Problem I)

a.) $\{(x_i, y_i)\}_{i=1}^{N}$ ; $y_i \in \mathbb{R}$ ; $x_i \in \mathbb{R}^d$ $d \geq N$

$y_i \overset{iid}{\sim} N(x_i^T w, \lambda^{-1})$ $w \sim N(0, \text{diag}(\alpha_1, \ldots \alpha_d)^{-1})$

$\alpha_k \sim \Gamma(a_0, b_0)$ $\lambda \sim \Gamma(e_0, f_0)$

Derive optimal $q(\cdot)$ of each distribution using

$$q(w, \alpha_1, \ldots, \alpha_d, \lambda) \approx p(w, \alpha_1, \ldots, \alpha_d, \lambda | y, x) = q(w) q(\lambda) \prod_{k=1}^{d} q(\alpha_k)$$

the optimal $q(\lambda)$

$q(\lambda) \propto \exp\left[ \underset{-q}{\mathbb{E}}\left[ \sum_{i=1}^{N} \ln p(y_i | w, \lambda, x) \right]\right] p(\lambda)$

$\propto \exp\left[ \underset{-q}{\mathbb{E}}\left[ \sum_{i=1}^{N} \frac{1}{2} \ln 2\pi + \frac{1}{2} \ln \lambda - \frac{\lambda}{2}(y_i - x_i^T w)^2 \right]\right] \times \frac{f_0}{\Gamma(e_0)} \lambda^{e_0 - 1} e^{-f_0 \lambda}$

$\propto \left(\frac{\lambda}{2\pi}\right)^{\frac{N}{2}} \exp\left[ -\frac{\lambda}{2} \mathbb{E}\left[ (y_i^2 - 2 y_i x_i^T w + w^T x_i x_i^T w) \right]\right] \frac{f_0}{\Gamma(e_0)} \lambda^{e_0 - 1} e^{-f_0 \lambda}$

$\propto \lambda^{N/2 + e_0 - 1} \exp\left[ -\frac{\lambda}{2}\left( \sum_{i=1}^{N} y_i^2 - 2 y_i x_i^T \mathbb{E}[w] + \mathbb{E}[\text{tr}(x_i x_i^T w w^T)]\right) - f_0 \lambda \right]$

Note that: $\mathbb{E}[w] = \mu$ $\mathbb{E}[w^T w] = \Sigma + \mu \mu^T$

$\text{tr}(x_i x_i^T \mathbb{E}[w w^T]) = \text{tr}(x_i x_i^T (\Sigma - \mu \mu^T))$

$q(\lambda) \propto \exp\left[ -\lambda \left( f_0 + \frac{1}{2}\left( \sum_{i=1}^{N} x_i^T \Sigma x_i + y_i^2 - 2 y_i x_i^T \mu + \mu^T x_i x_i^T \mu \right) \right)\right] \lambda^{e_0 + \frac{N}{2} - 1}$

$\propto \lambda^{e_0 + \frac{N}{2} - 1} \exp\left[ -\lambda\left[ f_0 + \frac{1}{2}\left( \sum_{i=1}^{N}(y_i - x_i^T \mu)^2 + x_i^T \Sigma x_i \right)\right]\right]$

Which is Gamma. Thus, $q(\lambda) \sim \Gamma(e_1, f_1)$

where $e_1 = e_0 + N/2$

$f_1 = f_0 + \frac{1}{2}\left( \underbrace{\sum_{i=1}^{N}(y_i - x_i^T \mu)^2 + x_i^T \Sigma x_i}_{\mathbb{E}[y_i - x_i^T w]} \right)$

and

$\underset{q(\lambda)}{\mathbb{E}}[\lambda] = \frac{e_1}{f_1}$

The optimal $q(\alpha_k)$

$$q(\alpha_K) \propto \exp\left[\mathbb{E}_{\underset{q}{}}\left[\underbrace{\sum_{i=1}^{N} \ln p(y_i|w, X, \lambda)}_{\text{can be ignored}} + \ln p(\lambda) + \ln p(w|\alpha_1,..,\alpha_K) + \sum_{k=1}^{d} \ln p(\alpha_1,..,\alpha_K)\right]\right]$$

$$\propto \alpha_k^{a_0-1} e^{-\alpha_k b_0} \times \exp\left[\mathbb{E}_{\underset{q}{}}\left[\frac{1}{2}\sum_{k=1}^{d} \ln \alpha_k - \frac{1}{2} w^T \text{diag}(\alpha_1,..,\alpha_d) w\right]\right]$$

$$\propto \alpha_k^{a_0-1} e^{-\alpha_k b_0} \exp\left[\mathbb{E}\left[\frac{1}{2}\ln \alpha_k - \frac{1}{2} w^T \text{diag}(\alpha_1,..,\alpha_d) w\right]\right]$$

$$\propto \alpha_k^{a_0+\frac{1}{2}-1} e^{-\alpha_k b_0} \exp\left[-\frac{1}{2}\text{trace}\left(\text{diag}(\alpha_1,..,\alpha_d) \mathbb{E}[ww^T]\right)\right]$$

$$\propto \alpha_k^{a_0-\frac{1}{2}} e^{-\alpha_k b_0} \exp\left[-\frac{1}{2}\text{trace}\left(\text{diag}(\alpha_1,..,\alpha_d)(\Sigma + \mu\mu^T)\right)\right]$$

$$\propto \alpha_k^{a_0-\frac{1}{2}} \exp\left[-\alpha_k b_0 - \frac{1}{2}\left(\alpha_k\left[\Sigma_{(k,k)} + \mu_k\mu_k^T\right]\right)\right]$$

Where $\Sigma_{(k,k)}$ denotes the $k^{th}$ row of $k^{th}$ column of $\Sigma$ and the $k^{th}$ element of $\mu$.

This is obviously a Gamma and so, $q(\alpha_k) \sim \Gamma(a_1, b_1)$
where $\quad a_1 = a_0 + \frac{1}{2}$
$$b_1 = b_0 + \frac{1}{2}\left(\Sigma_{(k,k)} + \mu_k\mu_k^T\right)$$
and $\mathbb{E}_{q(\alpha_k)}[\alpha_k] = \frac{a_1}{b_1}$

the optimal $q(w)$

$$q(w) \propto \exp\left[\mathbb{E}_{\underset{q}{}}\left[\sum_{i=1}^{N} \ln p(y_i|w, x_i, \lambda) + \ln p(\lambda) + \ln p(w|\alpha_1,..,\alpha_d) + \underbrace{\sum_{k=1}^{d} \ln p(\alpha_k)}_{\text{Don't depend on } w}\right]\right]$$

$$q(w) \propto \exp\left[\mathbb{E}_{\underset{q}{}}\left[\sum_{i=1}^{N}\left(\frac{1}{2}\ln\left(\frac{\lambda}{2\pi}\right) - \frac{\lambda}{2}(y_i - x_i^T w)^2\right) - \frac{d}{2}\ln 2\pi + \frac{1}{2}\sum_{k=1}^{d}\ln\alpha_k - \frac{1}{2}w^T\text{diag}(\alpha_1,..,\alpha_d)w\right]\right]$$

$$\propto \exp\left[\mathbb{E}_{\underset{q}{}}\left[-\frac{\lambda}{2}\sum_{i=1}^{N}(y_i - x_i^T w)^2 - \frac{1}{2}w^T\text{diag}(\alpha_1,..,\alpha_d)w\right]\right]$$

$$\propto \exp\left[-\frac{1}{2}\mathbb{E}\left[\lambda\sum_{i=1}^{N}(y_i^2 - 2y_i x_i^T w + w^T x_i x_i^T w) - w^T\text{diag}(\alpha_1,..,\alpha_d)w\right]\right]$$

$$\propto \exp\left[-\frac{1}{2}\mathbb{E}\left[\lambda\sum_{i=1}^{N}x_i x_i^T + \text{diag}(\alpha_1,..,\alpha_d))ww^T - 2\lambda\sum_{i=1}^{N}y_i x_i^T w\right]\right]$$

$$\propto \exp\left[-\frac{1}{2}\mathbb{E}\left[(\omega - (\lambda\sum_{i=1}^{N}\chi_i\chi_i^T + diag(\alpha_1,..,\alpha_d)(\lambda\sum_{i=1}^{N}y_i\chi_i)^T\right.\right.$$
$$\times(\lambda + \sum_{i=1}^{N}\chi_i\chi_i^T + diag(\alpha_1,..,\alpha_d))\times(\omega - (\lambda\sum_{i=1}^{N}\chi_i\chi_i^T + diag(\alpha_1,..,\alpha_d))^{-1}(\lambda\sum_{i=1}^{N}y_i\chi_i)]]$$

$$\propto \exp\left[-\frac{1}{2}\left(\mathbb{E}[\omega] - [\mathbb{E}[\lambda]\sum_{i=1}^{N}\chi_i\chi_i^T + diag(\mathbb{E}[\alpha_1],..,\mathbb{E}[\alpha_d]))^{-1}(\mathbb{E}[\lambda]\sum y_i\chi_i)\right)^T\right.$$
$$\times(\mathbb{E}[\lambda]\sum_{i=1}^{N}\chi_i\chi_i^T + diag(\mathbb{E}[\alpha_1],..,\mathbb{E}[\alpha_d]))$$
$$\times(\mathbb{E}[\omega] - (\mathbb{E}[\lambda]\sum_{i=1}^{N}\chi_i\chi_i^T) + diag(\mathbb{E}[\alpha_1],..,\mathbb{E}[\alpha_d]))^{-1}(\mathbb{E}[\lambda]\sum_{i=1}^{N}y_i\chi_i)]$$

Which is   Normal Thus,

$$q(\omega) \sim N(\mu_1, \Sigma_1) \quad \text{where} \quad \mu_1 = \Sigma_1\left(\mathbb{E}[\lambda]\sum_{i=1}^{N}y_i\chi_i\right)$$

$$\Sigma_1 = \left[\mathbb{E}[\lambda]\sum_{i=1}^{N}\chi_i\chi_i^T + diag(\mathbb{E}[\alpha_1],..,\mathbb{E}[\alpha_d])\right]$$

b) VE for Bayesian Linear Regression

Input: $\{(\chi_i, y_i)\}_{i=1}^{N}$   $y \in \mathbb{R}$; $\chi \in \mathbb{R}^d$
$q(\lambda) \sim \Gamma(\lambda|e_1, f_1); q(\omega) \sim N(\omega|\mu_1, \Sigma_1); q(\alpha) \sim \Gamma(a_1, b_1)$

Output: Values for $a_1,...,a_d, b_1,..,b_d, e_1, f_1, \mu_1, \Sigma_1$

1) Initialize parameters $a_{10},..,a_{d0}, b_{10},..,b_{d0}, e_0, f_0, \mu_0, \Sigma_0$ in some way.

2) for $t = 1,..,T$
   - Update $q(\lambda)$ with $e_{1t} = e_0 + \frac{N}{2}$ & $f_{1t} = f_0 + \frac{1}{2}[\sum_{i=1}^{N}(y_i - \chi_i^T\mu_{t-1})^2 + \chi_i^T\Sigma_{t-1}\chi_i]$
   - Update $q(\alpha_k)$ with $a_{1k_t} = a_{0k_t} + \frac{1}{2}$ & $b_{1k_t} = b_{0k_t} + \frac{1}{2}[\Sigma_{t-1}(k,k) + \mu_{k_{t-1}}\mu_{k_{t-1}}^T]$
   - Update $q(\omega)$ with

$$\mu_t = \Sigma_t\left(\frac{e_{1t}}{f_{1t}}\sum_{i=1}^{N}y_i\chi_i\right) \quad \& \quad \Sigma_t = \left(\frac{e_{1t}}{f_{1t}}\sum_{i=1}^{N}\chi_i\chi_i^T + diag\left(\frac{a_{11t}}{b_{11t}},..,\frac{a_{1dt}}{b_{1dt}}\right)\right)^{-1}$$

Then evaluate $\mathcal{L}(\vec{a}_{1t}, \vec{b}_{1t}, e_1, f_1, \mu_1, \Sigma_1)$ to assess convergence ($< \varepsilon$).

c.) Calculate the Variational objective function

$$\mathcal{L}(\vec{a}_1, \vec{b}_1, e_1, f_1, \mu_1, \Sigma_1) = \int q(\lambda) \ln p(\lambda) d\lambda$$

$$+ \int \cdots \int q(w) \prod_{k=1}^{d} q(\alpha_k) \ln p(w | \alpha_1, \ldots, \alpha_d) dw \, d\alpha_1, \ldots, d\alpha_d$$

$$+ \sum_{k=1}^{d} \int q(\alpha_k) \ln p(\alpha_k) d\alpha_k$$

$$+ \sum_{i=1}^{N} \int \cdots \int \prod_{k=1}^{d} q(\alpha_k) q(\lambda) q(w) \ln p(y_i | x_i, w, \lambda) dw \, d\lambda \, d\alpha_1, \ldots, d\alpha_d$$

$$- \int q(\lambda) \ln q(\lambda) d\lambda - \int q(w) \ln q(w) dw - \left( \sum_{k=1}^{d} \int q(\alpha_k) \ln q(\alpha_k) d\alpha_k \right)$$

$$\mathcal{L}(\vec{a}_1, \vec{b}_1, e_1, f_1, \mu_1, \Sigma_1) = \mathbb{E}[\ln p(\lambda)] + \mathbb{E}[\ln p(w | \alpha_1, \ldots, \alpha_d)] + \sum_{k=1}^{d} \mathbb{E}[\ln p(\alpha_k)]$$

$$+ \sum_{i=1}^{N} \mathbb{E}[\ln p(y_i | x_i, w, \lambda)] - \mathbb{E}[\ln q(\lambda)] - \mathbb{E}[\ln q(w)]$$

$$- \left( \sum_{k=1}^{d} \mathbb{E}[\ln q(\alpha_k)] \right)$$

Since these expectations are with respect to the $q(\cdot)$ parameters, we can see that 3 terms are just Entropies, with constant terms. Namely

$$\mathbb{E}_{q(w)}[\ln q(w)] = \frac{N}{2} + \frac{N}{2} \ln 2\pi + \frac{1}{2} \ln |\Sigma_1| = \text{constants} + \frac{1}{2} \ln |\Sigma_1|$$

$$\mathbb{E}_{q(\lambda)}[\ln q(\lambda)] = e_1 - \ln f_1 + \ln \Gamma(e_1) + (1 - e_1) \psi(e_1)$$

$$\mathbb{E}_{q(\alpha_k)}[\ln q(\alpha_k)] = a_{1k} - \ln b_{1k} + \ln \Gamma(a_{1k}) + (1 - a_{1k}) \psi(a_{1k})$$

Where $\psi(x)$ is the digamma function $\left( \frac{\partial \ln \Gamma(x)}{\partial x} \right)$ and the definitions of entropy have been used for the Multivariate Normal and the Gamma distributions.

$$\mathbb{E}[\ln p(w|\alpha_1,...,\alpha_d)] = \mathbb{E}\left[-\frac{d}{2}\ln 2\pi + \frac{1}{2}\sum_{k=1}^{d}\ln\alpha_i - \frac{1}{2}w^T \text{diag}(\alpha_1,...,\alpha_d)w\right]$$

$$= -\frac{d}{2}\ln 2\pi + \frac{1}{2}\sum_{k=1}^{d}\mathbb{E}[\ln\alpha_k] - \frac{1}{2}\text{trace}\left(\text{diag}[\mathbb{E}[\alpha_1],...,\mathbb{E}[\alpha_d]]\,\mathbb{E}[ww^T]\right)$$

$$= \text{Constants} + \frac{1}{2}\sum_{k=1}^{d}[\psi(a_{1k}) - \ln b_{1k}] - \frac{1}{2}\text{trace}\left(\text{diag}\left(\frac{a_{11}}{b_{11}},...,\frac{a_{1d}}{b_{1d}}\right)\left(\Sigma_1 + \mu_1\mu_1^T\right)\right)$$

$$= \text{Constants} + \frac{1}{2}\sum_{k=1}^{d}(\psi(a_{1k}) - \ln b_{1k}) - \frac{1}{2}\left[\text{trace}\left(\text{diag}\left(\frac{a_{11}}{b_{11}},...,\frac{a_{1d}}{b_{1d}}\right)\Sigma_1\right) + \mu_1^T\left(\text{diag}\left(\frac{a_{11}}{b_{11}},...,\frac{a_{1d}}{b_{1d}}\right)\mu_1\right)\right]$$

$$\mathbb{E}_{q(\lambda)}[\ln p(\lambda)] = e_0\ln f_0 - \ln\Gamma(e_0) + (e_0-1)\mathbb{E}_{q(\lambda)}[\ln\lambda] - f_0\,\mathbb{E}_{q(\lambda)}[\lambda]$$

$$= \text{Constants} + (e_0-1)[\psi(e_1) - \ln f_1] - f_0\frac{e_1}{f_1}$$

where we reused $\mathbb{E}[\ln\lambda] = $ Entropy of a Gamma distribution.

$$\mathbb{E}[\ln p(\alpha_k)] = a_{0k}\ln b_{0k} - \ln\Gamma(a_{0k}) + (a_{0k}-1)\mathbb{E}[\ln\alpha_k] - b_{0k}\mathbb{E}[\alpha_k]$$

$$= \text{Constants} + (a_{0k}-1)[\psi(a_{1k}) - \ln b_{1k}] - b_{0k}\frac{a_{1k}}{b_{1k}}$$

$$\mathbb{E}[\ln p(y_i|x_i,w,\lambda)] = -\frac{1}{2}\ln 2\pi + \frac{\mathbb{E}[\ln\lambda]}{2} - \frac{\mathbb{E}[\lambda]}{2}\mathbb{E}\left[\sum_{i=1}^{N}(y_i - x_i^T w)^2\right]$$

$$= \text{Constants} + \frac{1}{2}(\psi(e_1) - \ln f_1) - \frac{1}{2}\times\frac{e_1}{f_1}\left(\sum_{i=1}^{N}(y_i - x_i^T\mu_1)^2 + x_i^T\Sigma_1 x_i\right)$$

Combining all of the terms yields:

$$\mathcal{L}(a_1, b_1, e_1, f_1, \mu_1, \Sigma_1) = \mathbb{E}[\ln p(\lambda)]$$
$$+ \mathbb{E}[\ln p(\omega | \alpha_1, \ldots, \alpha_d)]$$
$$+ \sum_{k=1}^{d} \mathbb{E}[\ln p(\alpha_k)]$$
$$+ \sum_{i=1}^{N} \mathbb{E}[\ln p(y_i | x_i, \omega, \lambda)]$$
$$- \mathbb{E}[\ln q(\lambda)]$$
$$- \mathbb{E}[\ln q(\omega)]$$
$$- \sum_{k=1}^{d} \mathbb{E}[\ln q(\alpha_k)]$$

$$\mathcal{L}(a_1, b_1, e_1, f_1, \mu_1, \Sigma_1) = \text{constants} + (e_0 - 1)[\psi(e_1) - \ln f_1] - f_0 \times \frac{e_1}{f_1}$$
$$+ \text{constants} + \frac{1}{2} \sum_{k=1}^{d} (\psi(a_{1k}) - \ln b_{1k}) - \frac{1}{2}\left[\text{trace}\left(\text{diag}\left(\frac{a_{11}}{b_{11}}, \ldots, \frac{a_{1d}}{b_{1d}}\right) \Sigma_1\right) + \mu_1^T \text{diag}\left(\frac{a_{11}}{b_{11}}, \ldots, \frac{a_{1d}}{b_{1d}}\right) \mu_1\right]$$
$$+ \text{constants} + \sum_{k=1}^{d}\left[(a_{0k} - 1)[\psi(a_{1k}) - \ln b_{1k}] - b_{0k}\frac{a_{1k}}{b_{1k}}\right]$$
$$+ \left[\text{constants} + \frac{1}{2}(\psi(e_1) - \ln f_1) - \frac{1}{2} \times \frac{e_1}{f_1}\left(\sum_{i=1}^{N}(y_i - x_i^T \mu_1)^2 + x_i^T \Sigma_1 x_i\right)\right]$$
$$- \left[e_1 - \ln f_1 + \ln \Gamma(e_1) + (1 - e_1)\psi(e_1)\right]$$
$$- \left[\text{constants} + \frac{1}{2}\ln|\Sigma_1|\right]$$
$$- \sum_{k=1}^{d}\left[a_{1k} - \ln b_{1k} + \ln \Gamma(a_{1k}) + (1 - a_{1k})\psi(a_{1k})\right]$$

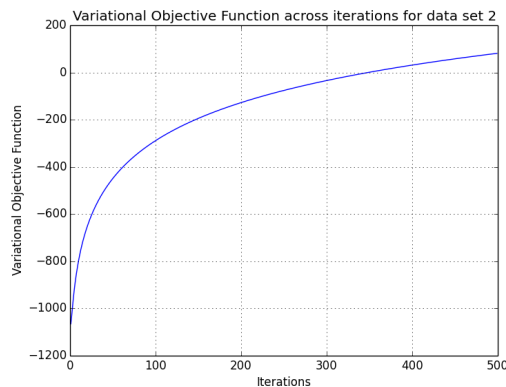Note, in the actual calculation of the objective function, all of the constants can be ignored.

# Problem 2

## a)



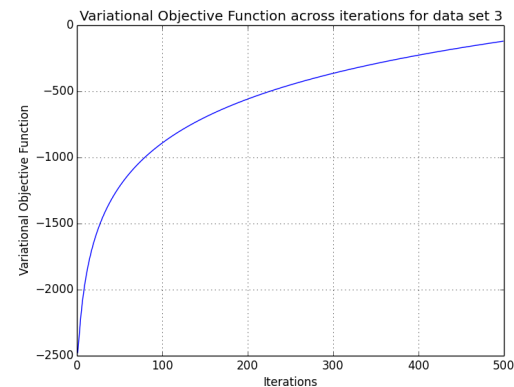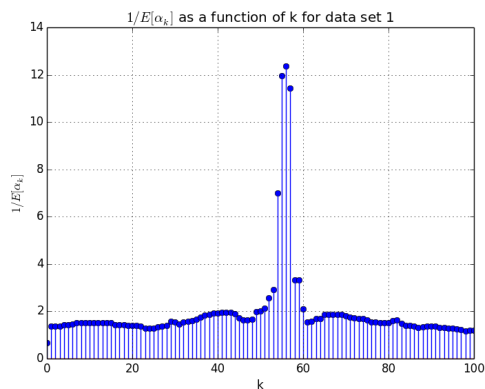Data Set 1 — Variational Objective Function across iterations for data set 1

Data Set 2 — Variational Objective Function across iterations for data set 2

Data Set 3 — Variational Objective Function across iterations for data set 3

## b)



Data Set 1 — $1/E[\alpha_k]$ as a function of k for data set 1

Data Set 2 — $1/E[\alpha_k]$ as a function of k for data set 2
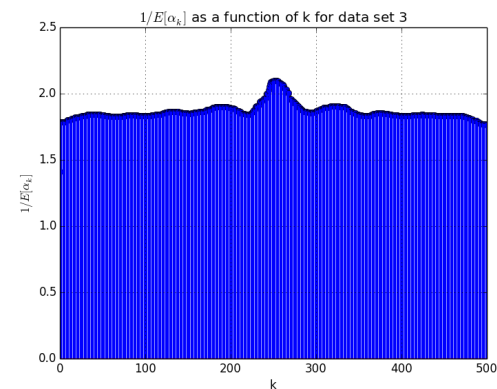
Data Set 3 — $1/E[\alpha_k]$ as a function of k for data set 3

## c)

**Data Set 1**

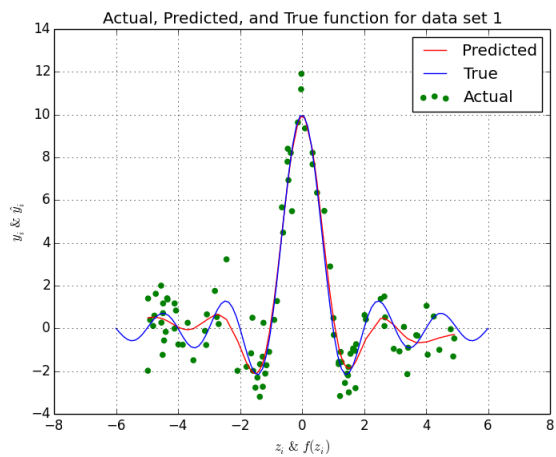$$1/\mathbb{E}_q[\lambda] = 3.67329015015$$

**Data Set 2**
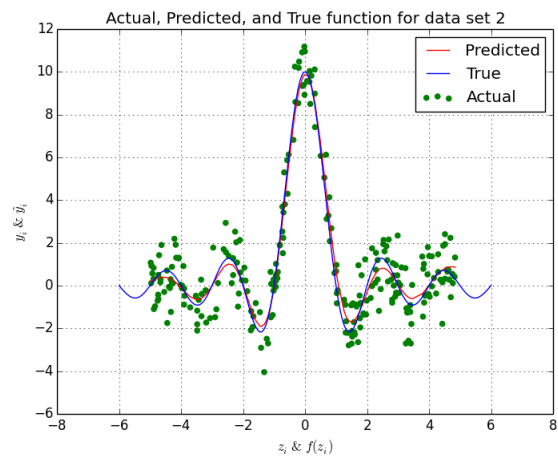
$$1/\mathbb{E}_q[\lambda] = 10.2601623179$$

**Data Set 3**

$$1/\mathbb{E}_q[\lambda] = 33.4137810438$$

d)

**Data Set 1**  **Data Set 2**  **Data Set 3**



Actual, Predicted, and True function for data set 1



Actual, Predicted, and True function for data set 2



Actual, Predicted, and True function for data set 3