

1 Binary Response Models

1.1 Theory

The outcome of interest in the Creditreform-dataset is of binary nature. Firms can either go bankrupt or not. This behavior is commonly modeled by Binary Response Models like the probit or the logit model. Binary response variables follow a Bernoulli probability function

$$f(y|x) = P(y = 1|x)^y (1 - P(y = 1|x))^{1-y}, \quad y \in \{0, 1\}, x \in \mathbb{R}^d, d \in \mathbb{N}, \quad (1)$$

where $P(y = 1|x)$ stands for the conditional probability of observing $y = 1$ given x . Both probit and logit models have in common that $P(y = 1|x)$ is modeled by a monotonic transformation of a linear function

$$P(y = 1|x) = G(x'\beta), \quad \beta \in \mathbb{R}^d, \quad (2)$$

where $x'\beta$ is the scalar product of x and β . Additionally we require that $0 \leq G(x'\beta) \leq 1$, since it denotes a probability.

For the probit model G will be the cumulative density function of the normal distribution

$$P(y = 1|x) = G(x'\beta) = \Phi(x'\beta) = \int_{-\infty}^{x'\beta} \frac{1}{\sqrt{2\pi}} \exp\left[-\left(\frac{t^2}{2}\right)\right] dt. \quad (3)$$

Here $\Phi(x'\beta)$ stands for the cumulative density function of the normal distribution.

For the logit model G will be replaced by the cumulative density function of the logistic distribution $\Lambda(x'\beta)$:

$$P(y = 1|x) = G(x'\beta) = \Lambda(x'\beta) = \frac{\exp(x'\beta)}{1 + \exp(x'\beta)} \quad (4)$$

[1]

The parameter vector β is obtained by the Maximum-Likelihood method. Given independent and identically distributed samples, the Likelihood function can be written as

$$\begin{aligned} L(\beta; y, x) &= \prod_{i=1}^n f(y_i|x_i) = \prod_{i=1}^n P(y_i = 1|x_i)^{y_i} (1 - P(y_i = 1|x_i))^{1-y_i} \\ &= \prod_{i=1}^n G(x_i'\beta)^{y_i} (1 - G(x_i'\beta))^{1-y_i} \end{aligned} \quad (5)$$

where we just take the product over all individual Bernoulli-functions.

The log-Likelihood can thus be written as

$$l = \log L(\beta; yx) = \sum_{i=1}^n y_i \log G(x_i \beta) + (1 - y_i) \log(1 - G(x_i \beta)) \quad (6)$$

The Maximum-Likelihood estimators $\beta_M L$ are calculated as

$$\beta_M L = \operatorname{argmax}(l) \quad (7)$$

and solve the first order conditions for a maximum.

$$\frac{\partial l}{\partial \beta} \stackrel{!}{=} 0 \quad (8)$$

In general, the resulting system of equations has no closed-form solution for $\beta_M L$ and numerical solutions are needed which can be obtained by iterative optimization techniques, one we will implement in the next section.[1]

1.2 Implementation

The architecture of the `brm`-class follows the general structure outlined in the chapter before. First, it generates a log-Likelihood-function, which it then optimizes using a Gradient-Descent-Algorithm. After training the model it is possible for the user to very easily obtain predictions by invoking the `predict()`-function, which has been augmented with a method for the `brm`-class.

1.2.1 Obtaining the Likelihood

The first task is to define a function that accepts a distribution and yet undefined data as it's input and first extracts all suitable variables, then expresses the Likelihood-function from 6 and finally returns another function which depends only on the weights β . This task is performed by the `get_likelihood()`-function. We outline the function pass of `get_likelihood()` first in pseudo-code followed by a look on the implementation in the R language.

The heavy lifting in this function is done by this R-snippet:

```

15  l = function(x){
16    - sum(apply(Xy_mat, 1,
17              function(X){
18                X[y_pos] * distr(t(X[-y_pos]) %% x,
19                                lower.tail = TRUE,
20                                log.p = TRUE) +
21                (1-X[y_pos]) * distr(t(X[-y_pos]) %% x,
22                                    lower.tail = FALSE,
23                                    log.p = TRUE))}))
24    # Return loglikelihood as a function of x (here 'x' stands for
25    the weights)
    return(l)

```

Algorithm 1 `get_likelhood()`

```
1: procedure SET UP AUXILIARY VARIABLES
2:    $grp \leftarrow$  unique labels
3:    $nums \leftarrow$  extract numeric columns
4:    $Xy\_mat \leftarrow$  bind numeric variables as a matrix
5:    $y\_pos \leftarrow$  cache position of the outcome variable
6: procedure SET UP LOG-LIKELIHOOD
7:    $l \leftarrow 0$ 
8:   for:  $x_i, y_i$  in data :
9:     calculate :  $j = y_i \log G(x_i \beta) + (1 - y_i) \log(1 - G(x_i \beta))$ 
10:    update :  $l \leftarrow l + j$ 
11:  return :  $l(\beta)$ 
```

The **for**-loop from the pseudo code is implemented as an **apply**-call to `Xy_mat`, which in turn is a matrix of numeric columns. The **apply**-function initially selects each row in `Xy_mat` and extracts the outcome-variable `Xy_mat[y_pos]` which corresponds to y_i from 6. It then computes the scalar product between the regressors of `Xy_mat`'s row, which plays the role of $x_i \beta$ in 6. The scalar product is wrapped in **dist** which represents the cumulative distribution function $G(x \beta)$ and is one of the arguments to `get_loglikelihood()`.

```
30   distr = switch (mode,
31                   "logit" = plogis,
32                   "probit" = pnorm
33   )
```

This way **dist** will point to the built-in functions for computing probabilities, depending if the user wishes to train a logit or a probit model. The arguments `lower.tail=TRUE` and `lower.tail=FALSE` stand for $G(x \beta)$ and $G(x \beta) = 1 - G(x \beta)$ respectively. Finally the resulting vector is summed up and multiplied by -1 . This is done because of the way we implemented the Gradient Descent algorithm. Currently `gradientDescentMinimizer()` can only find minima. However, Maximum-Likelihood estimation poses a maximization problem. Luckily, we can transform any maximization problem into a minimization problem by multiplying with minus one.

1.2.2 Gradient Descent

Since we now have a log-Likelihood function, the next step is to optimize it. To this effect we deploy a Gradient Descent algorithm. Theory tells us that in order to reach the minimum of a function $f(x)$ starting at a particular $x \in \mathbb{R}^d, d \in \mathbb{N}$ one needs to follow the negative gradient $\nabla f(x)$ of f evaluated at x . This leads to the iterative rule we can exploit

$$x_{t+1} = x_t - \eta \cdot \nabla f(x_t), \quad t \in \mathbb{N}, \eta \in \mathbb{R}^+ \quad (9)$$

where η is the learning rate. To this standard method of performing a Gradient Descent routine we will also make some minor modifications. First of all, we will approximate the gradients by taking finite differences, which is easier to implement albeit computationally inefficient. Finite differences are computed by

$$\nabla f(x_t) \approx \frac{f(x_{t+1}) - f(x_t)}{\epsilon}, \quad \epsilon > 0 \quad (10)$$

Secondly, we will before we initialize the algorithm, try a set of random points and chose the one that provides the lowest value of the objective function as a starting point for the Gradient Descent Routine. This also ensures to an extend that the algorithm, if it reaches convergence, finds the global minimum. The final modification will be to prune the gradients. When computing gradients using finite differences, it may happen that the gradient's values can become extremely high for large denominators and very small ϵ . In fact, they can become high enough for R to treat them as `Inf` which results in the gradients being treated as `NaN` (not a number). To counteract that, we will limit the gradients to the interval $[-100, 100]$.

Algorithm 2 `gradientDescentMinimizer()`

```

1: procedure SET UP AUXILIARY VARIABLES
2:   learn_rates  $\leftarrow$  descending sequence from learn to 0
3:   a  $\leftarrow$  matrix of 1000 randomly initialized points
4:   f_a  $\leftarrow$  vector of function values for each element of a
5:   update: a  $\leftarrow$  argmin(f_a)
6:   gradient  $\leftarrow$  compute gradient evaluated at a
7:   i  $\leftarrow$  0
8: procedure PERFORM GRADIENT DESCENT
9:   l  $\leftarrow$  0
10:  while: i  $\leq$  max_iter and any element of gradient  $>$  0 :
11:    update : a = a - learn_rates[i]  $\cdot$  gradient
12:    calculate : gradient = calculate gradient
13:    if i = max_iter then raise warning
14:  return : a
```

The `gradientDescentMinimizer()`-function accepts following arguments:

1. **obj:** an objective function, that accepts exactly one argument called 'x'.
2. **n_pars:** an integer specifying the dimensions of the objective.
3. **epsilon_step:** a float defining the stepwidth used for computing the finite differences.
4. **max_iter:** an integer for the maximum number of iteration before the algorithm aborts.

5. **precision**: a float defining the precision of the solution. All elements of the gradient have to be absolutely lower than **precision** for the algorithm to converge.
6. **learn**: a positive float representing the learning rate.
7. **verbose**: a boolean indicating if additional information during training is desired. The default is **FALSE**
8. **report_freq**: If **verbose** is **TRUE**, define how often to print the logstring. The default is 10 which corresponds to a console output being printed every 10 steps.

```

143 a = matrix(data = runif(1000 * n_pars,
144                        min = -10,
145                        max = 10),
146           ncol = n_pars)
147 f_a = apply(a, 1, obj)
148 a = a[which.min(f_a), ]

```

We begin by filling the matrix **a** with 1000 **n_pars**-dimensional points which we draw from the uniform distribution, making use of R's built-in **runif**-function. We draw random numbers within the range of $[-100, 100]$ to cover a wide part of the objective function's domain.

The workhorse in this routine is the **get_gradient()**-function, which computes the finite differences. First we need to compute the values of the objective function at the current and next step (lines 146 and 147). Then we can apply the current and next step as inputs to the objective function and compute the difference $f(x_{t+1}) - f(x_t)$. The current step is provided as the **x** argument to the function call. The next steps need to be inferred by the function. If $f(x_t)$ is multidimensional we need to perform an ϵ -step in each dimension of the vector, since we want to approximate the partial derivatives of $f(x)$ evaluated at x_t . I.e. first we want to increment just the first element of x and store the result, then just the second element, and repeat the process until we reach the last element. If we stack these vectors, we get a matrix of one-directional ϵ -steps that are essentially updates of the starting point x_t with which it is easy to compute the gradients as their element-wise difference, normalized by **epsilon_step**. The gradients are finally trimmed if necessary and returned.

```

130 get_gradient = function(x, d = n_pars,
131                        objective = obj,
132                        epsilon = epsilon_step){
133   init = matrix(data = x, nrow = d, ncol = d, byrow = TRUE)
134   steps = init + diag(x = epsilon, ncol = d, nrow = d)
135   f_steps = apply(steps, 1, objective)
136   f_comp = apply(init, 1, objective)
137   D = (f_steps - f_comp) / epsilon
138   D_trimmed = ifelse(abs(D) <= 100, abs(D), 100) * sign(D)
139   return(D_trimmed)}

```

Algorithm 3 `get_gradient()`

1: **procedure** SET UP AUXILIARY VARIABLES

2:

$$init \leftarrow \begin{bmatrix} x_1 & x_2 & \cdots \\ x_1 & x_2 & \cdots \\ x_1 & x_2 & \cdots \\ \vdots & \ddots & \cdots \end{bmatrix}$$

3:

$$steps \leftarrow \begin{bmatrix} x_1 + \epsilon & x_2 & \cdots \\ x_1 & x_2 + \epsilon & \cdots \\ x_1 & x_2 & \cdots \\ \vdots & \ddots & \cdots \end{bmatrix}$$

4: $f_comp \leftarrow$ **apply row-wise:** objective function to $init$ 5: $f_steps \leftarrow$ **apply row-wise:** objective function to $steps$ 6: **procedure** COMPUTE FINITE DIFFERENCES

7:

$$D \leftarrow \frac{f_steps - f_comp}{\epsilon}$$

8: **if** $anyd \in D \notin [-100; 100]$ **then replace** d by $100 \cdot \text{sign}(d)$

```
155 while(any(abs(gradient) >= precision) & i <= max_iter){
156     if(i %% report_freq == 0 & verbose) {
157         cat("\nStep:\t\t", i,
158             "\nx:\t\t", a,
159             "\ngradient:\t", gradient,
160             "\nlearn:\t", learn_rates[i],
161             "\n-----")
162
163         i = i + 1
164         a = a - learn_rates[i] * gradient
165         gradient = get_gradient(a)
166     }
167
168     cat("\nResults\n",
169         "\nIteration:\t", i,
170         "\nx:\t\t", a,
171         "\nf(x):\t\t", obj(a),
172         "\ndf(x):\t\t", gradient,
173         "\n")
174
175     if(i >= max_iter){
176         warning("Maximum number of iterations reached.")
177         return(a)
178     }
```

In each iteration the **while**-loop ensures that convergence has not been reached. This is implemented by a call to **any** wrapped around a vector of logi-

cal expressions. If any element of the gradient is still greater than the specified precision, the call to `any` will evaluate to `TRUE`. The second breaking criterion is a safeguard for the loop not to run infinite times. If the current iteration is larger than `max_iter` the algorithm will break and the user will receive a warning (lines 175-176). If the user wishes to receive information about the status of the algorithm during runtime, the optional argument `verbose` can be set to `TRUE` which will print a logstring to the console in regular intervals (lines 156-161).

1.3 Predictions

In order to facilitate making predictions based on the `brm`-class we augmented the built-in function `predict()` with a method that works on our custom class in a predefined way.

```

61 predict.brm = function(model, data){
62   Xb = as.matrix(cbind(constant = 1, data)) %*% model$weights
63   predictions = sapply(X = Xb,
64                        FUN = model$distribution,
65                        log.p = FALSE,
66                        lower.tail = TRUE)
67   return(predictions)
68 }

```

A call to `predict()` on a `brm`-model will add a column of ones to the provided `data` and multiply the matrix with the weights calculated during training of the model (line 62). Finally these scores of the index-function $X\beta$ will be applied to the correct distribution. The distribution is stored inside `model$distribution` which points to `pnorm` in case of `brm`-model of `mode` "probit" and a pointer to `plogis` if the `mode` is equal to "logit".

2 Linear Discriminant Analysis

2.1 Theory

Linear Discriminant Analysis (LDA) is a technique for dimensionality reduction that incorporates information on class-labels of the different observations. In contrast to Principal Component Analysis, which is a unsupervised dimensionality reduction technique, it finds the rotation that ensures the highest separability between classes. It accomplishes this goal by trying to maximize between class variance while simultaneously minimizing within class variance.

$$\max J_b(w) = w^T S_b w, \quad w \in \mathbb{R}^d, d \in \mathbb{N} \quad (11)$$

$$\min J_w(w) = w^T S_w w \quad (12)$$

This is done by maximizing the so called Raleigh coefficient

$$\max J = \frac{J_b(w)}{J_w(w)} = \frac{w'S_b w}{w'S_w w}. \quad (13)$$

The matrices for between and within class variance are defined as

$$S_b = \sum_{c=1}^C (\mu_c - \mu)(\mu_c - \mu)' \quad (14)$$

$$S_w = \sum_{c=1}^C \sum_{i \in c} (x_i - \mu_c)(x_i - \mu_c)' \quad (15)$$

where C is the number of classes, μ_c is the vector of sample means for each class respectively and μ is the vector of sample means for the full dataset. For identification purposes we can always chose weights w such that $w'S_w w = 1$, since J is constant with regards to rescalings. We can therefore replace w by αw which will result in the constant α canceling out. This way the initial optimization problem can be formulated as

$$\arg \min_w -\frac{1}{2} w'S_b w \quad s.t. \quad w'S_w w = 1 \quad (16)$$

with the lagrangian being

$$\mathcal{L} = -\frac{1}{2} w'S_b w + \frac{1}{2} \lambda (w'S_w w - 1). \quad (17)$$

The halves are added for more convenient matrix derivatives. The Karush-Kuhn-Tucker conditions imply that the solution to this maximization problem and subsequently the vector of weights we want to find needs to fulfill

$$S_b w = \lambda S_w w. \quad (18)$$

This is a generalized eigenvalue problem for which there exists a convenient R-solution in the form of the `geigen`-package.

2.2 Implementation

The result from 18, which is essentially the rotation of the underlying data's column space that ensures the highest separability, is implemented in the `lda`-class for the two-class case. The `lda()` function accepts two arguments:

1. **data:** a `data.frame` containing at least one column of factors indicating the class label.
2. **by:** a `character`-string equal to the column's name containing the class labels. Note that all other non-numerical columns will be ignored by the function, since LDA is only meaningful for continuous variables.

The function in a first step extracts the useable columns and the number of classes provided in **data**. It then performs a quick check if the prerequisites are met and then continues with the calculation of the class-means μ_1 and μ_2 , the overall mean μ , as well as the scatter-matrices S_b and S_w . With these we can solve the generalized eigenvalue problem from 18.

Algorithm 4 `lda()`

```

1: procedure SET UP AUXILIARY VARIABLES
2:   num
3:   classes
4:   mu
5:   x_bar
6:   S

```

3 Appendix

3.1 Unit Tests

References

- [1] Winkelmann, R., Boes, S. (2009): "*Analysis of Microdata*", 2nd edition.