

Introduction to Machine Learning, Fall 2023

Homework 2

(Due Tuesday Nov. 14 at 11:59pm (CST))

October 25, 2023

1. [10 points] [Convex Optimization Basics]

- (a) Proof any norm $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex. [2 points]
- (b) Determine the convexity (i.e., convex, concave or neither) of $f(x_1, x_2) = x_1^2/x_2$ on $\mathbb{R} \times \mathbb{R}_{>0}$. [2 points]
- (c) Determine the convexity of $f(x_1, x_2) = x_1/x_2$ on $\mathbb{R}_{>0}^2$. [2 points]
- (d) Recall Jensen's inequality $f(\mathbb{E}(X)) \leq \mathbb{E}(f(X))$ if f is convex for any random variable X . Proof the log sum inequality:

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left(\sum_{i=1}^n a_i \right) \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i}$$

where a_1, \dots, a_n and b_1, \dots, b_n are positive numbers. Hints: $f(x) = x \log x$ is strictly convex. [4 points]

Solution:

(a) for any norm we always have.

$$\begin{cases} \|x\| \geq 0, \|x\|=0 \iff x=0 & \textcircled{1} \\ \|cx\| = |c| \|x\| & \textcircled{2} \\ \|x+y\| \leq \|x\| + \|y\| & \textcircled{3} \end{cases}$$

$$\theta f(x) + (1-\theta)f(y) \geq f(\theta x + (1-\theta)y)$$

for $x, y \in \text{dom}$, we have $\theta x, (1-\theta)y$ also $\in \text{dom}$, $\theta \in [0, 1]$

with property $\textcircled{3}$ we have :

$$f(\theta x + (1-\theta)y) \leq f(\theta x) + f((1-\theta)y)$$

$\Rightarrow f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$ property $\textcircled{2}$

so according to definition of convex function, we proved.

$$(b) H = f'' = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} \frac{2}{x_2} & -2x_1 x_2^{-2} \\ -2x_1 x_2^{-2} & 2x_1^2 x_2^{-3} \end{bmatrix}$$

$$|H| = 4x_1^2 \cdot x_2^{-4} - 4x_1^2 x_2^{-4} = 0 \quad \frac{2}{x_2} > 0.$$

$\because H$ is symmetric and the determinant is 0 so is semi positive.
definite matrix So convex

$$c) \quad |H| = \begin{vmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{vmatrix} = \begin{vmatrix} 0 & -x_2^{-2} \\ -x_2^{-2} & 2x_1 x_2^{-3} \end{vmatrix} = -x_2^{-4} \leq 0.$$

so the Hessian matrix is not positive definite
which shows the function not convex

(d)

$$E(x) = \frac{b_i}{\sum_{i=1}^n b_i} \frac{a_i}{b_i} \quad \text{and} \quad f(x) = x \cdot \log x$$

we find: $f'(x) \geq 0$ $f(x)$ convex

$$\Rightarrow \text{let } \lambda_i = \frac{b_i}{\sum b_i}, \quad x_i = \frac{a_i}{b_i} \quad \therefore \sum_{i=1}^n \lambda_i = 1$$

\therefore we have $f(\sum \lambda_i x_i) \leq \sum \lambda_i f(x_i)$

$$\Rightarrow \sum \frac{b_i}{\sum b_i} \cdot \frac{a_i}{b_i} \cdot \log \sum \frac{b_i}{\sum b_i} \frac{a_i}{b_i} \leq \sum \frac{b_i}{\sum b_i} \cdot \frac{a_i}{b_i} \log \frac{a_i}{b_i}$$

$$\Rightarrow \sum a_i \cdot \log \sum \frac{a_i}{\sum b_i} \leq (\sum a_i) \log \frac{a_i}{\sum b_i}$$

$$\Rightarrow \sum a_i \cdot \log \frac{\sum a_i}{\sum b_i} \leq (\sum a_i) \log \frac{a_i}{\sum b_i}$$

PROVE

2. [10 points] [Linear Methods for Classification] Consider the “Multi-class Logistic Regression” algorithm. Given training set $\mathcal{D} = \{(x^i, y^i) \mid i = 1, \dots, n\}$ where $x^i \in \mathbb{R}^{p+1}$ is the feature vector and $y^i \in \mathbb{R}^k$ is a one-hot binary vector indicating k classes. We want to find the parameter $\hat{\beta} = [\hat{\beta}_1, \dots, \hat{\beta}_k] \in \mathbb{R}^{(p+1) \times k}$ that maximize the likelihood for the training set. Introducing the softmax function, we assume our model has the form

$$p(y_c^i = 1 \mid x^i; \beta) = \frac{\exp(\beta_c^\top x^i)}{\sum_{c'} \exp(\beta_{c'}^\top x^i)},$$

where y_c^i is the c -th element of y^i .

given x^i we should have $y_c^i = 1$
 $\begin{cases} 1 & c=t \\ 0 & c \neq t \end{cases} \rightarrow c$

- (a) Complete the derivation of the conditional log likelihood for our model, which is

$$\ell(\beta) = \ln \prod_{i=1}^n p(y_t^i \mid x^i; \beta) = \sum_{i=1}^n \sum_{c=1}^k \left[y_t^i (\beta_c^\top x^i) - y_t^i \ln \left(\sum_{c'} \exp(\beta_{c'}^\top x^i) \right) \right].$$

For simplicity, we abbreviate $p(y_t^i = 1 \mid x^i; \beta)$ as $p(y_t^i \mid x^i; \beta)$, where t is the true class for x^i . [4 points]

- (b) Derive the gradient of $\ell(\beta)$ w.r.t. β_1 , i.e.,

$$\nabla_{\beta_1} \ell(\beta) = \nabla_{\beta_1} \sum_{i=1}^n \sum_{c=1}^k \left[y_t^i (\beta_c^\top x^i) - y_t^i \ln \left(\sum_{c'} \exp(\beta_{c'}^\top x^i) \right) \right].$$

Remark: Log likelihood is always concave; thus, we can optimize our model using gradient ascent. (The gradient of $\ell(\beta)$ w.r.t. β_2, \dots, β_k is similar, you don't need to write them) [6 points]

Solution:

$$(a) L(\beta) = \prod_{i=1}^n \prod_{c=1}^k P(y_c^i \mid x^i, \beta) = \prod_{i=1}^n \prod_{c=1}^k \left[\frac{\exp(\beta_c^\top x^i)}{\sum_c \exp(\beta_c^\top x^i)} \right]^{y_c^i}$$

$$\Rightarrow \ell(\beta) = \ln(L(\beta)) = \ln \prod_{i=1}^n \prod_{c=1}^k \frac{\exp(\beta_c^\top x^i)}{\sum_c \exp(\beta_c^\top x^i)} = \sum_{i=1}^n \ln \prod_{c=1}^k \frac{\exp(\beta_c^\top x^i)}{\sum_c \exp(\beta_c^\top x^i)}$$

$$= \sum_{i=1}^n \sum_{c=1}^k (\ln \exp(\beta_c^\top x^i) \cdot y_c^i - y_c^i \ln (\sum_c \exp(\beta_c^\top x^i)))$$

$$= \sum_{i=1}^n \sum_{c=1}^k (y_c^i (\beta_c^\top x^i) - y_c^i \ln (\sum_c \exp(\beta_c^\top x^i)))$$

$$(b) \sum_{i=1}^n y_t^i (\beta_1^\top x^{i\top}) - y_t^i \cdot \frac{\exp(\beta_1^\top x^{i\top}) \cdot x^{i\top}}{\sum_c \exp(\beta_c^\top x^{i\top})}$$

3. [10 points] [Probability and Estimation] Suppose $\mathcal{D} = \{x_1, x_2, \dots, x_n\}$ are i.i.d. samples from exponential distribution with parameter $\lambda > 0$, i.e., $X \sim \text{Expo}(\lambda)$. Recall the PDF of exponential distribution is

$$p(x | \lambda) = \begin{cases} \lambda e^{-\lambda x}, & x > 0 \\ 0, & \text{otherwise} \end{cases}$$

- (a) To derive the posterior distribution of λ , we assume its prior distribution follows gamma distribution with parameters $\alpha, \beta > 0$, i.e., $\lambda \sim \text{Gamma}(\alpha, \beta)$ (since the range of gamma distribution is also $(0, +\infty)$, thus it's a plausible assumption). The PDF of λ is given by

$$p(\lambda | \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\lambda\beta},$$

where $\Gamma(\alpha) = \int_0^{+\infty} t^{\alpha-1} e^{-t} dt$, $\alpha > 0$. Show that the posterior distribution $p(\lambda | \mathcal{D})$ is also a gamma distribution and identify its parameters. Hints: Feel free to drop constants. [4 points]

- (b) Derive the maximum a posterior (MAP) estimation for λ under $\text{Gamma}(\alpha, \beta)$ prior. [3 points]

- (c) For exponential distribution $\text{Expo}(\lambda)$, $\sum_{i=1}^n x_i \sim \text{Gamma}(n, \lambda)$ and the inverse sample mean $\frac{1}{n} \sum_{i=1}^n x_i$ is the MLE for λ . Argue that whether $\frac{n-1}{n} \hat{\lambda}_{MLE}$ is unbiased ($\mathbb{E}(\frac{n-1}{n} \hat{\lambda}_{MLE}) = \lambda$). Hints: $\Gamma(z+1) = z\Gamma(z)$, $z > 0$. [3 points]

Solution:

$$(a) P(\lambda | D) = \frac{P(D | \lambda) P(\lambda)}{P(D)}$$

$$P(D) = \int_{-\infty}^{+\infty} P(D, \lambda) \cdot d\lambda$$

$$\Rightarrow P(\lambda | D) = \frac{\lambda^n e^{-\lambda(\sum_{i=1}^n x_i)} \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\lambda\beta}}{\int_{-\infty}^{+\infty} \lambda^n e^{-\lambda(\sum_{i=1}^n x_i)} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\lambda\beta} d\lambda}$$

$$= \frac{\lambda^{n+\alpha-1} e^{-\lambda(\sum_{i=1}^n x_i + \beta)}}{\int_{-\infty}^{+\infty} \lambda^{n+\alpha-1} e^{-\lambda(\sum_{i=1}^n x_i + \beta)} d\lambda}$$

$$\text{for } \int_{-\infty}^{+\infty} \lambda^{n+\alpha-1} e^{-\lambda(\sum_{i=1}^n x_i + \beta)} d\lambda$$

$$\Rightarrow \int_{-\infty}^{+\infty} \left(\frac{t}{\sum_{i=1}^n x_i + \beta} \right)^{n+\alpha-1} \cdot e^{-t} \cdot \frac{1}{\sum_{i=1}^n x_i + \beta} dt$$

$$\Rightarrow \left(\frac{1}{\sum_{i=1}^n x_i + \beta} \right)^{n+\alpha} \cdot \int_{-\infty}^{+\infty} t^{n+\alpha-1} \cdot e^{-t} \cdot dt$$

$$\frac{(\sum_{i=1}^n x_i + \beta)^{n+\alpha}}{\Gamma(n+\alpha)} \cdot \lambda^{n+\alpha-1} \cdot e^{-\lambda(\sum_{i=1}^n x_i + \beta)}$$

so for $P(\lambda | D)$ we have

$$\text{so } \alpha' = n + \alpha \quad \beta' = \sum_{i=1}^n x_i + \beta$$

we have $\text{Gamma}(\alpha', \beta')$

(b)

$$\begin{aligned}
 0 &= p(\lambda | D)' = \frac{(\beta + \sum_{i=1}^n x_i)^{\alpha+\alpha}}{\Gamma(n+\alpha)} \cdot (\lambda^{n+\alpha-1} \cdot (n+\alpha-1) \cdot e^{-(\beta + \sum_{i=1}^n x_i) \cdot \lambda} + \lambda^{n+\alpha-1} \cdot (-\beta - \sum_{i=1}^n x_i) \cdot e^{-(\beta + \sum_{i=1}^n x_i) \cdot \lambda}) \\
 &= (n+\alpha-1) \cdot \lambda^{n+\alpha-2} + \lambda^{n+\alpha-1} \cdot (-\beta - \sum_{i=1}^n x_i) = 0 \\
 \Rightarrow n+\alpha-1 &= \lambda \beta + \lambda \sum_{i=1}^n x_i \\
 \Rightarrow \lambda &= \frac{n+\alpha-1}{\beta + \sum_{i=1}^n x_i}
 \end{aligned}$$

(c) $\hat{\lambda} = \frac{\sum_{i=1}^n x_i}{n}$ = $E(\hat{\lambda}_{MLE})$

$$\because \sum_{i=1}^n x_i \sim \text{Gamma}(n, \lambda)$$

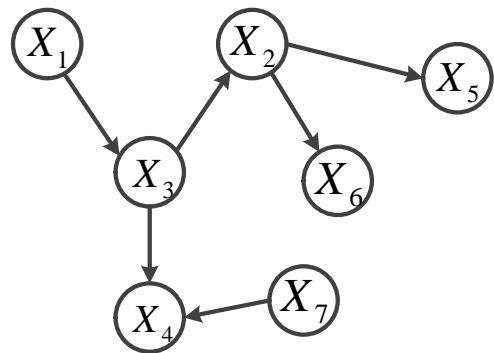
$$\begin{aligned}
 \therefore E\left[\sum_{i=1}^n x_i\right] &= \int_0^\infty \frac{\lambda^n}{\Gamma(n)} x^{n-1} e^{-\lambda x} \cdot x \cdot dx = \frac{\lambda \Gamma(n-1)}{\Gamma(n)} \int_0^\infty \frac{\lambda^{n-1}}{\Gamma(n-1)} x^{n-2} e^{-\lambda x} dx \\
 &= \frac{\lambda}{n-1}
 \end{aligned}$$

$$\therefore E\left[\frac{\sum_{i=1}^n x_i}{n}\right] = \frac{n\lambda}{n-1}$$

$$E\left(\frac{n-1}{n} \hat{\lambda}_{MLE}\right) = \frac{n-1}{n} E(\hat{\lambda}_{MLE}) = \frac{n-1}{n} \frac{n\lambda}{n-1} = \lambda$$

unbiased

4. [10 points] [Graphical Models] Given the following Bayesian Network,



answer the following questions.

- Factorize the joint distribution of X_1, \dots, X_7 according to the given Bayesian Network. [2 points]
- Justify whether $X_1 \perp X_5 | X_2$? [2 points]
- Justify whether $X_5 \perp X_7 | X_3, X_4$? [2 points]
- Justify whether $X_5 \perp X_7 | X_4$? [2 points]
- Write down the variables that are in the Markov blanket of X_3 . [2 points]

Solution:

$$(a) P(X_1, \dots, X_7) = P(X_1) P(X_2 | X_1) P(X_3 | X_1) P(X_4 | X_3, X_7) \\ P(X_5 | X_2) P(X_6 | X_2) P(X_7)$$

(b) $X_1 \rightarrow X_3 \rightarrow X_2$ is active triple. } \Rightarrow block
 $X_3 \rightarrow X_2 \rightarrow X_5$ is inactive triple.
 Yes

(c) $X_5 \perp X_7 | X_3, X_4$
 $X_3 \rightarrow X_2 \rightarrow X_5$ is active triple
 $X_4 \leftarrow X_3 \rightarrow X_2$ is inactive triple
 $X_3 \rightarrow X_4 \leftarrow X_7$ is active } \Rightarrow block
 Yes

(d) No
 $X_3 \rightarrow X_2 \rightarrow X_5$ active.
 $X_4 \leftarrow X_3 \rightarrow X_2$ active
 $X_3 \rightarrow X_4 \leftarrow X_7$ active } \Rightarrow unblock.

(e) $X_1 \ X_2 \ X_4 \ X_7$