

# Today

- ▶ Bayesian analysis of logistic regression
- ▶ Generalized linear mixed models
- ▶ CD on fixed and random effects
- ▶ HW 2 due February 28
- ▶ Case Studies [SSC 2014 Toronto](#)
- ▶ March/April: Semi-parametric regression (§10.7), generalized additive models, penalized regression methods (ridge regression, lasso); proportional hazards models (§10.8)

# Bayesian logistic regression

- ▶  $r_j \sim \text{Binom}(m_j, p_j)$
- ▶  $\log \frac{p_j}{1 - p_j} = \alpha + \beta x_j$
- ▶  $L(\alpha, \beta; \mathbf{y}) \propto \exp\{\alpha \sum y_j + \beta \sum (x_j y_j) - \sum m_j \log(1 + e^{\alpha + \beta x_j})\}$
- ▶  $\pi(\alpha, \beta | \mathbf{y}) \propto L(\alpha, \beta; \mathbf{y})\pi(\alpha, \beta)$
- ▶ flat prior  $\pi(\alpha, \beta) \propto 1$  popular for regression parameters  
proper posterior?
- ▶ implemented in the library `LearnBayes` via `logisticpost` Albert, 2009 *Bayesian Computation with R*

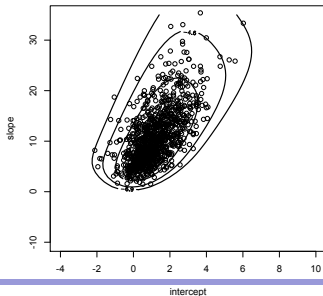
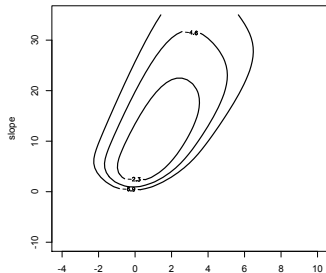
- ▶ bioassay data

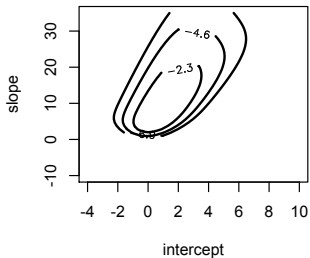
log(dose)	deaths	sample size
-0.86	0	5
-0.30	1	5
-0.05	3	5
0.73	5	5



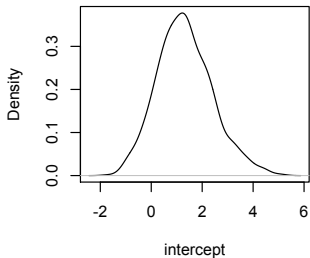
## ... Bayesian logistic regression

```
mycontour(logisticpost,c(-4,10,-10,35),bioassay)
## the limits were chosen using information in Gelman et al.,
## although they used 40 as the upper limit for beta, but I could not
> s <- simcontour(logisticpost, c(-4,10, -10, 35),m = 1000, data =bioassay)
> points(s) # just plotted 1000 points, otherwise plot is too black
> s <- simcontour(logisticpost, c(-4,10, -10, 35),m = 10000, data =bioassay)
# samples from posterior; more samples for getting quantiles
> quantile(s$x, c(0.025, 0.5, 0.975))
      2.5%      50%      97.5%
-0.6066507  1.2277272  3.7397231
> quantile(s$y, c(0.025, 0.5, 0.975))
      2.5%      50%      97.5%
 3.463158 10.770726 24.980196
> quantile(-s$x/s$y,c(.025,.5,.975))
      2.5%      50%      97.5%
-0.27589414 -0.11277051  0.09982482
```

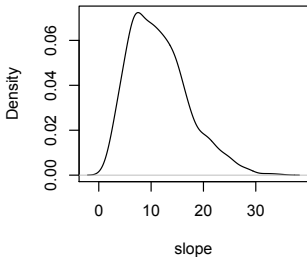




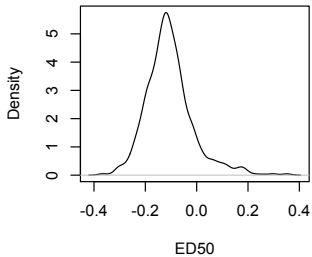
**Posterior for alpha**



**Posterior for beta**



**Posterior for  $-\alpha/\beta$**



## ... Bayesian logistic regression

		point estimate	lower 2.5% bound	upper 2.5% bound
$\alpha$	Wald	0.8466	-1.1510	2.844
	LRT		-0.8305	3.253
	Bayes		-0.5911	3.673
$\beta$	Wald	7.749	-1.8020	17.30
	LRT		1.7060	18.01
	Bayes		3.4213	25.30
ED50	Wald	-0.1092	-0.2963	0.0778
	Bayes		-0.2783	0.1067

```
> library(MCMCpack)
> posterior <- MCMClogit(y~x,data = databern)
> summary(posterior)
```

## ... Bayesian logistic regression

```
> library(MCMCpack)
> posterior <- MCMClogit(y~x,data = databern)
> summary(posterior)
```

```
Iterations = 1001:11000
Thinning interval = 1
Number of chains = 1
Sample size per chain = 10000
```

1. Empirical mean and standard deviation for each variable, plus standard error of the mean:

	Mean	SD Naive	SE	Time-series SE
(Intercept)	1.316	1.086	0.01086	0.03623
x	11.715	5.672	0.05672	0.20781

2. Quantiles for each variable:

	2.5%	25%	50%	75%	97.5%
(Intercept)	-0.63	0.5706	1.235	1.984	3.623
x	3.42	7.4827	10.731	15.003	24.931

## ... Bayesian logistic regression

		point estimate	lower 2.5% bound	upper 2.5% bound
$\alpha$	Wald	0.8466	-1.1510	2.844
	LRT		-0.8305	3.253
	Bayes		-0.5911	3.673
			-0.6300	3.623
$\beta$	Wald	7.749	-1.8020	17.30
	LRT		1.7060	18.01
	Bayes		3.4213	25.30
			3.4200	24.93
ED50	Wald	-0.1092	-0.2963	0.0778
	Bayes		-0.2783	0.1067
			-0.2749	0.1049

```
> library(MCMCpack)
> posterior <- MCMClogit(y~x,data = databern)
> summary(posterior)
```



# Posterior mode

```
s # samples from the posterior
      [,1]      [,2]
[1,]  0.078564218  4.590313
[2,] -0.130540858  6.144828
[3,]  1.113368385 14.986545
[4,] -0.567003218  6.761264
[5,]  0.551048901  5.926414
[6,]  1.563279919 14.031613
[7,]  0.294370457  3.165679
[8,]  1.869013672 14.637177
[9,]  0.247018100 11.806818
[10,] 2.018192523 16.877825
[11,] 1.751898117 11.932195
[12,] 3.013420092 20.085116
```

```
> post = density(s$y)
> which(post$y==max(post$y))
[1] 143
> post$x[143]
[1] 9.209531
> post2 = density(s$y, bw=1.5)
> which(post2$y == max(post2$y))
[1] 150
> post2$x[150]
[1] 8.867711
```

## Dependence through random effects

- ▶ Example: longitudinal data
- ▶  $Y_j = (Y_{j1}, \dots, Y_{jn_j})$  vector of observations on  $j$ th individual
- ▶ recall random effects model (normal theory):

$$Y_j = X_j\beta + Z_j b_j + \epsilon_j; \quad b_j \sim N(0, \sigma^2\Omega_b), \epsilon_j \sim N(0, \sigma^2\Omega_j)$$

- ▶ marginal distribution:

$$Y_j \sim N(X_j\beta, \sigma^2\Upsilon_j^{-1}) = N(X_j\beta, \sigma^2(\Omega_j + Z_j\Omega_b Z_j^T))$$

- ▶ sample of  $n$  i.i.d. such vectors leads to

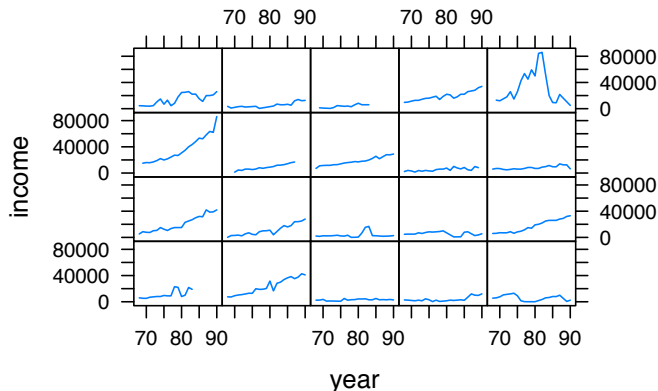
$$Y \sim N(X\beta, \sigma^2\Upsilon^{-1}),$$

- ▶  $\Omega = \text{diag}(\Omega_1, \dots, \Omega_m)$ ,  $\tilde{\Omega}_b = \text{diag}(\Omega_b, \dots, \Omega_b)$ ,

- ▶  $Z = \text{diag}(Z_1, \dots, Z_m)$ ,  $\sigma^2\Upsilon^{-1} = \Omega + Z\tilde{\Omega}_b Z^T$

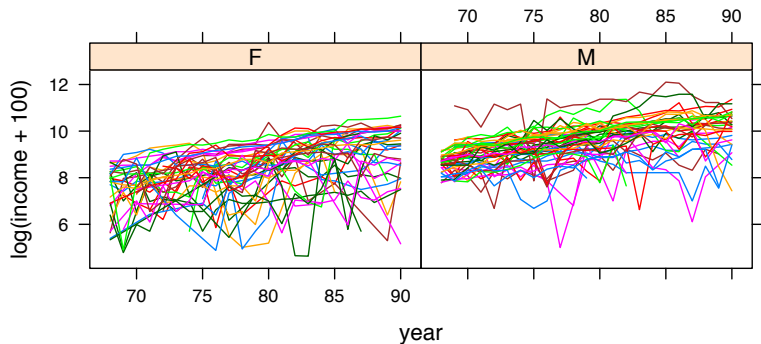
# Example: Panel Study of Income Dynamics

Faraway, §9.1



```
library(lattice)
xyplot(income ~ year | person, data = psid,
type="l", subset = (person < 21), strip = F)
```

# Example: Panel Study of Income Dynamics Faraway, §9.1



```
xyplot(log(income+100) ~ year | sex, data = psid,  
type="l", groups=person)
```

## ... PSID

```
> data(psid)
> head(psid)
  age educ sex income year person
1  31   12  M   6000   68     1
2  31   12  M   5300   69     1
3  31   12  M   5200   70     1
4  31   12  M   6900   71     1
5  31   12  M   7500   72     1
6  31   12  M   8000   73     1
> dim(psid)
[1] 1661     6
> library(lme4)
> psid$cyear = psid$year - 78
> mmod = lmer(log(income) ~ cyear*sex + age + educ +
+ (cyear | person), data=psid)
```

$$\begin{aligned}\log(\text{income})_{ij} &= \mu + \alpha \text{year}_i + \beta \text{sex}_j + (\alpha\beta) \text{year}_i \times \text{sex}_j \\ &\quad + \beta_2 \text{educ}_j + \beta_3 \text{age}_j + b_{0j} + b_{1j} \text{year}_i + \epsilon_{ij}, \\ \epsilon_{ij} &\sim N(0, \sigma^2), \quad b_j \sim N_2(0, \sigma^2 \Omega_b)\end{aligned}$$

## ... PSID

```
> mmod = lmer(log(income) ~ cyear*sex + age + educ +  
+ (cyear | person), data=psid)
```

$$\begin{aligned}\log(\text{income})_{ij} &= \mu + \alpha \text{year}_i + \beta \text{sex}_j + (\alpha\beta) \text{year}_i \times \text{sex}_j \\ &\quad + \beta_2 \text{educ}_j + \beta_3 \text{age}_j + b_{0j} + b_{1j} \text{year}_i + \epsilon_{ij}, \\ \epsilon_{ij} &\sim N(0, \sigma^2), \quad b_j \sim N_2(0, \sigma^2 \Omega_b)\end{aligned}$$

- ▶  $j$  indexes subjects,  $i$  indexes year
- ▶ variation in intercept between subjects  $b_{0j}$ ;  
in increase per year between subjects  $b_{1j}$
- ▶ year-to-year variation within subjects  $\epsilon_{ij}$

## ... PSID

$$\begin{aligned}\log(\text{income})_{ij} &= \mu + \alpha \text{year}_i + \beta \text{sex}_j + (\alpha\beta) \text{year}_i \times \text{sex}_j \\ &+ \beta_2 \text{educ}_j + \beta_3 \text{age}_j + b_{0j} + b_{1j} \text{year}_i + \epsilon_{ij}, \\ \epsilon_{ij} &\sim N(0, \sigma^2), \quad b_j \sim N_2(0, \sigma^2 \Omega_b)\end{aligned}$$

```
> summary(mmod)
Linear mixed model fit by REML ['lmerMod']
Formula: log(income) ~ cyear * sex + age + educ + (cyear | person)
Data: psid
```

```
REML criterion at convergence: 3819.776
```

```
Random effects:
```

Groups	Name	Variance	Std.Dev.	Corr
person	(Intercept)	0.2817	0.53071	
	cyear	0.0024	0.04899	0.19
Residual		0.4673	0.68357	

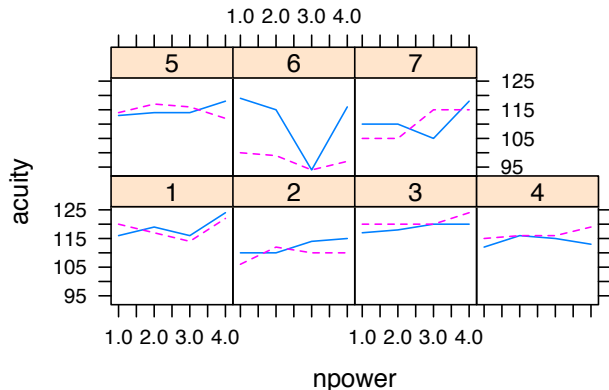
```
Number of obs: 1661, groups: person, 85
```

```
Fixed effects:
```

	Estimate	Std. Error	t value
(Intercept)	6.67420	0.54332	12.284
cyear	0.08531	0.00900	9.480
sexM	1.15031	0.12129	9.484
age	0.01093	0.01352	0.808
educ	0.10421	0.02144	4.861
cyear:sexM	-0.02631	0.01224	-2.150

# Example: Acuity of Vision

Faraway, §9.2



```
> xyplot(acuity ~ npower | subject, data=vision,  
+ type="l", groups=eye, lty=1:2, layout = c(4,2))
```



## ... vision

```
> head(vision)
  acuity power   eye subject npower
1    116   6/6  left      1        1
2    119  6/18  left      1        2
3    116  6/36  left      1        3
4    124  6/60  left      1        4
5    120   6/6  right     1        1
6    117  6/18  right     1        2
> eyemod <- lmer(acuity ~ power + (1 | subject) +
+ (1 | subject:eye), data = vision)
```

$$y_{ijk} = \mu + \rho_j + s_i + e_{ik} + \epsilon_{ijk}$$

$$s_i \sim N(0, \sigma_s^2), \quad e_{ik} \sim N(0, \sigma_e^2), \quad \epsilon_{ijk} \sim N(0, \sigma^2)$$

## ... vision

```
> summary(eyemod)
Linear mixed model fit by REML ['lmerMod']
Formula: acuity ~ power + (1 | subject) + (1 | subject:eye)
Data: vision
```

```
REML criterion at convergence: 328.7098
```

```
Random effects:
```

Groups	Name	Variance	Std.Dev.
subject:eye	(Intercept)	10.27	3.205
subject	(Intercept)	21.53	4.640
Residual		16.60	4.075

```
Number of obs: 56, groups: subject:eye, 14; subject, 7
```

```
Fixed effects:
```

	Estimate	Std. Error	t value
(Intercept)	112.6429	2.2349	50.40
power6/18	0.7857	1.5400	0.51
power6/36	-1.0000	1.5400	-0.65
power6/60	3.2857	1.5400	2.13

# Generalized linear mixed models



$$f(y_j | \theta_j, \phi) = \exp\left\{\frac{y_j \theta_j - b(\theta_j)}{\phi a_j} + c(y_j; \phi a_j)\right\}$$



$$b'(\theta_j) = \mu_j$$

- ▶ random effects

$$g(\mu_j) = \mathbf{x}_j^T \boldsymbol{\beta} + \mathbf{z}_j^T \mathbf{b}, \quad \mathbf{b} \sim N(\mathbf{0}, \Omega_b)$$

- ▶ likelihood

$$L(\boldsymbol{\beta}, \phi; \mathbf{y}) = \prod_{j=1}^n \int f(y_j | \boldsymbol{\beta}, \mathbf{b}, \phi) f(\mathbf{b}; \Omega_b) d\mathbf{b}$$

## ... generalized linear mixed models

- ▶ likelihood

$$L(\beta, \phi; \mathbf{y}) = \prod_{j=1}^n \int f(y_j | \beta, \mathbf{b}, \phi) f(\mathbf{b}; \Omega_b) d\mathbf{b}$$

- ▶ doesn't simplify unless  $f(y_j | \mathbf{b})$  is normal
- ▶ solutions proposed include
  - ▶ numerical integration, e.g. by quadrature
  - ▶ integration by MCMC
  - ▶ Laplace approximation to the integral – penalized quasi-likelihood
- ▶ reference: MASS library and book (§10.4):  
`glmmNQ`, `GLMMGibbs`, `glmmPQL`, **all in library(MASS)**  
`glmer` in `library(lme4)`

## Example: Balance experiment

Faraway, 10.1

- ▶ effects of surface and vision on balance; 2 levels of surface; 3 levels of vision
- ▶ surface: normal or foam
- ▶ vision: normal, eyes closed, domed
- ▶ 20 males and 20 females tested for balance, twice at each of 6 combinations of treatments
- ▶ auxiliary variables age, height, weight

Steele 1998, OzDASL

- ▶ linear predictor:  $\text{Sex} + \text{Age} + \text{Weight} + \text{Height} + \text{Surface} + \text{Vision} + \text{Subject} (?)$
- ▶ response measured on a 4 point scale; converted by Faraway to binary (stable/not stable)
- ▶ analysed using linear models at OzDASL

## ... balance

```
> balance <- glmer(stable ~ Sex + Age + Height + Weight + Surface + Vision +  
+ (1|Subject), family = binomial, data = ctsib)
```

```
# Subject effect is random
```

```
> summary(balance)
```

```
Generalized linear mixed model fit by maximum likelihood ['glmerMod']
```

```
...
```

```
Random effects:
```

Groups	Name	Variance	Std.Dev.
Subject	(Intercept)	8.197	2.863

Number of obs: 480, groups: Subject, 40

```
Fixed effects:
```

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	9.920750	13.358013	0.743	0.458
Sexmale	2.825305	1.762383	1.603	0.109
Age	-0.003644	0.080928	-0.045	0.964
Height	-0.151012	0.092174	-1.638	0.101
Weight	0.058927	0.061958	0.951	0.342
Surfacenorm	7.524423	0.888827	8.466	< 2e-16 ***
Visiondome	0.683931	0.530654	1.289	0.197
Visionopen	6.321098	0.839469	7.530	5.08e-14 ***

---

## ... balance

```
> library(MASS)

> balance2 <- glmmPQL(stable ~ Sex + Age + Height + Weight + Surface + Vision,
+ random = ~1 | Subject, family = binomial, data = ctsib)
> summary(balance2)
```

Random effects:

```
Formula: ~1 | Subject
      (Intercept) Residual
StdDev:    3.060712 0.5906232
```

Variance function:

```
Structure: fixed weights
Formula: ~invwt
```

Fixed effects: stable ~ Sex + Age + Height + Weight + Surface + Vision

	Value	Std.Error	DF	t-value	p-value
(Intercept)	15.571494	13.498304	437	1.153589	0.2493
Sexmale	3.355340	1.752614	35	1.914478	0.0638
Age	-0.006638	0.081959	35	-0.080992	0.9359
Height	-0.190819	0.092023	35	-2.073601	0.0455
Weight	0.069467	0.062857	35	1.105155	0.2766
Surfacenorm	7.724078	0.573578	437	13.466492	0.0000
Visiondome	0.726464	0.325933	437	2.228873	0.0263
Visionopen	6.485257	0.543980	437	11.921876	0.0000

## ... balance

```
> balance4 <- glmer(stable ~ Sex + Age + Height + Weight + Surface + Vision +
+ (1|Subject), family = binomial, data = ctsib, nAGQ = 9)
> summary(balance4)
```

Random effects:

Groups	Name	Variance	Std.Dev.
Subject	(Intercept)	7.8	2.793

Number of obs: 480, groups: Subject, 40

Fixed effects:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	13.551847	13.067369	1.037	0.2997
Sexmale	3.109307	1.724797	1.803	0.0714 .
Age	-0.001804	0.079161	-0.023	0.9818
Height	-0.175061	0.090239	-1.940	0.0524 .
Weight	0.065742	0.060606	1.085	0.2780
Surfacenorm	7.428046	0.872416	8.514	< 2e-16 ***
Visiondome	0.682509	0.527836	1.293	0.1960
Visionopen	6.210825	0.822012	7.556	4.17e-14 ***



- ▶ example: a clinical trial involves several or many centres
- ▶ an agricultural field trial repeated at a number of different farms, and over a number of different growing seasons
- ▶ a sociological study repeated in broadly similar form in a number of countries
- ▶ laboratory study uses different sets of analytical apparatus, imperfectly calibrated
- ▶ such factors are **non-specific**
- ▶ how do we account for them
  - ▶ on an appropriate scale, a parameter represents a shift in outcome
  - ▶ more complicated: the primary contrasts of concern vary across centres
  - ▶ i.e. treatment-center interaction

## ... non-specific effects

- ▶ suppose no treatment-center interaction
- ▶ example:

$$\text{logit}\{\Pr(Y_{ci} = 1)\} = \alpha_c + \mathbf{x}_{ci}^T \beta$$

- ▶ should  $\alpha_c$  be ?fixed? or ?random?
- ▶ effective use of a random-effects representation will require estimation of the variance component corresponding to the centre effects
- ▶ even under the most favourable conditions the precision achieved in that estimate will be at best that from estimating a single variance from a sample of a size equal to the number of centres
- ▶ very fragile unless there are at least, say, 10 centres and preferably considerably more

## ... non-specific effects

- ▶ if centres are chosen by an effectively random procedure from a large population of candidates, ... the random-effects representation has an attractive tangible interpretation. This would not apply, for example, to the countries of the EU in a social survey
- ▶ some general considerations in linear mixed models:
  - ▶ in balanced factorial designs, the analysis of treatment means is unchanged
  - ▶ in other cases, estimated effects will typically be 'shrunk', and precision improved
  - ▶ representation of the nonspecific effects as random effects involves independence assumptions which certainly need consideration and may need some empirical check

## ... non-specific effects

- ▶ if estimates of effect of important explanatory variables are essentially the same whether nonspecific effects are ignored, or are treated as fixed constants, then random effects model will be unlikely to give a different result
- ▶ it is important in applications to understand the circumstances under which different methods give similar or different conclusions
- ▶ in particular, if a more elaborate method gives an apparent improvement in precision, what are the assumptions on which that improvement is based, and are they reasonable?

## ... non-specific effects

- ▶ if there is an interaction between an explanatory variable [e.g. treatment] and a nonspecific variable
- ▶ i.e. the effects of the explanatory variable change with different levels of the nonspecific factor
- ▶ the first step should be to explain this interaction, for example by transforming the scale on which the response variable is measure or by introducing a new explanatory variable
  - ▶ example: two medical treatments compared at a number of centres show different treatment effects, as measured by an ratio of event rates
  - ▶ possible explanation: the difference of the event rates might be stable across centres
  - ▶ possible explanation: the ratio depends on some characteristic of the patient population, e.g. socio-economic status
- ▶ an important special application of random-effect models for interactions is in connection with overviews, that is, assembling of information from different studies of



**Happy Heart Day!**

*(Happy Valentine's Day, too!)*

**Replay**

© AGCM, Inc.