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

$$l(w) = \sum_{i=1}^{n} l_i(w) = \sum_{i=1}^{n} -y_i \log h(x_i) - (1-y_i)\log(1-h(x_i))$$

$$h(x) = g(w^T x) = \frac{1}{1+\exp(-w^T x)}$$

$$P(y=1 \mid x; w) = h(x)$$

$$\nabla l(w) = \sum_{i=1}^{n} -y_i \cdot \nabla \log h(x_i) - (1-y_i)\nabla \log(1-h(x_i))$$

$$= \sum_{i=1}^{n} -y_i \frac{\nabla h(x_i)}{h(x_i)} - (1-y_i)\frac{1-\nabla h(x_i)}{1-h(x_i)}$$

$$\nabla h(x_i) = \nabla\left(1+\exp(-w^T x_i)\right)^{-1} = -1\cdot(-x_i)\cdot\exp(-w^T x_i)\cdot\left(1+\exp(-w^T x_i)\right)^{-2}$$

$$= \frac{x_i \cdot \exp(-w^T x_i)}{\left(1+\exp(-w^T x_i)\right)^2} = x_i \cdot h(x)^2 \cdot \exp(-w^T x_i)$$

$$\nabla l(w) = \sum_{i=1}^{n} -y_i \cdot \frac{x_i \cdot h(x)^2 \cdot \exp(-w^T x_i)}{h(x_i)} - (1-y_i)\cdot \frac{-x_i h(x_i)^2 \cdot \exp(-w^T x_i)}{1-h(x_i)}$$

$$= \sum_{i=1}^{n} -y_i x_i \frac{\exp(-w^T x_i)}{1+\exp(-w^T x_i)} - (1-y_i)\frac{x_i}{1+\exp(-w^T x_i)}.$$

$$= \sum_{i=1}^{n} -y_i x_i(1-h(x_i)) - (1-y_i)x_i \cdot h(x_i)$$

$$= \sum_{i=1}^{n} 2 y_i x_i h(x_i) - y_i x_i - x_i h.$$

(b) $\nabla^2(L(w)) = \nabla X^T h(x)$

$\phantom{(b) \nabla^2(L(w))} = XSX^T$ , where $S$ is diagonal.

$h(x_i) > 0$ and $\exp(-w^T x_i) > 0$.

it is positive semidefinite and $L(w)$ is convex,

the global one.

**2.** (a) $f(x; \alpha, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\alpha)^2}{2\sigma^2}\right)$

$$L(\theta) = \prod_{i=1}^{n} f(x_i; \theta)$$

$$\log L(\alpha; x) = n \cdot \log(2\pi\sigma^2)^{\frac{1}{2}} - \frac{1}{2\sigma^2} \cdot \sum_{i=1}^{n}(x_i - \alpha)^2$$

$$\frac{\partial \log L(\alpha; x)}{\partial \alpha} = -\frac{1}{\sigma^2} \sum_{i=1}^{n}(x_i - \alpha) = 0$$

$$\boxed{\hat{\alpha} = \frac{1}{n}\sum_{i=1}^{n} x_i}$$

$$\frac{\partial \log L(\alpha; x)}{\partial(\sigma)} = -\frac{n}{2\sigma} + \frac{1}{2\sigma^4}\sum(x_i - \alpha)^2$$

$$\begin{cases} \frac{1}{\sigma^2}\sum(x_i - \hat{\alpha}) = 0 \\ \\ -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \cdot \sum(x_i - \alpha)^2 = 0 \end{cases}$$

$$\begin{cases} \hat{\sigma}^2 = \frac{1}{n}\sum_{i=1}^{n}(x - \hat{\alpha})^2 \\ \\ \hat{\alpha} = \frac{1}{n}\sum_{i=1}^{n} x_i \end{cases}$$

b) $f(x; \mu, \Sigma) = \dfrac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\tfrac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$

$L(\mu; x_1, x_2 \dots, x_n) = \dfrac{1}{(2\pi)^{np/2} |\Sigma|^{n/2}} \cdot \exp\left(-\tfrac{1}{2}\sum_{k=1}^{n}(x_k-\mu)^T \Sigma^{-1}(x_k-\mu)\right)$

$x_k, \quad k \in \{1, 2, \dots, n\}$

Let $S$ denotes $\quad S = \displaystyle\sum_{k=1}^{n.}(x_k - \bar{x})(x_k - \bar{x})^T$

$L(\mu; x_1, \dots, x_n) = \dfrac{1}{(2\pi)^{np/2} \cdot |\Sigma|^{n/2}} \cdot \exp\left(-\tfrac{1}{2}\operatorname{tr}(\Sigma^{-1}S) - \tfrac{1}{2}n(\bar{x}-\mu)^T \Sigma(\bar{x}-\mu)\right)$

$\log L(\mu; x_1 \dots, x_n) = \log \dfrac{1}{(2\pi)^{np/2} |\Sigma|^{\frac{n}{2}}} - \dfrac{n}{2}\log|\Sigma| - \tfrac{1}{2}\operatorname{tr}(\Sigma^{-1}S) - \tfrac{1}{2}n(\bar{x}-\mu)^T \Sigma^{-1}(\bar{x}-\mu)$

$\dfrac{\partial \log L(\mu; x_1 \dots, x_n)}{\partial \mu} = 0.$

$\therefore \bar{x} - \mu = 0$

$\therefore \hat{\mu} = \dfrac{1}{n}\displaystyle\sum_{i=1}^{n} x_i$

(a)

$$H(x) - H(x|Y) = -\int p(x)\,\ln(px)\,dx + \int p(x,y)\,\ln(x|Y)\,dy\,dx$$

$$= -\int \left(\int p(x,y)\,\ln(x|Y)\,dy\right) - p(x)\ln P(x)\,dx .$$

$$= -\int \left(\int p(x,y)\cdot(\ln P(y|x) + \ln(x) - \ln(y))\,dy\right) - p(x)\ln P(x)\,dx$$

$$\le \int \left(\int p(x,y)\,\ln P(y|x)\,dy + \int p(x,y)\,dy\cdot\ln P(x) - \int p(x,y)\,\ln P(y)\,dy\right.$$
$$\left. - p(x)\ln P(x)\right)dx$$

$$= -H(Y|x) + \iint p(x,y)\,\ln\frac{P(x)}{P(y)} - p(x)\ln P(x)\,dy\,dx .$$

$$= -H(Y|x) + \iint P(x,y)\,\ln P(y)\,dy\,dx$$

$$= -H(Y|x) + \int P(y)\ln P(y)\,dy .$$

$$= -H(Y|x) + H(Y)$$

$$= H(Y) - H(Y|x)$$

$$= I(x,Y)$$

(b) $I(x,Y) = H(x) - H(x|Y)$

$$= -\int P(x) \ln P(x)\, dx + \int P(x,Y) \ln(x|Y)\, dx\, dy\, Y.$$

$$= -\int P(f(Y)) \ln P(f(Y))\, df(Y) + \int P(f(Y),Y) \underbrace{\ln P(f(Y)|Y)}_{=1} df(Y)\, d$$

$$= -\int P(Y) \ln P(Y)\, dY$$

$$= H(Y)$$

$$I_{(x,Y)} = -\int P(f(Y)) \ln P(f(Y))\, df(Y)$$

$$= H(f(Y))$$

$$= H(x)$$

$$I(x,Y) = H(x) = H(Y)$$

(C) $\hat{P}(x) \triangleq \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}[x = x_i]$ — ①

$$\min_{\theta} D_{KL}(\hat{p} \| q) \triangleq \min_{\theta} - \int \hat{P}(x) \ln \frac{q(x|\theta)}{\hat{P}(x)} dx$$

$$= \min_{\theta} - \int \hat{P}(x) \ln q(x|\theta) dx + \int \hat{P}(x) \ln \hat{P}(x) dx.$$

$$\propto \min_{\theta} - \int \hat{P}(x) \ln q(x|\theta) dx.$$

plug in ① $\Rightarrow = \min_{\theta} - \int \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}[x = x_i] \ln q(x|\theta) dx$

$$= \min_{\theta} - \frac{1}{N} \sum_{i=1}^{N} \int \delta(x - x_i) \ln q(x|\theta) dx$$

$$= \min_{\theta} - \frac{1}{N} \sum_{i=1}^{N} \ln q(x|\theta) \propto \max_{\theta} q(x|\theta)$$

$\max_{\theta} q(x|\theta)$ is the miximax likelihood estimation given D

(d)    objective :     $\max \int_{-\infty}^{\infty} P(x) \ln P(x) \, dx$

constraints :
$$\begin{cases} \int_{-\infty}^{\infty} P(x) \, dx = 1 \\ \int_{-\infty}^{\infty} x P(x) \, dx = \mu \\ \int_{-\infty}^{\infty} (x-\mu)^2 P(x) \, dx = \delta^2 \end{cases}$$

$$-\int_{-\infty}^{\infty} P(x) \ln P(x) \, dx + \lambda_1 \left( \int_{-\infty}^{\infty} P(x) \, dx - 1 \right) + \lambda_2 \left( \int_{-\infty}^{\infty} x P(x) \, dx - \mu \right)$$
$$+ \lambda_3 \left( \int_{-\infty}^{\infty} (x-\mu)^2 P(x) \, dx - \delta^2 \right) = F(P(x))$$

$$\frac{\partial F(P(x))}{\partial x} = 0$$

$$\Rightarrow \hat{P}(x) = \frac{1}{\sqrt{2\pi}\delta} \exp\left( -\frac{(x-\mu)^2}{2\delta^2} \right)$$

4. (a) $w = \arg\min \sum c_i (y_i - \beta^T x_i - b)^2$

$= \arg\min \| c^{\frac{1}{2}} (y - xw) \|_2^2$

where $C$ is diagonal.

$\frac{\partial L(w)}{\partial w} = -2 (c^{\frac{1}{2}} x)^T c^{\frac{1}{2}} (y - xw) = 0$

$\hat{w} = ((cx)^T x)^{-1} (cx)^T y = (x^T c x)^{-1} x^T c y.$

$C_i = 1, \quad \hat{w} = (x^T x)^{-1} x^T y, \quad \hat{w} = \begin{bmatrix} b \\ \beta \end{bmatrix}$

$y_i = \beta^T x_i + b + \varepsilon_i$

$y = xw + \varepsilon$ . where $\varepsilon | x \sim N(0, \sigma^2 I)$

$\max P(y | w, x) = \max - (y - xw)^T (\sigma^2 I)^{-1} (y - xw)$

$= \min_w \| y - xw \|_2^2$

$\hat{w}_{MV} = (x^T x)^{-1} x^T y. = \hat{w}_{LS}$

(b) $y | x, w \sim N(Xw, \Sigma)$

$$\max P(y | x, w) = \min (y - Xw)^T \Sigma^{-1} (y - Xw)$$

$$= \min \| \Sigma^{-\frac{1}{2}} (y - Xw) \|_2^2$$

$$\hat{w}_{LS} = (X^T \Sigma X)^{-1} X^T \Sigma^{-1} y.$$

The MLE of $w$ with different noise variance

for each $i$ is equivalent to weight LS.

With matric $C = \Sigma^{-1}$,

1a) $\min\limits_{w,b} \quad \frac{1}{2}\|w\|_2^2 + C\sum\limits_{i=1}^{N}\varepsilon_i$

Subject to $\quad t^{(i)}(w^T x^{(i)}+b) \geq 1-\varepsilon_i$, where $\varepsilon_i \geq 0$

$$S_i \geq \max\left[0 - t^{(i)}(\underline{w}^T x^{(i)}+b)\right]$$

So $\min\left(\frac{1}{2}\|w\|^2 + C\sum\limits_{i=1}^{N} S_i\right)$ is equivalent to the

~~$S_i \geq \max$~~ $\min \frac{1}{2}\|w\|_2^2 + C\sum \max\left(0, 1-t^{(i)}(\underline{w}^T x^{(i)}+b)\right)$

---

1b) $\rho_{(w,b)} = \min d(x_i, H) = \min \dfrac{|w^T x_i + b|}{\|w\|_2}$

scale $w$ and $b$ by $\dfrac{1}{\min|w^T x_i + b|}$,

$\min \rho(\hat{w}, \hat{b}) = \dfrac{1}{\|\hat{w}\|_2}$

$y_i(w^T x_i + b) = 1 \implies \rho_{\min} = \dfrac{n}{\|w\|_2}$, where $n = \min|w^T x_i + b|$

$\rho_i = \dfrac{|w^T x_i + b|}{\|w\|_2} \qquad \forall i \qquad t^{(i)}(w^T x_i^{(i)}+b) \geq 1-\varepsilon_i$

$$\therefore \frac{\rho_i}{\rho_{\min}} = \varepsilon_i^* \quad \text{and} \quad \rho_i \propto \varepsilon_i^*$$

(C)

$$\min\left(\frac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\max\left(0, 1-t^{(i)}(w^Tx^{(i)}+b)\right)\right) - - - ①$$

when $C \to \infty$, ① is equivalent to $\min\left(\max\left(0, 1-t^{(i)}(w^Tx^{(i)}+b)\right)\right)$

$$t^{(i)}(w^Tx^{(i)}+b) = 1 \iff x^{(i)} \text{ is the closest point to the margin}.$$

The sum hard margin

$$\min\frac{1}{2}\|w\|_2^2 \quad \text{s.t.} \quad \text{the closest point to the margin}$$

these two are equivalent.