

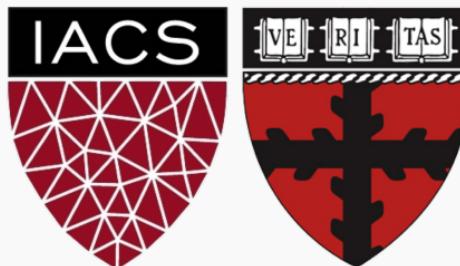
# Advanced Section #3: GLMs: Logistic Regression and Beyond

(Stat 139)  
(Stat 244)

Kevin Rader\*

\*special thanks to **Nick Stern** for help in original development

CS109A Introduction to Data Science  
Pavlos Protopapas, Kevin Rader, and Chris Tanner



# Outline

---

## 1. Motivation

- Limitations of linear regression

## 2. Anatomy

- Exponential Dispersion Family (EDF)
- Link function

## 3. Maximum Likelihood Estimation for GLM's

- Fischer Scoring



---

# Motivation



# Motivation

observation  $\rightarrow \mathbf{x}_i^T$  row in  $X$  matrix  
dot product

Linear regression framework:

$$\xrightarrow{\text{single observation}} y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i$$

$$\mathbf{x}_i = \begin{pmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \quad (p+1) \times 1$$

$$\xrightarrow{\text{all "n" observations}} \boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \quad (p+1) \times 1$$

Assumptions:

- Linearity: Linear relationship between expected value and predictors
- Normality: Residuals are normally distributed about expected value
- Homoskedasticity: Residuals have constant variance  $\sigma^2$
- Independence: Observations are independent of one another

check with plots. Need these to trust "inferences"  
(t-based)

# Motivation

Expressed mathematically...

- Linearity

$$\mathbb{E}[y_i] = \underline{x_i^T \beta}$$

- Normality

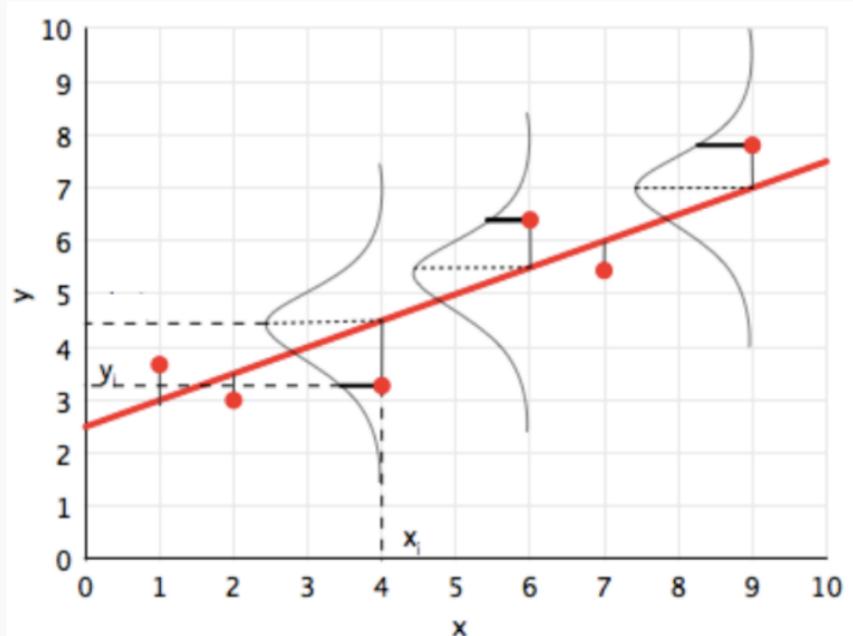
$$\underline{y_i} \sim \mathcal{N}(x_i^T \beta, \sigma^2)$$

- Homoskedasticity

$$\sigma^2 \text{ (instead of) } \underline{\sigma_i^2}$$

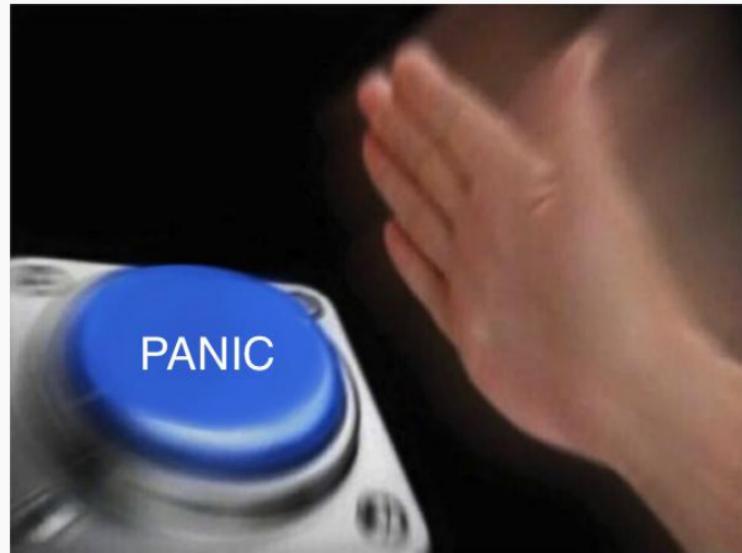
- Independence

$$p(y_i | y_j) = p(y_i) \text{ for } i \neq j$$



# Motivation

What happens when our assumptions break down?

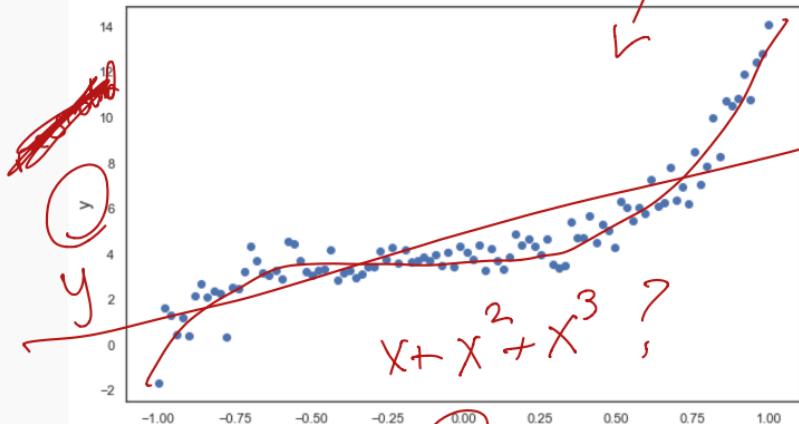


# Motivation

We have options within the framework of linear regression

bootstrap  
to do  
inferences!

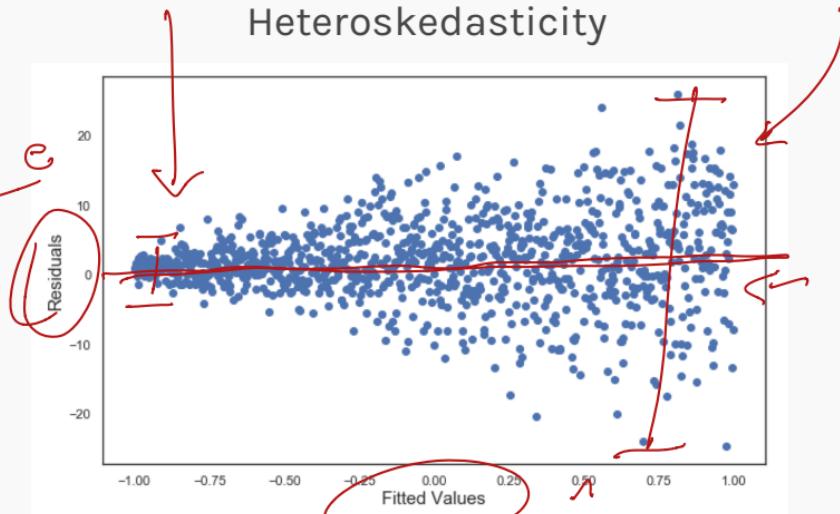
Nonlinearity



Transform X or Y

(Ex: Polynomial Regression)

Heteroskedasticity



Weight observations

(Ex: WLS Regression)



# Motivation

But assuming Normality can be pretty limiting...

Consider modeling the following random variables:

- Whether a coin flip is heads or tails (Bernoulli)  $\leftarrow$  bad
- Count of tropical storms in a given year (Poisson)  $\leftarrow$  linear is reasonable, but not best
- Time between stochastic events that occur w/ constant rate (gamma)
- Vote counts for multiple candidates in a poll (multinomial)  $\leftarrow$  bad

# Motivation

We can extend the framework for linear regression.

Enter:

## Generalized Linear Models

(GLMs)

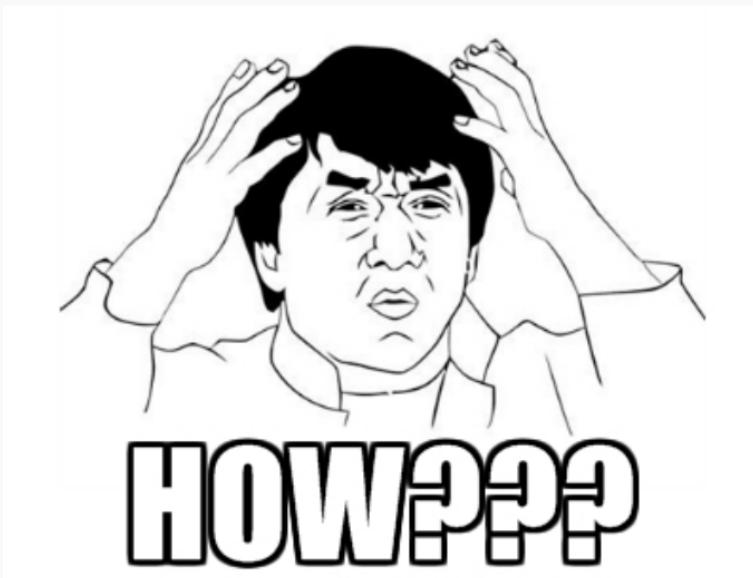
Relaxes:

- Normality assumption
- Homoskedasticity assumption

OLS  
Logistic Regression } special cases

# Motivation

---



---

# Anatomy of our R.V.'s Distribution

↑ γ



# Anatomy

Two adjustments must be made to turn LM into GLM

1. Assume response variable comes from a family of distributions called the exponential dispersion family (EDF).  
*odds; odds<sub>i</sub>*
2. The relationship between expected value and predictors is expressed through a link function.

$$Y_i \sim \text{Bern}(p_i)$$

$$E(Y_i) = p_i \leftarrow$$

logistic Regression is Link:

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots)}} = \frac{e^{\beta_0 + \beta_1 x_1 + \dots}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots}}$$



# Anatomy - EDF Family

The EDF family contains: Normal, Poisson, Gamma, and more!

The probability density function must follow this form:

$$f(y_i|\theta_i) = \exp\left(\frac{y_i\theta_i - b(\theta_i)}{\phi_i} + c(y_i, \phi_i)\right)$$

lef-over  
form of exponential family  
of distribution

Where

$\theta$  - “canonical parameter”  $\leftarrow \mu$

$\phi$  - “dispersion parameter”  $\leftarrow \sigma^2$

$b(\theta)$  - “cumulant function”

$c(y, \phi)$  - “normalization factor”



# Anatomy - EDF Family

Example: representing Bernoulli distribution in EDF form.

PDF of a Bernoulli random variable:

$$\log = \ln$$

$$f(y_i | p_i) = p_i^{y_i} (1 - p_i)^{1 - y_i}$$

Taking the log and then exponentiating (to cancel each other out) gives:

$$f(y_i | p_i) = \exp(y_i \log(p_i) + (1 - y_i) \log(1 - p_i))$$

Rearranging terms...

$$f(y_i | p_i) = \exp\left(y_i \log\left(\frac{p_i}{1 - p_i}\right) + \log(1 - p_i)\right)$$

framework  
of  
exponential family



# Anatomy - EDF Family

Comparing:

$$f(y_i | p_i) = \exp\left(y_i \log\left(\frac{p_i}{1-p_i}\right) + \log(1-p_i)\right) \quad \text{vs.} \quad f(y_i | \theta_i) = \exp\left(y_i \theta_i - b(\theta_i) + c(y_i, \phi_i)\right)$$

Choosing:

$$\theta_i = \log\left(\frac{p_i}{1-p_i}\right)$$



$$b(\theta_i) = \log(1 + e^{\theta_i})$$

$$c(y_i, \phi_i) = 0$$

And we recover the EDF form of the Bernoulli distribution

# Anatomy - EDF Family

The EDF family has some useful properties. Namely:

1.  $\mathbb{E}[y_i] \equiv \mu_i = b'(\theta_i)$
2.  $Var[y_i] = \phi_i b''(\theta_i)$

(the proofs for these identities are in the notes)

Plugging in the values we obtained for Bernoulli, we get back:

$$\mathbb{E}[y_i] = p_i, \quad Var[y_i] = p_i(1 - p_i)$$

# Anatomy - Link Function

---

Time to talk about the link function



# Anatomy - Link Function

Recall from linear regression that:

$$\mu_i = \underline{x_i^T \beta}$$

Does this work for the Bernoulli distribution?

why is this bad?

LHS:  $p_i \in (0, 1)$   $\mu_i = p_i = x_i^T \beta$   $\leftarrow$  unbounded R.H.S.

Solution: wrap the expectation in a function called the **link function**:

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

$$\ln\left(\frac{p}{1-p}\right) \leftarrow \text{un bounded}$$

\*For the Bernoulli distribution, the link function is the “logit” function (hence “logistic” regression)



# Anatomy - Link Function



Link functions are a choice, not a property. A good choice is: CDF

1. Differentiable (implies "smoothness")

2. Monotonic (guarantees invertibility)

1. Typically increasing so that  $\mu$  increases with  $\eta$

3. Expands the range of  $\mu$  to the entire real line

Example: Logit function for Bernoulli

$$g(\mu_i) = g(p_i) = \log\left(\frac{p_i}{1-p_i}\right)$$

interpretable  
to boot

for an  
unbounded  
R.V.  
→ {probit }  
regression  
uses  $\phi^{-1}$  instead  
of "logistic")

inverse CDF  
of standard  
normal  
distribution.



# Anatomy - Link Function

Logit function for Bernoulli looks familiar...

$$g(p_i) = \log\left(\frac{p_i}{1-p_i}\right) = \underline{\theta_i}$$

"canonical parameter" in the exponential family of distributions.

Choosing the link function by setting  $\theta_i = \eta_i$  gives us what is known as the **canonical link function**. Note:

$$\mu_i = b'(\theta_i) \rightarrow \theta_i = b'^{-1}(\mu_i)$$

(derivative of cumulant function must be invertible)

This choice of link, while not always effective, has some nice properties. Take STAT 149 to find out more!

↑  
for estimation



# Anatomy - Link Function

Here are some more examples (fun exercises at home)

Distribution $f(y_i \theta_i)$	Mean Function $\mu_i = b'(\theta_i)$	Canonical Link $\theta_i = g(\mu_i)$
Normal	$\theta_i$	$\mu_i$ ↪
Bernoulli/Binomial	$\frac{e^{\theta_i}}{1 + e^{\theta_i}}$	$\log\left(\frac{\mu_i}{1 - \mu_i}\right)$ ↪
Poisson	$e^{\theta_i}$	$\log(\mu_i)$ ↪
Gamma	$\frac{-1}{\theta_i}$	$\frac{-1}{\mu_i}$ ↪
Inverse Gaussian	$(-2\theta_i)^{-\frac{1}{2}}$	$\frac{-1}{2\mu_i^2}$ ↪

$\mu_i = \mu_1$   
 $\mu_i = p_i$   
 $\mu_i = \lambda_i$

---

# Maximum Likelihood Estimation



# Maximum Likelihood Estimation

Recall from linear regression – we can estimate our parameters,  $\theta$ , by choosing those that maximize the likelihood,  $L(y|\theta)$ , of the data, where:

$$L(y|\theta) = \prod_i^N p(y_i|\theta_i)$$

Joint P.D.F. of  $y | x, \theta$

In words: likelihood is the probability of observing a set of “N” independent datapoints, given our assumptions about the generative process.

# Maximum Likelihood Estimation

For GLM's we can plug in the PDF of the EDF family:

$$L(y|\theta) = \prod_{i=1}^N \exp\left(\frac{y_i\theta_i - b(\theta_i)}{\phi_i} + c(y_i, \phi_i)\right)$$

How do we maximize this? Differentiate w.r.t.  $\theta$  and set equal to 0.

Recall: taking the log first simplifies our life:

$$\log(L) = \ell(y|\theta) = \sum_{i=1}^N \frac{y_i\theta_i - b(\theta_i)}{\phi_i} + \sum_{i=1}^N c(y_i, \phi_i)$$

← this term drops out when differentiating

$$\ell'(\cdot|\theta) = \sum_{i=1}^N y_i - [b'(\theta)] \text{ set } 0$$

$b[\theta]$  is the  $E(Y)$

# Maximum Likelihood Estimation

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial \beta} \left( \frac{\partial \beta}{\partial \theta} \right)$$

we can obtain the  
 canonical  
 parameter

Through lots of calculus & algebra (see notes), we can find the following form for the derivative of the log-likelihood:

$$\ell'(y|\theta) = \sum_{i=1}^N \frac{1}{Var[y_i]} \frac{\partial \mu_i}{\partial \beta} (y_i - \mu_i)$$

← score equations for the mode

Setting this sum equal to 0 gives us the ~~generalized estimating equations~~:

$$\sum_{i=1}^N \frac{1}{Var[y_i]} \left[ \frac{\partial \mu_i}{\partial \beta} \right] (y_i - \mu_i) = 0$$

$\beta$  is a vector  
 $= 0$  ←  
a  $(p+1)$  of these,  
equations,



# Maximum Likelihood Estimation

When we use the canonical link, this simplifies to the **normal equations**:

$$\sum_{i=1}^N \frac{(y_i - \mu_i)x_i^T}{\phi_i} = 0$$

general

Let's attempt to solve the normal equations for the Bernoulli distribution. Plugging in  $\mu_i$  and  $\phi_i$  we get:

$$\sum_{i=1}^N \left( y_i - \frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}} \right) x_i^T = 0$$

solve for  
beta's to  
get your  
maximum likelihood  
estimates



# Maximum Likelihood Estimation

Sad news: we can't isolate  $\beta$  analytically.

Bernoulli  
and  
(it's canonical)  
link.



rely on iterative  
methods!

gradient descent  
newton's method  
Fisher's method

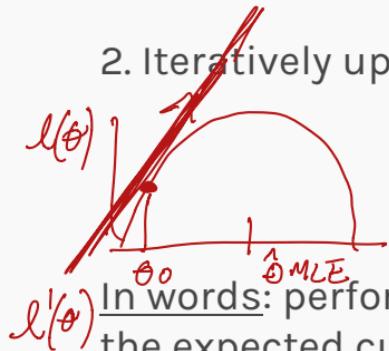
# Maximum Likelihood Estimation

Good news: we can approximate it numerically. One choice of algorithm is the **Fisher Scoring** algorithm.

In order to find the  $\theta$  that maximizes the log-likelihood,  $\ell(y|\theta)$ :

1. Pick a starting value for our parameter,  $\theta_0$ .

2. Iteratively update this value as follows:



In words: perform gradient ascent with a learning rate inversely proportional to the expected curvature of the function at that point.

$$\theta_{i+1} = \theta_i - \frac{\ell'(\theta_i)}{\mathbb{E}[\ell''(\theta_i)]}$$

normal equations  
"locally quadratic functions"  
how curved the likelihood is  
'Fisher's Information'

# Maximum Likelihood Estimation

---

Here are the results of implementing the Fisher Scoring algorithm for simple logistic regression in python:

DEMO



---

# Questions?

