

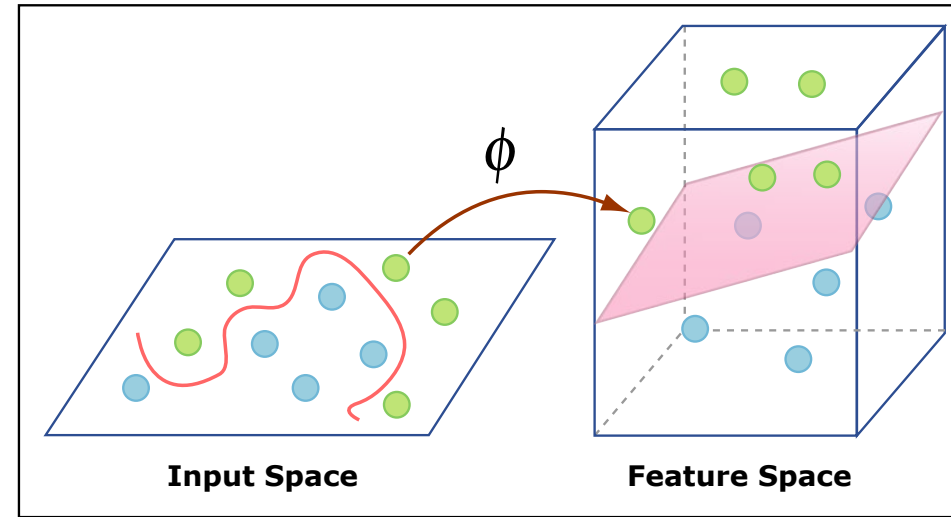
Kernels

We might want to consider something more complicated than a linear model:

Example 1: $[x^{(1)}, x^{(2)}] \rightarrow \Phi([x^{(1)}, x^{(2)}]) = [x^{(1)2}, x^{(2)2}, x^{(1)}x^{(2)}]$

Information unchanged, but now we have a **linear** classifier on the transformed points.

With the kernel trick, we just need kernel
 $k(\mathbf{a}, \mathbf{b}) = \Phi(\mathbf{a})^T \Phi(\mathbf{b})$



Example 4:

$$\begin{aligned}k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z} + c)^2 = \left(\sum_{j=1}^n x^{(j)} z^{(j)} + c \right) \left(\sum_{\ell=1}^n x^{(\ell)} z^{(\ell)} + c \right) \\&= \sum_{j=1}^n \sum_{\ell=1}^n x^{(j)} x^{(\ell)} z^{(j)} z^{(\ell)} + 2c \sum_{j=1}^n x^{(j)} z^{(j)} + c^2 \\&= \sum_{j,\ell=1}^n (x^{(j)} x^{(\ell)}) (z^{(j)} z^{(\ell)}) + \sum_{j=1}^n (\sqrt{2c} x^{(j)}) (\sqrt{2c} z^{(j)}) + c^2,\end{aligned}$$

and in $n = 3$ dimensions, one possible feature map is:

$$\Phi(\mathbf{x}) = [x^{(1)2}, x^{(1)}x^{(2)}, \dots, x^{(3)2}, \sqrt{2c}x^{(1)}, \sqrt{2c}x^{(2)}, \sqrt{2c}x^{(3)}, c]$$

and c controls the relative weight of the linear and quadratic terms in the inner product.

Even more generally, if you wanted to, you could choose the kernel to be any higher power of the regular inner product.

Dual representation, Sec 6.2

Primal problem: $\min_{\mathbf{w}} E(\mathbf{w})$

$$E = \frac{1}{2} \sum_n^N \{\mathbf{w}^T \mathbf{x}_n - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{t}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

$$\begin{aligned} \text{Solution } \mathbf{w} &= \mathbf{X}^+ \mathbf{t} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_M)^{-1} \mathbf{X}^T \mathbf{t} \\ &= \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_N)^{-1} \mathbf{t} = \mathbf{X}^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} = \mathbf{X}^T \mathbf{a} \end{aligned}$$

The kernel is $\mathbf{K} = \mathbf{X} \mathbf{X}^T$

Dual representation is : $\min_{\mathbf{a}} E(\mathbf{a})$

$$E = \frac{1}{2} \sum_n^N \{\mathbf{w}^T \mathbf{x}_n - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 = \|\mathbf{K}\mathbf{a} - \mathbf{t}\|_2^2 + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

Prediction

$$y = \mathbf{w}^T \mathbf{x} = \mathbf{a}^T \mathbf{X} \mathbf{x} = \sum_n^N a_n \mathbf{x}_n^T \mathbf{x} = \sum_n^N a_n k(\mathbf{x}_n, \mathbf{x})$$

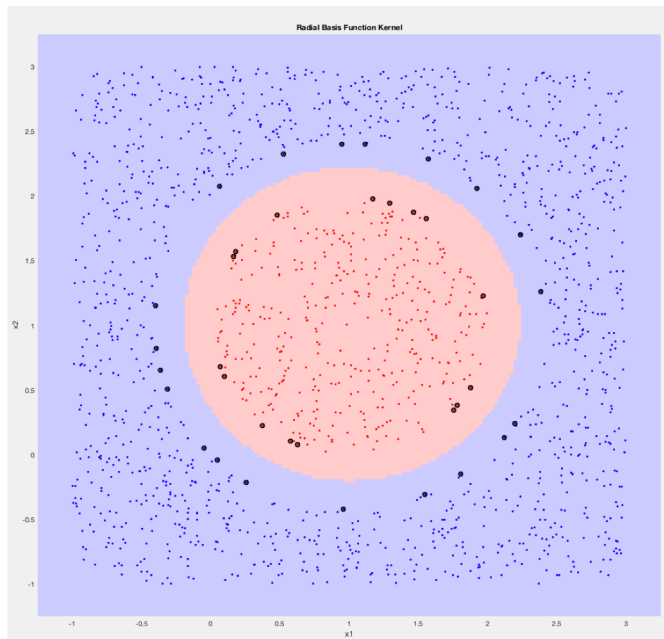
Dual representation, Sec 6.2

Prediction

$$y = \mathbf{w}^T \mathbf{x} = \mathbf{a}^T \mathbf{X} \mathbf{x} = \sum_n^N a_n \mathbf{x}_n^T \mathbf{x} = \sum_n^N a_n k(\mathbf{x}_n, \mathbf{x})$$

- Often **a is sparse** (... Support vector machines)
- We don't need to know **x** or $\varphi(\mathbf{x})$. ***Just the Kernel***

$$E(\mathbf{a}) = \|\mathbf{K} \mathbf{a} - \mathbf{t}\|_2^2 + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$



Gaussian Kernels

- Gaussian Kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{1}{2} (\mathbf{x} - \mathbf{x}')^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \mathbf{x}') \right)$$

Diagonal $\boldsymbol{\Sigma}$: (this gives ARD)

$$k(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{1}{2} \sum_i^N \frac{(x_i - x'_i)^2}{\sigma_i^2} \right)$$

Isotropic σ_i^2 gives an RBF


$$k(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2} \right)$$

Can be inner product in **infinite** dimensional space

Assume $x \in R^1$ and $\gamma > 0$.

$$\begin{aligned} e^{-\gamma\|x_i-x_j\|^2} &= e^{-\gamma(x_i-x_j)^2} = e^{-\gamma x_i^2 + 2\gamma x_i x_j - \gamma x_j^2} \\ &= e^{-\gamma x_i^2 - \gamma x_j^2} \left(1 + \frac{2\gamma x_i x_j}{1!} + \frac{(2\gamma x_i x_j)^2}{2!} + \frac{(2\gamma x_i x_j)^3}{3!} + \dots \right) \\ &= e^{-\gamma x_i^2 - \gamma x_j^2} \left(1 \cdot 1 + \sqrt{\frac{2\gamma}{1!}} x_i \cdot \sqrt{\frac{2\gamma}{1!}} x_j + \sqrt{\frac{(2\gamma)^2}{2!}} x_i^2 \cdot \sqrt{\frac{(2\gamma)^2}{2!}} x_j^2 \right. \\ &\quad \left. + \sqrt{\frac{(2\gamma)^3}{3!}} x_i^3 \cdot \sqrt{\frac{(2\gamma)^3}{3!}} x_j^3 + \dots \right) = \phi(x_i)^T \phi(x_j), \end{aligned}$$

where

$$\phi(x) = e^{-\gamma x^2} \left[1, \sqrt{\frac{2\gamma}{1!}} x, \sqrt{\frac{(2\gamma)^2}{2!}} x^2, \sqrt{\frac{(2\gamma)^3}{3!}} x^3, \dots \right]^T.$$


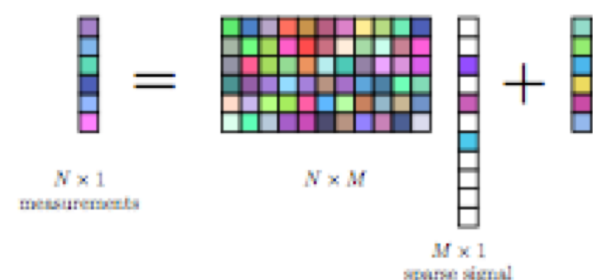
Sparse Bayesian Learning (SBL)

$$\text{Model : } \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}$$

$$\text{Prior : } \mathbf{x} \sim \mathcal{N}(\mathbf{x}; 0, \mathbf{\Gamma})$$

$$\mathbf{\Gamma} = \text{diag}(\gamma_1, \dots, \gamma_M)$$

$$\text{Likelihood : } p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \sigma^2 \mathbf{I}_N)$$



$$\text{Evidence : } p(\mathbf{y}) = \int_{\mathbf{x}} p(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x} = \mathcal{N}(\mathbf{y}; 0, \mathbf{\Sigma}_y)$$
$$\mathbf{\Sigma}_y = \sigma^2 \mathbf{I}_N + \mathbf{A}\mathbf{\Gamma}\mathbf{A}^H$$

$$\text{SBL solution : } \hat{\mathbf{\Gamma}} = \arg \max_{\mathbf{\Gamma}} p(\mathbf{y})$$
$$= \arg \min_{\mathbf{\Gamma}} \{ \log |\mathbf{\Sigma}_y| + \mathbf{y}^H \mathbf{\Sigma}_y^{-1} \mathbf{y} \}$$

M.E.Tipping, "Sparse Bayesian learning and the relevance vector machine," Journal of Machine Learning Research, June 2001.

Lecture 10

Support Vector Machines

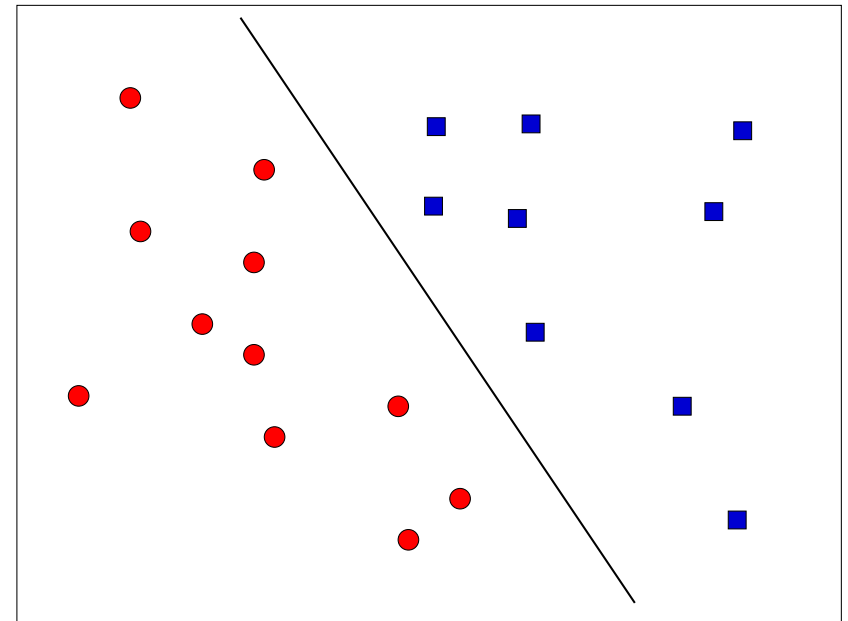
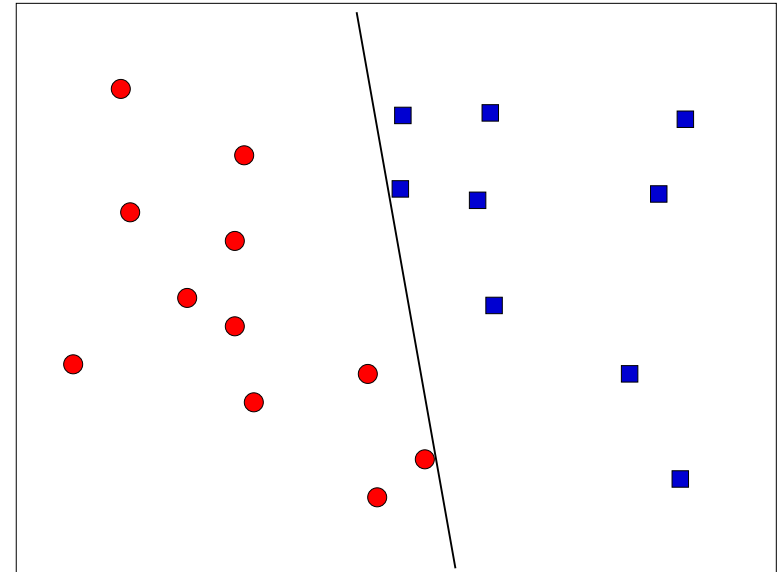
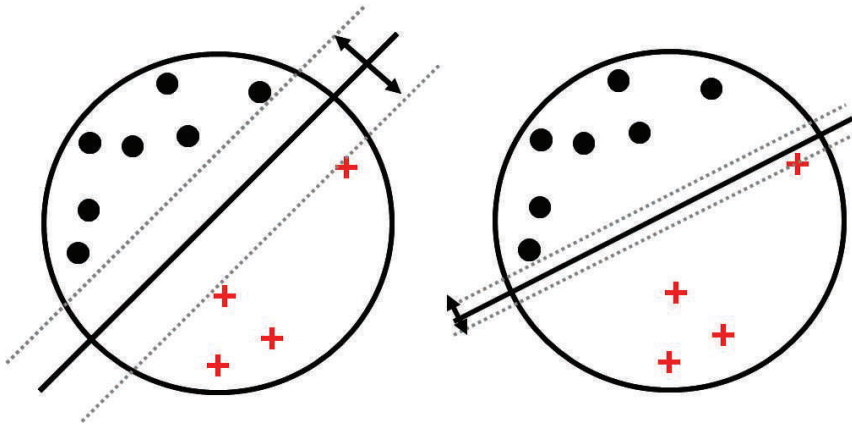
Non Bayesian!

Features:

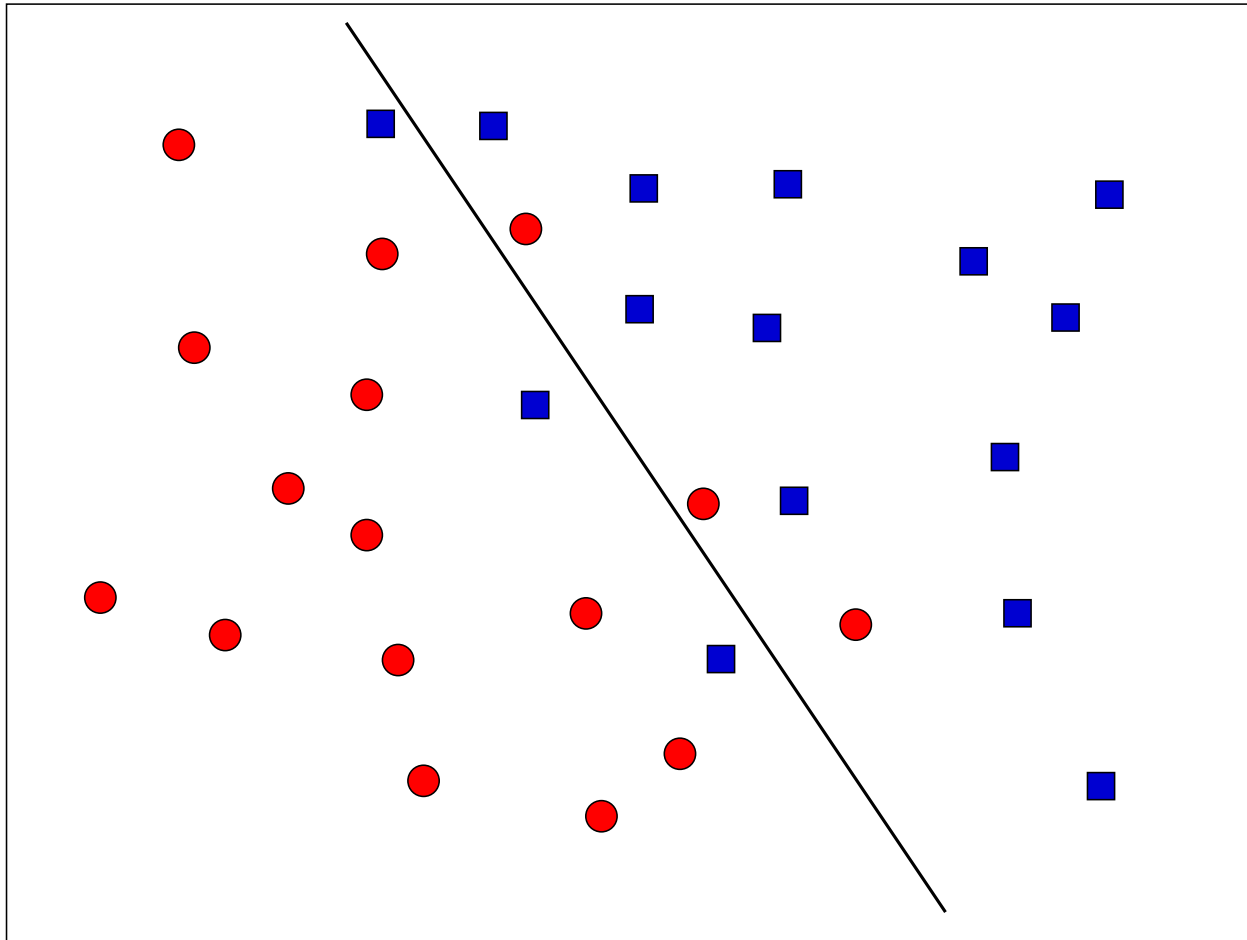
- Kernel
- Sparse representations
- Large margins

Regularize for plausibility

- Which one is best?
- We maximize the margin

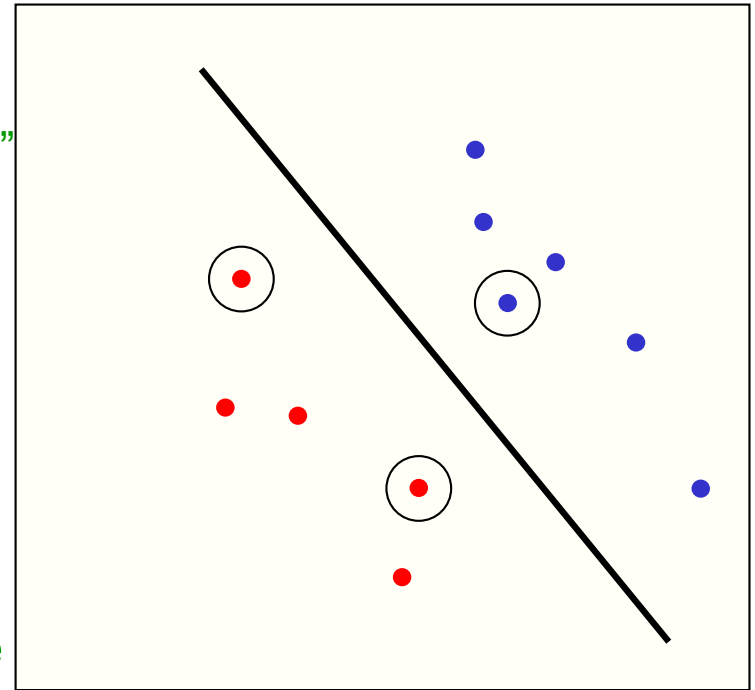


Regularize for plausibility



Support Vector Machines

- The line that maximizes the minimum margin is a good bet.
 - The model class of “hyper-planes with a margin m ” has a low VC dimension if m is big.
- This maximum-margin separator is determined by a subset of the datapoints.
 - Datapoints in this subset are called “support vectors”.
 - It is useful computationally if only few datapoints are support vectors, because the support vectors decide which side of the separator a test case is on.



The support vectors are indicated by the circles around them.

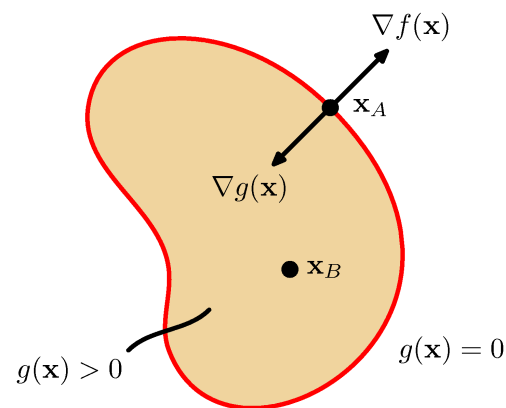
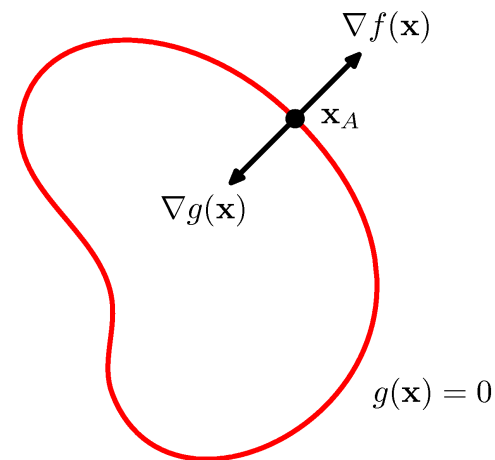
Lagrange multiplier (Bishop App E)

$$\max(f(x)) \text{ subject to } g(x) = 0$$

Taylor expansion

$$g(\mathbf{x} + \boldsymbol{\epsilon}) = g(\mathbf{x}) + \boldsymbol{\epsilon}^T \nabla g(\mathbf{x})$$

$$L(x, \lambda) = f(x) + \lambda g(x)$$



Lagrange multiplier (Bishop App E)

$\max(f(\mathbf{x}))$ subject to $g(\mathbf{x}) > 0$

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

Either $\nabla f(\mathbf{x}) = \mathbf{0}$

Then $g(\mathbf{x})$ is inactive, $\lambda=0$

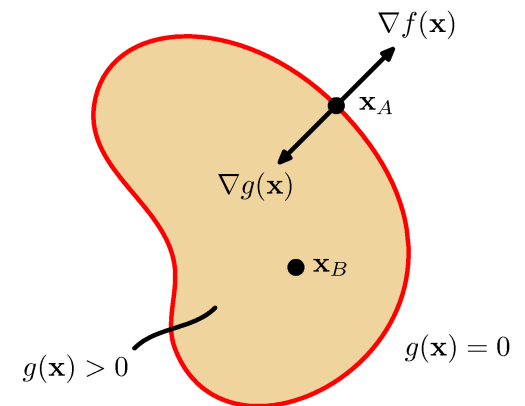
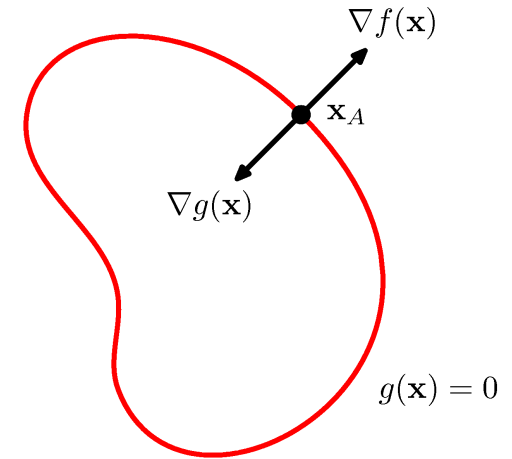
Or $g(\mathbf{x}) = 0$ but $\lambda > 0$

Thus optimizing $L(\mathbf{x}, \lambda)$ with the Karush-Kuhn-Trucker (KKT) equations

$$g(\mathbf{x}) \geq 0$$

$$\lambda \geq 0$$

$$\lambda g(\mathbf{x}) = 0$$



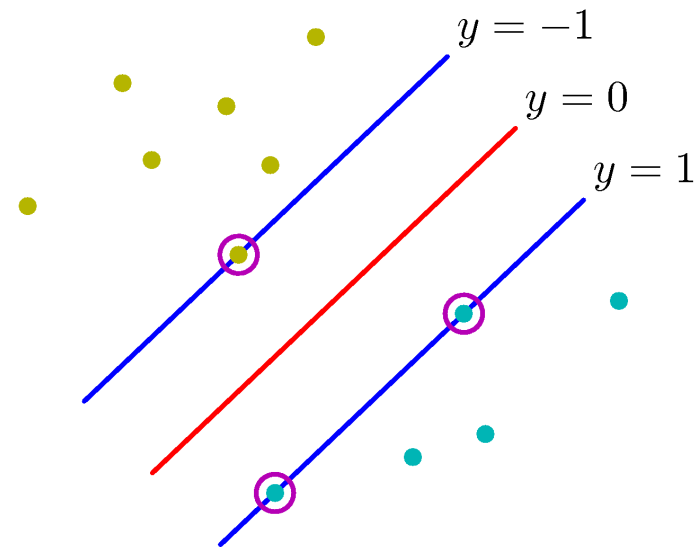
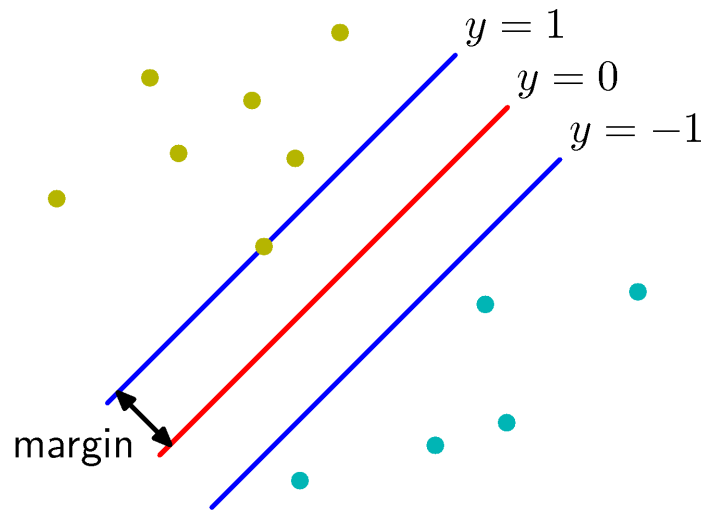
Testing a linear SVM

- The separator is defined as the set of points for which:

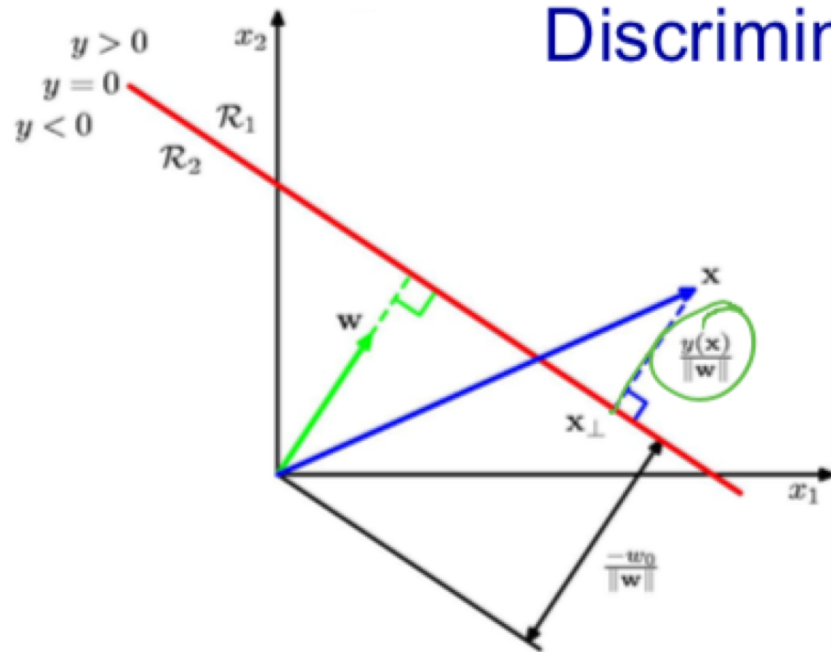
$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

so if $\mathbf{w} \cdot \mathbf{x}^c + b > 0$ say its a positive case

and if $\mathbf{w} \cdot \mathbf{x}^c + b < 0$ say its a negative case



Discriminant functions



The planar decision surface in data-space for the simple linear discriminant function:

$$\mathbf{w}^T \mathbf{x} + w_0 \geq 0$$

$$y = \mathbf{w}^T \mathbf{x} + w_0$$

$$y_{\perp} = \mathbf{w}^T \mathbf{x}_{\perp} + w_0 = 0 \quad \checkmark$$

$$\|\mathbf{w}\|_2^2 = \mathbf{w}^T \mathbf{w}$$

\mathbf{x} on plane $\Rightarrow y=0 \Rightarrow$

Distance from plane

$$r = \frac{y}{\|\mathbf{w}\|_2}$$

$$\begin{aligned} \mathbf{x} &= \mathbf{x}_{\perp} + r \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \\ \underline{\mathbf{w}^T \mathbf{x}} &= \mathbf{w}^T \mathbf{x}_{\perp} + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|_2} \\ &= -w_0 + r \|\mathbf{w}\|_2 \\ \mathbf{x} &= \mathbf{x}_{\perp} + \frac{y}{\|\mathbf{w}\|_2} \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \end{aligned}$$

Large margin

$$y = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$$

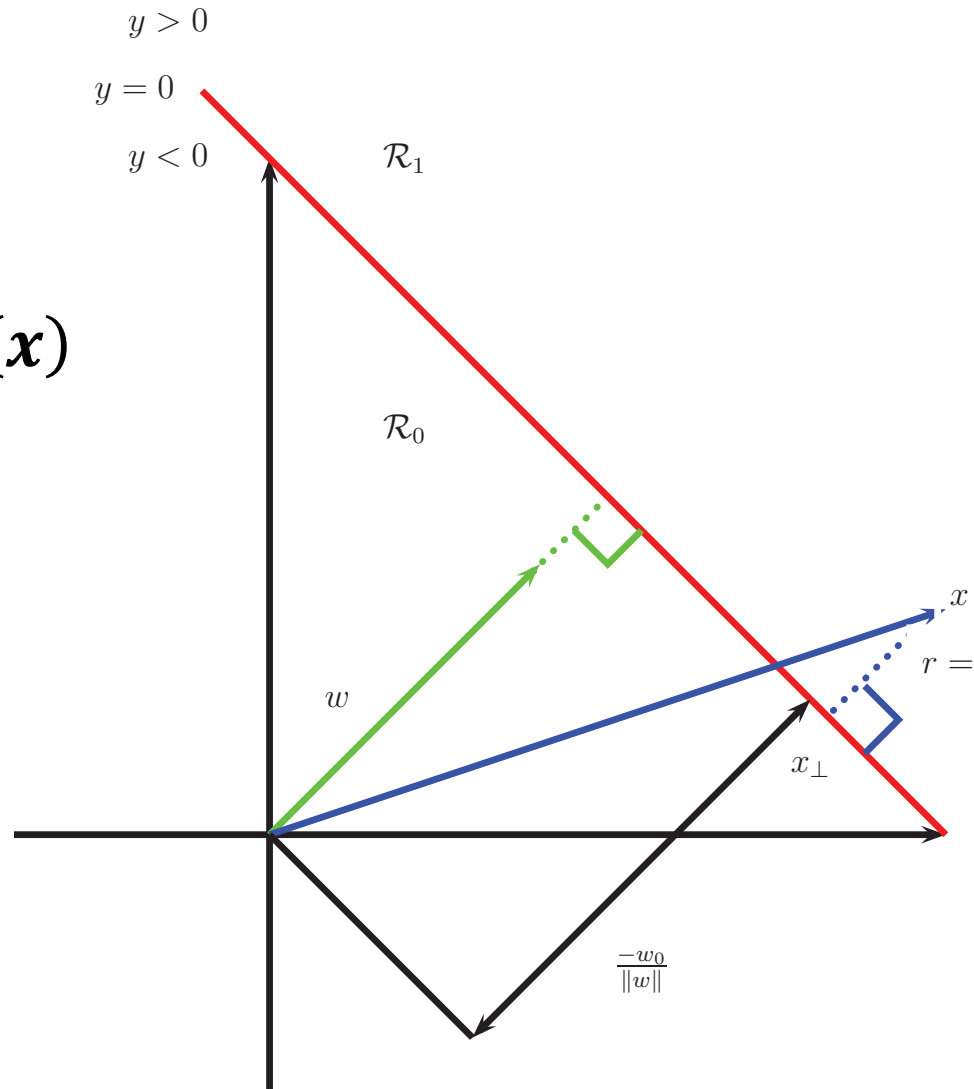
$$\mathbf{x}_n = \mathbf{x}_\perp + r_n \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

$$\mathbf{x} \text{ on plane } \Rightarrow y=0 \Rightarrow b = -\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

$$r_n = \frac{\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b}{\|\mathbf{w}\|} = \frac{y_n}{\|\mathbf{w}\|}$$

$$t_n y_n \geq 1$$

$$\max_{\mathbf{w}} \frac{1}{\|\mathbf{w}\|} \min_n t_n y_n$$



Maximum margin (Bishop 7.1)

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{Subject to} \quad t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1, \quad n = 1, \dots, N. \quad (7.5)$$

Lagrange function

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1\} \quad (7.7)$$

Differentiation

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \quad (7.8)$$

$$0 = \sum_{n=1}^N a_n t_n. \quad (7.9)$$

Dual representation

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \quad (7.10)$$

with respect to \mathbf{a} subject to the constraints

$$a_n \geq 0, \quad n = 1, \dots, N, \quad (7.11)$$

$$\sum_{n=1}^N a_n t_n = 0. \quad (7.12)$$

This can be solved with quadratic programming

Maximum margin (Bishop 7.1)

- KKT conditions

$$a_n \geq 0 \quad (7.14)$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0 \quad (7.15)$$

$$a_n \{t_n y(\mathbf{x}_n) - 1\} = 0. \quad (7.16)$$

either $a_n = 0$ or $t_n y(\mathbf{x}_n) = 1$.

- Solving for a_n

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \quad (7.8)$$

- Prediction

$$y(\mathbf{x}) = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b. \quad (7.13)$$

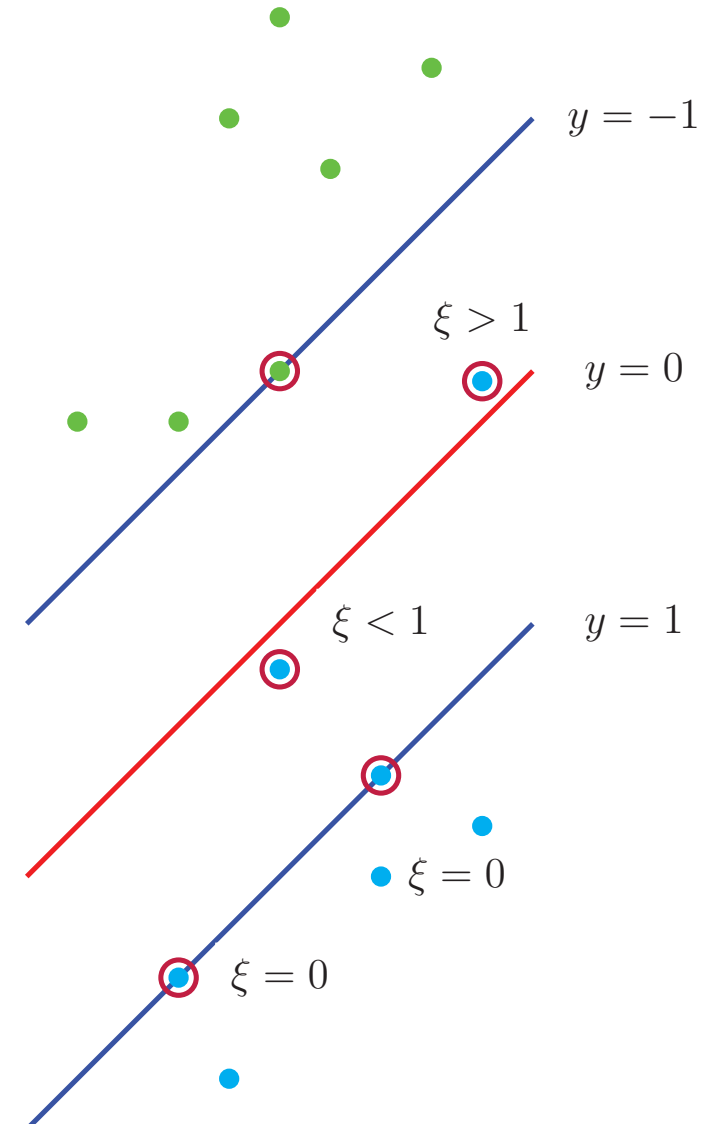
If there is no separating plane...

- Use a bigger set of features.
 - Makes the computation slow? “Kernel” trick makes the computation fast with many features.
- Extend definition of maximum margin to allow non-separating planes.
 - Use “slack” variables $\xi = |t_n - y(x_n)|$

$$t_n y(\mathbf{x}_n) \geq 1 - \xi_n, \quad n = 1, \dots, N \quad (7.20)$$

Objective function

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2 \quad (7.21)$$



SVM classification summarized--- Only kernels

- Minimize with respect to \mathbf{w}, w_0

$$C \sum_n^N \zeta_n + \frac{1}{2} \|\mathbf{w}\|^2 \quad (\text{Bishop 7.21})$$

- Solution found in dual domain with Lagrange multipliers
 - $a_n, n = 1 \cdots N$ and

- This gives the support vectors S

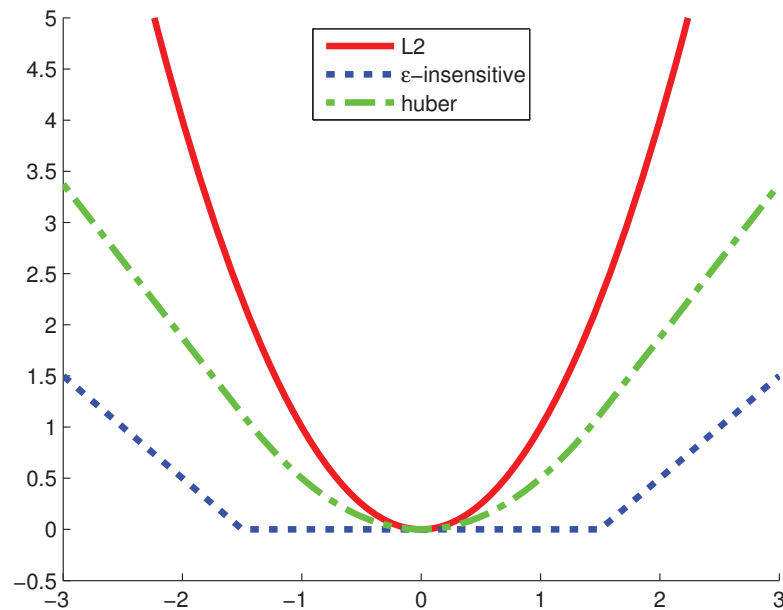
$$\hat{\mathbf{w}} = \sum_{n \in S} a_n t_n \boldsymbol{\varphi}(x_n) \quad (\text{Bishop 7.8})$$

- Used for predictions

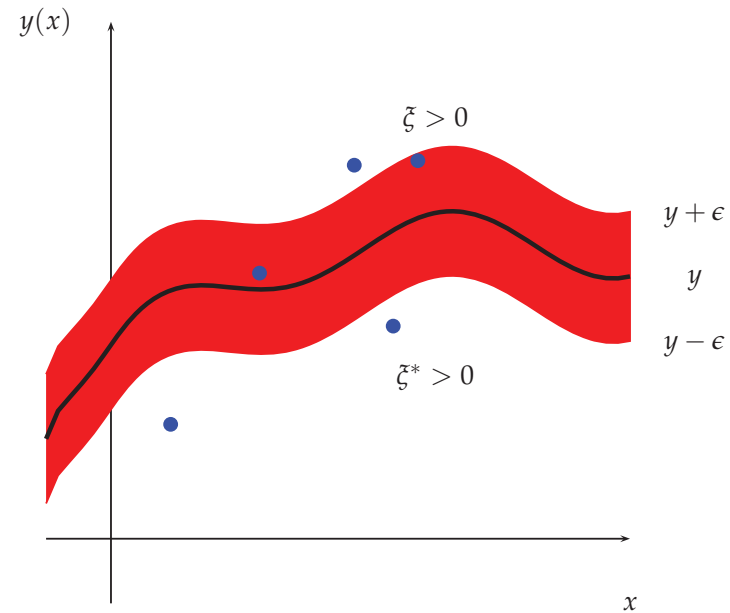
$$\hat{y} = w_0 + \mathbf{w}^T \boldsymbol{\varphi}(x) = w_0 + \sum_{n \in S} a_n t_n \boldsymbol{\varphi}(x_n)^T \boldsymbol{\varphi}(x)$$

$$= w_0 + \sum_{n \in S} a_n t_n k(x_n, x) \quad (\text{Bishop 7.13})$$

SVM for regression



(a)



(b)

Figure 14.10 (a) Illustration of ℓ_2 , Huber and ϵ -insensitive loss functions, where $\epsilon = 1.5$. Figure generated by `huberLossDemo`. (b) Illustration of the ϵ -tube used in SVM regression. Points above the tube have $\xi_i > 0$ and $\xi_i^* = 0$. Points below the tube have $\xi_i = 0$ and $\xi_i^* > 0$. Points inside the tube have $\xi_i = \xi_i^* = 0$. Based on Figure 7.7 of (Bishop 2006a).

SVMs are Perceptrons!

- SVM's use each training case, x , to define a feature $K(x, \cdot)$ where K is user chosen.
 - So the user designs the features.
- SVM do “feature selection” by picking support vectors, and learn feature weighting from a big optimization problem.
- \Rightarrow SVM is a clever way to train a standard perceptron.
 - What a perceptron cannot do, SVM cannot do.
- SVM DOES:
 - Margin maximization
 - Kernel trick
 - Sparse

SVM Code for classification (libsvm)

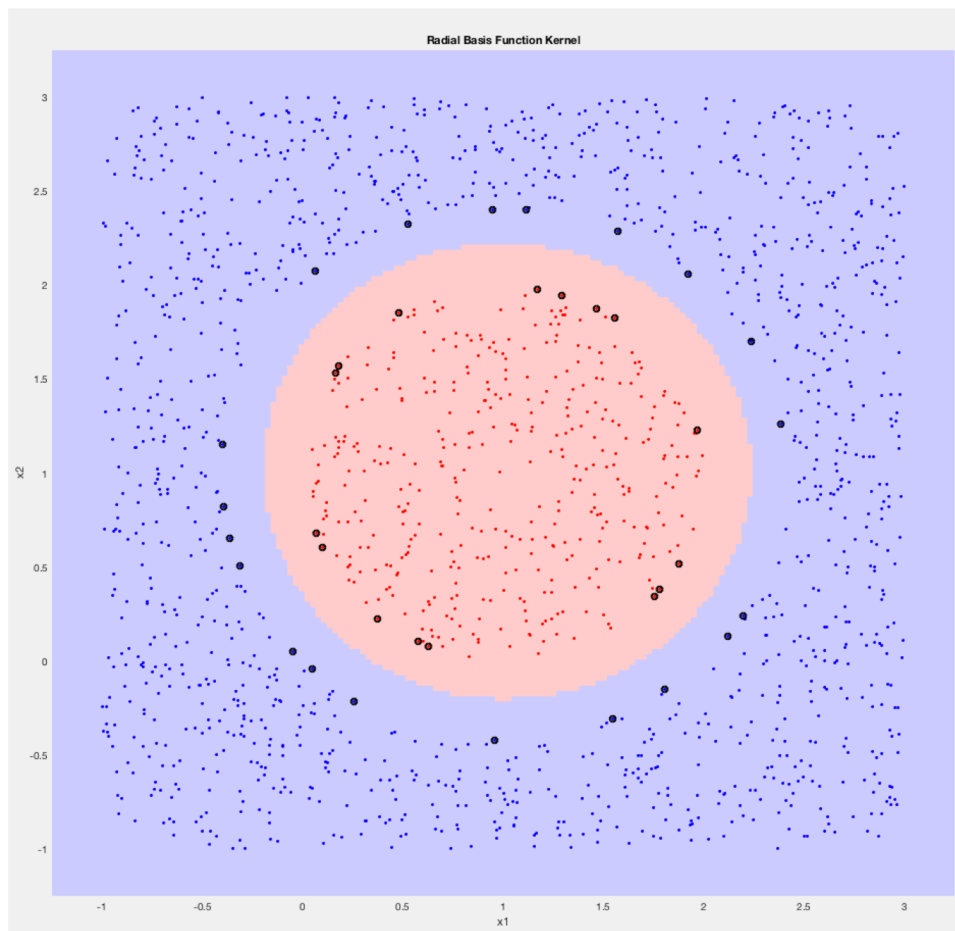
Part of ocean acoustic data set <http://noiselab.ucsd.edu/ECE285/SIO209Final.zip>
case 'Classify'

```
% train
```

```
model = svmtrain(Y, X,['-c 7.46 -g ' gamma ' -q ' kernel]);
```

```
% predict
```

```
[predict_label,~, ~] = svmpredict(rand([length(Y),1]), X, model,'-q');
```



```
>> modelmodel = struct with fields:  
Parameters: [5×1 double]  
nr_class: 2  
totalSV: 36  
rho: 8.3220  
Label: [2×1 double]  
sv_indices: [36×1 double]  
ProbA: [] ProbB: []  
nSV: [2×1 double]  
sv_coef: [36×1 double]  
SVs: [36×2 double]
```

Finding the Decision Function

libsvm

- \mathbf{w} : maybe **infinite** variables
- The **dual** problem

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C, i = 1, \dots, l \\ & \mathbf{y}^T \alpha = 0, \end{aligned}$$

Corresponds to
(Bishop 7.32)
With $\mathbf{y}=\mathbf{t}$

where $Q_{ij} = y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ and $\mathbf{e} = [1, \dots, 1]^T$

- At optimum

$$\mathbf{w} = \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i)$$

- A **finite** problem: #variables = #training data

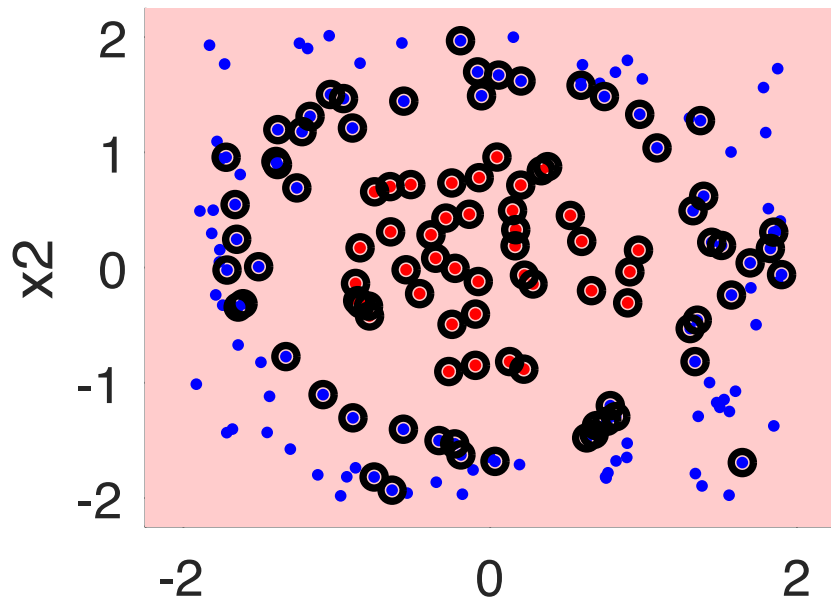


Navigation icons: back, forward, search, etc.

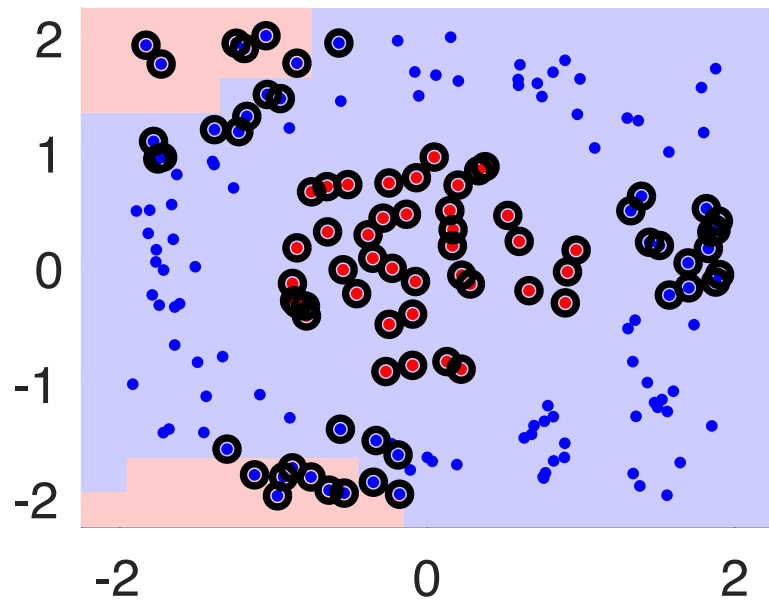
Using these results to eliminate \mathbf{w} , b , and $\{\xi_n\}$ from the Lagrangian, we obtain the dual Lagrangian in the form

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \quad (7.32)$$

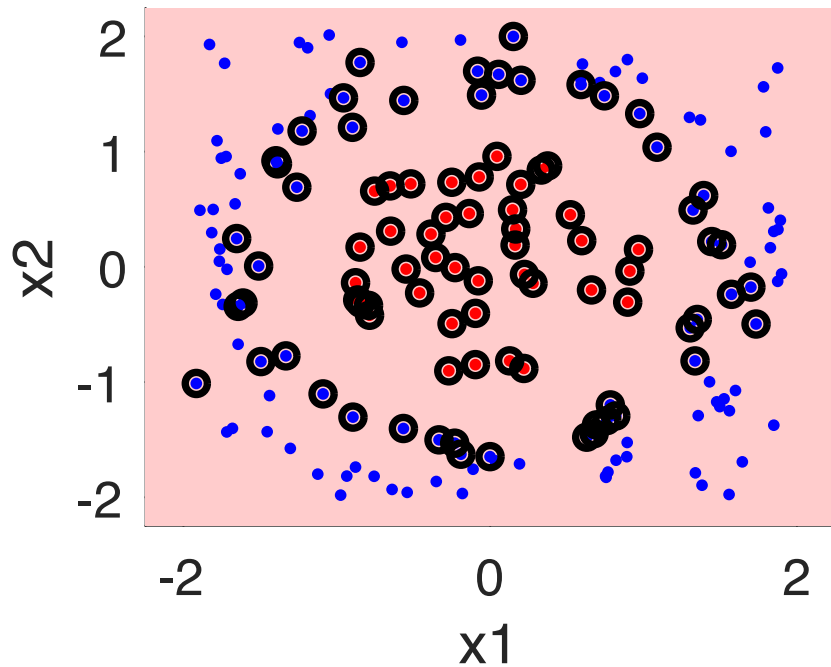
Linear Kernel



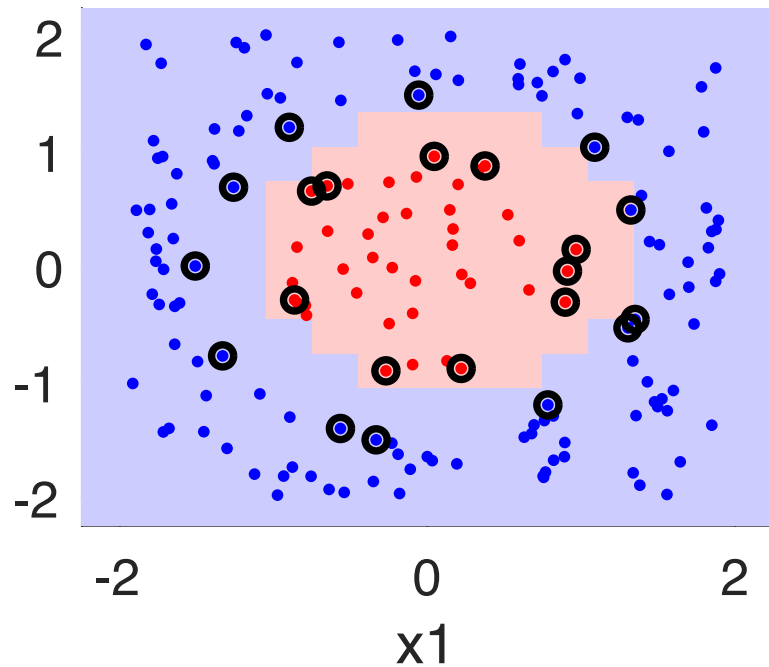
Sigmoid Function Kernel



Polynomial Kernel

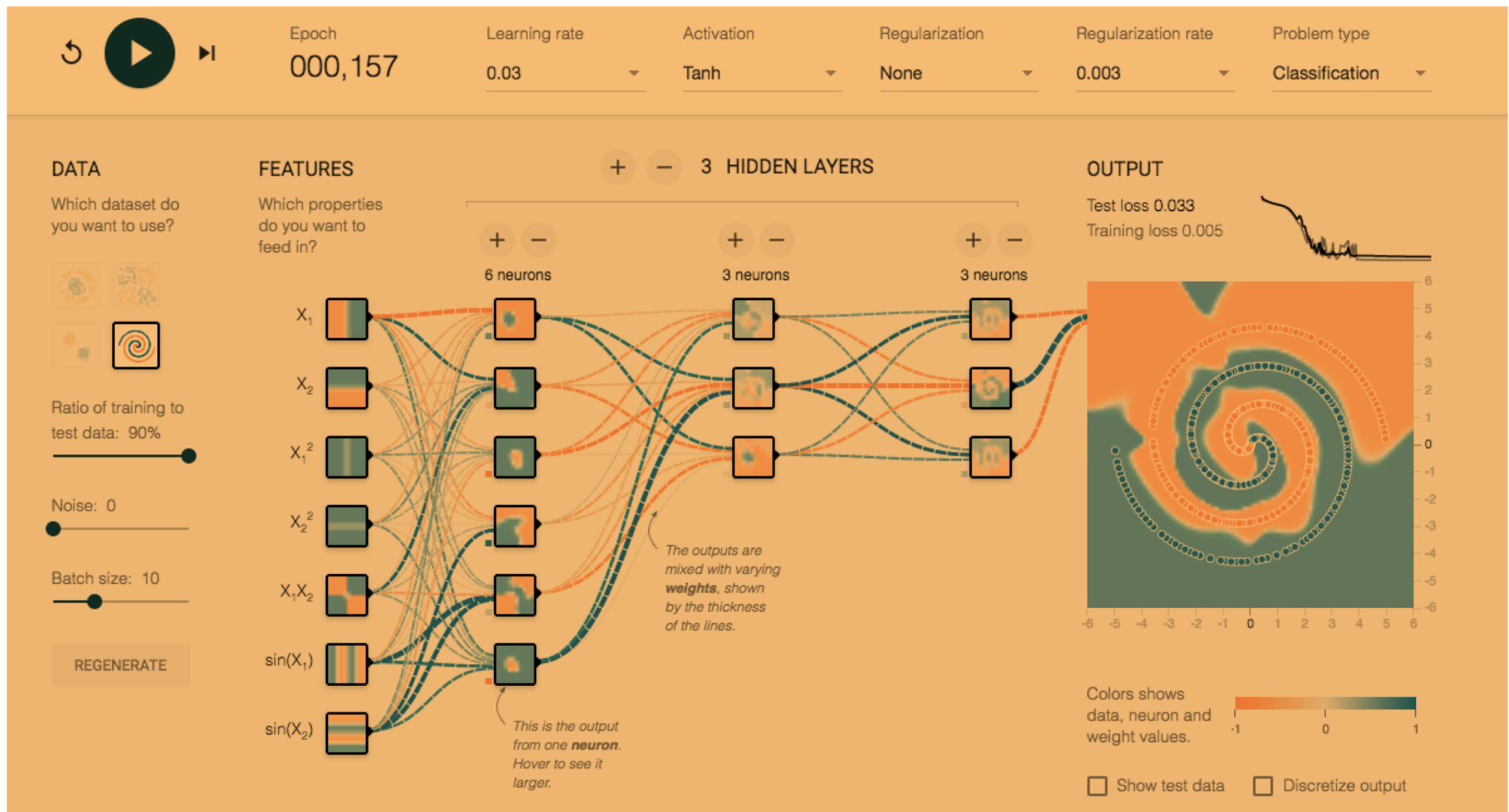


Radial Basis Function Kernel



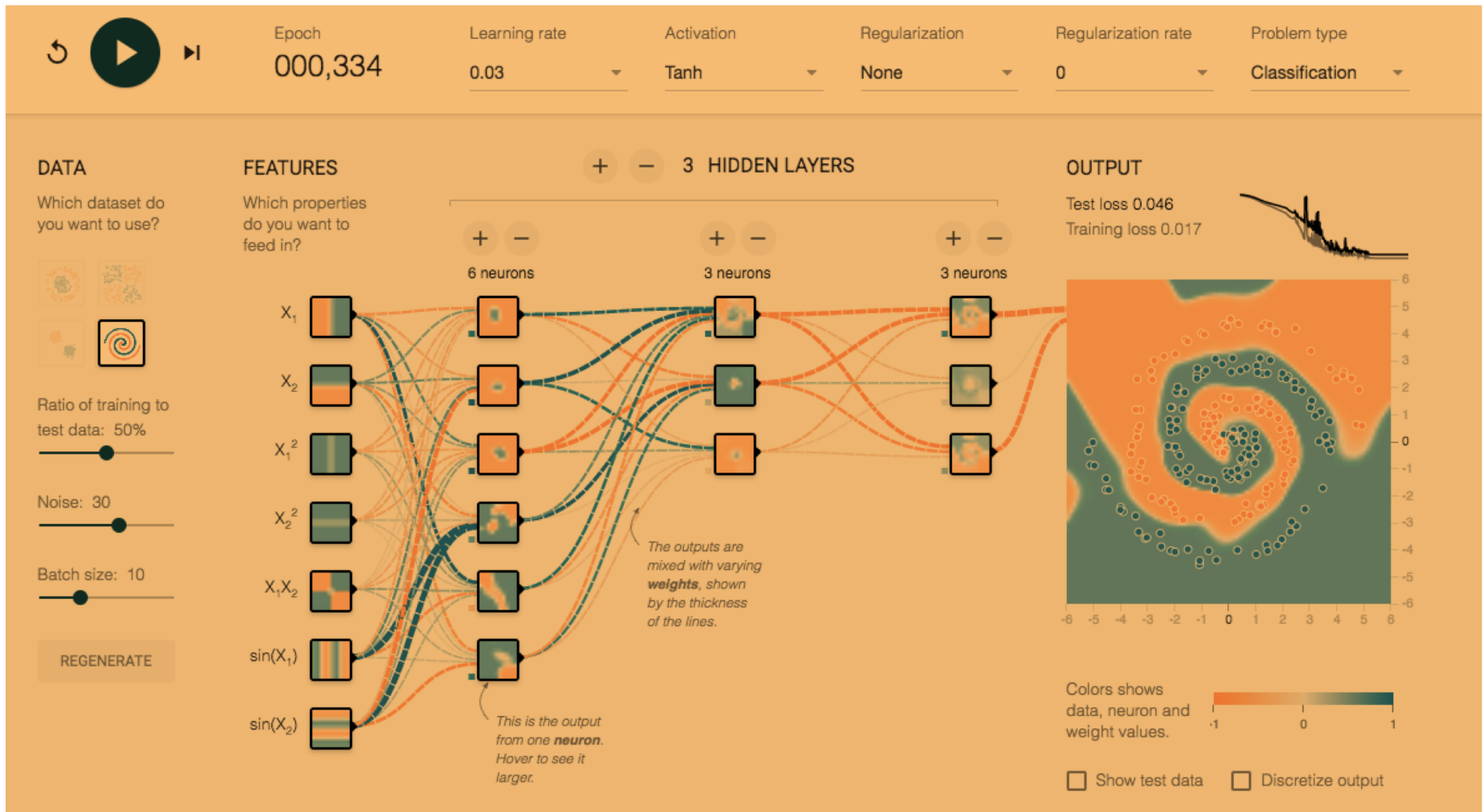
Tensorflow Playground

1. Fitting the spiral with default settings fail due to the small training set. The NN will fit to the training data which is not representative of the true pattern and the network will **generalize** poorly. Increasing the ratio of training to test data to 90% the NN finds the correct shape (1st image).



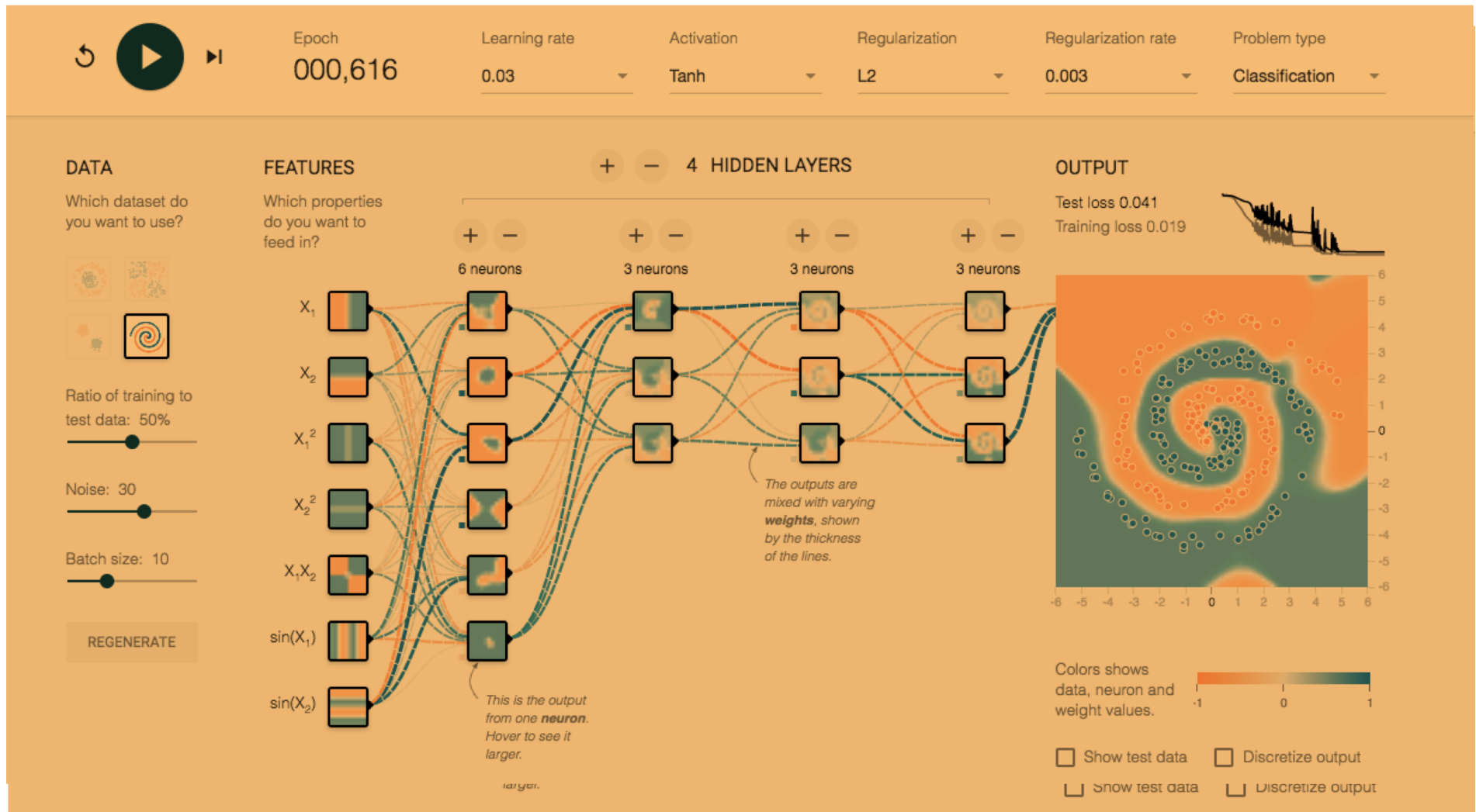
Tensorflow Playground

You can fix the generalization problem by adding **noise** to the data. This allows the small training set to generalize better as it reduce **overfitting** of the training data (2nd image).



Tensorflow Playground

Adding an additional hidden layer the NN fails to classify the shape properly. **Overfitting** once again becomes a problem even after you've added noise. This can be fixed by adding appropriate **L2 regularization** (third image).



- **NOT USED**