

Kernel Methods & Regularization

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

NetDB-ML, Spring 2014

1 What Learning Theory Taught Us?

- Generalization Performance
- Kernel Methods & Regularization

2 Extended Linear Models

- Regularized Nonlinear Regression
- Regularized Least Square Classification

3 Kernels and RKHS

1 What Learning Theory Taught Us?

- Generalization Performance
- Kernel Methods & Regularization

2 Extended Linear Models

- Regularized Nonlinear Regression
- Regularized Least Square Classification

3 Kernels and RKHS

1 What Learning Theory Taught Us?

- Generalization Performance
- Kernel Methods & Regularization

2 Extended Linear Models

- Regularized Nonlinear Regression
- Regularized Least Square Classification

3 Kernels and RKHS

Generalization Error

- Assuming a hypothesis class \mathcal{H} , let $h \in \mathcal{H}$ be the hypothesis trained from the dataset $\mathcal{X} = \{(\mathbf{x}^{(t)}, r^{(t)})\}_{t=1}^N$ by minimizing the empirical error:
$$R_{emp}[h] := \frac{1}{N} \sum_{t=1}^N l(h(\mathbf{x}^{(t)}), r^{(t)})$$

- Generalization error of h :

$$R[h] := \int p(\mathbf{x}, r) l(h(\mathbf{x}), r) d(\mathbf{x}, r) = E_{\mathcal{J} \times \mathcal{L}} [l(h(\mathbf{x}), r)]$$

- Let $h^* := \arg \inf_{g \in \mathcal{H}} R[g]$ and $R^* := \inf_{f: \mathcal{J} \rightarrow \mathcal{L}} R[f]$

- Instead, our ultimate goal is to have $R[h] \rightarrow R^*$!

- $R[h] - R^* = (R[h^*] - R^*) + (R[h] - R[h^*])$ Model太簡單，無法approximate

- $R[h^*] - R^*$ is called the **approximation error**: we need a complex model to reduce this

- $R[h] - R[h^*]$ is called the **estimation error**: we need a simple model

Model太複雜，會Overfit

Which Models to Assume?

- We can assume different models $\mathcal{H}_1, \mathcal{H}_2, \dots$ and compares them in the model selection process
- But there are too many different models
- This lecture introduces methods that allow the complexity of \mathcal{H} to be tuned *after* it is assumed
 - Simplifies the task of model assumption
 - Linear models suffice in most cases
- Does *not* mean \mathcal{H} will have the right complexity
 - Model selection is still required

1 What Learning Theory Taught Us?

- Generalization Performance
- Kernel Methods & Regularization

2 Extended Linear Models

- Regularized Nonlinear Regression
- Regularized Least Square Classification

3 Kernels and RKHS

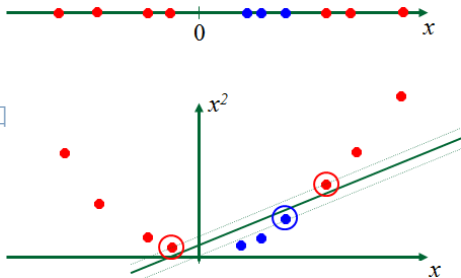
Kernel Methods

- To reduce approximation error, we need a complex model
- But complex model make the objective hard to solve
- One common approach, called **kernel methods**, is to map/lift/kernelize examples $\mathcal{X} = \{\mathbf{x}^{(t)}, r^{(t)}\}_{t=1}^N$ from original **input space** to a higher dimensional **feature space** $\Phi(\mathcal{X}) = \{\Phi(\mathbf{x}^{(t)}), r^{(t)}\}_{t=1}^N$
 - $k(\mathbf{a}, \mathbf{b}) := \langle \Phi(\mathbf{a}), \Phi(\mathbf{b}) \rangle$ is called the **kernel function**
不是讓hypothesis class更複雜
而是把data變複雜
- E.g., suppose $x \in \mathbb{R}$,
 $\Phi(x) = [x, x^2]^\top \in \mathbb{R}^2$
- Why does it work?
如果objective有做內積的可能
那麼feature space也要可以做內積

Kernel Methods

- To reduce approximation error, we need a complex model
- But complex model make the objective hard to solve
- One common approach, called **kernel methods**, is to map/lift/kernelize examples $\mathcal{X} = \{\mathbf{x}^{(t)}, r^{(t)}\}_{t=1}^N$ from original **input space** to a higher dimensional **feature space** $\Phi(\mathcal{X}) = \{\Phi(\mathbf{x}^{(t)}), r^{(t)}\}_{t=1}^N$
 - $k(\mathbf{a}, \mathbf{b}) := \langle \Phi(\mathbf{a}), \Phi(\mathbf{b}) \rangle$ is called the **kernel function**

- E.g., suppose $x \in \mathbb{R}$,
 $\Phi(x) = [x, x^2]^\top \in \mathbb{R}^2$
- Why does it work? 看右邊就知
- Note a linear h in the feature space is **not linear** in the input space anymore



Regularization

- To reduce estimation error, we need a simple model
- Note that in the same \mathcal{H} , there are still hypotheses that are more complex than the others
- We can add a term, called **regularization term**, in our objective that penalizes complex hypotheses in \mathcal{H} :

$$\arg \min_{g \in \mathcal{H}} \sum_{t=1}^N \overbrace{l(g(\mathbf{x}^{(t)}), r^{(t)})}^{\text{data term}} + \underbrace{\lambda \Omega(g)}_{\text{smoothness term}},$$

如果g太複雜 就提高「人造的」誤差

where

- Ω is a smoothness measure
- λ is a **hyperparameter** (fixed during the training process), which controls the trade-off between a) minimizing the empirical error; and b) maximizing function smoothness
- Why does it work? A “smoother” h allows unseen instances to learn values from those of nearby examples

1 What Learning Theory Taught Us?

- Generalization Performance
- Kernel Methods & Regularization

2 Extended Linear Models

- Regularized Nonlinear Regression
- Regularized Least Square Classification

3 Kernels and RKHS

1 What Learning Theory Taught Us?

- Generalization Performance
- Kernel Methods & Regularization

2 Extended Linear Models

- Regularized Nonlinear Regression
- Regularized Least Square Classification

3 Kernels and RKHS

Linear Regression

- Let $X = \begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ x_1^{(2)} & \cdots & x_d^{(2)} \\ \vdots & \ddots & \vdots \\ x_1^{(N)} & \cdots & x_d^{(N)} \end{bmatrix}$, $\mathbf{w} = [w_1, \dots, w_d]^\top$, and $\mathbf{r} = [r^{(1)}, r^{(2)}, \dots, r^{(N)}]^\top$

- The linear regression problem:

$$\arg \min_{\mathbf{w}, b} \left\| \mathbf{r} - \begin{bmatrix} \mathbf{1} & X \end{bmatrix} \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} \right\|^2$$

- b is called the bias term

- Solution: $\begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}^* = (\begin{bmatrix} \mathbf{1} & X \end{bmatrix}^\top \begin{bmatrix} \mathbf{1} & X \end{bmatrix})^{-1} \begin{bmatrix} \mathbf{1} & X \end{bmatrix}^\top \mathbf{r}$ if $\begin{bmatrix} \mathbf{1} & X \end{bmatrix}$ has full column rank

Closed form solution $(X^\top X)^{-1} X^\top \mathbf{r}$

- Objective:

$$\arg \min_{\mathbf{w}, b} \left\| \mathbf{r} - \begin{bmatrix} \mathbf{1} & \mathbf{X} \end{bmatrix} \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} \right\|^2,$$

$$\text{where } \mathbf{X} = \begin{bmatrix} \Phi(\mathbf{x}^{(1)})^\top \\ \vdots \\ \Phi(\mathbf{x}^{(N)})^\top \end{bmatrix}$$

- Note that the number of variables to solve now becomes (dimension of feature space + 1)

Regularization

- The *regularized least square problem*:

$$\arg \min_{\mathbf{w}, b} \|\mathbf{r} - [\mathbf{1} \ X] \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}\|^2 + \lambda \|\mathbf{w}\|^2,$$

- The bias term b is not regularized
- Why minimizing $\|\mathbf{w}\|^2$? A flat h learns from all examples (by the average of their label values)
- Can be expressed as an ordinary least squares:

$$\arg \min_{\mathbf{w}, b} \|\tilde{\mathbf{r}} - \tilde{\mathbf{X}} \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}\|^2, \text{ where}$$

$$\tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{1} & X \\ 0 & \mathbf{0} \\ \mathbf{0} & \sqrt{\lambda} \mathbf{I}_d \end{bmatrix} \in \mathbb{R}^{(N+d+1) \times (d+1)} \text{ and } \tilde{\mathbf{r}} = [\mathbf{r}, 0]^\top \in \mathbb{R}^{N+d+1}$$

[Proof]

- Solution: $\begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}^* = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{r}} =$

$$([\mathbf{1} \ X]^\top [\mathbf{1} \ X] + \lambda \begin{bmatrix} 0 & \mathbf{0} \\ \mathbf{0} & \sqrt{\lambda} \mathbf{I}_d \end{bmatrix})^{-1} [\mathbf{1} \ X]^\top \mathbf{r}$$

- $\tilde{\mathbf{X}}$ must be full column rank

The Bias Term

- With some particular kernel functions, we can simply set $b = 0$
- Simplified objective:

$$\arg \min_{\mathbf{w}} \|\mathbf{r} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

- Solution: $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{r} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{r}$
- In a very high (or infinite) dimensional feature space, this usually makes little difference in performance

在越高維度下，越不需要bias term也可以解釋資料，無限維下更是如此，因為可以透過各維度的線性組合，產生與bias term等效的結果，因此便少了一項

1 What Learning Theory Taught Us?

- Generalization Performance
- Kernel Methods & Regularization

2 Extended Linear Models

- Regularized Nonlinear Regression
- Regularized Least Square Classification

3 Kernels and RKHS

Regularized Least Square Classification

- $\mathbf{r} \in \{1, -1\}^N$
- Recall that in linear case, we cannot directly apply linear regression to the classification problem
 - Why? The linear model is obviously “too simple” such that the SSE will be large for all lines (large bias)
 - We therefore assumed the “logistic” model
- Can we directly apply regularized nonlinear regression to the classification problem?
 - Yes, as given a sufficiently high dimensional feature space, the nonlinear model will always be complex enough to produce low SSE
- This is called *regularized least square classification*

Questions?

- To maximize the flexibility of model complexity, we want a Φ that maps \mathbf{x} 's to a feature space of dimension as high as possible
 - Ideally, to an infinite dimensional feature space
 - Q1: How to obtain such an Φ ?
- Meanwhile, the number of variables to solve (in \mathbf{w}) increases as the dimension of feature space becomes higher
 - Q2: How to solve \mathbf{w} in an infinite dimensional feature space?

Outline

1 What Learning Theory Taught Us?

- Generalization Performance
- Kernel Methods & Regularization

2 Extended Linear Models

- Regularized Nonlinear Regression
- Regularized Least Square Classification

3 Kernels and RKHS

Common Kernel Functions

$$\Phi(\mathbf{a}) = 1, a, a^2 \quad \Phi(\mathbf{b}) = 1, b, b^2$$

$\langle \Phi(\mathbf{a}), \Phi(\mathbf{b}) \rangle = (\text{係數未標準化})$
 $1, a, a^2, b, ab, b^2 \dots$

- Kernel function: $k(\mathbf{a}, \mathbf{b}) := \langle \Phi(\mathbf{a}), \Phi(\mathbf{b}) \rangle$
- Polynomial kernel: $k(\mathbf{a}, \mathbf{b}) = (\mathbf{a}^\top \mathbf{b} / \alpha + \beta)^\gamma$

常用的
無限維
Kernel

- E.g., let $\alpha = 1$, $\beta = 1$, $\gamma = 2$ and $\mathbf{a} \in \mathbb{R}^2$, then
 $\Phi(\mathbf{a}) = [1, \sqrt{2}a_1, \sqrt{2}a_2, a_1^2, a_2^2, \sqrt{2}a_1a_2]^\top \in \mathbb{R}^6$

- Gaussian Radial Basis Function (RBF)¹ kernel: $k(\mathbf{a}, \mathbf{b}) = \exp(-\frac{\|\mathbf{a}-\mathbf{b}\|_2^2}{2\sigma^2})$
or $\exp(-\gamma\|\mathbf{a}-\mathbf{b}\|_2^2)$, $\gamma \geq 0$

- $k(\mathbf{a}, \mathbf{b}) = \exp(-\gamma\|\mathbf{a}-\mathbf{b}\|_2^2) = \exp(-\gamma\|\mathbf{a}\|^2 + 2\gamma\mathbf{a}^\top \mathbf{b} - \gamma\|\mathbf{b}\|^2) =$
 $\exp(-\gamma\|\mathbf{a}\|^2 - \gamma\|\mathbf{b}\|^2)(1 + \frac{2\gamma\mathbf{a}^\top \mathbf{b}}{1!} + \frac{(2\gamma\mathbf{a}^\top \mathbf{b})^2}{2!} + \dots)$

- Let $\mathbf{a} \in \mathbb{R}^2$, then $\Phi(\mathbf{a}) =$

$$\exp(-\gamma\|\mathbf{a}\|^2)[1, \sqrt{\frac{2\gamma}{1!}}a_1, \sqrt{\frac{2\gamma}{1!}}a_2, \sqrt{\frac{2\gamma}{2!}}a_1^2, \sqrt{\frac{2\gamma}{2!}}a_2^2, 2\sqrt{\frac{\gamma}{2!}}a_1a_2, \dots]^\top \in \mathbb{R}^\infty$$

- α , β , γ , and σ are hyperparameters in the objective

¹A radial basis function of \mathbf{a} and \mathbf{b} is a real-valued function whose values depend only on $\|\mathbf{a}-\mathbf{b}\|$

Questions Revisited

- Q1: How to obtain a feature mapping Φ whose range is infinite dimensional?
- Q2: How to solve w in an infinite dimensional feature space?
- Our previous definition of Φ over Gaussian RBF kernel answers Q1, but not Q2
- To answer Q2, we need a new perspective on Φ

Why Another Perspective?

- Recall that in regularized linear regression ($\arg \min_{\mathbf{w}} \|\mathbf{r} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$ with the bias term b omitted), we have $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{r}$
- Indeed, it can be shown that $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{r} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{r}$ [Homework]
- Letting $\mathbf{c} = (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{r}$ we see that $\mathbf{w}^* = \sum_{t=1}^N c_t \mathbf{x}^{(t)}$ is a linear combination of the examples 使empirical error最小的 \mathbf{w}^* 是來自於data的線性組合 \mathbf{c}_t 是 $d \times 1$
- Given any lifting Φ , if \mathbf{w} always admit the form $\mathbf{w} = \sum_{t=1}^N c_t \Phi(\mathbf{x}^{(t)})$ for some \mathbf{c} , we can instead solve:

看似要解無限維的 \mathbf{w}

但是 \mathbf{w} 必然是 \mathbf{c} 的線性組合

所以只要解 N 維的 \mathbf{c}

$$\arg \min_{\mathbf{c}} \|\mathbf{r} - \mathbf{K}\mathbf{c}\|^2 + \lambda \mathbf{c}^\top \mathbf{K}\mathbf{c}$$

$$\bullet \quad \|\mathbf{w}\|^2 = \mathbf{w}^\top \mathbf{w} = \sum_{i=1}^N \sum_{j=1}^N c^{(i)} \mathbf{x}^{(i)\top} \mathbf{x}^{(j)} c^{(j)} = \mathbf{c}^\top \mathbf{K}\mathbf{c}$$

- The number of variables to solve (in $\mathbf{c} \in \mathbb{R}^N$) now becomes independent with the dimension of feature space!

Definition

Given a function $k: \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}$, let $\mathbf{K} \in \mathbb{R}^{N \times N}$ be the kernel matrix where $K_{i,j} := k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \langle \Phi(\mathbf{x}^{(i)}), \Phi(\mathbf{x}^{(j)}) \rangle$. Then k is called a **kernel function** iff \mathbf{K} is positive semidefinite.

正半定才有Convex optimization可以用

- E.g., $k(\mathbf{a}, \mathbf{b}) = \mathbf{a}^\top \mathbf{b}$,
 - As $\mathbf{K} = \mathbf{X}\mathbf{X}^\top$ and for any \mathbf{v} , we have $\mathbf{v}^\top \mathbf{K} \mathbf{v} = \mathbf{v}^\top \mathbf{X}\mathbf{X}^\top \mathbf{v} = \|\mathbf{X}^\top \mathbf{v}\|^2 \geq 0$
- E.g., polynomial and Gaussian RBF [Proof]

Reproducing Kernel Hilbert Space (RKHS)

function space : 代一個數字，得到一個function

- We can define a lifting as

$\Phi(\mathbf{x}) = k(\mathbf{x}, \cdot)$, where k is a kernel function
lifting function是 \mathbf{x} 的function,
定義為 k 代入 \mathbf{x} 所得到的function

- And define a complete Hilbert space as the collection of all $\sum_{i=1}^n \alpha^{(i)} k(\mathbf{x}^{(i)}, \cdot)$

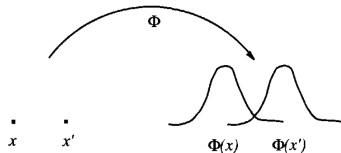
- n , $c^{(i)}$, and $\mathbf{x}^{(i)}$ are all arbitrary任意的
- Infinite dimensional

- With the inner product:² $\langle f, g \rangle := \sum_{i=1}^n \sum_{j=1}^m \alpha^{(i)} \beta^{(j)} k(\mathbf{x}^{(i)}, \mathbf{y}^{(j)})$ for any $f = \sum_{i=1}^n \alpha^{(i)} k(\mathbf{x}^{(i)}, \cdot)$ and $g = \sum_{j=1}^m \beta^{(j)} k(\mathbf{y}^{(j)}, \cdot)$

- Well-defined [Homework] $k(\mathbf{y})$ 其實是基底

- **Reproducing properties:** $\langle f, k(\mathbf{y}, \cdot) \rangle = f(\mathbf{y})$ and $\langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle = \langle k(\mathbf{x}, \cdot), k(\mathbf{y}, \cdot) \rangle = k(\mathbf{x}, \mathbf{y})$ [Proof]

²A vector space is called a Hilbert space if it is endowed with an inner product, and is complete if every Cauchy sequence (a sequence of points with decreasing distances) in it converges.



- $\Phi(\mathbf{x})$ is a **function**

內積有良好定義的空間，
可以稱為Hilbert space

Representer Theorem

Theorem

Let $L: \mathcal{X} \times \mathbb{R}^N \rightarrow \mathbb{R} \cup \{\infty\}$ be an arbitrary loss function over $\mathcal{X} = \{\mathbf{x}^{(t)}, r^{(t)}\}_{t=1}^N$, $C_k: \mathcal{X} \times \mathbb{R}^N \rightarrow \mathbb{R} \cup \{\infty\}$, $0 \leq i \leq K$, constraint functions, and $\Omega: [0, \infty) \rightarrow \mathbb{R}$ a strictly increasing function. Then each minimizer $h \in \mathcal{H}$ of the regularized risk functional:

$$\arg \min_{g \in \mathcal{H}} L((\mathbf{x}^{(1)}, r^{(1)}, g(\mathbf{x}^{(1)})), \dots, (\mathbf{x}^{(N)}, r^{(N)}, g(\mathbf{x}^{(N)}))) + \Omega(\|g\|_{\mathcal{RKHS}})$$

subject to $C_k((\mathbf{x}^{(1)}, r^{(1)}, g(\mathbf{x}^{(1)})), \dots, (\mathbf{x}^{(N)}, r^{(N)}, g(\mathbf{x}^{(N)}))) \leq 0, \forall k$

任何這種form

admits the form $h(\mathbf{x}) = \sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x})$. 都可以換成這種form

- L is more general than l seen previously as the former allows coupling between samples
- $\Omega(\|g\|_{\mathcal{RKHS}})$ can be written as $\tilde{\Omega}(\|g\|_{\mathcal{RKHS}}^2)$ without loss of generality
 - As the quadratic function is strictly increasing on $[0, \infty)$, hence Ω is strictly increasing iff $\tilde{\Omega}$ is so
 - In particular, $\tilde{\Omega}(\|g\|_{\mathcal{RKHS}}^2)$ can be $\lambda \|g\|_{\mathcal{RKHS}}^2$ for some $\lambda > 0$

Proof

- We consider $\tilde{\Omega}(\|g\|_{\mathcal{RKHS}}^2)$ for convenience
- By the fundamental theorem of linear algebra, we can decompose any $g \in \mathcal{H}$ into two vectors parallel and orthogonal to span of $k(\mathbf{x}^{(1)}, \cdot), \dots, k(\mathbf{x}^{(N)}, \cdot)$ respectively; i.e.,
 $g(\mathbf{x}) = \sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x}) + g_{\perp}(\mathbf{x})$ 任何向量都可以分別投影到兩個 orthogonal space, 表示成它們的和
反證: 假設不能寫成那個 form, 即 $g_{\perp} \neq 0$
- Since $\langle g_{\perp}, k(\mathbf{x}^{(i)}, \cdot) \rangle = 0$ for all $1 \leq i \leq N$, we have g_{\perp} 和任何 $k(\mathbf{x}^{(i)}, \cdot)$ 都垂直, 因此內積為 0
 $g(\mathbf{x}^{(i)}) = \langle g, k(\mathbf{x}^{(i)}, \cdot) \rangle = \sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x}^{(i)}) + \langle g_{\perp}, k(\mathbf{x}^{(i)}, \cdot) \rangle = \sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x}^{(i)})$
- Now suppose the minimizer h has the form $h - h_{\perp} = \sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x})$
 $h(\mathbf{x}) = \sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x}) + h_{\perp}(\mathbf{x})$, next we show that $h - h_{\perp}$ is always a better solution, which contradicts our assumption
一個更小的 h 也滿足, 所以是更好的 h , 與假設矛盾
- First, $h - h_{\perp}$ satisfies all constraints C_k 's, as $(h - h_{\perp})(\mathbf{x}^{(i)}) = h(\mathbf{x}^{(i)})$ for all $1 \leq i \leq N$
垂直的部分, 內積 = 0
- Due to the same reason, $h - h_{\perp}$ has the same loss score from L as h
- Furthermore, $\tilde{\Omega}(\|h - h_{\perp}\|^2) = \tilde{\Omega}(\|\sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x})\|_{\mathcal{RKHS}}^2) \leq \tilde{\Omega}(\|\sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x})\|_{\mathcal{RKHS}}^2 + \|h_{\perp}\|_{\mathcal{RKHS}}^2) = \tilde{\Omega}(\|h\|_{\mathcal{RKHS}}^2)$

- The minimizers of the problems with the form

$$\begin{aligned} \arg \min_{g \in \mathcal{H}} L((\mathbf{x}^{(1)}, r^{(1)}, g(\mathbf{x}^{(1)})), \dots, (\mathbf{x}^{(N)}, r^{(N)}, g(\mathbf{x}^{(N)}))) + \Omega(\|g\|_{\mathcal{RKHS}}) \\ \text{subject to } C_k((\mathbf{x}^{(1)}, r^{(1)}, g(\mathbf{x}^{(1)})), \dots, (\mathbf{x}^{(N)}, r^{(N)}, g(\mathbf{x}^{(N)}))) \leq 0, \forall k \end{aligned}$$

are called *kernel machines*

- The representer theorem tells us that although the RKHS is infinite dimensional, the solution will always be in a subspace spanned by $k(\mathbf{x}^{(1)}, \cdot), \dots, k(\mathbf{x}^{(N)}, \cdot)$
 - We don't need to search for the entire RKHS to obtain a solution
- For *any kernel*, we only need to solve N variables, *independent* with the dimension of feature space

Example: Regularized Nonlinear Regression

- Let $\Phi(\mathbf{x}) = k(\mathbf{x}, \cdot)$, we can write the objective ($b = 0$) as

$$\arg \min_{g: g(\Phi(\mathbf{x})) = \mathbf{w}^\top \Phi(\mathbf{x})} \sum_{t=1}^N (r^{(t)} - g(\Phi(\mathbf{x}^{(t)})))^2 + \lambda \|g\|_{\mathcal{RKHS}}^2$$

- Then for any kernel, we can always solve an alternative objective:

$$\arg \min_{\mathbf{c}} \|\mathbf{r} - \mathbf{K}\mathbf{c}\|^2 + \lambda \mathbf{c}^\top \mathbf{K} \mathbf{c}$$

- $\|g\|_{\mathcal{RKHS}}^2 = \langle g, g \rangle = \sum_{i=1}^N \sum_{j=1}^N c^{(i)} c^{(j)} k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \sum_{i=1}^N \sum_{j=1}^N c^{(i)} \mathbf{x}^{(i)\top} \mathbf{x}^{(j)} c^{(j)} = \mathbf{c}^\top \mathbf{K} \mathbf{c}$
- $\mathbf{c} \in \mathbb{R}^N$, so the number of variables to solve is independent with the dimension of feature space
- In the linear case, we have $\Phi(\mathbf{x}) = k(\mathbf{x}, \cdot) = \mathbf{x}^\top [\cdot]$
 - The representer theorem coincides with our previous observation in linear regression that $\mathbf{w}^* = \sum_{t=1}^N c_t \mathbf{x}^{(t)}$
 - Same kernel function $k(\mathbf{a}, \mathbf{b}) = \mathbf{a}^\top \mathbf{b}$, different feature spaces ($\Phi(\mathbf{x}) = \mathbf{x}$ vs. $\Phi(\mathbf{x}) = \mathbf{x}^\top [\cdot]$)

Semiparametric Representer Theorem

Theorem

Following the previous theorem, let $\tilde{g} := g + b\psi$, where $g \in \mathcal{H}$, $b \in \mathbb{R}$, and $\psi : \mathcal{J} \rightarrow \mathbb{R}$. Then each minimizer \tilde{h} of the regularized risk functional:

$$\begin{aligned} & \arg \min_{\tilde{g}} L((\mathbf{x}^{(1)}, r^{(1)}, \tilde{g}(\mathbf{x}^{(1)})), \dots, (\mathbf{x}^{(N)}, r^{(N)}, \tilde{g}(\mathbf{x}^{(N)}))) + \Omega(\|g\|_{\mathcal{R}\mathcal{K}\mathcal{H}\mathcal{S}}) \\ & \text{subject to } C_k((\mathbf{x}^{(1)}, r^{(1)}, \tilde{g}(\mathbf{x}^{(1)})), \dots, (\mathbf{x}^{(N)}, r^{(N)}, \tilde{g}(\mathbf{x}^{(N)}))) \leq 0, \forall k \end{aligned}$$

admits the form $\tilde{h}(\mathbf{x}) = \sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x}) + b\psi(\mathbf{x})$. [Proof]

- When $\psi(\mathbf{x}) = 1$, we have $\tilde{h}(\mathbf{x}) = \sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x}) + b$
- Applicable to the kernel machines with a unregularized bias term (i.e., $b \neq 0$)
 - What is the objective of regularized nonlinear regression with $b \neq 0$ after applying this theorem? [Homework]