

Rm 329

Statistics 315a.

in sequoia

HomeWork 3.

Dong-Bang Tsai

1)

For least squares ~~projections~~ and cubic smoothing splines,

$$\hat{f} = X(X^T X + \lambda I)^{-1} X^T y \quad \text{define } A = (X^T X + \lambda I)^{-1}$$

$$\Rightarrow \hat{f} = X A^T y$$

$$\therefore S = X A^T X^T$$

$$\Rightarrow \hat{f}(x_i) = x_i^T (X^T X + \lambda I)^{-1} X^T y = x_i^T A^T X^T y, \quad A = X^T X + \lambda I$$

$$\text{Now } \hat{f}_i(x_i) = x_i^T (X_{-i}^T X_{-i} + \lambda I)^{-1} X_{-i}^T y_{-i}$$

where X_{-i} is all the training set but with x_i removed.

$$\begin{aligned} \therefore \hat{f}_i(x_i) &= x_i^T (X^T X - x_i x_i^T + \lambda I)^{-1} (X^T y - x_i y_i) \\ &= x_i^T (A - x_i x_i^T)^{-1} (X^T y - x_i y_i) \end{aligned}$$

Using Mathematica, we can write $(A - x_i x_i^T)^{-1} = A^{-1} + \frac{A^{-1} x_i x_i^T A^{-1}}{1 - x_i^T A^{-1} x_i}$
I'll prove it later as well.

$$\hat{f}^{-1}(x_i) = x_i^T \left(A^T + \frac{A^T x_i x_i^T A^T}{1 - x_i^T A^T x_i} \right) (x^T y - x_i^T y_i)$$

$$= \left[x_i^T A^{-1} + \frac{(x_i^T A^T x_i) x_i^T A^T}{1 - (x_i^T A^T x_i)} \right] (x^T y - x_i^T y_i) \quad \text{Using } x_i^T A^T x_i = s_{ii}$$

$$= \left[x_i^T A^{-1} + \frac{s_{ii} x_i^T A^T}{1 - s_{ii}} \right] (x^T y - x_i^T y_i)$$

$$= x_i^T A^T x^T y - x_i^T A^{-1} x_i^T y_i + \left(\frac{s_{ii} x_i^T A^T}{1 - s_{ii}} \right) x^T y - \left(\frac{s_{ii} x_i^T A^T}{1 - s_{ii}} \right) x_i^T y_i$$

$$= \hat{f}(x_i) - s_{ii} y_i + \frac{s_{ii} \hat{f}(x_i)}{1 - s_{ii}} - \frac{y_i s_{ii}^2}{1 - s_{ii}}$$

$$= \frac{\hat{f}(x_i) - s_{ii} y_i + s_{ii}^2 y_i - y_i s_{ii}^2}{1 - s_{ii}}$$

$$= \frac{\hat{f}(x_i) - s_{ii} y_i}{1 - s_{ii}}$$

$$= \frac{\hat{f}(x_i) - y_i}{1 - s_{ii}} + y_i$$

$$\Rightarrow y_i - \hat{f}(x_i) = \frac{y_i - \hat{f}(x_i)}{1 - s_{ii}}$$

Note that

$$(A - x_i x_i^T) \left(A^T + \frac{A^T x_i x_i^T A^T}{1 - x_i^T A^T x_i} \right) = I \quad \therefore (A - x_i x_i^T)^{-1} = A^T + \frac{A^T x_i x_i^T A^T}{1 - x_i^T A^T x_i}$$

⑥ Since $(x_i^T x_i + \lambda I)^{-1}$ is positive-semidefinite,

$$\begin{aligned} \therefore x_i^T (x_i^T x_i + \lambda I)^{-1} x_i &= \left(x_i^T A^{-1} + \frac{s_{ii} x_i^T A^{-1}}{1 - s_{ii}} \right) x_i \\ &= x_i^T A^{-1} x_i + \cancel{\frac{s_{ii} x_i^T A^{-1} x_i}{1 - s_{ii}}} \\ &= s_{ii} + \frac{s_{ii}}{1 - s_{ii}} \geq 0 \end{aligned}$$

~~By symmetry~~ $\Rightarrow 0 \leq s_{ii} < 1$

$$\therefore |y_i - \hat{f}^{-1}(x_i)| = \left| \frac{y_i - \hat{f}(x_i)}{1 - s_{ii}} \right| \geq |y_i - \hat{f}(x_i)|$$

⑦ If S depends only on the x_i , and produce \hat{f}^{-1} from y doesn't depend on y itself.

 ~~$f(x_i) = \sum_j s_{ij} y_j = \sum_j s_{ij} \hat{f}^{-1}(x_i)$~~

$$\Rightarrow \hat{f}^{-1}(x_i) = \sum_{j \neq i} s_{ij} y_j + s_{ii} \hat{f}^{-1}(x_i)$$

$$\text{Since } \hat{f}(x_i) = \sum_j s_{ij} y_j$$

$$\Rightarrow \hat{f}(x_i) - \hat{f}^{-1}(x_i) = s_{ii} (y_i - \hat{f}^{-1}(x_i))$$

$$\therefore y_i - \hat{f}(x_i) = (1 - s_{ii})y_i - (-s_{ii})\hat{f}^{-1}(x_i)$$

$$\Rightarrow y_i - f^{-1}(x_i) = \frac{y_i - \hat{f}(x_i)}{1 - s_{ii}}$$

We get the same result.

2).

ESL 18.8

a) As state in ESL 18.7, we have a linear regression problem where $p \gg n$. There are infinitely many least-squares solutions all with zero residuals.

Now, the problem here is we use

$f(x) = \alpha + \beta^T x$ to classify input x into +1 and -1,
 $\therefore \hat{Y} = X\hat{\beta}$, we know $\hat{\beta}$ in \mathbb{R}^p ,

Using the results from ESL 18.7,

we can conclude that there are infinitely many directions defined by $\hat{\beta}$ in \mathbb{R}^p onto which data project to exactly two points.

(b)

The distance of the projected points to the origin is

$$\frac{\hat{B}^T X}{\|\hat{B}\|}$$

1. ~~the~~ the distance between the projected points

$$\text{is } \frac{1}{\|\hat{B}\|} - \frac{1}{\|\hat{B}\|} = \frac{2}{\|\hat{B}\|}.$$

(c)

From 18.9 part (c), when ~~or~~ $\lambda = 0$,

the solution $\hat{B}_0 = V D^T U^T y$ has residuals all equal to zero, and is unique in that it has the smallest Euclidean norm. However, it means that this \hat{B}_0 will ~~not~~ maximize the distance between the projected points defined by

$$\frac{2}{\|\hat{B}\|}, \quad \therefore \text{The maximal data piling direction}$$

$$\text{can be given by } \hat{B}_0 = V D^T U^T y = x^T y.$$

ESL 18.9.

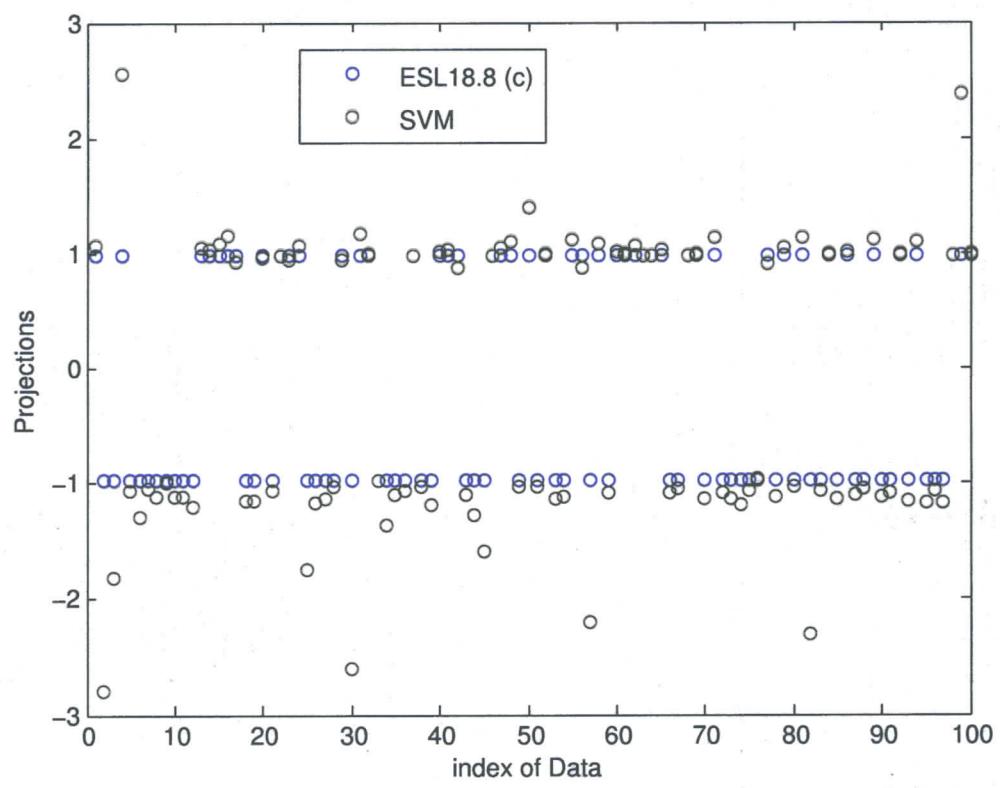
We want to compare the direction of the optimal separating hyperplane to ESL 18.8 (c) here.

Actually, I'll do a small experiment here,

$N=20$, $P=100$, randomly generate the sample and map y to -1 and 1.

We'll use 18.8 K1 to obtain the maximal data piling directions, and compare to this one obtained by solving optimal separation hyperplane using SVM.

We will find that the one obtained by SVM will have wider margin than maximal data piling method from ESL 18.8 (c). The code and figure are in the following pages.



```
1 % HW3 Q2, ESL 18.9
2 % Dong-Bang Tsai
3 clear;
4 N = 100;
5 p = 200;
6 X = randn(N, p);
7 Y = randn(N,1);
8 for i=1:N
9     if Y(i) > 0
10        Y(i) = 1;
11    else
12        Y(i) = -1;
13    end
14 end
15 [U,S,V] = svd(X, 'econ');
16 beta_p = (V/S)*U'*Y;
17 % The following is the projection distance obtained by ESL18.8 (c)
18 D_p = X*beta_p/sqrt(sum(beta_p.^2));
19
20 SVMStruct = svmtrain(X,Y);
21 Group = svmclassify(SVMStruct,X);
22 beta_s = SVMStruct.SupportVectors'*SVMStruct.Alpha;
23 % The following is the projection distance obtained by SVM
24 D_s = -X*beta_s/sqrt(sum(beta_s.^2));
25
26 x_indx = [1:N];
27 plot(x_indx, D_p, 'bo', x_indx, D_s, 'ro');
28 xlabel('index of Data'); ylabel('Projections');
29 legend('ESL18.8 (c)', 'SVM')
```

3). ESL 5.5

The source code is in the next page.

~~0.0530~~

For the raw data, there is no training error.

and the testing error will be 0.1583

Since in the raw data, we have 256 degree of freedom in each sample. therefore, in order to have 5 nodes. Different choice for the nodes, we use

$$\text{Knots} = [\underline{4}, \underline{8}, \underline{16}, \underline{32}, \underline{64}]$$

We find out that corresponding error will be

$$[\underline{0.0955}, \underline{0.0826}, \underline{0.0826}, \underline{0.0775}, \underline{0.0892}]$$

∴ We pick 32 knots in our case.

⇒ Full training error: ~~0.0530~~ 0.0530

Full testing error: 0.0831

As we expect, the test error decrease but training error increase when we apply smooth.

using cubic spline.

	Raw	filtered
Training error	0	0.0530
Testing error	0.1583	0.0831

```

1 % HW3 Q3, ESL 18.9
2 % Dong-Bang Tsai
3 clear;
4
5 % type of phonemes, 1)aa 2)ao 3)dcl 4)iy 5)sh
6 N_sample = 4509;
7 Y=zeros(N_sample,1);
8 Y_type = -1*ones(N_sample,1); % 0 is training set, 1 is testing set.
9 X=zeros(N_sample,256);
10
11 % This part of reading data takes most of mine time.
12 scanformat = '%d';
13 for i=1:256
14     scanformat = strcat(scanformat, ', %f');
15 end
16 scanformat = strcat(scanformat, ', %s');
17 fid = fopen('phoneme.data');
18 tline = fgetl(fid);
19 tline = fgetl(fid);
20 for i=1:N_sample
21     C = textscan(tline,scanformat);
22     for j=1:256
23         X(i,j) = C{j+1};
24     end
25     [tok,rem] = strtok(C{258},',');
26     if strcmp(tok, 'aa')
27         Y(i) = 1;
28     elseif strcmp(tok, 'ao')
29         Y(i) = 2;
30     elseif strcmp(tok,'dcl')
31         Y(i) = 3;
32     elseif strcmp(tok,'iy')
33         Y(i) = 4;
34     elseif strcmp(tok,'sh')
35         Y(i) = 5;
36     end
37     [tok,rem] = strtok(rem,',');
38     [tok,rem] = strtok(tok,'.');
39     if strcmp(tok,'train')
40         Y_type(i) = 0;
41     elseif strcmp(tok,'test')
42         Y_type(i) = 1;
43     end
44     tline = fgetl(fid);
45 end
46 fclose(fid);
47
48 indx_train = find(Y_type == 0);
49 Y_train = Y(indx_train);
50 X_train = X(indx_train,:);
51
52 indx_test = find(Y_type == 1);
53 Y_test = Y(indx_test);
54 X_test = X(indx_test,:);
55
56 Y_train_est = classify(X_train,X_train, Y_train,'quadratic');
57 misclassified_rate_train_raw = sum(Y_train_est ~= Y_train)/length(Y_train);
58
59 Y_test_est = classify(X_test, X_train, Y_train,'quadratic');
60 misclassified_rate_test_raw = sum(Y_test_est ~= Y_test)/length(Y_test);

```

```
61
62 knots = [50 100 150 200 250];
63 for i = 1:5
64     sp = spap2(knots(i),4,[1:256],X_train);
65     SpData{i} = sp.coefs;
66 end
67
68 error_knots = zeros(5,1);
69 for i=1:5
70     Y_sp = Y_train;
71     X_sp = SpData{i};
72     for j=1:10
73         Y_sp_test = Y_sp(334*j +1:334*(j+1));
74         Y_sp_train = Y_sp;
75         %Y_sp_train(334*j +1:334*(j+1)) = [];
76         %X_sp_test = X_sp(334*j +1:334*(j+1),:);
77         X_sp_train = X_sp;
78         %X_sp_train(334*j +1:334*(j+1),:) = [];
79         %Y_sp_test_est = classify(X_sp_test,X_sp_train, Y_sp_train,'quadratic');
80         %error_knots(i) = error_knots(i) + sum(Y_sp_test_est ~= Y_sp_test)/length
81         (Y_sp_test);
82     end
83 end
84
85
```

4) ESL 5.7

Integration by parts will give us

$$\int_a^b g''(x) h''(x) dx = [h'(x) g''(x)]_a^b - \int_a^b g'''(x) h'(x) dx$$

$$\left. \begin{array}{l} \text{ps: first term is zero,} \\ \text{and } g'''(x)=0 \text{ when } x \leq a \text{ or } x \geq b \end{array} \right\} = - \int_{x_1}^{x_n} g'''(x) h'(x) dx$$

$$= - \sum_{j=1}^{n-1} \int_{x_k}^{x_{j+1}} g'''(x) h'(x) dx$$

Since g is constructed by cubic spline, $\therefore g''(x) = \text{constant}$

$\therefore \int_a^b g''(x) h''(x) dx = - \sum_{j=1}^{n-1} g'''(x_j^+) \int_{x_k}^{x_{j+1}} h'(x) dx$ in the interval.

$$= - \sum_{j=1}^{n-1} g'''(x_j^+) [h(x_{j+1}) - h(x_j)]$$

Since it's interpolation, at node, $g(x_j) = \tilde{g}(x_j)$

$$\therefore h(x_j) = 0$$

$$\therefore \int_a^b g''(x) h''(x) dx = 0.$$

$$\textcircled{b} \quad \tilde{g} = h(x) + g(x)$$

$$\int_a^b \tilde{g}''(t)^2 dt = \int_a^b (h''(t) + g''(t))^2 dt.$$

$$= \int_a^b (h''(t))^2 + (g''(t))^2 + 2h''(t)g''(t) dt$$

$$= \int_a^b (h''(t))^2 + (g''(t))^2 dt$$

$$\geq \int_a^b g''(t)^2 dt.$$

\textcircled{c} If $\int_a^b h''(t)^2 dt = 0$, then the equality holds.

The condition will be

$$h''(t) = 0 \quad \forall t \in [a, b].$$

(c)

We want to find the function \hat{f} such that

$$\hat{f} = \arg \min_f \left\{ \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int_a^b f''(t)^2 dt \right\}$$

~~Now~~ Now, g is a natural cubic spline with knots at each of the x_i . Suppose that both g and \hat{f} interpolate the same sequence, we can write

$$\int_a^b \hat{f}''(t)^2 dt \geq \int_a^b g''(t)^2 dt \quad (\text{from } (b))$$

It imply that $h(t) = \hat{f}(t) - g(t) = 0 \forall t \in [a, b]$

since the loss term and penalty term with g will be never greater than the one with \hat{f} .

$$\Rightarrow \hat{f}(t) = g(t).$$

5) ESL 5.16.

From 5.45 of the textbook,

$$K(x, y) = \sum_{m=1}^M h_m(x) h_m(y) = \sum_{i=1}^{\infty} r_i \phi_i(x) \phi_i(y)$$

Multiply $\phi_j(x)$ on each side, and integrate it.

$$\sum_{m=1}^M \left(\int h_m(x) \phi_j(x) dx \right) h_m(y) = \sum_{i=1}^{\infty} r_i \left(\int \phi_i(x) \phi_j(x) dx \right) \phi_i(y)$$

Using the orthogonality of $\phi_i(x)$, i.e.,

$$\int \phi_i(x) \phi_k(x) dx = \begin{cases} 1, & \text{when } i=k \\ 0, & \text{when } i \neq k \end{cases}$$

$$\therefore \sum_{m=1}^M \left(\int h_m(x) \phi_j(x) dx \right) h_m(y) = \delta_j \phi_j(y)$$

Multiply with $\phi_k(y)$ and integrate again.

$$\sum_{m=1}^M \left[\int (h_m(x) \phi_j(x) dx) \right] \left[\int h_m(y) \phi_k(y) dy \right] = \delta_j \int \phi_j(y) \phi_k(y) dy$$

$$= r_j \delta_{j,k} \quad \text{where we apply orthogonality!}$$

matrix

Write down this equation in the ~~matrix~~ form.

then

$$\left(\int h(x) \phi(x) dx \right) \left(\int h(y) \phi(y) dy \right) = \text{diag}(r_1, r_2, \dots, r_m) = D_r.$$

define

~~$$V = \begin{pmatrix} \int h(x) \phi(x) dx \\ \vdots \\ \int h(x) \phi(x) dx \end{pmatrix}$$~~

$$V^T = D_r^{-\frac{1}{2}} \int h(x) \phi(x) dx - \mathbb{O}$$

$VTV = VV^T = I$, $\therefore V$ is an orthogonal matrix.

$$\therefore \left(\int h(x) \phi(x) dx \right) h(y) = D_r \phi(y)$$

and $\sum_{m=1}^M \left(\int h_m(x) \phi_j(x) dx \right) h_m(y) = \delta_j \phi_j(y)$

with ~~ϕ_j~~ eg ①.

$$D_r^{\frac{1}{2}} V^T h(y) = D_r \phi(y)$$

Using V is orthogonal, we can write.

$$h(x) = V D_r^{\frac{1}{2}} \phi(x)$$

from 5.63

$$\min_{\{\beta_m\}_1^M} \sum_{i=1}^N \left(y_i - \sum_{m=1}^M \beta_m h_m(x_i) \right)^2 + \gamma \sum_{m=1}^M \beta_m^2.$$

$$= \min_{\{\beta_m\}_1^M} \sum_{i=1}^N \|y_i - \beta^T h(x_i)\|^2 + \gamma \|\beta\|^2$$

$$= \min_{\beta} \sum_{i=1}^N \|y_i - \beta^T V D_r^{-\frac{1}{2}} \phi(x_i)\|^2 + \gamma \beta^T \beta \quad \left(\text{define } \beta = V D_r^{-\frac{1}{2}} c \right)$$

$$= \min_c \sum_{i=1}^N \|y_i - c^T \phi(x_i)\|^2 + \gamma (V D_r^{-\frac{1}{2}} c)^T (V D_r^{-\frac{1}{2}} c)$$

$$= \min_c \sum_{i=1}^N \|y_i - c^T \phi(x_i)\|^2 + \gamma c^T D_r^{-\frac{1}{2}} \underbrace{V^T V}_{I_r}^{-\frac{1}{2}} c$$

$$= \min_c \sum_{i=1}^N \|y_i - c^T \phi(x_i)\|^2 + \gamma c^T c D_r^{-1}$$

$$= \min_{\{c_j\}_1^\infty} \sum_{i=1}^N \left(y_i - \sum_{j=1}^\infty c_j \phi_j(x_i) \right)^2 + \gamma \sum_{j=1}^\infty \frac{c_j^2}{\sigma_j^2}$$

(b)

$$\therefore \hat{\beta} = (H^T H + \lambda I)^{-1} H^T y.$$

$$\therefore \hat{f} = H \hat{\beta} = I + (H^T H + \lambda I)^{-1} H^T y.$$

$$= \underbrace{(H^T)^{-1} H^T H^{-1}}_I + \underbrace{\lambda (H^T)^{-1} H^{-1}}_K y$$

$$= (I + \lambda K^{-1})^{-1} y.$$

$$\therefore \hat{f} = K(K + \lambda I)^{-1} y.$$

Note that $M < N$

$\therefore \cancel{K + \lambda I}$ are invertible

(c)

$$\text{Let } \hat{\alpha} = (K + \lambda I)^{-1} y.$$

$$\hat{f} = K \hat{\alpha}, \text{ from part (b)}$$

$$\therefore f(x) = h(x)^T \hat{\beta} = \sum_{i=1}^n k(x, x_i) \hat{\alpha}_i.$$

(d)

If $M < N$,

K is not invertible, but $\lambda \neq 0$, $K + \lambda I$ is invertible. i.e. The result still works

$$\hat{f} = H\hat{\beta} = K(K + \lambda I)^{-1}y$$

when $\lambda = 0$,

$$\hat{f} = H\hat{\beta} = H(H^T H)^{-1}H^T y = y$$