

HW1 solutions

Problem 1 (a) Consider the following optimization problem:

$$\begin{aligned} \max_f \text{Ent}(f) &:= \int_{-\infty}^{+\infty} f(x) \ln f(x) dx, \\ \text{s.t. } \int_{-\infty}^{+\infty} f(x) dx &= 1, \\ \int_{-\infty}^{+\infty} x f(x) dx &= 0, \\ \int_{-\infty}^{+\infty} x^2 f(x) dx &= 1. \end{aligned}$$

We set the Lagrangian as follows:

$$\mathcal{L}(f, \lambda, \mu, \nu) = \int_{-\infty}^{+\infty} f(x) \ln f(x) dx + \lambda \left(1 - \int_{-\infty}^{+\infty} f(x) dx \right) + \mu \left(- \int_{-\infty}^{+\infty} x f(x) dx \right) + \nu \left(1 - \int_{-\infty}^{+\infty} x^2 f(x) dx \right).$$

The functional derivative $\frac{\partial \mathcal{L}}{\partial f}$ is defined as: for any test function h , the following equality holds:

$$\int_{-\infty}^{+\infty} \frac{\partial \mathcal{L}}{\partial f} h dx = \left. \frac{\partial}{\partial \varepsilon} \right|_{\varepsilon=0} \mathcal{L}(f + \varepsilon h, \lambda, \mu, \nu).$$

Hence

$$\int_{-\infty}^{+\infty} \frac{\partial \mathcal{L}}{\partial f} h dx = \int_{-\infty}^{+\infty} (h \ln f + h) dx - \lambda \int_{-\infty}^{+\infty} h dx - \mu \int_{-\infty}^{+\infty} x h dx - \nu \int_{-\infty}^{+\infty} x^2 h dx.$$

Which implies that

$$0 = \frac{\partial \mathcal{L}}{\partial f} = \ln f + 1 - \lambda - \mu x - \nu x^2 \Rightarrow f(x) = e^{\lambda-1-\mu^2/(4\nu)} e^{\frac{(x+\mu/(2\nu))^2}{2/(2\nu)}}.$$

Since f is a probability density function, we can see that X follows normal distribution with mean 0 and variance 1.

(b) Similar to the previous problem, we consider the following optimization problem:

$$\begin{aligned} \max_f \text{Ent}(f) &:= \int_{-\infty}^{+\infty} f(x) \ln f(x) dx, \\ \text{s.t. } \int_{-\infty}^{+\infty} f(x) dx &= 1, \\ \int_{-\infty}^{+\infty} x f(x) dx &= m_1, \\ \int_{-\infty}^{+\infty} x^2 f(x) dx &= m_2, \\ \dots, \\ \int_{-\infty}^{+\infty} x^k f(x) dx &= m_k, \end{aligned}$$

We set the Lagrangian as follows:

$$\begin{aligned} \mathcal{L}(f, \lambda, \mu_1, \mu_2, \dots, \mu_k) &= \int_{-\infty}^{+\infty} f(x) \ln f(x) dx + \lambda \left(1 - \int_{-\infty}^{+\infty} f(x) dx \right) \\ &\quad + \mu_1 \left(m_1 - \int_{-\infty}^{+\infty} x f(x) dx \right) + \dots + \mu_k \left(m_k - \int_{-\infty}^{+\infty} x^k f(x) dx \right). \end{aligned}$$

Similarly, we maximize $\text{Ent}(f)$ by setting $\frac{\partial \mathcal{L}}{\partial f}$ to zero:

$$0 = \frac{\partial \mathcal{L}}{\partial f} = \ln f + 1 - \lambda - \mu_1 x - \dots - \mu_k x^k \Rightarrow f = e^{\lambda-1+\mu_1 x+\dots+\mu_k x^k}.$$

By plugging this into the constraint condition, we can solve $(\lambda, \mu_1, \dots, \mu_k)$.

Problem 2 The log-likelihood can be written as:

$$\ell(\theta|y) = yb(\theta) + c(\theta) + d(y).$$

The expectation of score function is 0:

$$0 = \mathbb{E} \left[\frac{\partial \ell(\theta|y)}{\partial \theta} \right] = \mathbb{E} [yb'(\theta) + c'(\theta)] \Rightarrow \mathbb{E} [y] = \frac{-c'(\theta)}{b'(\theta)}.$$

The property of Fisher information yields:

$$\begin{aligned} 0 &= \mathbb{E} \left[\frac{\partial^2 \ell(\theta|y)}{\partial \theta^2} \right] + \text{Var} \left[\frac{\partial \ell(\theta|y)}{\partial \theta} \right] = \mathbb{E} [yb''(\theta) + c''(\theta)] + (b'(\theta))^2 \text{Var} [y] \\ \Rightarrow \text{Var} [y] &= \frac{-b''(\theta)\mathbb{E} [y] - c''(\theta)}{(b'(\theta))^2} = \frac{\frac{c'(\theta)b''(\theta)}{b'(\theta)} - c''(\theta)}{(b'(\theta))^2} = \frac{\left(\frac{-c'(\theta)}{b'(\theta)} \right)'}{b'(\theta)}. \end{aligned}$$

(a) Denote $\mathbf{Y} = (y_1, \dots, y_n)^\top$, $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)^\top$. By the chain rule:

$$\begin{aligned} \mathbf{s}^\top &= \frac{\partial \ell(\boldsymbol{\theta}|\mathbf{Y})}{\partial \boldsymbol{\beta}^\top} = \frac{\partial \ell(\boldsymbol{\theta}|\mathbf{Y})}{\partial \boldsymbol{\theta}^\top} \cdot \frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\mu}^\top} \cdot \frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}^\top} \cdot \frac{\partial \boldsymbol{\eta}}{\partial \boldsymbol{\beta}^\top} \\ &= (y_1 b'(\theta_1) + c'(\theta_1) \quad \dots \quad y_n b'(\theta_n) + c'(\theta_n)) \begin{pmatrix} \left(\frac{-c'(\theta_1)}{b'(\theta_1)} \right)' & & \\ & \ddots & \\ & & \left(\frac{-c'(\theta_n)}{b'(\theta_n)} \right)' \end{pmatrix}^{-1} \\ &\quad \begin{pmatrix} \frac{\partial \mu_1}{\partial \eta_1} & & \\ & \ddots & \\ & & \frac{\partial \mu_n}{\partial \eta_n} \end{pmatrix} \mathbf{x} \\ &= \left(\frac{y_1 - \frac{-c'(\theta_1)}{b'(\theta_1)'}}{\left(\frac{-c'(\theta_1)}{b'(\theta_1)} \right)'}, \frac{\partial \mu_1}{\partial \eta_1}, \dots, \frac{y_n - \frac{-c'(\theta_n)}{b'(\theta_n)'}}{\left(\frac{-c'(\theta_n)}{b'(\theta_n)} \right)'}, \frac{\partial \mu_n}{\partial \eta_n} \right) \mathbf{x} \\ &= \left(\sum_{i=1}^n \frac{y_i - \mu_i}{\text{Var}[Y_i]} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot x_{i1}, \dots, \sum_{i=1}^n \frac{y_i - \mu_i}{\text{Var}[Y_i]} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot x_{in} \right) \end{aligned}$$

Which is the desired result.

(b) From (1)

$$\begin{aligned} \mathcal{I}_{jk} &= \mathbb{E} [s_j s_k] = \mathbb{E} \left[\left(\sum_{i=1}^n \frac{y_i - \mu_i}{\text{Var}[Y_i]} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot x_{ij} \right) \left(\sum_{i=1}^n \frac{y_i - \mu_i}{\text{Var}[Y_i]} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot x_{ik} \right) \right] \\ &= \sum_{1 \leq i, l \leq n} \mathbb{E} \left[\frac{(y_i - \mu_i)(y_l - \mu_l)}{\text{Var}[Y_i] \text{Var}[Y_l]} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \mu_l}{\partial \eta_l} \cdot x_{ij} x_{lk} \right] \\ &= \sum_{i=1}^n \mathbb{E} \left[\frac{(y_i - \mu_i)(y_i - \mu_i)}{\text{Var}[Y_i] \text{Var}[Y_i]} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot x_{ij} x_{ik} \right] \quad (\text{since } y_i \text{ s are i.i.d. sample}) \\ &= \sum_{i=1}^n \frac{x_{ij} x_{ik}}{\text{Var}[Y_i]} \cdot \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2. \end{aligned}$$

Problem 3 (a) We generate observations of Y following

$$\begin{aligned} p &= \frac{\exp(X^T \beta_0)}{1 + \exp(X^T \beta_0)}; \\ Y &\sim \text{Bernoulli}(p) \end{aligned}$$

The generated observations of Y are

```
array([0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 1, 1, 1,
       1, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0,
       0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1,
       , 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1,
       0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0,
       0, 0])
```

(b) IRLS algorithm for logistic regression can be found in lecture 2, thus you can use the update rules on the slides directly. The IRLS update rule is

$$\beta^{(k+1)} = \beta^{(k)} + \left(X^T W^{(k)} X\right)^{-1} \left(X^T (Y - p^{(k)})\right)$$

where

$$p^{(k)} = \frac{1}{1 + \exp(-X^T \beta^{(k)})}$$

$$W^{(k)} = \text{diag}\{p_1^{(k)}(1 - p_1^{(k)}), p_2^{(k)}(1 - p_2^{(k)}), \dots, p_n^{(k)}(1 - p_n^{(k)})\}$$

The MLE found by IRLS algorithm may differ with different initial value and stopping criterion. An example is $\hat{\beta} = (1.37086595; 0.66987777)^T$.

(c) The theoretical distribution of β is $N(\beta_0, \Sigma)$ where

$$\beta_0 = (-2, 1)^T, \quad \Sigma = (X^T W X)^{-1} = \begin{pmatrix} 0.19350715 & -0.0440933 \\ -0.0440933 & 0.10205911 \end{pmatrix}$$

The empirical mean and variance of 100 estimations of β is:

$$\hat{\mu} = \frac{1}{100} \sum_{i=1}^{100} \hat{\beta}^{(i)} = (2.19250614, 1.07360059)^T$$

$$\hat{\Sigma} = \frac{1}{99} \sum_{i=1}^{100} (\hat{\beta}^{(i)} - \hat{\mu})(\hat{\beta}^{(i)} - \hat{\mu})^T = \begin{pmatrix} 0.24955422 & -0.040046 \\ -0.040046 & 0.12947053 \end{pmatrix}$$

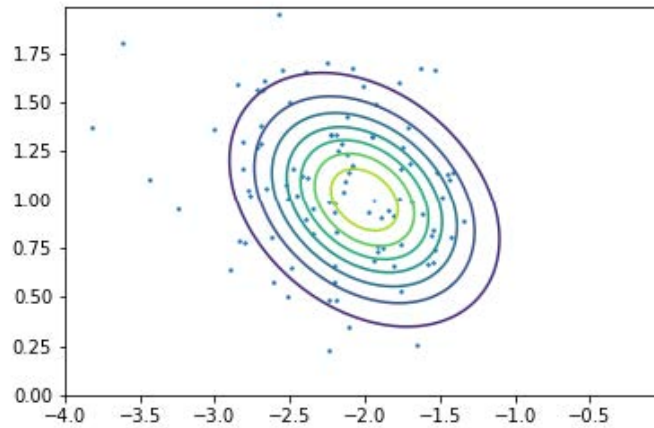


Figure 1: Contour plot for $n = 100$

(d) We repeat the procedure in (c) for $n = 10000$ and obtain the theoretical mean and variance

$$\beta_0 = (-2, 1)^T, \quad \Sigma = (X^T W X)^{-1} = \begin{pmatrix} 0.00174174 & -0.00053308 \\ -0.00053308 & 0.00096404 \end{pmatrix}$$

and the empirical mean and variance

$$\hat{\mu} = \frac{1}{100} \sum_{i=1}^{100} \hat{\beta}^{(i)} = (-2.00174993, 0.99698396)^T$$

$$\hat{\Sigma} = \frac{1}{99} \sum_{i=1}^{100} (\hat{\beta}^{(i)} - \hat{\mu})(\hat{\beta}^{(i)} - \hat{\mu})^T = \begin{pmatrix} 0.00186519 & -0.00062197 \\ -0.00062197 & 0.00114749 \end{pmatrix}$$

(All these results may differ with different settings.)

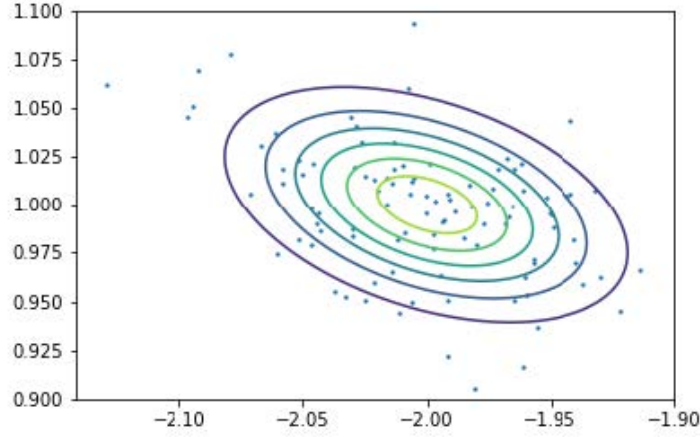


Figure 2: Contour plot for $n = 10000$

Problem 4 (a) It's obvious that $f(x, y) \geq 0$ and $f(x, y) = 0 \Rightarrow 1.5 - x + xy = 2.25 - x + xy^2 = 2.625 - x + xy^3 = 0 \Rightarrow (x, y) = (3, 0.5)$. Hence $(x^*, y^*) = (3, 0.5)$ is the unique point who reaches the global minimum.

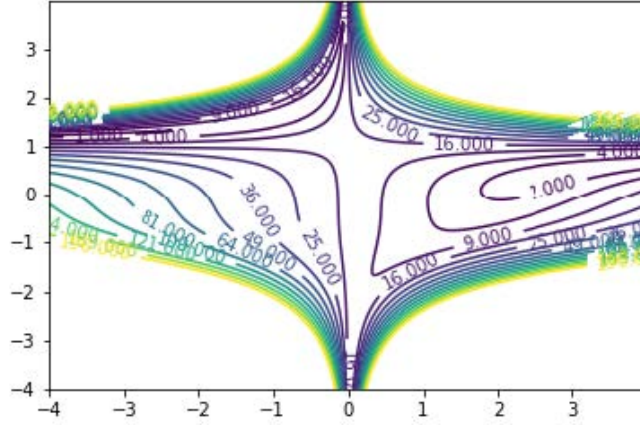


Figure 3: Contour plot for f

(b) Let $z = (x, y)^T$. The gradient of Beal function is

$$\begin{aligned} \frac{\partial f}{\partial x} &= (2y - 2)(1.5 - x + xy) + (2y^2 - 2)(2.25 - x + xy^2) + (2y^3 - 2)(2.625 - x + xy^3) \\ \frac{\partial f}{\partial y} &= 2x(1.5 - x + xy) + 4xy(2.25 - x + xy^2) + 6xy^2(2.625 - x + xy^3) \end{aligned}$$

For different optimization algorithms, the update schemes are as follows.

- gradient descent:

$$z^{(k+1)} = z^{(k)} - \alpha \nabla f(z^{(k)}).$$

- gradient descent with momentum:

$$\begin{aligned} m^{(k)} &= \mu m^{(k-1)} + (1 - \mu) \nabla f(z^{(k)}) \\ z^{(k+1)} &= z^{(k)} - \alpha m^{(k)} \end{aligned}$$

- Nesterov's acceleration:

$$\begin{aligned} \eta &= z^{(k)} + \frac{k-1}{k+2} (z^{(k)} - z^{(k-1)}) \\ z^{(k+1)} &= \eta - \alpha \nabla f(\eta) \end{aligned}$$

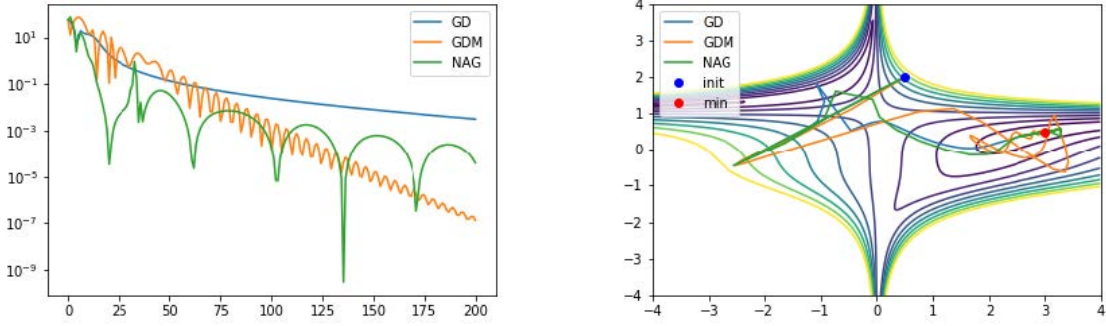


Figure 4: Stepsize = 0.0213 for gradient descent, Stepsize = 0.0812, $\mu = 0.9$ for gradient descent with momentum, Stepsize = 0.0271 for Nesterov's accelerated gradient descent.

(c) We can add random noise $N(0, 0.01)$ to obtain the stochastic gradient

$$g(z) = \nabla f(z) + \eta, \quad \eta \sim N(0, 1).$$

For vanilla SGD and its variants, the update schemes are:

- vanilla SGD:

$$z^{(k+1)} = z^{(k)} - \alpha g(z^{(k)}).$$

- AdaGrad:

$$\begin{aligned} G^{(k)} &= G^{(k-1)} + \nabla g(z^{(k)}) \odot g(z^{(k)}) \\ z^{(k+1)} &= z^{(k)} - \frac{\alpha}{\sqrt{G^{(k)} + \varepsilon}} \odot g(z^{(k)}). \end{aligned}$$

- RMSprop:

$$\begin{aligned} \mathbb{E}(g^2)_k &= \rho \mathbb{E}(g^2)_{k-1} + \rho g(z^{(k)}) \odot g(z^{(k)}) \\ z^{(k+1)} &= z^{(k)} - \frac{\alpha}{\sqrt{\mathbb{E}(g^2)_k + \varepsilon}} \odot g(z^{(k)}). \end{aligned}$$

- Adam:

$$\begin{aligned}
m^{(k)} &= \beta_1 m^{(k-1)} + (1 - \beta_1) g(z^{(k)}) \\
v^{(k)} &= \beta_2 v^{(k-1)} + (1 - \beta_2) g(z^{(k)}) \odot g(z^{(k)}) \\
\hat{m}^{(k)} &= \frac{m^{(k)}}{1 - \beta_1^k} \\
\hat{v}^{(k)} &= \frac{v^{(k)}}{1 - \beta_2^k} \\
z^{(k+1)} &= z^{(k)} - \frac{\alpha}{\sqrt{\hat{v}^{(k)} + \varepsilon}} \odot \hat{m}^{(k)}.
\end{aligned}$$

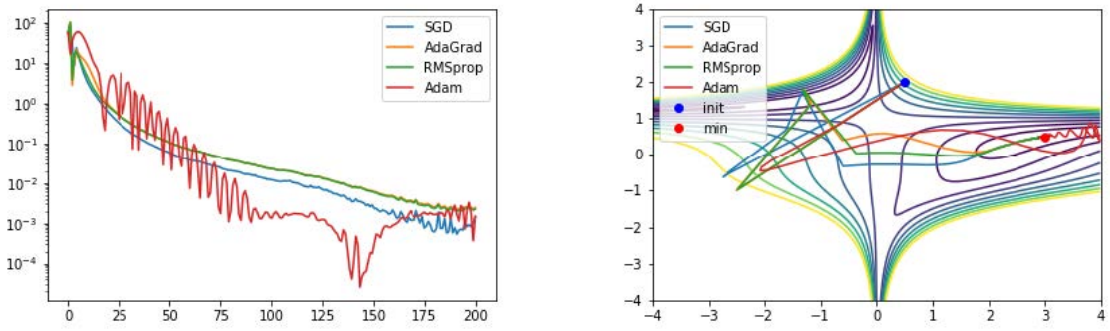


Figure 5: Stepsize = 0.029 for SGD, Stepsize = 2.99 for AdaGrad, Stepsize = 0.094, $\rho = 0.999$ for RMSprop, Stepsize = 1.152, $\beta_1 = 0.9$, $\beta_2 = 0.999$ for Adam.