

CS761 Notes

Young Wu

December 14, 2018

1 Lecture 1

1.1 CLASSIFY: like GG or NOT

$$y = \begin{cases} 1 & \text{if like GG} \\ 0 & \text{if not} \end{cases}$$

x_1 = number of stars for SS

x_2 = number of stars for CRA

$$\hat{y} = \begin{cases} +1 & \text{if } \frac{\hat{\mathbb{P}}\{y = 1|x_1, x_2\}}{\hat{\mathbb{P}}\{y = -1|x_1, x_2\}} \geq 1 \\ -1 & \text{otherwise} \end{cases}$$

1.2 Random variables

X and Y are random variables

Joint prob distribution: $p(x, y) = \mathbb{P}\{X = x \text{ and } Y = y\}$

Marginal: $p(x) = \sum_{y \in \text{all possible } y \text{ values}} p(x, y)$

Conditional dist: $p(y|x) = \frac{p(x, y)}{p(x)}$

Bayes Rule: $p(y|x) = \frac{p(x|y)p(y)}{p(x)}$

Expectation: $\mu = \mathbb{E}[X] = \sum_x xp(x), \hat{\mu} = \frac{1}{40} \sum_{i=1}^{40} x_i$

$$\mathbb{E}[f(X)] = \sum_x f(x)p(x)$$

$$\mathbb{E}[X^2] = \sum_x x^2 p(x)$$

$$\begin{aligned} \mathbb{E}[X + Y] &= \sum_x \sum_y (x + y) p(x, y) \\ &= \sum_x x \sum_y p(x, y) + \sum_y y \sum_x p(x, y) \\ &= \mathbb{E}[X] + \mathbb{E}[Y] \end{aligned}$$

$$\mathbb{E}[XY] = \sum_x \sum_y xy p(x, y)$$

Independence iff $p(x, y) = p(x)p(y)$

If $X \perp Y$,

$$\begin{aligned}\mathbb{E}[XY] &= \sum_x \sum_y xyp(x)p(y) \\ &= \sum_x xp(x) \sum_y yp(y) \\ &= \mathbb{E}[X] \mathbb{E}[Y]\end{aligned}$$

Variance:

$$\begin{aligned}\mathbb{V}[X] &= \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= \sum_x (x - \mathbb{E}[X])^2 p(x)\end{aligned}$$

Conditional Expectation: $\mathbb{E}[Y|X = x] = \sum_y yp(y|x)$

Sums of Independent random variables

x_1, x_2, \dots are indep random variables

$S_n = \sum_{i=1}^n x_i$, what is distribution of S_n

$$\begin{aligned}\mathbb{E}[S_n] &= \sum_{i=1}^n \mathbb{E}[X_i] \\ \mathbb{V}[S_n] &= \mathbb{E}[(S_n - \mathbb{E}[S_n])^2] \\ &= \mathbb{E}\left[\left(\sum_{i=1}^n (X_i - \mathbb{E}[X_i])\right)^2\right] \\ &= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])] \\ &= \sum_{i=1}^n \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] \\ &= \sum_{i=1}^n \mathbb{V}[X_i]\end{aligned}$$

$$\mathbb{E}[X_i - \mathbb{E}[X_i]] = 0$$

1.3 Example

$p = \mathbb{P}\{\text{unif randomly chosen student SS} = 3, \text{CRA} = 4\}$

$$\begin{aligned}\hat{p} &= \frac{1}{40} \sum_{i=1}^{40} \mathbb{1}_{\text{ith person says SS} = 3, \text{CRA} = 4} \\ &= \frac{1}{40} \sum_{i=1}^{40} \mathbb{1}_{i,3,4}\end{aligned}$$

$$\mathbb{E}[\hat{p}] = \frac{1}{40} \sum_{i=1}^{40} \mathbb{E}[\mathbb{1}_{i,3,4}] = p$$

Unbiased:

$$\begin{aligned} \mathbb{V}[\hat{p}] &= \mathbb{E} \left[\left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{i,3,4} - p \right)^2 \right] \\ &= \frac{1}{n^2} \mathbb{E} \left[\sum_{i=1}^n (\mathbb{1}_{i,3,4} - p)^2 \right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}[(\mathbb{1}_{i,3,4} - p)^2] \end{aligned}$$

and,

$$\begin{aligned} \mathbb{E}[(\mathbb{1}_{i,3,4} - p)^2] &= (1-p)^2 p + p^2 (1-p) = p(1-p) \\ \mathbb{E}[\hat{p}] &= p \\ \mathbb{V}[\hat{p}] &= \frac{p(1-p)}{n} \\ \text{std}(\hat{p}) &= \sqrt{\frac{p(1-p)}{n}} \end{aligned}$$

2 Lecture 2

2.1 Discrete Random Variables

Y random variable taking values in $\{a_1, \dots, a_m\}$.

$$\begin{aligned} p_j &= \mathbb{P}\{Y = a_j\}, j = 1, \dots, m \\ \sum_{j=1}^m p_j &= 1 \end{aligned}$$

Bernoulli:

$$\begin{aligned} Y &\in \{0, 1\} \\ p &= \mathbb{P}\{Y = 1\} \\ \mathbb{P}\{Y = y\} &= p^y (1-p)^{1-y} \\ \mathbb{E}[Y] &= 1 \cdot p + 0 \cdot (1-p) \\ \mathbb{V}[Y] &= \mathbb{E}[(Y - p)^2] = p(1-p) \end{aligned}$$

Binomial:

$$\begin{aligned} Y_1, Y_2, \dots, Y_n &\stackrel{iid}{\sim} \text{Be}(p) \\ \mathbb{P}\{Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n\} &= \prod_{i=1}^n \mathbb{P}\{Y_i = y_i\} \end{aligned}$$

$$= \prod_{i=1}^n \prod_{j=1}^m p^{y_i} (1-p)^{1-y_i}$$

Binomial random variable with params n, p

$$K := \sum_{i=1}^n Y_i \sim \text{Bi}(n, p)$$

$$\mathbb{P}\{K = k\} = \binom{n}{k} p^k (1-p)^{n-k}$$

Bin coef: $\frac{n!}{k!(n-k)!}$
Multinomial

$$Y \in \{a_1, \dots, a_m\}, i = 1, \dots, n, \text{ indep}$$

$$\mathbb{P}\{Y_1 = y_1, \dots, Y_n = y_n\} = \prod_{i=1}^n \mathbb{P}\{Y_i = y_i\}$$

$$= \prod_{i=1}^n \prod_{j=1}^m p_j^{\mathbb{1}_{y_i=j}}$$

$$K_j = \{\text{number times } a_j \text{ appears in } Y_1, \dots, Y_n\}$$

$$\mathbb{P}\{K_1 = k_1, \dots, K_m = k_m\} = \underbrace{\binom{n}{k_1, \dots, k_m}}_{\text{multinomial coef}} \prod_{j=1}^m p_j^{k_j}$$

Poisson

$$X \geq 0 \text{ integer-valued}$$

$$\mathbb{P}\{X = k\} = e^{-\lambda} \frac{\lambda^k}{k!}, \lambda > 0 \text{ param}$$

$$\mathbb{E}[X] = \lambda$$

$$\mathbb{V}[X] = \lambda$$

2.2 Optimal Binary Classification

feature X , label $Y \in \{0, 1\}$

$$(X, Y) \sim \mathbb{P}_{XY}$$

$$\{(X_i, Y_i)\}_{i=1}^n \stackrel{iid}{\sim} \mathbb{P}_{XY}$$

Classifier

$$f : X \rightarrow Y$$

$$\hat{y} = f(X)$$

Loss: 0/1 loss

$$Loss(\hat{y}, y) = \mathbb{1}_{\hat{y} \neq y} = \begin{cases} 1 & \text{if } \hat{y} \neq y \\ 0 & \text{otherwise} \end{cases}$$

Risk: expected loss

$$R(f) = \mathbb{E}[\mathbb{1}_{f(X) \neq Y}] = \mathbb{P}\{f(X) \neq Y\}$$

Bayes Risk:

$$R^* = \inf_f R(f)$$

min probability of error

Bayes Optimal Classifier

$$f^*(x) = \mathbb{1}_{\underbrace{\mathbb{P}\{Y=1|X=x\}}_{\eta(x)} = \frac{1}{2}}$$

pick label that is most probable given x

$$R^* = \mathbb{E}[\min\{\eta(x), 1 - \eta(x)\}]$$

$$\hat{y} = 1 \rightarrow p(\text{err}) = 1 - \eta(x)$$

$$\hat{y} = 0 \rightarrow p(\text{err}) = \eta(x)$$

Theorem 1. $R(f^*) = R^*$

Estimating $\mathbb{P}\{f(X) \neq Y\} =: p_f$
 labeled examples $\{(X_i, Y_i)\}_{i=1}^n \stackrel{iid}{\sim} \mathbb{P}_{XY}$

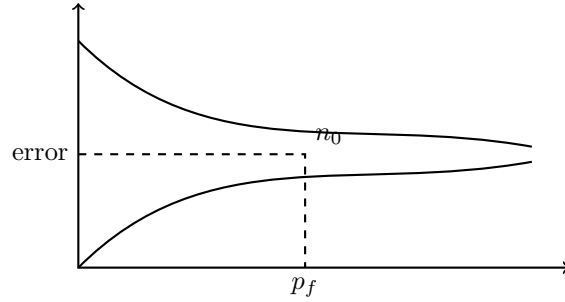
$$\hat{p}_f = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{1}_{f(X_i) \neq Y_i}}_{\text{iid Bernoulli random variable}}$$

$$\mathbb{E}[\hat{p}_f] = p_f$$

$$\mathbb{E}[|n\hat{p}_f - np_f|^2] = np_f(1 - p_f)$$

$$\mathbb{E}[|\hat{p}_f - p_f|^2] = \frac{p_f(1 - p_f)}{n} = \sigma^2$$

$$\sigma = \sqrt{\frac{p_f(1 - p_f)}{n}}$$



2.3 Analysis of Nearest Neighbor Classifier

Given $\{(X_i, Y_i)\}_{i=1}^n$

$$X_i \in \mathbb{R}^d$$

predict Y for new X

Theorem 2. (Cover and Hart 60's) $\mathbb{E}[R_n(X)] = \text{expected error of NN classifier}$

$$\lim_{n \rightarrow \infty} \mathbb{E}[R_n(x)] = \mathbb{E}[2\eta(x)(1 - \eta(x))] \leq 2R^*$$

as $n \rightarrow \infty$, NN of X , say $x' \in \{x_i\}_{i=1}^n$, $x \approx x'$

errs: $Y = 1, Y' = 0$ or $Y = 0, Y' = 1$

$$\eta(x)(1 - \eta(x')) + (1 - \eta(x))\eta(x') \approx 2\eta(x)(1 - \eta(x))$$

$$z := \min\{\eta(x), 1 - \eta(x)\}$$

$$\eta(x)(1 - \eta(x)) = z(1 - z)$$

$$\mathbb{E}[\eta(x)(1 - \eta(x))] = \mathbb{E}[z(1 - z)] = \mathbb{E}[z - z^2]$$

$$= \mathbb{E}[z] - \underbrace{\mathbb{E}[z^2]}$$

$\geq \mathbb{E}[z]^2$ Jensen's Inequality

$$\leq \mathbb{E}[z] - \mathbb{E}[z]^2$$

$$= \mathbb{E}[z](1 - \mathbb{E}[z])$$

$$\leq \mathbb{E}[z]$$

3 Lecture 3

3.1 Matrices

$$u_1, u_2, u_3$$

$$\|u_1\| = \|u_2\| = 1$$

$$u_i^T u_j = 0, i \neq j$$

$$U = [u_1 u_2 \dots u_m], u_i \in \mathbb{R}^n$$

Orthogonal matrix

$$U^T U = I_m$$

$$\begin{aligned} U^T U &= \begin{bmatrix} u_1^T \\ \vdots \\ u_m^T \end{bmatrix} \begin{bmatrix} u_1 & \dots & u_m \end{bmatrix} \begin{bmatrix} u_1^T u_1 & \dots & u_1^T u_m \\ \vdots & \ddots & \vdots \\ u_m^T u_1 & \dots & u_m^T u_m \end{bmatrix} \\ &= \begin{bmatrix} 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1 \end{bmatrix} \end{aligned}$$

$$\mathbb{R}^n = S \oplus S^\perp$$

$$\{\hat{u}_1, \hat{u}_2\} \in S, \{\hat{u}_3\} \in S^\perp$$

$$y = Ax, y \in \mathbb{R}^m, A \in M_{m \times n}, x \in \mathbb{R}^n$$

$$\begin{bmatrix} y_r \\ y_m \end{bmatrix} = \begin{bmatrix} \sigma_i & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_r \\ x_m \end{bmatrix}$$

$$y_r = \begin{bmatrix} y_1 \\ \vdots \\ y_r \end{bmatrix}$$

$$y_m = \begin{bmatrix} y_{r+1} \\ \vdots \\ y_m \end{bmatrix}, \text{ always zero}$$

$$\sigma_i = \begin{bmatrix} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_r \end{bmatrix}$$

$$x_r = \begin{bmatrix} x_1 \\ \vdots \\ x_r \end{bmatrix}$$

$$x_m = \begin{bmatrix} x_{r+1} \\ \vdots \\ x_n \end{bmatrix}, \text{ don't matter}$$

$$\mathbb{R}^m = \underbrace{N(A)}_{n-r} \oplus \underbrace{N(A)^\perp}_r$$

where $N(A)$ is the nullspace

$$\mathbb{R}^m = \underbrace{R(A)}_r \oplus \underbrace{R(A)^\perp}_{n-r}$$

3.2 SVD

$$\underbrace{A}_{m \times n} = \underbrace{U}_{m \times m} \underbrace{\Sigma}_{m \times n} \underbrace{V^T}_{n \times n}$$

$$U^T U = I$$

if U square, $U^T = U^{-1}$

$$\|Ux\|^2 = (Ux)^T (Ux) = x^T U^T U x = x^T x = \|x\|$$

$$(Ux)^T (Uy) = x^T U^T U y = x^T y$$

$$y = Ax = U \Sigma V^T x, V^T x = \begin{bmatrix} v_1^T x \\ \dots \\ v_n^T x \end{bmatrix}$$

$$V V^T x = (v_1 v_1^T + v_2 v_2^T + v_3 v_3^T) x = (v_1^T x) v_1 + (v_2^T x) v_2 + (v_3^T x) v_3$$

$$x = v_i$$

$$y = \sigma_i u_i$$

$$r = \text{rank}(A)$$

$$R(A) = \{u_1, \dots, u_r\}$$

$$R(A)^\perp = \{u_{r+1}, \dots, u_m\}$$

$$N(A) = \{v_{r+1}, \dots, v_n\}$$

$$N(A)^\perp = \{v_1, \dots, v_r\}$$

3.3 Matrix Identities

$$\begin{bmatrix} I & 0 \\ X & I \end{bmatrix}^{-1} = \begin{bmatrix} I & 0 \\ -X & I \end{bmatrix}$$

$$\begin{bmatrix} I & 0 \\ X & I \end{bmatrix} \begin{bmatrix} I & 0 \\ -X & I \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix}$$

$$\begin{bmatrix} X \\ Y \end{bmatrix} \sim \left(0, \begin{bmatrix} \Sigma_x & \Sigma_{x,y} \\ \Sigma_{y,x} & \Sigma_y \end{bmatrix} \right)$$

$$X|Y \sim (\dots, \Sigma_x - \Sigma_{x,y} \Sigma_y^{-1} \Sigma_{y,x})$$

$$\begin{aligned}
Ax_1 + Bx_2 &= y_1 \\
Cx_1 + Dx_2 &= y_2 \\
x_1 &= A^{-1}(y_1 - Bx_2) \\
CA^{-1}y_1 - CA^{-1}Bx_2 + Dx_2 &= y_2 \\
(D - CA^{-1}B)x_2 &= (y_2 - CA^{-1}y_1)
\end{aligned}$$

Matrix Iversion Lemma

$$(A - BD^{-1}C)^{-1} = A^{-1} + A^{-1}B(D - CA^{-1}B)^{-1}CA^{-1}$$

Other Identities:

$$\begin{aligned}
A(I + A)^{-1} &= I - (I + A)^{-1} \\
A &= (I + A) - I
\end{aligned}$$

3.4 Vector derivatives

$$\begin{aligned}
f(x) &= 0, f : \mathbb{R}^n \rightarrow \mathbb{R} \\
\frac{\partial f}{\partial x_1} &= 0 \text{ and } \frac{\partial f}{\partial x_2} = 0 \\
\frac{df}{dx} &= \text{gradient} = \nabla f \\
f &: \mathbb{R}^n \rightarrow \mathbb{R}^m \\
\frac{df}{dx} &\in \mathbb{R}^{m \times n} = \text{Jacobian} \\
c^T x &= c_1 x_1 + \dots + c_n x_n \\
\frac{dc^T x}{dx} &= \begin{bmatrix} c_1 \\ \dots \\ c_n \end{bmatrix} = c \\
\frac{dx^T Q x}{dx} &= (Q + Q^T)x
\end{aligned}$$

$$\begin{aligned}
\min_x \|Ax - b\|^2 \\
(Ax - b)^T (Ax - b) &= x^T (A^T A)x - 2b^T Ax + b^T b \\
\frac{d \text{ above}}{dx} &= 2A^T Ax - 2A^T b = 0 \\
A^T Ax - A^T b &= 0 \\
x_{opt} &= (A^T A)^{-1} A^T b
\end{aligned}$$

4 Lecture 4

4.1 Bayes Classifier

$$(X, Y) \sim \mathbb{P}_{XY}$$
$$\eta(x) := \mathbb{P}\{Y = 1 | X = x\}$$

$$f^*(x) = \begin{cases} 1 & \text{if } \eta(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbb{P}\{f^*(X) \neq Y\} = \mathbb{E}_x[\min(\eta(X), 1 - \eta(X))]$$
$$\{(X_i, Y_i)\}_{i=1}^n \stackrel{iid}{\sim} \mathbb{P}_{XY}$$

4.2 Nearest Neighbor Classifier

new unlabeled example x ,

$$i_{1nn}(x) = \arg \min_{i=1 \dots n} \|x - x_i\|$$
$$\hat{y} = y_{i_{1nn}(x)} \leftarrow f_{1nn}(x)$$

Theorem 3. *The following inequality holds,*

$$\lim_{n \rightarrow \infty} \mathbb{P}\{f_{1nn}(X) \neq Y\} \leq \mathbb{E}[2\eta(X)(1 - \eta(X))] \leq 2\mathbb{P}\{f^*(X) \neq Y\}$$

Theorem 4. *Let $R_k^\infty = \lim_{n \rightarrow \infty} \mathbb{P}\{f_{knn}(X) \neq Y\}$*

$$\mathbb{P}\{f^*(X) \neq Y\} \leq R_k^\infty \leq R_l^\infty \text{ for } l < k$$

4.3 KNN Classifier

Theorem 5. *(Stone 77) Let $k \rightarrow \infty$ and $\frac{k}{n} \rightarrow 0$ as $n \rightarrow \infty$, then,*

$$\lim_{n \rightarrow \infty} \mathbb{P}\{f_{1nn}(X) \neq Y\} = \mathbb{P}\{f^*(X) \neq Y\}$$

$$n \rightarrow k = \sqrt{n}$$

$$\mathbb{P}\{f_{1nn}(X) \neq Y\} - \mathbb{P}\{f^*(X) \neq Y\} \rightarrow 0 \text{ as } n \rightarrow \infty$$

4.4 Histogram Classifier

$$X_i \in [0, 1]^d, Y_i \in \{0, 1\}$$

$$d = 2$$

"bin", $\{B_i\}_{i=1}^M$ bins, $M = m^d$

$$\hat{p}_j = \frac{\sum_{i=1}^n \mathbb{1}_{x_i \in B_j, y_i = 1}}{\sum_{i=1}^n \mathbb{1}_{x_i \in B_j}}, j = 1, \dots, M$$

if $x \in B_j$ and $\hat{p}_j \geq \frac{1}{2}$ then label 1, otherwise 0.

$$\hat{\eta}_n(x) = \sum_{j=1}^M \hat{p}_j \mathbb{1}_{x \in B_j}$$

$$\hat{f}_n(x) = \begin{cases} 1 & \text{if } \hat{\eta}_n(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

Theorem 6. Let $M \rightarrow \infty$ and $\frac{n}{M} \rightarrow \infty$ as $n \rightarrow \infty$. Then,

$$\mathbb{P}\{\hat{f}_n(X) \neq Y\} \rightarrow \mathbb{P}\{f^*(X) \neq Y\}$$

Lemma 1. ,

$$\mathbb{P}\{\hat{f}_n(X) \neq Y\} - \mathbb{P}\{f^*(X) \neq Y\} \leq 2\mathbb{E}[|\hat{\eta}_n(x) - \eta(x)|]$$

$$p_j = \mathbb{P}\{Y = 1 | X \in B_j\}$$

$$p_j = \frac{\int_{B_j} \eta(x) p_x(x) dx}{\int_{B_j} p_x(x) dx}$$

$$\eta^{-x} = \sum_j^M p_j \mathbb{1}_{x \in B_j}$$

$$\mathbb{E}[|\eta(x) - \hat{\eta}_n(x)|] \leq \underbrace{\mathbb{E}[|\eta(x) - \eta^{-n(x)}|]}_{\text{Bias, } \rightarrow 0 \text{ as } M \rightarrow \infty} + \underbrace{\mathbb{E}[|\eta^{-x} - \hat{\eta}_n(x)|]}_{\text{Variance}}$$

5 Lecture 5

5.1 Bayes Classifier

$$\eta(x) = \mathbb{P}\{Y = 1 | X = x\}$$

$$f^{\star}(x) = \begin{cases} 1 & \text{if } \eta(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

Ideal Histogram Classifier

$$\begin{aligned} X &\in [0, 1]^d, Y = \{0, 1\} \\ M &= m^d \text{ "Bins" }, \{B_j\}_{j=1}^{\infty} \\ p_j &= \frac{\mathbb{E}[\mathbb{1}_{X \in B_j, Y=1}]}{\mathbb{E}[\mathbb{1}_{X \in B_j}]} \\ &= \frac{\int_{B_j} \eta(x) p(x) dx}{\int_{B_j} p(x) dx} \\ \bar{\eta}(x) &= \sum_{j=1}^M p_j \mathbb{1}_{x \in B_j} \end{aligned}$$

$$\bar{f}(x) = \begin{cases} 1 & \text{if } \bar{\eta}(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

Empirical Histogram Classifier

$$\begin{aligned} \hat{p}_j &= \frac{\sum_{i=1}^n \mathbb{1}_{x_i \in B_j, y_i=1}}{\sum_{i=1}^n \mathbb{1}_{x_i \in B_j}} \\ \{(x_i, y_j)\} &\stackrel{iid}{\sim} P_{xy} \\ \hat{\eta}(x) &= \sum_{j=1}^M \hat{p}_j \mathbb{1}_{x \in B_j} \end{aligned}$$

$$\hat{f}(x) = \begin{cases} 1 & \text{if } \hat{\eta}(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

A1 $p(x) \geq c > 0 \forall x$

A2 η is uniformly continuous

$$p(x) \geq c > 0 \text{ and } \frac{n}{M} \rightarrow \infty \Rightarrow N(x) \xrightarrow{as} \infty$$

Theorem 7. If $M \rightarrow \infty, \frac{n}{M} \rightarrow \infty$ as $n \rightarrow \infty$

$$\underbrace{\mathbb{P}\{\hat{f}(X) \neq Y\}}_{\mathbb{E}[\mathbb{1}_{\hat{f}(X) \neq Y}]} \rightarrow \underbrace{\mathbb{P}\{f^{\star}(X) \neq Y\}}_{\mathbb{E}[\mathbb{1}_{f^{\star}(X) \neq Y}]}$$

Proof. $\mathbb{P}\{\hat{f}(X) \neq Y\} - \mathbb{P}\{f^*(X) \neq Y\} \leq 2\mathbb{E}[|\eta(X) - \hat{\eta}(X)|]$

$$\begin{aligned}\mathbb{E}[|\eta(X) - \hat{\eta}(X)|] &= \mathbb{E}[|\eta(X) - \bar{\eta}(X) + \bar{\eta}(X) - \hat{\eta}(X)|] \\ &\leq \underbrace{\mathbb{E}[|\eta(X) - \bar{\eta}(X)|]}_{\text{deterministic error}} + \underbrace{\mathbb{E}[|\bar{\eta}(X) - \hat{\eta}(X)|]}_{\text{stochastic error} \rightarrow 0, as n \rightarrow \infty}\end{aligned}$$

where,

$$\begin{aligned}\mathbb{E}[|\eta(X) - \bar{\eta}(X)|] &= \int |\eta(x) - \bar{\eta}(x)| p(X) dx \\ &= \sum_{j=1}^M \int_{B_j} \underbrace{|\eta(x) - \bar{\eta}(x)|}_{\leq \varepsilon_m \rightarrow 0} p(x) dx \\ &\leq \varepsilon_m\end{aligned}$$

and,

$$\mathbb{E}[|\bar{\eta}(X) - \hat{\eta}(X)|] = \mathbb{E}\left[\mathbb{E}\left[\left|\bar{\eta}(x) - \frac{K(x)}{N(x)}\right| \middle| X = x\right]\right]$$

where,

x let $B(x)$ be its bin

$$\begin{aligned}K(x) &= \sum_{i=1}^n \mathbb{1}_{x_i \in B(x), y_i=1} \\ N(x) &= \sum_{i=1}^n \mathbb{1}_{x_i \in B(x)}\end{aligned}$$

$$K(x) | N(x) = n_x \sim \text{Binomial}(n_x, \bar{\eta}(x))$$

$$\begin{aligned}\mathbb{E}[K(x) | N(x) = n_x] &= n_x \bar{\eta}(x) \\ \mathbb{E}\left[\frac{K(x)}{N(x)} \middle| N(x) = n_x\right] &= \frac{n_x \bar{\eta}(x)}{n_x} = \bar{\eta}(x) \\ \mathbb{E}\left[\left(\bar{\eta}(x) - \frac{K(x)}{N(x)}\right)^2 \middle| N(x) = n_x\right] &= \mathbb{E}\left[\frac{1}{(n_x)^2} (n_x \bar{\eta}(x) - K(x))^2 \middle| N(x) = n_x\right] \\ &= \frac{1}{(n_x)^2} \mathbb{E}\left[(n_x \bar{\eta}(x) - K(x))^2 \middle| X = x\right] \\ &= \frac{1}{(n_x)^2} (n_x \bar{\eta}(x) (1 - \bar{\eta}(x))) \\ &= \frac{\bar{\eta}(x) (1 - \bar{\eta}(x))}{n_x}\end{aligned}$$

Recall, $\mathbb{E}[Z^2] \geq (\mathbb{E}[Z])^2$ Jensen's

$$\mathbb{E}[|\bar{\eta}(X) - \hat{\eta}(X)| | X = x] \leq \sqrt{\frac{\bar{\eta}(x) (1 - \bar{\eta}(x))}{n_x}}$$

□

$$N(x) \propto \frac{n}{M}$$

$$\mathbb{E} [|\bar{\eta}(X) - \hat{\eta}(X)|] = O\left(\sqrt{\frac{M}{n}}\right)$$

A3: η is 1-Lipschitz

$$|\eta(x) - \eta(x')| \leq \|x - x'\| \quad \forall x, x' \in [0, 1]^d$$

$$\text{if } x, x' \in B_j, \|x - x'\| \leq \frac{\sqrt{d}}{m} =: \varepsilon_m$$

$$\begin{aligned} \mathbb{P}\{\hat{f}(X) \neq Y\} - \mathbb{P}\{f^*(X) \neq Y\} &= O\left\{\max\left(\sqrt{\frac{m^d}{n}}, \frac{\sqrt{d}}{m}\right)\right\} \\ \sqrt{\frac{m^d}{n}} &= \frac{\sqrt{d}}{m} \\ \Rightarrow m^{d+2} &= dn \\ \Rightarrow m &= (dn)^{\frac{1}{d+2}} \\ \mathbb{P}\{\hat{f}(X) \neq Y\} - \mathbb{P}\{f^*(X) \neq Y\} &= O\left(\frac{\sqrt{d}}{(dn)^{\frac{1}{d+2}}}\right) \\ &= O\left(\frac{1}{n^{\frac{1}{d+2}}}\right) \\ &= O\left(n^{-\frac{1}{d+2}}\right) \end{aligned}$$

Curse of dim

5.2 Multivariate Normal Distribution

$$X \in \mathbb{R}^d$$

$$p(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$$\mu \in \mathbb{R}^d, \mathbb{E}[X] = \mu, \Sigma = \mathbb{E}\left[(X - \mu)(X - \mu)^T\right] \in \mathbb{R}^{d \times d}$$

$$\Sigma_{ij} = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)]$$

Covariation of X_i, X_j

$$X \sim N(\mu, \Sigma)$$

If $x \sim N(\mu, \Sigma)$, then $Ax + b \sim N(A\mu + b, A\Sigma A^T)$

Bivariate: $X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \mathbb{E}[X] = \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

$$\Sigma = \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

MVN Classifier

Suppose $X|Y = l \sim \underbrace{N(\mu_l, \Sigma_l)}_{\text{class-conditional distribution of } x}, l = 0, 1, \dots, k-1$

$$\mathbb{P}\{Y = l\} = \frac{1}{K}$$

$$\mathbb{P}\{X|Y = l\} \mathbb{P}\{Y = l\}$$

Bayes Classifier

$$f^*(x) = \arg \max_l p(x|y = l) \mathbb{P}\{Y = l\}$$

$$K = 2$$

$$f^*(x) = \begin{cases} 1 & \text{if } \underbrace{\log \left(\frac{p(x|y=1)}{p(x|y=0)} \right)}_{\text{log likelihood ratio}} \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Log LR

$$\begin{aligned} & \log \left(\frac{\frac{1}{\sqrt{(2\pi)^d |\Sigma_1|}} \exp \left(-\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right)}{\frac{1}{\sqrt{(2\pi)^d |\Sigma_0|}} \exp \left(-\frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) \right)} \right) \\ &= \underbrace{-\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0)}_{\text{quadratic in } x} + \text{const} \end{aligned}$$

If $\Sigma_0 = \Sigma_1 = \Sigma$,

$$\hat{y}(x) = \begin{cases} 1 & \text{if } \underbrace{2(\mu_1 - \mu_0)^T \Sigma^{-1} x}_w \geq \underbrace{\mu_0^T \Sigma \mu_0 - \mu_1^T \Sigma \mu_1}_t \\ 0 & \text{otherwise} \end{cases}$$

Linear classifier $w^T x > t$

$$\begin{aligned} \{x_i, y_i\}_{i=1}^n \\ \hat{\mu}_1 &= \frac{1}{n_1} \sum_{i \in Y_i=1} x_i, n_1 = \text{number of } Y_i = 1 \\ \hat{\Sigma}_1 &= \frac{1}{n_1} \sum_{i \in Y_i=1} (x_i - \hat{\mu}_1)(x_i - \hat{\mu}_1)^T \end{aligned}$$

6 Lecture 6

6.1 Generative Models

$Y = 0$ or $Y = 1$

$$\mathbb{P}\{Y = 0\} = 1 - \mathbb{P}\{Y = 1\}$$

$\mathbb{P}\{X|Y = 0\}, P\{X|Y = 1\}$, Class-conditional distributions

$$X|Y = 0 \sim N(\mu_0, \Sigma_0)$$

$$X|Y = 1 \sim N(\mu_1, \Sigma_1)$$

Region $\hat{y} = 0$ is R_0 , region $\hat{y} = 1$ is $R_1, x \in X$

$$\begin{aligned} X &= R_0 \cup R_1 \\ (X, Y) &\sim \mathbb{P}_{x,y} = \mathbb{P}\{X = x, Y = y\} = \mathbb{P}\{X|Y = y\} \mathbb{P}\{Y = y\} \\ \text{Total Err} &= c_{10} \mathbb{P}\{Y = 0, \hat{Y} = 1\} + c_{01} \mathbb{P}\{Y = 1, \hat{Y} = 0\}, c_{10}, c_{01} > 0 \\ &= c_{10} \int_{R_1} \underbrace{p(x|y=0)p(y=0)}_{\geq 0} dx + c_{01} \int_{R_0} \underbrace{p(x|y=1)p(y=1)}_{\geq 0} dx \end{aligned}$$

Aside,

X with density $p(x)$, A sub $X, \mathbb{P}\{X \in A\} = \int_A p(x) dx$

If $p(x|y=0)p(y=0) > p(x|y=1)p(y=1)$

Then put x in R_0

Conversely if $<$, then x in R_1

$$\begin{aligned} \frac{p(x|y=1)p(y=1)}{p(x|y=0)p(y=0)} &\stackrel{\geq \hat{y}=1}{\stackrel{< \hat{y}=0}{>}} 1 \\ \text{LR} &= \frac{p(x|y=1)}{p(x|y=0)} \stackrel{\geq \hat{y}=1}{\stackrel{< \hat{y}=0}{>}} \frac{c_{10}p(y=0)}{c_{01}p(y=1)} \end{aligned}$$

$$\text{LRT} = \frac{p(x|y=1)}{p(x|y=0)} \underset{\hat{y}=0}{\overset{\hat{y}=1}{\geq}} t$$

6.2 Constrain one type of error

$$\min \mathbb{P}\{\hat{y} = 0, Y = 1\} \text{ such that } \mathbb{P}\{\hat{y} = 1, Y = 0\} \leq \alpha < 1$$

Neyman-Pearson decision

Solution:

$$\frac{p(x|y=1)}{p(x|y=0)} \underset{\hat{y}=0}{\overset{\hat{y}=1}{\geq}} t_\alpha$$

Example

If $X|Y = j \sim N(\mu_j, \Sigma_j)$

$$p(x|y=j) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_j|}} e^{-\frac{1}{2}(x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j)}$$

$$\log(p(x|y=j)) = -\frac{1}{2} \underbrace{(x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j)}_{\text{Mahalanobis distance}} + \text{constant}$$

Special case:

$$\Sigma_0 = \Sigma_1 = \Sigma$$

Linear classifier

$\{(X_i, Y_i)\} \stackrel{iid}{\sim} \mathbb{P}_{X,Y}$ MVN cross conditionals

$$j = 0, 1$$

$$\Rightarrow \hat{\mu}_j = \frac{1}{\#\{i : y_i = j\}} \sum_{i: y_i = j} x_i$$

$$\hat{\Sigma}_j = \frac{1}{\#\{i : y_i = j\}} \sum_{i: y_i = j} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T$$

$$\mu_j = \mathbb{E}[X|Y = j]$$

$$\Sigma_j = \mathbb{E}\left[(X - \mu_j)(X - \mu_j)^T | Y = j\right]$$

Is there a natural notion of "distance" for general class-conditional distributions?

$$p_0(x) = \mathbb{P}\{X|Y = 0\}$$

$$p_1(x) = \mathbb{P}\{X|Y = 1\}$$

$$\log \text{LR} = \log \frac{P_1(x)}{P_0(x)} = \Lambda(x) \underset{\hat{y}=0}{\overset{\hat{y}=1}{\geq}} 0$$

if $X \sim q$ (q may be P_0 or P_1 or even something else)

What do we expect $\log \text{LR}$ to be?

$$\begin{aligned}\mathbb{E}_{q[\Lambda(X)]} &= \int q(x) \log \left(\frac{p_1(x)}{p_0(x)} \right) dx \\ &= \int q(x) \log \left(\frac{p_1(x)}{p_0(x)} \cdot \frac{q(x)}{q(x)} \right) dx \\ &= \int q(x) \log \left(\frac{q(x)}{p_0(x)} \right) dx - \int q(x) \log \left(\frac{q(x)}{p_1(x)} \right) dx \\ &= D(q\|p_0) - D(q\|p_1)\end{aligned}$$

KLD of p_i from q , Kullback-Leibler divergence

$$D(q\|p_0) \stackrel{\geq \hat{y}=1}{\stackrel{< \hat{y}=0}{\leq}} D(q\|p_1)$$

KL for MVN

$$D(N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)) = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - d + \log \left(\frac{|\Sigma_1|}{|\Sigma_0|} \right) \right)$$

$$D(q\|p) \geq 0$$

7 Lecture 7

7.1 Optimal Classification

$$(X, Y) \sim \mathbb{P}_{X,Y}$$

$p(x, y)$ be joint prob density or mass function

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

explicit notation

$$\begin{aligned}p(y|x) &= \mathbb{P}\{Y = y | X = x\} \\ \eta(x) &= \mathbb{P}\{Y = 1 | X = x\} \\ 1 - \eta(x) &= \mathbb{P}\{Y = 0 | X = x\}\end{aligned}$$

min prob of error classifier

$$f^*(x) = \begin{cases} 0 & \text{if } \eta(x) \geq \frac{1}{2} \\ 1 & \text{otherwise} \end{cases}$$

$$\eta(x) \geq \frac{1}{2} \Leftrightarrow \frac{\eta(x)}{1 - \eta(x)} \geq 1$$

$$\begin{aligned}
\frac{\eta(x)}{1-\eta(x)} &= \frac{\mathbb{P}\{Y=1|X=x\}}{\mathbb{P}\{Y=0|X=x\}} = \frac{p(y=1|x)}{p(y=0|x)} \\
&= \frac{\frac{p(x,y=1)}{p(x)}}{\frac{p(x,y=0)}{p(x)}} = \frac{p(x,y=1)}{p(x,y=0)} \\
&= \frac{p(x|y=1)p(y=1)}{p(x|y=0)p(y=0)} = \text{LR}
\end{aligned}$$

$$\begin{aligned}
p(x,y) &= p(y|x)p(x) \\
&= p(x|y)p(y)
\end{aligned}$$

class-conditional distribution

$$\begin{aligned}
p_0(x) &= p(x|y=0) \\
p_1(x) &= p(x|y=1) \\
\Lambda(x) &= \log\left(\frac{p_1(x)}{p_0(x)}\right) \\
\Lambda &\stackrel{\hat{y}=1}{\geq} \stackrel{\hat{y}=0}{\leq} 0
\end{aligned}$$

$X \sim q$ prob dist ("test")

$\Lambda(X)$ is a real-valued random variable

$$\begin{aligned}
\Lambda(X) &= \underbrace{\mathbb{E}[\Lambda(X)]}_{\text{deterministic, a number}} + \underbrace{(\Lambda(X) - \mathbb{E}[\Lambda(X)])}_{\text{zero-mean random variable}} \\
\mathbb{E}[\Lambda(X)] &= \int \Lambda(x) q(x) dx \\
&= \int q(x) \log\left(\frac{p_1(x)}{p_0(x)} \cdot \frac{q(x)}{q(x)}\right) dx \\
&= \underbrace{\int q(x) \log\left(\frac{q(x)}{p_0(x)}\right) dx}_{D(q\|p_0)} - \underbrace{\int q(x) \log\left(\frac{q(x)}{p_1(x)}\right) dx}_{D(q\|p_1)}
\end{aligned}$$

Kullback-Leibler (KL) divergences

Lemma 2. $D(q\|p) \geq 0$ for any q, p distributions, $D(q\|q) = 0$.

Proof. ,

$$\begin{aligned}
D(q\|p) &= \int q(x) \log\left(\frac{q(x)}{p(x)}\right) dx \\
&= - \int q(x) \log\left(\frac{p(x)}{q(x)}\right) dx \\
&= -\mathbb{E}_q\left[\log\left(\frac{p(x)}{q(x)}\right)\right]
\end{aligned}$$

$$\begin{aligned}
&\geq -\log \left(\mathbb{E} \left[\frac{p(x)}{q(x)} \right] \right) \\
&= -\log \left(\int q(x) \frac{p(x)}{q(x)} dx \right) \\
&= -\log(1) \\
&\geq 0
\end{aligned}$$

Jensen's Inequality. If f is convex, then $\mathbb{E}[f(Z)] \geq f(\mathbb{E}[Z])$ □

$$D(q\|p_0) - D(q\|p_1)$$

Case 1 : $q = p_1$

Case 2 : $q = p_0$

Example 1. $X|Y=0 \sim N(-\mu, 1), X|Y=1 \sim N(\mu, 1)$

$$\begin{aligned}
\Lambda(x) &= \log \left(\frac{\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}}}{\frac{1}{\sqrt{2\pi}} e^{-\frac{(x+\mu)^2}{2}}} \right) \\
&= -\frac{(x-\mu)^2}{2} + \frac{(x+\mu)^2}{2} \\
&= \frac{2\mu x + 2\mu x}{2} \\
&= 2\mu x \stackrel{\hat{y}=1}{\geq} \stackrel{\hat{y}=0}{<} \\
\Lambda(x) &= 2\mu x
\end{aligned}$$

$$\begin{aligned}
X &\sim p_1 \\
\Lambda(x) &= 2\mu X \sim N(2\mu^2, 4\mu^2), \sigma = 2\mu
\end{aligned}$$

μ bigger is better

MVN: $x|y=j \sim N(\mu_j, \Sigma), j=0,1$

$$\begin{aligned}
D(p_0\|p_1) &= D(p_1\|p_0) \\
&= \frac{1}{2} (\mu_1 - \mu_0)^T \Sigma^{-1} (\mu_1 - \mu_0) \\
D(p_1\|p_0) &= D(p_0\|p_1) \\
&= 2\mu^2
\end{aligned}$$

8 Lecture 8

8.1 Bayes Classifier (minimum prob of err)

$$f^*(x) = \begin{cases} +1 & \text{if } \mathbb{P}(y = 1|X = x) \geq \frac{1}{2} \\ -1 & \text{otherwise} \end{cases}$$

or

$$\begin{cases} +1 & \text{if } \log\left(\frac{p(x|y=1)}{p(x|y=-1)}\right) \geq \log\left(\frac{p(y=-1)}{p(y=+1)}\right) \\ -1 & \text{otherwise} \end{cases}$$

8.2 Nearest Neighbor Classifier

$$\begin{aligned} \{(X_i, Y_i)\} &\stackrel{iid}{\sim} \mathbb{P}_{XY} \\ f_{1nn}(x) &= y_{i_x}, i_x = \arg \min_i \|x - x_i\| \\ \lim_{n \rightarrow \infty} \mathbb{P}\{f_{1nn}(X) \neq Y\} &\leq 2\mathbb{P}\{f^*(X) \neq Y\} \end{aligned}$$

Example 2. $x \in \{-1, 1\}^d, y \in \{-1, 1\}, y = x_1$, and $x_2, \dots, x_d \stackrel{iid}{\sim} \pm 1$ with probability $\frac{1}{2}$.

Bayes Error = 0

$$n = 2, x = (+1, +1, \dots, +1)$$

$$d = 2, x = (1, 1)$$

Possible cases for $(x_1, +1)$ and $(x_2, -1)$,

1	-1	correct
-1	1	incorrect
1	1	correct
-1	-1	incorrect
1	-1	tie
-1	1	tie
1	-1	correct
-1	-1	incorrect

$$\frac{1}{4} \cdot \frac{1}{2} = \frac{1}{8}$$

$$\mathbb{P}\{f_{1nn}(X) \neq Y\} \rightarrow \frac{1}{2} \text{ as } d \rightarrow \infty$$

8.3 Generative Model Plug-in Classifier

$$p(y = +1) = p(y = -1)$$

$$X|Y = +1 \sim N(\theta, I)$$

$$X|Y = -1 \sim N(-\theta, I)$$

Bayes Classifier

$$-\frac{1}{2}\|x - \theta\|^2 + \frac{1}{2}\|x + \theta\|^2 - \frac{1}{2}\|x - \theta\|^2 + \frac{1}{2}\|x + \theta\|^2$$

$$f^*(x) = \begin{cases} +1 & \text{if } x^T \theta > 0 \\ -1 & \text{if } x^T \theta < 0 \end{cases}$$

Bayes Err Rate

$$\begin{aligned} \mathbb{P}\{f^*(X) \neq Y\} &= \frac{1}{2}\mathbb{P}\{x^T \theta > 0 | Y = -1\} + \frac{1}{2}\mathbb{P}\{x^T \theta < 0 | Y = +1\} \\ \mathbb{P}\{x^T \theta > 0 | Y = -1\}, x^T \theta &\sim N(-\|\theta\|^2, \|\theta\|^2) \\ \mathbb{P}\{Z > 0\} &= \mathbb{P}\{Z' > \|\theta\|^2\}, Z \sim N(-\|\theta\|^2, \|\theta\|^2), Z' \sim N(0, \|\theta\|^2) \\ &\leq \frac{\|\theta\|^2}{(\|\theta\|^4)} \\ &= \frac{1}{\|\theta\|^2} \end{aligned}$$

Use Markov,

$$\mathbb{P}\{Z > t\} \leq \mathbb{P}\{Z^2 > t^2\} \leq \frac{\mathbb{E}[Z^2]}{t^2}$$

Plug-in,

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n y_i x_i$$

$$\hat{f}(x) = \begin{cases} +1 & \text{if } x^T \hat{\theta} > 0 \\ -1 & \text{if } x^T \hat{\theta} < 0 \end{cases}$$

What is the distribution of $x^T \hat{\theta} | Y = -1$

$$X = -\theta + e_1, e_1 \sim N(0, I)$$

$$\hat{\theta} = \theta + e_2, e_2 \sim N\left(0, \frac{1}{n}I\right)$$

Ignoring constant factors,

$$\begin{aligned} \mathbb{P}\{x^T \hat{\theta} > 0 | Y = -1\} &\leq \frac{1}{\|\theta\|^2} + \frac{d^2}{n} \frac{1}{\|\theta\|^4} \\ \frac{d^2}{n} \frac{1}{\|\theta\|^4} &\approx \frac{1}{\|\theta\|^2} \\ n &\approx \left(\frac{\|\theta\|}{d}\right)^2 \end{aligned}$$

8.4 Maximum Likelihood Estimation

$$x_1, \dots, x_n \sim q$$
$$q \in \{p_\theta\}_{\theta \in \Theta}$$

Example 3. $x_1, \dots, x_n \stackrel{iid}{\sim} N(\theta, I)$ for some $\theta \in \mathbb{R}^d$

Two approaches

1. Method of Moments

$$\mu_f = \int f(x) q(x) dx, \text{ for any } f$$
$$\hat{\mu}_f = \frac{1}{n} \sum_{i=1}^n f(x_i)$$
$$\mu_f(\theta) = \int f(x) p_\theta(x) dx$$

find θ that minimizes,

$$|\hat{\mu}_f - \mu_f(\theta)| \rightarrow \hat{\theta}$$

for one or more functions f .

2. MLE

$$\hat{\theta} = \arg \max_{\theta} \prod_{i=1}^n p_\theta(x_i)$$
$$= \arg \min_{\theta} - \sum_{i=1}^n \log p_\theta(x_i)$$

9 Lecture 9

9.1 Maximum Likelihood

$$x_1, \dots, x_n \stackrel{iid}{\sim} q$$

Models $\{P_\theta\}_{\theta \in \Theta}$

$$L(\theta) = \log \left(\prod_{i=1}^n P_\theta(x_i) \right)$$

when viewed as function of θ is called the likelihood of θ .

$$\hat{\theta} = \arg \max_{\theta} L(\theta)$$

Example 4. $x|\theta \sim N(\theta, I), \theta \in \mathbb{R}^d$

$$\begin{aligned}
\log p(x|\theta) &= \log \left(\frac{1}{\sqrt{2\pi d}} \exp \left(-\frac{1}{2} (x_1 - \theta)^T (x_1 - \theta) \right) \right) \\
&= -\frac{1}{2} (x_1 - \theta)^T (x_1 - \theta) + \text{const} \\
&= -\frac{1}{2} (x_1^T x_1 - 2x_1^T \theta + \theta^T \theta) \\
\hat{\theta} &= \arg \max_{\theta \in \mathbb{R}^d} - \sum_{i=1}^n \frac{1}{2} (x_i - \theta)^T (x_i - \theta) \\
&\Rightarrow \sum_{i=1}^n (x_i - \theta) = 0 \\
&\Rightarrow \sum x_i = \sum \theta \\
&\Rightarrow \sum x_i = n\theta \\
&\Rightarrow \hat{\theta} = \frac{1}{n} \sum x_i
\end{aligned}$$

Example 5. $x|\theta \sim \text{Pois}(\theta), \theta > 0$

$$\begin{aligned}
p(x|\theta) &= e^{-\theta} \frac{\theta^x}{x!}, x = 1, 2, \dots \\
\log p(x|\theta) &= -\theta + x \log \theta - \log x! \\
&\quad x_1, x_2, \dots, x_n \\
&\quad \max_{\theta} \sum_i (-\theta + x_i \log \theta) \\
&\Rightarrow \sum_i (-\theta + x_i) = 0 \\
&\Rightarrow \hat{\theta} = \frac{1}{n} \sum x_i
\end{aligned}$$

MLE $\Rightarrow \hat{\theta}, P_{\hat{\theta}}$

Questions:

1. In what sense is $P_{\hat{\theta}}$ a good model for q
2. If $q = P_{\theta^*}$, then we hope that $\hat{\theta} \rightarrow \theta^*$ as $n \rightarrow \infty$

$$\begin{aligned}
\hat{\theta} &= \arg \max_{\theta} \sum_{i=1}^n \log p(x_i|\theta) \\
&= \arg \min_{\theta} - \sum_{i=1}^n \log p(x_i|\theta) + \log q(x_i) \\
&= \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \log \frac{q(x_i)}{p(x_i|\theta)} \\
&\quad \frac{1}{n} \sum_{i=1}^n \log \frac{q(x_i)}{p(x_i|\theta)} \rightarrow_{n \rightarrow \infty} D(q\|P_{\theta}) \text{ a.s.} \\
\theta^* &= \arg \min_{\theta} D(q\|p_{\theta})
\end{aligned}$$

$$\begin{aligned}
0 &\geq \frac{1}{n} \sum_{i=1}^n \log \frac{q(x_i)}{p(x_i|\hat{\theta})} - \sum_{i=1}^n \log \frac{q(x_i)}{p(x_i|\theta^*)} \\
&= D(q||p_{\hat{\theta}}) - D(q||p_{\theta^*}) \\
D(q||p_{\theta^*}) &\leq D(q||p_{\hat{\theta}}) \leq (\approx) D(q||p_{\theta^*})
\end{aligned}$$

Claim: $D(q||p_{\hat{\theta}}) \rightarrow D(q||p_{\theta^*})$ as $n \rightarrow \infty$

9.2 Central Limit Theorem

If z_1, \dots, z_n are zero mean, iid with variance σ^2 , then

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n z_i \right) \stackrel{\text{assmp dist}}{\sim} N(0, \sigma^2)$$

Theorem 8. $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} p_{\theta^*}$. Let $L(\theta) = \sum_{i=1}^n \log p(x_i|\theta)$. Assume $\frac{\partial L}{\partial \theta_j}$ and $\frac{\partial^2 L}{\partial \theta_j \partial \theta_k}$ exist $\forall j, k$. Then,

$$\sqrt{n} (\hat{\theta}_n - \theta^*) \stackrel{\text{assmp dist}}{\sim} N(0, I^{-1}(\theta^*))$$

where, Fisher Information matrix,

$$[I(\theta^*)]_{j,k} = -\mathbb{E} \left[\frac{\partial \log p(x|\theta)}{\partial \theta_j \partial \theta_k} \Big|_{\theta=\theta^*} \right]$$

Low curvature \Rightarrow high variance

High curvature \Rightarrow low variance in MLE

Proof. 1-dim case

$$L(\theta) = \sum_{i=1}^n \log p(x_i|\theta)$$

Taylor's Series (Mean Value Theorem):

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2} f''(\bar{x})(x - x_0)^2$$

\bar{x} is between x and x_0

$$\begin{aligned}
0 &= L'(\hat{\theta}) = L'(\theta^*) + L''(\bar{\theta})(\hat{\theta} - \theta^*) \\
&\Rightarrow (\hat{\theta} - \theta^*) = -\frac{L'(\theta^*)}{L''(\bar{\theta})}, \bar{\theta} \text{ is between } \hat{\theta} \text{ and } \theta^* \\
&\Rightarrow \sqrt{n}(\hat{\theta} - \theta^*) = -\frac{\frac{1}{\sqrt{n}} L'(\theta^*)}{\frac{1}{n} L''(\bar{\theta})} \\
\frac{1}{\sqrt{n}} L'(\theta^*) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial \log p(x_i|\theta)}{\partial \theta} \Big|_{\theta=\theta^*}
\end{aligned}$$

Then,

$$\begin{aligned}
\mathbb{E} \left[\frac{\partial \log p(x_i|\theta)}{\partial \theta} \Big|_{\theta=\theta^*} \right] &= \mathbb{E} \left[\frac{1}{p(x_i|\theta)} \frac{\partial p(x_i|\theta)}{\partial \theta} \Big|_{\theta=\theta^*} \right] \\
&= \int \frac{1}{p(x_i|\theta)} \frac{\partial p(x_i|\theta)}{\partial \theta} \Big|_{\theta=\theta^*} p(x_i|\theta^*) dx_i \\
&= \int 1 \frac{\partial p(x_i|\theta)}{\partial \theta} \Big|_{\theta=\theta^*} dx_i \\
&= \frac{\partial}{\partial \theta} \left(\int p(x_i|\theta) dx \right) \Big|_{\theta=\theta^*} \\
&= \frac{\partial 1}{\partial \theta} \Big|_{\theta=\theta^*} \\
&= 0
\end{aligned}$$

and,

$$\begin{aligned}
\frac{\partial^2 \log p(x|\theta)}{\partial \theta^2} &= \frac{\partial}{\partial \theta} \left(\frac{1}{p(x|\theta)} \frac{\partial p(x|\theta)}{\partial \theta} \right) \\
&= -\frac{1}{p^2(x|\theta)} \left(\frac{\partial p(x|\theta)}{\partial \theta} \right)^2 + \frac{1}{p(x|\theta)} \frac{\partial^2 p(x|\theta)}{\partial \theta^2}
\end{aligned}$$

In expectations given $\theta = \theta^*$,

$$\begin{aligned}
\mathbb{E} \left[\frac{\partial^2 \log p(x|\theta)}{\partial \theta^2} \right] &= \mathbb{E} \left[-\frac{1}{p^2(x|\theta)} \left(\frac{\partial p(x|\theta)}{\partial \theta} \right)^2 \right] + 0 \\
&= \int -\frac{1}{p^2(x|\theta)} \left(\frac{\partial p(x|\theta)}{\partial \theta} \right)^2 \Big|_{\theta=\theta^*} p(x|\theta^*) dx \\
&= -\mathbb{E} \left[\left(\frac{\partial \log p(x|\theta)}{\partial \theta} \right)^2 \Big|_{\theta=\theta^*} \right] \\
&= -\mathbb{V} \left[\frac{\partial \log p(x|\theta)}{\partial \theta} \Big|_{\theta=\theta^*} \right]
\end{aligned}$$

□

10 Lecture 10

10.1 Performance of MLE

$$\begin{aligned}
x_1, x_2, \dots, x_n &\stackrel{iid}{\sim} P_{\theta^*} \\
\hat{\theta}_n &= \arg \min_{\theta} - \sum_{i=1}^n \underbrace{\log p(x_i|\theta)}_{l(\theta)}
\end{aligned}$$

10.2 Fischer Info Matrix

$$I(\theta^*) = \left\{ \mathbb{E} \left[- \frac{\partial \log p(x|\theta)}{\partial \theta_j \partial \theta_k} \Big|_{\theta=\theta^*} \right] \right\}_{j,k}$$

$$\hat{\theta}_n \stackrel{asympt}{\sim} N \left(\theta^*, \underbrace{\frac{1}{n} I^{-1}(\theta^*)}_{\text{error covariance}} \right)$$

Assumptions:

1. $L'(\theta), L''(\theta)$ exist,
2. $\mathbb{E} \left[\frac{\partial \log p(x|\theta)}{\partial \theta} \right] = 0 \quad \forall \theta$

Example 6. $x_1, \dots, x_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$, σ^2 known, θ unknown, $\theta \in \mathbb{R}$

$$\hat{\theta} = \frac{1}{n} \sum x_i \sim N \left(\theta^*, \frac{\sigma^2}{n} \right)$$

$$\frac{\partial \log p(x|\theta)}{\partial \theta} = \frac{2}{\sigma} \left(-\frac{1}{2} \frac{(x-\theta)^2}{\sigma^2} + \text{const} \right) = \frac{x-\theta}{\sigma^2}$$

$$\frac{\partial^2 \log p(x|\theta)}{\partial \theta^2} = -\frac{1}{\sigma^2} \Rightarrow I(\theta^*) = \frac{1}{\sigma^2}$$

Example 7. $x_i \stackrel{iid}{\sim} \text{Unif}[0, \theta]$

$$\hat{\theta}_n = \max_{i=1, \dots, n} x_i$$

extremal stats

$$x_i \stackrel{iid}{\sim} N(0, 1)$$

$$\max_{i=1, \dots, n} x_i \stackrel{asympt}{\sim} \sqrt{2 \log n}$$

$$x_{i^*} \sim N(\mu, 1), \mu > 0$$

10.3 Sufficient Statistics

$$X \sim \text{Bernoulli}(\theta)$$

$$p(x|\theta) = \theta^x (1-\theta)^{1-x}, x \in \{0, 1\}$$

$$x_1, \dots, x_n \stackrel{iid}{\sim} \text{Be}(\theta)$$

$$S = \sum x_i \sim \text{Bi}(\theta, n)$$

$$p(s=k|\theta) = \binom{n}{k} \theta^k (1-\theta)^{n-k}$$

$$p(x_1, \dots, x_n, S|\theta) = \begin{cases} p(x_1, \dots, x_n|\theta) & \text{if } \sum x_i = S \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} p(x_1, \dots, x_n|s = k, \theta) &= \frac{p(x_1, \dots, x_n, s = k|\theta)}{p(s|\theta)} \\ &= \frac{\prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i}}{\binom{n}{k} \theta^k (1-\theta)^{n-k}} \\ &= \frac{\theta^k (1-\theta)^{n-k}}{\binom{n}{k} \theta^k (1-\theta)^{n-k}} \\ &= \frac{1}{\binom{n}{k}} \\ &= \max_{\theta} p(x_1, \dots, x_n, s|\theta) \\ &= \max_{\theta} p(x_1, \dots, x_n|s) p(s|\theta) \\ &= \max_{\theta} p \left(\underbrace{s}_{\text{sufficient statistics for } \theta} \mid \theta \right) \end{aligned}$$

10.4 Fischer-Neyman Factorization Thm

$$x_1, \dots, x_n \stackrel{iid}{\sim} P_{\theta}$$

$t(x_1, \dots, x_n)$ is a sufficient statistic for θ iff

$$p(x_1, \dots, x_n|\theta) = \underbrace{a(x_1, \dots, x_n)}_{\text{no } \theta \text{ dependence}} b(t, \theta)$$

MVN $x_i \stackrel{iid}{\sim} N(\mu, \Sigma)$, both unknowns, $\theta = (\mu, \Sigma)$, $x_i \in \text{real}^d$, $i = 1, \dots, n$, nd

Sufficient Statistics:

$$\begin{aligned} n\hat{\mu} &= \sum_i x_i \\ n\hat{\Sigma} &= \sum_i (x_i - \hat{\mu})(x_i - \hat{\mu})^T = t(x_1, \dots, x_n) \\ &\approx d + d^2 \text{ numbers} \end{aligned}$$

10.5 Rao-Blackwell Theorem

$x_1, \dots, x_n \stackrel{iid}{\sim} P_{\theta^*}, t = t(x_1, \dots, x_n)$, sufficient for θ ,

Let $\underbrace{f(x_1, \dots, x_n)}_{\hat{\theta}}$ be an estimator of θ

$$\underbrace{g([t(x_1, \dots, x_n)])}_{\text{just depends on } x\text{'s through sufficient statistics}} = \mathbb{E}[f(x_1, \dots, x_n) | t(x_1, \dots, x_n) = t]$$

$$\mathbb{E}[g(t(x_1, \dots, x_n))] = \mathbb{E}[\mathbb{E}[f(x_1, \dots, x_n) | t(x_1, \dots, x_n) = t]] = \mathbb{E}[f(x_1, \dots, x_n)]$$

Then,

$$\mathbb{E}[g(t(x_1, \dots, x_n))^2] \leq \mathbb{E}[(f(x_1, \dots, x_n) - \theta)^2]$$

Example 8. $n > 2, x_1, \dots, x_n \stackrel{iid}{\sim} N(\theta, 1), t_n = \frac{1}{n} \sum_{i=1}^n x_i$

$$f(x_1, \dots, x_n) = \frac{x_1 + x_2}{2}$$

10.6 Linear Models

predict y from $x \in \mathbb{R}^d$

$$\hat{y} = f(x^T w)$$

f non-linear, w = weights $\in \mathbb{R}^d$, "weighted sum of features"

$$\hat{y} = \sum_{j=1}^k w_j x_j$$

10.7 Linear Regression

$\{(x_i, y_i)\}_{i=1}^n$ training data

$$\hat{y}_i = x_i^T w$$

prediction errors

$$e_1 = y_1 - x_1^T w$$

$$e_2 = y_2 - x_2^T w$$

...

$$e_n = y_n - x_n^T w$$

Least squares

$$y = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix}$$

$$x = \begin{bmatrix} x_1^T \\ \dots \\ x_n^T \end{bmatrix}$$

$$w = \begin{bmatrix} w_1 \\ \dots \\ w_n \end{bmatrix}$$

$$\min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

$$\min_w \|y - Xw\|^2$$

$$\|y - Xw\|^2 = (y - Xw)^T (y - Xw)$$

$$= y^T y - 2y^T Xw + w^T X^T Xw$$

$$0 = \frac{\partial \text{above}}{\partial w} = 2X^T y - 2X^T Xw$$

$$X^T y = X^T Xw$$

normal equations

$\hat{w} = (X^T X)^{-1} X^T y$, assuming invertible

11 Lecture 11

11.1 Linear Models

$$e_1 = y_1 - x_1^T w$$

$$e_2 = y_2 - x_2^T w$$

11.2 Least Squares (LS)

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - x_i^T w)^2$$

$$(X^T X)^{-1} X^T y = \arg \min_w \|y - Xw\|_2^2$$

$$\sum_i (y_i - x_i^T w)^2 = \sum_i (y_i^2 - 2w^T x_i y_i + w^T x_i x_i^T w)$$

$$w^T \left(\sum \begin{bmatrix} x_{i1} \\ \dots \\ x_{id} \end{bmatrix} y_i \right) \rightarrow w^T \begin{bmatrix} \sum x_{i1} y_i \\ \dots \\ \sum x_{id} y_i \end{bmatrix}$$

If for example,

$$\begin{aligned}\sum x_{i1}y_i \text{ large } + &\Rightarrow w_1 > 0 \\ \sum x_{i2}y_i \text{ large } - &\Rightarrow w_2 < 0 \\ \sum x_{i3}y_i = 0 &\Rightarrow w_3 \approx 0\end{aligned}$$

11.3 MLE Perspective

$$\begin{aligned}y_i &= x_i^T w + v_i, v_i \stackrel{iid}{\sim} N(0, 1) \\ y_i - x_i^T w &= v_i \sim N(0, 1) \\ p(v_i) &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}v_i^2} \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y_i - x_i^T w)^2} \\ \log \prod_{i=1}^n p(v_i) &= \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y_i - x_i^T w)^2} \right) \\ &= -\frac{1}{2} \sum_{i=1}^n (y_i - x_i^T w)^2 + \text{const.}\end{aligned}$$

$$p(v_i) = \frac{1}{2} e^{-|v_i|} \text{ heavy tails}$$

$$\begin{aligned}\log \prod_{i=1}^n \frac{1}{2} e^{-|y_i - x_i^T w|} \\ &= -\frac{1}{2} \sum_{i=1}^n |y_i - x_i^T w| + \text{const} \\ &\Rightarrow \text{MLE } \max_w - \sum |y_i - x_i^T w| \\ &\Rightarrow \min_w \sum |y_i - x_i^T w| \\ &\Rightarrow \min_w \|y - x^T w\|_1\end{aligned}$$

$$z_1, z_2, \dots, z_n$$

$$\begin{aligned}\min_{\theta \in \mathbb{R}} \sum_{i=1}^n (z_i - \theta)^2 &\Rightarrow \theta \text{ is mean} \\ \min_{\theta \in \mathbb{R}} \sum_{i=1}^n |z_i - \theta| &\Rightarrow \theta \text{ is median}\end{aligned}$$

11.4 Binary Classification and Linear Prediction

$$y \in \{0, 1\}, \hat{y}_i = x_i^T w$$

$$\mathbb{P}\{y_i = 1\} = f(x_i^T w), f: \mathbb{R} \rightarrow [0, 1]$$

$$\hat{y}_i = \begin{cases} 1 & \text{if } \text{sign}(x_i^T w) > 0 \\ 0 & \text{otherwise} \end{cases}$$

11.5 Bernoulli Distribution

$$\begin{aligned} p(y_i) &= p_i^{y_i} (1 - p_i)^{1-y_i} \\ &= f(x_i^T w)^{y_i} (1 - f(x_i^T w))^{1-y_i} \\ &= \exp(y_i \log p_i + (1 - y_i) \log(1 - p_i)) \\ &= \exp\left(y_i \underbrace{\log\left(\frac{p_i}{1 - p_i}\right)}_{\theta_i} + \log(1 - p_i)\right) \end{aligned}$$

$$\text{Canonical } \theta_i = \log\left(\frac{p_i}{1 - p_i}\right)$$

$$\begin{aligned} \theta_i &= x_i^T w \\ y_i \theta_i &= y_i x_i^T w \\ &= w x_i y_i^T \\ e^{\theta_i} &= \frac{p_i}{1 - p_i} \\ e^{\theta_i} (1 - p_i) &= p_i \\ e^{\theta_i} &= (1 + e^{\theta_i}) p_i \\ p_i &= \frac{e^{\theta_i}}{1 + e^{\theta_i}} = \frac{1}{1 + e^{-\theta_i}} \\ f(x_i^T w) &= \frac{1}{1 + e^{-\theta_i}} \\ \max_{w \in \mathbb{R}^d} \prod_{i=1}^n \exp\left(y_i x_i^T w + \log\left(1 - \frac{1}{1 + e^{-x_i^T w}}\right)\right) \end{aligned}$$

11.6 Logistic Regression

$$\begin{aligned} \theta_i &= x_i^T w \\ p_i^{y_i} (1 - p_i)^{1-y_i} &= \exp(y_i \log p_i + (1 - y_i) \log(1 - p_i)) \\ p(y_i | \theta_i) &= \exp\left(y_i \log\left(\frac{1}{1 + e^{-\theta_i}}\right) + (1 - y_i) \log\left(\frac{1}{1 + e^{\theta_i}}\right)\right) \\ \log(p(y_i | \theta_i)) &= \begin{cases} \log\left(\frac{1}{1 + e^{-\theta_i}}\right) & \text{if } y_i = 1 \\ \log\left(\frac{1}{1 + e^{\theta_i}}\right) & \text{if } y_i = 0 \end{cases} \end{aligned}$$

$$z_i = 2y_i - 1 \in \{-1, 1\}$$

$$\sum_{i=1}^n \log p(z_i | \theta_i) = \sum_{i=1}^n \log \left(\frac{1}{1 + e^{-\theta_i z_i}} \right)$$

MLE, $\theta_i = w^T x_i$

$$\max_{w \in \mathbb{R}^d} \sum_{i=1}^n \log \left(\frac{1}{1 + e^{-w^T x_i z_i}} \right)$$

$$= \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \underbrace{\log \left(1 + e^{-w^T x_i z_i} \right)}_{\text{logistic loss function}}$$

LS,

$$\min_w \sum_{i=1}^n \underbrace{(y_i - x_i^T w)^2}_{\text{squared err loss function}}, z_i = 2y_i - 1$$

11.7 Multiclass Problems

$$y_i \in \{1, 2, \dots, k\}$$

k weight vectors, w_1, \dots, w_k

$$l \in \{1, \dots, k\}$$

$$p(y_i = l) = \frac{e^{w_l^T x_i}}{\sum_{j=1}^k e^{w_j^T x_i}} \in [0, 1]$$

$$\max_{w_1, \dots, w_k \in \mathbb{R}^d} \sum_{i=1}^n \log \left(\frac{e^{w_l^T x_i}}{\sum_{j=1}^k e^{w_j^T x_i}} \right)$$

Softmax.

Multinomial Logistic Regression

12 Lecture 12

12.1 GLMs

Models for $p(y|x)$ then depend on x only through a linear map $w^T x$. The weight $w \in \mathbb{R}^d$ can be fit by maximum likelihood.

$$p(y|x) = p(y|w^T x)$$

Gaussian: $y_i \in \mathbb{R}$,

$$p(y_i | w^T x_i) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (y_i - w^T x_i)^2\right)$$

Binomial: $y_i \in \{0, 1\}$

$$p(y_i | w^T x_i) = \exp\left(y_i \log\left(\frac{1}{1 + e^{-w^T x_i}}\right) + (1 - y_i) \log\left(\frac{1}{1 + e^{w^T x_i}}\right)\right)$$

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

$$p(y_i | w^T x_i) = \exp\left(\log\left(\frac{1}{1 + e^{-y_i x_i^T w}}\right)\right)$$

12.2 MLEs

Gaussian:

$$\begin{aligned} \min_{w \in \mathbb{R}^d} - \sum_{i=1}^n \underbrace{\frac{1}{2} (y_i - w^T x_i)^2}_{\text{squared error loss}} \\ = \min_w \|y - Xw\|^2 \end{aligned}$$

If X is full rank,

$$\hat{w} = (X^T X)^{-1} X^T y$$

Bernoulli:

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n \underbrace{\log(1 + \exp(-y_i x_i^T w))}_{\text{logistic loss}}$$

12.3 Binary Classification

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

sq err loss:

$$(y_i - w^T x_i)^2 = (1 - y_i x_i^T w)^2$$

logistic loss:

$$\log(1 + e^{-y_i x_i^T w})$$

hinge loss:

$$\max\{0, 1 - y_i x_i^T w\} = (1 - y_i x_i^T w)^+$$

0 – 1 loss:

$$\mathbb{1}_{y_i x_i^T w < 0}$$

Then,

$$\begin{aligned}\hat{y}_i &= \text{sign}(x_i^T w) \\ \{(x_i, y_i)\}_{i=1}^n &\stackrel{iid}{\sim} \mathbb{P}_{xy} \\ \mathbb{P}_{xy} \{y &\neq \text{sign}(\hat{w}^T x)\}\end{aligned}$$

12.4 Convex Losses

loss l is convex in w if for any w_0, w_1 , and $\lambda \in [0, 1]$

$$l(\lambda y_i x_i^T w_0 + (1 - \lambda) y_i x_i^T w_1) \leq \lambda l(y_i x_i^T w_0) + (1 - \lambda) l(y_i x_i^T w_1)$$

12.5 Least Square

$$\begin{aligned}\min_w \quad & \underbrace{\|y - Xw\|^2}_{(y-Xw)^T(y-Xw)} \\ \frac{\partial \text{ above}}{\partial w} &= -2X^T y + 2X^T Xw \\ &= -2X^T (y - Xw)\end{aligned}$$

Set it to 0,

$$\begin{aligned}X^T y &= X^T Xw \\ w^* &= (X^T X)^{-1} X^T y\end{aligned}$$

w_1 initial guess (often 0 or random)

$$\begin{aligned}w_{t+1} &= w_t - \gamma \nabla f(w_t), t = 1, 2, \dots \\ &= w_t + \underbrace{\gamma}_{\text{step size}} X^T (y - X^T w_t)\end{aligned}$$

13 Lecture 13

13.1 Gradient Descent and LS

$$\begin{aligned}\min_w \|y - Xw\|^2 &= \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 \\ f(w) &= \|y - Xw\|^2 \\ \nabla f(w) &= \frac{\partial}{\partial w} \left((y - Xw)^T (y - Xw) \right)\end{aligned}$$

$$\begin{aligned}
&= -X^T y + X^T X w \\
&= -X^T (y - X w) \\
w_{t+1} &= w_t - \gamma \nabla f(w_t), t = 1, 2, \dots, \gamma > 0 \\
&= w_t + \gamma X^T (y - X w_t) \\
&= w_t + \gamma (X^T - X^T X w_t) \\
&= w_t + \gamma X^T X \left((X^T X)^{-1} X^T y - w_t \right) \\
&= w_t - \gamma X^T X (w_t - w^*)
\end{aligned}$$

Goal:

$$\begin{aligned}
w_t - w^* &\rightarrow 0 \text{ as } t \rightarrow \infty \\
w_{t+1} - w^* &= w_t - w^* - \gamma X^T X (w_t - w^*) \\
&= (I - \gamma X^T X) \underbrace{(w_t - w^*)}_{v_t} \\
v_{t+1} &= (I - \gamma X^T X) v_t \\
&= (I - \gamma X^T X) (I - \gamma X^T X) v_{t-1} \\
&= (I - \gamma X^T X)^t v_1
\end{aligned}$$

Eigendecomposition of $(I - \gamma X^T X)$

$$I - \gamma X^T X = U D U^T$$

U : orthogonal: $U^T U = I$

D : diagonal $\begin{bmatrix} \gamma_1 & 0 \\ 0 & \gamma_2 \end{bmatrix}$

$$\begin{aligned}
D^2 &= \begin{bmatrix} \gamma_1^2 & 0 \\ 0 & \gamma_2^2 \end{bmatrix} \\
(UD D^T)^2 &= (UD U^T) (UD U^T) \\
&= U D I D U^T \\
&= U D^2 U^T
\end{aligned}$$

Convergence Condition:

$$\begin{aligned}
|\lambda_i| &< 1 \text{ for } i = 1, 2, \dots, d \\
D &= U^T U D U^T U \\
&= U^T (I - \gamma X^T X) U \\
&= (I - \gamma U^T X^T X U) \leftarrow \text{is diagonal}
\end{aligned}$$

$$U^T X^T X U = \text{diag}(\lambda_1(X^T X) > 0, \dots, \lambda_d(X^T X) > 0)$$

is positive definite.

$$\begin{aligned} |\lambda_i| &= |1 - \gamma \lambda_i(X^T X)| < 1 \\ \gamma \lambda_i(X^T X) &\leq 2, i = 1, \dots, d \\ \Rightarrow \gamma &< \frac{2}{\lambda_{\max}(X^T X)} \end{aligned}$$

LS GD:

$$w_{t+1} = w_t - \gamma \frac{\partial}{\partial w} \left(\frac{1}{2} \sum_{i=1}^n (y_i - x_i^T w)^2 \right) \Big|_{w=w_t}$$

Stochastic GD:

$$w_{t+1} = w_t - \gamma \frac{\partial}{\partial w} \frac{1}{2} (y_{i_t} - x_{i_t}^T w)^2 \Big|_{w=w_t}, i_t \sim \text{Unif}(1, \dots, n)$$

Then,

$$\begin{aligned} &\mathbb{E} \left[\frac{\partial}{\partial w} \frac{1}{2} (y_{i_t} - X_{i_t}^T w)^2 \right] \\ &= \mathbb{E} \left[\frac{\partial}{\partial w} \sum_{i=1}^n (y_i - X_i^T w) \mathbb{1}_{\{i=i_t\}} \right] \\ &= \frac{\partial}{\partial w} \frac{1}{n} \sum_{i=1}^n (y_i - X_i^T w)^2 \end{aligned}$$

13.2 SGD with Convex Loss

$$w^* = \arg \max_{w \in \mathbb{R}^d} = \frac{1}{T} \sum_{i=1}^T f_t(w), \text{ convex in } w$$

LS:

$$f_t(w) = (y_{i_t} - X_{i_t}^T w)^2$$

Logistic:

$$f_t(w) = \log(1 + \exp(-y_{i_t} x_{i_t}^T w))$$

$T \rightarrow \infty$, LS:

$$\begin{aligned} &\frac{1}{T} \sum_{i=1}^T \left(\underbrace{y_{i_t} - X_{i_t}^T w}_{Z_t \text{ iid}} \right)^2 \xrightarrow{as} \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T w)^2 \\ \mathbb{E}[Z_t] &= \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T w)^2 \end{aligned}$$

Convex iff f is above tangent plane
for any $w, u \in \mathbb{R}^d$

$$f(u) \geq f(w) + \nabla^T f(w)(u - w)$$

Example,

$$\begin{aligned} w^2 &= f(w) \\ \frac{\partial}{\partial x} w^2 &= 2w = 0 \\ w = 0, f(w) &= 0, \nabla f(0) = 0 \\ f(u) &\geq 0 + 0(u - w) \\ &\geq 0 \\ w = 1, f(1) &= 1, \nabla f(1) = 2 \\ f(u) &\geq 1 + 2(u - w) \end{aligned}$$

14 Lecture 14

14.1 Supervised ML

$\{(x_i, y_i)\}_{i=1}^n \stackrel{iid}{\sim} \mathbb{P}_{XY}$, unknown

learn to predict y from x

What would we do if we knew \mathbb{P}_{XY}
in Binary Classification

$$y \in \{0, 1\} \text{ or } \{-1, +1\}$$

$$\frac{\mathbb{P}\{Y = 1|X = x\}}{\mathbb{P}\{Y = -1|X = x\}} \underset{\text{class } -1}{\underset{1}{\geq}}$$

14.2 Bayes Rule

$$\begin{aligned} p(y|x) &= \frac{p(x|y)p(y)}{p(x)} \\ \frac{p(x|y = +1)p(y = +1)}{p(x|y = -1)p(y = -1)} &\underset{\text{class } -1}{\underset{1}{\geq}} \\ \frac{p(x|y = +1)}{p(x|y = -1)} &\underset{\text{class } -1}{\underset{1}{\geq}} \frac{p(y = -1)}{p(y = +1)} \end{aligned}$$

$$\eta(x) = p(y = +1|x)$$

$$\hat{\eta}(x) \underset{\text{class } -1}{\underset{1}{\geq}} \frac{1}{2}$$

"plug-in" approach

14.3 Empirical Risk Minimization

$\{(x_i, y_i)\}_{i=1}^n$, set of classifiers \mathcal{F}

$$f(x) = \pm 1$$

$$\min_{f \in \mathcal{F}} \sum_{i=1}^n \text{loss}(f(x_i), y_i)$$

0-1 loss: $\mathbb{1}_{\{f(x_i) \neq y_i\}} = \mathbb{1}_{\{y_i f(x_i) < 0\}}$

sq err loss: $(y_i - f(x_i))^2 = (1 - y_i f(x_i))^2$

abs err loss: $|1 - y_i f(x_i)|$

logistic loss: $\log(1 + \exp(-y_i f(x_i)))$

hinge: $\max\{0, 1 - y_i f(x_i)\}$

14.4 Linear Models

$$\begin{aligned} f(x_i) &= w^T x_i \\ &= \sum_{j=1}^n w_j x_{ij} \end{aligned}$$

14.5 GLMs

Parametric model $p(y_i|x_i) = p(y_i|w^T x_i)$

Binary Classification $p(y|x)$ is Bernoulli

$$\begin{aligned} y_i &\in \{-1, +1\} \\ p(y_i|w^T x_i) &= \frac{1}{1 + \exp(-y_i x_i^T w_i)} \\ \log(p(y_i|w^T x_i)) &= -\log(1 + \exp(-y_i x_i^T w_i)) \end{aligned}$$

MLE of w

$$\begin{aligned} &\max_w \prod_{i=1}^n p(y_i|w^T x_i) \\ &\Rightarrow \min_w - \sum_{i=1}^n \log p(y_i|x_i^T w) \\ &\Rightarrow \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w_i)) \end{aligned}$$

Example 9. $X|Y \sim N\left(y \frac{w}{2}, I\right)$

$$y \in \{-1, +1\}, p(y = +1) = p(y = -1)$$

opt classifier

$$\begin{aligned}
 p(y|x) &= \frac{p(x|y)p(y)}{p(x)} \\
 p(x) &= p(x, y = +1) + p(x, y = -1) \\
 p(x|y = +1) &= \exp\left(-\frac{1}{2}\left(x - \frac{w}{2}\right)^T \left(x - \frac{w}{2}\right)\right) \\
 &= \exp\left(-\frac{1}{2}x^T x\right) \exp\left(\frac{1}{8}w^T w\right) \exp\left(\frac{1}{2}w^T x\right)
 \end{aligned}$$

$$\begin{aligned}
 p(y = +1|X = x) &= \frac{p(x|y = +1)}{p(x|y = +1) + p(x|y = -1)} \\
 &= \frac{\exp\left(\frac{1}{2}w^T x\right)}{\exp\left(\frac{1}{2}w^T x\right) + \exp\left(-\frac{1}{2}w^T x\right)} \\
 &= \frac{1}{1 + \exp(-w^T x)}
 \end{aligned}$$

MLE of μ

$$\begin{aligned}
 \hat{\mu} &= \frac{1}{n} \sum_{i=1}^n y_i x_i \\
 \mathbb{E}[\hat{\mu}] &= \mu = \frac{w}{2} \\
 \Rightarrow \hat{w} &= 2\hat{\mu}
 \end{aligned}$$

15 Lecture 15

15.1 Optimization in ML

$$\begin{aligned}
 \min_{w \in \mathbb{R}^d} \sum_{t=1}^T f_t(w) \\
 f_t(w) &= (y_t - w^T x_t)^2 \\
 f_t(w) &= \log(1 - \exp(-y_t w^T x_t)) \\
 f_t(w) &= \max(0, 1 - y_t x_t^T w)
 \end{aligned}$$

$$(x_t, y_t) \stackrel{unif}{\sim} \{(x_i, y_i)\}_{i=1}^n$$

$$\frac{1}{T} \sum_{t=1}^T f_t(w) \xrightarrow{as} \frac{1}{n} \sum_{i=1}^n f_i(w)$$

15.2 Prob Model Approach

$$\begin{aligned}
 & p(y|w^T x) \\
 & \min_w \sum_{t=1}^n \underbrace{-\log p(y_t|w^T x_t)}_{f_t(w)} \\
 & \Rightarrow p(y_t|w^T x_t) = \exp(-f_t(w))
 \end{aligned}$$

What properties should f_t have?

1. f_t is convex function (opt)
2. f_t is non-negative (prob interpretation)

$$\begin{aligned}
 f_t &: \mathbb{R} \rightarrow [0, \infty) \\
 \exp(-f_t) &\in [0, 1]
 \end{aligned}$$

SGD $\{(x_i, y_i)\}_{i=1}^n$

$$\begin{aligned}
 (x_t, y_t) &\overset{unif}{\sim} \{(x_i, y_i)\}_{i=1}^n, t = 1, 2, \dots \\
 w_1 &\in \mathbb{R}^d \text{ arbitrary} \\
 w_{t+1} &= w_t - \gamma_t \nabla f_t(w_t), t = 1, 2, \dots
 \end{aligned}$$

γ_t is step size

Key fact

If f_t is convex:

$$f_t(w^*) \geq f_t(w) + (w^* - w)^T \nabla f_t(w)$$

Theorem 9. $\gamma_t = \gamma$ fixed, $\|\nabla f_t(w)\| \leq G \forall t, w$

Let $w_T^* = \arg \min_w \sum_{t=1}^T f_t(w)$

$$\frac{1}{T} \sum_{t=1}^T f_t(w_t) - f_t(w_T^*) \leq \frac{\|w_1 - w^*\|^2}{2\gamma T} + \frac{\gamma}{2} G^2$$

Remark: choose γ so,

$$\begin{aligned}
 \frac{1}{\gamma T} &\approx \gamma \Rightarrow \gamma = \frac{1}{\sqrt{T}} \\
 &O\left(\frac{1}{\sqrt{T}}\right)
 \end{aligned}$$

Proof. :

$$\|w_{t+1} - w^*\|^2 = \|w_t - \gamma \nabla f_t(w_t) - w^*\|^2$$

$$\begin{aligned}
&= \|w_t - w^\star\|^2 - 2\gamma (w_t - w^\star)^T \nabla f_t(w_t) + \underbrace{\gamma^2 \nabla f_t^T(f_t) \nabla f_t(w_t)}_{\|\nabla f_t(w_t)\|^2 \leq G^2} \\
&\Rightarrow (w_t - w^\star)^T \nabla f_t(w_t) \leq \frac{\|w_t - w^\star\|^2 - \|w_{t+1} - w^\star\|^2}{2\gamma} + \frac{\gamma G^2}{2} \\
&\Rightarrow f_t(w_t) - f_t(w^\star) \leq (w_t - w^\star)^T \nabla f_t(w_t) \leq \frac{\|w_t - w^\star\|^2 - \|w_{t+1} - w^\star\|^2}{2\gamma} + \frac{\gamma G^2}{2} \\
&\Rightarrow \frac{1}{T} \sum_{t=1}^T f_t(w_t) - f_t(w^\star) \leq \frac{1}{T} \sum_{t=1}^T \left(\frac{\|w_t - w^\star\|^2 - \|w_{t+1} - w^\star\|^2}{2\gamma} + \frac{\gamma G^2}{2} \right) \\
&\Rightarrow \frac{1}{T} \sum_{t=1}^T f_t(w_t) - f_t(w^\star) \leq \frac{\|w_1 - w^\star\|^2}{2\gamma} + \frac{\gamma G^2}{2}
\end{aligned}$$

where,

$$\begin{aligned}
&\sum_{t=1}^T \|w_t - w^\star\|^2 - \|w_{t+1} - w^\star\|^2 \\
&= \|w_1 - w^\star\|^2 - \|w_2 - w^\star\|^2 + \|w_2 - w^\star\|^2 - \|w_3 - w^\star\|^2 + \dots - \|w_{T+1} - w^\star\|^2 \\
&= \|w_1 - w^\star\|^2 - \|w_{T+1} - w^\star\|^2 \\
&\leq \|w_1 - w^\star\|^2
\end{aligned}$$

□

Theorem 10. $\gamma_t = \frac{1}{\sqrt{t}}, \|w_t\| \leq B, \|\nabla f_t(w)\| \leq G \forall t, w$

Let $w_T^\star = \arg \min_w \sum_{t=1}^T f_t(w)$

$$\frac{1}{T} \sum_{t=1}^T f_t(w_t) - f_t(w_T^\star) \leq \frac{2B^2 + G^2}{\sqrt{T}} \forall T$$

Proof. Similar,

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T f_t(w_t) - f_t(w^\star) &\leq \frac{1}{T} \sum_{t=1}^T \left(\frac{\|w_t - w^\star\|^2 - \|w_{t+1} - w^\star\|^2}{2\gamma_t} + \frac{\gamma_t G^2}{2} \right) \\
&= \frac{1}{T} \sum_{t=1}^T \left(\sqrt{t} \frac{\|w_t - w^\star\|^2}{2} - \sqrt{t} \frac{\|w_{t+1} - w^\star\|^2}{2} + \frac{1}{2\sqrt{t}} G^2 \right) \\
&= \frac{1}{T} \left(\frac{\|w_1 - w^\star\|^2}{2} - \frac{\|w_{T+1} - w^\star\|^2}{2} \right) + \frac{1}{T} \sum_{t=2}^T \frac{\|w_t - w^\star\|^2}{2} (\sqrt{t} - \sqrt{t-1}) + \frac{1}{T} \sum_{t=1}^T \frac{G^2}{2\sqrt{t}} \\
&\leq \frac{1}{T} \left(\frac{\|w_1 - w^\star\|^2}{2} - \frac{\|w_{T+1} - w^\star\|^2}{2} \right) + \frac{1}{T} \sum_{t=2}^T 2B^2 (\sqrt{t} - \sqrt{t-1}) + \frac{1}{T} \sum_{t=1}^T \frac{G^2}{2\sqrt{t}} \\
&\leq \frac{2B^2}{T} + \frac{2B^2}{\sqrt{T}} + \frac{G^2}{\sqrt{T}}
\end{aligned}$$

where,

$$\sum_{t=1}^T \frac{1}{\sqrt{t}} \approx \int_1^T \frac{1}{\sqrt{t}} dt \leq \sqrt{T}$$

□

16 Lecture 16

16.1 Regularization and Bayesian Inference

MLE GLE: $p(y_i | w^T x_i)$

$$\hat{w} = \arg \min_w \sum_{i=1}^n -\log p(y_i | w^T x_i)$$

Gaussian

$$y_i | w^T x_i \sim N(w^T x_i, I)$$

$$\min_w \sum_{i=1}^n (y_i - w^T x_i)^2$$

$$\min_w \|y - Xw\|^2$$

$$X^T y = X^T X \hat{w}$$

system of n equations in d unknowns

$$\{(x_i, y_i)\}_{i=1}^n, n < d$$

underdetermined

if $n < d$,

then exists $v \in \mathbb{R}^d$

such that

$$Xv = 0 \text{ and } v \neq 0$$

$$x_i^T v = 0, i = 1, \dots, n$$

$$X = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}, v = \begin{bmatrix} \alpha \\ -\alpha \\ \alpha \end{bmatrix}$$

$$d = 3, n = 2$$

$$\left(\begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, 1 \right), \left(\begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}, -1 \right), Xv = 0 \forall \alpha \in \mathbb{R}$$

$$X^T y = X^T X (\hat{w} + v) = X^T X \hat{w}$$

\hat{w} is a solution

then so is $\hat{w} + v$

What solution should we use?

$$\hat{w} = \arg \min_{w: x^T y = x^T x w} \|w\|$$

minimum norm solution

$$(x, y), x = x_i + \varepsilon$$

$$w^T x = w^T x_i + w^T \varepsilon \leq \|w\| \|\varepsilon\|$$

16.2 Computing Minimum Norm Solution

1. $\min_w \|w\|^2$ such that $X^T y = X^T X w$
2. $\min_w \|y - Xw\|^2 + \tau \|w\|^2$, tiny $\tau > 0$
3. $X^T X = U \Lambda U^T, \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, 0, 0), (X^T X)^{-1} = U \Lambda^{-1} U^T, \hat{w} = (X^T X)^{-1} X^T y$

$$\begin{aligned} \|y - Xw\|^2 &= \sum (y_i - x_i^T w)^2 \\ &= - \sum \log p(y_i | w^T x_i) \\ \tau \|w\|^2 &= \tau \sum_{j=1}^d w_j^2 \\ &= - \sum_{j=1}^d \log p(w_j) \\ &\Rightarrow - \log p(w_j) \propto \tau w_j^2 \\ p(w_j) &= \exp(-\log p(w_j)) \propto \exp(-\tau w_j^2) \end{aligned}$$

prior

Likelihood

$$y | Xw \sim N(Xw, I) = p(y | w)$$

$$w \sim N\left(0, \frac{1}{2\tau} I\right) = p(w)$$

prior prob for w

$$\begin{aligned} \min_w -\log p(y | w) - \log p(w) \\ = \max_w p(y | w) p(w) \end{aligned}$$

$$\begin{aligned}
&= \max_w \frac{p(y|w)p(w)}{p(y)} \\
&= \max_w p(w|y)
\end{aligned}$$

posterior distribution of w given y

Maximum a Posteriori estimator (MAP)

16.3 Bayesian Inference

$$\begin{aligned}
&x, \{p(x|\theta)\}_{\theta \in \Theta} \\
&x_1, x_2, \dots, x_n \stackrel{iid}{\sim} q
\end{aligned}$$

MLE:

$$\max_{\theta \in \Theta} \prod_{i=1}^n p(x_i|\theta)$$

$\hat{\theta}$ is MLE

roughly, $p_{\hat{\theta}}$ is density that approximately

$$\min_{\theta} \text{KL}(q, p_{\hat{\theta}})$$

We know or believe something a priori about θ

$p(\theta)$ is a weighting function on the Θ

$$\theta_{MAP} = \arg \max_{\theta \in \Theta} \prod_{i=1}^n p(x_i|\theta) p(\theta)$$

Example 10. $x_i \sim \text{Poisson}(\theta), \theta > 0$

$$p(x_i = k|\theta) = e^{-\theta} \frac{\theta^k}{k!}, k = 1, 2, \dots$$

x_i = number of times a "word" appears index i

$$\mathbb{E}[x_i] = \theta$$

MLE:

$$\begin{aligned}
&\min_{\theta} \underbrace{\sum_{i=1}^n (\theta - x_i \log \theta)}_{L(\theta)} \\
&\frac{\partial}{\partial \theta} L(\theta) = n - \frac{\sum x_i}{\theta} \\
&\hat{\theta}_{MLE} = \frac{1}{n} \sum x_i
\end{aligned}$$

unbiased

MAP:

$$p(\theta) = \alpha e^{-\alpha\theta}$$

exponential distribution

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} \prod_i p(x_i|\theta) p(\theta) \\ &= \arg \min_{\theta} \sum_i (\theta - x_i \log \theta + \alpha \theta) \\ &\Rightarrow \sum_i ((1 + \alpha) \theta - x_i \log \theta) \\ &\Rightarrow \hat{\theta}_{MAP} = \frac{1}{(1 + \alpha)n} \sum x_i = \frac{1}{1 + \alpha} \hat{\theta}_{MLE}\end{aligned}$$

iid Poiss(θ)

$$\begin{aligned}\sum_{x_i} \mathbb{E} \left[\sum x_i \right] &= n\theta \\ \mathbb{V} \left[\sum x_i \right] &= \sum \mathbb{V} [X_i] = n\theta \\ \mathbb{V} \left[\frac{1}{n} \sum x_i \right] &= \frac{\theta}{n}\end{aligned}$$

For MLE,

$$\begin{aligned}\mathbb{E} \left[\hat{\theta}_{MLE} \right] &= \theta \\ \mathbb{V} \left[\hat{\theta}_{MLE} \right] &= \frac{\theta}{n}\end{aligned}$$

For MAP,

$$\begin{aligned}\mathbb{E} \left[\hat{\theta}_{MAP} \right] &= \frac{1}{1 + \alpha} \theta \\ \mathbb{V} \left[\hat{\theta}_{MAP} \right] &= \left(\frac{1}{1 + \alpha} \right)^2 \frac{\theta}{n}\end{aligned}$$

MSE

$$\begin{aligned}\mathbb{E} \left[\left(\hat{\theta} - \theta \right)^2 \right] &= \text{Bias}^2 + \text{Variance} \\ &= \mathbb{E} \left[\left(\hat{\theta} - \mathbb{E} \left[\hat{\theta} \right] + \mathbb{E} \left[\hat{\theta} \right] - \theta \right)^2 \right] \\ &= \left(\mathbb{E} \left[\hat{\theta} \right] - \theta \right)^2 + \mathbb{E} \left[\left(\hat{\theta} - \mathbb{E} \left[\hat{\theta} \right] \right)^2 \right]\end{aligned}$$

$$\text{MSE} \left[\hat{\theta}_{MLE} \right] = 0 + \frac{\theta}{n} = \frac{\theta}{n}$$

$$\begin{aligned}\text{MSE} \left[\hat{\theta}_{MAP} \right] &= \left(1 - \frac{1}{1+\alpha} \right)^2 \theta^2 + \left(\frac{1}{1+\alpha} \right)^2 \frac{\theta}{n} \\ \alpha = 1 &\Rightarrow \frac{\theta^2}{4} + \frac{1}{4} \frac{\theta}{n} = \frac{\theta}{4} \left(\theta + \frac{1}{n} \right) \\ \frac{1}{4} \left(\theta + \frac{1}{n} \right) &\gtrless \frac{1}{n}\end{aligned}$$

17 Lecture 17

17.1 Bayesian Inference

prior distribution: $p(\theta)$

likelihood: $p(x|\theta)$

posterior distribution: $p(\theta|x) \propto p(x|\theta)p(\theta)$

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta \in \Theta} p(x|\theta)p(\theta) \\ &= \arg \min_{\theta} (-\log p(x|\theta) - \log p(\theta))\end{aligned}$$

17.2 Linear Bayesian Regression

$$\begin{aligned}\begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix} &= y = Xw + e \\ e &\sim N(0, \sigma^2 I) \\ \hat{\theta}_{LS,MLE} &= (X^T X)^{-1} X^T y \\ y|w &\sim N(Xw, \sigma^2 I) \\ p(y|w) &\propto \exp \left(-\frac{1}{2\sigma^2} \|y - Xw\|^2 \right)\end{aligned}$$

prior:

$$\begin{aligned}w &\sim N(0, \sigma_w^2 I) \\ \mathbb{E}[w] &= 0 \\ \mathbb{E}[\|w\|^2] &= \sigma_w^2 d \\ p(w) &\propto \exp \left(-\frac{1}{2\sigma_w^2} \|w\|^2 \right)\end{aligned}$$

$$\begin{aligned}\hat{w}_{MAP} &= \arg \max_w p(y|w)p(w) \\ &= \arg \max_w \exp \left(-\frac{1}{2\sigma^2} \|y - Xw\|^2 \right) \exp \left(-\frac{1}{2\sigma_w^2} \|w\|^2 \right) \\ &= \arg \min_w \frac{1}{2\sigma^2} \|y - Xw\|^2 + \frac{1}{2\sigma_w^2} \|w\|^2\end{aligned}$$

$$\begin{aligned}
&= \arg \min_w \frac{1}{\sigma^2} \left(\underbrace{\frac{1}{2} \|y - Xw\|^2 + \frac{1}{2} \lambda \|w\|^2}_{f_\lambda(w)} \right), \lambda = \frac{\sigma^2}{\sigma_w^2} > 0 \\
0 &= \frac{\partial f_\lambda(w)}{\partial w} = -X^T (y - Xw) + \lambda w \\
X^T y &= X^T Xw + \lambda w \\
&= (X^T X + \lambda I) w \\
\hat{w}_{MAP} &= (X^T X + \lambda I)^{-1} X^T y
\end{aligned}$$

ridge regression estimate

$$\begin{aligned}
y &= w + e \\
y|w &\sim N(w, \sigma^2) \\
w &\sim N(0, \sigma_w^2) \\
w|y &\sim N\left(\frac{\sigma_w^2}{\sigma^2 + \sigma_w^2} y, \frac{\sigma^2 \sigma_w^2}{\sigma^2 + \sigma_w^2}\right)
\end{aligned}$$

$$\begin{aligned}
p(w|y) &\propto p(y|w) p(w) \\
&\propto \exp\left(-\frac{1}{2\sigma^2} (y - w)^2 - \frac{1}{2\sigma_w^2} w^2\right) \\
&\propto \exp\left(-\frac{\sigma_w^2 (y - w)^2 + \sigma^2 w^2}{2\sigma^2 \sigma_w^2}\right) \\
&= \exp\left(-\frac{w^2 - 2\frac{\sigma_w^2}{\sigma^2 + \sigma_w^2} w + \frac{\sigma_w^2}{\sigma^2 + \sigma_w^2} y^2}{\frac{2\sigma^2 \sigma_w^2}{\sigma^2 + \sigma_w^2}}\right) \\
&\propto \exp\left(-\frac{\left(w - \frac{\sigma_w^2}{\sigma^2 + \sigma_w^2} y\right)^2}{2\frac{\sigma^2 \sigma_w^2}{\sigma^2 + \sigma_w^2}}\right)
\end{aligned}$$

$$\begin{aligned}
\hat{w}_{MAP} &= \frac{\sigma_w^2}{\sigma^2 + \sigma_w^2} y = \frac{\sigma_w^2}{\sigma^2 + \sigma_w^2} \underbrace{\hat{w}_{MLE}}_y \\
\frac{\sigma_w^2}{\sigma^2 + \sigma_w^2} &= \frac{1}{1 + \frac{1}{\lambda}} \\
\lambda &= \frac{\sigma_w^2}{\sigma^2} = \text{SNR}
\end{aligned}$$

signal to noise ratio

17.3 Gauss Markov Theorem

If x, y are jointly Gaussian vector

$$\begin{aligned} \begin{bmatrix} x \\ y \end{bmatrix} &\sim N \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right) \\ y|x &\sim N \left(\mu_y + \Sigma_{yx} \Sigma_{xx}^{-1} (x - \mu_x), \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \right) \\ \mathbb{E} \left[(x - \mu_x) (y - \mu_y)^T \right] &= \Sigma_{xy} \end{aligned}$$

Example 11. $y = Xw + e, e \sim N(0, \sigma^2 I)$

$$\begin{aligned} \Sigma_{ww} &= \mathbb{E} [ww^T] = \sigma_w^2 \\ \Sigma_{wy} &= \mathbb{E} [w (Xw + e)^T] \\ &= \mathbb{E} [ww^T] X^T \\ &= \sigma_w^2 X^T \\ \Sigma_{yy} &= \mathbb{E} [(Xw + e) (Xw + e)^T] \\ &= XX^T \sigma_w^2 + \sigma^2 I \\ w &\sim N(0, \sigma_w^2 I) \\ \begin{bmatrix} y \\ w \end{bmatrix} &\sim N \\ \mathbb{E} [w|y] &= \sigma_w^2 X^T (XX^T + \sigma^2 I)^{-1} y \\ \hat{w}_{MAP} &= \left(X^T X + \frac{\sigma_w^2}{\sigma^2} I \right)^{-1} X^T y \end{aligned}$$

17.4 Weiner Filter

$$X = By + e, e \sim N(0, \sigma^2 I)$$

y is an image

B is blur

e is noise

X blurry noise image

Likelihood,

$$\begin{aligned} X|y &\sim N(By, \sigma^2 I) \\ \hat{y}_{MLE} &= \arg \max_y \frac{1}{2} \|x - By\|^2 = B^{-1} X \end{aligned}$$

if B^{-1} exist,

$$\hat{y}_{MLE} = B^{-1} x = y + B^{-1} e$$

blow up noise

$$\begin{aligned}
y &\sim N(0, \Sigma_{yy}) \\
p(y|x) &\propto \exp\left(-\frac{1}{2\sigma^2}\|x - By\|^2\right) \exp\left(-\frac{1}{2}y^T \Sigma_{yy}^{-1}y\right) \\
\hat{y}_{MAP} &= \mathbb{E}[y|x] = \arg \min_y \underbrace{\frac{1}{2\sigma^2}\|x - By\|^2 + \frac{1}{2}y^T \Sigma_{yy}^{-1}y}_{f(y)} \\
0 &= \frac{\partial f(y)}{\partial y} = \frac{1}{\sigma^2}(-B^T(X - By) + \sigma^2 \Sigma_{yy}^{-1}y) \\
0 &= -B^T X + B^T By + \sigma^2 \Sigma_{yy}^{-1}y \\
B^T X &= (B^T B + \sigma^2 \Sigma_{yy}^{-1})y \\
\hat{y}_{MAP} &= (B^T B + \sigma^2 \Sigma_{yy}^{-1})^{-1} B^T X
\end{aligned}$$

If $\sigma^2 \ll 1$,

$$\hat{y}_{MAP} = (B^T B)^{-1} B^T X = B^{-1} B^{-T} B^T X = B^{-1} X$$

17.5 Linear Regression

$$\begin{aligned}
y &= Xw + e, e \sim N(0, \sigma^2 I) \\
y|w &\sim N(Xw, \sigma^2 I)
\end{aligned}$$

Ridge Prior

$$\begin{aligned}
p(w) &\propto \exp(-\lambda \|w\|^2) \\
&\Rightarrow \hat{w}_{MAP} = (X^T X + \lambda I)^{-1} X^T y
\end{aligned}$$

Prior knowledge: most w_j are 0, $j = 1, \dots, d$

i.e. most features not important

$$\begin{aligned}
&\sum_{j=1}^d \mathbb{1}_{\{w_j \neq 0\}} \text{ is small } \ll d \\
p(w) &\propto \exp\left(-\lambda \sum_{j=1}^d \mathbb{1}_{\{w_j \neq 0\}}\right), \lambda > 0 \\
p(w|y) &\propto \exp\left(-\frac{1}{2\sigma^2}\|y - Xw\|^2\right) \exp\left(-\lambda \sum_{j=1}^d \mathbb{1}_{\{w_j \neq 0\}}\right) \\
\hat{w}_{MAP} &= \arg \min_w \left(\frac{1}{2\sigma^2}\|y - Xw\|^2 + \lambda \sum_{j=1}^d \mathbb{1}_{\{w_j \neq 0\}}\right)
\end{aligned}$$

hard to optimize

$$p(w) \propto \exp\left(-\lambda \sum_{j=1}^d |w_j|\right)$$
$$\hat{w}_{MAP} = \arg \min_w \left(\frac{1}{2\sigma^2} \|y - Xw\|^2 + \lambda \|w\|_1 \right)$$

easy to optimize

LASSO

18 Lecture 18

18.1 Linear Prediction

$$\hat{y} = w^T x$$
$$\min_w \sum_{i=1}^n \underbrace{l(y_i, w^T x_i)}_{-\log \text{like}} + \underbrace{\lambda \text{pen}(w)}_{-\log \text{prior}}$$
$$p(y|w) \propto \exp(-l(y_i, w^T x_i))$$

18.2 Loss Functions

$$\text{sq err: } \|y - Xw\|_2^2 = \sum_{i=1}^n (y_i - w^T x_i)^2$$

$$\text{abs err: } \|y - Xw\|_1 = \sum_{i=1}^n |y_i - w^T x_i|$$

$$0-1 \text{ loss: } \sum_{i=1}^n \mathbb{1}_{\{y_i - w^T x_i < 0\}}$$

$$\text{logistic loss: } \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

$$\text{hinge: } \sum_{i=1}^n \max\{0, 1 - y_i w^T x_i\}$$

18.3 Penalties, Regularizers

$$2 \text{ norm: } \|w\|^2$$

$$1 \text{ norm: } \|w\|_1 = \sum_{j=1}^d |w_j|$$

$$\text{ideal pen: } \sum_{j=1}^d \mathbb{1}_{\{w_j \neq 0\}}$$

LASSO

$$\min_w \frac{1}{2} \|y - Xw\|^2 + \lambda \|w\|_1$$
$$d = n, X = I$$

$$\begin{aligned}
& \min_w \sum_{i=1}^n \frac{1}{2} (y_i - w_i)^2 + \lambda |w_i| \\
& \min_{w_i} \underbrace{\frac{1}{2} (y_i - w_i)^2 + \lambda |w_i|}_{f(w_i)} \\
& 0 = \frac{\partial f}{\partial w_i} = -y_i + w_i + \lambda \text{sign}(w_i), w_i \neq 0 \\
& w_i = y_i - \lambda \text{sign}(w_i) \\
& \Rightarrow \text{sign}(w_i) = \text{sign}(y_i)
\end{aligned}$$

$$\begin{aligned}
\hat{w}_i &= \begin{cases} y_i - \lambda \text{sign}(y_i) \\ 0 \end{cases} \\
f(\hat{w}_i) &= \begin{cases} \frac{\lambda^2}{2} + \lambda |y_i - \lambda \text{sign}(y_i)| \\ \frac{y_i^2}{2} \end{cases}
\end{aligned}$$

$$\begin{aligned}
|y_i| < \lambda &\Rightarrow \hat{w}_i = 0 \\
|y_i| > \lambda &\Rightarrow \hat{w}_i = y_i - \lambda \text{sign}(y_i) = \text{sign}(y_i) (|y_i| - \lambda)
\end{aligned}$$

soft threshold function

w_i unknown, $\varepsilon_i \sim N(0, 1)$

$$\begin{aligned}
y_i &\sim N(w_i, 1) \\
y_i &= w_i + \varepsilon_i
\end{aligned}$$

1. MLE

$$\begin{aligned}
& \min_w \frac{1}{2} \|y - w\|^2 \\
& \Rightarrow \hat{w}_i = y_i \\
& \mathbb{E} [\|\hat{w} - w\|^2] = n
\end{aligned}$$

2. G-MAP

$$\begin{aligned}
& \min_w \frac{1}{2} \|y - w\|^2 + \lambda \|w\|^2 \\
& \hat{w}_i = \frac{1}{1 + \lambda} y_i \\
& \mathbb{E} [\|\hat{w} - w\|^2] = \left(\frac{\lambda}{1 + \lambda} \right)^2 \|w\|^2 + \left(\frac{1}{1 + \lambda} \right)^2 n
\end{aligned}$$

3. Soft threshold

$$\hat{w}_i = \text{sign}(y_i) (|y_i| - \lambda)_+$$

4. Oracle, $\sigma^2 = 1$

$$\hat{w}_i = \begin{cases} 0 & \text{if } |w_i| \leq 1 \\ y_i & \text{if } |w_i| > 1 \end{cases}$$

$$\mathbb{E} [\|\hat{w}_o - w\|^2] = \sum_{i=1}^n \min \{w_i^2, 1\}$$

Theorem 11. Assume $y_i \sim N(w_i, 1)$, $\hat{w}_i = \text{sign}(y_i) (|y_i| - \lambda)_+$, and take $\lambda = \sqrt{2 \log n}$, then

$$\mathbb{E} [\|\hat{w} - w\|^2] \leq (2 \log n + 1) \left(1 + \sum_{i=1}^n \min \{w_i^2, 1\} \right)$$

Example 12. $k < n, w_i \neq 0$, then,

$$\begin{aligned} \sum_{i=1}^n \min \{w_i^2, 1\} &= k \\ (2 \log n + 1) \left(1 + \sum_{i=1}^n \min \{w_i^2, 1\} \right) &= O(k \log n) \end{aligned}$$

$$y \sim N(0, 1)$$

$$\begin{aligned} \mathbb{P} \{y \geq \lambda\} &= \frac{1}{\sqrt{2\pi}} \int_{\lambda}^{\infty} e^{-\frac{y^2}{2}} dy \\ &\leq \frac{1}{2} e^{-\frac{\lambda^2}{2}} \end{aligned}$$

$$\mathbb{P} \{|y_i| > \lambda\} \leq e^{-\frac{\lambda^2}{2}}$$

$$\begin{aligned} \mathbb{P} \{|y_{i,1}| \geq \lambda \text{ or } |y_{i,2}| \geq \lambda \text{ or } \dots |y_{i,n}| \geq \lambda\} &\leq n e^{-\frac{\lambda^2}{2}} \\ &= e^{-\frac{1}{2}(\lambda^2 - 2 \log n)} \end{aligned}$$

union bound

19 Lecture 19

$$\min_w \|y - Xw\|^2 + \lambda \underbrace{\text{pen}(w)}_{\|w\|^2 \text{ or } \|w\|_1}$$

$$\min_w (y - w)^2 + \lambda \text{pen}(w)$$

has a closed-form solution

$$\begin{aligned}\hat{w}_i &= y_i - \text{sign}(y_i) \min\{|y_i|, \lambda\} \\ L(w) &= \|y - Xw\|^2 + \lambda \text{pen}(w)\end{aligned}$$

iterates w_1, w_2, \dots

$$\begin{aligned}L(w_1) &\geq L(w_2) \geq L(w_3) \dots \\ L(w) &= \|y - Xw_k + Xw_k - Xw\|^2 + \lambda \text{pen}(w) \\ &= \|y - Xw_k\|^2 + \underbrace{\|X(w_k - w)\|^2 + 2(y - Xw_k)^T X(w_k - w) + \lambda \text{pen}(w)}_{\text{choose } w \text{ so this } \boxed{x} \leq \lambda \text{pen}(w_k)} \\ w_{k+1} &= \arg \min_w \boxed{x} \\ \boxed{x} &= \|X(w_k - w)\|^2 + 2(y - Xw_k)^T X(w_k - w) + \lambda \text{pen}(w) \\ &\leq \underbrace{\|X\|^2}_{\leq \gamma^{-1}} \|w_k - w\|^2 + 2 \underbrace{(y - Xw_k)^T X}_{\frac{1}{\gamma} V_k^T} (w_k - w) + \lambda \text{pen}(w), V_k = \gamma X^T (y - Xw_k) \\ \gamma \boxed{x} &\leq \|w_k - w\|^2 + 2V_k^T (w_k - w) + \gamma \lambda \text{pen}(w) \\ &= \|V_k + w_k - w\|^2 - \|V_k\|^2 + \gamma \lambda \text{pen}(w)\end{aligned}$$

then the minimization problem is

$$\begin{aligned}\min_w &\|V_k + w_k - w\|^2 + \gamma \lambda \text{pen}(w) \\ z_k &= w_k + \gamma X^T (y - Xw_k)\end{aligned}$$

GD iterate

$$\begin{aligned}\min_w &\|z_k - w\|^2 + \gamma \lambda \text{pen}(w) \\ \min_w &\sum_{j=1}^d \left((z_{k_j} - w_j)^2 + \gamma \lambda |w_j| \right).\end{aligned}$$

Goal:

$$\min_w \|y - Xw\|^2 + \lambda \|w\|_1$$

w_1 init

$$\begin{aligned}k &= 1, 2, \dots \\ z_k &= w_k + \gamma X^T (y - Xw_k) \\ \hat{w}_{k+1} &= z_k - \text{sign}(z_k) \min\{|z_k|, \gamma \lambda\} \\ \min_w L(w) &= \|y - Xw\|^2 + \lambda \|w\|^2 \\ \hat{w} &= (X^T X + \lambda I)^{-1} X^T y\end{aligned}$$

$$X_{n \times d}, w \in \mathbb{R}^d$$

$$\frac{\partial L}{\partial w} = -2X^T (y - Xw) + 2\lambda w = 0$$

$$X^T (y - Xw) = \lambda w$$

solution have form

$$w = \frac{1}{\lambda} X^T \underbrace{(y - Xw)}_{\lambda \alpha}$$

form of solution

$$w = X^T \alpha, \alpha \in \mathbb{R}^n$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \dots \\ x_n^T \end{bmatrix}$$

$$X^T = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}$$

$$w = \sum_{i=1}^n \alpha_i x_i \text{ for } \alpha_i \in \mathbb{R}$$

weights are linear combo of features

19.1 Dual Opt

$$\min_{\alpha \in \mathbb{R}^n} \|y - XX^T \alpha\|^2 + \lambda \|X^T \alpha\|^2$$

$$L(\alpha) = \|y - XX^T \alpha\|^2 + \lambda \alpha^T XX^T \alpha$$

$$\frac{\partial L}{\partial \alpha} = 2XX^T (y - XX^T \alpha) + 2\lambda XX^T \alpha$$

$$= 2XX^T (-y + XX^T \alpha + \lambda \alpha) = 0$$

$$y = (XX^T + \lambda I) \alpha$$

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

$$\hat{w} = X^T \hat{\alpha}$$

$$K = XX^T$$

kernel matrix

$$K_{ij} = K(X_i, X_j) = X_i^T X_j$$

$$\hat{w} = (X^T X + \lambda I)^{-1} X^T y$$

$$= X^T (XX^T + \lambda I)^{-1} y$$

$$X = UDV^T$$

19.2 Using nonlinear features

$$\begin{aligned}
 X &= \begin{bmatrix} x_1 \\ \dots \\ x_d \end{bmatrix} \in \mathbb{R}^d \\
 \Phi(X) &= \begin{bmatrix} \phi_1(x) \\ \dots \\ \phi_D(x) \end{bmatrix} \in \mathbb{R}^D \\
 \Phi &= \begin{bmatrix} \Phi(x_1)^T \\ \dots \\ \Phi(x_n)^T \end{bmatrix} \\
 \min_w &\|y - \Phi w\|^2 + \lambda \|w\|^2 \\
 \hat{w} &= (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y \\
 \hat{\alpha} &= \left(\underbrace{\Phi \Phi^T}_{K_\Phi} + \lambda \Phi \right)^{-1} y \\
 K_\Phi(i, j) &= \Phi(x_i)^T \Phi(x_j)
 \end{aligned}$$

Example 13. $x \in \mathbb{R}^2$

$$\begin{aligned}
 \Phi(x) &= \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} \in \mathbb{R}^3 \\
 K_\Phi(i, j) &= \Phi(x_i)^T \Phi(x_j) \\
 &= x_{i1}^2 x_{j1}^2 + 2x_{i1}x_{i2}x_{j1}x_{j2} + x_{i2}^2 x_{j2}^2 \\
 &= (x_{i1}x_{j1} + x_{i2}x_{j2})^2 \\
 &= (x_i^T x_j)^2
 \end{aligned}$$

Example 14. Other kernels,

$$\begin{aligned}
 K_\Phi(i, j) &= (x_i^T x_j + 1)^2 \\
 K_\Phi(i, j) &= \exp(-\lambda \|x_i - x_j\|^2) \\
 \hat{y} &= \sum_{i=1}^n \hat{\alpha}_i K_\Phi(x_i, x)
 \end{aligned}$$

20 Lecture 20

20.1 GM Thm

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

$$y|x \sim N(\mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(x - \mu_x), \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy})$$

If $\mu_x = \mu_y = 0$, MAP is $\sigma_{xy}\sigma_{xx}^{-1}x$

20.2 Kernels

$$\min_{w \in \mathbb{R}^d} \|y - Xw\|^2 + \lambda \|w\|^2$$

$$X = \begin{bmatrix} x_1^T \\ \dots \\ x_n^T \end{bmatrix}, n \times d$$

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$$

$$\phi(x) = \begin{bmatrix} \phi_1(x) \\ \dots \\ \phi_D(x) \end{bmatrix}$$

$$\min_{w \in \mathbb{R}^D} \|y - \Phi w\|^2 + \lambda \|w\|^2$$

$$\Phi = \begin{bmatrix} \phi(x_1)^T \\ \dots \\ \phi(x_n)^T \end{bmatrix}, n \times D$$

$$\Phi^T(y - \Phi w) = \lambda w$$

$$\hat{w} = \Phi^T \alpha, \alpha \in \mathbb{R}^n$$

$$\hat{w} = \sum_{i=1}^n \alpha_i \phi(x_i)$$

$$\min_{\alpha} \|y - \Phi(\Phi^T \alpha)\|^2 + \underbrace{\lambda \|\Phi^T \alpha\|^2}_{\lambda \alpha^T \Phi \Phi^T \alpha, K = \Phi \Phi^T}$$

$$\min_{\alpha} \|y - K\alpha\|^2 + \lambda \alpha^T K \alpha$$

$$K_{ij} = K(x_i, x_j) = \phi^T(x_i) \phi(x_j)$$

$$K = \Phi \Phi^T$$

$$= \begin{bmatrix} \phi^T(x_1) \\ \dots \\ \phi^T(x_n) \end{bmatrix} [\phi^T(x_1) \dots \phi^T(x_n)]$$

$$K^T = (\Phi \Phi^T)^T$$

$$= \Phi \Phi^T$$

$$2K(y - K\alpha) = 2\lambda K\alpha$$

$$y - K\alpha = \lambda \alpha$$

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

get new x , predict

$$\begin{aligned}
 \hat{y} &= \hat{w}^T \phi(x) \\
 &= (X^T \hat{\alpha})^T \phi(x) \\
 &= \left(\sum_{i=1}^n \hat{\alpha}_i \phi(x_i) \right)^T \phi(x) \\
 &= \sum_{i=1}^n \hat{\alpha}_i \underbrace{\phi^T(x_i) \phi(x)}_{K(x_i, x)}
 \end{aligned}$$

Example,

$$\begin{aligned}
 x &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2, \phi(x) = [1, x_1, x_2, x_1^2, x_1 x_2, x_2^2]^T \in \mathbb{R}^6 \\
 \phi^T(x) \phi(z) &= 1 + x_1 z_1 + x_2 z_2 + x_1^2 z_1^2 + x_1 x_2 z_1 z_2 + x_2^2 z_2^2 \\
 (1 + x^T z)^2 &= 1 + 2x^T z + (x^T z)^2 \\
 &= 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2 \\
 K(x_i, x_j) &= (1 + x_i^T x_j)^2 \\
 (1 + x^T z)^l &= \sum_{k=0}^l \binom{l}{k} (1)^{l-k} (x^T z)^k \\
 &= \sum_{k=0}^l \binom{l}{k} (x^T z)^k \\
 D &\gg d \\
 D &\gg n
 \end{aligned}$$

If $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$\begin{aligned}
 \phi(x) &= \{x_1^p x_2^q\}, p, q \geq 0, p + q \leq l \\
 \phi(x) &= \{x_1^{p_1} x_2^{p_2} \dots x_d^{p_d}\}, 0 \leq p_1 + p_2 + \dots p_d \leq l
 \end{aligned}$$

Gaussian kernel

$$\begin{aligned}
 K(x_i, x_j) &= \exp(-\beta \|x_i - x_j\|^2) \\
 &= \phi^T(x_i) \phi(x_j)
 \end{aligned}$$

infinite dim kernel space

$$\begin{aligned}
 e^x &= 1 + \frac{x}{1!} + \frac{x^2}{2!} + \dots \\
 e^{-x^2} &= 1 - \frac{x^2}{1!} + \frac{x^4}{2!} \\
 \hat{\alpha} &= (K + \lambda I)^{-1} y
 \end{aligned}$$

"nonparametric" learning

21 Lecture 21

$$\begin{aligned} & \{(x_i, y_i)\}_{i=1}^n, \phi : \mathbb{R}^d \rightarrow \mathbb{R}^D, x_i \in \mathbb{R}^d \\ & \min_{w \in \mathbb{R}^D} \|y - \Phi w\|^2 + \lambda \|w\|^2 \\ & \Phi = \begin{bmatrix} \phi^T(x_1) \\ \dots \\ \phi^T(x_n) \end{bmatrix} \\ & \hat{y} = \phi(x_i)^T w \\ & = w^T \phi(x_i) \\ & \Phi^T (y - \Phi w) = \lambda w \\ & w_\lambda = \frac{1}{\lambda} \Phi^T \underbrace{(y - \Phi w)}_{\alpha \in \mathbb{R}^n} \\ & w_\lambda = \Phi^T \alpha = \sum_{i=1}^n \alpha_i \phi(x_i) \end{aligned}$$

Dual:

$$\begin{aligned} & \min_{\alpha \in \mathbb{R}^n} \|y - \Phi \Phi^T \alpha\|^2 + \lambda \alpha^T \phi \phi^T \alpha \\ & (K + \lambda I) \alpha_\lambda = y \end{aligned}$$

Kernal trick:

$$\begin{aligned} \alpha_\lambda &= (K + \lambda I)^{-1} y \\ K &= \Phi \Phi^T, K_{ij} = \phi(x_i)^T \phi(x_j) \end{aligned}$$

21.1 Representer Theorem

$$\lambda > 0$$

$$w_\lambda = \arg \min_{w \in \mathbb{R}^D} \sum_{i=1}^n l(y_i w^T \phi(x_i)) - \lambda \|w\|^2$$

then

$$w_\lambda = \sum \alpha_i \phi(x_i)$$

Classifier: for a new x ,

$$\begin{aligned} \hat{y} &= \text{sign}(\hat{w} \phi(x)) \\ &= \text{sign}\left(\sum \alpha_i \phi(x_i)^T \phi(x)\right) \end{aligned}$$

$$\begin{aligned}
&= \text{sign} \left(\sum \alpha_i K(x_i, x) \right) \\
K(x, x') &= (x^T x' + 1)^k \\
\min_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n l \left(y_i \left(\sum_{j=1}^n \alpha_j \left(\phi(x_j)^T \right) \phi(x_i) \right) \right) &+ \lambda \sum \alpha_i \alpha_j \underbrace{\phi^T(x_i) \phi(x_j)}_{K(x_i, x_j)}
\end{aligned}$$

Proof. WLOG, $w_\lambda \in \mathbb{R}^D$

$$w_\lambda = \sum_{j=1}^n \alpha_j \phi(x_j) + u$$

where $u^T \phi(x_i) = 0, i = 1, \dots, n$

$$\begin{aligned}
\|w_\lambda\|^2 &= \left\| \sum \alpha_i \phi(x_j) + u \right\|^2 \\
&= \left\| \sum \alpha_i \phi(x_j) \right\|^2 + \|u\|^2 \\
&= \left\| \sum \alpha_i \phi(x_j) \right\|^2
\end{aligned}$$

□

21.2 Kernels

$$\begin{aligned}
K(x, x') &= (x^T x' + 1)^k \\
x \in \mathbb{R}^d, D &= \begin{bmatrix} d + K \\ K \end{bmatrix}
\end{aligned}$$

eg. $d = 10, k = 4, D = 1001$

$$\text{sign} \left(\sum \alpha_i K(x_i, x) \right)$$

what is this? multi-dim poly

21.3 Infinite Dim Feature Spaces

$$\begin{aligned}
l_p &= \left\{ (\beta_1, \beta_2, \dots) : \sum_{i \geq 1} \beta_i^p < \infty \right\} \\
p &= 2, l_2
\end{aligned}$$

Suppose we have a sequence of features $\{\phi_i(x)\}_{i \geq 1} \in l_2$

Define

$$\begin{aligned}
K(x, x') &= \sum_{i \geq 1} \phi_i(x) \phi_i(x') \\
&= \langle \phi(x), \phi(x') \rangle
\end{aligned}$$

l_2 inner product, symm

$$\begin{aligned} |K(x, x')|^2 &= |\langle \phi(x), \phi(x') \rangle|^2 \\ &\leq \|\phi(x)\|^2 \|\phi(x')\|^2 \\ &< \infty \end{aligned}$$

Cauchy-Schwarz

21.4 Taylor Series Kernel

$$\begin{aligned} f(z) &= \sum_{j=0}^{\infty} a_j z^j, z \in \mathbb{R} \\ e^z &= \sum_{j=0}^{\infty} \left(\frac{z^j}{j!} \right) < \infty \end{aligned}$$

Exp kernel:

$$K(x, x') = \exp(x^T x')$$

21.5 Gaussian Kernel

$$\begin{aligned} K(x, x') &= \exp\left(-\frac{\|x - x'\|^2}{\sigma^2}\right) \\ &= \exp\left(-\frac{\|x\|^2}{\sigma^2}\right) \exp\left(-\frac{\|x'\|^2}{\sigma^2}\right) \exp\left(-\frac{2x^T x'}{\sigma^2}\right) \end{aligned}$$

21.6 Kernels

$$\begin{aligned} \{\phi_j(x)\}_{j \geq 1} &\in l_2, K(x, x') = \langle \phi(x), \phi(x') \rangle \\ f(x) &= \sum_{i=1}^N \alpha_i K(x_i, x), x_i \in \mathbb{R}^d \\ \hat{y} &= \text{sign}(f(x)) \\ f \in \mathcal{H} &= \left\{ f : f(x) = \sum_{i=1}^N \alpha_i K(x_i, x) \right\} \end{aligned}$$

Hilbert Space

$$\begin{aligned} f, g &\in \mathcal{H} \\ \alpha f + \beta g &\in \mathcal{H} \\ \alpha, \beta &\in \mathbb{R} \end{aligned}$$

Lemma 3. $\mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a valid kernel iff it is symmetric and PSD for any $N \geq 1$ and any $x_1, \dots, x_N \in \mathbb{R}^d$. Gram matrix $K_{ij} = K(x_i, x_j)$ is symm PSD.

Proof. $x_1 \dots x_N$

$$K(x, x') = \phi^T(x) \phi(x')$$

$$K, N \times N$$

$$v \in \mathbb{R}^N$$

$$\begin{aligned} v^T K v &= \sum_{i,j=1}^N v_i v_j \underbrace{K(x_i, x_j)}_{\phi(x_i) \phi(x_j)} = \left(\sum_i v_i \phi(x_i) \right)^T \left(\sum_j v_j \phi(x_j) \right) \\ &= \left\| \sum v_i \phi(x_i) \right\|^2 \geq 0 \end{aligned}$$

□

Proof. K is PSD $\Rightarrow K(x, x')$ is an inner product

define $\phi(x)$ to be $X \rightarrow K(\cdot, x)$, "canonical feature map"

$$\mathcal{H} = \left\{ f : f(\cdot) = \sum_{i=1}^N \alpha_i K(\cdot, x_i) \right\}, x_i \in \mathbb{R}^d$$

□

Define inner product,

$$\begin{aligned} \langle f, g \rangle &= \langle \sum \alpha_i K(\cdot, x_i), \sum \beta_j K(\cdot, x_j) \rangle \\ &= \sum_{i,j=1}^N \alpha_i \beta_j K(x_i, x_j) \\ &= \alpha^T K \beta \end{aligned}$$

this is valid

$$\min_{w \in \mathbb{R}^D} \sum_{i=1}^n l(y_i w^T \phi(x_i)) + \lambda \|w\|^2$$

What if $D > n$?

$$\begin{aligned} K(x, x') &= \exp\left(-\frac{\|x - k'\|^2}{\sigma^2}\right) \\ \min_w \|y - \Phi w\|^2 + \lambda \|w\|^2 \\ \Phi &= \begin{bmatrix} \phi^T(x_1) \\ \dots \\ \phi^T(x_n) \end{bmatrix}, n \times D, D \gg n \\ w_\lambda &= (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y \\ w_0 &= \lim_{\lambda \rightarrow 0} w_\lambda = \underbrace{\Phi^T}_{D \times n} \underbrace{(\Phi \Phi^T)^{-1}}_{n \times n} y \end{aligned}$$

pseudo inverse

$$\Phi^T = \begin{bmatrix} \phi(x_1) & \dots & \phi(x_n) \end{bmatrix}$$

and if $D > n$,

$$\|y - \underbrace{\Phi w_0}_{\hat{y}}\|^2 = 0$$

min norm solution

$$w_0 = \sum_{i=1}^n \alpha_i \phi(x_i)$$

21.7 Laplacian Kernel

$$K(x, x') = \exp(-\beta \|x - x'\|)$$

$$\beta \sim n$$

22 Lecture 22

22.1 SVMs

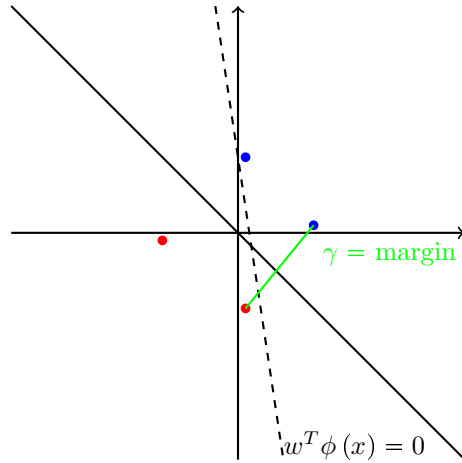
$$\min_w \sum_{i=1}^n (1 - y_i w^T \phi(x_i))_+ + \lambda \|w\|^2$$

$$\min_{\alpha} \sum_{i=1}^n \left(1 - y_i \sum_{j=1}^n \alpha_j K(x_i, x_j) \right)_+ \lambda \alpha^T K \alpha$$

Dual:

$$K(x, x') = \phi(x)^T \phi(x')$$

$$w^T \phi(x) = \sum_{i=1}^n \alpha_i K(x_i, x)$$



solid line is max margin separator
dashed line is separator but not max margin
sum of hinge losses

$$H(w) = \sum_{i=1}^n (1 - y_i w^T \phi(x_i))_+$$

linearly separable $\Rightarrow \exists w$ such that $H(w) = 0$

max-margin separator
 w with smallest norm
and $y_i w^T \phi(x_i) \geq 1$ for all i

22.2 max-margin opt

$$\min_w \|w\|^2 \text{ such that } y_i w^T \phi(x_i) \geq 1$$

$$\text{such that } \sum_{i=1}^n (1 - y_i w^T \phi(x_i))_+ = 0$$

Lagrangian form:

$$\min_w \sum_{i=1}^n (1 - y_i w^T \phi(x_i))_+ + \underbrace{\lambda}_{\text{Lagrange mult}} \|w\|^2$$

for $\lambda > 0$ tiny

22.3 Perceptron

$$\{(x_i, y_i)\} \hat{y}_i = \text{sign}(w^T x_i)$$

$$w_1 = \text{init}$$

$$w_{t+1} = w_t + \mu \underbrace{(y_{i_t} - \hat{y}_{i_t})}_{\text{error}} x_{i_t}, t \geq 1$$

$$y = 1, w^T x > 0$$

$$y - \hat{y} = \begin{cases} 0 & \text{if } y = 1, w^T x > 0 \\ 0 & \text{if } y = -1, w^T x < 0 \\ 2 & \text{if } y = 1, w^T x < 0 \\ -2 & \text{if } y = -1, w^T x > 0 \end{cases}$$

$$y - \hat{y} = \begin{cases} 0 & \text{if } y w^T x > 0 \\ 2y & \text{if } y w^T x < 0 \end{cases}$$

$$\begin{aligned} w_{t+1} &= w_t + \mu (2y_{i_t}) \mathbb{1}_{\{y_{i_t} w_t^T x_{i_t} < 0\}} x_{i_t} \\ &= w_t + 2\mu \mathbb{1}_{\{y_{i_t} w_t^T x_{i_t} < 0\}} y_{i_t} x_{i_t} \end{aligned}$$

SGD with loss function l with derivative l'

$$\begin{aligned} w_{t+1} &= w_t + \gamma (-l' (y_{i_t} w_t^T x_{i_t})) y_{i_t} x_{i_t} \\ &= w_t - \gamma \frac{\partial l}{\partial w} \Big|_{w=w_t} \\ \gamma &= 2\mu \end{aligned}$$

$$-l' (z) = \mathbb{1}_{\{z < 0\}}$$

$$\hat{y} = f (w^T x), \text{ single layer perceptron}$$

22.4 Multilayer Neural Network

$$\hat{y} = W_L f (W_{L-1} \dots (W_2 f (W_1 x + b_1) + b_2) + \dots + b_{L-1}) + b_L$$

Wx affine linear map

$f(\cdot)$ nonlinear coordinate-wise

$$f(v) = \begin{bmatrix} f(v_1) \\ \dots \\ f(v_n) \end{bmatrix}, v \in \mathbb{R}^n, \text{ "activation" function}$$

$$\min_{(w_j, b_j)_{j=1}^L} \sum_{i=1}^n l(y_i \hat{y}_i(\{w_j, b_j\}))$$

22.5 Two-Layer Neural Net

$$\hat{y} = W_2 f (W_1 x + b_1) + b_2$$

linear in W_2, b_2 , but nonlinear in w_1, b_1

22.6 Kernel Machine = 2 Layer Net

$$\hat{y} = \sum_{i=1}^n \alpha_i K(x_i, x)$$

linear in parameter

$$W_2 = \begin{bmatrix} \alpha_1 & \dots & \alpha_n \end{bmatrix}$$

$$K(x_i, x) = f(x_i^T + b_i)$$

Example 15. $(w_i^T x + 1)^k, f(\cdot) = (\cdot)^k$

Example 16. $\exp(x_i^T x), f(\cdot) = \exp(\cdot)$

Example 17. $\exp\left(-\frac{1}{2}\|x_i - x\|^2\right) = \exp\left(x_i^T x - \underbrace{\frac{1}{2}(\|x\|^2 + \|x_i\|^2)}_{b_1}\right)$

$$W_1^T = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}$$

$$\hat{y} = W_2 f(W_1 x + b_1)$$

$$\begin{bmatrix} x_1^T x \\ \dots \\ x_n^T x \end{bmatrix}$$

22.7 Difference

W_1, b_1 fixed for kernel machine

layer 1 width = n

Kernel machines are neural nets

Variable parameter for neural net

layer 2 width = anything

But note all neural nets are kernel machines

23 Lecture 23

23.1 Two-Layer NNs

$$y = W_2 f(W_1 x + b_1)$$

W_1, b_1 are fixed in kernel

23.2 Convolutions NNs

1 network layer (in general)

$$y = f(w_i^T x + b_i)$$

$$W_1 = \begin{bmatrix} w_1^T \\ w_2^T \\ \dots \\ w_n^T \end{bmatrix}$$

CNN

$$y_i = f(\|w_i \star x\|_{\infty} - b_i)$$

convolution

$$(w_i \star x)_k = \sum_{j=1}^d w_{ij} x_{k-j}$$

max pooling: max output of conv

$$\|w_i \star x\|_{\infty}$$

$$f(z) = \max\{0, z\}$$

23.3 Backprop = SGD

23.4 Stone-Weierstrauss Thm (informal)

any continuous function $[0, 1]^d$ can be approximated point-wise to arbitrary accuracy with a polynomial.

Theorem 12. *for any continuous f on $[0, 1]^d$ there exists a neural net*

$$g(x) = W_2 f(W_1 x)$$

$$W_1 \in \mathbb{R}^{D \times d}$$

$$f(u_i^T x) = (u_i^T + 1)^k$$

$$W_1 = \begin{bmatrix} u_1^T \\ \dots \\ u_n^T \end{bmatrix}$$

for k and n sufficiently large, $n = D$ and

$$u_i \in \mathbb{R}^d$$

$$u_i \stackrel{iid}{\sim} p$$

where p is any continuous density on $[0, 1]^d$

Example:

$$\begin{aligned}
g : \mathbb{R}^d &\rightarrow \mathbb{R} \\
g(x) &= v^T f(W, x) \\
&= v^T \begin{bmatrix} (v_1^T x + 1)^k \\ \dots \\ (v_n^T x + 1)^k \end{bmatrix} \\
&= v^T \begin{bmatrix} \phi(u_1)^T \phi(x) \\ \dots \\ \phi(u_n)^T \phi(x) \end{bmatrix} \\
\phi(x) &\in \mathbb{R}^D, D = \binom{d+k}{k} \\
g(x) &= \sum_{i=1}^n v_i \phi(u_i)^T \phi(x) \\
&= (\Phi^T v)^T \phi(x) \\
&= w^T \phi(x) \\
\Phi^T &= \begin{bmatrix} \phi(u_1) & \dots & \phi(u_D) \end{bmatrix}
\end{aligned}$$

general polynomial

If Φ^T is invertible

$$v = (\Phi^T)^{-1} w$$

Lemma 4. If $h : \mathbb{R}^d \rightarrow \mathbb{R}$ is a polynomial function $\neq 0$ and if $u \in \mathbb{R}^d$ is a continuous random point, then $\mathbb{P}\{h(u) = 0\} = 0$

Intuition:

$$\phi(u_1) \in \mathbb{R}^D$$

v_1, \dots, v_{D-1} be a basis for orthogonal subspace

$$\begin{aligned}
&\phi(u_2) \\
&\mathbb{P}\{v_j^T \phi(u_2) = 0\} = 0 \\
&\Rightarrow \phi(u_2) \text{ is linearly independent of } \phi(u_1)
\end{aligned}$$

24 Lecture 24

24.1 Probably Approx Correct Learning

\mathcal{X} = feature space, $\mathcal{X} = \mathbb{R}^d, \hat{y} = f(x), f \in \mathcal{F}$

\mathcal{Y} = label space, $\mathcal{Y} = \{-1, +1\}$

\mathcal{F} = hypothesis space, $\mathcal{F} = \{\text{linear classifiers}\}$

l = loss function, $l : \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{R}_+$

24.2 Empirical Risk Minimization (ERM)

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i))$$

$$\{(x_i, y_i)\}_{i=1}^n \stackrel{iid}{\sim} P$$

Emp Risk:

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i))$$

True Risk:

$$R(f) = \mathbb{E}_{(x_i, y_i) \sim P} [l(y_i, f(x_i))]$$

$$f^* = \arg \min_{f \in \mathcal{F}} R(f)$$

Want $\hat{R}(f) \approx R(f)$ for all f , uniformly close

$$\begin{aligned} \mathbb{P} \left\{ |\hat{R}(f) - R(f)| \geq t \right\} &\leq \mathbb{P} \left\{ |\hat{R}(f) - R(f)|^2 \geq t^2 \right\} \\ &\leq \frac{\mathbb{E} \left[\left(\hat{R}(f) - R(f) \right)^2 \right]}{t^2} \\ &= \frac{\mathbb{V} \left[\hat{R}(f) \right]}{t^2} \\ &\leq \frac{c^2}{4nt^2} \end{aligned}$$

Markov's Inequality

$$\hat{R}(f) = \frac{1}{n} \left(\sum_{i=1}^n \underbrace{l(y_i, f(x_i))}_{\text{bounded by } c} \right)$$

$$0 \leq l \leq c$$

$l = 0$ with probability $\frac{1}{2}$
 $l = c$ with probability $\frac{1}{2}$

$$\Rightarrow \mathbb{V}[l] \leq \frac{c^2}{4}$$

24.3 Chernoff's Bound

$$\mathbb{P} \left\{ \hat{R}(f) - R(f) \geq t \right\}$$

$$\begin{aligned}
\mathbb{P} \left\{ e^{\lambda(\hat{R}(f) - R(f))} \geq e^{\lambda t} \right\} &\leq e^{-\lambda t} \mathbb{E} \left[e^{\lambda(\hat{R}(f) - R(f))} \right] \\
\mathbb{P} \left\{ \hat{R}(f) - R(f) \geq t \right\} &\leq e^{-\frac{2nt^2}{c^2}} \\
\mathbb{P} \left\{ |\hat{R}(f) - R(f)| \geq t \right\} &= \mathbb{P} \left\{ \hat{R}(f) - R(f) \geq t \text{ or } R(f) - \hat{R}(f) \geq t \right\} \\
&\leq 2e^{-\frac{2nt^2}{c^2}}
\end{aligned}$$

by union bound

$$\mathbb{P} \left\{ \hat{R}(f_1) - R(f_1) \geq t \text{ or } \hat{R}(f_2) - R(f_2) \geq t \text{ or } \dots \right\} \leq \text{small bound}$$

Assume \mathcal{F} is finite.

$k = |\mathcal{F}|$ is the number of classifiers in \mathcal{F} .

24.4 Uniform Bound

$$\begin{aligned}
&\mathbb{P} \left\{ \hat{R}(f_1) - R(f_1) \geq t \text{ or } \hat{R}(f_2) - R(f_2) \geq t \text{ or } \dots \underbrace{\hat{R}(f_k) - R(f_k) \geq t}_{\text{bad things}} \right\} \\
&\leq \sum_{i=1}^k \left\{ |\hat{R}(f_i) - R(f_i)| \geq t \right\}, \text{ union bound} \\
&\leq \underbrace{|\mathcal{F}| 2e^{-\frac{2nt^2}{c^2}}}_{\delta} \\
&\delta = 2|\mathcal{F}| e^{-\frac{2nt^2}{c^2}} \\
\log \frac{2|\mathcal{F}|}{\delta} &= \frac{2nt^2}{c^2} \\
t &= \sqrt{\frac{c^2 \log \left(\frac{2|\mathcal{F}|}{\delta} \right)}{2n}}
\end{aligned}$$

with probability $\geq 1 - \delta$,

$$\max_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| \leq \sqrt{\frac{c^2 \log \left(\frac{2|\mathcal{F}|}{\delta} \right)}{2n}}$$

$$\begin{aligned}
f^* &= \arg \min_{f \in \mathcal{F}} \mathbb{E} [l(y, f(x))] \\
&= \arg \min_{f \in \mathcal{F}} R(f) \\
\hat{f} &= \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i))
\end{aligned}$$

$$\begin{aligned}
&= \arg \min_{f \in \mathcal{F}} \hat{R}(f) \\
R(\hat{f}) &\leq \hat{R}(\hat{f}) + \sqrt{\frac{c^2 \log\left(\frac{2|\mathcal{F}|}{\delta}\right)}{2n}} \text{ wp } \geq 1 - \delta \\
&\leq \hat{R}(f^*) + \sqrt{\frac{c^2 \log\left(\frac{2|\mathcal{F}|}{\delta}\right)}{2n}}, \text{ since } \hat{f} \text{ min } \hat{R} \\
&\leq R(f^*) + \sqrt{\frac{c^2 \log\left(\frac{2|\mathcal{F}|}{\delta}\right)}{n}}
\end{aligned}$$

Generalization Bound

decreases like $\frac{1}{\sqrt{n}}$
increases with $\log\left(\frac{1}{\delta}\right)$ and $\log |\mathcal{F}|$

$$\begin{aligned}
n &\geq \log |\mathcal{F}| \\
\Delta &= \min_{f \neq f^*} R(f) - R(f^*)
\end{aligned}$$

with prob $\geq 1 - \delta$

$$\begin{aligned}
R(\hat{f}) &\leq R(f^*) + \sqrt{\frac{2 \log\left(\frac{2|\mathcal{F}|}{\delta}\right)}{n}} \\
\mathbb{E}[R(\hat{f})] &\leq (1 - \delta) \left[R(f^*) + \sqrt{\frac{2 \log\left(\frac{2|\mathcal{F}|}{\delta}\right)}{n}} \right] + \delta c, \text{ take } \delta = \frac{1}{\sqrt{n}} \\
&\leq R(f^*) + \tilde{O}\left(\sqrt{\frac{\log |\mathcal{F}| + \log n}{n}}\right)
\end{aligned}$$

24.5 Infinite \mathcal{F}

$$\begin{aligned}
\mathcal{F} &= \left\{ \text{all linear classifiers on } [0, 1]^k \right\} \\
|\mathcal{F}_\varepsilon| &= o\left(\left(\frac{1}{\varepsilon}\right)^d\right) \\
\log |\mathcal{F}_\varepsilon| &\sim d \log \frac{1}{\varepsilon} \\
t &= \sqrt{\frac{c^2 \log\left(\frac{2|\mathcal{F}|}{\delta}\right)}{n}}
\end{aligned}$$

24.6 Hyperparameter Tuning

$$\hat{w}_\lambda = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (1 - y_i w^T x_i)_+ \lambda \|w\|^2, \lambda \in \Lambda$$

want "tune" λ

$$\begin{aligned} & \{(x_i, y_i)\}_{i=1}^n \\ & n = n_T + n_H \\ & w_\lambda = \arg \min \sum_{i=1}^{n_T} (1 - y_i w^T x_i)_+ + \lambda \|w\|^2 \end{aligned}$$

use hold out set $\{(x_i, y_i)\}_{i=n_T+1}^n$ validate, tune λ

$$\hat{R}(w_\lambda) = \frac{1}{n_H} = \sum_{i=n_T+1}^n \mathbb{1}_{\{y_i w_\lambda^T x_i < 0\}}$$

number of mistakes w_λ makes on hold out

$$\begin{aligned} \lambda^* &= \arg \min_{\lambda \in \Lambda} \mathbb{E} \left[\mathbb{1}_{\{y w_\lambda^T x < 0\}} \right] \\ \hat{\lambda} &= \arg \min_{\lambda \in \Lambda} \hat{R}(w_\lambda) \end{aligned}$$

Assume Λ is finite

$$\begin{aligned} \Lambda &= \{\lambda_1, \lambda_2, \dots, \lambda_k\} \\ |\hat{R}(w_\lambda) - R(w_\lambda)| &\leq \sqrt{\frac{2 \log \frac{2|\Lambda|}{\delta}}{n_H}}, \text{ wp } \geq 1 - \delta \end{aligned}$$

25 Lecture 25

Feature space \mathcal{X}

Label sapce \mathcal{Y}

Predictor $f : \mathcal{X} \rightarrow \mathcal{Y}, f \in \mathcal{F}$

Loss function $l : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$

$$\begin{aligned} & \{(x_i, y_i)\}_{i=1,2,\dots,n} \stackrel{iid}{\sim} P \\ & R(f) = \mathbb{E}_{(x,y) \sim P} [l(y, f(x))] \\ & \hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i)) \\ & f^* = \arg \min_{f \in \mathcal{F}} R(f) \\ & \hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}(f) \\ & \mathbb{P} \left\{ |R(\hat{f}) - R(f^*)| \geq t \right\} \end{aligned}$$

$$\leq 2|\mathcal{F}| \exp\left(-\frac{2nt^2}{c^2}\right)$$

$$t < \sqrt{\frac{\log |\mathcal{F}| + \log n}{n}}$$

25.1 Markov Ineq

$X > 0, \phi$ increasing

$$\mathbb{P}\{\phi(X) > \phi(t)\} = \mathbb{P}\{X \geq t\} \leq \frac{\mathbb{E}[X]}{t}$$

$$\mathbb{E}[X] \geq \mathbb{E}[X \mathbb{1}_{[X \geq t]}] \geq \mathbb{E}[t \mathbb{1}_{[X \geq t]}] = t \mathbb{P}\{X \geq t\}$$

25.2 Chebyshev Inequality

$$\mathbb{P}\{|X - \mathbb{E}[X]| \geq t\}$$

$$= \mathbb{P}\{|X - \mathbb{E}[X]|^2 \geq t^2\}$$

$$\leq \frac{\mathbb{V}[X]}{t^2}$$

25.3 Chernoeff Bound

$$\mathbb{P}\{X - \mathbb{E}[X] \geq t\}$$

$$= \mathbb{P}\{\exp(s(X - \mathbb{E}[X])) \geq e^{st}\}$$

$$\leq \frac{\mathbb{E}[\exp(s(X - \mathbb{E}[X]))]}{e^{st}}$$

$$\leq \min_{s>0} \frac{\mathbb{E}[e^{sX}]}{e^{st} e^{s\mathbb{E}[X]}}$$

25.4 Sub Gaussian Random Variables

Assume $\mathbb{E}[X] = 0$,

$$\mathbb{E}[e^{sX}] \leq e^{\frac{cs^2}{2}}, c > 0$$

$$\mathbb{P}\{|X| \geq t\} \leq \frac{e^{\frac{cs^2}{2}}}{e^{st}} = e^{\frac{cs^2}{2} - st}$$

$$\leq e^{-\frac{t^2}{2c}}$$

$$s = \frac{t}{c}$$

X is subGaussian with param c ,

$$S_n = \sum_{i=1}^n x_i$$

$$\begin{aligned}
\mathbb{E} [e^{sS_n}] &= \mathbb{E} \left[e^{s \sum_{i=1}^n x_i} \right] \\
&= \prod_{i=1}^n \mathbb{E} [e^{sX_i}] \leq \exp \left(\frac{ncs^2}{2} \right) \\
\mathbb{P} \{ |S_n| \geq t \} &\leq \exp \left(-\frac{t^2}{2nc} \right) \\
\mathbb{P} \left\{ \left| \frac{1}{n} \sum x_i \right| \geq t \right\} &\leq \exp \left(-\frac{nt^2}{2c} \right)
\end{aligned}$$

X is bounded, $x \in [a, b]$

$$\mathbb{E} [e^{sX}] \leq e^{\frac{(b-a)^2 s^2}{2}}$$

25.5 Shattering Coefficient of \mathcal{F}

$$\begin{aligned}
D_n &= \{(x_i, y_i)\}_{i=1, \dots, n} \\
\mathcal{S}(\mathcal{F}, n) &= \max_{x_1, \dots, x_n} |\{(f(x_1), f(x_2), \dots, f(x_n)), f \in \mathcal{F}\}|
\end{aligned}$$

Examples:

\mathcal{F} : 1-dim linear classifier

$$\begin{aligned}
\mathcal{F}_{D_n} &= \{(+, +, +), (-, -, -), (+, -, -), (-, -, +), (-, +, +), (+, +, -)\} \\
VC(\mathcal{F}) &\geq 2
\end{aligned}$$

\mathcal{F}' : 2-dim linear classifier

$$\begin{aligned}
\mathcal{F}_{D_n} &= \{-1, +1\}^3 \\
VC(\mathcal{F}) &\geq 3
\end{aligned}$$

The total number of classifiers is,

$$2 \binom{n}{d}$$

\mathcal{F} shatter D_n if \mathcal{F}_{D_n} contains all tuples of $\{+1, -1\}^n \Rightarrow |\mathcal{F}_{D_n}| = 2^n$.

$VC(\mathcal{F}) = k$ if k is the largest integer there exist D_k such that \mathcal{F} shatter D_k

VC of d dim linear classifiers is $d + 1$

25.6 FTSLT

$$\mathbb{P} \left\{ \sup_{f \in \mathcal{F}} |\hat{R}_n(f) - R(f)| \geq \varepsilon \right\} \leq 8\mathcal{S}(\mathcal{F}, n) e^{-\frac{n\varepsilon^2}{32}}$$

$$= O \left(e^{\frac{n\varepsilon^2}{c} + VC \log n} \right)$$

$$\mathcal{S}(\mathcal{F}, n) \leq (n+1)^{VC}$$

$$D_n = \{(x_i, y_i)\}_{i=1, \dots, n}$$

$$D'_n = \{(x'_i, y'_i)\}_{i=1, \dots, n}$$

$$\hat{R}'_n(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[f(x'_i) \neq y_i]}$$

Step 1:

$$\mathbb{P} \left\{ \sup_{f \in \mathcal{F}} |\hat{R}_n(f) - R(f)| \geq \varepsilon \right\} \leq 2\mathbb{P} \left\{ \sup_{f \in \mathcal{F}} |\hat{R}_n(f) - \hat{R}'_n(f)| > \frac{\varepsilon}{2} \right\}$$

$$\{\sigma_i\}_{i=1, \dots, n} = \begin{cases} +1 & \text{with probability } \frac{1}{2} \\ -1 & \text{with probability } \frac{1}{2} \end{cases}$$

$$\mathbb{P} \left\{ \sup_{f \in \mathcal{F}} |\hat{R}_n(f) - \hat{R}'_n(f)| > \frac{\varepsilon}{2} \right\}$$

$$= \mathbb{P} \left\{ \sum_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{f(x_i) \neq y_i\}} - \mathbb{1}_{\{f(x'_i) \neq y'_i\}} \right) \right| \geq \frac{\varepsilon}{2} \right\}$$

$$\leq \mathbb{P} \left\{ \sum_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{f(x_i) \neq y_i\}} \right) \right| \geq \frac{\varepsilon}{4} \right\} + \mathbb{P} \left\{ \sum_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{f(x'_i) \neq y'_i\}} \right) \right| \geq \frac{\varepsilon}{4} \right\}$$

$$\mathbb{P} \left\{ \sum_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{f(x_i) \neq y_i\}} \right) \right| \geq \varepsilon \right\}$$

$$= \mathbb{E} \left[\mathbb{1}_{\left\{ \sup_{f \in \mathcal{F}_{D_n}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i \left(\mathbb{1}_{\{f(x_i) \neq y_i\}} \right) \right| \geq \varepsilon \right\}} \right]$$

$$\leq \mathcal{S}(\mathcal{F}, n) \mathbb{E} \left[\sup_{f \in \mathcal{F}} \mathbb{P} \left\{ \frac{1}{n} \left| \sum_{i=1}^n \underbrace{\sigma_i \left(\mathbb{1}_{\{f(x_i) \neq y_i\}} \right)}_{iid, bounded by 1, expectation 0} \right| \geq \varepsilon \mid D_n \right\} \right]$$

$$\leq \mathcal{S}(\mathcal{F}, n) e^{-\frac{n\varepsilon^2}{32}}$$

26 Lecture 26

$$R(f) = \mathbb{P}\{f(X) \neq y\}$$

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{f(x_i) \neq y_i\}}$$

$$f \in \mathcal{F}, \hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}(f)$$

$$f^* = \arg \min_{f \in \mathcal{F}} R(f)$$

26.1 Uniform Derivation Bound

$$\sup_{f \in \mathcal{F}} |R(f) - \hat{R}_n(f)| \leq 6 \sqrt{\frac{\overbrace{d_{\mathcal{F}}} \log\left(\frac{n}{\delta}\right)}{VCdim(\mathcal{F}) n}} \text{ wp } \geq 1 - \delta$$

$$R(\hat{f}) \leq \hat{R}(\hat{f}) + 6 \sqrt{\frac{d_{\mathcal{F}} \log\left(\frac{n}{\delta}\right)}{n}}$$

$$\leq \hat{R}(f^*) + 6 \sqrt{\frac{d_{\mathcal{F}} \log\left(\frac{n}{\delta}\right)}{n}}$$

$$\leq R(\hat{f}) + 12 \sqrt{\frac{d_{\mathcal{F}} \log\left(\frac{n}{\delta}\right)}{n}}$$

generalization bound

26.2 VC Dimensions

VC dim of axis-aligned rectangles in $\mathbb{R}^d = 2d$

26.3 Sample Complexity

$$n \gg d_{\mathcal{F}}$$

26.4 Neural nets

VC dim of L layer ReLU = $O(L \cdot W)$

W is total number of weights $\neq d$

L is number of layers

26.5 Cross Validation

hold out $\frac{n}{10}$

$$\hat{f} = \arg \min_{\frac{9}{10}^n}$$

$$R(\hat{f}) \leq \hat{R}(\hat{f}) + \sqrt{\frac{2 \log\left(\frac{1}{\delta}\right)}{n}}$$

Chernoff bound

26.6 Convex Loss Functions

$$\begin{aligned} l(z) &= (1 - z)_+ \\ \mathbb{P}\{\text{err}\} &\leq R(f) = \mathbb{E}[\hat{l}(yf(x))] \\ \hat{R}(f) &= \frac{1}{n} \sum_{i=1}^n l(y_i f(x_i)) \\ \hat{f} &= \arg \min_{f \in \mathcal{F}} \hat{R}(f) \\ f^* &= \arg \min_{f \in \mathcal{F}} R(f) \end{aligned}$$

26.7 Linear Classifiers

$$\begin{aligned} f(x) &= w^T x \\ \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (1 - y_i w^T x_i)_+ \\ \sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| &\leq \mathbb{E} \left[\sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| \right] + \sqrt{\frac{2 \log \frac{1}{\delta}}{n}} \text{ wp } 1 - \delta \end{aligned}$$

McDiramid's Inequality

26.8 Rademacher Complexity

$$\begin{aligned} \mathbb{E} \left[\sup_f |R(f) - \hat{R}(f)| \right] &\leq \mathbb{E} \left[\sup_f |\hat{R}'(f) - \hat{R}(f)| \right] \\ &= \mathbb{E} \left[\sup_f \frac{1}{n} \left| \sum_{i=1}^n (l(y'_i f(x'_i)) - l(y_i f(x_i))) \right| \right] \\ &= \mathbb{E} \left[\sup_f \frac{1}{n} \sigma_i \left((l(y'_i f(x'_i)) - l(y_i f(x_i))) \right) \right] \\ &\leq 2 \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i l(y_i f(x_i)) \right| \right] \\ \sigma_i \text{ iid } \pm 1 \text{ wp } \frac{1}{2} \end{aligned}$$

Rademacher Complexity $\mathcal{R}_n(\mathcal{F})$

$$\|x_i\| \leq 1 \text{ for all } i$$

$$\begin{aligned} \mathcal{W}_M &= \{w \in \mathbb{R}^d : \|w\| \leq M\} \\ \min_{w \in \mathcal{W}_M} \sum_{i=1}^n (1 - y_i w^T x_i)_+ \\ \mathcal{R}_n(\mathcal{W}_M) &= 2\mathbb{E} \left[\sup_{w \in \mathcal{W}_n} \frac{1}{n} \left| \sum_{i=1}^n l(y_i w^T x_i) \right| \right] \\ &\leq 4\mathbb{E} \left[\sup_{w \in \mathcal{W}_n} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i y_i x_i \right| \right] \\ &\leq 4\mathbb{E} \left[\sup_{w \in \mathcal{W}_n} \frac{1}{n} \|M\| \left\| \sum_{i=1}^n \sigma_i x_i \right\| \right] \\ &\leq \frac{4M}{\sqrt{n}} \end{aligned}$$

Theorem 13. ('17,'18) \mathcal{F} = neural nets with L layers, M , using hinge or logistic

$$R(\hat{f}) \leq \hat{R}(\hat{f}) + O\left(\sqrt{\frac{LM^2}{n}}\right)$$

M is product of Frobenius norms of all weight matrix

27 Problem Set 1

27.1 Q1

No

27.2 Q2

No

27.3 Q3

Strategy 1: completely randomized,

$$\begin{aligned}\mathbb{P}\{X = \hat{X}, Y = \hat{Y}\} &= \mathbb{P}\{X = \hat{X}\} \mathbb{P}\{Y = \hat{Y}\} \\ &= \frac{1}{6} \frac{1}{6} \\ &= \frac{1}{36}\end{aligned}$$

Strategy 2: guess the other players number,

$$\begin{aligned}\mathbb{P}\{X = Y\} &= \frac{1}{6} \\ &< \frac{1}{36}\end{aligned}$$

This is minimal because,

$$\begin{aligned}\mathbb{P}\{X = \hat{X}, Y = \hat{Y}\} &\leq \mathbb{P}\{X = \hat{X}\} \\ &\leq \frac{1}{6}\end{aligned}$$

28 Problem Set 2

28.1 Q1

$$\begin{aligned}\mathbb{P}\{g(x) \neq Y|X = x\} &= \mathbb{P}\{Y = 1|X = x\} \mathbb{P}\{g(X) = 0|X = x\} + \mathbb{P}\{Y = 0|X = x\} \mathbb{P}\{g(X) = 1|X = x\} \\ &= \eta(x)(1 - \mathbb{P}\{g(X) = 1|X = x\}) + (1 - \eta(x)) \mathbb{P}\{g(X) = 1|X = x\} \\ \mathbb{P}\{g(x) \neq Y|X = x\} - \mathbb{P}\{f^*(x) \neq Y|X = x\} &= (2\eta(x) - 1)(\mathbb{P}\{f^*(x) = 1|X = x\} - \mathbb{P}\{g(x) = 1|X = x\}) \\ &\geq 0\end{aligned}$$

If $\eta(x) \geq \frac{1}{2}$, then,

$$\begin{aligned}2\eta(x) - 1 &\geq 0 \\ \mathbb{P}\{f^*(x) = 1|X = x\} - \mathbb{P}\{g(x) = 1|X = x\} &= 1 - \mathbb{P}\{g(x) = 1|X = x\} \geq 0\end{aligned}$$

If $\eta(x) < \frac{1}{2}$, then,

$$2\eta(x) - 1 < 0$$

$$\mathbb{P}\{f^*(x) = 1|X = x\} - \mathbb{P}\{g(x) = 1|X = x\} = 1 - \mathbb{P}\{g(x) = 1|X = x\} \leq 0$$

Therefore,

$$\begin{aligned}\mathbb{P}\{g(X) \neq Y\} &= \mathbb{E}[\mathbb{P}\{g(X) \neq Y\}|X] \\ &\geq \mathbb{E}[\mathbb{P}\{f^*(X) \neq Y\}|X] \\ &= \mathbb{P}\{f^*(X) \neq Y\}\end{aligned}$$

28.2 Q2

$$\begin{aligned}\mathbb{E}[Z] &= \mathbb{E}[Z|Z < t]\mathbb{P}\{Z < t\} + \mathbb{E}[Z|Z \geq t]\mathbb{P}\{Z \geq t\} \\ &\geq \mathbb{E}[Z|Z \geq t]\mathbb{P}\{Z \geq t\} \\ &\geq t\mathbb{P}\{Z \geq t\}\end{aligned}$$

28.3 Q3

Start with showing mean = 0,

$$\begin{aligned}\mathbb{E}[\mathbb{1}_{f(X_i) \neq Y_i} - p_f] &= \mathbb{P}\{f(X_i) \neq Y_i\} - p_f \\ &= 0\end{aligned}$$

Use Markov,

$$\begin{aligned}\mathbb{P}\{|\hat{p}_f - p_f| > \varepsilon\} &\leq \frac{\mathbb{E}[(\hat{p}_f - p_f)^2]}{\varepsilon^2} \\ &= \frac{1}{\varepsilon^2} \mathbb{E}\left[\left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{f(X_i) \neq Y_i} - p_f\right)^2\right] \\ &= \frac{1}{\varepsilon^2} \mathbb{V}\left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{f(X_i) \neq Y_i} - p_f\right] \\ &= \frac{1}{\varepsilon^2} \frac{1}{n^2} \sum_{i=1}^n \mathbb{V}[\mathbb{1}_{f(X_i) \neq Y_i} - p_f] \\ &= \frac{1}{(n\varepsilon)^2} \sum_{i=1}^n \mathbb{E}[\mathbb{1}_{f(X_i) \neq Y_i} - 2p_f \mathbb{1}_{f(X_i) \neq Y_i} + p_f^2] \\ &= \frac{1}{(n\varepsilon)^2} \sum_{i=1}^n (p_f - p_f^2) \\ &= \frac{p_f(1 - p_f)}{n\varepsilon^2}\end{aligned}$$

28.4 Q4

By convexity,

$$\begin{aligned} g(x) &\geq g(t) + g'(t)(x - t) \\ \mathbb{E}[g(x)] &\geq g(t) + g'(t)(\mathbb{E}[x] - t) \end{aligned}$$

With $t = \mathbb{E}[x]$ and $g(x) = x^2$

$$\mathbb{E}[x^2] \geq \mathbb{E}[x]^2$$

28.5 Q5

$$\begin{aligned} \mathbb{E}[X] &= \sum_{k=0}^{\infty} k e^{-\lambda} \frac{\lambda^k}{k!} \\ &= \lambda \sum_{k=1}^{\infty} e^{-\lambda} \frac{\lambda^{k-1}}{(k-1)!} \\ &= \lambda \end{aligned}$$

28.6 Q6

$$\begin{aligned} \mathbb{P}\{X_1 + X_2 = n\} &= \sum_{i=0}^n \mathbb{P}\{X_1 = i, X_2 = n - i\} \\ &= \sum_{i=0}^n e^{-\lambda_1} \frac{\lambda_1^i}{i!} e^{-\lambda_2} \frac{\lambda_2^{n-i}}{(n-i)!} \\ &= e^{-\lambda_1 - \lambda_2} \frac{1}{n!} \sum_{i=0}^n \binom{n}{i} \lambda_1^i \lambda_2^{n-i} \\ &= e^{-\lambda_1 - \lambda_2} \frac{(\lambda_1 + \lambda_2)^n}{n!} \\ \mathbb{P}\{X_1 = k | X_1 + X_2 = n\} &= \frac{\mathbb{P}\{X_1 = k, X_1 + X_2 = n\}}{\mathbb{P}\{X_1 + X_2 = n\}} \\ &= \frac{e^{-\lambda_1} \frac{\lambda_1^k}{k!} e^{-\lambda_2} \frac{\lambda_2^{n-k}}{(n-k)!}}{e^{-\lambda_1 - \lambda_2} \frac{(\lambda_1 + \lambda_2)^n}{n!}} \\ &= \binom{n}{k} \delta^k (1 - \delta)^{n-k}, \delta = \frac{\lambda_1}{\lambda_1 + \lambda_2} \end{aligned}$$

28.7 Q7

Define $\eta(x, k) = \mathbb{P}\{Y = k | X = x\}$ for $k = 1, \dots, m$, then

$$f^*(x) = \arg \max_k \eta(x, k)$$

28.8 Q8

The minimum error is $p_f = 1 - \mathbb{E}_x \left[\max_k \eta(x, k) \right]$

28.9 Q9

An estimator is,

$$\hat{p}_f = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{f(X_i) \neq Y_i}$$

28.10 Q10

Use the result from Q3,

$$\begin{aligned} \mathbb{P}\{p_f - \hat{p}_f > \varepsilon\} &\leq \frac{p_f(1-p_f)}{n\varepsilon^2} \leq \frac{0.25}{n\varepsilon^2} \\ \Rightarrow \mathbb{P}\{p_f > 0.05 + 0.05\} &\leq \frac{0.25}{1000 \cdot 0.05^2} \\ \Rightarrow p_f &< 0.1 \end{aligned}$$

29 Problem Set 3

29.1 Q1

$$\begin{aligned} \mathbb{P}\{f(x) \neq Y|X=x\} - \mathbb{P}\{f^*(x) \neq Y|X=x\} &= \eta(x)(1 - \mathbb{1}_{f(x)=1}) + (1 - \eta(x))\mathbb{1}_{f(x)=1} - \eta(x)(1 - \mathbb{1}_{f^*(x)=1}) - (1 - \eta(x))\mathbb{1}_{f^*(x)=1} \\ &= \eta(x)(2\mathbb{1}_{f^*(x)=1} - 2\mathbb{1}_{f(x)=1}) + \mathbb{1}_{f(x)=1} - \mathbb{1}_{f^*(x)=1} \\ &= (2\eta(x) - 1)(\mathbb{1}_{f^*(x)=1} - \mathbb{1}_{f(x)=1}) \end{aligned}$$

This is,

$$\begin{cases} 0 & \text{if } \eta(x) \geq \frac{1}{2} \text{ and } \tilde{\eta}(x) \geq \frac{1}{2} \\ 0 & \text{if } \eta(x) < \frac{1}{2} \text{ and } \tilde{\eta}(x) < \frac{1}{2} \\ 2\left(\eta(x) - \frac{1}{2}\right) \leq 2(\eta(x) - \tilde{\eta}(x)) & \text{if } \eta(x) \geq \frac{1}{2} \text{ and } \tilde{\eta}(x) < \frac{1}{2} \\ 2\left(\eta(x) - \frac{1}{2}\right) \leq 2(\eta(x) - \tilde{\eta}(x)) & \text{if } \eta(x) < \frac{1}{2} \text{ and } \tilde{\eta}(x) \geq \frac{1}{2} \end{cases}$$

Therefore,

$$\mathbb{P}\{f(x) \neq Y|X=x\} - \mathbb{P}\{f^*(x) \neq Y|X=x\} \leq 2|\eta(x) - \tilde{\eta}(x)|$$

29.2 Q2

1. The loss function is,

$$\begin{cases} c_{01} & \text{if } f(x) = 0, y = 1 \\ c_{10} & \text{if } f(x) = 1, y = 0 \end{cases}$$

Then the expect loss given f and f^* is,

$$\begin{aligned}
\mathbb{E}[l(f, X, Y) | X = x] - \mathbb{E}[l(f^*, X, Y) | X = x] &= c_{01}\eta(x)(1 - \mathbb{1}_{f(x)=1}) + c_{10}(1 - \eta(x))\mathbb{1}_{f(x)=1} - c_{01}\eta(x)(1 - \mathbb{1}_{f^*(x)=1}) \\
&= \eta(x)(-(c_{01} + c_{10})\mathbb{1}_{f(x)=1} + (c_{01} + c_{10})\mathbb{1}_{f^*(x)=1}) + c_{10}(\mathbb{1}_{f(x)=1} - \mathbb{1}_{f^*(x)=1}) \\
&= (\eta(x)(c_{01} + c_{10}) - c_{10})(\mathbb{1}_{f^*(x)=1} - \mathbb{1}_{f(x)=1})
\end{aligned}$$

To minimize the expected loss,

Set $f^*(x) = 1$ if and only if,

$$\begin{aligned}
\eta(x)(c_{01} + c_{10}) - c_{10} &\geq 0 \\
\eta(x) &\geq \frac{c_{10}}{c_{01} + c_{10}}
\end{aligned}$$

Therefore, the optimal classifier is,

$$\begin{cases} 1 & \text{if } \eta(x) \geq \frac{c_{10}}{c_{01} + c_{10}} \\ 0 & \text{otherwise} \end{cases}$$

2. No. $f^*(x)$ does minimize the average probability of error for any p .

1. The same as (a), π_i is already contained in $\eta(x) = \frac{p(x|Y=1)\pi_1}{p(x)}$.

29.3 Q3

1. The optimal classifier is $\hat{y} = x_1, R^* = 0$

1. All possible data points are $y = 1$ and $x = (1, 1)$ or $(1, -1)$ and $y = -1$ and $x = (-1, 1)$ or $(-1, -1)$, the probability of error is, assuming the original data are $(1, x_1)$ and $(-1, x_2)$, and the new data $(1, \delta)$ has label 1,

$$\begin{aligned}
\mathbb{P}\{f_n(x) \neq 1\} &= \frac{1}{2}\mathbb{P}\{\delta \neq x_1\}\mathbb{P}\{x_1 \neq x_2\} \\
&= \frac{1}{2} \frac{1}{2} \frac{1}{2} \\
&= \frac{1}{8}
\end{aligned}$$

2. Assume the new data has label 1, and the training data are $u = (1, u_2, \dots, u_d)$ and $v = (-1, v_2, \dots, v_d)$,

Let $U = \sum_j |u_j - x_j| \sim \text{Bin}\left(d-1, \frac{1}{2}\right)$ and $V = \sum_j |v_j - x_j| - 1 \sim \text{Bin}\left(d-1, \frac{1}{2}\right)$,

and note that $d-1+U-V \sim \text{Bin}\left(2d-2, \frac{1}{2}\right)$,

$$\begin{aligned}
\mathbb{P}\{f_n(x) \neq 1\} &= \mathbb{P}\{\|u - x\| < \|v - x\|\} + \frac{1}{2}\mathbb{P}\{\|u - x\| = \|v - x\|\} \\
&= \mathbb{P}\{U+1 < V\} + \frac{1}{2}\mathbb{P}\{1+U = V\} \\
&= \mathbb{P}\{d-1+V-U > d\} + \frac{1}{2}\mathbb{P}\{d-1+V-U = d\}
\end{aligned}$$

$$= \sum_{i=d+1}^{2d-2} \binom{2d-2}{i} \frac{1}{2^{2d-2}} + \frac{1}{2} \binom{2d-2}{d} \frac{1}{2^{2d-2}}$$

3. As $d \rightarrow \infty$, $\text{Bin} \left(2d-2, \frac{1}{2} \right)$ is approximately $N \left(d-1, \frac{d-1}{2} \right)$ or $N \left(d, \frac{d}{2} \right)$

$$\begin{aligned} \mathbb{P} \{ f_n(x) \neq 1 \} &= \mathbb{P} \{ d-1 + V - U > d \} + \frac{1}{2} \mathbb{P} \{ d-1 + V - U = d \} \\ &= \frac{1}{2} + 0 \\ &= \frac{1}{2} \end{aligned}$$

29.4 Q4

1. The MLE is,

$$\begin{aligned} \hat{y}(x) &= \arg \max_l p(y = l|x) \\ &= \arg \max_l p(x|y = l) p(y = l) \\ &= \arg \max_l \log p(x|y = l) + \log p(y = l) \\ &= \arg \max_l \frac{-1}{2} \log |\Sigma_l| - \frac{1}{2} (x - \mu_l)^T \Sigma_{l-1}^{-1} (x - \mu_l) + \log \pi_l \end{aligned}$$

Given data,

$$\hat{y}(x) = \arg \max_l \frac{-1}{2} \log |\hat{\Sigma}_l| - \frac{1}{2} (x - \hat{\mu}_l)^T \hat{\Sigma}_{l-1}^{-1} (x - \hat{\mu}_l) + \log \hat{\pi}_l$$

2. Use a common covariance matrix $\hat{\Sigma}$ instead of individual Σ_l ,

$$\begin{aligned} \hat{y}(x) &= \arg \max_l \frac{-1}{2} \log |\hat{\Sigma}| - \frac{1}{2} (x - \hat{\mu}_l)^T \left(\hat{\Sigma} \right)^{-1} (x - \hat{\mu}_l) + \log \hat{\pi}_l \\ &= \arg \max_l -\frac{1}{2} (x - \hat{\mu}_l)^T \hat{\Sigma}^{-1} (x - \hat{\mu}_l) + \log \hat{\pi}_l \\ &= \arg \max_l -x^T \left(\hat{\Sigma} \right)^{-1} x + 2\hat{\mu}_l^T \hat{\Sigma}^{-1} x - \hat{\mu}_l^T \hat{\Sigma}^{-1} \hat{\mu}_l + \log \hat{\pi}_l \\ &= \arg \max_l -x^T \left(\hat{\Sigma} \right)^{-1} x + 2\hat{\mu}_l^T \hat{\Sigma}^{-1} x - \hat{\mu}_l^T \hat{\Sigma}^{-1} \hat{\mu}_l + \log \hat{\pi}_l \\ &= \arg \max_l \left(2\hat{\mu}_l^T \hat{\Sigma}^{-1} \right) x - \left(\hat{\mu}_l^T \hat{\Sigma}^{-1} \hat{\mu}_l + \log \hat{\pi}_l \right) \end{aligned}$$

29.5 Q5

Use Markov's Inequality,

$$\mathbb{P} \{ X \geq t \} \leq \mathbb{P} \{ X^2 \geq t^2 \} \leq \frac{\mathbb{E}[X^2]}{t^2}$$

29.6 Q6

The event $\{X + Y > t\} \subseteq \left\{X > \frac{t}{2}\right\} \cup \left\{Y > \frac{t}{2}\right\}$

$$\mathbb{P}\{X + Y > t\} \leq \mathbb{P}\left\{X > \frac{t}{2}\right\} + \mathbb{P}\left\{Y > \frac{t}{2}\right\}$$

29.7 Q7

1. The optimal Bayes classifier is always,

$$f^*(x) = \begin{cases} 1 & \text{if } \eta(x) > \frac{1}{2} \\ -1 & \text{otherwise} \end{cases}$$

To simplify $\eta(x)$,

$$\begin{aligned} \eta(x) > \frac{1}{2} &\Leftrightarrow \mathbb{P}\{Y = 1|X = x\} > \frac{1}{2} \\ &\Leftrightarrow \mathbb{P}\{X = x|Y = 1\} \mathbb{P}\{Y = 1\} > \frac{1}{2} \mathbb{P}\{X = x\} \\ &\Leftrightarrow 2\mathbb{P}\{X = x|Y = 1\} \mathbb{P}\{Y = 1\} > \mathbb{P}\{X = x|Y = 1\} \mathbb{P}\{Y = 1\} + \mathbb{P}\{X = x|Y = -1\} \mathbb{P}\{Y = -1\} \\ &\Leftrightarrow \mathbb{P}\{X = x|Y = 1\} \mathbb{P}\{Y = 1\} > \mathbb{P}\{X = x|Y = -1\} \mathbb{P}\{Y = -1\} \\ &\Leftrightarrow \log \mathbb{P}\{X = x|Y = 1\} + \log \mathbb{P}\{Y = 1\} > \log \mathbb{P}\{X = x|Y = -1\} + \log \mathbb{P}\{Y = -1\} \\ &\Leftrightarrow -\frac{1}{2} \frac{1}{\sigma} (x - \theta)^T (x - \theta) + \log \pi_1 > -\frac{1}{2} \frac{1}{\sigma} (x + \theta)^T (x + \theta) + \log \pi_{-1} \\ &\Leftrightarrow -\frac{1}{2} \frac{1}{\sigma} (x^T x - 2x^T \theta + \theta^T \theta - x^T x - 2x^T \theta - \theta^T \theta) > \log \pi_{-1} - \log \pi_1 \\ &\Leftrightarrow \frac{2}{\sigma} x^T \theta > \log \pi_{-1} - \log \pi_1 \\ &\Leftrightarrow x^T \theta > \frac{\sigma}{2} (\log \pi_{-1} - \log \pi_1) \end{aligned}$$

Therefore, the optimal Bayes classifier is,

$$f^*(x) = \begin{cases} 1 & \text{if } x^T \theta > \frac{\sigma}{2} (\log \pi_{-1} - \log \pi_1) \\ -1 & \text{otherwise} \end{cases}$$

when $\pi_1 = \pi_{-1}$,

$$f^*(x) = \begin{cases} 1 & \text{if } x^T \theta > 0 \\ -1 & \text{otherwise} \end{cases}$$

2. The MLE is,

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} \sum_{i=1}^n \log \mathbb{P}\{x_i, y_i | \theta\} \\ &= \arg \max_{\theta} \sum_{i=1}^n \frac{-1}{2} \frac{1}{\sigma} (x_i - \theta y_i)^T (x_i - \theta y_i) + \log \pi_{y_i} \end{aligned}$$

$$\begin{aligned}
&= \arg \min_{\theta} \sum_{i=1}^n (x_i - \theta y_i)^T (x_i - \theta y_i) \\
&= \arg \min_{\theta} \sum_{i=1}^n (x_i^T x_i - 2y_i^T x_i + y_i^T \theta^T \theta) \\
&= \arg \min_{\theta} \sum_{i=1}^n (-2y_i^T x_i + \theta^T \theta) \\
&= \theta : \sum_{i=1}^n -2y_i x_i + 2\theta = 0 \\
&= \frac{1}{n} \sum_{i=1}^n x_i y_i
\end{aligned}$$

3. The plug-in classifier is

$$\hat{f}(x) = \begin{cases} 1 & x^T \hat{\theta} > \frac{\hat{\sigma}}{2} (\log \hat{\pi}_{-1} - \log \hat{\pi}_1) \\ -1 & \text{otherwise} \end{cases}$$

where,

$$\begin{aligned}
\hat{\sigma} &= \frac{1}{nd} (x_i - \theta y_i)^T (x_i - \theta y_i) \\
\hat{\pi}_j &= \frac{n_{y_i=j}}{n}
\end{aligned}$$

4. The error probability for x with true label -1 is,

$$\mathbb{P} \left\{ \tilde{f}(x) = 1 \right\} = \mathbb{P} \left\{ x^T \hat{\theta} > \frac{\hat{\sigma}}{2} (\log \hat{\pi}_{-1} - \log \hat{\pi}_1) \right\}$$

where

$$\begin{aligned}
x &\sim N(-\theta, \sigma^2 I) \\
\theta &\sim N\left(\theta, \frac{\sigma^2}{n} I\right)
\end{aligned}$$

5. The error probability from the previous part, and define $e_1 = x + \theta$ and $e_2 = \hat{\theta} - \theta$

$$\begin{aligned}
\mathbb{P} \left\{ \tilde{f}(x) = 1 \right\} &\leq \mathbb{P} \left\{ x^T \hat{\theta} > 0 \right\} \\
&= \mathbb{P} \left\{ (e_1 - \theta)^T (e_2 + \theta) > 0 \right\} \\
&= \mathbb{P} \left\{ (e_1^T e_2 + e_1^T \theta - e_2^T \theta - \theta^T \theta) > 0 \right\} \\
&\leq \mathbb{P} \left\{ (e_1^T - e_2^T) \theta > \frac{1}{2} \theta^T \theta \right\} + \mathbb{P} \left\{ e_1^T e_2 > \frac{1}{2} \theta^T \theta \right\} \\
&\leq \frac{\mathbb{E} \left[((e_1^T - e_2^T) \theta)^2 \right]}{\left(\frac{1}{2} \theta^T \theta \right)^2} + \frac{\mathbb{E} \left[(e_1^T e_2)^2 \right]}{\left(\frac{1}{2} \theta^T \theta \right)^2}
\end{aligned}$$

$$\begin{aligned}
&= \frac{\sigma^2 \left(1 + \frac{1}{n}\right) \theta^T \theta + \sigma^4 \frac{d^2}{n}}{\left(\frac{1}{2} \theta^T \theta\right)^2} \\
&= O \left(\max \left\{ \frac{\sigma^2}{\theta^T \theta}, \frac{\sigma^4 d^2}{n (\theta^T \theta)^2} \right\} \right)
\end{aligned}$$

To reduce the error, both of the following needs to hold,

$$\begin{aligned}
\sigma^2 &<< \theta^T \theta \\
\sigma^4 &<< \frac{n}{d^2} (\theta^T \theta)^2
\end{aligned}$$

or,

$$\sigma^2 << \min \left\{ \theta^T \theta, \frac{\sqrt{n}}{d} (\theta^T \theta) \right\}$$

30 Problem Set 4

30.1 Q1

1. Use formula,

$$\begin{aligned}
x &\sim N(\mu, \Sigma) \Rightarrow Ax + b \sim N(A\mu + b, A\Sigma A^T) \\
w^T x &\sim N(0, w^T \Sigma w)
\end{aligned}$$

2. Similarly,

$$\begin{aligned}
z = Ax &\sim N(0, A\Sigma A^T) \\
&\Rightarrow A\Sigma A^T = I \\
&\Rightarrow A = \Sigma^{-\frac{1}{2}} = U D^{-\frac{1}{2}} U^T
\end{aligned}$$

30.2 Q2

1. $p(x_i|\theta) = \frac{1}{\theta} \mathbb{1}_{x \leq \theta}$

$$\begin{aligned}
\hat{\theta}_n &= \arg \max_{\theta} p(x_1, \dots, x_n | \theta) \\
&= \arg \max_{\theta} \prod_{i=1}^n p(x_i | \theta) \\
&= \arg \max_{\theta} \frac{1}{\theta^n} \mathbb{1}_{\max_i (x_i) \leq \theta} \\
&= \max_i \{x_i\}
\end{aligned}$$

2. The CDF,

$$\begin{aligned}
 F_{\hat{\theta}_n}(x) &= \mathbb{P}\left\{\hat{\theta}_n \leq x\right\} \\
 &= \mathbb{P}\left\{\max_i \{x_i\} \leq x\right\} \\
 &= \prod_{i=1}^n \mathbb{P}\{x_i \leq x\} \\
 &= \left(\frac{x}{\theta}\right)^n
 \end{aligned}$$

Then the PDF,

$$\begin{aligned}
 f_{\hat{\theta}_n}(x) &= F'_{\hat{\theta}_n}(x) \\
 &= n \frac{x^{n-1}}{\theta^n}
 \end{aligned}$$

3. The MSE is,

$$\begin{aligned}
 \mathbb{E}\left[\left(\hat{\theta}_n - \theta\right)^2\right] &= \int_0^\theta (x - \theta)^2 n \frac{x^{n-1}}{\theta^n} dx \\
 &= n \int_0^\theta \frac{x^{n+1}}{\theta^n} - 2 \frac{x^n}{\theta^{n-1}} + \frac{x^{n-1}}{\theta^{n-2}} dx \\
 &= \frac{n}{\theta^n} \int_0^\theta x^{n+1} - 2x^n \theta + x^{n-1} \theta^2 dx \\
 &= \frac{n}{\theta^n} \left[\frac{x^{n+2}}{n+2} - 2 \frac{x^{n+1}}{n+1} \theta + \frac{x^n}{n} \theta^2 \right]_{x=0}^\theta \\
 &= \frac{n}{\theta^n} \left(\frac{\theta^{n+2}}{n+2} - 2 \frac{\theta^{n+2}}{n+1} + \frac{\theta^{n+2}}{n} \right) \\
 &= n \theta^2 \frac{(n+1)n - 2(n+2)n + (n+1)(n+2)}{n(n+1)(n+2)} \\
 &= \frac{2\theta^2}{(n+1)(n+2)} \\
 &\rightarrow 0 \text{ as } n \rightarrow \infty
 \end{aligned}$$

30.3 Q3

1. The first derivative test for e_1 is,

$$\begin{aligned}
 \sum_{i=1}^n \text{sign}(\hat{\theta}_1 - x_i) &= 0 \\
 \Rightarrow \hat{\theta}_1 &= \text{median}(x)
 \end{aligned}$$

The first derivative test for e_2 is,

$$\sum_{i=1}^n 2(\hat{\theta}_2 - x_i) = 0$$

$$\Rightarrow \hat{\theta}_2 = \frac{1}{n} \sum_{i=1}^n x_i$$

2. The data is (0.9, 1.0, 1.1),

$$\hat{\theta}_1 = 1$$

$$\hat{\theta}_2 = 1$$

3. The data is (0.9, 1.1, 100),

$$\hat{\theta}_1 = \frac{102}{3}$$

$$\hat{\theta}_2 = 1.1$$

30.4 Q4

Use the invariance property and find the MLE for θ first,

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} \prod_{i=1}^N p(\tau_i | \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^N (\log \theta - \tau_i \theta) \\ &= \theta : \frac{N}{\theta} - \sum_{i=1}^N \tau_i = 0 \\ &= \hat{\theta} = \frac{N}{\sum_{i=1}^N \tau_i} \end{aligned}$$

Then,

$$\begin{aligned} \hat{\mathbb{P}}\{\tau < 10\} &= \int_0^{10} \hat{\theta} e^{-\tau \hat{\theta}} d\tau \\ &= 1 - e^{-10\hat{\theta}} \\ &= 1 - e^{-10 \frac{N}{\sum_{i=1}^N \tau_i}} \end{aligned}$$

31 Problem Set 5

31.1 Q1

Use Factorization Theorem,

$$p(x|a, b) = \prod_{i=1}^N \frac{1}{b-a} \mathbb{1}_{a \leq x_i \leq b}$$

$$\begin{aligned}
&= \frac{1}{(b-a)^N} \mathbb{1}_{a \leq x_1 \leq x_2 \leq \dots \leq x_n \leq b} \\
&= \frac{1}{(b-a)^N} \mathbb{1}_{a \leq x_{(1)} \leq x_{(n)} \leq b}
\end{aligned}$$

Therefore, $x_{(1)}$ and $x_{(n)}$ are sufficient

31.2 Q2

1. The covariance matrix is,

$$\begin{aligned}
\Sigma &= \mathbb{E}[XX^T] - \mathbb{E}[X] \mathbb{E}[X]^T \\
&= Q - \mu^2 ee^T
\end{aligned}$$

Use Sherman-Morrison formula,

$$\Sigma^{-1} = Q^{-1} + \frac{\mu^2 Q^{-1} ee^T Q^{-1}}{1 - \mu^2 e^T Q^{-1} e}$$

Then,

$$\begin{aligned}
\log p(x|\mu) &= (X - \mu e)^T \Sigma^{-1} (X - \mu e) \\
&= X^T \Sigma^{-1} X + \mu^2 e^T \Sigma^{-1} e - 2\mu X^T \Sigma^{-1} e \\
&= X^T Q^{-1} X + \frac{\mu^2 X^T Q^{-1} ee^T Q^{-1} X}{1 - \mu^2 e^T Q^{-1} e} - 2\mu \left(X^T Q^{-1} e + \frac{\mu^2 X^T Q^{-1} ee^T Q^{-1} e}{1 - \mu^2 e^T Q^{-1} e} \right)
\end{aligned}$$

Therefore, by Factorization Theorem,

$$\begin{aligned}
t(X) &= X^T Q^{-1} e \\
&= X^T \left(\begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\
&= \frac{1}{\rho + 1} (X_1 + X_2)
\end{aligned}$$

Note that, $X_1 + X_2$ is another sufficient statistics.

2. $\mathbb{E}[X_1] = \mu$

1. First note that,

$$\mathbb{E}[X_1 | T_X = X_1 + X_2] = \frac{T_X}{2}$$

Then,

$$\mathbb{E}\left[\frac{T_X}{2}\right] = \mu$$

and,

$$\begin{aligned}
\mathbb{V} \left[\frac{T_X}{2} \right] &= \frac{1}{4} (\mathbb{V} [X_1] + \mathbb{V} [X_2] + 2Cov [X_1, X_2]) \\
&= \frac{1}{2} (1 + \rho) \\
&\leq \frac{1}{2} (1 + 1) \\
&= 1 \\
&= \mathbb{V} [X_1]
\end{aligned}$$

31.3 Q3

1. MLE is invariant,

$$\begin{aligned}
\hat{\phi} &= \hat{p}_1 - \hat{p}_2 \\
&= \frac{x_1}{n_1} - \frac{x_2}{n_2}
\end{aligned}$$

2. The Log-Like is for p_i is $x_i \log p_i + (n_i - x_i) \log (1 - p_i)$

$$\begin{aligned}
\frac{\partial \log (x|p)}{\partial p_i} &= \frac{x_i}{p_i} - \frac{n_i - x_i}{1 - p_i} \\
\frac{\partial^2 \log (x|p)}{\partial p_i^2} &= -\frac{x_i}{p_i^2} - \frac{n_i - x_i}{(1 - p_i)^2} \\
\mathbb{E} [x_i] &= n_i p_i \\
\mathbb{E} \left[-\frac{\partial^2 \log (x|p)}{\partial p_i^2} \right] &= \frac{n_i}{p_i} - \frac{n_i}{1 - p_i} \\
&= \frac{n_i}{p_i (1 - p_i)} \\
\mathbb{E} \left[-\frac{\partial \log (x|p)}{\partial p_i \partial p_j} \right] &= 0
\end{aligned}$$

Therefore,

$$I(p) = \begin{bmatrix} \frac{n_1}{p_1 (1 - p_1)} & 0 \\ 0 & \frac{n_2}{p_2 (1 - p_2)} \end{bmatrix}$$

3. Use MLE Asypt theorem,

$$\begin{aligned}
\hat{p} &\sim N(p, I) \\
\hat{\phi} &\sim N \left(p_1 - p_2, \frac{p_1 (1 - p_1)}{n_1} + \frac{p_2 (1 - p_2)}{n_2} \right)
\end{aligned}$$

4. Use results from (c),

$$\begin{aligned}\mathbb{P}\left\{|\hat{\phi} - \phi| > 0.01\right\} &< 0.05 \Rightarrow \mathbb{P}\left\{\frac{|\hat{\phi} - \phi|}{\sigma} > \frac{0.01}{\sigma}\right\} < 0.05 \\ &\Rightarrow \frac{0.01}{\frac{1}{\sqrt{2n}}} > 2 \\ &\Rightarrow n > 20000\end{aligned}$$

5. NO

31.4 Q4

1. For the estimator $\hat{N}_1 = \frac{2}{n} \left(\sum_{i=1}^n x_i \right) - 1$,

$$\begin{aligned}\mathbb{E}[x_i] &= \sum_{i=1}^N \frac{i}{N} \\ &= \frac{N+1}{2} \\ \mathbb{E}[x_i^2] &= \sum_{i=1}^N \frac{i^2}{N} \\ &= \frac{(N+1)(2N+1)}{6} \\ \mathbb{E}[\hat{N}_1] &= \frac{2}{n} \sum_{i=1}^n \mathbb{E}[x_i] - 1 \\ &= \frac{2}{n} n \frac{N+1}{2} \\ &= N \\ \mathbb{V}[\hat{N}_1] &= \frac{4}{n^2} \sum_{i=1}^n \mathbb{V}[x_i] \\ &= \frac{4}{n^2} n \left(\frac{(N+1)(2N+1)}{6} - \frac{(N+1)^2}{4} \right) \\ &= \frac{N^2 - 1}{3n} \\ \text{MSE}[\hat{N}_1] &= \frac{N^2 - 1}{3n}\end{aligned}$$

\hat{N}_1 is unbiased.

2. MLE,

$$\begin{aligned}\hat{N}_2 &= \arg \max_N p(x|N) \\ &= \arg \max_N \sum_{i=1}^n -\log(N) + \log(\mathbb{1}_{x_i \leq N})\end{aligned}$$

$$\begin{aligned}
&= \arg \max_N -n \log(N) + \log(\mathbb{1}_{x_{(n)} \leq N}) \\
&= x_{(n)} \\
\mathbb{E}[\hat{N}_2] &= \sum_{i=1}^N i \mathbb{P}\{\hat{N}_2 = i\} \\
&= \sum_{i=1}^N i \left(\left(\frac{i}{N} \right)^N - \left(\frac{i-1}{N} \right)^N \right) \\
&= \frac{1}{N^N} \left(\sum_{i=1}^N i i^N - \sum_{i=1}^{N-1} (i+1) i^N \right) \\
&= \frac{1}{N^N} \left(N^{N+1} + \sum_{i=1}^{N-1} i i^N - \sum_{i=1}^{N-1} (i+1) i^N \right) \\
&= N - \frac{1}{N^N} \sum_{i=1}^{N-1} i^N \\
\mathbb{E}[\hat{N}_2^2] &= \sum_{i=1}^N i^2 \mathbb{P}\{\hat{N}_2 = i\} \\
&= \sum_{i=1}^N i^2 \left(\left(\frac{i}{N} \right)^N - \left(\frac{i-1}{N} \right)^N \right) \\
&= \frac{1}{N^N} \left(\sum_{i=1}^N i^2 i^N - \sum_{i=1}^{N-1} (i+1)^2 i^N \right) \\
&= N^2 - \frac{1}{N^N} \sum_{i=1}^N (2i+1) i^N \\
\mathbb{V}[\hat{N}_2^2] &= N^2 - \frac{1}{N^N} \sum_{i=1}^N (2i+1) i^N - \left(N - \frac{1}{N^N} \sum_{i=1}^{N-1} i^N \right)^2 \\
\text{MSE}[\hat{N}_2] &= N^2 - 2N \mathbb{E}[\hat{N}_2] + \mathbb{E}[\hat{N}_2^2] \\
&= N^2 - 2N \left(N - \frac{1}{N^N} \sum_{i=1}^{N-1} i^N \right) + N^2 - \frac{1}{N^N} \sum_{i=1}^N (2i+1) i^N \\
&= \frac{2}{N} \sum_{i=1}^{N-1} i^N - \frac{1}{N^N} \sum_{i=1}^N (2i+1) i^N
\end{aligned}$$

3. MLE is biased. Use Rao-Blackwell, and define,

$$\hat{N}_3 = \mathbb{E}[\hat{N}_1 | x_{(n)}]$$

31.5 Q5

No

31.6 Q6

No

32 Sample Midterm

32.1 Q1

1. Guess and check,

$$w^{\star} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

2. Use formula,

$$\begin{aligned}\hat{w} &= (X^T X)^{-1} X^T y \\ &= \begin{bmatrix} \frac{4}{3} \\ \frac{3}{4} \\ \frac{4}{3} \end{bmatrix} \\ \hat{y} &= \begin{bmatrix} \frac{4}{3} \\ \frac{3}{2} \\ -\frac{2}{3} \\ -\frac{2}{3} \end{bmatrix} \\ \hat{y} &= \text{sign}(x^T \hat{w})\end{aligned}$$

32.2 Q2

1. Least Squares

$$\min_w \|y - Xw\|^2$$

2. A lot

$$\hat{w} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \hat{w} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

3. 0

1. No

1. 3

32.3 Q3

1. Sample means for each $y_i = j \in \{000, 001, 010, 011, \dots, 111\}$.

$$\hat{\mu}_j = \frac{1}{\#_{y_i=j}} \sum_{i:y_i=j} x_i$$

$$\hat{\Sigma}_j = \frac{1}{\#_{y_i=j}} \sum_{i:y_i=j} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T$$

2. Naive Bayes

$$x : 000 = \arg \max_j p(x|y = j)$$

3. Histogram

$$p(y = j) = \frac{\#_{y_i=j}}{n}$$

4. Marginal

$$p(x) = \sum_j p(x|y = j) p(y = j)$$

5. Marginal

$$\begin{aligned} p(x|y_i = 1) &= \frac{p(x, y_i = 1)}{p(y_i = 1)} \\ &= \frac{\sum_{y_i=1} p(x|y) p(y)}{\sum_{y_i=1} p(y)} \end{aligned}$$

6. Naive Bayes: has disease if $p(x|y = 000) < p(x|y \neq 000)$, where,

$$p(x|y \neq 0) = \frac{\sum_{y \neq 000} p(x|y) p(y)}{\sum_{y \neq 000} p(y)}$$

32.4 Q4

1. Algorithm 5

1. Chernoff bound for mean is,

$$\pm \sqrt{\frac{\log\left(\frac{10}{\delta}\right)}{2n}}$$

For $\delta = 0.05$ and $n = 50$, this is ± 0.23

The standard deviation is,

$$\sigma = n^{-\frac{1}{2}}$$

For $n = 50$, this is ± 0.14

2. No

32.5 Q5

1. Let the classes by ± 1 $p(y = 1|x) = \frac{\exp(w_1^T x)}{\exp(w_1^T x) + \exp(w_{-1}^T x)}$
 $= \frac{1}{1 + \exp((w_{-1} - w_1)^T x)}$
 Define $w = w_{-1} - w_1$

2. Given generative model,

$$\begin{aligned} p(y = k|x) &= \frac{p(x|y = k)p(y = k)}{p(x)} \\ &= \frac{e^{-\frac{1}{2}(x-\theta_k)^T(x-\theta_k)} \frac{1}{c}}{\sum_{j=1}^c e^{-\frac{1}{2}(x-\theta_j)^T(x-\theta_j)} \frac{1}{c}} \\ &= \frac{e^{-\theta_k^T x + \theta_k^T \theta_k}}{\sum_{j=1}^c e^{-\theta_j^T x + \theta_j^T \theta_j}} \\ &= \left(\frac{e^{-\theta_k^T x}}{\sum_{j=1}^c e^{-\theta_j^T x}} \right) \end{aligned}$$

since $\|\theta_k\| = \|\theta_j\|$ for each j .

3. The gradient is,

$$\begin{aligned} \frac{\partial}{\partial \theta_k} - \log(p(y|x)) &= \frac{\partial}{\partial \theta_k} \left(-\theta_y^T x + \log \left(\sum_{j=1}^c e^{\theta_j^T x} \right) \right) \\ &= -\mathbb{1}_{y=k} x + \frac{e^{-\theta_k^T x}}{\sum_{j=1}^c e^{-\theta_j^T x}} x \\ &= (-\mathbb{1}_{y=k} + p(y = k|x)) x \end{aligned}$$

32.6 Q6

1. Uniform sampling without replacement,

$$\mathbb{E} \left[\sum_{j \in S_i} \exp(w_j^T x_i) \right] = \frac{m-1}{c-1} \sum_{j \neq i} \exp(w_j^T x_i)$$

2. Use Popoviciu, and $-1 \leq w_j^T x_i \leq 1$

$$\mathbb{V} \left[\sum_{j \in S_i} \exp(w_j^T x_i) \right] \leq \frac{1}{4} (m-1) \left(e - \frac{1}{e} \right)^2$$

3. Use Markov,

$$\begin{aligned} & \mathbb{P} \left\{ \left| \sum_{j \in S_i} \exp(w_j^T x_i) - \mu \right| \geq \varepsilon \right\} \\ &= \mathbb{P} \left\{ \left(\sum_{j \in S_i} \exp(w_j^T x_i) - \mu \right)^2 \geq \varepsilon^2 \right\} \\ &\leq \frac{1}{\varepsilon^2} \mathbb{E} \left[\left(\sum_{j \in S_i} \exp(w_j^T x_i) - \mu \right)^2 \right] \\ &\leq \frac{1}{\varepsilon^2} \sigma^2 \end{aligned}$$

Use $\delta = \frac{1}{\varepsilon^2} \sigma^2$,

$$\varepsilon = \sqrt{\frac{\sigma^2}{\delta}}$$

Therefore,

$$\mathbb{P} \left\{ \sum_{j \in S_i} \exp(w_j^T x_i) - \mu \leq \sqrt{\frac{\sigma^2}{\delta}} \right\} = 1 - \delta$$

4. With probability $1 - \delta$,

$$\sum_{j \in S_i} \exp(w_j^T x_i) \leq \frac{m-1}{c-1} \sum_{j \neq i} \exp(w_j^T x_i) + \sqrt{\frac{\sigma^2}{\delta}}$$

Therefore,

$$\begin{aligned} & p(y_i = k | x_i) - \tilde{p}(y_i = k | x_i) \\ &\leq \frac{\exp(w_k^T x_i)}{\sum_{j=1}^c \exp(w_j^T x_i)} + \frac{\exp(w_k^T x_i)}{\exp(w_k^T x_i) + \frac{m-1}{c-1} \sum_{j \neq i} \exp(w_j^T x_i) + \sqrt{\frac{\sigma^2}{\delta}}} \end{aligned}$$

32.7 Q7

1. $X \sim \text{Binomial}\left(n, \frac{\pi}{4}\right)$, estimate by,

$$\hat{\pi} = \frac{4X}{n}$$

2. Mean and variance are,

$$\begin{aligned}\mathbb{E}[\hat{\pi}] &= \mathbb{E}\left[\frac{4X}{n}\right] \\ &= \frac{4}{n}n\frac{\pi}{4} \\ &= \pi \\ \mathbb{V}[\hat{\pi}] &= \left(\frac{4}{n}\right)^2 n\frac{\pi}{4}\left(1 - \frac{\pi}{4}\right) \\ &= \frac{1}{n}\pi(4 - \pi)\end{aligned}$$

3. Use Chebyshev,

$$\begin{aligned}\mathbb{P}\{|\hat{\pi} - \pi| \geq 0.001\} &\leq \frac{\mathbb{V}[\hat{\pi}]}{0.001^2} \\ &= \frac{1}{n}\pi(4 - \pi)1000000\end{aligned}$$

32.8 Q8

1. Factorization Theorem,

$$\begin{aligned}p(x_1, x_2 | \mu) &= c(\mu, \theta) \exp\left(-\frac{1}{2} \frac{(x_1 - \mu)^2}{\sigma_1^2} + \frac{(x_2 - \mu)^2}{\sigma_2^2}\right) \\ &= c(\mu, \theta) \exp\left(\mu \left(\frac{x_1}{\sigma_1^2} + \frac{x_2}{\sigma_2^2}\right)\right) \\ \Rightarrow T(x_1, x_2) &= \frac{x_1}{\sigma_1^2} + \frac{x_2}{\sigma_2^2}\end{aligned}$$

2. The expectation,

$$\mathbb{E}[T(x_1, x_2)] = \mu \left(\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}\right)$$

3. Unbiased estimator,

$$f(x_1, x_2) = \frac{1}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}} T(x_1, x_2)$$

4. Same

32.9 Q9

1. MLE is,

$$\hat{\alpha} = \arg \max_{\alpha} \sum_{i=1}^n \log(\alpha) - (1 + \alpha) \log(x_i)$$

$$\begin{aligned}
&= \alpha : \frac{n}{\alpha} - \sum_{i=1}^n \log(x_i) = 0 \\
&= \frac{n}{\sum_{i=1}^n \log(x_i)}
\end{aligned}$$

2. A sufficient statistic is,

$$T(x) = \sum_{i=1}^n \log(x_i)$$

3. Yes.

$$\hat{\alpha} = \frac{n}{T(x)}$$

32.10 Q10

$$\begin{aligned}
D(p\|q) &= \mathbb{E}_p \left[\frac{\log(p(x))}{\log(q(x))} \right] \\
&= \mathbb{E}_\theta \left[\frac{1}{2\sigma^2} ((x-\theta^*)^2 - (x-\theta)^2) \right] \\
&= \mathbb{E}_\theta \left[\frac{1}{2\sigma^2} (2x\theta - 2x\theta^* - \theta^2 + (\theta^*)^2) \right] \\
&= \mathbb{E}_\theta \left[\frac{1}{2\sigma^2} (2\theta^2 - 2\theta\theta^* - \theta^2 + (\theta^*)^2) \right] \\
&= \frac{1}{2\sigma^2} (\theta - \theta^*)^2 \\
\frac{\partial^2 D(\theta\|\theta^*)}{\partial \theta^2} &= \frac{1}{\sigma^2} \\
\mathbb{E} \left[-\frac{\partial^2 \log p(x|\theta)}{\partial \theta^2} \right] &= \mathbb{E} \left[-\frac{\partial^2}{\partial \theta^2} \left(\frac{1}{2\sigma^2} (x-\theta)^2 \right) \right] \\
&= \mathbb{E} \left[-\left(-\frac{1}{\sigma^2} \right) \right] \\
&= \frac{1}{\sigma^2}
\end{aligned}$$

32.11 Q11

1. Poisson distribution,

$$\begin{aligned}
p(y) &= \frac{1}{k!} \lambda^y e^{-\lambda} \\
&= \frac{1}{k!} \exp(y \log(\lambda) - \lambda) \\
&= b(y) \exp(\theta y - a(\theta)), \theta = \log(\lambda)
\end{aligned}$$

2. $\theta = \log(\lambda)$

1. The loglikelihood,

$$l(x|\theta) = \sum_{i=1}^n (y_i w^T x_i - \exp(w^T x_i) - \log(y_i!))$$

2. The derivative,

$$\frac{\partial}{\partial w} l(x|w) = \sum_{i=1}^n (y_i x_i - \exp(w^T x_i) x_i)$$

32.12 Q12

1. No

1. No

1. Second is incorrect.

1. Use bounds,

$$\gamma < \frac{2}{\lambda_{\max}(X^T X)}$$

where,

$$(\lambda - 5)(\lambda - 6) - 9 = 0$$

$$\lambda^2 - 11\lambda + 21 = 0$$

$$\lambda = \frac{11}{2} \pm \frac{1}{2}\sqrt{37} < \frac{17}{2}$$

Therefore,

$$\gamma = 0.2$$

$$< \frac{\frac{2}{17}}{\frac{2}{2}} = \frac{4}{17}$$

2. Use formula,

$$w_1 = w_0 + \gamma x^T (y - x^T w)$$

$$= \begin{bmatrix} -1 \\ 0.8 \\ -0.7 \end{bmatrix}$$

33 Midterm

33.1 Q1

1. No.

1. No.

1. One point is misclassified.

1. Empirical error rate is,

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\hat{y} \neq y}$$

Actual error rate is,

$$p = \mathbb{P}\{\hat{y} \neq y\}$$

Expected value and variance of empirical error rate is,

$$\begin{aligned} \mathbb{E}[\hat{p}] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\hat{y} \neq y}\right] \\ &= \frac{1}{16} \sum_{i=1}^{16} p \\ &= p \\ \mathbb{V}[\hat{p}] &= \mathbb{V}\left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\hat{y} \neq y}\right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{V}[\mathbb{1}_{\hat{y} \neq y}] \\ &= \frac{1}{n^2} np(1-p) \\ &= \frac{p(1-p)}{n} \end{aligned}$$

Use Markov Inequality (Chebyshev), with $n = 16$ in the diagram,

$$\begin{aligned} \mathbb{P}\left\{p - \hat{p} \geq \frac{1}{4}\right\} &\leq \mathbb{P}\left\{|\hat{p} - \mathbb{E}[\hat{p}]| \geq \frac{1}{16}\right\} \\ &= \mathbb{V}[\hat{p}] \cdot 16 \\ &= 16 \frac{p \cdot (1-p)}{n} \\ &\leq \frac{16}{4n} \\ &= \frac{1}{4} \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{P}\left\{p \leq \hat{p} + \frac{1}{4}\right\} &= 1 - \mathbb{P}\left\{p - \hat{p} \geq \frac{1}{4}\right\} \\ &\geq \frac{3}{4} \end{aligned}$$

33.2 Q2

1. No
1. Same as usual.
1. Same as usual.
1. Same.

33.3 Q3

1. The densities are,

$$p(x|y=0) = \mathbb{1}_{x \in [0,1]}$$

$$p(x|y=1) = \frac{1}{\theta} \mathbb{1}_{x \in [0,\theta]}$$

Therefore, the optimal classifier is,

$$\hat{y} = \begin{cases} 1 & \text{if } x \leq \theta \\ 0 & \text{if } x > \theta \end{cases}$$

The probability of error is,

$$\begin{aligned} \mathbb{P}\{\hat{y} \neq y\} &= \mathbb{P}\{y=0, x \leq \theta\} + \mathbb{P}\{y=1, x > \theta\} \\ &= \mathbb{P}\{x \leq \theta | y=0\} \cdot \mathbb{P}\{y=0\} + \mathbb{P}\{x > \theta | y=1\} \cdot \mathbb{P}\{y=1\} \\ &= \theta \cdot \frac{1}{2} + 0 \cdot \frac{1}{2} \\ &= \frac{\theta}{2} \end{aligned}$$

2. The usual MLE,

$$\begin{aligned} \hat{\theta} &= \max_{\theta} \prod_{i=1}^n \left(\frac{1}{\theta} \mathbb{1}_{x_i \in [0,\theta]} \right) \\ &= \max_{\theta} \left(\frac{1}{\theta^n} \mathbb{1}_{\max_i x_i \in [0,\theta]} \right) \\ &= \max_i x_i \end{aligned}$$

3. The probability is,

$$\begin{aligned} \mathbb{P}\{\hat{\theta} \leq (1-\varepsilon)\theta\} &= \mathbb{P}\left\{\max_i x_i \leq (1-\varepsilon)\theta\right\} \\ &= (\mathbb{P}\{x_i \leq (1-\varepsilon)\theta\})^n \\ &= \left(\frac{(1-\varepsilon)\theta}{\theta}\right)^n \\ &= (1-\varepsilon)^n \end{aligned}$$

4. The probability is given $\hat{\theta}$,

$$\begin{aligned}\mathbb{P}\{\hat{y} \neq y\} &= \mathbb{P}\{y = 0, x \leq \hat{\theta}\} + \mathbb{P}\{y = 1, x > \theta\} \\ &= \mathbb{P}\{x \leq \hat{\theta} | y = 0\} \cdot \mathbb{P}\{y = 0\} + \mathbb{P}\{x > \hat{\theta} | y = 1\} \cdot \mathbb{P}\{y = 1\} \\ &= \hat{\theta} \cdot \frac{1}{2} + \frac{\theta - \hat{\theta}}{\theta} \cdot \frac{1}{2}\end{aligned}$$

Then the probability that the plug in classifier is at most $\frac{\varepsilon}{2}$ greater than the minimum is,

$$\begin{aligned}\mathbb{P}\left\{\mathbb{P}\{\hat{y} \neq y\} \leq \mathbb{P}\{\hat{y} \neq y\} + \frac{\varepsilon}{2}\right\} &= \mathbb{P}\left\{\hat{\theta} \cdot \frac{1}{2} + \frac{\theta - \hat{\theta}}{\theta} \cdot \frac{1}{2} \leq \frac{\theta}{2} + \frac{\varepsilon}{2}\right\} \\ &= \mathbb{P}\left\{\theta\hat{\theta} + \theta - \hat{\theta} - \theta^2 - \varepsilon\theta \leq 0\right\} \\ &= \mathbb{P}\left\{\theta(1 - \varepsilon) - \hat{\theta} \leq \theta(\theta - \hat{\theta})\right\} \\ &\geq \mathbb{P}\left\{\theta(1 - \varepsilon) - \hat{\theta} \leq 0\right\} \\ &= \mathbb{P}\left\{\hat{\theta} \geq (1 - \varepsilon)\theta\right\}\end{aligned}$$

33.4 Q4

1. Use as test set,

$$\hat{p} = \frac{1}{m} \sum_{i=1}^n \mathbb{1}_{\hat{y}_i \neq y_i}$$

2. Only use data in that class,

$$\hat{p}_j = \frac{\sum_{i=1}^n \mathbb{1}_{\hat{y}_i \neq y_i, y_i = j}}{\sum_{i=1}^n \mathbb{1}_{y_i = j}}$$

3. Use the binomial probability again,

$$\begin{aligned}\mathbb{V}[\hat{p}] &= \frac{p(1-p)}{m} \\ \mathbb{V}[\hat{p}_j] &= \frac{p_j(1-p_j)}{n_j}\end{aligned}$$

4. Use the one with smallest test set error, use the one with smallest,

$$k^* = \arg \min_k \left(\max_j \hat{p}_{j,k} \right)$$

Can also use the confidence interval with Chernoff bounds,

$$[L_{j,k}, U_{j,k}] = \left[\hat{p}_{j,k} \pm \sqrt{\frac{\log\left(\frac{6l}{\delta}\right)}{2n_j}} \right]$$

For worse case best performance, use the one with smallest

$$k^* = \arg \min_k \left(\max_j \hat{p}_{j,k} + \sqrt{\frac{\log\left(\frac{6l}{\delta}\right)}{2n_j}} \right)$$

For a fixed δ , confident that k^* is the best if,

$$\hat{p}_{j,k^*} + \sqrt{\frac{\log\left(\frac{6l}{\delta}\right)}{2n_j}} \geq \hat{p}_{j,k} - \sqrt{\frac{\log\left(\frac{6l}{\delta}\right)}{2n_j}} \quad \forall k \neq k^*$$

34 Problem Set 6

34.1 Q1

1. Find the mean and the mode,

$$\begin{aligned} \arg \max_{\theta} p(\theta; \alpha) &= \arg \max_{\theta} \theta^{\alpha-1} (1-\theta)^{\alpha-1} \\ &= \arg \max_{\theta} (\alpha-1) \log(\theta) + (\alpha-1) \log(1-\theta) \\ &= \theta : \frac{1}{\theta} - \frac{1}{1-\theta} = 0 \\ &= \theta : 2\theta = 1 \\ &= \frac{1}{2} \end{aligned}$$

$$\begin{aligned} \mathbb{E}[X] &= \int_{\theta} \text{Beta}(\alpha, \alpha) \theta^{\alpha} (1-\theta)^{\alpha-1} \\ &= \int_{\theta} \text{Beta}(\alpha+1, \alpha) \theta^{\alpha} (1-\theta)^{\alpha-1} \end{aligned}$$

$$\begin{aligned} &= \frac{\Gamma(\alpha+1) \Gamma(2\alpha)}{\Gamma(2\alpha+1) \Gamma(\alpha)} \\ &= \frac{\alpha}{2\alpha} \\ &= \frac{1}{2} \end{aligned}$$

2. The posterior is,

$$p(\alpha|\theta) \propto p(\theta|\alpha) p(\alpha)$$

$$\begin{aligned}
&= \text{Beta}(\alpha, \alpha) \theta^{\alpha-1} (1-\theta)^{\alpha-1} \theta^s (1-\theta)^{N-s} \\
&\propto \theta^{\alpha-1+s} (1-\theta)^{\alpha-1+N-s}
\end{aligned}$$

This is the kernel of another Beta distribution.

3. The mean is,

$$\frac{\alpha + s}{2\alpha + n}$$

As $n \rightarrow \infty, s \rightarrow n\theta$

$$\frac{\alpha + s}{2\alpha + n} \rightarrow \theta$$

34.2 Q2

From Q1,

$$\hat{\theta} = \frac{\alpha + x}{2\alpha + n}$$

Since $x \sim \text{Bin}(\theta, n)$,

$$\begin{aligned}
M(\hat{\theta}) &= \mathbb{V}[\hat{\theta}] + \left(\theta - \mathbb{E}[\hat{\theta}]\right)^2 \\
&= \frac{n\theta(1-\theta)}{(2\alpha + n)^2} + \left(\theta - \frac{\alpha + n\theta}{2\alpha + n}\right)^2 \\
&= \frac{n\theta(1-\theta) + \alpha^2(2\theta-1)^2}{(2\alpha + n)^2}
\end{aligned}$$

Need the derivative to be 0,

$$\begin{aligned}
n(1-2\theta) + 4\alpha^2(2\theta-1) &= 0 \\
\Rightarrow 4\alpha^2 &= n \\
\Rightarrow \alpha &= \frac{1}{2}\sqrt{n}
\end{aligned}$$

Therefore, the minimax optimal estimator is,

$$\hat{\theta} = \frac{\sqrt{n} + 2x}{2\sqrt{n} + 2n}$$

34.3 Q3

No.

34.4 Q4

Maximize each coordinate separately,

$$\begin{aligned} & \min_{w_i} (y_i - w_i)^2 + \lambda \mathbb{1}_{w_i \neq 0} \\ & \Rightarrow -2(y_i - w_i) + \lambda = 0 \text{ or } w_i = 0 \\ & \Rightarrow w_i = y_i \mathbb{1}_{y_i^2 > \lambda} \end{aligned}$$

34.5 Q5

1. SGD step is,

$$w_{t+1} = w_t + \gamma_t x_i (y_i - x_i^T w_t)$$

2. SGD step is,

$$w_{t+1} = w_t + \gamma_t x_i (y_i - x_i^T w_t) - \lambda \gamma_t w_t$$

3. SGD step is,

$$w_{t+1} = w_t + \gamma_t x_i (y_i - x_i^T w_t) - \text{sign}(w_t) \gamma_t$$

4. SGD step is,

$$w_{t+1} = w_t + \gamma_t \frac{-y_i x_i \exp(-y_i w^T x_i)}{1 + \exp(-y_i w^T x_i)}$$

5. SGD step is,

$$w_{t+1} = w_t + \gamma_t y_i x_i \mathbb{1}_{y_i w_t^T x_i > 1}$$

6. SGD step is,

$$w_{t+1} = w_t + \gamma_t y_i x_i \mathbb{1}_{y_i w_t^T x_i > 1} - \lambda \gamma_t w_t$$

35 Problem Set 7

35.1 Q1

1. Show that $v^T K v \geq 0$ for all v ,

$$\begin{aligned} v^T K v &= v^T x x^T v \\ &= (v^T x)^2 \\ &\geq 0 \end{aligned}$$

2. Use sums and constant products preserve PSD,

$$(1 + x^T x')^p = \sum_{i=1}^n \binom{p}{i} (x^T x')^i$$

3. Same as (a)

35.2 Q2

1. Show that $v^T K v \geq 0$ for all v ,

$$\begin{aligned} v^T K v &= v^T (K_1 + K_2) v \\ &= v^T K_1 v + v^T K_2 v \\ &\geq 0 \end{aligned}$$

2. Use eigenvalue decomposition,

$$\begin{aligned} K_1 &= \sum_{i=1}^n \alpha_i a_i a_i^T \\ K_2 &= \sum_{i=1}^n \beta_i b_i b_i^T \\ v^T K v &= v^T (K_1 + K_2) v \\ &= v^T \left(\sum_{i=1}^n \alpha_i a_i a_i^T + \sum_{i=1}^n \beta_i b_i b_i^T \right) v \\ &= v^T \left(\sum_{i,j} \alpha_i \beta_j (a_i \cdot b_j) (a_i \cdot b_j)^T \right) v \\ &= \sum_{i,j} \alpha_i \beta_j (v^T a_i \cdot b_j)^2 \\ &\geq 0 \end{aligned}$$

3. From (a) and (b)

1. Taylor expansion and use (c)

35.3 Q3

1. If $f(x) = 1$,

$$\begin{aligned} \|\phi(x) - \mu_1\| &\leq \|\phi(x) - \mu_0\| \\ \Rightarrow \|\phi(x)\|^2 - 2\langle \mu_1, \phi(x) \rangle + \|\mu_1\|^2 &\leq \|\phi(x)\|^2 - 2\langle \mu_0, \phi(x) \rangle + \|\mu_0\|^2 \\ \Rightarrow \langle \mu_1 - \mu_0, \phi(x) \rangle &\geq \frac{1}{2} (\|\mu_1\|^2 - \|\mu_0\|^2) \\ \Rightarrow \langle w, \phi(x) \rangle + b &\geq 0 \end{aligned}$$

2. The first term is,

$$\begin{aligned}
& \langle w, \phi(x) \rangle = \langle \mu_1, \phi(x) \rangle - \langle \mu_0, \phi(x) \rangle \\
& = \langle \frac{1}{n_1} \sum_{y_i=1} \phi(x_i), \phi(x) \rangle - \langle \frac{1}{n_0} \sum_{y_i=0} \phi(x_i), \phi(x) \rangle \\
& = \sum_y \frac{1}{n_y} \sum_{y_i=y} \langle \phi(x_i), \phi(x) \rangle \\
& = \sum_y \frac{1}{n_y} \sum_{y_i=y} k(x_i, x)
\end{aligned}$$

3. Use representer theorem,

$$\begin{aligned}
L(\alpha) &= \sum_{i=1}^n \left(y_i - \langle \sum_{j=1}^n \alpha_j \phi(x_j), \phi(x_i) \rangle \right)^2 + \lambda \langle \sum_{i=1}^n \alpha_i \phi(x_i), \sum_{i=1}^n \alpha_i \phi(x_i) \rangle \\
&= \sum_{i=1}^n \left(\left(y_i - \sum_{j=1}^n \alpha_j K_{ij} \right)^2 + \lambda \sum_{j=1}^n \alpha_i \alpha_j K_{ij} \right) \\
\frac{\partial L}{\partial \alpha_j} &= \sum_{i=1}^n \left(2K_{ij} \left(y_i - \sum_{j=1}^n \alpha_j K_{ij} \right) + 2\lambda \alpha_i K_{ij} \right) \\
&= \sum_{i=1}^n 2K_{ij} \left(y_i - \sum_{k=1}^n \alpha_k K_{ik} \right) + \sum_{i=1}^n 2\lambda \alpha_i K_{ij}
\end{aligned}$$

With only row i_t from the data,

$$\begin{aligned}
\frac{\partial L}{\partial \alpha_j} &= 2K_{i_t j} \left(y_{i_t} - \sum_{k=1}^n \alpha_k K_{i_t k} \right) + 2\lambda \alpha_{i_t} K_{i_t j} \\
&= 2K_{i_t j} (y_{i_t} - K_{i_t} \alpha) + 2\lambda \alpha_{i_t} K_{i_t j}
\end{aligned}$$

Put the derivatives together for all j to get the gradient,

$$\nabla L = 2K_{i_t}^T (y_{i_t} - K_{i_t} \alpha) + 2\lambda K_{i_t} \alpha_{i_t}$$

Note that α and α_t are the same thing, and add this to the SGD step,

$$\begin{aligned}
\alpha_{t+1} &= \alpha_t - \gamma_t \nabla L \\
&= \alpha_t - \gamma_t (2K_{i_t}^T (y_{i_t} - K_{i_t} \alpha) + 2\lambda K_{i_t} \alpha_{i_t}) \\
&= (1 - 2\lambda \gamma_t K_{i_t} e_{i_t}^T) \alpha_t - 2\gamma_t K_{i_t}^T (K_{i_t} \alpha_t - y_{i_t})
\end{aligned}$$

where $e_{i_t}^T$ is $\left[0, 0, \dots, 0, \underbrace{1}_{\text{position } i_t}, 0, \dots, 0 \right]$.

Therefore, the SGD step is,

$$\alpha_{t+1} = \alpha_t - \gamma_t (2K_{i_t}^T (K_{i_t} \alpha - y_{i_t}) + 2\lambda K_{i_t} \alpha_t)$$

35.4 Q4

1. There are d terms in the form x_{ij} , d terms in the form x_{ij}^2 and $\binom{d}{2}$ terms in the form $x_{ij}x_{ik}$ and 1 constant term,

$$1 + 2d + \binom{d}{2} = \frac{1}{2} (d+1)(d+2)$$

Therefore the rank is at most,

$$\min \left\{ n, \frac{1}{2} (d+1)(d+2) \right\}$$

2. The rank is at most,

$$\min \left\{ n-1, \frac{1}{2} (d+1)(d+2) \right\}$$

35.5 Q5

All three are quadratic,

Left: $\text{sign}(x_1, x_2)$

Middle: $\text{sign}(ax_1 + b)$

Right: $\text{sign}(a(x_1 - h)^2 + b(x_1 - k)^2 + r)$

36 Problem Set 8

36.1 Q1

1. Larger than 0 parts.

$$y_1 = \max\{x_1, 0\}$$

$$y_2 = \max\{x_2, 0\}$$

2. The required mapping is,

$$\begin{aligned} y_1 &= a_{11}x_1 + a_{12}x_2 \\ &= a_{11} \max\{0, x_1\} - a_{11} \max\{0, -x_1\} + a_{12} \max\{0, x_2\} - a_{12} \max\{0, -x_2\} \end{aligned}$$

Similar for the other one.

3. The specific mapping is,

$$\begin{aligned} y_1 &= \frac{1}{2} (x_1 - x_2) \\ y_2 &= \frac{1}{2} (x_2 - x_1) \\ &= -y_1 \end{aligned}$$

4. The region is,

$$\begin{aligned}y_1 &\geq y_2 \geq 0 \\y_1 &= \max \{x_1, 0\} + \max \{x_2, 0\} \\y_2 &= \max \{x_1, 0\}\end{aligned}$$

36.2 Q2

1. A tree can be trained to have zero error on any training set, just split on every data point.

1. Error occurs if the curve pass through a cell: 11 cells out of 64 on both diagrams.

$$\frac{11}{64}$$

2. Suppose there are m cells, and $\frac{1}{2}$ of cells are on one side, then for both histogram and decision trees,

$$\mathbb{P}\{y = 1 | x_i \in B_i, y_i = 0\} = p \frac{1}{2} + (1 - p) \frac{1}{2} = \frac{1}{2}$$

3. The question is asking for the number of nodes in a tree with $2\sqrt{m}$ leaves. The fewest is a balanced tree,

$$4\sqrt{m}$$

4. Same as (d)

$$2O(m^{d-1}) = O(m^{d-1})$$

36.3 Q3

The collection of predictors are rankings,

$$\mathcal{F} = \{f_h : h \text{ is a ranking}\}$$

The size is $m! = O(m^m)$,

$$\begin{aligned}\log(m!) &= O(m \log m) \\ \mathbb{E}\left[R(\hat{f})\right] &\leq R(f^*) + O\left(\sqrt{\frac{\log |\mathcal{F}| + \log n}{n}}\right) \\ &= R(f^*) + \sqrt{\frac{m \log m + \log n}{n}}\end{aligned}$$

36.4 Q4

Hinge loss ignores the error in the positive half, the classifier is the mid point between the separation boundary, 5.11 and 6.1. Therefore, the classifier is,

$$\text{sign}\left(x - \frac{1}{2}(5.11 + 6.1)\right)$$

The support vector is the minimizer of the dual,

$$\alpha = (0, 0, 1, 1)$$

36.5 Q5

1. y is constant

$$\hat{y} = x_1$$

2. Solve the following system

$$w_1 + w_2 = 1$$

$$w_1 + w_3 = 1$$

$$w = (1 + c, -c, -c).$$

3. Minimize by choosing c ,

$$\begin{aligned}\min_w \|w\|^2 &= \min_c (1 + c)^2 + c^2 + c^2 \\ &= \min_c 3c^2 + 2c + 1 \\ &\Rightarrow 6c + 2 \\ &\Rightarrow c = -\frac{1}{3} \\ &\Rightarrow w = \left(\frac{2}{3}, -\frac{1}{2}, -\frac{1}{2}\right)\end{aligned}$$

4. The minimization is,

$$\begin{aligned}\min_{\alpha} \|y - XX^T \alpha\|^2 &= \min_{\alpha} \left\| \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} \right\|^2 \\ &= \min_{\alpha} \left\| \begin{bmatrix} 1 - 2\alpha_1 - \alpha_2 \\ 1 - \alpha_1 - 2\alpha_2 \end{bmatrix} \right\|^2 \\ &= \min_{\alpha} (1 - 2\alpha_1 - \alpha_2)^2 + (1 - \alpha_1 - 2\alpha_2)^2 \\ &\Rightarrow 1 - 2\alpha_1 - \alpha_2 = 1 - \alpha_1 - 2\alpha_2 = 0 \\ &\Rightarrow \alpha_1 = \alpha_2 = \frac{1}{3}\end{aligned}$$

$$\Rightarrow w = X^T \alpha = \left(\frac{2}{3}, -\frac{1}{2}, -\frac{1}{2} \right)$$

5. Use L_1 norm,

$$\min_c |1 + c| + |-c| + |-c| \Rightarrow w = (1, 0, 0)$$

36.6 Q6

1. Consider the linear model

$$\begin{aligned} y|w &\sim N(Xw, \sigma_v^2 I), w \sim N(0, \sigma_w^2 I) \\ \hat{w} &= \arg \max_w \log p(y|w) + \log p(w) \\ &= \arg \min_w \frac{1}{2\sigma_v^2} \|y - Xw\|^2 + \frac{1}{2\sigma_w^2} \|w\|^2 \\ &= \arg \min_w \|y - Xw\|^2 + \lambda \|w\|^2, \lambda = \frac{\sigma_v^2}{\sigma_w^2} \end{aligned}$$

2. Redefine the prior,

$$w \sim N(0, \text{diag}(\sqrt{2}\sigma_w, \sigma_w, \dots, \sigma_w))$$

3. Take the derivative,

$$\begin{aligned} -2X^T(y - Xw) + 2\lambda w &= 0 \\ \Rightarrow X^T y &= (X^T X + \lambda I)^{-1} w \\ \Rightarrow w &= (X^T X + \lambda I)^{-1} X^T y \\ \Rightarrow w &= \frac{1}{1 + \lambda} y \end{aligned}$$

For the other problem,

$$\begin{aligned} w &= \left(X^T X + \lambda \text{diag}\left(\frac{1}{2}, 1, \dots, 1\right) \right)^{-1} X^T y \\ \Rightarrow w &= \text{diag}\left(\frac{1}{1 + \frac{1}{2}\lambda}, \frac{1}{1 + \lambda}, \dots, \frac{1}{1 + \lambda}\right) y \end{aligned}$$

4. Use Bernoulli Normal,

$$\begin{aligned} y|w &\sim \text{Ber}(\mu_i), \mu_i = \frac{1}{1 + e^{-w^T x_i}}, w \sim N(0, \sigma_w^2 I) \\ \hat{w} &= \arg \max_w \log p(y|w) + \log p(w) \\ &= \arg \min_w \log \left(\mu_i^{y_i} (1 - \mu_i)^{1 - y_i} \right) + \lambda \|w\|^2 \\ &= \arg \min_w y_i \log \mu_i + (1 - y_i) \log (1 - \mu_i) + \lambda \|w\|^2 \end{aligned}$$

$$= \arg \min_w \log \left(1 + e^{-y_i w^T x_i} \right) + \lambda \|w\|^2$$

5. Use singular value decomposition,

$$X = UDV^T, \tilde{X} = U_{[1:r]}$$

36.7 Q7

1. The MLEs are,

$$p_i = \frac{n_i}{n}$$

$$p_{ij} = \frac{n_{ij}}{n_i}$$

2. Use x_k to predict,

$$\hat{j} = \arg \max_j p_{x_k j}$$

3. Max over both,

$$\hat{j} = \arg \max_{ij} p_{x_k i} p_{ij}$$

4. Buy stocks.

36.8 Q8

1. Use the one with minimum loss (error)

$$\hat{f} = \arg \min_{f_m} \frac{1}{n} \sum_{i=1}^n l(y_i, f_m(x_i))$$

2. The truly best is,

$$f^* = \arg \min_{f_m} \mathbb{E}[l_m]$$

$$\mathbb{P} \left\{ \hat{R}(f_m) - R(f_m) > t \right\} \leq \mathbb{E} \left[e^{\frac{\lambda}{n} \sum_{i=1}^n X_i} \right] e^{-\lambda t}$$

$$= \mathbb{E} [e^{\lambda X_i}] e^{-\lambda t}$$

$$= e^{\frac{\lambda^2}{8n} - \lambda t}$$

$$\leq e^{-2nt^2}, \lambda = 4nt$$

Then, union bound,

$$\mathbb{P} \left\{ |\hat{R}(f_m) - R(f_m)| > t \right\} \leq 2e^{-2nt^2}$$

Union bound again,

$$\begin{aligned} \mathbb{P} \left\{ \bigcup_{f_m} |\hat{R}(f_m) - R(f_m)| > t \right\} &\leq 2Me^{-2nt^2} \\ \mathbb{P} \left\{ R(\hat{f}) - \hat{R}(\hat{f}) < t \right\} &\leq 2Me^{-2nt^2} \\ \mathbb{P} \left\{ R(\hat{f}) - \hat{R}(f^*) < t \right\} &\leq 2Me^{-2nt^2} \\ \mathbb{P} \left\{ R(\hat{f}) - R(f^*) < 2t \right\} &\leq 2Me^{-2nt^2} \end{aligned}$$

3. Given $\hat{R}(\hat{f}) = 0.1$,

$$\begin{aligned} \mathbb{P} \{ 0.1 - R(f^*) < 2t \} &\leq 2Me^{-2nt^2} \\ \delta &= 2Me^{-2nt^2} \\ n &= \frac{1}{2t^2} \log \frac{2M}{\delta} \end{aligned}$$

36.9 Q9

1. Given $x|w \sim N(w, \sigma_v^2)$

$$\begin{aligned} \hat{w} &= \arg \max_w \frac{1}{\sigma_v^2} \|x - w\|^2 \\ &= \arg \max_w \|x - w\|^2 \\ &= x \\ \mathbb{E}[w] &= \mathbb{E}[x] \\ &= w \\ \mathbb{V}[w] &= \sigma_v^2 \end{aligned}$$

Same for MSE.

2. Use Gauss-Markov,

$$\begin{aligned} \hat{w} &= \mathbb{E}[w|x] \\ &= \mathbb{E}[w] + \frac{\text{Cov}[w, x]}{\mathbb{V}[X]} (x - \mathbb{E}[x]) \\ &= \frac{\sigma_\theta^2 x}{\sigma_\theta^2 + \sigma_v^2} \end{aligned}$$

3. This is Quiz 5

36.10 Q10

1. The empirical risk is,

$$\min_w \hat{R}(w) = \min_w \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2$$

2. Approximate by a partition of M

$$\begin{aligned} |w^T x - \tilde{w}^T x| &\leq \|w - \tilde{w}\| \|x\| \\ &\leq \frac{\sqrt{d}}{M} \end{aligned}$$

The size of the set is M^d

3. Same as Q8

1. Set δ

$$\begin{aligned} \delta &= 2M^d e^{-2nt^2} \\ M &= \left(\frac{1}{2} \delta e^{2nt^2} \right)^{\frac{1}{d}} \end{aligned}$$

2. Use a smaller space W_s the partitions with at most k non-zero entries.

1. The size of W_s is

$$|W_s| = \sum_{i=0}^k \binom{d}{i} (M-1)^i$$

Then same as Q8

- **Definitions :**

- Bayes Risk is $R^* = \inf_f R(f) = \inf_f \mathbb{E}[l(f, X, Y)] \stackrel{0-1 \text{ loss}}{=} \inf_f \mathbb{E}[\mathbb{1}_{f(X) \neq Y}] = \mathbb{P}\{f(X) \neq Y\}$.
- Optimal Bayes Classifier is $f^*(x) = 1$ if $\eta(x) \geq \frac{1}{2}$ and 0 otherwise; and equivalently $f^*(x) = 1$ if $\frac{\eta(x)}{1-\eta(x)} \geq 1$ and 0 otherwise; and equivalently $f^*(x) = 1$ if $\frac{p(x|Y=1)p(Y=1)}{p(x|Y=0)p(Y=0)} \geq 1$ and 0 otherwise.
- Log Likelihood Ratio: $\Lambda(x) = \log\left(\frac{p_1(x)}{p_0(x)}\right)$, where $p_j(x) = p(x|Y=j)$.
- Bayes Cost: $C = \sum_{i,j=0}^1 c_{i,j} \pi_j \mathbb{P}\{\text{decide } H_i | H_j\} = \sum_{i,j=0}^1 c_{i,j} \pi_j \int_{R_i} p_j(x) dx$, where $\pi_j = \mathbb{P}\{H_j\}$ and $R_j = \{x : \text{decide } H_j\}$.
- MLE Risk: $R_{MLE}(q, p_\theta) = \mathbb{E}[-\log p(x|\theta)]$, Excess Risk: $R_{MLE}(q, p_\theta) - R_{MLE}(q, q) = D(q||p_\theta) \geq 0$.
- Probably Approximately Correct: $\mathbb{P}\{R(\hat{f}) - R(f^*) \leq \varepsilon\} \geq 1 - \delta$ where $\varepsilon = \sqrt{\frac{2c^2 \log\left(\frac{2|\mathcal{F}|}{\delta}\right)}{n}}$, take $\delta = \frac{1}{\sqrt{n}}$, $\mathbb{E}[R(\hat{f})] \leq R(f^*) + O\left(\sqrt{\frac{\log|\mathcal{F}| + \log n}{n}}\right)$.

- **Estimators :**

- Empirical means and covariances: $\hat{\mu}_j = \frac{1}{\#\{y_i = j\}} \sum_{i: y_i = j} x_i$ and $\hat{\Sigma} = \frac{1}{n} \left(\sum_j \sum_{i: y_i = j} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T \right)$.
- Gaussian GLM: $p(y|x^T w) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - x^T w)^2\right)$, and $\hat{w} = (X^T X)^{-1} X^T y$.
- Binomial GLM: $p(y|x^T w) = \exp\left(y \log\left(\frac{1}{1 + e^{-x^T w}}\right) + (1 - y) \log\left(\frac{1}{1 + e^{x^T w}}\right)\right)$.
- Multinomial GLM: $p(y|x^T w) = \frac{\exp(x^T w_l)}{\sum_{j=1}^k \exp(x^T w_j)}$.
- Max a Posterior: $\theta_{MAP} = \max_{\theta} p(\theta|y) \propto p(y|\theta) p(\theta)$ to minimize loss $L(\theta, \hat{\theta}) = \mathbb{1}_{\{\|\hat{\theta} - \theta\| > \varepsilon\}}$.
- Bayesian minimum MSE estimator: $\hat{\theta} = \mathbb{E}[\theta|y]$ to minimize loss $L(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|^2$.
- Bayesian minimum MAE estimator: $\hat{\theta} = \text{median}[\theta|y]$ to minimize loss $L(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|_1$.
- Gaussian Penalty Function: $\min_w \log p(y|w^T x) + \lambda \|w\|^2$.
- Laplacian Penalty Function: $\min_w \log p(y|w^T x) + \lambda \|w\|_1$.
- Sparsity Penalty Function: $\min_w \log p(y|w^T x) + \lambda \|w\|_0$.
- Minimax Optimal Estimator: $\hat{\theta} = \arg \min_{\hat{\theta}} \sup_{\theta} R(\hat{\theta}, \theta)$.

- Support Vector Machine: $\min_w \sum_{i=1}^n (1 - y_i w^T \phi(x_i))_+ + \lambda \|w\|^2$, or $\min_w \|w\|$ such that $\sum_{i=1}^n (1 - y_i w^T \phi(x_i))_+ = 0$.
- Exponential Kernel: $K(x, x') = \exp(x^T x')$.
- Gaussian Kernel: $K(x, x') = \exp\left(-\frac{\|x^T - x'^T\|^2}{\sigma^2}\right)$.
- Laplacian Kernel: $K(x, x') = \exp(-\beta \|x^T - x'^T\|)$.
- **Distributions :**
- Multivariate Normal Distribution: $p(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$ or $\log p(x) \propto \log |\Sigma| - \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)$.
- Binomial Distribution: $p(x) = \binom{n}{x} p^x (1 - p)^{n-x}$.
- Hypergeometric Distribution: $p(x) = \frac{\binom{b}{x} \binom{N-b}{n-x}}{\binom{N}{n}}$.
- Multinomial Distribution: $p(x) = \binom{n}{x_1 x_2 \dots x_k} \prod_{i=1}^k p_i^{x_i}$.
- Exponential Distribution: $p(x) = \lambda e^{-\lambda x}$.
- Gamma Distribution: $p(x) = \frac{x^{a-1} \exp\left(-\frac{x}{b}\right)}{\Gamma(a) b^a}$.
- Beta Distribution: $p(x) = \frac{x^{a-1} (1-x)^{b-1}}{\text{Beta}(a, b)}$, where $\text{Beta}(a, b) = \frac{\Gamma(a) \Gamma(b)}{\Gamma(a+b)}$.
- Exponential Family: $p(y|\theta) = b(y) \exp(\theta^T T(y) - \alpha(\theta))$, θ is the natural parameter and $T(y)$ is the sufficient statistic. Canonical form is when $T(y) = y$, and $\log p(y|\theta) = \sum_{i=1}^n (w^T x_i y_i - \alpha(w^T x_i)) + \log b(y_i)$.
- **Other Statistics, Algebra :**
- Kullback-Leibler Divergence: $D(p_1 \| p_0) = \mathbb{E}_{1[\Lambda(X)]} = \int p_1(x) \log \frac{p_1(x)}{p_0(x)} dx$.
- Mahalanobis Distance: $(x - \mu)^T \Sigma^{-1} (x - \mu)$.
- Sufficiency: $t(X)$ is sufficient if $p(x|t, \theta) = p(x|t)$.
- Rao-Blackwellization: If f is an estimator and t is a sufficient statistic, then $\mathbb{E}[f(X)|t(X)]$ is the improved Rao-Blackwell estimator (in terms of MSE).
- Characteristic Equation of X is $\det(\lambda I - X) = 0$, where λ are the eigenvalues.

- Binomial Conjugate Prior: $\text{Binomial}(n, p) + \text{Beta}(\alpha, \beta) = \text{Beta}\left(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i\right)$, and Neg Binomial(r, p) + $\text{Beta}(\alpha, \beta) = \text{Beta}\left(\alpha + \sum_{i=1}^n x_i, \beta + rn\right)$.
- Poisson Conjugate Prior: $\text{Poisson}(\lambda) + \Gamma(k, \beta) = \text{Neg Binomial}\left(k + \sum_{i=1}^n x_i, \beta + n\right)$.
- Normal Conjugate Prior: $\text{Normal}(\mu, \sigma) + \text{Normal}(\mu_0, \sigma_0) = \text{Normal}\left(\frac{1}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}} \left(\frac{n}{\sigma^2} \bar{x} + \frac{\mu}{\sigma_0^2}\right), \frac{1}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}}\right)$
and multivariate version is, $\text{Normal}(\mu, \Sigma) + \text{Normal}(\mu_0, \Sigma_0) = \text{Normal}\left((\Sigma_0^{-1} + n\Sigma^{-1})^{-1} (\Sigma_0^{-1} \mu_0 + n\Sigma^{-1} \bar{x}), (\Sigma_0^{-1} + n\Sigma^{-1})^{-1}\right)$
- Uniform Conjugate Prior: $\text{Uniform}(0, \theta) + \text{Pareto}(x_m, k) = \text{Pareto}(\max\{x_1, \dots, x_n, x_m\}, k + n)$.
- Gamma Conjugate Prior: $\text{Gamma}(\alpha, \beta) + \text{Gamma}(\alpha_0, \beta_0) = \text{Gamma}\left(\alpha_0 + n\alpha, \beta_0 + \sum_{i=1}^n x_i\right)$.
- Dual Optimization: for $\min_w \sum_{i=1}^m l(1 - y_i x_i^T w) + \lambda \|w\|^2$ is $\min_{\alpha} \sum_{i=1}^m l\left(1 - y_i \sum_{j=1}^m \alpha_j x_i^T x_j\right) + \lambda \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j x_i^T x_j = \min_{\alpha} \sum_{i=1}^m l(1 - y_i K \alpha) + \lambda \alpha^T K \alpha$.
- Shattering Coefficient: for class \mathcal{F} is $\mathcal{S}(\mathcal{F}, n) = \max_{x \in X} |\{(f(x_1), \dots, f(x_n)) \in \{-1, +1\}^n, f \in \mathcal{F}\}|$.
- VC Dimension: largest integer k such that $\mathcal{S}(\mathcal{F}, k) = 2^k$.
- VC Dimension of Linear Hyperplane: $d + 1$.
- VC Dimension of Hyper-Rectangle: $2d$.
- VC Dimension of Neural Net: with sign activation: $O(|E| \log |E|)$, with sigmoid activation: $O(|E|^2)$.
- **Inequalities, Bounds :**
- Cauchy-Schwarz Inequality: $|\mathbb{E}[XY]| \leq \sqrt{\mathbb{E}[X^2] \mathbb{E}[Y^2]}$ or $|u^T v| \leq \|u\| \|v\|$.
- Holder's Inequality: For $\frac{1}{p} + \frac{1}{q} = 1$, $\mathbb{E}[|XY|] \leq (E[|X^p|])^{\frac{1}{p}} (E[|Y^q|])^{\frac{1}{q}}$.
- Markov's Inequality: For $X \geq 0$ and $a > 0$, $\mathbb{P}\{X > a\} \leq \frac{\mathbb{E}[X]}{a}$.
- Chebyshev's Inequality: For $t > 0$, $\mathbb{P}\{|X - \mathbb{E}[X]| \geq t\} \leq \frac{\mathbb{V}[X]}{t^2}$.
- Chebyshev-Cantelli Inequality: For $t \geq 0$, $\mathbb{P}\{X - \mathbb{E}[X] > t\} \leq \frac{\mathbb{V}[X]}{\mathbb{V}[X] + t^2}$.
- Jensen's Inequality: if f is convex, $\lambda f(x) + (1 - \lambda) f(y) \geq f(\lambda x + (1 - \lambda) y)$, then $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$.

- Association Inequality: if f and g are increasing, $\mathbb{E}[f(X)g(X)] \geq \mathbb{E}[f(X)]\mathbb{E}[g(X)]$, and if f is increasing, g is decreasing, then the following: $\mathbb{E}[f(X)g(X)] \leq \mathbb{E}[f(X)]\mathbb{E}[g(X)]$.

- Fourth Moment: $\mathbb{E}[|X|] \leq (\mathbb{E}[X^2])^{\frac{3}{2}} (\mathbb{E}[X^4])^{-\frac{1}{2}}$.

- Chernoff bound $(1 - \delta)$ confidence intervals for mean of $x_i \in [0, 1]$ in k dimensions: $\pm \sqrt{\frac{\log\left(\frac{2k}{\delta}\right)}{2n}}$, and for standard deviation: $\sigma = \frac{1}{\sqrt{n}}$. The minimum number of data to ensure ε error with δ probability is $n \geq \frac{1}{2\varepsilon^2} \log\left(\frac{2k}{\delta}\right)$.

- Popoviciu's Inequality: If $\mathbb{P}\{m \leq z \leq M\} = 1$, then $\mathbb{V}[Z] \leq \frac{1}{2}(M - m)^2$.

- Hoeffding's Inequality: $X_i \in [a_i, b_i]$, $S_n = \sum_{i=1}^n X_i$, for each $t > 0$, $\mathbb{P}\{|S_n - \mathbb{E}[S_n]| \geq t\} \leq 2 \exp\left(-2t^2 \left(\sum_{i=1}^n (b_i - a_i)^2\right)^{-1}\right)$.

- Corollary to Hoeffding's Inequality: $X_i \in [a, b]$, $c = (b - a)^2$, $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$, then $\mathbb{P}\{|\hat{\mu} - \mu| \geq t\} \leq 2 \exp\left(-\frac{2nt^2}{c}\right)$.

- Lemma to proof Hoeffding's Inequality: $\mathbb{E}[Z] = 0$, $Z \in [a, b]$, then $\mathbb{E}[e^{sZ}] \leq \exp\left(\frac{1}{8}(s^2(b - a)^2)\right)$.

- Sub Gaussian Tail Bound: if $\{Z_i\}_{i=1}^n$ are independent and $\mathbb{P}\{|Z_i - \mathbb{E}[Z_i]| \geq t\} \leq a \exp\left(-b \frac{t^2}{2}\right)$, then $\mathbb{P}\left\{\frac{1}{n} \sum_i Z_i - \mathbb{E}[Z] > \varepsilon\right\} \leq e^{-cn\varepsilon^2}$ and $\mathbb{P}\left\{\mathbb{E}[Z] - \frac{1}{n} \sum_i Z_i > \varepsilon\right\} \leq e^{-cn\varepsilon^2}$ with $c = \frac{b}{16a}$.

- Gaussian Tail Bound: $\frac{1}{\sqrt{2\pi}} \int_t^\infty \exp\left(-\frac{x^2}{2}\right) dx \leq \min\left\{\frac{1}{2} \exp\left(-\frac{t^2}{2}\right), \frac{1}{\sqrt{2\pi}t^2} \exp\left(-\frac{t^2}{2}\right)\right\}$

- Exponential Bound: If $\mathbb{P}\{|X_i - \mathbb{E}[X_i]| \geq t\} \leq a \exp\left(-\frac{bt^2}{2}\right)$, then $\mathbb{E}[e^{s(X_i - \mathbb{E}[X_i])}] \leq \exp\left(\frac{4as^2}{b}\right)$

- Empirical Risk Minimization: for loss functions bounded by c , $\mathbb{P}\{\hat{R}(f) - R(f) > t\} \leq \frac{c^2}{2nt^2}$ by Markov, $\leq 2 \exp\left(-\frac{2nt^2}{c^2}\right)$ by Chernoff and Union bound. If finite, $\mathbb{P}\{\hat{R}(f) - R(f) > t\} \leq \exp(-2nt^2)$. If multiple f , $\mathbb{P}\left\{\bigcup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| > t\right\} \leq 2|\mathcal{F}| \exp\left(-\frac{2nt^2}{c^2}\right)$.

- Empirical Risk Comparison: with probability δ , $R(\hat{f}) \leq \hat{R}(\hat{f}) + t \leq \hat{R}(f^*) + t \leq R(f^*) + 2t$, where

$$t = \sqrt{\frac{c^2 \log \frac{2|\mathcal{F}|}{\delta}}{2n}}.$$

- VC Dimension: $\mathbb{P}\left\{\sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| > \varepsilon\right\} \leq 8\mathcal{S}(\mathcal{F}, n) \exp\left(-\frac{n\varepsilon^2}{32}\right)$ and $\mathbb{E}\left[\sup_{f \in \mathcal{F}} |\hat{R}_n(f) - R(f)|\right] \leq 2\sqrt{\frac{\log \mathcal{S}(\mathcal{F}, n) + \log 2}{n}}$.

- Uniform Derivation Bound: $\mathbb{P} \left\{ \sup_{f \in \mathcal{F}} R(f) - \hat{R}_n(f) \geq 1 - \delta \right\} \leq 6 \sqrt{\frac{\text{VC}(\mathcal{F}) \log\left(\frac{n}{\delta}\right)}{n}}.$
- **Linear Algebra** :
- Singular Value Decomposition: $A = U \Sigma V^T$ satisfy $Av_i = \sigma_i u_i, A^T u_i = \sigma_i v_i$, where $U^T U = U U^T = V^T V = V V^T = I$.
- Schur Complements: $\begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} I & 0 \\ C A^{-1} & I \end{bmatrix} \begin{bmatrix} A & 0 \\ 0 & D - C A^{-1} B \end{bmatrix} \begin{bmatrix} I & A^{-1} B \\ 0 & I \end{bmatrix} = \begin{bmatrix} I & B D^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} A - B D^{-1} C & 0 \\ 0 & D \end{bmatrix} \begin{bmatrix} I & B D^{-1} \\ 0 & I \end{bmatrix}.$
- Matrix Inversion Lemma: $(A - B D^{-1} C)^{-1} = A^{-1} + A^{-1} B (D - C A^{-1} B)^{-1} C A^{-1}.$
- Sherman-Morrison Formula: $(A^{-1} u v^T)^{-1} = A^{-1} - \frac{A^{-1} u v^T A^{-1}}{1 + v^T A^{-1} u}.$
- Vector derivatives: $\frac{d c^T x}{d x} = c, \frac{d x^T x}{d x} = 2x, \frac{d x^T A x}{d x} = (A + A^T) x.$
- Diagonal Metrix: $(\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T = \Sigma^T (\Sigma^T \Sigma + \lambda I)^{-1}.$
- Representation Theorem: $\arg \min_w \sum_{i=1}^n l(y_i w^T \phi(x_i)) - \lambda \|w\|^2 = \sum_{i=1}^n \alpha_i \phi(x_i).$
- Kernel Requirement: A kernel is valid iff the Gram matrix $K_{ij} = K(x_i, x_j)$ is symmetric and PSD.
- **Statistics Theorems** :
- weak Law of Large Numbers: $\mathbb{E}[|X_i|] < \infty \Rightarrow \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} \mathbb{E}[X_i].$
- strong Law of Large Numbers: $\mathbb{E}[|X_i|] < \infty \Rightarrow \frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mathbb{E}[X_i]$ as $n \rightarrow \infty$.
- Central Limit Theorem: $\mathbb{E}[Z_i] = 0, \mathbb{V}[Z_i] = \sigma^2, \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n Z_i \right) \rightarrow^d N(0, \sigma^2).$
- Fisher-Neyman Factorization: $t(X)$ is sufficient iff $p(x|\theta) = a(x) b(t, \theta).$
- Rao-Blackwell Theorem: let $t(X)$ be sufficient and define $g(t(X)) = \mathbb{E}[f(X) | t(X)]$, then $\mathbb{E}[(g(t(X)) - \theta)^2] \leq \mathbb{E}[(f(X) - \theta)^2]$, equal iff $f(X) = g(t(X)).$
- Convergence of Log-Likelihood to KL: $\hat{\theta}_n = \arg \max_{\theta} p(x|\theta) = \arg \min_{\theta} \sum_{i=1}^n \log \frac{q(x_i)}{p(x_i|\theta)} \rightarrow \arg \min_{\theta} D(q||p_{\theta}).$
- Asyptotic Distribution of MLE: Let $\hat{\theta}_n = \arg \max_{\theta} p(x|\theta)$, and $\mathbb{E} \left[\frac{\partial \log p(x|\theta)}{\partial \theta} \right] = 0$, then $\hat{\theta}_n \stackrel{asympt}{\sim} N(\theta, n^{-1} I^{-1}(\theta^*))$ where information matrix $[I(\theta^*)]_{j,k} = -\mathbb{E} \left[\frac{\partial^2 \log p(x|\theta)}{\partial \theta_j \partial \theta_k} \Big|_{\theta=\theta^*} \right].$
- KL-Divergence Information Matrix identity: if $x|\theta \sim N(\theta, \sigma)$, then $\frac{\partial^2 D(p(x|\theta)||p(x|\theta^*))}{\partial \theta^2} \Big|_{\theta=\theta^*} = I(\theta^*).$

- Gauss-Markov Theorem: $\begin{bmatrix} x \\ y \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$, then $y|x \sim N(\mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(x - \mu_x), \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy})$,
or, $\mathbb{E}[w|x] = \mathbb{E}[w] + \frac{Cov(w, x)}{V(x)}(x - \mathbb{E}[X])$
- Direct Observation Model: $y = W + \varepsilon, \varepsilon \sim N(0, \sigma^2 I)$, and the soft-thresholding estimator $\hat{w}_i = \text{sign}(y_i) \max\{|y_i| - \lambda, 0\}$, oracle estimator $\hat{w}_i = y_i \mathbb{1}_{\{|w_i|^2 \geq \sigma^2\}}$, with $\lambda = \sqrt{2\sigma^2 \log n}$, then $\mathbb{E}[\|\hat{w} - w\|^2] \leq (2 \log n + 1) \left(\sigma^2 + \sum_{i=1}^n \min\{|w_i|^2, \sigma^2\} \right)$.
- **Optimization :**
- Step size choice: if $v_t = w_t - w^* = (I - \gamma X^T X) v_1$, then $v_t \rightarrow 0$ if the eigenvalues of $(I - \gamma X^T X) < 1 \Rightarrow \gamma < \frac{2}{\lambda_{\max}(X^T X)}$
- Constant step size: If $\|\nabla f_t(w)\| \leq G$ and $w^* = \arg \min_w \sum_{t=1}^T f_t(w)$, then gradient descent with $\gamma_t = \gamma$ starting at w_1 satisfies the following: $\frac{1}{T} \sum_{t=1}^T (f_t(w_t) - f_t(w^*)) \leq \frac{\|w_1 - w^*\|^2}{2\gamma T} + \frac{\gamma}{2} G^2$.
- Diminishing step size: If $\|\nabla f_t(w)\| \leq G, \|w^*\| \leq B$ and $w^* = \arg \min_w \sum_{t=1}^T f_t(w)$, then gradient descent with $\gamma_t = \frac{1}{\sqrt{t}}$ starting at w_1 satisfies: $\frac{1}{T} \sum_{t=1}^T (f_t(w_t) - f_t(w^*)) \leq \frac{2B^2 + G^2}{\sqrt{T}}$.
- Subgradients: for $\|w\|_1$ is $\text{sign}(w)$; for $\max\{0, x^T w\}$ is $x \mathbb{1}_{x^T w > 0}$.
- Proximal Gradient: $\min_w f(w) + c(w) \Rightarrow w_k = \text{prox}(w_{k-1} - t \nabla f(w_{k-1}))$, $\text{prox}(v) = \arg \min_u \left(\frac{1}{2} \|u - v\|^2 + t c(u) \right)$.
- 1 Norm Penalty: $\arg \min_w \|y - w\|^2 - \lambda \|w\|_1 = y - \text{sign}(y) \min\{|y|, \lambda\}$.
- 2 Norm Penalty: $\arg \min_w \|y - w\|^2 - \lambda \|w\|^2 = (X^T X + \lambda I)^{-1} X^T y$.
- General loss function: $w = w - 2\mu l'(y_i w^T x_i) y_i x_i$.
- Perception: $w = w + 2\mu \mathbb{1}_{\{y_i w^T x_i < 0\}} y_i x_i$.
- Backprop: given $y_i^{(l)} = f(z_i^{(l)})$, $\frac{\partial J}{\partial y_i^{(l)}} = \frac{\partial J}{\partial y_i^{(l)}} f'(z_i^{(l)}) y_j^{(l-1)}$, and $\frac{\partial J}{\partial y_i^{(l)}} = \sum_j \frac{\partial J}{\partial y_j^{(l+1)}} f'(z_j^{(l+1)}) w_{ji}^{(l+1)}$, $\frac{\partial J}{\partial y_i^{(L)}} = -(y_i - y_i^{(L)})$.
- **Other results, formulas :**
- Empirical Classifier Error: $\hat{p} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\hat{y}_i \neq y_i\}}$ with mean $\mathbb{E}[\hat{p}] = p$ and $\mathbb{V}[\hat{p}] = \frac{p(1-p)}{n}$, where $p = \mathbb{P}\{\hat{y} \neq y\} = \mathbb{E}[\mathbb{1}_{\{\hat{y} \neq y\}}]$ is the actual classifier error.
- KL-Divergence of Normal Distribution: With same variance: $X|Y \sim N(\mu_j, \Sigma)$ with common covariance is $D(p_0 \| p_1) = D(p_1 \| p_0) = \frac{1}{2} (\mu_1 - \mu_0)^T \Sigma^{-1} (\mu_1 - \mu_0)$. With different variances: $D(N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1))$ is $\frac{1}{2} \text{tr}(\Sigma_1^{-1} \Sigma_0) + \frac{1}{2} (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - d + \log \left(\frac{|\Sigma_1|}{|\Sigma_0|} \right)$.

- Optimal Bayes binary Classifier with common covariances and equal prior is $\hat{y}(x) = 1$ if $2(\mu_1 - \mu_0)^T \Sigma^{-1} x \geq \mu_0^T \Sigma \mu_0 - \mu_1^T \Sigma \mu_1$ is linear in x .
- Non-negative Expected Value: If $Y \geq 0$, then $\mathbb{E}[Y] = \sum_{i=1}^{\infty} \mathbb{P}\{Y \geq i\}$.
- Sum formulas: $\sum_{i=1}^n i = \frac{n(n+1)}{2}$; $\sum_{i=1}^n i^2 = \frac{n(n+1)(n+2)}{6}$.
- Bayesian Linear Regression with prior $w \sim N(0, \sigma_w^2 I)$ has $\hat{w} = (X^T X + \lambda I)^{-1} X^T y$, where $\lambda = \frac{\sigma^2}{\sigma_w^2}$.
- Minimax Optimal Estimator: if $\hat{\theta}_p = \arg \min_{\hat{\theta}} \int R(\hat{\theta}, \theta) p(\theta) d\theta$, and $\int R(\hat{\theta}_p, \theta) p(\theta) d\theta = \sup_{\theta} R(\hat{\theta}_p, \theta)$, then $\hat{\theta}_p$ is minimax optimal. In particular, if $R(\hat{\theta}_p, \theta)$ is constant, then it is minimax.
- Two-layer Neural Net: $\hat{y} = W_2 f(W_1 + b_1) + b_2$, $W_2 = \alpha$, $K(x_i, x) = f(x_i^T + b_i)$, $W_1 = x$ is a SVM.
- Random Feature: if $u_1, \dots, u_n, n \geq D$ have continuous density, then polynomial mapping $\Phi_n^T = [\phi(u_1), \dots, \phi(u_n)]$ has full rank D with probability 1.
- Sauer's Lemma: $\mathcal{S}(\mathcal{F}, n) \leq (n+1)^{VC(\mathcal{F})}$