

Problem Set 2. Solution

2020-12

1 Problem 1

- (1) Let X have the standard Laplace distribution. Use *importance sampling* based on 100000 draws from the standard normal as the proposal density to estimate $\mathbb{E}(X)$, $\text{Var}(X)$ and $\mathbb{P}(X > 2)$.

Solution. Let $X \sim \text{Laplace}(0, 1) \sim f(x)$, $q(x) \sim \mathcal{N}(0, 1)$, and $w(x) := \frac{f(x)}{q(x)} = \sqrt{\frac{\pi}{2}} \exp\left(\frac{1}{2}x^2 - |x|\right)$. Following $\mathbb{E}_f[h(X)] = \mathbb{E}_q\left[h(X)\frac{f(X)}{q(X)}\right] = \mathbb{E}_q[h(X)w(X)]$, we have

$$\widehat{\mathbb{E}}(X) = \frac{1}{n} \sum_{i=1}^n x_i w(x_i)$$

$$\widehat{\text{Var}}(X) = \frac{1}{n} \sum_{i=1}^n x_i^2 w(x_i) - (\widehat{\mathbb{E}}(X))^2$$

$$\widehat{\mathbb{P}}(X > 2) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(x_i > 2) w(x_i)$$

```

1 import numpy as np
2 import pandas as pd
3 from scipy import stats
4 from IPython.display import display
5
6 def is_est(size=10**5, reverse=False):
7     if reverse:
8         x = np.random.laplace(size=size)
9         w = stats.norm.pdf(x) / stats.laplace.pdf(x)
10    else:
11        x = np.random.normal(size=size)
12        w = stats.laplace.pdf(x) / stats.norm.pdf(x)
13    EX = (x * w).mean()
14    VarX = ((x ** 2.0) * w).mean()
15    P = ((x > 2.0) * w).mean()
16    return x, w, EX, VarX, P
17
18 np.random.seed(1234)
19 x, w, EX, VarX, P = is_est()
20
21 df = pd.DataFrame({
22     "E(X)" :    ["0", f"{EX:.4f}"],
23     "Var(X)":   ["2", f"{VarX:.4f}"],
24     "P(X>2)":  [f"{1-stats.laplace.cdf(2):.4f}", f"{P:.4f}"],
25     }, index=["Truth", "Estimate"])
26 )

```

```

27 display(df)
28 print(f"Weights: [{w.min():.2f}, {w.max():.2f}]")

```

	E(X)	Var(X)	P(X>2)
1 Truth	0	2	0.0677
3 Estimate	-0.0043	1.5628	0.0583
4 Weights:	[0.76, 127.48]		

■

- (2) Plot the logarithms of the importance weights as a histogram. Is the distribution of these log-importance-weights symmetric? Do you occasionally get extremely large importance weights? Extremely small ones? Which type of outliers is more worrisome?

Solution.

```

1 import plotly.graph_objects as go
2
3 fig = go.Figure()
4 fig.add_trace(go.Histogram(
5     x=w, xbins=dict(size=0.1), opacity=0.6
6 ))
7 fig.update_xaxes(type="log")
8 fig.update_layout(
9     xaxis_title_text='Importance Weights', # xaxis label
10    yaxis_title_text='Count',               # yaxis label
11    bargap=0.01,                           # gap between bars
12    width=360, height=240, template="seaborn",
13    margin=dict(l=5, r=5, t=5, b=5),
14 )
15 fig.show()
16 fig.write_image("01-02_importance-weight.pdf")

```

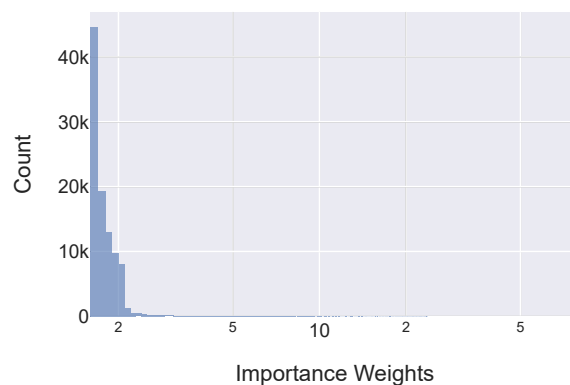


Figure 1.1: The distribution of log-importance-weights

The results show that the distribution of log-importance-weights is asymmetric. The outliers that have larger value, which lead to high variance, are more worrisome. ■

- (3) Compare your Monte Carlo estimates with the true values. Are there biases in the Monte Carlo estimates? How large are the Monte Carlo standard derivations? (Is there a theoretical formula for the standard deviation of your Monte Carlo estimator?)

Solution. The Monte Carlo estimates are unbiased. Standard derivations are given as follows.

```
1 print(f"{(x * w).std():.2f}, {(x ** 2.0) * w).std():.2f}, {(x > 2.0) * w).std():.2f}")
1 3.77, 13.78, 0.76
```

Let $I_n^{\text{IS}} := \widehat{\mathbb{E}}_f[h(X)] = \widehat{\mathbb{E}}_q[h(X)w(X)] = \frac{1}{n} \sum_{i=1}^n h(x_i)w(x_i)$, we have

$$\begin{aligned}\mathbb{E}_q[I_n^{\text{IS}}] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_q[h(X_i)w(X_i)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_f[h(X_i)] = \mathbb{E}_f[h(X)] =: I \\ \text{Var}(I_n^{\text{IS}}) &= \frac{1}{n} \text{Var}_q(h(X)w(X)) = \frac{1}{n} \left(\sqrt{\frac{\pi}{8}} \int_{\mathbb{R}} (h(x))^2 \exp\left(\frac{1}{2}x^2 - 2|x|\right) dx - I^2 \right).\end{aligned}$$

Then

$$\begin{aligned}\text{Var}(\widehat{\mathbb{E}}(X)) &= \frac{1}{n} \sqrt{\frac{\pi}{8}} \int_{\mathbb{R}} x^2 \exp\left(\frac{1}{2}x^2 - 2|x|\right) dx = \infty \\ \text{Var}(\widehat{\text{Var}}(X)) &= \frac{1}{n} \left(\sqrt{\frac{\pi}{8}} \int_{\mathbb{R}} x^4 \exp\left(\frac{1}{2}x^2 - 2|x|\right) dx - 4 \right) = \infty \\ \text{Var}(\widehat{\mathbb{P}}(X > 2)) &= \frac{1}{n} \left(\sqrt{\frac{\pi}{8}} \int_{x>2} \exp\left(\frac{1}{2}x^2 - 2|x|\right) dx - \frac{1}{4}e^{-4} \right) = \infty.\end{aligned}$$

- (4) Let Y have the standard normal distribution. Use *importance sampling* based on 100000 draws from the standard Laplace as the proposal density to estimate $\mathbb{E}(Y)$, $\text{Var}(Y)$ and $\mathbb{P}(Y > 2)$. Repeat parts (2) and (3).

Solution.

```
1 np.random.seed(1234)
2 x, w, EX, VarX, P = is_est(reverse=True)
3
4 df = pd.DataFrame({
5     "E(Y)" :    ["0", f"{EX:.4f}"],
6     "Var(Y)":   ["1", f"{VarX:.4f}"],
7     "P(Y>2)":  [f"{1-stats.norm.cdf(2):.4f}", f"{P:.4f}"],
8     }, index=["Truth", "Estimate"])
9 )
10 display(df)
11 print(f"Weights: [{w.min():.2f}, {w.max():.2f}]")
12
13 fig = go.Figure()
14 fig.add_trace(go.Histogram(
15     x=np.log(w), xbins=dict(size=1), opacity=0.6
16 ))
17 # fig.update_xaxes(type="log")
```

```

18 fig.update_layout(
19     xaxis_title_text='Importance Weights', # xaxis label
20     yaxis_title_text='Count',             # yaxis label
21     bargap=0.05,                          # gap between bars
22     width=360, height=240, template="seaborn",
23     margin=dict(l=5, r=5, t=5, b=5),
24 )
25 fig.show()
26 fig.write_image("01-02_importance-weight-reverse.pdf")
27 print(f"{{x * w}.std():.2f}, {{((x ** 2.0) * w).std():.2f}, {{(x > 2.0) * w).std():.2f}}")

```

	E(y)	Var(y)	P(y>2)
Truth	0	1	0.0228
Estimate	-0.0002	1.0007	0.0228
Weights: [0.00, 1.32]			
0.99, 1.08, 0.11			

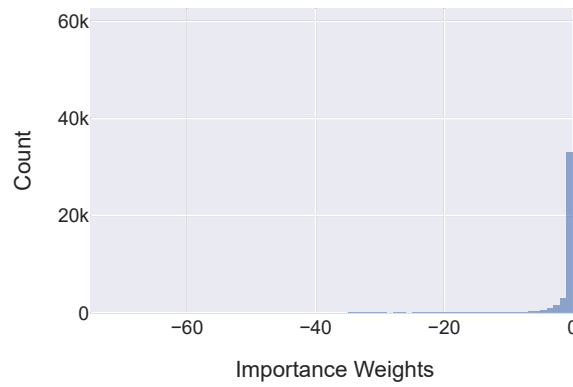


Figure 1.2: The distribution of log-importance-weights

$$\begin{aligned}
 \text{Var}(\widehat{\mathbb{E}}(X)) &= \frac{2}{n\pi} \int_{\mathbb{R}_+} x^2 \exp(-x^2 + x) dx \approx 0.0031^2 \approx \frac{0.99^2}{n} \\
 \text{Var}(\widehat{\text{Var}}(X)) &= \frac{1}{n} \left(\frac{2}{\pi} \int_{\mathbb{R}_+} x^4 \exp(-x^2 + x) dx - 1 \right) \approx 0.0034^2 \approx \frac{1.08^2}{n} \\
 \text{Var}(\widehat{\mathbb{P}}(X > 2)) &= \frac{1}{n} \left(\frac{1}{\pi} \int_{x>2} \exp(-x^2 + x) dx - 0.0228^2 \right) \approx 0.00035^2 \approx \frac{0.11^2}{n}.
 \end{aligned}$$

■

2 Problem 2

Consider the following Restricted Boltzmann Machine (RBM) with energy function

$$p_{\theta}(v) = \frac{1}{Z_{\theta}} \sum_h \exp(-E(v, h)), \quad E(v, h) = -b^T v - c^T h - h^T W v$$

Here the model parameters are $\theta = \{b, c, W\}$

(1) Show that $p(v \mid h) = \prod_{i=1}^n p(v_i \mid h)$, $p(h \mid v) = \prod_{i=1}^d p(h_i \mid v)$

Solution. It suffices to show the first identity since we have the second by symmetry.

$$\begin{aligned}
p(v, h) &= \frac{p(v, h)}{p(h)} \\
&= \frac{1}{p(h)} \frac{\exp(b^\top v + c^\top h + h^\top W v)}{\sum_u \exp(b^\top v + c^\top h + h^\top W u)} \\
&= \frac{1}{p(h)} \frac{\prod_i \exp(b_i v_i + c^\top h + h^\top W_{:,i} v_i)}{\sum_u \exp(b^\top v + c^\top h + h^\top W u)} \\
&= \frac{\prod_j \sum_{u_j} \exp(b_j u_j + c^\top h + h^\top W_{:,j} u_j)}{\sum_u \exp(b^\top v + c^\top h + h^\top W u)} \prod_i p(v_i | h) \\
&= \prod_i p(v_i | h)
\end{aligned}$$

■

(2) Derive the derivatives of the log-likelihood function w.r.t. the model parameters θ

Solution.

$$\begin{aligned}
L(\theta) &= \sum_i \log p(v^{(i)}) \\
\frac{\partial L}{\partial \theta} &= \sum_i \frac{\partial \log p(v^{(i)})}{\partial \theta} \\
&= \sum_i \frac{\partial}{\partial \theta} \left(\log \sum_h \exp(-E(v^{(i)}, h)) - \log \sum_{v, h} \exp(-E(v, h)) \right) \\
&= - \sum_i \frac{\sum_h \exp(-E(v^{(i)}, h)) \frac{\partial}{\partial \theta} E(v^{(i)}, h)}{\sum_h \exp(-E(v^{(i)}, h))} + \sum_i \frac{\sum_{v, h} \exp(-E(v, h)) \frac{\partial}{\partial \theta} E(v, h)}{\sum_{v, h} \exp(-E(v, h))} \\
&= - \sum_i \sum_h p(h | v^{(i)}) \frac{\partial}{\partial \theta} E(v^{(i)}, h) + \sum_i \sum_{v, h} p(v, h) \frac{\partial}{\partial \theta} E(v, h)
\end{aligned}$$

where

$$\begin{aligned}
\sum_h p(h | v) \frac{\partial E(v, h)}{\partial b_i} &= - \sum_h p(h | v) v_i = -v_i \\
\sum_h p(h | v) \frac{\partial E(v, h)}{\partial c_i} &= - \sum_h \sum_{k=1}^d p(h_k | v) h_i = - \sum_{h_i} p(h_i | v) h_i = -p(h_i = 1 | v) \\
\sum_h p(h | v) \frac{\partial E(v, h)}{\partial W_{i,j}} &= - \sum_h \sum_{k=1}^d p(h_k | v) h_i v_j = - \sum_{h_i} p(h_i | v) h_i v_j = -p(h_i = 1 | v) v_j,
\end{aligned}$$

then

$$\begin{aligned}
\frac{\partial L}{\partial b_i} &= \sum_k \left(v_i^{(k)} - \sum_v p(v) v_i \right) \\
\frac{\partial L}{\partial c_i} &= \sum_k \left(p(h_i = 1 | v^{(k)}) - \sum_v p(v) p(h_i = 1 | v) \right) \\
\frac{\partial L}{\partial W_{i,j}} &= \sum_k \left(p(h_i = 1 | v^{(k)}) v_j^{(k)} - \sum_v p(v) p(h_i = 1 | v) v_j \right)
\end{aligned}$$

■

(3) Use the following code to load the MNIST data set.

```

1 from sklearn.datasets import fetch_openml
2
3 X, y = fetch_openml('mnist_784', version=1, return_X_y=True)
4 X = (X/255).astype('float32')
5 X_train, X_test = X[:60000,:], X[60000:,:]

```

Train your RBM on the training data set using contrastive divergence ($k = 1$), with Gibbs sampling for the energy induced distribution. Report the reconstruction error $\|v - \tilde{v}\|^2$ on the training data and the test data as a function of the number of iterations, where \tilde{v} is the sample after $k = 1$ iteration of Gibbs sampling that starts at v . Do you have any interesting finding? Explain it.

Solution. Here we refer to the Learnergy framework proposed by [Roder *et al.* \(2020\)](#)¹.

```

1  import time
2  import torch
3  import torchvision
4  from torch.utils.data import DataLoader
5  from tqdm.notebook import tqdm
6  from learnergy.models.bernoulli import RBM
7  # import learnergy.utils.logging as l
8  # logger = l.get_logger(__name__)
9
10 class MyRBM(RBM):
11     def __init__(self, **kwargs):
12         super(MyRBM, self).__init__(**kwargs)
13
14     def train_func(
15         self, train_dataset, test_dataset, batch_size=64, epochs=10
16     ):
17         batches = DataLoader(train_dataset, batch_size=batch_size,
18                             shuffle=True, num_workers=0)
19         for epoch in range(epochs):
20             # logger.info('Epoch %d/%d', epoch+1, epochs)
21             print(f'Epoch {epoch}')
22             start = time.time()
23             mse = 0
24             pl = 0
25             for samples, _ in tqdm(batches):
26                 samples = samples.reshape(len(samples), self.n_visible)
27                 if self.device == 'cuda':
28                     samples = samples.cuda()
29                 _, _, _, visible_states = self.gibbs_sampling(samples)
30                 visible_states = visible_states.detach()
31                 cost = (
32                     torch.mean(self.energy(samples))
33                     - torch.mean(self.energy(visible_states))
34                 )
35                 self.optimizer.zero_grad()
36                 cost.backward()
37                 self.optimizer.step()
38                 batch_size = samples.size(0)
39                 batch_mse = torch.div(

```

¹See <https://github.com/gugarosa/learnergy/blob/master/learnergy/models/bernoulli/rbm.py>

```

40         torch.sum(torch.pow(samples - visible_states, 2)), batch_size).detach()
41         batch_pl = self.pseudo_likelihood(samples).detach()
42         mse += batch_mse
43         pl += batch_pl
44         mse /= len(batches)
45         pl /= len(batches)
46         mse_test, _ = self.reconstruct(test_dataset)
47         end = time.time()
48         self.dump(mse=mse.item(), mse_test=mse_test.item(),
49                 pl=pl.item(), time=end-start)
50         # logger.info('MSE: %f | MSE_test: %f | log-PL: %f', mse, mse_test, pl)
51
52 if __name__ == "__main__":
53     import plotly.graph_objects as go
54     # mnist data
55     train = torchvision.datasets.MNIST(
56         root='../dataset/mnist', train=True, download=True,
57         transform=torchvision.transforms.ToTensor()
58     )
59     test = torchvision.datasets.MNIST(
60         root='../dataset/mnist', train=False, download=True,
61         transform=torchvision.transforms.ToTensor()
62     )
63     # RBM
64     model = MyRBM(n_visible=784, n_hidden=128, steps=1, learning_rate=0.1,
65                 momentum=0, decay=0, use_gpu=True)
66     # Training
67     n_epochs = 500
68     model.train_func(train, test, batch_size=128, epochs=n_epochs)
69     # Saving
70     torch.save(model, 'model.pth')
71     mse = model.history["mse"]
72     mse_test = model.history["mse_test"]
73
74     fig = go.Figure()
75     fig.add_trace(go.Scatter(
76         y=mse, name="Training Error",
77         mode='lines+markers', opacity=0.6,
78     ))
79     fig.add_trace(go.Scatter(
80         y=mse_test, name="Testing Error",
81         mode='lines+markers', opacity=0.6,
82     ))
83     fig.update_layout(
84         xaxis_title_text='epoch',

```

```

85     yaxis_title_text='Reconstruction Error',
86     width=700, height=240, template="seaborn",
87     margin=dict(l=5, r=5, t=5, b=5),
88 )
89 fig.show()
90 fig.write_image("02-03_reconstruction-error.pdf")

```

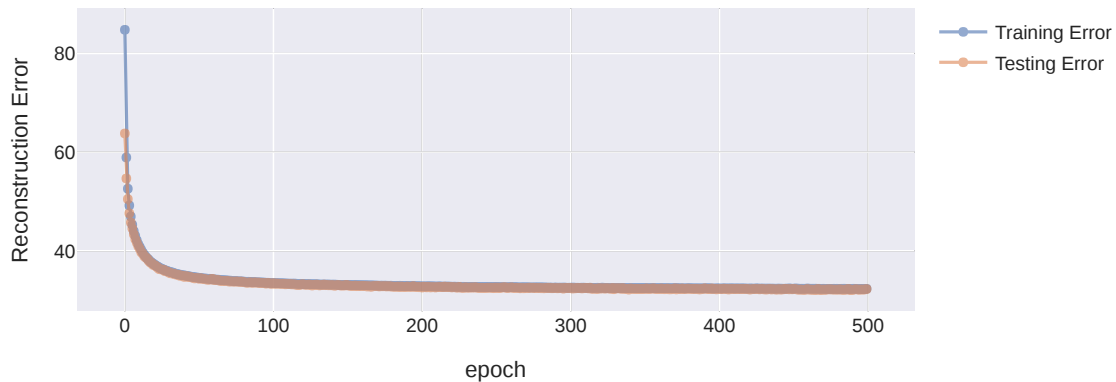


Figure 2.3: Reconstruction Errors

- (4) Generate samples from your trained RBM using Gibbs sampling with 10 independent chains, and each chain runs 200 iterations. Show the results in a 10×11 grid plot where columns correspond to samples at every 20 iterations.

Solution.

```

1  import torch
2  import plotly.graph_objs as go
3  from plotly.subplots import make_subplots
4
5  torch.manual_seed(1234)
6  rows, cols = 10, 11
7  batches = torch.FloatTensor(rows, 1, 28, 28).normal_(0, 256)
8
9  fig = make_subplots(
10     rows=rows, cols=cols,
11     shared_xaxes="all", shared_yaxes="all",
12 )
13  for row, data in zip(range(1, rows+1), batches):
14     v = data.to("cuda").reshape(1, model.n_visible)
15     pos_hidden_probs, pos_hidden_states = model.hidden_sampling(v)
16     neg_hidden_states = pos_hidden_states
17     fig.add_trace(go.Heatmap(
18         z=v.reshape(28, 28).detach().cpu().numpy(),
19         colorscale='Viridis', showscale = False
20     ), row=row, col=1,

```



```

21 )
22 for i in range(1, 201):
23     visible_probs, visible_states = model.visible_sampling(
24         neg_hidden_states, False
25     )
26     neg_hidden_probs, neg_hidden_states = model.hidden_sampling(
27         visible_states, False
28     )
29     if i % 20 == 0:
30         col = i//20+1
31         fig.add_trace(go.Heatmap(
32             z=visible_probs.reshape(28, 28).detach().cpu().numpy(),
33             colorscale='Viridis', showscale = False
34         ), row=row, col=col,
35         )
36         fig.update_xaxes(title_text=f"{i}", row=rows, col=col)
37         # fig.update_yaxes(title_text=f"{label.item()}", row=row, col=1)
38
39 fig.update_xaxes(showticklabels=False)
40 fig.update_yaxes(autorange='reversed', showticklabels=False)
41 fig.update_layout(
42     height=600, width=600,
43     template = "seaborn",
44     margin=dict(l=5, r=5, t=5, b=5),
45 )
46 fig.show()
47 fig.write_image("02-04_Gibbs-sampling.pdf")

```

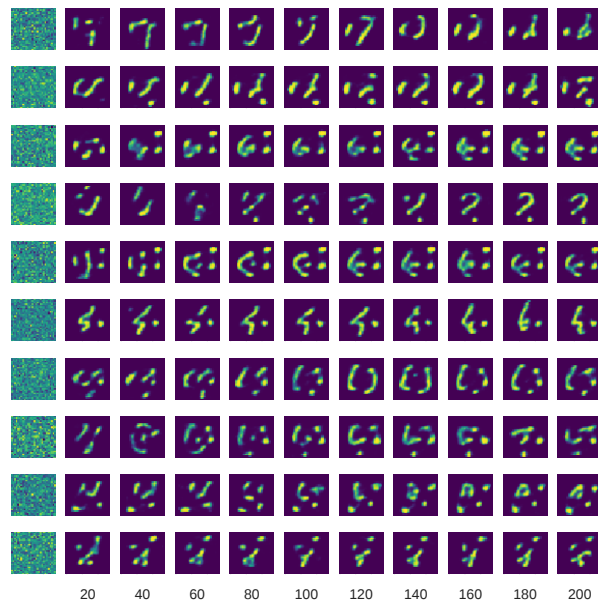


Figure 2.4: Gibbs sampling



3 Problem 3

Consider a logistic regression model with normal priors

$$y_i \sim \text{Bernoulli}(p_i), p_i = \frac{1}{1 + \exp(-x_i^\top \beta)}, \quad i = 1, \dots, n. \quad \beta \sim \mathcal{N}(0, \sigma_\beta^2)$$

where $\sigma_\beta = 1$. Download the data from the course website.

- (1). Implement a Hamiltonian Monte Carlo sampler to collect 500 samples (with 500 discarded as burn-in), show the scatter plot. Test the following two strategies for the number of leapfrog steps L : (1) use a fixed L ; (2) use a random one, say Uniform $(1, L_{\max})$. Do you find any difference? Explain it.

Solution. The posterior distribution of β is

$$p(\beta | y, X) \propto p(y | \beta, X) p(\beta) \propto \exp\left(-\frac{1}{2} \|\beta\|_2^2\right) \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}.$$

The potential energy function is

$$\begin{aligned} U(\beta) &= \frac{1}{2} \|\beta\|_2^2 - \log \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} + \text{constant} \\ &= \frac{1}{2} \|\beta\|_2^2 - \sum_{i=1}^n [y_i \log p_i + (1 - y_i) \log (1 - p_i)] + \text{constant} \\ &= \frac{1}{2} \|\beta\|_2^2 + \sum_{i=1}^n [(1 - y_i) x_i^\top \beta + \log (1 + \exp(-x_i^\top \beta))] + \text{constant} \\ &= \frac{1}{2} \|\beta\|_2^2 + (X\beta)^\top (1 - y) + 1^\top \log (1 + \exp(-X\beta)) + \text{constant}. \end{aligned}$$

The Hamiltonian is $H(\beta, r) = U(\beta) + K(r)$ where $K(r) = \frac{1}{2} \|r\|_2^2$. Then,

$$\begin{aligned} \frac{\partial H}{\partial \beta} &= \frac{\partial U}{\partial \beta} = \beta + \sum_{i=1}^n (p_i - y_i) x_i = \beta + X^\top \left(\frac{1}{1 + \exp(-X\beta)} - y \right) \\ \frac{\partial H}{\partial r} &= \frac{\partial K}{\partial r} = r \end{aligned}$$

```

1 import numpy as np
2 from scipy.special import expit as sigmoid
3
4 data = np.load("./mcs_hw2_p3_data.npy")
5 # X, y = data[:, :2], data[:, -1]
6
7 class HamiltonianLogisticRegression:
8     def __init__(self, data):
9         self.X, self.y = data[:, :2], data[:, -1]
10
11     def potential_energy(self, beta):
12         """ U(beta), potential energy
13         """
14         X, y = self.X, self.y
15         Xbeta = X@beta

```

```

16         return (
17             0.5*(beta@beta) + Xbeta@(1-y) - np.sum(np.log(sigmoid(Xbeta)))
18         )
19
20     def kinetic_energy(self, r):
21         """ K(r), Euclidean-Gaussian kinetic energy
22         """
23         return 0.5*(r@r)
24
25     def potential_energy_grad(self, beta):
26         X, y = self.X, self.y
27         return beta + X.T@(sigmoid(X@beta)-y)
28
29     def kinetic_energy_grad(self, r):
30         return r
31
32     def hamiltonian(self, beta, r):
33         return self.potential_energy(beta) + self.kinetic_energy(r)
34
35     def leapfrog(self, beta0, r0, n_leap_steps, leap_size):
36         for i in range(n_leap_steps):
37             r1 = r0 - 0.5 * leap_size * self.potential_energy_grad(beta0)
38             beta2 = beta0 + leap_size * self.kinetic_energy_grad(r1)
39             r2 = r1 - 0.5 * leap_size * self.potential_energy_grad(beta2)
40             beta0, r0 = beta2.copy(), r2.copy()
41         return beta0, r0
42
43     def hamiltonian_monte_carlo(
44         self, n_samples, n_discards, n_leapfrog_steps,
45         leap_size=0.1, fixed=True,
46     ):
47         dim = self.X.shape[-1]
48         n_leap_steps = n_leapfrog_steps
49         samples = np.zeros((n_samples + n_discards, dim))
50         beta = np.random.randn(dim)
51         for i in range(n_samples + n_discards):
52             if not fixed:
53                 n_leap_steps = np.random.randint(1, n_leapfrog_steps+1)
54                 r = np.random.normal(size=dim)
55                 beta0, r0 = self.leapfrog(
56                     beta0=beta, r0=r, n_leap_steps=n_leap_steps, leap_size=leap_size,
57                 )
58                 dH = self.hamiltonian(beta, r) - self.hamiltonian(beta0, r0)
59                 if np.random.random() < min(1, np.exp(dH)):
60                     beta = beta0.copy()

```

```

61     samples[i] = beta
62     return samples[n_discards:-1, :]
63
64 if __name__ == "__main__":
65     np.random.seed(1234)
66     hmc = HamiltonianLogisticRegression(data=data)
67     n_samples = n_discards = 500
68     sample_fix_steps = hmc.hamiltonian_monte_carlo(
69         n_samples, n_discards, leap_size=0.005, n_leapfrog_steps=4,
70     )
71     sample_rand_steps = hmc.hamiltonian_monte_carlo(
72         n_samples, n_discards, leap_size=0.005, n_leapfrog_steps=8, fixed=False
73     )
74
75     import plotly.graph_objects as go
76     fig = go.Figure()
77     fig.add_trace(go.Scatter(
78         x=sample_fix_steps[:,0], y=sample_fix_steps[:,1], name="L = 4",
79         mode='markers', opacity=0.9, marker_symbol="circle-open",
80     ))
81     fig.add_trace(go.Scatter(
82         x=sample_rand_steps[:,0], y=sample_rand_steps[:,1], name="L ~ U(1,8)",
83         mode='markers', opacity=0.9, marker_symbol="square-open",
84     ))
85     fig.update_layout(
86         width=700, height=480, template="plotly_white",
87         margin=dict(l=5, r=5, t=5, b=5),
88     )
89     fig.show()
90     fig.write_image("03-01_hamiltonian-monte-carlo.pdf")

```

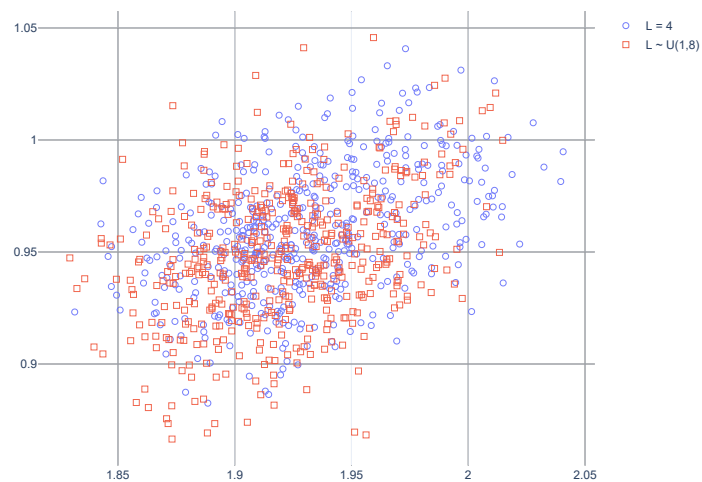


Figure 3.5: Hamiltonian Monte Carlo



- (2). Run HMC for 100000 iterations and discard the first 50000 samples as burn-in to form the ground truth. Implement stochastic gradient MCMC algorithms including SGLD, SGHMC and SGNHT. Show the convergence rate of different SGMCMC algorithms in terms of KL divergence to the ground truth as a function of iterations. You may want to use the ITE package <https://bitbucket.org/szzoli/ite-in-python/src/master/> to compute the KL divergence between two samples.

Solution.

```
1 class SGMCMC(HamiltonianLogisticRegression):
2     def __init__(self, **kwargs):
3         super(SGMCMC, self).__init__(**kwargs)
4
5     def potential_energy_sgrad(self, beta, batch_size):
6         n_data = self.X.shape[0]
7         indices = np.random.randint(0, n_data, size=batch_size)
8         X, y = self.X[indices], self.y[indices]
9         return beta + (n_data/batch_size)*X.T@ (sigmoid(X@beta)-y)
10
11     def sgld(
12         self, n_samples, batch_size=64,
13         truth=None, kl_divergence=None, interval=1,
14     ):
15         n_data, dim = self.X.shape[0], self.X.shape[1]
16         beta = np.random.randn(dim)
17         samples, kl_paths = [], {"i": [], "v": []}
18
19         for i in range(n_samples):
20             eps = 1/(((i//100)+10)*20)
21             grad = self.potential_energy_sgrad(beta, batch_size)
22             beta1 = beta - 0.5*eps*grad + np.random.normal(0, eps, dim)
23             dU = self.potential_energy(beta)-self.potential_energy(beta1)
24             if dU >= 0 or np.random.rand() < np.exp(dU):
25                 beta = beta1
26             samples.append(beta.copy())
27             if (i > 10) and (i % interval == 0):
28                 samples_arr = np.array(samples)
29                 kl_distance = kl_divergence(truth, samples_arr[:i])
30                 print([i, kl_distance])
31                 kl_paths["i"].append(i)
32                 kl_paths["v"].append(kl_distance)
33         return np.array(samples), kl_paths
34
35     def sghmc(
36         self, n_samples, n_leapfrog_steps=4, batch_size=64,
37         truth=None, kl_divergence=None, interval=1,
```

```

38     ):
39     n_data, dim = self.X.shape[0], self.X.shape[1]
40     beta = np.random.randn(dim)
41     samples, kl_paths = [], {"i":[], "v":[]}
42
43     for i in range(n_samples):
44         eps = 1/(((i//1000)+10)*20)
45         r = np.random.normal(0,1,2)
46         for j in range(n_leapfrog_steps):
47             beta += eps*r
48             grad = self.potential_energy_sgrad(beta, batch_size)
49             r += -eps*grad - eps*r+np.random.normal(0, 2*eps, 2)
50         samples.append(beta.copy())
51         if (i > 10) and (i % interval == 0):
52             samples_arr = np.array(samples)
53             kl_distance = kl_divergence(truth, samples_arr[:i])
54             print([i, kl_distance])
55             kl_paths["i"].append(i)
56             kl_paths["v"].append(kl_distance)
57     return np.array(samples), kl_paths
58
59 def sgnht(
60     self, n_samples, batch_size=64,
61     truth=None, kl_divergence=None, interval=1,
62     ):
63     n_data, dim = self.X.shape[0], self.X.shape[1]
64     beta, r = np.random.randn(dim), np.random.randn(dim)
65     samples, kl_paths = [], {"i":[], "v":[]}
66     A = 1
67     xi = A
68     for i in range(n_samples):
69         eps = 1/(((i//1000)+10)*20)
70         grad = self.potential_energy_sgrad(beta, batch_size)
71         r = r - eps*grad - eps*xi*r + np.sqrt(2*A)*np.random.normal(0, eps)
72         beta = beta + eps*r
73         xi = xi + eps*(r@r/2-1)
74         samples.append(beta)
75         if (i > 10) and (i % interval == 0):
76             samples_arr = np.array(samples)
77             kl_distance = kl_divergence(truth, samples_arr[:i])
78             print([i, kl_distance])
79             kl_paths["i"].append(i)
80             kl_paths["v"].append(kl_distance)
81     return np.array(samples), kl_paths
82

```

```

83
84 if __name__ == "__main__":
85     import os
86     import ite
87     co = ite.cost.BDKL_KnnK()
88     kl_divergence = co.estimate
89     np.random.seed(1234)
90
91     hmc = HamiltonianLogisticRegression(data=data)
92     n_samples, n_discards = 100000, 50000
93     if os.path.isfile(f"samples_hmc_{n_samples}.txt"):
94         sample_hmc = np.loadtxt(f"samples_hmc_{n_samples}.txt")
95     else:
96         sample_hmc = hmc.hamiltonian_monte_carlo(
97             n_samples, n_discards, leap_size=0.005, n_leapfrog_steps=4,
98         )
99         np.savetxt(f"samples_hmc_{n_samples}.txt", sample_hmc)
100
101     sgcmc = SGMCMC(data=data)
102     batch_size = 64
103     n_samples = 100000
104     sample_sgld, kl_sgld = sgcmc.sgld(
105         n_samples, batch_size=batch_size,
106         truth=sample_hmc, kl_divergence=kl_divergence, interval=5000,
107     )
108     sample_sghmc, kl_sghmc = sgcmc.sghmc(
109         n_samples, n_leapfrog_steps=3, batch_size=batch_size,
110         truth=sample_hmc, kl_divergence=kl_divergence, interval=5000,
111     )
112     sample_sgnht, kl_sgnht = sgcmc.sgnht(
113         n_samples, batch_size=batch_size,
114         truth=sample_hmc, kl_divergence=kl_divergence, interval=5000,
115     )
116
117     import plotly.graph_objects as go
118     fig = go.Figure()
119     fig.add_trace(go.Scatter(
120         x=kl_sgld["i"], y=kl_sgld["v"], name="SGLD",
121         mode='lines+markers', opacity=0.6,
122     ))
123     fig.add_trace(go.Scatter(
124         x=kl_sghmc["i"], y=kl_sghmc["v"], name="SGHMC",
125         mode='lines+markers', opacity=0.6,
126     ))
127     fig.add_trace(go.Scatter(

```

```

128     x=kl_sgnht["i"], y=kl_sgnht["v"], name="SGNHT",
129     mode='lines+markers', opacity=0.6,
130 ))
131 fig.update_layout(
132     xaxis_title_text='Iterations',
133     yaxis_title_text='KL-Divergence',
134     width=700, height=240, template="seaborn",
135     margin=dict(l=5, r=5, t=5, b=5),
136 )
137 fig.show()
138 fig.write_image("03-02_KL-Divergence.pdf")

```

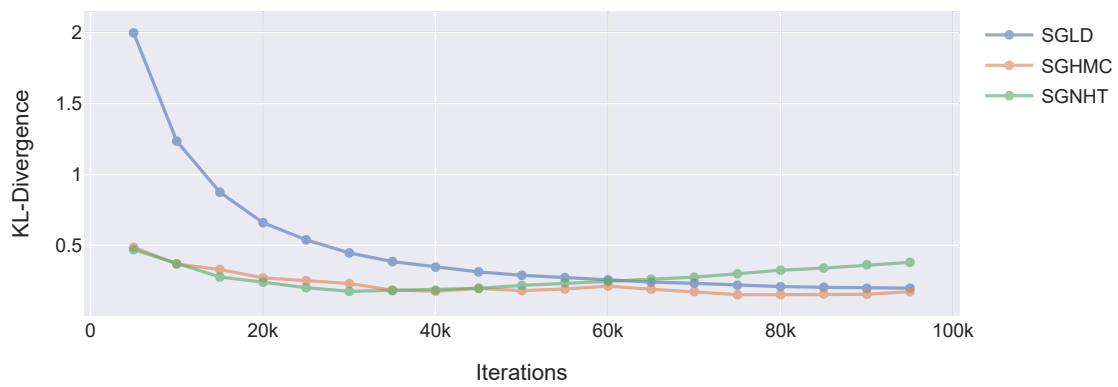


Figure 3.6: KL-Divergence

■

References

Roder, M., de Rosa, G. H. and Papa, J. P. (2020) Learnergy: Energy-based machine learners. 6