

Ex.1

See code and output at the end of  
this document .

Ex. 2

$$g_{\theta} = \theta^T x$$

$$[\theta_0, \theta_1, \theta_2] \begin{bmatrix} x_{bmi} \\ x_{ht} \\ 1 \end{bmatrix}$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{n=1}^N (y_n - g_{\theta}(x_n))^2$$

$$\Rightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmin}} \|y - X\theta\|^2$$

$$a) J(\theta) = \|y - X\theta\|^2 = \sqrt{(y - X\theta)^2}^2 = (y - X\theta)^2$$

$$\nabla_{\theta} J(\theta) = 0$$

$$2 X^T (y - X\hat{\theta}) = 0.$$

$$X^T y - X^T X \hat{\theta} = 0$$

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

To have a unique minimizer,  $\nabla_{\theta}^2 J(\theta)$  needs to be positive definite.

If  $X^T X$  is singular, we can use regularization & LASSO regression.

c) See code

d). Find  $\nabla \Sigma_{\text{train}}(\theta^k)$ .

$$\theta^{k+1} = \theta^k - \alpha^k \nabla \Sigma_{\text{train}}(\theta^k)$$

$$\nabla \Sigma_{\text{train}}(\theta^k) = \nabla_{\theta} \|y_n - x \theta^{(k)}\|^2$$

$$= 2 x^T (y_n - x \theta^{(k)})$$

$$= 2 x^T y - 2 x^T x \theta^k$$

$$= 2 (x^T y - x^T x \theta^k) = \underline{d}$$

$$\theta^{k+1} = \theta^k - \alpha^k (b - A \theta^k)$$

$$\theta^{k+1} = \theta^k - \alpha^k b + \alpha^k A \theta^k.$$

$$\min_{\alpha} J(\underline{\theta} + \alpha \underline{d}) = \|\underline{y} - x(\underline{\theta} + \alpha \underline{d})\|^2.$$

$$\begin{aligned} \nabla_{\alpha} J(\underline{\theta} + \alpha \underline{d}) &= \frac{d}{d\alpha} \left[ \|\underline{y}\|^2 + \|x(\underline{\theta} + \alpha \underline{d})\|^2 - 2 \underline{y}^T x(\underline{\theta} + \alpha \underline{d}) \right] \\ &= \frac{d}{d\alpha} \left[ \|\underline{y}\|^2 + (\underline{\theta} + \alpha \underline{d})^T x^T x (\underline{\theta} + \alpha \underline{d}) - 2 \underline{y}^T x (\underline{\theta} + \alpha \underline{d}) \right] \end{aligned}$$

$$= \frac{d}{d\alpha} \left[ (\underline{\theta} + \alpha \underline{d})^T x^T x (\underline{\theta} + \alpha \underline{d}) - 2 \underline{y}^T x \alpha \underline{d} \right].$$

$$= \frac{d}{d\alpha} \left[ \underline{\theta}^T x^T x (\underline{\theta} + \alpha \underline{d}) + \alpha \underline{d}^T x^T x (\underline{\theta} + \alpha \underline{d}) \right] - 2 \underline{y}^T x \underline{d}$$

$$= \underline{\theta}^T x^T x \underline{d} + \underline{d}^T x^T x \underline{\theta} + 2 \alpha \underline{d}^T x^T x \underline{d} - 2 \underline{y}^T x \underline{d}$$

$$= 2 \underline{\theta}^T x^T x \underline{d} + 2 \alpha \underline{d}^T x^T x \underline{d} - 2 \underline{y}^T x \underline{d}$$

$$= 2 (\underline{\theta}^T x^T x \underline{d} + \alpha \underline{d}^T x^T x \underline{d} - \underline{y}^T x \underline{d})$$

$$\nabla_{\alpha} J(\underline{\theta} + \alpha \underline{d}) = 0.$$

$$\underline{\theta}^T \underline{x}^T \underline{x} \underline{d} + \alpha \|\underline{x} \underline{d}\|^2 - \underline{y}^T \underline{x} \underline{d} = 0,$$

$$\alpha^k = - \frac{\underline{\theta}^T \underline{x}^T \underline{x} \underline{d}^k - \underline{y}^T \underline{x} \underline{d}^k}{\|\underline{x} \underline{d}^k\|^2}$$

$$= - \frac{(\underline{\theta}^T \underline{x}^T \underline{x} - \underline{y}^T \underline{x}) \underline{d}^k}{\|\underline{x} \underline{d}^k\|^2}$$

$$= - \frac{\frac{1}{2} \underline{d}^{kT} \cdot \underline{d}^k}{\underline{d}^{kT} \underline{x}^T \underline{x} \underline{d}^k}$$

$$= - \frac{1}{2} \frac{\underline{d}^{kT} \cdot \underline{d}^k}{\underline{d}^{kT} \cdot A \cdot \underline{d}^k}$$

$$4) \quad \hat{\theta}_\lambda = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2 \quad (3)$$

$$\hat{\theta}_\alpha = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|X\theta - y\|_2^2 \quad \text{s.t.} \quad \|\theta\|_2^2 \leq \alpha \quad (4)$$

$$\hat{\theta}_\epsilon = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|\theta\|_2^2 \quad \text{s.t.} \quad \|X\theta - y\|_2^2 \leq \epsilon \quad (5)$$

a) Looking at (3),

$$J(\hat{\theta}) = \|X\hat{\theta} - y\|^2 + \lambda \|\hat{\theta}\|^2 = (X\hat{\theta} - y)^T (X\hat{\theta} - y) + \lambda \hat{\theta}^T \hat{\theta}$$

$$\nabla_{\hat{\theta}} J(\hat{\theta}) = 0$$

$$2X^T(X\hat{\theta} - y) + 2\lambda \hat{\theta} = 0$$

$$2X^T X \hat{\theta} - 2X^T y + 2\lambda \hat{\theta} = 0$$

$$(X^T X + \lambda I) \hat{\theta} = X^T y$$

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

$$c) b). \hat{\theta}_\alpha = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|X\theta - y\|_2^2 \quad \text{s.t.} \quad \|\theta\|_2^2 \leq \alpha \quad (4)$$

$$\hat{\theta}_\epsilon = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|\theta\|_2^2 \quad \text{s.t.} \quad \|X\theta - y\|_2^2 \leq \epsilon \quad (5)$$

(i) For (3)

$$J(\theta) = \|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2$$

For (4) :

$$J(\theta, \gamma_\alpha) = \|X\theta - y\|_2^2 - \gamma_\alpha (\alpha - \|\theta\|_2^2)$$

For (5)

$$J(\theta, \gamma_\epsilon) = \|\theta\|_2^2 - \gamma_\epsilon (\epsilon - \|X\theta - y\|_2^2)$$

4 b ii)

(3) : stationarity :  $\nabla_{\theta} L = (X^T X + \lambda I) \hat{\theta} - X^T y = 0$

Primal feasibility : None because it has no constraint.

Dual feasibility : None . . . .

Complementary slackness : None . . . .

---

(4) . stationarity :  $\nabla_{\theta} L = 0$

$$X^T (X\theta - y) - \gamma_{\alpha} \theta = 0$$

Primal feasibility :  $(\alpha - \|\theta\|_2^2) \geq 0$

Dual feasibility :  $\gamma_{\alpha} \geq 0$

complementary slackness :  $\gamma_{\alpha} \cdot (\alpha - \|\theta\|_2^2) = 0$

---

(5) stationarity :  $\nabla_{\theta} L = 0$

$$\theta - \gamma_{\varepsilon} X^T (X\theta - y) = 0$$

Primal feasibility :  $(\varepsilon - \|X\theta - y\|_2^2) \geq 0$

Dual feasibility :  $\gamma_{\varepsilon} \geq 0$

complementary slackness :  $\gamma_{\varepsilon} (\varepsilon - \|X\theta - y\|_2^2) = 0$

$$4b) \text{ iii) } \hat{\theta}_\alpha = (X^T X + \alpha I)^{-1} X^T y.$$

$$X^T (X \hat{\theta}_\alpha - y) - \gamma_\alpha \hat{\theta}_\alpha = 0 \quad (\text{Stationarity})$$

$$X^T X \hat{\theta}_\alpha - X^T y - \gamma_\alpha \hat{\theta}_\alpha = 0$$

$$X^T X \hat{\theta}_\alpha - (X^T X + \lambda I) \hat{\theta}_\alpha - \gamma_\alpha \hat{\theta}_\alpha = 0.$$

$$(\cancel{X^T X} - \cancel{X^T X} - \lambda I - \gamma_\alpha I) \hat{\theta} = 0.$$

$$\boxed{\gamma_\alpha = \lambda} > 0 \quad (\text{Primal feasibility}),$$

$$\gamma_\alpha (\alpha - \|\hat{\theta}_\alpha\|^2) = 0 \quad (\text{complementary slackness})$$

$$\boxed{\alpha = \|(X^T X + \lambda I)^{-1} X^T y\|_2^2}$$



4 b) iv).

$$\hat{\theta}_\lambda - \gamma_\varepsilon x^T (x \hat{\theta}_\lambda - y) = 0.$$

$$\hat{\theta}_\lambda - \gamma_\varepsilon x^T x \hat{\theta}_\lambda - \gamma_\varepsilon x^T y = 0,$$

$$\hat{\theta}_\lambda - \gamma_\varepsilon x^T x \hat{\theta}_\lambda - \gamma_\varepsilon (x^T x + \lambda I) \hat{\theta}_\lambda = 0$$

$$\left[ I - \gamma_\varepsilon (x^T x - \cancel{x^T x} + \lambda I) \right] \hat{\theta}_\lambda = 0$$

$$[(1 - \gamma_\varepsilon \lambda) I] \hat{\theta}_\lambda = 0$$

$$\boxed{\gamma_\varepsilon = \frac{1}{\lambda}} > 0 \quad (\text{stationarity})$$

complementary slackness:  $\gamma_\varepsilon (\varepsilon - \|x \hat{\theta}_\lambda - y\|_2^2) = 0$

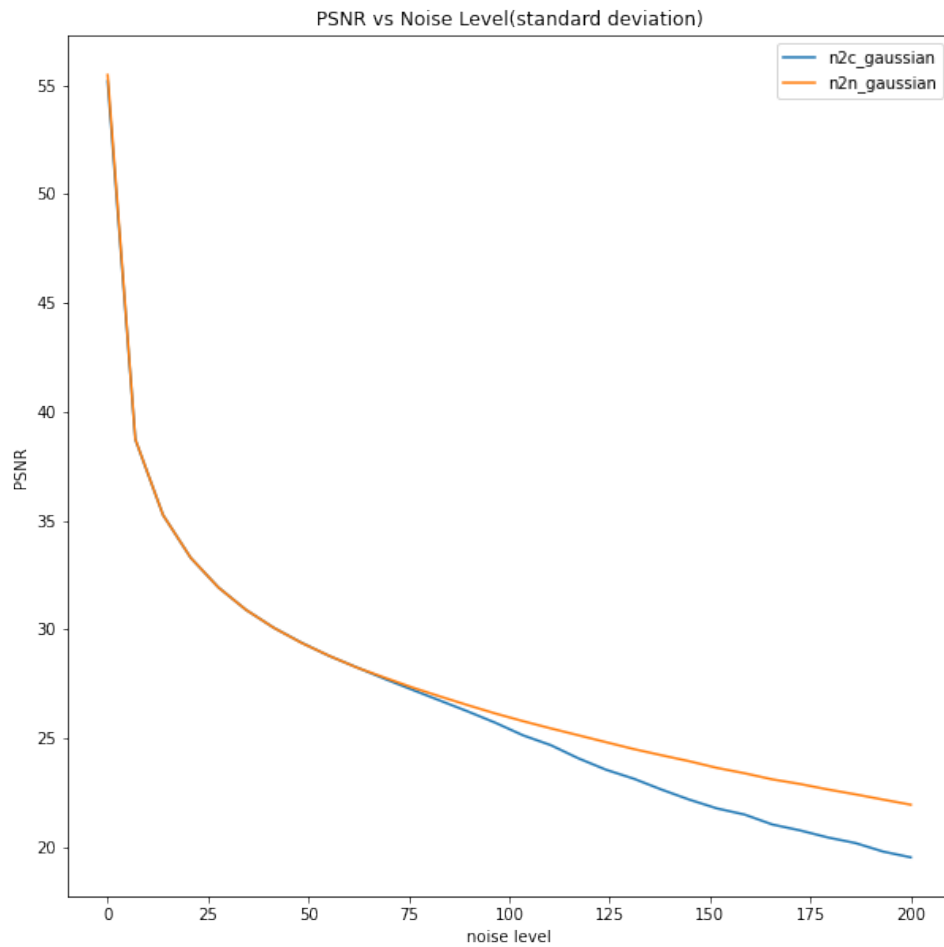
$$\varepsilon - \|x \hat{\theta}_\lambda - y\|_2^2 = 0$$

$$\varepsilon = \|x \hat{\theta}_\lambda - y\|_2^2$$

$$\boxed{\varepsilon = \|x (x^T x + \lambda I)^{-1} x^T y - y\|_2^2}$$

4 b) v). KKT conditions are necessary but not conclusive to guarantee a solution. However, since the problem we are optimizing is smooth, these conditions are sufficient for optimality, and we can claim that  $\hat{\theta}_n$  is the solution to (4).

5.



Noise levels are generated from 0 to 200 at normal distribution, then normalized by dividing by 255.

When the noise level is low, it's clear that both models perform similarly. However, as noise level gets higher, the model trained with noise data (n2n) has a higher average PSNR score than the one trained with clean data.