

STATISTICAL MODELS

Focus on REGRESSION MODELS: study the relationship between variables

In particular: 2 types of variables:

- RESPONSE / DEPENDENT variable Y
- one or more PREDICTORS / COVARIATES / INDEPENDENT variables X_1, X_2, \dots, X_p

GOAL: study how the response is influenced by the covariates (hence the relationship is not symmetric: the variables have different "roles")

- examples:
 - evaluate how the blood pressure is affected by a specific treatment, while also controlling for the individuals' characteristics (age, weight, ...)
 - predict the number of claims given the insurer's characteristics (age, past accidents, ...)

$\Rightarrow Y = g(X_1, \dots, X_p)$
and our goal is to study $g(\cdot)$

- As usual, the variables are observed on several individuals / statistical units (n is the number of observations)
- We consider only 1 response variable
- The number of predictors is $p \geq 1$.

The data can be organized into a matrix

DATA: matrix: - rows: individuals / statistical units
- columns: variables

statistical unit	response variable	1 st predictor	2 nd predictor	...	p-th predictor
1	y_1	x_{11}	x_{12}	...	x_{1p}
2	y_2	x_{21}	x_{22}	...	x_{2p}
...
i	y_i	x_{i1}	x_{i2}	...	x_{ip}
...
n	y_n	x_{n1}	x_{n2}	...	x_{np}

the submatrix of the covariates alone is called the "model matrix"

$$X = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}$$

When do we need statistics? In the applications we are interested in, the value of the response variable is not fixed, given the values of the predictors \rightarrow there is uncertainty: the relationship between Y and (X_1, \dots, X_p) is stochastic (not deterministic)

STATISTICAL MODELS: assume that the observations are REALIZATIONS of RANDOM VARIABLES \rightarrow goal is the study of whether and how the LAW of the response variable is affected by the independent variables: $Y \sim f(y; x_1, \dots, x_p)$.

Common assumption: the covariates are non-stochastic and measured without error. This is justified in experimental settings (e.g. fix the dose of the treatment and study the outcome). In observational studies this is not possible (e.g. demographic / economic / social studies). For simplicity, the hypothesis is maintained, with the interpretation that the analysis is performed CONDITIONALLY on the observed values of the covariates (i.e. $Y | X_1=x_1, \dots, X_p=x_p \sim f(y; x_1, \dots, x_p)$).
 (note: uppercase for r.v. lowercase for the observed value)

How do we actually build a model and perform the analysis?

THE FUNDAMENTAL STEPS:

- 1) MODEL SPECIFICATION
given the goal of the study and the available data, specify the model (also using past info, theories on the problem, ...)
- 2) ESTIMATION
estimate the model parameters (unknown quantities that define $g(\cdot)$) on the basis of the observed data
- 3) MODEL CHECKING / DIAGNOSIS
are the hypotheses underlying the model coherent with the observed data? YES: use the model
NO: go back to 1) and repeat

1. MODEL SPECIFICATION

1A) THE RANDOM COMPONENT

The type of model that we specify mainly depends on the nature of the response variable (remember that we are modeling the law of Y)

RESPONSE VARIABLE:

- QUANTITATIVE
 - \rightarrow continuous (support \mathbb{R}) \rightarrow Gaussian linear model; linear model via OLS (no Gaussian assumption)
 - \rightarrow discrete / counts (support \mathbb{N}_0) \rightarrow Poisson regression (GLM)
- QUALITATIVE (categorical)
 - nominal variables (no order in the levels)
 - \rightarrow binary (only 2 levels, e.g. presence/absence) \rightarrow Logistic regression (logit model), probit model (GLM)
 - \rightarrow more than 2 categories, not ordered (e.g. hair color) \rightarrow Logistic regression / multinomial model (GLM)
 - ordinal variables (the categories have an intrinsic ordering, e.g. rankings low/medium/high rates very unsatisfied/unsatisfied/satisfied/very satisfied) \rightarrow Cumulative logit/probit model (GLM)

The type of response variable drives the choice of the distribution $f(y; x_1, \dots, x_p)$.

1B) THE RELATIONSHIP between Y and X_1, \dots, X_p : $g(\cdot)$

it is deterministic: it is also called the SYSTEMATIC COMPONENT
 we will consider the case where $g(\cdot)$ is completely specified by a FINITE SET of (unknown) REAL PARAMETERS $\theta \in \Theta \subseteq \mathbb{R}^q$, $q \geq 1$ finite.
 The specific way each covariate enters the model depends on the type of variable (more details later...)

2. ESTIMATE

The estimate procedure consists in estimating the unknown parameters on the basis of the observed data. Once we estimate $\hat{\theta}$, the relationship between Y and x_1, \dots, x_p is completely known.

3. MODEL CHECKING

- Having uniquely defined the model, we need to check
- goodness of fit: does the model fit the observed data well?
 - do we need all the considered covariates or a more parsimonious model can be defined (without loss of fit)?
 - are the distributional assumptions satisfied?

If the model checking highlights some kind of problem, we have to go back to the model specification (and change, for example, the way the variables enter the model, the number of covariates, the assumptions on the law f) and repeat the procedure until step 3 gives good results.

Then, the model can be used for:

- inference on the parameters: understand the effect of each covariate
- prediction: given specific values of the covariates, what is the value of Y ? (careful with prediction at values of the X_j outside of the observed range, i.e. extrapolation)

So far, we have denoted the relationship between Y and (x_1, \dots, x_p) simply as $Y \sim f(y; x_1, \dots, x_p)$ meaning that the distribution of Y depends on the covariates.

ADDITIONAL ERROR TERM

The simplest way to introduce the stochastic component is to consider

$$Y = \underbrace{g(x_1, \dots, x_p)}_{\substack{\text{regression} \\ \text{function} \\ \text{deterministic}}} + \underbrace{\varepsilon}_{\substack{\text{error term} \\ \text{stochastic}}}$$

(notice: GLMs do not fall into this kind of specification)

Regression models can be classified based on:

1. the number of variables involved
2. the type of function linking Y to the x_j , $j=1, \dots, p$

1) NUMBER OF VARIABLES

1A. number of INDEPENDENT variables:

- "SIMPLE" regression: only 1 covariate $Y = g(x_1) + \varepsilon$
- "MULTIPLE" regression: $p > 1$ covariates $Y = g(x_1, \dots, x_p) + \varepsilon$

1B. number of DEPENDENT variables:

- univariate: only 1 response Y
- multivariate: the response is a vector $\underline{Y} = (Y_1, \dots, Y_m)$

2) TYPE OF FUNCTION $g(\cdot)$

2A. PARAMETRIC: g can be expressed using a FINITE number of parameters $\theta = (\theta_1, \dots, \theta_q) \in \Theta \subseteq \mathbb{R}^q$, q finite

- LINEAR: $g(\cdot)$ is a parametric function and it is LINEAR in the parameters. We denote the parameters with β .

Examples: $g(x) = \beta_1 x$
 $g(x) = \beta_1 x + \beta_2 x^2 + \beta_3 x^3$
 $g(x) = \beta_1 \log x - \beta_2 \sqrt{x}$
 $g(x_1, x_2, x_3) = \beta_1 x_1 + \beta_2 \log x_2 + \beta_3 e^{x_3 + x_2}$
 the parameters $\beta = (\beta_1, \beta_2, \beta_3)$ enter linearly.

Notice that the variables x_j need not be linear! We can transform them to better fit the data.

- LINEARIZABLE: the relation is not linear, but there is a transformation to make it so:

Example: the model $Y = \beta_1 \cdot x^2 + \varepsilon$ is not linear.
 But if we take the logarithm: $\log Y = \underbrace{\log \beta_1}_{\tilde{\beta}_1} + \beta_2 \underbrace{\log x^2}_{\tilde{x}} + \log \varepsilon$
 $\Rightarrow \tilde{Y} = \tilde{\beta}_1 + \beta_2 \cdot \tilde{x} + \tilde{\varepsilon}$ linear

- NON-LINEAR: it is parametric but it is not linear nor linearizable

Example: $Y = \frac{\beta_1 x}{\beta_2 + x} + \varepsilon$

2B. NONPARAMETRIC: the parameter space Θ is not a subset of \mathbb{R}^q (e.g. kernel regression, trees, RF, splines, nearest neighbors, ...)
 GP regression