# A Coefficient of Determination for GLMs

by Dabao Zhang

Rearranged by Jae Ho, Chang Spring, 2018

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の 0 0



#### OLS

 $R_{OLS}^2$ ; A Coefficient of Determination for OLS

Consider a linear model;

$$\boldsymbol{y}_{n \times 1} = \boldsymbol{X}_{n \times (p+1)} \boldsymbol{\beta}_{(p+1) \times 1} + \boldsymbol{\epsilon}_{n \times 1} , \quad \boldsymbol{\epsilon} \sim N_n(\boldsymbol{0}, \sigma^2 \mathbf{I}_n)$$

Then we have

$$SSE = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2, \quad SST = \sum_{i=1}^{n} (y_i - \bar{y}_i)^2$$
$$R_{OLS}^2 \stackrel{def}{=} 1 - \frac{SSE}{SST}$$

◆□ → ◆□ → ◆ 三 → ◆ 三 → ○ Q (や 2 / 42



# $\textcircled{O} R^2_{OLS} \text{ measures the goodness of fit}$

- **2** However, for GLMs?
- **③** Some GLMs with well-defined likelihood

 $l(\pmb{y},\mu(\pmb{X}))$  ; The log-likelihood of model  $E(\pmb{y}|\pmb{X})=\mu(\pmb{X})$  for observed data  $(\pmb{y},\pmb{X})$ 

 $l(\pmb{y},\mu(\pmb{1}_n))$  ; The log-likelihood of model  $E(\pmb{y})=\mu(\pmb{1}_n)=\mu\pmb{1}_n$ 

Magee(1990) ; realtionship between  $\mathbb{R}^2$  and the likelihood ratio statistics in the linear regression models

$$R_{LR}^2 = 1 - exp\left\{\frac{2}{n}l(\boldsymbol{y}, \hat{\boldsymbol{\mu}}(\boldsymbol{1}_n)) - \frac{2}{n}l(\boldsymbol{y}, \hat{\boldsymbol{\mu}}(\boldsymbol{X}))\right\}$$

where  $\hat{\mu}$ 's are obtained by MLE for corresponding model. For proof, refer to " $R^2$  Measures based on Wald and Likelihood Ratio Joint Significance Tests, Magee, L. ,1990"

For a logistic regression, perfectly fitted values result in  $l(\boldsymbol{y}, \hat{\mu}(\boldsymbol{X})) = 0$ . That is, samples can be perfectly separated by a linear function.

$$max(R_{LR}^2) = 1 - exp\left\{\frac{2}{n}l(\boldsymbol{y}, \hat{\mu}(\boldsymbol{1}_n))\right\}$$

For example, with balanced case-control data (Bernoulli-distributed with  $\pi = 0.5$ ),  $max(R_{LR}^2) = 0.75$ . That is,  $R_{LR}^2$  is bounded from above by  $l(\boldsymbol{y}, \hat{\mu}(\mathbf{1}_n))$  and will never attain value 1. This is incosistent with the existing concept of  $R^2$ .

Nagelkerke(1991); Suggested the corrected one.

$$R_N^2 \stackrel{def}{=} \frac{R_{LR}^2}{1 - exp\left\{\frac{2}{n}l(\boldsymbol{y}, \hat{\mu}(\mathbf{1}_n))\right\}} = \frac{R_{LR}^2}{max(R_{LR}^2)} \in [0, 1]$$

But still inconsistent with the classical definition of cofficient of determination.

Cameron and Windmeijer (1997); Use the KL divergence to quantify the uncertainty remaining in the response after accounting for predictors.

$$\hat{\theta}_{\boldsymbol{X}} \stackrel{def}{=} argmaxL(\theta|\boldsymbol{X})$$
; MLE given data  $\boldsymbol{X}$ .  
 $\hat{\theta}_{I_n} \stackrel{def}{=} argmaxL(\theta|I_n)$ ; MLE under saturated model.  
Saturated model example ;  $E(Y) = \boldsymbol{\beta}_{n \times 1}$ 

$$\begin{split} \hat{KL}(\theta, \boldsymbol{X}) \stackrel{def}{=} 2\{l(\hat{\theta}_{I_n}|I_n) - l(\hat{\theta}_{\boldsymbol{X}}|\boldsymbol{X})\}\\ R_{KL}^2 \stackrel{def}{=} 1 - \frac{\hat{KL}(\theta, \boldsymbol{X})}{\hat{KL}(\theta, \mathbf{1}_n)} \end{split}$$

Can be interpreted as the deviance reduction ratio due to the changes of predictors in X.

 OLS
 Specified Likelihood
 Limits
 GLM review
 Measurement Proposal
 Empirical Studies
 Real Data
 End

 00
 0000
 000000000
 0000000000
 0000000000
 0000
 000
 0

#### Limits

All the aforementioned generalized coefficients of determination are given through **completely specified likelihood function**.

However, these are not applicable for more general GLMs like quasi-models which specify only the **mean and variance** functions.

#### GLM review

Random Component

• Assume  $(y_i|X_i) \stackrel{iid}{\sim}$  a member of exponential family

2 Model 
$$\mu_i = E(y_i | X_i)$$

3 Link function g;  $g(\mu_i) = X_i \beta$  where  $X_i$  is an  $i^{th}$  row of data matrix.

#### Exponential family

Any random variable y in the exponential family has a probability density function of the form,

$$f(y, \theta, \phi) = \exp\left\{\frac{y\theta - b(\theta)}{\phi} - c(y, \phi)\right\}$$

with loglikelihood,

$$\ell(\theta, y, \phi) = \log\{f(y, \theta, \phi)\} = \frac{y\theta - b(\theta)}{\phi} - c(y, \phi)$$

 $\theta$ ; The canonical parameter of interest

 $\phi$ : A dispersion parameter which plays a role in the variance We use the Bartlett's Identities to derive a general expression for the variance function.

The first and second Bartlett results ensures that under suitable conditions (see Leibniz integral rule), for a density function dependent on  $\theta$ ,  $f_{\theta}(\cdot)$ ,

$$E_{\theta} \left[ \frac{\partial}{\partial \theta} \log(f_{\theta}(y)) \right] = 0$$
$$Var_{\theta} \left[ \frac{\partial}{\partial \theta} \log(f_{\theta}(y)) \right] + E_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \log(f_{\theta}(y)) \right] = 0$$

**Expected value of Y** : Taking the first derivative with respect to  $\theta$  of the log of the density in the exponential family form described above, we have

$$\frac{\partial}{\partial \theta} \log(f(y,\theta,\phi)) = \frac{\partial}{\partial \theta} \left[ \frac{y\theta - b(\theta)}{\phi} - c(y,\phi) \right] = \frac{y - b'(\theta)}{\phi}$$

. Then taking the expected value and setting it equal to zero leads to,

$$E_{\theta} \left[ \frac{y - b'(\theta)}{\phi} \right] = \frac{E_{\theta}[y] - b'(\theta)}{\phi} = 0$$
$$E_{\theta}[y] = b'(\theta)$$

◆□ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ <

Variance of Y: To compute the variance we use the second Bartlett identity,

$$\operatorname{Var}_{\theta}\left[\frac{\partial}{\partial\theta}\left(\frac{y\theta-b(\theta)}{\phi}-c(y,\phi)\right)\right] + \operatorname{E}_{\theta}\left[\frac{\partial^{2}}{\partial\theta^{2}}\left(\frac{y\theta-b(\theta)}{\phi}-c(y,\phi)\right)\right]$$

is 0. Then,

$$\operatorname{Var}_{\theta}\left[\frac{y-b'(\theta)}{\phi}\right] + \operatorname{E}_{\theta}\left[\frac{-b''(\theta)}{\phi}\right] = 0, \quad \operatorname{Var}_{\theta}\left[y\right] = b''(\theta)\phi$$

We have now a relationship between  $\mu$  and  $\theta$ , namely  $\mu = b'(\theta)$  and  $\theta = b'^{-1}(\mu)$ , which allows for a relationship between  $\mu$  and the variance,

 $V(\theta) = b''(\theta)$  = the part of the variance that depends on  $\theta$ 

$$V(\mu) = b''(b'^{-1}(\mu))$$

or,

$$V(\mu) = (b'' \circ b'^{-1})(\mu)$$

#### Variance Function - Bernoulli case

For example, let  $y \sim Bernoulli(p)$  then we express the density of the Bernoulli distribution in exponential family form,

$$f(y) = \exp\left(y\ln\frac{p}{1-p} + \ln(1-p)\right)$$

$$\theta = \ln \frac{p}{1-p} = logit(p),$$

which gives us  $p = \frac{e^{\theta}}{1+e^{\theta}}$ ,  $b(\theta) = \ln(1+e^{\theta})$  and  $b'(\theta) = \frac{e^{\theta}}{1+e^{\theta}} = p = \mu$ ,  $b''(\theta) = \frac{e^{\theta}}{1+e^{\theta}} - \left(\frac{e^{\theta}}{1+e^{\theta}}\right)^2$  OLSSpecified LikelihoodLimitsGLM reviewMeasurement ProposalEmpirical StudiesReal DataEnd000

#### Variance Function - Bernoulli case

This give us

$$V(\mu) = \mu(1 - \mu) = \mu - \mu^2.$$

In this case, dispersion parameter  $\phi = 1$ .

# Measuring Variation Changes Along the Variance Function

#### $\phi$ ; dispersion parameter

 $V(\ \cdot\ )$ ; known variance function in GLMs Consider a simpler measure of uncertainty, var, the variation.

$$var(y_i|X_i) \stackrel{def}{=} \phi V\{\mu(X_i)\} = \phi b''\{b'^{-1}(\mu(X_i))\}$$

where  $\mu(X_i) \stackrel{def}{=} E(y_i|X_i)$ . While the variance function describes the effect of the mean on the variation of the response variable besides the dispersion parameter, Jorgensen(1987) showed that the variance function  $V(\cdot)$  indeed characterizes the underlying exponential distributions.

# Measuring Variation Changes Along the Variance Function

For a response variable with its mean changing from a to b, its variation moves accordingly along the variance function from  $\phi V(a)$  to  $\phi V(b)$ .

Therefore, the variation change of the response variable should be measured using, instead of  $(a - b)^2$ ,

$$d_{V}(a,b) = \left[\int_{a}^{b} \sqrt{1 + \{V'(t)\}^{2}} dt\right]^{2}$$

which is an squared length of  $V(\cdot)$  between V(a) and V(b).



#### Variation change along the variance function





#### Variation change along the variance function



20 / 42



#### Variation change along the variance function



21/42

Variation

 $d_V(a, b)$  can differ dramatically from the Euclidean distance  $(a - b)^2$  when the underlying variance function is nonlinear.

As shown by Morris (1982, 1983), many popularly considered **exponential family distributions**, such as binomial, negative binomial, and gamma distributions, have quadratic variance functions. We assume a general case ;

$$V(\mu) = v_2 \mu^2 + v_1 \mu + v_0 , \ v_2 \neq 0$$

Variation

$$\int_{a}^{b} \sqrt{1 + V'(t)^{2}} dt = \int_{a}^{b} \sqrt{1 + (2v_{2}t + v_{1})^{2}} dt$$
$$= \frac{V'(t)\sqrt{1 + V'(t)^{2}} + \sinh^{-1}(V'(t))}{4v_{2}}\Big|_{a}^{b} = \sqrt{d_{V}(a, b)}$$

When  $v_2 = 0$ , that is, the variance function is linear or constant as in the case of Poisson distribution or normal distribution, we have  $d_V(a, b) = (1 + v_1^2)(b - a)^2$ . 
 OLS
 Specified Likelihood
 Limits
 GLM review
 Measurement Proposal
 Empirical Studies
 Real Data
 End

 00
 0000
 00000000
 00000000
 000000000
 000
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00
 00

$$R_V^2$$

The total variation in Y;  $\sum_{i=1}^{n} d_V\{y_i, \hat{\mu}_i(\mathbf{1}_n)\}$ The model with predictors  $\boldsymbol{X}$  reduces the unexplained variation in Y to  $\sum_{i=1}^{n} d_V\{y_i, \hat{\mu}_i(\boldsymbol{X})\}$ . Therefore, we define the coefficient of determination as

$$R_V^2 = 1 - \frac{\sum_{i=1}^n d_V\{y_i, \hat{\mu}_i(\mathbf{X})\}}{\sum_{i=1}^n d_V\{y_i, \hat{\mu}_i(\mathbf{1}_n)\}}.$$

・ロト ・ 回 ト ・ ヨ ト ・ ヨ ・ うらつ

Appropriate when only mean and variance functions can be specified(like the quasi-models). Therefore,  $\hat{\mu}(\mathbf{X})$  and  $\hat{\mu}(\mathbf{1}_n)$  may be derived from

quasi-likelihood estimators, other than MLE.

# Consistency of $R_V^2$ and its Extension

 $V'(\cdot)$  is **constant** for normal and Poisson distributions. That is,  $R_V^2 = R^2$  for OLS with normal distribution and log-linear model with Poisson distribution.

Similar to the coefficient of determination, the coefficient of **partial determination** is well-defined for linear models.

Measures the proportion of variation in the response variable not explained by a set of predictors that can be explained by an additional set of predictors.

#### Partial Determination

For example, considering two sets of predictors  $X_1$  and  $X_2$  in a linear regression model, we have

$$R^{2}(X_{2}|X_{1}) = 1 - \frac{SSE(X_{1}, X_{2})}{SSE(X_{1})} = \frac{R^{2}(X_{1}, X_{2}) - R^{2}(X_{1})}{1 - R^{2}(X_{1})}$$

(:) Recall that  $1 - R^2(X_1) = \frac{SSE(X_1)}{SST}$ .

measuring the proportion of remaining variation in the response, when including  $X_1$ , explained by  $X_2$ .

0 000000

Measurement Prop

l Empirical Studies 000000000

# Partial Determination

With our definition of  $R_V^2$ , we can easily extend it to a coefficient of partial determination for more general models.

$$R_V^2(X_2|X_1) = \frac{R_V^2(X_1, X_2) - R_V^2(X_1)}{1 - R_V^2(X_1)}$$

◆□▶ ◆□▶ ◆ ■▶ ◆ ■▶ ● ■ 少へで 27 / 42 
 OLS
 Specified Likelihood
 Limits
 GLM review
 Measurement Proposal
 Empirical Studies
 Real Data
 End

 00
 0000
 0
 0000000000
 0000000000
 0000
 000
 0
 0

# Adjusted $R_V^2$

 $R_V^2$  also suffers to increasing numbers of predictors as the classical  $R^2$ .

Therefore, averaged measures of the variation change along the variance function can be used to take consideration of effects caused by different numbers of predictors.

$$R_{V,adj}^2 \stackrel{def}{=} 1 - \frac{\sum_{i=1}^n d_V\{y_i, \hat{\mu}_i(\boldsymbol{X})\}/(n-p)}{\sum_{i=1}^n d_V\{y_i, \hat{\mu}_i(\boldsymbol{1}_n)\}/(n-1)}.$$

・ロト ・ 回 ト ・ ヨ ト ・ ヨ ・ うらつ

# Empirical Comparison

- Random Sampling for 100 times from 2 Populations
- from a member of Exponential family f(·) with means μ<sub>1</sub>, μ<sub>2</sub> respectively
- **3** 50 random samples for each sampling ;  $x_{1,1}, \cdots, x_{1,25} \stackrel{iid}{\sim} f(\cdot \mid \mu_1), \quad x_{2,1}, \cdots, x_{2,25} \stackrel{iid}{\sim} f(\cdot \mid \mu_2)$

 $X_1$ ; The population of the corresponding observation  $X_2$ ; From the standard normal distribution, independent of  $X_1$ and the response variable

I'm gonna draw a sketch for this data structure.



#### **Binomial Model**

Let

$$\mu_1 = \frac{e^{-\beta}}{1 + e^{-\beta}}, \quad \mu_2 = \frac{e^{\beta}}{1 + e^{\beta}}.$$

That is, we model the mean  $\mu(X_1) = \frac{e^{X_1\beta}}{1+e^{X_1\beta}}$ , with  $X_1 = 1$  or -1 indicating the two different populations. Here,  $\mu_1 + \mu_2 = 1$ .

The coefficients of determination are averaged over the 100 datasets for each  $\beta$  ranging from 0 to 5 with step=0.1.

# **Binomial Model**

When  $\beta = 0$ , we have  $\mu_1 = \mu_2$ .

Thus corresponding  $R_{\cdot}^2$ 's are supposed to report 0 (or very close to 0).

On the other hand,  $\mu_1 \to 0$ ,  $\mu_2 \to 1$  as  $\beta \to \infty$ . Which leads to the single-population based samples. Therefore, it is not surprising to observe that, when including the true predictor  $X_1, R_N^2, R_{KL}^2$ , and  $R_V^2$  approach 1.

See Fig.2.

#### Candidates

$$\begin{split} R_{LR}^2 &= 1 - exp \left\{ \frac{2}{n} l(\boldsymbol{y}, \hat{\mu}(\boldsymbol{1}_n)) - \frac{2}{n} l(\boldsymbol{y}, \hat{\mu}(\boldsymbol{X})) \right\} \\ R_N^2 &= \frac{R_{LR}^2}{max(R_{LR}^2)} \\ R_{KL}^2 &= 1 - \frac{\hat{KL}(\beta, \boldsymbol{X})}{\hat{KL}(\beta, \boldsymbol{1}_n)} \\ R_V^2 &= 1 - \frac{\sum_{i=1}^n d_V\{y_i, \hat{\mu}_i(\boldsymbol{X})\}}{\sum_{i=1}^n d_V\{y_i, \hat{\mu}_i(\boldsymbol{1}_n)\}}. \end{split}$$

#### Poisson & Gamma Model

Poisson Model ;  $\mu(X_1) = e^{X_1\beta}$ 

Gamma Model ;  $\mu(X_1) = \frac{100}{2+X_1\beta}$ , shape par.=100

The coefficients of determination are averaged over the 100 datasets for each  $\beta$  ranging from 0 to 5 with step=0.1.

See Fig.3, 4.

◆□▶ ◆□▶ ◆目▶ ◆目▶ 目 のへで 33 / 42

# Poisson & Gamma Model

When  $\beta \to \infty$ ,  $|\mu_1 - \mu_2|$  reaches the maximum in poisson model but is bounded by 50 in the gamma model.

Therefore, it is not surprising to observe that, when including the true predictor  $X_1, R_{KL}^2$  and  $R_V^2$  approach one in the Poisson model, but are barely bounded away from one in the gamma model.

# Why $R_N^2, R_{LR}^2$ falsely claim high $R^2$ ?

We consider the case  $\beta = 5$ . Denote  $X_{21}$  as the subset of  $X_2$  with  $E(X_1) = \mu_1$  and  $X_{22}$  as the subset of  $X_2$  with  $E(X_1) = \mu_2$ .

With the total sample size at 50, we have  $\bar{X}_{21} - \bar{X}_{22} \sim N(0, 0.08)$ . A large value of  $\bar{X}_{21} - \bar{X}_{22}$  implies falsely correlated  $X_1$  and  $X_2$ although  $X_1$  and  $X_2$  are truly independent. Plot of calculated  $R^2$  versus  $\bar{X}_{21} - \bar{X}_{22}$  in fig.5. It shows the robustness of different coefficients of determination when only a falsely correlated  $X_2$  is included.

# Why $R_N^2, R_{LR}^2$ falsely claim high $R^2$ ?

# $R_{LR}^2$ and $R_N^2$ ; Have tendency to severely overstate the variation proportion explained by the poisson or gamma model.

On the other hand, both  $R_{KL}^2$  and  $R_V^2$  are more robust to such false correlation.

 $X_1$  vs  $(X_1, X_2)$ 

 $R^2(X_1, X_2)$  is similar to  $R^2(X_1)$ 

. Though the former models always have slightly higher values.

We defined  $R_{V,adj}^2$  to adjust for increasing numbers of predictors. Since the K-L divergences used to define  $R_{KL}^2$  are indeed deviances, we can also address the issue of different degrees of freedom in the deviances by defining an adjusted version of  $R_{KL}^2$ as follows;

$$R_{KL,adj}^2 = 1 - \frac{\hat{KL}(\boldsymbol{y}, \hat{\boldsymbol{\mu}}(\boldsymbol{X}))/(n-p)}{\hat{KL}(\boldsymbol{y}, \hat{\boldsymbol{\mu}}(\boldsymbol{1}_n))/(n-1)}$$

◆□▶ ◆□▶ ◆目▶ ◆目▶ 目 のへで 37 / 42 OLSSpecified LikelihoodLimitsGLM reviewMeasurement ProposalEmpirical StudiesReal DataEnd0000000000000000000000000000000000000

 $X_1$  vs  $(X_1, X_2)$ 

As shown in Figure 6, both  $R_{V,adj}^2$  and  $R_{KL,adj}^2$  have lowered values when irrelevant predictors  $X_2$  and  $X_3$  are added to the model, with  $X_3$  independently simulated from a standard normal distribution.

# Real Data Analysis

We illustrate the different definitions of  $R^2$  by applying them to the data from a study of nesting horseshoe crabs included in Agresti (1996).

(C ; colors), (SC ; spine conditions), (CW ; carapace widths ), (W ; weights) of 173 female crabs, each with a male crab attached to her in her nest.

This study intended to investigate whether these factors affect the number of satellites, that is, any other males riding near a female crab.

#### Real Data Analysis

We consider binomial models ;

 $oldsymbol{y} = Whether \ a \ female \ crab \ had \ any \ satellites$  and Poisson models ;

#### $y = \# \ of \ satellites \ a \ female \ crab \ had$

As demonstrated in the previous section, both  $R_{LR}^2$  and  $R_N^2$ indeed severely overstate the variation proportion explained by the Poisson models.

#### Real Data Analysis

Quasi-binomial models and quasi-Poisson models can be fitted for the above cases, with  $R_V^2$  and partial  $R_V^2$  calculated.

 $\hat{\phi} = 1.0266$  for the binomial full model, and  $\hat{\phi} = 3.2354$  for the Poisson full model.

Use quasi-Poisson models to allow overdispersion (that is, use of likelihood-based  $\mathbb{R}^2$  is inappropriate)!

For Poisson or quasi-Poisson models we have same regression coefficients. This leads to the same  $R_V^2$  and partial  $R_V^2$ .

# Conclusions

 $\mathbb{R}^2$  ; Measures goodness of fit , also provides a measure of predictability.

 $R^2$  can be used to choose the optimal set of predictors when the model size, that is, the number of predictors, is fixed.

The *adj*. versions can be used to compare models including different numbers of predictors. For this reason,  $R_{adj}^2$  can be also used to help **model selection**, tuning parameter **selection**, etc.

Our extension  $R_{V,adj}^2$  makes all these possible when any statistical model with a well-defined variance function, such as quasi model.