

Introduction to Machine Learning (67577)

Exercise 2 - Linear Regression

Dor Roter

208772251

April 28, 2021

1 Theoretical Questions

1.1 Solutions of the Normal Equations

1. Let $X \in \mathbb{R}^{m \times d}$ be a matrix, then

$$\text{Ker}(X) = \{v \in \mathbb{R}^d \mid Xv = 0\} = \{v \in \mathbb{R}^d \mid X^T Xv = X^T \cdot 0\} = \{v \in \mathbb{R}^d \mid X^T Xv = 0\} = \text{Ker}(X^T X) \quad \square$$

2. Let $A \in \mathbb{R}^{n \times n}$.

firstly, let us show that for all $v \in \text{Im}(A^T)$ it holds that $v \in \text{Ker}(A)^\perp$:

$$v \in \text{Im}(A^T) \Rightarrow v = A^T x$$

$$\Rightarrow \forall u \in \text{Ker}(A) \quad \langle v | u \rangle = v^T \cdot u = (A^T x)^T \cdot u = x^T A \cdot u = x^T \cdot Au = x^T \cdot 0 = 0$$

$$\text{Therefore: } v \in \text{Im}(A^T) \Rightarrow \langle v | x \rangle = 0 \quad \forall x \in \text{Ker}(A) \Leftrightarrow v \in \text{Ker}(A)^\perp$$

as such, it holds that $\text{Im}(A^T) \subseteq \text{Ker}(A)^\perp$.

let us show that for all $v \in \text{Ker}(A)^\perp$ it holds that $v \in \text{Im}(A^T)$.

now, we assume by contradiction that $v \notin \text{Im}(A^T)$ and we will show that $v \notin \text{Ker}(A)^\perp$.

$$v \notin \text{Im}(A^T) \Rightarrow v \in \text{Im}(A^T)^\perp$$

$$\text{therefore: } \forall u \in \mathbb{R}^n \quad v^T \cdot (A^T u) = (Av) \cdot u = 0$$

as such, it holds that Av must be 0, which means $v \in \text{Ker}(A)$.

but, $v \in \text{Ker}(A) \Rightarrow v \notin \text{Ker}(A)^\perp$.

Therefore, it holds that $\text{Im}(A^T) \supseteq \text{Ker}(A)^\perp$.

$$\Rightarrow \text{Im}(A^T) = \text{Ker}(A)^\perp \quad \square$$

3. Let $y = Xw$ be a non-homogeneous system of linear equations, with X as square not invertible matrix.

As X is not invertible, the equations system could not have a single solution, Therefore there are either 0, or infinitely many solutions for the provided linear equations set.

There might be a solution to the equations iff $y \in Im(X)$ (could be written as a linear combination of X using the scalars specified in w).

Therefore:

$$y \in Im(X) \underset{q.2}{=} Ker(X^T)^\perp$$

$$\Rightarrow y \perp Ker(X^T) \quad \square$$

4. Let $X^T X w = X^T y$ be a normal linear system.

We will address the two following possibilities for $X^T X$ separately.

- (a) $X^T X$ **is invertible**: then $X^T X w = X^T y \Leftrightarrow w = (X^T X)^{-1} X^T y$ and we have found a single unique solution for the linear system.
- (b) $X^T X$ **is not invertible**: since $X^T X$ is a square matrix, by the lemma proved in q.3:

$$X^T y \perp Ker((X^T X)^T) = Ker(X^T X) \Leftrightarrow \text{"there are } \infty \text{ solutions for } X^T X w = X^T y \text{"}$$

Let us show $X^T y \perp Ker(X^T X)$:

$$\forall v \in Ker(X^T X) \quad \langle X^T y | v \rangle = (X^T y)^T v = y^T X v$$

since we have proved in q.1 that $Ker(X) = Ker(X^T X)$, it holds that $v \in Ker(X^T X) \Leftrightarrow v \in Ker(X) \Leftrightarrow X v = 0$ and thus:

$$\underbrace{\forall v \in Ker(X^T X) \quad \langle X^T y | v \rangle = y^T X v = y^T \cdot 0 = 0}_{\Downarrow X^T y \perp Ker(X^T X)}$$

1.2 Projection Matrices

5. Let $V \subseteq \mathbb{R}^d$, $dim(V) = k$, v_1, \dots, v_k be orthonormal basis of V , and $P = \sum_{l=1}^k v_l v_l^T$.

- (a) by P's definition $P_{i,j} = \sum_{l=1}^k (v_l)_i \cdot (v_l)_j = \sum_{l=1}^k (v_l)_j \cdot (v_l)_i = P_{j,i}$ and therefore P 's symmetric.
- (b) Since P is a symmetric matrix, it holds that P has an EVD decomposition $P = U D U^T$ where U is an orthogonal matrix and D is a diagonal matrix.

Therefore, $P^2 = U D U^T \cdot U D U^T = U D^2 U^T$.

By the claim proved in (d), $P^2 = P$, and so $U D^2 U^T = P^2 = P = U D U^T$.

As $diag(\lambda_1^2, \dots, \lambda_k^2) = D^2 = D = diag(\lambda_1, \dots, \lambda_k)$, it holds that for each $i \in [k]$ $\lambda_i = \lambda_i^2 \Leftrightarrow \lambda_i \in \{0, 1\}$.

Finally for all $i \in [k]$ it holds that:

$$P v_i = \left(\sum_{l=1}^k v_l v_l^T \right) \cdot v_i = \sum_{l=1}^k (v_l v_l^T \cdot v_i) = \sum_{l=1}^k v_l (v_l^T \cdot v_i) = \sum_{l=1}^k v_l \langle v_l | v_i \rangle \underbrace{=}_{\substack{v_1, \dots, v_k \text{ are orthonormal} \\ \langle v_l | v_i \rangle = \begin{cases} 1 & l = i \\ 0 & l \neq i \end{cases}}} \sum_{l=1}^k v_l \cdot \delta_{i,l} = v_i$$

and therefore v_1, \dots, v_k are the the eigenvectors corresponding to the eigenvalue 1.

- (c) Let $v \in V$, therefore there is a linear combination of the orthonormal basis of V , v_1, \dots, v_k and $\alpha_1, \dots, \alpha_k \in \mathbb{R}$ such as $v = \sum_{l=1}^k \alpha_l \cdot v_l$.

Therefore:

$$Pv = \left(\sum_{i=1}^k v_i v_i^T \right) \left(\sum_{j=1}^k \alpha_j v_j \right) = \sum_{i=1}^k \sum_{j=1}^k (v_i v_i^T \cdot \alpha_j v_j) = \sum_{i=1}^k \sum_{j=1}^k \alpha_j v_i \langle v_i | v_j \rangle = \sum_{i=1}^k \sum_{j=1}^k \alpha_j v_i \cdot \delta_{i,j} = \sum_{i=1}^k \alpha_i v_i = v$$

(d) Using P 's definition it holds that:

$$P^2 = \left(\sum_{i=1}^k v_i v_i^T \right) \cdot \left(\sum_{j=1}^k v_j v_j^T \right) = \sum_{i=1}^k \sum_{j=1}^k v_i v_i^T v_j v_j^T = \sum_{i=1}^k \sum_{j=1}^k v_i \langle v_i | v_j \rangle v_j^T = \sum_{i=1}^k \sum_{j=1}^k v_i \cdot \delta_{i,j} \cdot v_j^T = \sum_{i=1}^k v_i \cdot v_i^T = P$$

(e) Using the pervious lemmas:

$$(I - P)P = P - P^2 \stackrel{(d)}{=} P - P = 0$$

1.3 Least Squares

6. Let $X = U\Sigma V^T$ be the SVD decomposition of X .

Firstly, we will show $(X^T X)^{-1} = V(\Sigma^T \Sigma)^{-1} V^T$:

$$\begin{aligned} X^T X &= V \Sigma^T U^T \cdot U \Sigma V^T = V \Sigma^T \Sigma V^T \\ \Rightarrow \begin{cases} X^T X \cdot V(\Sigma^T \Sigma)^{-1} V^T = V \Sigma^T \underbrace{\Sigma V^T \cdot V}_{I} (\Sigma^T \Sigma)^{-1} V^T = V \underbrace{(\Sigma^T \Sigma) \cdot (\Sigma^T \Sigma)^{-1}}_I V^T = V \cdot V^T = I \\ V(\Sigma^T \Sigma)^{-1} V^T \cdot X^T X = V(\Sigma^T \Sigma)^{-1} \underbrace{V^T \cdot V}_{I} \Sigma^T \Sigma V^T = V \underbrace{(\Sigma^T \Sigma)^{-1} \cdot (\Sigma^T \Sigma)}_I V^T = V \cdot V^T = I \end{cases} \end{aligned}$$

And so

$$\hat{w} = (X^T X)^{-1} X^T y = V(\Sigma^T \Sigma)^{-1} \underbrace{V^T \cdot V}_{I} \Sigma^T U^T \cdot y = V(\Sigma^T \Sigma)^{-1} \Sigma^T U^T \cdot y$$

We will demonstrate now that $(\Sigma^T \Sigma)^{-1} \Sigma^T = \Sigma^\dagger$ and by so prove the general solution from the recitation equals to the one seen in class.

Since Σ is a diagonal matrix, $(\Sigma^T \Sigma)^{-1}_{i,j} = \begin{cases} \sigma_i^{-2} & \sigma_i \neq 0 \\ 0 & \sigma_i = 0 \end{cases}$ and so, $\left((\Sigma^T \Sigma)^{-1} \Sigma^T \right)_{i,j} = \begin{cases} \sigma_i^{-2} \sigma_i = \sigma_i^{-1} & \sigma_i \neq 0 \\ 0 \cdot \sigma_i = 0 & \sigma_i = 0 \end{cases} = \Sigma^\dagger$.

Therefore, $\hat{w} = (X^T X)^{-1} X^T y = V(\Sigma^T \Sigma)^{-1} \Sigma^T U^T \cdot y = V \Sigma^\dagger U^T \cdot y = X^\dagger y \square$

7. let $X \in \mathbb{R}^{m \times d}$, then $T_X : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is the respective linear transformation represented by X .

for $X^T X \in \mathbb{R}^{d \times d}$ it holds that:

$$X^T X \in \mathbb{R}^{d \times d} \text{ is invertible} \Leftrightarrow \text{Ker}(X^T X) = 0 \stackrel{q.1}{\Leftrightarrow} \text{Ker}(X) = 0 \Leftrightarrow \underbrace{\text{rank}(X) = \dim(\mathbb{R}^d)}_{\text{rank-nullity theorem}} - \text{Ker}(X) = d$$

And therefore:

$$\begin{aligned} \text{rank}(X) &= \dim(\text{Span}(x_1, \dots, x_m)) = d \\ \Leftrightarrow \text{Span}(x_1, \dots, x_m) &= \mathbb{R}^d \end{aligned}$$

As stated $X^T X \in \mathbb{R}^{d \times d}$ is invertible $\Leftrightarrow \text{Span}(x_1, \dots, x_m) = \mathbb{R}^d$. \square

8. Let us define $\hat{w} = X^\dagger y$, and let $\bar{w} \in \mathbb{R}^d$ be another solution of the linear equations set $X^T X w = X^T y$.

We have seen in recitation that $X = U \Sigma V^T = \begin{bmatrix} U_{\mathcal{R}} & U_{\mathcal{N}} \end{bmatrix} \begin{bmatrix} \mathcal{S} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_{\mathcal{R}}^T \\ V_{\mathcal{N}}^T \end{bmatrix} = U_{\mathcal{R}} \mathcal{S} V_{\mathcal{R}}^T = \tilde{U} \tilde{\Sigma} \tilde{V}^T$ where \mathcal{S} is a diagonal invertible $r \times r$ matrix.

Therefore, each w that solves $X^T X w = X^T y$ must be of form $w = \tilde{V} \tilde{\Sigma}^{-1} \hat{U}^T y$ where $\tilde{V} \tilde{\Sigma}^{-1} \hat{U}^T$ could each be padded into the original SVD decomposition dimentions.

Let us note that under this notation, $\hat{w} = V \Sigma^\dagger U^T y$ and $\bar{w} = \begin{bmatrix} V_{\mathcal{R}} & V_{\mathcal{N}} \end{bmatrix} \begin{bmatrix} \mathcal{S}^{-1} & & & 0 \\ & \sigma_{r+1} & & 0 \\ & & \ddots & 0 \\ & & & \sigma_d & 0 \end{bmatrix} \begin{bmatrix} U_{\mathcal{R}}^T \\ U_{\mathcal{N}}^T \end{bmatrix}$,

therefore also $\forall i \in [r] \hat{w}_i = \bar{w}_i$ and it holds that:

$$\|\bar{w}\|^2 \stackrel{\text{pythagorean theorem}}{=} \sum_{i=1}^d \bar{w}_i^2 = \sum_{i=1}^r \bar{w}_i^2 + \sum_{i=r+1}^d \bar{w}_i^2 \geq \sum_{i=1}^r \bar{w}_i^2 + \sum_{i=r+1}^d 0 = \sum_{i=1}^r \hat{w}_i^2 + \sum_{i=r+1}^d 0 = \|\hat{w}\|^2$$

$\Rightarrow \|\hat{w}\| \leq \|\bar{w}\|$ for each \bar{w} solution of $X^T X w = X^T y$. \square

2 Practical Questions

13. First, we can note the categorical features of the dataset include the following features: *waterfront*, *view*, *condition*, *grade*, *zipcode*, *lat* and *long*. While *waterfront*, *view*, *condition*, *grade* are simple boolean values or ordinal values, they are already already encoded as a needed for us, thus it is left to us to decide on the proper handling of *zipcode*, *lat* and *long*.

In the preprocessing process I have also clustered some non-categorical features, such as *yr_built* and *yr_renovated*, and transformed them into categorical dummy values in order to allow the model to better deal with their illinearity.

Next, I have also broken up "*sqft_basement*" into another binary column stating whether a basement exists, and using the passed time between the last renovation/building I've stated in "*is_new*" whether the building is relatively new (providing the linear model with the ability to take into account the combination of "*yr_built*" and "*yr_renovated*").

Lastly, in order to allow the model to "understand" geolocation, I have tupled a rough estimation of the latitude and longitude data into a single column, "*location*" which is in and of itself a categorical feature of an area (roughly 10km in diameter).

I have then encoded those new "*zipcode*", "*location*", "*yr_built*", "*yr_renovated*" features using one-hot encoding.

14. Scree-plot of non-categorical design matrix:

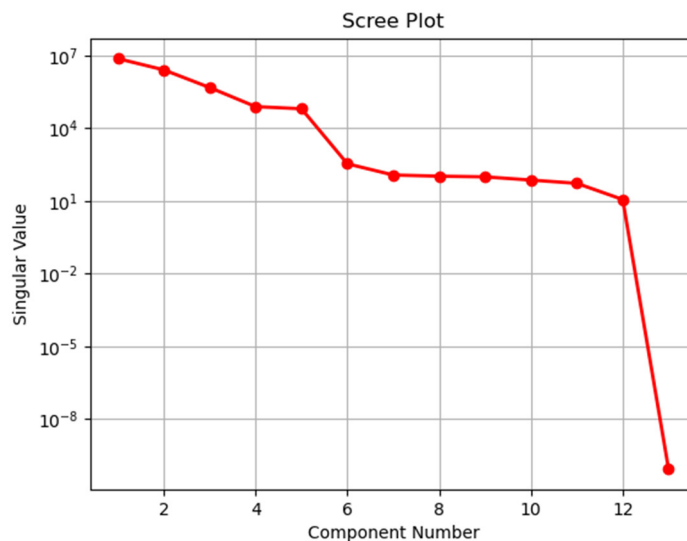


Figure 1: Scree-Plot of the design matrix

15. We note the data seems to have some close to linearly dependent features since part of the singular values are close to zero. This indicated the model would probably function well using a subset of the features available to us, and might generalize better if we do so (using just 5-10 features).

Furthermore, the high singular values are corresponding to a singular-vector which, as we have proven, is one of the features vectors, and thus those values correspond to important features for the model's prediction.

16. Plotting the training progress as a factor of p - the percent of the training dataset used we note that the MSE tends to be higher, and less stable when a smaller portion of the dataset is used to fit the model, and as p 's value surpasses 20% we can see the MSE settling towards its minimum point, with only minor improvements to the models accuracy.

Therefore, we can conclude that fitting this model requires a minimum of around 3.5 thousand samples in order to produce a somewhat credible estimator.

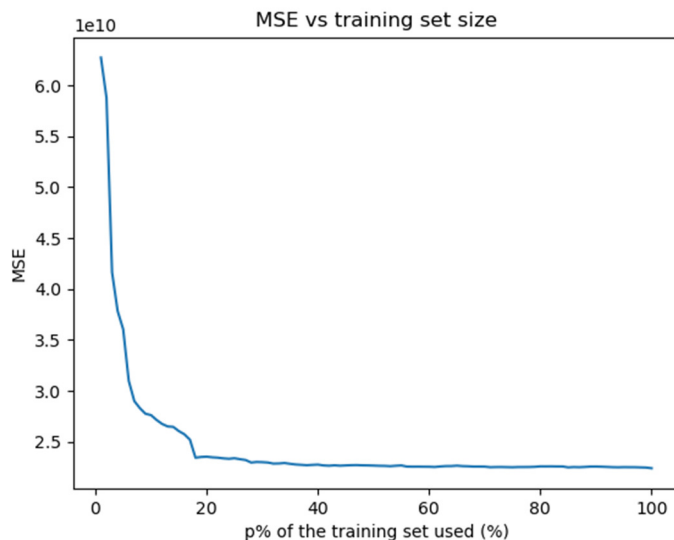


Figure 2: MSE plot by training set size

17. We notice by looking at the correlation coefficient of each feature in relation to the response vector that "*sqft_living*" and "*grade*" columns, with respective coefficients of 0.68, and 0.66, seem very beneficial to the model as their linear correlation with the response vector, described by their pearson correlation is close to 1, thus both rise and fall in unison with the response vector. This hypothesis is backed by the plots of both features against the response vector showing a rather clear linear relation.

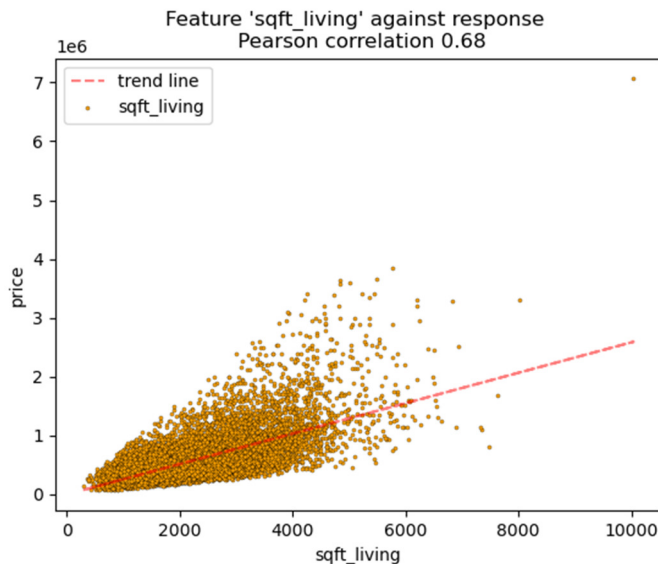


Figure 3: "*sqft_living*" vs "*price*"

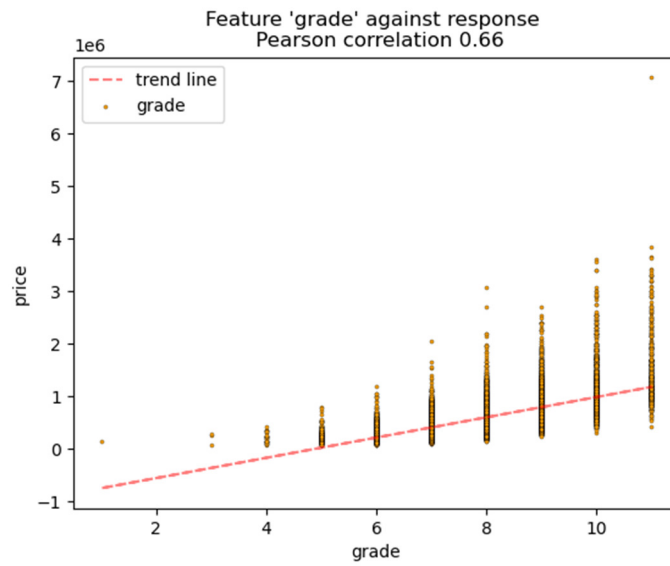


Figure 4: "grade" vs "price"

As opposed to these columns, the "condition" column has a pearson correlation coefficient of only 0.05, thus indicating there's not much of a linear relationship between the "condition" column and the response vector, and as such, this column seems less beneficial for the linear regression model's prediction

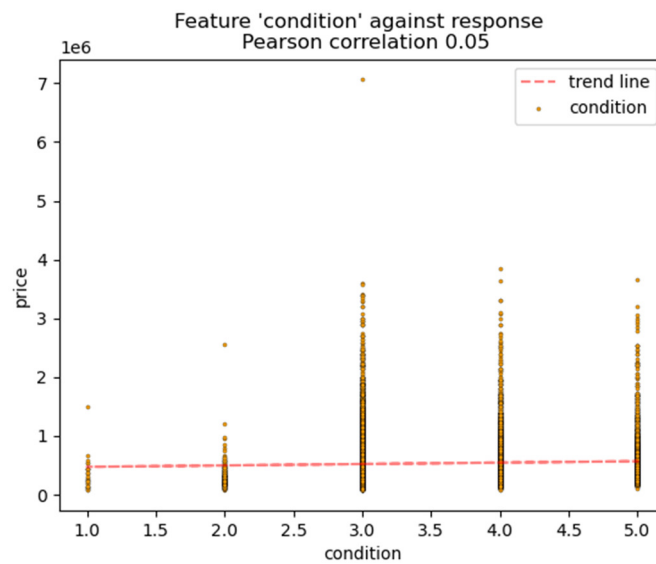


Figure 5: "condition" vs "price"