

1 Problem 1

1.1 a

$$(Ax)^T b = x^T A^T b = x^T (A^T A \hat{x}) = (x^T A^T)(A \hat{x}) = (Ax)^T (A \hat{x})$$

1.2 b

Using part a, we can substitute x for \hat{x} , giving us $(A\hat{x})^T b = \hat{x}^T A^T b = \|A\hat{x}\|^2$. $\frac{\|A\hat{x}\|^2}{\|Ax\|} = \|Ax\| \|b\|$, and therefore, $\frac{\|A\hat{x}\|^2}{\|Ax\| \|b\|} = \|Ax\|$. This means the equality holds as the two sides are equivalent to $\frac{\|A\hat{x}\|^2}{\|Ax\| \|b\|}$.

1.3 c

Since we know that $\cos(\theta_{\hat{x}}) = \frac{(A\hat{x})^T b}{\|A\hat{x}\| \|b\|}$, where $\theta_{\hat{x}}$ is the angle between $A\hat{x}$ and b . The same applies for $\cos(\theta_x) = \frac{(Ax)^T b}{\|Ax\| \|b\|}$, where θ_x is the angle between Ax and b . We know that $\theta_{\hat{x}} < \theta_x$ because TODO

2 Problem 2

2.1 a

$$\hat{x} = (A^T A)^{-1} A^T b = ((QR)^T QR)^{-1} (QR)^T b = (R^T Q^T QR)^{-1} R^T Q^T b = R^{-1} Q^{-1} (Q^T)^{-1} (R^T)^{-1} R^T Q^T b = R^{-1} (Q^T Q)^{-1} Q^T b = R^{-1} Q^T b, \text{ so } A\hat{x} = QRR^{-1}Q^T b = QQ^T b$$

2.2 b

$A\hat{x} = A(A^T A)^{-1} A^T b = AA^{-1} (A^T)^{-1} A^T b = b$, so $\|A\hat{x} - b\|^2 = \|b\|^2 + \|A\hat{x}\|^2 - 2\|A\hat{x}\| \|b\| = \|b\|^2 + \|A\hat{x}\|^2 - 2\|A\hat{x}\|^2 = \|b\|^2 - \|A\hat{x}\|^2$. From part a, we know that this is equal to $\|b\|^2 - \|QQ^T b\|^2 = \|b\|^2 - \|Q^T b\|^2$.

3 Problem 3

3.1 a

The weighted least squares is equal to $\sum_{i=1}^m w_i (\hat{a}_i^T x - b_i)^2$. The one-norm for vector x is equal to $\sum_{i=1}^m |x_i| = 1$, so we can express the weighted sum as the one-norm squared of $\sum_{i=1}^m (w_i (\hat{a}_i^T x - b_i))^2$, meaning if we have some diagonal matrix D with entries $\sqrt{w_1}, \sqrt{w_2}, \dots, \sqrt{w_m}$, then this is equal to $\|D(Ax - b)\|^2$.

3.2 b

For matrix B where $B = DA$, if A has linearly independent columns, then since D is a diagonal matrix, it also has linearly independent columns no w_i will be equal to 0. Therefore, if $Ax = 0$, then x must be 0 as it is not in the column

space of A . This means B is the product of 2 linearly independent matrices, so if $Bx = 0$, $DAx = 0$, and since D is also linearly independent, x must be 0. This means that if $Bx = 0$ only if x is 0, then B must be linearly independent.

3.3 c

We said that $B = DA$, so we can use matrix B in this weighted least squares solution, meaning that since $d = Db$, $\hat{d} = (B^T B)^{-1} B^T d = ((DA)^T DA)^{-1} (DA)^T Db = (A^T D^T DA)^{-1} A^T D^T Db = (A^T \text{diag}(w) A)^{-1} A^T \text{diag}(w) b = (A^T \text{diag}(w) A)^{-1} A^T \text{diag}(w) b$

4 Problem 4

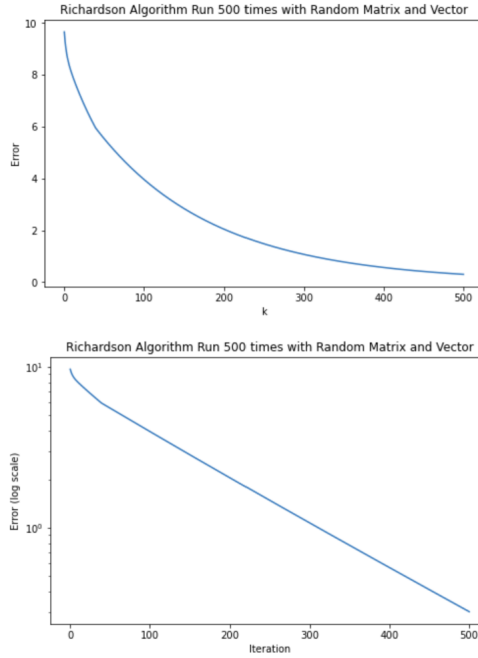
4.1 a

If $x^{k+1} = x^k$, then that means $\mu A^T (Ax^k - b) = 0$. \hat{x} is the solution of $A^T (Ax_k - b) = 0$. This means that $Ax_k = b$, but since all $x_{k+1} = x_k$, all $x_k = 0$. Therefore, $b = 0$, but since $\hat{x} = A^\dagger b$, $\hat{x} = 0$. Therefore, all $x_k = \hat{x}$.

4.2 b

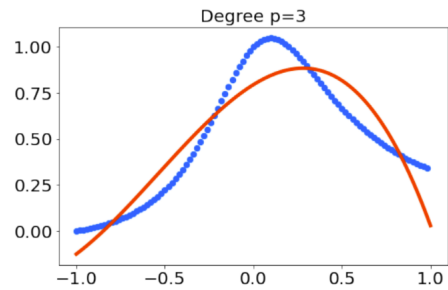
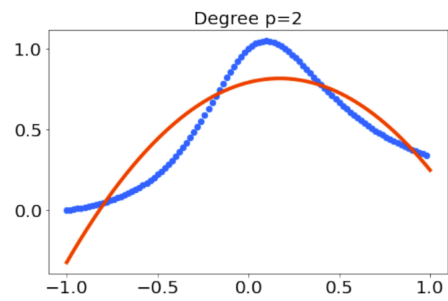
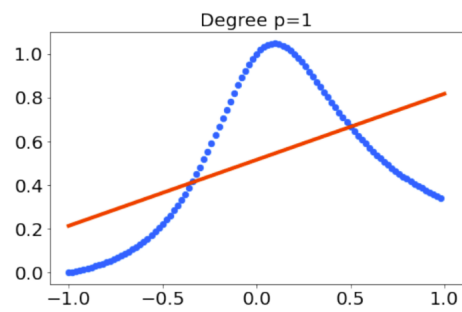
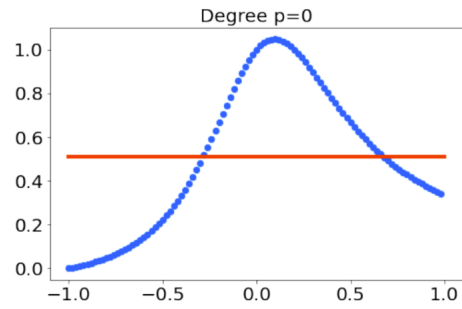
$x^{k+1} = x^k - \mu A^T (Ax^k - b) = x^k - \mu A^T Ax^k + \mu A^T b = x^k (I - \mu A^T A) + \mu A^T b$. Therefore, $x^{k+1} = Fx^k + g$ where $F = (I - \mu A^T A)$ and $g = \mu A^T b$.

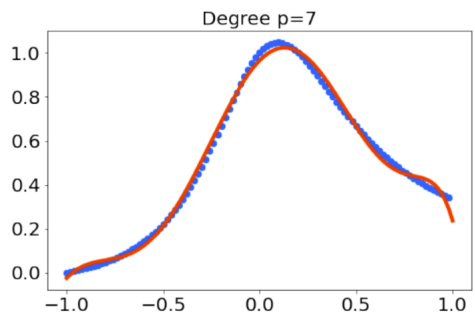
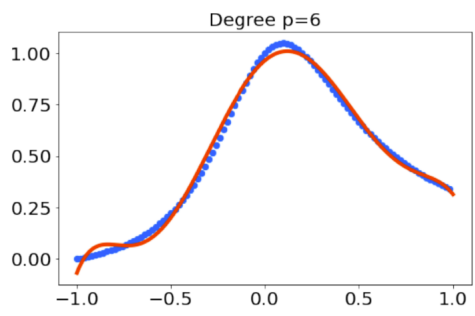
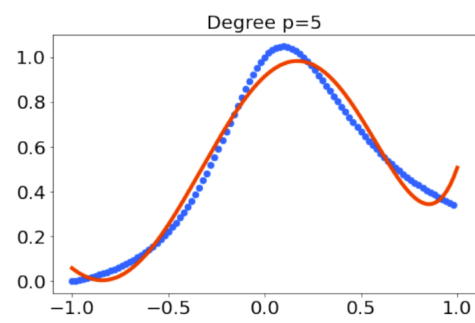
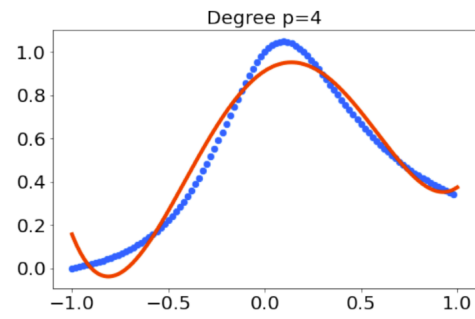
4.3 c

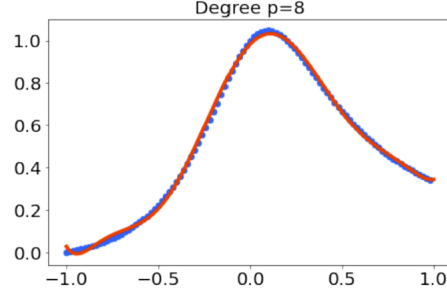


5 Problem 5

5.1 a

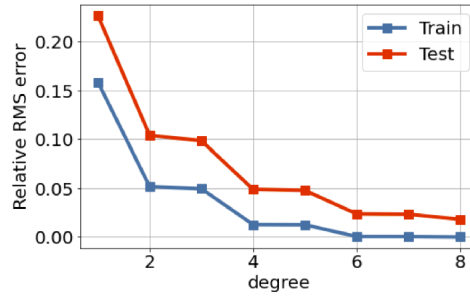






5.2 b

This tells us the reasonable degree fit is around 4-6, as the least difference in degree between train and test occurs among those values



6 Problem 6

$\|AX - Y\|^2 = \|X^T A^T - Y^T\|^2 = \|X^T W - Y^T\|^2$, meaning this is also equal to $\sum_{i=1}^m \|X^T w_i - y_i\|^2$, where w_i and y_i are the columns of the A^T and Y^T matrices, respectively. This means to minimize the summation, we need $X^T w_i = y_i$, or $w_i = (X^T)^{-1} y_i$. Since $W = A^T$, this means $A = W^T$, so $A = ((X^T)^{-1} Y^T)^T$. Since $X^{-1} X = I$, we can say $A = ((X^T)^{-1} X^{-1} X Y^T)^T = ((X X^T)^{-1} X Y^T)^T = Y X^T ((X X^T)^T)^{-1} = Y X^T (X^T X)^{-1} = Y X^\dagger$

7 Problem 7

7.1 a

If we let Y be a matrix of N rows and K columns that is populated with 1s and -1s. For any index Y_{nk} , if the value at $X^T B_n$ is in the group k , then $Y_{nk} = 1$. Otherwise, $Y_{nk} = -1$.

7.2 b

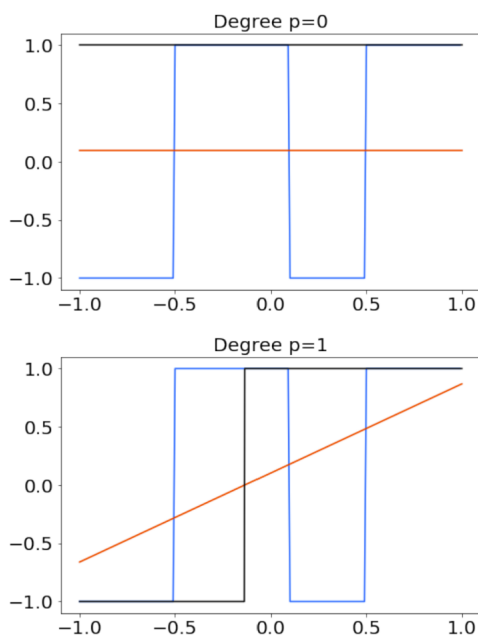
In order to minimize this, we need to set $Y = X^T \beta$, or ensure that for every column k in β and Y , $X^T \beta_k = Y_k$. By taking the norm of the solution, we see that this is then equal to $\|X^T \beta_1 - Y_1\|^2 + \|X^T \beta_2 - Y_2\|^2 \dots + \|X^T \beta_k - Y_k\|^2$.

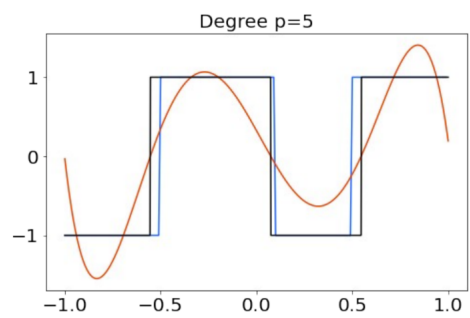
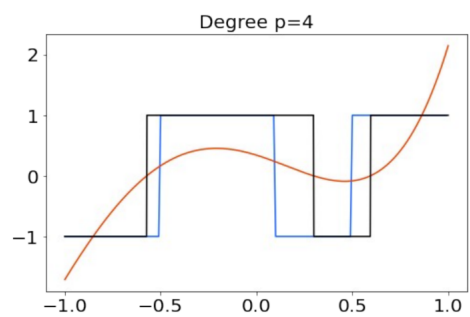
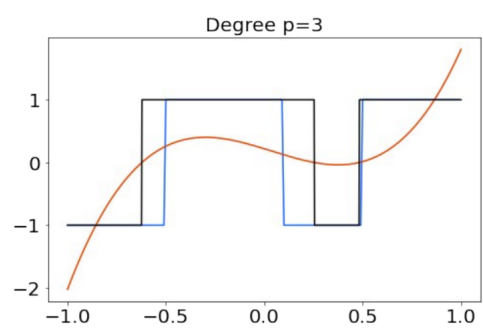
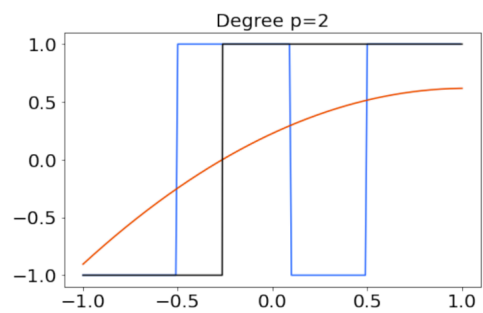
7.3 c

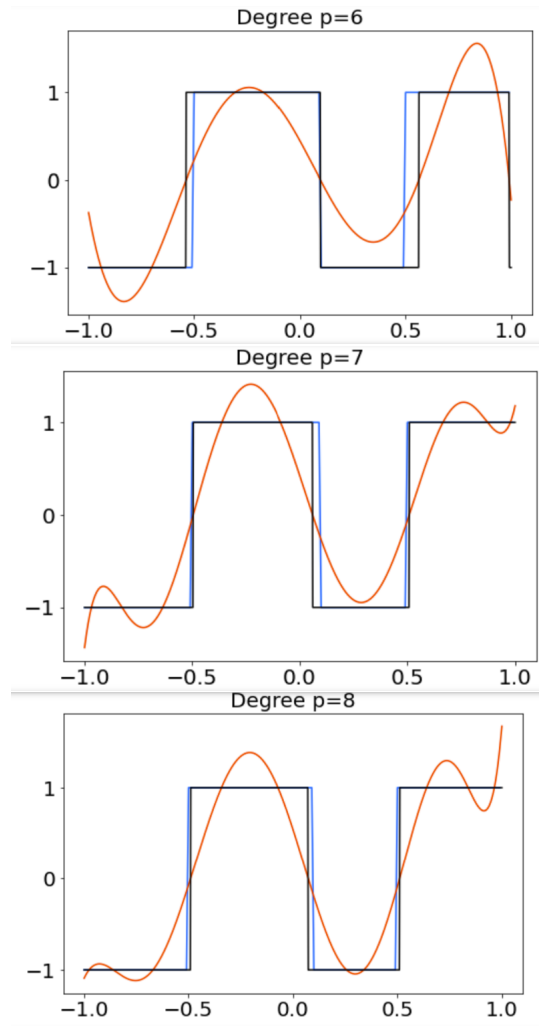
In order to minimize the least-squares equation, we need $Y = X^T \beta$, or $(X^T)^{-1} Y = \beta$. $(X^T)^\dagger Y = X(X^T X)^{-1} Y = X X^{-1} (X^T)^{-1} Y = (X^T)^{-1} Y$, so therefore, the two are equivalent.

8 Problem 8

8.1 a and b







8.2 c

The error rate decreases as the degree increases. This can be seen from the sign of the function matching the original function more and more as the degree increases.

The errors are:

Degree 0, Error: 0.0

Degree 1, Error: 55.000000000000001

Degree 2, Error: 55.000000000000001

Degree 3, Error: 49.5

Degree 4, Error: 55.000000000000001

Degree 5, Error: 55.000000000000001

Degree 6, Error: 55.00000000000001
Degree 7, Error: 55.00000000000001
Degree 8, Error: 55.00000000000001