

Homework 3 Bayesian Machine Learning

Problem 1. (50 points)

We have a data set of the form $\{(x_i, y_i)\}_{i=1}^N$, where $y \in \mathbb{R}$ and $x \in \mathbb{R}^d$. We assume d is large and not all dimensions of x are informative in predicting y . Consider the following regression model for this problem:

$$y_i \stackrel{iid}{\sim} \text{Normal}(x_i^T w, \lambda^{-1}), \quad w \sim \text{Normal}(0, \text{diag}(\alpha_1, \dots, \alpha_d)^{-1}),$$

$$\alpha_k \stackrel{iid}{\sim} \text{Gamma}(a_0, b_0), \quad \lambda \sim \text{Gamma}(e_0, f_0).$$

Use the density function $\text{Gamma}(\eta|\tau_1, \tau_2) = \frac{\tau_2^{\tau_1}}{\Gamma(\tau_1)} \eta^{\tau_1-1} e^{-\tau_2 \eta}$. In this homework, you will derive a variational inference algorithm for approximating the posterior distribution with

$$q(w, \alpha_1, \dots, \alpha_d, \lambda) \approx p(w, \alpha_1, \dots, \alpha_d, \lambda|y, x)$$

- a) Using the factorization $q(w, \alpha_1, \dots, \alpha_d, \lambda) = q(w)q(\lambda) \prod_{k=1}^d q(\alpha_k)$, derive the optimal form of each q distribution. Use these optimal q distributions to derive a variational inference algorithm for approximating the posterior.
- b) Summarize the algorithm derived in Part (a) using pseudo-code in a way similar to how algorithms are presented in the notes for the class.
- c) Using these q distributions, calculate the variational objective function. You will need to evaluate this function in the next problem to show the convergence of your algorithm.

ZOLTAN ONODI-SZUCS
202131

Data: $\{(x_i, y_i)\}_{i=1}^N$, $y_i \in \mathbb{R}$, $x_i \in \mathbb{R}^d$

$y_i \sim \mathcal{N}(x_i^T w, \lambda^{-1})$, $w \sim \mathcal{N}(0, \text{diag}(\lambda_1, \dots, \lambda_d)^{-1})$

$\lambda_k \sim \text{Gamma}(a_0, b_0)$, $\lambda \sim \text{Gamma}(e_0, f_0)$

$$a) \quad q(w, \lambda, \lambda_1, \dots, \lambda_d) = q(w) q(\lambda) \prod_{k=1}^d q(\lambda_k)$$

$$\mathcal{L}(w, \lambda, \lambda_1, \dots, \lambda_d) = \int \dots \int q(w, \lambda, \lambda_1, \dots, \lambda_d) \cdot \ln \frac{P(y, w, \lambda, \lambda_1, \dots, \lambda_d | x)}{q(w, \lambda, \lambda_1, \dots, \lambda_d)} dw d\lambda d\lambda_1 \dots d\lambda_d$$

$$\begin{aligned} P(y, w, \lambda, \lambda_1, \dots, \lambda_d) &= P(y | \lambda, w, \lambda, x) \cdot P(\lambda, w, \lambda | x) \\ &= P(y | w, \lambda, x) \cdot P(w | \lambda, \lambda, x) \cdot P(\lambda, \lambda | x) \\ &= P(y | w, \lambda, x) \cdot P(w | \lambda) \cdot P(\lambda) \cdot P(\lambda) \end{aligned}$$

$$P(y | w, \lambda, x) = \prod_{i=1}^N P(y_i | w, \lambda, x_i)$$

$$P(\lambda) = \prod_{k=1}^d P(\lambda_k)$$

EMSC - KMOO VARIATION
VARIATION

$$q(\lambda) \propto \exp \left\{ E_{q(\theta+\lambda)} [\ln P(y, \lambda, w, \lambda | x)] \right\}$$

$$\propto \exp \left\{ E_{q(\theta+\lambda)} [\ln P(y | w, \lambda, x) + \ln P(\lambda)] \right\}$$

$$\propto \exp \left\{ E_{q(w)} [\ln P(y | w, \lambda, x)] + \ln P(\lambda) \right\}$$

$$\propto \prod_{i=1}^N \lambda^{1/2} \exp \left\{ -\frac{1}{2} E_{q(w)} [(y_i - x_i^T w)^2] \right\} \cdot \lambda^{e_0-1} \cdot e^{-f_0 \lambda}$$

$$\propto \lambda^{N/2 + e_0 - 1} \cdot e^{-\lambda (f_0 + \frac{1}{2} \sum_{i=1}^N E_{q(w)} [(y_i - x_i^T w)^2])}$$

$$\sim \text{Gamma}(e', f') \quad e' = \frac{N}{2} + e_0$$

$$f' = f_0 + \frac{1}{2} \sum_{i=1}^N E_{q(w)} [(y_i - x_i^T w)^2]$$

$$q(w) \propto \exp \left\{ E_{q(\theta+w)} [\ln P(y, w, \lambda, \lambda | x)] \right\}$$

$$\propto \exp \left\{ E_{q(\theta+w)} [\ln P(y | w, \lambda, x) + \ln P(w | \lambda)] \right\}$$

$$\propto \exp \left\{ E_{q(\lambda)} [\ln P(y | w, \lambda, x)] + E_{q(w)} [\ln P(w | \lambda)] \right\}$$

$$\propto \prod_{i=1}^N \exp \left\{ -\frac{1}{2} E_{q(\lambda)} [\lambda] (y_i - x_i^T w)^2 \right\} \cdot \exp \left\{ -\frac{1}{2} E_{q(\lambda)} [\lambda] w^T w \right\}$$

$$\cdot \exp \left\{ -\frac{1}{2} E_{q(\lambda)} [w^T \cdot \text{diag}(\lambda) \cdot w] \right\}$$

$$\sim \mathcal{N}(\mu', \Sigma')$$

$$\Sigma' = (E_{q(\lambda)} [\lambda] \sum_{i=1}^N x_i x_i^T + E_{q(\lambda)} [\text{diag}(\lambda)])^{-1}$$

$$\mu' = \Sigma' E_{q(\lambda)} [\lambda] \sum_{i=1}^N y_i x_i$$

$$q(L_k) \propto E_{q(\theta+L_k)} [\ln P(Y, w, \lambda, L | X)]$$

$$\propto \exp \{ E_{q(\theta+L_k)} [\ln P(w | L) + \ln P(L_k)] \}$$

$$\propto \exp \{ E_{q(\theta+L_k)} \left[\sum_{k=1}^K \ln L_k \cdot e^{-\frac{1}{2} a_k^2 L_k} \right] + E_{q(L_k)} \left[\ln L_k \cdot e^{-\frac{1}{2} w_k^2 L_k} \right] +$$

$$+ \ln L_k \cdot e^{-\frac{1}{2} b_0^2 L_k} \}$$

$$\propto L_k^{a_0/2 - 1} \cdot e^{-\frac{1}{2} (b_0^2 + \frac{1}{2} E_{q(w)} [w_k^2]) L_k}$$

$$\sim \text{Gamma}(a_k', b_k')$$

$$a_k' = a_0 + \frac{1}{2}$$

$$b_k' = b_0 + \frac{1}{2} E_{q(w)} [w_k^2]$$

Expectations:

$$E_{q(\lambda)} [\lambda] = \frac{c}{4}$$

$$E_{q(w)} [(y_i - x_i^T w)^2] = (y_i - x_i^T \mu')^2 + x_i^T \Sigma' x_i$$

$$E_{q(w)} [w_k^2] = \Sigma'_{kk} + \mu_k'^2$$

$$E_{q(L)} [\text{diag}(L_k)] = \text{diag}\left(\frac{a_k'}{b_k'}\right)$$

Algorithm:

1) Initialize $a_0, b_0, e_0, f_0, \mu_0, \Sigma_0$ in some way

2) For iterations $t \dots T$:

a) Update $q(w) = \mathcal{N}(\mu_t, \Sigma_t)$ by setting

$$\Sigma_t = \left(\frac{e_{t-1}}{f_{t-1}} \sum_{i=1}^N x_i x_i^T + \text{diag} \left(\frac{a_{t-1}}{f_{t-1}} \right) \right)^{-1}, \quad \mu_t = \frac{1}{f_t} \sum_{i=1}^N y_i x_i$$

b) Update $q(\lambda) = \text{gamma}(e', f')$ by setting

$$e_t = \frac{N}{2} + e_0, \quad f_t = f_0 + \frac{1}{2} \sum_{i=1}^N \left((y_i - x_i^T \mu_t)^2 + x_i^T \Sigma_t x_i \right)$$

c) Update $q(a, b) = \text{gamma}(a', b')$ by setting

$$a_t = a_0 + \frac{1}{2}, \quad b_{b(t)} = b_0 + \frac{1}{2} \left(\Sigma_t + \mu_t^2 \right)$$

d) Evaluate the variational objective function \mathcal{L} to assess convergence

$$\mathcal{L} = \int \int q(w, \lambda, \lambda_1, \lambda_2) \ln \frac{p(y, w, \lambda, \lambda_1, \lambda_2 | x)}{q(w, \lambda, \lambda_1, \lambda_2)} dw d\lambda d\lambda_1 d\lambda_2$$

$$= E_q[\ln p(y, w, \lambda, \lambda_1, \lambda_2 | x)] - E_q[\ln q(w, \lambda, \lambda_1, \lambda_2)]$$

$$= E_q[\ln p(y | w, \lambda, x) + \ln p(w | \lambda) + \ln p(\lambda) + \ln p(\lambda_1)] - \\ - E_q[\ln q(w)] - E_q[\ln q(\lambda)] - E_q[\ln q(\lambda_1)]$$

$$= \frac{N}{2} E_{q(\lambda)}[\ln(\lambda)] - \frac{E_q[\lambda]}{2} \sum_{i=1}^N E_{q(w)}[(y_i - x_i^T w)^2] + \\ + \frac{1}{2} \sum_{k=1}^d E_{q(\lambda_k)}[\ln(\lambda_k)] - \frac{1}{2} E_{q(w)}[w^T \text{diag}(E_{q(\lambda_k)}[\lambda_k]) w] + \\ + (a_0 - 1) \sum_{k=1}^d E_{q(\lambda_k)}[\ln \lambda_k] - b_0 \sum_{k=1}^d E[\lambda_k] + \\ + (c_0 - 1) E_{q(\lambda)}[\ln(\lambda)] - f_0 E_{q(\lambda)}[\lambda] + \\ + \frac{1}{2} \ln(|\Sigma'|)$$

$$+ (e' - \ln f' + \ln(\Gamma(e')) + (1 - e') \Psi(e')) +$$

$$+ \sum_{k=1}^d [a'_k - \ln b'_k + \ln(\Gamma(a'_k)) + (1 - a'_k) \Psi(a'_k)]$$

$$+ \text{const.}$$

$$\begin{aligned}
L = & \frac{N}{2} (\Psi(e') - \ln(f')) - \frac{e'}{2f'} \sum_{i=1}^N (y_i - x_i^T \mu')^2 + x_i^T \Sigma^{-1} x_i \\
& + \frac{1}{2} \sum_{k=1}^d (\Psi(a'_k) - \ln(b'_k)) - \frac{1}{2} \sum_{k=1}^d (\Sigma_{kk}' + \mu_{kk}'^2) \frac{a'_k}{b'_k} \\
& + (a_0 - 1) \sum_{k=1}^d (\Psi(a'_k) - \ln(b'_k)) - b_0 \sum_{k=1}^d \frac{a'_k}{b'_k} \\
& + (e_0 - 1) (\Psi(e') - \ln(f')) - f_0 \frac{e'}{f'} \\
& + \frac{1}{2} \ln(|\Sigma'|) \\
& + (e' - \ln f' + \ln \Gamma(e') + (1 - e') \Psi(e')) \\
& + \sum_{k=1}^d [a'_k - \ln b'_k + \ln \Gamma(a'_k) + (1 - a'_k) \Psi(a'_k)]
\end{aligned}$$

Problem 2. (50 points)

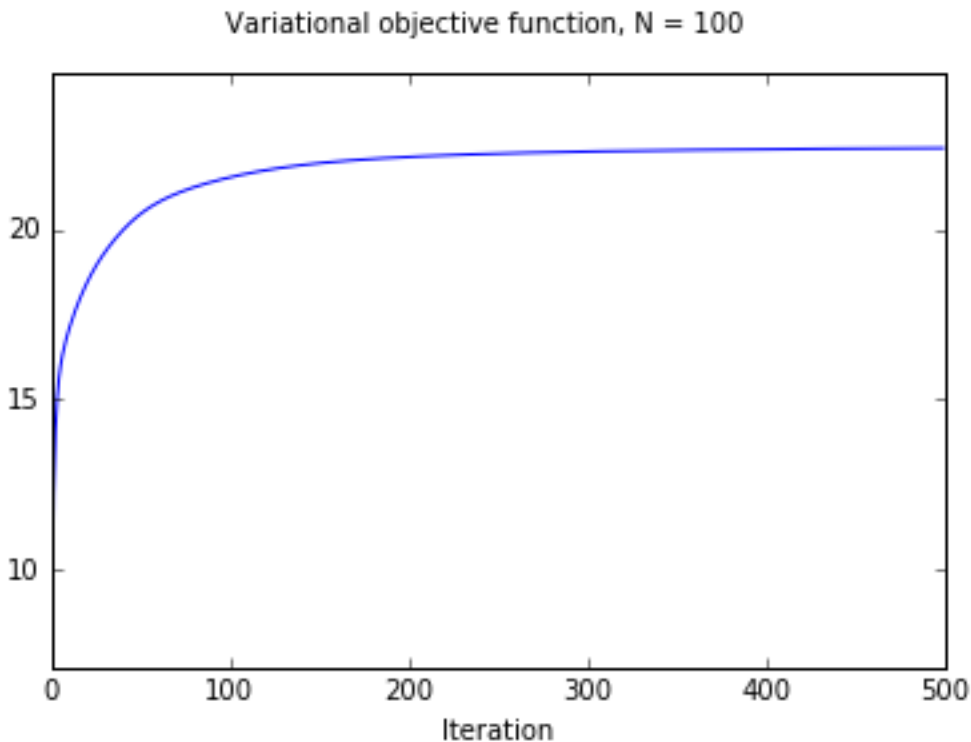
Implement the algorithm derived in Problem 1 and run it on the three data sets provided. Set the prior parameters $a_0 = b_0 = 10^{-16}$ and $e_0 = f_0 = 1$. We will not discuss sparsity-promoting “ARD” priors in detail in this course, but setting a_0 and b_0 in this way will encourage only a few dimensions of w to be significantly non-zero since many α_k should be extremely large according to $q(\alpha_k)$.

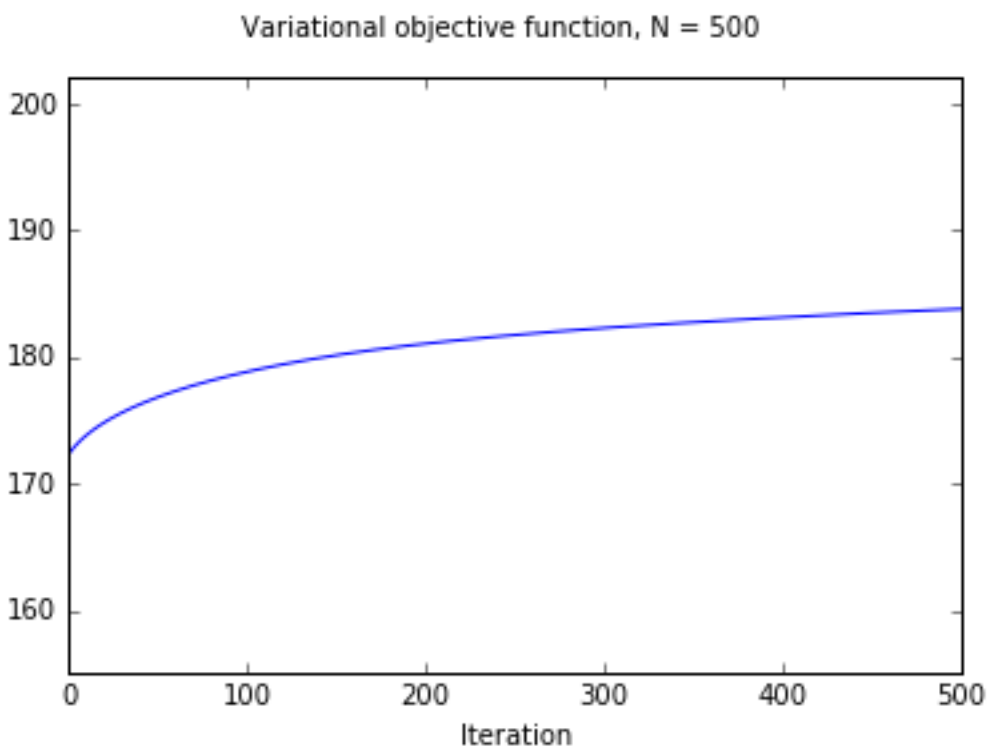
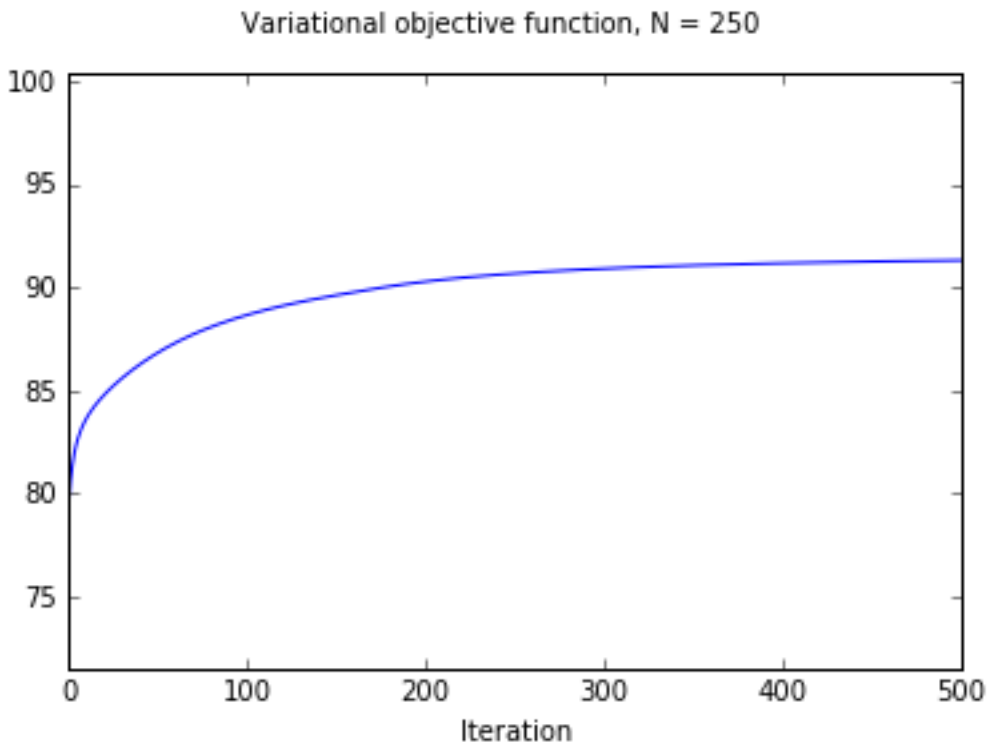
For each of the three data sets provided, show the following:

- Run your algorithm for 500 iterations and plot the variational objective function.
- Using the final iteration, plot $1/\mathbb{E}_q[\alpha_k]$ as a function of k .
- Give the value of $1/\mathbb{E}_q[\lambda]$ for the final iteration.
- Using $\hat{w} = \mathbb{E}_{q(w)}[w]$, calculate $\hat{y}_i = x_i^T \hat{w}$ for each data point. Using the z_i associated with y_i (see below), plot \hat{y}_i vs z_i as a solid line. On the same plot show (z_i, y_i) as a scatter plot. Also show the function $(z_i, 10 * \text{sinc}(z_i))$ as a solid line in a different color.

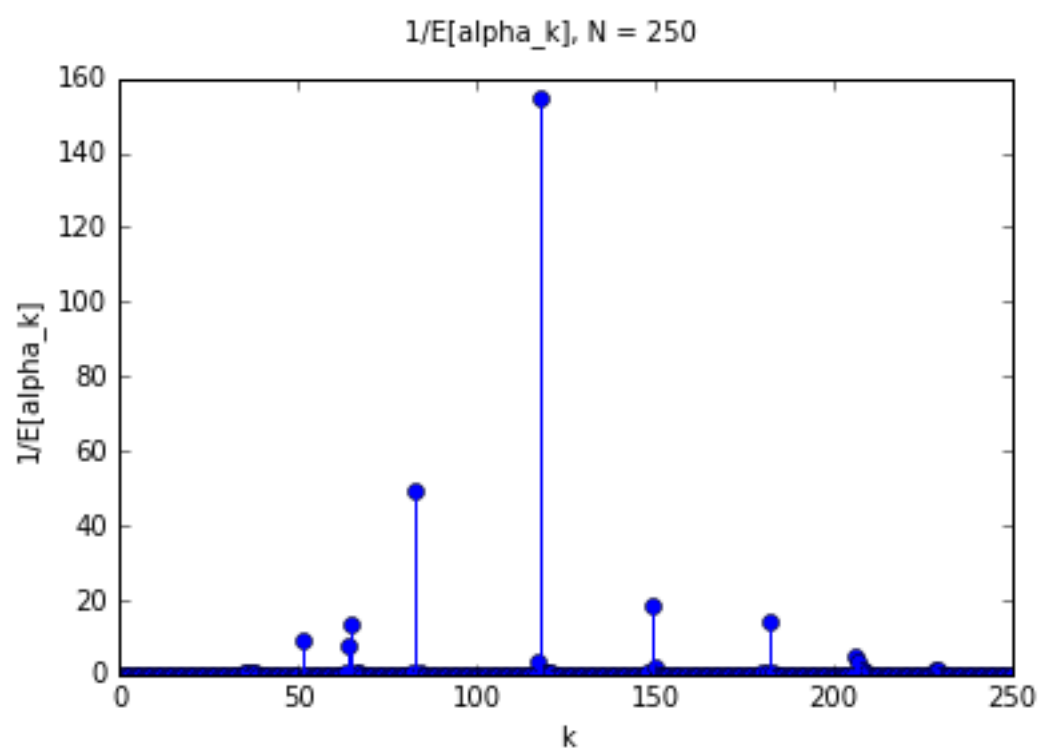
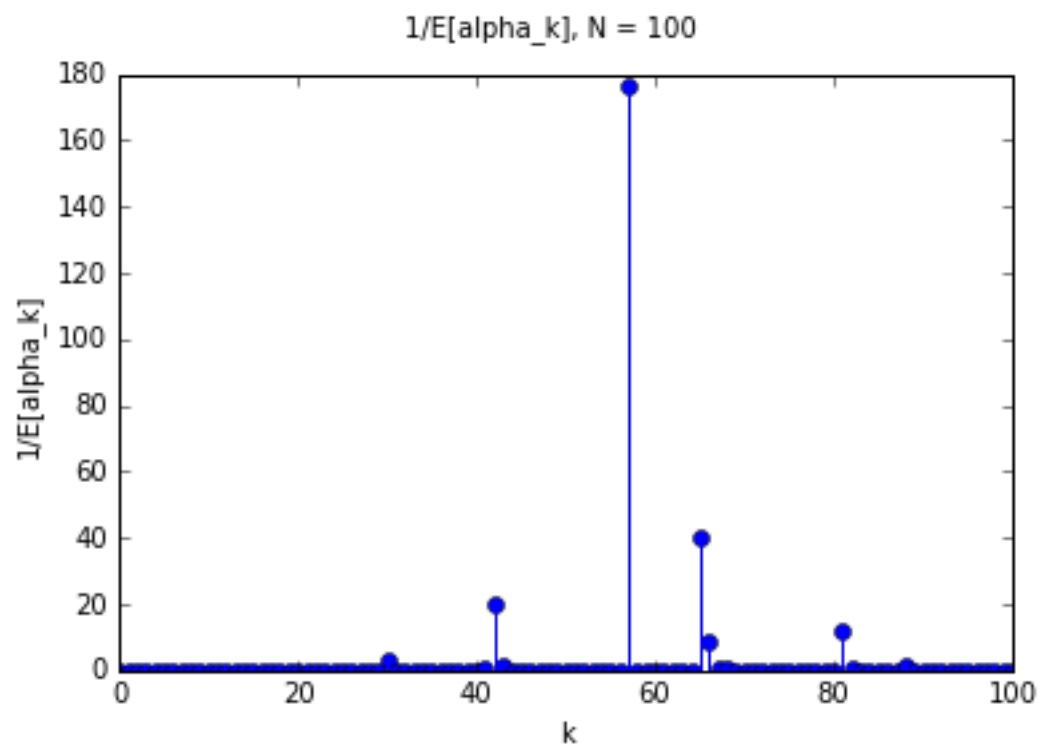
Hint about Part (d): z is the horizontal axis and y the vertical axis. Both solid lines should look like a function that smoothly passes through the data. The second line is ground truth.

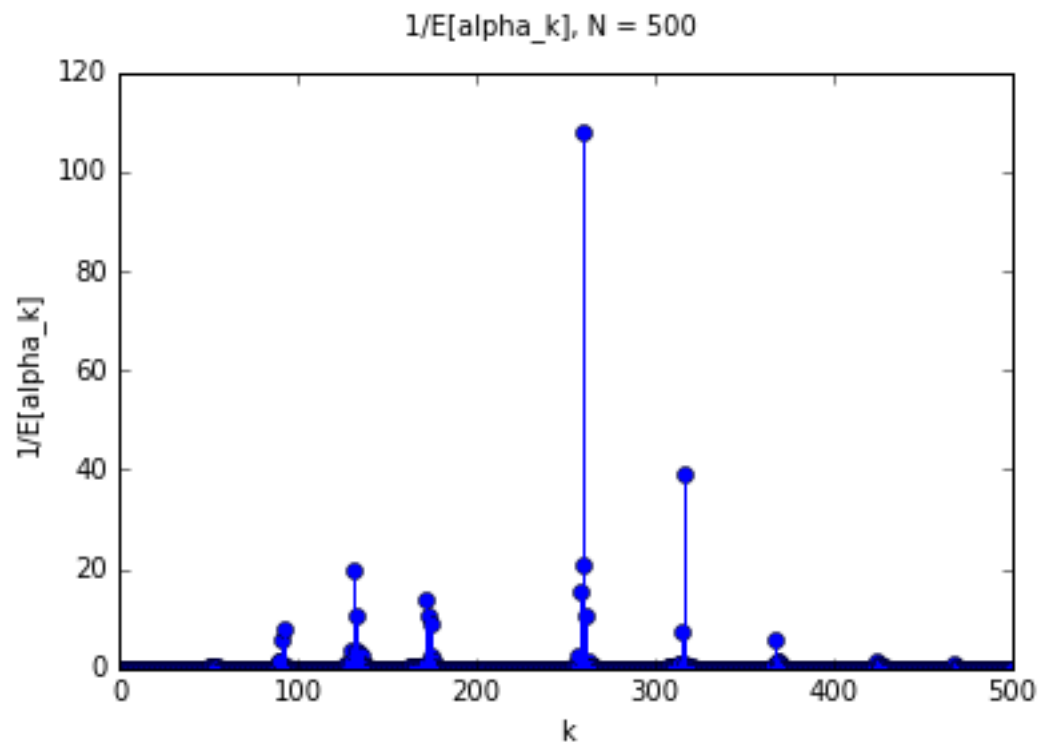
a, Note that the variational objective functions are positive, because I drop the constant terms. If I include them, it's the same shape but a negative number.





b,





c,
 $N = 100:$ $1/E[\lambda] = 1.07981360709$
 $N = 250:$ $1/E[\lambda] = 0.899449415713$
 $N = 500:$ $1/E[\lambda] = 0.978100451976$
d,

