

2/60

#### Overview Fitting and graphing

### Fitting and graphing: Overview

Object-oriented approach in R:



- Fit model (obj <- glm(...))  $\rightarrow$  a model object
- $\bullet$  print (obj) and summary (obj)  $\rightarrow$  numerical results
- ${\rm \circ}$  anova (obj) and Anova (obj)  $\rightarrow$  tests for model terms
- update(obj), add1(obj), drop1(obj) for model selection

#### Plot methods:

- plot (obj) often gives diagnostic plots
- Other plot methods:
  - Mosaic plots: mosaic (obj) for "loglm" and "glm" objects
  - Effect plots: plot (Effect (obj)) for nearly all linear models
  - Influence plots (car): influencePlot (obj) for "glm" objects

## Model-based methods: Overview

#### Structure

• Explicitly assume some probability distribution for the data, e.g., binomial, Poisson, ...

Model-based methods

- Distinguish between the systematic component— explained by the model— and a random component, which is not
- Allow a compact summary of the data in terms of a (hopefully) small number of parameters

#### Advantages

- Inferences: hypothesis tests and confidence intervals
- Methods for model selection: adjust balance between goodness-of-fit and parsimony
- Predicted values give model-smoothed summaries for plotting
- $\implies$  Interpret the fitted model graphically

#### Overview Objects and methods

### Objects and methods

#### How this works:

- Model objects have a "class" attribute:
  - loglm(): "loglm"
  - glm(): c("glm", "lm") inherits also from lm()
- Class-specific methods have names like method.class, e.g., plot.glm(), mosaic.loglm()
- Generic functions (print (), summary (), plot () ...) call the appropriate method for the class

arth.mod <- glm(Better ~ Age + Sex + Treatment, data=Arthritis)
class(arth.mod)</pre>

## [1] "glm" "lm"

Overview Modeling approaches

#### Objects and methods Methods for "glm" objects

[1] add1.glm\*

[3] anova.glm\* [5] asGnm.glm\*

[7] avPlot.glm\*

## [13] confint.glm\* [15] deviance.glm\* [17] dropterm.glm\* ## [19] extractAIC.glm\* [21] formula.glm\* [23] influence.glm\*

## [27] mmp.glm\* ## [29] modFit.glm ## [31] ncvTest.glm\* [33] predict.glm

[9] bootCase.glm\* [11] coeftest.glm\*

[25] linearHypothesis.

[37] residualPlot.glm\*

[35] profile.glm\*

[39] residuals.glm ## [41] rstudent.glm\* ## [43] sigmaHat.glm\* ## [45] summary.glm

## ##

##

## ##

##

## ##

##

##

##

library(MASS); library(vcdExtra) methods (class="glm")

# Objects and methods Some available plot () methods:

#### methods("plot")

			##	[1]	plot.acf*	plot.ACF*
	addterm.glm*		##	[3]	plot.augPred*	plot.coef.mer*
	Anova.glm*		##	[5]	plot.compareFits*	plot.correspondence*
	assoc.glm		##	[7]	plot.data.frame*	plot.decomposed.ts*
	Boot.glm*		##	[9]	plot.default	plot.dendrogram*
	ceresPlot.glm*		##	[11]	plot.density*	plot.ecdf
	confidenceEllipse.glm*		##	[13]	plot.eff*	plot.efflist*
	cooks.distance.glm*		##	[15]	plot.effpoly*	plot.factor*
	drop1.glm*		##	[17]	plot.formula*	plot.function
	effects.glm*		##	[19]	plot.gam*	plot.ggplot*
	family.glm*		##	[21]	plot.gls*	plot.gnm*
	gamma.shape.glm*		##	[23]	plot.goodfit*	plot.gtable*
	leveragePlot.glm*		##	[25]	plot.hclust*	plot.histogram*
glm*	logLik.glm*		##	[27]	plot.HLtest*	plot.HoltWinters*
	model.frame.glm*		##	[29]	<pre>plot.intervals.lmList*</pre>	plot.isoreg*
	mosaic.glm		##	[31]	plot.jam*	plot.lda*
	nobs.glm*		##	[33]	plot.lm*	plot.lme*
	print.glm*		##	[35]	plot.lmList*	plot.lmList.confint*
	qqPlot.glm*		##	[37]	plot.loddsratio*	plot.loglm*
	residualPlots.glm*		##	[39]	plot.mca*	plot.medpolish*
	rstandard.glm*		##	[41]	plot.merMod*	plot.mlm*
	sieve.glm		##	[43]	plot.mlm.efflist*	plot.nffGroupedData*
	summarise.glm*		# #	[45]	plot.nfnGroupedData*	plot.nls*
	vcov.glm*	5/60	##	[47]	plot.nmGroupedData*	plot.oddsratio*

Modeling approaches Overview

### Modeling approaches: Overview

### Association models Loglinear models (Contingency table form) [Admit][GenderDept] [AdmitDept][GenderDept] [AdmitDept][AdmitGender][GenderDept] Poisson GLMs (Frequency data frame) Freq ~ Admit + Gender\*Dept Freq ~ Admit\*Dept + Gender\*Dept Freq ~ Admit\*Dept + Admit\*Gender + Gender\*Dept Ordered variables Freq ~ right+left+Diag(right.left) Freq ~ right+left+Symm(right.left)

### Response models

- Binary response
- Categorical predictors: Logit models logit(Admit) ~ 1 logit(Admit) ~ Dept
- logit(Admit) ~ Dept + Gender
- Continuous/mixed predictors: Logistic regression models Pr(Admit) ~ Dept + Age + GRE
- Polytomous response
- Ordinal: proportional odds model Improve ~ Age + Sex + Treatment
- General: multinomial model
- WomenWork ~ Kids + HusbandInc

7/60

### Logistic regression models

#### **Response variable**

- Binary response: success/failure, vote: yes/no
- Binomial data: x successes in n trials (grouped data)
- Ordinal response: none < some < severe depression
- Polytomous response: vote Liberal, Tory, NDP, Green

#### **Explanatory variables**

- Quantitative regressors: age, dose
- Transformed regressors:  $\sqrt{age}$ , log(dose)
- Polynomial regressors:  $age^2$ ,  $age^3$ , ... (or better: splines)
- Categorical predictors: treatment, sex (dummy variables, contrasts)
- Interaction regessors: treatment  $\times$  age, sex  $\times$  age

This is exactly the same as in classical ANOVA, regression models

#### Examples

# Arthritis treatment data

Examples

Examples



- The response variable, Improved is ordinal: "None" < "Some" < "Marked"
- A binary logistic model can consider just Better = (Improved>"None")
- Other important predictors: Sex, Treatment
- Main Q: how does treatment affect outcome?
- How does this vary with Age and Sex?
- This plot shows the binary observations, with several model-based smoothings





- Admit/Reject can be considered a binomial response for each Dept and Gender
- Logistic regression here is analogous to an ANOVA model, but for log odds(Admit)
- (With categorical predictors, these are often called logit models)
- Every such model has an equivalent loglinear model form.
- This plot shows fitted logits for the main effects model, Dept + Gender

### Survival in the Donner Party

- Binary response: survived
- Categorical predictors: sex, family
- Quantitative predictor: age
- Q: Is the effect of age linear?
- Q: Are there interactions among predictors?
- This is a generalized pairs plot, with different plots for each pair



## Binary response: What's wrong with OLS?

Binary response

- For a binary response, Y ∈ (0, 1), want to predict π = Pr(Y = 1 | x)
- A linear probability model uses classical linear regression (OLS)
- Problems:

0/60

11/60

- Gives predicted values and CIs outside  $0 \le \pi \le 1$
- Homogeneity of variance is violated: V(π̂) = π̂(1 − π̂) ≠ constant
- Inferences, hypothesis tests are wrong!



#### Binary response



### Logistic regression: One predictor

For a single quantitative predictor, x, the simple linear logistic regression model posits a linear relation between the *log odds* (or *logit*) of Pr(Y = 1) and x,

$$\log[\pi(x)] \equiv \log\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x$$

- When β > 0, π(x) and the log odds increase as x increases; when β < 0 they decrease with x.</li>
- This model can also be expressed as a model for the probabilities  $\pi(x)$

$$\pi(x) = \text{logit}^{-1}[\pi(x)] = \frac{1}{1 + \exp[-(\alpha + \beta x)]}$$

### Logistic regression: One predictor

The coefficients of this model have simple interpretations in terms of odds and log odds:

• The odds can be expressed as a multiplicative model

$$\operatorname{podds}(Y=1) \equiv \frac{\pi(x)}{1-\pi(x)} = \exp(\alpha + \beta x) = e^{\alpha} (e^{\beta})^x$$
 (1)

Thus:

- $\beta$  is the change in the log odds associated with a unit increase in *x*.
- The odds are multiplied by  $e^{\beta}$  for each unit increase in *x*.
- $\alpha$  is log odds at x = 0;  $e^{\alpha}$  is the odds of a favorable response at this *x*-value.
- In R, use exp(coef(obj)) to get these values.
- Another interpretation: In terms of probability, the slope of the logistic regression curve is  $\beta \pi (1 \pi)$
- This has the maximum value  $\beta/4$  at  $\pi = \frac{1}{2}$

#### Binary response Logistic regression model

### Logistic regression models: Multiple predictors

- For a binary response,  $Y \in (0, 1)$ , let **x** be a vector of *p* regressors, and  $\pi_i$  be the probability,  $Pr(Y = 1 | \mathbf{x})$ .
- The logistic regression model is a linear model for the *log odds*, or *logit* that Y = 1, given the values in x,

$$\mathsf{logit}(\pi_i) \equiv \mathsf{log}\left(\frac{\pi_i}{1-\pi_i}\right) = \alpha + \mathbf{x}_i^\mathsf{T} \boldsymbol{\beta}$$
$$= \alpha + \beta_1 \mathbf{x}_{i1} + \beta_2 \mathbf{x}_{i2} + \dots + \beta_p \mathbf{x}_{ip}$$

 An equivalent (non-linear) form of the model may be specified for the probability, π<sub>i</sub>, itself,

$$\pi_i = \{1 + \exp(-[\alpha + \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta}])\}^{-1}$$

The logistic model is also a *multiplicative* model for the odds of "success,"

$$\frac{\pi_i}{1-\pi_i} = \exp(\alpha + \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta}) = \exp(\alpha) \exp(\alpha) \exp(\boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta})$$

Increasing  $x_{ij}$  by 1 increases logit( $\pi_i$ ) by  $\beta_j$ , and multiplies the odds by  $e^{\beta_j}$ .

```
17/60
```

### The summary () method gives details:

summary(arth.logistic)
##
## Call:

```
## glm(formula = Better ~ Age, family = binomial, data = Arthritis)
##
## Deviance Residuals:
     Min 1Q Median
##
                            30
                                       Max
## -1.5106 -1.1277 0.0794 1.0677 1.7611
##
## Coefficients:
##
           Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.6421 1.0732 -2.46 0.014 *
## Age
              0.0492
                       0.0194 2.54
                                         0.011 *
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
  (Dispersion parameter for binomial family taken to be 1)
##
##
      Null deviance: 116.45 on 83 degrees of freedom
## Residual deviance: 109.16 on 82 degrees of freedom
## AIC: 113.2
##
## Number of Fisher Scoring iterations: 4
```

#### Binary response Fittin

### Fitting the logistic regression model

Logistic regression models are the special case of generalized linear models, fit in R using glm(..., family=binomial) For this example, we define Better as any improvement at all:

```
data("Arthritis", package="vcd")
Arthritis$Better <- as.numeric(Arthritis$Improved > "None")
```

#### Fit and print:

arth.logistic <- glm(Better ~ Age, data=Arthritis, family=binomial)
arth.logistic</pre>

```
##
## Call: glm(formula = Better ~ Age, family = binomial, data = Arthritis)
##
Coefficients:
## (Intercept) Age
##
-2.6421 0.0492
##
## Degrees of Freedom: 83 Total (i.e. Null); 82 Residual
## Null Deviance: 116
## Residual Deviance: 109 AIC: 113
```

18/60

#### Binary response Fitting

### Interpreting coefficients

coef(arth.logistic)	<pre>exp(coef(arth.logistic))</pre>						
## (Intercept) Age ## -2.642071 0.049249	## (Intercept) Age ## 0.071214 1.050482						
	<pre>exp(10*coef(arth.logistic)[2]) ## Age ## 1.6364</pre>						
Interpretations:							
• log odds(Better) increase by $\beta = 0.0492$ for each year of age • odds(Better) multiplied by $e^{\beta} = 1.05$ for each year of age— a 5%							

over 10 years, odds(Better) are multiplied by exp(10 × 0.0492) = 1.64, a 64% increase.

• Pr(Better) increases by  $\beta/4 = 0.0123$  for each year (near  $\pi = \frac{1}{2}$ )

#### Binary response Multiple predictors

### Multiple predictors

Interpreting coefficients

cbind(coef=coef(arth.logistic2),

## TreatmentTreated 1.7598

-0.5781

0.0487

-1.4878

##

## (Intercept)

## I(Age - 50)

## SexMale

The main interest here is the effect of Treatment. Sex and Age are control variables. Fit the main effects model (no interactions):

 $\operatorname{logit}(\pi_i) = \alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_2 x_{i2}$ 

where  $x_1$  is Age and  $x_2$  and  $x_3$  are the factors representing Sex and Treatment, respectively. R uses dummy (0/1) variables for factors.

, J	0 if Female		v f	0	if Placebo		
$x_2 = \{$	1	if Male	$x_3 = \begin{cases} \\ \\ \\ \end{cases}$	1	if Treatment		

Multiple predictors

**OddsRatio**=exp(coef(arth.logistic2)), exp(confint(arth.logistic2)))

0.561 0.2647 1.132

1.050 1.0100 1.096

0.226 0.0652 0.689

5.811 2.1187 17.727

coef OddsRatio 2.5 % 97.5 %

- $\alpha$  doesn't have a sensible interpretation here. Why?
- $\beta_1$ : increment in log odds(Better) for each year of age.
- $\beta_2$ : difference in log odds for male as compared to female.
- $\beta_3$ : difference in log odds for treated vs. the placebo group

Binary response

### Multiple predictors: Fitting

Fit the main effects model. Use I (Age-50) to center Age, making  $\alpha$  interpretable.

coeftest() in Imtest gives just the tests of coefficients provided by
summary():

```
library(lmtest)
coeftest(arth.logistic2)
##
## z test of coefficients:
##
##
                   Estimate Std. Error z value Pr(>|z|)
## (Intercept)
                    -0.5781
                                0.3674 -1.57
                                                  0.116
## I(Age - 50)
                    0.0487
                                0.0207
                                          2.36
                                                  0.018 *
                    -1.4878
                                0.5948
                                         -2.50
                                                  0.012 *
## SexMale
## TreatmentTreated 1.7598
                                0.5365
                                                  0.001 **
                                         3.28
## ----
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

22/60

### Hypothesis testing: Questions

Overall test: How does my model, logit(π) = α + x<sup>T</sup>β compare with the null model, logit(π) = α?

Hypothesis tests

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_p = 0$$

• **One predictor**: Does *x<sub>k</sub>* significantly improve my model? Can it be dropped?

 $H_0: \beta_k = 0$  given other predictors retained

• Lack of fit: How does my model compare with a perfect model (saturated model)?

For ANOVA, regression, these tests are carried out using *F*-tests and *t*-tests. In logistic regression (fit by maximum likelihood) we use

- *F*-tests  $\rightarrow$  likelihood ratio  $G^2$  tests
- *t*-tests  $\rightarrow$  Wald *z* or  $\chi^2$  tests

- $\alpha = -0.578$ : At age 50, females given placebo have odds(Better) of  $e^{-0.578} = 0.56$ .
- $\beta_1 = 0.0487$ : Each year of age multiplies odds(Better) by  $e^{0.0487} = 1.05$ , a 5% increase.
- β<sub>2</sub> = -1.49: Males e<sup>-1.49</sup> = 0.26 × less likely to show improvement as females. (Or, females e<sup>1.49</sup> = 4.437 × more likely than males.)
- $\beta_3 = 1.76$ : Treated  $e^{1.76}$ =5.81  $\times$  more likely Better than Placebo

23/60

### Maximum likelihood estimation

Likelihood, L = Pr(data | model), as function of model parameters
For case *i*,

Hypothesis tests

$$\mathcal{L}_{i} = \begin{cases} p_{i} & \text{if } Y = 1 \\ 1 - p_{i} & \text{if } Y = 0 \end{cases} = p_{i}^{Y_{i}} (1 - p_{i}^{Y_{i}}) \quad \text{where} \quad p_{i} = 1/(1 + \exp(\mathbf{x}_{i}\boldsymbol{\beta}))$$

• Under independence, joint likelihood is the product over all cases

$$\mathcal{L} = \prod_{i}^{n} p_{i}^{Y_{i}} (1 - p_{i}^{Y_{i}})$$

•  $\implies$  Find estimates  $\hat{\beta}$  that maximize log  $\mathcal{L}$ . Iterative, but this solves the "estimating equations"

Hypothesis tests

 $\boldsymbol{X}^{\mathsf{T}}\boldsymbol{y} = \boldsymbol{X}^{\mathsf{T}}\widehat{\boldsymbol{p}}$ 

## **Overall test**

- Likelihood ratio test (G<sup>2</sup>)
  - Compare *nested* models, similar to incremental F tests in OLS
  - Let L<sub>1</sub> = maximized likelihood for our model logit(π<sub>i</sub>) = β<sub>0</sub> + x<sub>i</sub><sup>T</sup>β w/ k predictors

Hypothesis tests

- Let  $\mathcal{L}_0$  = maximized likelihood for **null** model logit( $\pi_i$ ) =  $\beta_0$  under  $H_0$ :  $\beta_1 = \beta_2 = \cdots = \beta_k = 0$
- Likelihood-ratio test statistic:

$$G^2 = -2\log\left(\frac{L_0}{L_1}\right) = 2(\log L_1 - \log L_0) \sim \chi_k^2$$

Hypothesis tests

26/60

### Wald tests and confidence intervals

- Analogous to *t*-tests in OLS
- $H_0: \beta_i = 0$

$$z = \frac{b_k}{s(b_k)} \sim \mathcal{N}(0,1) \quad \text{or} \quad z^2 \sim \chi_1^2$$
(Wald chi-square)

- Confidence interval:
  - $b_k \pm z_{1-\alpha/2} s(b_k)$

		Analys	is of	Maximum	Likelihood	Estimates		
e.g.,	Parameter		DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq	
	Intercept sex treat age	Female Treated	1 1 1	-4.5033 1.4878 1.7598 0.0487	1.3074 0.5948 0.5365 0.0207	11.8649 6.2576 10.7596 5.5655	0.0006 0.0124 0.0010 0.0183	

## LR, Wald and score tests

25/60

27/60

Testi	ng Global Null	Hypothesis:	BETA=0
Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	24.3859	3	<.0001
Score	22.0051	3	<.0001
Wald	17.5147	3	0.0006



# Different ways to measure departure from $H_0$ : $\beta = 0$

- $\bullet$  LR test: diff in log L
- Wald test:  $(\hat{\boldsymbol{\beta}} \boldsymbol{\beta}_0)^2$
- Score test: slope at  $\beta = 0$

#### Visualizing

#### Visualizing

### Plotting logistic regression data

Plotting a binary response together with a fitted logistic model can be difficult because the 0/1 response leads to much overplottting.

- Need to jitter the points
- Useful to show the fitted logistic curve
- Confidence band gives a sense of uncertainty
- Adding a non-parametric (loess) smooth shows possible nonlinearity
- NB: Can plot either on the response scale (probability) or the link scale (logit) where effects are linear



### Types of plots

- Conditional plots: Stratified plot of Y or logit(Y) vs. one X, conditioned by other predictors— only that subset is plotted for each
- Full-model plots: plots of fitted response surface, showing all effects; usually shown in separate panels
- Effect plots: plots of predicted effects for terms in the model, averaged over predictors not involved in a given term.



Visualizing Conditional plots

## Conditional plots with ggplot2



Conditional plots

Visualizing

Conditional plots with ggplot2 Plot of Arthritis treatment data, by Treatment (ignoring Sex)



## Conditional plot, faceted by Sex

gg + facet\_wrap(~ Sex)



The data is too thin for males to estimate each regression separately

Visualizing Full-model plots	Visualizing Full-model plots
Full-model plots	Plotting with ggplot2 package
Full-model plots show the fitted values on the logit scale or on the response scale (probability), usually with confidence bands. This often requires a bit of custom programming. Steps:	
<ul> <li>Obtain fitted values with predict (model, se.fit=TRUE) — type="link" (logit) is the default</li> <li>Can use type="response" for probability scale</li> <li>Join this to your data (cbind())</li> <li>Plot as you like: plot(), ggplot(),</li> </ul>	<pre>arth.fit2\$obs &lt;- c(-4, 4)[1+arth.fit2\$Better]  gg2 &lt;- ggplot(arth.fit2, aes(x=Age, y=fit, color=Treatment)) +   geom_line(size = 2) +   geom_ribbon(aes(ymin = fit - 1.96 * se.fit,</pre>
<pre>arth.fit2 &lt;- cbind(Arthritis,</pre>	<pre>geom_point(aes(y=obs), position=position_jitter(height=0.25, width=0)) gg2 + facet_wrap(~ Sex)</pre>
<pre>## ID Treatment Sex Age Improved Better fit se.fit ## 1 57 Treated Male 27 Some 1 -1.43 0.758 ## 2 46 Treated Male 29 None 0 -1.33 0.728 ## 3 77 Treated Male 30 None 0 -1.28 0.713 ## 4 17 Treated Male 32 Marked 1 -1.18 0.684</pre>	
33/60	34/60
Visualizing Full-model plots	Visualizing Full-model plots

# Full-model plots

Ploting on the logit scale shows the additive effects of age, treatment and sex



These plots show the data (jittered) as well as model uncertainty (confidence bands)

## Full-model plots

Ploting on the probability scale may be simpler to interpret



These plots show the data (jittered) as well as model uncertainty (confidence bands)



- One plot for each variable in the model
- Other variables: continuous— held fixed at median; factors— held fixed at most frequent value
- Partial residuals (r<sub>j</sub>): the coefficient β<sub>j</sub> in the full model is the slope of the simple fit of r<sub>j</sub> on x<sub>j</sub>.

#### General ideas Effect plots

#### Effect plots General ideas

### Effect plots: basic ideas

Show a given effect (and low-order relatives) controlling for other model effects.





Effect plots Examples

Sex

### Effect plots for generalized linear models: Details

- For simple models, full model plots show the complete relation between response and all predictors.
- Fox(1987)— For complex models, often wish to plot a specific main effect or interaction (including lower-order relatives) --- controlling for other effects
  - Fit full model to data with linear predictor (e.g., logit)  $\eta = X\beta$  and link
    - function  $g(\mu) = \eta \rightarrow \text{estimate } \boldsymbol{b}$  of  $\beta$  and covariance matrix  $\widehat{V}(\boldsymbol{b})$  of  $\boldsymbol{b}$ .
  - Construct "score data"
    - Vary each predictor in the term over its' range
    - Fix other predictors at "typical" values (mean, median, proportion in the data) → "effect model matrix," X\*
  - Use predict () on X\*
    - Calculate fitted effect values,  $\hat{\eta}^* = X^* b$ .
    - Standard errors are square roots of diag  $X^* \widehat{V(b)} X^{*T}$
  - Plot  $\hat{\eta}^*$ , or values transformed back to scale of response,  $g^{-1}(\hat{\eta}^*)$ .
- Note: This provides a general means to visualize interactions in all linear and generalized linear models.

Effect plots Examples

#### Full model plots:

arth.full <- Effect(c("Age", "Treatment", "Sex"), arth.logistic2)</pre> plot(arth.full, multiline=TRUE, ci.style="bands", colors = c("red "blue"), **lwd**=3, ...)



#### Age\*Treatment\*Sex effect plot

Plotting main effects:

0.2





Age





0.5

0.4

0.3

0.2

Placebo



Treatment

Treater

41/60

#### Effect plots Examples

#### Case studies Arrest

#### Model with interaction of Age x Sex

#### plot(allEffects(arth.logistic4), rows=1, cols=3)



• Only the high-order terms for Treatment and Sex\*Age need to be interpreted

B SECTION > TORONTO STAR < WEDNESDAY, DECEMBER 11, 2002 ★ thestar.com

- (How would you describe this?)
- The main effect of Age looks very different, averaged over Treatment and Sex

Case study: Arrests for Marijuana Possession Context & background

- In Dec. 2002, the *Toronto Star* examined the issue of racial profiling, by analyzing a data base of 600,000+ arrest records from 1996-2002.
- They focused on a subset of arrests for which police action was discretionary, e.g., simple possession of small quantities of marijuana, where the police could:
  - Release the arrestee with a summons—like a parking ticket
  - Bring to police station, hold for bail, etc.— harsher treatment
- Response variable: released Yes, No

45/60

• Main predictor of interest: skin-colour of arrestee (black, white)

The Toronto Star meets mosaic displays...

#### Race and Crime



# Man behind the numbers

Case studies Arrests ... Which got turned into this infographic: Same charge, different treatment Statistical analysis of single drug possession charges shows Degree of likelihood that blacks are much less likely to be released at the scene Much less likely to occur and much more likely to be held in custody for a bail hearing. Much more likely to occur Darker colours represent a stronger statistical link between skin colour and police treatment. More likely to occur Whites are more likely to be released at the scene 6,662 78% 14.5% 7.5% held charges laid released at the scene released at station for bail Blacks are much more likely to be held for bail hearings 2.446 64% 20% 16% held charges released at the scene released at station for bail

40

50

60

70

80

... Hey, they even spelled likelihood correctly!

20

30

10

SOURCE: Toronto police arrest records 1996-200

0%

100

90

### Arrests for Marijuana Possession: Data

#### Data Control variables:

- year, age, sex
- employed, citizen Yes, No
- checks Number of police data bases (previous arrests, previous convictions, parole status, etc.) in which the arrestee's name was found.

lik lik dat som	orary orary ta(Ari ne(Ari	(effects) (car) rests) rests)	# fo # fo	or Ar: or And	rests ova()	s data			
##		released	colour	year	age	sex	employed	citizen	checks
##	299	Yes	Black	2001	24	Male	Yes	Yes	3
##	766	Yes	White	2000	18	Male	Yes	Yes	2
##	1530	Yes	White	2000	14	Male	Yes	Yes	0
##	2367	Yes	White	1999	23	Male	Yes	Yes	0
##	2619	No	White	2001	22	Male	Yes	Yes	2
##	2664	No	Black	1998	38	Male	Yes	No	3
##	4202	No	White	2001	47	Female	No	Yes	1
##	4206	Yes	Black	1999	26	Male	No	Yes	0
##	4323	Yes	Black	1999	22	Male	No	Yes	6
# #	5102	Yes	White	2000	19	Male	Yes	Yes	0

### Arrests for Marijuana Possession: Model

To allow possibly non-linear effects of year, we treat it as a factor:

> Arrests\$year <- as.factor(Arrests\$year)</pre>

Logistic regression model with all main effects, plus interactions of colour:year and colour:age

```
> arrests.mod <- glm(released ~ employed + citizen + checks + colour *
+ year + colour * age, family = binomial, data = Arrests)</pre>
```

> Anova(arrests.mod)

Analysis of Deviance Table (Type II tests)

#### Response: released

49/60

51/60

	LR Chisq	DÍ	Pr(>Chisq)							
employed	72.673	1	< 2.2e-16	* * *						
citizen	25.783	1	3.820e-07	* * *						
checks	205.211	1	< 2.2e-16	* * *						
colour	19.572	1	9.687e-06	* * *						
year	6.087	5	0.2978477							
age	0.459	1	0.4982736							
colour:year	21.720	5	0.0005917	* * *						
colour:age	13.886	1	0.0001942	* * *						
Signif. code	s: 0 '**	**'	0.001 '**'	0.01	'*'	0.05	1.1	0.1	1	

50/60

1

Case studies Arrests

### Effect plots: colour

plot(Effect("colour", arrests.mod), ci.style="bands", ...)



- Effect plot for colour shows average effect controlling (adjusting) for *all* other factors simultaneously
- (The Star analysis, controlled for these one at a time.)
- ⇒ Evidence for different treatment of blacks and whites ("racial profiling")
- (Even Frances Nunziata could understand this.)
- NB: Effects smaller than claimed by *the Star*

Case studies Arrests

### Effect plots: Interactions

The story turned out to be more nuanced than reported by the *Toronto Star*, as shown in effect plots for interactions with colour.

> plot(effect("colour:year", arrests.mod), multiline = TRUE, ...)



- Up to 2000, strong evidence for differential treatment of blacks and whites
- Also evidence to support Police claim of effect of training to reduce racial effects in treatment

#### Case studies

### Effect plots: Interactions

The story turned out to be more nuanced than reported by the Toronto Star, as shown in effect plots for interactions with colour.



- Opposite age effects for blacks and whites-
- Young blacks treated more harshly than young whites
- Older blacks treated less harshly than older whites

### Effect plots: allEffects

All model effects can be viewed together using plot (allEffects (mod))

> arrests.effects <- allEffects(arrests.mod, xlevels = list(age = seg(15,</pre> 45, 5)))

> plot(arrests.effects, ylab = "Probability(released)")





#### Model diagnostics

### Model diagnostics

Probability(released)

0.90

0.85

As in regression and ANOVA, the validity of a logistic regression model is threatened when:

- Important predictors have been omitted from the model
- Predictors assumed to be linear have non-linear effects on Pr(Y = 1)
- Important interactions have been omitted

• A few "wild" observations have a large impact on the fitted model or coefficients

### Model specification: Tools and techniques

- Use non-parametric smoothed curves to detect non-linearity
- Consider using polynomial terms  $(X^2, X^3, ...)$  or regression splines (e.g., ns(X, 3)
- Use update (model, ...) to test for interactions—formula:  $. \sim .^2$

### Diagnostic plots in R

In R, plotting a glm object gives the "regression guartet" — basic diagnostic plots

Model diagnostics

arth.mod1 <- glm(Better ~ Age + Sex + Treatment, data=Arthritis, familv='binomial') plot(arth.mod1)



Better versions of these plots are available in the car package

55/60

#### Model diagnostics Leverage and influence

#### Model diagnostics Leverage and influence

### Unusual data: Leverage and Influence

- "Unusual" observations can have dramatic effects on estimates in linear models
  - Can change the coefficients for the predictors
  - Can change the predicted values for all observations

#### • Three archetypal cases:

- Typical X (low leverage), bad fit Not much harm
- Unusual X (high leverage), good fit Not much harm
- Unusual X (high leverage), bad fit BAD, BAD, BAD
- Influential observations: unusual in both X and Y
- Heuristic formula:

Influence =  $Leverage_X \times Residual_Y$ 

Model diagnostics Leverage and influence

![](_page_14_Figure_13.jpeg)

![](_page_14_Figure_15.jpeg)

#### Model diagnostics Leverage and influence

40

30

10 20 30 40

### Which cases are influential?

50

60

40

30

20

57/60

10 20 30 40

	ID	Treatment	Sex	Age	Better	StudRes	Hat	CookD
1	57	Treated	Male	27	1	1.922	0.08968	0.3358
15	66	Treated	Female	23	0	-1.183	0.14158	0.2049
39	11	Treated	Female	69	0	-2.171	0.03144	0.2626

70 80

![](_page_14_Figure_19.jpeg)

OL

58/60

50

60 70 80