### Kernels

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

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

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

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



Example 4:

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

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

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

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: 
$$\min_{\mathbf{w}} E(\mathbf{w})$$
  
 $E = \frac{1}{2} \sum_{n}^{N} \{\mathbf{w}^{T} \mathbf{x}_{n} - t_{n}\}^{2} + \frac{\lambda}{2} \|\mathbf{w}\|^{2} = \|\mathbf{X}\mathbf{w} - \mathbf{t}\|_{2}^{2} + \frac{\lambda}{2} \|\mathbf{w}\|^{2}$ 

Solution 
$$w = X^+ t = (X^T X + \lambda I_M)^{-1} X^T t$$
  
=  $X^T (XX^T + \lambda I_N)^{-1} t = X^T (K + \lambda I_N)^{-1} t = X^T a$ 

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

Dual representation is : 
$$\min_{a} E(a)$$
$$\mathcal{L} \in \mathcal{R}^{\mathcal{N}}$$
$$\mathcal{L} = \frac{1}{2} \sum_{n=1}^{N} \{ \mathbf{w}^{T} \mathbf{x}_{n} - t_{n} \}^{2} + \frac{\lambda}{2} \| \mathbf{w} \|^{2} = \| \mathbf{K} \mathbf{a} - \mathbf{t} \|_{2}^{2} + \frac{\lambda}{2} \mathbf{a}^{T} \mathbf{K} \mathbf{a}$$

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  $\varphi(x)$ . Just the Kernel  $E(a) = ||Ka - t||_2^2 + \frac{\lambda}{2}a^T Ka$



### **Gaussian Kernels**

Gaussian Kernel

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

Diagonal  $\Sigma$ : (this gives ARD)

$$k(x, x') = \exp\left(-\frac{1}{2}\sum_{i}^{N}\frac{\left(x_{i} - x_{i}'\right)^{2}}{\sigma_{i}^{2}}\right)$$

Isotropic  $\sigma_i^2$  gives an RBF

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

Can be inner product in infinite dimensional space Assume  $x \in R^1$  and  $\gamma > 0$ .

$$\underbrace{e^{-\gamma ||x_i - x_j||_{\iota}^2}}_{=e^{-\gamma (x_i - x_j)^2}} = e^{-\gamma x_i^2 + 2\gamma x_i x_j - \gamma x_j^2} = \underbrace{c^{-\gamma x_i^2}}_{=e^{-\gamma x_i^2 - \gamma x_j^2}} \left(1 + \frac{2\gamma x_i x_j}{1!} + \frac{(2\gamma x_i x_j)^2}{2!} + \frac{(2\gamma x_i x_j)^3}{3!} + \cdots\right) \stackrel{\text{(x)}}{=e^{-\gamma x_i^2 - \gamma x_j^2}} \left(1 \cdot 1 + \sqrt{\frac{2\gamma}{1!}} x_i \cdot \sqrt{\frac{2\gamma}{1!}} x_j + \sqrt{\frac{(2\gamma)^2}{2!}} x_i^2 \cdot \sqrt{\frac{(2\gamma)^2}{2!}} x_j^2 + \sqrt{\frac{(2\gamma)^2}{2!}} x_j^2 + \sqrt{\frac{(2\gamma)^2}{2!}} x_j^2 + \sqrt{\frac{(2\gamma)^3}{3!}} x_i^3 \cdot \sqrt{\frac{(2\gamma)^3}{3!}} x_j^3 + \cdots\right) = \phi(x_i)^T \phi(x_j),$$

where



#### Sparse Bayesian Learning (SBL)

 $\begin{array}{l} \mathsf{Model} : \ \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n} \\ \mathsf{Prior} : \ \mathbf{x} \sim \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{\Gamma}) \\ \mathbf{\Gamma} = \mathsf{diag}(\gamma_1, \dots, \gamma_M) \\ \mathsf{Likelihood} : \ p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \sigma^2 \mathbf{I}_N) \end{array}$ 



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

$$\begin{array}{l} \mathsf{SBL} \ \mathsf{solution} : \ \hat{\boldsymbol{\Gamma}} = \mathop{\arg\max}\limits_{\boldsymbol{\Gamma}} \, p(\mathbf{y}) \\ = \mathop{\arg\min}\limits_{\boldsymbol{\Gamma}} \, \big\{ \log |\boldsymbol{\Sigma}_{\mathbf{y}}| + \mathbf{y}^{H} \boldsymbol{\Sigma}_{\mathbf{y}}^{-1} \mathbf{y} \big\} \end{array}$$

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)) subject to g(x) = 0







## Lagrange multiplier (Bishop App E)

$$\max(f(\mathbf{x})) \text{ subject to } g(\mathbf{x}) \ge 0$$
  
$$\implies L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

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  $L(x, \lambda)$  with the Karesh-Kuhn-Trucker (KKT) equations

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





#### Testing a linear SVM

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

 $\mathbf{w}.\mathbf{x}+b=0$ so if  $\mathbf{w}.\mathbf{x}^{c}+b>0$  say its a positive case and if  $\mathbf{w}.\mathbf{x}^{c}+b<0$  say its a negative case





## The planar decision surface in data-space for the simple

linear discriminant function:  $\mathbf{w}^T \mathbf{x} + w_0 \ge 0$ 

= WTX+W.

 $y_1 = v x_1 + v y_2 = 0$   $||v x||_2^2 = v x_1 + v y_2$ 

X on plane => y=0 =>

 $X = X_1 + r \frac{1}{100 \text{ M}}$   $W = W \overline{Y} \times \tau r W \overline{W}$ Distance from plane  $X = X_{i} + \frac{-1}{4}$ 

Large margin  $y = \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}) + b$ y > 0 $x_n = x_\perp + r_n \frac{w}{\|w\|}$ y = 0y < 0th x **x** on plane => y=0 => $b = -w^T \phi(x)$  $\mathcal{R}_0$  $\underline{r_n} = \frac{\boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x_n}) + b}{\|\boldsymbol{w}\|} = \frac{\underline{y_n}}{\|\boldsymbol{w}\|}$ w $x_{\perp}$  $t_n y_n \ge 1$   $\checkmark$  $\frac{-w_0}{\|w\|}$  $\max_{\boldsymbol{w}} \frac{-}{\|\boldsymbol{w}\|} \min_{n} t_n y_n$ 

## Maximum margin (Bishop 7.1) Subject to $t_{p} \left( \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_{n}) + b \right) \ge 1, \qquad n = 1, \dots, N.$ (7.5)

Lagrange function  $L(\mathbf{v})$ 

 $\underset{\mathbf{w},b}{\operatorname{arg\,min}} \frac{1}{2} \|\mathbf{w}\|^2$ 

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

Differentiation

$$\mathbf{w} = \sum_{n=1}^{N} a_n t_n \phi(\mathbf{x}_n)$$
(7.8)  
$$0 = \sum_{n=1}^{N} a_n t_n.$$
(7.9)

**Dual representation** 

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

n=1

with respect to a subject to the constraints

$$a_n \ge 0, \qquad n = 1, \dots, N,$$
 (7.11)

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

### This can be solved with quadratic programming

## Maximum margin (Bishop 7.1)

• KKT conditions

$$a_n \ge 0$$
 (7.14)

$$t_n y(\mathbf{x}_n) - 1 \ge 0 \tag{7.15}$$

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

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

• Solving for a<sub>n</sub>

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

• Prediction

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

## If there is no separating plane...

y = -1

y = 0

y = 1

 $\xi >$ 

 $\mathcal{E} = 0$ 

 $\xi < 1$ 

 $\bigcirc$ 

- 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 
$$\xi = |t_n - y(x_n)|$$

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

# **Objective function** $C\sum_{n=1}^{\infty} \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$

(7.21)

C-700

## SVM classification summarized--- Only kernels

• Minimize with respect to w, w<sub>0</sub>

 $C\sum_{n=1}^{N}\zeta_{n} + \frac{1}{2}\|\boldsymbol{w}\|^{2} \qquad (\text{Bishop 7.21}) \leq 1$ 

Solution found in dual domain with Lagrange multipliers

 $-a_n$  ,  $n=1\cdots N$  and

• This gives the support vectors S

$$\widehat{\boldsymbol{w}} = \sum_{n \in S} a_n t_n \boldsymbol{\varphi}(xn)$$

(Bishop 7.8)

3)

Used for predictions

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

$$= w_0 + \sum_{n \in S} a_n t_n k(x_n, x)$$
 (Bishop 7.1)

### SVM for regression



**Figure 14.10** (a) Illustration of  $\ell_2$ , Huber and  $\epsilon$ -insensitive loss functions, where  $\epsilon = 1.5$ . Figure generated by huberLossDemo. (b) Illustration of the  $\epsilon$ -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, .) 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.
- =>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 <a href="http://noiselab.ucsd.edu/ECE285/SIO209Final.zip">http://noiselab.ucsd.edu/ECE285/SIO209Final.zip</a> 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

- w: maybe infinite variables
- The dual problem

• At

$$\begin{array}{ll} \min_{\boldsymbol{\alpha}} & \frac{1}{2} \boldsymbol{\alpha}^{T} Q \boldsymbol{\alpha} - \mathbf{e}^{T} \boldsymbol{\alpha} \end{array} \right] & \begin{array}{l} \text{Corresponds to} \\ \text{subject to} & 0 \leq \alpha_{i} \leq C, i = 1, \ldots, I \\ \mathbf{y}^{T} \boldsymbol{\alpha} = 0, \end{array} \\ \text{where } Q_{ij} = y_{i} y_{j} \phi(\mathbf{x}_{i})^{T} \phi(\mathbf{x}_{j}) \text{ and } \mathbf{e} = [1, \ldots, 1]^{T} \end{array} \\ \begin{array}{l} \text{With } \mathbf{y} = \mathbf{t} \\ \text{With } \mathbf{y} = \mathbf{t} \\ \text{At optimum} \end{array}$$

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

• A finite problem: #variables = #training data

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

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



#### **Polynomial Kernel**



#### **Sigmoid Function Kernel**



#### **Radial Basis Function Kernel**



## **Tensorflow Playground**

 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 (1<sup>st</sup> 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).

