# Linear Models and Regression II
## Introductory Lecture

ginsbourger@stat.unibe.ch

Spring Semester 2017

---

Introduction: beyond classical linear modelling
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## Foreword

Welcome to the first lecture of the course "Linear Models and Regression II"!

The course will take place from today until the end of the semester (ending on Friday June 2nd 2017). Note that the Easter holidays will take place from Friday 14th to Sunday 23rd April 2017.

Usual times for the lecture and exercises will be respectively Mondays from 8.15AM to 10AM and Fridays from 2.15PM to 4PM in room B1 (here). An exception is March 13-17 for which lecture/exercise days will be swapped.

The (oral) exam will take place a few weeks after the end of the semester, around end of June/beginning of July 2017. Regarding the exercises, given by Dr. Dario Azzimonti, things will work similarly as for LMRI.

---

Introduction: beyond classical linear modelling
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## Beyond classical linear modelling

In the last semester we have focused on statistical models of the form

$$Y = \sum_{j=1}^{p} \theta_j f_j(X) + \epsilon,$$

where $Y \in \mathcal{Y}, X \in \mathcal{X}$ are respectively a response and a covariable (both possibly multivariate), $p \in \mathbb{N} \setminus \{0\}$, $f_j : \mathcal{X} \to \mathcal{Y}$ are basis functions with coefficients $\theta_j \in \mathbb{R}$ ($1 \leq j \leq p$), and $\epsilon$ is a random error in $\mathcal{Y}$. We have typically assumed $\mathcal{Y} = \mathbb{R}, \mathcal{X} = \mathbb{R}^d$ and $\epsilon$ centred, Gaussian or at least $L^2$.

We have seen that this formulation encompasses a number of particular models (Simple and Multiple Linear Regression, Polynomial Regression ANOVA, ANCOVA, etc.) and also that under assumptions it is possible based on observed data $(X_i, Y_i)_{i=1,\ldots,n}$ ($n \in \mathbb{N} \setminus \{0\}$) to estimate $\theta$ and also to perform tests, derive exact and approximate confidence regions, etc.

---

Introduction: beyond classical linear modelling
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## Beyond classical linear modelling

In a number of situations however, classical models and/or methods from LMRI are **not suitable**. This concerns notably the following cases:

- It is not reasonable to assume $\epsilon$ Gaussian or related as
    - $Y$ is binary, categorical, integer-valued, or other scenarios for which an $\epsilon$ with distribution supported by the real-line does not make sense,
    - or $Y$ and $\epsilon$ may well be real-valued but another distribution is assumed (then it might be that linear models are suitable but not least-squares)

- The relationship between $Y$ and $\theta$ is non-linear,

- Beyond usual assumptions but possibly still in the linear set up, it might be that $\epsilon_i$'s corresponding to different observations are dependent, that the design matrix does not have full column-rank, etc.

**Introduction: beyond classical linear modelling**
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## A detour through non-linear least squares

In cases where responses $Y \in \mathbb{R}^n$ stem from a non-linear model of the form

$$Y = f(X, \boldsymbol{\theta}) + \epsilon$$

or in other words when $\boldsymbol{Y} = \boldsymbol{v}(\boldsymbol{\theta}) + \epsilon$ where $\boldsymbol{v}(\boldsymbol{\theta})$ has components $f(X_i, \boldsymbol{\theta})$ and $\epsilon$ is a centred squared-integrable noise, one common approach is to keep using the least-squares procedure.

A *Least Squares Estimator* $\widehat{\boldsymbol{\theta}}$ of $\boldsymbol{\theta}$ is then defined by its realizations —assuming their existence— minimizing given $\boldsymbol{y}$ the loss function

$$Q : \boldsymbol{\eta} \in \mathbb{R}^p \mapsto Q(\boldsymbol{\eta}) = ||\mathbf{y} - \boldsymbol{v}(\boldsymbol{\eta})|||^2$$
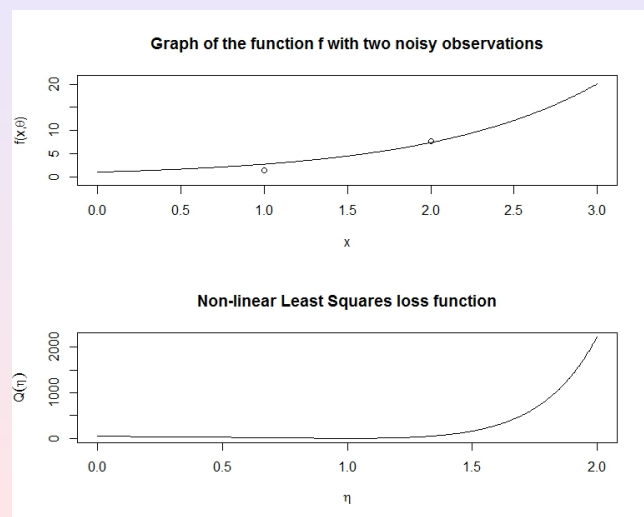
Conditions to ensure existence of the LS estimator —not discussed in detail here— classically involve $\Theta$ (e.g. assumed to be a subset of $\mathbb{R}^p$ with non-empty interior; compacity simplifies things but is not always assumed), $\boldsymbol{v}$ (typically assumed injective and twice continuously differentiable), and the $\epsilon_i$'s (typically assumed centred, squared-integrable and uncorrelated).

---

**Introduction: beyond classical linear modelling**
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## A detour through non-linear least squares

A simple and classical example follows (from "NonLinear Regression Modeling" (1983) by D.A. Ratkowsky, cited in "Regression non linéaire et applications" (1992) by A. Antoniadis, J. Berruyer and R. Carmona, a book that has largely inspired the first lectures of the current course).

Here $\mathcal{X} = \mathcal{Y} = \mathbb{R}$ and $Y_i = f(X_i, \theta) + \epsilon_i$ with $f(x; \theta) = \exp(\theta x)$, $\epsilon_i$ realizations of a centred Gaussian noise with given variance $\sigma^2$, and $X_1, X_2$ are also given. Here we take $\sigma = 2$, $X_1 = 1$ and $X_2 = 2$.

---

**Introduction: beyond classical linear modelling**
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## A detour through non-linear least squares

---

**Introduction: beyond classical linear modelling**
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## A detour through non-linear least squares

In this example the loss function is quite simple and it is even possible to calculate the (non-linear) LS estimator in closed form. In fact, we have

$$Q(\eta) = (y_1 - \exp(\eta))^2 + (y_2 - \exp(2\eta))^2$$

Noting that $Q$ is here convex differentiable, we get by annulling its derivative

$$2\exp(3\widehat{\theta}) - (2y_2 - 1)\exp(\widehat{\theta}) - y_1 = 0.$$

Posing $\gamma = \exp(\widehat{\theta})$ we can solve for $\gamma$ as real root of the third order polynomial $2\gamma^3 - (2y_2 - 1)\gamma - y_1$. With the (simulated) values $y_1 = 4.556237$ and $y_2 = 8.953329$, we get $\gamma = 2.7428935$ whereof $\widehat{\theta} = 1.009013$.

For more complicated cases (with larger $n$ and/or $p$, more complex functions, etc.), solving non-linear LS analytically is out of reach. Most of the time, non-linear LS calls for **iterative fitting procedures**. Similar procedures will also be useful in more general frameworks as discussed in the next sections.

**Introduction: beyond classical linear modelling**
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## A detour through non-linear least squares

Assuming $\Theta$ with a non-empty interior and uncorrelated centred noise with constant variance, most iterative procedures for non-linear least squares are based on the following result:

**Proposition**

Any LS estimator $\widehat{\theta}$ belonging to $\overset{\circ}{\Theta}$ is solution to the *normal equations*:

$$(\boldsymbol{Y} - \boldsymbol{v}(\eta))^T \nabla \boldsymbol{v}(\eta) = 0.$$

The *Gauss-Newton* algorithm addresses the issue of approaching solutions to this by relying on Taylor expansions of $\boldsymbol{v}$. Assuming that the current candidate parameter vector is $\theta^{(i)}$, we have at order one

$$\boldsymbol{v}(\eta) \approx \boldsymbol{v}(\theta^{(i)}) + \nabla \boldsymbol{v}(\theta^{(i)})^T (\eta - \theta^{(i)})$$

Replacing $\boldsymbol{v}$ by its local approximation in the normal equation, one gets

$$(\boldsymbol{y} - \boldsymbol{v}(\theta^{(i)}) - \nabla \boldsymbol{v}(\theta^{(i)})^T (\eta - \theta^{(i)}))^T \nabla \boldsymbol{v}(\theta^{(i)}) = 0$$

---

**Introduction: beyond classical linear modelling**
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## A detour through non-linear least squares

Expanding the former equation and assuming that the symmetric matrix $\nabla \boldsymbol{v}(\theta^{(i)})^T \nabla \boldsymbol{v}(\theta^{(i)})$ is positive definite (i.e., invertible, here), the principle of the Gauss-Newton algorithm is to define the next point as

$$\theta^{(i+1)} = \theta^{(i)} + \left( \nabla \boldsymbol{v}(\theta^{(i)})^T \nabla \boldsymbol{v}(\theta^{(i)}) \right)^{-1} \nabla \boldsymbol{v}(\theta^{(i)}) \left( \boldsymbol{y} - \boldsymbol{v}(\theta^{(i)}) \right).$$

Note that another popular approach consists in directly taking a Taylor expansion of $Q$'s gradient so as to get

$$\nabla Q(\eta) \approx \nabla Q(\theta^{(i)}) + \nabla^2 Q(\theta^{(i)})(\eta - \theta^{(i)})$$

A solution to the normal equation is then sought by iteratively setting

$$\theta^{(i+1)} = \theta^{(i)} - \nabla^2 Q(\theta^{(i)})^{-1} \nabla Q(\theta^{(i)}).$$

This method is known as the *Newton-Raphson* algorithm. Note that in both cases the choice of a good starting point may be crucial.

---

**Introduction: beyond classical linear modelling**
A brief overview of logistic and Poisson regression
Generalities about Maximum Likelihood Estimation

## On the way to Generalized Linear Models

In a number of situations, relations between observations and parameters $\theta$ are not as direct as previously.

This includes notably cases where $Y$ is binary, with a distribution $B(\rho)$ where the probability $\rho$ (say of $Y$ taking the value 1) is a function of $X$ and $\theta$. **Logistic regression** corresponds to a particular case of this.

This includes also cases where $Y$ is integer-valued, for example with a Poisson distribution $\mathcal{P}(\lambda)$ where the intensity $\lambda$ is a function of $X$ and $\theta$. **Poisson regression** corresponds to a particular case of this.

Before going into some overarching principles and investigating these models in more detail, let us first get a brief overview.

---

Introduction: beyond classical linear modelling
**A brief overview of logistic and Poisson regression**
Generalities about Maximum Likelihood Estimation

## Beyond linear models: overview of logistic regression

In logistic regression, the output $Y$ is binary, and has a probability $\rho$ of taking the values 0 or 1 that depends on covariables $X$ through a function parametrized by a vector $\theta \in \mathbb{R}^p$.

In the logistic regression model, the following form is assumed for $\rho$:

$$\rho(x) = \frac{\exp(\sum_{j=1}^p \theta_j f_j(x))}{1 + \exp(\sum_{j=1}^p \theta_j f_j(x))}$$

The function $h : t \in \mathbb{R} \to \exp(t)/(1 + \exp(t))$ is the *standard logistic function*. Its reciprocal bijection is $\text{logit} : \rho \in (0, 1) \to \text{logit}(\rho) = \log(\frac{\rho}{1-\rho})$. Reformulating the model in terms of the $\text{logit}$ function results in

$$\text{logit}(\rho(x)) = \sum_{j=1}^p \theta_j f_j(x)$$

It may seem quite close to a linear model but it is important to stress that least squares estimation does not apply here. Estimating $\theta$ based on binary observations hence calls for some specific procedure(s).

Introduction: beyond classical linear modelling
**A brief overview of logistic and Poisson regression**
Generalities about Maximum Likelihood Estimation

## Beyond linear models: overview of logistic regression

A standard approach to estimate $\theta$ in logistic regression is to maximize the so-called **likelihood function** of parameters given observations, i.e. to maximize the probability of observing $Y$ seen as a function of $\theta$.

Each $Y_i$ is assumed to have been generated independently from a Bernoulli distribution with parameter $\rho(X_i)$, i.e. has probability

$$\rho(X_i)^y (1 - \rho(X_i))^{1-y} = \frac{\exp(y \sum_{j=1}^{p} \theta_j f_j(X_i))}{1 + \exp(\sum_{j=1}^{p} \theta_j f_j(X_i))}$$

of being equal to $y \in \{0, 1\}$. The likelihood of $\eta$ given a vector of observations $(y_1, \ldots, y_n)$ is then defined as

$$L(\eta; (y_1, \ldots, y_n)) = \prod_{i=1}^{n} \rho(X_i)^{y_i} (1 - \rho(X_i))^{1-y_i} = \prod_{i=1}^{n} \frac{\exp(y_i \sum_{j=1}^{p} \eta_j f_j(X_i))}{1 + \exp(\sum_{j=1}^{p} \eta_j f_j(X_i))}$$
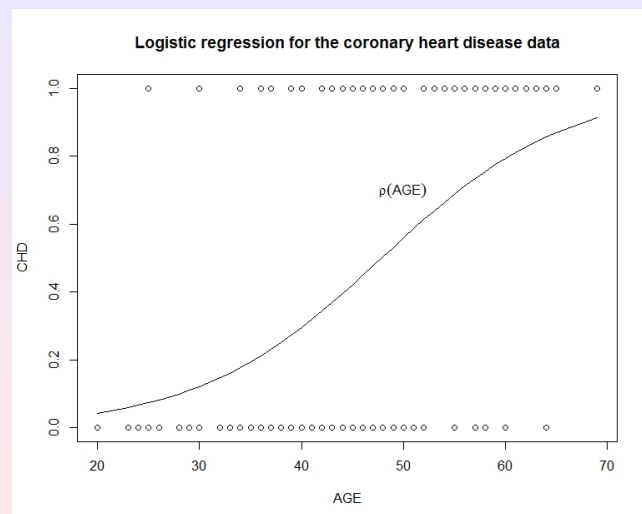
$\widehat{\theta}$ is then found by numerically maximizing $L$. We will study that in more detail in next lectures, but let us now illustrate the model with a first example.

Introduction: beyond classical linear modelling
**A brief overview of logistic and Poisson regression**
Generalities about Maximum Likelihood Estimation

## Beyond linear models: overview of logistic regression

Here we consider an example from a study on "Coronary Heart Disease" (Source: Hosmer, D.W., Lemeshow, S. and Sturdivant, R.X. (2013). Applied Logistic Regression: Third Edition. See `https://www.umass.edu/statdata/statdata/data/chdage.txt`).

The data set consists of $n = 100$ observations, with three variables: for each observation, an identification code, the age of the patient, and a binary variable indicating whether there is a coronary heart disease or not. Let us refer to the last one as $Y$, to the age as $X$ and investigate the relationship between $Y$ and $X$ assuming a simple logistic regression model ($d = 1$).

Introduction: beyond classical linear modelling
**A brief overview of logistic and Poisson regression**
Generalities about Maximum Likelihood Estimation

## Beyond linear models: overview of logistic regression



Logistic regression for the coronary heart disease data

Introduction: beyond classical linear modelling
**A brief overview of logistic and Poisson regression**
Generalities about Maximum Likelihood Estimation

## Beyond linear models: overview of logistic regression

For the sake of brevity we do not dive into the details of the fitted model here and much remains to be said about model interpretation, parameter estimation, confidence intervals and beyond. Along the next lectures, we will have the occasion to analyse the logistic regression model in more detail.

For now we however to a second model, Poisson regression, where the response is this time integer valued. As we will see later, these two seemingly distant examples actually share some essential properties.

Introduction: beyond classical linear modelling
**A brief overview of logistic and Poisson regression**
Generalities about Maximum Likelihood Estimation

## Beyond linear models: overview of Poisson regression

Let us now focus on *count data*, and more particularly on regression when the response of interest is integer valued, i.e. $Y \in \mathbb{N}$.

The Poisson distribution is a standard distribution over $\mathbb{N}$. A r.v. $Y$ is Poisson-distributed whenever its probability mass function is given by

$$P(Y = y) = \frac{\exp(-\lambda)}{y!}\lambda^y \text{ for } y \in \mathbb{N},$$

where the parameter $\lambda > 0$ is called *intensity*.
It is a well known fact that for $Y \sim \mathcal{P}(\lambda)$, $\mathbb{E}[Y] = \text{Var}[Y] = \lambda$.

In Poisson regression, $Y$ is assumed to follow a Poisson distribution which intensity depends on the covariables $X$. More precisely, the logarithm of the intensity function is supposed to depend linearly on basis functions of $X$:

$$\log(\lambda(x)) = \sum_{j=1}^{p} \theta_j f_j(x).$$

---

Introduction: beyond classical linear modelling
**A brief overview of logistic and Poisson regression**
Generalities about Maximum Likelihood Estimation

## Beyond linear models: overview of Poisson regression

Now, given observed $y_1, \ldots, y_n$ from a Poisson regression model with corresponding covariable instances $X_1, \ldots, X_n$, the parameters $\theta$ are typically estimated by maximizing the likelihood function

$$L(\eta; (y_1, \ldots, y_n)) = \prod_{i=1}^{n} \left\{ \exp\left( -\exp\left( \sum_{j=1}^{p} \eta_j f_j(X_i) \right) \right) \frac{\left( \exp(\sum_{j=1}^{p} \eta_j f_j(X_i)) \right)^{y_i}}{y_i!} \right\}$$

As we will see in the following, the maximization of $L$ is not as hard as it may seem since rewriting the problem and digging into some more theory allow practical simplifications as well as theoretical guarantees.

Practical examples will be provided at a later stage, when covering more specifically the topic of Poisson regression. For now let us focus on Maximum Likelihood Estimation in a much more general framework.

---

Introduction: beyond classical linear modelling
A brief overview of logistic and Poisson regression
**Generalities about Maximum Likelihood Estimation**

## About Maximum Likelihood Estimation

In the last two examples, the approach to fit model parameters has been to maximize so-called likelihood functions. We now review a few selected facts about this approach that is central in statistics and which domain of relevance goes well beyond logistic and Poisson regression.

Assume that $(E, \nu)$ is a measured space (often in our context $E = \mathbb{R}^n$ equipped with its Lebesgue measure) and that $(P_\theta)_{\theta \in \Theta}$ is a family of probability measures dominated by $\nu$. Let us denote by $p(\cdot; \theta)$ the corresponding probability density/mass functions, and assume further that $Y$ is a random element with distribution $P_\theta$ for some $\theta \in \Theta$.

The Maximum Likelihood Estimation (MLE) approach consists, given a realization $y$ of $Y$, in estimating $\theta$ by maximizing (*when possible*) the function

$$L(\cdot; y) : \eta \in \Theta \to L(\eta; y) := p(y; \eta).$$

---

Introduction: beyond classical linear modelling
A brief overview of logistic and Poisson regression
**Generalities about Maximum Likelihood Estimation**

## About Maximum Likelihood Estimation

As a motivating example, let us consider the particular case where there are only two alternative values for $\theta$, i.e. $\Theta = \{\theta_1, \theta_2\}$.

Assuming that one observes $Y$ which is assumed to be drawn from $P_\theta$ with $\theta \in \{\theta_1, \theta_2\}$, estimating $\theta$ then amounts to dividing the set $E$ into two disjoint subsets: the set $E_1$, with $\theta$ estimated to be $\theta_1$ when $Y \in E_1$, and $E_2$ with $\theta$ estimated to be $\theta_2$ when $Y \in E_2$. In other words,

$$\widehat{\theta} = \begin{cases} \theta_1 \text{ in case } Y \in E_1, \\ \theta_2 \text{ else.} \end{cases}$$

Now, nothing has been said yet on how to choose $E_1$.

The approach that we will follow here is to consider two kinds of error: the probability that $Y \in E_1$ when $\theta = \theta_2$ and the probability that $Y \in E_2$ when $\theta = \theta_1$, that is to say respectively $P_{\theta_2}(E_1)$ and $P_{\theta_1}(E_2)$.

**Introduction: beyond classical linear modelling**
**A brief overview of logistic and Poisson regression**
**Generalities about Maximum Likelihood Estimation**

## About Maximum Likelihood Estimation

Define now the function $C : E_1 \subset E \to C(E_1) = P_{\theta_2}(E_1) + P_{\theta_1}(E_2)$.

Noting that $P_{\theta_1}(E_2) = 1 - P_{\theta_1}(E_1)$ and remembering that $P_{\theta_i}$ are dominated by $\nu$ with densities $p(\cdot; \theta_i)$, we get that

$$C(E_1) = 1 + \int_{E_1} p(\boldsymbol{y}; \theta_2) - p(\boldsymbol{y}; \theta_1) \mathrm{d}\nu(\boldsymbol{y}).$$

From there, it becomes apparent that $C$ can be minimized by taking

$$E_1 = \{ \boldsymbol{y} \in E : p(\boldsymbol{y}; \theta_2) \leq p(\boldsymbol{y}; \theta_1) \},$$

or in other words, by estimating $\theta$ to be $\theta_1$ for $\boldsymbol{y}$ such that $p(\boldsymbol{y}; \theta_1) \geq p(\boldsymbol{y}; \theta_2)$ (note that what is done in the equality case does not change the value of $C$).

**Introduction: beyond classical linear modelling**
**A brief overview of logistic and Poisson regression**
**Generalities about Maximum Likelihood Estimation**

## About Maximum Likelihood Estimation

MLE is used well beyond the previous (didactic) case where possible values for the parameter $\theta$ were restricted to a set of two.

In a more general framework, another fact that supports MLE is the following. Assume that the $Y_i$'s are i.i.d. with a distribution possessing a density $p_0(\cdot; \theta)$ with respect to some measure $\nu_o$ on $\mathcal{Y}$, where $\theta \in \Theta$, and take $E = \mathcal{Y}^n$ equipped with the product measure $\nu = \nu_o^{\otimes n} = \nu_o \otimes \cdots \otimes \nu_o$.

Then the likelihood of $\eta$ given $\boldsymbol{y}$ is $L(\eta; \boldsymbol{y}) = p(\boldsymbol{y}; \eta) = \prod_{i=1}^n p_0(y_i; \eta)$. Now we observe that the so-called *log-likelihood*

$$\ell(\eta; \boldsymbol{y}) = \log \left( \prod_{i=1}^n p_0(y_i; \eta) \right) = \sum_{i=1}^n \log(p_0(y_i; \theta))$$

is —up to a normalization term $1/n$— a "sample analogue" of the function

$$\bar{\ell}_o : \eta \in \Theta \to \bar{\ell}_o(\eta) = \mathbb{E}[\log(p_0(Y; \eta))],$$

where $Y$ is assumed to follow $P_\theta$.

**Introduction: beyond classical linear modelling**
**A brief overview of logistic and Poisson regression**
**Generalities about Maximum Likelihood Estimation**

## About Maximum Likelihood Estimation

The crux is then that under regularity conditions $\bar{\ell}$ is actually maximal at the true parameter $\theta$. Indeed, considering $\eta \in \Theta$, we have formally

$$\bar{\ell}_o(\theta) - \bar{\ell}_o(\eta) = \mathbb{E}[\log(p_0(Y; \theta))] - \mathbb{E}[\log(p_0(Y; \eta))]$$

$$= \mathbb{E}\left[ -\log \left( \frac{p_0(Y; \eta)}{p_0(Y; \theta)} \right) \right]$$

$$\geq -\log \left( \mathbb{E}\left[ \frac{p_0(Y; \eta)}{p_0(Y; \theta)} \right] \right) = -\log(1) = 0.$$

This is an additional motivation for using MLE in the case of i.i.d. observations, especially when the number of replications is large.

**Introduction: beyond classical linear modelling**
**A brief overview of logistic and Poisson regression**
**Generalities about Maximum Likelihood Estimation**

## About Maximum Likelihood Estimation

Major questions that arise regarding MLE notably include

- for a given model, the existence of some likelihood maximizer(s),
- and, when existing, how to compute the MLE in practice.

Moreover, questions such as the consistency and further asymptotic properties of MLE are of utmost importance both in theory and in practice.

When the likelihood function is sufficiently regular (twice continuously differentiable over $\Theta \subset \mathbb{R}^p$ is often assumed/met here) and the true parameter is assumed to lie on the interior of the domain, the MLE is searched among the roots of the so-called *score function*:

$$U : \boldsymbol{\theta} \in \Theta \to U(\boldsymbol{\theta}) = \nabla \ell(\boldsymbol{\theta}; \boldsymbol{y}) = \left( \frac{\delta}{\delta \theta_j} \ell(\boldsymbol{\theta}; \boldsymbol{y}) \right)_{i=1,\dots,p}.$$

Now, finding solutions to $U(\boldsymbol{\theta}) = \boldsymbol{0}$ (called "normal equation") with a non-linear $U$ is often non-trivial and call for iterative procedures.

**Introduction: beyond classical linear modelling**
**A brief overview of logistic and Poisson regression**
**Generalities about Maximum Likelihood Estimation**

## About Maximum Likelihood Estimation

Another function that comes into play in maximum likelihood estimation is

$$\mathcal{J} : \boldsymbol{\theta} \in \Theta \to \mathcal{J}(\boldsymbol{\theta}) = -\nabla^2 \ell(\boldsymbol{\theta}; \boldsymbol{y}),$$

the so-called "observed information matrix".

$\mathcal{J}$ can notably be used to check if a solution to the normal equations actually is a local maximizer of $\ell$. Also, $\mathcal{J}$ is instrumental in numerical MLE procedures (Cf. Newton-Raphson algorithms). The "expected version" of $\mathcal{J}$ further plays a key role in asymptotic results regarding MLE's variance.

**Introduction: beyond classical linear modelling**
**A brief overview of logistic and Poisson regression**
**Generalities about Maximum Likelihood Estimation**

## About Maximum Likelihood Estimation

Note finally that the existence of solutions, the convergence of algorithms, asymptotic properties of MLE are model-dependent and require serious examination. Next, we will present a distribution class, Natural Exponential Families, for which MLE enjoys convenient properties and that will turn out to be relevant for logistic and Poisson regression.