

Homework 2: Bayesian Methods and Multiclass Classification

Introduction

This homework is about Bayesian methods and multiclass classification. In lecture we have primarily focused on binary classifiers trained to discriminate between two classes. In multiclass classification, we discriminate between three or more classes. We encourage you to first read the Bishop textbook coverage of these topic, particularly: Section 4.2 (Probabilistic Generative Models), Section 4.3 (Probabilistic Discriminative Models).

As usual, we imagine that we have the input matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$ (or perhaps they have been mapped to some basis Φ , without loss of generality) but our outputs are now “one-hot coded”. What that means is that, if there are c output classes, then rather than representing the output label y as an integer $1, 2, \dots, c$, we represent \mathbf{y} as a binary vector of length c . These vectors are zero in each component except for the one corresponding to the correct label, and that entry has a one. So, if there are 7 classes and a particular datum has label 3, then the target vector would be $C_3 = [0, 0, 1, 0, 0, 0, 0]$. If there are c classes, the set of possible outputs is $\{C_1 \dots C_c\} = \{C_k\}_{k=1}^c$. Throughout the assignment we will assume that output $\mathbf{y} \in \{C_k\}_{k=1}^c$.

The problem set has four problems:

- In the first problem, you will explore the properties of Bayesian estimation methods for the Bernoulli model as well as the special case of Bayesian linear regression with a simple prior.
- In the second problem, you will explore the properties of the softmax function, which is central to the method of multiclass logistic regression.
- In the third problem, you will dive into matrix algebra and the methods behind generative multiclass classifications. You will extend the discrete classifiers that we see in lecture to a Gaussian model.
- Finally, in the fourth problem, you will implement logistic regression as well as a generative classifier from close to scratch.

Problem 1 (Bayesian Methods, 10 pts)

This question helps to build your understanding of the maximum-likelihood estimation (MLE) vs. maximum a posterior estimator (MAP) and posterior predictive estimator, first in the Beta-Bernoulli model and then in the linear regression setting.

First consider the Beta-Bernoulli model (and see lecture 5.)

1. Write down the expressions for the MLE, MAP and posterior predictive distributions, and for a prior $\theta \sim \text{Beta}(4, 2)$ on the parameter of the Bernoulli, and with data $D = 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 1, 0$, plot the three different estimates after each additional sample.
2. Plot the posterior distribution (prior for 0 examples) on θ after 0, 4, 8, 12 and 16 examples. (Using whatever tools you like.)
3. Interpret the differences you see between the three different estimators.

Second, consider the Bayesian Linear Regression model, with data $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, $\mathbf{x}_i \in \mathbb{R}^m$, $y_i \in \mathbb{R}$, and generative model

$$y_i \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_i, \beta^{-1})$$

for (known) precision β (which is just the reciprocal of the variance). Given this, the likelihood of the data is $p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I})$. Consider the special case of an isotropic (spherical) prior on weights, with

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$$

4. Justify when you might use this prior in practice.
5. Using the method in lecture of taking logs, expanding and pushing terms that don't depend on \mathbf{w} into a constant, and finally collecting terms and completing the square, confirm that the posterior on weights after data D is $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\mathbf{m}_n, \mathbf{S}_n)$, where

$$\mathbf{S}_n = (\alpha\mathbf{I} + \beta\mathbf{X}^\top\mathbf{X})^{-1}$$

$$\mathbf{m}_n = \beta\mathbf{S}_n\mathbf{X}^\top\mathbf{y}$$

6. Derive the special case of the MAP estimator for this problem as the isotropic prior becomes arbitrarily weak. What does the MAP estimator reduce to?
7. What did we observe in lecture about this estimator for the case where the prior is neither weak nor strong?

Solution

1)

Bernoulli PMF:

$$p(x|\theta) = \prod_{i=1}^n \theta^x (1-\theta)^{1-x}$$

$\ln(p(x|\theta)) = n_1 \ln \theta + n_0 \ln(1-\theta)$ Setting the derivative (wrt θ) to 0, you get the expression for MLE (doesn't

make use of prior): $\frac{n_1}{n_0 + n_1}$

$p(\theta|D) = \text{Beta}(\theta|n_1 + \alpha, n_0 + \beta)$. Therefore, expression for MAP:

$$\theta_{\text{MAP}} = \frac{\alpha + n_1}{\alpha + \beta + n_1 + n_0 + 2}$$

Expression for Posterior Predictive Distribution:

$$p(x=1|D) = \frac{\alpha + n_1}{\alpha + \beta + n_1 + n_0 + 2}$$

See figure 1 for plot.

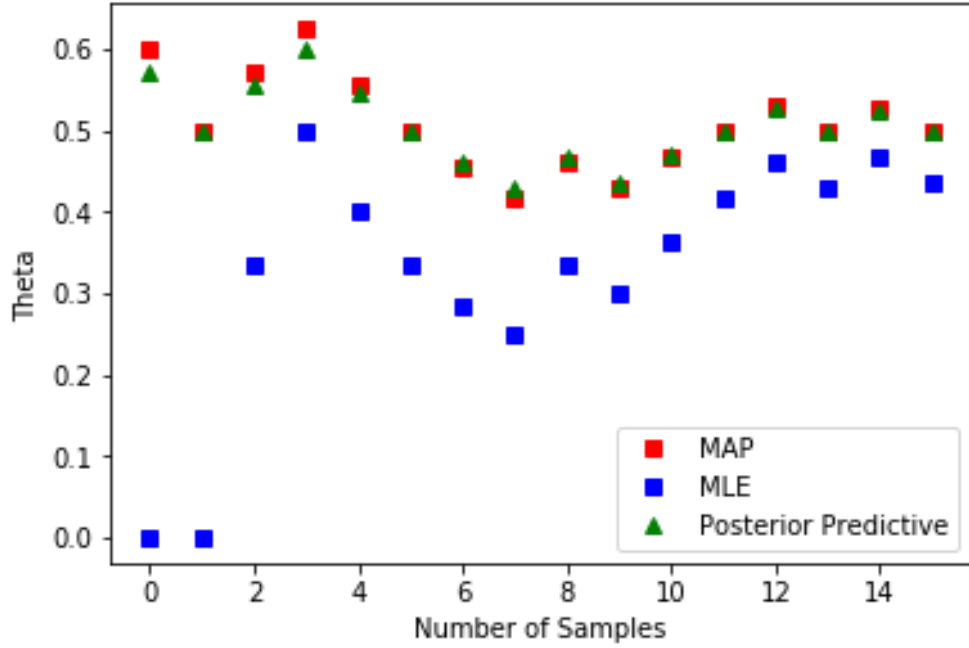


Figure 1: Plot of MLE, MAP, and posterior predictive distributions for Bernoulli with a Beta prior after various amounts of data.

2)

See Figure 2.

3)

The MLE is purely based on the data, so it starts off with bad estimates until it gets enough data. Since there are only 16 samples, the original beta prior retains an effect on the final estimate, causing the MAP and posterior predictive estimates to be larger than the MLE at all points. The MAP estimate starts off higher than the posterior predictive, because it subtracts the equivalent of 1/2 from the estimate, but after a few data points, the MAP and posterior predictive end up being essentially the same.

4)

One might use an isotropic (spherical) prior if it is expected that the weights should be centered around 0 (example, random noise).

5)

$$p(\mathbf{w}|D) \propto N(\mathbf{w}|0, \alpha^{-1}\mathbf{I})N(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}) = \frac{1}{|2\pi\alpha^{-1}\mathbf{I}|} \exp\left(-\frac{1}{2}\mathbf{w}^T\alpha^{-1}\mathbf{I}\mathbf{w}\right) \frac{1}{|2\pi\beta^{-1}\mathbf{I}|} \exp\left(-\frac{1}{2}(\mathbf{y}-\mathbf{X}\mathbf{w})^T\beta^{-1}\mathbf{I}(\mathbf{y}-\mathbf{X}\mathbf{w})\right)$$

$$\ln p(\mathbf{w}|D) = \text{Constant} - \frac{1}{2}[\mathbf{w}^T\alpha^{-1}\mathbf{I}\mathbf{w} + (\mathbf{y}-\mathbf{X}\mathbf{w})^T\beta^{-1}\mathbf{I}(\mathbf{y}-\mathbf{X}\mathbf{w})]$$

$$\ln p(\mathbf{w}|D) = \text{Constant} - \frac{1}{2}(\mathbf{w}^T(\alpha\mathbf{I} + \beta\mathbf{X}^T\mathbf{X})\mathbf{w} - 2\mathbf{w}^T(\beta\mathbf{X}^T\mathbf{y}))$$

Complete the square:

$$\ln p(\mathbf{w}|D) = \text{Constant} - \frac{1}{2}a(\mathbf{w} - b/a)^2, \text{ where } a = (\alpha\mathbf{I} + \beta\mathbf{X}^T\mathbf{X}), \text{ and } b = (\beta\mathbf{X}^T\mathbf{y})$$

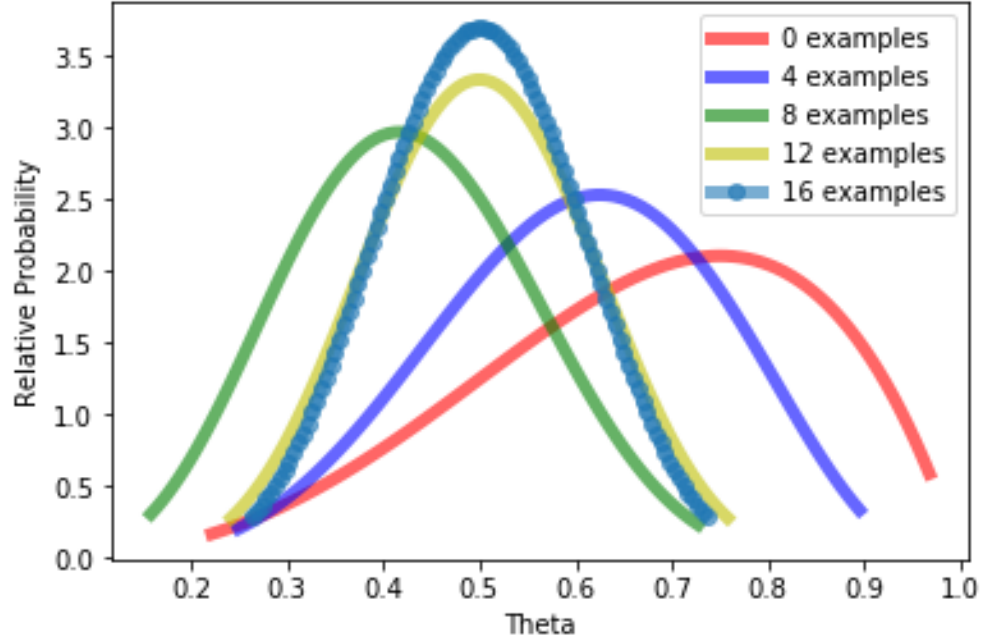


Figure 2: Plot of Beta posterior after various amounts of data.

Then $\mathbf{S}_n = a = (\alpha \mathbf{I} + \beta \mathbf{X}^T \mathbf{X})^{-1}$, and $\mathbf{m}_n = b/a = (\beta \mathbf{S}_n \mathbf{X}^T \mathbf{y})$.

This matches the normal distribution with mean of \mathbf{m}_n and variance of \mathbf{S}_n .

6)

If the prior becomes very weak, that means α is very large, so the MAP estimator, which is $\mathbf{m}_n = \frac{(\beta \mathbf{X}^T \mathbf{y})}{(\alpha \mathbf{I} + \beta \mathbf{X}^T \mathbf{X})}$, reduces to the MLE, which is $\frac{\mathbf{y}}{\mathbf{X}}$.

7)

If the prior is neither weak nor strong, it becomes the regularization constant of LASSO.

Problem 2 (Properties of Softmax, 8pts)

We have explored logistic regression, which is a discriminative probabilistic model over two classes. For each input \mathbf{x} , logistic regression outputs a probability of the class output y using the logistic sigmoid function.

The softmax transformation is an important generalization of the logistic sigmoid to the case of c classes. It takes as input a vector, and outputs a transformed vector of the same size,

$$\text{softmax}(\mathbf{z})_k = \frac{\exp(z_k)}{\sum_{\ell=1}^c \exp(z_\ell)}, \quad \text{for all } k$$

Multiclass logistic regression uses the softmax transformation over vectors of size c . Let $\{\mathbf{w}_\ell\} = \{\mathbf{w}_1 \dots \mathbf{w}_c\}$ denote the parameter vectors for each class. In particular, multiclass logistic regression defines the probability of class k as,

$$p(\mathbf{y} = C_k | \mathbf{x}; \{\mathbf{w}_\ell\}) = \text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_k = \frac{\exp(\mathbf{w}_k^\top \mathbf{x})}{\sum_{\ell=1}^c \exp(\mathbf{w}_\ell^\top \mathbf{x})}.$$

As above, we are using $\mathbf{y} = C_k$ to indicate the output vector that represents class k .

Assuming data $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, the negated log-likelihood can be written in the standard form, as

$$\mathcal{L}(\{\mathbf{w}_\ell\}) = - \sum_{i=1}^n \ln p(\mathbf{y}_i | \mathbf{x}_i; \{\mathbf{w}_\ell\})$$

Softmax is an important function in the context of machine learning, and you will see it again in other models, such as neural networks. In this problem, we aim to gain intuitions into the properties of softmax and multiclass logistic regression.

Show that:

1. The output of the softmax function is a vector with non-negative components that are at most 1.
2. The output of the softmax function defines a distribution, so that in addition, the components sum to 1.
3. Softmax preserves order. This means that if elements $z_k < z_\ell$, in \mathbf{z} , then $\text{softmax}(\mathbf{z})_k < \text{softmax}(\mathbf{z})_\ell$ for any k, ℓ .
4. Show that

$$\frac{\partial \text{softmax}(\mathbf{z})_k}{\partial z_j} = \text{softmax}(\mathbf{z})_k (I_{kj} - \text{softmax}(\mathbf{z})_j) \quad \text{for any } k, j$$

, where indicator $I_{kj} = 1$ if $k = j$ and $I_{kj} = 0$ otherwise.

5. Using your answer to the previous question, show that

$$\frac{\partial}{\partial \mathbf{w}_k} \mathcal{L}(\{\mathbf{w}_\ell\}) = \sum_{i=1}^n (p(\mathbf{y}_i = C_k | \mathbf{x}_i; \{\mathbf{w}_\ell\}) - y_{ik}) \mathbf{x}_i$$

By the way, this may be useful for Problem 3!

Solution

1)

Given the softmax is purely using exponential functions, all of these functions give non-negative outputs, so the softmax function output must be non-negative. Since, the numerator is a subset of the denominator, the output can only be at most 1.

2)

If you sum up the numerator of the softmax function, it will equal the denominator, which is a sum of all of the elements of the vector \mathbf{w} . Thus, the components sum to 1. This is a distribution, because it gives values corresponding to each datapoint.

3)

The denominator for each of the softmax functions for a particular vector of values will be the same, and if $z_k < z_l$, then the numerator term of $\text{softmax}(\mathbf{z})_k$ will be less than the numerator of $\text{softmax}(\mathbf{z})_l$, so $\text{softmax}(\mathbf{z})_k$ will be less than $\text{softmax}(\mathbf{z})_l$ for any k, l .

4)

$$\frac{\partial \exp(\mathbf{w}_k^T \mathbf{x})}{\partial z_j \sum_{l=1}^c \exp(\mathbf{w}_l^T \mathbf{x})} =$$

If $k = j$:

$$\begin{aligned} & \frac{(\mathbf{w}_k^T \mathbf{x}) * \exp(\mathbf{w}_k^T \mathbf{x}) * (\sum_{l=1}^c \exp(\mathbf{w}_l^T \mathbf{x})) - \exp(\mathbf{w}_k^T \mathbf{x}) * ((\mathbf{w}_k^T \mathbf{x}) * \exp(\mathbf{w}_k^T \mathbf{x}))}{(\sum_{l=1}^c \exp(\mathbf{w}_l^T \mathbf{x}))^2} \\ &= \text{softmax}(\mathbf{z})_k \left(\frac{\sum_{l=1}^c \exp(\mathbf{w}_l^T \mathbf{x}) - \exp(\mathbf{w}_l^T \mathbf{x})}{\sum_{l=1}^c \exp(\mathbf{w}_l^T \mathbf{x})} \right) = \text{softmax}(\mathbf{z})_k (I_{kj} - \text{softmax}(\mathbf{z})_k) \end{aligned}$$

If $k \neq j$:

$$\frac{-\exp(\mathbf{w}_k^T \mathbf{x}) * (\exp(\mathbf{w}_j^T \mathbf{x}))}{(\sum_{l=1}^c \exp(\mathbf{w}_l^T \mathbf{x}))^2} = \text{softmax}(\mathbf{z})_k * (-\text{softmax}(\mathbf{z})_j)$$

(noting that the components sum to 1).

5)

$$\frac{\partial}{\partial \mathbf{w}_k} L(\mathbf{w}_l) = - \sum_{i=1}^n \frac{p(y|x_i; \{\mathbf{w}_l\})(I_{kj} - p(y|x_i; \{\mathbf{w}_l\}))}{p(y|x_i; \{\mathbf{w}_l\})}$$

The indicator variable becomes a one-hot matrix over all classes and thus becomes $y_{ik} * \mathbf{x}_i$. Thus, the expression becomes:

$$= \sum_{i=1}^n (p(y|x_i; \{\mathbf{w}_l\}) - I_{kj}) = \sum_{i=1}^n (p(y|x_i; \{\mathbf{w}_l\}) - y_{ik}) \mathbf{x}_i$$

Problem 3 (Return of matrix calculus, 10pts)

Consider now a generative c -class model. We adopt class prior $p(\mathbf{y} = C_k; \boldsymbol{\pi}) = \pi_k$ for all $k \in \{1, \dots, c\}$ (where π_k is a parameter of the prior). Let $p(\mathbf{x}|\mathbf{y} = C_k)$ denote the class-conditional density of features \mathbf{x} (in this case for class C_k). Consider the data set $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ where as above $\mathbf{y}_i \in \{C_k\}_{k=1}^c$ is encoded as a one-hot target vector.

1. Write out the negated log-likelihood of the data set, $-\ln p(D; \boldsymbol{\pi})$.
2. Since the prior forms a distribution, it has the constraint that $\sum_k \pi_k - 1 = 0$. Using the hint on Lagrange multipliers below, give the expression for the maximum-likelihood estimator for the prior class-membership probabilities, i.e. $\hat{\pi}_k$. Make sure to write out the intermediary equation you need to solve to obtain this estimator. Double-check your answer: the final result should be very intuitive!

For the remaining questions, let the class-conditional probabilities be Gaussian distributions with the same covariance matrix

$$p(\mathbf{x}|\mathbf{y} = C_k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}), \text{ for } k \in \{1, \dots, c\}$$

and different means $\boldsymbol{\mu}_k$ for each class.

3. Derive the gradient of the negative log-likelihood with respect to vector $\boldsymbol{\mu}_k$. Write the expression in matrix form as a function of the variables defined throughout this exercise. Simplify as much as possible for full credit.
4. Derive the maximum-likelihood estimator for vector $\boldsymbol{\mu}_k$. Once again, your final answer should seem intuitive.
5. Derive the gradient for the negative log-likelihood with respect to the covariance matrix $\boldsymbol{\Sigma}$ (i.e., looking to find an MLE for the covariance). Since you are differentiating with respect to a *matrix*, the resulting expression should be a matrix!
6. Derive the maximum likelihood estimator of the covariance matrix.

[Hint: Lagrange Multipliers.] Lagrange Multipliers are a method for optimizing a function f with respect to an equality constraint, i.e.

$$\min_{\mathbf{x}} f(\mathbf{x}) \text{ s.t. } g(\mathbf{x}) = 0.$$

This can be turned into an unconstrained problem by introducing a Lagrange multiplier λ and constructing the Lagrangian function,

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x}).$$

It can be shown that it is a necessary condition that the optimum is a critical point of this new function. We can find this point by solving two equations:

$$\frac{\partial L(\mathbf{x}, \lambda)}{\partial \mathbf{x}} = 0 \quad \text{and} \quad \frac{\partial L(\mathbf{x}, \lambda)}{\partial \lambda} = 0$$

Cookbook formulas. Here are some formulas you might want to consider using to compute difficult gradients. You can use them in the homework without proof. If you are looking to hone your matrix calculus skills, try to find different ways to prove these formulas yourself (will not be part of the evaluation of this homework). In general, you can use any formula from the matrix cookbook, as long as you cite it. We opt for the following common notation: $\mathbf{X}^{-\top} := (\mathbf{X}^{\top})^{-1}$

$$\frac{\partial \mathbf{a}^{\top} \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-\top} \mathbf{a} \mathbf{b}^{\top} \mathbf{X}^{-\top}$$

$$\frac{\partial \ln |\det(\mathbf{X})|}{\partial \mathbf{X}} = \mathbf{X}^{-\top}$$

Solution

1)

This can be modeled in part as a multinomial distribution, so: $p(D; \boldsymbol{\pi}) = \prod_{i=1}^n \prod_{k=1}^c (\pi_k p(\mathbf{x}|\mathbf{y} = C_k))^{y_k}$

$$-\ln p(D; \boldsymbol{\pi}) = \sum_{i=1}^n \sum_{k=1}^c -(\ln(\pi_k) + \ln p(\mathbf{x}|\mathbf{y} = C_k)) * (y_k)$$

2)

From above, we know that the negative log likelihood of the prior is:

$$-\sum_{i=1}^n \sum_{k=1}^c (y_{ik} \ln \pi_k)$$

Optimizing the function with respect to the constraint gives:

$$L(\pi, \lambda) = -\sum_{i=1}^n \sum_{k=1}^c (y_{ik} \ln \pi_k) + \lambda(1 - \sum_{k=1}^c (\pi_k))$$

To attain the optimum of this function, we must solve the two equations setting the derivatives with respect to π_k and λ to 0:

$$\frac{\partial L(\pi, \lambda)}{\partial \pi_k} = 0 = -\sum_{i=1}^n \frac{y_{ik}}{\pi_k} - \lambda$$

$$\frac{\partial L(\pi, \lambda)}{\partial \lambda} = 0 = 1 - \sum_{k=1}^c (\pi_k)$$

Using these 2 equations, we see that:

$$\lambda = n \text{ and } \hat{\pi} = \frac{\sum_{i=1}^n (y_{ik})}{\sum_{i=1}^n \sum_{k=1}^c (y_{ik})}$$

This essentially means that the prior for each class is the total number of objects in that class divided by the total number of objects in all classes.

3)

$$-\ln p(\mathbf{x}|\mathbf{y} = C_k) = \text{Constant} + \frac{1}{2} \sum_{k=1}^c (\mathbf{x} - \mu_k)^T \Sigma^{-1} (\mathbf{x} - \mu_k)$$

$$= \frac{1}{2} \sum_{k=1}^c (\mu_k^T \Sigma^{-1} \mu_k + \mathbf{x}^T \Sigma^{-1} \mathbf{x} + 2\mathbf{x}^T \Sigma^{-1} \mu_k)$$

$$-\frac{\partial \ln p(\mathbf{x}|\mathbf{y})}{\partial \mu_k} = \sum_{k=1}^c (\mu_k^T \Sigma^{-1} + \mathbf{x}^T \Sigma^{-1})$$

4)

$-\frac{\partial \ln p(\mathbf{x}|\mathbf{y})}{\partial \mu_k} = 0 \Rightarrow \hat{\mu}_k = \frac{\sum_{i=1}^n \mathbf{x}_{ik}}{n_k}$ This means that the MLE for μ_k is the sum of the observations for a given class divided by the total number of observations in that class.

5/6)

$$-\ln p(\mathbf{x}|\mathbf{y} = C_k) = \text{Constant} + \frac{N}{2} \ln |\Sigma| + \frac{1}{2} \sum_{k=1}^c (\mathbf{x} - \mu_k)^T \Sigma^{-1} (\mathbf{x} - \mu_k)$$

$$-\ln p(\mathbf{x}|\mathbf{y} = C_k) = \text{Constant} + \frac{N}{2} \ln |\Sigma| + \frac{1}{2} \text{tr}(\Sigma^{-1} \sum_{k=1}^c (\mathbf{x} - \mu_k)^T (\mathbf{x} - \mu_k))$$

From matrix cookbook:

$$-\frac{\partial \ln p(\mathbf{x}|\mathbf{y})}{\partial \Sigma} = \frac{N}{2} \Sigma^{-T} + \frac{1}{2} \Sigma^T \sum_{k=1}^c (\mathbf{x} - \mu_k)^T (\mathbf{x} - \mu_k)$$

$$-\frac{\partial \ln p(\mathbf{x}|\mathbf{y})}{\partial \Sigma} = 0 \Rightarrow \hat{\Sigma} = \frac{1}{N} \sum_{k=1}^c (\mathbf{x} - \mu_k)^T (\mathbf{x} - \mu_k)$$

4. Classifying Fruit [15pts]

You're tasked with classifying three different kinds of fruit, based on their heights and widths. Figure 3 is a plot of the data. Iain Murray collected these data and you can read more about this on his website at http://homepages.inf.ed.ac.uk/imurray2/teaching/oranges_and_lemons/. We have made a slightly simplified (collapsing the subcategories together) version of this available as `fruit.csv`, which you will find in the Github repository. The file has three columns: type (1=apple, 2=orange, 3=lemon), width, and height. The first few lines look like this:

```
fruit,width,height
1,8.4,7.3
1,8,6.8
1,7.4,7.2
1,7.1,7.8
...
```

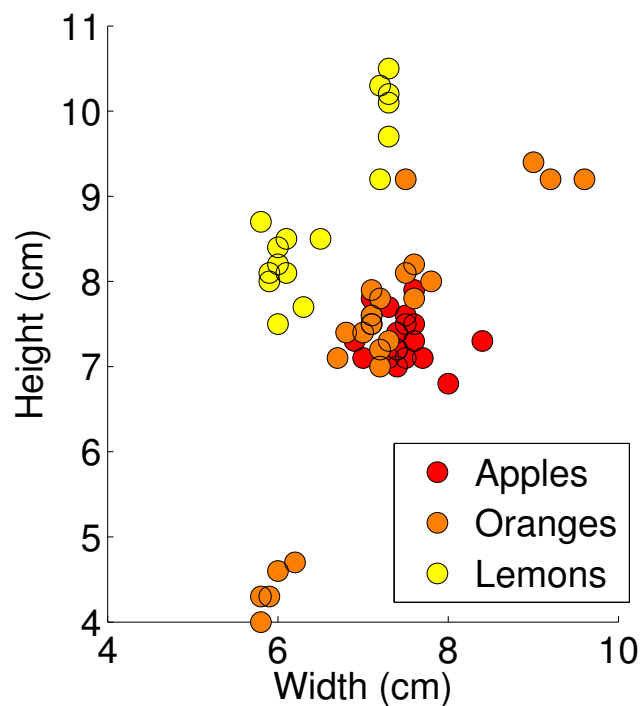


Figure 3: Heights and widths of apples, oranges, and lemons. These fruit were purchased and measured by Iain Murray: http://homepages.inf.ed.ac.uk/imurray2/teaching/oranges_and_lemons/.

Problem 4 (Classifying Fruit, 15pts)

You should implement the following:

- The three-class generalization of logistic regression, also known as softmax regression, for these data. You will do this by implementing gradient descent on the negative log likelihood. You will need to find good values for the learning rate η and regularization strength λ .
- A generative classifier with Gaussian class-conditional densities, as in Problem 3. In particular, make two implementations of this, one with a shared covariance matrix across all of the classes, and one with a separate covariance being learned for each class. Note that the staff implementation can switch between these two by the addition of just a few lines of code. In the separate covariance matrix case, the MLE for the covariance matrix of each class is simply the covariance of the data points assigned to that class, without combining them as in the shared case.

You may use anything in `numpy` or `scipy`, except for `scipy.optimize`. That being said, if you happen to find a function in `numpy` or `scipy` that seems like it is doing too much for you, run it by a staff member on Piazza. In general, linear algebra and random variable functions are fine. The controller file is `problem4.py`, in which you will specify hyperparameters. The actual implementations you will write will be in `LogisticRegression.py` and `GaussianGenerativeModel.py`.

You will be given class interfaces for `GaussianGenerativeModel` and `LogisticRegression` in the distribution code, and the code will indicate certain lines that you should not change in your final submission. Naturally, don't change these. These classes will allow the final submissions to have consistency. There will also be a few hyperparameters that are set to irrelevant values at the moment. You may need to modify these to get your methods to work. The classes you implement follow the same pattern as scikit-learn, so they should be familiar to you. The distribution code currently outputs nonsense predictions just to show what the high-level interface should be, so you should completely remove the given `predict()` implementations and replace them with your implementations.

- The `visualize()` method for each classifier will save a plot that will show the decision boundaries. You should include these in this assignment.
- Which classifiers model the distributions well?
- What explains the differences?

In addition to comparing the decision boundaries of the three models visually:

- For logistic regression, report negative log-likelihood loss for several configurations of hyperparameters. Why are your final choices of learning rate (η) and regularization strength (λ) reasonable? Plot loss during training for the best of these configurations, with iterations on the x-axis and loss on the y-axis (one way to do this is to add a method to the `LogisticRegression` Class that displays loss).
- For both Gaussian generative models, report likelihood. In the separate covariance matrix case, be sure to use the covariance matrix that matches the true class of each data point.

Solution

1)

The Gaussian generative model with separate co-variance matrices models the distributions the best. This is because it can form non-linear boundaries. The model with separate co-variance matrices is better than the one that assumes the same co-variance matrices, because the actual underlying distributions have different co-variances.

In Table 1 I report different negative log-likelihood losses for different hyperparameters. The final choices

Scores		
Eta	Lambda	Negative Log-Likelihood Score
0.1	0.1	4744.9
0.001	0.1	34.8
0.0001	0.1	25.8
0.0001	10	33.9
0.0001	1	22.3
0.0001	0.1	25.8
0.0001	0.01	28.2

Table 1: Negative log-likelihood scores for different hyperparameters.

are reasonable, because at a certain point the eta value will not affect much, but it slows down the algorithm to have a lower value. The lambda value I chose (.1) had slightly higher loss than 1, but this was due to the outliers of the blue points. Otherwise, it got closer to the red points, and would have been better for future predictions.

For the Gaussian Generative Model with separate covariance matrices, the likelihood was -158.1, and for the model with shared covariance matrices, the likelihood was -185.3.

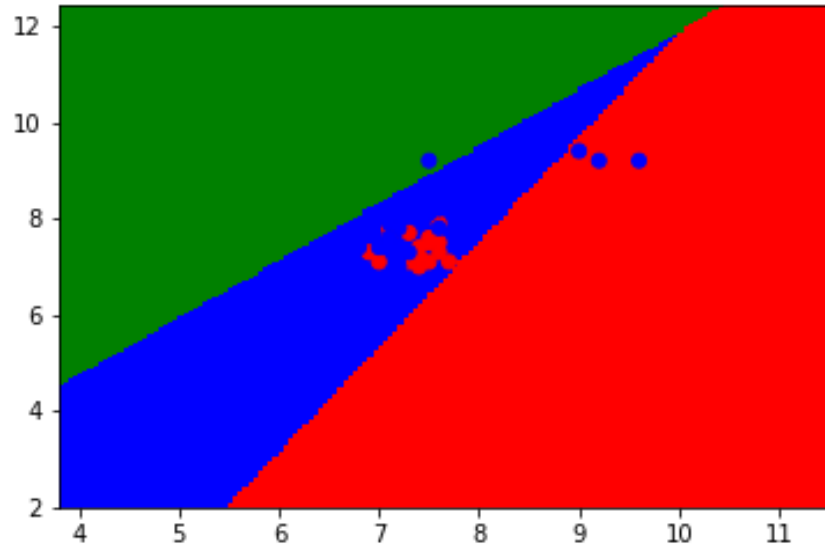


Figure 4: Plot of logistic regression classifier decision boundaries.

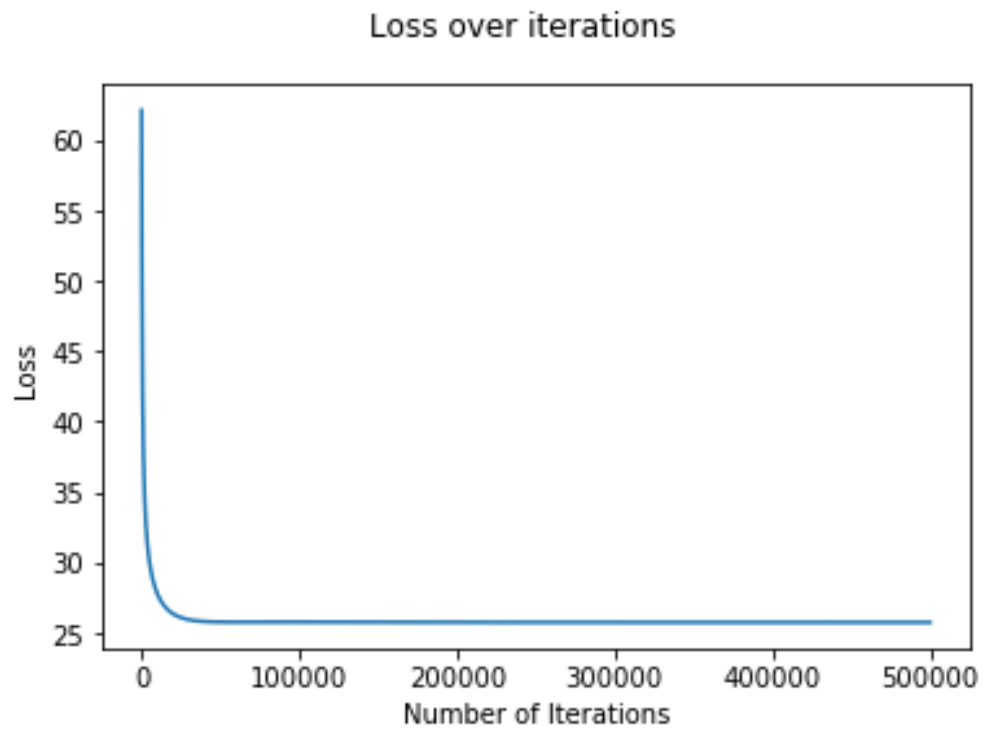


Figure 5: Plot of loss over many iterations in logistic regression classifier.

Calibration [1pt]

Approximately how long did this homework take you to complete? 20 hours

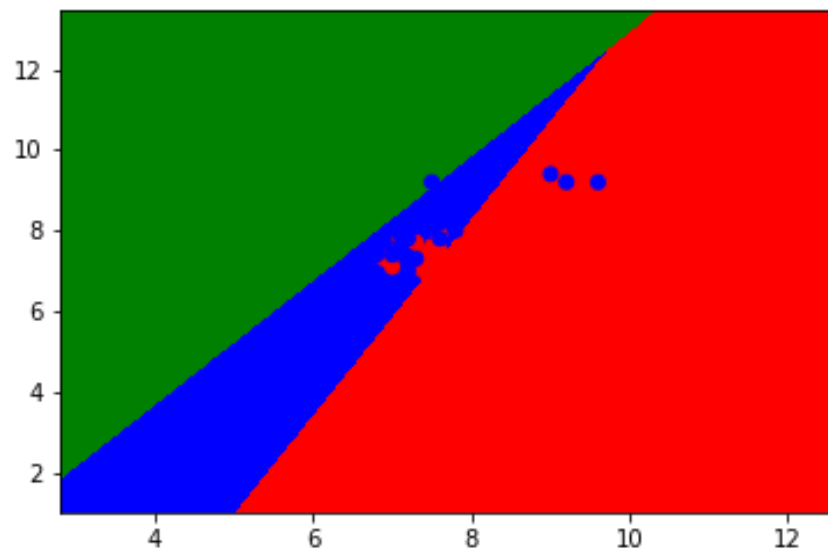


Figure 6: Plot of Gaussian Generative Model classifier decision boundaries with shared covariance matrix.

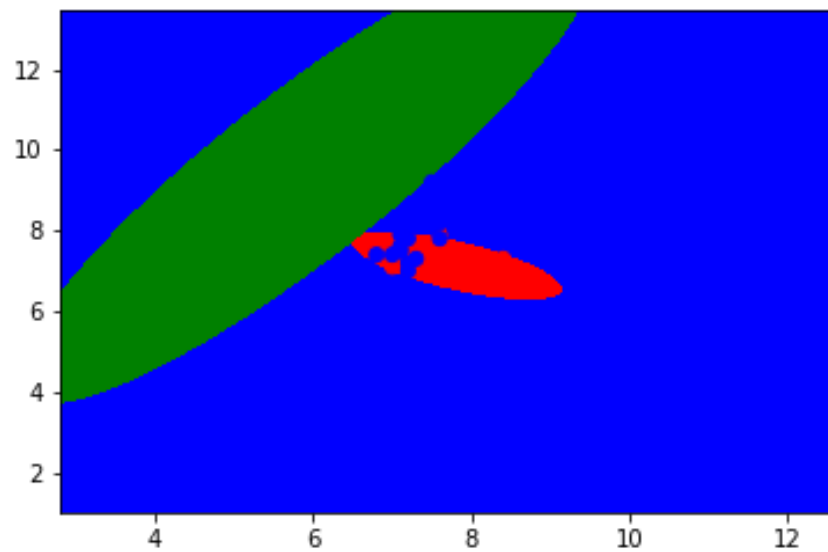


Figure 7: Plot of Gaussian Generative Model classifier decision boundaries with separate covariance matrix.