

MATH 342W / 650.4 Spring 2024 Homework #2

Osman Khan

Saturday 24th February, 2024

Problem 1

These are questions about Silver's book, chapter 2, 3. Answer the questions using notation from class (i.e. $t, f, g, h^*, \delta, \epsilon, e, z_1, \dots, z_t, \mathbb{D}, \mathcal{H}, \mathcal{A}, \mathcal{X}, \mathcal{Y}, X, y, n, p, x_1, \dots, x_p, x_1, \dots, x_n$, etc).

- (a) [harder] If one's goal is to fit a model for a phenomenon y , what is the difference between the approaches of the hedgehog and the fox? Connecting this to the modeling framework should really make you think about what Tetlock's observation means for political and historical phenomena.

Hedgehogs have simple \mathcal{H} , which are not complex enough for the events about which forecasts are made, while foxes have more complex \mathcal{H} . In addition, foxes have models that have done a better job at learning from \mathbb{D} , and also at model validation. Hedgehogs tend to spend more time thinking from a theoretical point of view, and almost ignoring \mathbb{D} .

- (b) [easy] Why did Harry Truman like hedgehogs? Are there a lot of people that think this way?

Harry Truman liked hedgehogs because they would have given him a clear answer/line of action. I think a lot of people are not like hedgehogs (or foxes) because not everyone have strong feelings and positions about things. I think most people have a goal of keeping as many people happy as possible, and are worried having a strong position might offend some.

- (c) [difficult] Why is it that the more education one acquires, the less accurate one's predictions become?

With more education, one tends to become more of an expert in one specific field, and have an increase in hubris, which would cloud that person's judgement in regards to model making for forecasts. It tends to confirm their biases.

- (d) [easy] Why are probabilistic classifiers (i.e. algorithms that output functions that return probabilities) better than vanilla classifiers (i.e. algorithms that only return the class label)? We will move in this direction in class soon.

It is more realistic to say a range of things can happen, rather than saying a specific thing would happen. For e.g., saying there is a high chance of 4-7 mm of rain happening tomorrow v. there is a high chance of 5.5 mm of rain happening tomorrow.

- (e) [easy] What algorithm that we studied in class is PECOTA most similar to?

PECOTA is most similar to the nearest neighbor model.

- (f) [easy] Is baseball performance as a function of age a linear model? Discuss.

According to James, baseball performance as a function of age was a quadratic model with a peak at 27. I agree with that as most athletes tend to peak at around 27, and their performances do follow the quadratic trajectory of getting better until 27, and then getting worse afterwards. There surely are exceptions to this, as some outliers do exist, i.e. Justin Verlander of the Astros.

- (g) [harder] How can baseball scouts do better than a prediction system like PECOTA?

Baseball scouts are able to see the soft skills that PECOTA cannot see. This is more clearly seen in minor-leagues (when players are younger). A good pitcher, with low speeds, might be getting lots of outs in the minor-league as the batters aren't good enough, but a good scout would know that pitcher would get thrashed in the major league.

- (h) [harder] Why hasn't anyone (at the time of the writing of Silver's book) taken advantage of Pitch f/x data to predict future success?

Pitch f/x was not available at all minor-league stadiums. Incidentally, it was replaced by Trackman in 2017.

Problem 2

These are questions about the SVM.

- (a) [easy] State the hypothesis set \mathcal{H} inputted into the support vector machine algorithm. Is it different than the \mathcal{H} used for \mathcal{A} = perceptron learning algorithm?

$$\mathcal{H} = \{1_{\vec{w} \cdot \vec{x} \geq 0 = w_0 + w_1 x_1 + w_2 x_2 \geq 0} : w_0 \in \mathbb{R}, w_1 \in \mathbb{R}, w_2 \in \mathbb{R} = \vec{w} \in \mathbb{R}^3\}.$$

No, it is the same Hypothesis Set (all lines: $\mathcal{H} = \{1_{x_2 \leq a + bx} : a \in \mathbb{R}, b \in \mathbb{R}\}$), just reparametrized. We reduce the dimension of the w vector and state the intercept separately as b , to get:

$$\mathcal{H} = 1_{\vec{w} \cdot \vec{x} - b \geq 0} : \vec{w} \in \mathbb{R}^p, b \in \mathbb{R}\}$$

- (b) [E.C.] Prove the max-margin linearly separable SVM converges. State all assumptions. Write it on a separate page.

- (c) [difficult] Let $\mathcal{Y} = \{-1, 1\}$. Rederive the cost function whose minimization yields the SVM line in the linearly separable case.

I) All $y=1$'s are above or equal to L_-U :

$$\forall i' s.t. y_i = 1, \vec{w}\vec{x}_i - (b + 1) \geq 0 \Rightarrow \vec{w}\vec{x}_i - b \geq 1 \Rightarrow y_i(\vec{w}\vec{x}_i - b) \geq 1$$

II) All $y=0$'s are below or equal to L_-L :

$$\forall i' s.t. y_i = -1, \vec{w}\vec{x}_i - (b - 1) \leq 0 \Rightarrow \vec{w}\vec{x}_i - b \leq -1 \Rightarrow y_i(\vec{w}\vec{x}_i - b) \geq 1$$

- (d) [easy] Given your answer to (c) rederive the cost function using the “soft margin” i.e. the hinge loss plus the term with the hyperparameter λ . This is marked easy since there is just one change from the expression given in class.

$$THE := \sum_{i=1}^n \max\{0, 1 - y_i(\vec{w}\vec{x}_i - b)\},$$

$$H_i = \max\{0, 1 - y_i(\vec{w}\vec{x}_i - b)\},$$

Vapnik Objective Function:

$$\underset{\vec{w}, b}{\operatorname{argmin}} \left\{ \frac{1}{n} THE + \lambda \|\vec{w}\|^2 \right\} \Rightarrow g = \mathcal{A}(\mathbb{D}, \mathcal{H}, \lambda)$$

Problem 3

These are questions are about the k nearest neighbors (KNN) algorithm.

- (a) [easy] Describe how the algorithm works. Is k a “hyperparameter”?

KNN looks at K nearest neighbors of a datapoint, and returns the mode of those K nearest neighbors as the prediction for a new data point. K is indeed a hyperparameter.

- (b) [difficult] [MA] Assuming $\mathcal{A} = \text{KNN}$, describe the input \mathcal{H} as best as you can.

For KNN, we will not be predicting data that have never occurred before. The prediction will always be a subset of the historical data. This means our \mathcal{H} could be a ‘multisubset’ of \mathbb{D} .

- (c) [easy] When predicting on \mathbb{D} with $k = 1$, why should there be zero error? Is this a good estimate of future error when new data comes in? (Error in the future is called *generalization error* and we will be discussing this later in the semester).

Errors are measured by $y - \hat{y}$. But, if $K=1$, then there is only one value/classification of that point, and it will be identical to \hat{y} . This is not a good estimate of future error as the new data does not necessarily have to belong to the same classification of y/\hat{y} .

Problem 4

These are questions about the linear model with $p = 1$.

- (a) [easy] What does \mathbb{D} look like in the linear model with $p = 1$? What is \mathcal{X} ? What is \mathcal{Y} ?

\mathbb{D} looks like $\langle \vec{x}, \vec{y} \rangle$, where both \mathcal{X} & \mathcal{Y} are column vectors.

- (b) [easy] Consider the line fit using the ordinary least squares (OLS) algorithm. Prove that the point $\langle \bar{x}, \bar{y} \rangle$ is on this line. Use the formulas we derived in class.

For $\mathcal{Y} = R$ or $\mathcal{Y} \subset R$, $p = 1$, our candidate sets are linear models. Then, for OLS, $g(x) = b_0 + b_1x$, where we derived $b_1 = r \frac{S_y}{S_x}$, $b_0 = \bar{y} - r \frac{S_y}{S_x} \bar{x}$. Now, if we plug in \bar{x} into $g(x)$ we get: $\bar{y} - r \frac{S_y}{S_x} \bar{x} + r \frac{S_y}{S_x} \bar{x} = \bar{y}$. ■

- (c) [harder] Consider the line fit using OLS. Prove that the average prediction $\hat{y}_i := g(x_i)$ for $x_i \in \mathbb{D}$ is \bar{y} .

As average is taking the sum and dividing by the number, we find that the 'average' prediction is: $\frac{\sum_{i=1}^n (\hat{y}_i)}{n} = \frac{\sum_{i=1}^n (b_0 + b_1 \hat{x}_i)}{n} = \frac{nb_0}{n} + \frac{nb_1 \sum_{i=1}^n (\hat{x}_i)}{n} = b_0 + b_1 \bar{x} = \bar{y}$. ■

- (d) [harder] Consider the line fit using OLS. Prove that the average residual e_i is 0 over \mathbb{D} .

$e = y - \hat{y}$. We know the average of y is \bar{y} , and from 4.d we know that the average prediction $\hat{y}_i = \bar{y} \therefore e_i = \bar{y} - \hat{y}_i = \bar{y} - \bar{y} = 0$. ■

- (e) [harder] Why is the RMSE usually a better indicator of predictive performance than R^2 ? Discuss in English.

RMSE is in the same unit as Y , and gives us a 95% confidence interval of containing the true Y . On the other hand, R^2 just gives us the relationship with the null model. If worse than the null model, R^2 is negative.

- (f) [harder] R^2 is commonly interpreted as “proportion of the variance explained by the model” and proportions are constrained to the interval $[0, 1]$. While it is true that $R^2 \leq 1$ for all models, it is not true that $R^2 \geq 0$ for all models. Construct an explicit example \mathbb{D} and create a linear model $g(x) = w_0 + w_1x$ whose $R^2 < 0$.

This would require us to have a terrible example, whereby our model is way off the data i.e. the null model has better performance.

We will let $\mathcal{D} = \langle \vec{X}, \vec{Y} \rangle$ where $\vec{X} = \langle 7, 17, 27, 37 \rangle$ while $\vec{Y} = \langle 37, 27, 17, 7 \rangle$. We will have $w_0 = -100$ and $w_1 = -1$, which will give terrible errors.

- (g) [difficult] You are given \mathbb{D} with n training points $\langle x_i, y_i \rangle$ but now you are also given a set of weights $[w_1 \ w_2 \ \dots \ w_n]$ which indicate how costly the error is for each of the i points. Rederive the least squares estimates b_0 and b_1 under this situation. Note that these estimates are called the *weighted least squares regression* estimates. This variant \mathcal{A} on OLS has a number of practical uses, especially in Economics. No need to simplify your answers like I did in class (i.e. you can leave in ugly sums).

When we reparametrized the Hypothesis space of line, we attached weights to features (except w_0 which was the intercept/bias term). Later on, while deriving least squared regression, at some point, we took the derivative w.r.t w_0 & w_1 . The weights given in this question, for Weighted Least Squares Regression, also has weights $[w_1, w_2, \dots w_n]$. We will denote the earlier weights as w_0^* & w_1^* to avoid confusion.

HW24g. $e_i^* = w_i e_i = w_i (y_i - \hat{y}_i)$, $SSE = \sum_{i=1}^n (e_i^*)^2 = \sum_{i=1}^n w_i^2 (y_i - \hat{y}_i)^2 = \sum_{i=1}^n w_i^2 (y_i - w_0^* - w_1^* x_i)^2$

$$= \sum w_i^2 y_i^2 + \sum w_i^2 w_0^{*2} + \sum w_i^2 w_1^{*2} x_i^2 - 2 \sum w_i^2 w_0^* y_i - 2 \sum w_i^2 w_1^* x_i y_i + 2 \sum w_i^2 w_0^* w_1^* x_i$$

$$= \sum w_i^2 y_i^2 + n w_0^{*2} \sum w_i^2 + n w_1^{*2} \sum w_i^2 x_i^2 - 2 n w_0^* \sum w_i^2 y_i - 2 n w_1^* \sum w_i^2 x_i y_i + 2 n w_0^* w_1^* \sum w_i^2 x_i$$

$$\frac{\partial [SSE]}{\partial w_0^*} = 2 n w_0^* \sum w_i^2 - 2 n \sum w_i^2 y_i + 2 n w_1^* \sum w_i^2 x_i \stackrel{\text{set}}{=} 0 \Rightarrow b_0 = \frac{\sum w_i^2 y_i - b_1 \sum w_i^2 x_i}{\sum w_i^2}$$

$$\frac{\partial [SSE]}{\partial w_1^*} = 2 n w_1^* \sum w_i^2 x_i^2 - 2 n \sum w_i^2 x_i y_i + 2 n w_0^* \sum w_i^2 x_i \stackrel{\text{set}}{=} 0 \Rightarrow b_1 = \frac{\sum w_i^2 x_i y_i - b_0 \sum w_i^2 x_i}{\sum w_i^2 x_i^2}$$

$$= b_1 \sum w_i^2 x_i^2 = \frac{(\sum w_i^2 x_i y_i) \sum w_i^2 - \sum w_i^2 x_i \sum w_i^2 y_i}{\sum w_i^2 x_i^2 \sum w_i^2 - (\sum w_i^2 x_i)^2}$$

$$= b_1 \sum w_i^2 x_i^2 \sum w_i^2 - b_1 \sum w_i^2 x_i \sum w_i^2 y_i = \sum w_i^2 x_i y_i \sum w_i^2 - \sum w_i^2 x_i \sum w_i^2 y_i$$

$$b_1 = \frac{(\sum w_i^2 x_i y_i) \sum w_i^2 - \sum w_i^2 x_i \sum w_i^2 y_i}{\sum w_i^2 x_i^2 \sum w_i^2 - (\sum w_i^2 x_i)^2}$$

$$b_1 = \frac{\sum w_i^2 x_i y_i - \sum w_i^2 x_i (\sum w_i^2 y_i / \sum w_i^2)}{\sum w_i^2 x_i^2 - \sum w_i^2 x_i (\sum w_i^2 y_i / \sum w_i^2)}$$

- (h) [harder] Interpret the ugly sums in the b_0 and b_1 you derived above and compare them to the b_0 and b_1 estimates in OLS. Does it make sense each term should be altered in this manner given your goal in the weighted least squares?

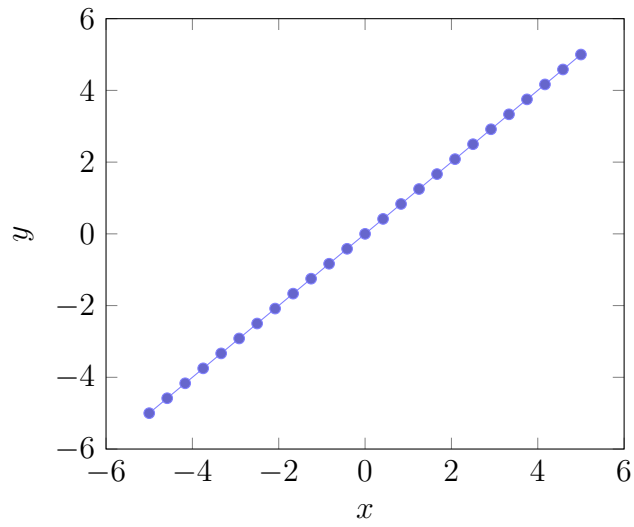
Yes, it makes sense, as our goal is to have the appropriate/estimate weight/impact of each error on our prediction.

- (i) [E.C.] In class we talked about $x_{raw} \in \{\text{red}, \text{green}\}$ and the OLS model was the sample average of the inputted x . Imagine if you have the additional constraint that x_{raw} is ordinal e.g. $x_{raw} \in \{\text{low}, \text{high}\}$ and you were forced to have a model where $g(\text{low}) \leq g(\text{high})$. Write about an algorithm \mathcal{A} that can solve this problem.

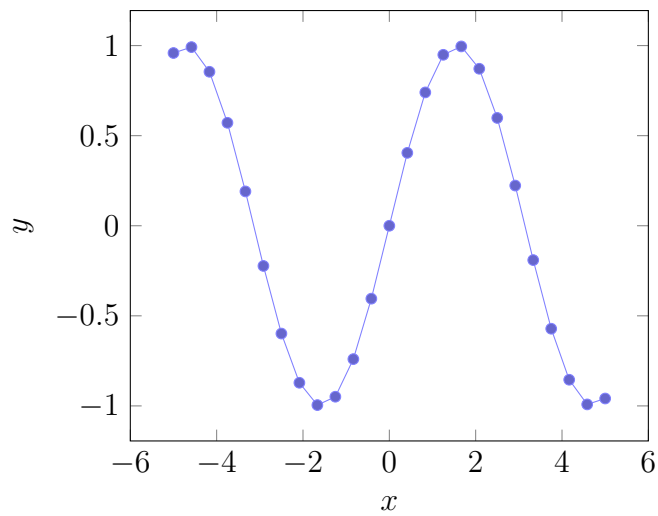
Problem 5

These are questions about association and correlation.

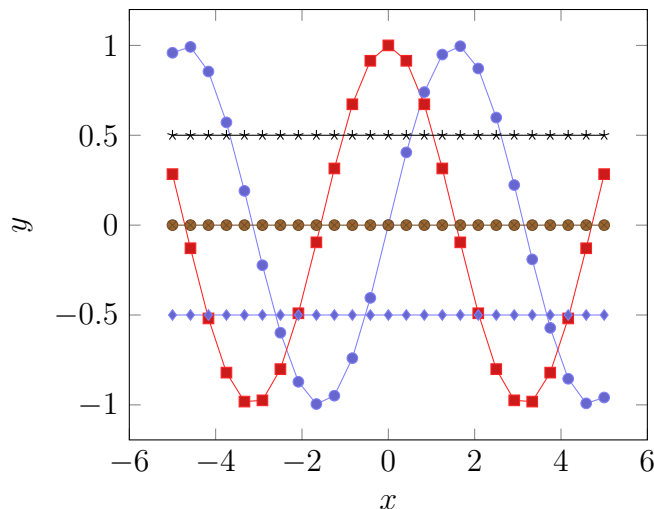
- (a) [easy] Give an example of two variables that are both correlated and associated by drawing a plot.



- (b) [easy] Give an example of two variables that are not correlated but are associated by drawing a plot.



- (c) [easy] Give an example of two variables that are not correlated nor associated by drawing a plot.



(d) [easy] Can two variables be correlated but not associated? Explain.

No, because correlation is a subset of association, as correlation is linear association.

Problem 6

These are questions about multivariate linear model fitting using the least squares algorithm.

(a) [difficult] Derive $\frac{\partial}{\partial \mathbf{c}} [\mathbf{c}^\top \mathbf{A} \mathbf{c}]$ where $\mathbf{c} \in \mathbb{R}^n$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ but *not* symmetric. Get as far as you can.

HW26a) $\mathbf{c} \in \mathbb{R}^n, \mathbf{A} \in \mathbb{R}^{n \times n}$, not symmetric, $\mathbf{c} = [c_1, c_2, \dots, c_n]^\top, \mathbf{c}^\top = [c_1, c_2, \dots, c_n]$, $\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$

$\mathbf{c}^\top \mathbf{A} \mathbf{c} = [c_1, c_2, \dots, c_n] \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} [c_1, c_2, \dots, c_n]^\top = [c_1, c_2, \dots, c_n] \begin{bmatrix} c_1 a_{11} + c_2 a_{12} + \dots + c_n a_{1n} \\ c_1 a_{21} + c_2 a_{22} + \dots + c_n a_{2n} \\ \vdots \\ c_1 a_{n1} + c_2 a_{n2} + \dots + c_n a_{nn} \end{bmatrix}$

$= c_1 \sum_{i=1}^n c_i a_{1i} + c_2 \sum_{i=1}^n c_i a_{2i} + \dots + c_n \sum_{i=1}^n c_i a_{ni} = \sum_{j=1}^n \sum_{i=1}^n c_j a_{ji} = \sum_{j=1}^n \left(\frac{\partial}{\partial c_j} (c_1 a_{1j} + c_2 a_{2j} + \dots + c_n a_{nj}) \right)$

$\frac{\partial}{\partial c_1} (c_1 a_{11} + c_2 a_{12} + \dots + c_n a_{1n}) = a_{11}$, $\frac{\partial}{\partial c_2} (c_1 a_{21} + c_2 a_{22} + \dots + c_n a_{2n}) = a_{21} + a_{22} + \dots + a_{2n}$, $\frac{\partial}{\partial c_k} (c_1 a_{k1} + c_2 a_{k2} + \dots + c_n a_{kn}) = a_{k1} + a_{k2} + \dots + a_{kn}$

(b) [easy] Given matrix $\mathbf{X} \in \mathbb{R}^{n \times (p+1)}$, full rank and first column consisting of the $\mathbf{1}_n$ vector, rederive the least squares solution \mathbf{b} (the vector of coefficients in the linear model shipped in the prediction function g). No need to rederive the facts about vector derivatives.

HW2.66) X , a matrix, $\in \mathbb{R}^{n \times (p+1)}$, full rank w/ 1st column = $\vec{1}_n$, so $X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ \vdots & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nn} \end{bmatrix}$

as $\hat{y}_1 = \vec{x}_1 \vec{\omega}$, $\hat{y}_2 = \vec{x}_2 \vec{\omega}$, \dots , $\hat{y} = X \vec{\omega} = \begin{bmatrix} \omega_0 + \omega_1 x_{11} + \omega_2 x_{12} + \dots + \omega_n x_{1n} \\ \omega_0 + \omega_1 x_{21} + \omega_2 x_{22} + \dots + \omega_n x_{2n} \\ \vdots \\ \omega_0 + \omega_1 x_{n1} + \omega_2 x_{n2} + \dots + \omega_n x_{nn} \end{bmatrix}$

$\vec{e} := \vec{y} - \hat{\vec{y}}$, $SSE = \sum_{i=1}^n e_i^2 = \vec{e}^T \vec{e} = (\vec{y} - \hat{\vec{y}})^T (\vec{y} - \hat{\vec{y}}) = (\vec{y}^T - \hat{\vec{y}}^T) (\vec{y} - \hat{\vec{y}})$

$= \vec{y}^T \vec{y} - \hat{\vec{y}}^T \vec{y} - \vec{y}^T \hat{\vec{y}} + \hat{\vec{y}}^T \hat{\vec{y}} = \vec{y}^T \vec{y} - 2 \hat{\vec{y}}^T \vec{y} + \hat{\vec{y}}^T \hat{\vec{y}}$

$= \vec{y}^T \vec{y} - 2 (X \vec{\omega})^T \vec{y} + (X \vec{\omega})^T X \vec{\omega} = \vec{y}^T \vec{y} - 2 \vec{\omega}^T X^T \vec{y} + \vec{\omega}^T X^T X \vec{\omega}$, so

$\frac{\partial SSE}{\partial \vec{\omega}} = \begin{bmatrix} \frac{\partial SSE}{\partial \omega_0} \\ \frac{\partial SSE}{\partial \omega_1} \\ \vdots \\ \frac{\partial SSE}{\partial \omega_n} \end{bmatrix} \stackrel{\text{set to } 0}{=} \text{solve for } b_0, b_1, \dots, b_p \text{ : } \vec{b} = (X^T X)^{-1} X^T \vec{y}$

- (c) [harder] Consider the case where $p = 1$. Show that the solution for \vec{b} you just derived in (b) is the same solution that we proved for simple regression. That is, the first element of \vec{b} is the same as $b_0 = \bar{y} - r \frac{s_y}{s_x} \bar{x}$ and the second element of \vec{b} is $b_1 = r \frac{s_y}{s_x}$.

HW26c) $p=1$ so $X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$, $X^T = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix}$, $X^T X = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$, $X^T y = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$

$= \begin{bmatrix} (1+1+\dots+1) (x_1+x_2+\dots+x_n) \\ (x_1+x_2+\dots+x_n) (x_1^2+x_2^2+\dots+x_n^2) \end{bmatrix}$, $\bar{y} = [y_1, y_2, \dots, y_n]$, $X^T y = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$

$= \begin{bmatrix} n\bar{y} \\ \sum x_i y_i \end{bmatrix}$, as $X^T X = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$, $(X^T X)^{-1} = \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{bmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{bmatrix}$

& $X^T y = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$, so $(X^T X)^{-1} X^T y$ would equal

$\begin{bmatrix} \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2} & \frac{-\sum x_i \sum y_i + n \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2} \end{bmatrix} = \begin{bmatrix} b_0 & b_1 \end{bmatrix}$

so $b_0 = \frac{\sum x_i^2 \bar{y} - \bar{x} \sum x_i y_i}{n \sum x_i^2 - n \bar{x}^2} = \frac{\sum x_i^2 \bar{y} - \bar{x} \sum x_i y_i}{\sum x_i^2 - n \bar{x}^2}$

$= \frac{\bar{y} \sum x_i^2}{S^2_x} - \frac{\bar{x} \sum x_i y_i}{S^2_x} + \frac{n \bar{x}^2 \bar{y}}{S^2_x} - \frac{n \bar{x}^2 \bar{y}}{S^2_x} - \frac{\bar{y} \sum x_i^2 - n \bar{x}^2 \bar{y}}{S^2_x} + \frac{\bar{x} \sum x_i y_i - n \bar{x} \bar{y} \sum x_i}{S^2_x}$

$= \frac{\bar{y} \sum x_i^2 - n \bar{x}^2}{S^2_x} - \frac{\bar{x} \sum x_i y_i - n \bar{x} \bar{y}}{S^2_x} = \frac{\bar{y} S^2_x - \bar{x} S_{xy}}{S^2_x}$

$= \bar{y} - \bar{x} r \frac{S_y}{S_x} = \bar{y} - \bar{x} r \frac{S_y}{S_x}$, &

$b_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} = \frac{n \sum x_i y_i - n \bar{x} \bar{y}}{n \sum x_i^2 - n \bar{x}^2} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$

$= \frac{S_{xy}}{S^2_x} = \frac{r S_x S_y}{S_x S_x} = r \frac{S_y}{S_x}$, $\therefore b_0 = \bar{y} - \bar{x} r \frac{S_y}{S_x}$, $b_1 = r \frac{S_y}{S_x}$

(d) [easy] If X is rank deficient, how can you solve for b ? Explain in English.

If X is rank deficient, then 1 or more of its columns are not linearly independent. This means they are not necessary. They must be removed so that we are left with columns that are linearly independent, and then we can solve for b .

(e) [difficult] Prove $\text{rank}[X] = \text{rank}[X^T X]$.

→ If $Xb = 0$ for some b , then $X^T Xb = 0$

← If $X^T Xb = 0$ for some b , then $b^T X^T Xb = 0 \Rightarrow Xb = 0 \therefore \text{rank}[X] = \text{rank}[X^T X]$

(f) [harder] [MA] If $p = 1$, prove $r^2 = R^2$ i.e. the linear correlation is the same as proportion of sample variance explained in a least squares linear model.

$$r^2 = \left(\frac{S_{xy}}{S_x S_y} \right)^2 = \frac{SST - SSE}{SST} = R^2$$

- (g) [harder] Prove that $g([1 \ \bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_p]) = \bar{y}$ in OLS.

$$\text{As } g(\vec{x}^*) = \hat{\mathbf{y}}^* = b_0 + b_1 x^*_1 + \dots + x^*_p$$

$$\text{then } g([1 \ \bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_p]) = b_0 + b_1 \bar{x}_1 + \dots + \bar{x}_n$$

$$\text{which} = \frac{b_0}{n} \sum (1) + \frac{b_1}{n} \sum (\bar{x}_1) + \dots + \frac{b_p}{n} \sum (\bar{x}_p)$$

$$\text{which} = \frac{1}{n} (b_0 \sum (1) + b_1 \sum (\bar{x}_1) + \dots + \sum (b_p \sum (\bar{x}_p)))$$

$$\text{which} = \frac{1}{n} = \frac{\hat{\mathbf{y}}}{n} = \underbrace{\bar{\mathbf{y}}}_{\text{from 4.c}} = \hat{\mathbf{y}}$$

- (h) [harder] Prove that $\bar{e} = 0$ in OLS.

$\bar{e} = \frac{1}{n} e_i$, as OLS is linear regression, the errors can either be positive (if the regression is above the data point) or negative (if the regression is below the data point). Add all these errors give us zero.

- (i) [difficult] If you model \mathbf{y} with one categorical nominal variable that has levels A, B, C , prove that the OLS estimates look like \bar{y}_A if $x = A$, \bar{y}_B if $x = B$ and \bar{y}_C if $x = C$. You can choose to use an intercept or not. Likely without is easier.

\vec{y} is one categorical nominal variable that has levels A, B, C, so \vec{x} would be a column vector of length n , while X would be an $n * 3$ matrix while X^T will be a $3 * n$ matrix. $\vec{b} = (X^T X)^{-1} X^T \vec{y}$, where $X^T X =$ a diagonal matrix of size three, which has nA , nB , & nC on its diagonal, so its inverse would have the reciprocal of those entries on the diagonal. $X^T \vec{y}$ would be a $3 * 1$ matrix with the entries being $\sum_{i=A} Y_i, \sum_{i=B} Y_i, \sum_{i=C} Y_i \therefore \vec{b}$ = a column vector of length 3: $[\bar{y}_A \bar{y}_B \bar{y}_C]^T$. As $x = A = g([100]), x = B = g([010]), x = C = g([001])$, their multiplication with \vec{b} would yield $\bar{y}_A, \bar{y}_B, \& \bar{y}_C$ respectively. ■

- (j) [harder] [MA] Prove that the OLS model always has $R^2 \in [0, 1]$.

OLS is always at least as good as the null model. Therefore, even when there is an error on all points, the worst it can be is the null model, which gives and R^2 of 0 as $R^2 = \frac{SST - SSE}{SST}$, while when the model is super good, SSE is almost 0 and R^2 is close to 100%.