

1. Kernel Ridge Regression (20 points)

This problem examines the similarities and differences between KRR with and without offset.

- (a) (3 points) If we just center the training data in the original input space (by subtracting off the mean of the feature vectors and responses) and then apply KRR without offset, is that equivalent to KRR with offset? Explain.

1) a

a)

with offset

$$\tilde{K}_w = \langle \bar{\Phi}(x_i), \bar{\Phi}(x_j) \rangle - \frac{1}{n} \sum_{i=1}^n \langle \bar{\Phi}(x_i), \bar{\Phi}(x_i) \rangle - \frac{1}{n} \sum_{j=1}^n \langle \bar{\Phi}(x_j), \bar{\Phi}(x_j) \rangle + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \langle \bar{\Phi}(x_i), \bar{\Phi}(x_j) \rangle$$

$$k(x) = \langle \bar{\Phi}(x_i), \bar{\Phi}(x) \rangle - \frac{1}{n} \sum_{i=1}^n \langle \bar{\Phi}(x_i), \bar{\Phi}(x_i) \rangle - \frac{1}{n} \sum_{j=1}^n \langle \bar{\Phi}(x_j), \bar{\Phi}(x) \rangle + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \langle \bar{\Phi}(x_i), \bar{\Phi}(x_j) \rangle$$

without offset

$$\tilde{K}_{w0} = \langle \Phi(x_i - \bar{x}), \Phi(x_j - \bar{x}) \rangle \neq \tilde{K}_w$$

$$k_{w0}(x) = \langle \Phi(x_i - \bar{x}), \Phi(x - \bar{x}) \rangle \neq k_w(x)$$

(b) (3 points) In KRR with offset, give a formula for the offset b using the kernel.

$$\begin{aligned}
 1b) \quad b &= \bar{y} - \hat{w}^T \bar{x} = \bar{y} - \left\{ \frac{1}{\mu} \left[X^T - X(\tilde{K} + \mu I)^{-1} \tilde{K} \right] \tilde{y} \right\}^T \bar{x} \\
 &= \bar{y} - \left\{ \frac{1}{\mu} X \left[I - (\tilde{K} + \mu I)^{-1} \tilde{K} \right] \tilde{y} \right\}^T \bar{x} \\
 &= \bar{y} - \frac{1}{\mu} \tilde{y}^T \left[I - (\tilde{K} + \mu I)^{-1} \tilde{K} \right] X \bar{x}
 \end{aligned}$$

In above term:

$$X \bar{x} = \begin{bmatrix} \langle \tilde{x}_1, \bar{x} \rangle \\ \vdots \\ \langle \tilde{x}_n, \bar{x} \rangle \end{bmatrix} = \begin{bmatrix} \langle x_1 - \bar{x}, \bar{x} \rangle \\ \vdots \\ \langle x_n - \bar{x}, \bar{x} \rangle \end{bmatrix}$$

The entries of this vector are:

$$\langle x_i - \bar{x}, \bar{x} \rangle = \frac{1}{n} \sum_r \langle x_i, x_r \rangle - \frac{1}{n^2} \sum_r \sum_s \langle x_r, x_s \rangle$$

For term $\left[I - (\tilde{K} + \mu I)^{-1} \tilde{K} \right]^T$:

$$\begin{aligned}
 \left[I - (\tilde{K} + \mu I)^{-1} \tilde{K} \right]^T &= I - \tilde{K} (\tilde{K} + \mu I)^{-1} = [\tilde{K} + \mu I - \tilde{K}] (\tilde{K} + \mu I)^{-1} \\
 &= \mu (\tilde{K} + \mu I)^{-1} \Rightarrow b = \bar{y} - \tilde{y}^T (\tilde{K} + \mu I)^{-1} \left\{ \frac{1}{n} \sum_r \langle x_i, x_r \rangle - \frac{1}{n^2} \sum_r \sum_s \langle x_r, x_s \rangle \right\}
 \end{aligned}$$

To kernelize, substitute $\langle \Phi(x), \Phi(x') \rangle$ for $\langle x, x' \rangle$.

$$b = \bar{y} - \tilde{y}^T (\tilde{K} + \mu I)^{-1} \left\{ \frac{1}{n} \sum_r \langle \Phi(x_i), \Phi(x_r) \rangle - \frac{1}{n^2} \sum_r \sum_s \langle \Phi(x_r), \Phi(x_s) \rangle \right\}$$

- (c) (5 points) This problem shows you how to compute the KRR with offset solution without a bunch of loops. The train-train kernel matrix is the $n \times n$ matrix $K = [k(\mathbf{x}_i, \mathbf{x}_j)]$. The “centered” train-train kernel matrix \tilde{K} is the $n \times n$ matrix whose entries are $\langle \tilde{\Phi}(\mathbf{x}_i), \tilde{\Phi}(\mathbf{x}_j) \rangle$, where $\tilde{\Phi}(\mathbf{x}) := \Phi(\mathbf{x}) - \frac{1}{n} \sum_{\ell=1}^n \Phi(\mathbf{x}_\ell)$ and Φ is a feature map corresponding to the kernel k . This is the matrix that arises in KRR with offset. Calculation of \tilde{K} is facilitated by the formula

$$\tilde{K} = K - KO - OK + OKO, \quad (1)$$

where O is a square matrix with all entries equal to $1/n$. You should verify this fact but you don’t need to turn your work in.

Now consider a test data set $\mathbf{x}'_1, \dots, \mathbf{x}'_m$, and let K' be the $n \times m$ train-test matrix with entries $k(\mathbf{x}_i, \mathbf{x}'_j)$, and let \tilde{K}' be the $n \times m$ centered train-test matrix, whose entries are $\langle \tilde{\Phi}(\mathbf{x}_i), \tilde{\Phi}(\mathbf{x}'_j) \rangle$. Determine a formula analogous to (1) for relating \tilde{K}' to K' . Also, determine a formula for computing the predicted outputs on all the test points. Your formula should yield a column vector of length m .

1(c)
page 1

Train-test

$$\begin{aligned} \tilde{K}' &= \left[\langle \tilde{\Phi}(x_i), \tilde{\Phi}(x_j) \rangle \right]_{i,j=1}^{n,m} = \left[\langle \Phi(x_i) - \frac{1}{n} \sum_{\ell} \Phi(x_{\ell}), \Phi(x_j) \rangle \right]_{i,j=1}^{n,m} \\ &= \left[\langle \Phi(x_i), \Phi(x_j) \rangle \right]_{i,j=1}^{n,m} - \frac{1}{n} \left[\langle \Phi(x_i), \sum_{\ell} \Phi(x_{\ell}) \rangle \right]_{i,j=1}^{n,m} - \frac{1}{n} \left[\langle \sum_{\ell} \Phi(x_{\ell}), \Phi(x_j) \rangle \right]_{i,j=1}^{n,m} \\ &\quad + \frac{1}{n^2} \left[\langle \sum_{\ell} \Phi(x_{\ell}), \sum_{\ell} \Phi(x_{\ell}) \rangle \right]_{i,j=1}^{n,m} \end{aligned}$$

$$[\langle \Phi(x_i), \Phi(x_j) \rangle] = K'$$

$$\begin{aligned} \frac{1}{n} \left[\langle \Phi(x_i), \sum_{\ell} \Phi(x_{\ell}) \rangle \right]_{i,j=1}^{n,m} &= \frac{1}{n} \langle \Phi(x_i), \sum_{\ell} \Phi(x_{\ell}) \rangle [1 \dots 1]_{1 \times m} \\ &= \frac{1}{n} \langle \Phi(x_i), [\Phi(x_1) \dots \Phi(x_n)] \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} \rangle [1 \dots 1]_{1 \times m} \end{aligned}$$

$$= [\langle \Phi(x_i), \Phi(x_j) \rangle]_{i,j=1}^{n,m} \times \frac{1}{n} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} [1 \dots 1]_{1 \times m}$$

$$= [\langle \Phi(x_i), \Phi(x_j) \rangle]_{i,j=1}^{n,m} \times \begin{bmatrix} \frac{1}{n} & \dots & \frac{1}{n} \\ \vdots & & \vdots \\ \frac{1}{n} & \dots & \frac{1}{n} \end{bmatrix}_{n \times m} = K \mathbf{O}_{nm}$$

$$\frac{1}{n} \left[\langle \sum_{\ell} \Phi(x_{\ell}), \Phi(x_j) \rangle \right]_{i,j=1}^{n,m} = \frac{1}{n} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} \left[\langle \sum_{\ell} \Phi(x_{\ell}), \Phi(x_j) \rangle \right]_{i,j=1}^{n,m}$$

$$= \frac{1}{n} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} \langle [1 \dots 1]_{1 \times n} \begin{bmatrix} \Phi(x_1) \\ \vdots \\ \Phi(x_n) \end{bmatrix}, \Phi(x_j) \rangle$$

$$\begin{aligned} &= \frac{1}{n} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} [1 \dots 1]_{1 \times n} [\langle \Phi(x_i), \Phi(x_j) \rangle]_{i,j=1}^{n,m} = \begin{bmatrix} \frac{1}{n} & \dots & \frac{1}{n} \\ \vdots & & \vdots \\ \frac{1}{n} & \dots & \frac{1}{n} \end{bmatrix}_{n \times n} [\langle \Phi(x_i), \Phi(x_j) \rangle]_{i,j=1}^{n,m} \\ &= \mathbf{O} K' \end{aligned}$$

7(c)
page 2

$$\begin{aligned}
 \frac{1}{h^2} \left\langle \sum_{i=1}^n \bar{\Phi}(x_i), \sum_{j=1}^n \bar{\Phi}(x_j) \right\rangle_{i,j=1}^n &= \frac{1}{h^2} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} \left\langle \sum_{i=1}^n \bar{\Phi}(x_i), \sum_{j=1}^n \bar{\Phi}(x_j) \right\rangle_{i,j=1}^n \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}_{1 \times n} \\
 &= \frac{1}{h^2} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} \left\langle \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}_{1 \times n} \begin{bmatrix} \bar{\Phi}(x_1) \\ \vdots \\ \bar{\Phi}(x_n) \end{bmatrix}_{n \times 1}, \begin{bmatrix} \bar{\Phi}(x_1) & \dots & \bar{\Phi}(x_n) \end{bmatrix}_{1 \times n} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} \right\rangle \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}_{1 \times n} \\
 &= \frac{1}{h^2} \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix}_{n \times n} \left\langle \bar{\Phi}(x_i), \bar{\Phi}(x_j) \right\rangle_{i,j=1}^n \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix}_{n \times n} \\
 &= \begin{bmatrix} \frac{1}{h} & \dots & \frac{1}{h} \\ \vdots & & \vdots \\ \frac{1}{h} & \dots & \frac{1}{h} \end{bmatrix}_{n \times n} \left\langle \bar{\Phi}(x_i), \bar{\Phi}(x_j) \right\rangle_{i,j=1}^n \begin{bmatrix} \frac{1}{h} & \dots & \frac{1}{h} \\ \vdots & & \vdots \\ \frac{1}{h} & \dots & \frac{1}{h} \end{bmatrix}_{n \times n} = O(K) \sigma_{nm}
 \end{aligned}$$

$$\tilde{K}' = K' - K \sigma_{nm} - O(K') + O(K) \sigma_{nm}$$

Let's revisit the body fat data which we have seen earlier. We saw that a linear fit was reasonable, but now let's try a nonlinear fit.

Use the first 150 examples for training, and the remainder for estimating the mean squared error.

You will be asked to implement two variants of kernel ridge regression. For the next two problems please turn in

- the mean squared error on the training data
- the mean squared error on the test inputs
- the offset b in part (d) using your formula from (b)

(d) (3 points) Implement kernel ridge regression *with* offset. Remember to report the offset in addition to the other requested items.

MSE_training= 28.7940, MSE_test= 34.1488, b=1.0314

Code for KRR with offset:

```
close all; clear all; clc
load bodyfat_data.mat
n=150;
lambda=3e-3;
d=2;
x_train=X(1:n,:);
y_train=y(1:n);
x_test=X(n+1:end,:);
y_test=y(n+1:end);
m=numel(y_test);
xbar=sum(x_train,1)/n;
ybar=sum(y_train)/n;
xtilde=x_train-repmat(xbar,[n 1]);
ytilde=y_train-ybar;
O=1/n*ones(n);
for i=1:n
    for j=1:n
        K(i,j)=gaus_ker(x_train(i),x_train(j));
    end
end
Ktilde=K-K*O-O*K+O*K*O;
yhat_train=ybar+ytilde'*(Ktilde+n*lambda*eye(n))^-1*Ktilde;
yhat_train=yhat_train';
e_train=y_train-yhat_train;
MSE_train=sum(e_train.^2)/98
for i=1:n
    for j=1:m
        Kprime(i,j)=gaus_ker(x_train(i),x_test(j));
    end
end
Onm=1/n*ones(n,m);
Ktildeprime=Kprime-K*Onm-O*Kprime+O*K*Onm;
yhat_test=ybar+ytilde'*(Ktilde+n*lambda*eye(n))^-1*Ktildeprime;
yhat_test=yhat_test';
e_test=y_test-yhat_test;
```

```
MSE_test=sum(e_test.^2)/98
b=ybar-ybar-ytilde'*(Ktilde+n*lambda*eye(n))^-1*(1/n*sum(K,2)-
(1/n^2)*sum(K)*ones(n,1))
```

Code for calculation of Gaussian kernel:

```
function k=gaus_ker(u,v)
sigma=1.5;
k=exp(-1/(2*sigma^2)*norm(u-v)^2);
end
```

- (e) (3 points) Implement kernel ridge regression *without* offset. Comment on any differences in performance with respect to part (b). This problem should be easy once you have solved (d)

MSE_training=28.8242, MSE_test=73.9540

Code for KRR without offset:

```
close all; clear all; clc
load bodyfat_data.mat

n=150;
lambda=3e-3;
d=2;

x_train=X(1:n,:);
y_train=y(1:n);

xbar=sum(x_train,1)/n;
ybar=sum(y_train)/n;

xtilde=x_train-repmat(xbar,[n 1]);
ytilde=y_train-ybar;

%%%%%%%%%% Ktilde %%%%%%%%%%
for i=1:n
    for j=i:n
        sum2=0;
        for r=1:n
            sum2=sum2+gaus_ker(x_train(i),x_train(r));
        end
        sum3=0;
        for s=1:n
            sum3=sum3+gaus_ker(x_train(s),x_train(j));
        end
        sum45=0;
        for r=1:n
            for s=1:n
                sum45=sum45+gaus_ker(x_train(r),x_train(s));
            end
        end
    end
end
```

```

        Ktilde(i,j)=gaus_ker(x_train(i),x_train(j))-1/n*sum2-
1/n*sum3+1/n^2*sum45;
        Ktilde(j,i)=Ktilde(i,j);
    end
end
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

x_test=X(n+1:end,:);
y_test=y(n+1:end);

for i=1:numel(y_test)
    yhat_test(i)=ybar+ytilde'*(Ktilde+n*lambda*eye(n))^-1*ktilde(x_test(i));
end
e=y_test-yhat_test;
MSE=sum(e.^2)/98

```

2. Support Vector Regression (20 points)

Support vector regression (SVR) is a method for regression analogous to the support vector classifier. Let $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$, $i = 1, \dots, n$ be training data for a regression problem.

In the case of *linear regression*, SVR solves

$$\begin{aligned}
 \min_{\mathbf{w}, b, \xi^+, \xi^-} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-) \\
 \text{s.t.} \quad & y_i - \mathbf{w}^T \mathbf{x}_i - b \leq \epsilon + \xi_i^+ \quad \forall i \\
 & \mathbf{w}^T \mathbf{x}_i + b - y_i \leq \epsilon + \xi_i^- \quad \forall i \\
 & \xi_i^+ \geq 0 \quad \forall i \\
 & \xi_i^- \geq 0 \quad \forall i
 \end{aligned}$$

where $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$, $\xi^+ = (\xi_1^+, \dots, \xi_n^+)^T$, and $\xi^- = (\xi_1^-, \dots, \xi_n^-)^T$.

Here $\epsilon > 0$ is fixed.

a. (5 points) Show that for an appropriate choice of λ , SVR solves

$$\min_{\mathbf{w}, b} \quad \frac{1}{n} \sum_{i=1}^n \ell_\epsilon(y_i, \mathbf{w}^T \mathbf{x}_i + b) + \lambda \|\mathbf{w}\|_2^2$$

where $\ell_\epsilon(y, t) = \max\{0, |y - t| - \epsilon\}$ is the so-called ϵ -insensitive loss, which does not penalize prediction errors below a level of ϵ .

2a
(page 1)

$$\min_{w, b, \xi_i^+, \xi_i^-} \frac{1}{2} \|w\|_2^2 + \frac{c}{n} \sum_i (\xi_i^+ + \xi_i^-)$$

SvR \rightarrow

$$\text{s.t. } y_i - w^T x_i - b \leq \varepsilon + \xi_i^+ \quad \forall i$$

$$w^T x_i + b - y_i \leq \varepsilon + \xi_i^- \quad \forall i$$

$$\xi_i^+ \geq 0, \quad \xi_i^- \geq 0 \quad \forall i$$

Constraints:

$$y_i - w^T x_i - b \leq \varepsilon + \xi_i^+ \rightarrow \xi_i^+ \geq y_i - w^T x_i - b - \varepsilon$$

$$w^T x_i + b - y_i \leq \varepsilon + \xi_i^- \rightarrow \xi_i^- \geq w^T x_i + b - y_i - \varepsilon$$

$$\left. \begin{array}{l} \xi_i^+ \geq y_i - w^T x_i - b - \varepsilon \\ \xi_i^+ \geq 0 \end{array} \right\} \rightarrow \xi_i^+ = \max \{0, y_i - w^T x_i - b - \varepsilon\}$$

$$\left. \begin{array}{l} \xi_i^- \geq w^T x_i + b - y_i - \varepsilon \\ \xi_i^- \geq 0 \end{array} \right\} \rightarrow \xi_i^- = \max \{0, w^T x_i + b - y_i - \varepsilon\}$$

Constrained optimization problem transforms to unconstrained optimization problem:

$$\min_{w, b} \frac{1}{2} \|w\|_2^2 + \frac{c}{n} \sum_i \max \{0, y_i - w^T x_i - b - \varepsilon\} + \max \{0, w^T x_i + b - y_i - \varepsilon\}$$

① $\left. \begin{array}{l} y_i - w^T x_i - b - \varepsilon \leq 0 \\ w^T x_i + b - y_i - \varepsilon \leq 0 \end{array} \right\} \rightarrow |y_i - w^T x_i - b| - \varepsilon \leq 0$
For this case, $\sum (0 + 0) = 0$

② $\left. \begin{array}{l} y_i - w^T x_i - b - \varepsilon < 0 \\ w^T x_i + b - y_i - \varepsilon > 0 \end{array} \right\} \rightarrow y_i - w^T x_i - b < -\varepsilon \rightarrow |y_i - w^T x_i - b| > \varepsilon$
For this case, $\sum 0 + (w^T x_i + b - y_i - \varepsilon) = \sum |y_i - w^T x_i - b| - \varepsilon$

③ $\left. \begin{array}{l} y_i - w^T x_i - b - \varepsilon > 0 \\ w^T x_i + b - y_i - \varepsilon < 0 \end{array} \right\} \rightarrow y_i - w^T x_i - b > \varepsilon \rightarrow |y_i - w^T x_i - b| > \varepsilon$
For this case, $\sum (y_i - w^T x_i - b - \varepsilon) + 0 = \sum |y_i - w^T x_i - b| - \varepsilon$

2a
(page 2)

$$\textcircled{4} \begin{cases} y_i - w^T x_i - b - \epsilon > 0 \\ w^T x_i + b - y_i - \epsilon > 0 \end{cases} \rightarrow \begin{cases} y_i - w^T x_i - b > \epsilon \\ y_i - w^T x_i - b < -\epsilon \end{cases} \xrightarrow{\epsilon > 0} \text{impossible}$$

Based on cases ①, ..., ④:

$$\begin{aligned} & \sum_i \max \{0, y_i - w^T x_i - b - \epsilon\} + \max \{0, w^T x_i + b - y_i - \epsilon\} \\ &= \sum_i \max \{0, |y_i - w^T x_i - b| - \epsilon\} \end{aligned}$$

$$\begin{aligned} & \Rightarrow \frac{1}{2} \|w\|_2^2 + \frac{C}{n} \sum_i \max \{0, |y_i - w^T x_i - b| - \epsilon\} + \max \{0, w^T x_i + b - y_i - \epsilon\} \\ &= \frac{1}{2} \|w\|_2^2 + \frac{C}{n} \sum_i \max \{0, |y_i - w^T x_i - b| - \epsilon\} \\ &= \frac{1}{2} \|w\|_2^2 + \frac{C}{n} \sum_i \ell_\epsilon(y_i, w^T x_i + b) \end{aligned}$$

objective function can be divided by any constant value.

Divide objective function by C :

$$\underbrace{\left(\frac{1}{2C}\right)}_r \|w\|_2^2 + \frac{1}{n} \sum_i \ell_\epsilon(y_i, w^T x_i + b)$$

- b. (8 points) The optimization problem is convex with affine constraints, and therefore strong duality holds. Use the KKT conditions to derive the dual optimization problem in a manner analogous to the support vector classifier. As in the SVC, you should eliminate the dual variables corresponding to the constraints $\xi_i^+ \geq 0$, $\xi_i^- \geq 0$.

2b

page 1)

SV C

$$\min_{w, b, \xi} \frac{1}{2} \|w\|_2^2 + \frac{c}{n} \sum_{i=1}^n \xi_i$$

$$\text{s.t. } y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$

$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|_2^2 + \frac{c}{n} \sum_{i=1}^n \xi_i$$

$$- \sum \alpha_i [y_i (w^T x_i + b) - (1 - \xi_i)] - \sum \beta_i \xi_i$$

$$\max_{\alpha \geq 0, \beta \geq 0} \min_{w, b, \xi} L(w, b, \xi, \alpha, \beta)$$

$$\frac{\partial L}{\partial w} = w - \sum \alpha_i y_i x_i = 0$$

$$\frac{\partial L}{\partial b} = -\sum \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial \xi_i} = \frac{c}{n} - \alpha_i - \beta_i = 0$$

SV R

$$\min_{w, b, \xi^+, \xi^-} \frac{1}{2} \|w\|_2^2 + \frac{c}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-)$$

$$\text{s.t. } y_i - w^T x_i - b \leq \varepsilon + \xi_i^+ \quad \forall i$$

$$w^T x_i + b - y_i \leq \varepsilon + \xi_i^- \quad \forall i$$

$$\xi_i^+ \geq 0, \quad \xi_i^- \geq 0$$

$$L(w, b, \xi^+, \xi^-, \alpha, \beta, \gamma, \lambda) =$$

$$\frac{1}{2} \|w\|_2^2 + \frac{c}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-)$$

$$+ \sum \alpha_i (y_i - w^T x_i - b - \varepsilon - \xi_i^+)$$

$$+ \sum \beta_i (w^T x_i + b - y_i - \varepsilon - \xi_i^-)$$

$$- \sum \gamma_i \xi_i^+ - \sum \lambda_i \xi_i^-$$

$$\max_{\alpha \geq 0, \beta \geq 0, \gamma \geq 0, \lambda \geq 0} \min_{w, b, \xi^+, \xi^-} L(w, b, \xi^+, \xi^-, \alpha, \beta, \gamma, \lambda)$$

$$\frac{\partial L}{\partial w} = w - \sum \alpha_i x_i + \sum \beta_i x_i = 0$$

$$\frac{\partial L}{\partial b} = -\sum \alpha_i + \sum \beta_i = 0$$

$$\frac{\partial L}{\partial \xi_i^+} = \frac{c}{n} - \alpha_i - \gamma_i = 0$$

$$\frac{\partial L}{\partial \xi_i^-} = \frac{c}{n} - \beta_i - \lambda_i = 0$$

2 b
(page 2)

SVC

$$w = \sum \alpha_i y_i x_i$$

$$\sum \alpha_i y_i = 0$$

$$\frac{c}{n} = \alpha_i + \beta_i$$

$$\frac{1}{2} \left(\sum_i \alpha_i y_i x_i \right)^T \left(\sum_j \alpha_j y_j x_j \right) + \sum (\alpha_i + \beta_i) \xi_i$$

$$- \left(\sum_i \alpha_i y_i x_i \right)^T \left(\sum_i \alpha_i y_i x_i \right) - b \sum \alpha_i y_i$$

$$+ \sum \alpha_i - \sum \alpha_i \xi_i - \sum \beta_i \xi_i$$

$$= -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum \alpha_i$$

$$\max_{\alpha, \beta} -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum \alpha_i$$

$$\text{s.t. } \sum \alpha_i y_i = 0 \quad \forall i$$

$$\alpha_i + \beta_i = \frac{c}{n} \quad \forall i$$

$$\alpha_i \geq 0, \beta_i \geq 0$$

SVR

$$w = -\sum (\beta_i - \alpha_i) x_i$$

$$\sum \alpha_i = \sum \beta_i$$

$$\frac{c}{n} = \alpha_i + \alpha_i = \lambda_i + \beta_i$$

$$\frac{1}{2} \left(\sum (\beta_i - \alpha_i) x_i \right)^T \left(\sum (\beta_j - \alpha_j) x_j \right)$$

$$+ \frac{c}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-)$$

$$+ \sum \alpha_i y_i + \sum (\beta_i - \alpha_i) x_i^T \alpha_j x_j$$

$$= \sum \alpha_i b - \sum \alpha_i \xi_i - \sum \alpha_i \xi_i$$

$$- \sum_i \left(\sum_j (\beta_j - \alpha_j) x_j \right)^T \beta_i x_i + \sum \beta_i b$$

$$- \sum \beta_i y_i - \sum \beta_i \xi_i - \sum \beta_i \xi_i$$

$$- \sum \alpha_i \xi_i^+ - \sum \lambda_i \xi_i^-$$

$$= -\frac{1}{2} \sum_{i,j=1}^n (\beta_i - \alpha_i) (\beta_j - \alpha_j) x_i^T x_j - \sum (\beta_i - \alpha_i) y_i$$

$$\max_{\alpha, \beta, \lambda} -\frac{1}{2} \sum_{i,j} (\beta_i - \alpha_i) (\beta_j - \alpha_j) x_i^T x_j - \sum (\beta_i - \alpha_i) y_i$$

$$\text{s.t. } \sum \alpha_i - \sum \beta_i = 0 \quad \forall i$$

$$\alpha_i + \alpha_i = \lambda_i + \beta_i = \frac{c}{n} \quad \forall i$$

$$\alpha_i \geq 0, \beta_i \geq 0, \lambda_i \geq 0$$

2b

(page 3)

SVCeliminate β_i

$$\max_{\alpha} -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j X_i^T X_j + \sum \alpha_i$$

$$\text{s.t.} \sum \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq \frac{1}{n} \quad \forall i$$

$$\text{1st KKT: } W^* + \sum_{i=1}^n \alpha_i^* (-y_i X_i) = 0$$

$$W^* = \sum_{i=1}^n \alpha_i^* y_i X_i$$

$$\text{Complementary Slackness} \quad \alpha_i^* [1 - \xi_i^* - y_i (W^{*T} X_i + b^*)] = 0$$

$$\beta_i^* \xi_i^* = 0$$

$$\alpha_i^* < \frac{1}{n} \rightarrow \beta_i^* = 0 \rightarrow \xi_i^*$$

$$i: 0 < \alpha_i^* < \frac{1}{n}$$

$$\xi_i^* = 0$$

$$1 - \xi_i^* - y_i (W^{*T} X_i + b^*) = 0$$

$$\rightarrow b^* = y_i - W^{*T} X_i$$

SVReliminate γ_i & λ_i

$$\max_{\alpha, \beta} -\frac{1}{2} \sum_{i,j} (\beta_i - \alpha_i)(\beta_j - \alpha_j) X_i^T X_j - \sum_i (\beta_i - \alpha_i) y_i$$

$$\text{s.t.} \sum \alpha_i - \sum \beta_i = 0 \quad \forall i$$

$$0 \leq \alpha_i \leq \frac{1}{n} \quad \forall i$$

$$0 \leq \beta_i \leq \frac{1}{n} \quad \forall i$$

$$W^* + \sum_{i=1}^n \alpha_i^* (-X_i) + \sum_{i=1}^n \beta_i^* X_i = 0$$

$$W^* = \sum_{i=1}^n (\alpha_i^* - \beta_i^*) X_i$$

$$\alpha_i^* [y_i - W^{*T} X_i - b^* - \epsilon - \xi_i^*] = 0$$

$$\beta_i^* [W^{*T} X_i + b^* - y_i - \epsilon - \xi_i^*] = 0$$

$$\gamma_i^* \xi_i^{*+} = 0 \quad \lambda_i^* \xi_i^{*-} = 0$$

$$\alpha_i^* < \frac{1}{n} \rightarrow \gamma_i^* = 0 \rightarrow \xi_i^{*+} = 0$$

$$\beta_i^* < \frac{1}{n} \rightarrow \lambda_i^* = 0 \rightarrow \xi_i^{*-} = 0$$

$$i: 0 < \alpha_i^* < \frac{1}{n}$$

$$\xi_i^* = 0$$

$$y_i - W^{*T} X_i - b^* - \epsilon - \xi_i^* = 0$$

$$\rightarrow b^* = y_i - W^{*T} X_i - \epsilon$$

$$j: 0 < \beta_j^* < \frac{1}{n}$$

$$\xi_j^* = 0$$

$$W^{*T} X_j + b^* - y_j - \epsilon - \xi_j^* = 0$$

$$\rightarrow b^* = y_j - W^{*T} X_j + \epsilon$$

- c. (4 points) Explain how to kernelize SVR. Be sure to explain how to determine b^* .

~~2.1 Page 12~~ SVC SVR

2(c) $f(x) = \text{sign}\{w^* x + b^*\}$

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

$$b^* = y_i - w^{*T} x_i$$

$$i: (0 < \alpha_i^* < \frac{1}{n})$$

$$f(x) = \text{sign}\left\{\sum_{i=1}^n \alpha_i^* y_i k(x, x_i) + b^*\right\}$$

where

$$b^* = y_i - \sum_{j=1}^n \alpha_j^* y_j k(x_j, x_i)$$

$$i: (0 < \alpha_i^* < \frac{1}{n})$$

$f(x) = w^{*T} x + b^*$

$$w^* = \sum_{i=1}^n (\alpha_i^* - \beta_i^*) x_i$$

$$b^* = y_i - w^{*T} x_i - \varepsilon$$

$$= y_i - w^{*T} x_i + \varepsilon$$

$$i: (0 < \beta_i^* < \frac{1}{n})$$

$$f(x) = \sum_{i=1}^n (\alpha_i^* - \beta_i^*) k(x, x_i) + b^*$$

where

$$b^* = y_i - \sum_{j=1}^n (\alpha_j^* - \beta_j^*) k(x_j, x_i)$$

$$i: (0 < \alpha_i^* < \frac{1}{n})$$

$$b^* = y_i - \sum_{j=1}^n (\alpha_j^* - \beta_j^*) k(x_j, x_i)$$

$$i: (0 < \beta_i^* < \frac{1}{n})$$

- d. (3 points) Argue that the final predictor will only depend on a subset of training examples, and characterize those training examples.

2(d)

Final classifier only depends on samples for which $0 < \alpha_i^* < \frac{1}{n}$

$$b^* = y_i - \sum_{j=1}^n (\alpha_j^* - \beta_j^*) k(x_j, x_i)$$

$$i: (0 < \beta_i^* < \frac{1}{n})$$

Final predictor depends only on samples for which

$$0 < \alpha_i^* < \frac{1}{n}$$

$$0 < \beta_i^* < \frac{1}{n}$$