# Model selection.

grouped → ungrouped data.

↓
contingency table.

### Marijuana use

| race | gender | Yes | No |
|------|--------|-----|-----|
| White | F | 420 | 620 |
| White | M | 483 | 579 |
| Others | F | 25 | 55 |
| Others | M | 32 | 62 |

$Y$ – Marijuana use          $(Y, X)$

$X$ – gender

$Z$ – race.

Test independence for $X$ and $Y$

$\mu_{ij}$ – expected frequencies          $\mu_{ij} = n \overline{\pi}_{ij}$

$\widehat{\mu_{ij}} = n \widehat{\pi}_{ij}$

$\overline{\pi}_{ij} = \overline{\pi}_{i+} \overline{\pi}_{+j}$

$G^2$ – test          Pearson. $\chi^2$ – test.

$$\text{logit}\,(P(Y=1)) = \alpha + \beta_1 x + \beta_2 z$$

↓          ↓

gender          race.

$X$, $Y$ are independent

$H_0$: $\beta_1 = 0$      Wald test

$H_1$: $\beta_1 \neq 0$      LRT- test

Goodness of fit test.

↙

comparison with a saturated model

↓

↘ Comparison with a null model.

↓

LRT → $G^2$ test

$G^2 = 2 \sum$ observed $\log \dfrac{observed}{fitted}$

score test. Pearson $\chi^2$ test

$$\chi^2 = \dfrac{\sum (observed - fitted)^2}{fitted} \sim \chi^2_{df}.$$

$\text{logit}\,(P(Y=1)) = \alpha + \beta_1 x + \beta_2 z$

$\text{logit}\,(P(Y=1)) = \alpha$

$H_0$: $\beta_1 = \beta_2 = 0$

$H_1$: at least one is not zero.

## Residuals.

Pearson residuals $e_i = \dfrac{\overset{y_i}{observed} - predicted}{\sqrt{Var(\hat{y_i})}}$

score    $\chi^2 = \sum\limits_{i=1}^{n} e_i^2$

Stand. res. $= \dfrac{y_i - n\,\hat{\pi}_i}{SE}$      $SE = \sqrt{Var(\hat{y_i})(1-h_i)}$

St. res. $\sim N(0,1)$ if the fit is good.

$$|st\ res| < 2.$$
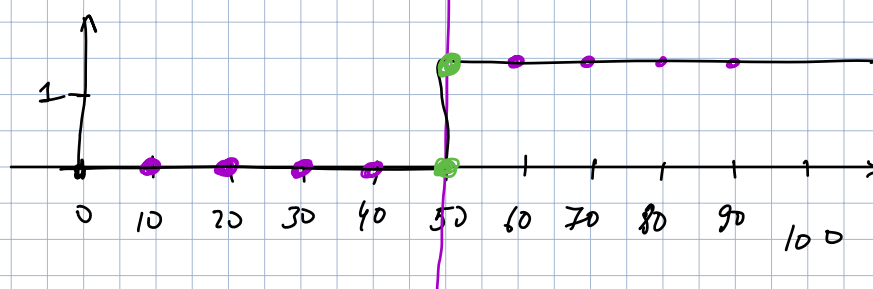$$< 3$$

Sparse data.

$x = 10, 20, 30, 40$      $y = 0$

$x = 60, 70, 80, 90$      $y = 1$



$\beta = \infty$      perfect separation.

quasi complete

$x = 10, 20, 30, 40, 50$      $y = 0$

$x = 50, 60, 70, 80, 90$      $y = 1$

Summary on model selection.

Model selection

→ backward step algorithm

↳ search all models and select the best one using AIC (BIC) criteria

↳ penalized regression (Lasso, Ridge)

# Goodness of fit.

| Method of goodness of fit or comparison of models. | Grouped data | Ungrouped data |
|---|---|---|
| 1. Comparison with a **saturated** model. | Yes $G^2$-test & RT Pearson $\chi^2$-test. | NO chi-sq. approx does not work |
| 2. Comparison of two nested model. | Yes LRT | Yes. LRT |
| 3. Comparison with the null model | Yes LRT | Yes LRT |
| 4. Residual analysis. | Yes | Yes. |
| 5 Correlation | No | Yes . |
| 6. ROC AUC | NO | Yes . |

# Alternative link functions.

$y \nearrow 1$
$\searrow 0$

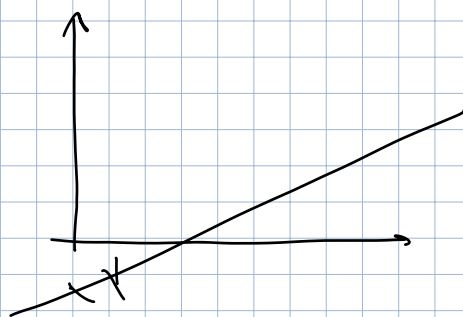$$logit\left( P(Y=1) \right) = \alpha + \beta_1 x_1 + \ldots + \beta_P x_P .$$

$$log\left( \frac{\pi}{1-\pi} \right) \leftarrow \text{logit lint} .$$

## Identity link function,

- $\pi = \alpha + \beta_1 x_1 + \ldots + \beta_P x_P .$

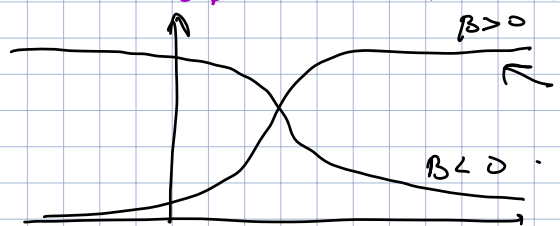$P(Y=1) = \alpha + \beta_1 x_1 + \ldots + \beta_P x_P .$

Problem ?

$\pi(x) = \alpha + \beta x$

$\begin{cases} \widehat{\pi} \text{ can be } < 0 . \\ \widehat{\pi} > 1 . \end{cases}$

lack of convergance
and giva an error
message .

## Probit model and norma latent variable model.

$\beta > 0$

$logit(\pi(x)) = \alpha + \beta x$

cdf of $\beta > 0 .$  S-shaped.

$\beta < 0 .$

some distribution ( logistic distr.)

$Y \sim N(0,1)$    $\Phi(y)$ — cdf of a st. normal distr.
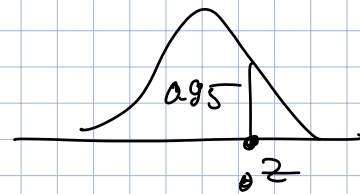
$\pi(x) = \Phi(\alpha + \beta x)$

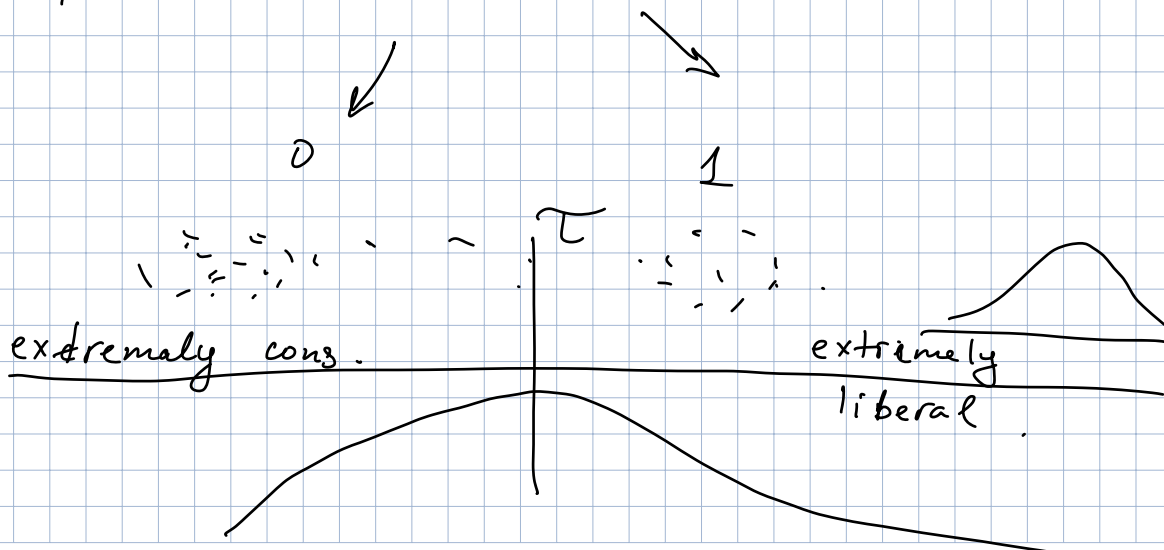$\Phi^{-1}(\pi(x)) = \alpha + \beta x$

probit$(\pi(x)) = \alpha + \beta x$

Probit link transforms
$P(Y=1)$         to    a   z-score.
probit$(0.5) = 0$
probit$(0.95) = 1.645$

Interpretation of the parameters in
probit model is simplest when we can
relate the model to a corresponding normal
linear model.

0                    1

$$y^* \qquad \text{if} \qquad y^* < \tau \qquad y = 0$$
$$y^* > \tau \qquad y = 1$$

$$y^* = \alpha + \beta_1 x_1 + \ldots + \beta_p x_p + \boxed{\varepsilon}$$
$$\varepsilon \sim N(0, \sigma^2)$$

$$\text{probit}\left( P(Y=1) \right) = \alpha + \beta x_1 + \ldots + \beta_p x_p .$$

We can interpret $\widehat{\beta_j}$ from probit model
fit as representing the estimated
changes in $\pounds y^*$ for 1 unit increase
in $x_j$ adjusting for other explanatory
variables.
For arbitrary value for $\text{var}(\varepsilon)$ $\widehat{\beta_j}$
is the estimated # of st. dev.
that the distr. $y^*$ shifts.

$\quad y^*$ - latent (unobserved variable)