

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 three 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 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 third 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

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

This prior makes sense when you have little prior information and do not know much about the relationship among features so you can simplify by assuming independence.

4. 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

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

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

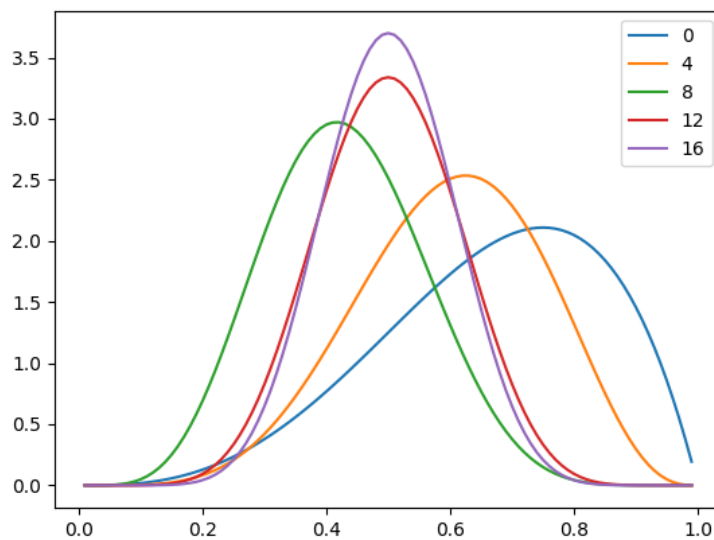
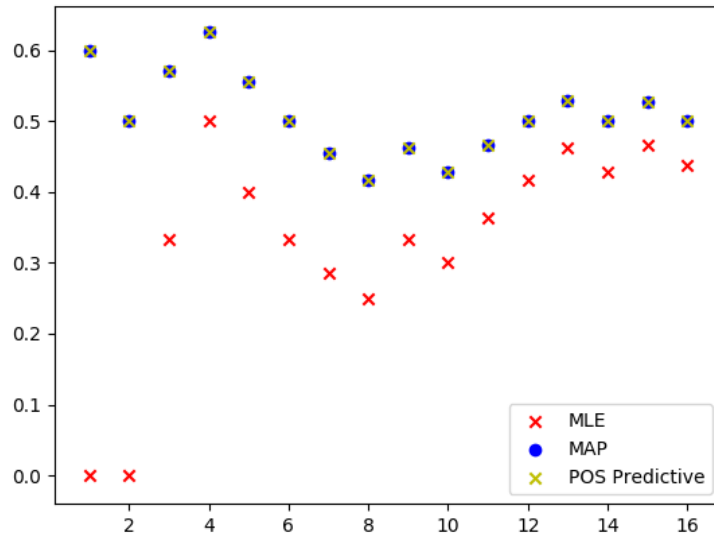
0.1 solution

1.

$$w_{MLE} = \operatorname{argmax}_w p(D|\mathbf{w})$$

$$w_{MAP} = \operatorname{argmax}_w p(\mathbf{w}|D) = \operatorname{argmax}_w p(D|\mathbf{w})p(w)$$

$$p(y|D, \mathbf{x}) = \int_{\mathbf{w}} p(y|\mathbf{w}, \mathbf{x})p(\mathbf{w}|D)d\mathbf{w}$$



2.

3. The MAP and posterior predictive follow the exact same points because the posterior mean is equal to the posterior mode, effectively making the -1 and the -2 in the numerator and denominator respectively. The MLE generally follows the same convergence path towards the 0.5 that is expected from the Bernoulli distribution, however at the beginning with two 0s, it's clear that MLE, without the beta prior is far off, and has also not converged to 0.5 by 16 points like both MAP and posterior predictive have.

4.

$$\begin{aligned}
& \log p(w, y) \propto \log p(y|w)p(w) \\
& = \log[(\det 2\pi\beta^{-1}I)^{-0.5} \exp(-1/2(y - Xw)^\top \beta^{-1}I(y - Xw) \cdot (\det 2\pi\alpha^{-1}I)^{-0.5} \exp(-1/2w^\top \alpha^{-1}Iw)] \\
& \quad \propto -(\beta/2)(y - Xw)^\top (y - Xw) - (\alpha/2)w^\top w \\
& = -(\alpha/2)w^\top w - (\beta/2)y^\top y - (\beta/2)(Xw)^\top (Xw) + \beta y^\top Xw \\
& = -1/2 \begin{pmatrix} w \\ y \end{pmatrix}^\top \begin{pmatrix} \alpha I + \beta X^\top X & -X^\top \\ -\beta X & \beta \end{pmatrix} \begin{pmatrix} w \\ y \end{pmatrix}
\end{aligned}$$

Using Gaussian conjugacy to find the conditional Gaussian from the joint (Murphy equation 4.69):

$$\begin{aligned}
\Sigma_{w|y} &= \Lambda_{w,w}^{-1} = (\alpha I + \beta X^\top X)^{-1} \\
\mu_{w|y} &= \Sigma_{w|y}(\Sigma_w^{-1}\mu + X^\top \Sigma_y^{-1}y) = \beta \Sigma_{x|y} X^\top y
\end{aligned}$$

Problem 2 (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: $\mathbf{X}^{-\top} := (\mathbf{X}^{\top})^{-1}$

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

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

0.2 Solution

1.

$$-\ln p(D|\pi) = -\sum_{i=1}^n \sum_{k=1}^c y_{i,k} [\ln p(x_i|y_i = C_k) + \ln \pi_k]$$

2. Adding the Lagrangian and differentiating:

$$\frac{\partial L}{\partial \pi} = \frac{\partial}{\partial \pi} \left[-\sum_{i=1}^n \sum_{k=1}^c y_{i,k} [\ln p(x_i|y_i = C_k) + \ln \pi_k] + \lambda \sum_k (\pi_k - 1) \right]$$

$$0 = 1/(\pi_k) \sum_{i=1}^n y_{i,k} + \lambda$$

$$\pi_k = (1/\lambda) \sum_{i=1}^n y_{i,k}$$

Since we differentiated with respect to λ :

$$\sum_k \pi_k - 1 = 0$$

$$\sum_k -(1/\lambda) \sum_{i=1}^n y_{i,k} - 1 = 0$$

$$\lambda = -N$$

Thus:

$$\pi_k = -(1/\lambda) \sum_{i=1}^n y_{i,k} = (1/N) \sum_{i=1}^n y_{i,k}$$

3. Using the normal PDF, our log-likelihood equation, calculating the gradient, and using line 81 of the Matrix cookbook:

$$\begin{aligned} -\ln p(D|\pi) &= -\sum_{i=1}^n \sum_{k=1}^c y_{i,k} [\ln p(x_i|y_i = C_k) + \ln \pi_k] \\ &= -\sum_{i=1}^n \sum_{k=1}^c y_{i,k} \left[\ln [(2\pi \det(\Sigma))^{-1/2} \exp(-(1/2)(x - \mu)^T \Sigma^{-1} (x - \mu))] + \ln \pi_k \right] \\ &= -\sum_{i=1}^n \sum_{k=1}^c y_{i,k} \left[-(1/2) \ln(2\pi) - (1/2) \ln(\det(\Sigma)) - (1/2)(x - \mu)^T \Sigma^{-1} (x - \mu) + \ln \pi_k \right] \\ \frac{\partial}{\partial \mu} \ln p(D|\pi) &= (1/2) \sum_{i=1}^n \sum_{k=1}^c y_{i,k} \frac{\partial}{\partial \mu} [(x - \mu)^T \Sigma^{-1} (x - \mu)] \\ &= (1/2) \sum_{i=1}^n \sum_{k=1}^c y_{i,k} \frac{\partial}{\partial \mu} [-x^T \Sigma^{-1} \mu - \mu^T \Sigma^{-1} x + \mu^T \Sigma^{-1} \mu] \\ &= (1/2) \sum_{i=1}^n \sum_{k=1}^c y_{i,k} [-2x^T \Sigma^{-1} + (\Sigma^{-1} + \Sigma^{-T}) \mu] \end{aligned}$$

4. Setting the gradient to zero and solving for μ :

$$0 = (1/2) \sum_{i=1}^n \sum_{k=1}^c y_{i,k} [-2x^T \Sigma^{-1} + (\Sigma^{-1} + \Sigma^{-T})\mu]$$

$$\mu = \frac{4 \sum_{i=1}^n \sum_{k=1}^c y_{i,k} x^T \Sigma^{-1}}{\sum_{i=1}^n \sum_{k=1}^c y_{i,k} (\Sigma^{-1} + \Sigma^{-T})}$$

5. We can use equation 61 from the Matrix Cookbook:

$$\begin{aligned} \frac{-\partial}{\partial \Sigma} &= - \sum_{i=1}^n \sum_{k=1}^c y_{i,k} \frac{\partial}{\partial \Sigma} [- (1/2) \ln(2\pi) - (1/2) \ln(\det(\Sigma)) - (1/2)(x - \mu)^T \Sigma^{-1}] \\ &= (1/2) \sum_{i=1}^n \sum_{k=1}^c y_{i,k} [\Sigma^{-T} - \Sigma^{-T}(x - \mu)(x - \mu)^T \Sigma^{-T}] \end{aligned}$$

3. Classifying Fruit [15pts]

You're tasked with classifying three different kinds of fruit, based on their heights and widths. Figure 1 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
...
```

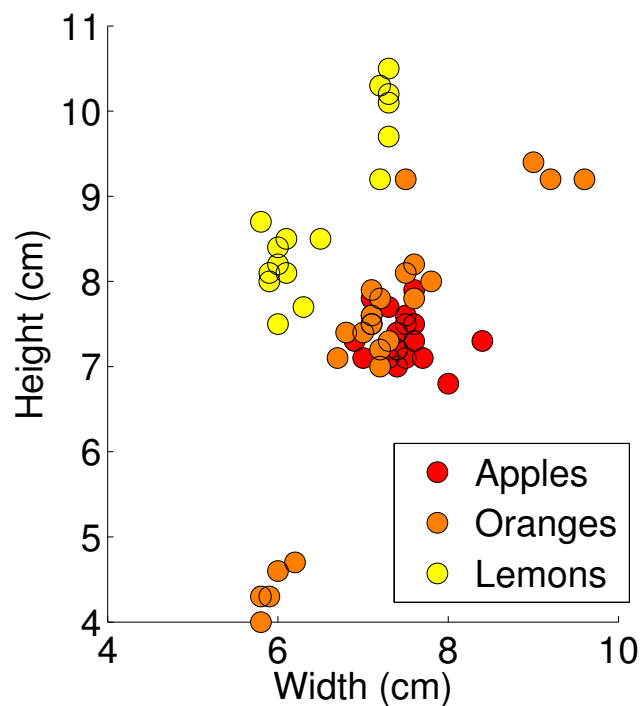


Figure 1: Heights and widths of apples, oranges, and lemons. These fruit were purchased and measured by Iain Murray: http://homepages.inf.ed.ac.uk/imurray2/teaching/oranges_and_lemons/.

Problem 3 (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 . See the third practice problem in the section 3 notes for information about multi-class logistic regression, softmax, and negative log likelihood.
- 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 `problem3.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?

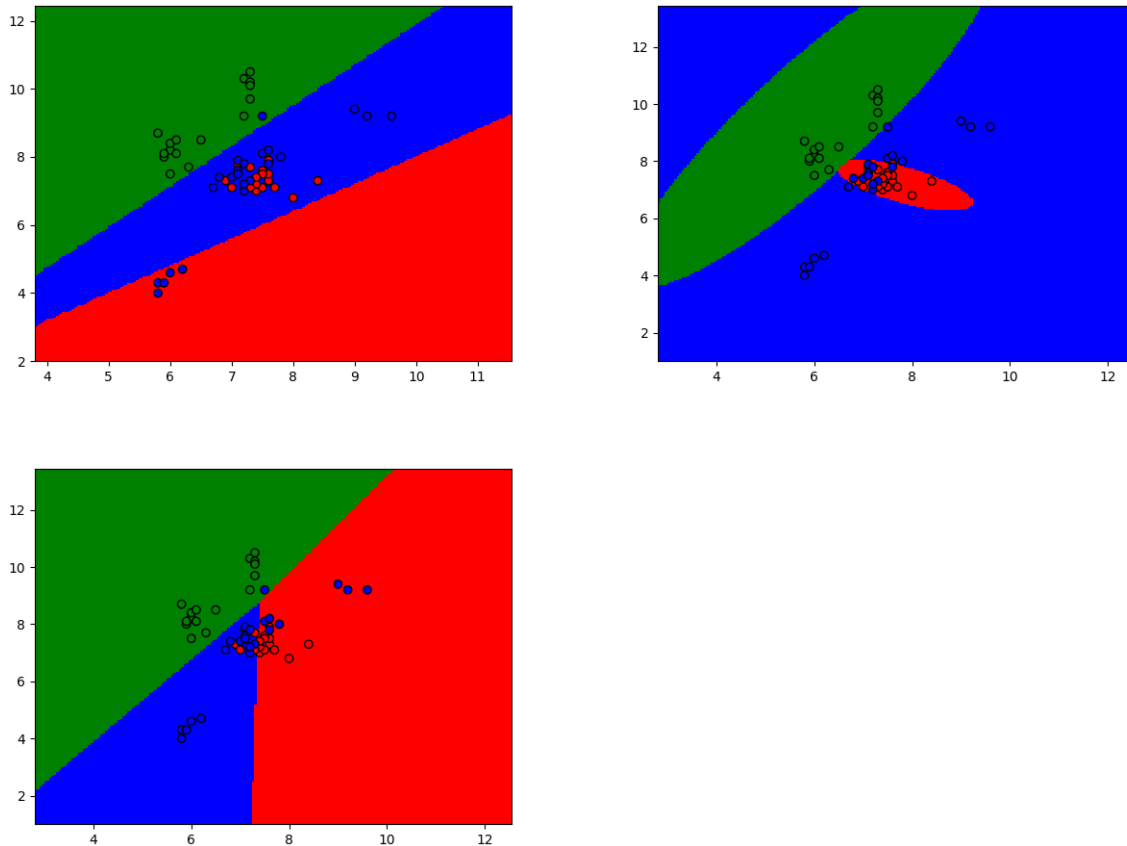
In addition to comparing the decision boundaries of the three models visually:

- For logistic regression, plot negative log-likelihood loss with iterations on the x-axis and loss on the y-axis for several configurations of hyperparameters. Note which configuration yields the best final loss. Why are your final choices of learning rate (η) and regularization strength (λ) reasonable? How does altering these hyperparameters affect convergence? Focus both on the ability to converge and the rate at which it converges (a qualitative description is sufficient).
- For both Gaussian generative models, report negative log likelihood. In the separate covariance matrix case, be sure to use the covariance matrix that matches the true class of each data point.

Finally, consider a fruit with width 4cm and height 11cm. To what class do each of the classifiers assign this fruit? What do these results tell you about the classifiers' ability to generalize to new data? Repeat for a fruit of width 8.5cm and height 7cm.

0.3 Solution

Plots of Logistic Regression, Gaussian without shared covariance, and Gaussian with shared covariance.

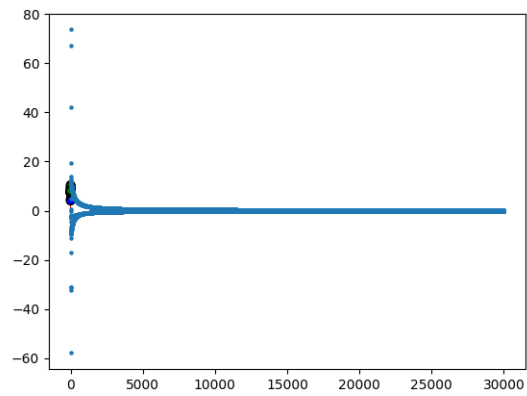


My model for the Logistic Regression isn't quite correct; the point where the two lines meet should be right in the middle of the center cluster of points, to separate out the blue points from the red, however the generative Gaussian model seems pretty good. With the current data, the Gaussian without shared covariance seems quite good, however the area which it would predict apple is quite small, and doesn't leave much variation for potential apple shapes. It is biased by the small cluster of red data points compared to the varied blue data points. The Gaussian with shared covariance seems a bit better to predict more variety of apple sizes, though it clearly leaves many of the blue points in the top right of the graph outside the prediction area.

For logistic regression, the separation lines have to be linear, which even with an accurate regression model would have limitations for this point. For Gaussian without shared covariance, since each class' variance is treated separately, this leads to prediction areas of varying sizes, since the variance per class is modeled separately. This is why the red prediction area is much smaller than the green and blue, and why the red area clearly has extended to cover the two outer points on the right side of the central cluster. The shared covariance graph is more evenly separated, since in the center the blue and red points are quite close together, which with a shared covariance means tight, separate prediction areas is not reasonable.

My graph below for the loss over iterations is currently not correct, and I've not had the time to plot more η or λ values. The reason for this error, and overall error in my logistic regression model, is that I had been combining the overall loss per iteration correctly. It's obvious from this, then, that the loss

jumps around too much, rather than monotonically increasing/decreasing towards the convergence point.



For the fruits of size 4x11cm and 8.5x7cm In the same order as above: LR: lemon, orange GSeparate: lemon, apple GShared: orange, apple

Calibration [1pt]

Approximately how long did this homework take you to complete? 15-20 hours

Name, Email, and Collaborators

Name: David Hughes

Email: davidralphhughes@college.harvard.edu

Collaborators: Alexander Muoz