

1 Linear Regression

1.1 Introduction

Suppose that we are trying to study two quantities x and y that we suspect are related – at least approximately – by a linear equation $y = ax + b$. Sometimes this linear relationship is predicted by theoretical considerations, and sometimes it is just an empirical hypothesis.

For example, if we are trying to determine the velocity of an object travelling towards us at constant speed, and we measure the distances d_1, d_2, \dots, d_n between us and the object at a series of times t_1, t_2, \dots, t_n , then since “distance equals rate times time” we have a theoretical foundation for the assumption that $d = rt + b$ for some constants r and b . On the other hand, because of unavoidable experimental errors, we can't expect that this relationship will hold exactly for the observed data; instead, we likely get a graph like that shown in fig. 1. We've drawn a line on the plot that seems to capture the true slope (and hence velocity) of the object.

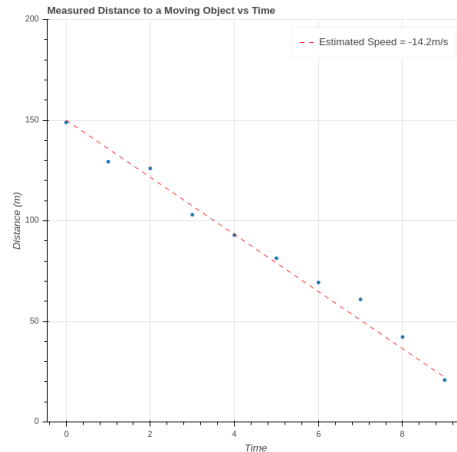


Figure 1: Physics Experiment

On the other hand, we might look at a graph such as fig. 2, which plots the gas mileage of various car models against their engine size (displacement), and observe a general trend in which bigger engines get lower mileage. In this situation we could ask for the best line of the form $y = mx + b$ that captures this relationship and use that to make general conclusions without necessarily having an underlying theory.

1.2 Least Squares (via Calculus)

In either of the two cases above, the question we face is to determine the line $y = mx + b$ that “best fits” the data $\{(x_i, y_i)_{i=1}^N\}$. The classic approach is to

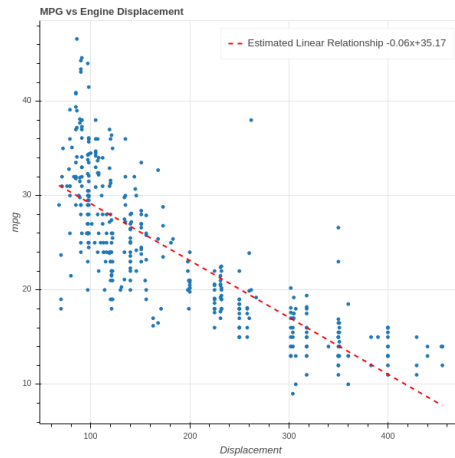


Figure 2: MPG vs Displacement ([1])

determine the equation of a line $y = mx + b$ that minimizes the “mean squared error”:

$$MSE(m, b) = \frac{1}{N} \sum_{i=1}^n (y_i - mx_i - b)^2$$

It’s worth emphasizing that the MSE is a function of two variables – the slope m and the intercept b – and that the data points $\{(x_i, y_i)\}$ are constants for these purposes. Furthermore, it’s a quadratic function in those two variables. Since our goal is to find m and b that minimize the MSE , we have a Calculus problem that we can solve by taking partial derivatives and setting them to zero.

To simplify the notation, let’s abbreviate MSE by E .

$$\frac{\partial E}{\partial m} = \frac{1}{N} \sum_1^N -2x_i(y_i - mx_i - b)$$

$$\frac{\partial E}{\partial b} = \frac{1}{N} \sum_1^N -2(y_i - mx_i - b)$$

We set these two partial derivatives to zero, so we can drop the -2 and regroup the sums to obtain two equations in two unknowns (we keep the $\frac{1}{N}$ because it is illuminating in the final result):

$$\begin{aligned} \frac{1}{N} \left(\sum_{i=1}^N x_i^2 \right) m + \frac{1}{N} \left(\sum_{i=1}^N x_i \right) b &= \frac{1}{N} \sum_{i=1}^N x_i y_i \\ \frac{1}{N} \left(\sum_{i=1}^N x_i \right) m + b &= \frac{1}{N} \sum_{i=1}^N y_i \end{aligned} \tag{1}$$

In these equations, notice that $\frac{1}{N} \sum_{i=1}^N x_i$ is the average (or mean) value of the x_i . Let's call this \bar{x} . Similarly, $\frac{1}{N} \sum_{i=1}^N y_i$ is the mean of the y_i , and we'll call it \bar{y} . If we further simplify the notation and write S_{xx} for $\frac{1}{N} \sum_{i=1}^N x_i^2$ and S_{xy} for $\frac{1}{N} \sum_{i=1}^N x_i y_i$ then we can write down a solution to this system using Cramer's rule:

$$\begin{aligned} m &= \frac{S_{xy} - \bar{x}\bar{y}}{S_{xx} - \bar{x}^2} \\ b &= \frac{S_{xx}\bar{y} - S_{xy}\bar{x}}{S_{xx} - \bar{x}^2} \end{aligned} \tag{2}$$

where we must have $S_{xx} - \bar{x}^2 \neq 0$.

1.2.1 Exercises

1. Verify that eq. 2 is in fact the solution to the system in eq. 1.
2. Suppose that $S_{xx} - \bar{x}^2 = 0$. What does that mean about the x_i ? Does it make sense that the problem of finding the "line of best fit" fails in this case?

1.3 Least Squares (via Geometry)

In our discussion above, we thought about our data as consisting of N pairs (x_i, y_i) corresponding to n points in the xy -plane \mathbf{R}^2 . Now let's turn that picture "on its side," and instead think of our data as consisting of *two* points in \mathbf{R}^n :

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \text{and} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Let's also introduce one other vector

$$E = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}.$$

First, let's assume that E and X are linearly independent. If not, then X is a constant vector (why?) which we already know is a problem from section 1.2, Exercise 2. Therefore E and X span a plane in \mathbf{R}^n .

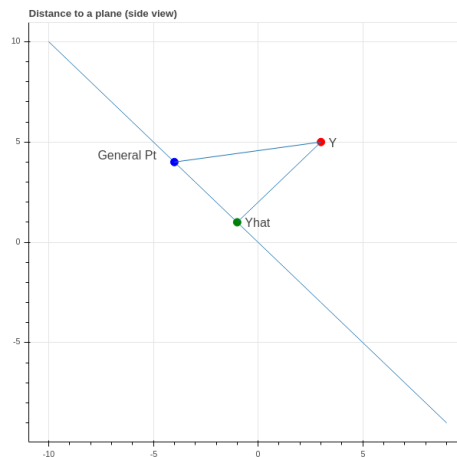


Figure 3: Distance to A Plane

Now if our data points (x_i, y_i) all *did* lie on a line $y = mx + b$, then the three vectors X , Y , and E would be linearly dependent:

$$Y = mX + bE.$$

Since our data is only approximately linear, that's not the case. So instead we look for an approximate solution. One way to phrase that is to ask:

What is the point \hat{Y} in the plane H spanned by X and E in \mathbf{R}^n which is closest to Y ?

If we knew this point \hat{Y} , then since it lies in H we would have $\hat{Y} = mX + bE$ and the coefficients m and b would be a candidate for defining a line of best fit $y = mx + b$. Finding the point in a plane closest to another point in \mathbf{R}^n is a geometry problem that we can solve.

Proposition: The point \hat{Y} in the plane spanned by X and E is the point such that the vector $Y - \hat{Y}$ is perpendicular to H .

Proof: See fig. 3 for an illustration – perhaps you are already convinced by this, but let’s be careful. $\hat{Y} = mX + bE$ such that

$$D = \|Y - \hat{Y}\|^2 = \|Y - mX - bE\|^2$$

is minimal. Using some vector calculus, we have

$$\frac{\partial D}{\partial m} = \frac{\partial}{\partial m}(Y - mX - bE) \cdot (Y - mX - bE) = -2(Y - mX - bE) \cdot X$$

and

$$\frac{\partial D}{\partial b} = \frac{\partial}{\partial b}(Y - mX - bE) \cdot (Y - mX - bE) = -2(Y - mX - bE) \cdot E.$$

So both derivatives are zero exactly when $\hat{Y} = (Y - mX - bE)$ is orthogonal to both X and E , and therefore every vector in H .

We also obtain equations for m and b just as in our first look at this problem.

$$\begin{aligned} m(X \cdot E) + b(E \cdot E) &= (Y \cdot E) \\ m(X \cdot X) + b(E \cdot X) &= (Y \cdot X) \end{aligned} \tag{3}$$

We leave it is an exercise below to check that these are the same equations that we obtained in eq. 2.

1.3.1 Exercises

1. Verify that eq. 2 and eq. 3 are equivalent.

1.4 The Multivariate Case (Calculus)

Having worked through the problem of finding a “line of best fit” from two points of view, let’s look at a more general problem. We looked above at a scatterplot showing the relationship between gas mileage and engine size (displacement). There are other factors that might contribute to gas mileage that we want to consider as well – for example: - a car that is heavy compared to its engine size may get worse mileage - a sports car with a drive train that gives fast acceleration as compared to a car with a transmission designed for long trips may have different mileage for the same engine size.

Suppose we wish to use engine displacement, vehicle weight, and acceleration all together to predict mileage. Instead of looking points (x_i, y_i) where x_i is the displacement of the i^{th} car model and we try to predict a value y from a corresponding x as $y = mx + b$ – let’s look at a situation in which our measured value y depends on multiple variables – say displacement d , weight w , and acceleration a with $k = 3$ – and we are trying to find the best linear equation

$$y = m_1d + m_2w + m_3a + b \tag{4}$$

But to handle this situation more generally we need to adopt a convention that will allow us to use indexed variables instead of d , w , and a . We will use the *tidy* data convention.

Tidy Data: A dataset is tidy if it consists of values x_{ij} for $i = 1, \dots, N$ and $j = 1, \dots, k$ so that:

- the row index corresponds to a *sample* – a set of measurements from a single event or item;
- the column index corresponds to a *feature* – a particular property measured for all of the events or items.

In our case,

- the *samples* are the different types of car models,
- the *features* are the properties of those car models.

For us, N is the number of different types of cars, and k is the number of properties we are considering. Since we are looking at displacement, weight, and acceleration, we have $k = 3$.

So the “independent variables” for a set of data that consists of N samples, and k measurements for each sample, can be represented by a $N \times k$ matrix

$$X = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Nk} \end{pmatrix}$$

and the measured dependent variables Y are a column vector

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}.$$

If m_1, \dots, m_k are “slopes” associated with these properties in eq. 4, and b is the “intercept,” then the predicted value \hat{Y} is given by a matrix equation

$$\hat{Y} = X \begin{bmatrix} m_1 \\ m_2 \\ \dots \\ m_k \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} b$$

and our goal is to choose these parameters m_i and b to make the mean squared error:

$$MSE(m_1, \dots, m_k, b) = \|Y - \hat{Y}\|^2 = \sum_{i=1}^N (y_i - \sum_{j=1}^k x_{ij}m_j - b)^2.$$

Here we are summing over the N different car models, and for each model taking the squared difference between the true mileage y_i and the “predicted” mileage $\sum_{j=1}^k x_{ij}m_j + b$. We wish to minimize this MSE.

Let’s make one more simplification. The intercept variable b is annoying because it requires separate treatment from the m_i . But we can use a trick to eliminate the need for special treatment. Let’s add a new feature to our data matrix (a new column) that has the constant value 1.

$$X = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} & 1 \\ x_{21} & x_{22} & \cdots & x_{2k} & 1 \\ \vdots & \vdots & \ddots & \vdots & 1 \\ x_{N1} & x_{N2} & \cdots & x_{Nk} & 1 \end{pmatrix}$$

Now our data matrix X is $N \times (k+1)$ and we can put our “intercept” $b = m_{k+1}$ into our vector of “slopes” m_1, \dots, m_k, m_{k+1} :

$$\hat{Y} = X \begin{bmatrix} m_1 \\ m_2 \\ \cdots \\ m_k \\ m_{k+1} \end{bmatrix}.$$

and our MSE becomes

$$MSE(M) = \|Y - XM\|^2$$

where

$$M = \begin{bmatrix} m_1 \\ m_2 \\ \cdots \\ m_k \\ m_{k+1} \end{bmatrix}.$$

Remark: Later on (see section 1.6) we will see that if we “center” our features about their mean, by subtracting the average value of each column of X from

that column; and we also subtract the average value of Y from the entries of Y , then the b that emerges from the least squares fit is zero. As a result, instead of adding a column of 1's, you can change coordinates to center each feature about its mean, and keep your X matrix $N \times k$.

The Calculus approach to minimizing the MSE is to take its partial derivatives with respect to the m_i and set them to zero. Let's first work out the derivatives in a nice form for later.

Proposition: The gradient of $MSE(M) = E$ is given by

$$\nabla E = \begin{bmatrix} \frac{\partial}{\partial M_1} E \\ \frac{\partial}{\partial M_2} E \\ \vdots \\ \frac{\partial}{\partial m_{M+1}} E \end{bmatrix} = -2X^\top Y + 2X^\top X M \quad (5)$$

where X^\top is the transpose of X .

Proof: First, remember that the ij entry of X^\top is the ji entry of X . Also, we will use the notation $X[j, :]$ to mean the j^{th} row of X and $X[:, i]$ to mean the i^{th} column of X . (This is copied from the Python programming language; the ':' means that index runs over all possibilities).

Since

$$E = \sum_{j=1}^N (Y_j - \sum_{s=1}^{k+1} X_{js} M_s)^2$$

we compute:

$$\begin{aligned} \frac{\partial}{\partial M_t} E &= -2 \sum_{j=1}^N X_{jt} (Y_j - \sum_{s=1}^{k+1} X_{js} M_s) \\ &= -2 \left(\sum_{j=1}^N Y_j X_{jt} - \sum_{j=1}^N \sum_{s=1}^{k+1} X_{jt} X_{js} M_s \right) \\ &= -2 \left(\sum_{j=1}^N X_{tj}^\top Y_j - \sum_{j=1}^N \sum_{s=1}^{k+1} X_{tj}^\top X_{js} M_s \right) \\ &= -2 \left(X^\top [t, :] Y - \sum_{s=1}^{k+1} \sum_{j=1}^N X_{tj}^\top X_{js} M_s \right) \\ &= -2 \left(X^\top [t, :] Y - \sum_{s=1}^{k+1} (X^\top X)_{ts} M_s \right) \\ &= -2 \left(X^\top [t, :] Y - (X^\top X) [t, :] M \right) \end{aligned} \quad (6)$$

Stacking up the different rows to make E yields the desired formula.

Proposition: Assume that $D = X^\top X$ is invertible (notice that it is a $(k+1) \times (k+1)$ square matrix so this makes sense). The solution M to the multivariate least squares problem is

$$M = D^{-1} X^\top Y \quad (7)$$

and the “predicted value” \hat{Y} for Y is

$$\hat{Y} = X D^{-1} X^\top Y. \quad (8)$$

1.5 The Multivariate Case (Geometry)

Let’s look more closely at the equation obtained by setting the gradient of the error, eq. 5, to zero. Remember that M is the unknown vector in this equation, everything else is known:

$$X^\top Y = X^\top X M$$

Here is how to think about this:

1. As M varies, the $N \times 1$ matrix $X M$ varies over the space spanned by the columns of the matrix X . So as M varies $X M$ is a general element of the subspace H of R^N spanned by the $k+1$ columns of X .
2. The product $X^\top X M$ is a $(k+1) \times 1$ matrix. Each entry is the dot product of the general element of H with one of the $k+1$ basis vectors of H .
3. The product $X^\top Y$ is a $(k+1) \times 1$ matrix whose entries are the dot product of the basis vectors of H with Y .

Therefore, this equation asks for us to find M so that the vector $X M$ in H has the same dot products with the basis vectors of H as Y does. The condition

$$X^\top \cdot (Y - X M) = 0$$

says that $Y - X M$ is orthogonal to H . This argument establishes the following proposition.

Proposition: Just as in the simple one-dimensional case, the predicted value \hat{Y} of the least squares problem is the point in H closest to Y – or in other words the point \hat{Y} in H such that $Y - \hat{Y}$ is perpendicular to H .

1.5.1 Orthogonal Projection

Recall that we introduced the notation $D = X^\top X$, and let’s assume, for now, that D is an invertible matrix. We have the formula (see eq. 8):

$$\hat{Y} = X D^{-1} X^\top Y.$$

Proposition: The matrix $P = XD^{-1}X^\top$ is an $N \times N$ matrix called the orthogonal projection operator onto the subspace H spanned by the columns of X . It has the following properties:

- PY belongs to the subspace H for any $Y \in \mathbf{R}^N$.
- $(Y - PY)$ is orthogonal to H .
- $P * P = P$.

Proof: First of all, $PY = XD^{-1}X^\top Y$ so PY is a linear combination of the columns of X and is therefore an element of H . Next, we can compute the dot product of PY against a basis of H by computing

$$X^\top PY = X^\top XD^{-1}X^\top Y = X^\top Y$$

since $X^\top X = D$. This equation means that $X^\top(Y - PY) = 0$ which tells us that $Y - PY$ has dot product zero with a basis for H . Finally,

$$PP = XD^{-1}X^\top XD^{-1}X^\top = XD^{-1}X^\top = P.$$

It should be clear from the above discussion that the matrix $D = X^\top X$ plays an important role in the study of this problem. In particular it must be invertible or our analysis above breaks down. In the next section we will look more closely at this matrix and what information it encodes about our data.

1.6 Centered coordinates

Recall from last section that the matrix $D = X^\top X$ is of central importance to the study of the multivariate least squares problem. Let's look at it more closely.

Lemma: The i, j entry of D is the dot product

$$D_{ij} = X[:, i] \cdot X[:, j]$$

of the i^{th} and j^{th} columns of X .

Proof: In the matrix multiplication $X^\top X$, the i^{th} row of X^\top gets "dotted" with the j^{th} column of X to product the i, j entry. But the i^{th} row of X^\top is the i^{th} column of X , as asserted in the statement of the lemma.

A crucial point in our construction above relied on the matrix D being invertible. The following Lemma shows that D fails to be invertible only when the different features (the columns of X) are linearly dependent.

Lemma: D is not invertible if and only if the columns of X are linearly dependent.

Proof: If the columns of X are linearly dependent, then there is a nonzero vector m so that $Xm = 0$. In that case clearly $Dm = X^\top Xm = 0$ so D is not invertible. Suppose D is not invertible. Then there is a nonzero vector m with

$Dm = X^\top X m = 0$. This means that the vector Xm is orthogonal to all of the columns of X . Since Xm belongs to the span H of the columns of X , if it is orthogonal to H it must be zero.

In fact, the matrix D captures some important statistical measures of our data, but to see this clearly we need to make a slight change of basis. First recall that $X[:, k+1]$ is our column of all 1, added to handle the intercept. As a result, the dot product $X[:, i] \cdot X[:, k+1]$ is the sum of the entries in the i^{th} column, and so if we let μ_i denote the average value of the entries in column i , we have

$$\mu_i = \frac{1}{N}(X[:, i] \cdot X[:, k+1])$$

Now change the matrix X by elementary column operations to obtain a new data matrix X_0 by setting

$$X_0[:, i] = X[:, i] - \frac{1}{N}(X[:, i] \cdot X[:, k+1])X[:, k+1] = X[:, i] - \mu_i X[:, k+1]$$

for $i = 1, \dots, k$.

In terms of the original data, we are changing the measurement scale of the data so that each feature has average value zero, and the subspace H spanned by the columns of X_0 is the same as that spanned by the columns of X . Using X_0 instead of X for our least squares problem, we get

$$\hat{Y} = X_0 D_0^{-1} X_0^\top Y$$

and

$$M_0 = D_0^{-1} X_0^\top Y$$

where $D_0 = X_0^\top X_0$.

Proposition: The matrix D_0 has a block form. Its upper left block is a $k \times k$ symmetric block with entries

$$(D_0)_{ij} = (X[:, i] - \mu_i X[:, k+1]) \cdot (X[:, j] - \mu_j X[:, k+1])$$

Its $(k+1)^{\text{st}}$ row and column are all zero, except for the $(k+1), (k+1)$ entry, which is N .

Proof: This follows from the fact that the last row and column entries are (for $i \neq k+1$):

$$(X[:, i] - \mu_i X[:, k+1]) \cdot X[:, k+1] = (X[:, i] \cdot X[:, k+1]) - N\mu_i = 0$$

and for $i = k+1$ we have $X[:, k+1] \cdot X[:, k+1] = N$ since that column is just N 1's.

Proposition: If the x coordinates (the features) are centered so that they have mean zero, then the intercept b is

$$\bar{Y} = \frac{1}{N} \sum y_i.$$

Proof: By centering the coordinates, we replace the matrix X by X_0 and D by D_0 . and we are trying to minimize $\|Y - X_0 M_0\|^2$. Use the formula from eq. 7 to see that

$$M_0 = D_0^{-1} X_0^T Y.$$

The b value we are interested in is the last entry m_{k+1} in M_0 . From the block form of D_0 , we know that D_0^{-1} has bottom row and last column zero except for $1/N$ in position $(k+1) \times (k+1)$. Also X_0^T has last row consisting entirely of 1. So the bottom entry of $X_0^T Y$ is $\sum_{i=1}^N y_i$, and the bottom entry b of $D_0^{-1} X_0^T Y$ is

$$\mu_Y = \frac{1}{N} \sum_{i=1}^N y_i.$$

as claimed.

Corollary: If we make a further change of coordinates to define

$$Y_0 = Y - \mu_Y \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

then the associated b is zero. As a result we can forget about the extra column of 1's that we added to X to account for it and reduce the dimension of our entire problem by 1.

Just to recap, if we center our data so that $\mu_Y = 0$ and $\mu_i = 0$ for $i = 1, \dots, k$, then the least squares problem reduces to minimizing

$$E(M) = \|Y - XM\|^2$$

where X is the $N \times k$ matrix with j^{th} row $(x_{j1}, x_{j2}, \dots, x_{jk})$ for $j = 1, \dots, N$ and the solutions are as given in eq. 7 and eq. 8.

1.7 Caveats about Linear Regression

1.7.1 Basic considerations

Reflecting on our long discussion up to this point, we should take note of some of the potential pitfalls that lurk in the use of linear regression.

1. When we apply linear regression, we are explicitly assuming that the variable Y is associated to X via linear equations. This is a big assumption!

2. When we use multilinear regression, we are assuming that changes in the different features have independent effects on the target variable y . In other words, suppose that $y = ax_1 + bx_2$. Then an increase of x_1 by 1 increases y by a , and an increase of x_2 by 1 increases y by b . These effects are independent of one another and combine to yield an increase of $a + b$.
3. We showed in our discussion above that linear regression problem has a solution when the matrix $D = X^\top X$ is invertible, and this happens when the columns of D are linearly independent. When working with real data, which is messy, we could have a situation in which the features we are studying are, in fact, dependent – but because of measurement error, the samples that we collected aren't. In this case, the matrix D will be “close” to being non-invertible, although formally still invertible. In this case, computing D^{-1} leads to numerical instability and the solution we obtain is very unreliable.

1.7.2 Simpson's Effect

Simpson's effect is a famous phenomenon that illustrates that linear regression can be very misleading in some circumstances. It is often a product of “pooling” results from multiple experiments. Suppose, for example, that we are studying the relationship between a certain measure of blood chemistry and an individual's weight gain or less on a particular diet. We do our experiments in three labs, the blue, green, and red labs. Each lab obtains similar results – higher levels of the blood marker correspond to greater weight gain, with a regression line of slope around 1. However, because of differences in the population that each lab is studying, some populations are more susceptible to weight gain and so the red lab sees a mean increase of almost 9 lbs while the blue lab sees a weight gain of only 3 lbs on average.

The three groups of scientists pool their results to get a larger sample size and do a new regression. Surprise! Now the regression line has slope -1.6 and increasing amounts of the marker seem to lead to *less* weight gain!

This is called Simpson's effect, or Simpson's paradox, and it shows that unknown factors (confounding factors) may cause linear regression to yield misleading results. This is particularly true when data from experiments conducted under different conditions is combined; in this case, the differences in experimental setting, called *batch effects*, can throw off the analysis very dramatically. See fig. 4 .

1.7.3 Exercises

1. When proving that D is invertible if and only if the columns of X are linearly independent, we argued that if $X^\top X m = 0$ for a nonzero vector m , then $X m$ is orthogonal to the span of the columns of X , and is also an element of that span, and is therefore zero. Provide the details: show that

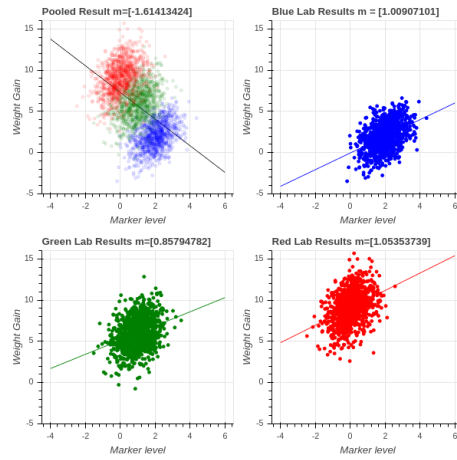


Figure 4: Simpson's Effect

if H is a subspace of \mathbf{R}^N , and x is a vector in H such that $x \cdot h = 0$ for all $h \in H$, then $x = 0$.

References

- [1] U.C. IRVINE ML REPOSITORY. Auto MPG Dataset. Available at <http://archive.ics.uci.edu/ml/datasets/Auto+MPG>.