
SOLUTION - Do Not Distribute

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 topics, 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 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

Problem 1.1

The expressions for MLE, MAP, and the posterior predictive distribution are as follows:

$$p(\theta|D) = \theta_{\text{MLE}} = \frac{n_1}{n_0 + n_1}$$

$$p(\theta|D) = \text{Beta}(\theta|n_1 + \alpha, n_0 + \beta) = \theta_{\text{MAP}} = \frac{\alpha + n_1 - 1}{\alpha + \beta + n_1 + n_0 - 2}$$

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

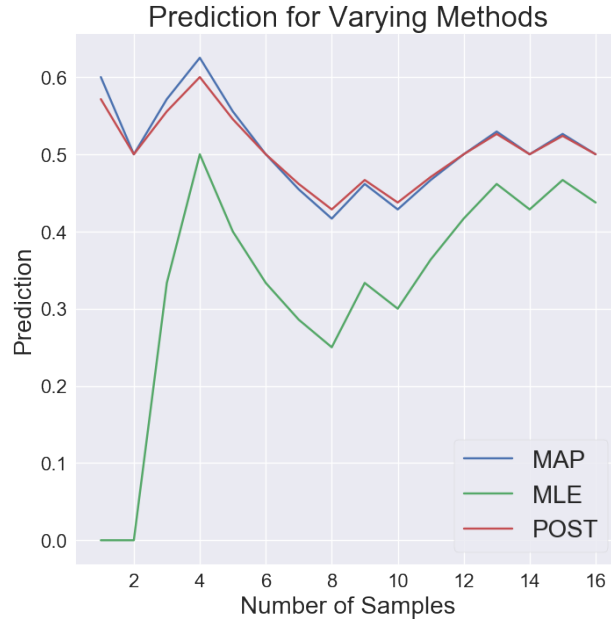


Figure 1: Plot for the three different estimates after each additional sample.

Problem 1.2

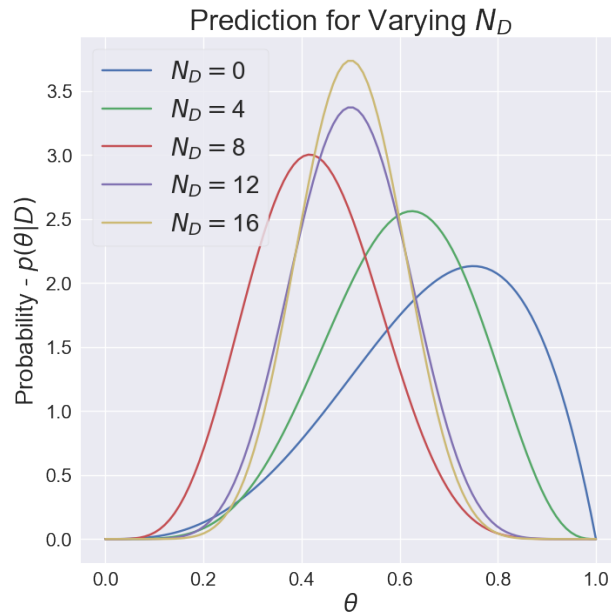


Figure 2: Posterior distribution with pripro on θ after 0, 4, 8, 12, and 16 examples.

Problem 1.3

As we can see for the MAP and posterior estimators they follow closely to each other. On the other hand, the MLE follows a much lower path compared to MAP and the posterior estimators. This comes from the fact that we are not using a prior. Another effect of this, is that we can see that at the beginning, it is at 0, because it does not consider any distribution effects, but rather only the beginning data, in which it was all 0's. As we increase the number of samples, we also see that the MLE estimator starts converging to the other two solutions.

Problem 1.4

This prior could be useful in practice when you do not have any information on the distribution of the sample beforehand. This would be a good initial prior as it is uniform and imitates a unit, multivariate, and isotropic Gaussian distribution. Also, this does not introduce artificial bias to the model. Furthermore, it is important to note that this introduces a prior that has a mean of 0 and with independent features as the non-diagonal elements of the covariance are 0.

Problem 1.5

Given that we have the following:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I})$$

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

and Bayes rule gives us

$$p(\mathbf{w}|D) \propto p(\mathbf{w})p(\mathbf{y}|\mathbf{X}, \mathbf{w})$$

such that

$$\ln(p(\mathbf{w}|D)) \propto \underbrace{\ln(p(\mathbf{w}))}_{(1)} + \underbrace{\ln(p(\mathbf{y}|\mathbf{X}, \mathbf{w}))}_{(2)}.$$

$$(1) \Rightarrow \ln(\mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})) = c - \frac{1}{2}(\mathbf{w}^\top \alpha \mathbf{I} \mathbf{w})$$

$$(2) \Rightarrow \ln(\mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I})) = c - \frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^\top \beta \mathbf{I} (\mathbf{y} - \mathbf{X}\mathbf{w}) = c - \frac{\beta}{2}(\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\mathbf{w} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w})$$

constant

Putting (1) and (2) together we obtain

$$\ln(p(\mathbf{w}|D)) \propto c - \frac{\alpha \mathbf{I}}{2}(\mathbf{w}^\top \mathbf{w}) - \frac{\beta}{2}(-2\mathbf{y}^\top \mathbf{X}\mathbf{w} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w})$$

$$\ln(p(\mathbf{w}|D)) \propto c - \frac{1}{2}(\alpha \mathbf{I} \mathbf{w}^\top \mathbf{w} - 2\beta \mathbf{y}^\top \mathbf{X}\mathbf{w} + \beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w})$$

$$\ln(p(\mathbf{w}|D)) \propto c - \frac{1}{2}(\mathbf{w}^\top (\alpha \mathbf{I} + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{y}^\top \mathbf{X}\mathbf{w})$$

which using completing the square yields the form

$$\ln(\mathcal{N}(\mathbf{w}|\mathbf{m}_n, \mathbf{S}_n)) = c - \frac{1}{2} (\mathbf{w}^\top \mathbf{S}_n^{-1} \mathbf{w} - 2\mathbf{m}_n^\top \mathbf{S}_n^{-1} \mathbf{w} + \mathbf{m}_n^\top \mathbf{S}_n^{-1} \mathbf{m}_n)$$

where

$$\begin{aligned}\mathbf{S}_n &= (\alpha \mathbf{I} + \beta \mathbf{X}^\top \mathbf{X})^{-1} \\ \mathbf{m}_n &= \beta \mathbf{S}_n \mathbf{X}^\top \mathbf{y}.\end{aligned}$$

and ultimately simplified to

$$\ln(\mathcal{N}(\mathbf{w}|\mathbf{m}_n, \mathbf{S}_n)) = c - \frac{1}{2} (\mathbf{w} - \mathbf{m}_n)^\top \mathbf{S}_n^{-1} (\mathbf{w} - \mathbf{m}_n).$$

■

Problem 1.6

If the isotropic prior becomes arbitrarily weak, then $\alpha \rightarrow 0$, meaning that the MAP estimator becomes

$$\ln(\mathcal{N}(\mathbf{w}|\mathbf{m}_n, \mathbf{S}_n)) = c - \frac{1}{2} (\mathbf{w} - \mathbf{m}_n)^\top \mathbf{S}_n^{-1} (\mathbf{w} - \mathbf{m}_n)$$

where

$$\begin{aligned}\mathbf{S}_n &= (\beta \mathbf{X}^\top \mathbf{X})^{-1} \\ \mathbf{m}_n &= \beta \mathbf{S}_n \mathbf{X}^\top \mathbf{y}.\end{aligned}$$

which is simply the MLE distribution.

Problem 1.7

If the prior is neither weak nor strong, then we see that it behaves moderately as a regularization of the model and we recover essentially the ridge regression once again. This can be seen by bringing together the α and the β into one fractional parameter. This is really useful, as we can see the similarities between the functions of the different methods.

End Solution

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!

Problem 2.1

We know that $e^x > 1$ for $x \in \mathbb{R}$ and that

$$\sum_{l=1}^c \exp z_l \geq \exp z_k \Rightarrow \frac{\exp z_k}{\sum_{l=1}^c \exp z_l} \leq 1.$$

This means that the softmax function is a vector with non-negative components that are at most 1. ■

Problem 2.2

$$\sum_{k=1}^c \text{softmax}(\mathbf{z})_k = \frac{\sum_{k=1}^c \exp z_k}{\sum_{l=1}^c \exp z_l} = 1$$
■

Problem 2.3

If $z_k < z_l$, then

$$\exp(z_k) < \exp(z_l).$$

Also,

$$\text{softmax}(\mathbf{z})_k = \frac{\exp z_k}{\sum_{i=1}^c \exp z_i} \quad \text{and} \quad \text{softmax}(\mathbf{z})_l = \frac{\exp z_l}{\sum_{i=1}^c \exp z_i}$$

then,

$$\frac{\exp z_k}{\sum_{i=1}^c \exp z_i} < \frac{\exp z_l}{\sum_{i=1}^c \exp z_i}.$$

Thus, this ultimately means that

$$\text{softmax}(\mathbf{z})_k < \text{softmax}(\mathbf{z})_l$$
■

Problem 2.4

$$\begin{aligned}
 \frac{\partial \text{softmax}(\mathbf{z})_k}{\partial z_j} &= \frac{\partial}{\partial z_j} \left(\frac{\exp z_k}{\sum_{l=1}^c \exp z_l} \right) \\
 &= \frac{\partial}{\partial z_j} \left(\exp z_k \frac{1}{\sum_{l=1}^c \exp z_l} \right) \\
 &= \underbrace{\frac{\partial \exp z_k}{\partial z_j} \cdot \frac{1}{\sum_{l=1}^c \exp z_l}}_{(1)} + \underbrace{\exp z_k \cdot \frac{\partial}{\partial z_j} \left(\frac{1}{\sum_{l=1}^c \exp z_l} \right)}_{(2)}
 \end{aligned}$$

If we now take the individual derivatives, in which (2) we use the chain rule, we get

$$\begin{aligned}
 (1) &= \mathbf{I}_{kj} \frac{\exp z_k}{\sum_{l=1}^c \exp z_l} \\
 (2) &= \exp z_k \frac{\partial}{\partial u} \frac{1}{u} \quad \text{where} \quad u = \left(\sum_{l=1}^c \exp z_l \right) \quad \text{and} \quad \frac{\partial}{\partial u_j} = \exp z_j \\
 &= - \exp z_k \frac{1}{(\sum_{l=1}^c \exp z_l)^2} \exp z_j
 \end{aligned}$$

Putting (1) and (2) together, this simplifies to

$$\boxed{\text{softmax}(\mathbf{z})_k (\mathbf{I}_{kj} - \text{softmax}(\mathbf{z})_j)}$$

■

Problem 2.5

$$p(\mathbf{y}_i | \mathbf{x}_i, \{\mathbf{w}_\ell\}) = \prod_{k=1}^c p(\mathbf{y}_i = C_k | \mathbf{x}_i \{\mathbf{w}_\ell\})$$

$$\begin{aligned}
 \mathcal{L}(\{\mathbf{w}_\ell\}) &= - \ln \sum_{i=1}^N p(\mathbf{y}_i | \mathbf{x}_i, \{\mathbf{w}_\ell\}) \\
 &= - \sum_{i=1}^N \sum_{j=1}^c \ln \left(\prod_{k=1}^c p(\mathbf{y}_i = C_k | \mathbf{x}_i \{\mathbf{w}_\ell\}) \right)^{y_{ij}} \\
 &= - \sum_{i=1}^N \sum_{j=1}^c y_{ij} \ln(\text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_j)
 \end{aligned}$$

Taking the derivative w.r.t \mathbf{w}_k we see that

$$\begin{aligned}
\frac{\partial}{\partial \mathbf{w}_k} (\mathcal{L}(\{\mathbf{w}_\ell\})) &= - \sum_{i=1}^N \sum_{j=1}^c \frac{\partial}{\partial \mathbf{w}_k} (y_{ij} \ln(\text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_j)) \\
&= - \sum_{i=1}^N \sum_{j=1}^c y_{ij} \frac{\partial [\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top}{\partial \mathbf{w}_k} \frac{\text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_j (\mathbf{I}_{jk} - \text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_k)}{\text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_j} \\
&= - \sum_{i=1}^N \sum_{j=1}^c \mathbf{x}_i y_{ij} (\mathbf{I}_{jk} - \text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_k) \\
&= - \sum_{i=1}^N \sum_{j=1}^c \mathbf{x}_i y_{ij} \mathbf{I}_{jk} - \sum_{i=1}^N \sum_{j=1}^c \mathbf{x}_i y_{ij} \text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_k \\
&= - \sum_{i=1}^N \sum_{j=1}^c \mathbf{x}_i y_{ij} \mathbf{I}_{jk} + \sum_{i=1}^N \sum_{j=1}^c \mathbf{x}_i y_{ij} \text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_k \\
&= - \sum_{i=1}^N \mathbf{x}_i y_{ik} + \sum_{i=1}^N \mathbf{x}_i \text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_k \quad \text{where all } i^{th} \text{ row of } y_{ij} \text{ will sum to 1} \\
&= \sum_{i=1}^N x_i (\text{softmax}([\mathbf{w}_1^\top \mathbf{x} \dots \mathbf{w}_c^\top \mathbf{x}]^\top)_k - y_{ik}) \\
&= \boxed{\sum_{i=1}^N (p(\mathbf{y}_i | \mathbf{x}_i, \{\mathbf{w}_\ell\}) - y_{ik}) x_i}
\end{aligned}$$

■

End Solution

Problem 3 (Return of matrix calculus, 10pts)

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

1. Write out the negated log-likelihood of the data set, $-\ln p(D; \pi)$.
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

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

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:

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

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: $\mathbf{X}^{-\top} := (\mathbf{X}^\top)^{-1}$

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

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

Problem 3.1

If we want to find the negative log likelihood of the provided data set ($-\ln p(D; \pi) = -\ln p(\{\mathbf{x}_i, \mathbf{y}_i\}; \pi)$, which is a joint distribution), then we must recall

$$\prod_{i=1}^n p(\mathbf{x}_i, \mathbf{y}_i; \pi) = \prod_{i=1}^n p(\mathbf{y}_i = C_k; \pi) p(\mathbf{x}_i | \mathbf{y}_i = C_k)$$

Thus, to obtain the likelihood we multiply it out for each term

$$p(\mathbf{y}_i = C_k; \pi) = \pi_k$$

$$\begin{aligned} -\ln p(\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n, \pi) &= -\ln \prod_{i=1}^n \prod_{k=1}^c \pi_k^{y_{ik}} \\ &= -\sum_{i=1}^n \sum_{k=1}^c y_{ik} \ln \pi_k \end{aligned}$$

Such that at the end, because of the ln's we obtain the sum of the prior and the posterior

$$-\sum_{i=1}^n \sum_{k=1}^c y_{ik} (\ln(\pi_k) + \ln(p(\mathbf{x}_i | \mathbf{y}_i = C_k)))$$

Problem 3.2

To obtain the expression for the maximum-likelihood estimation for the prior class membership we must minimize the negative log likelihood expressed above as

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

with the constraint

$$\sum_{k=1}^c \pi_k - 1 = 0.$$

We must solve the coupled equations using the Lagrange multiplier approach. To do so we end up with the equation

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

Now, as according to the Lagrange Multipliers, we need to optimize w.r.t. π and also w.r.t. λ .

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

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

Putting the above 2 together we obtain

$$\hat{\pi}_k = \frac{\sum_{i=1}^n y_{ik}}{\sum_{k=1}^c \sum_{i=1}^n y_{ik}} = \frac{\mathbf{Y}^\top \mathbf{1}}{\mathbf{1} \mathbf{Y} \mathbf{1}}$$

Problem 3.3

So the negative likelihood is given as

$$\begin{aligned} -\ln \prod_{i=1}^n \prod_{k=1}^c p(\mathbf{x}_i | y_i = C_k)^{y_{ik}} &= \sum_{i=1}^n \sum_{k=1}^c -\ln(\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}))^{y_{ik}} \\ &= \sum_{i=1}^n \sum_{k=1}^c -y_{ik} \ln \left((2\pi)^{-1/2} (|\boldsymbol{\Sigma}|)^{-1/2} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i) \right) \right) \\ &= \sum_{i=1}^n \sum_{k=1}^c \left(\frac{1}{2} \ln(2\pi) + \frac{1}{2} \ln(|\boldsymbol{\Sigma}|) + \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i) \right) y_{ik} \\ &= \sum_{i=1}^n \sum_{k=1}^c \left(c + \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i) \right) y_{ik} \end{aligned}$$

Taking the gradient w.r.t. $\boldsymbol{\mu}_k$

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(-\ln \prod_{i=1}^n \prod_{k=1}^c p(\mathbf{x}_i | y_i = C_k)^{y_{ik}} \right) &= \sum_{i=1}^n \sum_{k=1}^c \frac{\partial}{\partial \boldsymbol{\mu}_i} \left(c + \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i) \right) y_{ik} \\ &= \frac{1}{2} \sum_{i=1}^n -(\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}^{-\top})(\mathbf{x}_i - \boldsymbol{\mu}_i) y_{ik} \quad \text{using chain rule} \\ &= \frac{1}{2} \sum_{i=1}^n -(2\boldsymbol{\Sigma}^{-1})(\mathbf{x}_i - \boldsymbol{\mu}_i) y_{ik} \\ &= \boldsymbol{\Sigma}^{-1} \sum_{i=1}^n -(\mathbf{x}_i - \boldsymbol{\mu}_i) y_{ik} \\ &= \boldsymbol{\Sigma}^{-1} \left(\sum_{i=1}^n \mathbf{x}_i y_{ik} - \boldsymbol{\mu}_i \sum_{i=1}^n y_{ik} \right) \end{aligned}$$

Problem 3.4

Here we must solve for the gradient above and set it to 0 to obtain the estimator

$$\begin{aligned}
 \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(-\ln \prod_{i=1}^n \prod_{k=1}^c p(\mathbf{x}_i | y_i = C_k)^{y_{ik}} \right) &= 0 \\
 &= \boldsymbol{\Sigma}^{-1} \left(\sum_{i=1}^n \mathbf{x}_i y_{ik} - \boldsymbol{\mu}_i \sum_{i=1}^n y_{ik} \right) \\
 &= \sum_{i=1}^n \mathbf{x}_i y_{ik} - \boldsymbol{\mu}_i \sum_{i=1}^n y_{ik} \quad \nearrow n_k
 \end{aligned}$$

$$\hat{\boldsymbol{\mu}}_k = \frac{1}{n_k} \sum_{i=1}^n \mathbf{x}_i y_{ik}$$

Problem 3.5

From the previous problem we know that

$$-\ln \prod_{i=1}^n \prod_{k=1}^c p(\mathbf{x}_i | y_i = C_k)^{y_{ik}} = \sum_{i=1}^n \sum_{k=1}^c \left(\frac{1}{2} \ln(2\pi) + \frac{1}{2} \ln(|\boldsymbol{\Sigma}|) + \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i) \right) y_{ik}$$

Now if we take the gradient of this w.r.t. $\boldsymbol{\Sigma}$ using the formulas from the cookbook we get

$$\begin{aligned}
 \frac{\partial}{\partial \boldsymbol{\Sigma}} \left(-\ln \prod_{i=1}^n \prod_{k=1}^c p(\mathbf{x}_i | y_i = C_k)^{y_{ik}} \right) &= \sum_{i=1}^n \sum_{k=1}^c \frac{\partial}{\partial \boldsymbol{\Sigma}} \left(\frac{1}{2} \ln(2\pi) + \frac{1}{2} \ln(|\boldsymbol{\Sigma}|) + \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i) \right) y_{ik} \\
 &= \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^c (\boldsymbol{\Sigma}^{-\top} - \boldsymbol{\Sigma}^{-\top} (\mathbf{x}_i - \boldsymbol{\mu}_i) (\mathbf{x}_i - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}^{-\top}) y_{ik} \\
 &= \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^c \boldsymbol{\Sigma}^{-1} (\mathbf{I} - \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i) (\mathbf{x}_i - \boldsymbol{\mu}_i)^\top) y_{ik} \\
 &= \frac{1}{2} \boldsymbol{\Sigma}^{-1} \sum_{i=1}^n \sum_{k=1}^c (\mathbf{I} - \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i) (\mathbf{x}_i - \boldsymbol{\mu}_i)^\top) y_{ik}
 \end{aligned}$$

Problem 3.6

To derive the maximum likelihood, we take the derivative and set it to zero and solve for Σ

$$\begin{aligned}\frac{\partial}{\partial \Sigma} \left(-\ln \prod_{i=1}^n \prod_{k=1}^c p(\mathbf{x}_i | y_i = C_k)^{y_{ik}} \right) &= 0 \\ &= \frac{1}{2} \Sigma^{-1} \sum_{i=1}^n \sum_{k=1}^c (\mathbf{I} - \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i)(\mathbf{x}_i - \boldsymbol{\mu}_i)^\top) y_{ik} \\ &= \sum_{i=1}^n \sum_{k=1}^c (\mathbf{I} - \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_i)(\mathbf{x}_i - \boldsymbol{\mu}_i)^\top) y_{ik}\end{aligned}$$

Where we can now continue solving for Σ

$$\begin{aligned}\sum_{i=1}^n \sum_{k=1}^c \mathbf{I} y_{ik} &= \Sigma^{-1} \sum_{i=1}^n \sum_{k=1}^c (\mathbf{x}_i - \boldsymbol{\mu}_i)(\mathbf{x}_i - \boldsymbol{\mu}_i)^\top y_{ik} \\ \mathbf{I} \sum_{i=1}^n \sum_{k=1}^c y_{ik} &= \Sigma^{-1} \sum_{i=1}^n \sum_{k=1}^c (\mathbf{x}_i - \boldsymbol{\mu}_i)(\mathbf{x}_i - \boldsymbol{\mu}_i)^\top y_{ik} \\ \Sigma &= \frac{1}{n_i} \sum_{i=1}^n \sum_{k=1}^c (\mathbf{x}_i - \boldsymbol{\mu}_i)(\mathbf{x}_i - \boldsymbol{\mu}_i)^\top y_{ik}\end{aligned}$$

End Solution

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:

```
1 fruit,width,height
2 1,8.4,7.3
3 1,8,6.8
4 1,7.4,7.2
5 1,7.1,7.8
6 ...
```

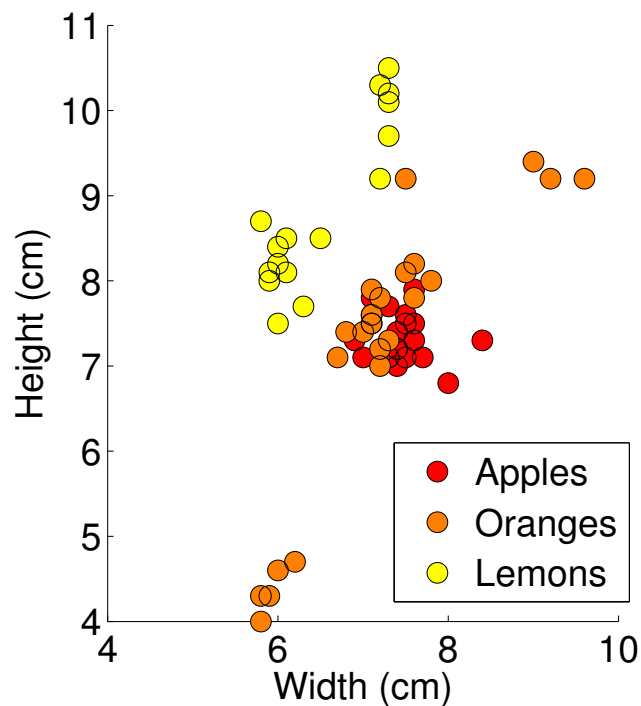


Figure 3: Heights and widths of apples, oranges, and lemons. These fruit were purchased and measured by Iain Murray: http://homepages.inf.ed.ac.uk/imurray2/teaching/oranges_and_lemons/.

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

The `visualize()` method for each classifier will save a plot that will show the decision boundaries. You should include these in this assignment.

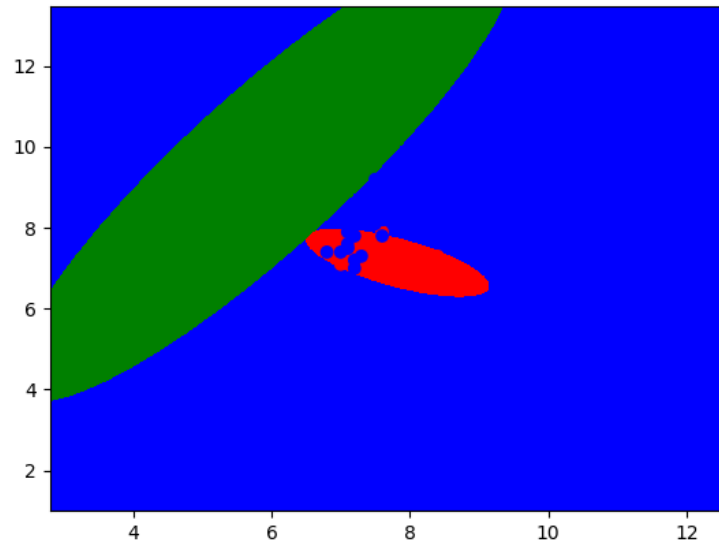


Figure 4: Separate Covariances Generative Model Plot

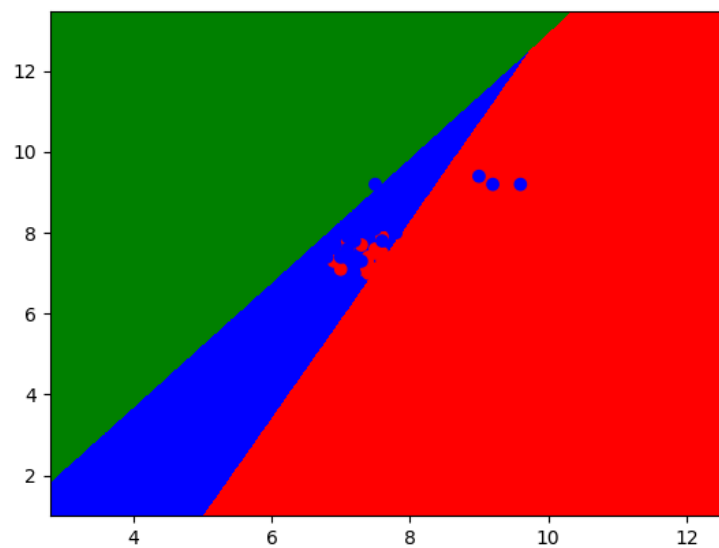


Figure 5: Separate Covariances Generative Model Plot

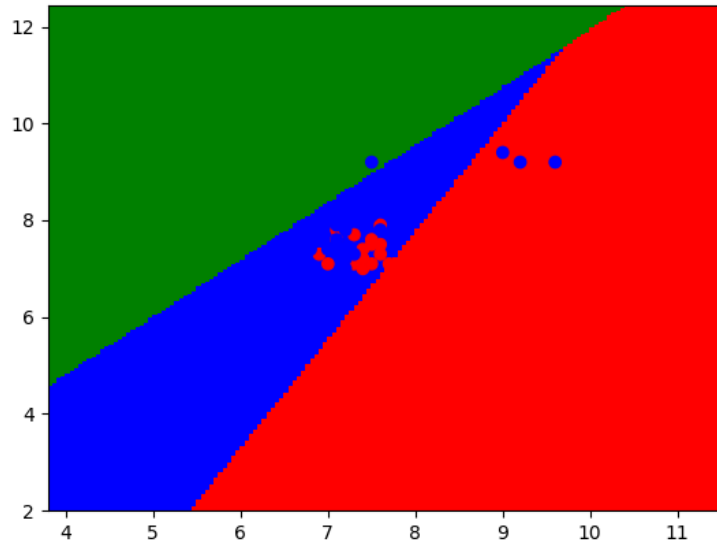


Figure 6: Logistic Regression Plot

Which classifiers model the distributions well?

Since the data is extremely difficult to classify well, given that some of the values of the fruits overlay, purely based on the appearances of the plots, I believe the Generative model with Separate Covariances did the best out of all of the 3 models.

What explains the differences?

The differences between all of the plots is that the Logistic Regression and the Shared Covariance Generative model have linear separators, whereas the Separate Covariance Generative model behaves more like multivariate Gaussians. This can be explained largely by the fact that we are incorporating all of the classification into one shared common matrix, which will limit the 'degrees of freedom' of the model, not allowing it to behave like the one with the Separate Covariance Generative models.

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).

	$\eta = 0.00001$	$\eta = 0.0001$	$\eta = 0.001$
$\lambda = 0.0$	34.55	31.20	36.36
$\lambda = 0.0001$	34.56	31.20	36.38
$\lambda = 0.001$	34.58	31.26	36.64
$\lambda = 1.0$	43.48	42.25	52.26

Table 1: Table of Different Hyperparameters Configuration, reporting value of loss after 50,000 iterations

The final choices of the learning rate and the regularization strength are reasonable because for one, they minimize the loss as we have seen in the above table. Also, if we consider a regularization term of 0, we see that the plot greatly changes even though the value of the loss changes marginally. This is

attributed to possibly over-fitting some of the data and having wights that are too large and is completely a numerical artifact. Indeed, it is normal to see the loss increase as we increasing the regularization term, as we know we are increasing the value of the loss directly by adding the term.

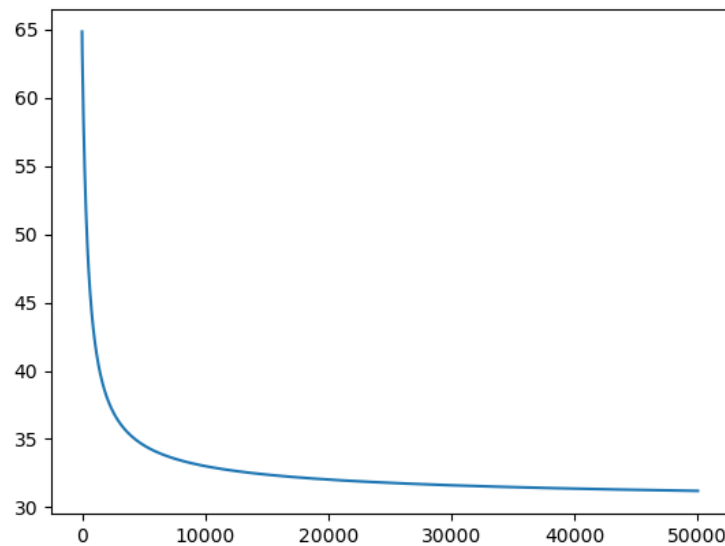


Figure 7: Logistic Regression Loss Plot for the best Hyperparameters, as highlighted in yellow in Table 1.

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.

Non-shared covariance log likelihood is: 291.065036278

Shared covariance log likelihood is: 318.265266347

End Solution

Calibration [1pt]

Approximately how long did this homework take you to complete?

Answer

Approx 24 hours

Problem4.py

```
1 # Don't change these imports. Note that the last two are the
2 # class implementations that you will implement in
3 # LogisticRegression.py and GaussianNaiveBayes.py
4 import matplotlib.pyplot as plt
5 import pandas as pd
6 from LogisticRegression import LogisticRegression
7 from GaussianGenerativeModel import GaussianGenerativeModel
8
9
10 ## These are the hyperparameters to the classifiers. You may need to
11 # adjust these as you try to find the best fit for each classifier.
12
13 # Logistic Regression parameters
14 eta = 1E-4
15 lambda_parameter = 0.0001
16 iterations=50000
17
18 # Do not change anything below this line!!
19 # -----
20
21 # Read from file and extract X and Y
22 df = pd.read_csv("fruit.csv")
23 X = df[['width', 'height']].values
24 Y = (df['fruit'] - 1).values
25
26 nb1 = GaussianGenerativeModel(isSharedCovariance=False)
27 nb1.fit(X,Y)
28 nb1.visualize("generative_result_separate_covariances.png")
29 nb1.likelihood()
30
31 nb2 = GaussianGenerativeModel(isSharedCovariance=True)
32 nb2.fit(X,Y)
33 nb2.visualize("generative_result_shared_covariances.png")
34 nb2.likelihood()
35
36 lr = LogisticRegression(eta=eta, lambda_parameter=lambda_parameter, iterations=iterations)
37 lr.fit(X,Y)
38 lr.visualize('logistic_regression_result.png')
```

Listing 1: Problem 4.py Main File

GaussianGenerativeModel.py

```
1 from scipy.stats import multivariate_normal
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import matplotlib.colors as c
5 import math
6 import pandas as pd
7
8 # Please implement the fit and predict methods of this class. You can add additional private
9 # methods
10 # by beginning them with two underscores. It may look like the _dummyPrivateMethod below.
11 # You can feel free to change any of the class attributes, as long as you do not change any of
12 # the given function headers (they must take and return the same arguments), and as long as you
13 # don't change anything in the .visualize() method.
14 class GaussianGenerativeModel:
```

```

14 def __init__(self, isSharedCovariance=False):
15     self.isSharedCovariance = isSharedCovariance
16     self.numberOfClassess=3
17     self.iterations=100
18
19 # Just to show how to make 'private' methods
20 def __dummyPrivateMethod(self, input):
21     return None
22
23 # TODO: Implement this method!
24 def fit(self, X, Y):
25     self.X = X
26     self.Y = Y
27     isSharedCovariance=self.isSharedCovariance
28     numberOfClasses = self.numberOfClassess
29     iterations=self.iterations
30
31     mu=np.zeros((2,3))
32     pik=[];sigma=[];valuesAll=[]
33
34     if isSharedCovariance: sigma=np.zeros((2,2))
35     for k in xrange(numberOfClassess): # looping over class K
36         values=X[Y==k,:]
37         valuesAll.append(values)
38         YY=Y[Y==k]
39
40         pik.append(float(len(YY))/float(len(Y)))
41         mu[0,k]=np.mean([values[:,0]])
42         mu[1,k] = np.mean([values[:,1]])
43         if not isSharedCovariance:
44             sigma.append(np.cov(values,rowvar=0))
45         if isSharedCovariance:
46             sigma = sigma+(float(len(values))/float(len(X)))*(np.cov(values, rowvar=0))
47
48
49     self.mu=mu # pass back the values of the mean
50     self.sigma=sigma # pass back the values of the variance
51     self.pik=pik
52     return
53
54 def likelihood(self):
55     mu=self.mu
56     sigma=self.sigma
57     Y=self.Y
58     X=self.X
59     C = np.array(pd.get_dummies(Y))
60     pik=self.pik
61     isSharedCovariance=self.isSharedCovariance
62
63     lkh=0
64
65     for k in xrange(C.shape[1]):
66         for i in xrange(C.shape[0]):
67
68             mult=np.asmatrix(X[i,:]-mu[:,k])
69             if isSharedCovariance:
70                 sigmaInv = np.linalg.inv(sigma)
71                 lkh += (np.log(2.*math.pi)+(1./2.)*np.log(np.linalg.det(sigma))+(1./2.)*np.
72                 matmul(np.matmul(mult,sigmaInv),mult.T))*C[i,k]-np.log(pik[k])
73             else:
74                 sigmaInv = np.linalg.inv(sigma[k])
75                 lkh += (np.log(2. * math.pi) + (1. / 2.) * np.log(np.linalg.det(sigma[k])) +
76                 (1. / 2.) * np.matmul(
77                     np.matmul(mult, sigmaInv), mult.T)) * C[i, k] - np.log(pik[k])
78             if isSharedCovariance:

```

```

77         print 'shared covariance log likelihood is: '+str(lkh[0,0])+'\n'
78     else:
79         print 'non-shared covariance log likelihood is: '+str(lkh[0,0])+'\n'
80
81
82     # TODO: Implement this method!
83     def predict(self, X_to_predict):
84         # The code in this method should be removed and replaced! We included it just so that the
            distribution code
85         # is runnable and produces a (currently meaningless) visualization.
86         isSharedCovariance=self.isSharedCovariance
87         mu=self.mu
88         pik=self.pik
89         sigma=self.sigma
90         predict=np.zeros((X_to_predict.shape[0],3))
91         for k in xrange(3):
92             if not isSharedCovariance:
93                 mv=multivariate_normal(mu[:,k],sigma[k]) #create a multivariate distribution for
each class with shared variance
94             if isSharedCovariance:
95                 mv = multivariate_normal(mu[:, k], sigma)
96                 predict[:,k]=mv.pdf(X_to_predict)*pik[k]# obtain the PDF for a given class for each
of the vales of X
97             Y=np.argmax(predict,axis=1)
98
99         return np.array(Y)
100
101     # Do not modify this method!
102     def visualize(self, output_file, width=3, show_charts=False):
103         X = self.X
104
105         # Create a grid of points
106         x_min, x_max = min(X[:, 0] - width), max(X[:, 0] + width)
107         y_min, y_max = min(X[:, 1] - width), max(X[:, 1] + width)
108         xx,yy = np.meshgrid(np.arange(x_min, x_max, .005), np.arange(y_min,
109         y_max, .005))
110
111         # Flatten the grid so the values match spec for self.predict
112         xx_flat = xx.flatten()
113         yy_flat = yy.flatten()
114         X_topredict = np.vstack((xx_flat, yy_flat)).T
115
116         # Get the class predictions
117         Y_hat = self.predict(X_topredict)
118         Y_hat = Y_hat.reshape((xx.shape[0], xx.shape[1]))
119
120         cMap = c.ListedColormap(['r','b','g'])
121
122         # Visualize them.
123         plt.figure()
124         plt.pcolormesh(xx,yy,Y_hat, cmap=cMap)
125         plt.scatter(X[:, 0], X[:, 1], c=self.Y, cmap=cMap)
126         plt.savefig(output_file)
127         if show_charts:
128             plt.show()

```

Listing 2: Gaussian Generative Model File

LogisticRegression.py


```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.colors as c
4 import pandas as pd
5 from scipy.misc import logsumexp
6
7 # Please implement the fit and predict methods of this class. You can add additional private
  methods
8 # by beginning them with two underscores. It may look like the __dummyPrivateMethod below.
9 # You can feel free to change any of the class attributes, as long as you do not change any of
10 # the given function headers (they must take and return the same arguments), and as long as you
11 # don't change anything in the .visualize() method.
12 class LogisticRegression:
13     def __init__(self, eta, lambda_parameter, iterations):
14         self.eta = eta
15         self.lambda_parameter = lambda_parameter
16         self.numberOfClassess=3
17         self.iterations=iterations
18
19     # Just to show how to make 'private' methods
20     def __dummyPrivateMethod(self, input):
21         return None
22
23
24     # TODO: Implement this method!
25     def fit(self, X, C):
26         self.X = X
27         self.C = C
28
29         lambda_parameter=self.lambda_parameter
30         eta = self.eta
31         iterations=self.iterations
32         numberOfClassess=self.numberOfClassess
33
34         C=np.array(pd.get_dummies(C)) # convert all of the values to one-hot vectors
35
36         Bias=np.ones((C.shape[0])) #create a bias column
37         X=np.column_stack((Bias,X)) # stack the bias column
38
39         w=np.ones((X.shape[1],numberOfClassess)) # initialize a weight matrix
40         loss=[]; itAll=[]
41         for it in xrange(iterations):
42             wx=np.matmul(X,np.transpose(w)) ### double check this to make sure its correct
43             sf= np.exp(wx) # convert the wx to the exponential form for the softmax function
44             for i in xrange(wx.shape[0]): # for all the classes
45                 sumSf=np.sum(sf[i,:]) # create a sum vector for the soft max
46                 sf[i,:]= sf[i,:]/sumSf # divide each class by the sum vector for that respective
47
48         class
49
50             L=np.matmul(np.transpose((sf-C)),X)
51             loss.append(-np.sum((C*np.log(sf)))+lambda_parameter*np.sum(w*w))
52
53             w=w-eta*(L+2*lambda_parameter*w)
54             itAll.append(it)
55
56         plt.figure()
57         plt.plot(itAll,loss)
58         plt.savefig('logistic_regression_loss_new.png')
59
60         print 'the final loss for the logistic regression model is: '+str(loss[-1])+'\n'
61
62         self.weights=w
63         return
64
65     # TODO: Implement this method!
66     def predict(self, X_to_predict):

```

```

64     # The code in this method should be removed and replaced! We included it just so that the
distribution code
65     # is runnable and produces a (currently meaningless) visualization.
66     w=self.weights
67
68     Bias = np.ones((X_to_predict.shape[0])) # create a bias column
69     X_to_predict = np.column_stack((Bias, X_to_predict)) # stack the bias column
70
71     Y=np.argmax(np.matmul(X_to_predict ,np.transpose(w)),axis=1)
72
73     return Y
74
75 def visualize(self, output_file , width=2, show_charts=False):
76     X = self.X
77
78     # Create a grid of points
79     x_min, x_max = min(X[:, 0] - width), max(X[:, 0] + width)
80     y_min, y_max = min(X[:, 1] - width), max(X[:, 1] + width)
81     xx,yy = np.meshgrid(np.arange(x_min, x_max, .05), np.arange(y_min,
82     y_max, .05))
83
84     # Flatten the grid so the values match spec for self.predict
85     xx_flat = xx.flatten()
86     yy_flat = yy.flatten()
87     X_topredict = np.vstack((xx_flat ,yy_flat)).T
88
89     # Get the class predictions
90     Y_hat = self.predict(X_topredict)
91     Y_hat = Y_hat.reshape((xx.shape[0], xx.shape[1]))
92
93     cMap = c.ListedColormap(['r','b','g'])
94
95     # Visualize them.
96     plt.figure()
97     plt.pcolormesh(xx,yy,Y_hat, cmap=cMap)
98     plt.scatter(X[:, 0], X[:, 1], c=self.C, cmap=cMap, lw=1)
99     plt.savefig(output_file)
100     if show_charts:
101         plt.show()

```

Listing 3: Logistic Regression Model File