

1.) a.) $\sigma_X = \frac{1}{1+e^{-X}}$

$$\sigma'_X = -(1+e^{-X})^{-2}(-e^{-X})$$

$$= \frac{e^{-X}}{(1+e^{-X})^2}$$

$$= \frac{1}{1+e^{-X}} \left(\frac{e^{-X}}{1+e^{-X}} \right) = \frac{1}{1+e^{-X}} \left(\frac{1+e^{-X}}{1+e^{-X}} - \frac{1}{1+e^{-X}} \right)$$

$$= \frac{1}{1+e^{-X}} \left(1 - \frac{1}{1+e^{-X}} \right) = \boxed{\sigma(X)[1-\sigma(X)]}$$

b.) $NLL(\theta) = -\sum_{i=1}^N [y_i \log \mu_i + (1-y_i) \log(1-\mu_i)]$ 1.9 $\mu(X) = \sigma(\theta^T \vec{x}_i)$

$$= -\sum_{i=1}^N [y_i \log(\sigma(\theta^T \vec{x}_i)) + (1-y_i) \log(1-\sigma(\theta^T \vec{x}_i))]$$

$$\nabla_\theta NLL(\theta) = -\sum_{i=1}^N y_i \frac{1}{\sigma(\theta^T \vec{x}_i)} (\nabla_\theta \sigma)(\theta^T \vec{x}_i) + (1-y_i) \frac{1}{1-\sigma(\theta^T \vec{x}_i)} (\nabla_\theta \sigma)(\theta^T \vec{x}_i)$$

$$= -\sum_{i=1}^N y_i \underbrace{\frac{1}{\sigma(\theta^T \vec{x}_i)} (\sigma(1-\sigma)) \theta^T \vec{x}_i}_{(1-\sigma(\theta^T \vec{x}_i))} \vec{x}_i + (1-y_i) \underbrace{\frac{1}{1-\sigma(\theta^T \vec{x}_i)} (\sigma(1-\sigma)) (\theta^T \vec{x}_i)}_{\sigma(\theta^T \vec{x}_i)} (-\vec{x}_i)$$

$$= -\sum_{i=1}^N y_i (1-\sigma(\theta^T \vec{x}_i)) \vec{x}_i - (1-y_i) \sigma(\theta^T \vec{x}_i) \vec{x}_i$$

$$= -\sum_{i=1}^N y_i \vec{x}_i - y_i \sigma(\theta^T \vec{x}_i) \vec{x}_i - \sigma(\theta^T \vec{x}_i) \vec{x}_i + y_i \sigma(\theta^T \vec{x}_i) \vec{x}_i$$

$$= -\sum_{i=1}^N (y_i - \underbrace{\sigma(\theta^T \vec{x}_i)}_{\mu_i}) \vec{x}_i = \sum_{i=1}^N (\mu_i - y_i) \vec{x}_i$$

$$= \boxed{X^T(\vec{\mu} - \vec{y})}$$
 This can't be set to zero to find critical pt (MLE)

c.) Hessian

$$H_\theta = \nabla_\theta (\nabla_\theta NLL(\theta))^T \quad \begin{matrix} \xrightarrow{\text{Hessian = full matrix of 2nd}} \\ \text{partial derivatives s.t. } (H_f)_{i,j} = \frac{\partial^2 f}{\partial x_i \partial x_j} \end{matrix}$$

$$= \sum_{i=1}^N \nabla_\theta (x_i^T (\vec{\mu}_i - \vec{y}_i))^T = \sum_{i=1}^N \nabla_\theta ((\vec{\mu}_i - \vec{y}_i)^T \vec{x}_i) = \sum_{i=1}^N (\nabla_\theta (\vec{\mu}_i - \vec{y}_i)^T \vec{x}_i)$$

$$= \sum_{i=1}^N (\nabla_\theta \sigma(\theta^T \vec{x}_i)^T)^T \vec{x}_i = \sum_{i=1}^N (\nabla_\theta \sigma(\vec{x}_i^T \theta))^T \vec{x}_i$$

$$= \sum_{i=1}^N \vec{x}_i^T \sigma(\theta^T \vec{x}_i) (1-\sigma(\theta^T \vec{x}_i)) \vec{x}_i = X^T \text{diag}(\mu(1-\mu)) X$$

$$= \boxed{X^T S X} \quad S \equiv \text{diag}(\mu(1-\mu))$$

quadratic form of summation
uses diagonal matrix by definition

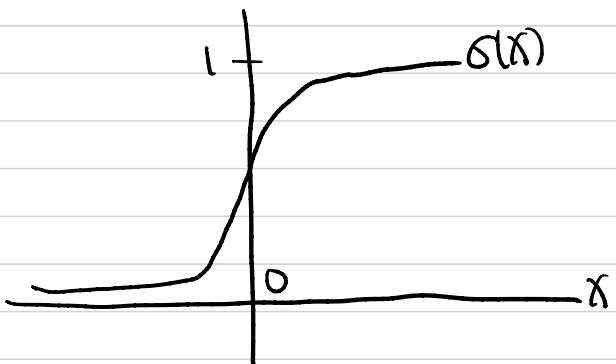
H_θ is positive semidefinite if $x^T H_\theta x \geq 0 \quad \forall x \in \mathbb{R}^n$

$$x^T H_\theta x = x^T X^T S X x = x^T X^T S(XX) = x^T X^T S^{1/2} S^{1/2}(XX) = \|S^{1/2} XX\|_2^2$$

Since $\|XX\|_2^2 \geq 0$, we now check if $\|S^{1/2}\|_2^2 \geq 0$. This is true if all values of $\mu_i(1-\mu_i) \geq 0$ because S is a diagonal matrix.

(positive semidefinite if matrix is Hermitian + all eigenvalues are real & non-negative)

$$\mu_i(1-\mu_i) = \sigma(\theta^T \vec{x}_i)(1-\sigma(\theta^T \vec{x}_i))$$



$$\sigma(x) \quad \forall x, 0 < \sigma(x) < 1.$$

Thus, μ_i must be between 0 & 1.
So, H_θ must be positive semidefinite.

$$2) P(x; \sigma^2) = \frac{1}{2} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

$$\int_{-\infty}^{\infty} P(x; \sigma^2) dx = 1$$

$$= \int_{-\infty}^{\infty} \frac{1}{2} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx$$

$$Z = \int_{-\infty}^{\infty} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx$$

$$Z^2 = \int_{y=-\infty}^{\infty} \int_{x=-\infty}^{\infty} \exp\left(-\frac{x^2+y^2}{2\sigma^2}\right) dx dy$$

Switch to polar coordinates $x = r \cos \theta, y = r \sin \theta$

$$Z^2 = \int_{\theta=0}^{2\pi} \int_{r=0}^{\infty} \exp\left(-\frac{r^2}{2\sigma^2}\right) r dr d\theta \quad u = e^{-r^2/2\sigma^2}$$

$$du = e^{-r^2/2\sigma^2} \left(-\frac{1}{2\sigma^2}(2r)\right) dr$$

$$= 2\pi \int_{r=0}^{\infty} \exp\left(-\frac{r^2}{2\sigma^2}\right) r dr$$

$$= 2\pi \left(-\sigma^2 \exp\left(-\frac{r^2}{2\sigma^2}\right)\right) \Big|_{r=0}^{\infty}$$

$$= -2\pi\sigma^2(0 - 1) = 2\pi\sigma^2$$

$$Z = \sqrt{2\pi\sigma^2} = \boxed{\sqrt{2\pi}\sigma}$$

normal distribution

$$N(x|\mu,\sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

3.) a.) $P(W) = \prod_j N(w_j | 0, \tau^2)$

$$\operatorname{argmax}_W \underbrace{\sum_{i=1}^N \log N(y_i | w_0 + W^T x_i, \sigma^2)}_{X_1, \dots, X_N} + \underbrace{\sum_{j=1}^D \log N(w_j | 0, \tau^2)}_{\text{Data}}$$

$$\theta_{MAP} = \operatorname{argmax}_{\theta} p(\theta|D) = \operatorname{argmax}_{\theta} \underbrace{p(D|\theta)}_{p(D)} p(\theta)$$

ignore w₀ in argmax

$$= \operatorname{argmax} (\log p(D|\theta) + \log p(\theta))$$

$$= \operatorname{argmax}_W \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) \exp \left(-\frac{(y_i - w_0 + W^T x_i)^2}{2\sigma^2} \right) + \sum_{j=1}^D \log \frac{1}{\sqrt{2\pi}\tau} \exp \left(-\frac{w_j^2}{2\tau^2} \right)$$

$$= \operatorname{argmax}_W \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) + \left(-\frac{(y_i - w_0 + W^T x_i)^2}{2\sigma^2} \right) + \sum_{j=1}^D \log \left(\frac{1}{\sqrt{2\pi}\tau} \right) + \frac{-w_j^2}{2\tau^2}$$

remove constants from argmax

$$= \operatorname{argmax}_W \sum_{i=1}^N \frac{1}{2\sigma^2} (-(y_i - w_0 + W^T x_i)^2) + \sum_{j=1}^D \frac{1}{2\tau^2} (-w_j)^2$$

scale by $2\sigma^2$

$$= \operatorname{argmax}_W \sum_{i=1}^N -(y_i - w_0 + W^T x_i)^2 + \sum_{j=1}^D \frac{\sigma^2}{\tau^2} (-w_j)^2$$

this expression is equivalent to argmin of its negative:

$$= \operatorname{argmin}_W \sum_{i=1}^N (y_i - w_0 + W^T x_i)^2 + \frac{\sigma^2}{\tau^2} \sum_{j=1}^D w_j^2 \quad \text{Define } \frac{\sigma^2}{\tau^2} \equiv \lambda$$

$$= \operatorname{argmin}_W \sum_{i=1}^N (y_i - w_0 + W^T x_i)^2 + \lambda \|W\|_2^2$$

We can scale by a factor of $\frac{1}{N}$ (will not change argmin)

$$= \boxed{\operatorname{argmin}_W \frac{1}{N} \sum_{i=1}^N (y_i - w_0 + W^T x_i)^2 + \lambda \|W\|_2^2}$$

turns it
into average
squared
error

$$b) \text{ minimize: } \|Ax - b\|_2^2 + \|\Gamma x\|_2^2$$

$$\nabla_x \|Ax - b\|_2^2 + \|\Gamma x\|_2^2 = 0$$

$$= \nabla_x [(Ax - b)^T (Ax - b) + (\Gamma x)^T (\Gamma x)]$$

$$= \nabla_x [(x^T A^T - b^T)(Ax - b) + (x^T \Gamma^T)(\Gamma x)]$$

$$= \nabla_x [x^T A^T A x - x^T A^T b - \underbrace{b^T A x}_{(x^T A^T b)^T} + b^T b + x^T \Gamma^T \Gamma x]$$

$$= 2A^T A x - 2A^T b + 2\Gamma^T \Gamma x = 0$$

$$(2A^T A + 2\Gamma^T \Gamma)x = 2A^T b$$

$$\boxed{x = (A^T A + \Gamma^T \Gamma)^{-1} A^T b}$$

$$d) \text{ minimize: } \|Ax + b1 - y\|_2^2 + \|\Gamma x\|_2^2$$

$1 = \text{ones matrix}$

$$\nabla_x [(Ax + b1 - y)^T (Ax + b1 - y) + (\Gamma x)^T (\Gamma x)]$$

$$= \nabla_x [(x^T A^T + b1^T - y^T)(Ax + b1 - y) + (x^T \Gamma^T \Gamma x)]$$

$$= \nabla_x [x^T A^T A x + b1^T A x - y^T A x + x^T A^T b1 + \underbrace{b1^T b1}_{b1^T b1} - y^T b1 - x^T A^T y - b1^T y + y^T y + x^T \Gamma^T \Gamma x]$$

$$= \nabla_x [x^T A^T A x + 2b1^T A x - 2y^T A x - 2b1^T y + b1^T b1 + y^T y + x^T \Gamma^T \Gamma x]$$

$$= 2A^T A x + 2b1^T A - 2y^T A + 2\Gamma^T \Gamma x = 0$$

Also find ∇_b and set equal to 0 to solve for optimal b .

$$\nabla_b [f] = 21^T A x - 21^T y + 2b n = 0$$

$$b = \frac{1^T A x - 1^T y}{-n} = \boxed{\frac{-1^T (Ax - y)}{n}}$$

Solve for x :

$$(2A^T A + 2\Gamma^T \Gamma)x + 2\left(\frac{-1^T (Ax - y)}{n}\right) 1^T A - 2y^T A = 0$$

$$(A^T A + \Gamma^T \Gamma - \frac{1^T A (1^T A)}{n})x = -\frac{1^T y (1^T A)}{n} + y^T A$$

$\frac{1^T A^T 1}{n}$

$$X = [A^T A + \Gamma^T \Gamma + \frac{1}{n} \mathbf{1} \mathbf{1}^T A^T A]^{-1} [\frac{1}{n} \mathbf{1} - \mathbf{1}^T A Y + Y^T A]$$