

CS 446 / ECE 449 — Homework 4

sehsu2, dcku2 (partner)

Version 1.0

Instructions.

- Homework is due **Tuesday, April 4, at noon CST**; no late homework accepted.
- Everyone must submit individually at gradescope under **hw4** and **hw4code**.
- The “written” submission at **hw4** **must be typed**, and submitted in any format gradescope accepts (to be safe, submit a PDF). You may use L^AT_EX, Markdown, Google Docs, MS Word, whatever you like; but it must be typed!
- When submitting at **hw4**, gradescope will ask you to mark out boxes/pages around each of your answers; please do this precisely!
- Please make sure your NetID is clear and large on the first page of the homework.
- Your solution **must** be written in your own words. Please see the course webpage for full academic integrity information. Briefly, you may have high-level discussions with at most 3 classmates, whose NetIDs you should place on the first page of your solutions, and you should cite any external reference you use; despite all this, your solution must be written in your own words.
- We reserve the right to reduce the auto-graded score for **hw4code** if we detect funny business (e.g., your solution lacks any algorithm and hard-codes answers you obtained from someone else, or simply via trial-and-error with the autograder).
- Coding problems come with suggested “library routines”; we include these to reduce your time fishing around APIs, but you are free to use other APIs.
- When submitting to **hw4code**, only upload the two python files **hw4.py** and **hw4_utils.py**. Don’t upload a zip file or additional files.

Version history.

1. Initial version.

1. Linear Regression/SVD.

Throughout this problem let \mathbf{X} be the $n \times d$ matrix with the feature vectors $(\mathbf{x}_i)_{i=1}^n$ as its rows. Suppose we have the singular value decomposition $\mathbf{X} = \sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^\top$.

Solution.

- (a) The empirical risk minimization problem with squared loss is given by:

$$\mathbf{R}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^d \sum_{j=1}^{n_i} (\mathbf{w}^\top \mathbf{x}_i - y_{i_j})^2.$$

Substituting the training set, we get:

$$\mathbf{R}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^d \sum_{j=1}^{n_i} (\mathbf{w}^\top \mathbf{e}_i - y_{i_j})^2.$$

Setting the gradient with respect to \mathbf{w} equal to zero, we get:

$$\frac{\partial \mathbf{R}(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{n} \sum_{i=1}^d \sum_{j=1}^{n_i} 2(\mathbf{w}^\top \mathbf{e}_i - y_{i_j}) \mathbf{e}_i = 0.$$

To minimize the empirical risk, the gradient is equal to the zero vector as $\mathbf{w}^\top \mathbf{e}_i = w_i$:

$$\sum_{i=1}^d \sum_{j=1}^{n_i} (w_i - y_{i_j}) \mathbf{e}_i = 0.$$

For each $k \in (1, \dots, d)$ of the gradient, we obtain:

$$\sum_{j=1}^{n_k} (w_k - y_{k_j}) = 0.$$

The above equation is satisfied when:

$$w_k = \frac{1}{n_k} \sum_{j=1}^{n_k} y_{k_j}.$$

- (b) We will use the SVD and compute the empirical risk. Since $\mathbf{y} = \sum_{i=1}^r a_i \mathbf{u}_i$ and a_i are some random constants, the SVD form of \mathbf{X} is $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$, and the pseudoinverse of \mathbf{X} is obtained from $\mathbf{X}^+ = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^\top$:

$$\hat{\mathbf{w}}_{\text{ols}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^\top \sum_{i=1}^r a_i \mathbf{u}_i = \mathbf{V} \mathbf{\Sigma}^\top \sum_{i=1}^r a_i \mathbf{e}_i$$

$$\mathbf{X} \hat{\mathbf{w}}_{\text{ols}} = \mathbf{X} \mathbf{X}^+ \mathbf{y} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top \mathbf{V} \sum_{i=1}^r a_i \mathbf{e}_i = \mathbf{U} \sum_{i=1}^r a_i \mathbf{e}_i = \sum_{i=1}^r a_i \mathbf{u}_i$$

Since $\mathbf{X} \hat{\mathbf{w}}_{\text{ols}} = \mathbf{y}$, this implies that the empirical risk is zero when the label vector \mathbf{y} is a linear combination of the $\{\mathbf{u}_i\}_{i=1}^r$. On the other hand, the label vector \mathbf{y} is not a linear combination of the $\{\mathbf{u}_i\}_{i=1}^r$ if the empirical risk is nonzero. The empirical risk is shown as:

$$\mathbf{R}(\hat{\mathbf{w}}_{\text{ols}}) = \frac{1}{n} \|\mathbf{X} \hat{\mathbf{w}}_{\text{ols}} - \mathbf{y}\|^2$$

The above equation shows that the empirical risk must be nonzero because $\mathbf{P}_\mathbf{X} \mathbf{y} \neq \mathbf{y}$ and we know that $|\mathbf{P}_\mathbf{X} \mathbf{y} - \mathbf{y}|^2 > 0$.

- (c) If $\mathbf{X}^\top \mathbf{X}$ is invertible, all of its eigenvalues are nonzero. This implies that all the singular values of \mathbf{X} are nonzero. Since the number of nonzero singular values equals the rank of \mathbf{X} , the rank of \mathbf{X} is d . Therefore, $(\mathbf{x}_i)_{i=1}^n$ spans \mathbb{R}^d . If $(\mathbf{x}_i)_{i=1}^n$ spans \mathbb{R}^d , the rank of \mathbf{X} is d . This means that \mathbf{X} has d nonzero singular values. Since the squared singular values are the eigenvalues of $\mathbf{X}^\top \mathbf{X}$, all eigenvalues of $\mathbf{X}^\top \mathbf{X}$ are nonzero. Thus, $\mathbf{X}^\top \mathbf{X}$ is invertible. Consequently, $\mathbf{X}^\top \mathbf{X}$ is invertible if and only if $(\mathbf{x}_i)_{i=1}^n$ spans \mathbb{R}^d . This characterizes when linear regression has a unique solution due to the normal equation.
- (d) Let's consider a 3×2 matrix \mathbf{X} given by:

$$\mathbf{X} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

Now, let's compute $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{X} \mathbf{X}^\top$:

$$\mathbf{X}^\top \mathbf{X} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mathbf{X} \mathbf{X}^\top = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

To check whether a matrix is invertible or not, we calculate its determinant to see if it is nonzero. If the determinant is nonzero, then it is invertible. Otherwise, it is not invertible.

$$\det(\mathbf{X}^\top \mathbf{X}) = 1$$

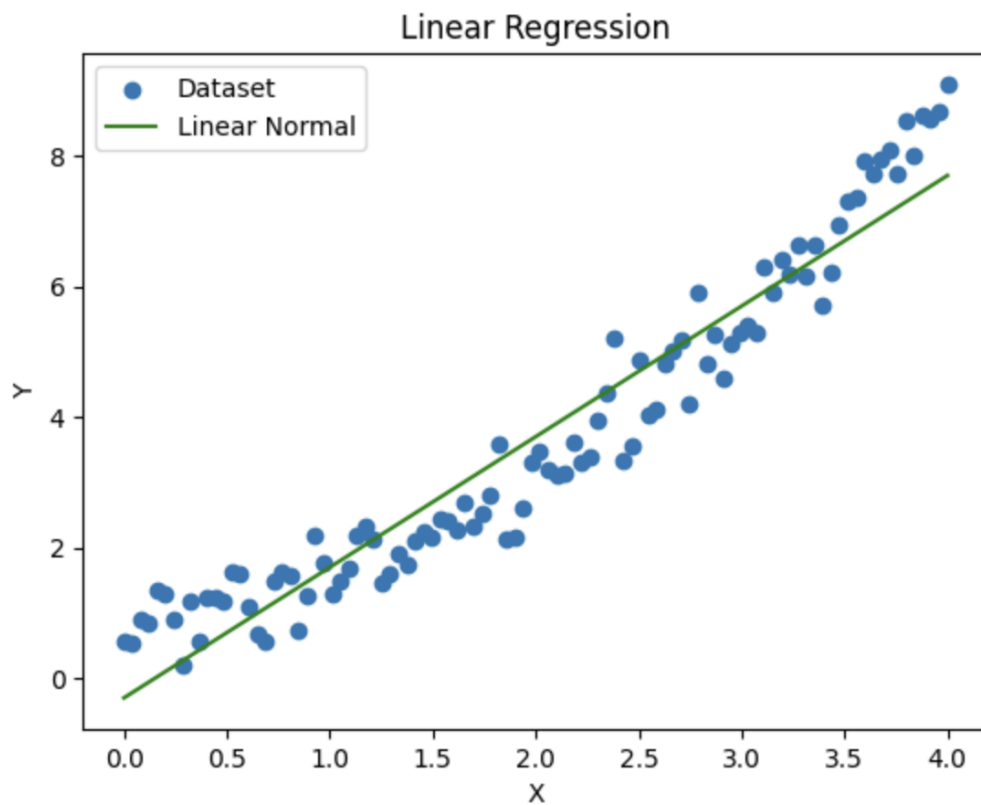
$$\det(\mathbf{X} \mathbf{X}^\top) = 0$$

The above equations show that $\mathbf{X}^\top \mathbf{X}$ is invertible and $\mathbf{X} \mathbf{X}^\top$ is not invertible.

2. Linear Regression.

Recall that the empirical risk in the linear regression method is defined as $\hat{\mathcal{R}}(\mathbf{w}) := \frac{1}{2n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$, where $\mathbf{x}_i \in \mathbb{R}^d$ is a data point and y_i is an associated label.

(c)



3. Polynomial Regression.

In Problem 3 you constructed a linear model $\mathbf{w}^\top \mathbf{x} = \sum_{i=1}^d x_i w_i$. In this problem you will use the same setup as in the previous problem, but enhance your linear model by doing a quadratic expansion of the features. Namely, you will construct a new linear model $f_{\mathbf{w}}$ with parameters

$$(w_0, w_{01}, \dots, w_{0d}, w_{11}, w_{12}, \dots, w_{1d}, w_{22}, w_{23}, \dots, w_{2d}, \dots, w_{dd})^\top,$$

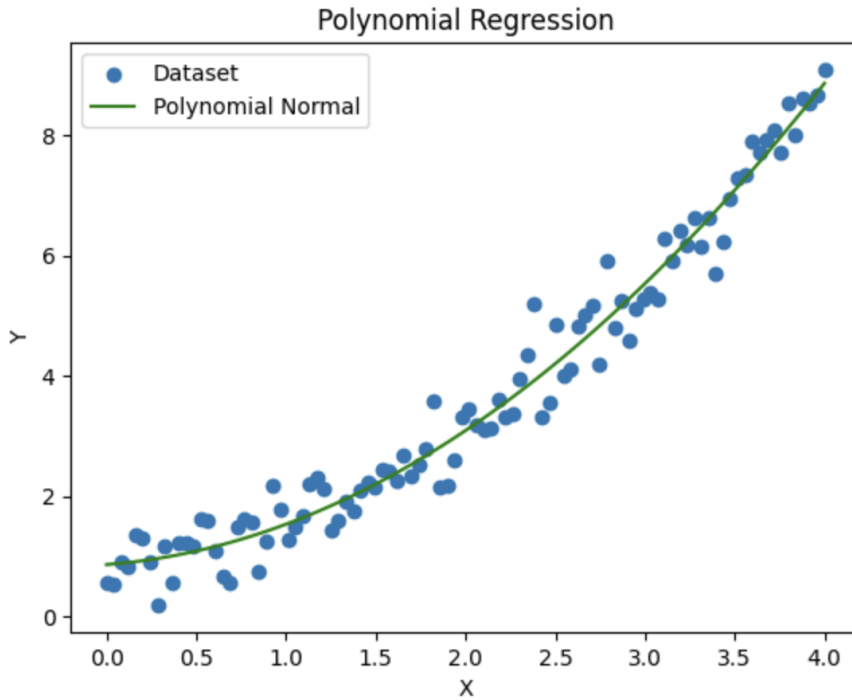
defined by

$$f_{\mathbf{w}}(x) = \mathbf{w}^\top \phi(\mathbf{x}) = w_0 + \sum_{i=1}^d w_{0i} x_i + \sum_{i \leq j}^d w_{ij} x_i x_j.$$

(a)

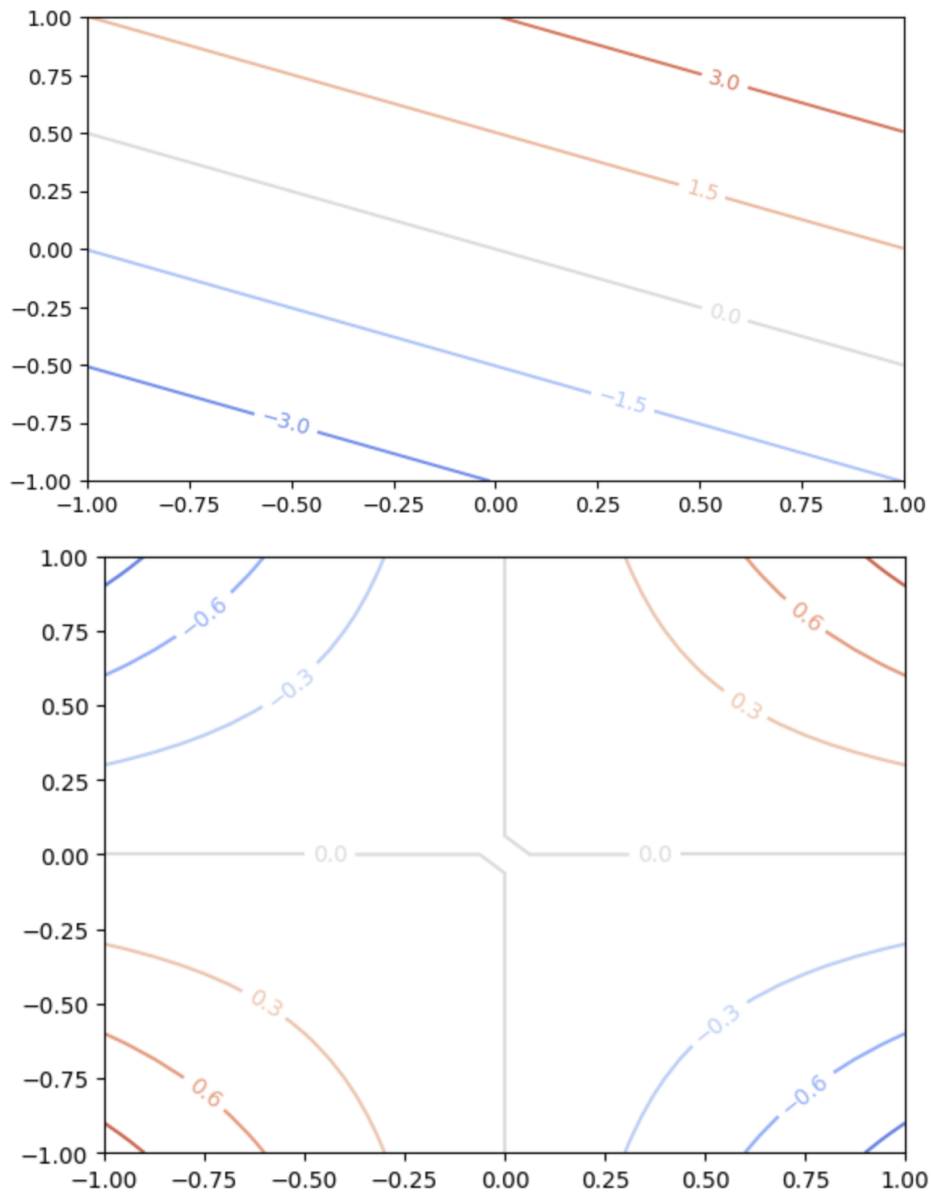
$$\phi(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_3 \end{bmatrix}$$

(d)



The polynomial model approximates the data better compared to the linear model. According to the polynomial regression figure, the polynomial normal seems to fit the dataset better.

(e)



According to the figures, a linear model (top figure) cannot correctly classify all points, as the problem is not linearly separable. A polynomial model (bottom figure), such as a second-degree polynomial model, can correctly classify all points because it can create a non-linear decision boundary.