# Lecture 14. Kernel Methods

## COMP90051 Statistical Machine Learning

Semester 2, 2015
Lecturer: Andrey Kan

THE UNIVERSITY OF
MELBOURNE

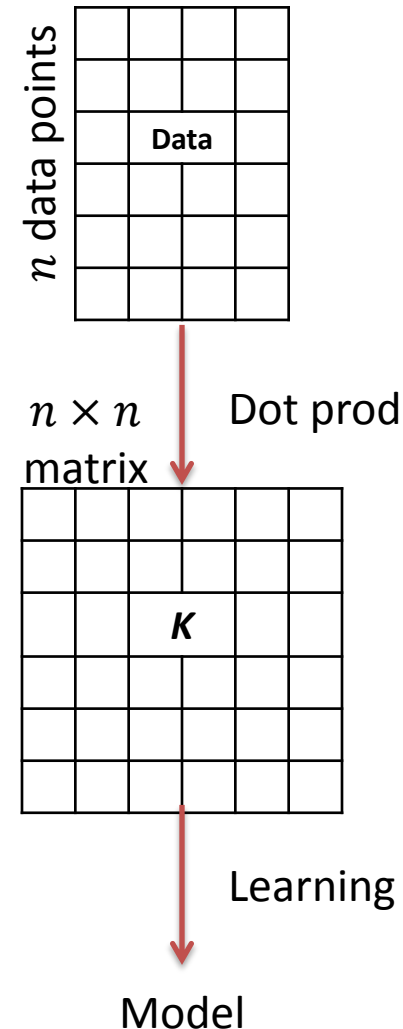POSTERA CRESCAM LAUDE

# Kernel Methods

*Very general family of linear techniques that can be made non-linear in a multitude of ways*

# Overview

- ## Kernel matrix $K$
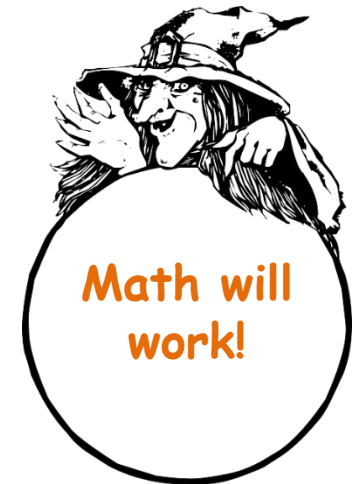  - Square $n$ x $n$ matrix measuring pairwise similarity
  - Entries $K_{ij} = \boldsymbol{x}_i \cdot \boldsymbol{x}_j$

- ## Kernel methods
  - Linear methods relying on training data only through $\boldsymbol{K}$
  - **Make non-linear** by running linear approach in new feature space; but don't need to map the data, just need $\boldsymbol{K}$

- ## **Modular**: learning algorithm, feature space, decouple
  - Can design general learning algorithms
  - Can design general feature mappings/kernels

$n$ data points

Data

$n \times n$ matrix    Dot prod

$K$

Learning

Model

# Dot products, dot products, …

- How does the SVM depend on the data?

$$\min_{\boldsymbol{w}} \left( \frac{1}{2} \|\boldsymbol{w}\|^2 + \frac{C}{n} \sum_{i=1}^{n} l(1 - y_i \boldsymbol{w} \cdot \boldsymbol{x}_i - b) \right)$$
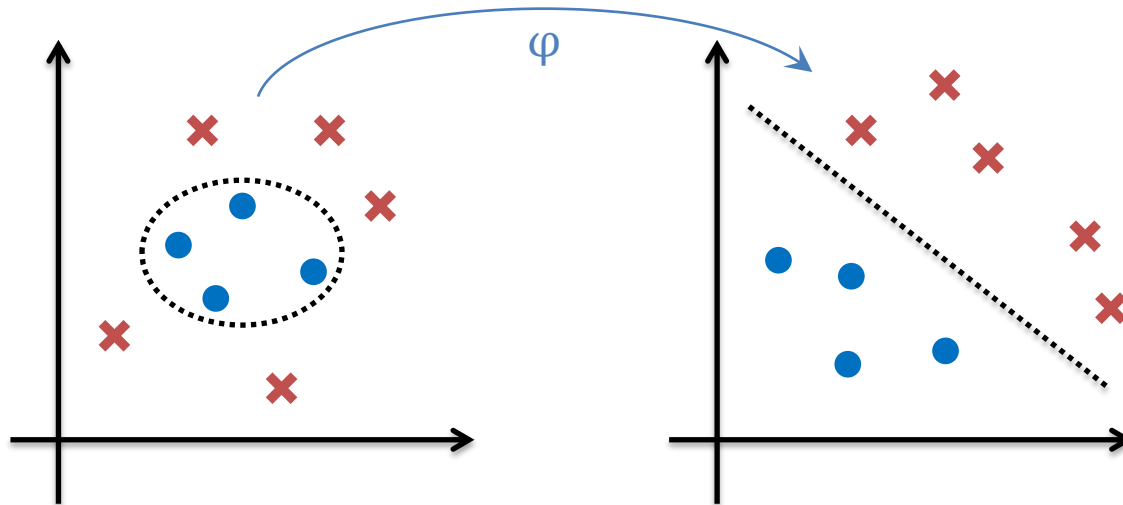
**Math will work!**

- Famous result: "Representer Theorem"
  - ∗ Solution in span of data!!  $\boldsymbol{w}^* = \sum_{i=1}^{n} \alpha_i y_i \boldsymbol{x}_i$
  - ∗ Predictions become $f(\boldsymbol{x}) = \boldsymbol{w} \cdot \boldsymbol{x} = \sum_{i=1}^{n} \alpha_i y_i \boldsymbol{x}_i \cdot \boldsymbol{x}$
  - ∗ Support vectors are those $\boldsymbol{x}_i$ with $\alpha_i \neq 0$
    → why the SVM is non-parametric!
  - ∗ Finding the $\alpha_i$ involves only dot products between data

4

# Non-linear SVM

- Map data into a new feature space

    * Example: original features and products of pairs
    $$\varphi(x_1, \ldots, x_d) = (x_1, \ldots, x_d, x_1 x_1, x_1 x_2, \ldots, x_d x_d)$$

- Run linear SVM in new space (use kernel matrix)

- Decision boundary is non-linear in original space

# Flexibility of non-linear SVM



www.kernel-methods.net

# Blessing of dimensionality

- Kernels implicitly map data to a very high dimensional features space

- Potentially dangerous because it looks like we now have many more parameters than data points
  * That is, $\boldsymbol{w}$ for the transformed features space is high dimensional
  * Curse of dimensionality
  * Danger of overfitting

- Representer theorem: $f(\boldsymbol{x}) = \sum_{i=1}^{n} \alpha_i y_i \boldsymbol{x}_i \cdot \boldsymbol{x}$
  * The number of parameters is at most $n$, independent of dimensionality
  * Support vectors are those data points with non-zero $\alpha_i$

- Usually the number of non-zero $\alpha_i$ is smaller than $n$
  * Sparse kernel machines

# *Danger!!*   Potential problem

**Stop!**

- Training linear SVM cubic in dimension *d*

- Non-linear → more dimensions → intractable?

- Representer Theorem to the rescue
  - ∗ Need only find $\alpha_i$'s which determine weight vector
  - ∗ "Only" *n* of them, independent of *d*
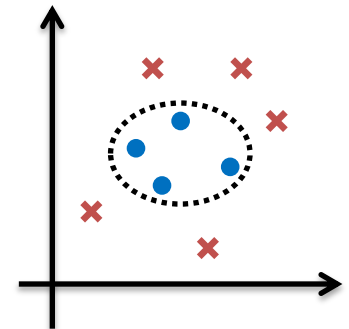    *…* problematic if "big data"

# But wait… what about the kernel matrix?

**Stop!**

- *Computing* kernel matrix naïvely expensive
  * Map data to $d'$-dim feature space, then dot product; takes $O(d'n^2)$

- Computing kernels *directly* can be **cheap as**

- Example: $p$-degree polynomial kernel
  * $\varphi(x_1, \ldots, x_d) = (x_1, \ldots, x_d, x_1 x_1, x_1 x_2, \ldots, x_d x_d, \ldots)$
  * $d' = O(d^p)$ is pretty yuck $\rightarrow$ $O(d^p)$ per matrix entry
  * Trick to cut down to $O(d)$:  $\varphi(\boldsymbol{u}) \cdot \varphi(\boldsymbol{v}) = (1 + \boldsymbol{u} \cdot \boldsymbol{v})^p$

- In fact, don't bother with feature spaces just make up matrix
  * Mercer's Theorem: provides a tool for identifying valid kernels

9

# Some popular kernels

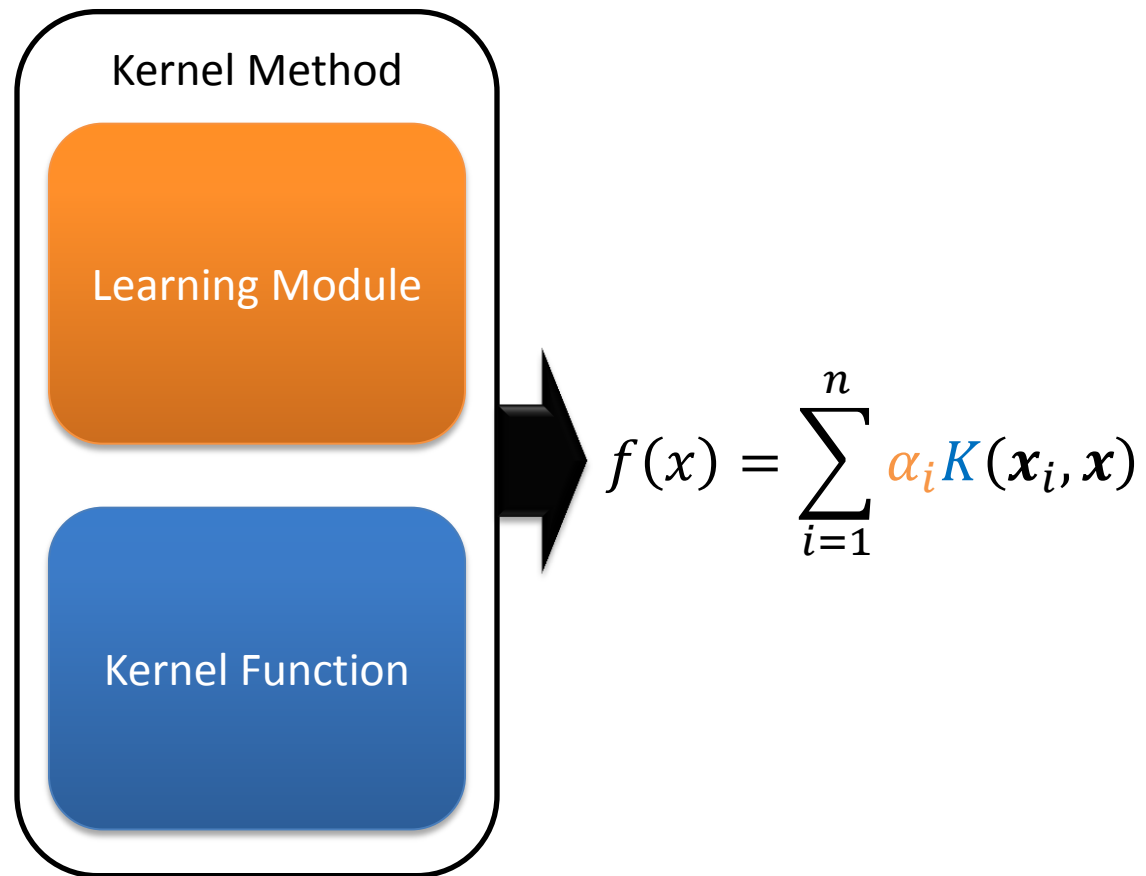| Kernel | $k(u, v)$ | Parameters | Mapping $x$ to features |
|---|---|---|---|
| Linear | $u \cdot v$ | - | Identity |
| Polynomial | $(1 + u \cdot v)^p$ | Integer degree p>1 | Degree $p$ polynomial in the $x_i$'s |
| Gaussian aka Radial Basis Function (RBF) | $\exp\left(\dfrac{-\|u - v\|^2}{2\sigma^2}\right)$ | Width $\sigma$ | Infinite dimensional, nothing explicit! |

# Many more…

www.kernel-methods.net

# Modular learning

- More sophisticated learning algorithms, and kernels, exist

- Representer + Mercer Theorems imply their design decouples

Kernel Method

Learning Module

Kernel Function

$$f(x) = \sum_{i=1}^{n} \alpha_i K(\boldsymbol{x}_i, \boldsymbol{x})$$

# ...and the algorithms

www.kernel-methods.net

# Checkpoint

- Which of the following statements is true?

    Any method that uses a feature space transformation $\varphi(\boldsymbol{x})$ is a kernel method

    Support vectors are points from the training set

    Feature mapping $\varphi(\boldsymbol{x})$ makes data linearly separable

art: OpenClipartVectors at
pixabay.com (CC0)

# Kernelised SVM

*Historically first kernel method*

# Feature transformations and kernels

- A feature transformation $\varphi(x)$ can help transforming data into a linearly separable form

- The kernel trick can be used to perform the transformation implicitly, without actually computing the transformation for each point

- One <u>does not have to use kernels</u>: feature transformation $\varphi(x)$ can be also use explicitly
  - * If the number of resulting features is reasonable

# Solutions to the SVM problem

- Recall the soft-margin SVM problem statement:

$$\min_{\boldsymbol{w}} \left( \frac{1}{2} \|\boldsymbol{w}\|^2 + \frac{C}{2} \sum_{i=1}^{n} \xi_i \right)$$

Subject to $\xi_i \geq 0$, $y_i(\boldsymbol{w}\boldsymbol{x}_i + b) \geq 1 - \xi_i$ for $i = 1, \dots, n$

- This optimization problem (and many others) can be solved in two ways, called the primal formulation and dual formulation

- The primal approach is to solve the problem as it is stated here
  * If you wish, you can first apply $\varphi(\boldsymbol{x})$ here directly

- The complexity of training is $O(d^3)$, where $d$ is the dimensionality of data
  * Recall that [solving optimization = training]
  * If you have applied $\varphi(\boldsymbol{x})$ then $d$ is the dimensionality of resulting space

# The dual formulation

- The idea of the dual formulation is to combine equations for the objective function and constraints into a single new objective function without constraints

  * This can be analytically convenient

  * Also this can lead to a different perspective on the problem

- This method is called *Lagrange multipliers*, and the corresponding new objective function is called *Lagrangian*

18

# Formulating the dual SVM problem

- Soft-margin SVM original problem:

$$\min_{\boldsymbol{w},b} \left( \frac{1}{2} \|\boldsymbol{w}\|^2 + \frac{C}{n} \sum_{i=1}^{n} \xi_i \right)$$

  * Subject to $\xi_i \geq 0$, $y_i(\boldsymbol{w}\boldsymbol{x}_i + b) \geq 1 - \xi_i$ for $i = 1, \dots, n$

- Lagrangian

$$L_P \stackrel{\text{def}}{=} \frac{1}{2} \|\boldsymbol{w}\|^2 + C \sum_{i=1}^{n} \xi_i - \sum_{i=1}^{n} \alpha_i [y_i(\boldsymbol{w}\boldsymbol{x}_i + b) \geq 1 - \xi_i] - \sum_{i=1}^{n} \mu_i \xi_i$$

- Fix $\alpha_i, \mu_i$ and minimise Lagrangian with respect to $\boldsymbol{w}, b, \xi$:

  * Set partial derivatives $\frac{\partial L_P}{\partial w_i}, \frac{\partial L_P}{\partial b}, \frac{\partial L_P}{\partial \xi_i}$ to zero

# Solving the SVM problem

- Setting the partial derivatives to zero gives

  * $w^* = \sum_{i=1}^{n} \alpha_i y_i \boldsymbol{x}_i$

  * $0 = \sum_{i=1}^{n} \alpha_i y_i$

  * $\alpha_i = C/n - \mu_i$

  * $\alpha_i, \mu_i, \xi_i \geq 0$

- Substituting these back into $L_P$ gives

  * $L_D = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \boldsymbol{x}_i \boldsymbol{x}_j$

**Proofs are outside the scope**

- $L_D$ gives a lower bound on the original solution

- Therefore we maximize $L_D$ w.r.t. $\alpha_i$, s.t. $0 = \sum_{i=1}^{n} \alpha_i y_i$ and $0 \leq \alpha_i \leq \frac{C}{n}$

- In addition, Karush-Kuhn-Tucker conditions complete unique characterization of solution

# About those alpha's

- Transform

$$\min_{\boldsymbol{w}} \frac{1}{2} \|\boldsymbol{w}\|^2 + \frac{C}{2} \sum_{i=1}^{n} \xi_i$$

into "dual problem"

$$\max_{\boldsymbol{\alpha}} \ \boldsymbol{\alpha} \cdot \boldsymbol{1} - \frac{1}{2} \boldsymbol{\alpha}^T G \boldsymbol{\alpha}$$

subject to $\ 0 \le \alpha_i \le \frac{C}{n}$

where $G_{ij} = y_i y_j \boldsymbol{x}_i \cdot \boldsymbol{x}_j = y_i y_j K_{ij}$
Gram matrix                                          Kernel matrix

*Want different non-linear mapping? Swap out K*

- Quadratic program but in *n*, not *d* variables

- See how regularisation restricts data influence?
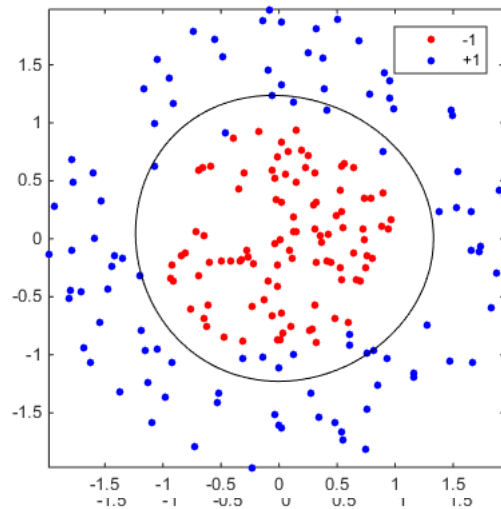
21

# Solutions to the SVM problem

- The complexity of the primal solution is $O(d^3)$, where $d$ is the dimensionality of data (after applying $\varphi$)

- The complexity of the dual solution is $O(n^3)$

- For "big data", and finite-dimensional explicit transformation $\varphi(\boldsymbol{x})$ the primal form can be preferable

# SVM: decisions to make

- Regularisation parameter $C$

- Choice of a kernel

- Note that some kernels introduce additional parameters, e.g., $p$ in polynomial kernel $(1 + \boldsymbol{u} \cdot \boldsymbol{v})^p$


- SVM is a sparse kernel machine

- Regularisation also helps avoid the curse of dimensionality

- In practice: everything is set using cross-validation

# Varying regularization parameter

Low $C$, tendency to underfit

High $C$, tendency to overfit



Polynomial kernel is used in all cases                                    24

# SVMs for multiclass classification

- Multiclass classification has to be reduced to binary classification

- **One-versus-all**: $K$ binary classifications are made for each of the $K$ classes. Each time, classification is "class $k$" vs "the rest". A new instance is assigned to class for which it has the strongest confidence (furthest away from the boundary).

- **One-versus-one**: $K(K-1)/2$ classifications are made, covering all possible pairs of classes. A new instance is assigned the most voted class (class that gets most number of "wins").

# SVM vs Neural Networks

| Kernelised SVM | Feed-forward ANN |
|---|---|
| Linear SVM applied to non-linear data using kernelisation | Non-linear model by design (non-linear transfer functions) |
| Variable number of parameters, up to N (data points): predictions are made using support vectors | Parametric model: fixed number of parameters |
| Inherently batch training (see dual formulation, Kernel matrix) | Naturally amenable to online training (e.g., stochastic gradient descent) |
| SVM training always finds a global minimum (convex problem) | Global minimum is not guaranteed |
| Single dimensional output | Naturally adaptable for multiclass classification and multidimensional output |

# Summary

- **Kernel methods**
  - ∗ Modular decoupling of linear learner and feature mapping
  - ∗ Built on Representer Theorem, Mercer's Theorem
  - ∗ Tonnes of kernel functions
  - ∗ Tonnes of kernel methods (SVM, PCA)