

CS 229, Summer 2019

Problem Set #1 Solutions

YOUR NAME HERE (YOUR SUNET HERE)

Due Monday, July 15 at 11:59 pm on Gradescope.

Notes: (1) These questions require thought, but do not require long answers. Please be as concise as possible. (2) If you have a question about this homework, we encourage you to post your question on our Piazza forum, at <http://piazza.com/stanford/summer2019/cs229>. (3) If you missed the first lecture or are unfamiliar with the collaboration or honor code policy, please read the policy on the course website before starting work. (4) For the coding problems, you may not use any libraries except those defined in the provided `environment.yml` file. In particular, ML-specific libraries such as scikit-learn are not permitted. (5) To account for late days, the due date is Monday, July 15 at 11:59 pm. If you submit after Monday, July 15 at 11:59 pm, you will begin consuming your late days. If you wish to submit on time, submit before Monday, July 15 at 11:59 pm.

All students must submit an electronic PDF version of the written questions. We highly recommend typesetting your solutions via \LaTeX . All students must also submit a zip file of their source code to Gradescope, which should be created using the `make_zip.py` script. You should make sure to (1) restrict yourself to only using libraries included in the `environment.yml` file, and (2) make sure your code runs without errors. Your submission may be evaluated by the auto-grader using a private test set, or used for verifying the outputs reported in the writeup.

1. [40 points] Linear Classifiers (logistic regression and GDA)

In this problem, we cover two probabilistic linear classifiers we have covered in class so far. First, a discriminative linear classifier: logistic regression. Second, a generative linear classifier: Gaussian discriminant analysis (GDA). Both the algorithms find a linear decision boundary that separates the data into two classes, but make different assumptions. Our goal in this problem is to get a deeper understanding of the similarities and differences (and, strengths and weaknesses) of these two algorithms.

For this problem, we will consider two datasets, along with starter codes provided in the following files:

- `src/linearclass/ds1_{train,valid}.csv`
- `src/linearclass/ds2_{train,valid}.csv`
- `src/linearclass/logreg.py`
- `src/linearclass/gda.py`

Each file contains n examples, one example $(x^{(i)}, y^{(i)})$ per row. In particular, the i -th row contains columns $x_0^{(i)} \in \mathbb{R}$, $x_1^{(i)} \in \mathbb{R}$, and $y^{(i)} \in \{0, 1\}$. In the subproblems that follow, we will investigate using logistic regression and Gaussian discriminant analysis (GDA) to perform binary classification on these two datasets.

(a) [10 points]

In lecture we saw the average empirical loss for logistic regression:

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n \left(y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right),$$

where $y^{(i)} \in \{0, 1\}$, $h_{\theta}(x) = g(\theta^T x)$ and $g(z) = 1/(1 + e^{-z})$.

Find the Hessian H of this function, and show that for any vector z , it holds true that

$$z^T H z \geq 0.$$

Hint: You may want to start by showing that $\sum_i \sum_j z_i x_i x_j z_j = (x^T z)^2 \geq 0$. Recall also that $g'(z) = g(z)(1 - g(z))$.

Remark: This is one of the standard ways of showing that the matrix H is positive semi-definite, written “ $H \succeq 0$.” This implies that J is convex, and has no local minima other than the global one. If you have some other way of showing $H \succeq 0$, you’re also welcome to use your method instead of the one above.

Answer:

Note that $g'(z) = g(z)(1 - g(z))$, and thus for $h(x) = g(\theta^T x)$, we have $\partial h(x)/\partial \theta_k = h(x)(1 - h(x))x_k$. We will use the latter fact in the following derivations.

Letting $h_{\theta}(x^{(i)}) = g(\theta^T x^{(i)}) = 1/(1 + \exp(-\theta^T x^{(i)}))$, we have

$$\begin{aligned} \frac{\partial \log h_{\theta}(x^{(i)})}{\partial \theta_k} &= \frac{1}{h_{\theta}(x^{(i)})} \frac{\partial h_{\theta}(x^{(i)})}{\partial \theta_k} = \frac{1}{h_{\theta}(x^{(i)})} \frac{\partial h_{\theta}(x^{(i)})}{\partial \theta_k} = (1 - h_{\theta}(x^{(i)}))x_k^{(i)} \\ \frac{\partial \log(1 - h_{\theta}(x^{(i)}))}{\partial \theta_k} &= \frac{1}{1 - h_{\theta}(x^{(i)})} \frac{\partial(1 - h_{\theta}(x^{(i)}))}{\partial \theta_k} = -h_{\theta}(x^{(i)})x_k^{(i)}. \end{aligned}$$

Substituting into our equation for $J(\theta)$, we have

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_k} &= -\frac{1}{n} \sum_{i=1}^n y^{(i)}(1 - h_\theta(x^{(i)}))x_k^{(i)} - (1 - y^{(i)})h_\theta(x^{(i)})x_k^{(i)} \\ &= -\frac{1}{n} \sum_{i=1}^n (y^{(i)} - h_\theta(x^{(i)}))x_k^{(i)}.\end{aligned}$$

Consequently, the (k, l) entry of the Hessian is given by

$$\begin{aligned}H_{kl} &= \frac{\partial^2 J(\theta)}{\partial \theta_k \partial \theta_l} = -\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta_l} (y^{(i)} - h_\theta(x^{(i)}))x_k^{(i)} \\ &= \frac{1}{n} \sum_{i=1}^n h_\theta(x^{(i)})(1 - h_\theta(x^{(i)}))x_l^{(i)}x_k^{(i)}.\end{aligned}$$

Using the fact that $X = xx^T$ if and only if $X_{ij} = x_i x_j$, we have

$$H = \frac{1}{n} \sum_{i=1}^n h_\theta(x^{(i)})(1 - h_\theta(x^{(i)}))x^{(i)}x^{(i)T}.$$

To prove that H is positive semi-definite, we show $z^T H z \geq 0$ for all $z \in \mathbb{R}^d$. Letting $z \in \mathbb{R}^d$, we have

$$\begin{aligned}z^T H z &= \frac{1}{n} z^T \left(\sum_{i=1}^n h_\theta(x^{(i)})(1 - h_\theta(x^{(i)}))x^{(i)}x^{(i)T} \right) z \\ &= \frac{1}{n} \sum_{i=1}^n h_\theta(x^{(i)})(1 - h_\theta(x^{(i)}))(z^T x^{(i)})(x^{(i)T} z) \\ &= \frac{1}{n} \sum_{i=1}^n h_\theta(x^{(i)})(1 - h_\theta(x^{(i)}))(z^T x^{(i)})^2 \geq 0.\end{aligned}$$

The last inequality holds because $0 \leq h_\theta(x^{(i)}) \leq 1$ so $h_\theta(x^{(i)})(1 - h_\theta(x^{(i)})) \geq 0$, and trivially $(z^T x^{(i)})^2 \geq 0$. Therefore H is positive semi-definite.

- (b) [5 points] **Coding problem.** Follow the instructions in `src/linearclass/logreg.py` to train a logistic regression classifier using Newton's Method. Starting with $\theta = \vec{0}$, run Newton's Method until the updates to θ are small: Specifically, train until the first iteration k such that $\|\theta_k - \theta_{k-1}\|_1 < \epsilon$, where $\epsilon = 1 \times 10^{-5}$. Make sure to write your model's predicted probabilities on the validation set to the file specified in the code.

Include a plot of the **validation data** with x_1 on the horizontal axis and x_2 on the vertical axis. To visualize the two classes, use a different symbol for examples $x^{(i)}$ with $y^{(i)} = 0$ than for those with $y^{(i)} = 1$. On the same figure, plot the decision boundary found by logistic regression (i.e, line corresponding to $p(y|x) = 0.5$).

Answer:

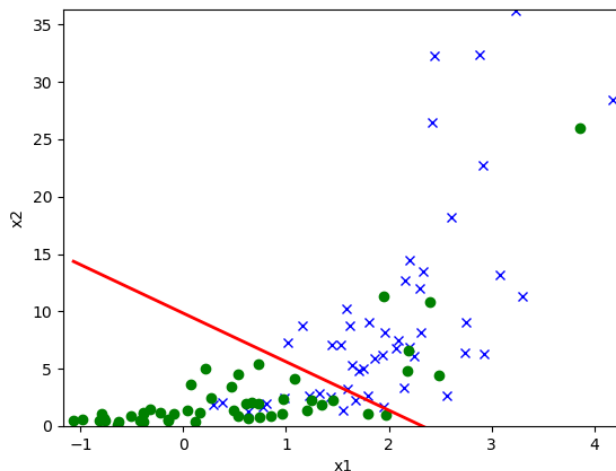


Figure 1: Separating hyperplane for logistic regression on Dataset 1

- (c) [5 points] Recall that in GDA we model the joint distribution of (x, y) by the following equations:

$$\begin{aligned}
 p(y) &= \begin{cases} \phi & \text{if } y = 1 \\ 1 - \phi & \text{if } y = 0 \end{cases} \\
 p(x|y=0) &= \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right) \\
 p(x|y=1) &= \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right),
 \end{aligned}$$

where ϕ , μ_0 , μ_1 , and Σ are the parameters of our model.

Suppose we have already fit ϕ , μ_0 , μ_1 , and Σ , and now want to predict y given a new point x . To show that GDA results in a classifier that has a linear decision boundary, show the posterior distribution can be written as

$$p(y=1 | x; \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(-(\theta^T x + \theta_0))},$$

where $\theta \in \mathbb{R}^d$ and $\theta_0 \in \mathbb{R}$ are appropriate functions of ϕ , Σ , μ_0 , and μ_1 .

Answer: For shorthand, we let $\mathcal{H} = \{\phi, \Sigma, \mu_0, \mu_1\}$ denote the parameters for the problem.

Since the given formulae are conditioned on y , use Bayes rule to get:

$$\begin{aligned}
 p(y = 1 | x; \phi, \Sigma, \mu_0, \mu_1) &= \frac{p(x|y = 1; \phi, \Sigma, \mu_0, \mu_1)p(y = 1; \phi, \Sigma, \mu_0, \mu_1)}{p(x; \phi, \Sigma, \mu_0, \mu_1)} \\
 &= \frac{p(x|y = 1; \mathcal{H})p(y = 1; \mathcal{H})}{p(x|y = 1; \mathcal{H})p(y = 1; \mathcal{H}) + p(x|y = 0; \mathcal{H})p(y = 0; \mathcal{H})} \\
 &= \frac{\exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right) \phi}{\exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right) \phi + \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right) (1 - \phi)} \\
 &= \frac{1}{1 + \frac{1-\phi}{\phi} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0) + \frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right)} \\
 &= \frac{1}{1 + \exp\left(\log\left(\frac{1-\phi}{\phi}\right) - \frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0) + \frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right)}.
 \end{aligned}$$

Now, we expand and rearrange the difference of quadratic terms in the preceding expression, finding that

$$\begin{aligned}
 (x - \mu_0)^T \Sigma^{-1}(x - \mu_0) - (x - \mu_1)^T \Sigma^{-1}(x - \mu_1) &= x^T \Sigma^{-1}x - \mu_0^T \Sigma^{-1}x - x^T \Sigma^{-1}\mu_0 + \mu_0^T \Sigma^{-1}\mu_0 - x^T \Sigma^{-1}x + \mu_1^T \Sigma^{-1}x + x^T \Sigma^{-1}\mu_1 - \mu_1^T \Sigma^{-1}\mu_1 \\
 &= -2\mu_0^T \Sigma^{-1}x + \mu_0^T \Sigma^{-1}\mu_0 + 2\mu_1^T \Sigma^{-1}x - \mu_1^T \Sigma^{-1}\mu_1 \\
 &= 2(\mu_1 - \mu_0)^T \Sigma^{-1}x + \mu_0^T \Sigma^{-1}\mu_0 - \mu_1^T \Sigma^{-1}\mu_1.
 \end{aligned}$$

Thus, we have

$$p(y = 1 | x; \mathcal{H}) = \frac{1}{1 + \exp\left(\log\left(\frac{1-\phi}{\phi}\right) + \frac{1}{2}\mu_1^T \Sigma^{-1}\mu_1 - \frac{1}{2}\mu_0^T \Sigma^{-1}\mu_0 + (\mu_0 - \mu_1)^T \Sigma^{-1}x\right)}.$$

and setting

$$\theta = -\Sigma^{-1}(\mu_0 - \mu_1) \quad \text{and} \quad \theta_0 = \frac{1}{2}(\mu_0^T \Sigma^{-1}\mu_0 - \mu_1^T \Sigma^{-1}\mu_1) - \log \frac{1-\phi}{\phi}$$

gives that

$$p(y = 1 | x; \phi, \Sigma, \mu_0, \mu_1) = \frac{1}{1 + \exp(-(\theta^T x + \theta_0))}.$$

- (d) [7 points] Given the dataset, we claim that the maximum likelihood estimates of the parameters are given by

$$\begin{aligned}
 \phi &= \frac{1}{n} \sum_{i=1}^n 1\{y^{(i)} = 1\} \\
 \mu_0 &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 0\}x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}} \\
 \mu_1 &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\}x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}} \\
 \Sigma &= \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T
 \end{aligned}$$

The log-likelihood of the data is

$$\begin{aligned}\ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^n p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi).\end{aligned}$$

By maximizing ℓ with respect to the four parameters, prove that the maximum likelihood estimates of ϕ , μ_0 , μ_1 , and Σ are indeed as given in the formulas above. (You may assume that there is at least one positive and one negative example, so that the denominators in the definitions of μ_0 and μ_1 above are non-zero.)

Answer:

First, derive the expression for the log-likelihood of the training data:

$$\begin{aligned}\ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi) \\ &= \sum_{i=1}^n \log p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^n \log p(y^{(i)}; \phi) \\ &\simeq \sum_{i=1}^n \left[\frac{1}{2} \log \frac{1}{|\Sigma|} - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) + \frac{y^{(i)}}{2} \log \phi + \frac{1 - y^{(i)}}{2} \log(1 - \phi) \right]\end{aligned}$$

where constant terms independent of the parameters have been ignored in the last expression.

Now, the likelihood is maximized by setting the derivative (or gradient) with respect to each of the parameters to zero.

$$\begin{aligned}\frac{\partial \ell}{\partial \phi} &= \sum_{i=1}^n \left[\frac{y^{(i)}}{2\phi} - \frac{1 - y^{(i)}}{2(1 - \phi)} \right] \\ &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\}}{\phi} - \frac{n - \sum_{i=1}^n 1\{y^{(i)} = 1\}}{1 - \phi}\end{aligned}$$

Setting this equal to zero and solving for ϕ gives the maximum likelihood estimate.

For μ_0 , take the gradient of the log-likelihood, and then use the same kinds of tricks as were used to analytically solve the linear regression problem.

$$\begin{aligned}\nabla_{\mu_0} \ell &= -\frac{1}{2} \sum_{i:y^{(i)}=0} \nabla_{\mu_0} (x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0) \\ &= -\frac{1}{2} \sum_{i:y^{(i)}=0} \nabla_{\mu_0} [\mu_0^T \Sigma^{-1} \mu_0 - x^{(i)T} \Sigma^{-1} \mu_0 - \mu_0^T \Sigma^{-1} x^{(i)}] \\ &= -\frac{1}{2} \sum_{i:y^{(i)}=0} \nabla_{\mu_0} \text{tr} [\mu_0^T \Sigma^{-1} \mu_0 - x^{(i)T} \Sigma^{-1} \mu_0 - \mu_0^T \Sigma^{-1} x^{(i)}] \\ &= -\frac{1}{2} \sum_{i:y^{(i)}=0} [2\Sigma^{-1} \mu_0 - 2\Sigma^{-1} x^{(i)}]\end{aligned}$$

The last step uses matrix calculus identities (specifically, those given in page 8 of the lecture notes), and also the fact that Σ (and thus Σ^{-1}) is symmetric.

Setting this gradient to zero gives the maximum likelihood estimate for μ_0 . The derivation for μ_1 is similar to the one above.

For Σ , we find the gradient with respect to $S = \Sigma^{-1}$ rather than Σ just to simplify the derivation (note that $|S| = \frac{1}{|\Sigma|}$). You should convince yourself that the maximum likelihood estimate S_n found in this way would correspond to the actual maximum likelihood estimate Σ_n as $S_n^{-1} = \Sigma_n$.

$$\begin{aligned}\nabla_S \ell &= \sum_{i=1}^n \nabla_S \left[\frac{1}{2} \log |S| - \frac{1}{2} \underbrace{(x^{(i)} - \mu_{y^{(i)}})^T}_{b_i^T} S \underbrace{(x^{(i)} - \mu_{y^{(i)}})}_{b_i} \right] \\ &= \sum_{i=1}^n \left[\frac{1}{2|S|} \nabla_S |S| - \frac{1}{2} \nabla_S b_i^T S b_i \right]\end{aligned}$$

But, we have the following identities:

$$\nabla_S |S| = |S| (S^{-1})^T$$

$$\nabla_S b_i^T S b_i = \nabla_{str} (b_i^T S b_i) = \nabla_{str} (S b_i b_i^T) = b_i b_i^T$$

In the above, we again used matrix calculus identities, and also the commutativity of the trace operator for square matrices. Putting these into the original equation, we get:

$$\begin{aligned}\nabla_S \ell &= \sum_{i=1}^n \left[\frac{1}{2} S^{-1} - \frac{1}{2} b_i b_i^T \right] \\ &= \frac{1}{2} \sum_{i=1}^n [\Sigma - b_i b_i^T]\end{aligned}$$

Setting this to zero gives the required maximum likelihood estimate for Σ .

- (e) [5 points] **Coding problem.** In `src/linearclass/gda.py`, fill in the code to calculate ϕ , μ_0 , μ_1 , and Σ , use these parameters to derive θ , and use the resulting GDA model to make predictions on the validation set. Make sure to write your model's predictions on the validation set to the file specified in the code.

Include a plot of the **validation data** with x_1 on the horizontal axis and x_2 on the vertical axis. To visualize the two classes, use a different symbol for examples $x^{(i)}$ with $y^{(i)} = 0$ than for those with $y^{(i)} = 1$. On the same figure, plot the decision boundary found by GDA (i.e, line corresponding to $p(y|x) = 0.5$).

Answer:

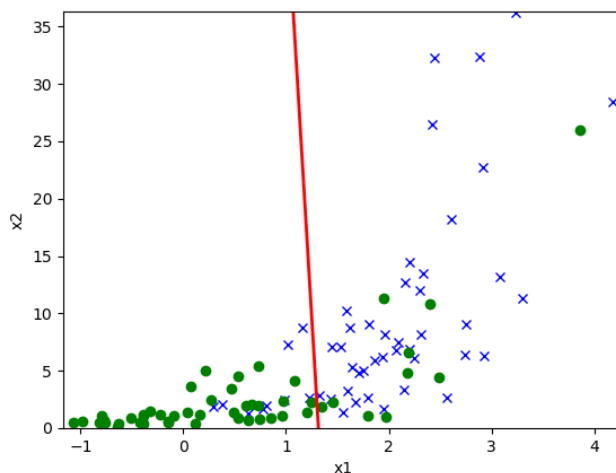


Figure 2: Separating hyperplane for GDA on Dataset 1

- (f) [2 points] For Dataset 1, compare the validation set plots obtained in part (b) and part (e) from logistic regression and GDA respectively, and briefly comment on your observation in a couple of lines.

Answer: The decision boundary with logistic regression appears more reasonable than GDA on Dataset 1. The accuracy is also higher with the logistic regression model on Dataset 1.

- (g) [5 points] Repeat the steps in part (b) and part (e) for Dataset 2. Create similar plots on the **validation set** of Dataset 2 and include those plots in your writeup.

On which dataset does GDA seem to perform worse than logistic regression? Why might this be the case?

Answer:

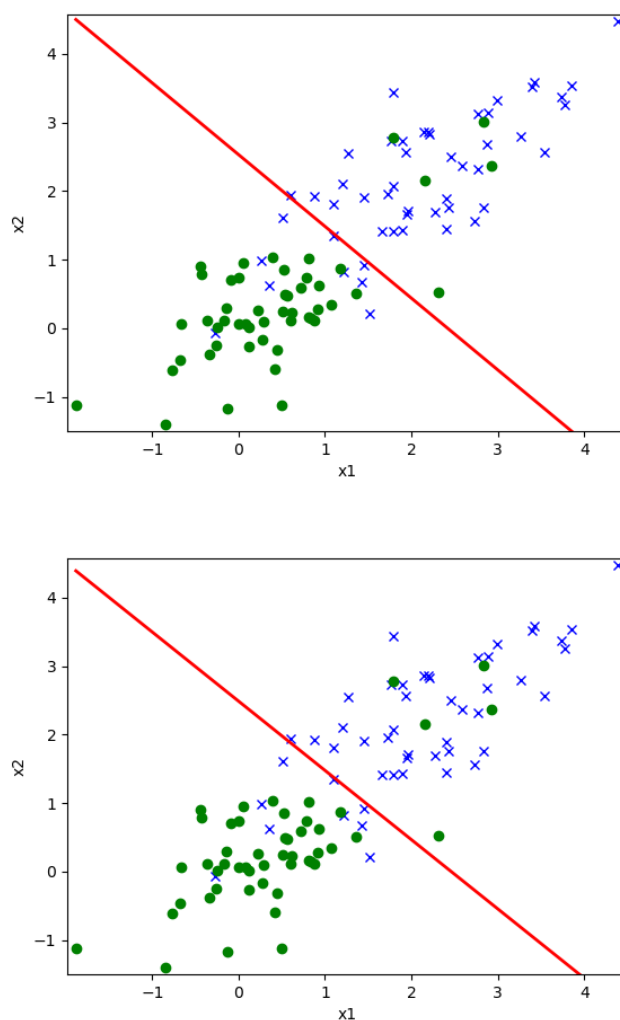


Figure 3: Separating hyperplanes for logistic regression (top) and GDA (bottom) on Dataset 2.

GDA seems to perform worse than logistic regression on Dataset 1. This is probably because the $x^{(i)}$'s are not Gaussian. We saw in the notes that logistic regression is more robust than GDA when the underlying dataset is not drawn from a multivariate Gaussian.

- (h) [1 points] For the dataset where GDA performed worse in parts (f) and (g), can you find a transformation of the $x^{(i)}$'s such that GDA performs significantly better? What might this transformation be?

Answer: If you transform dataset 2 by setting $x_2 := \log x_2$, then the $x^{(i)}$'s become Gaussian. Thus GDA does much better with this transformation.

2. [30 points] Incomplete, Positive-Only Labels

In this problem we will consider training binary classifiers in situations where we do not have full access to the labels. In particular, we consider a scenario, which is not too infrequent in real life, where we have labels only for a subset of the positive examples. All the negative examples and the rest of the positive examples are unlabelled.

We formalize the scenario as follows. Let $\{(x^{(i)}, t^{(i)})\}_{i=1}^n$ be a standard dataset of i.i.d distributed examples. Here $x^{(i)}$'s are the inputs/features and $t^{(i)}$ are the labels. Now consider the situation where $t^{(i)}$'s are not observed by us. Instead, we only observe the labels of some of the positive examples. Concretely, we assume that we observe $y^{(i)}$'s that are generated by

$$\begin{aligned}\forall x, \quad p(y^{(i)} = 1 \mid t^{(i)} = 1, x^{(i)} = x) &= \alpha, \\ \forall x, \quad p(y^{(i)} = 0 \mid t^{(i)} = 1, x^{(i)} = x) &= 1 - \alpha \\ \forall x, \quad p(y^{(i)} = 1 \mid t^{(i)} = 0, x^{(i)} = x) &= 0, \\ \forall x, \quad p(y^{(i)} = 0 \mid t^{(i)} = 0, x^{(i)} = x) &= 1\end{aligned}$$

where $\alpha \in (0, 1)$ is some unknown scalar. In other words, if the unobserved “true” label $t^{(i)}$ is 1, then with α chance we observe a label $y^{(i)} = 1$. On the other hand, if the unobserved “true” label $t^{(i)} = 0$, then we always observe the label $y^{(i)} = 0$.

Our final goal in the problem is to construct a binary classifier h of the true label t , with only access to the partial label y . In other words, we want to construct h such that $h(x^{(i)}) \approx p(t^{(i)} = 1 \mid x^{(i)})$ as closely as possible, using only x and y .

Real world example: Suppose we maintain a database of proteins which are involved in transmitting signals across membranes. Every example added to the database is involved in a signaling process, but there are many proteins involved in cross-membrane signaling which are missing from the database. It would be useful to train a classifier to identify proteins that should be added to the database. In our notation, each example $x^{(i)}$ corresponds to a protein, $y^{(i)} = 1$ if the protein is in the database and 0 otherwise, and $t^{(i)} = 1$ if the protein is involved in a cross-membrane signaling process and thus should be added to the database, and 0 otherwise.

For the rest of the question, we will use the dataset and starter code provided in the following files:

- `src/posonly/{train,valid,test}.csv`
- `src/posonly/posonly.py`

Each file contains the following columns: x_1 , x_2 , y , and t . As in Problem 1, there is one example per row. The $y^{(i)}$'s are generated from the process defined above with some unknown α .

(a) [5 points] Coding problem: ideal (fully observed) case

First we will consider the hypothetical (and uninteresting) case, where we have access to the true t -labels for training. In `src/posonly/posonly.py`, write a logistic regression classifier that uses x_1 and x_2 as input features, and train it using the t -labels. We will ignore the y -labels for this part. Output the trained model's predictions on the **test set** to the file specified in the code.

Create a plot to visualize the test set with x_1 on the horizontal axis and x_2 on the vertical axis. Use different symbols for examples $x^{(i)}$ with true label $t^{(i)} = 1$ than those with $t^{(i)} = 0$. On the same figure, plot the decision boundary obtained by your model (i.e, line

corresponding to model's predicted probability = 0.5) in red color. Include this plot in your writeup.

Answer:

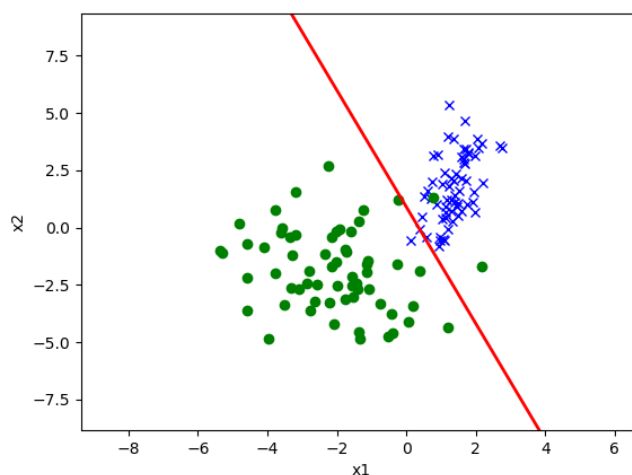


Figure 4: Separating hyperplane for logistic regression on t -labels

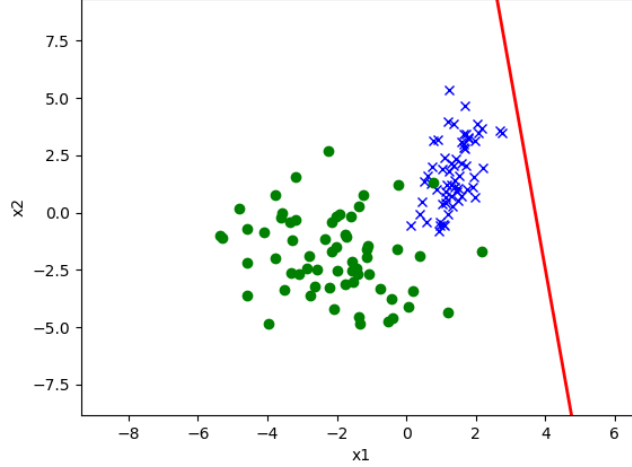
(b) [5 points] **Coding problem: The naive method on partial labels**

We now consider the case where the t -labels are unavailable, so you only have access to the y -labels at training time. Extend your code in `src/posonly/posonly.py` to re-train the classifier (still using x_1 and x_2 as input features), but using the y -labels only. Output the predictions on the **test set** to the appropriate file (as described in the code comments).

Create a plot to visualize the test set with x_1 on the horizontal axis and x_2 on the vertical axis. Use different symbols for examples $x^{(i)}$ with true label $t^{(i)} = 1$ (even though we only used the $y^{(i)}$ labels for training, use the true $t^{(i)}$ labels for plotting) than those with $t^{(i)} = 0$. On the same figure, plot the decision boundary obtained by your model (i.e, line corresponding to model's predicted probability = 0.5) in red color. Include this plot in your writeup.

Note that the algorithm should learn a function $h(\cdot)$ that approximately predicts the probability $p(y^{(i)} = 1 \mid x^{(i)})$. Also note that we expect it to perform poorly on predicting the probability of interest, namely $p(t^{(i)} = 1 \mid x^{(i)})$.

Answer:

Figure 5: Separating hyperplane for logistic regression on y -labels naively trained

In the following sub-questions we will attempt to solve the problem with only partial observations. That is, we only have access to $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$, and will try to predict $p(t^{(i)} = 1 | x^{(i)})$.

(c) [5 points] **Warm-up with Bayes rule**

Show that under our assumptions, for any i ,

$$p(t^{(i)} = 1 | y^{(i)} = 1, x^{(i)}) = 1 \quad (1)$$

That is, observing a positive partial label $y^{(i)} = 1$ tells us for sure the hidden true label is 1. Use Bayes rule to derive this (an informal explanation will not earn credit).

Answer:

$$\begin{aligned}
 p(t^{(i)} = 1 | y^{(i)} = 1, x^{(i)}) &= \frac{p(y^{(i)} = 1 | t^{(i)} = 1, x^{(i)})p(t^{(i)} = 1 | x^{(i)})}{p(y^{(i)} = 1 | x^{(i)})} \\
 &= \frac{p(y^{(i)} = 1 | t^{(i)} = 1, x^{(i)})p(t^{(i)} = 1 | x^{(i)})}{p(y^{(i)} = 1 | t^{(i)} = 0, x^{(i)})p(t^{(i)} = 0 | x^{(i)}) + p(y^{(i)} = 1 | t^{(i)} = 1, x^{(i)})p(t^{(i)} = 1 | x^{(i)})} \\
 &= \frac{p(y^{(i)} = 1 | t^{(i)} = 1, x^{(i)})p(t^{(i)} = 1 | x^{(i)})}{0 \cdot p(t^{(i)} = 0 | x^{(i)}) + p(y^{(i)} = 1 | t^{(i)} = 1, x^{(i)})p(t^{(i)} = 1 | x^{(i)})} \\
 &= \frac{p(y^{(i)} = 1 | t^{(i)} = 1, x^{(i)})p(t^{(i)} = 1 | x^{(i)})}{p(y^{(i)} = 1 | t^{(i)} = 1, x^{(i)})p(t^{(i)} = 1 | x^{(i)})} \\
 &= 1
 \end{aligned}$$

- (d) [5 points] Show that for any example, the probability that true label $t^{(i)}$ is positive is $1/\alpha$ times the probability that the partial label is positive. That is, show that

$$p(t^{(i)} = 1 | x^{(i)}) = \frac{1}{\alpha} \cdot p(y^{(i)} = 1 | x^{(i)}) \quad (2)$$

Note that the equation above suggests that if we know the value of α , then we can convert a function $h(\cdot)$ that approximately predicts the probability $h(x^{(i)}) \approx p(y^{(i)} = 1 | x^{(i)})$ into a function that approximately predicts $p(t^{(i)} = 1 | x^{(i)})$ by multiplying the factor $1/\alpha$.

Answer:

We have

$$\begin{aligned} p(y^{(i)} = 1 | x^{(i)}) &= p(y^{(i)} = 1 \wedge t^{(i)} = 1 | x^{(i)}) \quad (\text{from previous sub-question}) \\ &= p(t^{(i)} = 1 | x^{(i)})p(y^{(i)} = 1 | t^{(i)} = 1, x^{(i)}) \\ &= p(t^{(i)} = 1 | x^{(i)})\alpha. \end{aligned}$$

Dividing both sides by α , we have the desired result.

(e) [5 points] **Estimating α**

The solution to estimate $p(t^{(i)} | x^{(i)})$ outlined in the previous sub-question requires the knowledge of α which we don't have. Now we will design a way to estimate α based on the function $h(\cdot)$ that approximately predicts $p(y^{(i)} = 1 | x^{(i)})$ (which we obtained in part b).

To simplify the analysis, let's assume that we have magically obtained a function $h(x)$ that perfectly predicts the value of $p(y^{(i)} = 1 | x^{(i)})$, that is, $h(x^{(i)}) = p(y^{(i)} = 1 | x^{(i)})$.

We make the crucial assumption that $p(t^{(i)} = 1 | x^{(i)}) \in \{0, 1\}$. This assumption means that the process of generating the "true" label $t^{(i)}$ is a noise-free process. This assumption is not very unreasonable to make. Note, we are NOT assuming that the observed label $y^{(i)}$ is noise-free, which would be an unreasonable assumption!

Now show that

$$\alpha = \mathbb{E}[h(x^{(i)}) | y^{(i)} = 1] \quad (3)$$

The above result motivates the following algorithm to estimate α by estimating the RHS of the equation above using samples: Let V_+ be the set of labeled (and hence positive) examples in the validation set V , given by $V_+ = \{x^{(i)} \in V | y^{(i)} = 1\}$.

Then we use

$$\alpha \approx \frac{1}{|V_+|} \sum_{x^{(i)} \in V_+} h(x^{(i)}).$$

to estimate α . (You will be asked to implement this algorithm in the next sub-question. For this sub-question, you only need to show equation (3). Moreover, this sub-question may be slightly harder than other sub-questions.)

Answer:

Solution 1: This is a solution that is easier to understand but it strongly relies on the assumption that we made here. The Solution 2 below is less intuitive but can work for more general settings than the setting in this sub-question.

We fix i , and let $A = \{x : p(t^{(i)} = 1 | x^{(i)} = x) = 1\}$ be the set of all "positive examples" and $B = \{x : p(t^{(i)} = 0 | x^{(i)} = x) = 1\}$ be the set of all negative examples. By the assumption $p(t^{(i)} = 1 | x^{(i)}) \in \{0, 1\}$, we know that A and B covers the whole input space. We first claim that $p(x^{(i)} \notin A, y^{(i)} = 1) = 0$. This is because when $x^{(i)} \notin A$, we have that $x^{(i)} \in B$ which implies that $p(t^{(i)} = 0 | x^{(i)}) = 1$, which implies that $p(t^{(i)} = 1 | x^{(i)}) = 0$. Therefore, the event $x^{(i)} \notin A$ implies the event $t^{(i)} = 0$, which implies that $y^{(i)} = 0$. Thus $p(x^{(i)} \notin A, y^{(i)} = 1) = 0$. It follows that $p(x^{(i)} \notin A | y^{(i)} = 1) = 0$, which implies that $p(x^{(i)} \in A | y^{(i)} = 1) = 1$.

By 2d, for $x \in A$, we have that $h(x) = \alpha p(t^{(i)} = 1 \mid x^{(i)} = x) = \alpha$. That is, $p(h(x^{(i)}) = \alpha \mid x^{(i)} = x) = 1$. This implies that $p(h(x^{(i)}) = \alpha \mid x^{(i)} = x, y^{(i)} = 1) = 1$. (Note that for any event E , and any event F with positive probability, $p(E) = 1$ implies that $p(E \mid F) = 1$. This is because $p(E^c) = 0$ which implies $p(E^c \mid F) = p(E^c \cap F)/p(F) = 0$.)

$$\begin{aligned}
 p(h(x^{(i)}) = \alpha \mid y^{(i)} = 1) &= p(h(x^{(i)}) = \alpha \mid y^{(i)} = 1, x^{(i)} \in A) \cdot p(x^{(i)} \in A \mid y^{(i)} = 1) \quad (4) \\
 &= 1 \cdot p(x^{(i)} \in A \mid y^{(i)} = 1) \\
 &\quad \text{(by } p(h(x^{(i)}) = \alpha \mid x^{(i)} = x, y^{(i)} = 1) = 1) \\
 &= 1 \cdot 1 \quad \text{(because } p(x^{(i)} \in A \mid y^{(i)} = 1) = 1) \\
 &= 1 \quad (5)
 \end{aligned}$$

Therefore $\mathbb{E}[h(x^{(i)}) \mid y^{(i)} = 1] = \alpha$.

Solution 2

Let $g(x^{(i)}) = p(t^{(i)} = 1 \mid x^{(i)})$. By d), we have that

$$h(x^{(i)}) = \alpha g(x^{(i)})$$

. Then, we have that

$$\mathbb{E}[h(x^{(i)}) \mid y^{(i)} = 1] = \mathbb{E}[h(x^{(i)}) \mid y^{(i)} = 1] \quad (6)$$

$$= \frac{\mathbb{E}[h(x^{(i)}) \mathbf{1}\{y^{(i)} = 1\}]}{p(y^{(i)} = 1)} \quad \text{(by def of conditional expectation)}$$

$$= \frac{\mathbb{E}[\mathbb{E}[\mathbf{1}\{y^{(i)} = 1\} \mid x^{(i)}] \cdot h(x^{(i)})]}{\mathbb{E}[p(y^{(i)} = 1 \mid x^{(i)})]} \quad \text{((by law of expectation))}$$

$$= \frac{\mathbb{E}[h(x^{(i)}) \cdot h(x^{(i)})]}{\mathbb{E}[h(x^{(i)})]} \quad (7)$$

$$= \alpha \frac{\mathbb{E}[g(x^{(i)}) \cdot g(x^{(i)})]}{\mathbb{E}[g(x^{(i)})]} \quad (8)$$

$$= \alpha \quad \text{(since } g(x^{(i)}) \in \{0, 1\})$$

(f) [5 points] **Coding problem.**

Using the validation set, estimate the constant α by averaging your classifier's predictions over all labeled examples in the validation set:¹

$$\alpha \approx \frac{1}{|V_+|} \sum_{x^{(i)} \in V_+} h(x^{(i)}).$$

Add code in `src/posonly/posonly.py` to rescale your predictions $h(y^{(i)} = 1 \mid x^{(i)})$ of the classifier that is obtained from part b, using the equation (2) obtained in part (d) and using the estimated value for α .

Finally, create a plot to visualize the test set with x_1 on the horizontal axis and x_2 on the vertical axis. Use different symbols for examples $x^{(i)}$ with true label $t^{(i)} = 1$ (even though

¹There is a reason to use the validation set, instead of the training set, to estimate the α . However, for the purpose of this question, we sweep the subtlety here under the rug, and you don't need to understand the difference between the two for this question.

we only used the $y^{(i)}$ labels for training, use the true $t^{(i)}$ labels for plotting) than those with $t^{(i)} = 0$. On the same figure, plot the decision boundary obtained by your model (i.e, line corresponding to model's **adjusted** predicted probability = 0.5) in red color. Include this plot in your writeup.

Answer:

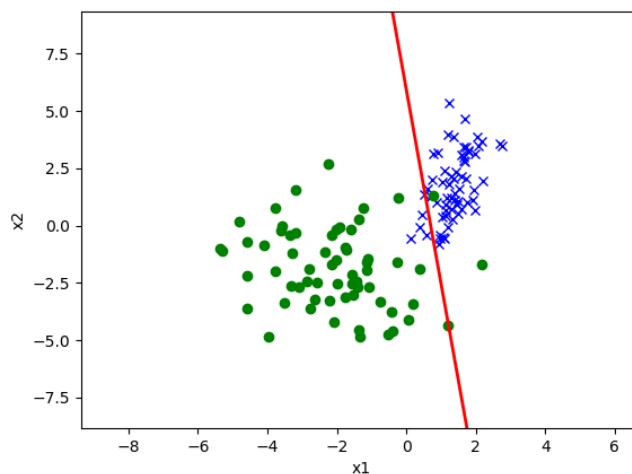


Figure 6: Separating hyperplane for logistic regression on y -labels with correction.

Remark: We saw that the true probability $p(t \mid x)$ was only a constant factor away from $p(y \mid x)$. This means, if our task is to only rank examples (*i.e.* sort them) in a particular order (e.g, sort the proteins in order of being most likely to be involved in transmitting signals across membranes), then in fact we do not even need to estimate α . The rank based on $p(y \mid x)$ will agree with the rank based on $p(t \mid x)$.

3. [25 points] Poisson Regression

In this question we will construct another kind of a commonly used GLM, which is called Poisson Regression. In a GLM, the choice of the exponential family distribution is based on the kind of problem at hand. If we are solving a classification problem, then we use an exponential family distribution with support over discrete classes (such as Bernoulli, or Categorical). Similarly, if the output is real valued, we can use Gaussian or Laplace (both are in the exponential family). Sometimes the desired output is to predict counts, for e.g., predicting the number of emails expected in a day, or the number of customers expected to enter a store in the next hour, etc. based on input features (also called covariates). You may recall that a probability distribution with support over integers (i.e. counts) is the Poisson distribution, and it also happens to be in the exponential family.

In the following sub-problems, we will start by showing that the Poisson distribution is in the exponential family, derive the functional form of the hypothesis, derive the update rules for training models, and finally using the provided dataset train a real model and make predictions on the test set.

- (a) [5 points] Consider the Poisson distribution parameterized by λ :

$$p(y; \lambda) = \frac{e^{-\lambda} \lambda^y}{y!}.$$

(Here y has positive integer values and $y!$ is the factorial of y .) Show that the Poisson distribution is in the exponential family, and clearly state the values for $b(y)$, η , $T(y)$, and $a(\eta)$.

Answer: Rewrite the distribution function as:

$$\begin{aligned} p(y; \lambda) &= \frac{e^{-\lambda} e^{y \log \lambda}}{y!} \\ &= \frac{1}{y!} \exp(y \log \lambda - \lambda) \end{aligned}$$

Comparing with the standard form for the exponential family:

$$\begin{aligned} b(y) &= \frac{1}{y!} \\ \eta &= \log \lambda \\ T(y) &= y \\ a(\eta) &= e^\eta. \end{aligned}$$

- (b) [3 points] Consider performing regression using a GLM model with a Poisson response variable. What is the canonical response function for the family? (You may use the fact that a Poisson random variable with parameter λ has mean λ .)

Answer: The canonical response function for the GLM model will be:

$$\begin{aligned} g(\eta) &= E[y; \eta] \\ &= \lambda \\ &= e^\eta. \end{aligned}$$

- (c) [7 points] For a training set $\{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$, let the log-likelihood of an example be $\log p(y^{(i)}|x^{(i)}; \theta)$. By taking the derivative of the log-likelihood with respect to θ_j , derive the stochastic gradient ascent update rule for learning using a GLM model with Poisson responses y and the canonical response function.

Answer: The log-likelihood of an example $(x^{(i)}, y^{(i)})$ is defined as $\ell(\theta) = \log p(y^{(i)}|x^{(i)}; \theta)$. To derive the stochastic gradient ascent rule, use the results in part (a) and the standard GLM assumption that $\eta = \theta^T x$.

$$\begin{aligned}
 \frac{\partial \ell(\theta)}{\partial \theta_j} &= \frac{\partial \log p(y^{(i)}|x^{(i)}; \theta)}{\partial \theta_j} \\
 &= \frac{\partial \log \left(\frac{1}{y^{(i)}!} \exp(\eta^T y^{(i)} - e^\eta) \right)}{\partial \theta_j} \\
 &= \frac{\partial \log \left(\exp((\theta^T x^{(i)})^T y^{(i)} - e^{\theta^T x^{(i)}}) \right)}{\partial \theta_j} + \frac{\partial \log \left(\frac{1}{y^{(i)}!} \right)}{\partial \theta_j} \\
 &= \frac{\partial \left((\theta^T x^{(i)})^T y^{(i)} - e^{\theta^T x^{(i)}} \right)}{\partial \theta_j} \\
 &= \frac{\partial \left((\sum_k \theta_k x_k^{(i)}) y^{(i)} - e^{\sum_k \theta_k x_k^{(i)}} \right)}{\partial \theta_j} \\
 &= x_j^{(i)} y^{(i)} - e^{\sum_k \theta_k x_k^{(i)}} x_j^{(i)} \\
 &= (y^{(i)} - e^{\theta^T x^{(i)}}) x_j^{(i)}.
 \end{aligned}$$

Thus the stochastic gradient ascent update rule should be:

$$\theta_j := \theta_j + \alpha \frac{\partial \ell(\theta)}{\partial \theta_j},$$

which reduces here to:

$$\theta_j := \theta_j + \alpha (y^{(i)} - e^{\theta^T x^{(i)}}) x_j^{(i)}.$$

- (d) [10 points] **Coding problem**

Consider a website that wants to predict its daily traffic. The website owners have collected a dataset of past traffic to their website, along with some features which they think are useful in predicting the number of visitors per day. The dataset is split into train/valid sets and the starter code is provided in the following files:

- `src/poisson/{train,valid}.csv`
- `src/poisson/poisson.py`

We will apply Poisson regression to model the number of visitors per day. Note that applying Poisson regression in particular assumes that the data follows a Poisson distribution whose natural parameter is a linear combination of the input features (*i.e.*, $\eta = \theta^T x$). In `src/poisson/poisson.py`, implement Poisson regression for this dataset and use *full batch gradient ascent* to maximize the log-likelihood of θ . For the stopping criterion, check if the change in parameters has a norm smaller than a small value such as 10^{-5} .

Using the trained model, predict the expected counts for the **validation set**, and create a scatter plot between the true counts vs predicted counts (on the validation set). In the scatter plot, let x-axis be the true count and y-axis be the corresponding predicted expected count. Note that the true counts are integers while the expected counts are generally real values.

Answer:

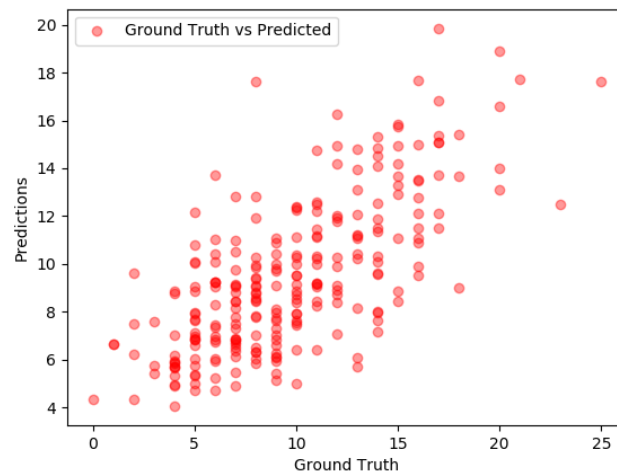


Figure 7: Ground Truth vs Prediction plot on the validation set

4. [15 points] Convexity of Generalized Linear Models

In this question we will explore and show some nice properties of Generalized Linear Models, specifically those related to its use of Exponential Family distributions to model the output.

Most commonly, GLMs are trained by using the negative log-likelihood (NLL) as the loss function. This is mathematically equivalent to Maximum Likelihood Estimation (*i.e.*, maximizing the log-likelihood is equivalent to minimizing the negative log-likelihood). In this problem, our goal is to show that the NLL loss of a GLM is a *convex* function w.r.t the model parameters. As a reminder, this is convenient because a convex function is one for which any local minimum is also a global minimum, and there is extensive research on how to optimize various types of convex functions efficiently with various algorithms such as gradient descent or stochastic gradient descent.

To recap, an exponential family distribution is one whose probability density can be represented

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta)),$$

where η is the *natural parameter* of the distribution. Moreover, in a Generalized Linear Model, η is modeled as $\theta^T x$, where $x \in \mathbb{R}^d$ are the input features of the example, and $\theta \in \mathbb{R}^d$ are learnable parameters. In order to show that the NLL loss is convex for GLMs, we break down the process into sub-parts, and approach them one at a time. Our approach is to show that the second derivative (*i.e.*, Hessian) of the loss w.r.t the model parameters is Positive Semi-Definite (PSD) at all values of the model parameters. We will also show some nice properties of Exponential Family distributions as intermediate steps.

For the sake of convenience we restrict ourselves to the case where η is a scalar. Assume $p(Y|X; \theta) \sim \text{ExponentialFamily}(\eta)$, where $\eta \in \mathbb{R}$ is a scalar, and $T(y) = y$. This makes the exponential family representation take the form

$$p(y; \eta) = b(y) \exp(\eta y - a(\eta)).$$

- (a) [5 points] Derive an expression for the mean of the distribution. Show that $\mathbb{E}[Y; \eta] = \frac{\partial}{\partial \eta} a(\eta)$ (note that $\mathbb{E}[Y; \eta] = \mathbb{E}[Y|X; \theta]$ since $\eta = \theta^T x$). In other words, show that the mean of an exponential family distribution is the first derivative of the log-partition function with respect to the natural parameter.

Hint: Start with observing that $\frac{\partial}{\partial \eta} \int p(y; \eta) dy = \int \frac{\partial}{\partial \eta} p(y; \eta) dy$.

Answer: The expectation of an exponential family distribution is first derivative of the log-partition function $a(\eta)$ w.r.t the natural parameters η .

$$\begin{aligned} \frac{\partial}{\partial \eta} p(y; \eta) &= p(y; \eta) \left(y - \frac{\partial}{\partial \eta} a(\eta) \right) \\ \Rightarrow \int \frac{\partial}{\partial \eta} p(y; \eta) dy &= \int p(y; \eta) \left(y - \frac{\partial}{\partial \eta} a(\eta) \right) dy \\ \Rightarrow \frac{\partial}{\partial \eta} \int p(y; \eta) dy &= \int y p(y; \eta) dy - \int p(y; \eta) \frac{\partial}{\partial \eta} a(\eta) dy \\ &= \frac{\partial}{\partial \eta} 1 = \mathbb{E}[Y; \eta] - \frac{\partial}{\partial \eta} a(\eta) \int p(y; \eta) dy \\ &\Rightarrow \mathbb{E}[Y; \eta] = \frac{\partial}{\partial \eta} a(\eta) \end{aligned}$$

- (b) [5 points] Next, derive an expression for the variance of the distribution. In particular, show that $\text{Var}(Y; \eta) = \frac{\partial^2}{\partial \eta^2} a(\eta)$ (again, note that $\text{Var}(Y; \eta) = \text{Var}(Y|X; \theta)$). In other words, show that the variance of an exponential family distribution is the second derivative of the log-partition function w.r.t. the natural parameter.

Hint: Building upon the result in the previous sub-problem can simplify the derivation.

Answer:

$$\begin{aligned}
 \frac{\partial^2}{\partial \eta^2} a(\eta) &= \frac{\partial}{\partial \eta} \mathbb{E}[Y|X; \theta] \\
 &= \frac{\partial}{\partial \eta} \int yp(y; \eta) dy \\
 &= \int yp(y; \eta) \left(y - \frac{\partial}{\partial \eta} a(\eta) \right) dy \\
 &= \int y^2 p(y; \eta) dy - \frac{\partial}{\partial \eta} a(\eta) \int yp(y; \eta) dy \\
 &= \mathbb{E}[Y^2; \eta] - \mathbb{E}[Y; \eta]^2 \\
 &= \text{Var}(Y; \eta).
 \end{aligned}$$

- (c) [5 points] Finally, write out the loss function $\ell(\theta)$, the NLL of the distribution, as a function of θ . Then, calculate the Hessian of the loss w.r.t θ , and show that it is always PSD. This concludes the proof that NLL loss of GLM is convex.

Hint 1: Use the chain rule of calculus along with the results of the previous parts to simplify your derivations.

Hint 2: Recall that variance of any probability distribution is non-negative.

Answer:

$$\begin{aligned}
 \ell(\theta) &= -\log [p(y; \eta)] \\
 &= -\log [b(y) \exp(\eta y - a(\eta))] \\
 &= a(\eta) - \eta y + C \\
 &= a(\theta^T x) - \theta^T xy + C
 \end{aligned}$$

Gradient of the loss w.r.t. θ is

$$\begin{aligned}
 \nabla_{\theta} \ell(\theta) &= \frac{\partial}{\partial \eta} a(\eta) \nabla_{\theta} \eta - yx \\
 &= \frac{\partial}{\partial \eta} a(\eta) x - yx
 \end{aligned}$$

Hessian of the loss w.r.t. θ is

$$\begin{aligned}
 \nabla_{\theta}^2 \ell(\theta) &= \nabla_{\theta} (\nabla_{\theta} \ell(\theta)) \\
 &= \nabla_{\theta} \left(\frac{\partial}{\partial \eta} a(\eta) x - y x \right) \\
 &= x \nabla_{\theta} \left(\frac{\partial}{\partial \eta} a(\eta) \right) \\
 &= x \frac{\partial}{\partial \eta} \left(\frac{\partial}{\partial \eta} a(\eta) \right) \nabla_{\theta} \eta \\
 &= \frac{\partial^2}{\partial \eta^2} a(\eta) x x^T \\
 &= \text{Var}(Y; \eta) x x^T
 \end{aligned}$$

Observe that $x x^T$ is always PSD because for any vector $z \in \mathbb{R}^d$ we have $z^T (x x^T) z = (z^T x)(x^T z) = (z^T x)^2 \geq 0$. Also, for any value of θ , $\text{Var}(Y; \eta)$ is always positive. Therefore the Hessian $\nabla_{\theta}^2 \ell(\theta) = \text{Var}(Y; \eta) x x^T$ of GLM's NLL loss is PSD, and hence convex. (Note that just saying elements of H is non-negative is not sufficient in showing it is a PSD)

Remark: The main takeaways from this problem are:

- Any GLM model is convex in its model parameters.
- The exponential family of probability distributions are mathematically nice. Whereas calculating mean and variance of distributions in general involves integrals (hard), surprisingly we can calculate them using derivatives (easy) for exponential family.

5. [25 points] **Linear regression: linear in what?**

In the first two lectures, you have seen how to fit a linear function of the data for the regression problem. In this question, we will see how linear regression can be used to fit non-linear functions of the data using feature maps. We will also explore some of its limitations, for which future lectures will discuss fixes.

(a) [5 points] **Learning degree-3 polynomials of the input**

Suppose we have a dataset $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$ where $x^{(i)}, y^{(i)} \in \mathbb{R}$. We would like to fit a third degree polynomial $h_\theta(x) = \theta_3 x^3 + \theta_2 x^2 + \theta_1 x^1 + \theta_0$ to the dataset. The key observation here is that the function $h_\theta(x)$ is still linear in the unknown parameter θ , even though it's not linear in the input x . This allows us to convert the problem into a linear regression problem as follows.

Let $\phi : \mathbb{R} \rightarrow \mathbb{R}^4$ be a function that transforms the original input x to a 4-dimensional vector defined as

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} \in \mathbb{R}^4 \quad (9)$$

Let $\hat{x} \in \mathbb{R}^4$ be a shorthand for $\phi(x)$, and let $\hat{x}^{(i)} \triangleq \phi(x^{(i)})$ be the transformed input in the training dataset. We construct a new dataset $\{(\phi(x^{(i)}), y^{(i)})\}_{i=1}^n = \{(\hat{x}^{(i)}, y^{(i)})\}_{i=1}^n$ by replacing the original inputs $x^{(i)}$'s by $\hat{x}^{(i)}$'s. We see that fitting $h_\theta(x) = \theta_3 x^3 + \theta_2 x^2 + \theta_1 x^1 + \theta_0$ to the old dataset is equivalent to fitting a linear function $h_\theta(\hat{x}) = \theta_3 \hat{x}_3 + \theta_2 \hat{x}_2 + \theta_1 \hat{x}_1 + \theta_0$ to the new dataset because

$$h_\theta(x) = \theta_3 x^3 + \theta_2 x^2 + \theta_1 x^1 + \theta_0 = \theta_3 \phi(x)_3 + \theta_2 \phi(x)_2 + \theta_1 \phi(x)_1 + \theta_0 = \theta^T \hat{x} \quad (10)$$

In other words, we can use linear regression on the new dataset to find parameters $\theta_0, \dots, \theta_3$. Please write down 1) the objective function $J(\theta)$ of the linear regression problem on the new dataset $\{(\hat{x}^{(i)}, y^{(i)})\}_{i=1}^n$ and 2) the update rule of the batch gradient descent algorithm for linear regression on the dataset $\{(\hat{x}^{(i)}, y^{(i)})\}_{i=1}^n$.

Terminology: In machine learning, ϕ is often called the feature map which maps the original input x to a new set of variables. To distinguish between these two sets of variables, we will call x the input **attributes**, and call $\phi(x)$ the **features**. (Unfortunately, different authors use different terms to describe these two things. In this course, we will do our best to follow the above convention consistently.)

Answer:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)})^2$$

Differentiating this objective, we get the update rule:

$$\begin{aligned} \nabla_\theta J(\theta) &= \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)}) (\hat{x}^{(i)}) \\ &\propto \lambda \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)}) (\hat{x}^{(i)}) \end{aligned}$$

Where λ is the learning rate. The update rule is:

$$\theta := \theta - \lambda \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)}) (\hat{x}^{(i)})$$

(b) [5 points] **Coding question: degree-3 polynomial regression**

For this sub-question question, we will use the dataset provided in the following files:

`src/featuremaps/{train,valid,test}.csv`

Each file contains two columns: x and y . In the terminology described in the introduction, x is the attribute (in this case one dimensional) and y is the output label.

Using the formulation of the previous sub-question, implement linear regression with **normal equations** using the feature map of degree-3 polynomials. Use the starter code provided in `src/featuremaps/featuremap.py` to implement the algorithm.

Create a scatter plot of the training data, and plot the learnt hypothesis as a smooth curve over it. Submit the plot in the writeup as the solution for this problem.

Remark: Suppose \hat{X} is the design matrix of the transformed dataset. You may sometimes encounter a non-invertible matrix $\hat{X}^T \hat{X}$. For a numerically stable code implementation, always use `np.linalg.solve` to obtain the parameters directly, rather than explicitly calculating the inverse and then multiplying it with $\hat{X}^T y$.

Answer:

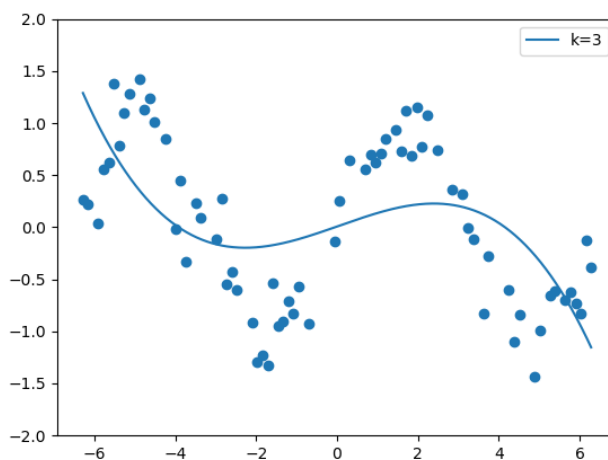


Figure 8: Polynomial regression with degree 3

(c) [5 points] **Coding question: degree- k polynomial regression**

Now we extend the idea above to degree- k polynomials by considering $\phi : \mathbb{R} \rightarrow \mathbb{R}^{k+1}$ to be

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^k \end{bmatrix} \in \mathbb{R}^{k+1} \quad (11)$$

Follow the same procedure as the previous sub-question, and implement the algorithm with $k = 3, 5, 10, 20$. Create a similar plot as in the previous sub-question, and include the hypothesis curves for each value of k with a different color. Include a legend in the plot to indicate which color is for which value of k .

Submit the plot in the writeup as the solution for this sub-problem. Observe how the fitting of the training dataset changes as k increases. Briefly comment on your observations in the plot.

Answer:

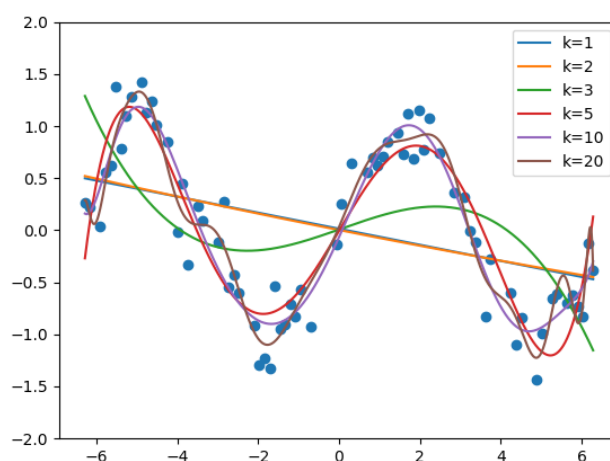


Figure 9: Polynomial regression with kernel sizes 1,2,3,5,10 and 20

Generally higher degree polynomials tend to fit the data better, though very high degree polynomials can be numerically unstable.

(d) [5 points] **Coding question: other feature maps**

You may have observed that it requires a relatively high degree k to fit the given training data, and this is because the dataset cannot be explained (i.e., approximated) very well by low-degree polynomials. By visualizing the data, you may have realized that y can be approximated well by a sine wave. In fact, we generated the data by sampling from $y = \sin(x) + \xi$, where ξ is noise with Gaussian distribution. Please update the feature map ϕ to include a sine transformation as follows:

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^k \\ \sin(x) \end{bmatrix} \in \mathbb{R}^{k+2} \quad (12)$$

With the updated feature map, train different models for values of $k = 0, 1, 2, 3, 5, 10, 20$, and plot the resulting hypothesis curves over the data as before.

Submit the plot as a solution to this sub-problem. Compare the fitted models with the previous sub-question, and briefly comment about noticeable differences in the fit with this feature map.

Answer:

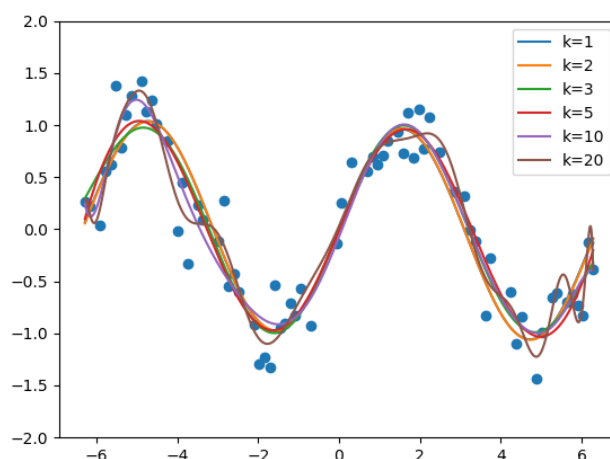


Figure 10: Polynomial regression with other features with kernel sizes 1,2,3,5,10 and 20

In the presence of the $\sin(x)$ feature, the models seem to fit the data better / more robustly. However the numerical instability that comes with high degree polynomials remain.

(e) [5 points] **Overfitting with expressive models and small data**

For the rest of the problem, we will consider a small dataset (a random subset of the dataset you have been using so far) with much fewer examples, provided in the following file:

`src/featuremaps/small.csv`

We will be exploring what happens when the number of features start becoming bigger than the number of examples in the training set. Run your algorithm on this small dataset

using the following feature map

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^k \end{bmatrix} \in \mathbb{R}^{k+1} \quad (13)$$

with $k = 1, 2, 5, 10, 20$.

Create a plot of the various hypothesis curves (just like previous sub-questions). Observe how the fitting of the training dataset changes as k increases. Submit the plot in the writeup and comment on what you observe.

Remark: The phenomenon you observe where the models start to fit the training dataset very well, but suddenly “goes wild” is due to what is called *overfitting*. The intuition to have for now is that, when the amount of data you have is small relative to the expressive capacity of the family of possible models (that is, the hypothesis class, which, in this case, is the family of all degree k polynomials), it results in overfitting.

Loosely speaking, the set of hypothesis function is “very flexible” and can be easily forced to pass through all your data points especially in unnatural ways. In other words, the model explains the noises in the training dataset, which shouldn’t be explained in the first place. This hurts the predictive power of the model on test examples. We will describe overfitting in more detail in future lectures when we cover learning theory and bias-variance tradeoffs.

Answer:

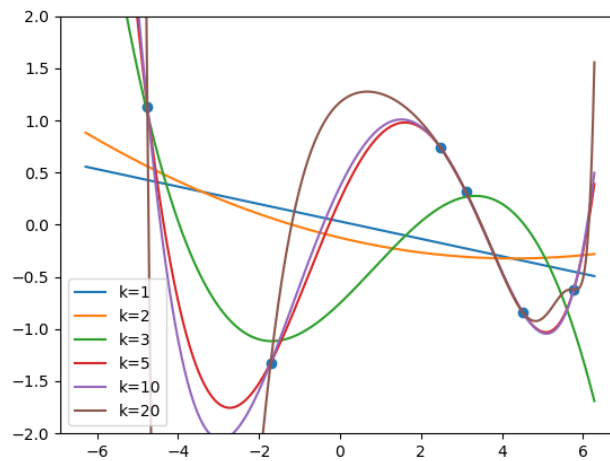


Figure 11: Polynomial regression with kernel sizes 1,2,3,5,10 and 20 on small dataset

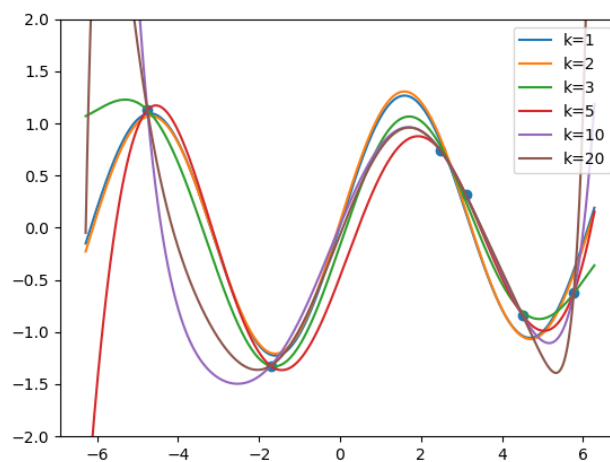


Figure 12: Regression with other features with kernel sizes 1,2,3,5,10 and 20 on small dataset

We see that when the dataset is small, higher degree polynomials tend to pass through all the points, but qualitatively seem like a poor fit. Numerical instability with high degree polynomials remain a problem even with small data, with or without $\sin(x)$.