

Quiz 6 occurs December 2, in the discussion section. The quiz will be based upon the problems below.

Quiz 6 Problems

Exercise 1. Suppose we are presented with data points $(x_1, y_1), \dots, (x_n, y_n)$. We would like to find the line $y = mx + b$ which lies “closest” to all of these data points. Such a line is known as a **linear regression**. There are many ways to define the “closest” such line. The standard method is to use **least squares minimization**. A line which lies close to all of the data points should make the quantities $(y_i - mx_i - b)$ all very small. We would like to find numbers m, b such that the following quantity is minimized:

$$f(m, b) = \sum_{i=1}^n (y_i - mx_i - b)^2.$$

Using the second derivative test (or a convexity argument), show that the minimum value of f is achieved when

$$m = \frac{(\sum_{i=1}^n x_i)(\sum_{j=1}^n y_j) - n(\sum_{k=1}^n x_k y_k)}{(\sum_{i=1}^n x_i)^2 - n(\sum_{j=1}^n x_j^2)} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{j=1}^n (x_j - \bar{x})^2}.$$

$$b = \frac{1}{n} \left(\sum_{i=1}^n y_i - m \sum_{j=1}^n x_j \right) = \bar{y} - m\bar{x}.$$

Briefly explain why this is actually the minimum value of $f(m, b)$.

Exercise 2. I believe that the number of home runs hit by an MLB baseball player in a single season is linearly related to the number of strikeouts they have in a single season.

Here the data can be found from:

<http://www.seanlahman.com/baseball-archive/statistics>

I recommend using the 2020 Version, comma delimited version. The data is in a zip file, and home run data can be found in `Core` then `batting.csv` then the columns `HR` and `SO`.

That is, $(x_1, y_1), (x_2, y_2), \dots$ is your data, where x_i is the number of home runs hit by the i^{th} person (or i^{th} row) in the data file, and y_i is the number of strikeouts of the i^{th} person (or i^{th} row) in the data file. And we need to find the slope m and intercept b of the line

$$y = mx + b$$

that best fits the data. That is, use the least squares linear regression line.

Plot the data, together with the best fit line. Does the line fit the data well (visually)?

As a test of goodness of fit, compute the quantity

$$\frac{1}{n} \sum_{i=1}^n \phi(y_i - [mx_i + b])$$

where $\phi(t) = 1 - e^{-t^2/10}$ for all $t \in \mathbf{R}$, and n is the number of data points. (Note that $\phi(0) = 0$ and $\lim_{t \rightarrow \pm\infty} \phi(t) = 1$.) (This quantity is therefore between 0 and 1.) Is this quantity close to 0?

Do all above steps again for the following different statement:

I also believe that the number of stolen bases an MLB baseball player in a single season is linearly related to the number of hits they have in a single season (see the columns SB and H).

Finally, do all above steps again for the following statement:

I also believe that the number of hits of an MLB baseball player in a single season is approximately a constant plus the square of the number of doubles they have in a single season (see the columns 2B and H). That is, $(x_1, y_1), (x_2, y_2), \dots$ is your data, where y_i is the number of hits by the i^{th} person (or i^{th} row) in the data file, and x_i is the number of doubles of the i^{th} person (or i^{th} row) in the data file. And we need to find the parameters m, b of the parabola

$$y = mx^2 + b$$

that best fits the data. That is, use the least squares linear regression of m, b .

As a test of goodness of fit, compute the quantity

$$\frac{1}{n} \sum_{i=1}^n \phi(y_i - [mx_i^2 + b])$$

where $\phi(t) = 1 - e^{-t^2/10}$ for all $t \in \mathbf{R}$, and n is the number of data points. (Note that $\phi(0) = 0$ and $\lim_{t \rightarrow \pm\infty} \phi(t) = 1$.) (This quantity is therefore between 0 and 1.) Is this quantity close to 0?

Exercise 3. Suppose you have three vegetarian turkeys with numerical “quality” values of

$$Y_i = \beta_1 + \varepsilon_i, \quad \forall 1 \leq i \leq 3$$

Suppose you have three vegetarian turkeys with numerical “quality” values of

$$Y_i = \beta_2 + \varepsilon_i, \quad \forall 4 \leq i \leq 6$$

And suppose you have three vegetarian turkeys with numerical “quality” values of

$$Y_i = \beta_3 + \varepsilon_i, \quad \forall 7 \leq i \leq 9.$$

Here $\varepsilon_1, \dots, \varepsilon_9$ are i.i.d. Gaussians with mean zero and unknown variance $\sigma^2 > 0$, and $\beta_1, \beta_2, \beta_3 \in \mathbf{R}$ are unknown.

In order to test the null hypothesis that $\beta_1 = \beta_2 = \beta_3$, we perform the F test, i.e. we evaluate

$$F := \sup_{c_1, \dots, c_3 \in \mathbf{R}: \sum_{i=1}^3 c_i = 0} \frac{\left(\sum_{j=1}^3 c_j \bar{Y}_j - \sum_{j=1}^3 c_j \beta_j \right)^2}{S^2 \sum_{j=1}^3 \frac{c_j^2}{3}}$$

Report a p -value for F for the observation that:

$$Y_1 = 5, Y_2 = 6, Y_3 = 7$$

$$Y_4 = 5, Y_5 = 5, Y_6 = 5$$

$$Y_7 = 6, Y_8 = 7, Y_9 = 8.$$

Do you have confidence in accepting the null hypothesis?

Exercise 4. Let

$$h(x) := \frac{1}{1 + e^{-x}}, \quad \forall x \in \mathbf{R}.$$

Fix $x \in \mathbf{R}$ and $y \in [0, 1]$. Define $t: \mathbf{R}^2 \rightarrow \mathbf{R}$ by

$$t(a, b) := \log \left([h(ax + b)]^y [1 - h(ax + b)]^{1-y} \right), \quad \forall a, b \in \mathbf{R}.$$

Show that t is concave. Conclude that t has at most one global maximum.

Exercise 5. Consider the following table with turkey data. We have 8 (vegetarian) turkeys, with various temperatures x (Fahrenheit), and the status y of each turkey is cooked (corresponding to a value of $y = 1$) or not cooked (corresponding to a value of $y = 0$). Using logistic regression, find $a, b \in \mathbf{R}$, i.e. find a function

$$h(ax + b)$$

that best fits your data, where $h(t) = 1/(1 + e^{-t})$ for all $t \in \mathbf{R}$.

That is, given a temperature x , $h(ax + b)$ should be close to 1 when the turkey is cooked, and $h(ax + b)$ should be close to 0 when the turkey is not cooked.

Turkey	Temperature	Done? Yes or no.
1	140	no
2	145	no
3	150	no
4	155	yes
5	160	no
6	165	yes
7	170	yes
8	175	yes

Exercise 6 (Optional). Suppose $\beta \in \mathbf{R}^m$ is an unknown vector, and A is a known $m \times n$ real matrix. Let $\varepsilon \in \mathbf{R}^n$ be a vector of i.i.d. standard Gaussian random variables. Our observation is $Y := A\beta + \varepsilon$, and the goal is to recover the unknown vector w . In linear least squares regression, we try to determine the best linear relationship w between the rows of A

and the observation Y . Assume that $n \leq m$ and the matrix A has full rank (so that $A^T A$ is invertible). Show that the vector $x \in \mathbf{R}^m$ that minimizes the quantity

$$\|Y - Ax\|^2 := \sum_{i=1}^n (y_i - (Ax)_i)^2$$

is

$$x := (A^T A)^{-1} A^T Y. \quad (*)$$

Equivalently, show that x minimizes

$$\mathbf{E} \|Y - x\|^2$$

over all choices of vectors $x \in \mathbf{R}^m$ such that $x = By$ for some $n \times m$ real matrix B , and such that $\mathbf{E}x = w$. (Since ε is the only random variable here, \mathbf{E} denotes expected value with respect to ε .)