

## Logistic Regression

$$\text{Odds} = \frac{\pi}{1-\pi}$$

$$\text{logit}(\pi_i) = \log\left(\frac{\pi_i}{1-\pi_i}\right) = \alpha + \beta_i x_{ji}$$

Estimate  $\alpha$  and  $\beta$  maximum likelihood

$$\eta_i = \alpha + \beta x_i$$

$$\frac{\hat{\pi}_i}{1-\hat{\pi}_i} = \exp(\eta_i)$$

$$\hat{\pi}_i = \frac{\exp(\eta_i)}{1 + \exp(\eta_i)} = (1 + \exp(-\eta_i))^{-1}$$

## Generalized Linear Models

$$g(\mu_i) = \eta_i = \beta^T x_i$$

$$\text{- Linear component: } \eta_i = \beta^T x_i$$

- **Link function:** let univariate function  $g$  be monotonic (strictly decreases or increases) and differentiable

- **Random component:**  $y_i$  are independent and from an exponential family, which impose the variance of  $y_i$  depends on  $\mu_i$  through a variance function  $\text{var}(y_i) = \phi V(\mu_i)$  where  $\phi$  is called the *dispersion parameter*.

Classic linear regression assumes normal dist.

Logistic regression assumes  $g(\mu) = \log\left[\frac{\mu}{1-\mu}\right]$  and  $y_i$  is a Bernoulli trial ( $\mathbb{P}(y) = \mu(1-\mu)$ ). Other possible links for binary responses include

- Probit:

$$g(\mu) = \Phi^{-1}(\mu), \text{ where } \Phi \text{ is the cumulative standard normal distribution}$$

- Complementary log-log:

$$g(\mu) = \log(-\log(1-\mu))$$

## Maximum Likelihood Estimation of Proportions

Suppose we draw a random sample of size  $n = 5$  from some population, send them an offer, and  $x = 2$  respond. What is our best guess of the response probability?

$$L(\pi) = \binom{5}{2} \pi^2 (1-\pi)^3$$

- Let  $l(\pi) = \log(L(\pi))$ , *log-likelihood*

$$l(\pi) = \log(10) + 2 \log(\pi) + 3 \log(1-\pi)$$

$$\frac{dl(\pi)}{d\pi} = \frac{2}{\pi} - \frac{3}{1-\pi} = 0 \Rightarrow \hat{\pi} = \frac{2}{5}$$

## Maximum Likelihood Estimation of Means

Suppose we draw a random sample of size  $n = 3$  from a normal population with unknown mean  $\mu$  and known variance  $\sigma^2 = 5$ . The observed values are  $x_1 = 4, x_2 = 5, x_3 = 6$ . What is the best guess of  $\mu$ ?

Pick  $\mu$  so that the probability of observing the three values given  $\mu$  is maximized:

$$f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(x-\mu)^2\right]$$

$$L(\mu) = \prod_{i=1}^3 \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(x_i - \mu)^2\right]\right)$$

$$l(\mu) = \sum_{i=1}^3 \left[-\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(x_i - \mu)^2\right]$$

$$= -k_1 - k_2 \sum_{i=1}^3 (x_i - \mu)^2$$

$$\frac{dl(\mu)}{d\mu} = 2k_2 \sum_{i=1}^3 (x_i - \mu) = 0 \Rightarrow \mu = \frac{1}{3} \sum_{i=1}^3 x_i$$

$$\text{Note: } \frac{d^2 l(\mu)}{d\mu^2} = -6k_2 < 0 \text{ thus maximum}$$

## Likelihood Function for Logistic Regression

$$\log\left(\frac{\pi_i}{1-\pi_i}\right) = \alpha + \beta x_i$$

$$\text{Probability distribution: } f_i(y_i) = \pi_i^{y_i} (1-\pi_i)^{1-y_i}$$

$$\Rightarrow f(y_1, \dots, y_n) = \prod_{i=1}^n \pi_i^{y_i} (1-\pi_i)^{1-y_i}$$

$$\log(f) = l = \log(L) = \text{Log-Likelihood}$$

$$\sum_{i=1}^n y_i (\alpha + \beta x_i) - \sum_{i=1}^n \log[1 + \exp(\alpha + \beta x_i)]$$

Maximize this with respect to  $\alpha$  and  $\beta$

## Residual versus null deviance

**Residual deviance:** deviance for full model

**Null deviance:** deviance for the intercept-only model (think of as SST)

Difference between them measures how much variation is explained by the model and plays the role of extra sums of squares. Can test overall significance of the model b/c it has a **Chi-squared distribution**.

```
> fit$null.deviance - fit$deviance
[1] 1095.99
```

## Likelihood-Ratio Test

We can test  $H_0: \beta_{p+1} = \dots = \beta_{p+q} = 0$  with the test statistic  $D_1 - D_2$ , which has a chi square distribution with  $q$  degrees of freedom (the D's represent deviance)

## Predictive Accuracy of Classifiers

- GLMs minimize deviance

-  $AIC = deviance + 2p$  common penalized measure

$$\text{- Classification rate: } \frac{TN+TP}{TN+FP+FN+TP}$$

$$\text{- Recall: } \frac{TP}{FN+TP}$$

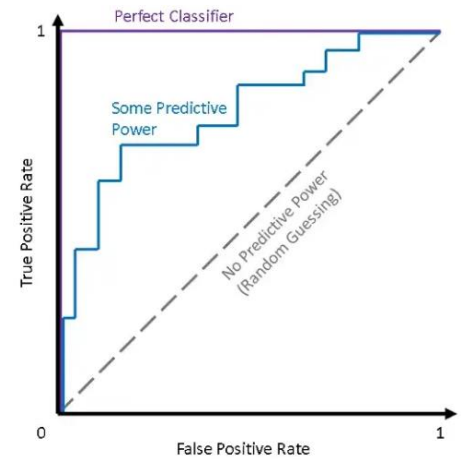
$$\text{- TPR (True Positive Rate) (Sensitivity): } \frac{TP}{TP+FN}$$

$$\text{- FPR (False Positive Rate) (1-Specificity): } \frac{FP}{TN+FP}$$

$$\text{- Precision: } \frac{TP}{FP+TP}$$

$$\text{- F1: } \frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}} = 2 \frac{P \cdot R}{P+R}$$

- **AUC** (Area under curve): ROC curve is produced by calculating and plotting the true positive rate against the false positive rate for a single classifier at a variety of thresholds



## Scoring Model:

- **Proxy behavior:** behavior that has been observed in the past that is similar to future behavior you would like to predict, e.g., response to similar offer sent yesterday

Warning: your model may not work because this is only a proxy behavior. Seasonality, the state of the economy, what competition is doing, etc. usually all affect response

- **Target period:** time period when proxy offer was active

- **Base period:** a period of time prior to the target period. Information from the base period will be used to predict proxy behavior

Performance usually assessed with a *gains table*

1. Find quantiles of predicted values  $\hat{y}$
2. Compute number of responders and revenue by quantile, also averages
3. Compute cumulative counts and revenues by quantile, also averages and lifts

A	B	C	D	E	F	G	H	I	J	K	L	M
Quantile of $\hat{y}$	Num Resp	Rev Amt	Resp Rate	Avg Amt	Num Resp	Num Resp	Rev Amt	Resp Rate	Avg Amt	Resp Rate	Rev Amt	Lift
1	10569	1681	15.901	0.159	15.1	10569	1681	15.901	0.159	15.1	2.98	3.18
2	10569	477	4.0241	0.0451	3.81	21138	2158	19.9742	0.102	9.45	1.91	1.99
3	10568	307	2.3716	0.0290	2.24	31706	2465	22.3458	0.0777	7.05	1.45	1.49
4	10569	201	1.5484	0.0190	1.46	42275	2066	23.8942	0.0631	5.65	1.18	1.19
5	10569	159	1.1608	0.0150	1.10	52844	2825	25.0549	0.0535	4.74	1	1

- **Columns A and B:** Quantiles of  $\hat{y}$

- **Columns C and D:** number of responders and total revenue by quantile

- **Columns E and F:** response rate and average revenue in quantile, e.g.,  $\frac{1681}{10569} = 15.9\%$  and  $\frac{96728}{10569} = \$15.1$  per contact

- **Column G:** Depth of contacts or cumulative counts, e.g.,  $10,569 + 10,569 = 21,138$

- **Columns H and I:** cumulative responders and revenue by quantile:

Row 2:  $1,681 + 444 = 2,158$  responders and  $159,501 + 40,241 = 199,742$  revenue at 40%

Last row: 2,825 total responders and 250,549 total revenue

- **Columns J and K:** cumulative response rate and average revenue in quantile:

Row 2:  $\frac{2,158}{21,138} = 9.45\%$  and  $\frac{199,742}{21,138} = \$9.45$  per contact at 40%

Last row: **contacting at random** gives 5.35% respond rate, \$4.74 per contact

- **Columns L and M:** lift of model over random guessing, e.g.  $\frac{15.1\%}{5.35\%} = 2.98$  indicates the response rate from using model to pick best 20% of the names is improved by 57% over picking names at random. Revenue more than tripled (3.93)!

### Generalized Logistic Regression Model

Let  $Y$  in  $\{1, 2, \dots, K\}$  be a multinomial r.v. for the outcome with  $K$  values. The probability that observation  $i$  comes from class  $k$  is  $\pi_{ik} = P(Y_i = k)$ , where  $\pi_{i1} + \dots + \pi_{iK} = 1$

#### Approaches:

- *One vs all:* fit  $K$  separate logistic regression models

- *Generalized logit:* fit  $K - 1$  models. Pick class 1 as the base category, but, as with the binary logit, this choice is arbitrary (models equivalent with a different base category)

#### Generalized logit model math:

- Model ( $k = 2, \dots, K$ ):

$$\log\left(\frac{\pi_k}{\pi_1}\right) = \beta_k^T \mathbf{x} \quad (k = 2, \dots, K)$$

- Solving for  $\pi_k$ :

$$\pi_k = \pi_1 \exp(\beta_k^T \mathbf{x}) \quad (k = 2, \dots, K)$$

- The probabilities must sum to 1:

$$1 = \sum_{j=1}^K \pi_j = \pi_1 + \sum_{(j=2)}^K \pi_1 \exp(\beta_j^T \mathbf{x}) = \pi_1 [1 + \sum_{j=2}^K \exp(\beta_j^T \mathbf{x})]$$

- Solving for  $\pi_1$ :

$$\pi_1 = \frac{1}{1 + \sum_{j=2}^K \exp(\beta_j^T \mathbf{x})}$$

- Substituting back into formula for  $\pi_k$ :

$$\pi_k = \frac{\exp(\beta_k^T \mathbf{x})}{1 + \sum_{j=2}^K \exp(\beta_j^T \mathbf{x})} \quad (k = 2, \dots, K)$$

#### Generalized logit estimation:

- Likelihood function:

$$L(\mathbf{B}) = \prod_{i=1}^n \prod_{k=1}^K \pi_{ik}^{z_{ik}}$$

Where  $z_{ik} = 1$  when  $y_i = k$  and  $z_{ik} = 0$  otherwise

- The log-likelihood is the cost function

$$\log(L(\mathbf{B})) = l(\mathbf{B}) = \sum_{i=1}^n \sum_{k=1}^K z_{ik} \log(\pi_{ik})$$

- Denote optimal values of  $\beta_k$  by  $\mathbf{b}_k$ . Compute estimated probabilities that each observation is in class  $k$  as:

$$p_{i1} = \frac{1}{1 + \sum_{j=2}^K \exp(\mathbf{b}_j^T \mathbf{x}_i)} \text{ and } p_{ik} = \exp(\mathbf{b}_k^T \mathbf{x}^T) p_{i1}$$

- The *maximum likelihood classifier* assigns  $i$  to the class with the largest  $p_{ik}$

Base category is  $y = 1$

$$\log\left(\frac{\pi_2}{\pi_1}\right) = -8.69 + 1.42x_1 + 0.99x_2$$

$$\log\left(\frac{\pi_3}{\pi_1}\right) = -7.64 + 0.66x_1 + 1.52x_2$$

All four models give similar results

- There are three possible base categories and all give equivalent models

Call: multinom(formula = y ~ x1 + x2, data = train)

Coefficients:  
(Intercept) x1 x2  
2 -8.693704 1.4200155 0.9893357  
3 -7.644341 0.6568461 1.5195375  
Residual Deviance: 68.61044

multinom(formula = factor(y, levels = c(2, 1, 3)) ~ x1 + x2, data = train)

Coefficients:  
(Intercept) x1 x2  
1 8.691775 -1.4197191 -0.9891887  
3 1.047763 -0.7629366 0.5302861  
Residual Deviance: 68.61043

multinom(formula = factor(y, levels = 3:1) ~ x1 + x2, data = train)

Coefficients:  
(Intercept) x1 x2  
2 -1.047771 0.7629246 -0.5302758  
1 7.643941 -0.6567638 -1.5194762  
Residual Deviance: 68.61043

- All three models have the same deviance (cost function value)

- Let  $\beta_{ij}$  be slopes for  $\log\left(\frac{\pi_i}{\pi_j}\right) = \beta_{ij}^T \mathbf{x}$ . Note that

$$\frac{\pi_i}{\pi_j} = e^{\beta_{ij}^T \mathbf{x}}$$

- Clearly,  $\beta_{ij} = -\beta_{ji}$ , e.g.,

$$\log\left(\frac{\pi_3}{\pi_1}\right) = \beta_{31}^T \mathbf{x} \Rightarrow \frac{\pi_3}{\pi_1} = e^{(\beta_{31}^T \mathbf{x})} \Rightarrow \frac{\pi_1}{\pi_3} = e^{-\beta_{31}^T \mathbf{x}}$$

-  $\beta_{23} = \beta_{21} - \beta_{31}$ , e.g.,

$$0.7629246 = 1.4200155 - 0.6568461$$

$$\frac{\pi_2}{\pi_3} = \frac{\pi_2}{\pi_3} * \frac{\pi_1}{\pi_1} = \frac{\pi_2}{\pi_1} * \frac{\pi_1}{\pi_3} = e^{\beta_{21}^T \mathbf{x}} * e^{-\beta_{31}^T \mathbf{x}} = e^{(\beta_{21} - \beta_{31})^T \mathbf{x}}$$

#### Making predictions

- **Binary logit:** individual chooses between two options and selects the one that provides greater utility

- **Multinomial logit:** individual chooses among more than two alternatives and selects the one that provides the greatest utility

- **Ordered logit:** individual reveals the strength of his or her preferences with respect to a single outcome

- **Conditional logit:** allows variables that vary across alternatives and possibly across the individuals as well, e.g., choice of mode of transportation (e.g., train, bus, car). Characteristics or attributes of these include time waiting, how long it takes to get to work, and cost.

- **Log-linear analysis:** class of models that subsumes (includes or absorbs) the logit models and more.

#### Poisson

- The PMF of a Poisson distribution:

$$P(Y = y) = \frac{\mu^y e^{-\mu}}{y!}, \mu > 0, y = 0, 1, 2, \dots$$

- Suppose we observe a random sample of ordered pairs  $(x_i, y_i)$ ,  $i = 1, \dots, n$  from where  $y_i$  has a Poisson distribution with mean  $\mu_i$ . More generally,  $x_i$  could be a vector of predictors. Assume (log link function)  $\log(\mu_i) \rightarrow$  rate ratio

$$\ln(RR_{\text{source}}) = \ln\left(\frac{PATE(y_{\text{source}})}{PATE(y_{\text{no source}})}\right) = \ln(PATE|y_{\text{source}}) - \ln(PATE|y_{\text{no source}}) = [b_0 + b_1(1) + b_2 AGE] - [b_0 + 1] = b_1$$

$$s.t. RR = e^{\ln(RR)} = e^{b_1}$$

$$\log(\mu_i) = \eta_i = \alpha + \beta x_i \Leftrightarrow \mu_i = e^{\eta_i}$$

- The likelihood is:

$$L = \prod_{i=1}^n \frac{e^{-\mu_i} \mu_i^{y_i}}{y_i!}$$

- The log-likelihood of a Poisson model is:

$$l = \log(L) = \sum_{i=1}^n [y_i \log(\mu_i) - \mu_i - \log(y_i!)]$$

- The log-likelihood of a saturated Poisson model is:

$$\hat{\mu}_i \equiv y_i \Rightarrow l_s = \sum_{i=1}^n [y_i \log(y_i) - y_i - \log(y_i!)]$$

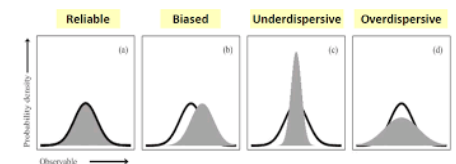
- Let  $\hat{\mu}_i \equiv \exp(\hat{\alpha} + \hat{\beta} x_i)$  be the MLEs. The deviance is:

$$D^2 = -2(l - l_s) = -2 \left[ \sum_{i=1}^n y_i \log\left(\frac{\hat{\mu}_i}{y_i}\right) + \sum_{i=1}^n (y_i - \hat{\mu}_i) \right]$$

Where  $y_i \log\left(\frac{\hat{\mu}_i}{y_i}\right) = 0$  for  $y_i = 0$

#### Beyond Poisson

- Recall that for Poisson  $Y$ ,  $E(Y) = V(Y) = \mu$ , but in practice we may find  $E(Y) < V(Y)$ , called the problem of **overdispersion** (or **underdispersion**  $E(Y) > V(Y)$ )



- **Negative binomial distribution (NBD) regression** models allow for overdispersion

- Beyond overdispersion, we often observe too many zero values for either Poisson or NBD. What to do?

*Zero-inflated Poisson (ZIP)* assumes a mixture distribution

$$P(Y = 0) = \pi + (1 - \pi)e^{-\mu},$$

$$P(Y = y) = (1 - \pi) \frac{e^{-\mu} \mu^y}{y!}, \mu > 0, y = 1, 2, \dots$$

Where  $\pi \in [0, 1]$  is the probability of extra (structural) zeros

*Zero-inflated negative binomial (ZINB)*

*Hurdle models*

- GLMs can accommodate other distributions including exponential and gamma

Gamma regression

- The PMF of a gamma distribution is:

f(y) = 1 / (beta^alpha \* Gamma(alpha)) \* y^(alpha-1) \* e^(-y/beta), alpha > 0, beta > 0, y > 0

- Where alpha is the shape parameter and beta is the scale parameter. The mean of a gamma distribution is mu = alpha\*beta, so beta = mu/alpha. The pdf can be written in terms of mu:

f(y) = (alpha/mu)^alpha \* Gamma(alpha)^-1 \* y^(alpha-1) \* e^(-y\*alpha/mu)

- The log pdf is given by:

log(f(y)) = alpha\*log(alpha) - alpha\*log(mu) - log[Gamma(alpha)] + (alpha-1)\*log(y) - y\*alpha/mu = alpha\*(-y/mu - log(mu)) - log(Gamma(alpha)) + alpha\*log(alpha) - log(y)

- The log-likelihood is:

l = sum\_{i=1}^n [alpha\*(-y\_i/mu\_i - log(mu\_i)) - log(Gamma(alpha)) + alpha\*log(alpha\*y\_i) - log(y\_i)]

- The log-likelihood of the saturate model is:

l\_s = sum\_{i=1}^n [alpha\*(-1 - log(y\_i)) - log(Gamma(alpha)) + alpha\*log(alpha\*y\_i) - log(y\_i)]

- The deviance is -2(l - l\_s), but notice how the last three terms of l and l\_s are identical and thus cancel out. Thus the deviance is:

-2(l - l\_s) = -2\*alpha\*sum\_{i=1}^n (-y\_i/mu - log(mu\_i) + 1 + log(y\_i)) = -2\*alpha\*sum\_{i=1}^n (log(-y\_i/mu\_i) - y\_i/mu\_i)

Regression Terms and Symbols:

Table with 4 columns: Term, ACT, JWHT, Other. Rows include Sum of squared errors, Total sum of squares, Mean squared error, and Residual Standard Error.

SSE = sum\_{i=1}^n [y\_i - f(X\_i)]^2

MSE = SSE/n

R^2 = 1 - SSE/SST = 1 - RSS/((n-1)\*var(y))

Penalized Estimates

R^2 = 1 - (SSE/(n-p-1)) / (SST/(n-1)) or AIC = deviance + 2p

Multicollinearity

Pipe: e.g., x1 -> x2 -> y. If you are studying x1 -> y then do not control for x2

Latent construct: predictors manifestations of common, underlying latent construct, e.g., w -> x1 and w -> x2. Often, estimate w and use it instead of x1 and x2

Back door confound (fork): include control to block back-door path, e.g., if w -> y and w -> x -> y then control for w to study x -> y

Collider: usually do not control for colliders, e.g., if x -> w and y -> w, then do not control for collider w when studying x -> y

Extras

F = (SST - SSE) / (SSE / (n - p - 1)) = (DeltaSSE / Delta df) / (bottom from full model)

- 1. (9 points) The Poisson distribution has the following PMF: P(X = x) = (lambda^x \* e^-lambda) / x! ...
- 2. (30 points) Data from 37 patients receiving a non-depleted allogeneic bone marrow transplant were examined to see which variables were associated with the development of acute graft-versus-host disease (GVHD) ...
- 3. (8 points) Complete the following table giving the Kaplan-Meier estimates of the survival function and retention rate.
- 4. (6 points) Multiple logistic regression was used to construct a prognostic index to predict significant coronary artery disease from data on 348 patients with valvular heart disease ...
- 5. The plot below is from a KM model of the newspaper data, stratifying by the publication type (print+digital versus digital only). What does the plot tell you?

