

Generalized Linear Models

Want a model for $(Y | X_1, \dots, X_p)$.

1. Specify a **family of distributions**.

2. Specify a **link function $g()$** such that:

$$g\{E(Y|X_1, \dots, X_p)\} = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p.$$

Note: $\eta = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$ referred to as **linear predictor**.

> `glm(y~x, family=..., link=...)`

1

Three primary examples:

$$Y \sim N(\mu, \sigma^2), \mu = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p.$$

$$Y \sim \text{Bernoulli}(p), \text{logit}(p) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p.$$

Note: Binomial data can be treated as Bernoulli data. (Some issues lurking here - what is sample size? also deviance.)

$$Y \sim \text{Poisson}(\lambda), \log \lambda = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p.$$

Note: Have given the ‘default’ link function in each case. This is more than convention, tied up with “exponential family” theory - the **canonical link function** is that which equates the **natural parameter of the distribution** with the **linear predictor**.

Other links sometimes used in Bernoulli case: probit $\Phi^{-1}(p)$, complementary log-log: $\log\{-\log(1-p)\}$.

2

Computation (say ϕ known):

Distribution + link = fully-specified prob. model. Gives log-likelihood function $l(\beta)$.

NR algorithm to maximize $l()$, requires evaluation of $l'()$ and $l''()$.

Fisher-scoring modification: Replace $l''()$ with $E_{Y|X}\{l''()\}$.

Turns out we have an **iteratively reweighted least squares** (IRLS) algorithm.

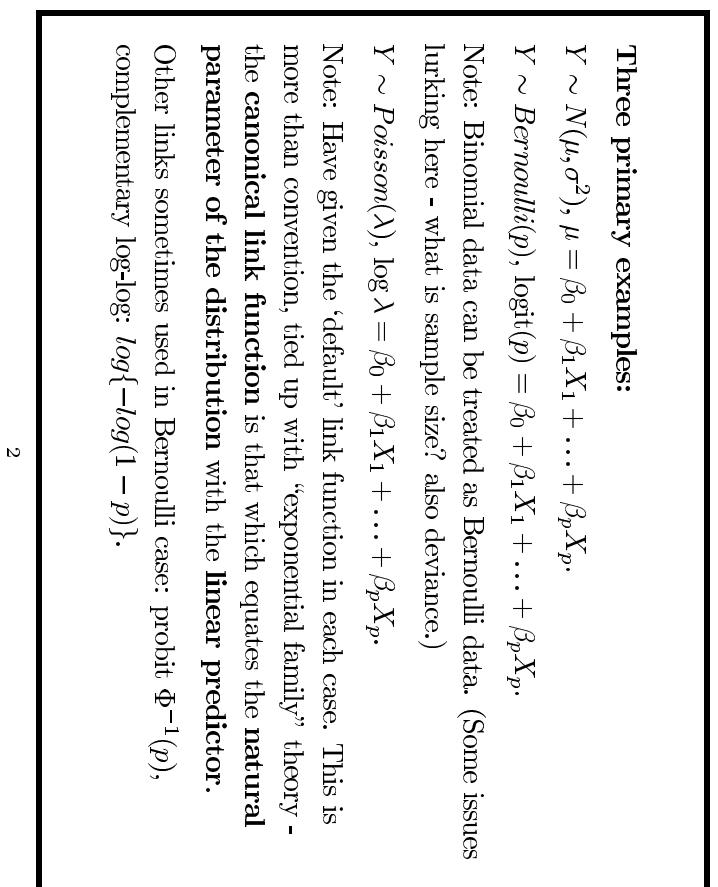
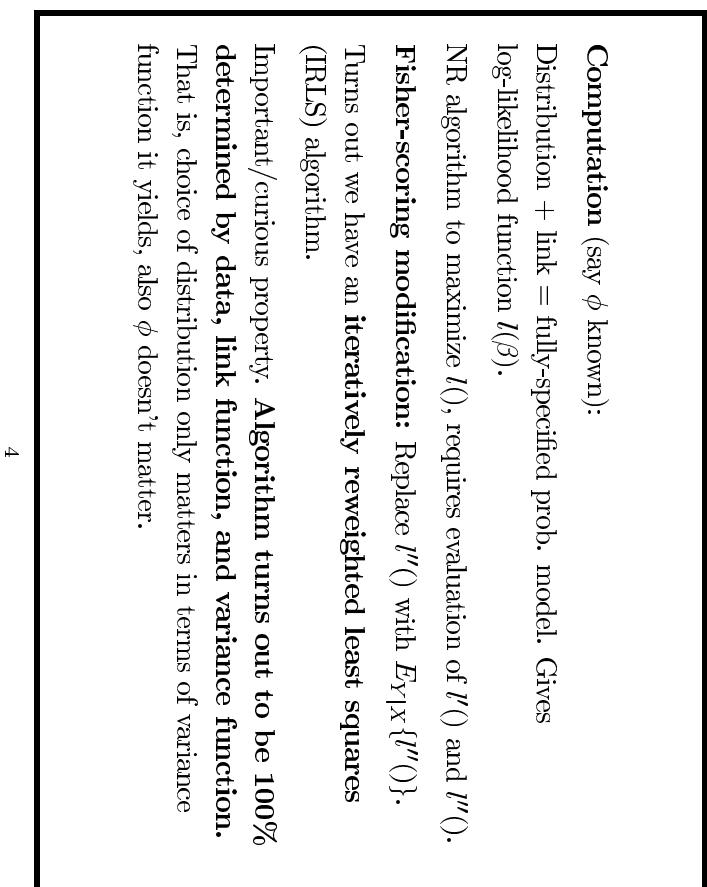
Important/curious property. **Algorithm turns out to be 100% determined by data, link function, and variance function.** That is, choice of distribution only matters in terms of variance function it yields, also ϕ doesn’t matter.

Normal: $v(\mu) = 1, \phi$ unknown

Bernoulli: $v(\mu) = \mu(1 - \mu), \phi = 1$.

Poisson: $v(\mu) = \mu, \phi = 1$.

3



4

So, have a general regression strategy for multiple types of Y variable.

Point estimation.

Gives strategy for add-on estimation of ϕ when needed. Pretend $\phi = 1$ and get $\hat{\beta}$, then **deviance** (recall defn, D_M in text). Then estimate $\hat{\phi} = (n - p)^{-1}D_M$.

Rationale: real/scaled/residual deviance is $\phi^{-1}D_M$.

Also, the role of the variance function leads to **quasi-likelihood** ideas: choose a variance function as one sees fit, don't worry if it corresponds to a real probability model or not.

5

SE / Interval estimation.

Hypothesis testing: nested, q predictors versus p predictors.

Poisson regression example

Deaths from coronary heart disease in a (famous) cohort study...

	Smokers		Non-smokers		
	age	deaths	person-years	deaths	person-years
35-44	32	52407	2	18790	
45-54	104	43248	12	10673	
55-64	206	28612	28	5710	
65-74	186	12663	28	2585	
75-84	102	5317	31	1462	
Residual-based diagnostics					

6

7

```
glm(formula = deaths ~ I(log(pyyears)) + age.grp + smoke,
family = poisson, data = dat)
```

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-21.4739	2.2556	-9.520	< 2e-16 ***
I(log(pyyears))	2.4352	0.2270	10.727	< 2e-16 ***
age.grp	1.7702	0.1542	11.478	< 2e-16 ***
smoke	-1.6991	0.3548	-4.789	1.68e-06 ***

Null deviance: 644.269 on 9 degrees of freedom
Residual deviance: 25.576 on 6 degrees of freedom

9

```
glm(formula = deaths ~ offset(log(pyyears)) + age.grp +
smoke,
family = poisson, data = dat)
```

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-7.28250	0.12076	-60.304	< 2e-16 ***
age.grp	0.83583	0.02904	28.777	< 2e-16 ***
smoke	0.40637	0.10720	3.791	0.00015 ***

Null deviance: 935.067 on 9 degrees of freedom
Residual deviance: 69.182 on 7 degrees of freedom

10

```
glm(formula = deaths ~ offset(log(pyyears)) + age.grp +
smoke,
family = poisson, data = dat)
```

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-8.61296	0.29172	-29.524	< 2e-16 ***
age.grp	1.98113	0.16025	12.363	< 2e-16 ***
I(age.grp^2)	-0.19768	0.02737	-7.223	5.08e-13 ***
smoke	1.13342	0.28077	4.037	5.42e-05 ***
age.grp:smoke	-0.30755	0.09704	-3.169	0.00153 **

Null deviance: 935.067 on 9 degrees of freedom
Residual deviance: 12.176 on 6 degrees of freedom

11

```
glm(formula = deaths ~ offset(log(pyyears)) + age.grp +
smoke,
family = poisson, data = dat)
```

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-8.61296	0.29172	-29.524	< 2e-16 ***
age.grp	1.98113	0.16025	12.363	< 2e-16 ***
I(age.grp^2)	-0.19768	0.02737	-7.223	5.08e-13 ***
smoke	1.13342	0.28077	4.037	5.42e-05 ***
age.grp:smoke	-0.30755	0.09704	-3.169	0.00153 **

Null deviance: 935.067 on 9 degrees of freedom
Residual deviance: 1.6354 on 5 degrees of freedom

12

