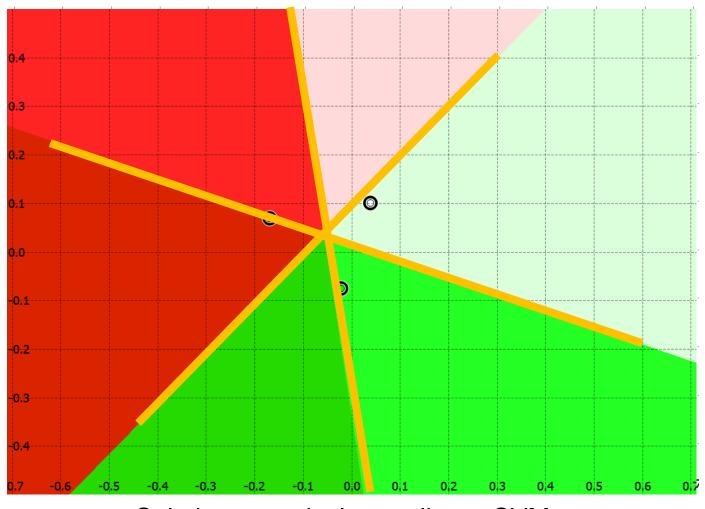


Draw the decision boundary (assume RBF kernel)?



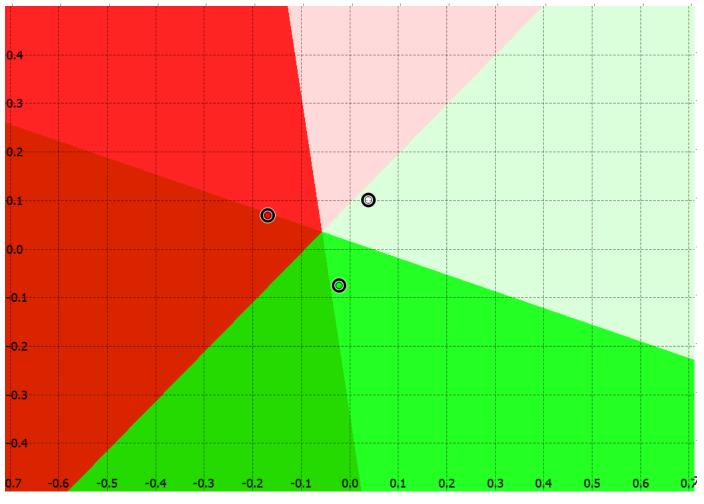






Solution ~ equivalent to linear SVM





Solution ~ equivalent to linear SVM



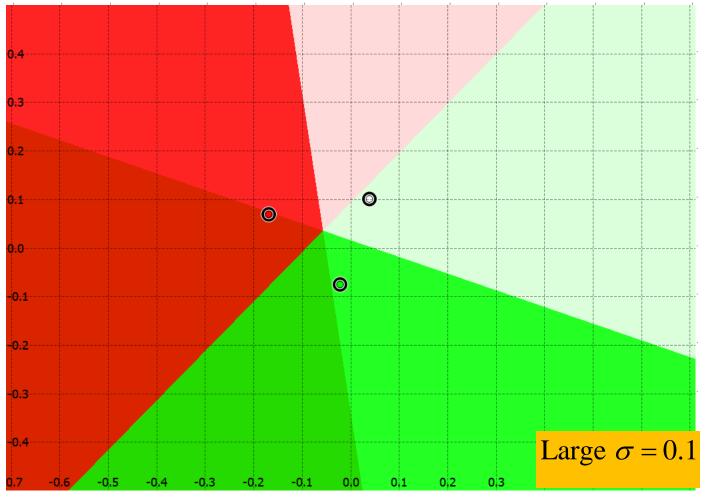
Discuss effect of kernel width σ

Compute the class label in a winner-take-all approach:

$$\mathbf{j} = \underset{j=1,\dots,K}{\operatorname{arg\,max}} \left(\sum_{i=1}^{M} y^{i} \alpha_{i}^{j} \underbrace{k(x, x^{i})}_{\sim 0 \text{ when } ||x-x^{i}|| \to \infty} + b^{j} \right) = \underset{j=1,\dots,K}{\operatorname{arg\,max}} (b^{j})$$

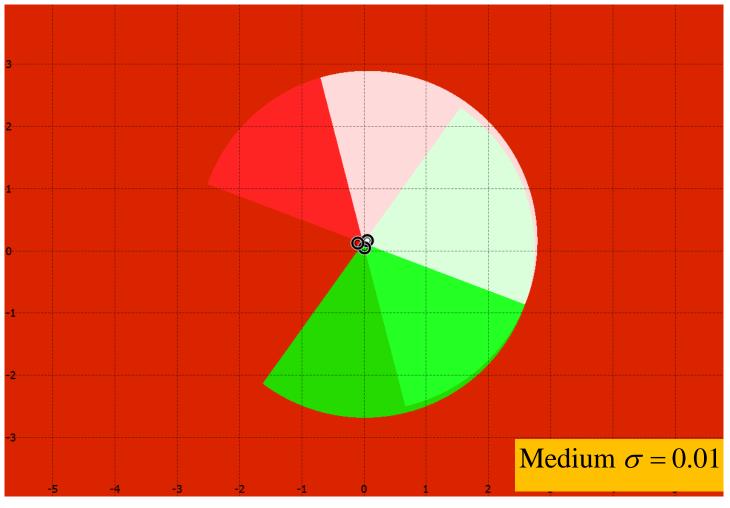
The kernel width does not affect this boundary except when very far from the datapoints where numerically the RBF function becomes zero and prediction is based on b





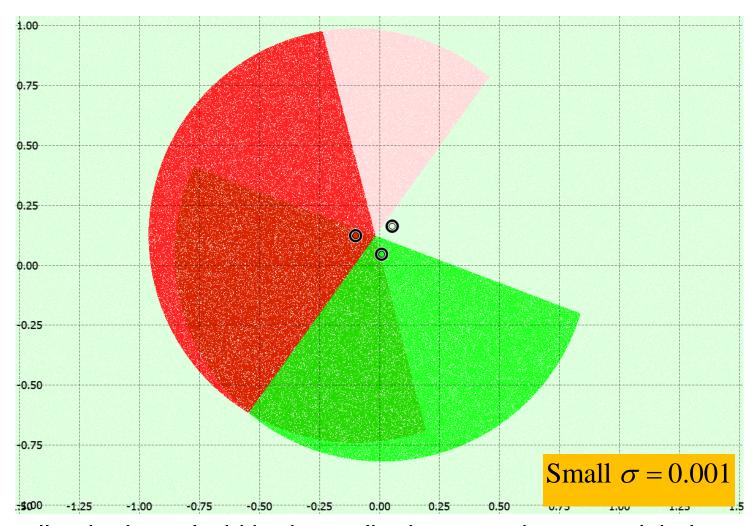
Solution ~ equivalent to linear SVM





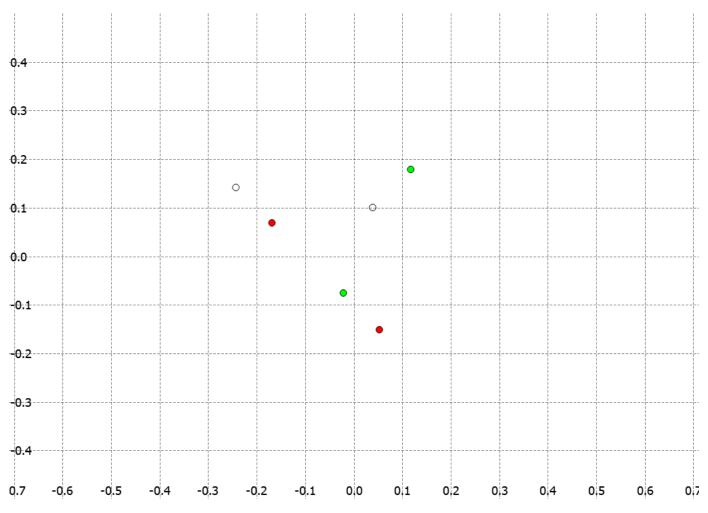
Far from datapoints, the red class labels wins.





The smaller the kernel width, the earlier in space the wrong label appears, as the RBF vanishes faster.





Draw the decision boundary (assume RBF kernel)?



