

1) Dual expression computation Ridge Regression
 Soln

Using the slack variable $\xi_i = y_i - w^T x_i$

First we note that for maximization

$$\Rightarrow \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \quad \text{which is subject}$$

$$\text{to } y_i (w^T x_i) \geq \xi_i \quad \forall i$$

\Rightarrow In this case we only have 2 opt variables
 w, ξ_i

\Rightarrow The optimization problem then becomes,

1) Dual expression computation Ridge Regression
 Soln

Using the slack variable $\xi_i = y_i - w^T x_i$

First we note that for maximization

$$\Rightarrow \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \quad \text{which is subject}$$

$$\text{to } y_i (w^T x_i) \geq \xi_i \quad \forall i$$

\Rightarrow In this case we only have 2 opt variables
 w, ξ_i

\Rightarrow The optimization problem then becomes,

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

Subject to

$$1 - \xi_i - y_i (w^T x_i) \leq 0 \quad \forall i$$

Using the slack variable

First we note that for maximization
 $\Rightarrow \frac{1}{2}(\|\omega\|^2 + C \sum_{i=1}^n \xi_i)$ which is subject

to $y_i (\omega^T x_i) \geq \xi_i \quad \forall i$

\Rightarrow In this case we only have 2 opt variables
 ω, ξ_i

\Rightarrow The optimization problem then becomes

$$\min_{\omega, \xi} \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i$$

Subject to

$$1 - \xi_i - y_i (\omega^T x_i) \leq 0 \quad \forall i$$

$$-\xi_i \leq 0 \quad \forall i$$

The lagrangian of the above becomes

$$L(\omega, \xi, \mu, \lambda) = \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \mu_i (1 - \xi_i - y_i (\omega^T x_i))$$

$$\Rightarrow \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i$$

to $y_i (\omega^T x_i) \geq \xi_i \quad \forall i$

\Rightarrow In this case we only have 2 opt variables
 ω, ξ_i

\Rightarrow The optimization problem then becomes

$$\min_{\omega, \xi} \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i$$

Subject to

$$1 - \xi_i - y_i (\omega^T x_i) \leq 0 \quad \forall i$$

$$-\xi_i \leq 0 \quad \forall i$$

The lagrangian of the above becomes

$$L(\omega, \xi, \mu, \lambda) = \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \mu_i (1 - \xi_i - y_i (\omega^T x_i)) - \sum_{i=1}^n \lambda_i \xi_i$$

We then take the derivative of L w.r.t ξ_i

\Rightarrow In the case we only
 w, ξ_i
 \Rightarrow The optimization problem then becomes

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

Subject to

$$1 - \xi_i - y_i (w^T x_i) \leq 0 \quad \forall i$$

$$-\xi_i \leq 0 \quad \forall i$$

The Lagrangian of the above becomes

$$L(w, \xi, \mu, \lambda) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \mu_i (1 - \xi_i - y_i (w^T x_i)) - \sum_{i=1}^n \lambda_i \xi_i$$

We then take the derivative of L w.r.t ξ_i to form the dual form, then

Subject to

$$1 - \xi_i - y_i (w^T x_i) \leq 0 \quad \forall i$$

$$-\xi_i \leq 0 \quad \forall i$$

The Lagrangian of the above becomes

$$L(w, \xi, \mu, \lambda) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \mu_i (1 - \xi_i - y_i (w^T x_i)) - \sum_{i=1}^n \lambda_i \xi_i$$

We then take the derivative of L w.r.t ξ_i to form the dual form, then

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \mu_i \lambda_i = C \quad \forall i$$

$$\Rightarrow 1 - \xi_i - y_i (w^T x_i) \leq 0 \quad \forall i$$

This implies that

$$\min_{w, q} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

Subject to

$$1 - \xi_i - y_i (w^T x_i) \leq 0 \quad \forall i$$

$$-\xi_i \leq 0 \quad \forall i$$

The Lagrangian of the above becomes

$$L(w, q, \mu, \lambda) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \mu_i (1 - \xi_i - y_i (w^T x_i))$$

$$- \sum_{i=1}^n \lambda_i \xi_i$$

We then take the derivative of L w.r.t ξ_i to form the dual form, then

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \mu_i \lambda_i = C \quad \forall i$$

$$\Rightarrow 1 - \xi_i - y_i (w^T x_i) \leq 0 \quad \forall i$$

\Rightarrow This implies that

$$\mu_i (1 - \xi_i - y_i (w^T x_i)) = 0 ; \lambda_i \xi_i = 0 \quad \forall i$$

\Rightarrow Thus for the complementary slackness condition we have

$$1 - y_i (w^T x_i) = 0$$

\Rightarrow Finding the value for all μ_i the dual function takes the form

$$Q(\mu, \lambda) = \inf_{w, q} L(w, q, \mu, \lambda)$$

but from the Lagrangian we had that the term

$$\xi_i (C - \mu_i - \lambda_i) \xi_i$$

Taking the infimum w.r.t ξ_i

→ Thus for the complementary slackness condition we have

$$1 - \xi_i (\omega^\top x_i) = 0$$

⇒ Finding the value for all μ_i the dual function takes the form

$$\mathcal{L}(\alpha, \lambda) = \inf_{\omega, \xi} L(\omega, \xi, \alpha, \lambda)$$

but from the Lagrangian we had that the term

$$\sum_i (c - \mu_i - \lambda_i) \xi_i$$

=) Finally taking the infimum w.r.t ξ_i and imposing $c - \mu_i - \lambda_i \geq 0$ yields the dual problem

$$\max_{\lambda} q(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n \lambda_i \alpha_i^\top y_j y_j^\top \alpha_j$$

with primal problem given by

takes the form $\max_{\omega, \xi} L(\omega, \xi, \alpha, \lambda)$ dual function

$$\mathcal{L}(\alpha, \lambda) = \inf_{\omega, \xi} L(\omega, \xi, \alpha, \lambda)$$

but from the Lagrangian we had that the term

$$\sum_i (c - \mu_i - \lambda_i) \xi_i$$

=) Finally taking the infimum w.r.t ξ_i and imposing $c - \mu_i - \lambda_i \geq 0$ yields the dual problem

$$\max_{\lambda} q(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n \lambda_i \alpha_i^\top y_j y_j^\top \alpha_j$$

with primal problem given by

$$\frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i$$

Subject to

$$y_i (\omega^\top x_i) \geq 1 - \xi_i$$

as $C \rightarrow \infty$

but from the Lagrangian we had that the term

$$\sum_i (c - \alpha_i - \lambda_i) \xi_i$$

=> Finally taking the infimum w.r.t ξ_i and imposing $c - \alpha_i - \lambda_i \geq 0$ yields the dual problem

$$\max_{\lambda} q(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j x_i^T x_j$$

with primal problem given by

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

Subject to

$$y_i (w^T x_i) \geq 1 - \xi_i \quad \text{as } C \rightarrow \infty$$

Dual expression

2) Dual expression computation for SVM

solve

using the slack variable

$$\xi_i = y_i - w^T x_i - b$$

For the minimization $\frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i$

Subject to

$$y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

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

Opt variables in the case w, b, ξ_i

2) Dual expression computation for SVM

Using the slack variable

$$\xi_i = y_i - w^T x_i - b$$

$$\text{For the minimization } \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i$$

Subject to

$$y_i (w^T x_i + b) \geq 1 - \xi_i, \quad \forall i$$

$$\text{with } \xi_i \geq 0 \quad \forall i$$

Opt variables in the case w, b, ξ_i

\Rightarrow Note that for all possible deviations for change in $w^T x$ and b we have,

$$w^T x + b = 1$$

$$w^T x + b = 0$$

$$w^T x + b = -1$$

\Rightarrow We then consider $w^T x + b = 0$, and let the

Subject to

$$y_i (w^T x_i + b) \geq 1 - \xi_i, \quad \forall i$$

$$\text{with } \xi_i \geq 0 \quad \forall i$$

Opt variables in the case w, b, ξ_i

\Rightarrow Note that for all possible deviations for change in $w^T x$ and b we have,

$$w^T x + b = 1$$

$$w^T x + b = 0$$

$$w^T x + b = -1$$

\Rightarrow We then consider $w^T x + b = 0$, and let the optimization problem be given as

$$\min_{w, b, \xi_i} \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i \quad \text{ch. note } w^T w = \|w\|^2$$

Subject to

\Rightarrow Note that for all possible deviations for change in $w^T x$ and b we have.

$$w^T x + b = 1$$

$$w^T x + b = 0$$

$$w^T x + b = -1$$

\Rightarrow We then consider $w^T x + b = 0$, and let the optimization problem be given as

$$\min_{w, b, \xi_i} \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i \quad \text{subject to } w^T w = 1 \quad \text{and} \quad \begin{aligned} & 1 - \xi_i - y_i (w^T x_i + b) \leq 0 \quad \forall i \\ & -\xi_i \leq 0 \quad \forall i \end{aligned}$$

Subject to

$$1 - \xi_i - y_i (w^T x_i + b) \leq 0 \quad \forall i$$

$$-\xi_i \leq 0 \quad \forall i$$

\Rightarrow Thus the Lagrangian becomes

\Rightarrow Note that for all possible deviations for change in $w^T x$ and b we have.

$$w^T x + b = 1$$

$$w^T x + b = 0$$

$$w^T x + b = -1$$

\Rightarrow We then consider $w^T x + b = 0$, and let the optimization problem be given as

$$\min_{w, b, \xi_i} \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i \quad \text{subject to } w^T w = 1 \quad \text{and} \quad \begin{aligned} & 1 - \xi_i - y_i (w^T x_i + b) \leq 0 \quad \forall i \\ & -\xi_i \leq 0 \quad \forall i \end{aligned}$$

Subject to

$$1 - \xi_i - y_i (w^T x_i + b) \leq 0 \quad \forall i$$

$$-\xi_i \leq 0 \quad \forall i$$

\Rightarrow Thus the Lagrangian becomes

$$L(w, b, \xi, \mu, \lambda) = \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i + \dots +$$

$$+ \sum_{i=1}^n \mu_i (1 - \xi_i - y_i (\omega^T x_i + b)) - \sum_{i=1}^n \lambda_i \xi_i$$

\Rightarrow Taking derivatives of L w.r.t b and ξ_i respectively we have

$$\nabla_b L = 0 \Rightarrow \omega^* = \sum_{i=1}^n \mu_i x_i y_i x_i$$

$$\Rightarrow \frac{\partial L}{\partial b} = 0 \Rightarrow \sum \mu_i x_i y_i = 0$$

$$\Rightarrow \frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \mu_i x_i \lambda_i = C \quad \forall i$$

$$+ \sum_{i=1}^n \mu_i (1 - \xi_i - y_i (\omega^T x_i + b)) - \sum_{i=1}^n \lambda_i \xi_i$$

\Rightarrow Taking derivatives of L w.r.t b and ξ_i respectively we have

$$\nabla_b L = 0 \Rightarrow \omega^* = \sum_{i=1}^n \mu_i x_i y_i x_i$$

$$\Rightarrow \frac{\partial L}{\partial b} = 0 \Rightarrow \sum \mu_i x_i y_i = 0$$

$$\Rightarrow \frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \mu_i x_i \lambda_i = C \quad \forall i$$

$$\Rightarrow 1 - \xi_i - y_i (\omega^T x_i + b) \leq 0 ; \quad \xi_i \geq 0 \quad \forall i$$

$$\Rightarrow \mu_i (1 - \xi_i - y_i (\omega^T x_i + b)) = 0 ; \quad \lambda_i \xi_i = 0 , \forall i$$

\Rightarrow For the complementary slackness condition, we have

$$\nabla_{\omega} L = 0 \Rightarrow \omega^* = \sum_{i=1}^n \mu_i x_i y_i z_i$$

$$\Rightarrow \frac{\partial L}{\partial b} = 0 \Rightarrow \sum \mu_i y_i z_i = 0$$

$$\Rightarrow \frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \mu_i z_i = c \quad \forall i$$

$$\Rightarrow 1 - \xi_i - y_i (\omega^T x_i + b) \leq 0 ; \quad \xi_i \geq 0 \quad \forall i$$

$$\Rightarrow \mu_i (1 - \xi_i - y_i (\omega^T x_i + b)) = 0 ; \quad \lambda_i \xi_i = 0 , \forall i$$

\Rightarrow For the complementary slackness condition, we have
that

$$1 - y_i (\omega^T x_i + b) = 0 \quad \text{---(1)}$$

$$b = y_i - \vec{x}_i^T \vec{\omega} , \quad i \text{ such that } 0 < \mu_i < c \quad \text{---(2)}$$

\Rightarrow Thus finding value for all μ_i , the dual function

$$\Rightarrow \frac{\partial L}{\partial b} = 0 \Rightarrow \sum \mu_i y_i z_i = 0$$

$$\Rightarrow \frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \mu_i z_i = c \quad \forall i$$

$$\Rightarrow 1 - \xi_i - y_i (\omega^T x_i + b) \leq 0 ; \quad \xi_i \geq 0 \quad \forall i$$

$$\Rightarrow \mu_i (1 - \xi_i - y_i (\omega^T x_i + b)) = 0 ; \quad \lambda_i \xi_i = 0 , \forall i$$

\Rightarrow For the complementary slackness condition, we have
that

$$1 - y_i (\omega^T x_i + b) = 0 \quad \text{---(1)}$$

$$b = y_i - \vec{x}_i^T \vec{\omega} , \quad i \text{ such that } 0 < \mu_i < c \quad \text{---(2)}$$

\Rightarrow Thus finding value for all μ_i , the dual function
becomes

$$L(\mu_i) = \inf_{w, b, \xi} L(w, b, \xi, \mu, \lambda)$$

$$\Rightarrow 1 - \xi_i - y_i (\omega^T x_i + b) \leq 0 ; \quad \xi_i \geq 0 \quad \forall i$$

$$\Rightarrow u_i (1 - \xi_i - y_i (\omega^T x_i + b)) = 0 ; \quad u_i \xi_i = 0 , \forall i$$

\Rightarrow For the complementary slackness condition, we have
that

$$1 - y_i (\omega^T x_i + b) = 0 \quad \text{--- (1)}$$

$$b = y_j - x_j^T \omega, \quad j \text{ such that } 0 < u_j < c \quad \text{--- (2)}$$

\Rightarrow Thus finding value for all u_i , the dual function becomes

$$Q(u) = \inf_{\omega, b, \xi} L(\omega, b, \xi, u, \lambda)$$

but from lagrangian we have the term

$$\sum_i (c - u_i - \lambda_i) \xi_i$$

\Rightarrow Now taking the infimum w.r.t ξ_i and
imposing $c = u_i + \lambda_i \quad \forall i$, the dual problem
becomes

$$\max_u Q(u) = \sum_{i=1}^n u_i - \frac{1}{2} \sum_{i,j=1}^n u_i u_j y_i y_j x_i^T x_j$$

Subject to

$$0 \leq u_i \leq C \quad i = 1, \dots, n \quad \sum_{i=1}^n y_i u_i = 0$$

and for the primal problem we have

$$\min_w L(w, b, \xi, u, \lambda)$$

\Rightarrow Now taking the infimum w.r.t q_i and
imposing $c = u_i - \lambda_i$ the dual problem
becomes

$$\max_u q(u) = \sum_{i=1}^n u_i - \frac{1}{2} \sum_{i,j=1}^n u_i u_j y_i y_j x_i^T x_j$$

Subject to

$$0 \leq u_i \leq C \quad i=1, \dots, n \quad \sum_{i=1}^n y_i u_i = 0$$

and for the primal problem we have

$$\text{minimize}_w \frac{1}{2} w^T w + c \sum_{i=1}^n q_i$$

Subject to

$$y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

as $C = N \rightarrow \infty$

How to obtain prediction for x ?

$$\max_u q(u) = \sum_{i=1}^n u_i - \frac{1}{2} \sum_{i,j=1}^n u_i u_j y_i y_j x_i^T x_j$$

Subject to

$$0 \leq u_i \leq C \quad i=1, \dots, n \quad \sum_{i=1}^n y_i u_i = 0$$

and for the primal problem we have

$$\text{minimize}_w \frac{1}{2} w^T w + c \sum_{i=1}^n q_i$$

Subject to

$$y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

as $C = N \rightarrow \infty$

(Q) How to obtain prediction for x ?

Soln

\Rightarrow We let K be a $n \times n$ matrix - with

$K_{ij} = k(x_i, x_j)$ such that K is positive
semi-definite if n and all data sets
 $\{x_1, \dots, x_n\}$

and for the primal problem we have

$$\text{minimize } \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i$$

Subject to

$$y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

as $C = N \rightarrow \infty$

3 (10) How to obtain prediction for \hat{x}

Solve

\Rightarrow We let K be a $n \times n$ matrix with
- $K_{ij} = K(x_i, x_j)$ such that K is positive
semi-definite $\forall n$ and all data sets
 $\{x_1, \dots, x_n\}$

\Rightarrow Thus, given any n , and feature vectors x_1, \dots, x_n
we will have that for all scalars c_1, \dots, c_n

$$\sum_{i,j=1}^n c_i c_j K(x_i, x_j) = 0 \quad \leftarrow (1)$$

minimize $\frac{1}{2} \sum_{i=1}^n c_i^2$

Subject to

$$y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

as $C = N \rightarrow \infty$

3 (10) How to obtain prediction for \hat{x}

Solve

\Rightarrow We let K be a $n \times n$ matrix with
- $K_{ij} = K(x_i, x_j)$ such that K is positive
semi-definite $\forall n$ and all data sets
 $\{x_1, \dots, x_n\}$

\Rightarrow Thus, given any n , and feature vectors x_1, \dots, x_n
we will have that for all scalars c_1, \dots, c_n

$$\sum_{i,j=1}^n c_i c_j K(x_i, x_j) = 0 \quad \leftarrow (1)$$

\Rightarrow Now for the kernel function given below

$$K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2 / 2\sigma^2)$$

$$\text{Subject to } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \epsilon_i \quad \text{as } \epsilon = N \rightarrow \infty$$

3 (Q) How to obtain prediction for \mathbf{x}_t
Solve

\Rightarrow We let K be a $n \times n$ matrix with
 $k_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ such that K is positive
 semi-definite $\forall i, j$ and all data sets
 $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$

\Rightarrow Thus, given any n , and feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$
 we will have that for all scalars c_1, \dots, c_n

$$\sum_{i,j=1}^n c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) = 0 \quad (1)$$

\Rightarrow Now for the Kernel function given below

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2)$$

Any data values \mathbf{x}_i and \mathbf{x}_j feed forward
 into the function (1) will yield the result of
 the predicted \mathbf{x}_t values

\Rightarrow This can be generally classified by the function
 function $f(\mathbf{x})$ for \mathbf{x}_t prediction as

$$\begin{aligned} f(\mathbf{x}) &= \sum_{i=1}^n y_i k(\mathbf{x}_i, \mathbf{x}_t) + b \\ &= \sum_{i=1}^n y_i e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_t\|^2}{2\sigma^2}} + b \end{aligned}$$

\Rightarrow This can be generally classified by the function
function $f(x)$ for x prediction as

$$f(x) = \sum_{i=1}^n w_i y_i k(x_i, x_j) + b$$
$$= \sum_{i=1}^n w_i y_i e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}} + b$$
$$=$$

- 4) Maximum likelihood Normal distribution
 \Rightarrow Compute MLE to the MLE estimation
(MLE) as $n \rightarrow \infty$

Solv
Given $x_1, x_2, x_3, \dots, x_n$ for estimation we know

that the p.d.f for Normal distribution is generally given by

$$= \sum_{i=1}^n w_i y_i e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}} + b$$
$$=$$

- 4) Maximum likelihood Normal distribution
 \Rightarrow Compute MLE to the MLE estimation
(MLE) as $n \rightarrow \infty$

Solv

Given $x_1, x_2, x_3, \dots, x_n$ for estimation we know

that the p.d.f for Normal distribution is generally given by

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

If we apply the equation to n -number of observations we have

$$\prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

4) Maximum likelihood Normal distribution
 \Rightarrow Compute MLE to the MLE estimation
 (MLE) as $n \rightarrow \infty$

Sols

Given $x_1, x_2, x_3, \dots, x_n$ for estimation we know

that the p.d.f. for normal distribution is generally given by

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

\Rightarrow If we apply the equation to n-number of observations we have.

$$f(x, \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

\Rightarrow Compute MLE to the MLE estimation
 (MLE) as $n \rightarrow \infty$

Sols

Given $x_1, x_2, x_3, \dots, x_n$ for estimation we know

that the p.d.f. for normal distribution is generally given by

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

\Rightarrow If we apply the equation to n-number of observations we have.

$$f(x, \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$\Rightarrow L(\mu, x_i, \sigma^2) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}\sigma^2} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

\Rightarrow multiplying this we obtain

$$\begin{aligned}\Rightarrow \frac{1}{\sigma\sqrt{2\pi}} \cdot \frac{1}{\sigma\sqrt{2\pi}} \cdots \frac{1}{\sigma\sqrt{2\pi}} &= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \\ \Rightarrow e^{-\frac{1}{2\sigma^2}(x_1-\mu)^2} \cdot e^{-\frac{1}{2\sigma^2}(x_2-\mu)^2} \cdots e^{-\frac{1}{2\sigma^2}(x_n-\mu)^2} &= e^{\frac{1}{2\sigma^2} + (\dots + e^{\frac{1}{2\sigma^2} \cdot e^{\frac{1}{2\sigma^2} \cdots e^{\frac{1}{2\sigma^2} \cdot e^{-\frac{1}{2\sigma^2}(x_i-\mu)^2}}})^n} \\ = e^{\frac{1}{2\sigma^2} \cdot \sum_{i=1}^n (x_i-\mu)^2} &= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \cdot e^{-\frac{1}{2\sigma^2}(x_i-\mu)^2}\end{aligned}$$

Taking the log of the above function we have

\Rightarrow multiplying this we obtain

$$\begin{aligned}\Rightarrow \frac{1}{\sigma\sqrt{2\pi}} \cdot \frac{1}{\sigma\sqrt{2\pi}} \cdots \frac{1}{\sigma\sqrt{2\pi}} &= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \\ \Rightarrow e^{-\frac{1}{2\sigma^2}(x_1-\mu)^2} \cdot e^{-\frac{1}{2\sigma^2}(x_2-\mu)^2} \cdots e^{-\frac{1}{2\sigma^2}(x_n-\mu)^2} &= e^{\frac{1}{2\sigma^2} + (\dots + e^{\frac{1}{2\sigma^2} \cdot e^{\frac{1}{2\sigma^2} \cdots e^{\frac{1}{2\sigma^2} \cdot e^{-\frac{1}{2\sigma^2}(x_i-\mu)^2}}})^n} \\ = e^{\frac{1}{2\sigma^2} \cdot \sum_{i=1}^n (x_i-\mu)^2} &= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \cdot e^{-\frac{1}{2\sigma^2}(x_i-\mu)^2}\end{aligned}$$

\Rightarrow Taking the log of the above function we have

$$\log\left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n + \log\left(e^{\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2}\right) =$$

$$\Rightarrow \frac{1}{\sigma\sqrt{2\pi}} \cdot \frac{1}{\sigma\sqrt{2\pi}} \cdots \frac{1}{\sigma\sqrt{2\pi}} = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n$$

$$\Rightarrow e^{\frac{-1}{2\sigma^2}(x_1-\mu)^2} \cdot e^{\frac{-1}{2\sigma^2}(x_2-\mu)^2} \cdots e^{\frac{-1}{2\sigma^2}(x_n-\mu)^2}$$

$$= e^{\frac{1}{2\sigma^2} \cdot \sum_{i=1}^n (x_i-\mu)^2} = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \cdot e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2}$$

\Rightarrow Taking the log of the above function we have

$$\log \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^n + \log \left(e^{\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2} \right)$$

$$\Rightarrow n \log \left(\frac{1}{\sigma\sqrt{2\pi}} + -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2 \right) = n \log(n) - n \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2$$

$$\Rightarrow e^{\frac{-1}{2\sigma^2}(x_1-\mu)^2} \cdot e^{\frac{-1}{2\sigma^2}(x_2-\mu)^2} \cdots e^{\frac{-1}{2\sigma^2}(x_n-\mu)^2}$$

$$= e^{\frac{1}{2\sigma^2} \cdot \sum_{i=1}^n (x_i-\mu)^2} = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \cdot e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2}$$

\Rightarrow Taking the log of the above function we have

$$\log \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^n + \log \left(e^{\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2} \right)$$

$$\Rightarrow n \log \left(\frac{1}{\sigma\sqrt{2\pi}} + -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2 \right) = n \log(n) - n \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2$$

$$\leq -n \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2$$

$$= \underline{\partial L} = \underline{\partial L} = -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2 \cdot 2 \cdot -1$$

$$= e^{-\frac{1}{2\sigma^2} \cdot \sum_{i=1}^n (x_i - \mu)^2} = \left(\frac{1}{\sigma\sqrt{2\pi}}\right) \cdot e$$

\Rightarrow Taking the log of the above function we have

$$\log\left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n + \log\left(e^{-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2}\right)$$

$$\Rightarrow n \log\left(\frac{1}{\sigma\sqrt{2\pi}} + -\frac{1}{2\sigma^2} \sum (x_i - \mu)^2\right) = n \log(n) - n \log(\sigma(2\pi)) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

$$\leq -n \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

$$= \frac{\partial L}{\partial \mu} = \frac{\partial L}{\partial \mu} = -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu) \cdot 2 \cdot -1$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

$$\log\left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n + \log\left(e^{-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2}\right)$$

$$\Rightarrow n \log\left(\frac{1}{\sigma\sqrt{2\pi}} + -\frac{1}{2\sigma^2} \sum (x_i - \mu)^2\right) = n \log(n) - n \log(\sigma(2\pi)) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

$$\leq -n \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

$$= \frac{\partial L}{\partial \mu} = \frac{\partial L}{\partial \mu} = -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu) \cdot 2 \cdot -1$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

$$\Rightarrow \frac{1}{\sigma^2} \sum x_i - \frac{1}{\sigma^2} \mu = 0 = \frac{1}{\sigma^2} \mu = \frac{1}{\sigma^2} \sum x_i$$

$$\frac{n\bar{x}}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^n x_i \Rightarrow n\bar{x} = \sigma^2 \frac{1}{\sigma^2} \sum_{i=1}^n x_i$$

$$\Rightarrow \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow \text{estimation}$$

S) Decision Trees

a) Discussion

In this problem we gonna be discussing our decision tree information gain based on Entropy which is used to determine or calculate the purity of a substance.

$$\Rightarrow \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow \text{estimation}$$

S) Decision Trees

a) Discussion

In this problem we gonna be discussing our decision tree information gain based on Entropy which is used to determine or calculate the purity of a substance.

In other words we going to use Information gain to help in getting and understanding its role in building decision tree by getting the total average for entropy based on each criterion split in the below example

5) Decision Trees

a) Discussion

In this problem we gonna be discussing our decision tree information gain based on Entropy which is used to determine of calculate the purity of a substance.

→ In other words we going to use Information gain to help in getting and understanding its role in building decision tree by getting the total average for entropy based on each specific split in the below example.

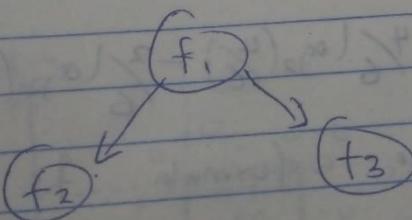
e.g., let's say we are splitting a function f_1 as shown below.

$$(f_1)$$

In this problem we gonna be discussing our decision tree information gain based on Entropy which is used to determine of calculate the purity of a substance.

→ In other words we going to use Information gain to help in getting and understanding its role in building decision tree by getting the total average for entropy based on each specific split in the below example.

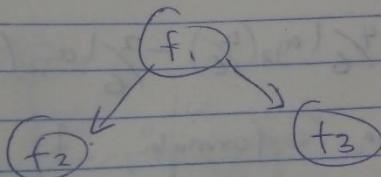
e.g., let's say we are splitting a function f_1 as shown below.



The function f_1 is given by

\Rightarrow In other words, we are going to use information gain to help in getting and understanding its role in building decision trees by getting the total usage for entropy based on each specific split in the below example.

Ques. Let's say we are splitting a function f_i as shown below:

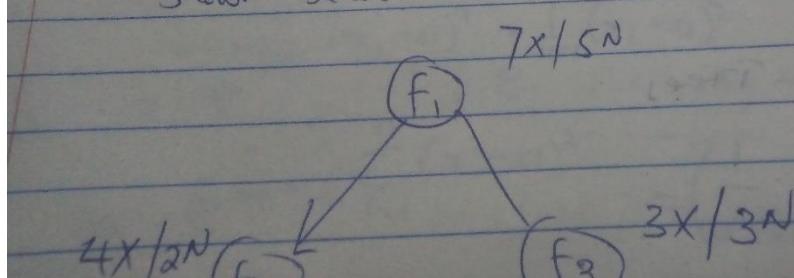


Note that Gain function is given by

$$\text{Gain}(S, A) = H(S) - \sum_{s \in \text{Val}} \frac{|S_s|}{|S|} H(s)$$

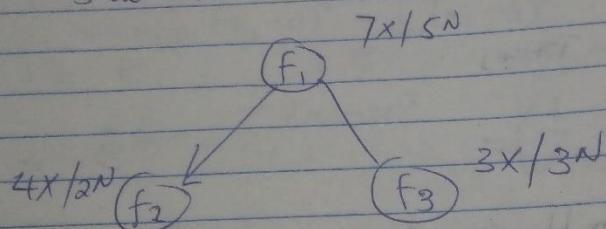
where S_s is the subset after splitting and S is the total subset.

\Rightarrow Now, let's say, we have to split f_1 and we are splitting it based on f_2 and f_3 as shown below:



whose S_V is the subset after splitting and
 S is the total sub set

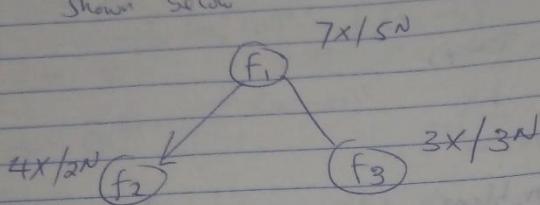
\Rightarrow Now, let's say we have to split (f_1) and we
are splitting it based on (f_2) and (f_3) as
shown below



Assume also we have a $7x/5N$
 $\Rightarrow 7$ -Values and $5N$ nodes total
and distributed as shown above

\Rightarrow Now, the first step is to compute $H(S)$ which
is basically for the subset $(f_1) = H(f_1) \approx H(S)$

\Rightarrow Now, let's say we have
are splitting it based on (f_2) and (f_3) as
shown below



Assume also we have a $7x/5N$
 $\Rightarrow 7$ -Values and $5N$ nodes total
and distributed as shown above

\Rightarrow Now, the first step is to compute $H(S)$ which
is basically for the subset $(f_1) = H(f_1) \approx H(S)$
which is the root node

\Rightarrow Then applying the formula

$$H(S) = -P_1 \log_2 P_1 - P_2 \log_2 P_2$$

$$4x/2^n \quad (f_2)$$

$$3x/3^n \quad (f_3)$$

Assume also we have $7x/5^n$
 $\Rightarrow 7$ -Values and 5 nodes total
 and distributed as shown above

\Rightarrow Now the first step is to compute $H(S)$ which
 is basically for the subset $(f_1) = H(f_1) \approx H(S)$
 which is the root node

\Rightarrow Then applying the formula

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \log_2 \left(\frac{2}{6}\right) \approx 0.87$$

\Rightarrow We use the same formula to obtain $H(f_2) = 0$
 and $H(f_3) = 4$

\Rightarrow MD is because we have values ranging
 from $0 \rightarrow 1$

$$\text{and distributed as}$$

\Rightarrow Now the first step is to compute $H(S)$ which
 is basically for the subset $(f_1) = H(f_1) \approx H(S)$
 which is the root node

\Rightarrow Then applying the formula

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{2}{6} \log_2 \left(\frac{2}{6}\right) \approx 0.87$$

\Rightarrow We use the same formula to obtain $H(f_2) = 0.73$
 and $H(f_3) = 4$

\Rightarrow MD is because we have values ranging
 from $0 \rightarrow 1$

$\Rightarrow 1 \Rightarrow$ for completely impurity

$0 \Rightarrow$ Full complete purity

\Rightarrow since we are computing $H(s)$ from root node, we have

$$\text{Gain}(s, f_1) = H(s) - \sum_{12} H(f_2) - \sum_{12} H(f_3)$$

$$= 0.87 - \sum_{12} 0.73 - \sum_{12} 1 = 0.005$$

\Rightarrow So basically if the above split gives the highest information gain value then the above structure will actually be used to split the particular dataset itself in order to construct the decision tree.

\Rightarrow since we are computing $H(s)$ from root node, we have

$$\text{Gain}(s, f_1) = H(s) - \sum_{12} H(f_2) - \sum_{12} H(f_3)$$

$$= 0.87 - \sum_{12} 0.73 - \sum_{12} 1 = 0.005$$

\Rightarrow So basically if the above split gives the highest information gain value then the above structure will actually be used to split the particular dataset itself in order to construct the decision tree.

\Rightarrow The higher the information gain the higher chance that splitting criterion will be used for constructing the decision tree.

\Rightarrow So in case we use conditional gain, then minimizing the conditional gain makes sense in building a decision tree.

$$Gain(s, f_1) = H(s) - \sum_{i=1}^2 H(f_{1i}) - \sum_{i=1}^2 H(f_{2i})$$

$$= 0.87 - \sum_{i=1}^2 \frac{6}{12} \times 0.73 - \sum_{i=1}^2 \frac{6}{12} \times 1 = 0.005$$

\Rightarrow So basically if the above split gives the highest information gain value then the above structure will actually be used to split the particular dataset itself in order to construct the decision tree.

\Rightarrow The higher the information gain the higher chance that splitting criterion will be used for constructing the decision tree.

\Rightarrow So in case we use conditional gain then minimizing the conditional gain makes sense to greedily build a decision tree.

b) With the given data points to draw the decision boundary we have that

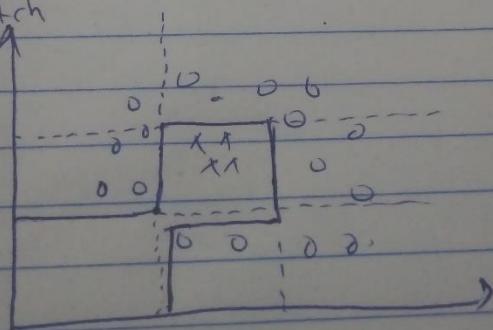
\Rightarrow the particular dataset itself in order to construct the decision tree.

The higher the information gain the higher chance that splitting criterion will be used for constructing the decision tree.

\Rightarrow So in case we use conditional gain then minimizing the conditional gain makes sense to greedily build a decision tree.

b) With the given data points to draw the decision boundary we have that

a sketch



= Understanding how decision tree works is that decision boundaries will always be perpendicular to X and Y axis as far as for that case we made the above estimation

=> So, basically the above sketch of the model can be said to be overfitted and this happens when we increase the depth of the tree.
This tells us that the model above is prone to overfitting.

= Understanding how decision tree works is that decision boundaries will always be perpendicular to X and Y axis as far as for that case we made the above estimation

=> So, basically the above sketch of the model can be said to be overfitted and this happens when we increase the depth of the tree.
=> This tells us that the model above is prone to overfitting

=> Thus, we see that if no limit on the depth of the algorithm, then this will be prone to overfitting implying that we are going to have a complex decision boundary accommodating

- => Understanding how decision tree works, we make the decision boundaries will always be perpendicular to X and Y axis as far as for that case we made the above estimation
- => So, basically, the above sketch of the model can be said to be overfitted and this happens when we increase the depth of the tree.
- => This tells us that the model above is prone to overfitting.
- => Thus, we see that if no limit on the depth of the algorithm, then this will be prone to overfitting implying that we going to have a complex decision boundary accommodating all similar data points.

D) (C) Increase the value of max depth may overfit the data

- => So, basically, the above sketch of the model can be said to be overfitted and this happens when we increase the depth of the tree.
- => This tells us that the model above is prone to overfitting.
- => Thus, we see that if no limit on the depth of the algorithm, then this will be prone to overfitting implying that we going to have a complex decision boundary accommodating all similar data points.

(C) (C) Increase the value of max depth may overfit the data

[END]