

Logistic Regression

Edps 590BAY

Carolyn J. Anderson

Department of Educational Psychology



©Board of Trustees, University of Illinois

Spring 2018

I Overview

- ▶ Logistic regression
- ▶ MLE estimation: GMAT data
- ▶ Bayesian estimation: GMAT data
- ▶ Various computational options:
 - ▶ thinning.
 - ▶ multiple cores (faster).
 - ▶ miscellaneous
- ▶ Quasi and complete separation.
- ▶ HSB data?

I As a GLM

- ▶ Random component: The response variable is **dicotomous**, coded $Y_i = 1$ or 0 . The distribution of Y_i is **Binomial** and we are interesting in probability that $Y_i = 1$, i.e., $Pr(x_i)$.
- ▶ Systematic component: A linear predictor such as

$$b_0 + b_1x_{1i} + \dots + b_kx_{ki}$$

The explanatory or predictor variables may be quantitative (continuous), qualitative (discrete or ordinal), or both (mixed).

- ▶ Link Function: The log of the odds that an event occurs, otherwise known as the **logit**:

$$\text{logit}(Pr(y_i = 1)) = \log\left(\frac{y_i = 1}{1 - Pr(y_i = 1)}\right)$$

Putting this all together, the logistic regression model is

$$\text{logit}(\pi(x_i)) = \log\left(\frac{\pi(x_i)}{1 - \pi(x_i)}\right) = b_0 + b_1x_{1i} + \dots + b_kx_{ki}$$

I Some Uses of Logistic Regression

- ▶ Identify “risk” factors for certain conditions (diabetes, divorce, etc.)
- ▶ Descriptive discriminant analysis—describe differences between groups.
- ▶ Adjust for “bias” (propensity score analysis/matching).
- ▶ Predict probabilities.
- ▶ Classify individuals.
- ▶ Discrete choice.
- ▶ Social network analysis.
- ▶ Pseudo-likelihood estimation.
- ▶ ... and many others

I Interpreting logistic regression models

The model

$$\text{logit}(Pr(y_i = 1)) = b_0 + b_1 x_i$$

Equivalently,

$$Pr(y_i = 1) = \frac{\exp[b_0 + b_1 x_i]}{1 + \exp[b_0 + b_1 x_i]} = \frac{1}{1 + \exp[-(b_0 + b_1 x_i)]}$$

We'll be using GMAT data (from Johnson & Wichern text) where y is admission to medical school as a function of GMAT scores and college GPA.

- ▶ Response: $Y_i = 1$ if the student is admitted, and 0 if the student was denied to wait listed (borderline).
- ▶ Explanatory: x_{1i} = student GMAT score and x_{2i} is college GPA.

I MLE Estimation

```
logreg.mle ← glm(admit ~ gmat, data=med, family=binomial)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.5996	-0.5439	-0.2154	0.4038	2.4350

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-17.397141	3.758928	-4.628	3.69e-06 ***
gmat	0.033609	0.007386	4.551	5.35e-06 ***

Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1

Dispersion parameter for binomial family taken to be 1

Null deviance: 111.533 on 84 degrees of freedom

Residual deviance: 60.106 on 83 degrees of freedom

AIC: 64.106

Number of Fisher Scoring iterations: 6

I Interpretation

The odds of admission given GMAT score = $gmat_i$ are

$$\frac{Pr(admission|gmat_i)}{Pr(deny|gmat_i)} = \exp[(-17.3971 + 0.0336(gmat_i))]$$

and when $gmat_i + 1$,

$$\frac{Pr(admission|(gmat_i + 1))}{Pr(deny|(gmat_i + 1))} = \exp[(-17.3971 + 0.0336(gmat_i + 1))]$$

The ratio is the odds ratio:

$$\begin{aligned} & \frac{Pr(admit|gmat_i + 1)/Pr(not\ admit|gmat_i + 1)}{Pr(admit|gmat_i)/Pr(not\ admit|gmat_i)} \\ &= \frac{\exp[(-17.3971 + 0.0336(gmat_i + 1))]}{\exp[(-17.3971 + 0.0336(gmat_i))]} = \exp(0.0336) \end{aligned}$$

I Interpretation

Our odds ratio equals $\exp(0.0336) = 1.034171$

the odds admission given $gmat_i + 1 = 1.0342$ odds given $gmat_i$

In words, the **odds** of admission is 0.0336 times the odds of admission for one point score lower on GMAT. Note that the standard deviation of GMAT is 81.52235. Using this,

the odds admission given $gmat_i + 81.52235 = 15.47384 \exp[\text{the odds given } gmat_i]$

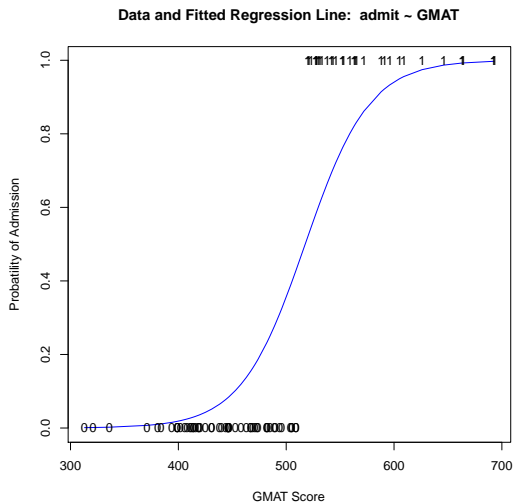
In words, the **odds** of admission is 15.47 times the odds of admission for one standard deviation lower on GMAT.

Can also say “the probability of admission is higher the greater a student’s GMAT score”

or “it is more likely than a student will be admitted if they have a higher GMAT score”

You should **NOT** say “the probability of admission is 1.03 times higher the greater a student’s GMAT score”

I Fitted Probability of Admission



I DataList

```
dataList ← list(  
  y=med$admit,  
  gmat=med$gmat,  
  n=length(med$admit)  
)
```

I JAGS Model

```
logreg1 ← ‘‘model {  
  for (i in 1:n) {  
    y[i] ~ dbern(p[i])  
    p[i] ← 1/(1 + exp(-eta[i]))  
    eta[i] ← b0 + b1*gmat[i]  
  }  
  b0 ~ dnorm(0,1/1000)  
  b1 ~ dnorm(0,1/1000)  
}’’
```

```
writeLines(logreg1,con="logreg1.txt")
```

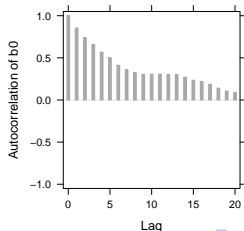
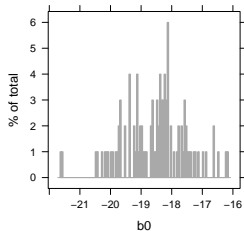
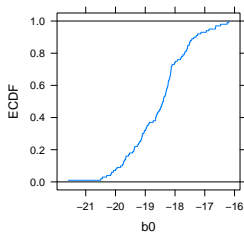
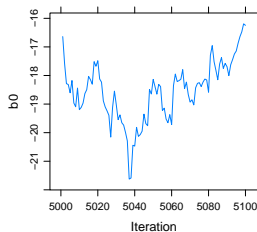
I runjags to check model set up

```
start1 ← list("b0"=1,"b1"=-.04)

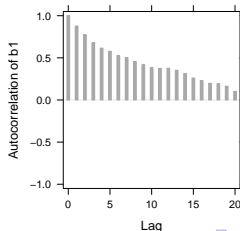
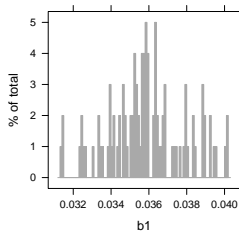
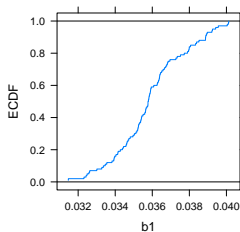
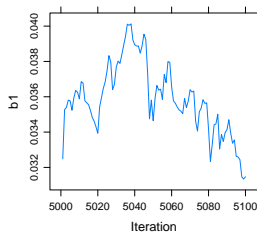
logreg1.chk ← run.jags(
  model=logreg1,
  sample=100,
  data=dataList,
  inits=start1,
  monitor=c("b0","b1", "dic"),
  n.chains=1
)

print(logreg1.chk)
plot(logreg1.chk)
```

I From Checking Set-Up



I From Checking Set-Up



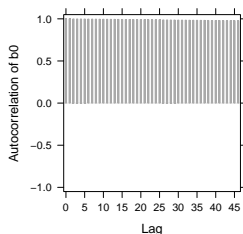
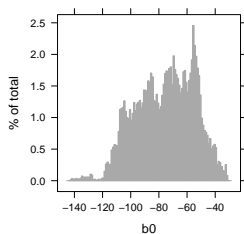
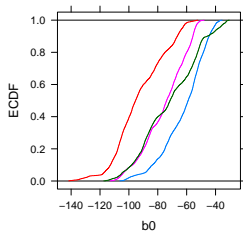
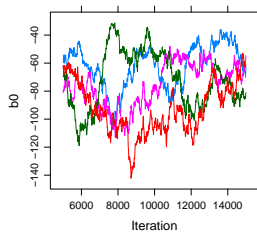
I Now Run to Convergence

```
start ← list(  
  list("b0"=1,"b1"=-.04),  
  list("b0"=-7,"b1"=.04),  
  list("b0"=5,"b1"=-.02),  
  list("b0"=-1,"b1"=-0.2)  
)
```

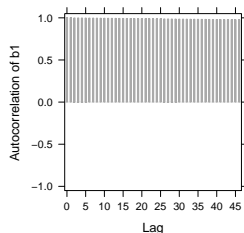
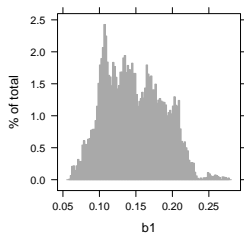
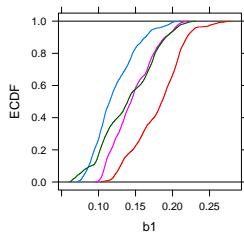
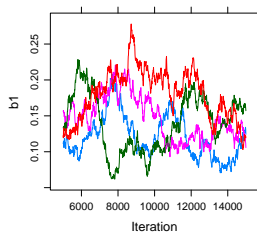
I Now Run to Convergence

```
logreg1.jags ← run.jags(  
  model=logreg1,  
  data=dataList,  
  inits=start,  
  monitor=c("b0","b1","dic"),  
  n.chains=4  
)  
  
print(logreg1.jags)  
plot(logreg1.jags)
```


I Looks Pretty Bad



I Looks Pretty Bad



I Fix & Speed up

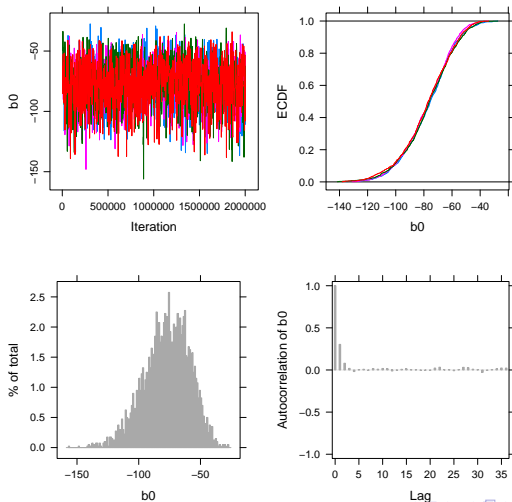
```
# So I don't get warning message (if I want to use parallel)
```

```
start1 ← list("b0"= 1,"b1"=-.04, "b2"=0,  
             .RNG.name="base::Wichmann-Hill", .RNG.seed=23)  
start2 ← list("b0"= 7,"b1"= .04, "b2"=-.5,  
             .RNG.name="base::Marsaglia-Multicarry",  
             .RNG.seed=97)  
start3 ← list("b0"= 5,"b1"=-.02, "b2"=0.5,  
             .RNG.name="base::Super-Duper", .RNG.seed=3351)  
start4 ← list("b0"=-1,"b1"=-0.2, "b2"=.2,  
             .RNG.name="base::Mersenne-Twister",  
             .RNG.seed=7531)
```

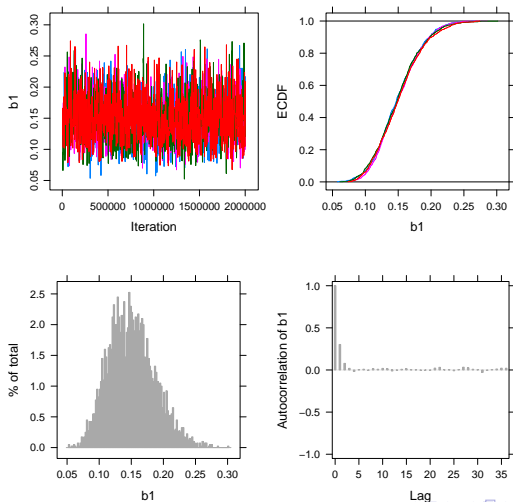
I Fix (and speed up)

```
logreg1.faster leftarrow run.jags (  
  method="parallel",  
  summarise=FALSE,  
  plots=FALSE,  
  model=logreg1,  
  data=dataList,  
  inits=list(start1,start2,start3,start4),  
  sample=100000,  
  monitor=c("b0","b1"),  
  n.chains=4,  
  thin=20  
)  
add.summary(logreg1.faster)  
plot(logreg1.faster)
```

I Better



I Better



I Statistics

JAGS model summary statistics from 400000 samples (thin = 20; chains = 4; adapt+burnin = 5000):

	Lower95	Median	Upper95	Mean	SD	Mode
b0	-114.49	-76.158	-42.598	-77.383	18.763	-
b1	0.082886	0.14817	0.22246	0.15056	0.036439	-

	MCerr	MC%ofSD	SSeff	AC.200	psrf
b0	0.38809	2.1	2338	0.88931	1.0006
b1	0.00075426	2.1	2334	0.88931	1.0008

Total time taken: 6.9 minutes

... It took a long time because of large number of iterations

I Model Evaluation for Logistic Regression

These are my favorites, which can be used with Bayesian as well as MLE: Receiver Operating Characteristic (ROC), Classification tables, and the Concordance index.

Suppose we have a simple model

$$\text{logit}(\widehat{Pr}(y_i = 1)) = \hat{b}_0 + \hat{b}_1 x_i$$

Let π_o be a cut-point or cut-score and $\widehat{Pr}(y_i = 1)$ be a predicted probability of the model. The predicted response is

$$\hat{y}_i = \begin{cases} 1 & \text{if } \widehat{Pr}(y_i = 1) > \pi_o \\ 0 & \text{if } \widehat{Pr}(y_i = 1) \leq \pi_o \end{cases}$$

Classification Table:

		Predicted	
		$\hat{y}_i = 1$	$\hat{y}_i = 0$
Actual	$y = 1$	correct	incorrect (false negative)
	$y = 0$	incorrect (false positive)	correct

I Classification Table

We're more interested in conditional proportions and probabilities:

		Predicted		
		$\hat{y}_i = 1$	$\hat{y}_i = 0$	
Actual	$y = 1$	n_{11}/n_{1+}	n_{12}/n_{1+}	n_{1+}
	$y = 0$	n_{21}/n_{2+}	n_{22}/n_{2+}	n_{2+}

$$n_{11}/n_{1+} = \text{proportion } (\hat{y} = 1|y = 1) = \text{"sensitivity"}$$

$$n_{22}/n_{2+} = \text{proportion } (\hat{y} = 0|y = 0) = \text{"specificity"}$$

$$\begin{aligned} p(\text{correct}) &= p(\hat{y} = 1 \ \& \ y = 1) + p(\hat{y} = 0 \ \& \ y = 0) \\ &= p(\hat{y} = 1|y = 1)p(y = 1) + (\hat{y} = 0|y = 0)p(y = 0) \\ &= (\text{sensitivity})p(y = 1) + (\text{specificity})p(y = 0) \end{aligned}$$

I GMAT Example

- ▶ Let the cut-score equal $\pi_o = .50$.
- ▶ Compare $p(\hat{y} = 1)$ and classify i as $Y = 1$ if $p(\hat{y} = 1) > \pi_o$, otherwise classify i as $Y = 0$.
- ▶ Tabulate the results

		Predicted		
		$\hat{y} = 1$	$\hat{y} = 0$	
Actual	$y = 1$	25	6	31
	$y = 0$	6	48	52

- ▶ The Conditional proportions:

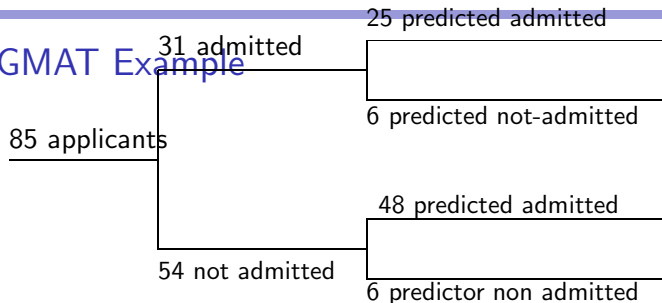
$$\text{Sensitivity} = 25/31 = .81$$

$$\text{Specificity} = 48/52 = .92$$

- ▶ The proportion correct = $(25 + 48)/85 = .86$

Sensitivity = .81%

I GMAT Example



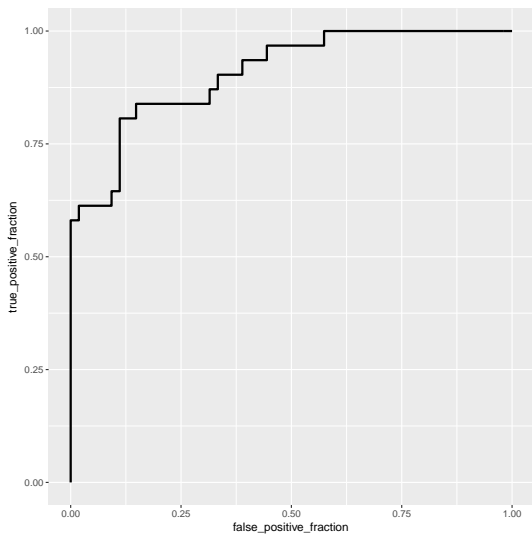
Specificity = .92%

Percent correct = .86%

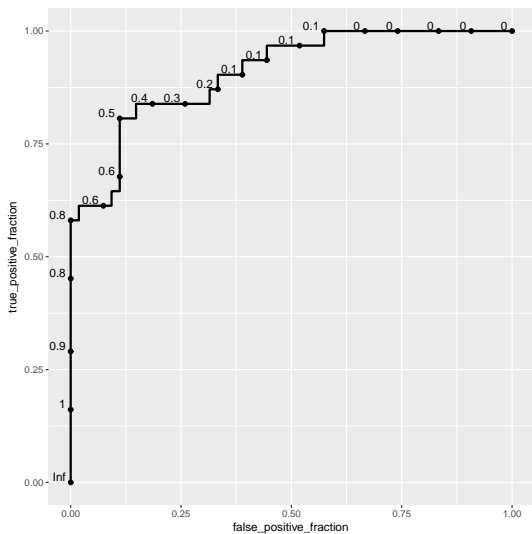
I Sensitivity, Specificity & $p(\text{Correct})$

- ▶ For every cut-score you will get a different result.
- ▶ Do this for lots of possible cut-scores and plot the results \rightarrow ROC curve.

I ROC curve



I ROC curve



I Area Under ROC Curve

Concordance: Take two cases i and j where $y_i = 1$ and $y_j = 0$ ($i \neq j$),

If $\hat{P}_r(y_i = 1) > \hat{P}_r(y_j = 1)$, then the pair is concordant

If $\hat{P}_r(y_i = 1) < \hat{P}_r(y_j = 1)$, then the pair is discordant

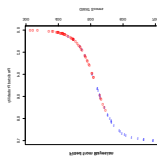
If $\hat{P}_r(y_i = 1) = \hat{P}_r(y_j = 1)$, then the pair is tie

The area under the ROC curve equals the concordance index. The concordance index is an estimate of the probability that predictions and outcomes are concordant.

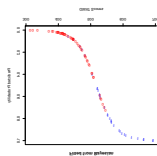
Our example, $c = 0.91$

This also provides a way to compare models, the solid dots in the next figure are from a model with more predictors. For example, . . .

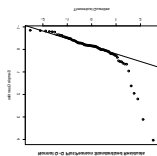
I Fitted Probabilities



I Fitted Probabilities & Data



I QQ Plot of Standardized Residuals



I admit \sim GMAT + GPA

```
logreg2 ← "model {  
  for (i in 1:n) {  
    y[i] ~ dbern(p[i])  
    p[i] ← 1/(1 + exp(-eta[i]))  
    eta[i] ← b0 + b1*gmat[i] + b2*gre[i]  
  }  
hspace.5in b0 ~ dnorm(0,1/1000)  
b1 ~ dnorm(0,1/1000)  
b2 ~ dnorm(0,1/1000)  
}"
```

I admit \sim GMAT + GPA

```
logreg2.runjags ← run.jags (  
  method="parallel",  
  summarise=FALSE,  
  plots=FALSE,  
  sample=500000,  
  model=logreg2,  
  data=dataList,  
  inits=start,  
  monitor=c("b0", "b1", "b2"),  
  n.chains=4,  
  thin=20  
)
```

I admit \sim GMAT + GPA

Calculating summary statistics...

Calculating the Gelman-Rubin statistic for 3 variables...

. JAGS model summary statistics from 2,000,000 samples (thin = 20; chains = 4; adapt+burnin = 5000):

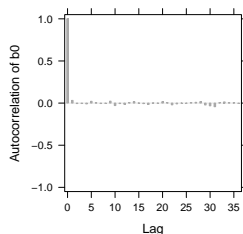
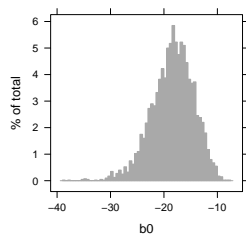
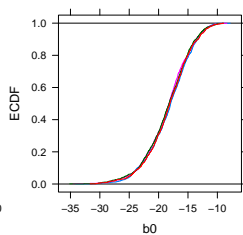
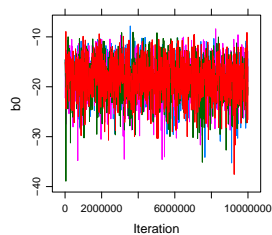
Note: Because I ran 4 chains with thin=20, the number of iterations was actually $500,000 \times 20 = 10,000,000$ per chain!

	Lower95	Median	Upper95	Mean	SD	Mode
b0	-26.101	-18.029	-11.082	-18.311	3.8781	-
b1	-3.8507	0.26955	2.9303	0.11768	1.7769	-
b2	-2.8994	-0.23269	3.8824	-0.082296	1.777	-

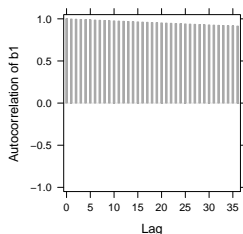
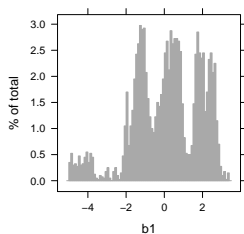
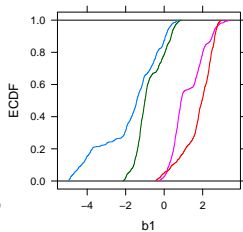
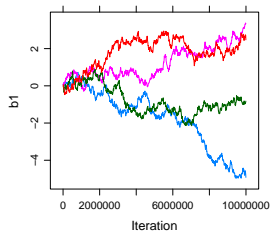
	MCerr	MC%fSD	SSeff	AC.200	psrf
b0	0.014642	0.4	70155	0.34361	1
b1	0.51228	28.8	12	0.99987	2.3686
b2	0.51484	29	12	0.99987	2.3685

Total time taken: 41.6 minutes

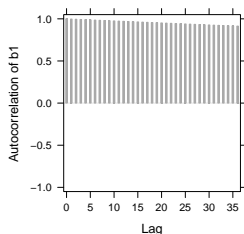
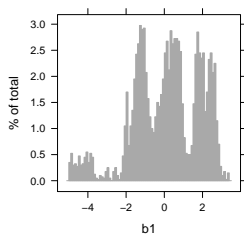
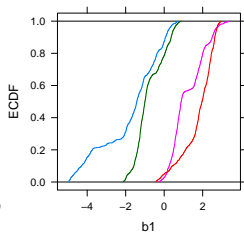
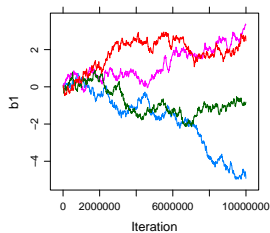
I Trace, Density and Auto-Correlation



I Really Bad



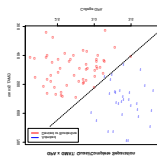
I Really bad



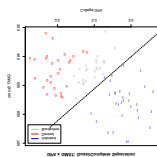
I What's Going On

- ▶ $r(gmat, gpa) = .46$
- ▶ From glm:
"Warning message:
glm.fit: fitted probabilities numerically 0 or 1 occurred "
- ▶ Could centering help? Not here.
- ▶ Case of quasi/complete separation?

I This is a “Problem”



I This is a “Problem”



I How Long Things Took

Note: I was typing and doing other things while models ran, so times are larger than would take if I wasn't typing and surfing.

	time	Notes
Check	6.1 sec	No errors, terrible fit
Increase iterations		Not good enough
thin=20	2.4 min	Better
+parallel	40-45 sec	OK
+drop extra	40-43 sec	OK
2 predictors		
lots of iterations	≈ 50 mintues	Terrible

- ▶ Parallel really speeds things up, but tied up computer.
- ▶ Adding iterations and emp.new really slows things down.
- ▶ Better starting values requires shorter chains.
- ▶ Expect long times as increase number of parameters and more complex models.
- ▶ Stan is slower than JAGS, but has other advantages.
- ▶ Illinois campus cluster: <https://campuscluster.illinois.edu>, which does have latest version of R.