

**Problem 1** (A Classic on the Gaussian Algebra, 10pts)

Let  $X$  and  $Y$  be independent univariate Gaussian random variables. In the previous problem set, you likely used the closure property that  $Z = X + Y$  is also a Gaussian random variable. Here you'll prove this fact.

- (a) Suppose  $X$  and  $Y$  have mean 0 and variances  $\sigma_X^2$  and  $\sigma_Y^2$  respectively. Write the pdf of  $X + Y$  as an integral.
- (b) Evaluate the integral from the previous part to find a closed-form expression for the pdf of  $X + Y$ , then argue that this expression implies that  $X + Y$  is also Gaussian with mean 0 and variance  $\sigma_X^2 + \sigma_Y^2$ . Hint: what is the integral, over the entire real line, of

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right),$$

i.e., the pdf of a univariate Gaussian random variable?

- (c) Extend the above result to the case in which  $X$  and  $Y$  may have arbitrary means.
- (d) Univariate Gaussians are supported on the entire real line. Sometimes this is undesirable because we are modeling a quantity with positive support. A common way to transform a Gaussian to solve this problem is to exponentiate it. Suppose  $X$  is a univariate Gaussian with mean  $\mu$  and variance  $\sigma^2$ . What is the pdf of  $e^X$ ?

**Solution**

- (a) Let  $X = x$ . Then  $Z = z \iff Y = z - x$ . Thus

$$\begin{aligned} p_Z(z) &= \int_{-\infty}^{\infty} p_X(x)p_Y(z-x)dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{x^2}{2\sigma_X^2}\right) \cdot \frac{1}{\sqrt{2\pi}\sigma_Y} \exp\left(-\frac{(z-x)^2}{2\sigma_Y^2}\right) dx \end{aligned}$$

- (b)

$$\begin{aligned} &\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{x^2}{2\sigma_X^2}\right) \cdot \frac{1}{\sqrt{2\pi}\sigma_Y} \exp\left(-\frac{(z-x)^2}{2\sigma_Y^2}\right) dx \\ &= \frac{1}{\sqrt{2\pi(\sigma_X^2 + \sigma_Y^2)}} \exp\left(-\frac{z^2}{2(\sigma_X^2 + \sigma_Y^2)}\right) \cdot \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\frac{\sigma_X\sigma_Y}{\sqrt{\sigma_X^2 + \sigma_Y^2}}}} \exp\left(\frac{x - \frac{z\sigma_X^2}{\sigma_X^2 + \sigma_Y^2}}{2\left(\frac{\sigma_X\sigma_Y}{\sqrt{\sigma_X^2 + \sigma_Y^2}}\right)^2}\right) dx. \end{aligned}$$

(Disclaimer: I got Mathematica to do the evaluation above, since I assumed a high probability of my making some small algebra error.)

But now the term under the integral is itself a normal distribution, and so it evaluates to 1, leaving us with

$$\begin{aligned}
& \frac{1}{\sqrt{2\pi(\sigma_X^2 + \sigma_Y^2)}} \exp\left(-\frac{z^2}{2(\sigma_X^2 + \sigma_Y^2)}\right) \cdot \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi} \frac{\sigma_X \sigma_Y}{\sqrt{\sigma_X^2 + \sigma_Y^2}}} \exp\left(\frac{x - \frac{z\sigma_X^2}{\sigma_X^2 + \sigma_Y^2}}{2\left(\frac{\sigma_X \sigma_Y}{\sqrt{\sigma_X^2 + \sigma_Y^2}}\right)^2}\right) \\
&= \frac{1}{\sqrt{2\pi(\sigma_X^2 + \sigma_Y^2)}} \exp\left(-\frac{z^2}{2(\sigma_X^2 + \sigma_Y^2)}\right) \cdot 1
\end{aligned}$$

and now this expression is easily identifiable as a Gaussian

$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{z^2}{2\sigma^2}\right) \quad \text{with variance } \sigma^2 = \sigma_X^2 + \sigma_Y^2 \text{ and mean } 0.$$

(c) Given arbitrary  $\mu_X, \mu_Y$  the means of  $X$  and  $Y$ , respectively, we have:

$$\begin{aligned}
p_Z(z) &= \int_{-\infty}^{\infty} p_X(x) p_Y(z-x) dx \\
&= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left(-\frac{(x-\mu_X)^2}{2\sigma_X^2}\right) \cdot \frac{1}{\sqrt{2\pi}\sigma_Y} \exp\left(-\frac{((z-x)-\mu_Y)^2}{2\sigma_Y^2}\right) dx \\
&= \frac{1}{\sqrt{2\pi(\sigma_X^2 + \sigma_Y^2)}} \exp\left(-\frac{(z-(\mu_X + \mu_Y))^2}{2(\sigma_X^2 + \sigma_Y^2)}\right) \cdot \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi} \frac{\sigma_X \sigma_Y}{\sqrt{\sigma_X^2 + \sigma_Y^2}}} \exp\left(\frac{x - \frac{(z-\mu_Y)\sigma_X^2 + \mu_X \sigma_Y^2}{\sigma_X^2 + \sigma_Y^2}}{2\left(\frac{\sigma_X \sigma_Y}{\sqrt{\sigma_X^2 + \sigma_Y^2}}\right)^2}\right) \\
&= \frac{1}{\sqrt{2\pi(\sigma_X^2 + \sigma_Y^2)}} \exp\left(-\frac{(z-(\mu_X + \mu_Y))^2}{2(\sigma_X^2 + \sigma_Y^2)}\right) \text{ again since everything under the integral is also a Gaussian}
\end{aligned}$$

where thanks again must be extended to Mathematica for help with the integration.

And thus, again, we get another Gaussian

$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(z-\mu)^2}{2\sigma^2}\right) \quad \text{with variance } \sigma^2 = \sigma_X^2 + \sigma_Y^2 \text{ and mean } \mu = \mu_X + \mu_Y.$$

(d) Let  $Y = e^X$ . Then  $\forall x \in X, y \in Y, x = \log(y)$ . Moreover,  $\frac{dx}{dy} = \frac{d}{dx} \log(y) = \frac{1}{|y|}$ .

Then by the change of variables formula (Murphy, pg 50), we have that

$$\begin{aligned}
p_Y(y) &= p_X(x) \cdot \left|\frac{dx}{dy}\right| \\
&= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \cdot \frac{1}{|y|} \\
&= \frac{1}{|y|\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\log(y)-\mu)^2}{2\sigma^2}\right)
\end{aligned}$$

**Problem 2** (Regression, 13pts)

Suppose that  $X \in \mathbb{R}^{n \times m}$  with  $n \geq m$  and  $Y \in \mathbb{R}^n$ , and that  $Y \sim \mathcal{N}(Xw, \sigma^2 I)$ . You learned in class that the maximum likelihood estimate  $\hat{w}$  of  $w$  is given by

$$\hat{w} = (X^T X)^{-1} X^T Y$$

- (a) Why do we need to assume that  $n \geq m$ ?
- (b) Define  $H = X(X^T X)^{-1} X^T$ , so that the “fitted” values  $\hat{Y} = X\hat{w}$  satisfy  $\hat{Y} = HY$ . Show that  $H$  is an orthogonal projection matrix that projects onto the column space of  $X$ , so that the fitted y-values are a projection of  $Y$  onto the column space of  $X$ .
- (c) What are the expectation and covariance matrix of  $\hat{w}$ ?
- (d) Compute the gradient with respect to  $w$  of the log likelihood implied by the model above, assuming we have observed  $Y$  and  $X$ .
- (e) Suppose we place a normal prior on  $w$ . That is, we assume that  $w \sim \mathcal{N}(0, \tau^2 I)$ . Show that the MAP estimate of  $w$  given  $Y$  in this context is

$$\hat{w}_{MAP} = (X^T X + \lambda I)^{-1} X^T Y$$

where  $\lambda = \sigma^2/\tau^2$ . (You may employ standard conjugacy results about Gaussians without proof in your solution.)

[Estimating  $w$  in this way is called *ridge regression* because the matrix  $\lambda I$  looks like a “ridge”. Ridge regression is a common form of *regularization* that is used to avoid the overfitting (resp. underdetermination) that happens when the sample size is close to (resp. higher than) the output dimension in linear regression.]

- (f) Do we need  $n \geq m$  to do ridge regression? Why or why not?
- (g) Show that ridge regression is equivalent to adding  $m$  additional rows to  $X$  where the  $j$ -th additional row has its  $j$ -th entry equal to  $\sqrt{\lambda}$  and all other entries equal to zero, adding  $m$  corresponding additional entries to  $Y$  that are all 0, and then computing the maximum likelihood estimate of  $w$  using the modified  $X$  and  $Y$ .

(a) If  $n < m$ , then  $X^T X$  is singular and thus non-invertible.

(b) Consider that by construction,  $\hat{Y}$  is a projection of  $Y$  onto a subspace of its basis vectors (indeed, this is essentially by definition what regression is). Considering this, if  $\hat{Y} = HY$ , then  $H$  *must* be the projection map  $H : \mathbb{R}^n \rightarrow \mathbb{R}^m$  of  $Y$  onto the basis of  $X$ .

(c) Let  $N \sim \mathcal{N}(0, \sigma^2 I)$  be the noise in our observed  $Y$ , and  $Y = Xw + N$ . Then

$$\begin{aligned} \mathbb{E}\hat{w} &= \mathbb{E}((X^T X)^{-1} X^T Y) \\ &= \mathbb{E}((X^T X)^{-1} X^T (Xw + N)) \\ &= \mathbb{E}((X^T X)^{-1} (X^T X)w) + \mathbb{E}((X^T X)^{-1} X^T N) \\ &= I \cdot \mathbb{E}w + (X^T X)^{-1} X^T \mathbb{E}N \\ &= I \cdot w + 0 \quad \text{since by assumption, the noise is normally distributed with mean 0} \\ &= w \end{aligned}$$

(d) We have the log likelihood

$$\begin{aligned}
\log(y|X, w, \sigma^2 I) &= \sum_{i=1}^n \log p(y_i | \mathbf{x}_i, w) \\
&= \sum_{i=1}^n \log \mathcal{N}(y_i | w^\top \mathbf{x}_i, \sigma^2 I) \\
&= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - w^\top \mathbf{x}_i)^2
\end{aligned}$$

and thus taking the gradient with respect to  $w$ , we have

$$\begin{aligned}
\frac{\partial}{\partial w} \log(y|X, w, \sigma^2 I) &= \frac{\partial}{\partial w} \left( -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - w^\top \mathbf{x}_i)^2 \right) \\
&= \sum_{i=1}^n (y_i - w^\top \mathbf{x}_i) \mathbf{x}_i^\top
\end{aligned}$$

(e) If  $w \sim \mathcal{N}(0, \tau^2 I)$ , then  $p(w) = \prod_{j=1}^m \mathcal{N}(0w_j | 0, \tau^2 I)$  and thus the MAP estimate is given by

$$\begin{aligned}
\arg \max_w \log p(w|y, X) &= \arg \max_w \log p(y, X|w)p(w) \\
&= \arg \max_w \left( \sum_{i=1}^n \log \mathcal{N}(y_i | w^\top \mathbf{x}_i, \sigma^2) + \sum_{j=1}^m \log \mathcal{N}(w_j | 0, \tau^2 I) \right)
\end{aligned}$$

and now differentiating with respect to  $w$  and setting the partial to 0,

$$\begin{aligned}
\frac{\partial}{\partial w} \sum_{i=1}^n \log \mathcal{N}(y_i | w^\top \mathbf{x}_i, \sigma^2) + \sum_{j=1}^m \log \mathcal{N}(w_j | 0, \tau^2 I) &= 0 \\
\implies \sum_{i=1}^N (y_i - \hat{w}^\top \mathbf{x}_i) \mathbf{x}_i^\top + \frac{\sigma^2}{\tau^2} \sum_{j=1}^m \hat{w}_j &= 0 \\
\implies X^\top Y - X^\top X \hat{w} = \lambda \hat{w} \quad \text{where } \lambda = \frac{\sigma^2}{\tau^2} \\
\implies (X^\top X + \lambda I) \hat{w} &= X^\top Y \\
\implies \hat{w} = (X^\top X + \lambda I)^{-1} X^\top Y \quad \text{where } \lambda = \frac{\sigma^2}{\tau^2}
\end{aligned}$$

(f) We don't need  $n \geq m$  for ridge regression, since if  $\lambda \neq 0$ ,  $X^\top X + \lambda I$  is invertible even if  $X^\top X$  is not. Thinking about it from a Bayesian point of view, this says that if we have some prior on what  $w$  is, then we can get away with having fewer data points, since we have a prior to inform what  $w$  might be.

(g) Let  $X' = \begin{pmatrix} X \\ \sqrt{\lambda}I \end{pmatrix}$ ,  $y' = \begin{pmatrix} y \\ \mathbf{0} \end{pmatrix}$ . Then

$$\begin{aligned}
\frac{\partial}{\partial w} \log(y'|X', w, \sigma^2 I) &= \sum_{i=1}^n (y'_i - w^\top \mathbf{x}'_i) \mathbf{x}'_i^\top \\
&= X'^\top y' - X'^\top X' w \\
&= \begin{pmatrix} X \\ \sqrt{\lambda}I \end{pmatrix}^\top \begin{pmatrix} y \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} X \\ \sqrt{\lambda}I \end{pmatrix}^\top \begin{pmatrix} X \\ \sqrt{\lambda}I \end{pmatrix} w \\
&= X^\top Y - (X^\top X + \sqrt{\lambda}\sqrt{\lambda}I)w \\
&= X^\top Y - (X^\top X + \lambda I)w
\end{aligned}$$

and now setting the derivative to 0, we get

$$w = (X^\top X + \lambda I)^{-1} X^\top Y,$$

just like in ridge regression.

**Problem 3** (The Dirichlet and Multinomial Distributions, 12pts)

The Dirichlet distribution over  $K$  categories is a generalization of the beta distribution. It has a shape parameter  $\alpha \in \mathbb{R}^K$  with non-negative entries and is supported over the set of  $K$ -dimensional positive vectors whose components sum to 1. Its density is given by

$$f(\theta_{1:K}|\alpha_{1:K}) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \theta_k^{\alpha_k-1}$$

(Notice that when  $K = 2$ , this reduces to the density of a beta distribution.) For the rest of this problem, assume a fixed  $K \geq 2$ .

- (a) Suppose  $\theta$  is Dirichlet-distributed with shape parameter  $\alpha$ . Without proof, state the value of  $E(\theta)$ . Your answer should be a vector defined in terms of either  $\alpha$  or  $K$  or potentially both.
- (b) Suppose that  $\theta \sim \text{Dir}(\alpha)$  and that  $X \sim \text{Cat}(\theta)$ , where  $\text{Cat}$  is a Categorical distribution. That is, suppose we first sample a  $K$ -dimensional vector  $\theta$  with entries in  $(0, 1)$  from a Dirichlet distribution and then roll a  $K$ -sided die such that the probability of rolling the number  $k$  is  $\theta_k$ . Prove that the posterior  $p(\theta|X)$  also follows a Dirichlet distribution. What is its shape parameter?
- (c) Now suppose that  $\theta \sim \text{Dir}(\alpha)$  and that  $X^{(1)}, X^{(2)}, \dots \stackrel{iid}{\sim} \text{Cat}(\theta)$ . Show that the posterior predictive after  $n - 1$  observations is given by,

$$P(X^{(n)} = k | X^{(1)}, \dots, X^{(n-1)}) = \frac{\alpha_k^{(n)}}{\sum_k \alpha_k^{(n)}}$$

where for all  $k$ ,  $\alpha_k^{(n)} = \alpha_k + \sum_{i=1}^{n-1} \mathbf{1}\{X^{(i)} = k\}$ . (Bonus points if your solution does not involve any integrals.)

- (d) Consider the random vector  $Z_k = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X^{(i)} = k\}$  for all  $k$ . What is the mean of this vector? What is the distribution of the vector? (If you're not sure how to rigorously talk about convergence of random variables, give an informal argument. Hint: what would you say if  $\theta$  were fixed?) What is the marginal distribution of a single class  $p(Z_k)$ ?
- (e) Suppose we have  $K$  distinct colors and an urn with  $\alpha_k$  balls of color  $k$ . At each time step, we choose a ball uniformly at random from the urn and then add into the urn an additional new ball of the same color as the chosen ball. (So if at the first time step we choose a ball of color 1, we'll end up with  $\alpha_1 + 1$  balls of color 1 and  $\alpha_k$  balls of color  $k$  for all  $k > 1$  at the start of the second time step.) Let  $\rho_k^{(n)}$  be the fraction of all the balls that are of color  $k$  at time  $n$ . What is the distribution of  $\lim_{n \rightarrow \infty} \rho_k^{(n)}$ ? Prove your answer.

(a)

$$\mathbb{E}(\theta) = \left(\frac{\alpha_1}{\alpha}, \dots, \frac{\alpha_K}{\alpha}\right).$$

(b) Let  $p(X|\theta) = \prod_{k=1}^K \theta_k^{N_k}$  be the likelihood of  $X$  given  $\theta$  such that  $X$  is categorical, and  $N_k$  be the number of times  $k$  was rolled. Then

$$\begin{aligned}
p(\theta|X) &\propto p(X|\theta)p(\theta) \\
&\propto \prod_{k=1}^K \theta_k^{N_k} \theta_k^{\alpha_k-1} \\
&= \prod_{k=1}^K \theta_k^{\alpha_k+N_k-1}
\end{aligned}$$

which is given by another Dirichlet distribution

$$Dir(\theta|\alpha + N)$$

with shape parameter  $\alpha + N$ , where  $N = (N_1, N_2, \dots, N_K)$ .

(c) Since we want to have the sum of the  $\theta_k$  sum to 1, so that it is a probability distribution, we can use Lagrange multipliers to find the posterior predictive without integrals. Then, taking the log of the likelihood above, we get that our Lagrangian is

$$\begin{aligned}
L &= \log\left(\prod_{k=1}^K \theta_k^{\alpha_k+N_k-1}\right) \\
&= \sum_{k=1}^K (\alpha_k + N_k - 1) \log \theta_k + \lambda(1 - \sum_{k=1}^K \theta_k).
\end{aligned}$$

Now, differentiating with respect to  $\lambda$  and setting the partial to 0, we get

$$\begin{aligned}
\frac{\partial}{\partial \lambda} \sum_{k=1}^K (\alpha_k + N_k - 1) \log \theta_k + \lambda(1 - \sum_{k=1}^K \theta_k) \\
= 1 - \sum_{k=1}^K \theta_k = 0.
\end{aligned}$$

Now doing the same with respect to  $\theta_k$ , we have

$$\begin{aligned}
\frac{\partial}{\partial \theta_k} \sum_{k=1}^K (\alpha_k + N_k - 1) \log \theta_k + \lambda(1 - \sum_{k=1}^K \theta_k) \\
= \frac{\alpha_k + N_k - 1}{\theta_k} - \lambda = 0 \\
\implies \frac{\alpha_k + N_k - 1}{\lambda} = \theta_k
\end{aligned}$$

Then we have that

$$\begin{aligned}
\sum_{k=1}^K \theta_k &= \sum_{k=1}^K \frac{\alpha_k + N_k - 1}{\lambda} = 1 \\
\implies \lambda &= \sum_{k=1}^K \alpha_k + \sum_{k=1}^K N_k - K.
\end{aligned}$$

Thus

$$\begin{aligned}\hat{\theta}_k &= \frac{\alpha_k + N_k - 1}{\sum_{k=1}^K \alpha_k + \sum_{k=1}^K N_k - K} \\ &= \frac{\alpha_k + N_{k-1}}{\sum_{k=1}^K \alpha_k + \sum_{k=1}^K N_{k-1}}\end{aligned}$$

Then in order to find the posterior predictive, we can simply consider that the posterior

$$p(X^{(n)} = k | X^{(1)}, \dots, X^{(n)}) = \mathbb{E}(\theta_k | X^{(1)}, \dots, X^{(n)}).$$

And now we have that

$$\begin{aligned}p(X^{(n)} = k | X^{(1)}, \dots, X^{(n)}) &= \mathbb{E}(\theta_k | X^{(1)}, \dots, X^{(n)}) \\ &= \frac{\alpha_k + N_{k-1}}{\sum_{k=1}^K \alpha_k + \sum_{k=1}^K N_{k-1}} \\ &= \frac{\alpha_k^{(n)}}{\sum_k \alpha_k^{(n)}}.\end{aligned}$$

(d) In the infinite limit, the distribution assumes that of the MLE, since the information from the prior vanishes as  $n \rightarrow \infty$ . In particular, we have that each  $\hat{\theta}_k = \frac{N_k}{\sum_k N_k}$ , which is just the fraction of times that we see side  $k$ .

(e)



# Physicochemical Properties of Protein Tertiary Structure

In the following problems we will code two different approaches for solving linear regression problems and compare how they scale as a function of the dimensionality of the data. We will also investigate the effects of linear and non-linear features in the predictions made by linear models.

We will be working with the regression data set Protein Tertiary Structure: <https://archive.ics.uci.edu/ml/machine-learning-databases/00265/> (download CASP.csv). This data set contains information about predicted conformations for 45730 proteins. In the data, the target variable  $y$  is the root-mean-square deviation (RMSD) of the predicted conformations with respect to the true properly folded form of the protein. The RMSD is the measure of the average distance between the atoms (usually the backbone atoms) of superimposed proteins. The features  $\mathbf{x}$  are physico-chemical properties of the proteins in their true folded form. After downloading the file CASP.csv we can load the data into python using

```
>>> import numpy as np
>>> data = np.loadtxt("CASP.csv", delimiter = ",", skiprows = 1)
```

We can then obtain the vector of target variables and the feature matrix using

```
>>> y = data[:, 0]
>>> X = data[:, 1:]
```

We can then split the original data into a training set with 90% of the data entries in the file CASP.csv and a test set with the remaining 10% of the entries. Normally, the splitting of the data is done at random, but here **we ask you to put into the training set the first 90% of the elements from the file CASP.csv** so that we can verify that the values that you will be reporting are correct. (This should not cause problems, because the rows of the file are in a random order.)

We then ask that you **normalize** the features so that they have zero mean and unit standard deviation in the training set. This is a standard step before the application of many machine learning methods. After these steps are done, we can concatenate a **bias feature** (one feature which always takes value 1) to the observations in the normalized training and test sets.

We are now ready to apply our machine learning methods to the normalized training set and evaluate their performance on the normalized test set. In the following problems, you will be asked to report some numbers and produce some figures. Include these numbers and figures in your assignment report. **The numbers should be reported with up to 8 decimals.**

## Problem 4 (7pts)

Assume that the targets  $y$  are obtained as a function of the normalized features  $\mathbf{x}$  according to a Bayesian linear model with additive Gaussian noise with variance  $\sigma^2 = 1.0$  and a Gaussian prior on the regression coefficients  $\mathbf{w}$  with *precision* matrix  $\Sigma^{-1} = \tau^{-2}\mathbf{I}$  where  $\tau^{-2} = 10$ . Code a routine using the **QR decomposition** (see Section 7.5.2 in Murphy's book) that finds the Maximum a Posteriori (MAP) value  $\hat{\mathbf{w}}$  for  $\mathbf{w}$  given the normalized training data

- Report the value of  $\hat{\mathbf{w}}$  obtained.
- Report the root mean squared error (RMSE) of  $\hat{\mathbf{w}}$  in the normalized test set.

## Solution

The value of  $\hat{\mathbf{w}}$  found was such that

$$\hat{\mathbf{w}}^\top = [1.34308416e^{-3}, 1.88519687e^{-3}, 1.26299471e^1, -9.76437769e^{-2}, -1.98338279e^{-6}, \\ -2.47087737e^{-2}, -8.93816049e^{-5}, 1.62440374e^{-2}, 1.37538211e^{-1}, 0.00000000e^0]$$

The RMSE of  $\hat{\mathbf{w}}$  was 1116.4433156731677.

### Problem 5 (14pts)

L-BFGS is an iterative method for solving general nonlinear optimization problems. For this problem you will use this method as a black box that returns the MAP solution by sequentially evaluating the objective function and its gradient for different input values. The goal of this problem is to use a built-in implementation of the L-BFGS algorithm to find a point estimate that maximizes our posterior of interest. Generally L-BFGS requires your black box to provide two values: the current objective and the gradient of the objective with respect to any parameters of interest. To use the optimizer, you need to first write two functions: (1) to compute the loss, or the *negative* log-posterior and (2) to compute the gradient of the loss with respect to the weights  $w$ .

As a preliminary to coming work in the class, we will use the L-BFGS implemented in PyTorch. [Warning: For this assignment we are using a small corner of the PyTorch world. Do not feel like you need to learn everything about this library.]

There are three parts to using this optimizer:

1. Create a vector of weights in NumPy, wrap in a pytorch **Tensor** and **Variable**, and pass to the optimizer.

```
from torch import Tensor
from torch.autograd import Variable

# Construct a PyTorch variable array (called tensors).
weights = Variable(Tensor(size))

# Initialize an optimizer of the weights
optimizer = torch.optim.LBFGS([weights])

...
```

2. Write a python function that uses the current weights to compute the log-posterior **and** sets `weights.grad` to be the gradient of the log-posterior with respect to the current weights.

```
def black_box():
    # Access the value of the variable as a numpy array.
    weights_data = weights.data.numpy()

    ...

    # Set the gradient of the variable.
    weights.grad = Tensor({numpy})

    return {objective}
```

3. Repeatedly call `optimizer.step(black_box)` to optimize.

[If you are feeling adventurous, you might find it useful to venture into the land of autograd and check your computation with PyTorch's `torch.autograd.gradcheck.get_numerical_jacobian`.]

- After running for 100 iterations, report the value of  $\hat{\mathbf{w}}$  obtained.
- Report the RMSE of the predictions made with  $\hat{\mathbf{w}}$  in the normalized test set.

**Problem 6** (14pts)

Linear regression can be extended to model non-linear relationships by replacing the original features  $\mathbf{x}$  with some non-linear functions of the original features  $\phi(\mathbf{x})$ . We can automatically generate one such non-linear function by sampling a random weight vector  $\mathbf{a} \sim \mathcal{N}(0, \mathbf{I})$  and a corresponding random bias  $b \sim \mathcal{U}[0, 2\pi]$  and then making  $\phi(\mathbf{x}) = \cos(\mathbf{a}^T \mathbf{x} + b)$ . By repeating this process  $d$  times we can generate  $d$  non-linear functions that, when applied to the original features, produce a non-linear mapping of the data into a new  $d$  dimensional space. We can encode these  $d$  functions into a matrix  $\mathbf{A}$  with  $d$  rows, each one with the weights for each function, and a  $d$ -dimensional vector  $\mathbf{b}$  with the biases for each function. The new mapped features are then obtained as  $\phi(\mathbf{x}) = \cos(\mathbf{A}\mathbf{x} + \mathbf{b})$ , where  $\cos$  applied to a vector returns another vector whose elements are the result of applying  $\cos$  to the individual elements of the original vector.

Generate 4 sets of non-linear functions, each one with  $d = 100, 200, 400, 600$  functions, respectively, and use them to map the features in the original normalized training and test sets into 4 new feature spaces, each one of dimensionality given by the value of  $d$ . After this, for each value of  $d$ , find the MAP solution  $\hat{\mathbf{w}}$  for  $\mathbf{w}$  using the corresponding new training set and the method from problem 4. Use the same values for  $\sigma^2$  and  $\tau^{-2}$  as before. You are also asked to record the time taken by the method QR to obtain a value for  $\hat{\mathbf{w}}$ . In python you can compute the time taken by a routine using the time package:

```
>>> import time
>>> time_start = time.time()
>>> routine_to_call()
>>> running_time = time.time() - time_start
```

Next, compute the RMSE of the resulting predictor in the normalized test set. Repeat this process with the method from problem 5 (L-BFGS).

- Report the test RMSE obtained by each method for each value of  $d$ .

You are asked to generate a plot with the results obtained by each method (QR and L-BFGS) for each value of  $d$ . In this plot the  $x$  axis should represent the time taken by each method to run and the  $y$  axis should be the RMSE of the resulting predictor in the normalized test set. The plot should contain 4 points in red, representing the results obtained by the method QR for each value of  $d$ , and 4 points in blue, representing the results obtained by the method L-BFGS for each value of  $d$ . Answer the following questions:

- Do the non-linear transformations help to reduce the prediction error? Why?
- What method (QR or L-BFGS) is faster? Why?
- (Extra Problem, Not Graded) Instead of using random  $\mathbf{A}$ , what if we treat  $\mathbf{A}$  as another parameter for L-BFGS to optimize? You can do this by wrapping it as a variable and passing to the constructor. Compute its gradient as well in *black\_box* either analytically or by using PyTorch *autograd*.