

# Homework #2

CS1810-S25

Due: February 28, 2025 at 11:59 PM

---

## Classification and Bias-Variance Trade-offs

### Introduction

This homework is about classification, bias-variance trade-offs, and uncertainty quantification.

The datasets that we will be working with relate to astronomical observations and loan applicants. The first dataset, found at `data/planet-obs.csv`, contains information on whether a planet was observed (as a binary variable) at given points in time. This will be used in Problem 1. The second dataset, available at `data/hr.csv`, details different loan applicants and their measured debt to income ratio and credit score. You will work with this data in Problem 3.

As a general note, for classification problems we imagine that we have the input matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  (or perhaps they have been mapped to some basis  $\Phi$ , without loss of generality) with outputs now “one-hot encoded.” This means that if there are  $K$  output classes, rather than representing the output label  $y$  as an integer  $1, 2, \dots, K$ , we represent  $\mathbf{y}$  as a “one-hot” vector of length  $K$ . A “one-hot” vector is defined as having every component equal to 0 except for a single component which has value equal to 1. For example, if there are  $K = 7$  classes and a particular data point belongs to class 3, then the target vector for this data point would be  $\mathbf{y} = [0, 0, 1, 0, 0, 0, 0]$ . We will define  $C_1$  to be the one-hot vector for the 1st class,  $C_2$  for the 2nd class, etc. Thus, in the previous example  $\mathbf{y} = C_3$ . If there are  $K$  total classes, then the set of possible labels is  $\{C_1 \dots C_K\} = \{C_k\}_{k=1}^K$ . Throughout the assignment we will assume that each label  $\mathbf{y} \in \{C_k\}_{k=1}^K$  unless otherwise specified. The most common exception is the case of binary classification ( $K = 2$ ), in which case labels are the typical integers  $y \in \{0, 1\}$ .

### Resources and Submission Instructions

We encourage you to read CS181 Textbook’s Chapter 3 for more information on linear classification, gradient descent, and classification in the discriminative setting. Read Chapter 2.8 for more information on the trade-offs between bias and variance.

In problems 1 and 3, you may use `numpy` or `scipy`, but not `scipy.optimize` or `sklearn`. Example code is given in the provided notebook. **We highly recommend that you use Google Colab for problems 1 and 3 to avoid numerical stability issues.**

Please type your solutions after the corresponding problems using this L<sup>A</sup>T<sub>E</sub>X template, and start each problem on a new page.

Please submit the **writeup PDF to the Gradescope assignment ‘HW2’**. Remember to assign pages for each question. Please submit your **L<sup>A</sup>T<sub>E</sub>X file and code files to the Gradescope assignment ‘HW2 - Supplemental’**. **You must include your plots in your writeup PDF**. The supplemental files will only be checked in special cases, e.g. honor code issues, etc.

**Problem 1** (Exploring Bias-Variance and Uncertainty)

In this problem, we will explore the bias and variance of a few different model classes when it comes to logistic regression and investigate two sources of predictive uncertainty in a synthetic (made-up) scenario.

We are using a powerful telescope in the northern hemisphere to gather measurements of some planet of interest. At certain times however, our telescope is unable to detect the planet due to its positioning around its star. The data in `data/planet-obs.csv` records the observation time in the “Time” column and whether the planet was detected in the “Observed” column (with the value 1 representing that it was observed). These observations were taken over a dark, clear week, which is representative of the region. Since telescope time is expensive, we would like to build a model to help us schedule and find times when we are likely to detect the planet.

1. Split the data into 10 mini-datasets of size  $N = 30$  (i.e. dataset 1 consists of the first 30 observations, dataset 2 consists of the next 30, etc. This has already been done for you). Consider the three bases  $\phi_1(t) = [1, t]$ ,  $\phi_2(t) = [1, t, t^2]$ , and  $\phi_3(t) = [1, t, t^2, t^3, t^4, t^5]$ . For each of these bases, fit a logistic regression model using  $\text{sigmoid}(\mathbf{w}^\top \phi(t))$  to each dataset by using gradient descent to minimize the negative log-likelihood. This means you will be running gradient descent 10 times for each basis, once for each dataset.

Use the given starting values of  $\mathbf{w}$  and a learning rate of  $\eta = 0.001$ , take 1,000 update steps for each gradient descent run, and make sure to average the gradient over the data points at each step. These parameters, while not perfect, will ensure your code runs reasonably quickly.

2. After consulting with a domain expert, we find that the probability of observing the planet is periodic as the planet revolves around its star—we are more likely to observe the planet when it is in front of its star than when it is behind it. In fact, the expert determines that observation follows the generating process  $y \sim \text{Bern}(f(t))$ , where  $f(t) = 0.4 \times \cos(1.1t + 1) + 0.5$  for  $t \in [0, 6]$  and  $y \in \{0, 1\}$ . Note that we, the modelers, do not usually see the true data distribution. Knowledge of the true  $f(t)$  is only exposed in this problem to allow for verification of the true bias.

Use the given code to plot the true process versus your learned models. Include your plots in your solution PDF.

**In no more than 5 sentences**, explain how bias and variance reflected in the 3 types of curves on the graphs. How do the fits of the individual and mean prediction functions change? Keeping in mind that none of the model classes match the true generating process exactly, discuss the extent to which each of the bases approximates the true process.

### Problem 1 (cont.)

3. If we were to increase the size of each dataset drawn from  $N = 30$  to a larger number, how would the bias and variance change for each basis? Why might this be the case? You may experiment with generating your own data that follows the true process and plotting the results, but this is **not** necessary. **Your response should not be longer than 5 sentences.**
4. Consider the test point  $t = 0.1$ . Using your models trained on basis  $\phi_3$ , report the predicted probability of observation of the *first* model (the model trained on the first 30 data points). How can we interpret this probability as a measure of uncertainty? Then, compute the variance of the classification probability over your 10 models at the same point  $t = 0.1$ . How does this measurement capture another source of uncertainty, and how does this differ from the uncertainty represented by the classification probability? Repeat this process (reporting the first model's classification probability and the variance over the 10 models) for the point  $t = 3.2$ .

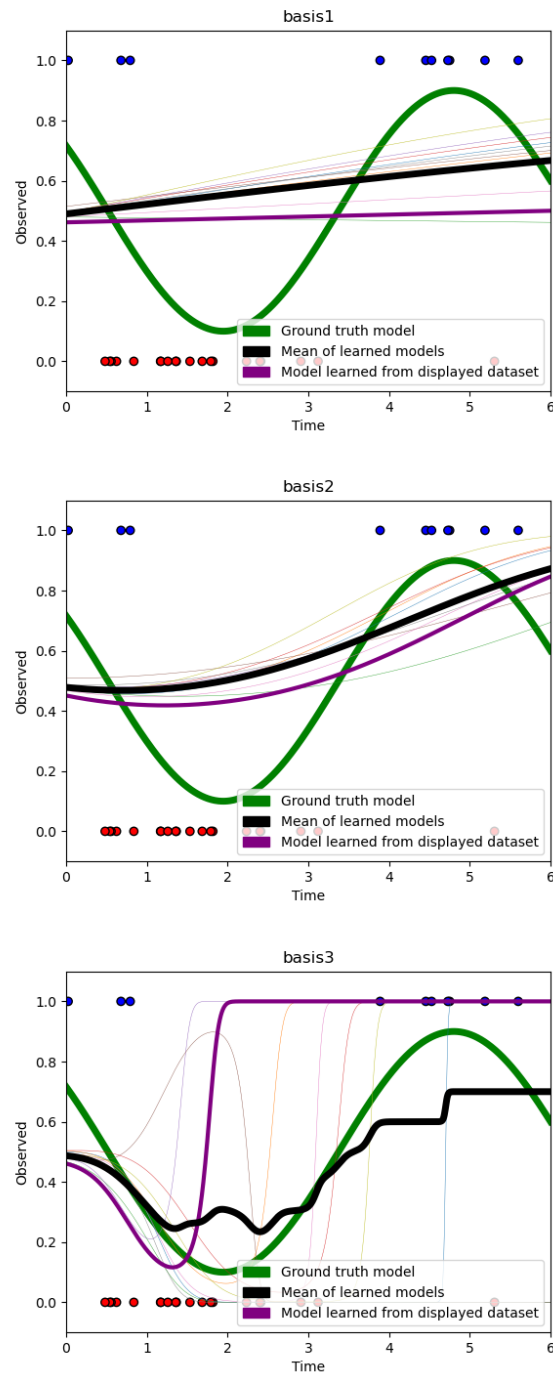
Compare the uncertainties and their sources at times  $t = 0.1$  and  $t = 3.2$ .

5. We now need to make some decisions about when to request time on the telescope. The justifications of your decisions will be sent to your funding agency, which will determine whether you will be allocated funds to use the telescope for your project. **In no more than 10 lines**, answer the following questions.
  - To identify the ideal time, which model(s) would you use and why?
  - What time would you request, and why?
  - Your funding agency suggests using a different telescope in a humid area near the equator. Can you still use your model to determine when the planet is likely to be visible? Why? Are there adaptations that may be necessary?
  - You seek out a team that has used the alternative telescope for observing this planet, and they provide you their observation file `data/planet-obs-alternate.csv`. Compare the observations from your telescope to theirs. What seems to be happening? What might be an appropriate model for this? Your funding agency asks you to refit your models on these new data. Do you think this is a reasonable ask, and if so, how will it help you make better decisions about when to request viewing time? If not, why do you think the additional modeling will not help? You do *not* need to do any modeling for this question!

In these questions, we are looking for your reasoning; there may be more than one valid answer.

Solution:

2. We have the following:



Note that as we go from basis1 to basis3, the curves get progressive more crazy while adhering to the curve itself better (bias going down while variance goes up). We can also watch the progression from underfit to overfit by observing the mean of learned models - the mean of the learned models almost

has the sinusoidal shape of the true process in basis3, but there are many more fluctuations in the curve.

Basis1 is pretty trash, with wild underfitting, and basis2 is a tiny bit better but also quite bad.

- 3) We implement this, and we see that the bias plateaus while the variance continues to decrease. This makes sense since we're trying to use multiple polynomial functions to approximate a sinusoidal function, so there's a fundamental issue with trying to predict a sinusoidal function using a polynomial basis.

- 4) We have the following:

$t = 0.1$  :  $p = 0.49926$ , variance = 0.0027834

$t = 3.2$  :  $p = 0.99999$ , variance = 0.249999

For  $t = 0.1$ , probability is 0.5, which means that the model is still very uncertain (equally likely to predict 0 or 1). The variance is very low, meaning that models are all quite similar and agree on the prediction.

For  $t = 3.2$ , probability is 1, which means that the model is very certain that the planet will be observed. However, the variance is relatively higher, meaning that models are quite different and disagree on the prediction.

- 5) To identify the ideal time, I would use the basis3 model. It's better to have low bias and high variance, since you can continually add data in order to decrease the variance, but decreasing bias is very difficult if your model is fundamentally flawed. I would request at time  $t = 3.2$  since it showed the highest probability of observation.

If we got a new telescope, we can still use our approach, but we would have to train on different data.

Comparing the observations from the other telescope, we can see that there's no clear sinusoidal wave because everything is rather uniformly spread on 1 and 0 so there is no clear decision boundary. It wouldn't be reasonable to refit my models on the data because the polynomial fit would either result in underfitting with high bias, or extreme overfitting because the data is so sporadic.

**Problem 2** (Maximum likelihood in classification)

Consider now a generative  $K$ -class model. We adopt class prior  $p(\mathbf{y} = C_k; \boldsymbol{\pi}) = \pi_k$  for all  $k \in \{1, \dots, K\}$  (where  $\pi_k$  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  $\mathbf{y}_i \in \{C_k\}_{k=1}^K$  is encoded as a one-hot target vector and the data are independent.

1. Write out the 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. Briefly state why your final answer is intuitive.

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

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

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

3. Derive the gradient of the 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  $\hat{\boldsymbol{\mu}}_k$  for vector  $\boldsymbol{\mu}_k$ . Briefly state why your final answer is intuitive.
5. Derive the gradient for the 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  $\hat{\boldsymbol{\Sigma}}$  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  $\lambda$  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. First, we identify the likelihood function as the product of  $p(\mathbf{x}_i, \mathbf{y}_i \mid \boldsymbol{\pi})$  over all  $i$ : we are given that the data are independent, so we can express this as a product:

$$\begin{aligned} p(D; \boldsymbol{\pi}) &= \prod_{i=1}^n p(\mathbf{x}_i, \mathbf{y}_i; \boldsymbol{\pi}) \\ &= \prod_{i=1}^n p(\mathbf{x}_i | \mathbf{y}_i) p(\mathbf{y}_i; \boldsymbol{\pi}) \\ &= \prod_{i=1}^n \prod_{k=1}^K p(\mathbf{x}_i | \mathbf{y} = C_k)^{I(\mathbf{y}_i = C_k)} p(\mathbf{y}_i = C_k; \boldsymbol{\pi}) \\ &= \prod_{i=1}^n \prod_{k=1}^K p(\mathbf{x}_i | \mathbf{y} = C_k)^{I(\mathbf{y}_i = C_k)} \pi_k \end{aligned}$$

Then, we take the log of the likelihood function to get

$$\ln p(D; \boldsymbol{\pi}) = \sum_{i=1}^n \sum_{k=1}^K I(\mathbf{y}_i = C_k) \ln p(\mathbf{x}_i | \mathbf{y} = C_k) + \sum_{i=1}^n \sum_{k=1}^K I(\mathbf{y}_i = C_k) \ln \pi_k.$$

2. Using the hint from below, we want to optimize the log-likelihood function subject to the constraint  $\sum_k \pi_k - 1 = 0$ . We can do this by constructing the Lagrangian function:

$$L(\boldsymbol{\pi}, \lambda) = \sum_{i=1}^n \sum_{k=1}^K I(\mathbf{y}_i = C_k) \ln p(\mathbf{x}_i | \mathbf{y} = C_k) + \sum_{i=1}^n \sum_{k=1}^K I(\mathbf{y}_i = C_k) \ln \pi_k + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right).$$

Taking the derivative of  $L$  with respect to  $\pi_k$  and setting it to zero, we get

$$\frac{\partial L}{\partial \pi_k} = \sum_{i=1}^n I(\mathbf{y}_i = C_k) \frac{1}{\pi_k} + \lambda = 0.$$

Solving for  $\pi_k$ , we get

$$\hat{\pi}_k = -\frac{1}{\lambda} \sum_{i=1}^n I(\mathbf{y}_i = C_k).$$

Note that we must have the sum of all  $\pi_k$  equal to 1, so we can substitute  $\lambda = -n$  to get

$$\hat{\pi}_k = \frac{1}{n} \sum_{i=1}^n I(\mathbf{y}_i = C_k).$$

This is intuitive because the maximum likelihood estimator for the prior class-membership probabilities is the proportion of data points that belong to class  $C_k$ .

3. We start by writing the log-likelihood function as

$$\ln p(D; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^n \ln p(\mathbf{x}_i | \mathbf{y}_i; \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \ln p(\mathbf{y}_i; \boldsymbol{\pi}).$$

Note that we can express

$$\ln p(\mathbf{x}_i | \mathbf{y}_i; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) = -\frac{1}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k).$$

We can then take the gradient of the log-likelihood with respect to  $\boldsymbol{\mu}_k$  (while ignoring constants like  $|\Sigma|$ ) to get

$$\nabla_{\boldsymbol{\mu}_k} \ln p(D; \boldsymbol{\mu}, \Sigma) = - \sum_{i=1}^n I(\mathbf{y}_i = C_k) \nabla_{\boldsymbol{\mu}_k} \left( \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right).$$

Taking the gradient, we get

$$\nabla_{\boldsymbol{\mu}_k} \ln p(D; \boldsymbol{\mu}, \Sigma) = \sum_{i=1}^n I(\mathbf{y}_i = C_k) \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k).$$

4. To find the maximum likelihood estimator  $\widehat{\boldsymbol{\mu}}_k$ , we set the gradient equal to zero

$$\sum_{i=1}^n I(\mathbf{y}_i = C_k) \Sigma^{-1} (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_k) = 0.$$

$$\widehat{\boldsymbol{\mu}}_k = \frac{\sum_{i=1}^n I(\mathbf{y}_i = C_k) \mathbf{x}_i}{\sum_{i=1}^n I(\mathbf{y}_i = C_k)}.$$

This makes sense because the MLE is just the mean for each class  $k$ .

5. First, write the log-likelihood

$$\ln p(D; \boldsymbol{\mu}, \Sigma) = \sum_{i=1}^n \ln p(\mathbf{x}_i | \mathbf{y}_i; \boldsymbol{\mu}, \Sigma) + \ln p(\mathbf{y}_i; \boldsymbol{\pi}).$$

We can express

$$\ln p(\mathbf{x}_i | \mathbf{y}_i; \boldsymbol{\mu}, \Sigma) = \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma) = -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k).$$

We can then take the gradient of the log-likelihood with respect to  $\Sigma$  (while ignoring constants like  $|\Sigma|$ ) to get

$$\nabla_{\Sigma} \ln p(D; \boldsymbol{\mu}, \Sigma) = -\frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K I(\mathbf{y}_i = C_k) \nabla_{\Sigma} \left( (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right).$$

Taking the gradient, we get

$$\nabla_{\Sigma} \ln p(D; \boldsymbol{\mu}, \Sigma) = -\frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K I(\mathbf{y}_i = C_k) \left( \Sigma^{-1} - \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \Sigma^{-1} \right).$$

6. To find the maximum likelihood estimator  $\widehat{\Sigma}$ , we set the gradient equal to zero

$$-\frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K I(\mathbf{y}_i = C_k) \left( \Sigma^{-1} - \Sigma^{-1} (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_k) (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_k)^\top \Sigma^{-1} \right) = 0.$$

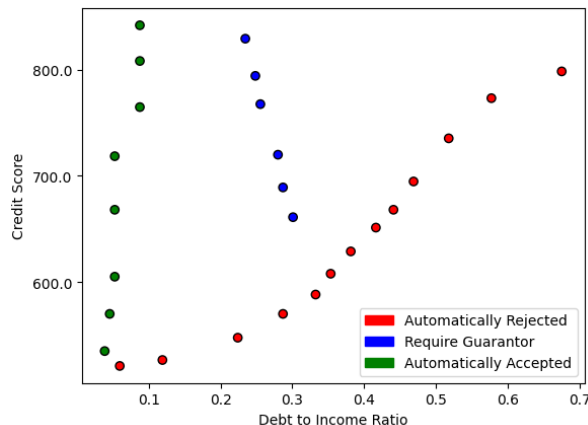
$$\widehat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K I(\mathbf{y}_i = C_k) (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_k) (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_k)^\top.$$





### Problem 3 (Classifying Loan Applicants)

In this problem, you will code up three different classifiers to classify different types of loan applicants. The file `data/hr.csv` contains data on debt to income ratio measured in tenths of a percent and credit score. The data can be plotted on these two axes:



Please implement the following classifiers in the `SoftmaxRegression` and `KNNClassifier` classes. For this problem, apply the following transformation to all data:

$$\phi(\mathbf{x}) = \left[ x_1 \cdot \frac{200}{7} - 7.5, \frac{x_2 - 500}{140} + 0.5 \right]^\top$$

where  $x_1$  and  $x_2$  represent the values for debt to income ratio and credit score, respectively. This transformation has been applied to the data for you in the notebook.

- A generative classifier with Gaussian class-conditional densities with a *shared covariance matrix*** across all classes. Feel free to re-use your Problem 2 results.
- Another generative classifier with Gaussian class-conditional densities , but now with a *separate covariance matrix*** learned for each class. (Note: The staff implementation can switch between the two Gaussian generative classifiers with just a few lines of code.)
- A multi-class logistic regression classifier** using the softmax activation function. In your implementation of gradient descent, **make sure to include a bias term and use L2 regularization** with regularization parameter  $\lambda = 0.001$ . Limit the number of iterations of gradient descent to 200,000, and set the learning rate to be  $\eta = 0.001$ .
- Another multi-class logistic regression classifier** with the additional feature map:

$$\phi(\mathbf{x}) = [\ln(x_1 + 10), x_2^2]^\top$$

where  $x_1$  and  $x_2$  represent the values for debt to income ratio and credit score, respectively.

- A kNN classifier** in which you classify based on the  $k = 1$  and  $k = 5$  nearest neighbors and the following distance function:

$$\text{dist}(\text{loan\_app}_1, \text{loan\_app}_2) = (\text{debt}_1 - \text{debt}_2)^2 / 9 + (\text{credit}_1 - \text{credit}_2)^2$$

where nearest neighbors are those with the smallest distances from a given point.

Note 1: When there are more than two labels, no label may have the majority of neighbors. Use the label that has the most votes among the neighbors as the choice of label.

Note 2: The grid of points for which you are making predictions should be interpreted as our test space. Thus, it is not necessary to make a test point that happens to be on top of a training point ignore itself when selecting neighbors.

### Problem 3 (cont.)

After implementing the above classifiers, complete the following exercises:

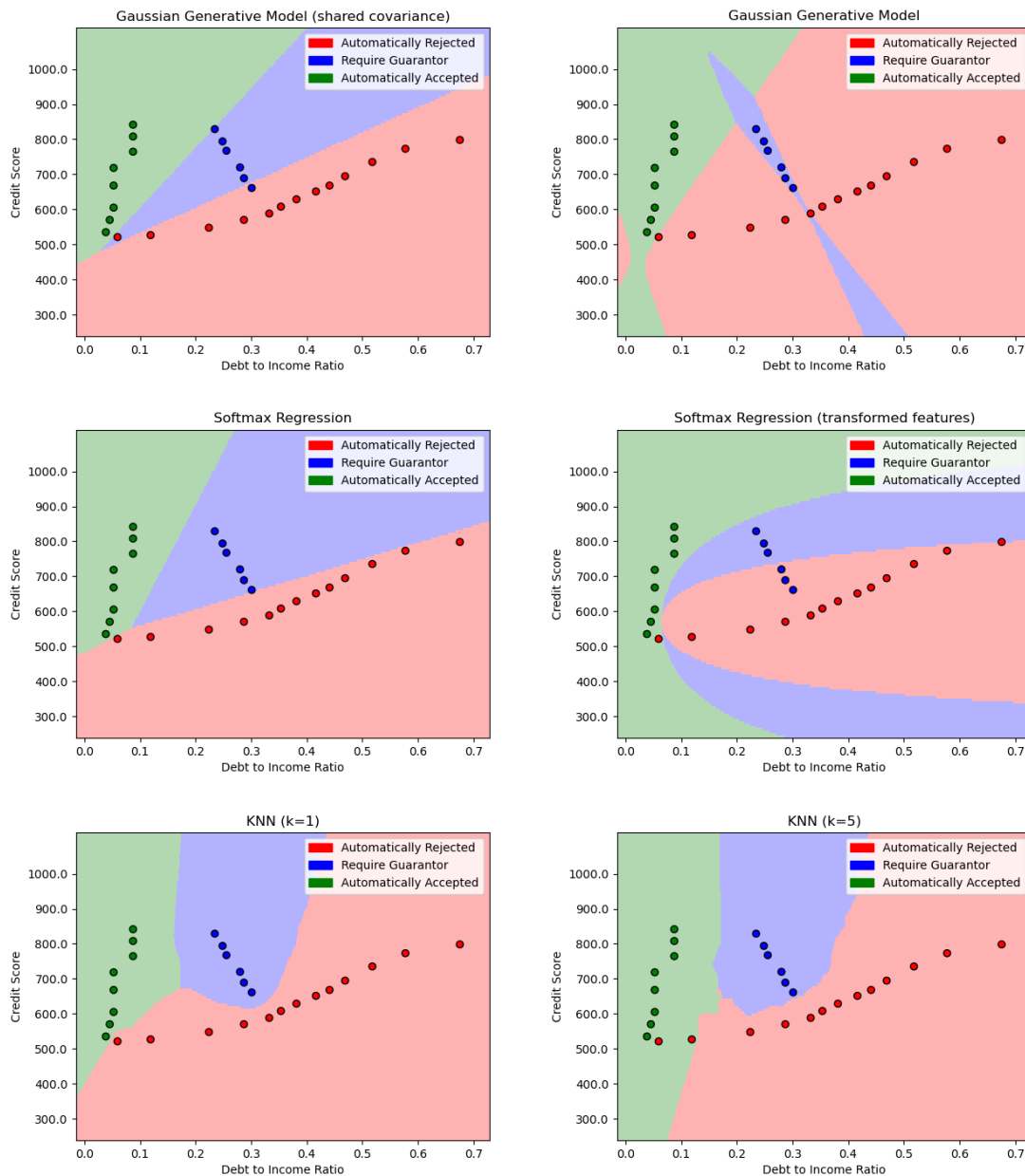
1. Plot the decision boundaries generated by each classifier for the dataset. Include them in your PDF. Identify the similarities and differences among the classifiers. What explains the differences—in particular, which aspects or properties of each model dictate the shape of its decision boundary?
2. Consider a loan applicant with Debt to Income Ratio 0.32 and Credit Score 350. To which class does each classifier assign this applicant? Report the classification probabilities of this applicant for models (c) and (d).

Interpret how each model makes its classification decision. What else should we, the modelers, be aware of when making predictions on a point “far” from our training data? **Your response should no be longer than 5 sentences.**

3. Can you think of any ethical problem that might arise from using this classifier to make loan decisions? You may approach this from any angle you like. For instance, can you think of someone who might have a low credit score and high debt-to-income ratio that you believe should nonetheless be offered a loan? Are there other variables that should be accounted for to ensure fair decisions? Are credit scores and debt-to-income ratio good bases for loan decisions? More generally, is using a classifier trained on past decisions to determine loan eligibility problematic in any way?

## Solution:

1) Implementation is in the Python notebook.



From top to bottom, note that when covariance is shared, our decision boundaries are linear. When covariance is separate for each class, the decision boundaries are quadratic.

The softmax regression without transformed features has linear boundaries, which makes sense because of the setup of softmax. However, our transformation basis is not linear, which results in nonlinear decision boundaries for the transformed features softmax regression.

The kNN classifier with  $k = 1$  and  $k = 5$  both have very interesting boundaries, but their shapes are

roughly the same. You can see traces of overfitting between the automatically rejected and automatically accepted regions, with every single point in our training data being classified correctly.

2. Implementation is in the Python Notebook: we have that every model classifies this individual as 0 (rejected) except for softmax regression with transformed features, which classifies the individual as 1 (requires guarantor).

The probabilities are as follows: Softmax regression:  $[1.00000000e+00, 9.93305871e-24, 4.05542149e-34]$

Softmax regression with transformed features:  $[2.16186498e-01, 7.83763958e-01, 4.95440650e-05]$

We can see that the softmax regression is very confident, but softmax regression with transformed features is only 80% confident. When a piece of data is far from the training data, we are making extrapolations that are not really backed by any confidence. This is especially true for the softmax regression with transformed features, which is not as confident as the other models (since it uses non-linear decision boundaries).

3. An ethical problem here is that a high debt-to-income ratio is very multi-dimensional. For example, someone who just graduated college might have a lot of debt and a relatively low income, yet they are more trustworthy than someone who might have no debt because they are always getting helped by others, yet also have no income so essentially no way of paying back the loan. A classifier with only two parameters cannot reasonably predict something with as many confounding variables as loan trustworthiness.

In addition to credit score + debt-to-income ratio, we can also consider things like net worth, current career average salary, as well as averages over the past few years to get a better idea of the applicant's current financial situation.

**Name:** Zihongbo Wang

**Collaborators and Resources:** Ossimi Ziv, Ryan Jiang, and the CS181 Textbook (as well as a youtube video about trolls 2 that clarified naive bayes)