

Lecture 3c: Logistic Regression for Binary Classification Problems

Lecturer: Jeffrey Varner

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications.*

The key concepts covered in this lecture include:

- **Logistic regression** is a statistical method used for binary classification that models the relationship between a dependent categorical variable (label) and one or more independent variables (features) by estimating probabilities through the logistic function.
- **Maximum likelihood estimation (MLE)** is a statistical technique to estimate the parameters of a probability distribution by maximizing the likelihood function, thereby determining the parameter values that make the observed data most probable.
- **Gradient descent** is an optimization algorithm used to minimize a function by iteratively adjusting parameters in the opposite direction of the gradient. Iteration continues until a local minimum of the function is found.
- **Alternatives to gradient descent** include heuristic optimization algorithms such as the Nelder-Mead Simplex Algorithm, Simulated Annealing, Genetic Algorithms, and Particle Swarm Optimization, which can estimate model parameters without relying on the gradient.

1 Introduction

In this lecture, we will introduce logistic regression for binary classification problems. Logistic regression is a statistical method used for binary classification that models the relationship between a dependent categorical variable (label) and one or more independent variables (features) by estimating probabilities through the logistic function. We will start by discussing the logistic model and then explore the maximum likelihood estimation of the model parameters. Unfortunately, the likelihood function of the logistic regression model does not have a closed-form analytical solution, so we must estimate the model parameters using numerical optimization algorithms. To address this challenge, we will introduce the gradient descent algorithm. This numerical optimization algorithm minimizes a function by iteratively adjusting the parameters in the opposite direction of the gradient. In addition to gradient descent, we will discuss alternative heuristic optimization algorithms that can estimate the model parameters without relying on the gradient. These alternative methods might be considered when the objective function is non-differentiable or has many local minima. Let's start by discussing the logistic regression model.

2 Logistic Regression

Logistic regression is a statistical method used for binary classification problems, where the dependent variable (label) is a binary categorical variable (e.g., ± 1 , etc), and the independent variables (features) are continuous or categorical. Unlike the Perceptron model, which outputs the class label directly, logistic regression models the probability that a given input belongs to a particular class based on the input features. The logistic regression model estimates the probability that a given feature vector belongs to a particular class based on the input features. In particular, logistic regression is a discriminative model, which means

it directly models the conditional probability of the label given the features, i.e., $P(y|\mathbf{x})$. This is in contrast to generative models, e.g., Naive Bayes, which we'll explore later, which model the joint probability of the features and the label, i.e., $P(y, \mathbf{x}) = P(\mathbf{x}|y) \cdot P(y)$. The logistic regression model uses the logistic function to model the probability of the binary label $y \in \{-1, +1\}$ given the (augmented) feature vector $\hat{\mathbf{x}} = (x_1, x_2, \dots, x_m, 1)$:

$$P_\theta(y|\hat{\mathbf{x}}) = \frac{1}{1 + e^{-y \cdot (\hat{\mathbf{x}}^\top \theta)}} \quad (1)$$

where $\theta \in \mathbb{R}^p$ ($p = m + 1$) is an (unknown) parameter vector (that we need to estimate somehow), and e is the base of the natural logarithm. The logistic function is a sigmoid function that maps the input, i.e., $-y \cdot (\hat{\mathbf{x}}^\top \theta)$ to the range $[0, 1]$. Thus, it is suitable for modeling probabilities. The logistic model *parameterizes* a probability distribution $P_\theta(y|\hat{\mathbf{x}}) : \mathbb{R}^p \times \mathbb{R} \rightarrow [0, 1]$ which is given by the logistic function. Given that we have a binary decision, i.e., a label $y \in \{-1, +1\}$, we can write the probability of the label y given the feature vector $\hat{\mathbf{x}}$ as:

$$P_\theta(y = 1|\hat{\mathbf{x}}) + P_\theta(y = -1|\hat{\mathbf{x}}) = 1 \quad (2)$$

where $P_\theta(y = 1|\hat{\mathbf{x}})$ and $P_\theta(y = -1|\hat{\mathbf{x}})$ are the probabilities of the positive and negative classes, respectively. We model these probabilities using the logistic function as follows:

$$P_\theta(y|\hat{\mathbf{x}}) = \begin{cases} \sigma(\hat{\mathbf{x}}^\top \theta) & \text{if } y = 1 \\ 1 - \sigma(\hat{\mathbf{x}}^\top \theta) & \text{if } y = -1 \end{cases} \quad (3)$$

where $\sigma(z) = 1/(1 + e^{-z})$.

2.1 Decision Boundary

The logistic regression model predicts the label $y \in \{-1, +1\}$ for a given feature vector $\hat{\mathbf{x}}$ by computing the probability $P_\theta(y|\hat{\mathbf{x}})$ and comparing it to a threshold, e.g., 0.5. The model predicts the positive class if the probability exceeds the threshold ($y = 1$). Otherwise, it predicts the negative class ($y = -1$). In the binary case, this is equivalent to selecting the class label with the highest probability. Suppose we have computed the probability of each binary label for n examples. We can store this data in a matrix $\mathbf{P} \in \mathbb{R}^{n \times 2}$, each row corresponds to an instance, and each column holds the probabilities of the positive and negative classes for a given example. Then, for each example, we can predict the label \hat{y} by comparing the probabilities:

$$i = \arg \max_y \{P_\theta(y = 1|\hat{\mathbf{x}}_i), P_\theta(y = -1|\hat{\mathbf{x}}_i)\} \quad \text{for } i = 1, 2, \dots, n \quad (4)$$

where the label is then given by:

$$\hat{y}_i = \begin{cases} 1 & \text{if } i = 1 \\ -1 & \text{otherwise} \end{cases} \quad (5)$$

Of course, these predictions are based on the model parameters θ , which we need to estimate from the training data. The following section will discuss the maximum likelihood estimation of the logistic regression model parameters.

2.2 Maximum Likelihood Estimation (MLE)

Maximum likelihood estimation (MLE) is a technique to estimate the parameters of a probability distribution by maximizing the likelihood function. In logistic regression, MLE estimates the parameters of the logistic regression model that make the observed label conditioned on the features the most probable. Given a set

of training examples $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, where $\mathbf{x}_i \in \mathbb{R}^m$ is an m -dimensional feature vector and $y_i \in \mathbb{R}$ is the binary (scalar) label, the likelihood function is defined as:

$$L(\theta) = \prod_{i=1}^n P_{\theta}(y_i | \hat{\mathbf{x}}_i) \quad (6)$$

where $P_{\theta}(y_i | \hat{\mathbf{x}}_i)$ is the probability of observing the label y_i given the feature vector $\hat{\mathbf{x}}_i$ and the model parameters θ . It's hard to maximize the likelihood function directly (because of the product), so we take the logarithm of the likelihood function to simplify the optimization:

$$\begin{aligned} \log L(\theta) &= \sum_{i=1}^n \log P_{\theta}(y_i | \mathbf{x}_i) \quad (\text{substituting } P_{\theta}(y_i | \mathbf{x}_i)) \\ &= - \sum_{i=1}^n \log \left(1 + e^{-y_i \cdot (\mathbf{x}_i^{\top} \theta)} \right) \end{aligned}$$

The maximum likelihood estimation (MLE) of the logistic regression model parameters θ^* is obtained by maximizing the log-likelihood function $\log L(\theta)$:

$$\theta^* = \arg \max_{\theta} \left[- \sum_{i=1}^n \log \left(1 + e^{-y_i \cdot (\hat{\mathbf{x}}_i^{\top} \theta)} \right) \right] \quad (7)$$

or alternatively by minimizing the negative log-likelihood function $\mathcal{L}(\theta) = -\log L(\theta)$:

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^n \log \left(1 + e^{-y_i \cdot (\hat{\mathbf{x}}_i^{\top} \theta)} \right) \quad (8)$$

Unfortunately, the optimization problem has no closed-form analytical solution, regardless of the perspective we take. Thus, we must estimate the model parameters using some numerical algorithm, such as gradient descent algorithm (or one of many other approaches). These approaches iteratively update the parameters θ to maximize the log-likelihood function (or minimize the negative log-likelihood function) until a stopping criterion is met, e.g., the parameter change is below a threshold or the maximum number of iterations is reached.

2.3 Gradient Descent

Gradient descent is a numerical search algorithm that minimizes a function by iteratively adjusting the parameters in the opposite direction of the gradient. Suppose there exists an objective function $\mathcal{L}(\theta)$ that we want to minimize with respect to the parameter θ , i.e., the negative log-likelihood function. In general, an objective function measures the difference between the predicted values and the observed values in some way, e.g., the mean squared error (MSE), the cross-entropy loss, or the negative log-likelihood. In logistic regression, the objective function is the negative log-likelihood function, which measures the difference between the predicted probabilities and the observed labels. However, whatever form the objective function takes, we assume that it is differentiable and that we can compute the gradient, i.e., $\nabla l(\theta)$ for the negative log-likelihood function, which points in the direction of the steepest increase of the function. This gives us a way to update the parameters to minimize the objective function using the update rule:

$$\theta_{k+1} = \theta_k - \alpha(k) \cdot \nabla_{\theta} \mathcal{L}(\theta_k) \quad \text{where } k = 0, 1, 2, \dots$$

The (hyper) parameter $\alpha(k) > 0$ is the learning rate (which can be a function of the iteration count k), and $\nabla_{\theta}\mathcal{L}(\theta)$ is the gradient of the negative log-likelihood function with respect to the parameters. We iterate until a stopping criterion is met, i.e., $\|\theta_{k+1} - \theta_k\| \leq \epsilon$, the maximum number of iterations is reached, or some other stopping criterion is met. Pseudo-code for a naive gradient descent algorithm (for a fixed learning rate) is shown in Algorithm 1.

Algorithm 1 Naive Gradient Descent for Negative Log-Likelihood $\mathcal{L}(\theta)$

```

1: Input: Initial parameters  $\theta_0$ , learning rate  $\alpha$ , stopping criterion  $\epsilon$ , maximum iterations  $N$ 
2: Output: Optimal parameter estimates  $\theta$ 
3: Initialize  $\theta \leftarrow \theta_0$  ▷ Initialize parameters to the initial guess
4:  $k \leftarrow 0$  ▷ Initialize iteration counter
5: while  $k \leq N$  or  $\|\theta_{k+1} - \theta_k\| \leq \epsilon$  do
6:    $\mathbf{d} \leftarrow \nabla\mathcal{L}(\theta_k)$  ▷ Compute gradient using analytical or numerical method, evaluate at  $\theta_k$ 
7:    $\theta_{k+1} \leftarrow \theta_k - \alpha \cdot \mathbf{d}$  ▷ Update parameters using the gradient direction  $\mathbf{d}$ 
8:    $k \leftarrow k + 1$ 
9: end while
10: return  $\theta$ 

```

Choose the learning Rate $\alpha(k)$

The choice of the learning rate α is crucial, as a too-large value can cause the algorithm to diverge, while a too-small value can slow down convergence. Choosing an appropriate learning rate is often challenging and may require tuning through hyperparameter optimization techniques. In practice, we may use adaptive learning rate methods, such as the Adam optimizer, which adjusts the learning rate based on the gradient's magnitude and the second moment of the gradient (1). Alternatively, Adagrad and RMSprop (unpublished) are other adaptive learning rate methods that can improve the convergence of the gradient descent algorithm (2).

2.4 Alternatives to Gradient Descent

The central issue with gradient descent is computing the gradient, which can be challenging for complex models or large datasets. Further, the choice of the learning rate α is also crucial, as a too-large value can cause the algorithm to diverge, while a too-small value can slow down convergence. In these cases, alternative heuristic optimization algorithms can be used to estimate the model parameters. The central theme of these approaches is that they directly evaluate the objective function at different points in the parameter space, update the model parameters using a heuristic, and then continue exploring the parameter space. New search points are generated randomly or with some heuristic, and the objective function is evaluated at these points to determine the next search direction. Let's walk through some of the alternatives to gradient descent.

There are many alternatives to gradient descent. The Nelder-Mead algorithm (3), also known as the simplex algorithm, is a direct search optimization algorithm that does not require the gradient of the objective function. It maintains a simplex (a geometric shape with $n+1$ vertices in an n -dimensional parameter space) that evolves through geometric operations such as reflection, expansion, contraction, and shrinkage. Thus, the Nelder-Mead algorithm can optimize non-differentiable or noisy objective functions, making it suitable for many optimization problems. However, the Nelder-Mead algorithm may struggle with high-dimensional problems or functions with many local minima, as it does not leverage gradient information to guide the search. Genetic algorithms (GAs), popularized by John Holland in the 1970s (4), are adaptive heuristic search techniques inspired by natural selection and genetics principles. They are designed to solve optimization

and search problems by iteratively evolving a population of candidate solutions through selection, crossover, and mutation. GAs aim to improve solution quality over generations by mimicking the evolutionary process. This makes them particularly effective for complex problems that may be discontinuous, non-differentiable, or highly nonlinear. However, GAs have several hyperparameters that need to be tuned, such as the population size, crossover rate, mutation rate, and the number of generations. Further, GAs can be computationally expensive, especially for large search spaces, and may require significant computational resources to converge to the optimal solution. This will be especially true for complex objective functions that are expensive to evaluate, as the algorithm will need to evaluate the objective function for each candidate solution in the population. Particle Swarm Optimization (PSO), developed by Kennedy and Eberhart in the mid-1990s (5) is an example of a meta-heuristic optimization algorithm inspired by the social behavior of birds and fish, which utilizes a population of candidate solutions, referred to as particles that move through the search space to find optimal solutions. Each particle adjusts its position based on its own experience and the collective knowledge of the swarm, allowing for efficient exploration and exploitation of the solution space to address complex optimization problems across various fields.

Simulated annealing, developed by Kirkpatrick et al. (6), is a probabilistic optimization algorithm inspired by the physical process of heating and then slowly cooling (annealing) materials to minimize defects. Simulated annealing works by iteratively exploring the solution space. First, a random (candidate) solution is generated, and then the objective function is evaluated. The difference in the objective function between the current and candidate solutions is computed, and the candidate solution is accepted if it improves the objective function. Alternatively, the candidate solution may be accepted with a certain probability even if it worsens the objective function. Thus, simulated annealing can escape local minima and explore the solution space more thoroughly than gradient descent, at least when the system is hot. As the number of iterations increases, the algorithm gradually decreases the temperature, i.e., the probability of accepting worse solutions, allowing it to converge to the global optimum. This method effectively solves complex optimization problems with large search spaces, where traditional techniques may struggle to find the global optimum. However, simulated annealing can be computationally expensive. It may require tuning of hyperparameters, particularly the annealing schedule, i.e., the decrease in temperature at each iteration, which impacts how willing the algorithm is to accept worse solutions to perform well. The pseudocode for the simulated annealing algorithm is shown in Algorithm 2.

Algorithm 2 Simulated Annealing

```

1: Initialize current solution  $x \leftarrow x_0$ 
2: Initialize temperature  $T \leftarrow T_0$ 
3: Define cooling schedule  $T \leftarrow \alpha \cdot T$ , where  $0 < \alpha < 1$ 
4: Define maximum iterations  $N$ 
5: Initialize best solution  $x_{\text{best}} \leftarrow x$ 
6: for  $i = 1$  to  $N$  do
7:   Generate a new candidate solution  $x_{\text{new}}$  in the neighborhood of  $x$ 
8:   Compute  $\Delta E \leftarrow f(x_{\text{new}}) - f(x)$ 
9:   if  $\Delta E < 0$  then
10:    Accept  $x_{\text{new}}$ :  $x \leftarrow x_{\text{new}}$ 
11:   else
12:    Accept  $x_{\text{new}}$  with probability  $P \leftarrow e^{-\Delta E/T}$ 
13:    Draw a random number  $r \in [0, 1]$ 
14:    if  $r < P$  then
15:       $x \leftarrow x_{\text{new}}$ 
16:    end if
17:   end if
18:   if  $f(x) < f(x_{\text{best}})$  then
19:     Update  $x_{\text{best}} \leftarrow x$ 
20:   end if
21:   Update temperature:  $T \leftarrow \alpha \cdot T$ 
22:   if  $T$  is below a threshold  $T_{\text{min}}$  then
23:     break
24:   end if
25: end for
26: return  $x_{\text{best}}$ 

```

3 Summary and Conclusions

In this lecture, we introduced logistic regression for binary classification problems. Logistic regression is a powerful tool for binary classification and is widely applied in numerous fields, including healthcare, finance, and marketing. Like the Perceptron model, logistic regression uses a linear model to describe the relationship between the features and the label. However, using the logistic function, it estimates the probability that a given input belongs to a particular class. We discussed the maximum likelihood estimation (MLE) technique, which maximizes the likelihood of the observed labels given the features to estimate the parameters in the logistic model. However, unlike our previous discussion of linear regression, the logistic regression model's likelihood function does not have a closed-form analytical solution. Thus, we must estimate the model parameters using numerical optimization algorithms, such as the gradient descent algorithm. Gradient descent is an optimization algorithm that minimizes a function by iteratively adjusting parameters in the opposite direction of the gradient. However, we also discussed other (heuristic) optimization algorithms that can estimate model parameters without relying on the gradient. For instance, simulated annealing, genetic algorithms, and particle swarm optimization are all methods that can be utilized to estimate the model parameters. These methods can be particularly useful when the objective function is non-convex, non-differentiable, or has many local minima.

References

1. Kingma DP, Ba J. Adam: A Method for Stochastic Optimization; 2017. Available from: <https://arxiv.org/abs/1412.6980>.
2. Duchi J, Hazan E, Singer Y. Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *J Mach Learn Res*. 2011;12(null):2121–2159.
3. Nelder JA, Mead R. A Simplex Method for Function Minimization. *The Computer Journal*. 1965;7(4):308–313. doi:10.1093/comjnl/7.4.308.
4. Holland JH. *Adaptation in natural and artificial systems: an introductory analysis with applications to biology, control, and artificial intelligence*. Ann Arbor: University of Michigan Press; 1975.
5. Kennedy J, Eberhart R. Particle swarm optimization. In: *Proceedings of ICNN'95 - International Conference on Neural Networks*. vol. 4; 1995. p. 1942–1948 vol.4.
6. Kirkpatrick S, Gelatt CD Jr, Vecchi MP. Optimization by simulated annealing. *Science*. 1983;220(4598):671–80. doi:10.1126/science.220.4598.671.