1.  What is **gradient descent**? Why is it important?

2.  What is **Logistic Regression**? What is "logistic"? What are we "regressing"?

3.  Bob tries to gather information about this year's apple harvest and ran a search in his favorite online news outlet. He retrieved a number of articles but found that a large portion of the retrieved articles are about the Apple laptops and computers -- and hence irrelevant to his search. He wants to build a logistic regression classifier, which uses the counts of selected words in the news articles to predict the class of the news article (fruit vs. computer). He built the following data set of 5 training instances and 1 test instance. Develop a logistic regression classifier to predict label $\hat{y} = 1$ (fruit) and $\hat{y} = 0$ (computer).

| ID | apple | ibm | lemon | sun | | CLASS |
|----|-------|-----|-------|-----|---|-------|
| TRAINING INSTANCES | | | | | | |
| A | 1 | 0 | 1 | 5 | 1 | FRUIT |
| B | 1 | 0 | 1 | 2 | 1 | FRUIT |
| C | 2 | 0 | 0 | 1 | 1 | FRUIT |
| D | 2 | 2 | 0 | 0 | 0 | COMPUTER |
| E | 1 | 2 | 1 | 7 | 0 | COMPUTER |
| TEST INSTANCES | | | | | | |
| T | 1 | 2 | 1 | 5 | | ? |

For the moment, we assume that we already have an estimate of the model parameters, i.e., the weights of the 4 features (and the bias $\theta_0$) is $\hat{\theta} = [\theta_0, \theta_1, \theta_2, \theta_3, \theta_4] = [0.2, 0.3, -2.2, 3.3, -0.2]$.

  (i).  Explain the intuition behind the model parameters, and their meaning in relation to the features

  (ii).  Predict the test label.

  (iii).  Recall the conditional likelihood objective

$$\log \mathcal{L}(\theta) = -\sum_{i=1}^{n} y_i \log\big(\sigma(x_i; \theta)\big) + (1 - y_i) \log\big(1 - \sigma(x_i; \theta)\big)$$

We want to make sure that the Loss (the negative log likelihood) our model, is lower when its prediction the correct label for test instance T, than when it's predicting a wrong label.

Compute the negative log-likelihood of the test instance (1) assuming that the true label y = 1 (fruit), i.e., our classifier made a mistake; and (2) assuming the true label as y = 0 (computer), i.e., our classifier predicted correctly.

4. For the model created in question 4, compute a single gradient descent update for parameter $\theta_1$ given the training instances given above. Recall that for each feature j, we compute its weight update as

$$\theta_j \leftarrow \theta_j - \eta \sum_i (\sigma(x_i; \theta) - y_i) x_{ij}$$

Summing over all training instances $i$. We will compute the update for $\theta_j$ assuming the current parameters as specified above, and a learning rate $\eta = 0.1$.

5. [OPTIONAL] What is the relation between "odds" and "probability"?

6. [OPTIONAL] (a) What is **Regression**? How is it similar to **Classification**, and how is it different?

(b) Come up with one typical classification task, and one typical regression task. Specify the range of valid values of $y$ (results) and possible valid values for $x$ (attributes).