

STA414: Statistical Methods for Machine Learning II

Homework 3

Tianyu (Shin) Ren

1002379188

University of Toronto

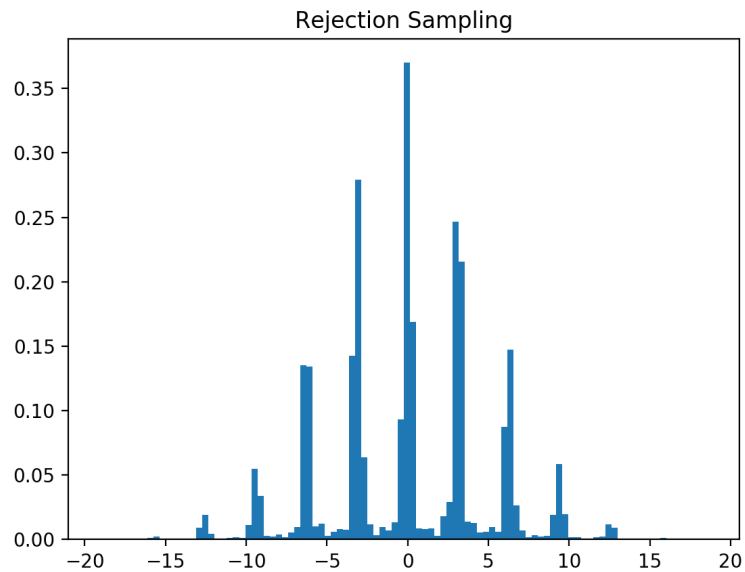
April 2019

1. Nature's rejection sampler

- (a) **Estimate $\int p(x, g = 0 | \theta = 0) dx$ by summing over 10,000 values of x , equally spaced from -20 to 20**

The estimate of $p(g = 0 | \theta = 0)$: 0.8

- (b) **Implement a rejection sampler to sample from $p(x | g = 1, \theta = 0)$**

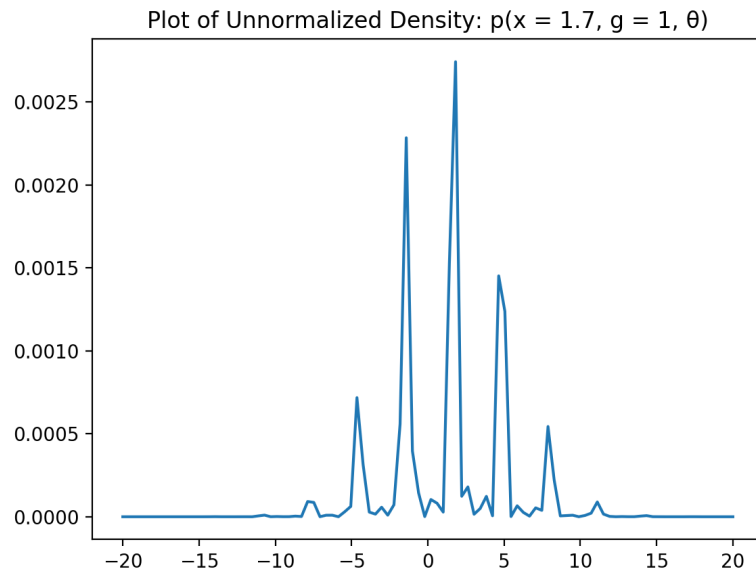


The fraction of accepted samples: 0.2

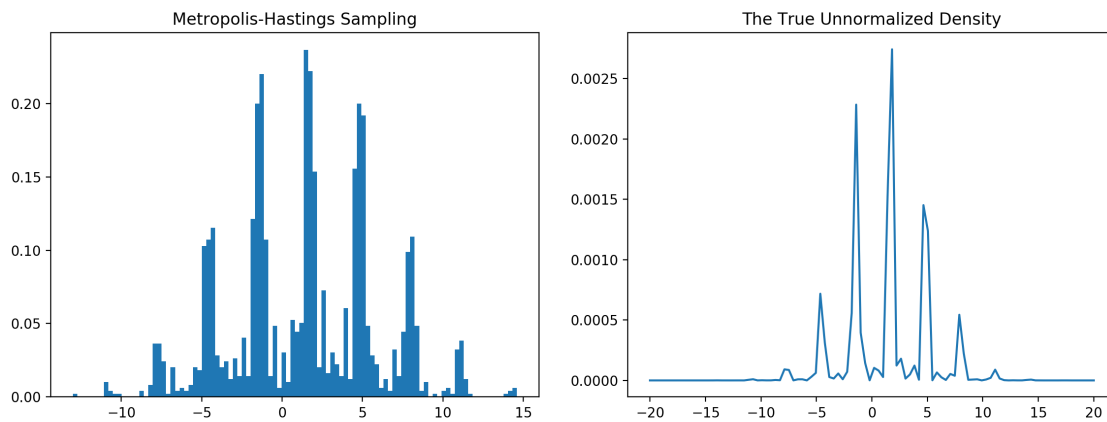
- (c) **Implement a self-normalized importance sampler to estimate $p(g = 0 | \theta = 0)$. Use $p(x | \theta = 0)$ as a proposal distribution.**

The estimate of the fraction of photons that are absorbed: 0.79

- (d) **Plot the unnormalized density $p(x = 1.7, g = 1, \theta)$, as a function of θ**



- (e) Write a Metropolis-Hastings sampler to sample from $p(\theta|x, a = 1)$. For the proposal distribution, use a Gaussian centered around the current sample



- (f) Use samples from your MH chain to estimate the posterior probability $p(-3 < \theta < 3|x = 1.7, g = 1)$

Estimate of the posterior probability: 0.47

2. Gradient estimators

- (a) Prove $\mathbb{E}_{p(x|\theta)} [\nabla_{\theta} \log p(x|\theta)] = 0$

$$\begin{aligned}
\mathbb{E}_{p(x|\theta)} [\nabla_{\theta} \log p(x|\theta)] &= \int \nabla_{\theta} \log p(x|\theta) p(x|\theta) dx \\
&= \int \frac{\nabla_{\theta} p(x|\theta)}{p(x|\theta)} p(x|\theta) dx \\
&= \int \nabla_{\theta} p(x|\theta) dx \\
&= \nabla_{\theta} \int p(x|\theta) dx \\
&= \nabla_{\theta} 1 \\
&= 0
\end{aligned}$$

(b) **Show that** $\mathbb{E}_{p(b|\theta)} \left[f(b) \frac{\partial}{\partial \theta} \log p(b|\theta) \right] = \frac{\partial}{\partial \theta} \mathbb{E}_{p(b|\theta)} [f(b)]$

$$\begin{aligned}
\mathbb{E}_{p(b|\theta)} \left[f(b) \frac{\partial}{\partial \theta} \log p(b|\theta) \right] &= \int f(b) \frac{\partial}{\partial \theta} \log p(b|\theta) \cdot p(b|\theta) db \\
&= \int f(b) \frac{\frac{\partial}{\partial \theta} p(b|\theta)}{p(b|\theta)} p(b|\theta) db \\
&= \int f(b) \frac{\partial}{\partial \theta} p(b|\theta) db \\
&= \frac{\partial}{\partial \theta} \int f(b) p(b|\theta) db \\
&= \frac{\partial}{\partial \theta} \mathbb{E}_{p(b|\theta)} [f(b)]
\end{aligned}$$

(c) **Show that** $\mathbb{E}_{p(b|\theta)} \left[[f(b) - c] \frac{\partial}{\partial \theta} \log p(b|\theta) \right] = \frac{\partial}{\partial \theta} \mathbb{E}_{p(b|\theta)} [f(b)]$, **for any fixed** c

$$\begin{aligned}
\mathbb{E}_{p(b|\theta)} \left[[f(b) - c] \frac{\partial}{\partial \theta} \log p(b|\theta) \right] &= \int [f(b) - c] \frac{\partial}{\partial \theta} \log p(b|\theta) \cdot p(b|\theta) db \\
&= \int f(b) \frac{\partial}{\partial \theta} p(b|\theta) db - c \int \frac{\partial}{\partial \theta} p(b|\theta) db \\
&= \frac{\partial}{\partial \theta} \mathbb{E}_{p(b|\theta)} [f(b)] - c \frac{\partial}{\partial \theta} \int p(b|\theta) db \\
&= \frac{\partial}{\partial \theta} \mathbb{E}_{p(b|\theta)} [f(b)] - c \frac{\partial}{\partial \theta} 1 \\
&= \frac{\partial}{\partial \theta} \mathbb{E}_{p(b|\theta)} [f(b)]
\end{aligned}$$

(d) **Give an example where** $\mathbb{E}_{p(b|\theta)} \left[[f(b) - c(b)] \frac{\partial}{\partial \theta} \log p(b|\theta) \right] \neq \frac{\partial}{\partial \theta} \mathbb{E}_{p(b|\theta)} [f(b)]$ **for some function** $c(b)$, **and show that it is biased.**

$$\text{Let } c(b) = f(b) \implies \mathbb{E}_{p(b|\theta)} \left[[f(b) - c(b)] \frac{\partial}{\partial \theta} \log p(b|\theta) \right] = 0$$

But $\frac{\partial}{\partial \theta} \mathbb{E}_{p(b|\theta)} [f(b)] \neq 0$ in general. So, it is biased.

3. Comparing variances of gradient estimators

(a) Compute $\mathbb{V} [\hat{L}_{MC}]$ as a function of D

$$\mathbb{V} [\hat{L}_{MC}] = \mathbb{V} \left[\sum_{d=1}^D x_d \right] = \sum_{d=1}^D \mathbb{V} [x_d] = D \cdot \mathbb{V} [x_d] = D$$

(b) Derive a closed form for $\hat{g}^{\text{SF}}[f] = [f(x) - c(\theta)] \frac{\partial}{\partial \theta} \log p(x|\theta)$

$$\begin{aligned} \hat{g}^{\text{SF}}[f] &= \left[\sum_{d=1}^D x_d - \sum_{d=1}^D \theta_d \right] \frac{\partial}{\partial \theta} \left[\log \left(\frac{1}{\sqrt[p]{2\pi}} \right) - \frac{1}{2} (x - \theta)^T (x - \theta) \right] \\ &= \left[\sum_{d=1}^D (x_d - \theta_d) \right] \frac{\partial}{\partial \theta} \left[\log \left(\frac{1}{\sqrt[p]{2\pi}} \right) - \frac{1}{2} (x^T x - 2\theta^T x + \theta^T \theta) \right] \\ &= \left[\sum_{d=1}^D (x_d - \theta_d) \right] (x - \theta) \\ &= (\mathbb{1}^T \epsilon) \epsilon \end{aligned}$$

$$\text{where } \epsilon = x - \theta = \begin{pmatrix} x_1 - \theta_1 \\ \vdots \\ x_d - \theta_d \end{pmatrix}, \quad \mathbb{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \quad \epsilon \sim \mathcal{N}(0, I)$$

(c) Derive the scalar value $\mathbb{V} [\hat{g}_1^{\text{SF}}]$ as a function of D

By Hint: $\mathbb{E}(\epsilon^3) = 0$, $\mathbb{E}(\epsilon^4) = 3$, $\mathbb{E}(\epsilon^2) = \mathbb{V}(\epsilon) + \mathbb{E}(\epsilon) = 1$, $\mathbb{E}(\epsilon_i \epsilon_j) = \text{COV}(\epsilon_i, \epsilon_j) + \mathbb{E}(\epsilon_i) \mathbb{E}(\epsilon_j) = 0$

$$\begin{aligned} \mathbb{V} [\hat{g}_1^{\text{SF}}] &= \mathbb{V}((\mathbb{1}^T \epsilon) \epsilon_1) \\ &= \mathbb{E}([(\mathbb{1}^T \epsilon) \epsilon_1]^2) - \mathbb{E}((\mathbb{1}^T \epsilon) \epsilon_1)^2 \\ &= \mathbb{E}(\epsilon_1^2 (\sum_d \epsilon_d^2 + 2 \sum_i \sum_j \epsilon_i \epsilon_j)) - \mathbb{E}((\mathbb{1}^T \epsilon) \epsilon_1)^2 \\ &= \mathbb{E}(\epsilon_1^2 \sum_d \epsilon_d^2) + 2 \sum_i \sum_j \mathbb{E}(\epsilon_i \epsilon_j) - \mathbb{E}((\mathbb{1}^T \epsilon) \epsilon_1)^2 \\ &= \mathbb{E}(\epsilon_1^4 + \epsilon_1^2 \epsilon_2^2 + \cdots + \epsilon_1^2 \epsilon_D^2) - \mathbb{E}(\epsilon_1^2 + \epsilon_1 \epsilon_2 + \cdots + \epsilon_1 \epsilon_D)^2 \\ &= \mathbb{E}(\epsilon_1^4) + \mathbb{E}(\epsilon_1^2) \mathbb{E}(\epsilon_2^2 + \cdots + \epsilon_D^2) - [\mathbb{E}(\epsilon_1^2) + \mathbb{E}(\epsilon_1) \mathbb{E}(\epsilon_2 + \cdots + \epsilon_D)]^2 \\ &= 3 + 1 \cdot (D - 1) - 1^2 \\ &= D + 1 \end{aligned}$$

- (d) **Derive this gradient estimator for $\nabla_{\theta} L(\theta)$, and give $\mathbb{V} [\hat{g}_1^{\text{REPARAM}}]$ as a function of D**

$$\begin{aligned}
 \nabla_{\theta} L(\theta) &= \nabla_{\theta} \mathbb{E}_{p(b|\theta)} [f(x)] \\
 &= \nabla_{\theta} \mathbb{E}_{p(\epsilon)} \left[\sum_{d=1}^D (\theta_d + \epsilon_d) \right] \\
 &= \mathbb{E}_{p(\epsilon)} \left[\nabla_{\theta} \sum_{d=1}^D (\theta_d + \epsilon_d) \right] \\
 &= \mathbb{E}_{p(\epsilon)} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\
 &= (1, \dots, 1)^T
 \end{aligned}$$

$$\mathbb{V} [\hat{g}_1^{\text{REPARAM}}] = \mathbb{V} \left[\frac{\partial f}{\partial x} \frac{\partial x}{\partial \theta} \right] = \mathbb{V} \left[\frac{\partial f}{\partial x} (1) \right] = \mathbb{V} [(1, \dots, 1)^T] = 0$$

Appendix

Python Code

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.stats import norm
4 np.random.seed(414)
5
6
7 # numerical integration
8 def p_x_given_theta(x, theta):
9     return norm.pdf(x, loc=theta, scale=4)
10
11 def f(x, theta):
12     return np.sin(5 * (x - theta)) ** 2 / (25 * np.sin(x - theta) ** 2)
13
14 def p_g_given_x_theta(x, theta, g=1):
15     return f(x, theta) ** g * (1 - f(x, theta)) ** (1 - g)
16
17 def numerical_integration(g, theta):
18     x = np.linspace(-20, 20, 10000)
19     prod = p_g_given_x_theta(x, theta, g) * p_x_given_theta(x, theta)
20     area = np.sum(prod) * 40 / len(x)
21     # f1 = plt.figure()
22     # ax1 = f1.add_subplot(111)
23     # plt.plot(x, prod)
24     return area
25
26 # rejection sampling
27 def rejection_sampling(iteration=1000):
28     samples = []
29     i, j = 0, 0
30     M = 1
31     while i < iteration:
32         u = np.random.uniform(0, 1)
33         x = np.random.normal(0, 4) / 0.8
34         if u < p_g_given_x_theta(x, theta=0) / M:
35             samples.append(x)
36             i += 1
37         j += 1
38     print("fraction of accepted samples: " + str(round(i/j, 2)))
39     return samples
40
41 # importance sampling
42 def importance_sampling(iteration=1000):
43     i, summation = 0, 0
44     while i < iteration:
45         x = np.random.normal(0, 4)
46         summation += p_g_given_x_theta(x, theta=0, g=0)
47         i += 1
48     return summation / i
49
50 # Q1(d)
51 def p_x_given_theta(x, theta):
52     return 1 / (4 * np.sqrt(2 * np.pi)) * np.exp(-(x - theta) ** 2 / 32)
53

```

```

54 def cauchy(theta):
55     return 1 / (10 * np.pi * (1 + (theta / 10) ** 2))
56
57 def p_x_g_theta(theta):
58     """
59     the unnormalized density p(x = 1.7, g = 1, theta), as a function of theta
60     """
61     return f(1.7, theta) * p_x_given_theta(1.7, theta) * cauchy(theta)
62
63 # Metropolis-Hastings sampling
64 def metropolis_hastings(iteration=10000):
65     accepted, rejected = [], []
66     i, x = 0, 0
67     while i < iteration:
68         u = np.random.uniform(0, 1)
69         x_new = np.random.normal(x, 100)
70         func = lambda a, b : p_g_given_x_theta(1.7, a) * p_x_given_theta(1.7, a) *
71             cauchy(a) * norm.pdf(b, loc=a, scale=100)
72         fraction = func(x_new, x) / func(x, x_new)
73         if u < fraction:
74             x = x_new
75             accepted.append(x)
76         else:
77             rejected.append(x)
78         i += 1
79     return accepted, rejected
80
81 if __name__ == '__main__':
82     # Q1(a) estimate the fraction of photons that get absorbed on average
83     print("estimate of p(g = 0|theta = 0):" + str(round(numerical_integration(g=0,
84         theta=0), 2)))
85
86 # Q1(b) plot rejection sampling
87 x = rejection_sampling(iteration=10000)
88 f2 = plt.figure()
89 ax2 = f2.add_subplot(111)
90 count, bins, ignored = plt.hist(x, 100, density=True)
91 plt.title("Rejection Sampling")
92
93 # Q1(c) estimate p(g = 0|theta = 0)
94 print("the estimate (importance sampling): " + str(round(importance_sampling(), 2)
95     ))
96
97 # Q1(d) plot unnormalized density p(x = 1.7, g = 1, theta)
98 theta = np.linspace(-20,20,100) # 100 linearly spaced numbers
99 y = p_x_g_theta(theta)
100 f3 = plt.figure()
101 ax3 = f3.add_subplot(111)
102 plt.plot(theta, y)
103 plt.title("Plot of Unnormalized Density: p(x = 1.7, g = 1, theta)")
104
105 # Q1(e) plot Metropolis-Hastings sampling
106 theta_accepted = metropolis_hastings(iteration=100000)[0]
107 f4 = plt.figure()
108 ax4 = f4.add_subplot(111)
109 count, bins, ignored = plt.hist(theta_accepted, 100, density=True)
110 plt.title("Metropolis-Hastings Sampling")

```



```
109 # plot the true unnormalized density
110 f5 = plt.figure()
111 ax5 = f5.add_subplot(111)
112 theta = np.linspace(-20,20,100)
113 y = p-g-given-x_theta(1.7, theta) * p-x-given-theta(1.7, theta) * cauchy(theta)
114 plt.plot(theta, y)
115 plt.title("The True Unnormalized Density")
116
117 # Q1(f) estimate the posterior probability
118 within = 0
119 for theta in theta_accepted:
120     if -3 < theta < 3:
121         within +=1
122 estimate_posterior = within / len(theta_accepted)
123 print("estimate of the posterior probability: " + str(round(estimate_posterior, 2)
124       ))
124 plt.show()
```