
DS-GA 1003: Machine Learning (Spring 2020)

Homework 5: Probabilistic models

Due: Wednesday, April 15, 2020 at 11:59pm

In this homework we'll be investigating conditional probability models, with a focus on various interpretations of logistic regression, with and without regularization. Along the way we'll discuss the calibration of probability predictions, both in the limit of infinite training data and in a more bare-hands way. On the Bayesian side, we'll recreate from scratch the Bayesian linear gaussian regression example we discussed in lecture. We'll also have several optional problems that work through many basic concepts in Bayesian statistics via one of the simplest problems there is: estimating the probability of heads in a coin flip. Later we'll extend this to the probability of estimating click-through rates in mobile advertising. Along the way we'll encounter empirical Bayes and hierarchical models.

Instructions. You should upload your code and plots to Gradescope. Please map the Gradescope entry on the rubric to your name and NetId. You must follow the policies for submission detailed in Homework 0.

1 Logistic Regression

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

1.1 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 $\mathcal{F}_{\text{score}} = \{x \mapsto x^T w \mid w \in \mathbf{R}^d\}$. 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 $\ell_{\text{logistic}}(y, w) = \log(1 + \exp(-y w^T x))$. You should see that it is a margin-based loss function and an upper bound of the 0-1 loss. Given the logistic loss, we can now minimize the empirical risk on our dataset 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) of the data.

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

1.2 Linearly Separable Data

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

1. 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 .
2. 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 increase the likelihood of the data by increasing a scalar c on \hat{w} unboundedly, 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.]

1.3 Regularized Logistic Regression

As we've shown in Section 1.2, when the data is linearly separable, MLE for logistic regression may end up with very large weights, which is a sign of overfitting. In this part, we will apply regularization to fix the problem.

The ℓ_2 regularized logistic regression objective function can be defined as

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

1. Prove that the objective function $J_{\text{logistic}}(w)$ is convex. You may use any facts mentioned in the [convex optimization notes](#).
2. 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.
3. 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.
4. 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.
5. [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 Bayesian Logistic Regression with Gaussian Priors

Let’s continue with logistic regression in the Bayesian setting, where we introduce a prior $p(w)$ on $w \in \mathbf{R}^d$.

1. For the dataset \mathcal{D} described in Section 1, give an expression for the posterior density $p(w | \mathcal{D})$ in terms of the negative log-likelihood function $\text{NLL}_{\mathcal{D}}(w)$ and the prior density $p(w)$ (up to a proportionality constant is fine).
2. Suppose we take a prior on w of the form $w \sim \mathcal{N}(0, \Sigma)$. Is this a conjugate prior to the likelihood given by logistic regression?
3. Find a covariance matrix Σ such that MAP estimate for w after observing data \mathcal{D} is the same as the minimizer of the regularized logistic regression function defined in Section 1.3 (and prove it). [Hint: Consider minimizing the negative log posterior of w . Also, remember you can drop any terms from the objective function that don’t depend on w . You may freely use results of previous problems.]
4. In the Bayesian approach, the prior should reflect your beliefs about the parameters before seeing the data and, in particular, should be independent on the eventual size of your dataset. Following this, you choose a prior distribution $w \sim \mathcal{N}(0, I)$. For a dataset \mathcal{D} of size n , how should you choose λ in our regularized logistic regression objective function so that the minimizer is equal to the mode of the posterior distribution of w (i.e. is equal to the MAP estimator).

3 Coin Flipping with Partial Observability

Consider flipping a biased coin where $p(z = H | \theta_1) = \theta_1$. 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 $p(x = H | z = H, \theta_2) = \theta_2$ and $p(x = T | z = T) = 1$.

1. Show that $p(x = H | \theta_1, \theta_2) = \theta_1 \theta_2$.
2. 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 . Can we estimate θ_1 and θ_2 using MLE? Explain your judgment.

3. We additionally obtained a set of clean results \mathcal{D}_c of size N_c , where x is directly observed without the reporter in the middle. Given that there are c_h heads and c_t tails, estimate θ_1 and θ_2 by MLE. Feel free to directly apply previous results. Note that the likelihood is $L(\theta_1, \theta_2) = p(\mathcal{D}_r, \mathcal{D}_c \mid \theta_1, \theta_2)$.
4. Since the clean results are expensive, we only have a small number of those. Let's put a prior distribution on θ_1 : $\text{Beta}(h, t)$. Derive the MAP estimates for θ_1 and θ_2 . Feel free to directly apply previous results.