

Problem c. Stochastic Gradient Descent

Haoyang Chen
hc2812

for $i = 0, \dots, N$ do:

 initialize d_i randomly and $t=1$

 While ($t \leq T$ & Stopping condition is not True) do:

 np.random.shuffle(X)

 for $k = 0, \dots, X.shape[0]-1$ do:

$$D_k = X[k]^T X[k] - \sum_{j=0}^{i-1} \lambda_j d_j d_j^T$$

$$y = d_i - \nabla_{d_i} (-d_i^T D_k d_i)$$

$$d_i = \frac{y}{\|y\|}$$

$$t = t + 1$$

$$\lambda_i = d_i^T X^T X d_i$$

Neural Networks and Deep Learning
Haoyang Chen hc2812 HW1

Theory Part:

Problem d:

(i) Y is uniformly distributed when $x \in [0, 255]$

$$\therefore y = \text{cdf}(x) = \int_0^x p(u) du = \int_0^x \frac{f(u)}{\int_0^{255} f(t) dt} du \approx \int_0^x f(u) du = \int_0^x \frac{1}{256} e^{-\frac{(x-u)^2}{256}} du = \Phi\left(\frac{x-u}{8}\right) = \Phi\left(x - 127.5\right)$$

Thus the mapping function $g(x)$ is the cumulative distribution function $\Phi(x - 127.5)$

$$(ii) P(X=x) = \int_0^\infty \int_0^\infty p(x=y, Y=z, Z=z) dy dz = \int_0^1 \int_0^1 8xyz dy dz = 2x \cdot [y^2]_0^1 \cdot [z^2]_0^1 = 2x, \quad x \in [0, 1]$$

$$P(Y=y) = \int_0^\infty \int_0^\infty p(x=y, Y=y, Z=z) dx dz = \int_0^1 \int_0^1 8xyz dx dz = 2y \cdot [x^2]_0^1 \cdot [z^2]_0^1 = 2y, \quad y \in [0, 1]$$

$$P(Z=z) = \int_0^\infty \int_0^\infty p(x=y, Y=y, Z=z) dx dy = \int_0^1 \int_0^1 8xyz dx dy = 2z \cdot [x^2]_0^1 \cdot [y^2]_0^1 = 2z, \quad z \in [0, 1]$$

$$\therefore P(x) = \begin{cases} 2x, & x \in [0, 1] \\ 0, & \text{otherwise} \end{cases}, \quad P(y) = \begin{cases} 2y, & y \in [0, 1] \\ 0, & \text{otherwise} \end{cases}, \quad P(z) = \begin{cases} 2z, & z \in [0, 1] \\ 0, & \text{otherwise} \end{cases}$$

$$\begin{aligned} E(XYZ) &= \int_0^\infty \int_0^\infty \int_0^\infty xyz p(x, y, z) dx dy dz = \int_0^1 \int_0^1 \int_0^1 xyz \cdot 8xyz dx dy dz = 8 \left[\frac{x^3}{3}\right]_0^1 \cdot \left[\frac{y^3}{3}\right]_0^1 \cdot \left[\frac{z^3}{3}\right]_0^1 \\ &= \frac{8}{27} \end{aligned}$$

x and y are conditionally independent given z .

$$\text{Proof: } P(x, y|z) = \frac{P(x, y, z)}{P(z)} = \frac{8xyz}{2z} = 4xy, \quad \text{for } x, y, z \in [0, 1]$$

$$P(x|z) = \frac{P(x, z)}{P(z)} = \frac{\int_0^\infty p(x, y, z) dy}{P(z)} = \frac{\int_0^1 8xyz dy}{2z} = \frac{4xz}{2z} = 2x, \quad \text{for } x, z \in [0, 1]$$

$$P(y|z) = \frac{P(y, z)}{P(z)} = \frac{\int_0^\infty p(x, y, z) dx}{P(z)} = \frac{\int_0^1 8xyz dx}{2z} = \frac{4yz}{2z} = 2y, \quad \text{for } y, z \in [0, 1]$$

$$\therefore P(x, y|z) = \begin{cases} 4xy, & \text{for } x, y, z \in [0, 1] \\ 0, & \text{otherwise} \end{cases} \quad P(x|z) = \begin{cases} 2x, & \text{for } x, z \in [0, 1] \\ 0, & \text{otherwise} \end{cases}, \quad P(y|z) = \begin{cases} 2y, & \text{for } y, z \in [0, 1] \\ 0, & \text{otherwise} \end{cases}$$

$$\therefore P(x, y|z) = P(x|z) P(y|z) \quad \therefore x \text{ and } y \text{ are conditionally independent given } z \quad \#$$

Problem e.

$$\begin{aligned}
 1. P(\underline{\mathcal{U}} | \underline{X}, \Sigma) &\propto P(\underline{X} | \underline{\mathcal{U}}, \Sigma) P(\underline{\mathcal{U}} | \underline{\mathcal{U}}_0, \Sigma_0) \propto \prod_{i=1}^m \exp\left(-\frac{1}{2}(\underline{X}^{(i)} - \underline{\mathcal{U}})^T \Sigma^{-1} (\underline{X}^{(i)} - \underline{\mathcal{U}})\right) \cdot \exp\left(-\frac{1}{2}(\underline{\mathcal{U}} - \underline{\mathcal{U}}_0)^T \Sigma_0^{-1} (\underline{\mathcal{U}} - \underline{\mathcal{U}}_0)\right) \\
 &= \exp\left(-\frac{1}{2}(\underline{\mathcal{U}} - \underline{\mathcal{U}}_0)^T \Sigma_0^{-1} (\underline{\mathcal{U}} - \underline{\mathcal{U}}_0) + \sum_{i=1}^m (\underline{X}^{(i)} - \underline{\mathcal{U}})^T \Sigma^{-1} (\underline{X}^{(i)} - \underline{\mathcal{U}})\right) \\
 &\propto \exp\left(-\frac{1}{2}(\underline{\mathcal{U}} - \underline{\mathcal{U}}_0)^T \Sigma_0^{-1} (\underline{\mathcal{U}} - \underline{\mathcal{U}}_0)\right) = N(\underline{\mathcal{U}} | \underline{\mathcal{U}}_m, \Sigma_m)
 \end{aligned}$$

$$\text{where } \underline{\mathcal{U}}_m = (\Sigma_0^{-1} + m\Sigma^{-1})^{-1} (\Sigma_0^{-1}\underline{\mathcal{U}}_0 + m\Sigma^{-1}\bar{X}), \quad \bar{X} = \frac{1}{m} \sum_{i=1}^m \underline{X}^{(i)}$$

$$\Sigma_m^{-1} = \Sigma_0^{-1} + m\Sigma^{-1}$$

Thus $\underline{\mathcal{U}}_{MAP} = \underline{\mathcal{U}}_m$ which maximize $P(\underline{\mathcal{U}} | \underline{X}, \Sigma)$

I guess professor means $\Sigma_{MAP} = \Sigma_m$, but in bayesian statistics, it's wired. Σ_m is related to the posterior distribution of $\underline{\mathcal{U}}$, not the posterior distribution of Σ . In bayesian statistics, if Σ is unknown, $\Sigma_m^{-1} = \Sigma_0^{-1} + m\Sigma_{MAP}^{-1}$, and $\Sigma_{MAP} = \arg \max_{\Sigma} P(\Sigma | \underline{X})$.

$$2. E(\underline{\mathcal{U}}_{MAP}) = E((\Sigma_0^{-1} + m\Sigma^{-1})^{-1} (\Sigma_0^{-1}\underline{\mathcal{U}}_0 + m\Sigma^{-1}\bar{X})) \neq E(\bar{X}) = \underline{\mathcal{U}}$$

$$E(\Sigma_{MAP}) = E((\Sigma_0^{-1} + m\Sigma^{-1})^{-1}) \neq \Sigma$$

Thus it's biased estimators

$$3. L(\underline{X}, \underline{\mathcal{U}}, \Sigma) = \prod_{i=1}^m \frac{1}{(2\pi)^{\frac{d}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\underline{X}^{(i)} - \underline{\mathcal{U}})^T \Sigma^{-1} (\underline{X}^{(i)} - \underline{\mathcal{U}})\right)$$

$$\propto |\Sigma|^{-\frac{m}{2}} \exp\left(-\frac{1}{2} \sum_{i=1}^m (\underline{X}^{(i)} - \underline{\mathcal{U}})^T \Sigma^{-1} (\underline{X}^{(i)} - \underline{\mathcal{U}})\right)$$

$$l(\underline{X}, \underline{\mathcal{U}}, \Sigma) = \log L(\underline{X}, \underline{\mathcal{U}}, \Sigma) \propto -\frac{m}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^m (\underline{X}^{(i)} - \underline{\mathcal{U}})^T \Sigma^{-1} (\underline{X}^{(i)} - \underline{\mathcal{U}})$$

$$\frac{\partial l(\underline{X}, \underline{\mathcal{U}}, \Sigma)}{\partial \underline{\mathcal{U}}} = 0 \Rightarrow \underline{\mathcal{U}}_{MLE} = \bar{X} = \frac{1}{m} \sum_{i=1}^m \underline{X}^{(i)}$$

$$\frac{\partial l(\underline{X}, \underline{\mathcal{U}}, \Sigma)}{\partial \Sigma} = 0 \Rightarrow \Sigma_{MLE} = \frac{1}{m} \sum_{i=1}^m (\underline{X}^{(i)} - \underline{\mathcal{U}}_{MLE})(\underline{X}^{(i)} - \underline{\mathcal{U}}_{MLE})^T = \frac{1}{m} \sum_{i=1}^m (\underline{X}^{(i)} - \bar{X})(\underline{X}^{(i)} - \bar{X})^T$$

Comparing the MLE & MAP, Σ_{MLE} is an unbiased estimator while Σ_{MAP} is biased