

## 1. [15 points] Logistic Regression: Training stability

In this problem, we will be delving deeper into the workings of logistic regression. The goal of this problem is to help you develop your skills debugging machine learning algorithms (which can be very different from debugging software in general).

We have provided a implementation of logistic regression in `src/p01_lr.py`, and two labeled datasets  $A$  and  $B$  in `data/ds1_a.txt` and `data/ds1_b.txt`

Please do not modify the code for the logistic regression training algorithm for this problem. First, run the given logistic regression code to train two different models on  $A$  and  $B$ . You can run the code by simply executing `python p01_lr.py` in the `src` directory.

- (a) [2 points] What is the most notable difference in training the logistic regression model on datasets  $A$  and  $B$ ? *training dataset B does not converge. (slowly).*
- (b) [5 points] Investigate why the training procedure behaves unexpectedly on dataset  $B$ , but not on  $A$ . Provide hard evidence (in the form of math, code, plots, etc.) to corroborate your hypothesis for the misbehavior. Remember, you should address why your explanation does *not* apply to  $A$ . *B can be perfectly separated by a boundary. Many solutions can be used to fit training set.*
- Hint:** The issue is not a numerical rounding or over/underflow error.
- (c) [5 points] For each of these possible modifications, state whether or not it would lead to the provided training algorithm converging on datasets such as  $B$ . Justify your answers.
- Using a different constant learning rate.
  - Decreasing the learning rate over time (e.g. scaling the initial learning rate by  $1/t^2$ , where  $t$  is the number of gradient descent iterations thus far).
  - Linear scaling of the input features.
  - Adding a regularization term  $\|\theta\|_2^2$  to the loss function.
  - Adding zero-mean Gaussian noise to the training data or labels.
- (d) [3 points] Are support vector machines, which use the hinge loss, vulnerable to datasets like  $B$ ? Why or why not? Give an informal justification.

**Hint:** Recall the distinction between functional margin and geometric margin.

This can cause  
the ~~loss~~ a very  
small learning  
rate which falls  
in the range of  
 $\uparrow 1e-15$

should be  
yes.

- i. the reason for slow convergence is caused by small grad, changing learning rate does not ~~defix~~ it.
- ii. No.
- iii. No. *can have the same affect as scaling weight.*
- iv. Yes. ~~more~~ then loss can't be zero.
- v. Yes. then training set would not be perfectly separatable.

*(d) No. SVM uses geometric margin. When dataset is perfectly separatable, it would only be one answer.*

2. (a).  $\ell(\theta) = \frac{1}{n} \sum_{i=1}^n P(y^{(i)} | x^{(i)}; \theta)$

$\ell(\theta) = \sum_{i=1}^n y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log (1 - h(x^{(i)}))$

$\frac{\partial}{\partial \theta_j} \ell(\theta) = \sum_{i=1}^n \left( y^{(i)} \frac{1}{h(x^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - h(x^{(i)})} \right) h(x^{(i)}) (1 - h(x^{(i)})) x_j$

$= \sum_{i=1}^n \left( y^{(i)} (1 - h(x^{(i)})) - (1 - y^{(i)}) h(x^{(i)}) \right) x_j$

$= \sum_{i=1}^n (y^{(i)} - h(x^{(i)})) x_j$  We know  $\sum_{i=1}^n x_j = 0$ .

$\frac{\partial}{\partial \theta_0} = \sum_{i=1}^n (y^{(i)} - h(x^{(i)})) = 0$ . set it to 0.

$\sum_{i=1}^n y^{(i)} = \sum_{i=1}^n h(x^{(i)})$

because  $i \in (0, 1) \rightarrow 1 \leq i \leq n$ .

b). ~~No. we can only get the probability of a  $P(y^{(i)} = 1 | x^{(i)}; \theta)$ .~~

~~When  $P \rightarrow 0.5$ , the probability for predict mistakenly is pretty large.~~

~~The converse does not hold because when a model achieves perfect accuracy,~~

~~Then  $P(y^{(i)} | x^{(i)}; \theta) = 1$ .~~

No, it is saying the sum of

ci). if perfectly accuracy is respect to training set, (infinitely large).

~~Yes. because the interval  $I$  can be infinitely small. When shrinking size to 1~~

~~then every prediction (probability) can be perfect accuracy.~~

~~The converse does not hold. because when interval only has one point.~~

if ~~set~~ test set is infinitely large.

No. perfectly calibrated  $\rightarrow$  perfectly accuracy. it only implies the average is

- (c) [5 points] **Coding problem.** We will now tune the hyperparameter  $\tau$ . In `src/p05c_tau.py`, find the MSE value of your model on the validation set for each of the values of  $\tau$  specified in the code. For each  $\tau$ , plot your model's predictions on the validation set in the format described in part (b). Report the value of  $\tau$  which achieves the lowest MSE on the **valid** split, and finally report the MSE on the **test** split using this  $\tau$ -value.

this is wrong

0.017

because I did not notice the interval  $I$  is open on both sides.

(b) perfectly calibrated.  $\rightarrow$  perfect accuracy.

when shrink the size of interval to 1. then  $P(y^{(i)}=1|x^{(i)}; \theta) = \mathbb{1}\{y^{(i)}=1\}$ .

since rhs is either 1 or 0, the lhs must also be 1 or 0.

Then for  $\forall x^{(i)}$ ,  $P(y^{(i)}|x^{(i)}; \theta) = 1$  or 0. If it is 1, and  $y^{(i)}=1$ , the prediction matches the  $\boxed{y^{(i)}=1}$ . When it is 0,  $P(y^{(i)}|x^{(i)}; \theta) = 0$ . and there is no example. so. it achieves perfect accuracy.

perfect accuracy  $\rightarrow$  perfectly calibrated.

$\hookrightarrow$  for  $\forall x^{(i)}$  the prediction matches the label.  $h_\theta(x^{(i)}) = y^{(i)}$ , then  $h_\theta(x^{(i)})$  is 1 or 0.

for any interval  $(a, b)$ . suppose there are  $c$   $y^{(i)}=1$  labels. and  $d$   $y^{(i)}=0$  labels.

where  $c, d \geq 0$ .  $\sum_{i \in (a,b)} P(y^{(i)}=1|x^{(i)}; \theta) = \sum_{i \in (a,b)} h_\theta(x^{(i)}) = c \times 1 + d \times 0 = c$ .

$\sum_{i \in I(a,b)} \mathbb{1}(y^{(i)}=1) = c \times 1 + d \times 0 = c$ . Thus  $\sum_{i \in I(a,b)} \mathbb{1}(y^{(i)}=1) = \sum_{i \in I(a,b)} P(y^{(i)}=1|x^{(i)}; \theta)$

(c). from  $\frac{\partial}{\partial \theta_0} \ell(\theta) = \sum_{i=1}^n (y^{(i)} - h(x^{(i)})) + 2\theta_0 = 0$ . then  $\sum_{i=1}^n h(x^{(i)}) = \sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\} + 2\theta_0$

it shifts the makes the model not well calibrated.

wrong!

2.(b). Assume the model achieves perfect accuracy.

And pick the  $I = (a, b) = (0.5, 1)$ .

$$\forall i \in I_{a,b} \quad 0.5 < P(y^{(i)} = 1 | x^{(i)}; \theta) < 1$$

$$\text{Then } \sum_{i \in I_{a,b}} P(y^{(i)} = 1 | x^{(i)}; \theta) < \sum_{i \in I_{a,b}} 1(y^{(i)} = 1).$$

Since all predictions must be 1 to achieve perfect accuracy.

This violates the well-calibrated property

Thus the model is not well-calibrated.

Assume the model is well-calibrated. Then.

$$\frac{\sum P(y^{(i)} = 1 | x^{(i)}; \theta)}{|\{i \in I_{a,b}\}|} = \frac{\sum 1(y^{(i)} = 1)}{|\{i \in I_{a,b}\}|}$$

pick  $I = (0.5, 1)$ .

Because the model output is  $0 < P(y^{(i)} = 1 | x^{(i)}; \theta) < 1$ .

The average  $< \frac{1}{n} \sum P(y^{(i)} = 1 | x^{(i)}; \theta) < 1$ .

Because when  $0 < p < 1$ , the model would predicts 1.

But the total sum of  $y^{(i)} = 1$  is not equal to the total num of examples.

So the model is not perfectly accurate!

Rewrite

$$3. (a). \theta_{MAP} = \arg \max_{\theta} p(\theta | x, y)$$

$$= \arg \max_{\theta} p(\theta | x, y) \cdot p(y | x) \quad \text{since } p(y | x) \text{ is constant.}$$

$$= \propto \frac{p(\theta, x, y)}{p(x, y)} \cdot \frac{p(y, x)}{p(x)}.$$

$$= \propto \frac{p(y | \theta, x) \cdot p(\theta | x)}{p(x)}.$$

$$= \propto \frac{p(y | \theta, x) \cdot p(\theta | x) \cdot p(x)}{p(x)}.$$

$$= \arg \max_{\theta} p(y | \theta, x) p(\theta).$$

$$(b). \theta_{MAP} = \arg \max_{\theta} p(\theta | x, y) = \arg \max_{\theta} p(y | x, \theta) p(\theta)$$

$$\text{Since } \theta \sim \mathcal{N}(0, \eta^2 I) \quad = \arg \max_{\theta} p(y | x, \theta) \frac{1}{\eta \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{\theta - 0}{\eta} \right)^2} \quad \text{since } \frac{1}{\eta \sqrt{2\pi}} \text{ is constant.}$$

$$= \arg \max_{\theta} p(y | x, \theta) e^{-\frac{1}{2} \left( \frac{\theta^T \theta}{\eta^2 I} \right)}$$

$$= \arg \max_{\theta} \log p(y | x, \theta) + \left( -\frac{1}{2} \cdot \frac{1}{\eta^2 I} \cdot \|\theta\|_2^2 \right).$$

$$= \arg \min_{\theta} -\log p(y | x, \theta) + \frac{1}{2\eta^2 I} \|\theta\|_2^2$$

$$\text{where } \lambda = \frac{1}{2\eta^2 I}$$

$$\begin{aligned}
(c). \quad \theta_{\text{MAP}} &= \arg \min_{\theta} -\log \left( \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta) \right) + \lambda \|\theta\|_2^2 \\
&= \arg \min_{\theta} - \left( \sum_{i=1}^m \log \left( \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right) \right) \right) + \lambda \|\theta\|_2^2 \\
&= \arg \min_{\theta} - \sum_{i=1}^m \log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \left( \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right) + \lambda \|\theta\|_2^2 \\
&= \arg \min_{\theta} \sum_{i=1}^m -\log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) + \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} + \lambda \|\theta\|_2^2 \\
&\quad \text{Since } -\log \left( \frac{1}{\sqrt{2\pi}\sigma} \right), \sigma^2 \text{ are constant.} \\
&= \arg \min_{\theta} \left( \sum_{i=1}^m \frac{1}{2} (y^{(i)} - \theta^T x^{(i)})^2 \right) + \lambda \|\theta\|_2^2
\end{aligned}$$

$$\begin{aligned}
(d). \quad \theta_{\text{MAP}} &= \arg \max_{\theta} P(y|x, \theta) P(\theta) \quad \sum_{i=1}^n \\
&= \arg \max_{\theta} \log(P(y|x, \theta)) + \log \left( \frac{1}{2bI} \exp \left( -\frac{|\theta|}{bI} \right) \right) \\
&= \arg \max_{\theta} \log P(y|x, \theta) + \sum_{i=1}^n \log \left( \frac{1}{2bI} \right) - \frac{|\theta|}{bI} \\
&= \arg \max_{\theta} \log \left( \prod_{i=1}^n P(y^{(i)} | x^{(i)}, \theta) \right) - \frac{n}{bI} \|\theta\|_1, \quad \gamma = \frac{n}{bI}. \\
&= \arg \min_{\theta} - \sum_{i=1}^n \log \left( \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right) \right) + \gamma \|\theta\|_1 \\
&= \arg \min_{\theta} \left( \sum_{i=1}^n -\log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) + \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right) + \gamma \|\theta\|_1 \\
&= \arg \min_{\theta} \sum_{i=1}^n \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} + \gamma \|\theta\|_1 \\
&= \arg \min_{\theta} \|X\theta - \vec{y}\|_2^2 + \gamma \|\theta\|_1
\end{aligned}$$

Rewrite 3C.

$$p(\vec{y} | X; \theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$\log(p(\vec{y} | X; \theta)) = \sum_{i=1}^m \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) + \sum_{i=1}^m \frac{-(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$$

$$= -m \log(\sqrt{2\pi}\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

$$= -m \left( \frac{1}{2} \log(2\pi) + \log(\sigma) \right) - \frac{1}{2\sigma^2} \|\vec{y} - X\theta\|_2^2$$

$$= \underbrace{-\frac{1}{2}m \log(2\pi)}_{\text{constant}} - \underbrace{m \log(\sigma)}_{\text{constant}} - \frac{1}{2\sigma^2} \|\vec{y} - X\theta\|_2^2$$

$$\theta_{\text{MAP}} = \arg \min_{\theta} -\log p(\vec{y} | X; \theta) + \frac{1}{2\eta^2 I} \|\theta\|_2^2$$

$$= \arg \min_{\theta} \frac{1}{2\sigma^2} \|\vec{y} - X\theta\|_2^2 + \frac{1}{2\eta^2 I} \|\theta\|_2^2$$

$$\nabla_{\theta} = \frac{1}{\sigma^2} (-X^T \vec{y} + X^T X \theta) + \frac{1}{\eta^2 I} \theta = 0$$

$$\frac{1}{\sigma^2} X^T X \theta - \frac{1}{\sigma^2} X^T \vec{y} + \frac{1}{\eta^2 I} \theta = 0$$

$$\left( \frac{1}{\sigma^2} X^T X + \frac{1}{\eta^2 I} \right) \theta = \frac{1}{\sigma^2} X^T \vec{y}$$

$$(X^T X + \frac{\sigma^2}{\eta^2 I}) \theta = X^T \vec{y}$$

$$\theta = (X^T X + \frac{\sigma^2}{\eta^2 I})^{-1} X^T \vec{y}$$

$$\begin{aligned} \nabla_{\theta} \|\vec{y} - X\theta\|_2^2 &= \nabla_{\theta} (\vec{y} - X\theta)^T (\vec{y} - X\theta) \\ &= \nabla_{\theta} (\vec{y}^T - \theta^T X^T) (\vec{y} - X\theta) \\ &= \nabla_{\theta} \vec{y}^T \vec{y} - \underbrace{\vec{y}^T X\theta - \theta^T X^T \vec{y}}_{2\vec{y}^T X\theta} + (\theta^T X^T) (X\theta) \\ &= -2X^T \vec{y} + 2X^T X\theta \end{aligned}$$

Rewrites 3 (d).  $\theta \sim \mathcal{L}(0, bI)$ .

$$P(\theta) = \prod_{i=1}^n \frac{1}{2b} \exp\left(-\frac{|y_i - 0|}{b}\right).$$

$$= \left(\frac{1}{2b}\right)^n \exp\left(-\frac{\|\theta\|_1}{b}\right)$$

$$\log P(\theta) = \log\left(\frac{1}{(2b)^n}\right) - \frac{\|\theta\|_1}{b}$$

$$= -n \log(2b) - \frac{\|\theta\|_1}{b} \quad \text{since } -n \log(2b) \text{ is constant.}$$

$$\theta_{\text{MAP}} = \arg \max_{\theta} \log P(y|x, \theta) - \frac{\|\theta\|_1}{b} - n \log(2b)$$

$$= \arg \min_{\theta} -\log P(y|x, \theta) + \frac{\|\theta\|_1}{b}$$

$$= \arg \min_{\theta} \frac{1}{2\sigma^2} \|\vec{y} - X\theta\|_2^2 + \frac{1}{b} \|\theta\|_1.$$

$$= \arg \min_{\theta} \|\vec{y} - X\theta\|_2^2 + \frac{2\sigma^2}{b} \|\theta\|_1,$$

$$\text{where } \gamma = \frac{2\sigma^2}{b}$$



(a). ~~where~~ where  $x$  is the  $i$ th entry and  $z$  is the  $j$ th entry.

$$K_{ij} = K(x, z) = K_1(x, z) + K_2(x, z) \quad \left\{ \begin{array}{l} \text{since } K_1, K_2 \text{ are kernels.} \end{array} \right.$$

$$K_{ji} = K(z, x) = K_1(z, x) + K_2(z, x) \quad \left\{ \begin{array}{l} K_{1,ij} = K_{1,ji} \Rightarrow K_1(x^{(i)}, x^{(j)}) = K_1(x^{(j)}, x^{(i)}) \\ \text{same for } K_2. \end{array} \right.$$

Since R side,  $K_{ij} = K_{ji}$ .

~~Thus~~ Thus the ~~kernel~~ kernel matrix is symmetric.

let  $z$  be any vector,  $z \in \mathbb{R}^n$ .

$$\begin{aligned} z^T K z &= \sum_i \sum_j z_i K_{ij} z_j = \sum_i \sum_j z_i (K_1(x^{(i)}, x^{(j)}) + K_2(x^{(i)}, x^{(j)})) z_j \\ &= \sum_i \sum_j z_i K_1(x^{(i)}, x^{(j)}) z_j + \sum_i \sum_j z_i K_2(x^{(i)}, x^{(j)}) z_j \\ &= \sum_i \sum_j z_i K_{1,ij} z_j + \sum_i \sum_j z_i K_{2,ij} z_j \\ &= z^T K_1 z + z^T K_2 z \end{aligned}$$

since  $K_1, K_2$  are semi-definite, then both terms  $\geq 0$ .

Thus,  $K$  is semi-definite. Thus,  $K$  is a valid kernel.

$$(b). K_{ij} = K(x^{(i)}, x^{(j)}) = K_1(x^{(i)}, x^{(j)}) - K_2(x^{(i)}, x^{(j)}) \quad \begin{array}{l} \text{since } K_1(x^{(i)}, x^{(j)}) = K_1(x^{(j)}, x^{(i)}) \\ K_2 \quad \quad \quad = K_2 \quad \dots \end{array}$$

$$K_{ji} = K(x^{(j)}, x^{(i)}) = K_1(x^{(j)}, x^{(i)}) - K_2(x^{(j)}, x^{(i)})$$

Thus  $K_{ij} = K_{ji}$ . symmetry.

$$\begin{aligned} z^T K z &= \sum_i \sum_j z_i K_{ij} z_j = \sum_i \sum_j z_i \cancel{K_{1,ij}} K_{1,ij} z_j - \sum_i \sum_j z_i K_{2,ij} z_j \\ &= \cancel{z^T K_1 z} z^T K_1 z - z^T K_2 z \end{aligned}$$

When ~~for~~ for  $z \in \mathbb{R}^n$ , <sup>when</sup>  $z^T K_1 z < z^T K_2 z$ ,  $K(x, z)$  is not a valid kernel.

(c).  $K_{ij} = K(x^{(i)}, x^{(j)}) = \alpha K_1(x^{(i)}, x^{(j)})$  since  $K_1(x^{(i)}, x^{(j)}) = K_1(x^{(j)}, x^{(i)})$

$K_{ji} = K(x^{(j)}, x^{(i)}) = \alpha K_1(x^{(j)}, x^{(i)})$   $K_{ij} = K_{ji}$  thus  $K$  is symmetric.

$Z^T K Z = \sum_i \sum_j Z_i^T K_{ij} Z_j = \sum_i \sum_j Z_i^T \alpha K_{ij} Z_j$

$= \alpha \sum_i \sum_j Z_i^T K_{ij} Z_j$  since  $K_1$  is kernel.

$\sum_i \sum_j Z_i^T K_{ij} Z_j \geq 0$ . then  $Z^T K Z \geq 0$ .

(d).  
Wrong

$K_{ij} = -\alpha K_1(x^{(i)}, x^{(j)})$

since  $K_1(x^{(i)}, x^{(j)}) = K_1(x^{(j)}, x^{(i)})$

$K_{ji} = -\alpha K_1(x^{(j)}, x^{(i)})$

$K_{ij} = K_{ji}$  thus  $K$  is symmetric.

$Z^T K Z = \sum_i \sum_j Z_i^T K_{ij} Z_j = -\alpha \sum_i \sum_j Z_i^T K_{ij} Z_j$  since  $K_1$  is kernel.

$\sum_i \sum_j Z_i^T K_{ij} Z_j \geq 0$ . thus  $Z^T K Z \geq 0$ .

$Z^T K Z = Z^T (-\alpha K_1) Z$   
 $= -\alpha Z^T K_1 Z$   
 $\leq 0$ .

(e).  $K_{ij} = K_1(x^{(i)}, x^{(j)}) K_2(x^{(i)}, x^{(j)})$

$K_{ji} = K_1(x^{(j)}, x^{(i)}) K_2(x^{(j)}, x^{(i)})$

$= K_1(x^{(i)}, x^{(j)}) K_2(x^{(i)}, x^{(j)})$   $K_1, K_2$  are kernels

$= K_{ij}$ .

Thus matrix  $K$  is symmetric.

For an arbitrary  $Z \in \mathbb{R}^n$

$Z^T K Z = \sum_i \sum_j Z_i K_{ij} Z_j = \sum_i \sum_j (Z_i K_{ij} Z_j) (K_{2ij}) = \sum_i \sum_j \sum_a \sum_b (Z_i K_{ij} Z_j) (c^T K_{ab} d)$  when  $a=i$   $b=j$

or  $(0 \cdot K_{2ab} \cdot 0)$  when  $a \neq i$  or  $b \neq j$

(1). we know  $K_{1ij} = K_{1ji}$  and  $K_{2ij} = K_{2ji}$

Thus  $Z_j K_{1ji} Z_i K_{2ji} = Z_i K_{1ij} Z_j K_{2ij}$

Thus  $c^T K_2 d = d^T K_2 c$

Then  $\sum_i \sum_j (Z_i K_{ij} Z_j) (c^T K_2 d) = \sum_i \sum_j (Z_i K_{ij} Z_j) \left( \sum_{i \neq j} (c^T K_2 d)^2 + \sum_{i=j} c^T K_2 d \right)$

since  $K_2$  is kernel  $c^T K_2 d \geq 0$  for  $c=d$

also,  $(c^T K_2 d)^2 \geq 0$ . also  $\sum_i \sum_j Z_i K_{ij} Z_j \geq 0$ .

Thus  $Z^T K Z \geq 0$ . Thus  $K$  is positive semidefinite

$= \sum_i \sum_j (Z_i K_{ij} Z_j) (c^T K_2 d)$  while  $c$  is  $\mathbb{R}^n$  vector with  $i$ th term being 1 and others being 0.

$d$  is  $\mathbb{R}^n$  vector with  $j$ th term being 1 and others being 0.

(f).  $k_{ij} = K(x^{(i)}, x^{(j)}) = f(x^{(i)}) f(x^{(j)}) = f(x^{(j)}) f(x^{(i)}) = K(x^{(j)}, x^{(i)}) = k_{ji}$   $K$  is symmetric.

$$Z^T K Z = \sum_i \sum_j z_i k_{ij} z_j = \sum_i \sum_j z_i f(x^{(i)}) f(x^{(j)}) z_j$$

since  $\sum_i z_i f(x^{(i)}) f(x^{(j)}) z_j = \sum_i \sum_j \left( z_i f(x^{(i)}) f(x^{(j)}) z_j \right)^2 + \sum_i \sum_{j \neq i} z_i f(x^{(i)}) f(x^{(j)}) z_j$

Thus  $Z^T K Z \geq 0$ . Thus  $K$  is positive semidefinite.

(g).  $k_{ij} = K_2(\phi(x^{(i)}), \phi(x^{(j)})) = K_2(\phi(x^{(j)}), \phi(x^{(i)})) = k_{ji}$   $K$  is symmetric.

$$Z^T K Z = \sum_i \sum_j z_i k_{ij} z_j = \sum_i \sum_j z_i K_2(\phi(x^{(i)}), \phi(x^{(j)})) z_j \quad \text{since } K_2 \text{ is a kernel then } \sum_i \sum_j z_i k_{ij} z_j \geq 0.$$

Thus  $Z^T K Z \geq 0$ .  $K$  is positive semidefinite.

(h).  $k_{ij} = p(k_1(x^{(i)}, x^{(j)})) = p(k_1(x^{(j)}, x^{(i)})) = k_{ji}$   $K$  is symmetric.

$$Z^T K Z = \sum_i \sum_j z_i p(k_1(x^{(i)}, x^{(j)})) z_j = \sum_i \sum_j z_i (a_1 k_1(x^{(i)}, x^{(j)}) + a_2 k_1(x^{(j)}, x^{(i)}) + \dots + a_p k_1(x^{(i)}, x^{(j)})) z_j$$

$$= \sum_i \sum_j \sum_p a_p \sum_p z_i k_1(x^{(i)}, x^{(j)}) z_j = \sum_p a_p \sum_i \sum_j z_i k_1(x^{(i)}, x^{(j)}) z_j$$

Since  $k_1$  is kernel  $\sum_i \sum_j z_i k_1(x^{(i)}, x^{(j)}) z_j \geq 0$ .

Since  $a_p \geq 0$ , the whole term  $\geq 0$ .  $Z^T K Z \geq 0$ ,  $K$  is positive semidefinite.

There is a missing power in  $k_1(x^{(i)}, x^{(j)})^k$

5. (a) (i). Use a map to represent infinite dimensional vector. index: value.

where index represent the position of the entry, non-recorded entries are 0.

$\theta^{(0)}$  is the empty map.

$$(ii). \theta^{(i+1)}(x^{(i+1)}) = \sum_{key \in \theta^{(i)} \text{ map and } \phi(x^{(i+1)}) \text{ map}} \theta^{(i)}[key] \cdot \phi(x^{(i+1)})[key].$$

$$(iii). \theta^{(i+1)} := \theta^{(i)} + \lambda (y^{(i+1)} - \text{sign}(\sum_{key} \theta^{(i)}[key] \phi(x^{(i+1)})[key])) \phi(x^{(i+1)})$$

---

(a) (i). use a zero vector whose dimension is same as  $x^i$ . to represent  $\theta^{(i)}$

$\theta^{(i)}$  is represented in the same manner st.  $\theta^{(i)} = \phi(\theta_{true}^{(i)})$ .

(ii). since  $\theta^{(i)T} \phi(x^{(i+1)})$  is an inner product. we can represent it as kernel trick.

Using  $k(\theta_{true}^{(i)}, x^{(i+1)})$  to represent it.

$$(iii). \theta_{true}^{(i+1)} := \theta_{true}^{(i)} + \lambda (y^{(i+1)} - k(\theta_{true}^{(i)}, x^{(i+1)})) x^{(i+1)}.$$

---

(a) (i). We know that  $\theta^{(i)}$  is a linear combination of  $\phi(x_1), \phi(x_2), \dots, \phi(x_i)$

we can record the coefficients of  $\phi(x)$  and number  $i$

the  $x_k$  entry can be got by  $i \% n$ . where  $n$  is total number of  $x$ .

The coefficients can be got as ~~vector~~ vector  $C[i]$ . also a learning rate.

$$\theta^{(i)} = \lambda \sum_j b^{(j)} \phi(x^{(j)}) \text{ since } \phi \text{ is a known mapping. recording } x^{(j)} \text{ to } x^{(i)} \text{ will get}$$

sufficient to ~~compute~~ compute  $\phi(x^{(i)})$  to  $\phi(x^{(j)})$ .

$$\theta^{(0)} =$$

5. (a). (i). ~~the~~  $\theta^{(i)}$  can be represented as a linear combination of  $\phi(x_1), \phi(x_2), \dots, \phi(x_i)$ .

$$\theta^{(i)} = \sum_{k=1}^i b^{(k)} \phi(x^{(k)}). \quad \text{Since we know the } \phi \text{ mapping will not change.}$$

$\phi(x^k)$  can be represented as  $x^k$ . Assume there are  $n$  inputs.

we can record  $k$  and use it as  $k \% n$  to represent index of  $x$ .

We also need to record  $b^{(k)}$  as an array of coefficients.

Knowing  $\vec{b}$ ,  $i$  are enough to represent  $\theta^{(i)}$ .

$\theta^{(0)}$  can be represented as  $i=0$   $\vec{b}=0$ .

$$\begin{aligned} \text{(ii). } h_{\theta^{(i)}}(x^{(i+1)}) &= g(\theta^{(i)\top} \phi(x^{(i+1)})) = g\left(\sum_{k=1}^i b^{(k)} \phi(x^{(k)}) \cdot \phi(x^{(i+1)})\right) \\ &= g\left(\sum_{k=1}^i b^{(k)} [\phi(x^{(k)}) \cdot \phi(x^{(i+1)})]\right) \\ &= g\left(\sum_{k=1}^i b^{(k)} K(x^{(k)}, x^{(i+1)})\right). \end{aligned}$$

$$\text{(iii). } b^{i+1} = y^{(i+1)} - h_{\theta^{(i)}}(x^{(i+1)}).$$

Put  $b^{i+1}$  to the end of  $\vec{b}$ , and  $i := i+1$

(c). the dot product kernel performs poorly because the data is not linearly separable in two-dim. rbf raises it to a higher dimension which makes it separable.