

Bhishan Poudel

HW Assignment 8 (Due by 10:30am on Dec 7)

1 Theory (150 points)

1. [Kernel Nearest Neighbor, 50 points]

The nearest-neighbour classifier 1-NN assigns a new input vector \mathbf{x} to the same class as that of the nearest input vector \mathbf{x}_n from the training set, where in the simplest case, the distance is defined by the Euclidean metric $\|\mathbf{x} - \mathbf{x}_n\|^2$. By expressing this rule in terms of scalar products and then making use of kernel substitution, formulate the nearest-neighbour classifier for a general nonlinear kernel.

2. [Distance-Weighted Nearest Neighbor, 50 points]

We have seen how to use kernels to formulate a distance-weighted nearest neighbor algorithm, when the labels are binary. Formulate a kernel-based, distance-weighted nearest neighbor that works for K classes, where $K \geq 2$.

3. [Naive Bayes, 50 points]

The Naive Bayes algorithm for text categorization presented in class treats all sections of a document equally, ignoring the fact that words in the title are often more important than words in the text in determining the document category. Describe how you would modify the Naive Bayes algorithm for text categorization to reflect the constraint that words in the title are K times more important than the other words in the document for deciding the category, where K is an input parameter (include pseudocode).

4. [Logistic Regression (*), 50 points]

Assume that a binary feature x_i is equal to 1 for all training examples \mathbf{x} belonging to a particular class C_k , and zero otherwise (i.e. x_i perfectly separates examples from class C_k from all other examples). Show that in this case the magnitude of the ML solution for \mathbf{w}_k goes to infinity, thus motivating the use of a prior over the parameters (Hint: use the fact that the gradient on slide 24 must vanish at the solution).

2 Submission

Turn in a hard copy of your homework report at the beginning of class on the due date. On this theory assignment, **clear and complete explanations and proofs of your results are as important as getting the right answer.**

Q11 K-Nearest Neighbor (kernel based classifier)

a) 1 nearest neighbor

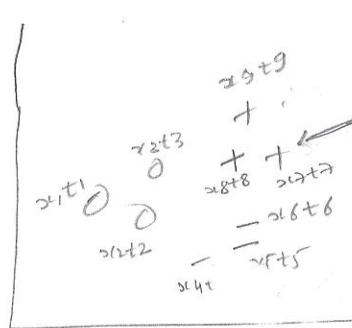
training dataset : $(x_1, t_1), (x_2, t_2), \dots, (x_n, t_n)$

test dataset : x

required output : \hat{y} (label of x across of x)

t_1, t_2, \dots, t_n are classes from c_1, c_2, \dots, c_m

for example { circle plus minus }
 $c_1 \quad c_2 \quad c_m$



x (test)
 predicted label = \hat{y}
 true label = y

distance

$$d(x, x') = \|x - x'\| = \sqrt{\|x - x'\|^2}$$

(Euclidean distance)

How to find \hat{y} ?

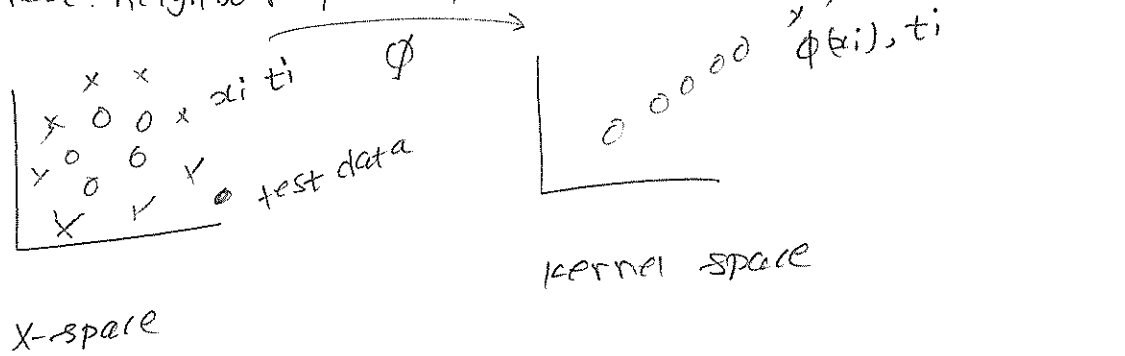
\Rightarrow For 1-nearest neighbor, calculate the distance (maybe Euclidean, cosine, edit, etc) from test point \vec{x} to all the data points $x_i, 1 \leq i \leq n$ then find the minimum distance data point x' .

Assign label of x' as the predicted label to \vec{x} .

$$\hat{y} = \text{class of } x'$$

= class of x_i such that $d(x, x_i)$ is minimum

① 1-nearest neighbor kernel space,



Here, when calculating distance, we use kernel method,

X-space:

$$d(x, x') = \|x - x'\| = \sqrt{\|x - x'\|^2} = \sqrt{(x - x')^T (x - x')}$$

(Euclidean distance)

kernel-space:

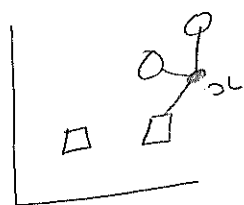
$$d(x, x') = \sqrt{(\phi(x) - \phi(x'))^T (\phi(x) - \phi(x'))}$$

$$= \sqrt{K(x, x')}$$

→ we compute the distance of test x to the all the training examples x_i using kernel method then choose the minimum distance training point x' and then,

$$\text{label of } x = \text{label of } x'$$

③ k-nearest neighbor kernel method.



3-nearest neighbor
two are circle
 \therefore label of x is circle.

kernel \Rightarrow distances are calculated using
kernel trick

$$d(x_i, x_j) = \sqrt{k(x_i, x_j)}$$

kernel can be polynomial kernel,
gaussian kernel and so on.

Let x_1, x_2, \dots, x_k are nearest of x

then,

$$y(x) = \operatorname{argmax}_{t \in T} \sum_{i=1}^k \delta_t(x_i)$$

Here, $t_i = \{0, 1\}$ ^{circle square}

$k=3$

for x_1 and x_2 $t = t_i = 0$ number = 2

for x_3 $t = t_i = 1$ number = 1

$$\operatorname{argmax} = 0$$

Summary

input $(x_1, t_1), (x_2, t_2), \dots, (x_n, t_n)$ $t_n \in \{c_1, c_2, \dots, c_K\}$

test $x, t = ?$

map input from x -space to kernel-space

x -space input $\phi(x_1), t_1$ $\phi(x_2), t_2$ $\phi(x_n), t_n$

kernel space test $\phi(x), t = ?$

compute similarity $k(\phi(x_i))$ for all $x_i, 1 \leq i \leq n$

(compute distance $= \sqrt{k(\phi(x_i))}$ is redundant since
distance $<$ similarity)

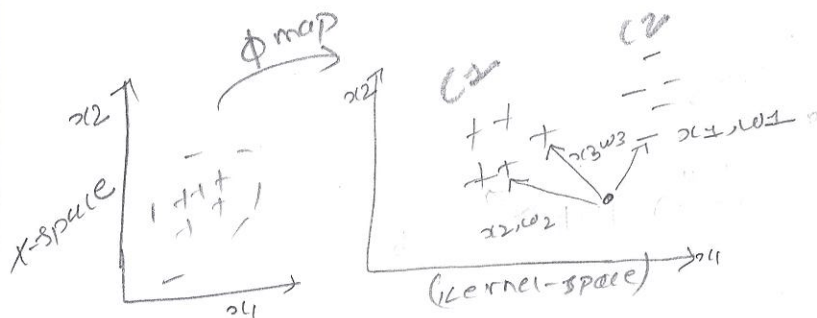
then, predicted class of x is,

$$y(x) = \underset{t \in T}{\operatorname{argmax}} \sum_{i=1}^n f_t(t_i)$$

PN2

distance-weighted Nearest Neighbor for classification

a) kernel space distance-weighted KNN binary classification



$$t_i \in \{+1, -1\}$$

3- nearest neighbor
contribution from + : $w_2 + w_3$

contribution from - : w_1

suppose $w_2 + w_3 > w_1$

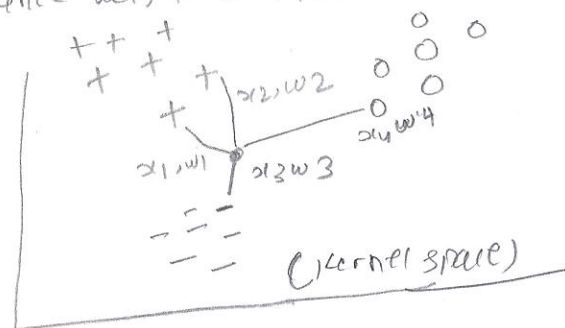
then label of $x = +$

- x is test example
- x_1, x_2, \dots, x_k are the nearest points to x , calculated using Euclidean distance in kernel space.
- $k(x, x_i)$ is the measure of similarity between point x and x_i in kernel space.

$$\hat{y} = \text{sign} \sum_{i=1}^k k(x, x_i) \delta_t(t_i)$$

b) kernel-based distance weighted KNN for multiclass classification

Note: nearer points have more contribution



4- nearest neighbor
contribution from class $c_1 = +$
 $= w_1 + w_2$

contribution from class $c_2 = -$
 $= w_3$

contribution from class $c_3 = 0$
 $= w_4$

$$\hat{y} = \text{argmax}_{t \in T} \sum_{i=1}^k k(x, x_i) \delta_t(t_i)$$

suppose $w_1 + w_2 > w_3$ and $w_1 + w_2 > w_4 \Rightarrow \hat{y} = +$

Ans

Q3 Naïve Bias

① simple text classification using Naïve Bias

- each document is an example vector
- a document x has n words w_1, w_2, \dots, w_n
- from all the documents $D = \{x_1, x_2, \dots, x_m\}$ we create vocabulary file (this is feature vector for whole dataset D)
- vocabulary $V = \{w_1, w_2, \dots, w_{|V|}\}$, there are $|V|$ number of words in vocabulary.
- each document x belongs to one of the categories $C_k = \{c_1, c_2, \dots, c_K\}$

for example $D = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} \text{line one} \\ \text{line two} \\ \text{line three} \end{bmatrix}$ (input data)

$\begin{bmatrix} \text{news} \\ \text{sports} \\ \text{religion} \end{bmatrix}$ (target / category label)

we want \Rightarrow given a document x find its category $p(c_k | x)$ probability of document x belonging to category c_k (eg news)

Naïve Bias gives \Rightarrow $\begin{cases} \text{① } P(c_k) & \text{probability of each categories (priors)} \\ \text{② } P(w_i | c_k) & \text{conditional probability of data } x \text{ w.r.t. category } c_k \end{cases}$

then we use Bayes's theorem to get $p(c^* | x)$.

$$p(c^* | x) = \frac{p(x | c^*) P(c^*)}{p(x)}$$

$$= P(c^*) \prod_{i=1}^n P(w_i | c^*)$$

where c^* corresponds to the $p(c_k | x)$ that has maximum probability.

$$c^* = \underset{c_k}{\operatorname{argmax}} p(c_k | x)$$

NB for this consideration when all words have equal probability

(A)

D : documents x_1, x_2, \dots, x_n
 $\{c_1, c_2, \dots, c_K\}$: classes c_1, c_2, \dots, c_K $\{w_1, w_2, \dots, w_V\}$
 V : vocabulary (unique words in D - stopwords - below threshold words)
 D_K : subset examples having class c_K
 n_K : # of words in D_K

1. for each category c_K :

2. D_K = subset with category c_K

3. $P(c_K) = \frac{|D_K|}{|D|}$ (prior)

4. n_K = # words in D_K

5. for each word $w_i \in V$:

6. let $n_{ki} = \# \text{ of } w_i \text{ in } D_K$

7. set $P(w_i | c_K) = \frac{n_{ki} + 1}{n_K + N}$ ← Laplace smoothing
 cond prob

8. Return priors $P(c_K)$
 and cond prob $P(w_i | c_K)$

(B)

Then we find the posterior probability for each class c_K ,

$$P(c_K | x) = \frac{P(c_K) \cdot P(x | c_K)}{P(x)}$$

(C)

$$\propto P(c_K) P(x | c_K)$$

$$\propto P(c_K) \cdot \prod_{i=1}^n P(w_i | c_K)$$

w_i are words in test document x

$$\propto \ln [P(c_K) \cdot \prod_{i=1}^n P(w_i | c_K)]$$

$$\propto \ln P(c_K) + \sum_{i=1}^n \ln (P(w_i | c_K))$$

(to prevent underflow)

Testing

$$c^* = \underset{c_K}{\operatorname{argmax}} P(c_K | x)$$

$$c^* = \underset{c_K}{\operatorname{argmax}} P(c_K) \prod_{i=1}^n P(w_i | c_K)$$

$$\propto \underset{c_K}{\operatorname{argmax}} [\ln P(c_K) + \sum_{i=1}^n \ln (P(w_i | c_K))]$$

(D)

Example of Naive Bias

train	doc	words	class	total
	1	chinese Beijing chinese	$c_1 = c$	$n_1 = 8$ words $d_1 = 3$ examples $n_1 + V = 8 + 6 = 14$
	2	chinese chinese shanghai	c_1	
	3	chinese macao	c_1	
	4	Tokyo Japan chinese	c_2	$n_2 = 3$ words $d_2 = 1$ example $n_2 + V = 3 + 6 = 9$ \Rightarrow total words 6 unique words
test	1	chinese chinese chinese Tokyo Japan (c_2)	?	$ V = 6$

	w_1	w_2	w_3	w_4	w_5	w_6	
vocabulary: {	chinese,	Beijing,	shanghai,	macao,	Tokyo,	Japan}	$ V =6$
word frequency:	5	1	1	1	0	0	
word freq in c_2 :	1	0	0	0	1	1	
Laplace smoothing:	1	1	1	1	1	1	
$\frac{n_{ki} + 1}{n_k + V }$	6, 2 14, 9	2, 2 14, 9	2, 2 14, 9	2, 2 14, 9	1, 2 14, 9	1, 2 14, 9	

category $c_2 = j$

category $c_1 = c$

- prior $P(c_1) = |D_1| / |D| = 3/4$ (3 documents belongs to class c_1)
- number of words in class, $n_1 = 8$
- conditional probabilities ($P(w_i | c_k) = \frac{n_{ki} + 1}{n_k + |V|}$)

chinese $P(w_1 | c_1) = \frac{5+1}{8+6} = \frac{6}{14} = \frac{3}{7}$ \Rightarrow Laplace

Beijing $P(w_2 | c_1) = \frac{1+1}{8+6} = \frac{2}{14} = \frac{1}{7}$

shanghai $P(w_3 | c_1) = \frac{1+1}{8+6} = \frac{2}{14} = \frac{1}{7}$

macao $P(w_4 | c_1) = \frac{1+1}{8+6} = \frac{2}{14} = \frac{1}{7}$

Tokyo $P(w_5 | c_1) = \frac{0+1}{8+6} = \frac{1}{14}$

Japan $P(w_6 | c_1) = \frac{0+1}{8+6} = \frac{1}{14}$

category $c_2 = j$

- prior $P(c_2) = |D_2| / |D| = 1/4$ (1 document belongs to class c_2)
- $n_2 = 3$
- conditional probabilities

$P(w_1 | c_2) = \frac{1+1}{3+6} = \frac{2}{9}$

$P(w_2 | c_2) = \frac{0+1}{3+6} = \frac{1}{9}$

$P(w_3 | c_2) = \frac{0+1}{3+6} = \frac{1}{9}$

$P(w_4 | c_2) = \frac{0+1}{3+6} = \frac{1}{9}$

$P(w_5 | c_2) = \frac{1+1}{3+6} = \frac{2}{9}$

$P(w_6 | c_2) = \frac{1+1}{3+6} = \frac{2}{9}$

③ choosing class

prior chinese Tokyo Japan

$P(c_1 | x) \propto P(c_1) \prod_{i \in c_1} P(w_i | c_1) = \frac{3}{4} \cdot \left(\frac{3}{7}\right)^3 \cdot \frac{1}{14} \cdot \frac{1}{14} = 0.0003$

$P(c_2 | x) \propto P(c_2) \prod_{i \in c_2} P(w_i | c_2) = \frac{1}{4} \cdot \left(\frac{2}{9}\right)^3 \cdot \frac{2}{9} \cdot \frac{2}{9} = 0.0001$

$c^* = \underset{c_k}{\operatorname{argmax}} P(c_k) \prod_{j=1}^n P(w_j | c_k) = \underset{c_k}{\operatorname{argmax}} \{0.0003, 0.0001\} = \text{Label of } 0.0003 = c_1$

P.T.O.

⑤ give more importance to title of doc

① for each category c_k :

② binclass prior $P(c_k) = |D_k| / |D|$ (where D_k is subset of D with class c_k)

③ # loop through title words
let n_k = number of title words in class c_k all examples

for each title word $w_i \in V_1$:

let n_{ki} = number of words in title of all class c_k examples

find conditional prob of each title words

$$P_1(w_i | c_k) = \frac{n_{ki} + 1}{n_k + |V_1|} \times K \leftarrow \text{multiplication factor for title words}$$

④ # loop through body words

let n_k = number of body words

for each body word $w_i \in V_2$:

let n_{ki} = number of words in body of all class c_k examples (for example c_k is class 1)

find conditional prob of each body words

$$P_2(w_i | c_k) = \frac{n_{ki} + 1}{n_k + |V_2|}$$

\Rightarrow output $P(c_k)$, $P_1(w_i | c_k)$, $P_2(w_i | c_k)$

choosing class:

we again suppose title words and body words are independent,

$$P(c_k | x) \propto \underbrace{P_1(x | c_k) P(c_k)}_{\text{for title}} \cdot \underbrace{P_2(x | c_k) P(c_k)}_{\text{for body}}$$