

1. (Projections) Are the following statements true or false? Prove that they are true or provide a counterexample.

(a) The projection of a vector on a subspace  $\mathcal{S}$  is equal to

$$\mathcal{P}_{\mathcal{S}} x = \sum_{i=1}^n \langle x, b_i \rangle b_i$$

for any basis  $b_1, \dots, b_d$  of  $\mathcal{S}$ . False Consider  $\mathbf{b}_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$  and  $\mathbf{b}_2 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ , they form a basis of  $\mathbf{R}^2$ . When using the definition  $\mathcal{P}_{\mathcal{S}} x = \sum_{i=1}^n \langle x, b_i \rangle b_i$  we would expect that  $\mathcal{P}_{\mathcal{S}} b_1 = b_1$ . However  $\mathcal{P}_{\mathcal{S}} b_1 = \begin{bmatrix} 2 \\ 5 \end{bmatrix} \neq b_1$ .

(b) The orthogonal complement of the orthogonal complement of a subspace  $\mathcal{S} \subseteq \mathbb{R}^n$  is  $\mathcal{S}$ . True Let  $S^{\perp} = \{x | \langle x, y \rangle = 0, \forall y \in S\}$  a subspace of an inner product space  $X$ , then  $S^{\perp\perp} = \{x | \langle x, y \rangle = 0, \forall y \in S^{\perp}\}$ . The inner product being symmetric,  $S \subseteq S^{\perp\perp}$ . Since for any vector  $x \in X$ , we have  $x = y + z$  where  $y \in S, z \in S^{\perp}$ , using Gram-schmidt orthonormalization process, we can find a basis of  $S$  and  $S^{\perp}$  which express any vector of  $X$  as a linear combination of these two basis and combining these two basis together forms a new basis for  $X$  so  $\dim X = \dim S + \dim S^{\perp}$ . If  $\dim X = n$  and  $\dim S = m$  then  $\dim S^{\perp} = n - m$ . Similarly  $\dim S^{\perp\perp} = n - (n - m) = m$  so  $\dim S^{\perp\perp} = \dim S$ , so  $S^{\perp\perp} \subseteq S$  and since the dimension of a space or subspace is the cardinality of its basis, thus  $S = S^{\perp\perp}$ .

(c) Replacing each entry of a vector in  $\mathbb{R}^n$  by the average of all its entries is equivalent to projecting the vector onto a subspace. True consider  $\mathbf{v} =$

$$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}, \text{ we want } \mathbf{w} = \begin{bmatrix} \frac{\sum_{i=1,n} v_i}{n} \\ \vdots \\ \frac{\sum_{i=1,n} v_i}{n} \end{bmatrix}. \text{ The orthogonal projection of } \mathbf{v} \text{ onto}$$

the vector  $\mathbf{b}$  is defined as  $\frac{\mathbf{v} \cdot \mathbf{b}}{\|\mathbf{b}\|^2}$ , take  $\mathbf{b} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$ .

2. (Eigen decomposition) The populations of deer and wolfs in Yellowstone are well approximated by

$$d_{n+1} = \frac{5}{4}d_n - \frac{3}{4}w_n, \quad (1)$$

$$w_{n+1} = \frac{1}{4}d_n + \frac{1}{4}w_n, \quad n = 0, 1, 2, \dots, \quad (2)$$

where  $d_n$  and  $w_n$  denote the number of deer and wolfs in year  $n$ . Assuming that there are more deer than wolfs to start with ( $w_0 < d_0$ ), what is the proportion between the numbers of deer and wolfs as  $n \rightarrow \infty$ ?

Rewriting the problem in a matrix form:

$$\begin{pmatrix} d_{n+1} \\ w_{n+1} \end{pmatrix} = \frac{1}{4} \begin{pmatrix} 5 & -3 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} d_n \\ w_n \end{pmatrix}$$

Let  $A = \frac{1}{4} \begin{pmatrix} 5 & -3 \\ 1 & 1 \end{pmatrix}$ ,  $v_{n+1} = \begin{pmatrix} d_{n+1} \\ w_{n+1} \end{pmatrix}$ ,  $v_0 = \begin{pmatrix} d_0 \\ w_0 \end{pmatrix}$  then  $v_{n+1} = Av_n = A^n v_0$ . We are looking to find the eigen decomposition so we can understand the behavior of  $v_n$  as  $n \rightarrow \infty$ .  $\det(A - \lambda I) = \frac{1}{2}(2\lambda^2 - 3\lambda + 1)$ , we find for eigenvalues  $\lambda_1 = \frac{1}{2}$  and  $\lambda_2 = 1$  with corresponding eigenvectors

$w_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ ,  $w_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}$ . Since A is diagonalizable the vectors  $\{w_1, w_2\}$  forms a basis of  $\mathbf{R}^2$  and we can express  $v_0$  in this basis as  $v_0 = \alpha w_1 + \beta w_2$  for some  $\alpha, \beta \in \mathbf{R}$ , thus  $v_{n+1} = \alpha A^n w_1 + \beta A^n w_2 = \alpha \lambda_1^n w_1 + \beta \lambda_2^n w_2 = \alpha (\frac{1}{2})^n w_1 + \beta w_2$ . Then taking the  $n \rightarrow \infty$ , the first term goes to zero and  $v_{n+1} \sim \beta w_2$ . So asymptotically  $\frac{d_{n+1}}{w_{n+1}} \sim 3$  which verifies the initial condition:  $w_0 < d_0$ .

3. Function approximation) In this problem we will work in the real inner product space  $L^2[-1, 1]$  given by

$$L^2[-1, 1] = \left\{ f : [-1, 1] \rightarrow \mathbb{R} \mid \int_{-1}^1 f(x)^2 dx < \infty \right\}.$$

On this space, the inner product is given by

$$\langle f, g \rangle = \int_{-1}^1 f(x)g(x) dx.$$

In the following exercises, you may use a computer to perform the integral calculations.

- (a) The functions  $\{1, x, x^2\}$  form a basis for the 3-dimensional subspace  $P_2$  of  $L^2[-1, 1]$  consisting of the polynomials of degree at most 2. Give the orthonormal basis for  $P_2$  obtained by applying Gram-Schmidt to this set of functions.

Using Gram-Schmidt orthonormalization process, we find

$$\begin{aligned} v_1 &= 1 \\ v_2 &= x - \langle x, 1 \rangle \frac{1}{\langle 1, 1 \rangle} \\ &= x \\ v_3 &= x^2 - \langle x^2, v_2 \rangle \frac{v_2}{\langle v_2, v_2 \rangle} - \langle x^2, v_1 \rangle \frac{v_1}{\langle v_1, v_1 \rangle} \\ &= x^2 - \frac{1}{3} \end{aligned}$$

Then we normalize each of these vectors to obtain:

$$\begin{aligned}w_1 &= \frac{v_1}{\|v_1\|} = \frac{\sqrt{2}}{2} \\w_2 &= \frac{v_2}{\|v_2\|} = \sqrt{\frac{3}{2}} x \\w_3 &= \frac{v_3}{\|v_3\|} = \sqrt{\frac{45}{8}} \left(x^2 - \frac{1}{3}\right)\end{aligned}$$

(b) Compute the orthogonal projection of  $f(x) = \cos(\pi x/2)$  onto  $P_2$ .

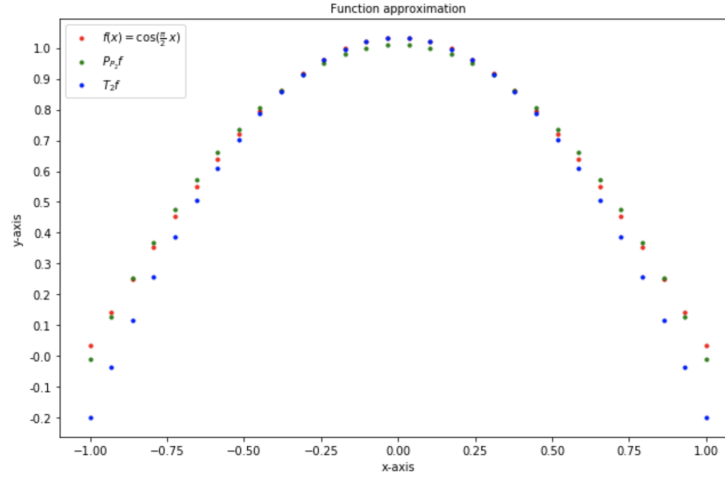
The projection of  $f(x) = \cos(\frac{\pi}{2} x)$  in the orthonormal basis  $\{w_1, w_2, w_3\}$  is:  $\sum_{i=1,3} \langle f, w_i \rangle w_i$ , where:

$$\begin{aligned}\langle f, w_1 \rangle &= \int_{-1}^1 \cos\left(\frac{\pi}{2} x\right) \frac{\sqrt{2}}{2} dx \\&= \frac{4}{\pi\sqrt{2}} \sim 0.9 \\ \langle f, w_2 \rangle &= \int_{-1}^1 \cos\left(\frac{\pi}{2} x\right) \frac{\sqrt{3}}{2} x dx \\&= 0 \\ \langle f, w_3 \rangle &= \int_{-1}^1 \cos\left(\frac{\pi}{2} x\right) \sqrt{\frac{45}{8}} \left(x^2 - \frac{1}{3}\right) dx \\&= 2\sqrt{10} \frac{\pi^2 - 12}{\pi^3} \sim -0.43\end{aligned}$$

(c) Plot  $f(x) = \cos(\pi x/2)$ ,  $\mathcal{P}_{P_2} f$ , and  $T_2 f$  on the same axis. Here  $\mathcal{P}_{P_2} f$  is the projection computed in the previous part, and  $T_2 f$  is the quadratic Taylor polynomial for  $f$  centered at  $x = 0$ :

$$T_2 f(x) = f(0) + f'(0)x + \frac{f''(0)}{2}x^2.$$

Include this plot in your submitted homework document.



- (d) The plot from the previous part shows that  $\mathcal{P}_{P_2}f$  is a better approximation than  $T_2f$  over most of  $[-1, 1]$ . Explain why this is the case.

$\mathcal{P}_{P_2}f$  is the orthogonal projection of  $f(x)$  over the subspace of polynomials of degree 2:  $\{w_1, w_2, w_3\}$ , like the Taylor expansion  $T_2f$ . The difference is that the Taylor polynomial is a polynomial expansion of  $f$  at 0. So in a neighborhood of 0, there is almost no differences between  $f$  and  $T_2f$ , but as we move away the approximation given by  $T_2f$  is worst than  $\mathcal{P}_{P_2}f$ .

#### 4. Scalar linear approximation

- (a) First we write  $E[(ax + b - y)^2] = E[((ax - y) - (-b))^2]$ , we know that the best mean-squared error minimizer of a random variable is its mean so  $-b = E[ax - y] = aE[x] - E[y] = a\mu_x - \mu_y$ . Substituting  $b$  in the expression we want to minimize gives us:

$$\begin{aligned} E[(ax + b - y)^2] &= E[(ax - y - (a\mu_x - \mu_y))^2] \\ &= E[\{a(\mu_x - x) - (y - \mu_y)\}^2] \\ &= a^2 E[(x - \mu_x)^2] + E[(y - \mu_y)^2] - 2aE[(x - \mu_x)(y - \mu_y)] \\ &= a^2 \sigma_x^2 + \sigma_y^2 - 2a \text{Cov}(x, y) \end{aligned}$$

Let  $f(a) = a^2 \sigma_x^2 + \sigma_y^2 - 2a \text{Cov}(x, y)$ , then  $f'(a) = 2(\sigma_x^2 a - \text{Cov}(x, y))$  and  $f''(a) = 2\sigma_x^2$ . The function is strictly convex, and its second derivative is positive, thus its minimizer is  $a = \frac{\text{Cov}(x, y)}{\sigma_x^2} = \rho_{x, y} \frac{\sigma_y}{\sigma_x}$ .

- (b) Applying the result from the previous question, the best linear estimate of  $y$  given  $x$  is  $y = \rho_{x, y} \frac{\sigma_y}{\sigma_x} (x - \mu_x) + \mu_y$ . Notice that  $\text{Var}(x) = \text{Var}(y \mid z) = E[y^2 \mid z^2] - E[y \mid z]^2 = E[y^2]E[z^2] - E[y]^2 E[z]^2 = (\sigma_y^2 + \mu_y^2)\sigma_z^2$  where we have used that  $a$  and  $z$  are independent and  $z$  has zero-mean. And  $E[x] = E[y \mid z] = E[y] = 0$ . Thus the best linear estimate of  $y$  given  $x$  is:  $\rho_{x, y} \frac{\sigma_y}{\sigma_z \sqrt{\sigma_y^2 + \mu_y^2}} x + \mu_y$ .

- (c) If in the expression above of  $y$  given  $x$ ,  $z$  is normally distributed with  $\sigma_z = 1$  then  $y$  is perfectly estimated from  $x$ .

## 5. Gradients

- (a) Compute the gradient of  $f(x) = b^T x$  where  $b \in \mathbf{R}^d$  and  $f : \mathbf{R}^d \rightarrow \mathbf{R}$ .  
 $\frac{\partial f(x)}{\partial x_j} = \sum_i b_i \frac{\partial x_i}{\partial x_j} = b_j$ , thus  $\nabla f(x) = b$ .
- (b) Compute the gradient of  $f(x) = x^T A x$  where  $A \in \mathbf{R}^{d \times d}$  and  $f : \mathbf{R}^d \rightarrow \mathbf{R}$ .  $f(x) = x^T A x = \sum_{i=1}^d \sum_{j=1}^d a_{ij} x_i x_j$ , then

$$\begin{aligned}
 \frac{\partial f}{\partial x_k} &= \sum_{i=1}^d \sum_{j=1}^d a_{ij} \frac{\partial x_i x_j}{\partial x_k} \\
 &= \sum_{i=1}^d \sum_{j=1}^d a_{ij} (x_j \delta_{ik} + x_i \delta_{jk}) \\
 &= \sum_{i=1}^d \sum_{j=1}^d a_{ij} x_j \delta_{ik} + \sum_{i=1}^d \sum_{j=1}^d a_{ij} x_i \delta_{jk} \\
 &= \sum_{j=1}^d a_{kj} x_j + \sum_{i=1}^d a_{ik} x_i \\
 &= (Ax)_k + (Ax)_k^T
 \end{aligned}$$

thus  $\nabla f(x) = (A + A^T)x$ .