

CS:E4830 Kernel Methods in Machine Learning

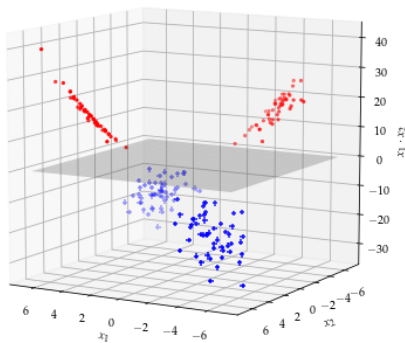
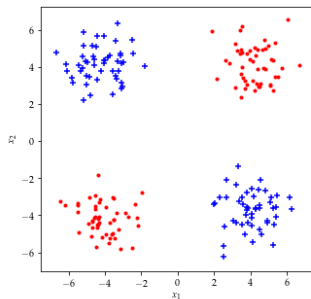
Lecture 12 : Course Review

Rohit Babbar

3rd April, 2019

Explicit Feature Mapping

- Dataset in 2-D (left), which is not linearly separable can be separated by a plane in 3-D (third feature is the product x_1x_2)



Kernel Methods - Motivation

- Most learning algorithms such as Support Vector Machines, and Logistic regression (classification part used in deep networks) can be written in the form of Inner/dot product between vectors in the feature space $\phi(x_i)$ s, i.e. $\langle \phi(x_i), \phi(x_j) \rangle$.
- The prediction function has the following form :

$$f(x) = \text{Some function of } \left(\sum_{i=1}^n \langle \phi(x_i), \phi(x) \rangle \right)$$

- Kernels are functions which give us the dot product $\langle \phi(x_i), \phi(x_j) \rangle$ directly without explicitly computing the feature expansion $\phi(\cdot)$

Properties of Kernels

- **Positive Scalar Multiple** - For any $\alpha > 0$, if $k(.,.)$ is a kernel, then $\alpha k(.,.)$ is also a kernel.
- **Conic Sum of Kernels** - For kernels $(k_j)_{j=1}^K$, and $(\alpha_j)_{j=1}^K > 0$, $\sum_{j=1}^K \alpha_j k_j$ is also a kernel
- **Difference of Kernels** is not necessarily a kernel
- **Product** - product of kernels is also a kernel
- **Mappings** - For an arbitrary function $f : \mathcal{X} \mapsto \mathbb{R}$, and a kernel $k(.,.)$, $\hat{k}(x, x') = f(x)k(x, x')f(x')$ is also a kernel
- Used above properties to prove that polynomial, exponential, and Gaussian kernels are valid kernel functions.

Positive Definite Functions

Definition - Positive definite functions

A symmetric function $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ is positive definite if $\forall n \geq 1, \forall (a_1, \dots, a_n) \in \mathbb{R}^n, \forall (x_1, \dots, x_n) \in \mathcal{X}^n$,

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j k(x_i, x_j) \geq 0$$

Moore-Aronszajn Theorem

A function $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ is a kernel if and only if it is symmetric and positive definite.

The kernel matrix

- A **kernel matrix** (also called the **Gram matrix**), is an $N \times N$ matrix of pairwise similarity values is used:

$$K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_N) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \dots & k(x_N, x_N) \end{bmatrix}$$

- Each entry is an inner product between two data points
 $k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$, where $\phi(\cdot)$ is a feature map in vector form
- Since an inner product is symmetric, therefore K is a symmetric matrix
- In addition, K is positive definite

Definition (RKHS)

Let \mathcal{H} be a Hilbert space of real-valued **functions** on the input \mathcal{X} . Then $\mathcal{H}(\subset \mathcal{R}^{\mathcal{X}})$ is defined to an **Reproducing kernel Hilbert Space (RKHS)** with $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ as the reproducing kernel, if the following conditions are satisfied

- $\forall x \in \mathcal{X}, k(., x) \in \mathcal{H}$ i.e., the space \mathcal{H} contains all functions of the form $k(., x)$ for every element x in the input space \mathcal{X} ,
- $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}$, the following property holds : $f(x) = \langle f, k(., x) \rangle_{\mathcal{H}}$ (it is called the reproducing property of the kernel).

Definition (RKHS)

Let \mathcal{H} be a Hilbert space of real-valued **functions** on the input \mathcal{X} . Then $\mathcal{H}(\subset \mathcal{R}^{\mathcal{X}})$ is defined to be an **Reproducing kernel Hilbert Space (RKHS)** if and only if, for any element x in the input space \mathcal{X} , the following function F , which takes a function f from the Hilbert Space \mathcal{H} , and maps it to its value $f(x) \in \mathbb{R}$

$$\begin{aligned} F : \quad \mathcal{H} &\mapsto \mathbb{R} \\ f &\mapsto f(x) \end{aligned}$$

is continuous

We saw the equivalence between these two definitions

RKHS norm controls smoothness

RKHS norm and smoothness

$$\begin{aligned}|f(x) - f(x')| &= |\langle f, k(x, \cdot) \rangle - \langle f, k(x', \cdot) \rangle| \quad (\text{reproducing property applied to } f) \\ &= |\langle f, k(x, \cdot) - k(x', \cdot) \rangle| \quad (\text{linearity of dot product}) \\ &\leq \|k(\cdot, x) - k(\cdot, x')\|_{\mathcal{H}} \|f\|_{\mathcal{H}} \quad (\text{by Cauchy-Schwarz inequality})\end{aligned}$$

- $\|f\|_{\mathcal{H}}$ controls how much the values at two points x and x' differ compared to their distance
- Larger value of $\|f\|_{\mathcal{H}}$ allows higher variations (potentially non-smooth functions)

Smaller RKHS norm \implies Smooth functions

- The same happens in finite dimensions when we add regularization $\|w\|^2$ for linear regression and SVM

Notion of Generalization

It is desired that the error of our classifier is close to that of Bayes classifier. However, another desirable quality in machine learning algorithms is

Generalization

- Let f_n be a classifier obtained by some algorithm (such as deep net or SVM or Random forest) which is based on a finite training sample of size n .
- The classifier f_n generalizes well if the difference between empirical and expected of f_n is low, i.e.,

$$|R(f_n) - R_{emp}(f_n)| \approx 0$$

- Note that having low generalization gap does imply low expected or test error, it just means that **empirical error is a good indicator of expected error**

Large vs Small Function class

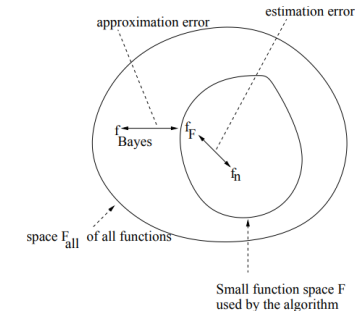


Figure: Pictorial depiction of the components of classification error

- The space F_{all} contains all possible functions that may be implemented using SVM, Deep nets, Random Forest and everything else
- **Estimation error** - $(R(f_n) - R(f_F))$ - **finiteness of training data**
- **Approximation error** - $(R(f_F) - R(f_{Bayes}))$ - **choice of function class**
- For example - If someone is claiming that using a deep net on a certain ML problem works better than SVM, which of the two errors is actually going down?

Large vs Small Function class

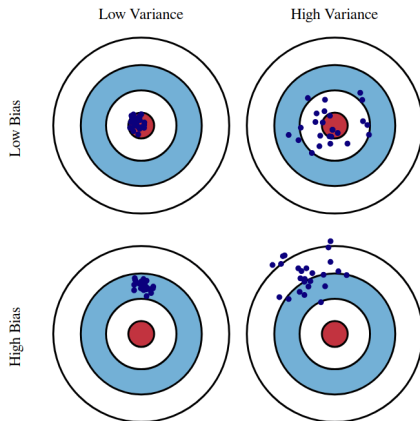


Figure: Pictorial depiction of the components of classification error

- **Estimation error** - $(R(f_n) - R(f_{\mathcal{F}}))$ - corresponds to **Variance**
- **Approximation error** - $(R(f_{\mathcal{F}}) - R(f_{Bayes}))$ - corresponds to **Bias**

Error variation with Function class capacity

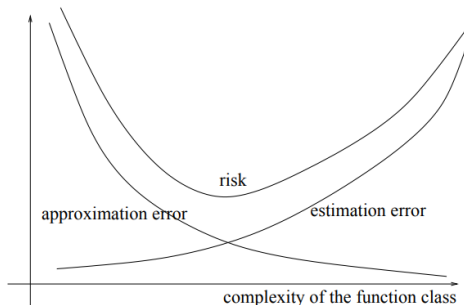


Figure: Variation of error components with the complexity of function class (tutorial by Von Luxburg and Schoelkopf)

- To the left with low complexity function class -
 - Linear classifiers or kernel classifier with high variance
- To the right with high complexity function class -
 - Deep neural networks

Empirical Risk Minimization

In practice, learning algorithms (do not have access to the underlying data generating distribution P over $\mathcal{X} \times \mathcal{Y}$) are based on minimizing error on the training data. Formally, this is given as follows :

Principle of ERM

The idea behind the principle of Empirical Risk Minimization is to find a classifier in a pre-defined function class which minimizes the empirical risk. That is

$$f_n := \arg \min_{f \in \mathcal{F}} R_{emp}(f)$$

- We want to check if the classifier (function) f_n that we learn from ERM is consistent or not

$$P(R(f_n) - R(f_{\mathcal{F}}) > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty$$

Uniform Convergence

- **Uniform Convergence** is a condition over a function class which ensures consistency of ERM, and is given by $|R_{emp}(f) - R(f)| < \epsilon, \forall f \in \mathcal{F}$ for some finite sample size n
- Alternatively, the condition of Uniform Convergence can be stated $\sup_{f \in \mathcal{F}} |R_{emp}(f) - R(f)| < \epsilon$

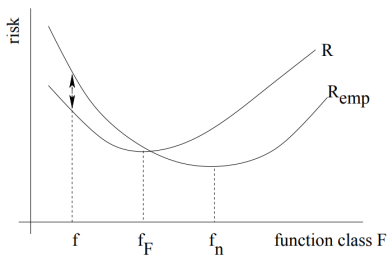


Figure: Under Uniform Convergence, the difference between the two curves becomes arbitrarily small for some large but finite sample size n

NASC for consistency of ERM

- Uniform convergence is a sufficient condition for the consistency of ERM
- Is it also necessary?

Theorem by Vapnik and Chervonenkis

Uniform convergence, i.e.,

$$\mathbb{P}(\sup_{f \in \mathcal{F}} |R(f) - R_{emp}(f)| > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty$$

$\forall \epsilon > 0$ is a necessary and sufficient condition for consistency of ERM with respect to the function class \mathcal{F} .

Capacity of Function Class

The main quantity of interest from the previous theorem is the following :

$$\mathbb{P}(\sup_{f \in \mathcal{F}} |R(f) - R_{emp}(f)| > \epsilon)$$

- Can we study the above quantity in the non-asymptotic regime, i.e. when the sample size n is finite
 - Practically, this also matters more since we normally have finite data size
- In bounding the quantity $\mathbb{P}(\sup_{f \in \mathcal{F}} |R(f) - R_{emp}(f)| \geq \epsilon)$, there are two challenges :
 - Infinitely many functions, due to **continuous nature of the function class**
 - The expected risk $R(f)$, which depends on the underlying probability distribution, and **cannot be computed from training data**
- To get a handle on this, we need the following three concepts :
 - Union bound - $\mathbb{P}(\sup_{f \in \mathcal{F}} |R(f) - R_{emp}(f)| \geq \epsilon) \leq 2m \exp(-2n\epsilon^2)$
 - Symmetrization -
 $\mathbb{P}(\sup_{f \in \mathcal{F}} |R(f) - R_{emp}(f)| > \epsilon) \leq 2\mathbb{P}(\sup_{f \in \mathcal{F}} |R_{emp}(f) - R'_{emp}(f)| > \epsilon/2)$
 - Shattering - $\mathbb{P}(\sup_{f \in \mathcal{F}} |R(f) - R_{emp}(f)| \geq \epsilon) \leq 2\mathcal{N}(\mathcal{F}, 2n) \exp(-n\epsilon^2/4)$

Representer Theorem

- For the following optimization

$$f_{\mathcal{H}} := \arg \min_f \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda \theta(\|f\|_{\mathcal{H}})$$

where $\theta : [0, \infty) \mapsto \mathbb{R}$ is non-decreasing function

- Even though the above problem is potentially an infinite dimensional optimization problem, **Representer Theorem** states its solution can be expressed in the following form

$$f_{\mathcal{H}} = \sum_{i=1}^n \alpha_i k(\cdot, x_i)$$

where $\alpha_i \in \mathbb{R}$

- Infinite to finite dimensional problem

Solving Kernel Ridge Regression

- Let's denote by
 - $y \in \mathbb{R}^n$, the label vector denoting the true values for the inputs
 - The kernel matrix K , where $K_{ij} = K(x_i, x_j)$
 - $\alpha \in \mathbb{R}^n$, the co-efficients we want to find
- For the input instance, the prediction by the desired function can be written as follows :

$$(\hat{f}(x_1), \dots, \hat{f}(x_n))^T = K\alpha$$

- We also know that

$$\|f\|_{\mathcal{H}}^2 = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j) = \alpha^T K \alpha$$

- Solving Kernel Ridge Regression involves solving

$$\arg \min_{\alpha \in \mathbb{R}^n} \frac{1}{n} (K\alpha - y)^T (K\alpha - y) + \lambda \alpha^T K \alpha$$

Solving Kernel Logistic Regression

- By representer theorem, any solution to kernel logistic regression is given by

$$\hat{f}(x) = \sum_{i=1}^n \alpha_i k(x_i, x)$$

- Also, we have the following :

- For the input instance, the prediction by the desired function can be written as follows :

$$(\hat{f}(x_1), \dots, \hat{f}(x_n))^T = K\alpha$$

- We also know that

$$\|f\|_{\mathcal{H}}^2 = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j) = \alpha^T K \alpha$$

- Therefore, we need to solve the following :

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i [K\alpha]_i)) + \frac{\lambda}{2} \alpha^T K \alpha$$

Convex sets

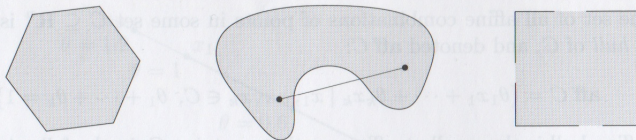
- A **line segment** between $x_1 \in \mathbb{R}^n$ and $x_2 \in \mathbb{R}^n$ is defined as all points that satisfy

$$x = \theta x_1 + (1 - \theta)x_2, 0 \leq \theta \leq 1$$

- A **convex set** contains the line segment between any two distinct points in the set

$$x_1, x_2 \in C, 0 \leq \theta \leq 1 \Rightarrow \theta x_1 + (1 - \theta)x_2 \in C$$

- Below: Convex and non-convex sets. Q: Which ones are convex?



Convex functions

- A function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is convex if (i) the domain of f is a convex set and (ii) for all x, y , and $0 \leq \theta \leq 1$, we have

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y).$$

- Geometrical interpretation: the graph of the function lies below the line segment from $(x, f(x))$ to $(y, f(y))$
- A function f is
 - strictly convex if strict inequality holds above
 - concave if $-f$ is convex.



Duality: Lagrangian

- Consider the primal optimisation problem

$$\begin{aligned} \min_{x \in \mathcal{D}} \quad & f_0(x) \\ \text{s.t.} \quad & f_i(x) \leq 0, i = 1, \dots, m \\ & h_i(x) = 0, i = 1, \dots, p \end{aligned}$$

with variable $x \in \mathbb{R}^n$

- Augment the objective function with the weighted sum of the constraint functions to form the **Lagrangian** of the optimization problem:

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

- $\lambda_i, i = 1, \dots, m$ and $\nu_i, i = 1, \dots, p$ (ν is the greek letter 'nu') are called the **Lagrange multipliers** or **dual variables**

Lagrange dual function

- The **Lagrange dual function** $g : \mathbb{R}^m \times \mathbb{R}^p \mapsto \mathbb{R}$ is the minimum value of the Lagrangian over x :

$$g(\lambda, \nu) = \inf_x L(x, \lambda, \nu) = \inf_x \{f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)\}$$

- Intuitively:
 - Fixing coefficients (λ, ν) corresponds to certain level of penalty,
 - The infimum returns the optimal x for that level of penalty
 - $g(\lambda, \nu)$ is the corresponding value for the Lagrangian
 - $g(\lambda, \nu)$ is a concave function as a pointwise infimum of a family of affine functions of (λ, ν)

The Lagrange dual problem

- For each pair (λ, ν) , $\lambda \geq 0$, the Lagrange dual function gives a lower bound on the optimal value of p^* .
- What is the tightest lower bound that can be achieved? We need to find the maximum
- This gives us a optimization problem

$$\begin{aligned} \max_{\lambda, \nu} \quad & g(\lambda, \nu) \\ \text{s.t.} \quad & \lambda \geq 0 \end{aligned}$$

- It is called the **Lagrange dual problem** of the original optimization problem.
- It is a convex optimisation problem, since it is equivalent to minimising $-g(\lambda, \nu)$ which is a convex function

SVM Problem Formulation

- Using Representer theorem, the problem can be reformulated as

$$\min_{\alpha \in \mathbb{R}^n} \left\{ \frac{1}{n} \sum_{i=1}^n \ell_{\text{hinge}}(y_i [K\alpha]_i) + \lambda \alpha^T K \alpha \right\}$$

- The above optimization problem is convex
- However, it is non-smooth optimization problem

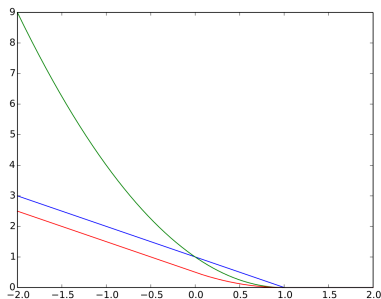


Figure: $z = yf(x)$ in the above graph

Rewriting in terms of Primal Variables

SVM Primal Formulation

$$\min_{\alpha \in \mathbb{R}^n, \xi \in \mathbb{R}^n} \left\{ \frac{1}{n} \sum_{i=1}^n \xi_i + \lambda \alpha^T K \alpha \right\}$$

such that

$$\begin{cases} 1 - y_i [K\alpha]_i - \xi_i \leq 0 & \text{for } i = 1, \dots, n \\ -\xi_i \leq 0 & \text{for } i = 1, \dots, n \end{cases}$$

SVM Dual Formulation

$$\max_{\alpha \in \mathbb{R}^n} 2 \sum_{i=1}^n \alpha_i y_i - \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x_i, x_j)$$

such that

$$0 \leq y_i \alpha_i \leq \frac{1}{2\lambda n} \text{ for } i = 1, \dots, n$$

Pictorial Depiction for α values

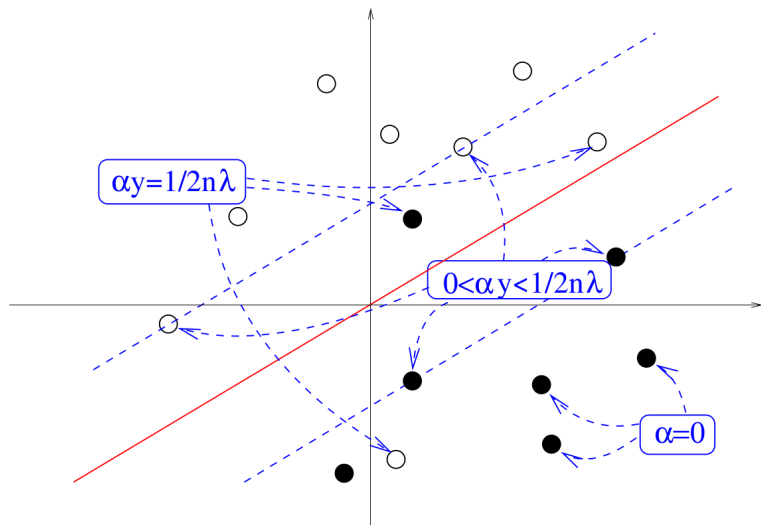


Figure: Training points with different values of α , (Picture : Julien Mairal)

Summary

- Broadly topics covered in this course
 - Basics about Kernels and RKHS
 - Foundational learning theory
 - Algorithms - Supervised and Unsupervised
 - Convex optimization overview
 - Hands-on with various algorithms

Summary

- Broadly topics covered in this course
 - Basics about Kernels and RKHS
 - Foundational learning theory
 - Algorithms - Supervised and Unsupervised
 - Convex optimization overview
 - Hands-on with various algorithms
- Some of the topics not covered
 - Structured prediction and learning with graphs
 - Other variants of SVM such as Multi-class SVM, RankSVM
- However, the course provides fundamentals and mathematical foundations to further understand these topics as desired

Summary

- Broadly topics covered in this course
 - Basics about Kernels and RKHS
 - Foundational learning theory
 - Algorithms - Supervised and Unsupervised
 - Convex optimization overview
 - Hands-on with various algorithms
- Some of the topics not covered
 - Structured prediction and learning with graphs
 - Other variants of SVM such as Multi-class SVM, RankSVM
- However, the course provides fundamentals and mathematical foundations to further understand these topics as desired
- Exam on 12.4 and 29.5, please register for either one of the dates

Thank you for following the course
and
Good Luck for the exam!