# CSC412/2056 Assignment #1

**Problem 1** (Variance and covariance, 6 points)

Let $X$ and $Y$ be two continuous independent random variables.

(a) Starting from the definition of independence, show that the independence of $X$ and $Y$ implies that their covariance is zero.

(b) For a scalar constant $a$, show the following two properties, starting from the definition of expectation:

$$\mathbb{E}(X + aY) = \mathbb{E}(X) + a\mathbb{E}(Y)$$
$$\text{var}(X + aY) = \text{var}(X) + a^2\text{var}(Y)$$

---

**Problem 2** (Densities, 5 points)

Answer the following questions:

(a) Can a probability density function (pdf) ever take values greater than 1?

(b) Let $X$ be a univariate normally distributed random variable with mean 0 and variance $1/100$. What is the pdf of $X$?

(c) What is the value of this pdf at 0?

(d) What is the probability that $X = 0$?

---

**Problem 3** (Calculus, 4 points)

Let $x, y \in \mathbb{R}^m$ and $A \in \mathbb{R}^{m \times m}$. Please answer the following questions, writing your answers in vector notation.

(a) What is the gradient with respect to $x$ of $x^T y$?

(b) What is the gradient with respect to $x$ of $x^T x$?

(c) What is the gradient with respect to $x$ of $x^T A$?

(d) What is the gradient with respect to $x$ of $x^T A x$?

---

**Problem 4** (Linear Regression, 10pts)

Suppose that $X \in \mathbb{R}^{n \times m}$ with $n \geq m$ and $Y \in \mathbb{R}^n$, and that $Y \sim \mathcal{N}(X\beta, \sigma^2 I)$. In this question you will derive the result that the maximum likelihood estimate $\hat{\beta}$ of $\beta$ is given by

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

(a) What are the expectation and covariance matrix of $\hat{\beta}$, for a given true value of $\beta$?

(b) Show that maximizing the likelihood is equivalent to minimizing the squared error $\sum_{i=1}^n (y_i - x_i\beta)^2$. [Hint: Use $\sum_{i=1}^n a_i^2 = a^T a$]

(c) Write the squared error in vector notation, (see above hint), expand the expression, and collect like terms. [Hint: Use $\beta^T x^T y = y^T x \beta$ (why?) and $x^T x$ is symmetric ]

(d) Take the derivative of this expanded expression with respect to $\beta$ to show the maximum likelihood estimate $\hat{\beta}$ as above. [Hint: Use results 3.c and 3.d for derivatives in vector notation.]

**Problem 5** (Ridge Regression, 10pts)

Suppose we place a normal prior on $\beta$. That is, we assume that $\beta \sim \mathcal{N}(0, \tau^2 I)$.

(a) Show that the MAP estimate of $\beta$ given $Y$ in this context is

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

where $\lambda = \sigma^2 / \tau^2$.

Estimating $\beta$ 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 that happens when the sample size is close to the output dimension in linear regression.

(b) 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 and then computing the maximum likelihood estimate of $\beta$ using the modified $X$ and $Y$.

**Problem 6** (Gaussians in high dimensions, 10pts)

In this question we will investigate how our intuition for samples from a Gaussian may break down in higher dimensions. Consider samples from a $D$-dimensional unit Gaussian $\mathbf{x} \sim \mathcal{N}(\mathbf{0}_D, \mathbb{I}_D)$ where $\mathbf{0}_D$ indicates a column vector of $D$ zeros and $\mathbb{I}_D$ is a $D \times D$ identity matrix.

1. Starting with the definition of Euclidean norm, quickly show that the distance of $\mathbf{x}$ from the origin is $\sqrt{\mathbf{x}^\mathsf{T}\mathbf{x}}$

2. In low-dimensions our intuition tells us that samples from the unit Gaussian will be near the origin. Draw 10000 samples from a $D = 1$ Gaussian and plot a normalized histogram for the distance of those samples from the origin. Does this confirm your intuition that the samples will be near the origin?

3. Draw 10000 samples from $D = \{1, 2, 3, 10, 100\}$ Gaussians and, on a single plot, show the normalized histograms for the distance of those samples from the origin. As the dimensionality of the Gaussian increases, what can you say about the expected distance of the samples from the Gaussian's mean (in this case, origin).

4. From Wikipedia, if $\mathbf{x}_i$ are $k$ independent, normally distributed random variables with means $\mu_i$ and standard deviations $\sigma_i$ then the statistic $Y = \sqrt{\sum_{i=1}^{k}(\frac{x_i-\mu_i}{\sigma_i})^2}$ is distributed according to the $\chi$-distribution. On the previous normalized histogram, plot the probability density function (pdf) of the $\chi$-distribution for $k = \{1, 2, 3, 10, 100\}$.

5. Taking two samples from the $D$-dimensional unit Gaussian, $\mathbf{x}_a, \mathbf{x}_b \sim \mathcal{N}(\mathbf{0}_D, \mathbb{I}_D)$ how is $\mathbf{x}_a - \mathbf{x}_b$ distributed? Using the above result about $\chi$-distribution, how is $||\mathbf{x}_a - \mathbf{x}_b||_2$ distributed? (Hint: start with a $\mathcal{X}$-distributed random variable and use the change of variables formula.) Plot the pdfs of this distribution for $k = \{1, 2, 3, 10, 100\}$. How does the distance between samples from a Gaussian behave as dimensionality increases? Confirm this by drawing two sets of 1000 samples from the $D$-dimensional unit Gaussian. On the plot of the $\chi$-distribution pdfs, plot the normalized histogram of the distance between samples from the first and second set.

6. In lecture we saw examples of interpolating between latent points to generate convincing data. Given two samples from a gaussian $\mathbf{x}_a, \mathbf{x}_b \sim \mathcal{N}(\mathbf{0}_D, \mathbb{I}_D)$ the linear interpolation between them $x_\alpha$ is defined as a function of $\alpha \in [0, 1]$

$$\text{lin\_interp}(\alpha, \mathbf{x}_a, \mathbf{x}_b) = \alpha \mathbf{x}_a + (1 - \alpha)\mathbf{x}_b$$

   For two sets of 1000 samples from the unit gaussian in $D$-dimensions, plot the average log-likelihood along the linear interpolations between the pairs of samples as a function of $\alpha$. (i.e. for each pair of samples compute the log-likelihood along a linear space of interpolated points between them, $\mathcal{N}(x_\alpha|0, I)$ for $\alpha \in [0, 1]$. Plot the average log-likelihood over all the interpolations.) Do this for $D = \{1, 2, 3, 10, 100\}$, one plot per dimensionality. Comment on the log-likelihood under the unit Gaussian of points along the linear interpolation. Is a higher log-likelihood for the interpolated points necessarily better? Given this, is it a good idea to linearly interpolate between samples from a high dimensional Gaussian?

7. Instead we can interpolate in polar coordinates: For $\alpha \in [0, 1]$ the polar interpolation is

$$\text{polar\_interp}(\alpha, \mathbf{x}_a, \mathbf{x}_b) = \sqrt{\alpha}\mathbf{x}_a + \sqrt{(1 - \alpha)}\mathbf{x}_b$$

   This interpolates between two points while maintaining Euclidean norm. On the same plot from the previous question, plot the probabilitiy density of the polar interpolation between pairs of samples from two sets of 1000 samples from $D$-dimensional unit Gaussians for $D = \{1, 2, 3, 10, 100\}$. Comment on the log-likelihood under the unit Gaussian of points along the polar interpolation. Give an intuative explanation for why polar interpolation is more suitable than linear interpolation for high dimensional Gaussians. **For 6. and 7. you should have one plot for each $D$ with two curves on each**.

8. (**Bonus 5pts**) In the previous two questions we compute the average loglikelihood of the linear and polar interpolations under the unit gaussian. Instead, consider the norm along the interpolation, $\sqrt{\mathbf{x}_\alpha^\mathsf{T}\mathbf{x}_\alpha}$. As we saw previously, this is distributed according to the $\mathcal{X}$-distribution. Compute and plot the average log-likelihood of the norm along the two interpolations under the the $\mathcal{X}$-distribution for $D = \{1, 2, 3, 10, 100\}$, i.e. $\mathcal{X}_D(\sqrt{\mathbf{x}_\alpha^\mathsf{T}\mathbf{x}_\alpha})$. There should be one plot for each $D$, each with two curves corresponding to log-likelihood of linear and polar interpolations. How does the log-likelihood along the linear interpolation compare to the log-likelihood of the true samples (endpoints)? Using your answer for questions 3 and 4, provide geometric intuition for the log-likelihood along the linear and polar interpolations. Use this to further justify your explanation for the suitability of polar v.s. linear interpolation.