

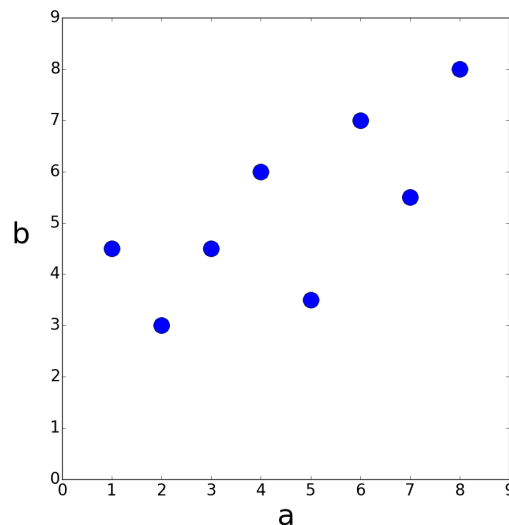
**This homework is due April 12, 2016, at Noon.**

**1. Homework process and study group**

Who else did you work with on this homework? List names and student ID's. (In case of hw party, you can also just describe the group.) How did you work on this homework?

Working in groups of 3-5 will earn credit for your participation grade.

**2. Mechanical: Linear Least Squares**



a	1	2	3	4	5	6	7	8
b	4.5	3	4.5	6	3.5	7	5.5	8

(a) Consider the above data points. Find the linear model of the form

$$b = xa \quad (1)$$

that best fits the data, i.e. find  $x$  to minimize

$$\left\| \begin{bmatrix} b_1 \\ \vdots \\ b_8 \end{bmatrix} - \begin{bmatrix} a_1 \\ \vdots \\ a_8 \end{bmatrix} x \right\|^2 \quad (2)$$

**Do not use Python for this calculation and show your work. (A calculator is okay).** Once you've computed  $x$ , compute the squared error between your model's prediction and the actual  $b$  values as shown in Equation (2). Plot the best fit line along with the data points to examine the quality of the fit. (If you're plotting by hand, it is okay if your plot of  $b = xa$  is approximate.)

- (b) You will notice from your graph that you can get a better fit by adding a  $b$ -intercept. That is we can get a better fit for the data by assuming a linear model of the form

$$b = x_1 a + x_2 \quad (3)$$

In order to do this, we need to augment our  $A$  matrix for the least squares calculation with a column of 1's (do you see why?) so that it has the form

$$A = \begin{bmatrix} a_1 & 1 \\ \vdots & \vdots \\ a_8 & 1 \end{bmatrix} \quad (4)$$

Find  $x_1$  and  $x_2$  to minimize

$$\left\| \begin{bmatrix} b_1 \\ \vdots \\ b_8 \end{bmatrix} - \begin{bmatrix} a_1 & 1 \\ \vdots & \vdots \\ a_8 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right\|^2 \quad (5)$$

**Again, do not use python for this calculation and show your work. A calculator is okay but take the inverse by hand using the formula for a  $2 \times 2$  inverse.**

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \quad (6)$$

Compute the squared error between your model's prediction and the actual  $b$  values as shown in Equation (5). Plot your new linear model. Is it a better fit for the data?

- (c) Let  $\vec{\hat{x}}$  be the solution to a general linear least squares problem,

$$\vec{\hat{x}} = \underset{\vec{x}}{\operatorname{argmin}} \left\| \vec{b} - A\vec{x} \right\|^2 \quad (7)$$

Show that the error vector  $\vec{b} - A\vec{\hat{x}}$  is orthogonal to the columns of  $A$  by direct manipulation. (*i.e. plug in the formula for the linear-least-squares estimate into the error vector. And then see if this vector is such that  $A^T$  times it is zero.*)

### 3. Retail Store Marketing

**Intro** The retail store EehEeh Sixteen would like to create a smart system where it decides which promotion to give to its customers when they checkout, depending on things they may be interested in. The promotion is supposed to be printed alongside the receipt and be used during their next purchase. The problem is, the customers don't disclose what their interests are when they checkout, and the only data the retail store can use are their current purchase data.

**The Setting** The store uses the following set of attributes in their decision making process: interest in party products, interest in family products, interest in student products and interest in office products. These attributes are used to describe each of the promotions the store offers. More concretely, the store attaches

to each promotion  $A$ , a "score" vector  $\vec{s}_A \in \mathbb{R}^4$  such that  $\vec{s}_A = \begin{bmatrix} \text{party-related score} \\ \text{family-related score} \\ \text{student-related score} \\ \text{office-related score} \end{bmatrix}$  which describes the

Interest Category	Spending Category			
	Food	Movies	Art	Books & Supplies
Party	40%	33%	22%	5%
Family	70%	10%	10%	10%
Student	20%	10%	15%	55%
Office	5%	2%	20%	73%

Table 1: The distribution of spending of people in each category.

ideal target customer. Therefore, the store would like to infer these same attributes about each customer at time of checkout so that they can print a promotion tailored to that customer on the receipt.

The data that the algorithm is allowed to use are the subtotals (in the current purchase) in the following four categories: Food, movies, art, and books & supplies.

**The Goal** EehEeh Sixteen hired the same intern from the Framingham heart study to devise an algorithm that takes a customer's purchase subtotals in the four categories listed above (food, movies, art and books & supplies), and decides which promotion to print on the receipt. The intern is lost and given the awesomeness of your help last time, he needs your help again. In this problem, you will walk him through a possible design of such an algorithm.

- (a) Assuming we somehow have the interests of a customer  $c$  in a vector  $\vec{x}_c = \begin{bmatrix} c_{\text{party}} \\ c_{\text{family}} \\ c_{\text{student}} \\ c_{\text{office}} \end{bmatrix}$  and a set of promotions  $A_1, A_2, \dots, A_N$ , with their attached vectors of scores  $\vec{s}_{A_1}, \vec{s}_{A_2}, \dots, \vec{s}_{A_N}$ . We would like to select which promotion is best aligned with the preferences of the customer. Assuming we have a function  $\text{sim}(\vec{x}_c, \vec{s}_A)$  which outputs a similarity score (higher score means more similar) between the customer  $c$  and the promotion  $A$ , how can we select which promotion to print to the customer on her receipt?
- (b) Would  $\text{sim}_1(\vec{x}_c, \vec{s}_A) = \|\vec{x}_c - \vec{s}_A\|$  be a good similarity measure? Why? What about  $\text{sim}_2(\vec{x}_c, \vec{s}_A) = \frac{1}{\|\vec{x}_c - \vec{s}_A\|}$ ? Why? What about  $\text{sim}_3(\vec{x}_c, \vec{s}_A) = \langle \vec{x}_c, \vec{s}_A \rangle$ ? Why? What about  $\text{sim}_4(\vec{x}_c, \vec{s}_A) = \left\langle \vec{x}_c, \frac{\vec{s}_A}{\|\vec{s}_A\|} \right\rangle$ ? Why?
- (c) The intern hands you research that the EehEeh Sixteen research division conducted, which calculated the distribution of spending in the store for people who are purely interested in only one category. The results are depicted in Table 1. Use this information to devise a system of linear equations, such that solving this system will result in the customer's preferences given her spending.
- (d) Combine these results into a complete algorithm.
- (e) Run the algorithm on a customer, Jane Doe, that spent \$6 on food, \$4 on movies, \$1 on art and \$5

on books. With promotions  $A_1, A_2, A_3$  and  $A_4$  targeted at customers with preferences  $\vec{s}_{A_1} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$ ,

$$\vec{s}_{A_2} = \begin{bmatrix} \frac{2}{3} \\ -\frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{3} \end{bmatrix}, \vec{s}_{A_3} = \begin{bmatrix} -\frac{1}{2} \\ \frac{1}{2} \\ \frac{5}{2} \\ -\frac{1}{2} \end{bmatrix} \text{ and } \vec{s}_{A_4} = \begin{bmatrix} 0 \\ \frac{1}{2} \\ 0 \\ \frac{1}{2} \end{bmatrix}$$

- (f) Will there ever be a customer for which the system devised in part (c) will yield no solutions or infinite solutions?

#### 4. Labelling patients using gene expression data

Least-squares techniques are useful for many different kinds of prediction problems. The core ideas we learned in class have been extensively further developed. These ideas are commonly used in machine learning for finance, healthcare, advertising, image processing and many other fields. Here we'll explore how least squares can be used for classification of data in a medical context.

Gene expression data of patients, along with other factors such as height, weight, age, family history, is often used to understand the likelihood that a patient might develop certain common diseases such as diabetes. Gene expression profiles can be read using DNA microarray technology, which uses tissue samples from a patient. This data, along with the patient specific characteristics above, can be combined into a vector to get a set of features that describe each patient.

Many scientific studies look at models in mice to understand how gene expression relates to diabetes. Previous studies have shown that the expression of the *tomosin2* and *ts1* genes are correlated to the onset of diabetes in mice. How can we predict whether or not a mouse will develop diabetes based on data about this expression as well as other factors of the mouse? We will use some (fake) data to explore this.

We are given information about the age and weight of the mouse, and additionally have access to data about whether the genes *tomosin2*, *ts1* and *chn1* (a third gene) were expressed or not. The gene expression data is captured using vectors that are  $+1$  if the gene is expressed and  $-1$  if the gene is not expressed. Similarly, whether or not a mouse has diabetes is also captured using a  $+1, -1$  vector, where  $+1$  indicates that the mouse has diabetes. Using this data we would like to develop a linear model that predicts whether or not a mouse will have diabetes.

$$\alpha_1(\text{age}) + \alpha_2(\text{weight}) + \alpha_3(\text{tomosin2}) + \alpha_4(\text{ts1}) + \alpha_5(\text{chn1}) \quad (8)$$

We would like the above expression to be positive if the mouse has diabetes and negative if the mouse does not have diabetes.

- (a) In problems such as this, it is common to use some *training* data to generate a model. Turns out, a good heuristic for this can be developed using a least squares technique. Set up a linear model for the problem in a format we have used for least-squares problems  $A\vec{x} = \vec{b}$ . Here,  $\vec{b}$  will be a vector with  $+1, -1$  entries.  $\alpha_i$ 's are your unknowns.
- (b) Using the (fake) *training* data `diabetes_train.npy`, generate the linear model using the least squares technique, i.e. find the unknown model parameters for the given data set. Include the unknown parameter values in the writeup of your homework. Use the provided ipython notebook file.
- (c) Now it is time to use the model you have developed to make some predictions! It is interesting to note here that we are not looking for a real number to model whether each mouse has diabetes or not, we are looking for a binary label. So we will use the *sign* of the expression above to assign a  $\pm 1$  value to each mouse. Predict whether each mouse with the characteristics in the *test* data set `diabetes_test.npy` will get diabetes. There are four mice in the test data set. Include the  $\pm 1$  vector that indicates whether or not they have diabetes in your writeup.

#### 5. Image analysis

Applications in medical imaging often require an analysis of images based on the pixels of the image. For instance, we might want to count the number of cells in a given sample. One way to do this is to “take a

picture” of the cells and use the pixels to determine the locations and thus the number of cells. Alternatively, automatic detection of shape is useful in image classification as well (e.g. consider a robot autonomously trying to find out where a mug is in it’s field of vision).

Let us focus back on the medical imaging scenario. You are interested in finding the exact position and shape of a cell in an image, so you want to find the equation of the ellipse that bounds the cell relative to a given coordinate system that is represented by the image. Your collaborator uses edge detection techniques to find a bunch of points that are approximately along the edge of the cell. We assume that the origin is in the center of the image with standard axes and collect the following points:  $(.3, -.7)$ ,  $(.5, .91)$ ,  $(.9, -.99)$ ,  $(1, 1.01)$ ,  $(1.2, -.93)$ ,  $(1.5, .8)$ ,  $(2, 0)$ . Submit your code for all parts of this problem, feel free to add it to the original iPython notebook.

Recall that a quadratic equation of the form

$$ax^2 + bxy + cy^2 + dx + ey = 1. \quad (9)$$

can be used to represent an ellipse (if  $b^2 - 4ac < 0$ ), and a quadratic equation of the form

$$a(x^2 + y^2) + dx + ey = 1 \quad (10)$$

is a circle if  $d^2 + e^2 - 4a > 0$ . The circle has fewer parameters.

- How can you find the equation of a circle that surrounds the cell? First, provide a setup and formulate a set of matrix equations to do this, i.e. an equation of the form  $A \cdot \vec{x} = \vec{b} + \vec{e}$ , where  $\vec{b}$  represents your observations and  $\vec{e}$  represents the unknown errors.
- How can you find the equation of an ellipse that surrounds the cell? Provide a setup and formulate a set of matrix equations to do this as above.
- Write a short program in iPython to fit a circle using these points. If you model your system of equations as  $A\vec{x} = \vec{b} + \vec{e}$  where  $\vec{e}$  is the error vector and the number of data points is  $N$ , what is  $\frac{\|\vec{e}\|}{N}$ ? Plot your points and the best fit circle in iPython.
- Write a short program in iPython to fit an ellipse using these points. If you model your system of equations as  $A\vec{x} = \vec{b} + \vec{e}$  where  $\vec{e}$  is the error vector and the number of data points is  $N$ , what is  $\frac{\|\vec{e}\|}{N}$ ? Plot your points and the best fit ellipse in iPython. How does this error compare to the one in the previous subpart? Which technique is better?

## 6. GPS Receivers

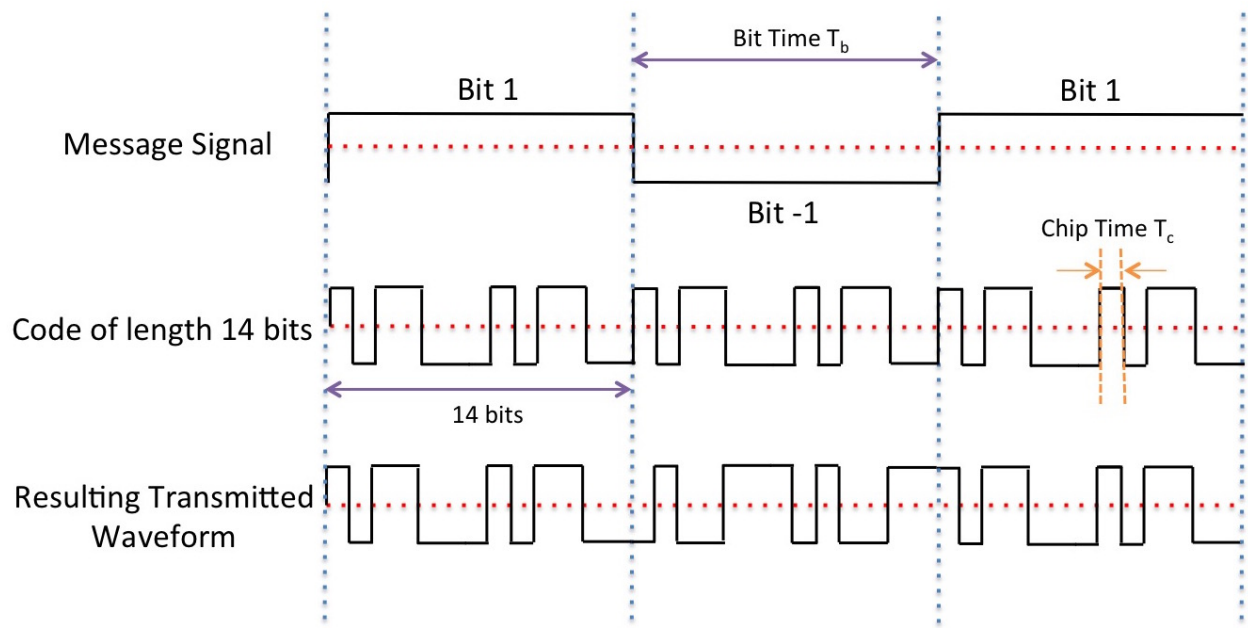
The Global Positioning System (GPS) is a space-based satellite navigation system that provides location and time information in all weather conditions, anywhere on or near the Earth where there is an unobstructed line of sight to four or more GPS satellites. Gold All satellites use the same carrier frequency to transmit the signals. To permit this without undue interference between the users, the satellites employ “spread-spectrum” technology (very similar to CDMA) and a special coding scheme where each transmitter is assigned a code that serves as a “signature”. An example is depicted in the figure below. The *message signal* to be transmitted changes at a much slower timescale than the timescale of the signature code, which is very fast.

In the example figure, we are showing a message signal that is a stream of  $+1$  or  $-1$  values. The signature code is also a stream of  $+1$  or  $-1$  values of length 14 bits. The signature code is multiplied by the appropriate *message signal* to get the final transmitted waveform. Let  $T_b$  be the “bit time”, i.e. time for each message bit and  $T_c$  be the “chip time” which is the time for each new symbol of the code to be generated. In the figure,  $T_b = 14T_c$ .

The GPS satellites use “Gold codes” which are 1023 bits long. So, for our problem,  $T_b = 1023T_c$ . (In reality,  $T_b = 20 \times 1023 \times T_c$ .) The Gold codes have special properties:

- Their auto-correlation of a Gold code (correlation with itself) is very high.
- The cross-correlation between different codes is very low.

These codes are generated using a linear feedback shift register (LFSR). You can read more about this if you are interested but for the problem you don’t need to know how this works. The take-away is that the Gold codes are vectors of  $+1$  and  $-1$  values that are “almost orthogonal”.



For the purpose of this question we only consider 24 GPS satellites. Download the file `prob10.ipynb` and the corresponding data files for the following questions:

- Auto-correlate the Gold code of satellite 10 with itself and plot it. Python has functions for this. What do you observe?
- Cross-correlate the Gold code of satellite 10 with satellite 13 and plot it. What do you observe?
- Now, consider a random signal, i.e. a signal that is not generated due to a specific code but is a random  $\pm 1$  sequence. Cross-correlate it with the Gold code of satellite 10. What do you observe? How does this compare to the cross-correlation of satellite 10 and satellite 13? What does this mean about our ability to identify satellites?
- The signals received by a receiver include signals from the satellites as well as an additional noise term. This is often modeled as a Gaussian noise term. You don’t need to understand Gaussians here but we use them since they form a good model for how the transmitted signal might be perturbed (large perturbations are very unlikely, and small perturbations are more likely).

Use the Gaussian noise generator to generate a random vector of length 1023, and cross correlate this with the Gold code of satellite 10. What do you observe?

For the next subparts of this problem, the signal is corrupted by Gaussian noise. Use the observation from this subpart for solving the rest of the question.

- (e) Now, assume that signals from multiple satellites are added at the receiver. So the signatures of multiple different satellites are present in the code. In addition, noise might be added to the signal. What are the satellites present in `data1.npy`?
- (f) Let's assume that you can hear only one satellite, Satellite A, at the location you are in (though this never happens in reality). Let's also assume that this satellite is transmitting a length 5 sequence of  $+1$  and  $-1$  after modulating it onto the 1023 bit Gold code corresponding to Satellite A. Find out from `data2.npy` which satellite it is and what sequence of  $\pm 1$  it is transmitting.
- (g) For the purpose of this problem, we'll assume that all the satellites transmit the same unique sequence of  $+1$ s and  $-1$ s. These are transmitted using the procedure described in the figure (called modulation, which we will learn more about soon.)

Signals from different transmitters arrive at the receiver with different propagation delays. So effectively the signals from different satellites are superimposed on each other with different offsets at the start. This propagation delay is used to find out how far the satellite is from the receiver. To find out exactly what the offset due to propagation delay is, you want to figure out the starting point of the signal transmission, and you can do this by cross-correlating the signature codes with the received signal at different offsets. What do you expect to observe when you cross-correlate the signature for a particular satellite (say Satellite A) with the received signal at the offset corresponding to the propagation delay of Satellite A?

What satellites are you able to see in `data3.npy` and what are the relative delays assuming that the message signal that was being sent was exactly  $[1 \ 1 \ -1 \ -1 \ -1]$ .

- 7. Your Own Problem** Write your own problem related to this week's material and solve it. You may still work in groups to brainstorm problems, but each student should submit a unique problem. What is the problem? How to formulate it? How to solve it? What is the solution?