

Homework 4: Probabilistic models

Due: Friday, March 25, 2022 at 11:59PM EST

Instructions: Your answers to the questions below, including plots and mathematical work, should be submitted as a single PDF file. It's preferred that you write your answers using software that typesets mathematics (e.g. LaTeX, LyX, or MathJax via iPython), though if you need to you may scan handwritten work. You may find the minted package convenient for including source code in your LaTeX document. If you are using LyX, then the listings package tends to work better.

1 Logistic Regression

Consider a binary classification setting with input space $\mathcal{X} = \mathbb{R}^d$, outcome space $\mathcal{Y}_{\pm} = \{-1, 1\}$, and a dataset $\mathcal{D} = ((x^{(1)}, y^{(1)}), \text{Dots}, (x^{(n)}, y^{(n)}))$.

Equivalence of ERM and probabilistic approaches

In the lecture we derived logistic regression using the Bernoulli response distribution. In this problem you will show that it is equivalent to ERM with logistic loss.

ERM with logistic loss.

Consider a linear scoring function in the space $\mathcal{F}_{\text{score}} = \{x \mapsto x^T w \mid w \in \mathbb{R}^d\}$. A simple way to make predictions (similar to what we've seen with the perceptron algorithm) is to predict $\hat{y} = 1$ if $x^T w > 0$, or $\hat{y} = \text{sign}(x^T w)$. Accordingly, we consider margin-based loss functions that relate the loss with the margin, $y x^T w$. A positive margin means that $x^T w$ has the same sign as y , i.e. a correct prediction. Specifically, let's consider the **logistic loss** function $\ell_{\text{logistic}}(y, w) = \log(1 + \exp(-y w^T x))$. This is a margin-based loss function that you have now encountered several times. Given the logistic loss, we can now minimize the empirical risk on our dataset \mathcal{D} to obtain an estimate of the parameters, \hat{w} .

MLE with a Bernoulli response distribution and the logistic link function.

As discussed in the lecture, given that $p(y = 1 \mid x; w) = 1/(1 + \exp(-x^T w))$, we can estimate w by maximizing the likelihood, or equivalently, minimizing the negative log-likelihood (NLL $_{\mathcal{D}}(w)$ in short) of the data.

1. Show that the two approaches are equivalent, i.e. they will produce the same solution for w .

Linearly Separable Data

In this problem, we will investigate the behavior of MLE for logistic regression when the data is linearly separable.

2. Show that the decision boundary of logistic regression is given by $\{x: x^T w = 0\}$. Note that the set will not change if we multiply the weights by some constant c .
3. Suppose the data is linearly separable and by gradient descent/ascent we have reached a decision boundary defined by \hat{w} where all examples are classified correctly. Show that we can always increase the likelihood of the data by multiplying a scalar c on \hat{w} , which means that MLE is not well-defined in this case. (Hint: You can show this by taking the derivative of $L(c\hat{w})$ with respect to c , where L is the likelihood function.)

Regularized Logistic Regression

As we've shown in above, when the data is linearly separable, MLE for logistic regression may end up with weights with very large magnitudes. Such a function is prone to overfitting. In this part, we will apply regularization to fix the problem.

The ℓ_2 regularized logistic regression objective function can be defined as

$$\begin{aligned} J_{\text{logistic}}(w) &= \hat{R}_n(w) + \lambda \|w\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp \left(-y^{(i)} w^T x^{(i)} \right) \right) + \lambda \|w\|^2. \end{aligned}$$

4. Prove that the objective function $J_{\text{logistic}}(w)$ is convex. You may use any facts mentioned in the convex optimization notes.
5. Complete the `f_objective` function in the skeleton code, which computes the objective function for $J_{\text{logistic}}(w)$. (Hint: you may get numerical overflow when computing the exponential literally, e.g. try e^{1000} in Numpy. Make sure to read about the log-sum-exp trick and use the numpy function `logaddexp` to get accurate calculations and to prevent overflow.
6. Complete the `fit_logistic_regression_function` in the skeleton code using the `minimize` function from `scipy.optimize`. Use this function to train a model on the provided data. Make sure to take the appropriate preprocessing steps, such as standardizing the data and adding a column for the bias term.
7. Find the ℓ_2 regularization parameter that minimizes the log-likelihood on the validation set. Plot the log-likelihood for different values of the regularization parameter.
8. [Optional] It seems reasonable to interpret the prediction $f(x) = \phi(w^T x) = 1/(1 + e^{-w^T x})$ as the probability that $y = 1$, for a randomly drawn pair (x, y) . Since we only have a finite sample (and we are regularizing, which will bias things a bit) there is a question of how well “calibrated” our predicted probabilities are. Roughly speaking, we say $f(x)$ is well calibrated if we look at all examples (x, y) for which $f(x) \approx 0.7$ and we find that close to 70% of those examples have $y = 1$, as predicted... and then we repeat that for all predicted probabilities in $(0, 1)$. To see how well-calibrated our predicted probabilities are, break the predictions on the validation set into groups based on the predicted probability (you can play with the size of the groups to get a result you think is informative). For each group, examine the percentage of positive labels. You can make a table or graph. Summarize the results. You may get some ideas and references from scikit-learn’s discussion.

2 Coin Flipping with Partial Observability

Consider flipping a biased coin where $p(z = H \mid \theta_1) = \theta_1$. However, we cannot directly observe the result z . Instead, someone reports the result to us, which we denote by x . Further, there is a chance that the result is reported incorrectly *if it's a head*. Specifically, we have $p(x = H \mid z = H, \theta_2) = \theta_2$ and $p(x = T \mid z = T) = 1$.

9. Show that $p(x = H \mid \theta_1, \theta_2) = \theta_1 \theta_2$.
10. Given a set of reported results \mathcal{D}_r of size N_r , where the number of heads is n_h and the number of tails is n_t , what is the likelihood of \mathcal{D}_r as a function of θ_1 and θ_2 .

11. Can we estimate θ_1 and θ_2 using MLE? Explain your judgment.