

**Imperial College
London**

COURSEWORK

PROBABILISTIC INFERENCE (CO-493)

Logistic Regression and MCMC

Logistic Regression

The goal of this coursework is to familiarise yourself with Logistic Regression. We will also see two different approaches to solve the regression task.

For simplicity we will only consider the binary classification case, in which target variables are $y \in \{0, 1\}$.

In logistic regression, the probability of a data point x being of class 1 is given by

$$p(y = 1 | x, \theta) = \sigma(x^\top \theta),$$

where, $\sigma(z) = 1/(1 + \exp(-z))$ is the *sigmoid* function.

Combining this with a Bernoulli likelihood and summing over all data points $\{x_i, y_i\}_{i=1}^N$ we end up with a negative log-likelihood function that looks like this:

$$-\log p(y|X, \theta) = - \sum_i (y_i \log \sigma(x_i^\top \theta) + (1 - y_i) \log(1 - \sigma(x_i^\top \theta)))$$

You will see this expression in many other classification problems, especially in deep learning, where it is known as the *cross-entropy loss*.

Your goal in this tutorial is to learn how to perform inference over the parameters θ in logistic regression, including point estimates θ_{ML} and θ_{MAP} and approximations to the posterior $p(\theta|X, y)$.

Provided to you are **logistic_regression_assignment.py**, a skeleton file in which you will provide solutions, and **logistic_regression_extra.ipynb**, a jupyter notebook that can be used for visualisation purposes. The notebook provides plots of the regression and might help the understanding.

Submit your final version to CATE via the LabTS system.

Task 1: 10 marks

Complete the definition of the function **predict**. This function takes data points X and a parameter vector θ and returns $p(y = 1 | x, \theta)$ for $x \in X$.

Task 2: 10 marks

Complete the definition of the function **log_likelihood**. This function takes data points X , labels y and a parameter vector θ and returns $\log p(y|X, \theta)$.

Task 3: 10 marks

The likelihood function does not have a closed form solution and similar to GP coursework we will use numerical optimisation to fit parameters.

In the previous coursework you saw how to derive the gradients and optimise likelihood with them. For this coursework you can do it that way, or any other way you want. The optimisation is convex, so numerical difference [default in scipy-optimizers] should also work and give you same answer. With correct gradients you should converge in 1 or 2 iterations.

Complete the function **max_lik_estimate** to optimise the log-likelihood function you've written above and obtain θ_{ML} .

Bayesian logistic regression

To move on to Bayesian inference on the parameters θ , we put a prior on the parameters.

More specifically, we use a Gaussian prior parametrised by a mean m and a covariance S :

$$\theta \sim \mathcal{N}(m, S)$$

Similarly to θ_{ML} , θ_{MAP} does not have an analytical solution and we are required to use numerical optimisation methods.

Maximum A Posteriori (MAP)

First, we build a MAP estimate of m and S

Task 4: 10 marks

Complete the definition of the function **neg_log_posterior**. This function takes data points X , labels y , prior mean m and covariance S as well as a parameter vector θ and returns the negative log posterior.

Task 5: 10 marks

Complete the function **map_estimate** to optimise the negative log posterior you've written above and obtain θ_{MAP} .

The Laplace approximation

As we have hinted above, in logistic regression the posterior distribution over θ doesn't have an analytical solution. This is an example of *approximate Bayesian inference*: The exact posterior is analytically intractable so that we have to approximate it using one of various techniques. The one we'll use in this part of the coursework is called the "Laplace approximation".

In brief, the Laplace approximation is a Gaussian centred at the peak of the pdf of interest with the same curvature. Let us make this a bit more rigorous below.

Let us say we have a probability distribution $p(\mathbf{z})$ we want to approximate. The distribution $p(\mathbf{z})$ is of the form

$$p(\mathbf{z}) = \frac{1}{Z} \tilde{p}(\mathbf{z}) ,$$

where $\tilde{p}(\mathbf{z})$ is an unnormalised distribution that we can evaluate easily, but Z is unknown. Formally, the Laplace approximation results from a second-order Taylor expansion of $\log \tilde{p}(\mathbf{z})$ around \mathbf{z}_0 :

$$\log \tilde{p}(\mathbf{z}) \approx \log \tilde{p}(\mathbf{z}_0) + \frac{d}{d\mathbf{z}} \log \tilde{p}(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} (\mathbf{z} - \mathbf{z}_0) + \frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^\top \frac{d^2}{d\mathbf{z}^2} \log \tilde{p}(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} (\mathbf{z} - \mathbf{z}_0)$$

Now let us evaluate this expression at the mode of $p(\mathbf{z})$, which is the same as the mode of $\tilde{p}(\mathbf{z})$. We define the mode \mathbf{z}^* such that

$$\frac{d}{d\mathbf{z}} \tilde{p}(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}^*} = \mathbf{0} .$$

At this point, the $\mathcal{O}(\mathbf{z})$ term of the expansion vanishes and we are left with

$$\log \tilde{p}(\mathbf{z}) \approx \log \tilde{p}(\mathbf{z}^*) - \frac{1}{2} (\mathbf{z} - \mathbf{z}^*)^\top \mathbf{A} (\mathbf{z} - \mathbf{z}^*)$$

Or, equivalently,

$$\tilde{p}(\mathbf{z}) \approx \tilde{p}(\mathbf{z}^*) \exp \left(-\frac{1}{2} (\mathbf{z} - \mathbf{z}^*)^\top \mathbf{A} (\mathbf{z} - \mathbf{z}^*) \right) ,$$

where

$$\mathbf{A} = -\frac{d^2}{d\mathbf{z}^2} \log \tilde{p}(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}^*} .$$

Because this distribution is Gaussian, we can easily calculate the normalisation constant. By inspection, we can identify the mean and the covariance, and write down the Laplace approximation of $p(\mathbf{z})$ as

$$q(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{z}^*, \mathbf{A}^{-1})$$

Task 6: 10 marks

As an example, let us use the unnormalized distribution $\tilde{p}(\mathbf{z}) = x e^{-x/2}$. When normalized properly, this is in fact the χ^2 distribution with $k = 4$ degrees of freedom.

Complete the function **laplace_q** that takes an input vector \mathbf{z} and returns the value of the Laplace approximation $q(\mathbf{z})$ of $\tilde{p}(\mathbf{z}) = x e^{-x/2}$.

Task 7: 20 marks

So far, we have obtained the mode (peak) of the posterior through the MAP estimate above. Here we are extending this to a posterior distribution over θ . However, as we mentioned above the posterior doesn't have an analytical form, so we will use the **Laplace transform**.

Complete the function **get_posterior**, based on your previous code, that will calculate the Laplace approximation $q(\theta)$ of the true posterior $p(\theta|X, y)$ and return the mean and variance of q .

The jupyter notebook offers code to visualise the behaviour and the diversity of q by drawing a number $j = 1, \dots, J$ of samples $\theta_j \sim q(\theta)$. For each sample it plots its predicted class probabilities $\sigma(x\theta_j)$.

Hint: The extension of the Laplace approximation to multivariate distributions is straightforward, and in this case the variance of the Gaussian is the inverse Hessian of the negative log posterior $A = -\nabla_{\theta} \nabla_{\theta} \log p(\theta|X, y)$.

MCMC Sampling

The Laplace approximation is part of a family of methods known as *deterministic approximate inference*. In addition, there's another set of methods known as *stochastic approximate inference* which, includes most of the sampling techniques mentioned in the course.

Here we will implement the Metropolis algorithm that uses a proposal distribution $q(x'|x^{(t)})$ that depends on the last sample $x^{(t)}$. We then perform the Metropolis algorithm:

1. Generate proposal $x' \sim q(x'|x^{(t)})$

2. If

$$\frac{q(x^{(t)}|x')\tilde{p}(x')}{q(x'|x^{(t)})\tilde{p}(x^{(t)})} \geq u, \quad u \sim U[0, 1]$$

accept the sample $x^{(t+1)} = x'$. Otherwise set $x^{(t+1)} = x^{(t)}$.

Task 8: 20 marks

Complete the function **metropolis_hastings_sample** that returns *nb_iter* samples of the true posterior $p(\theta | X, y)$ with prior $\theta \sim \mathcal{N}(m, S)$ using the Metropolis algorithm.