1. Shown that in regularized linear regression, $\boxed{\boldsymbol{w}^* \underset{①}{=} (\boldsymbol{X}^\top \boldsymbol{X} + \lambda \boldsymbol{I}_d)^{-1} \boldsymbol{X}^\top \boldsymbol{r}} \underset{②}{=} \boldsymbol{X}^\top (\boldsymbol{X}\boldsymbol{X}^\top + \lambda \boldsymbol{I}_N)^{-1} \boldsymbol{r}$, a linear combination of the examples.

- Regularized linear regression can be written as follow :

  arg min $\frac{1}{2}\|r - Xw\|^2 + \frac{1}{2}\lambda\|w\|^2$ .

- Applying first derivative (to w), we get $X^\top Xw - X^\top r + \lambda w$

- Minimum occurs when $X^\top Xw + \lambda I_d w = (X^\top X + \lambda I_d)w = X^\top r$ ,

  and w* $= (X^\top X + \lambda I_d)^{-1} X^\top r$

1. Shown that in regularized linear regression, $\boldsymbol{w}^* \underset{①}{=} (\boldsymbol{X}^\top \boldsymbol{X} + \lambda \boldsymbol{I}_d)^{-1} \boldsymbol{X}^\top \boldsymbol{r} \underset{②}{=} \boldsymbol{X}^\top (\boldsymbol{X}\boldsymbol{X}^\top + \lambda \boldsymbol{I}_N)^{-1}\boldsymbol{r}$, a linear combination of the examples.

- $(A+BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1}+DA^{-1}B)^{-1}DA^{-1}$    $A=\lambda I_d$ , $B=X^\top$ , $C=I_N$ , $D=X$

- $(X^\top X+\lambda I_d)^{-1}X^\top = (1/\lambda)I_d - (1/\lambda)I_d X^\top(I_N+X\cancel{(1/\lambda)I_d}X^\top)^{-1}X\cancel{(1/\lambda)I_d}X^\top$

  $= (1/\lambda)I_d X^\top - (1/\lambda)I_d X^\top(\lambda I_N+XX^\top)^{-1}XX^\top$

  $= (1/\lambda) I_d X^\top(I_N - \boxed{(\lambda I_N+XX^\top)^{-1}XX^\top \boxed{- (\lambda I_N+XX^\top)^{-1}(\lambda I_N)}} + \boxed{(\lambda I_N+XX^\top)^{-1}(\lambda I_N)})$

  $= (1/\lambda) I_d X^\top(I_N - \boxed{(\lambda I_N+XX^\top)^{-1}(\lambda I_N+XX^\top)} + (\lambda I_N+XX^\top)^{-1}(\lambda I_N))$

  $= (1/\lambda) I_d X^\top(I_N - I_N + (\lambda I_N+XX^\top)^{-1}(\lambda I_N))$

  $= X^\top(\lambda I_N+XX^\top)^{-1}$

- 因此， $(X^\top X+\lambda I_d)^{-1}X^\top r = X^\top(\lambda I_N+XX^\top)^{-1}r$

2. Show that in an RKHS, the inner product $\langle f, g \rangle := \sum_{i=1}^{n} \sum_{j=1}^{m} \alpha^{(i)} \beta^{(j)} k(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(j)})$ for any $f = \sum_{i=1}^{n} \alpha^{(i)} k(\boldsymbol{x}^{(i)}, \cdot)$ and $g = \sum_{j=1}^{m} \beta^{(j)} k(\boldsymbol{y}^{(j)}, \cdot)$ is well-defined; i.e., it satisfies

(a) *symmetry*: $\langle f, g \rangle = \langle g, f \rangle$;

- $\sum_{i=1}^{n} \sum_{j=1}^{m}$ 是Symmetric，$\alpha^{(i)} \beta^{(j)}$ 也是Symmetric，因此只要$k$是symmetric，那麼 $\sum_{i=1}^{n} \sum_{j=1}^{m} \alpha^{(i)} \beta^{(j)} k(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(j)})$ 就是Symmetric。事實如此，因為對於原始資料**x**，它 的Lifting function Φ 對於**x** 的每個維度的feature都是對稱的。

(b) *linearity*: $\langle af + bg, h \rangle = a\langle f, h \rangle + b\langle g, h \rangle$ for any $a, b \in \mathbb{R}$; and

- 關於純量積，a$\sum_{i=1}^{n} \alpha^{(i)} k(\boldsymbol{x}^{(i)}, \cdot)$ = $\sum_{i=1}^{n}$ a$\alpha^{(i)} k(\boldsymbol{x}^{(i)}, \cdot)$ 因此如此定義滿足乘法分配律。

- 而af+bg= $\sum_{i=1}^{n}$ aα$^{(i)}$ k (x$^{(i)}$,·) + $\sum_{j=1}^{m}$ bβ$^{(j)}$ k (y$^{(j)}$,·) ，定義h=$\sum_{q=1}^{p}$ γ$^{(q)}$ k (z$^{(q)}$,·)

- ⟨ af+bg,h ⟩ = $\sum_{q=1}^{p}$ ( $\sum_{i=1}^{n}$ aα$^{(i)}$ γ$^{(q)}$ k (x$^{(i)}$,z$^{(j)}$) + $\sum_{j=1}^{m}$ b β$^{(j)}$ γ$^{(q)}$ k (y$^{(j)}$, z$^{(j)}$) )

  = $\sum_{i=1}^{n}$ $\sum_{q=1}^{p}$ aα$^{(i)}$ γ$^{(q)}$ k (x$^{(i)}$,z$^{(j)}$) + $\sum_{j=1}^{m}$ $\sum_{q=1}^{p}$ b β$^{(j)}$ γ$^{(q)}$ k (y$^{(j)}$, z$^{(j)}$) = a ⟨ f,h ⟩ + b ⟨ g,h ⟩

2. Show that in an RKHS, the inner product $\langle f, g \rangle := \sum_{i=1}^{n} \sum_{j=1}^{m} \alpha^{(i)} \beta^{(j)} k(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(j)})$ for any $f = \sum_{i=1}^{n} \alpha^{(i)} k(\boldsymbol{x}^{(i)}, \cdot)$ and $g = \sum_{j=1}^{m} \beta^{(j)} k(\boldsymbol{y}^{(j)}, \cdot)$ is well-defined; i.e., it satisfies

(c) *positive definiteness*: $\langle f, f \rangle \geq 0$ with equality holds iff $f(\cdot) = 0(\cdot)$.　自己內積自己恆正

- $\langle f, f \rangle = \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha^{(i)} \alpha^{(j)} k(x^{(i)}, x^{(j)})$

　$= \sum_{i=1}^{n} \sum_{j=1}^{n} k(\alpha^{(i)} x^{(i)}, \alpha^{(j)} x^{(j)})$

　$= \sum_{i=1}^{n} \sum_{j=1}^{n} \langle \alpha^{(i)} \Phi(x^{(i)}), \alpha^{(j)} \Phi(x^{(j)}) \rangle$

　$= \langle \sum_{i=1}^{n} \alpha^{(i)} \Phi(x^{(i)}), \sum_{j=1}^{n} \alpha^{(j)} \Phi(x^{(j)}) \rangle$

　$= \langle \sum_{i=1}^{n} \alpha^{(i)} \Phi(x^{(i)}), \sum_{i=1}^{n} \alpha^{(i)} \Phi(x^{(i)}) \rangle$

　$= \| \sum_{i=1}^{n} \alpha^{(i)} \Phi(x^{(i)}) \|^2$

　$\geq 0$ ， $\alpha^{(i)} > 0$，因此等號成立在 $f(\cdot) = 0(\cdot)$

根據定義

純量積(乘法分配律)

根據定義

加法結合律

就是這樣

3. Show that in a large-margin linear classifier, the margin between the hyperplanes $\{\boldsymbol{x} : \boldsymbol{w}^\top \boldsymbol{x} - b = 1\}$ and $\{\boldsymbol{x} : \boldsymbol{w}^\top \boldsymbol{x} - b = -1\}$ is $2/\|\boldsymbol{w}\|$.

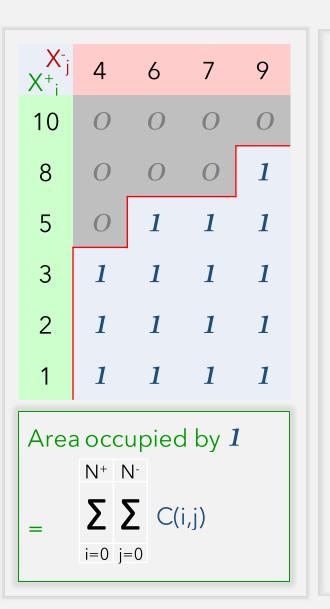- 推導兩平面間的距離公式的過程通常很冗長。因此假設我們已經有平行的觀念，亦即在N+1維空間中的2個N維超平面互相平行，即代表它們有相同的法向量，同時此兩個超平面的距離在任意處皆相等

- $w^\top x - b - 1 = 0$ 和 $w^\top x - b + 1 = 0$ 這兩個超平面互相平行，因為它們的法向量皆為w

- 對於超平面$w^\top x - 1 = b$上任一點$x_0$，$x_0$到超平面$w^\top x + 1 = b$的最短路徑，必然是沿著它們的法向量$w/\|w\|$走$\delta$的距離，因此可以記為$w^\top x_0 - 1 = b$，$w^\top(x_0 + \delta w/\|w\|) + 1 = b$

- 因此，$w^\top x_0 - 2 = w^\top(x_0 + \delta/\|w\|)$，即$w^\top(\delta/\|w\|)w = (\delta/\|w\|)w^\top w = \delta\|w\|^{(-1+2)} = \delta\|w\| = -2$

- 最後，$|\delta| = 2/\|w\|$，我們得到這兩個超平面距離為$2/\|w\|$

4. Prove the Semiparametric Representer theorem.

- 改寫Representer Theorem的證明

1. 將$\tilde{g}$中的g分解為<span style="color:blue">平行於</span>和<span style="color:green">垂直於</span>空間span(k(x$^{(1)}$,·),…, k(x$^{(N)}$,·))的部分，得到$\tilde{g}$ = g$_{//}$ + g$_{\perp}$ + b$\psi$ ，其中 g$_{//}$ = $\sum_{t=1}^{N} c_t k(\boldsymbol{x}^{(t)}, \boldsymbol{x})$

2. Loss function可以寫成 L( (**x**$^{(i)}$,**r**$^{(i)}$, g$_{//}$(**x**$^{(i)}$)+b$\psi$(**x**$^{(i)}$))$_{i=1,…,N}$ )，因為g$_{\perp}$(x$^{(i)}$)=0

3. ||g$_{//}$|| ≤ ||g|| ，因此$\Omega$(||g$_{//}$||$_{RKHS}$) ≤ $\Omega$(||g||$_{RKHS}$)，而將g代換成g$_{//}$不會對loss function和constraints造成任何影響，所以我們應該將g代換成g$_{//}$

4. 因此，若得以最小化loss function，必須有g$_{\perp}$ =0，此時$\tilde{g}$ = g$_{//}$ + b$\psi$

4. Prove the Semiparametric Representer theorem.

5. 由於 $g_\parallel \in \text{span}(k(\mathbf{x}^{(1)},\cdot),\dots, k(\mathbf{x}^{(N)},\cdot))$，我們可以將 $\tilde{g}^{(i)}_{i=1,\dots,N}$ 改寫成

$$\tilde{g}^{(i)} = \sum_{t=1}^{N} c_t\, k(\mathbf{x}^{(t)},\mathbf{x}^{(i)}) + b\psi(\mathbf{x}^{(i)})$$

6. 因此，$\tilde{h}(\mathbf{x})$ 便有 $\sum_{t=1}^{N} c_t\, k(\mathbf{x}^{(t)},\mathbf{x}) + b\psi(\mathbf{x})$ 的形式

6. Show that the *Area Under the ROC Curve* (AUC) is equal to the probability that a classifier ranks a randomly chosen positive instance higher than a randomly chosen negative one.

**table (6-a)**

| Label | Rank |
|-------|------|
| +1 | 1 |
| +1 | 2 |
| +1 | 3 |
| -1 | 4 |
| +1 | 5 |
| -1 | 6 |
| -1 | 7 |
| +1 | 8 |
| -1 | 9 |
| +1 | 10 |

**table (6-b)**

| $X^+_i$ \ $X^-_j$ | 4 | 6 | 7 | 9 |
|---|---|---|---|---|
| 10 | O | O | O | O |
| 8 | O | O | O | 1 |
| 5 | O | 1 | 1 | 1 |
| 3 | 1 | 1 | 1 | 1 |
| 2 | 1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 1 | 1 |

Area occupied by *1*

$$= \sum_{i=0}^{N^+} \sum_{j=0}^{N^-} C(i,j)$$

- 在任何一個資料集，對於這個Classifier而言，隨機選擇一個Positive instance($X^+_i$),$1\leq i\leq N^+$ 和一個Negative instance($X^-_j$),$1\leq j\leq N^-$，那麼 Rank($X^+_i$)<Rank($X^-_j$)的機率，即為：

$$\sum_{i=0}^{N^+} \sum_{j=0}^{N^-} C(i,j) / (N^+N^-), \quad C(i,j) = \begin{cases} 1, & \text{if } Rank(X^+_i) < Rank(X^-_j) \\ 0, & \text{if } Rank(X^+_i) > Rank(X^-_j) \end{cases}$$

如table (6-b)

而這恰好就是「左方表格中，被*1*所佔據的面積」 table (6-b)

(表格中*O*和*1*所在的區域其長寬皆為1，而每小格面積均等)，

而表格中紅色折線即為ROC Curve