

Geometric Foundations of Data Analysis I

Joshua Maglione

September 23, 2024

Contents

1	Introduction	1
2	Least squares fitting	2
2.1	Build up	2
2.2	Line of best fit	3
2.3	In class exercises pt. I	4
2.4	Plane of best fit	5
2.5	Hyperplane of best fit	7
2.6	Why Equation (2.4) works	7
2.7	In class exercises pt. II	8
2.8	Nonlinear fittings	8
2.9	Coefficient of determination (R^2 values)	9
2.10	In class exercises pt. III	11

1 Introduction

Data analysis and more broadly Data Science is a vast and important field within Computer Science and Mathematics. At the core, the goal is to make sense of data, which can be measurements, survey results, behavior patterns, etc. Often this data comes to us in a very “high dimension”. That is, there are so many variables that it is impossible to visualize, and even in low dimensions, it may not be clear what the best conclusion is based on the analysis.

A few references seem to agree that the *total data* on all computers is something like 10^{23} bytes or about 100 zettabytes. While all of this data is not concentrated in one organization, we still require highly sophisticated tools to make sense of a huge amount of data. My goal with this course is that you will have a solid foundation with some standard tools. From this bedrock one could explore more sophisticated methods of data analysis more easily.

We will consider four key topics:

1. Least Squares Fitting,

2. Principal Component Analysis,
3. k -Means and Hierarchical Clustering,
4. Nearest Neighbors and the Johnson–Lindenstrauss Theorem

We will be working with the assumption that **the data we care about is preserved by orthogonal and linear transformations**. This is not true with all data—for example, one should not take (proper) linear combinations of people. However, for data like grams of different kinds of food, this is completely plausible. This assumption will not always be necessary, but we will just keep this in mind.

Half of our time will be spent bringing these ideas to life and getting our hands dirty. We will be working with Jupyter Notebooks to build familiarity with the concepts we will discuss. This will be done using Python and standard data analysis packages like Pandas.

2 Least squares fitting

This method of analysis is both simple and powerful. Like with most things in mathematics, without some good guiding examples, we can get lost in the formulas.

2.1 Build up

Suppose we have the following data as seen in Figure 2.1. (Maybe a company produces parts once per month in lots that vary in size depending on demand. We write i for the production run, x_i for the lot size, and y_i for the hours of labor.)

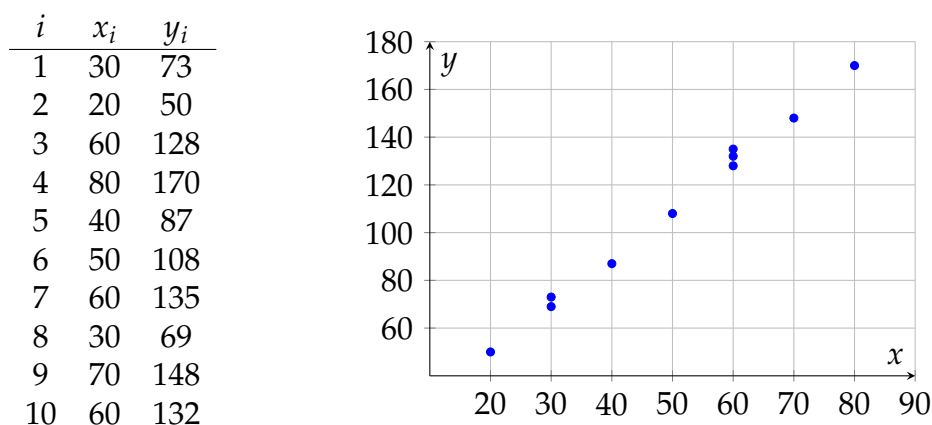


Figure 2.1: Data points exhibiting a linear relationship.

It seems clear from the plot that the data fits a geometric pattern—there is a linear phenomenon. If we try to find the line that is somehow closest to all the data points, we might draw something like in Figure 2.2.

Although the line is not a perfect fit, it seems to tell us something about the relationship between the lots size and the amount of hours.

Questions:

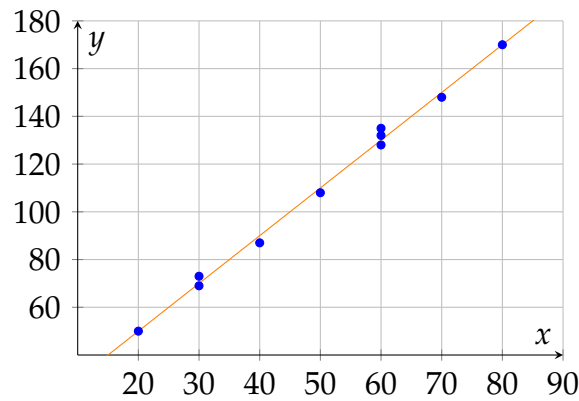


Figure 2.2: Line of best fit with data points.

- What makes this line “better” than alternatives?
- How are we quantifying “better”?
- Why are we using a line?

We will answer the first question later, probably. Let us consider the question of quantifying “better”. Least squares fitting is all about minimizing the squares of differences between the line and the actual data points. We will make this precise very soon.

2.2 Line of best fit

We know that the equation of a non-vertical line has the form

$$y = b_0 + b_1x.$$

If we had n data points of the form (x_i, y_i) , we could choose b_0 and b_1 to minimize the following sum

$$S(b_0, b_1) = \sum_{i=1}^n (y_i - (b_0 + b_1x_i))^2.$$

We can even solve for these values. Since S is a function in terms of b_0 and b_1 , all possible minima occur when the partial derivatives of S are 0. In other words, the minima arise as values (b_0, b_1) such that

$$\begin{aligned} \frac{\partial S}{\partial b_0} &= -2 \sum_{i=1}^n (y_i - (b_0 + b_1x_i)) = 0, \\ \frac{\partial S}{\partial b_1} &= -2 \sum_{i=1}^n (x_i y_i - x_i(b_0 + b_1x_i)) = 0. \end{aligned}$$

These equations are linear equations, so we can solve for these with techniques from linear algebra. This means, we need to solve two equations in the unknown

b_0 and b_1 :

$$\begin{aligned} nb_0 + b_1 \sum x_i &= \sum y_i, \\ b_0 \sum x_i + b_1 \sum x_i^2 &= \sum x_i y_i. \end{aligned} \quad (2.1)$$

Using the data from our example in Section 2.1, we have

$$\sum x_i = 500, \quad \sum y_i = 1100, \quad \sum x_i^2 = 28400, \quad \sum x_i y_i = 61800.$$

Thus, the equations we need to solve are

$$\begin{aligned} 10b_0 + 500b_1 &= 1100, \\ 500b_0 + 28400b_1 &= 61800, \end{aligned}$$

which yield $b_0 = 10$ and $b_1 = 2$. Going back to the context of the initial problem: this solution tells us that by increasing the lot size by one, we expect to increase the labor hours by two.

2.2.1 Written as matrices

Let us write the equations in (2.1) with matrices. This might seem like overkill at this stage, but it will set us up nicely to generalize. Let

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \quad Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad B = \begin{pmatrix} b_0 \\ b_1 \end{pmatrix}. \quad (2.2)$$

Therefore, the equations in (2.1) are equivalent to the single matrix equation:

$$X^t X B = X^t Y. \quad (2.3)$$

If $X^t X$ is invertible, then $B = (X^t X)^{-1} X^t Y$.

2.3 In class exercises pt. I

1. (a) With X , Y , and B as defined in Equation (2.2), show that

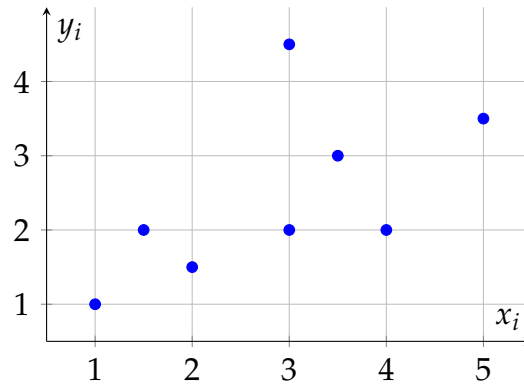
$$X^t X = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}, \quad X^t Y = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}.$$

- (b) Show that $X^t X B = X^t Y$ is equivalent to Equation (2.1):

$$\begin{aligned} nb_0 + b_1 \sum x_i &= \sum y_i, \\ b_0 \sum x_i + b_1 \sum x_i^2 &= \sum x_i y_i. \end{aligned}$$

2. Find a least squares fitting line to the following data and draw in the line:

i	x_i	y_i
1	1.0	1.0
2	2.0	1.5
3	3.0	2.0
4	1.5	2.0
5	3.5	3.0
6	3.0	4.5
7	4.0	2.0
8	5.0	3.5



(Round b_0 and b_1 to the nearest half integer.)

Week 1

2.4 Plane of best fit

We consider two independent variables and one dependent variable now. Consider the following data points as given in Figure 2.3.

i	x_{i1}	x_{i2}	y_i
0	278	36	287
1	252	31	256
2	344	35	300
3	134	33	182
4	215	35	248
5	261	40	271
6	131	39	149
7	463	43	411
8	167	46	214
9	298	42	291
10	230	60	314
11	293	67	352
12	290	37	298
13	271	31	252
14	385	63	439
15	354	36	328

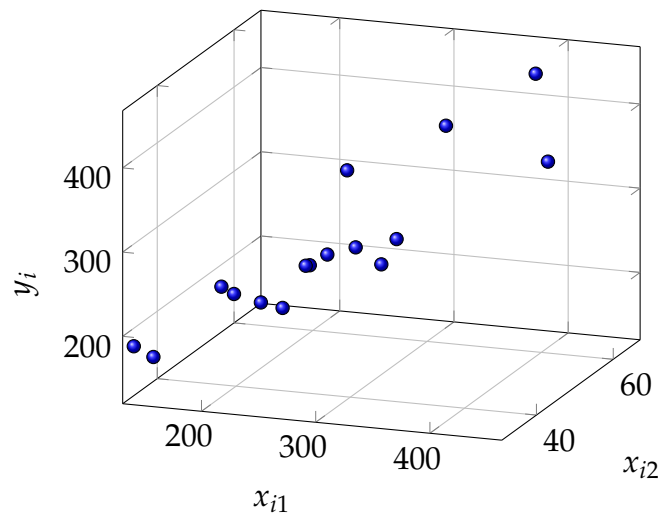


Figure 2.3: Data points in \mathbb{R}^3 .

We can put some meaning to these data. For example, suppose a company is selling a product, and we have 16 populations of people labeled 0 through 15. The values x_{i1} are the population sizes in 100s of people; the values x_{i2} are the average yearly income in €1000 per capita; and the values y_i are the number of

sales of the product. (There might be dependencies between population size and average income, but our model treats them as independent.)

It looks like though there is a plane of best fit for the data—thanks to the suggestive viewing angle. Our goal is to find a plane, given by

$$y = b_0 + b_1x_1 + b_2x_2.$$

We can just do what we did last time. That is, for

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} \end{pmatrix}, \quad Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad B = \begin{pmatrix} b_0 \\ b_1 \\ b_2 \end{pmatrix},$$

we need to solve for B in the equation

$$X^tXB = X^tY.$$

Thus, if X^tX is invertible, there is a unique B , which is equal to $(X^tX)^{-1}X^tY$. For our example, we have

$$X^tX = \begin{pmatrix} 16 & 4366 & 674 \\ 4366 & 1309480 & 187024 \\ 674 & 187024 & 30330 \end{pmatrix}, \quad X^tY = \begin{pmatrix} 4592 \\ 1343400 \\ 200571 \end{pmatrix}.$$

Therefore, the plane of best fit is approximately

$$y = -11.3 + 0.7x_1 + 2.6x_2.$$

Putting all the data together we have a plane of best fit as seen in Figure 2.4.

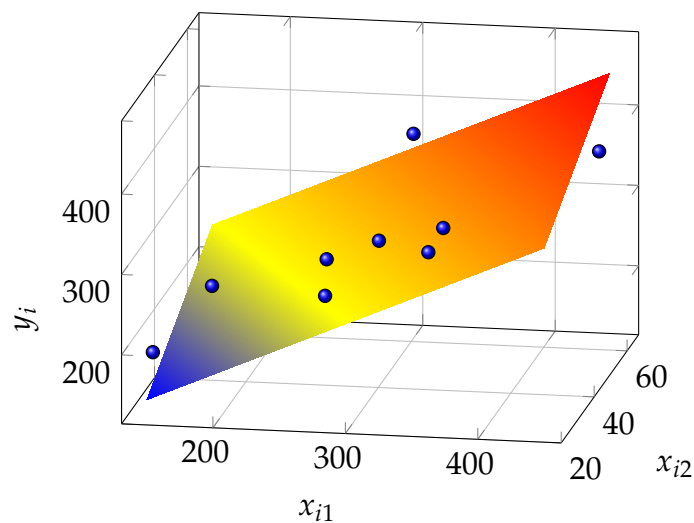


Figure 2.4: Data points together with the plane of best fit.

2.5 Hyperplane of best fit

Now we go to the general case. Suppose we have $p - 1$ independent variables and 1 dependent variable, where $p \geq 2$. We assume we have n data points of the form

$$(x_{i1}, x_{i2}, \dots, x_{i,p-1}, y_i) \in \mathbb{R}^p.$$

The least squares fitting for these data is a hyperplane of the form

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_{p-1}x_{p-1}.$$

To solve for the values b_i , we do as we did before. We define matrices

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1,p-1} \\ 1 & x_{21} & x_{22} & \cdots & x_{2,p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{n,p-1} \end{pmatrix}, \quad Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad B = \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_{p-1} \end{pmatrix}.$$

As before, the values we want are given by the equation

$$X^t X B = X^t Y. \quad (2.4)$$

2.6 Why Equation (2.4) works

The heart of least squares is (Euclidean) distance. The distance between two points $x = (x_1, \dots, x_p)$ and $y = (y_1, \dots, y_p)$ in \mathbb{R}^p is

$$d(x, y) = \|x - y\| = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_p - y_p)^2}.$$

For a vector $v = (v_1, \dots, v_p) \in \mathbb{R}^p$, the **length** of v is

$$\|v\| = d(0, v) = \sqrt{v_1^2 + v_2^2 + \dots + v_p^2}.$$

Recall that the dot product of two (column) vectors u and v is

$$u \cdot v = u^t v = u_1 v_1 + u_2 v_2 + \dots + u_p v_p.$$

Thus, the length of v is $\|v\| = \sqrt{v \cdot v}$; in other words $\|v\|^2 = v \cdot v$. In addition, if $u \cdot v = 0$, we say that u and v are **orthogonal** (or perpendicular).

The goal of least squares is to *minimize distance*; more specifically to minimize $\|Y - XB\|$. Note that the column vector Y has entries that are the *actual* y_i values, and the column vector

$$XB = \begin{pmatrix} B \cdot (1, x_{11}, x_{12}, \dots, x_{1,p-1}) \\ B \cdot (1, x_{21}, x_{22}, \dots, x_{2,p-1}) \\ \vdots \\ B \cdot (1, x_{n1}, x_{n2}, \dots, x_{n,p-1}) \end{pmatrix} = \begin{pmatrix} b_0 + b_1 x_{11} + b_2 x_{12} + \dots + b_{p-1} x_{1,p-1} \\ b_0 + b_1 x_{21} + b_2 x_{22} + \dots + b_{p-1} x_{2,p-1} \\ \vdots \\ b_0 + b_1 x_{n1} + b_2 x_{n2} + \dots + b_{p-1} x_{n,p-1} \end{pmatrix}$$

Therefore, $\|Y - XB\|$ is the square root of a sum of squares of the form

$$y_i - b_0 + b_1 x_{i1} + b_2 x_{i2} + \dots + b_{p-1} x_{i,p-1}.$$

Hence minimizing $\|Y - XB\|$ is the same as minimizing $\|Y - XB\|^2$, which is a sum of *squares*.

Proposition 2.1. *The minimal distance $\|Y - XB\|$ is achieved by solving for B in*

$$X^t XB = X^t Y.$$

Week 2

Proof. Consider the subspace $U = \{Xu \mid u \in \mathbb{R}^p\}$ of \mathbb{R}^n , and observe that our desired solution XB is contained in U . Since $\|Y - XB\|$ is minimal, we must have that the vector $Y - XB$ is orthogonal to all vectors contained in U .¹ That is, $(Xu) \cdot (Y - XB) = 0$ for all $u \in \mathbb{R}^p$. In other words, we have for all $u \in \mathbb{R}^p$,

$$\begin{aligned} 0 &= (Xu)^t(Y - XB) = u^t X^t(Y - XB) \\ &= u^t (X^t Y - X^t XB). \end{aligned}$$

Because $u^t (X^t Y - X^t XB) = 0$ for all $u \in \mathbb{R}^p$, it follows that $X^t Y - X^t XB = 0$. \square

2.7 In class exercises pt. II

1. Determine $X^t X$ and $X^t Y$ with

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1,p-1} \\ 1 & x_{21} & x_{22} & \cdots & x_{2,p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{n,p-1} \end{pmatrix}, \quad Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}.$$

2. Using (1) and by taking partial derivatives of

$$S(b_0, \dots, b_{p-1}) = \sum_{i=1}^n (y_i - (b_0 + b_1 x_{i1} + b_2 x_{i2} + \cdots + b_{p-1} x_{i,p-1}))^2, \quad (2.5)$$

show that the hyperplane of best fit is obtained by solving $X^t XB = X^t Y$. (You could try this for $p = 3$ first.)

2.8 Nonlinear fittings

Although all of our examples so far have been linear fittings, we will demonstrate that least squares fittings works in the nonlinear case. What is important is that we have a candidate equation to fit. In the linear cases, we tried to fit

$$y = b_0 + b_1 x_1 + b_2 x_2 + \cdots + b_{p-1} x_{p-1}.$$

Suppose we have the following data as given in Figure 2.5. Instead of trying to fit the line $y = b_0 + b_1 x$, we could try to fit the parabola:

$$y = b_0 + b_1 x + b_2 x^2.$$

¹To see why this is true, see Section 6.3.1 of [3], which is all about orthogonal decompositions.

We can treat this the same way as before. Of course the quantities x and x^2 are *not* independent, but we can ignore this. Set

$$x_{i1} = x_i, \quad x_{i2} = x_i^2.$$

Therefore, the hyperplane of best fit for the data (x_{i1}, x_{i2}, y_i) will give us the parabola of best fit. *Try this on your own!*

x_i	y_i
2.27	2.50
5.06	-16.13
1.45	4.23
5.89	-22.46
0.48	1.37
-0.22	0.86
1.44	11.85
-1.77	-14.71
2.45	9.42
-1.54	-14.07
7.55	-55.62
1.76	4.45
5.16	-19.56
3.26	-2.79
3.23	5.20
0.85	8.09

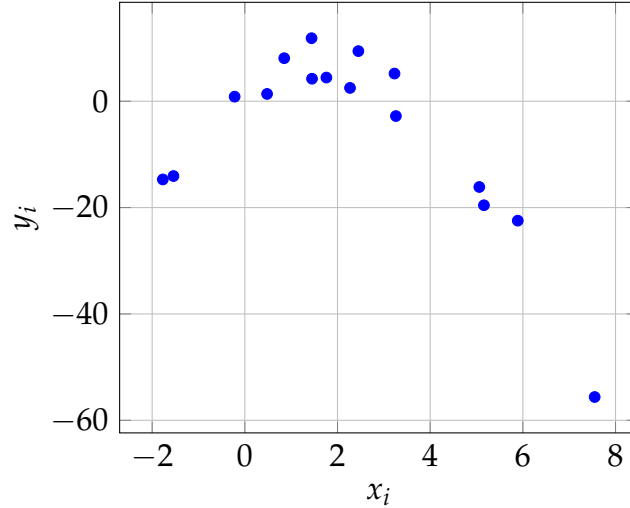


Figure 2.5: Data points demonstrating a nonlinear relationship.

So one can fit any hypersurface $y = f(x_1, \dots, x_{p-1})$ to the given data. The function f in this case is called the **regression function**. This general method of analysis is known as **regression analysis**. A few questions arise:

- Which surface is “best”?
- How can we quantify “best”?
- Even in the line case ($p = 2$), how can we quantify how well data fits our line?

2.9 Coefficient of determination (R^2 values)

We are going to make precise how well our hyperplane fits our data. Recall that hyperplanes can be replaced by hypersurfaces; see Section 2.8. First we establish some notation. Suppose we have n data points $(x_{i1}, x_{i2}, \dots, x_{i,p-1}, y_i) \in \mathbb{R}^p$. Then we define

$$\begin{aligned} \text{(Fitted value)} \quad \hat{y}_i &= b_0 + b_1 x_{i1} + b_2 x_{i2} + \dots + b_{p-1} x_{i,p-1}, \\ \text{(Residual)} \quad e_i &= y_i - \hat{y}_i, \end{aligned}$$

$$\text{(Sample mean)} \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

These yield vectors in \mathbb{R}^n as follows

$$\hat{Y} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{pmatrix} = XB, \quad E = \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix} = Y - \hat{Y}, \quad \bar{Y} = \begin{pmatrix} \bar{y} \\ \vdots \\ \bar{y} \end{pmatrix} = \bar{y} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}.$$

From our n data points, we have three points in \mathbb{R}^n given by Y , \hat{Y} , and \bar{Y} . Three points always lie on a plane, so the three points determine a triangle on such a plane. What does this triangle look like? If it is a triangle (and not a line or a single point), then the next lemma proves it must be a right triangle.

Lemma 2.2. *The vectors $E = Y - \hat{Y}$ and $\hat{Y} - \bar{Y}$ are orthogonal.*

Proof. Suppose $X^t X B = X^t Y$. We need to prove two equations. For the first,

$$0 = X^t(Y - XB) = X^t(Y - \hat{Y}) = X^t E.$$

Hence, $X^t E = 0$. For the second,

$$\begin{aligned} \bar{Y}^t E &= \bar{y} \sum_{i=1}^n (y_i - \hat{y}_i) \\ &= \bar{y} \sum_{i=1}^n (y_i - (b_0 + b_1 x_{i1} + b_2 x_{i2} + \cdots + b_{p-1} x_{i,p-1})) \\ &= -\frac{\bar{y}}{2} \cdot \frac{\partial S}{\partial b_0} = 0, \end{aligned}$$

where S is defined in Equation (2.5), so $\bar{Y}^t E = 0$. Thus, we have

$$(\hat{Y} - \bar{Y}) \cdot E = (XB)^t E - \bar{Y}^t E = 0. \quad \square$$

Remark 2.3. One can simplify the proof for Lemma 2.2 by applying an isometry to the data, so that $\bar{y} = 0$. That is, one only needs to prove that E and \hat{Y} are orthogonal.

The lengths of the differences of the vectors are important and have names:

$$\begin{aligned} \text{(Sums of Squares Total – SST)} &: \|Y - \bar{Y}\|^2, \\ \text{(Sums of Squares Error – SSE)} &: \|Y - \hat{Y}\|^2, \\ \text{(Sums of Squares Regression – SSR)} &: \|\hat{Y} - \bar{Y}\|^2. \end{aligned}$$

The value SST measures the *total variability* of the data set. For example, $\sqrt{SST} = \|Y - \bar{Y}\|$ is the distance from the actual data Y to the sample mean \bar{Y} . Using the same ideas, we can see that SSE measures the error of our regression and that SSR measures the distance from our regression to the sample mean.

Proposition 2.4.

$$SST = SSE + SSR.$$

Proof. Apply Lemma 2.2 and the Pythagorean Theorem:

$$\|Y - \bar{Y}\|^2 = \|Y - \hat{Y}\|^2 + \|\hat{Y} - \bar{Y}\|^2. \quad \square$$

Now we can describe a quantity that measures how good our regression fits the given data.

Definition 2.5. The **coefficient of determination** (also known as the R^2 -value) is

$$R^2 = \frac{SSR}{SST} = \frac{\|\hat{Y} - \bar{Y}\|^2}{\|Y - \bar{Y}\|^2}.$$

Proposition 2.6. $0 \leq R^2 \leq 1$.

Proof. Since each SST and SSR are squares, they are nonnegative. By Proposition 2.4, we have $0 \leq SSR \leq SST$. \square

2.9.1 What do the extremes means?

The one case where R^2 is meaningless is when $SST = 0$. This implies both $SSR = SSE = 0$. Moreover, $Y = \bar{Y} = \bar{y}\mathbb{1}$, where $\mathbb{1}$ is the all ones column vector. Hence, every data point y_i is the same and, therefore, equal to the mean. Let's never return to this case.

We can have $SSR = 0$, which is equivalent to $R^2 = 0$. This implies that $\|\hat{Y} - \bar{Y}\|^2 = 0$, so that $\hat{Y} = \bar{Y}$. In other words, our prediction \hat{y}_i is just simply the mean. This means we have not found any relationship between the independent variables and the dependent variables.

In the other extreme we have $SSR = SST$, which is equivalent to $R^2 = 1$. This implies that $Y = \hat{Y}$, so the given data lies (exactly) on the surface given by $y = f(x_1, \dots, x_{p-1})$. That is, the regression function exactly predicts the data.

To summarize, when $R^2 = 0$, we cannot deduce any relationship between the independent and dependent variables, and when $R^2 = 1$, we understand completely the relationship between the independent and dependent variables. Very roughly speaking, the R^2 can be thought of as the ratio of how well the regression fits the data.

2.10 In class exercises pt. III

1. Prove the following.

- (a) $\|\bar{Y}\|^2 = n\bar{y}$.
- (b) $Y \cdot \bar{Y} = \hat{Y} \cdot \bar{Y} = \|\bar{Y}\|^2$.
- (c) $Y \cdot \hat{Y} = \|\hat{Y}\|^2$.

2. Use (1) to show that

(a) $SST = \|Y\|^2 - \|\bar{Y}\|^2,$

(b) $SSE = \|Y\|^2 - \|\hat{Y}\|^2,$

(c) $SSR = \|\hat{Y}\|^2 - \|\bar{Y}\|^2.$

3. What are the R^2 values for the examples above?

References

- [1] Christos Sotiriou, Soek-Ying Neo, Lisa M. McShane, Edward L. Korn, Philip M. Long, Amir Jazaeri, Philippe Martiat, Steve B. Fox, Adrian L. Harris, and Edison T. Liu, *Breast cancer classification and prognosis based on gene expression profiles from a population-based study*, Proceedings of the National Academy of Sciences **100** (2003), no. 18, 10393–10398.
- [2] Google, *What is Clustering?* (2023), date accessed: 16 Oct. 2023. <https://developers.google.com/machine-learning/clustering/overview>.
- [3] Dan Margalit and Joseph Rabinoff, *Interactive Linear Algebra*, 2019, <https://textbooks.math.gatech.edu/ila/>.
- [4] Jonathon Shlens, *A tutorial on principal component analysis*, 2014, [arXiv:1404.1100](https://arxiv.org/abs/1404.1100).