

Linear regression model: classical approach

The observed data: (i). Any single observation is given by (Y, \mathbf{z})

Y is the observed variable

$\mathbf{z} = (z_0, \dots, z_k)$ is the covariate

$z_0 \in \{0, 1, 2\}$ is assumed to be a genotype

The distribution of Y is assumed to be a normal; any element z_i of the covariate \mathbf{z} can be

nominal (categorical, factor)

ordinal

quantitative (numeric)

Covariates:

the covariate $z_0 \in \{0, 1, 2\}$ (genotype) is assumed to be of special ordinal/nominal type (object "genotype")

type of each other covariate should be specified by user or selected by default

By default

a binary categorical covariate (2 levels) is assumed to be a nominal, and should be converted to 0 and 1 levels;

a categorical covariate having 3-8 levels is assumed to be an ordinal;

other covariates are assumed to be numeric.

(ii). The observed data is a sample (\mathbf{Y}, \mathbf{z}) :

$\mathbf{Y} = (Y_1, \dots, Y_n)'$ is $n \times 1$ -vector

$\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)$, $\mathbf{z}_i = (z_{i,0}, \dots, z_{i,k})$ are the covariates

Y_1, \dots, Y_n are conditionally independent given \mathbf{z}

$\text{Var}(Y|\mathbf{z}) = \sigma^2$ for any \mathbf{z} ($\text{Var}(Y_1) = \dots = \text{Var}(Y_n)$)

Linear regression model: (i). Regressor $\mathbf{X} : \mathbb{R}^k \rightarrow \mathbb{R}^m$; $\mathbf{X}(\mathbf{z})$ is defined for any observation

(ii). Linear regression model for single observation

$$\mathbb{E}_\theta(Y|\mathbf{z}) = \mathbf{X}(\mathbf{z})'\boldsymbol{\beta}$$

$\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)'$ is $m \times 1$ -vector of parameters

$Y \sim \mathcal{N}(\mathbf{X}(\mathbf{z})'\boldsymbol{\beta}, \sigma)$, where σ is the nuisance variance parameter

$\boldsymbol{\theta} = (\beta_1, \dots, \beta_m, \sigma)$ is the full parameter of the model

The linear regression model is determined by formula

$$Y \sim g(z_0) + f(\mathbf{z}^*), \quad \mathbf{z}^* = (z_1, \dots, z_k),$$

where $g \in \{CD, D, R, A\}$ is the genetic model and f is the covariate term

$$f(\mathbf{z}^*) = f_1(\mathbf{z}^*) + \dots + f_r(\mathbf{z}^*),$$

$$f_i(\mathbf{z}^*) = z_1^{r_{1i}} * \dots * z_k^{r_{ki}}, \quad r_{1i}, \dots, r_{ki} \in \{0, 1, \dots\} \text{ is of multiple regression type}$$

warning: values $r_{ki} \geq 2$ are not available for factors

The parameter

The intercept parameter is included into the model by default

The covariate z_0 generates

single parameter if $g \in \{D, R, A\}$ (difference with the baseline level $z_0 = 0$)

vector of length two if $g = CD$ (main effect with respect to the baseline level $z_0 = 0$)

Term $f(\mathbf{z}^*)$ should be transformed to $\sum_{\kappa \in K} z_\kappa$, where κ is a multi-index; K is a (ordered) set of subsets of $\{1, \dots, k\}$

the term $f_i(\mathbf{z}^*) = z_1^{r_{1i}} * \dots * z_k^{r_{ki}}$ generates K_i is the set of non empty subsets of $\{r_{k,i} : r_{k,i} = 1\}$
then $K = \cup_{i=1}^r K_i$
the elements of K should be sorted (κ 's of smaller length should be sorted before κ 's of larger length)

We associate a number d_κ to each index $\kappa \in K$

$d_i = 1$ if z_i is numeric or ordinal ($\kappa = i$)
 $d_i = r$ for z_i^r
 $d_i + 1$ is a number of levels of the covariate z_i , if z_i is a factor ($\kappa = i$)
in general case $d_\kappa = \prod_{i \in \kappa} d_i$.

The total length of the nuisance (part corresponding to f) parameter is $d = \sum_{\kappa \in K} d_\kappa$.

The regressor

The regressor is a vector (x_0, x_K) such that

$x_0 = (1, z_0)$ if $g \in \{D, R, A\}$
 $x_0 = (1, \mathbb{I}_{\{z_0=1\}}, \mathbb{I}_{\{z_0=2\}})$ if $g = "CD"$ (Model for TEST)
 $x_0 = (1, z_0)$ if $g = "CD"$ (Model for QAS)
 $x_0 = 1$ (Model for base)

We use the ordered set K to create regressors for nuisance parameter

$x_i = z_i$ if z_i is numeric or ordinal ($\kappa = i$)
 $x_i = (z_i, z_i^2, \dots, z_i^r)$ for z_i^r
 x_i is a vector of length d_i of $\{0, 1\}$ with 1 on j -th position if $z_i = a_{j+1}$, a_j is the j -th level of z_i (for example $(0, 0, 1, 0, 0)$).
in general case $x_\kappa = \text{as.vector}(\otimes_{i \in \kappa} x_i)$, where \otimes is the outer product (vector should be sorted in a proper order).
the regressor x_K , which corresponds to the nuisance parameter is obtained by merging x_κ over $\kappa \in K$ in the proper order.

(iii). Linear regression model for the observed data

$$\mathbb{E}(\mathbf{Y}|\mathbf{z}) = \mathbf{X}'\boldsymbol{\beta}$$

$\mathbf{X} = (\mathbf{X}(\mathbf{z}_1), \dots, \mathbf{X}(\mathbf{z}_n))$ is the $m \times n$ -matrix of regressors

$\boldsymbol{\beta}$ is the $1 \times m$ -vector of parameters

assumption $\mathbf{Y} \sim \mathcal{N}(\mathbf{X}'\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$, \mathbf{I}_n is the identity matrix

the matrix $\mathbf{X} = (\mathbf{X}'_g, \mathbf{X}^{*'})'$, where \mathbf{X}_g contains the regressors related to g and \mathbf{X}^* are the last d rows of the matrix \mathbf{X}

Nuisance parameter reduction

the parameter related to $(\mathbf{z}_1, \dots, \mathbf{z}_k)$ is not a subject of our interests

If the matrix \mathbf{X}^* , which contains the last d rows of the matrix \mathbf{X} , is not of the full rank d , we reduce the matrix \mathbf{X}^* to \mathbf{X}° by choosing only basis vectors of the linear space generated by rows of the matrix \mathbf{X}^*

The reduced matrix \mathbf{X} is determined by changing block \mathbf{X}^* to \mathbf{X}° in the matrix \mathbf{X}

QAS and P-value. Under $\mathbf{X}\mathbf{X}'$ is of full rank (positively definite), the least square estimator (LSE) is given by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}\mathbf{X}')^{-1}\mathbf{X}'\mathbf{Y}$$

The sum of squares

$$S(\boldsymbol{\beta}) = \|\mathbf{Y} - \mathbf{X}'\boldsymbol{\beta}\|^2 = (\mathbf{Y} - \mathbf{X}'\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}'\boldsymbol{\beta})$$

Notatons ("Model for QAS" is equal to "Model for TEST" if $g = A, D, R$):

$\mathbf{X}/\mathbf{X}_q/\mathbf{X}_0$ is the regressors in the "Model for TEST"/"Model for QAS"/"Baseline model" respectively
 $\hat{\beta}/\hat{\beta}_q/\hat{\beta}_0$ is the LSE from the "Model for TEST"/"Model for QAS"/"Baseline model" respectively
 $S(\hat{\beta})/S_q(\hat{\beta}_q)/S_0(\hat{\beta}_0)$ is the sums of squares from the "Model for TEST"/"Model for QAS"/"Baseline model" respectively

The algorithm

The input data contain: "phenotype" (numeric vector), "genotype", "covariates" (matrix), "covariate types" (vector) and "formula"

- (i) Create regressors related to covariates \mathbf{X}^*
- (ii) Check that the rank of the matrix \mathbf{X}^* is full
 if the rank of matrix \mathbf{X}^* is not equal to d , apply nuisance parameter reduction $\mathbf{X}^* := \mathbf{X}^\circ$.
- (iii) Loop by genotypes
 Create regressors \mathbf{X}_g for "Model for TEST", for "Model for QAS" and for "Baseline model"
 Check that the rank of the matrix \mathbf{X}_g related to g is maximal (2 if $g = A, D, R$ or 3 if $g = CD$)
 if the rank of the matrix \mathbf{X}_g in the "Model for TEST" is not maximal, report $gas[i] = NA$ and $pv[i] = NA$
 otherwise, continue
 Create matrix \mathbf{X} by merging \mathbf{X}_g and \mathbf{X}^*
 If the rank of matrix \mathbf{X} is not equal to its number of rows use nuisance parameter reduction to whole \mathbf{X} (\mathbf{X}_g should not be changed)
 Fit the "Model for QAS"
 assign $gas[i]$ the coefficient, related to \mathbf{z}_0 (second element of the vector $\hat{\beta}$)
 Fit the "Model for TESTS" and the "Baseline model"
 $SS_e = S(\hat{\beta})$; $\overline{SS}_e = SS_e/\text{df}$, where $\text{df} = n - \text{rk}(\mathbf{X})$
 $SS_h = S_0(\hat{\beta}_0) - SS_e$; $\overline{SS}_h = SS_h/\text{df}_h$, where $\text{df}_h = 2$ if $g = "CD"$ and $\text{df}_h = 1$ if $g = "A, D, R"$
 $\mathbb{F} = \overline{SS}_h/\overline{SS}_e$ is the F -statistic, which has under null hypothesis $F_{\text{df}_h, \text{df}}$ -distribution
 assign $pv[i] := 1 - F_{\text{df}_h, \text{df}}(\mathbb{F})$, where $F_{\text{df}_h, \text{df}}$ is the distribution function of the Fisher-Snedecor distribution $F_{\text{df}_h, \text{df}}$
- (iii) Return the vectors pv and gas .