

Logistic regression and neural networks

Eirik Nordgård
Geophysical Institute, University of Oslo

2019
October

Abstract

Here you can write your abstract

All material for this project may be found on <https://github.com/eirikngard/Project2>

1 Introduction Change number 2

Introduction to the subject

HER SKRIVER JEG EN ENDRING

2 Theory

2.1 Logistic regression

Logistic regression is used to solve classification problems. These are problems concerning outcomes, y_i , in form of discrete variables. Commonly the classification problems in question has two possible outcomes, true or false. Using the credit card data as an example, the two outcomes would be if a client would pay his/hers debt or not. This type of outcome is often called binary outcome, and can easily be programmed as 1 or 0.

The dependent variables (outcomes) y_i are discrete, ranging from $k = 0, \dots, K - 1$ where K is classes. The main goal is to predict the output classes from a designmatrix $\hat{X} \in R^{n \times p}$ (change R to real numbers R) made of n samples carrying p features/predictors. The

simplest output is perhaps a binary output, only having values 0 or 1 meaning yes/no, true/false, pay/dont pay etc. From there it is desired to identify the classes to which new unseen samples belong.

A discrete output can be obtained in several ways, but the simplest may be to have a sign function, which maps the output of a linear regressor to values 0, 1, $f(a) = \text{sign}(a) = 1$ if $a \geq 0$ and 0 if otherwise. Although this model is very simple, in many cases it might be convenient to know the probability of an output belonging a given category rather than a single value. This is done using the *logistic function*, and is often called a "soft classifier". The classifier outputs the probability of x_i belonging to a category $y_i = \{0, 1\}$. In this case the classifier is given by a *Sigmoid* function, often referred to as a likelihood function. The Sigmoid function $p(y)$ is given by[2]

$$p(y) = \frac{1}{1 + e^{-y}} = \frac{e^y}{1 + e^y} \quad (1)$$

Like other likelihoods, eq. (1) here we have that the likelihood $p(y_i = 0) = 1 - p(y_i = 1)$, always resulting in a summed likelihood of 1. Assuming two classes y_i being either 0 or 1, we have

$$y_i = \beta_0 + \beta_1 x_i + \dots + \beta_n x_i \quad (2)$$

Used in in eq.(1), this gives

$$p(y_i = 1 | x_i, \hat{\beta}) = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \quad (3)$$

and

$$p(y_i = 0 | x_i, \hat{\beta}) = 1 - p(y_i = 1 | x_i, \hat{\beta}) \quad (4)$$

Eq. (3) and eq. (4) can then be used to find the maximum likelihood of an event. The product of all individual probabilities of a specific outcome y_i is used to obtain a log-likelihood function, which in turn gives the cost function

$$C(\hat{\beta}) = \sum_{i=1}^n (y_i \log p(y_i = 1 | x_i, \hat{\beta}) + (1 - y_i) \log [1 - p(y_i = 1 | x_i, \hat{\beta})])$$

Eq. (2.1) is called the *cross entropy*. Being convex, a local minimizer of this function will also be a global minimizer. Hence, minimizing eq. (2.1) with respect to each β gives

$$\frac{\partial \hat{\beta}}{\partial \beta_0} = - \sum_{i=1}^n (y_i - \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}) \quad (5)$$

and

$$\frac{\partial \hat{\beta}}{\partial \beta_1} = - \sum_{i=1}^n (y_i x_i - x_i \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}) \quad (6)$$

It is worth mentioning that far from all of these functions are convex. Many cost functions used in machine learning are in fact non-convex and of high dimensionality. Computing eq. (5) and eq. (6) we get the following expressions for the first and the second derivative:

$$\frac{\partial C(\hat{\beta})}{\partial \hat{\beta}} = -\hat{X}^T (\hat{y} - \hat{p}) \quad (7)$$

$$\frac{\partial^2 C(\hat{\beta})}{\partial \hat{\beta} \partial \hat{\beta}^T} = \hat{X}^T \hat{W} \hat{X} \quad (8)$$

where \hat{y} is a vector with n elements y_i , the $n \times p$ matrix \hat{X} containing the x_i values, the vector \hat{p} is the fitted probabilities $p(y_i | x_i, \hat{\beta})$ and \hat{W} being a diagonal matrix with elements $p(y_i | x_i, \hat{\beta})(1 - p(y_i | x_i, \hat{\beta}))$ (skal $y = 1$ i denne?). These derivatives are used to create the Hessian matrix, later used to solve the minimization problem. Unlike in linear regression, solving these for $\hat{\beta}$ now requires a gradient descent method[2].

2.2 Gradient Descent Method

The idea of a gradient descent (GD) method is that a function $F(x)$, with $x = x_1, \dots, x_n$, decreases fastest if one goes from x in the direction of the negative gradient $-\text{gradient}F(x)$. If

$$x_{v+1} = x_v - \gamma_v \nabla F(x_v) \quad (9)$$

where γ_v is called learning rate or step length, in this case $\gamma_v > 0$. For small enough γ_v we will then have $F(x_{v+1}) \leq F(x_v)$, meaning we are always moving towards smaller values and eventually a minimum of F ¹. Choosing appropriate learning rate is important. If chosen too small the method will converge very slowly, and chosen too large may result in unpredictable behaviour. The first step doing the GD is to make a guess x_0 for a minimum of F . Note that the minima is a global minima only if the function is a convex function.

Maby include a paragraph on convex functions, as it is desired to have such functions to proceed with GD?

Logistic regression er som et neuralt nettverk uten layers.

The cost function I actually used:

¹this paragraph until here is almost copy of a paragraph in slides of Splines. Might wanna find another introduction to avoid copying

$$C(\hat{\beta}) = - \sum_{i=1}^n (y_i(\beta_0 + \beta_1 x_i) - \log(1 + e^{\beta_0 + \beta_1 x_i})) \quad (10)$$

Eq. (10) is known as the cross entropy, and provide the gradient used for the standard gradient descent. The gradient is the derivative of eq. (10) with respect to $\hat{\beta}$, thus

$$\frac{\partial C(\hat{\beta})}{\partial \hat{\beta}} = -\hat{X}^T(\hat{y} - \hat{p}) \quad (11)$$

² [2].

2.2.1 Standard/Steepest Gradient Descent

I have used this as the first gradient descent. Iterate 1000, calculated the gradient, minimizing the cost function. This iteration outputs betas, used to make prediction of y. Thereafter y is used in the sigmoid function to classify the data as default or not. Making a threshold of 0.5, where sigmoid output of less than 0.5 is 0, and above is 1. Then I count how many of the predicted values match with the values of the data y (last column in data set) to get the accuracy score for this gradient descent.

2.2.2 Stochastic Gradient Descent

Write about epochs and minibatches. Why we use them, how we use them and maybe how to optimize the number of these?

Anyhow, I have used stochastic gradient descent with linear regression as well, and found it to be around 1

2.3 Artificial Neural Networks

Artificial Neural Networks (ANN) is a computational system that by training on examples can learn to perform tasks. Increasing the

number of examples the ANN can learn from will increase the accuracy of the tasks performed. It is called artificial neural network because it is supposed to mimic a biological system such as a brain, where neurons interact by sending signals in the form of mathematical functions. The ANN is constructed in layers containing an arbitrary number of neurons, which in turn is connected by a weight variable. Similarly to an actual brain, the neurons in a ANN can only communicate with each other if the incoming signal exceeds a certain threshold value. If the signal is "strong" enough, an output is sent forward, but if this activation threshold is not reached the neuron remains inactive, providing zero output[2].

The neurons described above may be on different forms. One type of artificial neuron is called the *perceptron*. This works by taking n binary inputs x_1, x_2, \dots, x_n and producing one single binary output. The importance of each binary input is represented by n weights, w_1, w_2, \dots, w_n . The binary output is determined by if the weighted sum of inputs is less or greater than a threshold (bias) value. Thus,

$$\text{output} = 0 \text{ if } \sum_i w_i x_i \leq \text{threshold}$$

$$\text{output} = 1 \text{ if } \sum_i w_i x_i \geq \text{threshold}$$

This is the equivalent to whether the neuron is activated or not. Described in another way, the perceptron makes a decision (0 or 1) by weighing up arguments pro or con the decision. The network may also have several layers of neurons, each taking the output of the "previous" layers as input. These layers are called *hidden layers*, and allows for very complicated decision making networks. For understanding, the next paragraph states a simple example.

²Make sure p (sigmoid) is defined somewhere above

Imagine you want to visit the cabin. Two inputs may be if the weather is nice or not, and if you get to see yo grandmother or not. The output is simply to go or not to go to the cabin. Then the question is what weighs most, the weather or you being able to visit your grandmother. Say the threshold is 3. The weight for nice weather is 2 and the weight for visiting the grandmother is 2. Then you actually require both nice weather and that you get to see you grandmother to go to the cabin. If the weather is nice, but you cant see your grandmother, we have $1 * 2 + 0 * 2 = 2 < 3$ and the threshold is not reached. Similarly if the weather is bad, but you can see your grandmother we have $0 * 2 + 1 * 2 = 2 < 3$. But as both requirements are fulfilled we have $1 * 2 + 1 * 2 = 4 > 3$ and the threshold is reached, giving the binary output of 1, which in this case means you are going to the cabin. IF the threshold were lowered to 1, you would go to the cabin even though only one of the inputs were fulfilled. Thus, changing the weights and the threshold yields different decision making.

If the perceptron where to make a wrong classification, a problem arises. The desired path of action would be to make small changes is the biases or weights of the network and check if it then made the correct classification. If not, you would repeat this process making small changes until the classification succeeds. However, making small changes in biases or weights is actually resulting in very large changes in the output. This problem is solved using another type of neuron, namely the *sigmoid*.

This neuron functions in the same way as the perceptron, but with different input and output. Now the inputs can be any value between 0 and 1. And instead of having output of 0 or 1, the sigmoid has output $\sigma(w * x + b)$, where

b is the bias ³. and σ is the sigmoid function defines eq. (1). Hence the out of the sigmoid will be

$$output = \frac{1}{1 + e^{-\sum w_i x_i - b}} \quad (12)$$

[4] Rydd opp litt i denne teoridelen. Den likningen som kommer under burde være med, fordi den setes jo egnetlig inn i sigmoid, som gir/er likningen ovenfor. INCLUDE SOME SORT OF SCECH OG SIMOID?

WHAT FOLLOWS should possibly be after the first paragraph in this section: *A simple neuron model can hence be written*

$$y = f\left(\sum_{i=1}^n w_i x_i + b_i\right) = f(u) \quad (13)$$

[2] In eq. (13) x_i is the input signals, weighted by w_i . y is the output, equal to the value of its activation function. b_i is the bias/threshold. This can in turn be used to separate different type of neural networks. Such variations may be for instance Recurrent Neural Networks or Convolutional Neural Networks⁴. UNTIL HERE.

2.3.1 Backpropagation (ANN)

Backpropagation is an algorithm computing the gradient of the cost function. This quantity tells us what the partial derivatives with respect to both the bias and the weights are, and may be very useful if you want to know how fast the cost is changing when the bias or weights are changing. This actually allows us to get information on how the behaviour of the entire network changes when the bias or weights are changed. Like before, we have

$$a^l = \sigma(w^l a^{l-1} + b^l) \quad (14)$$

³include this notation for output further up as well?

⁴Further reading on the differnt types of neural network may be found on: URL

where a is the activation (referred to as output later), b is the bias, w is the weights, l is the layer number and σ is a function to every element in the vector $w^l a^{l-1} + b^l$, hereafter called z^l . This way the activation in one layer is related to the activation in the previous layer. Thus we have $a^l = \sigma(z^l)$. (Sigma is the same as f above i think. SJEKK⁵). z^l is called the weighted input to the neurons in layer l .

Then, for backpropagation to work, it is necessary to assume the the cost function can be written as an average $C = \frac{1}{n} \sum_x C_x$ over cost functions C_x for individual training points, x . Also, the cost function must be able to be written as a function of the outputs from the neural network (SIDE 42 NIELSEN).

While computing the derivatives of the cost $\frac{\partial C}{\partial b^l}$ and $\frac{\partial C}{\partial w^l}$, an error will occur. This error is just referred to as *error*, and is written δ_j^l in the j -th neuron in the l -th layer. For later notation, the k -th neuron is belongs to the $(l-1)$ -th layer. The quantity $\frac{\partial C}{\partial z_j^l}$ provides a good estimate of this error (EXPLANATION for why IN page 44 Nielsen), and from here

$$\delta_j^l = \frac{\partial C}{\partial z_j^l} \quad (15)$$

of neuron j in layer l . Backpropagation then relates δ to the derivatives of the cost. Rewriting eq. (15) gives

$$\delta_j^l = \frac{\partial C}{\partial a_j^l} \sigma(z_j^l) \quad (16)$$

which vectorized is

$$\delta^l = \nabla_a C \odot \sigma(z^l) \quad (17)$$

In eq. (17) $\nabla_a C$ expresses the rate of change of C with respect to the activation output. Further rewriting gives

$$\delta^l = (a^l - y) \odot \sigma'(z^l) \quad (18)$$

⁵sjekk om f og sigma er det samme

which in terms of the error in the next layer, δ^{l+1} , is

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l) \quad (19)$$

The name *backpropagation* arises from eq. (19). Here the error is moved backward through the activation function in layer l , giving the error in the weighted input to layer l .

Eq. (19) finally allows us to obtain equations for the partial derivatives with respect to the bias and the weights, respectively

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l \quad (20)$$

simplified to

$$\frac{\partial C}{\partial b} = \delta \quad (21)$$

and

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \quad (22)$$

simplified to

$$\frac{\partial C}{\partial w} = a_{in} \delta_{out} \quad (23)$$

Eq.(22) has a very useful consequence. When the activation a is small, the derivative will also be small. This means that the gradient changes fairly little during a gradient descent, meaning the weight *learns slowly*.

The resulting backpropagation algorithm then looks something like this:

1. **Input x :** Find corresponding activation a^l for the input layer.
2. **Feedforward:** For $l = 2, 3, \dots, L$ compute z^l (EQ HASENT NUMBER YET) and a^l (eq.(14)).
3. **Output error δ^L :** Compute vector $\delta^L = \nabla_a C \odot \sigma'(z^L)$

4. **Backpropagate the error:** Compute $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$ for each l .
5. **Output:** Gradient of the cost function: $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$ and $\frac{\partial C}{\partial b_j^l} = \delta_j^l$

[4]

2.3.2 Measuring performance

Maybe include some words on how i measure the performance of the models. Have used accuracy score. Possibly use area under curve and F1 score.

2.4 Data description

The data set used in this analysis is credit card data from credit card holders in Taiwan. In total the dataset contains 23 variables, which is used to employ a response variable, *default payment* with value $1 = \text{default} = \text{not pay}$, $0 = \text{not default} = \text{pay}$. The dataset contains the following information:

- X1: Amount of the given credit (NT dollar): it includes both the individual consumer credit and his/her family (supplementary) credit.
- X2: Gender (1 = male; 2 = female).
- X3: Education (1 = graduate school; 2 = university; 3 = high school; 4 = others).
- X4: Marital status (1 = married; 2 = single; 3 = others).
- X5: Age (year).
- X6–X11: History of past payment. We tracked the past monthly payment records (from April to September, 2005) as follows: X6 = the repayment status in September, 2005; X7 = the repayment

status in August, 2005;...;X11 = the repayment status in April, 2005. The measurement scale for the repayment status is: 1 = pay duly; 1 = payment delay for one month; 2 = payment delay for two months; ...; 8 = payment delay for eight months; 9 = payment delay for nine months and above.

- X12–X17: Amount of bill statement (NT dollar). X12 = amount of bill statement in September, 2005; X13 = amount of bill statement in August, 2005;...;X17 = amount of bill statement in April, 2005.
- X18–X23: Amount of previous payment (NT dollar). X18 = amount paid in September, 2005; X19 = amount paid in August, 2005;...;X23 = amount paid in April, 2005.[5]

All of these variables determines the probability of whether a client defaults or not on his/her debt the following month.

3 Method

3.1 Data Preparation

Not processed in any way, the dataset is roughly 30000X23 data points. Within these values are several illegal ones, not mentioned in the data description. Thus, the following data are removed from the set: $EDUACTION = 0, 5, 6$. $MARRIGE = 0$. $PAY_X = -2$. $BILL_{AMTX} < 0$. $PAY_{AMTX} < 0$. Additionally, according to Vladimir G. Drugov [6] 86,5 of PAY_X has the illegal value 0. Removing these would mean losing too much of the total data, so these values are not removed. We are then left with 22455X23 data points.⁶ Containing a lot more not default (0) than default (1), the data is very skewed. This may

⁶ONEHOTENCODING. Why encode those we encode.

result in models gaining better performance in predicting zeros than ones. To correct this skewed distribution the data is down-sampled so that the model can be trained on equally many zeros and ones. This leaves $12977x$ ⁷.

4 Conclusion

Something here

5 Conclusion

And we'd like the network to learn weights and biases so that the output from the network correctly classifies the digit.

For example, suppose the network was mistakenly classifying an image as an "8" when it should be a "9". We could figure out how to make a small change in the weights and biases so the network gets a little closer to classifying the image as a "9". And then we'd repeat this, changing the weights and biases over and over to produce better and better output. The network would be learning.

6 REMEMBER

- Contro all references, that they are written correctly
- Add correct superscripts to the part of backpropogation
- Just include formulas you actually use in the theory part.
- Maybe include illustrative figures on how neural network work.
- Write about Softmax somewhere. You are using this in your code.

⁷insert correct number according to onehot here.

Example of multilined eq:

$$\frac{dC(\beta)}{d\beta_j} = \frac{d}{d\beta} \left[\frac{1}{n} \sum_{i=0}^{n-1} (y_i - \beta_0 x_{i,0} - \beta_1 x_{i,1} - \dots - \beta_{n-1} x_{i,n-1})^2 \right] = 0 \quad (24)$$

7 Appendix

8 References

References

- [1] Hastie, Trevor. Tibshirani, Robert. Friedman, Jerome. *The Elements of Statistical Learning. Data Mining, Interference, and Prediction*. Second Edition. Springer, 2009. Chapter 3, Chapter 7
- [2] M. Hjorth-Jensen Lecture Notes in FYS-STK4155. *Data Analysis and Machine Learning: Linear Regression and more Advanced Regression Analysis*. URL: <https://compphysics.github.io/MachineLearning/doc/pub/Regression/html/Regression.html> Unpublished, 2019.
- [3] Wikipedia: Bias-Variance tradeoff URL: https://en.wikipedia.org/wiki/Bias_variance_tradeoff Read: 23.09.2019
- [4] Nilsen BLABLABLA URL
- [5] The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. I-Cheng Y. and Che-hui L. URL: Det er et problem med å få lagt til linken
- [6] URL: