

Announcements

- Drop lowest HW
 - Exams out soon
 - Final 3:30-5:30p Here
 - Wed Off rescheduled to 3-4p
-

Logistic Regression pt II

To perform inference for Logistic Regression, we rely on asymptotic results (i.e. large sample approximations).

I Wald Test for Indiv. Coefficients

Main Idea: When $n \rightarrow \infty$, the MLE

$$\hat{\beta}_k \sim \underline{N}(\beta_k, \underline{V_k}) \text{ where } V_k \text{ is defined as follows:}$$

Let G denote the negative Hessian matrix

Fisher Information \rightarrow

$$G = - \begin{bmatrix} \frac{\partial^2 l}{\partial \beta_0^2} & \frac{\partial^2 l}{\partial \beta_0 \partial \beta_1} & \dots & \frac{\partial^2 l}{\partial \beta_0 \partial \beta_{p-1}} \\ \vdots & \frac{\partial^2 l}{\partial \beta_1^2} & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 l}{\partial \beta_0 \partial \beta_{p-1}} & \dots & \dots & \frac{\partial^2 l}{\partial \beta_{p-1}^2} \end{bmatrix}$$

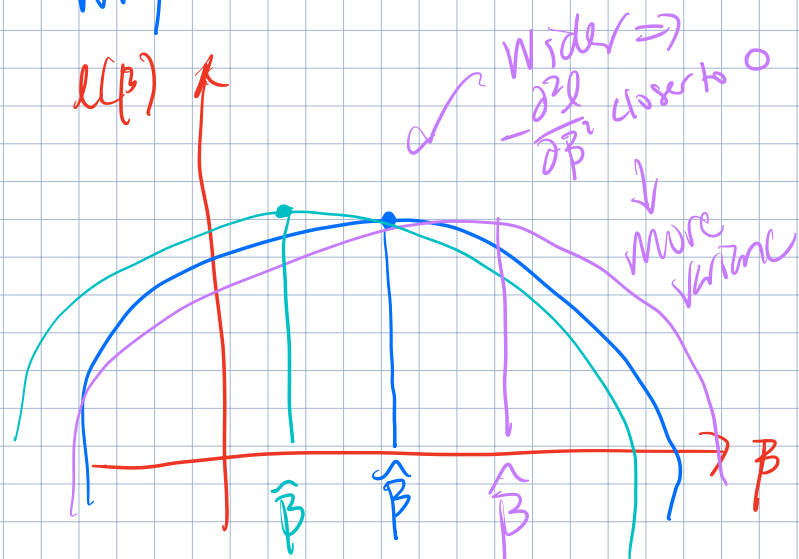
Then the variance of $\hat{\beta}_k$, is

$$\text{Var}(\hat{\beta}) = G^{-1} \Big|_{\beta = \hat{\beta}} \Rightarrow \text{Var}(\hat{\beta}_k) = \left[G^{-1} \Big|_{\beta = \hat{\beta}} \right]_{k+1, k+1}$$

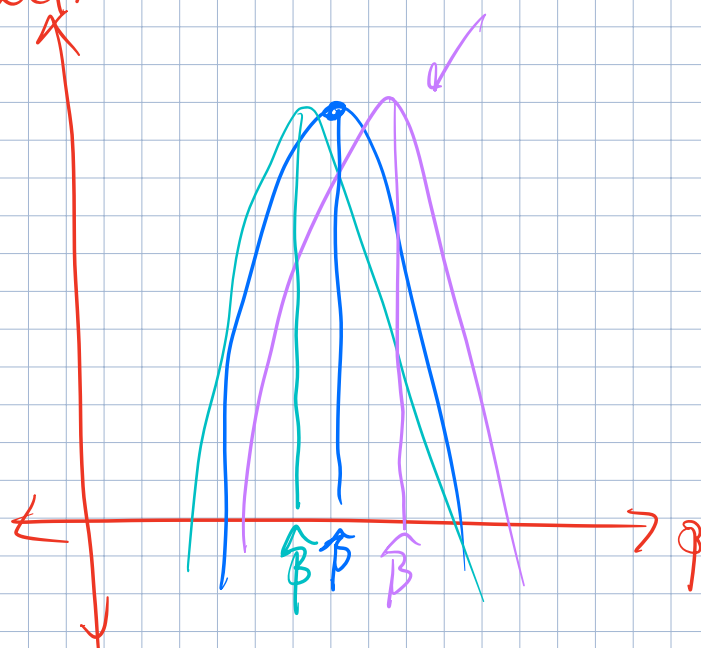
\uparrow
 $(k+1)^{\text{st}}$ diag b/c of β_0

Why does this work?

$l(\beta)$



$l(\beta)$



ok! Now we know that

$$\frac{\hat{\beta}_k - \beta_k}{\sqrt{V_k}} = \frac{\hat{\beta}_k - \bar{\beta}_k}{SE(\hat{\beta}_k)} \sim N(0, 1)$$

So an asymptotic test for β_k can be constructed as:

$$H_0: \beta_k = 0 \quad \text{vs} \quad H_1: \beta_k \neq 0$$

The test statistic:

$$Z = \frac{\hat{\beta}_k}{SE(\hat{\beta}_k)} \quad \text{leads to the decision rule:}$$

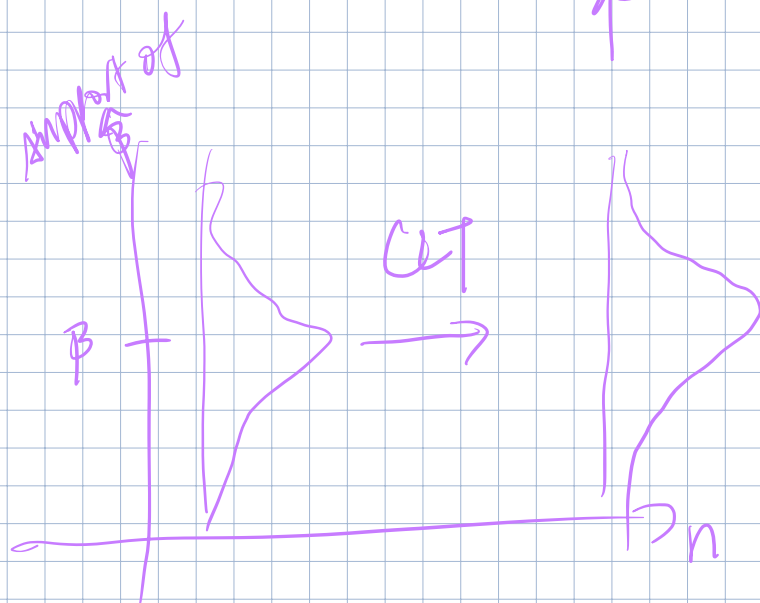
If $|Z| > Z_{1-\frac{\alpha}{2}}^*$ then reject H_0 .

⇒ We think the X_k covariate is predictive of the probability of "success".

$$\lceil \ell(\hat{\beta}) \approx \ell(\beta) + (\hat{\beta} - \beta)\ell'(\beta) + \dots \rceil$$

$$\underline{\underline{(\hat{\beta} - \beta)}} = \frac{1}{\underline{\underline{\ell'(\beta)}}} [\ell(\hat{\beta}) - \ell(\beta)] \leftarrow N(0, 1)$$

$$\ell(\beta) = \sum_{i=1}^n \underbrace{\log f(y_i | \beta)}_{\text{iid RV}} \rightarrow N(0, \underline{\underline{\quad}})$$



II Deviance & Likelihood Ratio Tests for Reduced & Full Models

Suppose we want to compare two models w/ a diff # of predictors, for ex:

$$H_0: \beta_r = \beta_{r+1} = \dots = \beta_{p-1} = 0$$

vs

H_1 : at least one $\{\beta_j\}_{j=r}^{p-1}$ is not zero



$$H_0: \text{Reduced Model} \Rightarrow \logit(\hat{\pi}_i) = \beta_0 + \sum_{j=r}^{p-1} \beta_j X_{ji}$$

vs

$$H_1: \text{Full Model} \Rightarrow \logit(\hat{\pi}_i) = \beta_0 + \sum_{j=1}^{p-1} \beta_j X_{ji}$$

We will need to construct a test based on loglikelihood \rightarrow compare loglikelihoods across the fitted models under H_0 & H_1

Deviance

Def: $\text{dev} = -2\ell(\beta)$

For ex:

In the OLS regression model, what is deviance?

$$\begin{aligned}\mathcal{L}(\beta) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ji})^2}{2\sigma^2}} \\ &= (2\pi\sigma^2)^{-n/2} e^{-\sum_i \frac{(y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ji})^2}{2\sigma^2}}\end{aligned}$$

\Rightarrow

$$\ell(\beta) = -\frac{n}{2} \log(2\pi\sigma^2) - \sum_i \frac{(y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ji})^2}{2\sigma^2}$$

$$\text{Dev} = -2\ell(\beta) = \underbrace{n \log(2\pi\sigma^2)} + \frac{\sum_i (y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ji})^2}{\sigma^2}$$

$$\text{Dev}(\hat{\beta}) = -2\ell(\hat{\beta}) = k_1 + k_2(\text{SSE})$$

So Deviance is like a generalization of the familiar SSE from OLS.

With this idea, I can consider a statistic based on deviance differences:

$$\Delta D = \text{Dev}(\text{reduced}) - \text{Dev}(\text{full})$$

Decision Rule:

$$\text{If } \Delta D > \chi^2_{df, 1-\alpha}, \text{ reject } H_0$$

$$\text{Here } df = p_{\text{full}} - p_{\text{red}} = p - v$$

If we reject H_0 , we have evidence that the full model is significantly better fit

than the reduced model according to likelihood
 \Rightarrow include the extra predictors!!

Recall, for logistic regression, the loglikelihood is:

$$l(\underline{\pi}) = \sum_{i=1}^n (y_i \log(\pi_i / (1 - \pi_i)) + \log(1 - \pi_i))$$

\therefore the deviance of a fitted model is just

$$\text{Dev}(\hat{\underline{\beta}}) = -2l(\hat{\underline{\beta}}) = -2 \left[\sum_{i=1}^n (y_i \log(\hat{\pi}_i / (1 - \hat{\pi}_i)) + \log(1 - \hat{\pi}_i)) \right]$$

$\begin{matrix} p-1 & \text{if full} \\ v-1 & \text{if reduced} \end{matrix}$

where

$$\log(\hat{\pi}_i / (1 - \hat{\pi}_i)) = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j x_{ij}$$

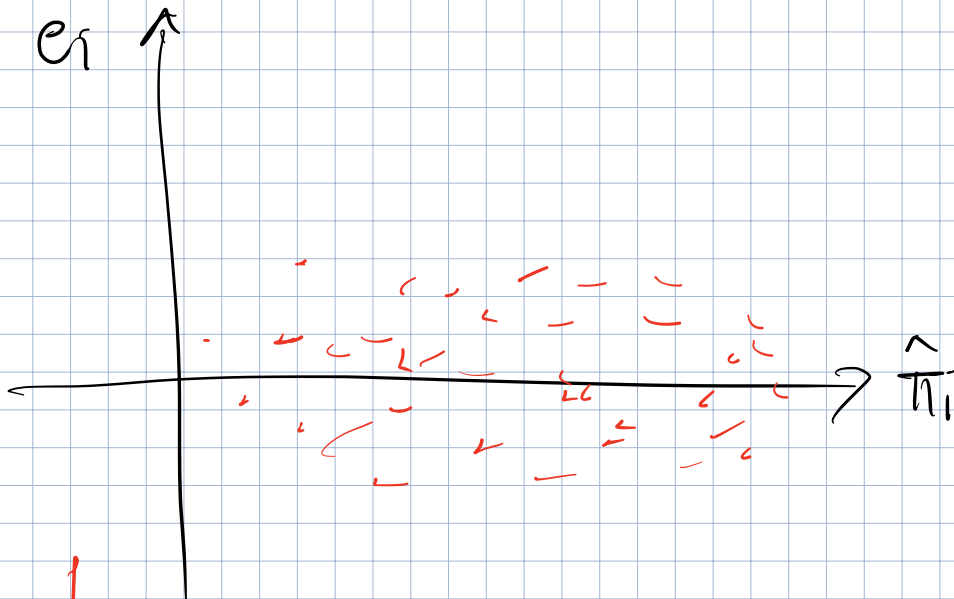
III Deviance Versions of Residuals

① Pearson Residuals:

- feel like standardized residuals in OLS:

$$e_i = \frac{y_i - \hat{\pi}_i}{\hat{SE}(y_i)} = \frac{y_i - \hat{\pi}_i}{\sqrt{\hat{\pi}_i (1 - \hat{\pi}_i)}}$$

We could make a "Pearson residual plot"
& would hope for no pattern in resids:



↳ interpret like OLS residuals

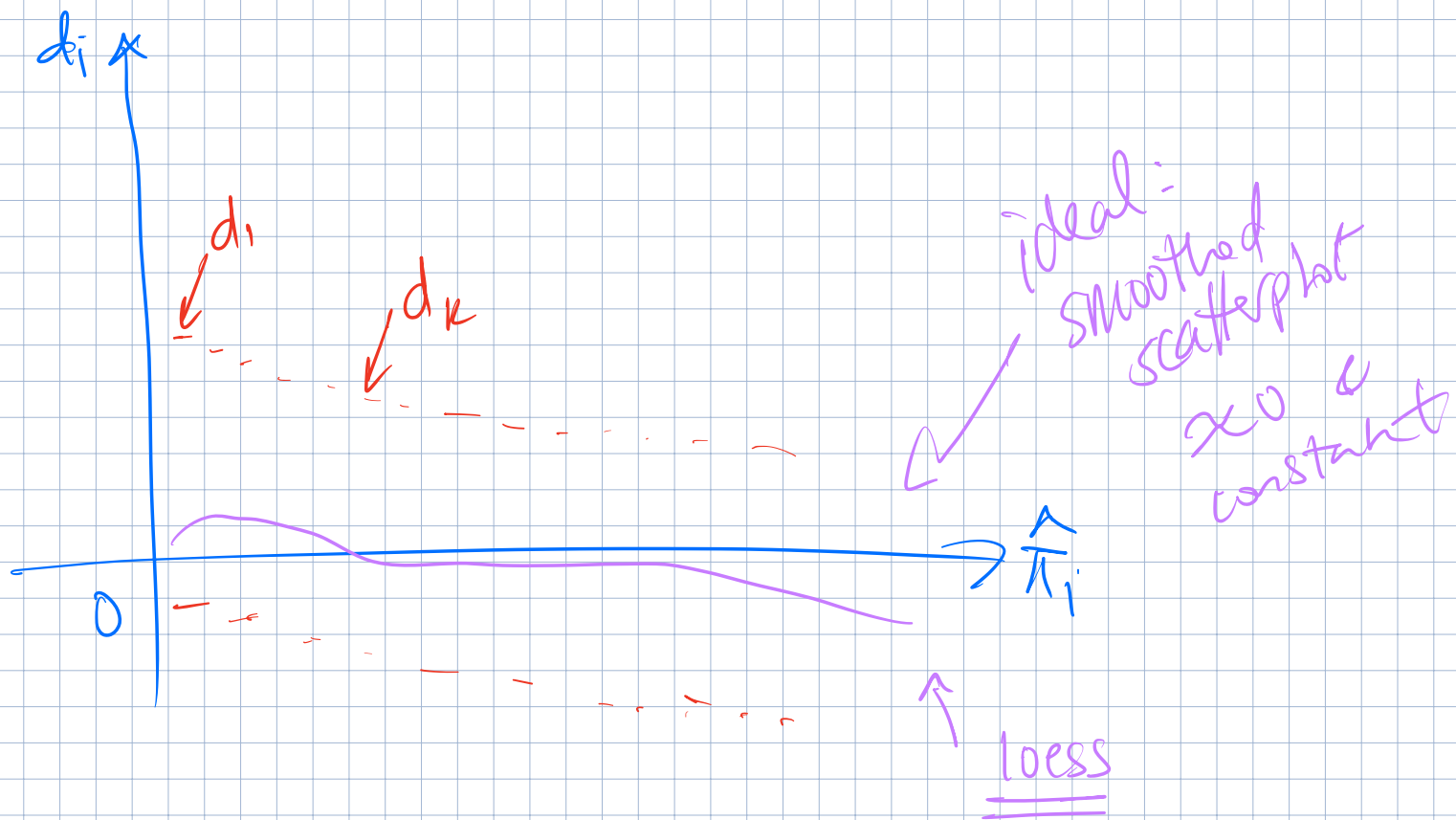
② Deviance Reside

Defin deviance residuals as:

$$d_i = \begin{cases} \sqrt{-2 (y_i \log(\hat{\pi}_i / (1 - \hat{\pi}_i)) + \log(1 - \hat{\pi}_i))}, & \text{if } y_i = 1 \\ -\sqrt{-2 (y_i \log(\hat{\pi}_i / (1 - \hat{\pi}_i)) + \log(1 - \hat{\pi}_i))}, & \text{if } y_i = 0 \end{cases}$$

The square of each deviance residual d_i^2 is basically how much deviance is attributable to the i th obs when considering the deviance of the full sample.

For diagnostics, we can look @ scatterplots:-



Pseudo R^2

Because there is no OLS principle,
the regular R^2 doesn't exist for
logistic regression to "explain variance".

Some alternative "pseudo R^2 " stats have
been proposed to assess goodness-of-fit:

Wron:
$$\text{pseudo } R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

McFadden's
$$\text{pseudo } R^2 = 1 - \frac{\ell(\text{Full})}{\ell(\text{Null})}$$

Model Selection for Log Reg:

☒ C_p or $Adj R^2$

☒ AIC/BIC ok \rightarrow b/c they're
likelihood based
criteria