Kernel Methods for Learning

John Shawe-Taylor

Department of Computer Science University College London jst@cs.ucl.ac.uk

International Summer School on Machine Learning Benicasim, Spain

June 2013

Aim:

The course is intended to give an overview of the kernel approach to pattern analysis. This will cover:

- Why linear pattern functions?
- Why kernel approach?
- How to plug and play with the different components of a kernel-based pattern analysis system?

1

What won't be included:

- Other approaches to Pattern Analysis
- Complete History
- Bayesian view of kernel methods
- Most recent developments

OVERALL STRUCTURE

- Part 1: Introduction to the Kernel methods approach.
- Part 2: Projections and subspaces in the feature space.
- **Part 3:** Stability of Pattern Functions with the example of Support Vector Machines.
- **Part 4:** Other learning algorithms: novelty detection, boosting and multiple kernel learning.
- Part 5: Kernel design strategies.

Part 1

- Kernel methods approach
- Worked example of kernel Ridge Regression
- Properties of kernels.

Kernel methods

Kernel methods (re)introduced in 1990s with Support Vector Machines

- Linear functions but in high dimensional spaces equivalent to non-linear functions in the input space
- Statistical analysis showing large margin can overcome curse of dimensionality
- Extensions rapidly introduced for many other tasks other than classification

Kernel methods approach

- Data embedded into a Euclidean feature (or Hilbert) space
- Linear relations are sought among the images of the data
- Algorithms implemented so that only require inner products between vectors
- Embedding designed so that inner products of images of two points can be computed directly by an efficient 'short-cut' known as the kernel.

Worked example: Ridge Regression

Consider the problem of finding a homogeneous real-valued linear function

$$g(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle = \mathbf{x}' \mathbf{w} = \sum_{i=1}^{n} w_i x_i,$$

that best interpolates a given training set

$$S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$$

of points \mathbf{x}_i from $X \subseteq \mathbb{R}^n$ with corresponding labels y_i in $Y \subseteq \mathbb{R}$.

Possible pattern function

 Measures discrepancy between function output and correct output – squared to ensure always positive:

 $f_g((\mathbf{x}, y)) = (g(\mathbf{x}) - y)^2$

Note that the pattern function f_g is not itself a linear function, but a simple functional of the linear functions g.

• We introduce notation: matrix \mathbf{X} has rows the m examples of S. Hence we can write

 $\xi = \mathbf{y} - \mathbf{X}\mathbf{w}$

for the vector of differences between $g(\mathbf{x}_i)$ and y_i .

Optimising the choice of g

Need to ensure flexibility of g is controlled – controlling the norm of w proves effective:

 $\min_{\mathbf{w}} \mathcal{L}_{\lambda}(\mathbf{w}, S) = \min_{\mathbf{w}} \lambda \|\mathbf{w}\|^2 + \|\xi\|^2,$

where we can compute

$$\|\xi\|^2 = \langle \mathbf{y} - \mathbf{X}\mathbf{w}, \mathbf{y} - \mathbf{X}\mathbf{w} \rangle$$

= $\mathbf{y}'\mathbf{y} - 2\mathbf{w}'\mathbf{X}'\mathbf{y} + \mathbf{w}'\mathbf{X}'\mathbf{X}\mathbf{w}$

Setting derivative of $\mathcal{L}_{\lambda}(\mathbf{w}, S)$ equal to 0 gives

 $\mathbf{X}'\mathbf{X}\mathbf{w} + \lambda\mathbf{w} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_n)\mathbf{w} = \mathbf{X}'\mathbf{y}$

ISSML, June 2013

Primal solution

We get the primal solution weight vector:

$$\mathbf{w} = \left(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_n\right)^{-1}\mathbf{X}'\mathbf{y}$$

and regression function

$$g(\mathbf{x}) = \mathbf{x}'\mathbf{w} = \mathbf{x}' \left(\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_n\right)^{-1} \mathbf{X}'\mathbf{y}$$

Dual solution

A dual solution should involve only computation of inner products – this is achieved by expressing the weight vector as a linear combination of the training examples:

$$\mathbf{X}'\mathbf{X}\mathbf{w} + \lambda\mathbf{w} = \mathbf{X}'\mathbf{y} \quad \text{implies}$$
$$\mathbf{w} = \frac{1}{\lambda}\left(\mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\mathbf{w}\right) = \mathbf{X}'\frac{1}{\lambda}\left(\mathbf{y} - \mathbf{X}\mathbf{w}\right) = \mathbf{X}'\alpha,$$

where

$$\alpha = \frac{1}{\lambda} \left(\mathbf{y} - \mathbf{X} \mathbf{w} \right) \tag{1}$$

or equivalently

$$\mathbf{w} = \sum_{i=1}^{m} \alpha_i \mathbf{x}_i$$

ISSML, June 2013

Dual solution

Substituting $\mathbf{w} = \mathbf{X}' \alpha$ into equation (1) we obtain:

 $\lambda \alpha = \mathbf{y} - \mathbf{X}\mathbf{X}' \alpha$

implying

$$(\mathbf{X}\mathbf{X}' + \lambda \mathbf{I}_m) \, \alpha = \mathbf{y}$$

This gives the dual solution:

$$\alpha = \left(\mathbf{X}\mathbf{X}' + \lambda\mathbf{I}_m\right)^{-1}\mathbf{y}$$

and regression function

$$g(\mathbf{x}) = \mathbf{x}'\mathbf{w} = \mathbf{x}'\mathbf{X}'\alpha = \sum_{i=1}^{m} \alpha_i \langle \mathbf{x}, \mathbf{x}_i \rangle$$

ISSML, June 2013

Key ingredients of dual solution

Step 1: Compute

$$\alpha = \left(\mathbf{K} + \lambda \mathbf{I}_m\right)^{-1} \mathbf{y}$$

where $\mathbf{K} = \mathbf{X}\mathbf{X}'$ that is $\mathbf{K}_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$

Step 2: Evaluate on new point x by

$$g(\mathbf{x}) = \sum_{i=1}^{m} \alpha_i \langle \mathbf{x}, \mathbf{x}_i \rangle$$

Important observation: Both steps only involve inner products

ISSML, June 2013

Applying the 'kernel trick'

Since the computation only involves inner products, we can substitute for all occurrences of $\langle \cdot, \cdot \rangle$ a kernel function κ that computes:

 $\kappa(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$

and we obtain an algorithm for ridge regression in the feature space F defined by the mapping

 $\phi: \mathbf{x} \longmapsto \phi(\mathbf{x}) \in F$

Note if ϕ is the identity this remains in the input space.

ISSML, June 2013

A simple kernel example

The simplest non-trivial kernel function is the quadratic kernel:

 $\kappa(\mathbf{x}, \mathbf{z}) = \langle \mathbf{x}, \mathbf{z} \rangle^2$

involving just one extra operation. But surprisingly this kernel function now corresponds to a complex feature mapping:

$$\kappa(\mathbf{x}, \mathbf{z}) = (\mathbf{x}'\mathbf{z})^2 = \mathbf{z}'(\mathbf{x}\mathbf{x}')\mathbf{z}$$
$$= \langle \operatorname{vec}(\mathbf{z}\mathbf{z}'), \operatorname{vec}(\mathbf{x}\mathbf{x}') \rangle$$

where vec(A) stacks the columns of the matrix A on top of each other. Hence, κ corresponds to the feature mapping

 $\phi: \mathbf{x} \longmapsto \operatorname{vec}(\mathbf{xx'})$

ISSML, June 2013

Implications of the kernel trick

- Consider for example computing a regression function over 1000 images represented by pixel vectors say $32 \times 32 = 1024$.
- By using the quadratic kernel we implement the regression function in a 1,000,000 dimensional space
- but actually using less computation for the learning phase than we did in the original space.

Implications of kernel algorithms

- Can perform linear regression in very highdimensional (even infinite dimensional) spaces efficiently.
- This is equivalent to performing non-linear regression in the original input space: for example quadratic kernel leads to solution of the form

$$g(\mathbf{x}) = \sum_{i=1}^{m} \alpha_i \langle \mathbf{x}_i, \mathbf{x} \rangle^2$$

that is a quadratic polynomial function of the components of the input vector \mathbf{x} .

• Using these high-dimensional spaces must surely come with a health warning, what about the curse of dimensionality?

ISSML, June 2013

Part 2

- Simple classification algorithm
- Principal components analysis.
- Kernel canonical correlation analysis.

Simple classification algorithm

 Consider finding the centres of mass of positive and negative examples and classifying a test point by measuring which is closest

$$h(\mathbf{x}) = \operatorname{sgn} \left(\|\phi(\mathbf{x}) - \phi_{S_{-}}\|^{2} - \|\phi(\mathbf{x}) - \phi_{S_{+}}\|^{2} \right)$$

 we can express as a function of kernel evaluations

$$h(\mathbf{x}) = \operatorname{sgn}\left(\frac{1}{m_{+}}\sum_{i=1}^{m_{+}}\kappa(\mathbf{x},\mathbf{x}_{i}) - \frac{1}{m_{-}}\sum_{i=m_{+}+1}^{m_{+}}\kappa(\mathbf{x},\mathbf{x}_{i}) - b\right),$$

where

$$b = \frac{1}{2m_+^2} \sum_{i,j=1}^{m_+} \kappa(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{2m_-^2} \sum_{i,j=m_++1}^m \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

ISSML, June 2013

Simple classification algorithm

 equivalent to dividing the space with a hyperplane perpendicular to the line half way between the two centres with vector given by

$$\mathbf{w} = \frac{1}{m^+} \sum_{i=1}^{m^+} \phi(\mathbf{x}_i) - \frac{1}{m^-} \sum_{i=m^++1}^{m^-} \phi(\mathbf{x}_i)$$

- Function is the difference in likelihood of the Parzen window density estimators for positive and negative examples
- We will see some examples of the performance of this algorithm in a moment.

ISSML, June 2013

Variance of projections

 Consider projections of the datapoints φ(x_i) onto a unit vector direction v in the feature space: average is given by

$$\mu_{\mathbf{v}} = \hat{\mathbf{E}} \left[\left\| P_{\mathbf{v}}(\phi(\mathbf{x})) \right\| \right] = \hat{\mathbf{E}} \left[\mathbf{v}' \phi(\mathbf{x}) \right] = \mathbf{v}' \phi_S$$

of course this is 0 if the data has been centred.

• average squared is given by

 $\hat{\mathbf{E}}\left[\|P_{\mathbf{v}}(\phi(\mathbf{x}))\|^{2}\right] = \hat{\mathbf{E}}\left[\mathbf{v}'\phi(\mathbf{x})\phi(\mathbf{x})'\mathbf{v}\right] = \frac{1}{m}\mathbf{v}'\mathbf{X}'\mathbf{X}\mathbf{v}$

ISSML, June 2013

Variance of projections

• Now suppose v has the dual representation $v = X'\alpha$. Average is given by

$$\mu_{\mathbf{v}} = \frac{1}{m} \alpha' \mathbf{X} \mathbf{X}' \mathbf{j} = \frac{1}{m} \alpha' \mathbf{K} \mathbf{j}$$

• average squared is given by

$$\frac{1}{m}\mathbf{v'}\mathbf{X'}\mathbf{X}\mathbf{v} = \frac{1}{m}\alpha'\mathbf{X}\mathbf{X'}\mathbf{X}\mathbf{X'}\alpha = \frac{1}{m}\alpha'\mathbf{K}^2\alpha$$

• Hence, variance in direction ${\bf v}$ is given by

$$\sigma_{\mathbf{v}}^2 = \frac{1}{m} \alpha' \mathbf{K}^2 \alpha - \frac{1}{m^2} (\alpha' \mathbf{K} \mathbf{j})^2$$

ISSML, June 2013

Fisher discriminant

• The Fisher discriminant is a thresholded linear classifier:

 $f(\mathbf{x}) = \operatorname{sgn}(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b$

where \mathbf{w} is chosen to maximise the quotient:

$$J(\mathbf{w}) = \frac{(\mu_{\mathbf{w}}^{+} - \mu_{\mathbf{w}}^{-})^{2}}{(\sigma_{\mathbf{w}}^{+})^{2} + (\sigma_{\mathbf{w}}^{-})^{2}}$$

 As with Ridge regression it makes sense to regularise if we are working in high-dimensional kernel spaces, so maximise

$$J(\mathbf{w}) = \frac{(\mu_{\mathbf{w}}^+ - \mu_{\mathbf{w}}^-)^2}{(\sigma_{\mathbf{w}}^+)^2 + (\sigma_{\mathbf{w}}^-)^2 + \lambda \|\mathbf{w}\|^2}$$

ISSML, June 2013

Fisher discriminant

- Using the results we now have we can substitute dual expressions for all of these quantities and solve using lagrange multipliers.
- The resulting classifier has dual variables

$$\alpha = (\mathbf{B}\mathbf{K} + \lambda \mathbf{I})^{-1}\mathbf{y}$$

where $\mathbf{B} = \mathbf{D} - \mathbf{C}$ with

$$\mathbf{C}_{ij} = \begin{cases} 2m^{-}/(mm^{+}) & \text{if } y_i = 1 = y_j \\ 2m^{+}/(mm^{-}) & \text{if } y_i = -1 = y_j \\ 0 & \text{otherwise} \end{cases}$$

ISSML, June 2013

and

$$\mathbf{D} = \begin{cases} 2m^{-}/m & \text{if } i = j \text{ and } y_{i} = 1\\ 2m^{+}/m & \text{if } i = j \text{ and } y_{i} = -1\\ 0 & \text{otherwise} \end{cases}$$

and $b = 0.5 \alpha \mathbf{Kt}$ with

$$\mathbf{t}_i = \left\{ \begin{array}{ll} 1/m^+ & \text{if } y_i = 1 \\ 1/m^- & \text{if } y_i = -1 \\ 0 & \text{otherwise} \end{array} \right.$$

giving a decision function

$$f(\mathbf{x}) = \operatorname{sgn}\left(\sum_{i=1}^{m} \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}) - b\right)$$

ISSML, June 2013

Preprocessing

- Corresponds to feature selection, or learning the feature space
- Note that in kernel methods the feature space is only determined up to orthogonal transformations (change of basis):

 $\hat{\phi}(\mathbf{x}) = \mathbf{U}\phi(\mathbf{x})$

for some orthogonal transformation \mathbf{U} ($\mathbf{U}'\mathbf{U} = \mathbf{I} = \mathbf{U}\mathbf{U}'$), then

 $\hat{\kappa}(\mathbf{x}, \mathbf{z}) = \langle \mathbf{U}\phi(\mathbf{x}), \mathbf{U}\phi(\mathbf{z}) \rangle = \phi(\mathbf{x})'\mathbf{U}'\mathbf{U}\phi(\mathbf{z}) = \phi(\mathbf{x})'\phi(\mathbf{z}) = \kappa(\mathbf{x}, \mathbf{z})$

• so feature selection in a kernel defined feature space is eqivalent to subspace projection

ISSML, June 2013

Subspace methods

- Principal components analysis: choose directions to maximise variance in the training data
- Canonical correlation analysis: choose directions to maximise correlations between two different views of the same objects
- Gram-Schmidt: greedily choose directions according to largest residual norms (not covered)
- Partial least squares: greedily choose directions with maximal covariance with the target (not covered)

In all cases we need kernel versions in order to apply these methods in high-dimensional kernel defined feature spaces

Principal Components Analysis

- PCA is a subspace method that is it involves projecting the data into a lower dimensional space.
- Subspace is chosen to ensure maximal variance of the projections:

 $\mathbf{w} = \operatorname{argmax}_{\mathbf{w}:\|\mathbf{w}\|=1} \mathbf{w}' \mathbf{X}' \mathbf{X} \mathbf{w}$

 This is equivalent to maximising the Raleigh quotient:

 $\frac{\mathbf{w'X'Xw}}{\mathbf{w'w}}$

ISSML, June 2013

Principal Components Analysis

• We can optimise using Lagrange multipliers in order to remove the contraints:

$$L(\mathbf{w},\lambda) = \mathbf{w}'\mathbf{X}'\mathbf{X}\mathbf{w} - \lambda\mathbf{w}'\mathbf{w}$$

taking derivatives wrt \mathbf{w} and setting equal to 0 gives:

 $\mathbf{X}'\mathbf{X}\mathbf{w} = \lambda\mathbf{w}$

implying \mathbf{w} is an eigenvector of $\mathbf{X}'\mathbf{X}$.

• Note that

$$\lambda = \mathbf{w}' \mathbf{X}' \mathbf{X} \mathbf{w} = \sum_{i=1}^{m} \langle \mathbf{w}, \mathbf{x}_i \rangle^2$$

ISSML, June 2013

Principal Components Analysis

- So principal components analysis performs an eigenvalue decomposition of $\mathbf{X}'\mathbf{X}$ and projects into the space spanned by the first k eigenvectors
- Captures a total of

$$\sum_{i=1}^k \lambda_i$$

of the overall variance:

$$\sum_{i=1}^m \|\mathbf{x}_i\|^2 = \sum_{i=1}^n \lambda_i = \operatorname{tr}(\mathbf{K})$$

ISSML, June 2013

Kernel PCA

- We would like to find a dual representation of the principal eigenvectors and hence of the projection function.
- Suppose that $\mathbf{w}, \lambda \neq 0$ is an eigenvector/eigenvalue pair for $\mathbf{X'X}$, then \mathbf{Xw}, λ is for $\mathbf{XX'}$:

$$(\mathbf{X}\mathbf{X}')\mathbf{X}\mathbf{w} = \mathbf{X}(\mathbf{X}'\mathbf{X})\mathbf{w} = \lambda\mathbf{X}\mathbf{w}$$

• and vice versa $\alpha, \lambda \to \mathbf{X}' \alpha, \lambda$

$$(\mathbf{X}'\mathbf{X})\mathbf{X}'\alpha = \mathbf{X}'(\mathbf{X}\mathbf{X}')\alpha = \lambda\mathbf{X}'\alpha$$

Note that we get back to where we started if we do it twice.

ISSML, June 2013

Kernel PCA

• Hence, 1-1 correspondence between eigenvectors corresponding to non-zero eigenvalues, but note that if $\|\alpha\| = 1$

$$\|\mathbf{X}'\alpha\|^2 = \alpha'\mathbf{X}\mathbf{X}'\alpha = \alpha'\mathbf{K}\alpha = \lambda$$

so if $\alpha^i, \lambda_i, i = 1, ..., k$ are first k eigenvectors/values of **K** $\frac{1}{\sqrt{\lambda_i}} \alpha^i$

are dual representations of first k eigenvectors $\mathbf{w}^1, \ldots, \mathbf{w}^k$ of $\mathbf{X}'\mathbf{X}$ with same eigenvalues.

• Computing projections:

$$\langle \mathbf{w}^i, \phi(\mathbf{x}) \rangle = \frac{1}{\sqrt{\lambda_i}} \langle \mathbf{X}' \alpha^i, \phi(\mathbf{x}) \rangle = \frac{1}{\sqrt{\lambda_i}} \sum_{j=1}^m \alpha_j^i \kappa(\mathbf{x}_i, \mathbf{x})$$

ISSML, June 2013

Kernel CCA

• Canonical correlation analysis finds correlations between different views of the same object:

 $\mathbf{x} \longmapsto (\phi_a(\mathbf{x}), \phi_b(\mathbf{x}))$

- Examples:
 - documents in two languages
 - image and its caption
 - brain scan and corresponding activity description
- Seek \mathbf{w}_a creating feature $x_a = \mathbf{w}'_a \phi_a(\mathbf{x})$ and \mathbf{w}_b creating feature $x_b = \mathbf{w}'_b \phi_b(\mathbf{x})$ that maximise:

$$\rho = \frac{\hat{\mathbb{E}}[x_a x_b]}{\sqrt{\hat{\mathbb{E}}[x_a^2]\hat{\mathbb{E}}[x_b^2]}} = \frac{\hat{\mathbb{E}}[\mathbf{w}_a' \phi_a(\mathbf{x})\phi_b(\mathbf{x})'\mathbf{w}_b]}{\sqrt{\hat{\mathbb{E}}[(\mathbf{w}_a' \phi_a(\mathbf{x}))^2]\hat{\mathbb{E}}[(\mathbf{w}_a' \phi_a(\mathbf{x}))^2]}}$$

ISSML, June 2013

Kernel CCA

• Using our standard notation

$$\rho = \frac{\mathbf{w}_a' \mathbf{X}_a' \mathbf{X}_b \mathbf{w}_b}{\sqrt{\mathbf{w}_a' \mathbf{X}_a' \mathbf{X}_a \mathbf{w}_a \mathbf{w}_b' \mathbf{X}_b' \mathbf{X}_b \mathbf{w}_b}}$$

• Since invariant to rescalings we can set two factors in denominator equal to 1. Using Lagrange multipliers obtain $L(\mathbf{w}_a, \mathbf{w}_b, \lambda_a, \lambda_b)$ as

$$\mathbf{w}_a' \mathbf{X}_a' \mathbf{X}_b \mathbf{w}_b - rac{\lambda_a}{2} \mathbf{w}_a' \mathbf{X}_a' \mathbf{X}_a \mathbf{w}_a - rac{\lambda_b}{2} \mathbf{w}_b' \mathbf{X}_b' \mathbf{X}_b \mathbf{w}_b$$

• Gives coupled equations

 $\mathbf{X}_{a}'\mathbf{X}_{b}\mathbf{w}_{b} = \lambda_{a}\mathbf{X}_{a}'\mathbf{X}_{a}\mathbf{w}_{a} \text{ and } \mathbf{w}_{a}'\mathbf{X}_{a}'\mathbf{X}_{b} = \lambda_{b}\mathbf{w}_{b}'\mathbf{X}_{b}'\mathbf{X}_{b}$

ISSML, June 2013

Kernel CCA

- Simple to verify that $\lambda_a = \lambda_b$
- Resulting in generalised eigenvalue problem:

$$\begin{pmatrix} \mathbf{0} & \mathbf{X}'_{a}\mathbf{X}_{b} \\ \mathbf{X}'_{b}\mathbf{X}_{a} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{w}_{a} \\ \mathbf{w}_{b} \end{pmatrix} = \lambda \begin{pmatrix} \mathbf{X}'_{a}\mathbf{X}_{a} & \mathbf{0} \\ \mathbf{0} & \mathbf{X}'_{b}\mathbf{X}_{b} \end{pmatrix} \begin{pmatrix} \mathbf{w}_{a} \\ \mathbf{w}_{b} \end{pmatrix}$$

 We would like to make a kernel version of this procedure – but must ensure that the flexibility is controlled to rule out spurious correlations:

 $\rho = \max_{\mathbf{w}_a, \mathbf{w}_b} \mathbf{w}_a' \mathbf{X}_a' \mathbf{X}_b \mathbf{w}_b$ subject to: $(1 - \tau) \mathbf{w}_a' \mathbf{X}_a' \mathbf{X}_a \mathbf{w}_a + \tau \mathbf{w}_a' \mathbf{w}_a = 1$ and $(1 - \tau) \mathbf{w}_b' \mathbf{X}_b' \mathbf{X}_b \mathbf{w}_b + \tau \mathbf{w}_b' \mathbf{w}_b = 1$

ISSML, June 2013
ISSML, June 2013

Kernel CCA

• We now dualise by letting $\mathbf{w}_a = \mathbf{X}'_a \alpha$ and $\mathbf{w}_b = \mathbf{X}'_b \beta$ to obtain:

 $\rho = \max_{\alpha,\beta} \alpha' \mathbf{K}_a \mathbf{K}_b \beta$ subject to: $(1 - \tau) \alpha' \mathbf{K}_a^2 \alpha + \tau \alpha' \mathbf{K}_a \alpha = 1$ and $(1 - \tau) \beta' \mathbf{K}_b^2 \beta + \tau \beta' \mathbf{K}_b \beta = 1$

• giving the generalised eigenvalue problem

$$\begin{pmatrix} \mathbf{0} & \mathbf{K}_{a}\mathbf{K}_{b} \\ \mathbf{K}_{b}\mathbf{K}_{a} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$
$$= \lambda \begin{pmatrix} (1-\tau)\mathbf{K}_{a}^{2} + \tau\mathbf{K}_{a} & \mathbf{0} \\ \mathbf{0} & (1-\tau)\mathbf{K}_{b}^{2} + \tau\mathbf{K}_{b} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

37

ISSML, June 2013

Part 3

- Statistical analysis of the stability of patterns.
- Rademacher complexity.
- Generalisation of SVMs
- Support Vector Machine Optimisation

Generalisation of a learner

- Assume that we have a learning algorithm A that chooses a function A_F(S) from a function space F in response to the training set S.
- From a statistical point of view the quantity of interest is the random variable:

 $\epsilon(S, \mathcal{A}, \mathcal{F}) = \mathbb{E}_{(\mathbf{x}, y)} \left[\ell(\mathcal{A}_{\mathcal{F}}(S), \mathbf{x}, y) \right],$

where ℓ is a 'loss' function that measures the discrepancy between $\mathcal{A}_{\mathcal{F}}(S)(\mathbf{x})$ and y.

Generalisation of a learner

- For example, in the case of classification ℓ is 1 if the two disagree and 0 otherwise, while for regression it could be the square of the difference between $\mathcal{A}_{\mathcal{F}}(S)(\mathbf{x})$ and y.
- We refer to the random variable $\epsilon(S, \mathcal{A}, \mathcal{F})$ as the generalisation of the learner.

Example of Generalisation I

- We consider the Breast Cancer dataset from the UCI repository.
- Use the simple Parzen window classifier described in Part 2: weight vector is



where w^+ is the average of the positive training examples and w^- is average of negative training examples. Threshold is set so hyperplane bisects the line joining these two points.

Example of Generalisation II

• Given a size *m* of the training set, by repeatedly drawing random training sets *S* we estimate the distribution of

 $\epsilon(S, \mathcal{A}, \mathcal{F}) = \mathbb{E}_{(\mathbf{x}, y)} \left[\ell(\mathcal{A}_{\mathcal{F}}(S), \mathbf{x}, y) \right],$

by using the test set error as a proxy for the true generalisation.

 We plot the histogram and the average of the distribution for various sizes of training set – initially the whole dataset gives a single value if we use training and test as the all the examples, but then we plot for training set sizes:

342, 273, 205, 137, 68, 34, 27, 20, 14, 7.

ISSML, June 2013

Example of Generalisation III

• Since the expected classifier is in all cases the same:

$$\mathbb{E} \left[\mathcal{A}_{\mathcal{F}}(S) \right] = \mathbb{E}_{S} \left[\mathbf{w}_{S}^{+} - \mathbf{w}_{S}^{-} \right]$$
$$= \mathbb{E}_{S} \left[\mathbf{w}_{S}^{+} \right] - \mathbb{E}_{S} \left[\mathbf{w}_{S}^{-} \right]$$
$$= \mathbb{E}_{y=+1} \left[\mathbf{x} \right] - \mathbb{E}_{y=-1} \left[\mathbf{x} \right],$$

we do not expect large differences in the average of the distribution, though the non-linearity of the loss function means they won't be the same exactly.

Error distribution: full dataset



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013

Observations

- Things can get bad if number of training examples small compared to dimension (in this case input dimension is 9)
- Mean can be bad predictor of true generalisation i.e. things can look okay in expectation, but still go badly wrong
- Key ingredient of learning keep flexibility high while still ensuring good generalisation

Controlling generalisation

- The critical method of controlling generalisation for classification is to force a large margin on the training data
- Equivalent to minimising the norm while keeping the separation fixed (at say $\pm 1)$
- Support Vector Machines implement this strategy

Controlling generalisation

- Now consider using an SVM on the same data and compare the distribution of generalisations
- SVM distribution in red



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013



ISSML, June 2013

Expected versus confident bounds

- For a finite sample the generalisation $\epsilon(S, \mathcal{A}, \mathcal{F})$ has a distribution depending on the algorithm, function class and sample size m.
- Traditional statistics as indicated above has concentrated on the mean of this distribution – but this quantity can be misleading, eg for low fold cross-validation.

Expected versus confident bounds cont.

- Statistical learning theory has preferred to analyse the tail of the distribution, finding a bound which holds with high probability.
- This looks like a statistical test significant at a 1% confidence means that the chances of the conclusion not being true are less than 1% over random samples of that size.
- This is also the source of the acronym PAC: probably approximately correct, the 'confidence' parameter δ is the probability that we have been misled by the training set.

Concentration inequalities

- Statistical Learning is concerned with the reliability or stability of inferences made from a random sample.
- Random variables with this property have been a subject of ongoing interest to probabilists and statisticians.

Concentration inequalities cont.

• As an example consider the mean of a sample of m 1-dimensional random variables X_1, \ldots, X_m :

$$S_m = \frac{1}{m} \sum_{i=1}^m X_i.$$

• Hoeffding's inequality states that if $X_i \in [a_i, b_i]$

$$P\{|S_m - \mathbb{E}[S_m]| \ge \epsilon\} \le 2\exp\left(-\frac{2m^2\epsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}\right)$$

Note how the probability falls off exponentially with the distance from the mean and with the number of variables.

ISSML, June 2013

Concentration for SLT

- We are now going to look at deriving SLT results from concentration inequalities.
- Perhaps the best known form is due to McDiarmid (although he was actually representing previously derived results):

McDiarmid's inequality

Theorem 1. Let X_1, \ldots, X_n be independent random variables taking values in a set A, and assume that $f: A^n \to \mathbb{R}$ satisfies

 $\sup_{x_1,\ldots,x_n,\hat{x}_i\in A} |f(x_1,\ldots,x_n) - f(x_1,\ldots,\hat{x}_i,x_{i+1},\ldots,x_n)| \le c_i,$

for $1 \leq i \leq n$. Then for all $\epsilon > 0$,

$$P\left\{f\left(X_{1},\ldots,X_{n}\right)-\mathbb{E}f\left(X_{1},\ldots,X_{n}\right)\geq\epsilon\right\}\leq\exp\left(\frac{-2\epsilon^{2}}{\sum_{i=1}^{n}c_{i}^{2}}\right)$$

• Hoeffding is a special case when $f(x_1, \ldots, x_n) = S_n$

ISSML, June 2013

Using McDiarmid

• By setting the right hand side equal to δ , we can always invert McDiarmid to get a high confidence bound: with probability at least $1 - \delta$

$$f(X_1, \dots, X_n) < \mathbb{E}f(X_1, \dots, X_n) + \sqrt{\frac{\sum_{i=1}^n c_i^2}{2} \log \frac{1}{\delta}}$$

• If $c_i = c/n$ for each *i* this reduces to

$$f(X_1,\ldots,X_n) < \mathbb{E}f(X_1,\ldots,X_n) + \sqrt{\frac{c^2}{2n}\log\frac{1}{\delta}}$$

ISSML, June 2013
Rademacher complexity

• Rademacher complexity is a new way of measuring the complexity of a function class. It arises naturally if we rerun the proof using the double sample trick and symmetrisation but look at what is actually needed to continue the proof:

ISSML, June 2013

Rademacher proof beginnings

For a fixed $f \in \mathcal{F}$ we have

$$\mathbb{E}\left[f(\mathbf{z})\right] \leq \hat{\mathbb{E}}\left[f(\mathbf{z})\right] + \sup_{h \in \mathcal{F}} \left(\mathbb{E}[h] - \hat{\mathbb{E}}[h]\right).$$

where \mathcal{F} is a class of functions mapping from Z to [0,1] and $\hat{\mathbb{E}}$ denotes the sample average.

We must bound the size of the second term. First apply McDiarmid's inequality to obtain ($c_i = 1/m$ for all *i*) with probability at least $1 - \delta$:

$$\sup_{h\in\mathcal{F}} \left(\mathbb{E}[h] - \hat{\mathbb{E}}[h]\right) \le \mathbb{E}_S \left[\sup_{h\in\mathcal{F}} \left(\mathbb{E}[h] - \hat{\mathbb{E}}[h]\right)\right] + \sqrt{\frac{\ln(1/\delta)}{2m}}.$$

ISSML, June 2013

Deriving double sample result

• We can now move to the ghost sample by simply observing that $\mathbb{E}[h] = \mathbb{E}_{\tilde{S}}\left[\hat{\mathbb{E}}[h]\right]$:

$$\mathbb{E}_{S}\left[\sup_{h\in\mathcal{F}}\left(\mathbb{E}[h] - \hat{\mathbb{E}}[h]\right)\right] = \mathbb{E}_{S}\left[\sup_{h\in\mathcal{F}}\mathbb{E}_{\tilde{S}}\left[\frac{1}{m}\sum_{i=1}^{m}h(\tilde{\mathbf{z}}_{i}) - \frac{1}{m}\sum_{i=1}^{m}h(\mathbf{z}_{i})\middle|S\right]\right]$$

ISSML, June 2013

Deriving double sample result cont.

Since the sup of an expectation is less than or equal to the expectation of the sup (we can make the choice to optimise for each \tilde{S}) we have

$$\mathbb{E}_{S}\left[\sup_{h\in\mathcal{F}}\left(\mathbb{E}[h] - \hat{\mathbb{E}}[h]\right)\right] \leq \mathbb{E}_{S}\mathbb{E}_{\tilde{S}}\left[\sup_{h\in\mathcal{F}}\frac{1}{m}\sum_{i=1}^{m}\left(h(\tilde{\mathbf{z}}_{i}) - h(\mathbf{z}_{i})\right)\right]$$

ISSML, June 2013

Adding symmetrisation

Here symmetrisation is again just swapping corresponding elements – but we can write this as multiplication by a variable σ_i which takes values ± 1 with equal probability:

$$\mathbb{E}_{S}\left[\sup_{h\in\mathcal{F}}\left(\mathbb{E}[h]-\hat{\mathbb{E}}[h]\right)\right] \leq \\ \leq \mathbb{E}_{\sigma S\tilde{S}}\left[\sup_{h\in\mathcal{F}}\frac{1}{m}\sum_{i=1}^{m}\sigma_{i}\left(h(\tilde{\mathbf{z}}_{i})-h(\mathbf{z}_{i})\right)\right] \\ \leq 2\mathbb{E}_{S\sigma}\left[\sup_{h\in\mathcal{F}}\frac{1}{m}\sum_{i=1}^{m}\sigma_{i}h(\mathbf{z}_{i})\right] \\ = R_{m}\left(\mathcal{F}\right),$$

assuming \mathcal{F} closed under negation $f \mapsto -f$.

ISSML, June 2013

Rademacher complexity

The Rademacher complexity provides a way of measuring the complexity of a function class \mathcal{F} by testing how well on average it can align with random noise:

$$R_m(\mathcal{F}) = \mathbb{E}_{S\sigma} \left[\sup_{f \in \mathcal{F}} \frac{2}{m} \sum_{i=1}^m \sigma_i f(\mathbf{z}_i) \right].$$

is known as the Rademacher complexity of the function class \mathcal{F} .

ISSML, June 2013

Main Rademacher theorem

The main theorem of Rademacher complexity: with probability at least $1 - \delta$ over random samples *S* of size *m*, every $f \in \mathcal{F}$ satisfies



- Note that Rademacher complexity gives the expected value of the maximal correlation with random noise – a very natural measure of capacity.
- Note that the Rademacher complexity is distribution dependent since it involves an expectation over the choice of sample – this might seem hard to compute.

ISSML, June 2013

Empirical Rademacher theorem

• Since the empirical Rademacher complexity

$$\hat{R}_{m}(\mathcal{F}) = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{2}{m} \sum_{i=1}^{m} \sigma_{i} f(\mathbf{z}_{i}) \middle| \mathbf{z}_{1}, \dots, \mathbf{z}_{m} \right]$$

is concentrated, we can make a further application of McDiarmid to obtain with probability at least $1-\delta$

$$\mathbb{E}_{\mathcal{D}}\left[f(\mathbf{z})\right] \leq \hat{\mathbb{E}}\left[f(\mathbf{z})\right] + \hat{R}_m(\mathcal{F}) + 3\sqrt{\frac{\ln(2/\delta)}{2m}}.$$

ISSML, June 2013

Application to large margin classification

• Rademacher complexity comes into its own for Boosting and SVMs.

ISSML, June 2013

Application to Boosting

• We can view Boosting as seeking a function from the class (*H* is the set of weak learners)

$$\left\{\sum_{h\in H} a_h h(\mathbf{x}) : a_h \ge 0, \sum_{h\in H} a_h \le B\right\} = \operatorname{conv}_B(H)$$

by minimising some function of the margin distribution (assume H closed under negation).

- Adaboost corresponds to optimising an exponential function of the margin over this set of functions.
- We will see how to include the margin in the analysis later, but concentrate on computing the Rademacher complexity for now.

ISSML, June 2013

Rademacher complexity of convex hulls

Rademacher complexity has a very nice property for convex hull classes:

$$\hat{R}_{m}(\operatorname{conv}_{B}(H)) = \frac{2}{m} \mathbb{E}_{\sigma} \left[\sup_{\substack{h_{j} \in H, \sum_{j} a_{j} \leq B}} \sum_{i=1}^{m} \sigma_{i} \sum_{j} a_{j} h_{j}(\mathbf{x}_{i}) \right] \\
\leq \frac{2}{m} \mathbb{E}_{\sigma} \left[\sup_{\substack{h_{j} \in H, \sum_{j} a_{j} \leq B}} \sum_{j} a_{j} \sum_{i=1}^{m} \sigma_{i} h_{j}(\mathbf{x}_{i}) \right] \\
\leq \frac{2}{m} \mathbb{E}_{\sigma} \left[\sup_{\substack{h_{j} \in H}} B \sum_{i=1}^{m} \sigma_{i} h_{j}(\mathbf{x}_{i}) \right] \\
\leq B \hat{R}_{m}(H).$$

ISSML, June 2013

Rademacher complexity of convex hulls cont.

• Hence, we can move to the convex hull without incurring any complexity penalty for B = 1!

ISSML, June 2013

Rademacher complexity for SVMs

• The Rademacher complexity of a class of linear functions with bounded 2-norm:

$$\begin{cases} \mathbf{x} \to \sum_{i=1}^{m} \alpha_{i} \kappa(\mathbf{x}_{i}, \mathbf{x}) : \alpha' \mathbf{K} \alpha \leq B^{2} \\ \\ \subseteq \{ \mathbf{x} \to \langle \mathbf{w}, \phi(\mathbf{x}) \rangle : \| \mathbf{w} \| \leq B \} \\ \\ = \mathcal{F}_{B}, \end{cases}$$

where we assume a kernel defined feature space with

 $\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \kappa(\mathbf{x}, \mathbf{z}).$

ISSML, June 2013

Rademacher complexity of \mathcal{F}_B The following derivation gives the result

$$\hat{R}_{m}(\mathcal{F}_{B}) = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}_{B}} \left| \frac{2}{m} \sum_{i=1}^{m} \sigma_{i} f(\mathbf{x}_{i}) \right| \right] \\
= \mathbb{E}_{\sigma} \left[\sup_{\|\mathbf{w}\| \leq B} \left| \left\langle \mathbf{w}, \frac{2}{m} \sum_{i=1}^{m} \sigma_{i} \phi(\mathbf{x}_{i}) \right\rangle \right| \right] \\
\leq \frac{2B}{m} \mathbb{E}_{\sigma} \left[\left\| \sum_{i=1}^{m} \sigma_{i} \phi(\mathbf{x}_{i}) \right\| \right] \\
= \frac{2B}{m} \mathbb{E}_{\sigma} \left[\left(\left\langle \sum_{i=1}^{m} \sigma_{i} \phi(\mathbf{x}_{i}), \sum_{j=1}^{m} \sigma_{j} \phi(\mathbf{x}_{j}) \right\rangle \right)^{1/2} \right] \\
\leq \frac{2B}{m} \left(\mathbb{E}_{\sigma} \left[\sum_{i,j=1}^{m} \sigma_{i} \sigma_{j} \kappa(\mathbf{x}_{i}, \mathbf{x}_{j}) \right] \right)^{1/2} = \frac{2B}{m} \sqrt{\sum_{i=1}^{m} \kappa(\mathbf{x}_{i}, \mathbf{x}_{i})}$$

ISSML, June 2013

ISSML, June 2013

Support Vector Machines (SVM)

 SVM seeks linear function in a feature space defined implicitly via a kernel κ:

 $\kappa(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$

that optimises a bound on the generalisation.

• The first step is to introduce a loss function which upper bounds the discrete loss

 $P(y \neq \operatorname{sgn}(g(\mathbf{x}))) = \mathbb{E}\left[\mathcal{H}(-yg(\mathbf{x}))\right],$

where $\mathcal H$ is the Heaviside function.

ISSML, June 2013

Margins in SVMs

• Critical to the bound will be the margin of the classifier

 $\gamma(\mathbf{x},y) = yg(\mathbf{x}) = y(\langle \mathbf{w},\phi(\mathbf{x})\rangle + b)$:

positive if correctly classified, and measures distance from the separating hyperplane when the weight vector is normalised.

• The margin of a linear function g is

$$\gamma(g) = \min_{i} \gamma(\mathbf{x}_i, y_i)$$

though this is frequently increased to allow some 'margin errors'.

ISSML, June 2013



ISSML, June 2013

Applying the Rademacher theorem

• Consider the loss function $\mathcal{A}:\mathbb{R} \to [0,1],$ given by

$$\mathcal{A}(a) = \left\{ egin{array}{ll} 1, & ext{if } a > 0; \ 1+a/\gamma, & ext{if } -\gamma \leq a \leq 0; \ 0, & ext{otherwise}. \end{array}
ight.$$

 By the Rademacher Theorem and since the loss function A dominates H, we have that

$$\mathbb{E} \left[\mathcal{H}(-yg(\mathbf{x})) \right] \leq \mathbb{E} \left[\mathcal{A}(-yg(\mathbf{x})) \right] \\ \leq \hat{\mathbb{E}} \left[\mathcal{A}(-yg(\mathbf{x})) \right] + \\ \hat{R}_m(\mathcal{A} \circ \mathcal{F}) + 3\sqrt{\frac{\ln(2/\delta)}{2m}}.$$

ISSML, June 2013

ISSML, June 2013

Empirical loss and slack variables

• But the function $\mathcal{A}(-y_i g(\mathbf{x}_i)) \leq \xi_i / \gamma$, for $i = 1, \ldots, m$, and so

$$\mathbb{E}\left[\mathcal{H}(-yg(\mathbf{x}))\right] \leq \frac{1}{m\gamma} \sum_{i=1}^{m} \xi_i + \hat{R}_m(\mathcal{A} \circ \mathcal{F}) + 3\sqrt{\frac{\ln(2/\delta)}{2m}}.$$

- The final missing ingredient to complete the bound is to bound $\hat{R}_m(\mathcal{A} \circ \mathcal{F})$ in terms of $\hat{R}_m(\mathcal{F})$.
- This can be obtained in terms of the maximal slope of the function \mathcal{A} : $\hat{R}_m(\mathcal{A} \circ \mathcal{F}) \leq \frac{2}{\gamma} \hat{R}_m(\mathcal{F})$.

ISSML, June 2013

Final SVM bound

• Assembling the result we obtain:

$$P(y \neq \operatorname{sgn}(g(\mathbf{x}))) = \mathbb{E}\left[\mathcal{H}(-yg(\mathbf{x}))\right]$$
$$\leq \frac{1}{m\gamma} \sum_{i=1}^{m} \xi_i + \frac{4}{m\gamma} \sqrt{\sum_{i=1}^{m} \kappa(\mathbf{x}_i, \mathbf{x}_i)} + 3\sqrt{\frac{\ln(2/\delta)}{2m}}$$

• Note that for the Gaussian kernel this reduces to

$$P(y \neq \operatorname{sgn}(g(\mathbf{x}))) \leq \frac{1}{m\gamma} \sum_{i=1}^{m} \xi_i + \frac{4}{\sqrt{m\gamma}} + 3\sqrt{\frac{\ln(2/\delta)}{2m}}$$

ISSML, June 2013

Using a kernel

- Can consider much higher dimensional spaces using the kernel trick
- Can even work in infinite dimensional spaces, eg using the Gaussian kernel:

$$\kappa(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2}\right)$$

ISSML, June 2013

Error distribution: dataset size: 342



ISSML, June 2013

Error distribution: dataset size: 273



ISSML, June 2013

Applying to 1-norm SVMs

We take the following formulation of the 1-norm SVM to optimise the bound:

$$\begin{array}{ll} \min_{\mathbf{w},b,\gamma,\xi} & -\gamma + C \sum_{i=1}^{m} \xi_i \\ \text{subject to} & y_i \left(\langle \mathbf{w}, \phi \left(\mathbf{x}_i \right) \rangle + b \right) \geq \gamma - \xi_i, \, \xi_i \geq 0, \\ & i = 1, \dots, m, \, \text{and} \, \|\mathbf{w}\|^2 = 1. \end{array}$$

$$\begin{array}{l} \text{(2)} \\ \xi_i = (\gamma - y_i g(\mathbf{x}_i))_+, \end{array}$$

where $g(\cdot) = \langle \mathbf{w}, \phi(\cdot) \rangle + b$.

ISSML, June 2013

Forming the Lagrangian $L(\mathbf{w}, b, \gamma, \xi, \alpha, \beta, \lambda)$:

$$-\gamma + C \sum_{i=1}^{m} \xi_{i} - \sum_{i=1}^{m} \alpha_{i} \left[y_{i} (\langle \phi \left(\mathbf{x}_{i} \right), \mathbf{w} \rangle + b) - \gamma + \xi_{i} \right] - \sum_{i=1}^{m} \beta_{i} \xi_{i} + \lambda \left(\| \mathbf{w} \|^{2} - 1 \right)$$

with $\alpha_i \geq 0$ and $\beta_i \geq 0$.

ISSML, June 2013

Taking derivatives gives:

$$\frac{\partial L(\mathbf{w}, b, \gamma, \xi, \alpha, \beta, \lambda)}{\partial \mathbf{w}} = 2\lambda \mathbf{w} - \sum_{i=1}^{m} y_i \alpha_i \phi(\mathbf{x}_i) = \mathbf{0},$$
$$\frac{\partial L(\mathbf{w}, b, \gamma, \xi, \alpha, \beta, \lambda)}{\partial \xi_i} = C - \alpha_i - \beta_i = 0,$$
$$\frac{\partial L(\mathbf{w}, b, \gamma, \xi, \alpha, \beta, \lambda)}{\partial b} = \sum_{i=1}^{m} y_i \alpha_i = 0,$$
$$\frac{\partial L(\mathbf{w}, b, \gamma, \xi, \alpha, \beta, \lambda)}{\partial \gamma} = 1 - \sum_{i=1}^{m} \alpha_i = 0.$$

ISSML, June 2013

$$L(\alpha, \lambda) = -\frac{1}{4\lambda} \sum_{i,j=1}^{m} y_i y_j \alpha_i \alpha_j \kappa \left(\mathbf{x}_i, \mathbf{x}_j\right) - \lambda,$$

which, again optimising with respect to λ , gives

$$\lambda^* = \frac{1}{2} \left(\sum_{i,j=1}^m y_i y_j \alpha_i \alpha_j \kappa \left(\mathbf{x}_i, \mathbf{x}_j \right) \right)^{1/2}$$

ISSML, June 2013

equivalent to maximising

$$L(\alpha) = -\sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j \kappa \left(\mathbf{x}_i, \mathbf{x}_j\right),$$

subject to the constraints

$$0 \le \alpha_i \le C, \quad \sum_{i=1}^m \alpha_i = 1 \quad \sum_{i=1}^m y_i \alpha_i = 0$$

to give solution

$$\alpha_i^*, i = 1, \dots, m$$

ISSML, June 2013

This is a convex quadratic programme: minimising a convex quadratic objective subject to linear constraints: convex if Hessian G is positive semidefinite:

 $\mathbf{G}_{ij} = y_i y_j \kappa \left(\mathbf{x}_i, \mathbf{x}_j \right)$

Matrix psd iff $\mathbf{u}'\mathbf{Gu} \ge 0$ for all \mathbf{u} :

$$\mathbf{u}'\mathbf{G}\mathbf{u} = \sum_{i,j=1}^{m} u_i u_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$
$$= \left\langle \sum_{i=1}^{m} u_i y_i \phi(\mathbf{x}_i), \sum_{j=1}^{m} u_j y_j \phi(\mathbf{x}_j) \right\rangle$$
$$= \left\| \left\| \sum_{i=1}^{m} u_i y_i \phi(\mathbf{x}_i) \right\|^2 \ge 0$$

ISSML, June 2013

Kuhn-Tucker conditions:

$$\alpha_{i} \left[y_{i}(\langle \phi(\mathbf{x}_{i}), \mathbf{w} \rangle + b) - \gamma + \xi_{i} \right] = 0$$

$$\beta_{i}\xi_{i} = 0$$

These imply:

• $\alpha_i \neq 0$ only if

$$y_i(\langle \phi(\mathbf{x}_i), \mathbf{w} \rangle + b) = \gamma - \xi_i$$

these correspond to support vectors – their margins are less than or equal to γ .

• $\xi_i \neq 0$ only if $\beta_i = 0$ implying that $\alpha_i = C$, i.e. for $0 < \alpha_i < C$ margin is exactly γ .

ISSML, June 2013

The solution can then be computed as:

choose

i, j such that $-C < \alpha_i^* y_i < 0 < \alpha_j^* y_j < C$

$$b^{*} = -0.5 \left(\sum_{k=1}^{m} \alpha_{k}^{*} y_{k} \kappa \left(\mathbf{x}_{k}, \mathbf{x}_{i} \right) + \sum_{k=1}^{m} \alpha_{k}^{*} y_{k} \kappa \left(\mathbf{x}_{k}, \mathbf{x}_{j} \right) \right)$$
$$f(\cdot) = \operatorname{sgn} \left(\sum_{j=1}^{m} \alpha_{j}^{*} y_{j} \kappa \left(\mathbf{x}_{j}, \cdot \right) + b^{*} \right);$$

ISSML, June 2013

We can compute the margin as follows:

$$\lambda^* = \frac{1}{2} \left(\sum_{i,j=1}^m y_i y_j \alpha_i^* \alpha_j^* \kappa \left(\mathbf{x}_i, \mathbf{x}_j \right) \right)^{1/2}$$
$$\gamma^* = (2\lambda^*)^{-1} \left(\sum_{k=1}^m \alpha_k^* y_k \kappa \left(\mathbf{x}_k, \mathbf{x}_j \right) + b^* \right)$$

Similarly we can compute

$$\sum_{i=1}^{m} \xi_i = \frac{-2\lambda^* + \gamma^*}{C}$$

if we wish to compute the value of the bound.

ISSML, June 2013

Decision boundary and γ margin for 1-norm svm with a gaussian kernel:



ISSML, June 2013

- Have introduced a slightly non-standard version of the SVM but makes ν-SVM very simple to define.
- Consider expressing $C = 1/(\nu m)$:
 - implies $0 \le \alpha_i \le 1/(\nu m)$
 - if $\xi > 0$ then $\alpha_i = 1/(\nu m)$, but $\sum_{i=1}^m \alpha_i = 1$ so at most νm inputs can have this hold.
 - equally at least νm inputs have $\alpha_i \neq 0$
- Hence, v can be seen as the fraction of 'support vectors', a natural measure of the noise in the data.

ISSML, June 2013
Alternative form of the SVM problem

Note more traditional form of the dual SVM optimisation:

$$L(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j \kappa \left(\mathbf{x}_i, \mathbf{x}_j\right).$$

with constraints

$$0 \le \alpha_i \le C, \qquad \sum_{i=1}^m y_i \alpha_i = 0$$

ISSML, June 2013

Alternative form of the SVM problem

- Arises from considering renormalising so that output at margin is 1 and minimising the weight vector.
- The values of the regularisation parameter *C* do not correspond.
- Has advantage of simple kernel adatron algorithm if we consider the case of fixing b = 0 which removes the constraint $\sum_{i=1}^{m} \alpha_i y_i = 0$, so can perform gradient descent on individual α_i independently.
- SMO algorithm performs the update on pairs of *α_i*, *α_j* to ensure constraints remain satisfied.

ISSML, June 2013

Part 4

- Novelty detection
- Boosting
- Multiple Kernel Learning: bounds and algorithms

ISSML, June 2013

We can also motivate novelty detection by a similar analysis as that for SVM: consider a hypersphere centred at c of radius r and the function g:

$$g\left(\mathbf{x}\right) = \begin{cases} 0, & \text{if } \|\mathbf{c} - \phi(\mathbf{x})\| \le r; \\ (\|\mathbf{c} - \phi(\mathbf{x})\|^2 - r^2)/\gamma, & \text{if } r^2 \le \|\mathbf{c} - \phi(\mathbf{x})\|^2 \le r^2 + \gamma; \\ 1, & \text{otherwise.} \end{cases}$$

with probability at least $1-\delta$

$$\mathbb{E}[g(\mathbf{x})] \le \hat{\mathbb{E}}[g(\mathbf{x})] + \frac{6R^2}{\gamma\sqrt{m}} + 3\sqrt{\frac{\ln(2/\delta)}{2m}}$$

Note that tension is between creating a tight bound and defining a small sphere.

ISSML, June 2013

Let

$$\xi_i = (\|\mathbf{c} - \phi(\mathbf{x})\|^2 - r^2)_+$$

so that

$$\hat{\mathbb{E}}[g(\mathbf{x})] \le \frac{1}{\gamma m} \|\xi\|_1$$

Treating γ as fixed we minimise the bound by minimising $\|\xi\|_1$ and r:



ISSML, June 2013

Again introducing the Lagrangian $L(\mathbf{c}, r, \alpha, \xi)$

$$r^{2} + C \sum_{i=1}^{m} \xi_{i} + \sum_{i=1}^{m} \alpha_{i} \left[\|\phi(\mathbf{x}_{i}) - \mathbf{c}\|^{2} - r^{2} - \xi_{i} \right] - \sum_{i=1}^{m} \beta_{i} \xi_{i}.$$

Differentiating with respect to the primal variables gives:



ISSML, June 2013

The final equation implies that $\alpha_i \leq C$. Substituting, we obtain

$$L(\mathbf{c}, r, \alpha, \xi) = r^{2} + C \sum_{i=1}^{m} \xi_{i}$$

+
$$\sum_{i=1}^{m} \alpha_{i} \left[\|\phi(\mathbf{x}_{i}) - \mathbf{c}\|^{2} - r^{2} - \xi_{i} \right] - \sum_{i=1}^{m} \beta_{i} \xi_{i}$$

=
$$\sum_{i=1}^{m} \alpha_{i} \langle \phi(\mathbf{x}_{i}) - \mathbf{c}, \phi(\mathbf{x}_{i}) - \mathbf{c} \rangle$$

=
$$\sum_{i=1}^{m} \alpha_{i} \kappa \left(\mathbf{x}_{i}, \mathbf{x}_{i}\right) - \sum_{i,j=1}^{m} \alpha_{i} \alpha_{j} \kappa \left(\mathbf{x}_{i}, \mathbf{x}_{j}\right),$$

ISSML, June 2013

Hence optimisation maximise

$$W(\alpha) = \sum_{i=1}^{m} \alpha_i \kappa \left(\mathbf{x}_i, \mathbf{x}_i \right) - \sum_{i,j=1}^{m} \alpha_i \alpha_j \kappa \left(\mathbf{x}_i, \mathbf{x}_j \right)$$

subject to $\sum_{i=1}^{m} \alpha_i = 1$ and $0 \le \alpha_i \le C, i = 1, \dots, m$. with final novelty test being:

$$f(\cdot) = \mathcal{H}\left[\kappa\left(\cdot, \cdot\right) - 2\sum_{i=1}^{m} \alpha_{i}^{*}\kappa\left(\mathbf{x}_{i}, \cdot\right) + D\right]$$

where

$$D = \sum_{i,j=1}^{m} \alpha_i^* \alpha_j^* \kappa \left(\mathbf{x}_i, \mathbf{x}_j \right) - \left(r^* \right)^2 - \gamma$$

ISSML, June 2013



ISSML, June 2013

Final Boosting bound

 Applying a similar strategy for Boosting with the 1-norm of the slack variables we arrive at Linear programming boosting that minimises

$$\sum_{h} a_h + C \sum_{i=1}^{m} \xi_i,$$

where $\xi_i = (1 - y_i \sum_h a_h h(\mathbf{x}_i))_+$.

• with corresponding bound:

$$P(y \neq \operatorname{sgn}(g(\mathbf{x}))) = \mathbb{E} \left[\mathcal{H}(-yg(\mathbf{x})) \right]$$
$$\leq \frac{1}{m} \sum_{i=1}^{m} \xi_i + \hat{R}(H) \sum_h a_h + 3\sqrt{\frac{\ln(2/\delta)}{2m}}$$

ISSML, June 2013

Linear programming machine

- Controversy of why boosting works and relation to bagging.
- The previous bound suggests an optimisation similar to that of SVMs.
- seeks linear function in a feature space defined explicitly.
- For example using the 1-norm it seeks w to solve

 $\min_{\mathbf{w},b,\xi} \quad \|\mathbf{w}\|_1 + C \sum_{i=1}^m \xi_i$ subject to $y_i \left(\langle \mathbf{w}, \mathbf{x}_i \rangle + b \right) \ge 1 - \xi_i, \, \xi_i \ge 0,$ $i = 1, \dots, m.$

ISSML, June 2013

Linear programming boosting

 Very slight generalisation considers the features as a set H_{ij} of 'weak' learners (and include the constant function as one weak learner and negative of each weak learner):

$\min_{\mathbf{a}, \boldsymbol{\xi}}$	$\ \mathbf{a}\ _1 + C \sum_{i=1}^m \xi_i$
subject to	$y_i \mathbf{H}_i \mathbf{a} \ge 1 - \xi_i, \ \xi_i \ge 0, \ a_i \ge 0$ i = 1,, m.

ISSML, June 2013

Alternative version

• Can explicitly optimise margin with 1-norm fixed:

$\max_{\rho, \mathbf{a}, \xi}$	$ ho - D \sum_{i=1}^{m} \xi_i$
subject to	$y_i \mathbf{H}_i \mathbf{a} \ge \rho - \xi_i, \ \xi_i \ge 0, a_j \ge 0$
	$\sum_{j=1}^{n} a_j = 1.$

• Dual has the following form:

$$\begin{array}{ll} \min_{\beta,\mathbf{u}} & \beta \\ \text{subject to} & \sum_{\substack{i=1\\m \\ \sum_{i=1}^m u_i = 1, \ 0 \le u_i \le D. \end{array} \end{array} \\ \end{array}$$

ISSML, June 2013

Column generation

Can solve the dual linear programme using an iterative method:

- 1 initialise $u_i = 1/m, i = 1, \dots, m, \beta = \infty, J = \emptyset$
- 2 choose j^* that maximises $f(j) = \sum_{i=1}^{m} u_i y_i \mathbf{H}_{ij}$
- 3 if $f(j^*) \leq \beta$ solve primal restricted to J and exit

$$4 \quad J = J \cup \{j^\star\}$$

- 5 Solve dual restricted to set J to give u_i, β
- 6 Go to 2
- Note that u_i is a distribution on the examples
- Each j added acts like an additional weak learner
- f(j) is simply the weighted classification accuracy
- Hence gives 'boosting' algorithm with previous weights updated satisfying error bound
- Guaranteed convergence and soft stopping criteria

ISSML, June 2013

Multiple kernel learning

 MKL puts a 1-norm constraint on a linear combination of kernels:

$$\left\{\kappa(\mathbf{x}, \mathbf{z}) = \sum_{t=1}^{N} z_t \kappa_t(\mathbf{x}, \mathbf{z}) : z_t \ge 0, \sum_{t=1}^{N} z_t = 1\right\}$$

and trains an SVM while optimizing z_t – a convex problem, c.f. group Lasso.

• obtain corresponding bound:

 $P(y \neq \operatorname{sgn}(g(\mathbf{x})))$ $\leq \frac{1}{m\gamma} \sum_{i=1}^{m} \xi_i + \frac{1}{\gamma} \hat{R}_m \left(\bigcup_{t=1}^{N} \mathcal{F}_t\right) + 3\sqrt{\frac{\ln(2/\delta)}{2m}}$

ISSML, June 2013

Bounding MKL

Need a bound on

$$\hat{R}_m\left(\mathcal{F}=\bigcup_{t=1}^N\mathcal{F}_t\right)$$

where $\mathcal{F}_t = \{ \mathbf{x} \to \langle \mathbf{w}, \phi_t(\mathbf{x}) \rangle : \| \mathbf{w} \| \le 1 \}.$

• First note further applications of McDiarmid gives with probability $1 - \delta_0$ of a random selection of σ^* :

$$\hat{R}_{m}(\mathcal{F}) \leq \frac{2}{m} \sup_{f \in \mathcal{F}} \sum_{i=1}^{m} \sigma_{i}^{*} f(\mathbf{x}_{i}) + 4\sqrt{\frac{\ln(1/\delta_{t})}{2m}}$$
$$\frac{2}{m} \sup_{f \in \mathcal{F}_{t}} \sum_{i=1}^{m} \sigma_{i}^{*} f(\mathbf{x}_{i}) \leq \hat{R}_{m}(\mathcal{F}_{t}) + 4\sqrt{\frac{\ln(1/\delta_{t})}{2m}}$$

with probability $1 - \delta_t$

ISSML, June 2013

and

Bounding MKL

• Hence taking $\delta_t = \delta/2(N+1)$ for $t = 0, \dots, N$

$$\hat{R}_{m}\left(\mathfrak{F}=\bigcup_{t=1}^{N}\mathfrak{F}_{t}\right)$$

$$\leq \frac{2}{m}\sup_{f\in\mathfrak{F}}\sum_{i=1}^{m}\sigma_{i}^{*}f(\mathbf{x}_{i})+4\sqrt{\frac{\ln(2(N+1)/\delta)}{2m}}$$

$$\leq \frac{2}{m}\max_{1\leq t\leq N}\sup_{f\in\mathfrak{F}_{t}}\sum_{i=1}^{m}\sigma_{i}^{*}f(\mathbf{x}_{i})+4\sqrt{\frac{\ln(2(N+1)/\delta)}{2m}}$$

$$\leq \frac{2}{m}\max_{1\leq t\leq N}\hat{R}_{m}(\mathfrak{F}_{t})+8\sqrt{\frac{\ln(2(N+1)/\delta)}{2m}}$$

with probability $1 - \delta/2$.

ISSML, June 2013

Bounding MKL

 This gives an overall bound on the generalisation of MKL of

$$P(y \neq \operatorname{sgn}(g(\mathbf{x}))) \leq \frac{1}{m\gamma} \sum_{i=1}^{m} \xi_i + \frac{2}{\gamma m} \max_{1 \leq t \leq N} \operatorname{tr}(\mathbf{K}_t) + \frac{8\sqrt{\frac{\ln(2(N+1)/\delta)}{2m}} + 3\sqrt{\frac{\ln(4/\delta)}{2m}}$$

where \mathbf{K}_t is the t-th kernel matrix.

• Bound gives only a logarithmic (additive) dependence on the number of kernels.

ISSML, June 2013

Linear Programming MKL

• Column generation gives efficient MKL if we can pick the best weak learner in each \mathcal{F}_t efficiently:

$$\sup_{f \in \mathcal{F}_t} \sum_{i=1}^m u_i y_i f(\mathbf{x}_i) = \sup_{\mathbf{w}: \|\mathbf{w}\| \le 1} \sum_{i=1}^m u_i y_i \langle \mathbf{w}, \phi_t(\mathbf{x}_i) \rangle$$
$$= \sup_{\mathbf{w}: \|\mathbf{w}\| \le 1} \left\langle \mathbf{w}, \sum_{i=1}^m u_i y_i \phi_t(\mathbf{x}_i) \right\rangle$$
$$= \left\| \sum_{i=1}^m u_i y_i \phi_t(\mathbf{x}_i) \right\|$$
$$= \sqrt{\mathbf{u}' \mathbf{Y} \mathbf{K}_t \mathbf{Y} \mathbf{u}} =: N_t$$

easily computable from the kernel matrices (note that \mathbf{u} is sparse after first iteration and can also be chosen sparse at the start).

ISSML, June 2013

Linear Programming MKL

• The optimal weak learner from \mathcal{F}_t is realised by the weight vector that achieves the supremum

$$\mathbf{w} = \frac{\sum_{i=1}^{m} u_i y_i \phi_t(\mathbf{x}_i)}{\left\|\sum_{i=1}^{m} u_i y_i \phi_t(\mathbf{x}_i)\right\|}$$

which has dual representation:

$$\alpha_i = \frac{1}{N_t} u_i y_i$$

• Hence, can use the linear programming boosting approach to implement multiple kernel learning.

ISSML, June 2013

Part 5

- Kernel design strategies.
- Kernels for text and string kernels.
- Kernels for other structures.
- Kernels from generative models.

ISSML, June 2013

- Already seen some properties of kernels:
 - symmetric:

 $\kappa(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \langle \phi(\mathbf{z}), \phi(\mathbf{x}) \rangle = \kappa(\mathbf{z}, \mathbf{x})$

- kernel matrices psd:

$$\mathbf{u}'\mathbf{K}\mathbf{u} = \sum_{i,j=1}^{m} u_i u_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$
$$= \left\langle \sum_{i=1}^{m} u_i \phi(\mathbf{x}_i), \sum_{j=1}^{m} u_j \phi(\mathbf{x}_j) \right\rangle$$
$$= \left\| \left\| \sum_{i=1}^{m} u_i \phi(\mathbf{x}_i) \right\|^2 \ge 0$$

ISSML, June 2013

- These two properties are all that is required for a kernel function to be valid: symmetric and every kernel matrix is psd.
- Note that this is equivalent to all eigenvalues nonnegative – recall that eigenvalues of the kernel matrix measured the sum of the squares of the projections onto the eigenvector.
- If we have uncountable domains should also have continuity, though there are exceptions to this as well.

ISSML, June 2013

Proof outline:

• Define feature space as class of functions:

$$\mathcal{F} = \left\{ \sum_{i=1}^{m} \alpha_i \kappa(\mathbf{x}_i, \cdot) : m \in \mathbb{N}, \mathbf{x}_i \in X, \alpha_i \in \mathbb{R}, i = 1, \dots, m \right\}$$

- Linear space
- embedding given by

 $\mathbf{x} \mapsto \kappa(\mathbf{x}, \cdot)$

ISSML, June 2013

• inner product between

$$f(\mathbf{x}) = \sum_{i=1}^{m} \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}) \text{ and } g(\mathbf{x}) = \sum_{i=1}^{n} \beta_i \kappa(\mathbf{z}_i, \mathbf{x})$$

defined as

$$\langle f,g\rangle = \sum_{i=1}^{m} \sum_{j=1}^{n} \alpha_i \beta_j \kappa(\mathbf{x}_i, \mathbf{z}_j) = \sum_{i=1}^{m} \alpha_i g(\mathbf{x}_i) = \sum_{j=1}^{n} \beta_j f(\mathbf{z}_j),$$

- well-defined
- $\langle f, f \rangle \geq 0$ by psd property.

ISSML, June 2013

• so-called reproducing property:

$$\langle f, \phi(\mathbf{x}) \rangle = \langle f, \kappa(\mathbf{x}, \cdot) \rangle = f(\mathbf{x})$$

 implies that inner product corresponds to function evaluation – learning a function corresponds to learning a point being the weight vector corresponding to that function:

 $\langle \mathbf{w}_f, \phi(\mathbf{x}) \rangle = f(\mathbf{x})$

ISSML, June 2013

Kernel constructions

For κ_1, κ_2 valid kernels, ϕ any feature map, **B** psd matrix, $a \ge 0$ and f any real valued function, the following are valid kernels:

- $\kappa(\mathbf{x}, \mathbf{z}) = \kappa_1(\mathbf{x}, \mathbf{z}) + \kappa_2(\mathbf{x}, \mathbf{z}),$
- $\kappa(\mathbf{x}, \mathbf{z}) = a\kappa_1(\mathbf{x}, \mathbf{z}),$
- $\kappa(\mathbf{x}, \mathbf{z}) = \kappa_1(\mathbf{x}, \mathbf{z})\kappa_2(\mathbf{x}, \mathbf{z}),$
- $\kappa(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})f(\mathbf{z}),$
- $\kappa(\mathbf{x}, \mathbf{z}) = \kappa_1(\phi(\mathbf{x}), \phi(\mathbf{z})),$
- $\kappa(\mathbf{x}, \mathbf{z}) = \mathbf{x}' \mathbf{B} \mathbf{z}.$

ISSML, June 2013

Kernel constructions

Following are also valid kernels:

- $\kappa(\mathbf{x}, \mathbf{z}) = p(\kappa_1(\mathbf{x}, \mathbf{z}))$, for *p* any polynomial with positive coefficients.
- $\kappa(\mathbf{x}, \mathbf{z}) = \exp(\kappa_1(\mathbf{x}, \mathbf{z})),$
- $\kappa(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} \mathbf{z}\|^2 / (2\sigma^2)).$

Proof of third: normalise the second kernel:

$$\frac{\exp(\langle \mathbf{x}, \mathbf{z} \rangle / \sigma^2)}{\sqrt{\exp(\|\mathbf{x}\|^2 / \sigma^2)} \exp(\|\mathbf{z}\|^2 / \sigma^2)} = \exp\left(\frac{\langle \mathbf{x}, \mathbf{z} \rangle}{\sigma^2} - \frac{\langle \mathbf{x}, \mathbf{x} \rangle}{2\sigma^2} - \frac{\langle \mathbf{z}, \mathbf{z} \rangle}{2\sigma^2}\right)$$
$$= \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2}\right).$$

ISSML, June 2013

Subcomponents kernel

For the kernel $\langle \mathbf{x}, \mathbf{z} \rangle^s$ the features can be indexed by sequences

$$\mathbf{i} = (i_1, \dots, i_n), \sum_{j=1}^n i_j = s$$

where

$$\phi_{\mathbf{i}}(\mathbf{x}) = x_1^{i_1} x_2^{i_2} \dots x_n^{i_n}$$

A similar kernel can be defined in which all subsets of features occur:

 $\phi: \mathbf{x} \mapsto (\phi_A(\mathbf{x}))_{A \subseteq \{1, \dots, n\}}$

where

$$\phi_A(\mathbf{x}) = \prod_{i \in A} x_i$$

ISSML, June 2013

Subcomponents kernel

So we have

$$\kappa_{\subseteq}(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$$

=
$$\sum_{A \subseteq \{1, \dots, n\}} \phi_A(\mathbf{x}) \phi_A(\mathbf{y})$$

=
$$\sum_{A \subseteq \{1, \dots, n\}} \prod_{i \in A} x_i y_i = \prod_{i=1}^n (1 + x_i y_i)$$

Can represent computation with a graph:



Each path in the graph corresponds to a feature.

ISSML, June 2013

Can also represent polynomial kernel

 $\kappa(\mathbf{x}, \mathbf{y}) = \left(\langle \mathbf{x}, \mathbf{y} \rangle + R\right)^d = \left(x_1 y_1 + x_2 y_2 + \dots + x_n y_n + R\right)^d$



ISSML, June 2013

The ANOVA kernel is represented by the graph:



ISSML, June 2013

Features are all the combinations of exactly d distinct features, while computation is given by recursion:

$$\begin{aligned} \kappa_0^m(\mathbf{x}, \mathbf{z}) &= 1, \text{ if } m \ge 0, \\ \kappa_s^m(\mathbf{x}, \mathbf{z}) &= 0, \text{ if } m < s, \\ \kappa_s^m(\mathbf{x}, \mathbf{z}) &= (x_m z_m) \kappa_{s-1}^{m-1}(\mathbf{x}, \mathbf{z}) + \kappa_s^{m-1}(\mathbf{x}, \mathbf{z}) \end{aligned}$$

While the resulting kernel is given by

 $\kappa_d^n(\mathbf{x},\mathbf{z})$

in the bottom right corner of the graph.

ISSML, June 2013

- Initialise DP(1) = 1;
- for each node compute

$$DP(i) = \sum_{j \to i} \kappa_{\left(u_j \to u_i\right)}(\mathbf{x}, \mathbf{z}) DP(j)$$

• result given at output node s: $\kappa(\mathbf{x}, \mathbf{z}) = DP(s)$.

ISSML, June 2013

Kernels for text

 The simplest representation for text is the kernel given by the feature map known as the vector space model

 $\phi: d \mapsto \phi(d) = (\mathrm{tf}(t_1, d), \mathrm{tf}(t_2, d), \dots, \mathrm{tf}(t_N, d))'$

where t_1, t_2, \ldots, t_N are the terms occurring in the corpus and tf(t, d) measures the frequency of term t in document d.

• Usually use the notation **D** for the document term matrix (cf. **X** from previous notation).

ISSML, June 2013

Kernels for text

• Kernel matrix is given by

 $\mathbf{K} = \mathbf{D}\mathbf{D}'$

wrt kernel

$$\kappa(d_1, d_2) = \sum_{j=1}^N \operatorname{tf}(t_j, d_1) \operatorname{tf}(t_j, d_2)$$

 despite high-dimensionality kernel function can be computed efficiently by using a linked list representation.

ISSML, June 2013
- The standard representation does not take into account the importance or relationship between words.
- Main methods do this by introducing a 'semantic' mapping S:

 $\hat{\kappa}(d_1, d_2) = \phi(d_1)' \mathbf{S} \mathbf{S}' \phi(d_2)$

ISSML, June 2013

• Simplest is diagonal matrix giving term weightings (known as inverse document frequency – tfidf):

$$w(t) = \ln \frac{m}{\mathrm{df}(t)}$$

• Hence kernel becomes:

$$\kappa(d_1, d_2) = \sum_{j=1}^{N} w(t_j)^2 \mathrm{tf}(t_j, d_1) \mathrm{tf}(t_j, d_2)$$

ISSML, June 2013

- In general would also like to include semantic links between terms with off-diagonal elements, eg stemming, query expansion, wordnet.
- More generally can use co-occurrence of words in documents:

 $\mathbf{S}=\mathbf{D}'$

SO

$$(\mathbf{SS'})_{ij} = \sum_d \operatorname{tf}(i,d) \operatorname{tf}(j,d)$$

ISSML, June 2013

• Information retrieval technique known as latent semantic indexing uses SVD decomposition:

$$\mathbf{D}' = \mathbf{U} \mathbf{\Sigma} \mathbf{V}'$$

so that

$$d \mapsto \mathbf{U}'_k \phi(d)$$

which is equivalent to peforming kernel PCA to give latent semantic kernels:

 $\tilde{\kappa}(d_1, d_2) = \phi(d_1)' \mathbf{U}_k \mathbf{U}'_k \phi(d_2)$

ISSML, June 2013

String kernels

• Consider the feature map given by

 $\phi_u^p(s) = |\{(v_1, v_2) : s = v_1 u v_2\}|$

for $u \in \Sigma^p$ with associated kernel

$$\kappa_p(s,t) = \sum_{u \in \Sigma^p} \phi_u^p(s) \phi_u^p(t)$$

ISSML, June 2013

String kernels

• Consider the following two sequences:

s ="statistics" t ="computation"

The two strings contain the following substrings of length 3:

"sta", "tat", "ati", "tis",
"ist", "sti", "tic", "ics"
"com", "omp", "mpu", "put",
"uta", "tat", "ati", "tio", "ion"

and they have in common the substrings "tat" and "ati", so their inner product would be $\kappa(s,t) = 2$.

ISSML, June 2013

Trie based p-spectrum kernels

- Computation organised into a trie with nodes indexed by substrings – root node by empty string *ε*.
- Create lists of substrings at root node:

 $L_s(\epsilon) = \{(s(i:i+p-1), 0): i = 1, |s|-p+1\}$

Similarly for t.

- Recursively through the tree: if $L_s(v)$ and $L_t(v)$ both not empty: for each $(u,i) \in L_*(v)$ add (u,i+1) to list $L_*(vu_{i+1})$
- At depth p increment global variable kern initialised to 0 by $|L_s(v)||L_t(v)|$.

ISSML, June 2013

Gap weighted string kernels

 Can create kernels whose features are all substrings of length p with the feature weighted according to all occurrences of the substring as a subsequence:

ϕ	са	ct	at	ba	bt	cr	ar	br
cat	λ^2	λ^3	λ^2	0	0	0	0	0
car	λ^2	0	0	0	0	λ^3	λ^2	0
bat	0	0	λ^2	λ^2	λ^3	0	0	0
bar	0	0	0	λ^2	0	0	λ^2	λ^3

• This can be evaluated using a dynamic programming computation over arrays indexed by the two strings.

ISSML, June 2013

Tree kernels

• We can consider a feature mapping for trees defined by

 $\phi: T \longmapsto (\phi_S(T))_{S \in I}$

where *I* is a set of all subtrees and $\phi_S(T)$ counts the number of co-rooted subtrees isomorphic to the tree *S*.

- The computation can again be performed efficiently by working up from the leaves of the tree integrating the results from the children at each internal node.
- Similarly we can compute the inner product in the feature space given by all subtrees of the given tree not necessarily co-rooted.

ISSML, June 2013

Probabilistic model kernels

- There are two types of kernels that can be defined based on probabilistic models of the data.
- The most natural is to consider a class of models index by a model class *M*: we can then define the similarity as

$$\kappa(x,z) = \sum_{m \in M} P(x|m)P(z|m)P_M(m)$$

also known as the marginalisation kernel.

 For the case of Hidden Markov Models this can be again be computed by a dynamic programming technique.

ISSML, June 2013

Probabilistic model kernels

- Pair HMMs generate pairs of symbols and under mild assumptions can also be shown to give rise to kernels that can be efficiently evaluated.
- Similarly hidden tree generating models of data, again using a recursion that works upwards from the leaves.

ISSML, June 2013

Fisher kernels

Fisher kernels are an alternative way of defining kernels based on probabilistic models.

• We assume the model is parametrised according to some parameters: consider the simple example of a 1-dim Gaussian distribution parametrised by μ and σ :

$$M = \left\{ P(x|\theta) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) : \theta = (\mu,\sigma) \in \mathbb{R}^2 \right\}$$

• The Fisher score vector is the derivative of the log likelihood of an input *x* wrt the parameters:

$$\log \mathcal{L}_{(\mu,\sigma)}(x) = -\frac{\left(x-\mu\right)^2}{2\sigma^2} - \frac{1}{2}\log\left(2\pi\sigma\right).$$

ISSML, June 2013

Fisher kernels

• Hence the score vector is given by:

$$\mathbf{g}(\theta^{0}, x) = \left(\frac{(x - \mu_{0})}{\sigma_{0}^{2}}, \frac{(x - \mu_{0})^{2}}{\sigma_{0}^{3}} - \frac{1}{2\sigma_{0}}\right).$$

• Taking $\mu_0 = 0$ and $\sigma_0 = 1$ the feature embedding is given by:

ISSML, June 2013



ISSML, June 2013

Fisher kernels

Can compute Fisher kernels for various models including

- ones closely related to string kernels
- Hidden Markov Models

ISSML, June 2013

Conclusions

Kernel methods provide a general purpose toolkit for pattern analysis

- kernels define flexible interface to the data enabling the user to encode prior knowledge into a measure of similarity between two items – with the proviso that it must satisfy the psd property.
- composition and subspace methods provide tools to enhance the representation: normalisation, centering, kernel PCA, kernel Gram-Schmidt, kernel CCA, etc.
- algorithms well-founded in statistical learning theory enable efficient and effective exploitation of the high-dimensional representations to enable good off-training performance.

ISSML, June 2013

Where to find out more

```
Web Sites: www.support-vector.net (SV Machines)
```

www.kernel-methods.net (kernel methods)

www.kernel-machines.net (kernel Machines)

www.neurocolt.com (Neurocolt: lots of TRs)

www.pascal-network.org

References

 N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive Dimensions, Uniform Convergence, and Learnability. *Journal of the ACM*, 44(4):615–631, 1997.

ISSML, June 2013

- [2] M. Anthony and P. Bartlett. Neural Network Learning: Theoretical Foundations. Cambridge University Press, 1999.
- [3] M. Anthony and N. Biggs. *Computational Learning Theory*, volume 30 of *Cambridge Tracts in Theoretical Computer Science*. Cambridge University Press, 1992.
- [4] M. Anthony and J. Shawe-Taylor. A result of Vapnik with applications. *Discrete Applied Mathematics*, 47:207–217, 1993.
- [5] K. Azuma. Weighted sums of certain dependent random variables. *Tohoku Math J.*, 19:357–367, 1967.
- [6] P. Bartlett and J. Shawe-Taylor. Generalization performance of support vector machines and other pattern classifiers. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods* — *Support Vector Learning*, pages 43–54, Cambridge, MA, 1999. MIT Press.
- [7] P. L. Bartlett. The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network.

ISSML, June 2013

IEEE Transactions on Information Theory, 44(2):525–536, 1998.

- [8] P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: risk bounds and structural results. *Journal of Machine Learning Research*, 3:463– 482, 2002.
- [9] S. Boucheron, G. Lugosi, , and P. Massart. A sharp concentration inequality with applications. *Random Structures and Algorithms*, pages vol.16, pp.277–292, 2000.
- [10] O. Bousquet and A. Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- [11] N. Cristianini and J. Shawe-Taylor. *An introduction to Support Vector Machines*. Cambridge University Press, Cambridge, UK, 2000.
- [12] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. In *Computational Learning Theory: Eurocolt* '95, pages 23–37. Springer-Verlag, 1995.

ISSML, June 2013

- [13] W. Hoeffding. Probability inequalities for sums of bounded random variables. J. Amer. Stat. Assoc., 58:13–30, 1963.
- [14] M. Kearns and U. Vazirani. An Introduction to Computational Learning Theory. MIT Press, 1994.
- [15] V. Koltchinskii and D. Panchenko. Rademacher processes and bounding the risk of function learning. *High Dimensional Probability II*, pages 443 – 459, 2000.
- [16] J. Langford and J. Shawe-Taylor. PAC bayes and margins. In *Advances in Neural Information Processing Systems 15*, Cambridge, MA, 2003. MIT Press.
- [17] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer, 1991.
- [18] C. McDiarmid. On the method of bounded differences. In 141 London Mathematical Society Lecture Notes Series, editor, *Surveys in Combinatorics 1989*, pages 148–188. Cambridge University Press, Cambridge, 1989.
- [19] R. Schapire, Y. Freund, P. Bartlett, and W. Sun Lee. Boosting the margin: A new explanation for the effectiveness of voting methods. *Annals of Statistics*,

ISSML, June 2013

1998. (To appear. An earlier version appeared in: D.H. Fisher, Jr. (ed.), Proceedings ICML97, Morgan Kaufmann.).

- [20] J. Shawe-Taylor, P. L. Bartlett, R. C. Williamson, and M. Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 44(5):1926–1940, 1998.
- [21] J. Shawe-Taylor and N. Cristianini. On the generalisation of soft margin algorithms. *IEEE Transactions on Information Theory*, 48(10):2721–2735, 2002.
- [22] J. Shawe-Taylor and N. Cristianini. *Kernel Methods* for Pattern Analysis. Cambridge University Press, Cambridge, UK, 2004.
- [23] J. Shawe-Taylor, C. Williams, N. Cristianini, and J. S. Kandola. On the eigenspectrum of the gram matrix and its relationship to the operator eigenspectrum. In *Proceedings of the 13th International Conference on Algorithmic Learning Theory (ALT2002)*, volume 2533, pages 23–40, 2002.
- [24] M. Talagrand. New concentration inequalities in product

ISSML, June 2013

spaces. Invent. Math., 126:505–563, 1996.

- [25] V. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.
- [26] V. Vapnik and A. Chervonenkis. Uniform convergence of frequencies of occurrence of events to their probabilities. *Dokl. Akad. Nauk SSSR*, 181:915 – 918, 1968.
- [27] V. Vapnik and A. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–280, 1971.
- [28] Tong Zhang. Covering number bounds of certain regularized linear function classes. *Journal of Machine Learning Research*, 2:527–550, 2002.

ISSML, June 2013