

CS 5875 Applied Machine Learning

Homework 1

Weiming Zhang, Yan Jiang

I. OBJECTIVE

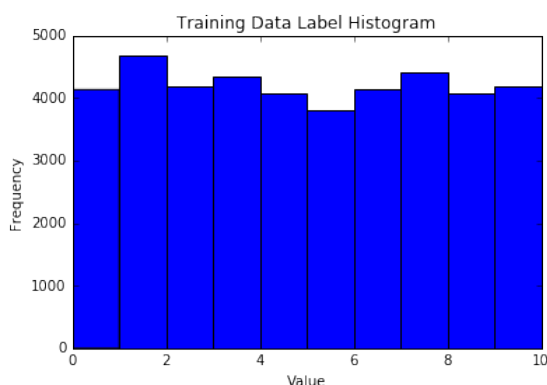
Working on KNN classification method and logistic regression.

II. PROBLEM 1: DIGIT RECOGNIZER

- b) Write a function to display an MNIST digit. Display one of each digit

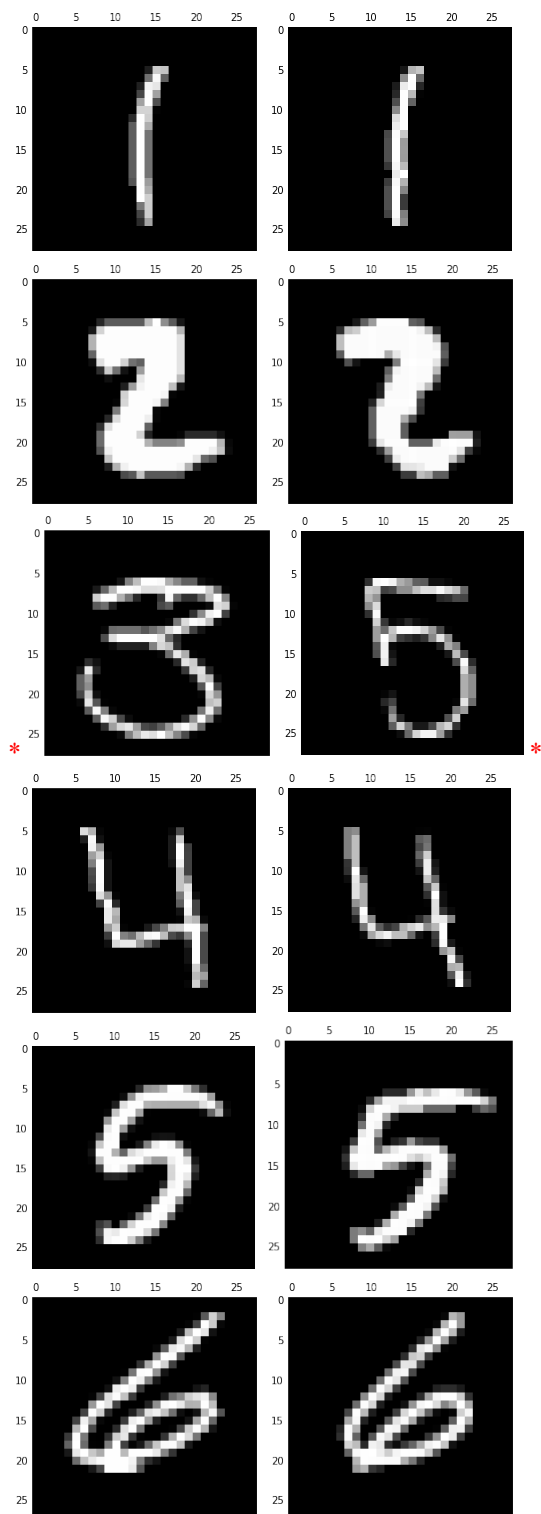
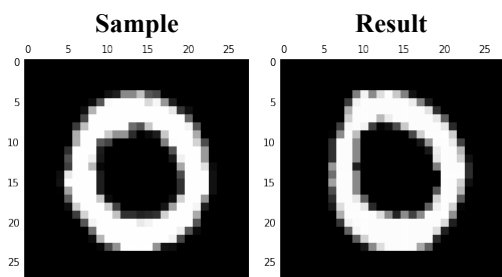
See source code function `displayDigit()`.

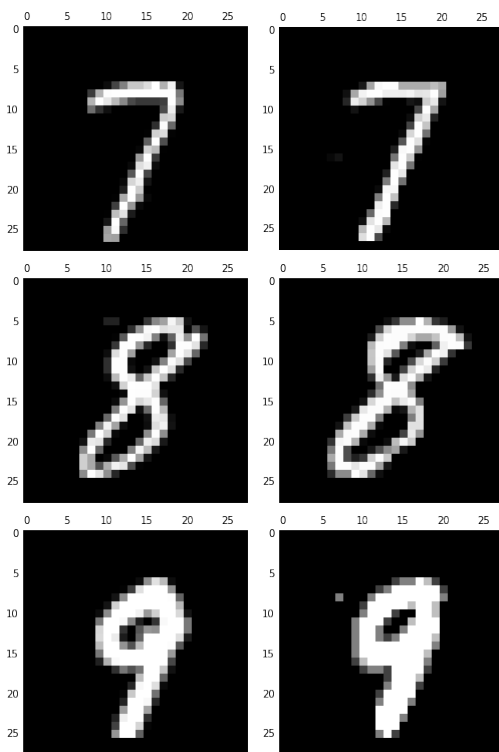
- c) Examine the prior probability of the classes in the training data. Is it uniform across the digits? Display a normalized histogram of digit counts. Is it even?



The distribution of the prior probability of the classes in the training data is uniform across the digits. It's evenly distributed. See code for this part in `showLabelDistribution()`.

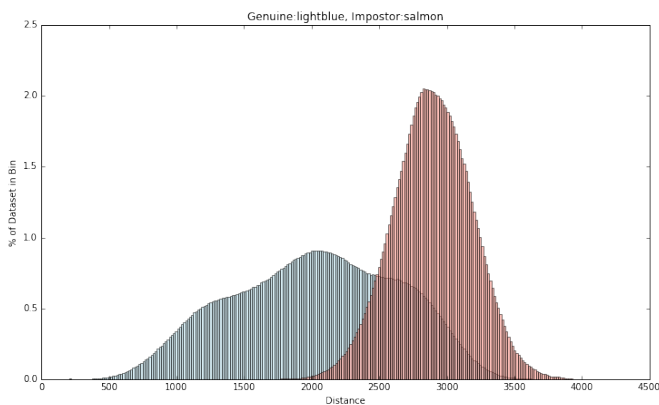
- d) Pick one example of each digit from your training data. Then, for each sample digit, compute and show the best match (nearest neighbor) between your chosen sample and the rest of the training data. Use L2 distance between the two images' pixel values as the metric. This probably won't be perfect, so add an asterisk next to the erroneous examples.





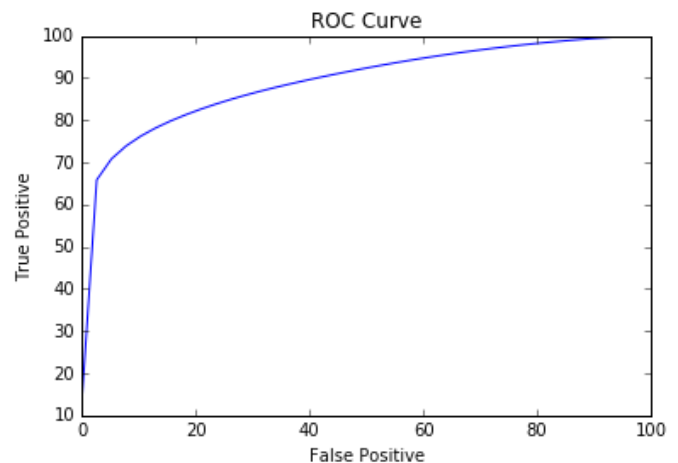
Note for digit 3 the prediction is incorrect (marked with an asterisk). See code for this part in `classify()` in the first panel.

- e) Consider the case of binary comparison between the digits 0 and 1. Ignoring all the other digits, compute the pairwise distances for all genuine matches and all impostor matches, again using the L2 norm. Plot histograms of the genuine and impostor distances on the same set of axes.



See code for this part in `genuineAndImpostor()`.

- f) Generate an ROC curve from the above sets of distances. What is the equal error rate? What is the error rate of a classifier that simply guesses randomly?



Equal error rate is around 18%

For random guess equal error rate is 50%

See code for this part in `rocCurve()`.

- g) Implement a K-NN classifier.

See code for this implementation in the **second panel** of the **Question1.ipynb**

- h) Using the training data for all digits, perform 3 fold cross-validation on your K-NN classifier and report your average accuracy.

K = 3

Average: 96.55%

See code for this part in `threefold()`.

- i) Generate a confusion matrix (of size 10×10) from your results. Which digits are particularly tricky to classify?

[1408	0	0	0	0	1	2	0	0	0]
[0	1573	2	0	1	0	0	0	0	1]
[14	10	1291	3	0	2	2	23	6	2]
[2	6	13	1385	0	13	0	7	10	6]
[0	13	0	0	1260	0	7	1	1	34]
[5	2	0	23	1	1208	17	0	5	13]
[8	2	0	0	1	4	1343	0	0	0]
[1	12	3	0	3	0	0	1440	0	18]
[7	25	4	12	5	19	6	4	1272	15]
[7	3	2	10	10	4	0	24	4	1359]

From the matrix we see that digits such as 3, 5 and 9 are relatively tricky to classify.

See code at `confusionMatrix()`.

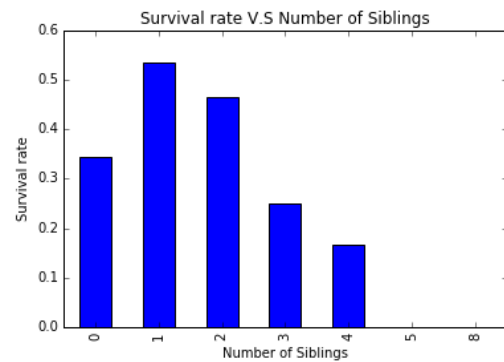
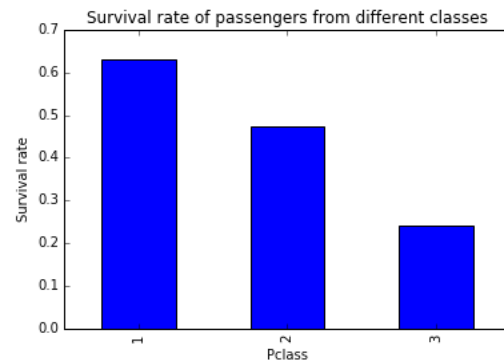
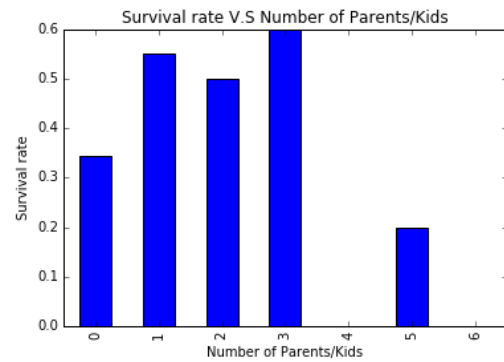
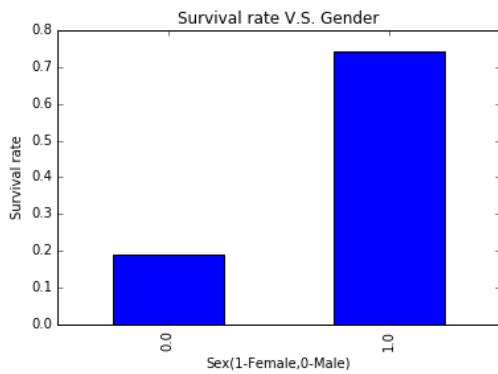
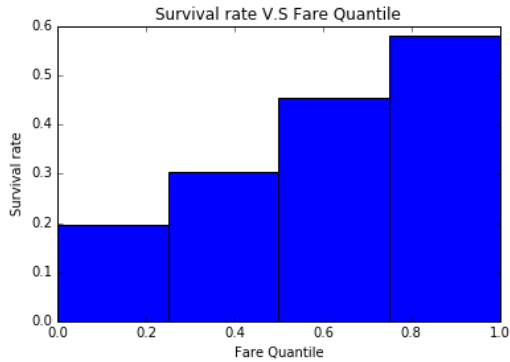
- j) Train your classifier with all of the training data, and test your classifier with the test data. Submit your results to Kaggle.

Submitted with accuracy: 96.929%

III. PROBLEM 2: THE TITANIC DISASTER

- b) Using logistic regression, try to predict whether a passenger survived the disaster. You can choose the features (or combinations of features) you would like to use or ignore, provided you justify your reasoning.

Here are some plots that help us determine the features to choose for classification.



Cross validation accuracy = 80.6%

See code in **Question2.ipynb**

- c) Train your classifier using all of the training data, and test it using the testing data. Submit your results to Kaggle.

Submitted with accuracy: 75.1%

Applied Machine Learning

Homework 1 - Written Exercises

Yan Jiang, Weiming Zhang

September 14, 2016

Exercise 1. Variance of a sum. Show that the variance of a sum is $\text{var}[X + Y] = \text{var}[X] + \text{var}[Y] + 2\text{cov}[X, Y]$, where $\text{cov}[X, Y]$ is the covariance between random variables X and Y .

Proof.

$$\begin{aligned}\text{var}[X + Y] &= E((X + Y - E(X + Y))^2) \\ &= E(X + Y)^2 - (E(X + Y))^2 \\ &= E(X^2 + 2XY + Y^2) - (E(X + Y))^2 \\ \text{var}[X] &= E(X^2) - (E(X))^2 \\ \text{var}[Y] &= E(Y^2) - (E(Y))^2 \\ \text{cov}[X + Y] &= E((X - E(X))(Y - E(Y))) \\ &= E(XY - E(X)Y - XE(Y) + E(X)E(Y)) \\ \text{var}[X] + \text{var}[Y] + 2\text{cov}[X + Y] &= E(X^2) - (E(X))^2 + E(Y^2) - (E(Y))^2 \\ &\quad + 2E(XY) - 2E(X)E(Y) - 2E(X)E(Y) + 2E(X)E(Y) \\ &= E(X^2) + 2E(XY) + E(Y^2) - ((E(X))^2 + (E(Y))^2 + 2E(X)E(Y)) \\ &= E(X^2 + 2XY + Y^2) - (E(X) + E(Y))^2 \\ &= E(X^2 + 2XY + Y^2) - (E(X + Y))^2 \\ &= \text{var}[X + Y]\end{aligned}$$

□

Exercise 2. Bayes rule for medical diagnosis (Source: Koller) After your yearly checkup, the doctor has bad news and good news. The bad news is that you tested positive for a serious disease, and that the test is 99% accurate (i.e., the probability of testing positive given that you have the disease is 0.99, as is the probability of testing negative given that you do not have the disease). The good news is that this is a rare disease, striking only one in 10,000 people. What are the chances that you actually have the disease? (Show your calculations as well as giving the final result.)

Solution:

Let θ indicate the case that a person do have this disease, and let x represent the case that a person is tested positive for this disease.

Here we have:

$$\begin{aligned}P(\theta) &= 0.0001 \\P(x | \theta) &= 0.99\end{aligned}$$

Thus,

$$\begin{aligned}P(\theta | x) &= \frac{P(\theta) P(x | \theta)}{P(x)} \\&= \frac{P(\theta) P(x | \theta)}{P(x | \theta) P(\theta) + P(x | \bar{\theta}) P(\bar{\theta})} \\&= \frac{0.0001 \times 0.99}{0.0001 \times 0.99 + 0.01 \times 0.9999} \\&= 0.98\end{aligned}$$

Exercise 3. Gradient and Hessian of log-likelihood for logistic regression.

(a) Let $\sigma(a) = \frac{1}{1+e^{-a}}$ be the sigmoid function. Show that $\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$

Solution:

$$\begin{aligned}\frac{d\sigma(a)}{da} &= (-1) \cdot \frac{1}{(1 + e^{-a})^2} \cdot e^{-a} \cdot (-1) \\&= \frac{e^{-a}}{(1 + e^{-a})^2} \\ \sigma(a)(1 - \sigma(a)) &= \frac{1}{1 + e^{-a}} \left(1 - \frac{1}{1 + e^{-a}}\right) \\&= \frac{1}{1 + e^{-a}} \cdot \frac{e^{-a}}{1 + e^{-a}} \\&= \frac{e^{-a}}{(1 + e^{-a})^2}\end{aligned}$$

Thus,

$$\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$$

- (b) Using the previous result and the chain rule of calculus, derive the expression for the gradient of the log likelihood given in HTF Eqn. 4.21.

Solution:

The likelihood of N observation is

$$L(\beta) = \prod_i^N p(x_i; \beta)^{y_i} (1 - p(x_i; \beta))^{1-y_i}$$

The log likelihood is

$$l(\beta) = \sum_{i=1}^N [y_i \log(x_i; \beta) + (1 - y_i) \log(1 - p(x_i; \beta))]$$

$$\frac{dl(\beta)}{d\beta} = \sum_{i=1}^N \left[y_i \frac{1}{p(x_i; \beta)} \cdot \frac{dp(x_i; \beta)}{d\beta} + (1 - y_i) \cdot \frac{1}{1 - p(x_i; \beta)} \cdot (-1) \cdot \frac{dp(x_i; \beta)}{d\beta} \right]$$

Since we have:

$$\frac{dp(x_i; \beta)}{d\beta} = p(x_i; \beta) (1 - p(x_i; \beta)) x_i$$

Thus,

$$\begin{aligned} \frac{dl(\beta)}{d\beta} &= \sum_{i=1}^N \left[y_i \frac{1}{p(x_i; \beta)} p(x_i; \beta) (1 - p(x_i; \beta)) x_i + (1 - y_i) \frac{1}{1 - p(x_i; \beta)} (-1) p(x_i; \beta) (1 - p(x_i; \beta)) x_i \right] \\ &= \sum_{i=1}^N [y_i (1 - p(x_i; \beta)) x_i + (1 - y_i) (-1) p(x_i; \beta) x_i] \\ &= \sum_{i=1}^N x_i [y_i - p(x_i; \beta)] \end{aligned}$$

- (c) As noted in HTF Eqn. 4.25, the Hessian matrix for the log likelihood can be written (up to a sign) as $\mathbf{X}^T \mathbf{W} \mathbf{X}$. Prove that this matrix is positive definite.

Proof. Since \mathbf{W} is a $\mathbf{N} \times \mathbf{N}$ diagonal matrix of weights with i th diagonal element $p(x_i; \beta) (1 - p(x_i; \beta))$, for any non zero vector

$$\mathbf{X} = [x_1, x_2, \dots, x_i, \dots]^T$$

,

$$\mathbf{X}^T \mathbf{W} \mathbf{X} = \sum_{i=1}^N [x_i^2 p(x_i; \beta) (1 - p(x_i; \beta))] > 0$$

Thus, $\mathbf{X}^T \mathbf{W} \mathbf{X}$ is positive definite.

□