

MGLM package vignette

Yiwen Zhang Hua Zhou

July 25, 2013

The analysis of multivariate count data arises in numerous fields including genomics, image analysis, text mining, and sports analytics. The multinomial logit model is limiting due to its restrictive mean-variance structure. Moreover, it assumes that counts of different categories are negatively correlated. Models that allow over-dispersion and possess more flexible positive and/or negative correlation structures offer more realism. We implement four models in the R package **MGLM**: multinomial logit (MN), Dirichlet multinomial (DM), generalized Dirichlet multinomial (GDM), and negative multinomial (NegMN). Distribution fitting, regression, hypothesis testing, and variable selection are treated in a unified framework.

The simulated data we plan to analyze is multivariate count data with d categories.

```
> require(MGLM)
```

1 Distribution fitting

The function **MGLMfit** fits various multivariate discrete distributions and outputs a list with the maximum likelihood estimate (MLE) and relevant statistics.

When fitting distributions, i.e. no covariates involved, MN is a sub-model of DM, and DM is a sub-model of GDM. **MGLMfit** outputs the p-value of the likelihood ratio test (LRT) for comparing the fitted model with the most commonly used multinomial model. NegMN model does not have a nesting relationship with any of the other three models. Therefore no LRT is performed when fitting a NegMN distribution.

1.1 Multinomial (MN)

We first generate data from a multinomial distribution. Note the multinomial parameter (must be positive) supplied to the ‘rmn’ function is automatically scaled to be a probability vector.

```
> set.seed(123)
> n <- 200
> d <- 4
> alpha <- rep(1,d)
> m <- 50
> Y <- rmn(m, alpha, n)
```

Fitting the DM distribution to the multinomial data shows no advantage.

```
> mnFit <- MGLMfit(Y, dist="DM")
> print(mnFit)
```

	estimate	SE
alpha_1	5270901	786946676
alpha_2	5063596	755995891
alpha_3	5030755	751092797
alpha_4	5160065	770398732

```
Distribution: Dirichlet Multinomial
Log-likelihood: -1457.788
BIC: 2936.769
AIC: 2923.576
LRT test p value: 1
Iterations: 35
```

The DM parameter estimates and their standard errors are both big, indicating possible overfitting by the DM model. This is confirmed by the fact that the p-value of the LRT for comparing MN to DM is close to 1.

1.2 Dirichlet-multinomial (DM)

DM is a Dirichlet mixture of multinomials and allows over-dispersion. Same as MN model, it assumes that the counts of any two different categories are negatively correlated. We generate the data from the DM model and fit distribution.

```
> set.seed(123)
> n <- 200
> d <- 4
> alpha <- rep(1, d)
> m <- 50
> Y <- rdirm(m, alpha, n)

> dmFit <- MGLMfit(Y, dist="DM")
> print(dmFit)
```

	estimate	SE
alpha_1	0.9766705	0.07658856
alpha_2	0.9951423	0.07925470
alpha_3	1.0061205	0.08003311
alpha_4	0.9045003	0.07254733

```
Distribution: Dirichlet Multinomial
Log-likelihood: -2011.225
BIC: 4043.644
AIC: 4030.451
LRT test p value: 0
Iterations: 4
```

The estimate is very close to the true value with small standard errors. The LRT shows that the DM model is significantly better than the MN model.

1.3 Generalized Dirichlet-multinomial (GDM)

GDM model uses $d - 2$ more parameters than the DM model and allows both positive and negative correlations among categories. DM is a sub-model of GDM. Here we fit a GDM model to the above DM sample.

```
> gdmFit <- MGLMfit(Y, dist="GDM")
> print(gdmFit)
```

	estimate	SE
alpha_1	1.1584741	0.12340343
alpha_2	0.9932931	0.43723944
alpha_3	0.8399666	0.10637444
beta_1	3.7068630	0.22641418
beta_2	1.9793891	0.09464475
beta_3	0.7596409	0.08440347

```
Distribution: Generalized Dirichlet Multinomial
Log-likelihood: -2007.559
BIC: 4046.907
AIC: 4027.117
LRT test p value: 0
Iterations: 19
```

GDM yields a slightly larger log-likelihood value but a larger BIC, suggesting DM as a preferred model. Now we simulate data from GDM and fit the GDM distribution.

```
> set.seed(124)
> n <- 200
> d <- 4
> alpha <- rep(1, d-1)
> beta <- rep(1, d-1)
> m <- 50
> Y <- rgdirm(m, alpha, beta, n)
> gdmFit <- MGLMfit(Y, dist="GDM")
> print(gdmFit)
```

	estimate	SE
alpha_1	1.0198346	0.10348010
alpha_2	0.8261251	0.10931634
alpha_3	0.7743869	0.09256818
beta_1	1.0621114	0.09556607
beta_2	0.8462160	0.10872545
beta_3	0.9245594	0.13500769

```
Distribution: Generalized Dirichlet Multinomial
Log-likelihood: -1820.616
BIC: 3673.021
AIC: 3653.231
LRT test p value: 0
Iterations: 18
```

1.4 Negative multinomial (NegMN)

NegMN model is a multivariate extension to the negative binomial model. It assumes positive correlation among the counts. To generate data from NegMN model and fit the NegMN distribution,

```
> set.seed(1220)
> n <- 100
> d <- 4
> p <- 5
> prob <- rep(0.2, d)
> beta <- 10
> Y <- rnegmn(prob, beta, n)
> negmnFit <- MGLMfit(Y, dist="NegMN")
> print(negmnFit)
```

	estimate	SE
p_1	0.1881512	0.009583840
p_2	0.1943109	0.009837429
p_3	0.1915110	0.009722206
p_4	0.1961775	0.009914205
phi	12.3139348	2.266096911

```
Distribution: Negative Multinomial
Log-likelihood: -1104.579
BIC: 2232.184
AIC: 2219.158
LRT test p value:
Iterations: 4
```

2 Regression

In regression, the $n \times p$ covariate matrix X is similar to that used in the `glm` function. The response should be a $n \times d$ count matrix. Unlike estimating a parameter vector β in GLM, we need to estimate a parameter matrix B when the responses are multivariate. The dimension of the parameter matrix depends on the model:

- MN: $p \times (d - 1)$
- DM: $p \times d$
- GDM: $p \times 2(d - 1)$
- NegMN: $p \times (d + 1)$

The GDM model provides the most flexibility, but also requires most parameters. When using function `MGLMreg` to run regression, we can pick the model by specifying the option `dist="MN"`, `"DM"`, `"GDM"` or `"NegMN"`.

The rows $B_{j\cdot}$ of the parameter matrix correspond to covariates. By default, the function output the Wald test statistics and p-values for testing $H_0 : B_{j\cdot} = \mathbf{0}$

vs $H_a : B_{j\cdot} \neq \mathbf{0}$. If specifying the option `LRT=TRUE`, the function also outputs LRT statistics and p-values.

Next we demonstrate that model mis-specification results in failure in hypothesis testing. We simulate response data from the GDM model. Covariates X_1 and X_2 have no effect and X_3, X_4, X_5 have impact on the response.

```
> set.seed(1234)
> n <- 200
> p <- 5
> d <- 4
> X <- matrix(runif(p*n), n, p)
> alpha <- matrix(c(0.6, 0.8, 1), p, d-1, byrow=TRUE)
> alpha[c(1,2),] <- 0
> Alpha <- exp(X%%alpha)
> beta <- matrix(c(1.2, 1, 0.6), p, d-1, byrow=TRUE)
> beta[c(1,2),] <- 0
> Beta <- exp(X%%beta)
> m <- runif(n, min=0, max=25) + 25
> Y <- rgdirm(m, Alpha, Beta)
```

We fit various regression models and test significance of covariates.

2.1 Multinomial regression

```
> mnReg <- MGLMreg(Y~0+X, dist="MN")
> print(mnReg)
```

Call: `MGLMreg(formula = Y ~ 0 + X, dist = "MN")`

Coefficients:

	Col_1	Col_2	Col_3
X1	0.2770632	-0.1827597	-0.1232039
X2	0.5430639	0.4227301	0.2465230
X3	0.3332517	0.5176055	0.2218513
X4	0.3568425	0.4867224	0.5654272
X5	-0.3024545	0.1925076	0.3237132

Hypothesis test:

	wald value	Pr(>wald)
X1	24.63244	1.842859e-05
X2	21.99680	6.533133e-05
X3	23.10908	3.832310e-05
X4	25.07475	1.489470e-05
X5	49.37327	1.086326e-10

Distribution: Multinomial

Log-likelihood: -2194.448

BIC: 4468.371

AIC: 4418.896

Iterations: 5

The Wald test shows that all predictors are significantly different from 0, including the null predictors X_1 and X_2 .

2.2 Dirichlet-multinomial regression

```
> dmReg <- MGLMreg(Y~0+X, dist="DM")
> print(dmReg)
```

Call: MGLMreg(formula = Y ~ 0 + X, dist = "DM")

Coefficients:

	Col_1	Col_2	Col_3	Col_4
X1	0.1541366	-0.1182637	-0.1883392	-0.01317227
X2	0.1832643	0.1420339	-0.1833942	-0.33388595
X3	1.1431455	1.2548275	1.0926350	0.81125857
X4	0.3927026	0.5454212	0.5900143	0.13113141
X5	0.2263497	0.6601082	0.9395989	0.48703441

Hypothesis test:

	wald value	Pr(>wald)
X1	3.349794	5.010817e-01
X2	7.845338	9.741080e-02
X3	25.497385	3.995531e-05
X4	8.735126	6.807201e-02
X5	23.136043	1.189431e-04

Distribution: Dirichlet Multinomial

Log-likelihood: -1683.961

BIC: 3473.889

AIC: 3407.922

Iterations: 6

Again, Wald test declares all predictors as significant.

2.3 Generalized Dirichlet-multinomial Regression

```
> gdmReg <- MGLMreg(Y~0+X, dist="GDM")
> print(gdmReg)
```

Call: MGLMreg(formula = Y ~ 0 + X, dist = "GDM")

Coefficients:

	alpha_Col_1	alpha_Col_2	alpha_Col_3	beta_Col_1	beta_Col_2	beta_Col_3
X1	-0.1602454	-0.2515410	0.3157057	-0.1540508	-0.1247925	0.4675919
X2	-0.2773313	0.4707074	0.0144257	-0.3115543	0.2541531	-0.1714098
X3	1.2432626	1.3371596	1.1717863	1.7556099	1.4185048	0.9160326
X4	0.3099390	0.6977423	0.8090867	0.6693109	0.7894918	0.3058891
X5	0.6649190	0.1160624	1.2812194	1.6628817	0.4753575	0.8534991

Hypothesis test:

	wald value	Pr(>wald)
--	------------	-----------

```

X1  2.078755 9.123148e-01
X2  3.945574 6.840417e-01
X3  35.526568 3.406618e-06
X4  17.732994 6.935276e-03
X5  55.903886 3.044167e-10

```

```

Distribution: Generalized Dirichlet Multinomial
Log-likelihood: -1658.646
BIC: 3476.242
AIC: 3377.292
Iterations: 20

```

When using the correct model, Wald test is able to differentiate the null effects from the significant ones. GDM regression yields the highest log-likelihood and smallest BIC.

2.4 Negative multinomial regression

```

> negReg <- MGLMreg(Y~0+X, dist="NegMN", regBeta=FALSE)
> print(negReg)

```

```

Call: MGLMreg(formula = Y ~ 0 + X, dist = "NegMN", regBeta = FALSE)

```

Coefficients:

```

$alpha
      Col_1      Col_2      Col_3      Col_4
X1  0.24359303 -0.21636491 -0.15652287 -0.03138873
X2  0.06183009 -0.05594709 -0.23269725 -0.47983600
X3 -0.16090208  0.02269929 -0.27122053 -0.49512992
X4 -0.17621570 -0.04386905  0.03473824 -0.53015978
X5 -0.60796591 -0.11582733  0.01293004 -0.31587585

```

\$phi

```

13.77597

```

Hypothesis test:

```

      wald value      Pr(>wald)
X1   14.74648 1.150232e-02
X2   26.10450 8.516656e-05
X3   35.24829 1.342363e-06
X4   28.34379 3.117744e-05
X5   51.20854 7.838372e-10

```

```

Distribution: Negative Multinomial
Log-likelihood: -2908.896
BIC: 5929.056
AIC: 5859.792
Iterations: 15

```

Again, the Wald test declares all predictors to be significant.

2.5 Prediction

We can use the fitted model to make prediction. The output of the prediction function is the probabilities of the d categories. This helps answer questions such as whether certain features increase the probability of observing category j . Take the fitted GDM model as an example:

```
> newX <- matrix(runif(1*p), 1, p)
> pred <- predict(gdmReg, newX)
> pred
```

	Col_1	Col_2	Col_3	Col_4
[1,]	0.2023061	0.3336253	0.33736	0.1267085

3 Sparse regression

Regularization is an important tool for model selection and improving the risk property of the estimates. In the package, we implemented three types of penalties on the parameter matrix B :

- select by entries
- select by rows
- select by rank

The function `MGLMtune` finds the optimal tuning parameter with smallest BIC and outputs the estimate using the chosen tuning parameter. The output from `MGLMtune` is a list containing the solution path and the final estimate. Users can either provide a vector of tuning parameters with option `lambdas` or specify the number of grid points via option `ngridpt` and let the function decide the default tuning parameters. The function `MGLMsparsereg` computes the regularized estimate at a given tuning parameter value `lambda`.

We generate the data from the DM model, with row sparsity, and show how each penalty type works.

```
> n <- 100
> p <- 10
> d <- 5
> m <- runif(n, min=0, max=25) + 25
> set.seed(134)
> X <- matrix(rnorm(n*p), n, p)
> alpha <- matrix(0, p, d)
> alpha[c(1,3, 5), ] <- 1
> Alpha <- exp(X%*%alpha)
> Y <- rdirn(size=m, alpha=Alpha)
```

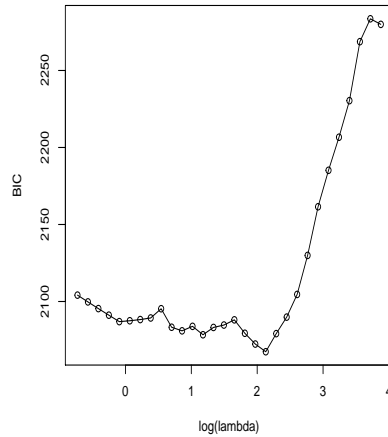



Figure 1: Variable selection by entries

3.1 Select by entries

Figure 1 displays the trace of BIC along the solution path and the regularized estimate at optimal tuning parameter value.

```
> sweep <- MGLMtune(Y~0+X, dist="DM", penalty="sweep", ngridpt=30)
> print(sweep$select)
```

```
Call: MGLMtune(formula = Y ~ 0 + X, dist = "DM", penalty = "sweep",
  ngridpt = 30)
```

```
Distribution: Dirichlet Multinomial
Log-likelihood: -978.428
BIC: 2067.38
AIC: 2004.856
Degrees of freedom: 24
Lambda: 8.416109
Iterations: 33
```

3.2 Select by rows

Since the rows of the parameter matrix correspond to predictors, selecting by rows performs variable selection at the predictor level. Figure 2 displays the trace of BIC along the solution path and the regularized estimate at optimal tuning parameter value.

```
> group <- MGLMtune(Y~0+X, dist="DM", penalty="group", ngridpt=30)
> print(group$select)
```

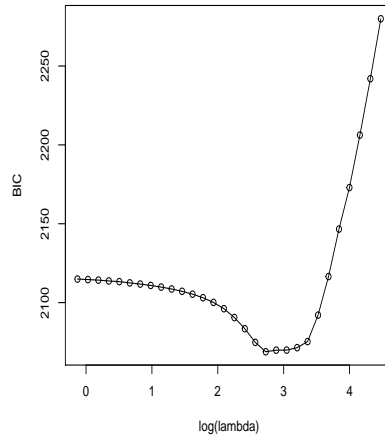


Figure 2: Variable selection by groups

```
Call: MGLMtune(formula = Y ~ 0 + X, dist = "DM", penalty = "group",
  ngridpt = 30)
```

```
Distribution: Dirichlet Multinomial
Log-likelihood: -971.2706
BIC: 2068.865
AIC: 1997.403
Degrees of freedom: 27.43079
Lambda: 15.30747
Iterations: 35
```

3.3 Select by singular values

Nuclear norm regularization encourages low rank in the regularized estimate. Figure 3 displays the trace of BIC along the solution path and the regularized estimate at optimal tuning parameter value.

```
> nuclear <- MGLMtune(Y~0+X, dist="DM", penalty="nuclear", ngridpt=30)
> print(nuclear$select)
```

```
Call: MGLMtune(formula = Y ~ 0 + X, dist = "DM", penalty = "nuclear",
  ngridpt = 30)
```

```
Distribution: Dirichlet Multinomial
Log-likelihood: -987.2953
BIC: 2056.726
AIC: 2010.261
Degrees of freedom: 17.83546
Lambda: 36.06196
Iterations: 19
```

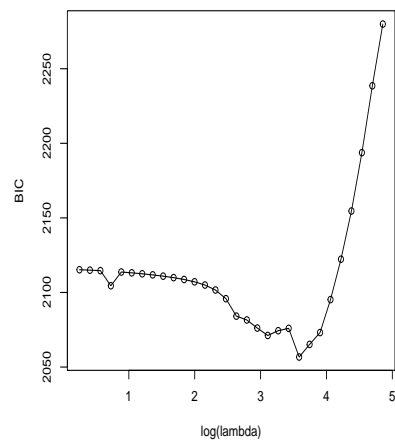


Figure 3: Variable selection by singular values