

Disciplina:

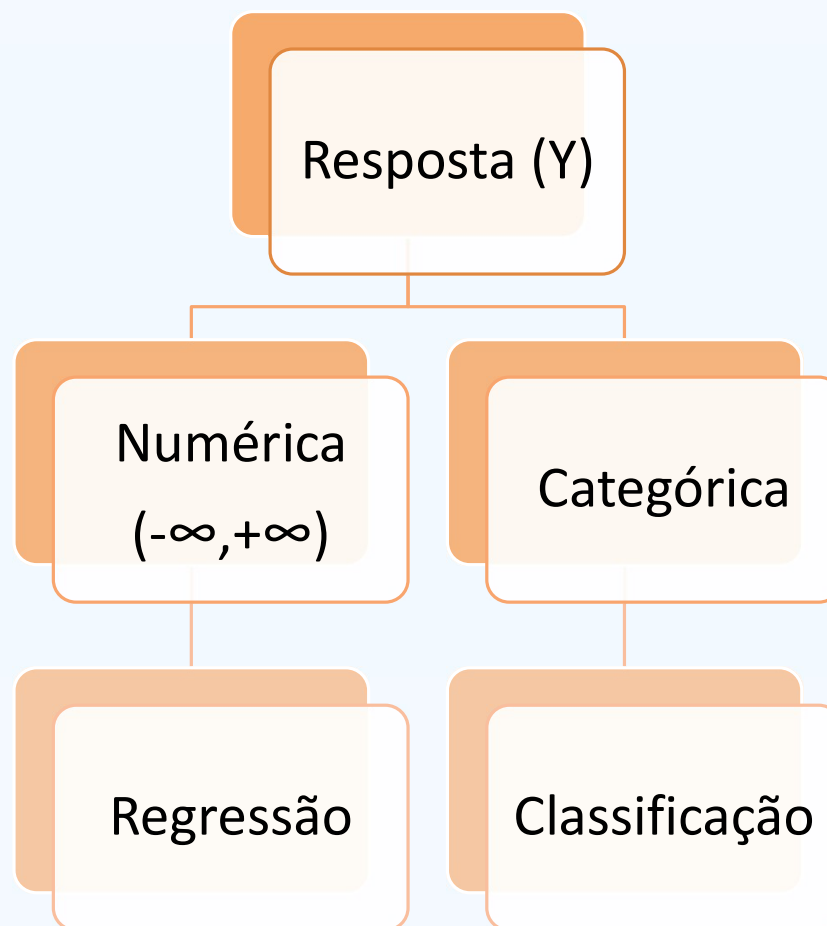
MODELOS ESTATÍSTICOS

Professor: Anaíle Mendes Rabelo

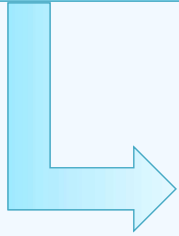


Regressão Logística

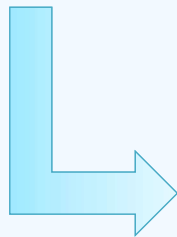
Modelos de Machine Learning



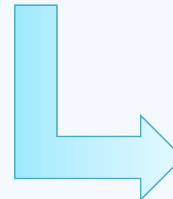
**Variáveis
Categóricas**



Binárias



**Estudar a
probabilidade de
ocorrência**



**Obtemos uma
resposta numérica
(probabilidade)**

**Probabilidade
Valores estão
entre 0 e 1**

Regressão Logística

Quando utilizar a Regressão logística

VARIÁVEL DEPENDENTE COM DISTRIBUIÇÃO BINOMIAL

REPROVAÇÃO NO TESTE DE HOMOCEDASTICIDADE

RESÍDUOS NÃO TEM DISTRIBUIÇÃO NORMAL

Regressão logística

- Técnica de classificação usada para prever uma resposta qualitativa

Exemplos

- Fraudes –
Probabilidade da transação ser fraudulenta ou não.
Classificação usada para prever uma resposta qualitativa (binária)
- Marketplace -
probabilidade do cliente comprar ou não a mercadoria

Pressupostos da Regressão Logística

A variável resposta precisa ser qualitativa, dicotômica ou binária(modelo tradicional)

As preditoras podem ser quantitativas ou categóricas (transformadas em binárias ou Dummy)

Assume que as observações são independentes, que uma não afeta a outra

Por que não uma regressão linear?

Problema da ordenação:

Não podemos utilizar uma regressão linear para prever eventos categóricos.

Ex: Condição médica dos pacientes:

- AVC
- Parada Cardíaca
- Overdose

Podemos simplesmente ordenar e realizar a regressão linear?

1 = AVC

2 = Parada Cardíaca

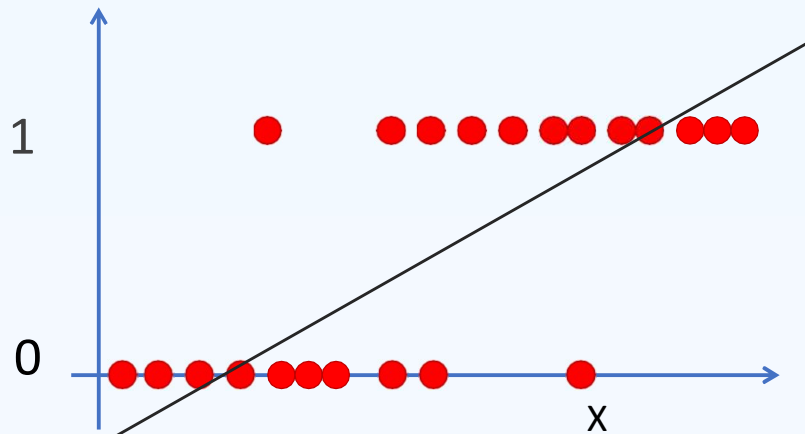
3 = Overdose

Devemos utilizar as variáveis dummy e transformar as variáveis.

Regressão Logística - Definição Teórica

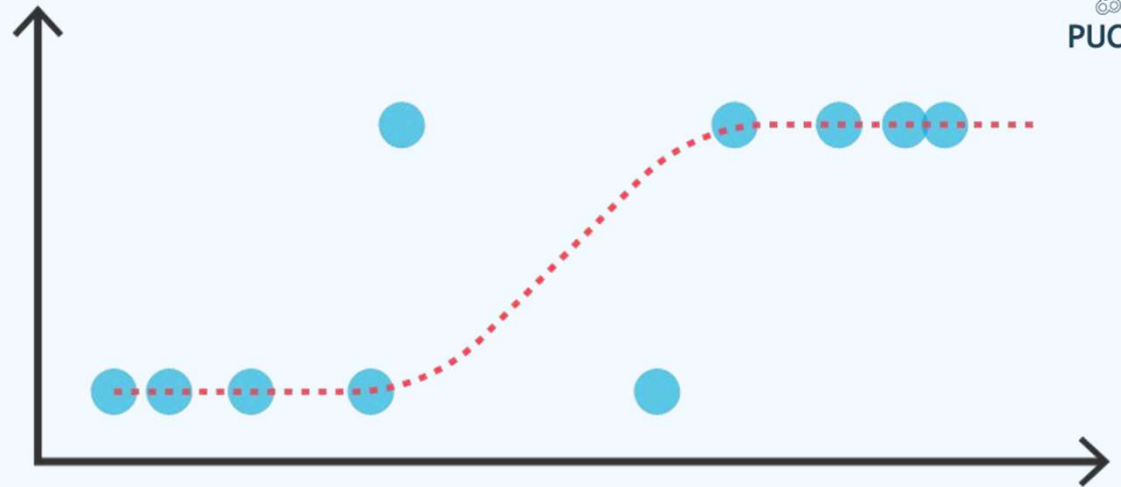
Permite **estimar a probabilidade** associada à **ocorrência de determinado evento** em vista de um conjunto de variáveis preditoras.

- Probabilidade de sucesso (1)
- Probabilidade de fracasso (0)



Ao interpretar Y como probabilidade, temos que realizar transformações para que a resposta de nossa regressão esteja entre 0 e 1

Regressão Logística



Função Logística

- Retorna os valores entre 0 e 1
- Tem formato de “S”

Suponha que o modelo tenha a seguinte forma:

$$Y = X'B + e$$

Em que $X' = [1, X_{i1}, X_{i2}, \dots, X_{in}]$, $B = [\beta_0, \beta_1, \beta_2, \dots, \beta_n]$, e a variável resposta entre 0 e 1.

Assumimos que a variável resposta é uma variável aleatória de Bernoulli, com função de Probabilidade:

y_i	Probabilidade
1	$P(y_i = 1) = p_i$
0	$P(y_i = 0) = 1 - p_i$

Recapitulando - Distribuição de Bernoulli

A variável aleatória Y tem distribuição de Bernoulli se apresenta apenas **dois resultados possíveis**, representados por 0 (fracasso ou negativo) e 1 (sucesso ou positivo). O parâmetro $0 < p < 1$ é a **probabilidade de sucesso**. Dessa forma, a função de probabilidade é

$$p(y) = \begin{cases} 1 - p & \text{se “fracasso” ou } y = 0 \\ p & \text{se “sucesso” ou } y = 1, \end{cases}$$
$$= p^y \cdot (1 - p)^{1-y}, \quad y \in \{0, 1\}.$$

Denotamos por $Y \sim \text{Ber}(p)$.

Com isso, Y apresenta:

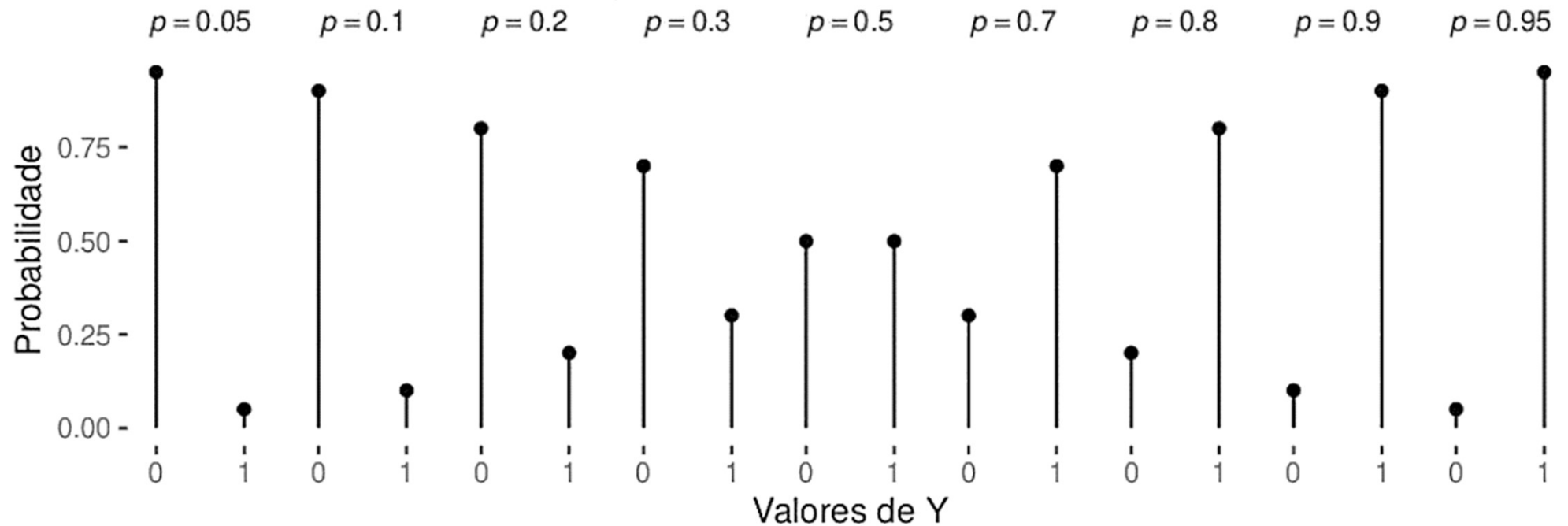
- ▶ $\mu = E(Y) = p$.
- ▶ $\sigma^2 = V(Y) = p \cdot (1 - p)$.

Distribuição de Bernoulli

Onde,

$P \rightarrow$ Probabilidade de sucesso

$1 - p \rightarrow$ Probabilidade de fracasso



Regressão Logística

- Uma vez que $E(y_i) = 1(p_i) + 0(1 - p_i) = p_i$
- Temos que:

$$E(y_i) = x'_i \beta = p_i$$

Logo, a **resposta encontrada na regressão logística** sempre será a **probabilidade de sucesso (1)**

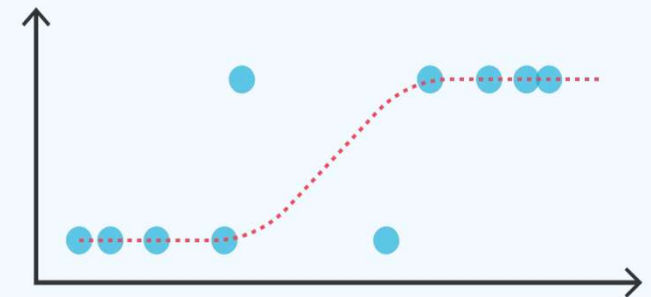
Regressão logística

Resposta (Y)
entre 0 e 1

Requer
transformação
da nossa
função

Função Logit

$$E(y) = \frac{e^{x'\beta}}{1 + e^{x'}} = \frac{1}{1 + e^{-x'\beta}}$$



Conseguimos transformar essa linha em forma de s, em uma linha reta (linearização), para que possa ser possível o cálculo dos coeficientes?

Regressão Logística

A regressão logística pode ser linearizada:

$$\eta = x'\beta$$

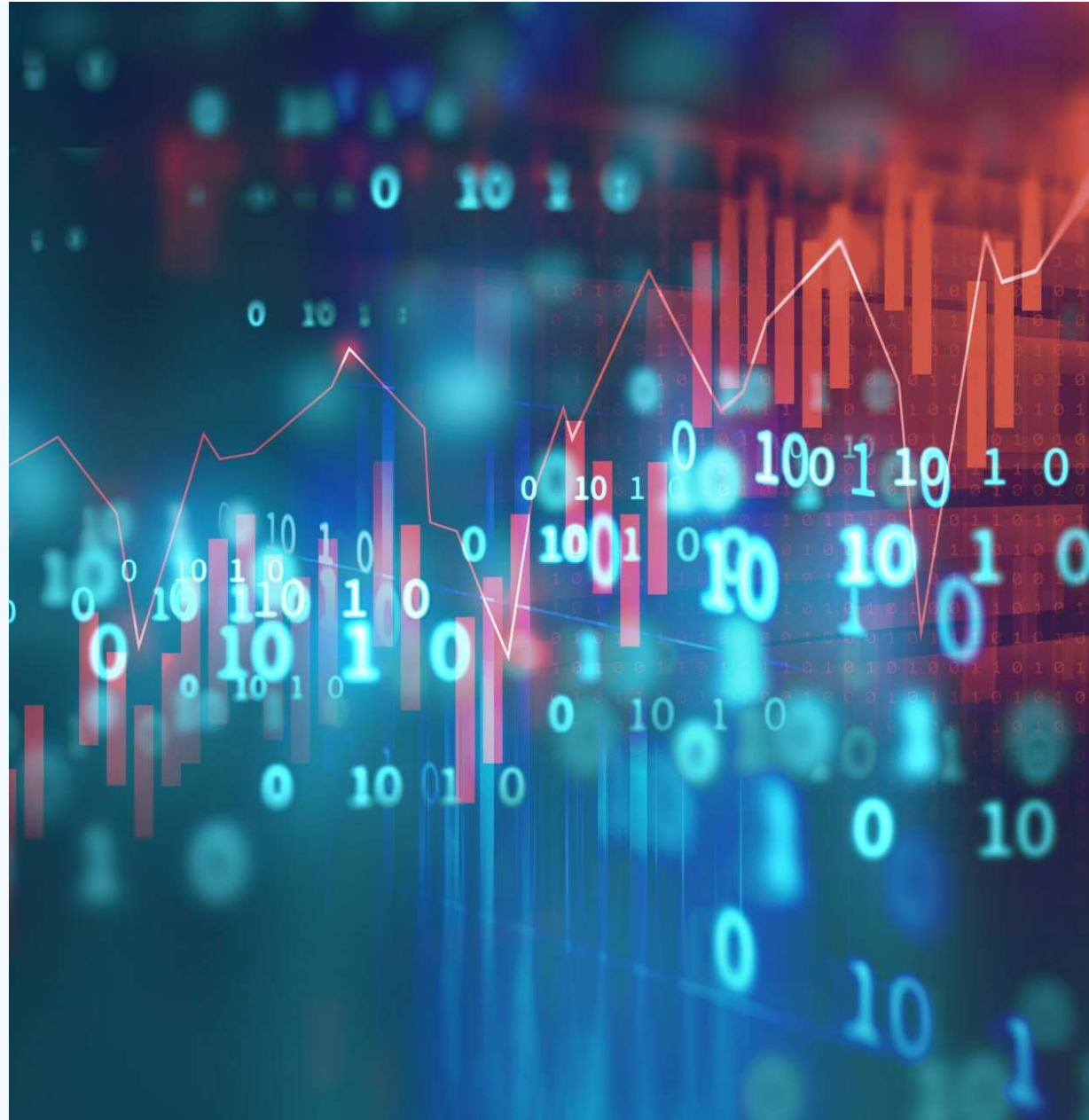
Ser o preditor linear, onde η é definido pela transformação.

$$\eta = \ln \frac{p}{p-1}$$

Essa transformação é frequentemente chamada de **transformação logit** da probabilidade p e a razão $\frac{p}{p-1}$ é chamada de chance odds .



Função ligação que associa os valores esperados da função resposta aos preditores lineares do modelo.



Regressão Linear – Analisando o erro

- Como temos uma resposta binária (0 e 1), temos que os termos de erro só podem ter dois valores:

$$\begin{aligned}\varepsilon_i &= 1 - x'_i\beta, & y_i &= 1 \\ \varepsilon_i &= -x'_i\beta, & y_i &= 0\end{aligned}$$

Logo os **erros não podem ser normais**, e a **variância não é constante**.

$$\begin{aligned}E(\sigma^2_{yi}) &= E\{y_i - E(y_i)\} = (1 - p_i)^2 p_i + (0 - p_i)^2 (1 - p_i) = p_i(1 - p_i) \\ \sigma^2_{yi} &= E(y_i)[1 - E(y_i)]\end{aligned}$$

Estimação de Parâmetros

- A estimação dos parâmetros de $x'_i\beta$ é realizada a partir do método de máxima verossimilhança;
- Como nossos dados seguem a distribuição de Bernoulli, então a distribuição de probabilidade é dada por:

$$f_i(y_i) = p_i^{y_i} [1 - E(p_i)]^{1 - y_i}, \quad i = 1, 2, 3, \dots, n$$

- E cada observação assume o valor de 0 e 1.
- Logo a função de verossimilhança para v.a. independentes pode ser dada por:

$$L(y_1, y_2, y_3, \dots, y_n, \beta) = \prod_{i=1}^n f_i(y_i) = \prod_{i=1}^n p_i^{y_i} [1 - E(p_i)]^{1 - y_i}$$

Estimação de Parâmetros

A forma geral de um modelo de regressão logística é

$$y_i = E(y_i) + \varepsilon_i$$

Em que as observações são variáveis aleatórias independentes de Bernoulli com valores esperados

$$E(y) = p_i = \frac{e^{x'\beta}}{1 + e^{x'\beta}}$$

Interpretação dos Parâmetros

- Considere o caso em que o preditor linear tem apenas uma variável preditora, de forma que o valor ajustado do preditor linear em um valor particular de x , x_i é:

$$\hat{\eta}(x_i) = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

O valor ajustado em $x_i + 1$ é:

$$\begin{aligned}\hat{\eta}(x_i + 1) &= \hat{\beta}_0 + \hat{\beta}_1 (x_i + 1) \\ \hat{\eta}(x_i + 1) &= \hat{\beta}_0 + \hat{\beta}_1 (x_i) + \hat{\beta}_1\end{aligned}$$

- E a diferença dos valores previstos é:

$$\hat{\eta}(x_i + 1) - \hat{\eta}(x_i) = \hat{\beta}_1$$

$$\hat{\eta}(x_i + 1) - \hat{\eta}(x_i) = \ln(odds_{x_i+1}) - \ln(odds_{x_i}) = \ln \frac{odds_{x_i+1}}{odds_{x_i}} = \hat{\beta}_1$$

$$\widehat{O_R} = \frac{odds_{x_i+1}}{odds_{x_i}} = e^{\hat{\beta}_1}$$

- Para as variáveis preditoras binárias, podemos realizar a seguinte análise:

$$\ln \frac{p}{1-p} = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Para $x_i = 0$, temos :

$$\ln \frac{p_0}{1-p_0} = \hat{\beta}_0$$

Para $x_i = 1$, temos :

$$\ln \frac{p_1}{1-p_1} = \hat{\beta}_0 + \hat{\beta}_1$$

Logo:

$$\ln \frac{p_1}{1-p_1} - \ln \frac{p_0}{1-p_0} = \hat{\beta}_0 + \hat{\beta}_1 - \hat{\beta}_0$$

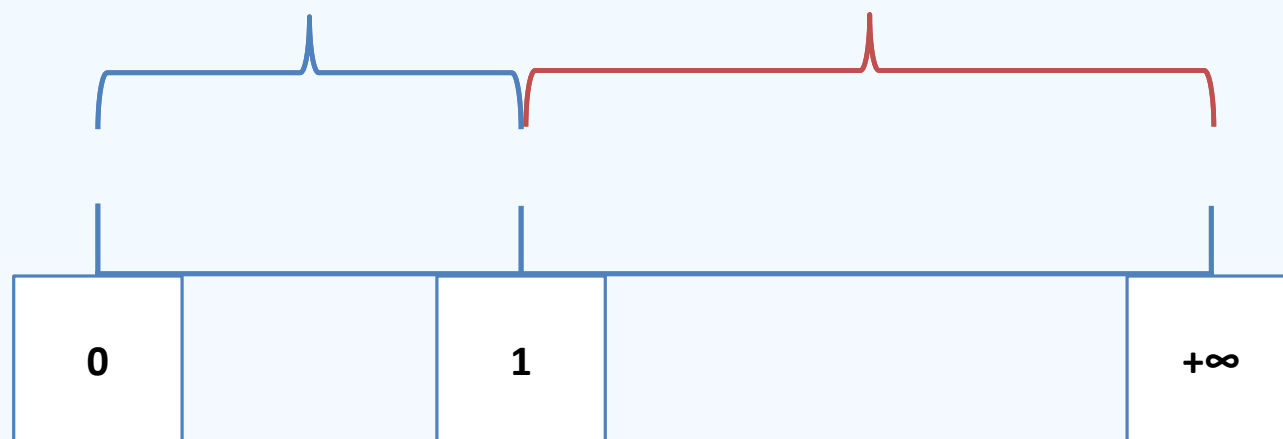
$$\ln \frac{\frac{p_1}{1-p_1}}{\frac{p_0}{1-p_0}} = \hat{\beta}_1 \rightarrow \frac{\frac{p_1}{1-p_1}}{\frac{p_0}{1-p_0}} = e^{\hat{\beta}_1}$$

$$\widehat{O}_R = \frac{odds_{x_i+1}}{odds_{x_i}} = e^{\widehat{\beta}_1}$$

ODDS

Reduz a probabilidade de ocorrência

Aumenta a probabilidade de ocorrência



Exemplo de Interpretação:

Suponha que você esteja estudando a probabilidade de uma pessoa comprar um produto online com base em duas variáveis: idade e gênero. Após ajustar um modelo de regressão logística, você obtém os seguintes resultados:

- Para a variável idade, o odds ratio é 1.05.
- Para a variável gênero (sendo 1 para masculino e 0 para feminino), o odds ratio é 0.8.

Interpretação:

- Para a idade: A cada aumento de uma unidade na idade, as chances de comprar o produto aumentam em 5%.
- Para o gênero: As chances de comprar o produto são 20% menores para homens em comparação com mulheres.

Tipos de resíduos nos GLMs

Tipos de Resíduos nos GLMs

O **resíduo de Pearson** é definido como (r_i^P) é definido como:

$$r_i^P = \frac{y_i - \hat{\mu}_i}{\sqrt{V(\hat{\mu}_i)}},$$

Componente da estatística de Pearson generalizada X_p^2 .

Desvantagem: Muito assimétrico para modelos não normais

Resíduo de Pearson Padronizado

O resíduo de Pearson tem uma versão padronizada, com média 0 e variância aproximadamente 1, definido por:

$$r_i^{P'} = \frac{y_i - \hat{\mu}_i}{\sqrt{\phi V(\hat{\mu}_i)(1 - \hat{h}_{ii})}},$$

Em que \hat{h}_{ii} é o i-ésimo elemento da diagonal da matriz H (de pesos e observações)

Os resíduos de Pearson padronizados apresentam propriedade razoáveis de segunda ordem, mas podem ter distribuições muito divergentes da normal

Resíduo Componente da Deviance

São as raízes quadradas dos componentes da *deviance* com sinal igual a $y_i - \hat{\mu}_i$. Sabe-se que a *deviance* é dada por:

$$D_p = 2 \sum_{i=1}^n [y_i(\tilde{\theta}_i - \hat{\theta}_i)] + b(\hat{\theta}_i) - b(\tilde{\theta}_i)$$

Então temos:

$$r_i^D = \text{sinal}(y_i - \hat{\mu}_i) \cdot \sqrt{2y_i(\tilde{\theta}_i - \hat{\theta}_i) + b(\hat{\theta}_i) - b(\tilde{\theta}_i)}$$

O resíduo r_i^D representa a distância da observação y_i ao valor ajustado $\hat{\mu}_i$, medida na escala da verossimilhança.

Resíduo Componente da Deviance

Quando r_i^D  temos que a i-ésima observação está mal ajustada pelo modelo

O resíduo componente da *deviance* possui uma versão padronizada que é dada por:

$$r_i^{D'} = \frac{r_i^D}{\sqrt{\phi (1 - \widehat{h}_{ii})}},$$

Em que h_{ii} é o i-ésimo elemento da diagonal da matriz H(pesos e observações)

Resíduo Quantílico

O **Resíduo Quantílico** aleatorizado pode ser definido por:

$$r_i^q = \phi^{-1} [F(y_i; \mu_i; \phi)]$$

De forma que se o modelo estiver corretamente especificados, estes resíduos seguiram uma distribuição normal padrão.

Se y_i fo uma v.a. discreta, então $F(y_i; \mu_i; \phi)$ é uma com fda com “saltos” em cada valor de y_i .

Análise de Resíduos

Os resíduos de Pearson e componente da Deviance geralmente **não possuem boas aproximações com a distribuição normal**, ainda que o **modelo ajustado esteja correto**.

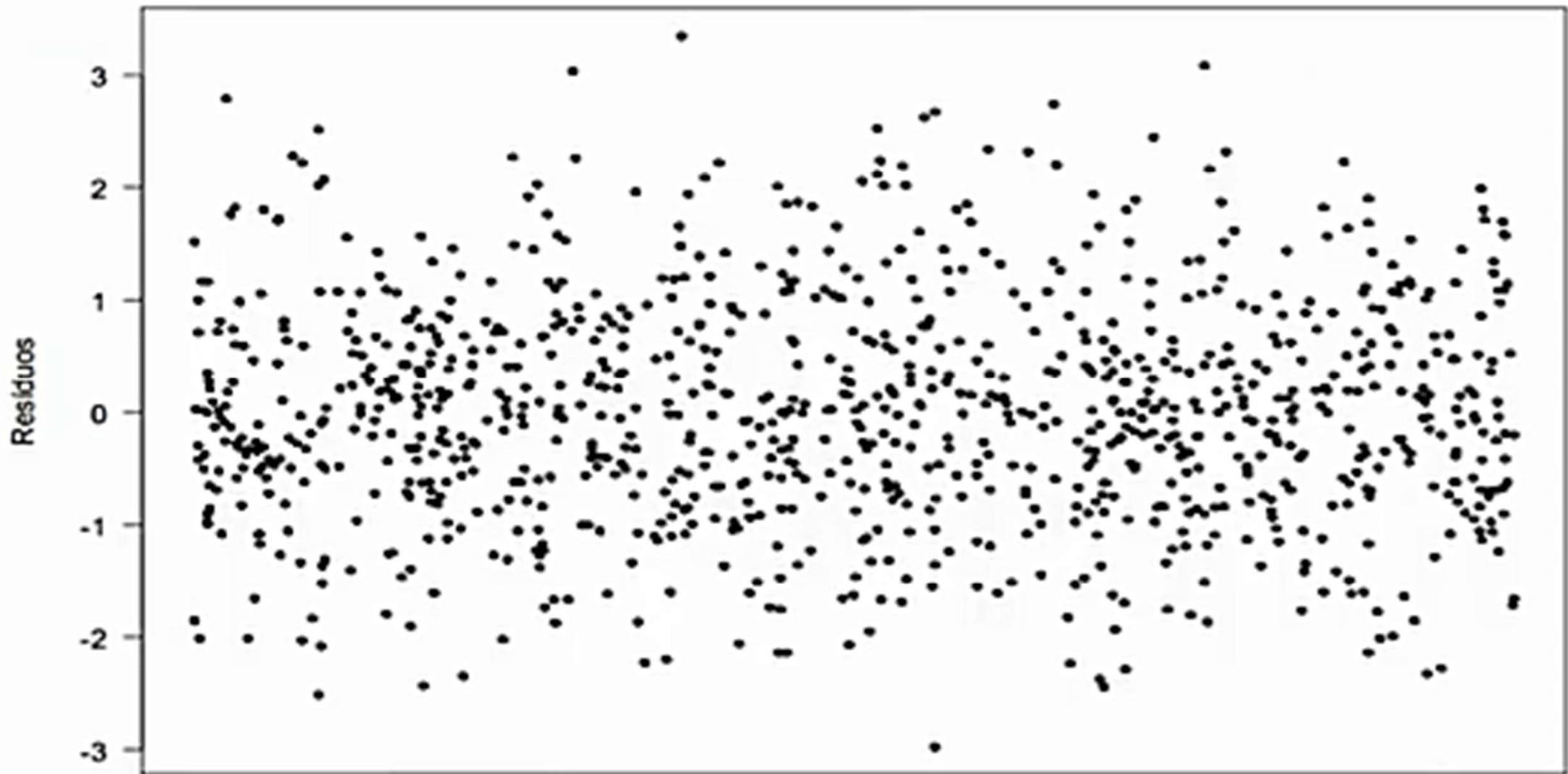
A avaliação da qualidade do ajuste baseado em gráficos de probabilidade (QQ-plot), podem não ser adequados.

Um tipo de resíduo que, por construção, tem **distribuição normal caso o modelo ajustado esteja correto** é o **resíduo quantílico aleatorizado**

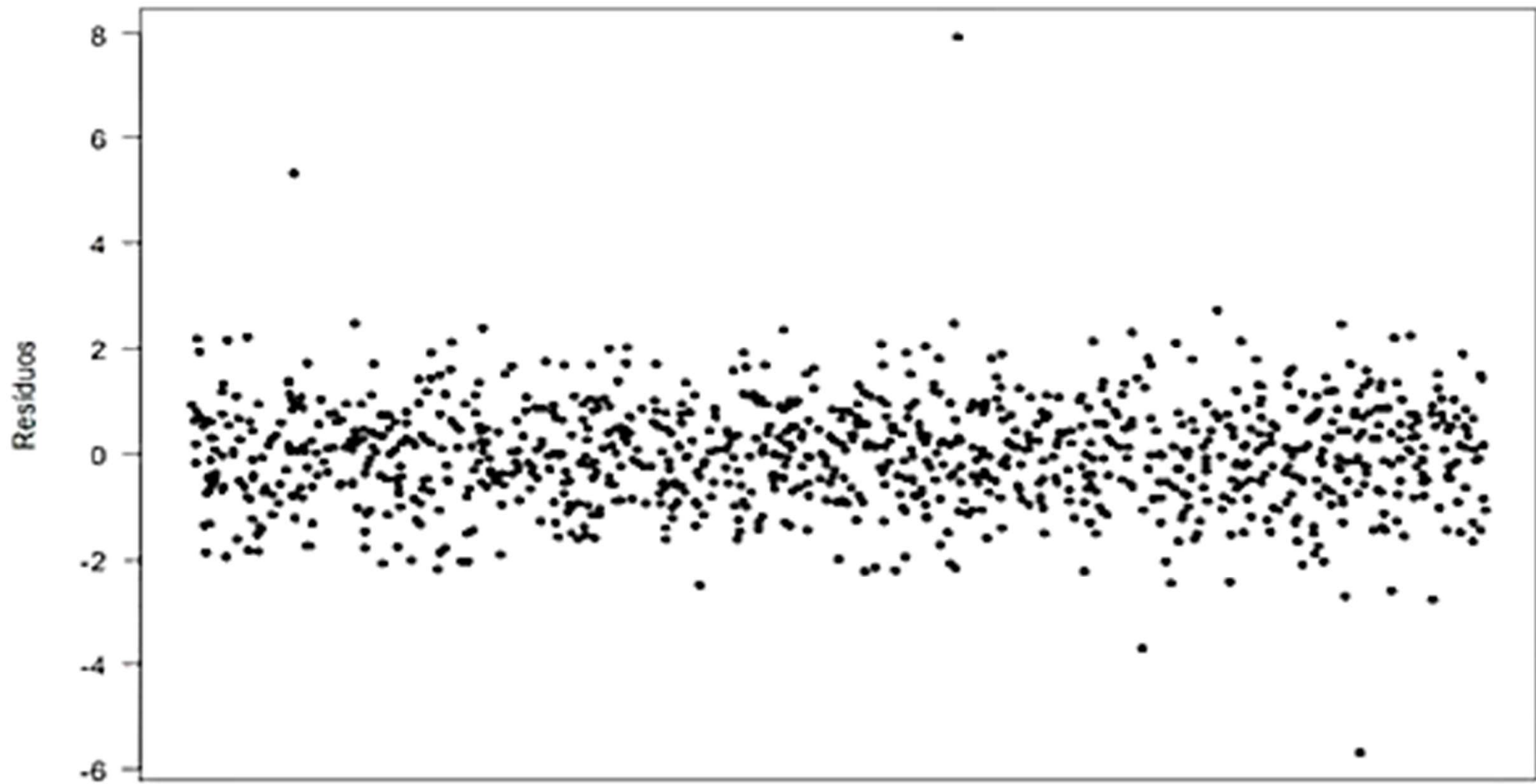
Análise de Resíduos

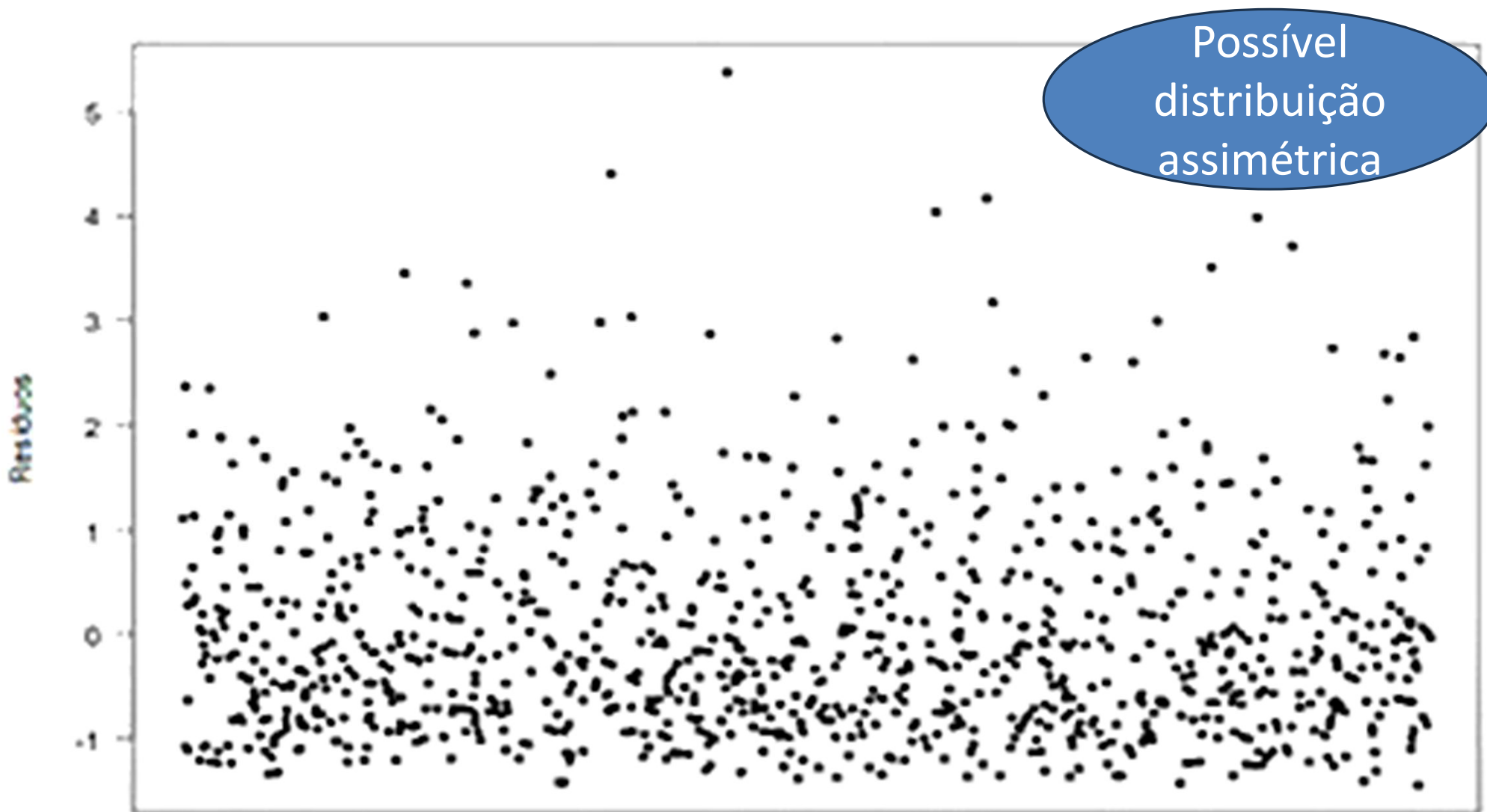
- **Resíduos vs valores ajustados:** Para um modelo bem ajustado, deve-se observar a dispersão aleatória dos pontos, centrada em zero, com média e variâncias constantes e sem valores extremos.
- **Resíduos vs variáveis incluídas no modelo:** padrões não aleatórios indicam que a variável não está inserida corretamente no modelo;
- **Resíduos vs variáveis não incluídas no modelo:** Padrões não aleatórios sinalizam a necessidade (e a forma) de inclusão da variável no modelo;
- **Resíduos vs ordem de coleta de dados:** Padrões não aleatórios indicam dependência das observações gerada pela ordem de coleta (no tempo , no espaço, ...)

Resíduos vs Valores ajustados

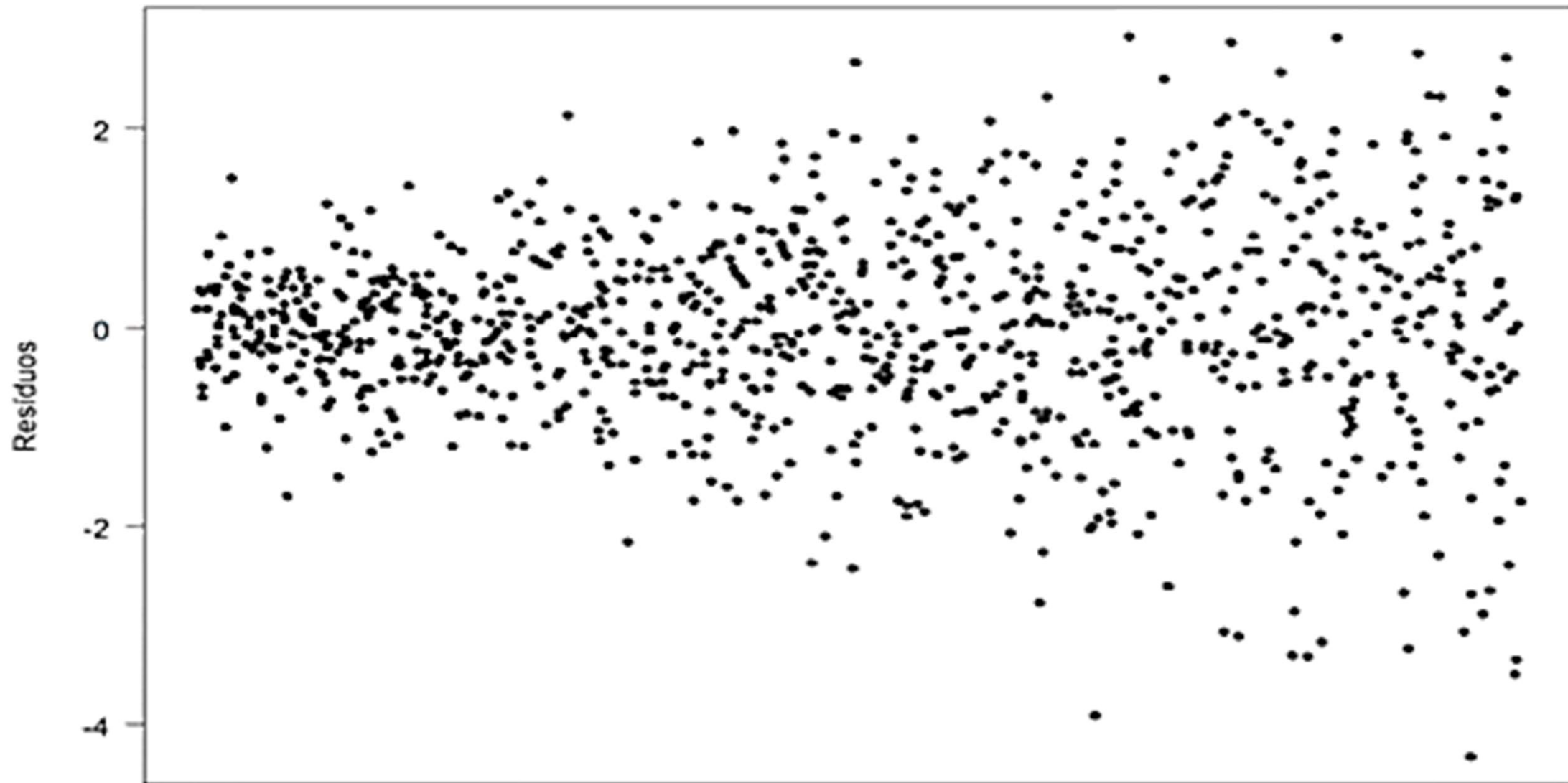


Resíduos vs Valores ajustados

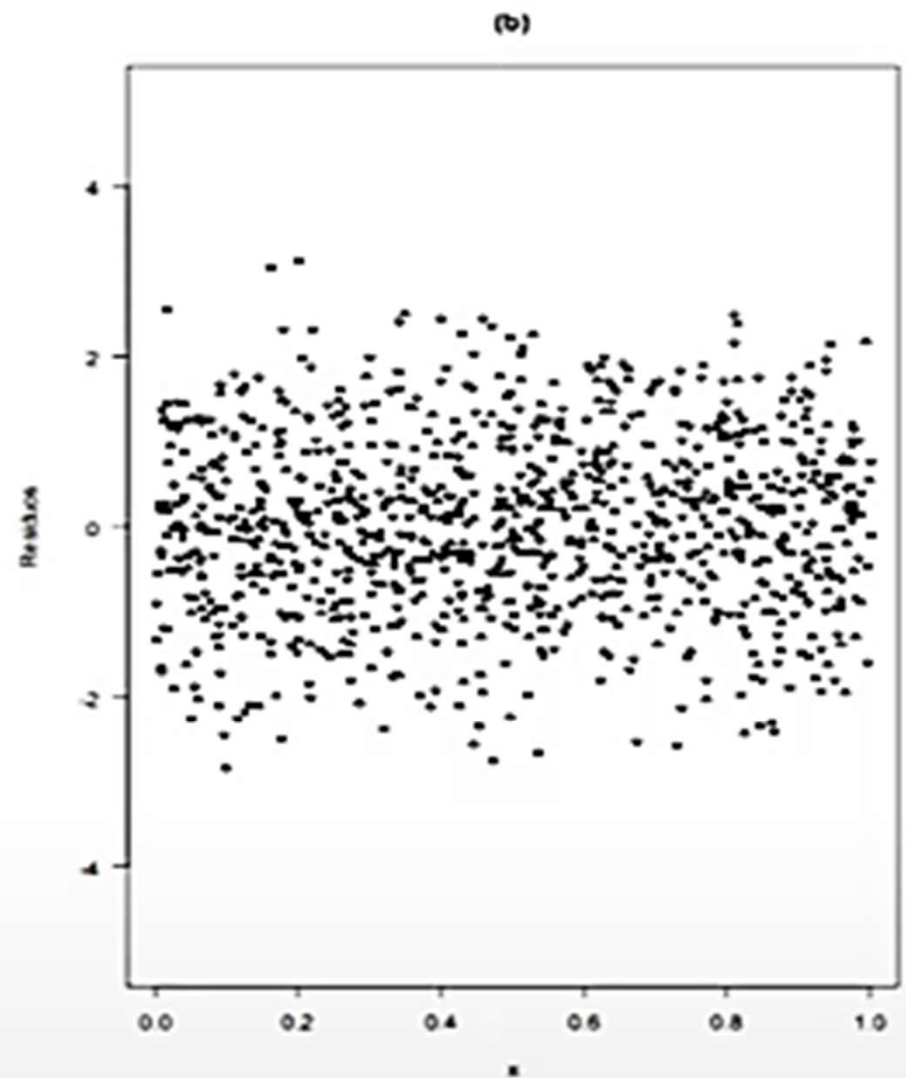
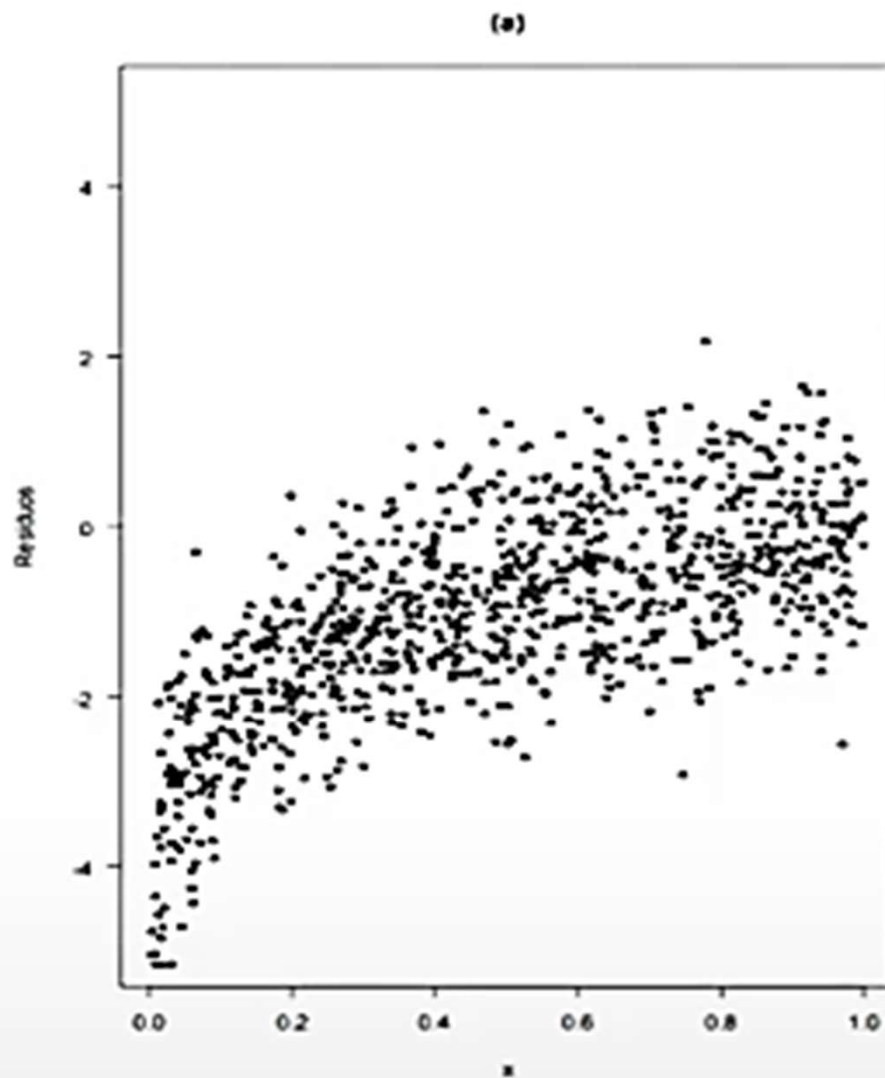




Resíduos vs Valores ajustados



Resíduos vs Variáveis incluídas no modelo



Multicolinearidade

- **Preditores correlacionados com outros preditores**, resulta quando você tem fatores que são, de certa forma, um pouco **redundantes**.
- Ou seja, quando **duas ou mais variáveis independentes em um modelo de regressão encontram-se altamente correlacionadas**
- Examinar a matriz de correlação das variáveis independentes.
 - 0,70 Altamente correlacionadas
 - 0,80 Alerta
- O valor do fator de inflação da variância (VIF), **que mede quanto a variância do coeficiente estimado para uma variável é inflada devido à multicolinearidade** com as outras variáveis independentes.
- VIFs maiores que 10 indicam alta multicolinearidade, enquanto valores entre 5 e 10 podem ser preocupantes.
- A maneira mais simples de lidar com a multicolinearidade é **excluir a variável multicolinear**

PARÂMETROS DOS MODELOS

- **Verificar a significância das variáveis do modelo**
- Teste de hipótese para determinar se a variável preditora do modelo é significativamente relacionada com variável resposta do modelo
 - Teste de Wald
 - Teste de Razão de verossimilhança



Teste de Razão de Verossimilhança

- Compara valores observados x preditos , com e sem determinadas variáveis.
- A comparação dos valores observados x valores preditos é baseada na log verossimilhança
- Pode-se pensar que o valor observado de Y, é um valor predito resultante de um modelo saturado.

$$D = - 2 \ln \left(\frac{\textit{Verossimilhança modelo ajustado}}{\textit{Verossimilhança modelo saturado}} \right)$$

Devaince: para a reg. Logística irá desempenhar o mesmo papel que o SQR na reg. Linear.

Modelo Saturado -> Modelo que se ajusta perfeitamente os dados

Teste de Razão de Verossimilhança

- Utilizando a *deviance* para comparação de modelos que não sejam saturados.

$$G = D(\text{modelo sem a variável}) - D(\text{modelo com a variável})$$



Modelo saturado é o mesmo para os 2.

$$G = -2 \ln \left(\frac{\text{Verossimilhança sem a variável}}{\text{Verossimilhança com a variável}} \right)$$

Teste de Razão de Verossimilhança

H_0 : *O modelo nulo (mais simples) é verdadeiro, ou seja, a inclusão das variáveis adicionais no modelo completo não melhora significativamente o ajuste do modelo.*

H_1 : *O modelo completo é significativamente melhor que o modelo nulo.*

H_0 : *Verossimilhança do modelo Nulo = Verossimilhança do Modelo Completo*

H_1 : *Verossimilhança do modelo Completo > Verossimilhança do Modelo Nulo*

Teste Wald

Obtido por comparação entre a estimativa de máxima verossimilhança do parâmetro $\hat{\beta}_j$ e a estimativa de seu erro padrão.

$$H_0 : \hat{\beta}_1 = 0$$

$$H_1 : \hat{\beta}_1 \neq 0$$

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

A estatística do Teste Wald para a regressão logística é dada por:

$$W_j = \frac{\hat{\beta}_1}{SE(\hat{\beta}_1)}$$

Se não rejeitarmos H_0 , temos que a variável X não explica a variável resposta.

Medidas de qualidade do ajuste do modelo

- Para analisar o desempenho geral do modelo ajustado podemos utilizar vários tipos de Testes de qualidade de ajuste.
- Teste que necessitam dados replicados (múltiplas observações com os mesmos valores para todos os preditores):
 - χ^2 de Pearson
 - Deviance
- O teste hosmer-lemeshow é útil para conjuntos de dados não replicados ou que contém apenas algumas observações replicadas.
 - As observações são agrupadas com base em suas probabilidades estimadas.

Deviance

- Pequenos valores de *Deviance* (ou elevado valor p) implicam que o modelo fornece um ajuste satisfatório aos dados , enquanto grandes valores de *deviance* implicam que o modelo atual não é adequado.
- Podemos dividir o *deviance* pelo graus de liberdade.
 - Se $\frac{D}{n-k} \gg 1 \rightarrow O \text{ Modelo não é adequado aos dados}$
 - Onde $n - k$ é o graus de liberdade. k é o número de parâmetro do modelo
 - D é dado por:

$$D = 2 \sum_{i=1}^n \left[y_i \ln \left(\frac{y_i}{n_i \hat{p}_i} \right) + (n_i - y_i) \ln \left(\frac{n_i - y_i}{n_i (1 - \hat{p}_i)} \right) \right]$$

χ^2 de Pearson

- Compara as probabilidade de sucesso e fracasso observadas e esperadas em cada grupo de observações
 - Nº esperado de sucesso : $n_i \hat{p}_i$
 - Nº esperado de fracasso: $n_i(1 - \hat{p}_i)$
- A estatística de de Pearson é dada por:

$$\chi^2_{n-k} = \sum_{i=1}^n \frac{(y_i - n_i \hat{p}_i)^2}{n_i \hat{p}_i (1 - \hat{p}_i)}$$

- **Valores pequenos da estatística de teste ou um grande valor de p , implica que o modelo fornece um ajuste satisfatório aos dados**

Exemplos