



Applied Microeconometrics

Fundamentals

Contents

Motivation	3
1 Statistical models	4
1.1   ⊃ Microeconometrics models	4
1.2   ⊃ Regression models	4
1.3   ⊃ Non-/Semi-/Parametric models	5
2 Data	7
2.1   Types of observational data	7
2.2   Sampling procedures	7
3 Statistical inference [under a frequentist approach]	8
3.1   Frequentist vs Bayesian inference	8
3.2   Estimation	9
3.2.1   Regression analysis	9
3.2.2   Estimators	9
3.2.3   Estimator properties	11
3.2.4   Uncertainty in the estimate: computing SEs & CIs	11
3.3   Hypothesis testing	14
3.3.1   Statistical tests	14
3.3.2   Null Hypothesis Significance Testing (NHST) paradigm	14
3.3.3   Type I/II errors, size and power	15
3.3.4   Criticisms of the NHST and ‘statistical significance’	16
4 Statistical inference [under a Bayesian approach]	17
4.1   Steps of Bayesian inference	17
4.2   Choosing $\theta$ ’s prior distribution	18
4.3   Estimating $\theta$ ’s posterior distribution	18
5 Prediction	20
6 Model comparison	21
6.1   Comparing nested models: $F$ tests	21
6.2   Comparing non-nested models: IC, CV	21
7 Other branches of statistical modeling	24
7.1   Statistical Inference Using Agent-based models (ABMs)	24
Key ideas [one pager]	26
References	27

<b>Appendix A A small library of regression models</b>	<b>28</b>
A.1 Expanding from the CLRM . . . . .	28
A.2 Limited outcome models . . . . .	29
A.3 Nonparametric models . . . . .	32
A.4 Multilevel models . . . . .	34

*er: Sections and lines in brown correspond to content which is **very much** ‘under construction’.*

# Motivation

Research questions related to the goal of sustainable development bring together social and natural systems, and are therefore particularly conducive to interdisciplinary work. The social system part demands some training in the social sciences, and in effect interdisciplinary researchers may have an economics background.

Applied microeconomics work in recent years has largely concerned the identification of causal relationships between variables, such that the current dominant methods and terminology are largely fitted to that goal. In applied work from other disciplines, one is likely to encounter alternative types of models, estimation methods, terminology, and even ultimate goals of the statistical analysis (e.g., predictive inference vs causal inference). If nothing else, an applied interdisciplinary researcher should be able to communicate with these different academic disciplines. This means notably understanding what a given method does in statistical terms, in other words: where it fits in the ‘family tree’ of statistical approaches. This will enable them to both: choose the most appropriate method given the problem at hand (when understanding what the method is doing, the empowered researcher need not resort only to the most common method in a given discipline), and justify that choice in front of the different disciplinary communities.

The purpose of this document is therefore twofold:

1. To detail the typical methods of applied microeconomics, which are our reference base. This includes defining and distinguishing common notions that may be conflated (*a model, an equation, a regression, a specification, an estimation method...*);
2. To put those into context, i.e., place them in the greater ‘family tree’ or space of statistical methods, and delineate a few other branches of that tree that may be relevant for empirical interdisciplinary research.

Let us start by defining microeconometrics:

**Econometrics** = (originally) the application of statistical methods to economic data, in order to measure the relationships of economic theory, i.e., obtain estimates that can be given a structural interpretation.

**Microeconometrics** = the use of these statistical methods to study microdata pertaining to individuals, households, and firms.

Ultimately, applied economics is a specific area of applied statistics. A distinguishing feature is the emphasis placed on causal modeling.

# 1 Statistical models

**A model** is a formal representation of a theory about a system, to ultimately describe that system.

**A statistical model** is a mathematical model of the data generating process (DGP)<sup>a</sup> of the sample  $\{y_i, X_i\}_{i=1}^n$ .

- What distinguishes it from other mathematical models is that it is non-deterministic: some variables are stochastic or “random”, they have probability distributions.<sup>b</sup>
- It is written as relationships between these *random* variables and some non-random variables, to study the **variation** of random variables. Specifically, it can serve 3 purposes: description (summarizing a sample); extraction of information; prediction.

<sup>a</sup>Formally, it combines the set of possible observations or “sample space”  $\mathcal{S}$  and a collection of joint probability distributions on  $\mathcal{S}$  (which ideally would include the “true” probability distribution induced by the DGP; but it doesn’t need to, we accept that are models are false).

<sup>b</sup>Indeed, the task of statistics can be described as quantifying evidence and reasoning under *uncertainty*.

## 1.1 ⊃ Microeconometrics models

All empirical investigations in *microeconometrics* aim to uncover important relationships to understand microeconomic behavior. They can broadly be separated into two types of approaches, depending on the extent to which they rely on microeconomic theory:

- **Structural analysis** heavily depends on economic theory. Model specifications are derived from specifications of the economic behavior. The goal is to analyze structural relationships for interdependent microeconomic variables (e.g., to estimate structural parameters that characterize individual preferences or technological relationships).

$$g(y, X, e|\theta) = 0, \quad \theta = \text{structural parameters}$$

- **Reduced form analysis** makes much less use of economic theory. The goal is to uncover associations among variables, by using regression models.

The **reduced form** of a system of structural equations is the result of solving the system for the dependent (i.e., nonlagged and endogenous) variables. This **gives the dependent variables as functions of the independent variables (exogenous variables or lags of the dependent)**.

$$y = h(X, e|\pi), \quad \pi = \text{reduced form parameters that are functions of } \theta$$

## 1.2 ⊃ Regression models

**A regression model** is a statistical model which models a *dependent variable*  $y$  as a function of *independent variables*  $X$ .

It has 4 components:  $y, X$ , unknown parameters  $\beta$  and an error term  $e$ .

The variables  $\{y, x_1, \dots, x_k\}$  have an unknown joint distribution and complicated covariance structure. Instead of looking at the full joint distribution, regression models simplify the problem by **focusing on the conditional distribution<sup>1</sup> of  $y$ , given  $X$** .

<sup>1</sup>Different regression models will look at different parts of the distribution, and specify them differently. Ex: classical linear regression model:  $\mathbb{E}[y|X] = f(X) = X\beta$ ; quantile regression model:  $\mathbb{Q}[y|X] = f(X) \dots$

Generally speaking, we use regression models for three major functions: estimation, hypothesis testing, and prediction.

**Writing a regression model** means that we consider that a sample  $\{y_i, X_i\}_{i=1}^n$  is generated by the process described by that model. We can write the model interchangeably:

- as a system of  $n$  equations:  $y_i = f(X_i, e_i|\beta)$ ,  $\forall i = 1, \dots, n$
- using matrix notation:  $y = f(X, e|\beta)$ , where the error term  $e$  is a vector of  $n$  random variables, with an  $n \times n$  symmetric covariance matrix.

In Bayesian inference, parameters are considered random variables, therefore the data generating process is written conditional on  $\beta$ :  $y = f(X, e|\beta)$ . In frequentist inference, the parameters are considered fixed, therefore we write  $y = f(X, e, \beta)$ .

Example: The classical linear regression model assumes a linear conditional expectation function and an additive error term:

$$y_i = X_i' \beta + e_i = \beta_0 + \sum_{j=1}^k \beta_j x_{ij} + e_i, \quad e_i \stackrel{\text{iid}}{\sim} (0, \sigma^2), \quad i = 1, \dots, n$$

$$y = X\beta + e, \quad e \sim (0, \sigma^2 \mathbf{I}_n)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

**Choosing a model specification** To carry out regression analysis, one must first choose a model specification: select which independent variables to include and an appropriate functional form of  $f$ .

Specification error occurs when either the functional form or the choice of independent variables poorly represents relevant aspects of the true DGP. Though “correct specification” is, in practice, unrealistic, as we do not observe the true DGP, we try to avoid the three basic types of misspecification:

- using an inappropriate functional form;
- including an  $x$  that is theoretically *irrelevant* (it has no partial effect on  $y$ )  $\rightarrow$  *overspecified* model;
- excluding an  $x$  that is theoretically *relevant* (it may cause  $y$ )  $\rightarrow$  *underspecified* model.

### 1.3 $\supset$ Non-/Semi-/Parametric models

The specification of a statistical model can be:

- **parametric or “finite-dimensional”**: the model is a family of distributions that has a *finite* number of parameters.<sup>2</sup> We assume that the data come from a population that can be adequately modeled by a probability distribution with a *fixed* set of parameters.
  - For *regression* models, it means that the distribution of the error term is fully characterized.

---

<sup>2</sup>Recall that a statistical model is a collection  $\mathcal{P}$  of probability distributions on some sample space  $\mathcal{S}$ . We can write it as  $\mathcal{P} = \{P_\theta | \theta \in \Theta\}$ , where  $\Theta$  is the parameter space. Hence we can write a parametric model as  $\mathcal{P} = \{P_\theta | \theta \in \Theta \subseteq \mathbb{R}^k\}$ .

- When the parameters uniquely specify the distribution,<sup>3</sup> we say that they are “identifiable”.

*Ex: The Poisson family of distributions is parametrized by a single number  $\lambda > 0$ ; the normal family is parametrized by two numbers  $\{\mu, \sigma\}$ .*

- **non-parametric:** the model makes no assumptions about a parametric distribution, it determines it from data.<sup>4</sup> The model has parameters, but their number and nature aren’t fixed in advance.
  - For *regression* models, it means that no parametric form is assumed for the relationship between the dependent and the independent variables. *Ex: Kriging; LOESS.*
- **semi-parametric:** the model combines parametric and nonparametric models.
 

*Ex: Only a few moments are specified:  $\mathbb{E}[e] = 0$  and  $\mathbb{V}[e] = \mathbb{E}[ee'] = \Omega$ .*

Why care about parametrization? Because what we are interested in is the class of probability distributions (as this will be our postulated model for observed data), and the parameter describes an integral feature of the probability distribution, s.t. knowledge about the parameter translates easily to knowledge about the distribution.

### Identification in parametric models

**Identification of a parameter** = its unique determination, given sufficient observations. *Assuming we had enough observations, could we determine the parameter?*

The model being “well-identified”, i.e., the identification of all its parameters, is required for consistent estimation — and thus for meaningful statistical inference. It can be obtained through the functional form (by the parametrization of the error distribution) or from exclusion, inequality and covariance restrictions.

Example of non-identification: in the linear regression  $y = X\beta + e$ , perfect collinearity between regressors means we can’t identify  $\beta$ .

<sup>3</sup>I.e., the correspondence of each distribution in  $\mathcal{P}$  with a  $\theta$  is 1-1, s.t.  $P_{\theta_1} = P_{\theta_2} \Rightarrow \theta_1 = \theta_2$ .

<sup>4</sup>Nonparametric regression requires larger sample sizes than regression based on parametric models, because the data must supply the model structure in addition to the model estimates. Nonparametric models also usually contain strong assumptions about independencies.

## 2 Data

Empirical studies can be separated into two classes, based on the type of data collected:

Study	Data collection
Experimental	The researcher records data about subjects while applying treatments and controlling conditions (active participation).
Observational	The researcher records data about subjects without applying a treatment (passive participation). If the goal is to uncover characteristics of a population, they may: <ul style="list-style-type: none"> <li>• inspect the entire population: perform a <b>census</b>;</li> <li>• inspect a subset: take <b>sample data</b> <math>S_t</math> from the population probability distribution <math>F(W_t   \theta_t)</math>.</li> </ul>

### 2.1 Types of observational data

Observational data can be grouped into 3 categories, based on the dimensions: units (N) and time (T):

- **Cross-sectional** [N]: observations for several units, at one point in time;
- **Time series** [T]: observations for a single unit, at repeated points in time;
- **Longitudinal** [N × T]: observations for several units, at repeated points in time.

When *the same units* are observed over time, we have **panel data**.<sup>5</sup> The panel can be:

- balanced: all observed units  $i$  have data across all periods  $t$ ;
- unbalanced: some units have more observations than others.

Variation *between* units at one point in time is called *between*-variation, while variation *within* one unit across time is called *within*-variation. The total variance of observed variables can be split into within- and between-variation.

One of the strengths of longitudinal data is its potential for supporting causal relationships because of its ability to deal with observable and unobservable effects.

### 2.2 Sampling procedures

**Random sampling** ensures the *data* probability distribution is the same as the *population* distribution. If sampling isn't random, it is **biased**: the data distribution differs from the population distribution.

Common random sampling procedures include:

- **Simple random sampling** — the assumption on which statistical inference theory is based.
- **Stratified random sampling**: the population is divided into L subgroups or “strata”, of  $N_1 \neq N_2, \dots, N_L$  units. Simple random samples of sizes  $n_1, n_2, \dots, n_L$  are drawn independently.
  - **Proportionate stratified random sampling**  
Ex: in a “10% sample, stratified across subgroups”, the same fraction is applied on each subgroup.

---

<sup>5</sup>“Panel data” and “longitudinal data” are often used interchangeably, as most often it is the same units that are observed over time. However keeping the distinction, as delineated in [Mertens et al. \(2017\)](#), can be useful.

### 3 Statistical inference [under a frequentist approach]

**Inferential statistics** or **statistical inference** consists in *inferring* properties of a population,<sup>a</sup> by calculating statistics from a sample drawn from the population.

It contrasts with descriptive statistics, which is solely concerned with properties of the observed data, not a larger population.

<sup>a</sup>*Population, DGP, and underlying probability distribution* could be used interchangeably. The data observed are of random variables, and we want to estimate parameters  $\theta$  of their joint probability distribution. Making statistical inferences = deducing properties of (conditional) probability distributions.

Statistical inference combines data and (explicit or implicit) prior assumptions,<sup>6</sup> and generally involves:

- **Estimation** — 1. Estimating the value (point estimation) or potential range of values (confidence interval estimation) of an unknown parameter  $\theta$  that characterizes the probability distribution of some feature of interest in the population; 2. Assessing the uncertainty around that estimate.
- **Hypothesis testing** — Testing for a specific value of the unknown parameter  $\theta$ .

#### 3.1 Frequentist vs Bayesian inference

There are two main paradigms for inference, whose difference is rooted in their definition of probability. Consider a parameter  $\theta$  of unknown true value  $\theta_0$ , and an *event*  $\theta = \tilde{\theta}$  (i.e.,  $\theta$  taking this value  $\tilde{\theta}$ ).

Frequentist approach	Bayesian approach
Definition of <i>probability</i> $\mathcal{P}$	
$\mathcal{P} \equiv$ the frequency of occurrence of an event; hence only repeatable events have $\mathcal{P}$ s (ex: coin flips).	$\mathcal{P} \equiv$ one's belief in an event; hence any event, incl. non-repeatable, can have a $\mathcal{P}$ .
Implication regarding $\theta$	
$\implies \theta$ is <i>fixed</i> . We can't assign $\mathcal{P}$ s to events such as $\theta \leq \tilde{\theta}$ . We handle our uncertainty in the value of $\theta$ by limiting error rates (over imaginary experiments).	$\implies \theta$ is a <i>random variable</i> . We can assign a $\mathcal{P}$ distribution over possible values of $\theta$ , to represent our uncertainty/belief in the value of $\theta$ .
Estimating $\theta$ using data	
1. Collect sample data, estimate the value (point $\hat{\theta}$ ) or potential range of values (confidence interval $\text{CI}[\hat{\theta}]$ ) of $\theta$ that is most consistent with the data.  Result: a conclusion, summary of data, in the form of: – a “true/false” statement from a significance test, expected to be correct ...% of the time; or – a confidence interval, expected to cover the true value ...% of the time. (“time” = number of possible samples from the pop.)	1. Define a $\mathcal{P}$ distribution over possible values of $\theta$ 2. Collect sample data and update this distribution, by applying Bayes' theorem to each possible value: $P(\tilde{\theta} \text{data}) = \frac{P(\text{data} \tilde{\theta}) \times P(\tilde{\theta})}{P(\text{data})}$ Result: a <i>posterior</i> $\mathcal{P}$ distribution for $\theta$ . We can compute a 95% credible interval, s.t. “after seeing the data, there is a 95% chance that this CI contains the true $\theta$ .”
Prediction	
Use the point estimate $\hat{\theta}$ as the most likely value of $\theta$ , and its CI.	Use the full posterior $\mathcal{P}$ distribution of $\hat{\theta}$ , which allows for taking into account the uncertainty in $\hat{\theta}$ .

<sup>6</sup>E.g., in Bayesian inference, an accurate prior (an assumption) will pull our estimates toward the true value. In frequentist inference, assuming a particular error distribution (i.e., parametric inference techniques) lends us power.



The sections below describe the *ABC* of statistical inference in the context of regression analysis, and under a frequentist approach, which is the classical approach in econometrics.

## 3.2 Estimation

### 3.2.1 Regression analysis

**Regression analysis** = a set of statistical processes for **estimating the relationship between a dependent variable  $y$  and independent variables  $X$** :<sup>a</sup>  $y = f(X, e|\theta)$ .

It is a way of summarizing and drawing inferences from data. It can have two purposes:

- prediction (interest is in  $\hat{y}$ ): the **prediction** of the conditional distribution of  $y$ , given  $X$ ;
- comparison (interest is in  $\hat{\beta}$ ): comparing groups (which differ in  $X$ ) or estimating causal effects.<sup>b</sup>

<sup>a</sup>Recall the definition of a [regression model](#).

<sup>b</sup>Regression coefficient estimates  $\hat{\beta}$  should be interpreted as “effects” only in causal inference. Otherwise, the safest interpretation is as a comparison, using the word “differences” rather than the words “effects” or “changes”. E.g., “the average difference in  $y$ , comparing two individuals that differ in  $x$  by one unit, is  $\hat{\beta} = 0.29$ ” or “adding 1 unit to  $x$  corresponds to an increase of  $\hat{\beta} = 0.29$  in an individual’s predicted  $y$ ”.

△ Regressions calculate the *distribution of values* of the relation between  $y$  and  $X$ . The output is a conditional *distribution*  $f_{y|X}$ . We can then choose to focus on its conditional mean  $\mathbb{E}[y|X]$ , its conditional quantiles  $\mathbb{Q}_{y|X}()$ ...

### 3.2.2 Estimators

We have a set of observations  $x_1, \dots, x_n$ , i.e., a realization of the sample of random variables  $X_1, \dots, X_n$ .

**An estimand  $\theta$**  is a quantity of interest that we want to estimate, e.g., a parameter or some summary of the data. *Ex: the population mean  $\mu_X$ .*

**An estimator  $\hat{\theta}_n$**  of an estimand is a sample statistic, i.e., a function of the random sample (and therefore a random variable):  $T_n = t(X_1, \dots, X_n)$ . Its values will vary sample to sample.

*Ex: the sample mean  $\bar{X}_n$  is an estimator for the population mean  $\mu_X$ .*

**An estimate** is a realization of that r.v.  $\hat{\theta}_n$ , calculated for our specific sample:  $t_n = t(x_1, \dots, x_n)$ .

The most common estimators in microeconometrics are extremum estimators: they solve a min/max problem.

- **Maximum Likelihood (ML)**<sup>7</sup>

We want to find the value of  $\theta$  that makes the observed data most likely. The likelihood function in a

<sup>7</sup>The ML estimator is just a type of statistic, and can be conceptualized under either inference approach. From the vantage point of Bayesian inference, ML is a special case of ‘maximum a posteriori’ estimation that assumes a uniform prior distribution of the parameters. In frequentist inference, ML is a special case of extremum estimation, where the objective function is the likelihood.

regression model is the probability density of the data given the parameters and predictors:

$$\begin{aligned}\mathcal{L}(y | X, \theta) &= f(X_1, \dots, X_n, \theta) \\ &= f(X_1, \theta) \dots f(X_n, \theta) \\ &= \prod_{i=1}^n f(X_i, \theta) \\ \log \mathcal{L}(y | X, \theta) &= \sum_{i=1}^n \log f(X_i = x_i, \theta)\end{aligned}$$

We compute  $\hat{\theta}_{\text{ML}} \equiv \operatorname{argmax}_{\theta} \mathcal{L}(y | X, \theta) = \operatorname{argmax}_{\theta} \log \mathcal{L}(y | X, \theta)$

### • Least Squares (LS)

The fit of a model  $y = g(X, e)$  to each data point  $i$  is measured by its residual  $r_i \equiv y_i - g(x_i, \hat{\beta})$ . We are interested in the values of the parameters that best fit the data, i.e., that minimize the sum of the squares of (eventually a function  $k()$  of) the residuals.<sup>8</sup>

$$\hat{\theta}_{\text{LS}} \equiv \operatorname{argmin}_{\theta} \sum_{i=1}^n k(r_i)^2$$

When the model is linear, i.e., a linear combination of the parameters  $g(X, \beta) = \sum_j \beta_j h_j(X)$ , Least Squares is a **Linear Least Squares (LLS)**.

#### \* Ordinary Least Squares (OLS)

The OLS estimator has an exact closed-form solution:

$$\hat{\beta}_{\text{OLS}} \equiv \operatorname{argmin}_{\beta} \sum_{i=1}^n r_i^2 = (X'X)^{-1} X'y$$

In the simple case of the univariate regression model ( $y = \alpha + \beta x + e$ ), the estimand is  $\beta_{\text{OLS}} = \frac{\operatorname{cov}[x, y]}{\operatorname{V}[x]}$  and the estimator (its sample analog)  $\hat{\beta}_{\text{OLS}} = \frac{\frac{1}{n} \sum_i (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n} \sum_i (x_i - \bar{x})^2}$ . In a multivariate regression model, the coefficient on each  $x_k$  is  $\beta_{\text{OLS}}^k = \frac{\operatorname{cov}[\tilde{x}_k, y]}{\operatorname{V}[\tilde{x}_k]}$  where  $\tilde{x}_k$  is the residual from the regression of  $x_k$  on all the other covariates.

#### \* Weighted Least Squares (WLS)

When errors are heteroscedastic, i.e., each has variance  $\sigma_i$ , OLS won't be efficient among linear unbiased estimators. For least squares to give us the most *efficient* linear unbiased estimator, we minimize a *weighted* sum of squared residuals, using weights  $w_i \propto \frac{1}{\sigma_i}$ .

$$\hat{\beta}_{\text{WLS}} \equiv \operatorname{argmin}_{\beta} \sum_{i=1}^n w_i r_i^2$$

#### \* Generalized Least Squares (GLS)

When errors are heteroscedastic or correlated, i.e., when  $x_1, \dots, x_n \stackrel{iid}{\sim} f(x|\theta)$  doesn't hold (the

---

<sup>8</sup>Indeed, let  $e \equiv y - \hat{y}$  be the unobserved error,  $L(e)$  the loss. We want to minimize the expected loss  $\mathbb{E}[L(e)|X]$ . So we look for the fit  $g(X, \hat{\beta})$  that minimizes the mean of that function  $L()$  of the residuals. For a squared error loss function  $L(e) = e^2$ , it means minimizing the sum of squared residuals  $\sum_i r_i^2$ . That fit is the **conditional mean**:  $g(X, \hat{\beta}_{\text{LS}}) \equiv \operatorname{argmin}_{g(\cdot)} \mathbb{E}[(y - g(X, \beta))^2] = \dots = \mathbb{E}[y|X]$ .

covariance matrix  $\Omega \equiv \text{cov}[e|X]$  is not diagonal with values  $\sigma^2$ , OLS will again be inefficient. We minimize instead the squared *Mahalanobis length*<sup>9</sup> of the residuals:

$$\hat{\beta}_{\text{GLS}} \equiv \underset{\beta}{\text{argmin}} \sum_{i=1}^n \overrightarrow{d_M}^2(r_i)$$

When the model is a linear combination of the parameters, the GLS estimator has an exact closed-form solution:  $\hat{\beta}_{\text{GLS}} = \underset{\beta}{\text{argmin}} (y - X\beta)' \Omega^{-1} (y - X\beta) = \dots = (X' \Omega^{-1} X)^{-1} X' \Omega^{-1} y$

\* **Two-Stage Least Squares (2SLS)**

When regressors are correlated with the errors, we need a matrix of instruments  $Z$  s.t.  $\mathbb{E}[z_i e_i] = 0$ .

$$\hat{\beta}_{\text{2SLS}} = (X' Z (Z' Z)^{-1} Z' X)^{-1} X' Z (Z' Z)^{-1} Z' y$$

• **Least (symmetric) absolute error**

We are interested in minimizing a different loss function: the absolute error loss,  $L(e) = |e|$ . The corresponding estimator will be more robust to outliers. The optimal fit, i.e., the least absolute deviations fit, is the **conditional median**:  $g(X, \hat{\beta}_{\text{LSA}}) \equiv \underset{g(\cdot)}{\text{argmin}} \mathbb{E}[|y - g(X, \beta)|] = \dots = \text{med}_{y|X}$ .

• **Least asymmetric absolute error**

We can generalize to an asymmetric loss function:  $L_\alpha(e) \equiv \begin{cases} (1-\alpha)|e| & \text{if } e < 0 \\ \alpha|e| & \text{if } e \geq 0 \end{cases} = (\alpha - \mathbb{1}\{e < 0\}) \times e$ , which places a different penalty on overprediction and underprediction. The optimal fit is the **conditional quantile**  $g(X, \hat{\beta}_{\text{LAA}, \alpha}) \equiv \underset{g(\cdot)}{\text{argmin}} \mathbb{E}[L_\alpha(y - g(X, \beta))] = \dots = \mathbb{Q}_{y|X}(\alpha)$ .

*Note: We have phrased all of the above in terms of the objective of finding the best fit. We could have also phrased it with the objective of making predictions about a specific part of the outcome distribution:*

Objective: fit	Objective: prediction	Optimal estimator/predictor
$\min L(e) \equiv e^2$	predict $\mathbb{E}[y X]$	$\hat{\beta}_{\text{LS}} \equiv \underset{\beta}{\text{argmin}} \sum_i (y_i - g(X_i, \beta))^2$
$\min L(e) \equiv  e $	predict $\text{med}_{y X}$	$\hat{\beta}_{\text{LSA}} \equiv \underset{\beta}{\text{argmin}} \sum_i  y_i - g(X_i, \beta) $
$\min L_\alpha(e) \equiv (\alpha - \mathbb{1}\{e < 0\})e$	predict $\mathbb{Q}_{y X}(\alpha)$	$\hat{\beta}_{\text{LAA}, \alpha} \equiv \underset{\beta}{\text{argmin}} \sum_i L_\alpha(y_i - g(X_i, \beta))$

### 3.2.3 Estimator properties

See section 2 in <https://clairepalandri.github.io/docs/CLRM&estimators.pdf>.

### 3.2.4 Uncertainty in the estimate: computing SEs & CIs

#### i. The uncertainty in any sample statistic can be captured by its SE & CI<sub>95%</sub>

Samples are not unique. Many different samples could have been taken from the population. Any sample statistic (sample mean, slope parameter estimates...) will vary from sample to sample, hence it is a random variable, with a *sampling* probability distribution.

<sup>9</sup>The Mahalanobis distance is a measure of the distance between a point P and a distribution D. It is a multi-dimensional generalization of the idea of measuring how many standard deviations away P is from the mean of D. It is unitless and scale-invariant, and takes into account the correlations of the data set.

We are interested in the population parameter  $\theta$ , and have computed an estimate  $\hat{\theta}$  from our sample. As different samples would have lead to different  $\hat{\theta}$ s,  $\hat{\theta}$  has a sampling distribution. If the distribution is rather condensed, i.e., the standard deviation is low *relative to the estimate*, it means we have high certainty about our estimate. We would quantify this certainty by computing  $\text{SD}[\hat{\theta}]$  – and then use it to construct confidence intervals and test statistics. As we do not observe the sampling distribution (we haven’t taken all the possible samples), we cannot observe  $\text{SD}[\hat{\theta}]$ . However, we can estimate it, and we’ll call that estimate a “standard error”  $\text{SE}[\hat{\theta}]$ .

For any sample statistic  $\hat{\theta}$ , estimated with  $n - k$  degrees of freedom:

- **Standard Error  $\text{SE}[\hat{\theta}]$**  = an estimate of the standard deviation of its distribution.
- The **95% Confidence Interval  $\text{CI}_{95\%}[\theta]$**  = the range of values s.t. “*I have a 95% confidence level that the true  $\theta$  is in that range.*”

**Correctly interpreting the CI** This confidence interval is based on the *sampling* distribution; the confidence refers to our uncertainty about the *sampling* method. The CI is therefore correctly interpreted in terms of repeated samples: “*Imagine we drew all possible random samples of size  $n$ . This interval would contain the true  $\theta$  in 95% of the samples.*”<sup>10</sup> I.e., we believe the 95% CI contains the true value, with the understanding that we’ll be wrong 5% of the time. Another — maybe more adequate — name suggested for such intervals is “compatibility intervals”, as they give a range of parameter values that are most compatible with our data and model/assumptions (Gelman and Greenland, 2019).

## ii. Traditional approach: asymptotic theory

Consider a parameter of interest  $\theta$ , and its estimator with  $n - k$  degrees of freedom  $\hat{\theta}$ .

### 1. Standard Error $\text{SE}[\hat{\theta}]$

- Example 1:  $\theta$  is the population mean  $\mu_x$ ,  $\hat{\theta}$  is the sample mean  $\bar{x}$ .
  - Population:  $X$ ’s mean  $\mu_x$  and variance  $\sigma_x^2$  are unobserved.
  - Sample: We measure the sample mean  $\bar{x}$ . Its variance  $\mathbb{V}[\bar{x}] = \frac{\sigma_x^2}{n}$  is unobserved, as the population variance  $\sigma_x^2$  is unobserved. A reasonable estimate for  $\sigma_x^2$  that we do observe is the *sample* variance  $s_x^2$ .<sup>11</sup> We can thereby estimate  $\mathbb{V}[\bar{x}]$  by  $\hat{\mathbb{V}}[\bar{x}] \equiv \frac{s_x^2}{n}$ , and  $\text{SD}[\bar{x}]$  by  $\text{SE}[\bar{x}] \equiv \frac{s_x}{\sqrt{n}}$ .
- Example 2:  $\theta$  is a regression slope  $\beta_{\text{OLS}}$  in the multivariate linear regression model.
  - Population: parameter  $\beta$  and error variance  $\sigma^2$  are unobserved.
  - Sample: We measure the parameter estimate  $\hat{\beta} \sim (\beta, \mathbb{V}[\hat{\beta}])$ . The formula of  $\mathbb{V}[\hat{\beta}]$  is known but unobserved — as it is notably a function of  $\sigma$ .  
For simplicity, consider the simple case of normal errors:  $e|X \sim \mathcal{N}(0, \sigma^2 I)$ . Then  $\mathbb{V}[\hat{\beta}] = \sigma^2(X'X)^{-1}$ . We can consistently estimate the population variance  $\sigma^2$  by the bias-adjusted

<sup>10</sup>This is a probability statement about the interval, not the population parameter. It says  $P(\beta \in \text{CI} \mid \beta) = 95\%$ . This is different from saying “*there is a 95% probability that the true  $\beta$  lies within this range*”, i.e.,  $P(\beta \in \text{CI} \mid \text{CI}) = 95\%$ . CIs are a frequentist concept, and this second erroneous interpretation contradicts the frequentist interpretation of probability. In the strict frequentist paradigm, the parameter is unobserved but it is set, so a probability statement on its value does not make sense. The probability applies to the interval, not to the true parameter value.

<sup>11</sup> $\triangle$  The standard deviation of the sample  $s$  has nothing to do with the standard error of the estimate  $\text{SE}[\hat{\theta}]$ . The first converges to the standard deviation of the population  $\sigma$  as  $n \rightarrow \infty$ , the second to 0.

sample variance  $s^2 \equiv \frac{1}{n-k} \sum_i r_i^2$ .<sup>12</sup> We can thereby estimate  $\mathbb{V}[\hat{\beta}]$  by  $\widehat{\mathbb{V}}[\hat{\beta}] \equiv s^2 (X'X)^{-1} = \frac{1}{n-k} \sum_i r_i^2 (X'X)^{-1}$ , and  $\text{SD}[\hat{\beta}]$  by  $\text{SE}[\hat{\beta}] \equiv \sqrt{\frac{1}{n-k} \sum_i r_i^2 (X'X)^{-1}}$ .

## 2. Confidence Intervals $\text{CI}[\theta]$

We want to give a range of estimates for the unknown parameter  $\theta$ . Consider *the estimate of the centered and standardized estimate*  $\frac{\hat{\theta}-\theta}{\text{SE}[\hat{\theta}]}$ .<sup>13</sup>

- Example 1:  $\hat{\theta} \equiv \bar{x}$   
 $\frac{\bar{x}-\mu}{\sqrt{s^2/n}}$  has a  $\mathcal{T}$  distribution with  $n-1$  degrees of freedom. Hence, by definition:

$$\begin{aligned} & \text{P}\left(q_{\mathcal{T}_{n-1}}(0.025) \leq \frac{\bar{x}-\mu}{\sqrt{s^2/n}} \leq q_{\mathcal{T}_{n-1}}(0.975)\right) = 0.95 \\ \iff & \text{P}\left(\bar{x} - q_{\mathcal{T}_{n-1}}(0.975) \frac{s}{\sqrt{n}} \leq \mu \leq \bar{x} - q_{\mathcal{T}_{n-1}}(0.025) \frac{s}{\sqrt{n}}\right) = 0.95 \end{aligned}$$

where  $q_{\mathcal{T}_{n-1}}()$  is the quantile function of the  $\mathcal{T}_{n-1}$  distribution. We can thus define the 95% CI:

$$\text{CI}_{95\%}[\mu] \equiv \left[ \bar{x} - q_{\mathcal{T}_{n-1}}(0.975) \frac{s}{\sqrt{n}} ; \bar{x} - q_{\mathcal{T}_{n-1}}(0.025) \frac{s}{\sqrt{n}} \right] = \left[ \bar{x} \pm q_{\mathcal{T}_{n-1}}(0.975) \frac{s}{\sqrt{n}} \right]$$

- Example 2:  $\hat{\theta} \equiv \hat{\beta}_{\text{OLS}}$   
*If errors are normally distributed*, then the sampling distribution of  $\frac{\hat{\beta}-\beta}{\text{SE}[\hat{\beta}]}$  is a Student's  $\mathcal{T}$  distribution with  $n-k$  degrees of freedom. Hence, by definition:

$$\begin{aligned} & \text{P}\left(q_{\mathcal{T}_{n-k}}(0.025) \leq \frac{\hat{\beta}-\beta}{\text{SE}[\hat{\beta}]} \leq q_{\mathcal{T}_{n-k}}(0.975)\right) = 0.95 \\ \iff & \text{P}\left(\hat{\beta} - q_{\mathcal{T}_{n-k}}(0.975) \text{SE}[\hat{\beta}] \leq \beta \leq \hat{\beta} - q_{\mathcal{T}_{n-k}}(0.025) \text{SE}[\hat{\beta}]\right) = 0.95 \end{aligned}$$

We can thus define the 95% CI:

$$\text{CI}_{95\%}[\beta] \equiv \left[ \hat{\beta} - q_{\mathcal{T}_{n-k}}(0.975) \text{SE}[\hat{\beta}] ; \hat{\beta} - q_{\mathcal{T}_{n-k}}(0.025) \text{SE}[\hat{\beta}] \right] = \left[ \hat{\beta} \pm q_{\mathcal{T}_{n-k}}(0.975) \text{SE}[\hat{\beta}] \right]$$

The larger its degrees of freedom, the closer a  $\mathcal{T}$  distribution gets to the standard normal distribution. Therefore, in both examples, when  $n-k$  is sufficiently large, we can simply use the normal distribution.<sup>14</sup>

$$\text{CI}_{95\%}[\theta] \simeq \left[ \hat{\theta} \pm q_{\mathcal{N}}(0.975) \text{SE}[\hat{\theta}] \right] = \left[ \hat{\theta} \pm 1.96 \text{SE}[\hat{\theta}] \right]$$

### iii. Simulation approach: Bootstrap

The traditional approach relies on the assumed *asymptotic* sampling distribution of the statistic. This distribution rests on asymptotic theory (that usually leads to limit normal and  $\chi_2$  sampling distributions). When our sample size is small (making this asymptotic approximation incorrect), or when analytical expressions for the uncertainty of the particular statistic are complicated, i.e., when conventional analytic approximations fail, we can create an alternative sampling approximation of the finite-sample distribution of interest by “**Bootstrap**”.

The Bootstrap procedure is a way to estimate the sampling distribution of the sample statistic, by resampling with replacement from the current sample to generate multiple “resamples”.<sup>15</sup> Supposing 100 Bootstrap resamples, we can obtain 100 estimates and estimate  $\text{SE}[\hat{\theta}]$  by their standard deviation.

Advantages and limits:

<sup>12</sup>The bias here refers to that from the reduced degrees of freedom stemming from estimating the sample means.

<sup>13</sup>The standardized estimate is  $\frac{\hat{\theta}-\theta}{\text{SD}[\hat{\theta}]}$ . The scaling term  $\text{SD}[\hat{\theta}]$  is unknown, therefore we replace it by its estimate  $\text{SE}[\hat{\theta}]$ .

<sup>14</sup> $1.96 \simeq 2$ , therefore it is common to read that statistically significant estimates are at least two standard errors from zero.

<sup>15</sup>Of course, we sample with replacement, to get samples of the same size  $n$ .

- + It does not assume any underlying distribution of the data.
- + It can be applied to any sample statistic.
- + Bootstrap CIs are asymptotically consistent (though we can't know the true CI) and more accurate than the traditional intervals.
- Inference still relies on an appropriately drawn sample; and assumes independent resamples. Therefore with structured models, one must think carefully about the design of the resampling procedure (e.g. with clusters: should we sample within or across clusters?).
- Simple but time-consuming.

### 3.3 Hypothesis testing

#### 3.3.1 Statistical tests

**A statistical test** is a method of verifying a statistical hypothesis.

**A statistical hypothesis** is a hypothesis on the probability distribution of  $T$ , where  $T$  is a **test-statistic** computed from the data, whose probability distribution is connected to our research question.

The general approach to conducting a statistical test consists of the following steps:

1. Write the null hypothesis  $H_0$  — the hypothesis to nullify.
2. Design a test statistic  $T$  that summarizes the deviation of the data from what would be expected under  $H_0$ , and has a specific distribution under  $H_0$ . Ex:
  - a t-test is a test in which the test statistic has a Student's  $\mathcal{T}$  distribution under  $H_0$ ;<sup>16</sup>
  - an F-test is a test in which the test statistic has an  $\mathcal{F}$  distribution under  $H_0$ .
3. Compute the realized value of  $T$  for our data:  $T_{\text{obs}}$ .
4. Look whether it falls in the tails of the distribution. That would mean it is very unlikely given  $H_0$ . Therefore we can reasonably reject  $H_0$ .

#### 3.3.2 Null Hypothesis Significance Testing (NHST) paradigm

Our goal is to statistically test the **hypothesis of a relationship between  $y$  and  $x_j$** , i.e., that  $\beta_j \neq 0$ . Null hypothesis testing proceeds by *reductio ad absurdum*: a hypothesis is assumed valid if its counterclaim is highly implausible. We'll test whether  $\beta_j = 0$  is highly implausible.

---

<sup>16</sup>t-tests are commonly applied for test statistics that would follow a normal distribution if the value of the scaling term were known. When the scaling term is unknown and is replaced by an estimate based on the data, these test statistics (under certain conditions) follow a Student's  $\mathcal{T}$  distribution.

1. Write  $H_0$       we define the null hypothesis  $H_0: \beta = 0$
2. Design  $T$       we define the  $t$ -statistic  $T \equiv \frac{\hat{\beta} - \beta_0}{\text{SE}[\hat{\beta}]} = \frac{\hat{\beta} - 0}{\text{SE}[\hat{\beta}]}$ . If errors are normal,  $T \underset{H_0}{\sim} \mathcal{T}_{n-k}$ .<sup>17</sup>
3. Compute  $T_{\text{obs}}$        $T_{\text{obs}} \equiv T(\text{observed data})$
4. Interpret      we define the 2-tailed<sup>18</sup>  $p$ -value  $\equiv \text{P}(|T| \geq |T_{\text{obs}}| \mid H_0)$ , i.e., the probability of observing data as extreme as that actually observed, assuming  $H_0$ .<sup>19</sup>

$p$ -value small  $\iff T_{\text{obs}}$  falls in the tail of the Student's  $\mathcal{T}$ -distribution  
 $\iff$  observing our  $T_{\text{obs}}$  is highly unlikely under  $H_0$   
 $\implies$  reject  $H_0$   
 $\implies$  there is a relationship between  $y$  and  $x$ .

In econometrics, the standard approach is to dichotomize the evidence using a  $p$ -value threshold, usually the *significance level*  $\alpha = 5\%$ .  $\hat{\beta}$  is “statistically significant” iff  $p \leq 0.05$ , i.e., there is less than a 5% chance of observing the effect size that was observed if there was in fact no effect.

### 3.3.3 Type I/II errors, size and power

A test can lead to two types of mistakes:

- **Type 1 error** or *false positive*:  $\{- \mid H_0\}$  reject  $H_0$  when shouldn't... (*overconfident*)
- **Type 2 error** or *false negative*:  $\{+ \mid H_0\}$  don't reject  $H_0$  when should (*overcautious*)

We define a test's:

- **size**  $\alpha_T$  = probability of erroneously rejecting  $H_0$   $\equiv \text{P}(\text{type 1 error}) = \text{P}(- \mid H_0)$
- **power**  $\kappa_T$  = probability of correctly rejecting  $H_0$   $\equiv 1 - \text{P}(\text{type 2 error}) = \text{P}(- \mid H_0)$

Intuitively, we would like to minimize the size and maximize the power of our test. To guarantee  $\alpha_T \leq 0.05$ , we simply set the significance level  $\alpha = 0.05$ . To guarantee  $\kappa_T \geq 0.80$ , we need a sufficiently large sample size  $N$ , or the “Minimum Detectable Effect” will be very high.<sup>20</sup>

<sup>17</sup>This is a very strong assumption! And it means that if errors are far from normal, the result of the  $t$ -test has no interpretation...

<sup>18</sup>We can actually use the test statistic  $T$  to carry out two different tests:

- A two-tailed test: if we want to test for the possibility of the relationship in both directions.  $H_0: \beta_j = 0, H_a: \beta_j \neq 0$ . Both tails of  $T$ 's distribution constitute therefore the “critical region”, each containing  $\frac{\alpha}{2}$  of the values. By default, statistical packages report the two-tailed  $p$ -values.
- A one-tailed test: to test for the possibility of the relationship only in one direction. E.g.:  $H_0: \beta_j = 0, H_a: \beta_j > 0$ . Only one tail of  $T$ 's distribution makes the critical region, containing  $\alpha$  of the values. Only  $z$ - and  $t$ -tests can accommodate one-tailed tests.  $F$ -tests,  $\chi^2$ -tests... cannot as their distributions are not symmetric.

<sup>19</sup>⚠ The  $p$ -value is often misinterpreted to be the probability of the null hypothesis, whereas it is the probability of the data, given the null.  $p\text{-value} = \text{P}(\text{obs} \mid \text{hyp}) \neq \text{P}(\text{hyp} \mid \text{obs})$ . To quantify the probability of the null (which would arguably be more intuitive and interesting), we would need to turn to Bayesian inference.

<sup>20</sup>In hypothesis testing in econometrics, we typically want at least 80% power and a maximum size of 5%. I.e., we accept to incorrectly reject the null a maximum of 5% of the time, and to correctly reject it at least 80% of the time (i.e., 80% of studies conducted with a given sample size will correctly reject the null). 95% > 80%: econometrics is more focused on avoiding overconfidence than worried about being overcautious. Note also that having a high sample size  $N$  is not sufficient to have higher statistical power — empirical studies have actually found zero or weak correlations between the two. The power of a study depends indeed on the sample size, the true size of the effect, measurement variance, and the number of comparisons performed. Note finally that studies of small effects, although potentially important, are unlikely to be statistically significant

**Power calculations** Having adequate power means that if there really is an effect, the empirical strategy and data will enable the test to detect it. Low powered studies will instead “miss” the effect.<sup>21</sup> Post-estimation, it is useful to perform a retrospective design analysis and ask: “*Was my study sufficiently powered?*”, especially if we found a statistically significant non-null effect. But it must be done correctly:

△ To estimate the power one must first postulate a ‘true’ effect size, which can be thought of as that observed in an infinitely large sample. That effect size should be determined from a literature review, not the effect size observed in one’s study! The latter is noisy, and generally overestimated (publication bias), and would therefore lead to overestimates of power.

### 3.3.4 Criticisms of the NHST and ‘statistical significance’

The 2-way binary approach to statistical hypothesis testing, based on the NHST falsificationist paradigm (where the underlying truth is  $H_0$  “no effect” or  $H_a$  “effect”) and the measured outcome is a binary statement of ‘statistical significance’ from a p-value threshold, is heavily criticized. It is argued that:

- The underlying reality is not a simple Yes or No: in social sciences, the null hypothesis of zero effect (i.e., conditional independence of  $y$  and  $T$  given  $X$ ) is generally implausible — there are virtually no true zeros — and thus uninteresting. The null model is very false, so we are very likely to reject it with enough data.
  - ↪ Instead, we need to find alternatives to thinking in terms of conditional independence in order to study causality. The idea would be to estimate these dependences directly, rather than modeling the world in terms of conditional independence and estimating this structure through the testing of null effects.
- Interpreting p-values dichotomously loses a lot of information.
  - ↪ Instead, one could interpret p-values continuously, the strength of evidence for  $H_0$  being a continuous function of the p-value.
- Interpreting p-values dichotomously may induce selection bias: to be publishable, estimates must be ‘significant’, i.e., more than two standard errors away from 0; which selects for overestimates.

---

because they have insufficient power to detect the magnitudes of effects.

<sup>21</sup>Lacasse et al. (2020) is a good example of this. Rephrasing their specific independent and dependent variables as generic  $x$  and  $y$ : “*Because enrollment in the trial was stopped before we had reached our proposed sample size, the trial was underpowered, with the consequence of a wide confidence interval around the point estimate. [...] The data that were accrued could not rule out benefit or harm from  $x$ .*” As summarized in the abstract: “*Our underpowered trial provides no indication that  $x$  has a positive or negative effect on  $y$ .*”



## 4 Statistical inference [under a Bayesian approach]

As aforementioned, one of the main distinguishing features of Bayesian inference is the expression of all information, including uncertainty, using probability. From this paradigm stem new possibilities at every step of inference.

### 4.1 Steps of Bayesian inference

We start with the same situation: consider a population parameter of unknown true value  $\theta$ , and an *event*  $\theta = \tilde{\theta}$  (i.e.,  $\theta$  taking this specific value). We are interested in estimating  $\theta$  using data.

- (0) Definition of  $\theta$ : whereas in frequentist inference  $\theta$  is considered fixed, now it is considered a *random variable*. This means that at all times it has a probability distribution  $\mathcal{P}$  which represents our state of belief about its actual value.
- (1) Prior to observing data, this  $\mathcal{P}$  is  $\theta$ 's “prior distribution” and represents our prior belief about  $\theta$ .
- (2) Estimation: as more evidence (data) becomes available, we use Bayes' Theorem to update probability statements about  $\theta$ , which results in a *posterior* distribution:

$$\text{posterior density} \rightarrow f(\theta | \text{data}) = \frac{\text{prior density} \cdot f(\theta) \cdot \text{likelihood}^{22}}{\text{scaling factor or "evidence"} \cdot f(\text{data})}$$

In the canonical setup of a regression of  $y$  on  $x$ , the estimation step precisely consists of combining the model, data, and prior through Bayes' theorem, which is applied to each possible value of  $\theta$  to compute a posterior distribution of  $\theta$ :

- $y = \{y_1, \dots, y_n\}$  the observed sampled values of the outcome variable of interest  $Y$
- $\theta$  a parameter of  $y$ 's distribution
- $\alpha$  a hyperparameter of  $\theta$ 's distribution

$$\text{posterior} \rightarrow P(\tilde{\theta} | y, \alpha) = \frac{P(\tilde{\theta}, y | \alpha)}{P(y | \alpha)} = \frac{P(y | \tilde{\theta}, \alpha) P(\tilde{\theta} | \alpha)}{P(y | \alpha)} \propto P(y | \tilde{\theta}, \alpha) P(\tilde{\theta} | \alpha)$$

likelihood ← prior

- (3) Inference: we can summarize our updated belief about  $\theta$  from the posterior distribution, by reporting:
  - a measurement of central tendency (e.g., the mean or median);
  - a 95% credible interval, s.t. “after seeing the data, there is a 95% chance that this CI contains the true  $\theta$ .”

If the goal were prediction rather than inference, we would proceed similarly up to the last step; then:

- (3') Prediction: we propagate the uncertainty in  $\theta$  into the predictions of new data points  $\{\tilde{y}\}$  by using simulations: we repeatedly draw from the posterior a value of  $\theta$  and compute a new data point  $\tilde{y}$ , thereby creating a predictive distribution.

<sup>22</sup>Note that the likelihood is now  $f(D|\theta)$  instead of  $f(D)$ , as  $\theta$  is no longer fixed. Or, equivalently, the (conditional) probability model is  $f(y|\theta, x)$  instead of  $f(y|x)$ .

## Advantages/Distinctions of Bayesian inference w.r.t. frequentist inference

- *Intuitive interpretation of findings*

Because uncertainty is encoded probabilistically, our uncertainty about  $\theta$  after observing the data is represented by a distribution of values. We can effortlessly compute a 95% credible interval from the posterior, with an intuitive interpretation: “There is a 95% chance that this CI contains the true  $\theta$ .” In comparison, the frequentist 95% confidence interval refers to our uncertainty about the *sampling method* — not  $\theta$  — and is thereby interpreted in terms of repeated samples: “Imagine we drew all possible random samples of size  $n$ . This interval would contain the true  $\theta$  in 95% of the samples.” More generally, the Bayesian framework enables us to actually answer questions like “What is the probability that  $\theta = 4$ ?” Whereas the Frequentist framework produces convoluted estimates: the probability of the data assuming that  $\theta = 4$ .

- *Including prior information*

In Bayesian inference emerges a compromise between prior information and data. More generally, it is a way to include multiple sources of information.

- *Making predictions is facilitated by the computation being ~~optimization~~-simulation-based*

A Bayesian will argue that one should do predictions based on the whole posterior distribution of possible coefficient values, while prediction based on point estimates disregards all information about how imprecise the point estimate is. Because the uncertainty about  $\theta$  is encoded in a probability distribution, we can simply propagate this uncertainty into predictions of a new data point  $\tilde{y}$ , by using simulations. We draw a value from  $\theta$ 's posterior distribution and make a *probabilistic* prediction of a new data point  $\tilde{y}$  for this value of  $\theta$ , and repeat this simulation  $S$  times. The  $S$  resulting values make the *posterior predictive distribution*  $f(\tilde{y}|Y, \alpha) = \int_S f(\tilde{y}|\theta) f(\theta|Y, \alpha) d\theta$ .

Bayesian inference is the discipline of updating our belief about the world based on further observation of the world. Whereas frequentist inference is focused on summarizing the information in the data. These summaries of data have known statistical properties but have limited value as predictions.

## 4.2 Choosing $\theta$ 's prior distribution

We include additional information using a prior distribution

- Using an uninformative or “flat” prior (the uniform distribution) results in the posterior distribution being equal to the product of the likelihood and a mere constant, s.t. the mode of the posterior distribution is the ML or LS estimator.
- Weakly Informative Priors: “What you should be doing when you think you want to use noninformative priors.” [https://statmodeling.stat.columbia.edu/2009/05/24/handy\\_statistic/](https://statmodeling.stat.columbia.edu/2009/05/24/handy_statistic/). Ex: The R function `rstanarm::stan_glm()` adjusts the default priors based on the scale of the variables in the model.
- “Conjugate” prior probability distributions (for the ... distribution): the posterior distributions  $f(\theta|x)$  are in the same family as the prior probability distribution  $f(\theta)$ .
- Bayesian inference is a compromise between prior and data, where each has a weight proportional to the inverse square of its s.e.  $\rightarrow \text{SE}_{\text{Bayes}} < \text{both } \text{SE}_{\text{prior}} \text{ and } \text{SE}_{\text{data}}$

## 4.3 Estimating $\theta$ 's posterior distribution

- When the likelihood has an analytical expression, we can combine it with the prior to derive the posterior analytically.

- Most of the time, there is no such analytical expression. To estimate the posterior distribution, we can use Markov Chain Monte-Carlo (MCMC) algorithms: a family of iterative sampling algorithms<sup>23</sup> that yield approximate draws from the posterior distribution:
  - “*Monte-Carlo*” refers to the practice of estimating the properties of a distribution by examining random samples from the distribution. Ex: instead of finding the mean of a normal distribution by directly calculating it from the distribution’s equations, we would draw random samples from the normal distribution and calculate the sample mean.
  - “*chain*” means that the random samples are generated by a special sequential process: each random sample is used as a stepping stone to generate the next random sample. Note that this means that the draws are not independent.
  - The “*Markov*” property of the chain is that, while each new sample depends on the one before it, new samples do not depend on any samples before the previous one.
- Most of the time, **fitting a Bayesian model = generating a set of posterior simulations** (representing different possible values of the parameter vector  $\theta$ ), which we typically summarize using its median, its median absolute deviation (a more robust estimator of scale than the standard deviation), and uncertainty intervals.

---

<sup>23</sup>For example, Stan uses as inference algorithms two MCMC algorithms: the Hamiltonian Monte Carlo algorithm and its adaptive variant the “no-U-turn sampler”.

## 5 Prediction

Prediction isn't part of statistical inference, but it can be the ultimate research goal, motivating the initial statistical inference step. Whether the ultimate goal is inference or prediction,<sup>24</sup> both first require finding a model that describes the relationship between the independent variables and the outcome in our data. The use of the resulting model then differs:

- Inference: Use the model to learn about the data generating process.
- Prediction: Use the model to predict the outcomes for new data points.

---

<sup>24</sup>Note: In machine learning, the term inference is sometimes used instead to mean “making a prediction, by evaluating an already trained model”. In this context, inferring properties of the model is referred to as training or learning (rather than inference), and using a model for prediction is referred to as inference (instead of prediction).

## 6 Model comparison

Learning from data has generally one of two ultimate objectives: inference or prediction. Model comparison should proceed in line with the objective. After a brief paragraph on *nested* model discrimination, this section focuses on model comparison for prediction, our objective will therefore be predictive performance.<sup>25</sup> Much of this section is taken from [Gelman et al. \(2014\)](#).

### 6.1 Comparing nested models: $F$ tests

If two models are *nested*, i.e., one represents a special case of the other, we can easily discriminate between them using a standard hypothesis test of the parametric restrictions on the nested one.

The key questions are: (1) is the improvement in fit large enough to justify the additional difficulty in fitting, and in a Bayesian context (2) is the prior distribution on the additional parameters reasonable?

### 6.2 Comparing non-nested models: IC, CV

We want to know which model gives the best predictions of new data generated from the true DGP. Ideally, we would measure the model's out-of-sample predictive accuracy or error, for such new data produced from the true DGP. After describing exactly what the quantity we would like to measure is, we will describe methods for estimating an *approximation* of it, given the data we have.

There are different ways of defining a model's predictive accuracy or error:

- If one is predicting a *point*, predictive accuracy can be defined using an error measure, such as the absolute error or the squared error. Individual errors are aggregated and averaged to obtain a summary measure of predictive accuracy, such as the Mean Absolute Error (MAE) or the Root Mean Squared Error (RMSE):<sup>26</sup>

$$MAE \equiv \frac{1}{N} \sum_i |\hat{y}_i - y_i|, \quad RMSE = \sqrt{MSE} \equiv \sqrt{\frac{1}{N} \sum_i (\hat{y}_i - y_i)^2}$$

- A more general<sup>27</sup> summary is the *log likelihood* or *log predictive density* (LPD). For any data  $y=y_1, \dots, y_m$  produced from the true DGP, i.e., taken from the *unknown* data distribution  $f$ ,  $LPD(y) \equiv \ln P(y|\theta) = \ln \prod_i P(y_i|\theta)$ .

Therefore for *out-of-sample* data:

If inference for $\theta$ is summarized by a point estimate $\hat{\theta}(y)$	If inference for $\theta$ is summarized by a posterior distribution $p_{post,\theta}()$
<p>▷ For a new data point <math>\tilde{y}_i \sim f</math>:</p> <p><math>LPD(\tilde{y}_i) = \ln P(\tilde{y}_i \hat{\theta})</math></p>	<p><math>LPD(\tilde{y}_i) = \ln p_{post,y}(\tilde{y}_i) \equiv \int P(\tilde{y}_i \theta) p_{post,\theta}(\theta) d\theta</math></p>
<p>▷ As new data points are themselves unknown, the expectation:</p> <p><math>ELPD \equiv \mathbb{E}_f[LPD(\tilde{y}_i)] = \mathbb{E}_f[\ln P(\tilde{y}_i \hat{\theta})]</math></p>	<p><math>ELPD \equiv \mathbb{E}_f[LPD(\tilde{y}_i)] = \mathbb{E}_f[\ln p_{post,y}(\tilde{y}_i)]</math></p>

<sup>25</sup>In classical econometrics focused on inference, especially when the goal is causal inference, the research design drives the model specification such that there isn't so much need for model comparison and selection.

<sup>26</sup>The RMSE is the standard deviation of the residuals, i.e., of the unexplained variation. It is an absolute measure of fit of the model to the data. (Whereas  $R^2$  is a relative measure of fit. Note that one should absolutely not select a model based on  $R^2$ , as this would favor overfitting.) Note that: (1) RMSE is scale-dependent (it has the same unit as  $y$ ), therefore it can only be compared across models in the same units; (2) compared to the MAE, the RMSE penalizes large errors more.

<sup>27</sup>It is proportional to the MSE if the model is normal.

In practice,  $f$  and  $\theta$  are unknown, so we can't compute ELPD. We will try to approximate it, using existing data (hence knowing that any method will be correct at best only in expectation...).

- **Adjusted within-sample predictive accuracy:** a natural estimate of the expected log predictive density for *new* data is the log predictive density for *existing* data. **Information criteria** such as AIC and WAIC give approximately unbiased estimates of ELPD by correcting for how much the fitting of  $k$  parameters increases predictive accuracy, by chance alone. These are scoring methods from information theory.
- **Cross-validation:** the model is fit to a training set, then the fit evaluated on a holdout set.

Both methods are based on adjusting the log predictive density of the observed data by subtracting an approximate bias correction. The measures differ in their starting points (how they measure the log predictive density) and their adjustments. Asymptotically, AIC is equal to LOO-CV computed using ML estimation, and Bayesian LOO-CV is equal to WAIC.

### Information Criteria (IC)

Goal: we want the best model fit (maximized likelihood), but we penalize model complexity (to not overfit the data). Most IC are expressed on the deviance scale; the model with smallest IC is preferred.

Let  $k$  be the number of parameters,  $n$  the sample size.

- **Akaike information criterion (AIC)**

- starting point: the log predictive density, conditional on a point estimate:  $\ln \hat{\mathcal{L}} \equiv \ln P(y|\hat{\theta})$ ;
- adjustment for overfitting: uses the simplest bias correction, based on the asymptotic normal posterior distribution, for which<sup>28</sup> simply subtracting  $k$  corrects for the number of parameters:

$$AIC \equiv -2 (\widehat{\text{ELPD}}_{\text{AIC}}) = -2 (\ln \hat{\mathcal{L}} - k) = -2 \ln \hat{\mathcal{L}} + 2k$$

$AIC_c$  is the AIC corrected for small samples:  $AIC_c = -2 \ln \hat{\mathcal{L}} + 2k \frac{n}{n-k-1} \xrightarrow{n \rightarrow +\infty} AIC$

*Limit: when we go beyond linear models with flat priors, e.g., models with hierarchical structures or informative priors, the number of effective parameters isn't  $k$  so we can't simply subtract  $k$ .*

- **Watanabe-Akaike information criterion (WAIC)**

- starting point: the log predictive density, averaging over the posterior distribution  $p_{\text{post}}(\theta) = P(\theta|y)$  (i.e., a fully Bayesian approach);
- adjustment for overfitting: corrects for the *effective* number of parameters.

### Cross-validation (CV)

Cross-validation consists in partitioning the data into a training set  $y_t$  and a validation set  $y_v$ , fitting the model to the training set, and evaluating this predictive accuracy (fit) using the validation set. It is based on the log predictive density, but can use any starting point (i.e., either averaging over the posterior distribution  $p_{\text{post}}(\theta)$  or conditioning on a point estimate  $\hat{\theta}$ ).

In Bayesian CV, fitting the model to  $y_t$  yields a posterior distribution for  $\theta$ :  $p_{\text{post}}(\theta) \equiv P(\theta|y_t)$ . We assume we can summarize it by  $S$  simulation draws  $\theta^1, \dots, \theta^S$ . We can then compute the log predictive density for  $y_v$  as:  $\text{LPD}(y_v) \equiv \ln P(y_v|\theta^{\text{post}}) \equiv \frac{1}{S} \sum_{s=1}^S \ln P(y_v|\theta^s)$

The CV process is repeated using different partitions, and the resulting log predictive densities are averaged into a single estimate of out-of-sample predictive accuracy.

---

<sup>28</sup>This is also true in the special case of a normal linear model with a uniform prior distribution.

- **K-fold CV**

The data are randomly partitioned into  $K$  equal-sized sets.  $K = 10$  is commonly used. The CV process is repeated  $K$  times, each time using one subsample for validation — such that each observation is used for validation exactly once — and the  $K$  results are averaged into one estimate:

$$\text{LPD}_{K\text{-CV}} = \sum_{k=1}^K \ln \left( \frac{1}{S} \sum_{s=1}^S P(y_k | \theta^s) \right)$$

- **‘Leave-one-out’ CV = n-fold CV**

In the extreme case of  $n$  partitions, each validation set represents a single data point:

$$\text{LPD}_{\text{LOO-CV}} = \sum_{i=1}^n \ln \left( \frac{1}{S} \sum_{s=1}^S P(y_i | \theta^s) \right)$$

In any CV process, each prediction is conditioned on  $n - v$  data points instead of  $n$ , which causes underestimation of the predictive fit. We can correct for this bias by estimating how much better predictions would be obtained if conditioning on  $n$  data points (Gelman et al., 2014).

**Conclusion** Neither cross-validation nor information criteria are perfect. AIC does not work in settings with strong prior information, WAIC relies on a data partition unamenable to structured models such as for spatial or network data, cross-validation is computationally expensive as getting a stable estimate requires many data partitions and fits. Gelman et al. (2014)’s preferred choice is “cross-validation, with WAIC as a fast and computationally convenient alternative. WAIC is fully Bayesian (using the posterior distribution rather than a point estimate) [...]. A useful goal of future research would be a bridge between WAIC and cross-validation with much of the speed of the former and robustness of the latter.”

**TO ADD: Model Shrinkage Methods, and other methods to deal with highly correlated predictors**

- LASSO (Least Absolute Shrinkage and Selection Operator)
- PCA

## 7 Other branches of statistical modeling

### 7.1 Statistical Inference Using Agent-based models (ABMs)

**Agent-Based Models** are computational models<sup>a</sup> that simulate the actions and interactions of autonomous agents within a system, to assess their effects on the system as a whole. The goal is to re-create and predict the emergence<sup>b</sup> of higher-level system properties from simple agent-level behaviors, taking a “bottom-up” approach.

ABMs are generally composed of 3 elements:

1. many **agents** with assigned attributes;
2. simple **rules** about: their individual decision-making process, how they interact, how they learn and adapt—these rules can be deterministic or probabilistic;
3. an **environment**.

<sup>a</sup>Computational models are mathematical models that study the behavior of a system by computer simulation. The system studied is often a complex nonlinear system for which simple analytical solutions are not available. Experimentation is therefore done by modifying the model’s parameters, and comparing outcomes. Examples include weather forecasting models, flight simulator models, neural network models, and ABMs.

<sup>b</sup>The process of *emergence* can be expressed as “the whole is greater than the sum of its parts”.

**Goal of ABMs** ABMs allow us to observe how the behaviors of individual agents affect the system as a whole and if any emergent structure develops within the system. They show how small-scale changes can affect large-scale outcomes within the system.

At a formal level, an ABM is just a statistical model. But agent-based modeling differs from other types of statistical modeling because it describes only the behavior of the agents in a system, rather than global properties of the system.

#### Use in different fields

- In economics: ABMs can describe the microeconomic actions of adaptive agents, which give rise to emergent behavior in the form of macroeconomic structures; which, in turn, influence agent decisions. Ex: we can represent the economy as a complex system, with crashes and booms that emerge from non-linear responses to small changes.
- In ecology: ABMs are often called individual-based models (IBMs), and are used to study population dynamics, plant-animal interactions...
- In epidemiology: epidemiological ABMs now complement traditional compartmental models (such as the deterministic SIR — Susceptible/Infectious/Recovered — model) which they have tended to surpass in terms of prediction accuracy to model the spread of epidemics.

#### Statistical inference

1. Model validation and selection, uncertainty quantification, and fitting ABMs to data: **There does not seem to be (yet) formal guidelines and procedures from the statistical literature, for: fitting ABMs to data, for making quantified statements of uncertainty about the outputs, e.g., calculating confidence intervals on predictions, nor for testing whether a specific parameter (rule) is needed in an ABM. See Banks and Hooten (2021); Heard et al. (2015).**
2. Statistical inference  
Because of the variety of input rules and the complexity of outputs, the likelihood function of an ABM



is generally intractable. One must hence perform likelihood-free inference. [Heard et al. \(2015\)](#) suggest that two main tools allow that: emulators and approximate Bayesian computation (ABC).

## Key ideas

Statistics is about reasoning under uncertainty, and therefore **probability distributions**. Inferential statistics proceeds by learning from data, it asks: *given sample data, what are we able to infer about the population?*

In microeconometrics, inference is usually conducted under a frequentist approach:

Steps	Options
1. Choose & write a model, the one we think is closest to the true and unobserved DGP.	<i>(linear regression model w. normal errors, logistic regression model, SEM...)</i>
★ <b>Bring in data</b> ★	
2. Estimate the model, i.e., estimate the conditional distribution. When the specification is parametric, it means estimating parameters. a. estimation $\implies$ " $\hat{\beta} = \dots$ " b. hypothesis testing $\implies$ " $\hat{\beta}$ is/isn't statistically significant"	<i>(OLS, 2SLS, ML,...)</i>
3. Validate & compare the model.	

In frequentist statistics, we trust that the results given by these statistical tools (estimators, tests...) give us relevant indications about the population, because of the tools' asymptotic properties (which stem from laws of large numbers (LLNs) and central limit theorems (CLTs)).

## References

- Bafumi, J. and Gelman, A. (2007). Fitting multilevel models when predictors and group effects correlate. *SSRN Electronic Journal*, ISSN: 1556-5068, DOI: [10.2139/ssrn.1010095](https://doi.org/10.2139/ssrn.1010095).
- Banks, D. L. and Hooten, M. B. (2021). Statistical Challenges in Agent-Based Modeling. *Am. Stat.*, pages 1–8, DOI: [10.1080/00031305.2021.1900914](https://doi.org/10.1080/00031305.2021.1900914).
- Bell, A., Fairbrother, M., and Jones, K. (2019). Fixed and random effects models: making an informed choice. *Quality & quantity*, 53(2):1051–1074, ISSN: 0033-5177, DOI: [10.1007/s11135-018-0802-x](https://doi.org/10.1007/s11135-018-0802-x).
- Gelman, A. and Greenland, S. (2019). Are confidence intervals better termed “uncertainty intervals”? *BMJ*, 366, DOI: [10.1136/bmj.15381](https://doi.org/10.1136/bmj.15381).
- Gelman, A. and Hill, J. (2006). *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge University Press, ISBN: [9781139460934](https://doi.org/10.1017/9781139460934).
- Gelman, A., Hwang, J., and Vehtari, A. (2014). Understanding predictive information criteria for Bayesian models. *Stat. Comput.*, 24(6):997–1016, DOI: [10.1007/s11222-013-9416-2](https://doi.org/10.1007/s11222-013-9416-2).
- Heard, D., Dent, G., Schifeling, T., and Banks, D. (2015). Agent-based models and microsimulation. *Annual Review of Statistics and Its Application*, 2(1):259–272, DOI: [10.1146/annurev-statistics-010814-020218](https://doi.org/10.1146/annurev-statistics-010814-020218).
- Horrace, W. C. and Oaxaca, R. L. (2006). Results on the bias and inconsistency of ordinary least squares for the linear probability model. *Econ. Lett.*, 90(3):321–327, DOI: [10.1016/j.econlet.2005.08.024](https://doi.org/10.1016/j.econlet.2005.08.024).
- Lacasse, Y., Sériès, F., Corbeil, F., Baltzan, M., Paradis, B., Simão, P., Abad Fernández, A., Esteban, C., Guimarães, M., Bourbeau, J., Aaron, S. D., Bernard, S., and Maltais, F. (2020). Randomized trial of nocturnal oxygen in chronic obstructive pulmonary disease. *New England Journal of Medicine*, 383(12):1129–1138, DOI: [10.1056/NEJMoa2013219](https://doi.org/10.1056/NEJMoa2013219).
- Mertens, W., Pugliese, A., and Recker, J. (2017). Analyzing Longitudinal and Panel Data. In *Quantitative Data Analysis: A Companion for Accounting and Information Systems Research*, pages 73–98. Springer International Publishing, Cham, ISBN: [978-3-319-42700-3](https://doi.org/10.1007/978-3-319-42700-3).
- Wooldridge, J. M. (2013). *Introductory Econometrics: A Modern Approach*. Cengage Learning, fifth edition, ISBN: [9781111531041](https://doi.org/10.1111/9781111531041).

## A A small library of regression models

The textbook Classical Linear Regression Model (CLRM) can be generalized in various dimensions, such as:

- the power of the independent variables ( $\rightarrow$  polynomial regression);
- the link function relating the linear predictor  $X\beta$  to the outcome  $\mathbb{E}[y|X]$  ( $\rightarrow$  generalized linear model);
- the number of levels in the data ( $\rightarrow$  multilevel model)...

This section presents some of the models resulting from these generalizations.

**Notation** Recall that a statistical model is the combination of a sample space and a collection of *joint probability distributions* on that space; the goal being to represent the specific distribution induced by the DGP. Rather than look at the full joint distribution, regression models simplify the problem and focus on the *conditional distribution* of  $y|X$ .<sup>29</sup> All regression models are therefore first *conditional distributions*, and can be written as such. Based on the properties of each distribution, we can then also write them in a *conditional mean*  $+/ \times$  *error* form.

### A.1 Expanding from the CLRM

- Classical Linear Regression Model

$$\begin{aligned} y|X &\sim \mathcal{F}(X\beta, \sigma^2\mathbf{I}) \\ \iff y &= \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + e, \quad e \stackrel{\text{iid}}{\sim} \mathcal{F}(0, \sigma^2\mathbf{I}) \\ \iff y &= \mathbb{E}[y|X] + e, \quad \mathbb{E}[y|X] = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k, \quad e \stackrel{\text{iid}}{\sim} \mathcal{F}(0, \sigma^2\mathbf{I}) \end{aligned}$$

- Polynomial Regression

- Ex: LOESS (locally estimated scatterplot smoothing) is a nonparametric regression algorithm, in which  $\mathbb{E}[y|X]$  at each data point  $i$  is estimated using a weighted low-degree polynomial regression model that gives higher weights to the neighboring points (along  $X$ ).

$$\mathbb{E}[y|X] = \beta_0 + \beta_1 X_1 + \beta_2 X^2 + \dots + \beta_p X^p, \quad e \stackrel{\text{iid}}{\sim} \mathcal{F}(0, \sigma^2\mathbf{I})$$

- Generalized linear model (GLM)

GLMs are often used to predict outcomes of bounded or discrete form (outcomes that cannot be fit well with normally distributed additive errors). A GLM consists of three elements: a probability distribution we assume the outcome to be generated from  $\mathcal{F}()$ , a linear predictor  $X\beta$ , and an invertible link function  $g()$  that relates  $\mathbb{E}[y|X]$  to  $X\beta$ .

- Ex: the linear regression model is a GLM with normal data and identity link
- Ex: the logistic regression model is a GLM with Bernoulli data and logit link
- Ex: the Poisson regression model is a GLM with Poisson data and log link

$$y|X \sim \mathcal{F}_{\text{ExpFamily}}(\dots), \quad g(\mathbb{E}[y|X]) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

- Generalized additive model (GAM)

GAMs generalize further to allow for  $g(\mathbb{E}[y|X])$  to be a *nonlinear* smooth function of each predictor. The space of functions of which  $h$  is an element is the “basis”.

$$y|X \sim \mathcal{F}_{\text{ExpFamily}}(\dots), \quad g(\mathbb{E}[y|X]) = \beta_0 + h_1(X_1) + \dots + h_k(X_k)$$

---

<sup>29</sup>For simplicity, we assume we adopt a frequentist approach, therefore we need not write distributions of  $y$  as conditional on  $\theta$ , as  $\theta$  is fixed. If we adopted a Bayesian approach, we’d make it explicit that distributions of  $y$  are conditional on  $\theta$ .

GAMs penalize the complexity of the model to prevent overfitting the data, by adding a penalty for the size of the coefficients associated with the basis functions.

- Multilevel or “hierarchical” models

The lowest-level model is a regression, higher-level models model coefficients of the model immediately below them. These higher-level models can be regressions or distributions. All models are fitted simultaneously.

- Ex: 2-level, varying-intercept model; the group-level model is a regression:

$$\begin{cases} y_i \sim \mathcal{F}(\alpha_{j[i]} + \beta X_i, \sigma_y^2) & \forall i = 1, \dots, N, \quad j = 1, \dots, J \\ \alpha_j \sim \mathcal{F}(\gamma_0 + \gamma_1 W_j, \sigma_\alpha^2) & \forall j = 1, \dots, J \end{cases}$$

- Ex: 2-level, varying-intercept & slope model; the group-level models are distributions:

$$\begin{cases} y_i \sim \mathcal{F}(\alpha_{j[i]} + \beta_{j[i]} X_i, \sigma_y^2) & \forall i = 1, \dots, N, \quad j = 1, \dots, J \\ \alpha_j \sim \mathcal{N}(\mu_\alpha, \sigma_\alpha^2) & \forall j = 1, \dots, J \\ \beta_j \sim \Gamma(\gamma_\beta, \delta_\beta) & \forall j = 1, \dots, J \end{cases}$$

- Nonparametric models

Use large numbers of parameters to allow essentially arbitrary curves for the predicted value of  $y|X$ .

Ex: splines

- Incomplete data models

- Missing data. For some problems, we can set up a model specifically to handle the missingness mechanism. Ex censored data: extensions of ML / Bayesian regression include the censoring into the likelihood.

- Measurement error in the predictors  $x$ :<sup>30</sup> we observe  $x^* = x + \eta$ . If we can estimate the variance of the measurement errors, we can either just apply a bias correction on the raw estimate from the regression of  $y$  on  $x^*$ , or directly fit the full “simultaneous-equation model” using a marginal likelihood or Bayesian approach. Same maths as in IV.

## A.2 Limited outcome models

A *limited* dependent variable  $y$ , i.e., a  $y$  that is categorical or constrained to fall in a certain range, often arises in econometrics. With such data, linear regression is not an appropriate estimation method, as it does not take into account the constraint on possible values of the dependent variable. The strategy is to transform the limited  $y$  into a continuous, real-valued variable  $y' \in (-\infty, \infty)$ , that we can then model as  $y' = X\beta + \varepsilon$ , using a link function  $g(y)$ .

Limited $y$	Appropriate regression models
binary: $y \in \{0, 1\}$	probit regression, logit regression
count: $y \in \{0, 1, 2, 3, \dots\}$	Poisson regression, negative binomial regression
interval: $y \in [0, 1]$	<b>fractional response</b>
censored	censored regression, e.g., Tobit

<sup>30</sup>Measurement error in  $y$  does not pose a problem besides imprecision, as it just goes into the error term. It is measurement error in  $x$  that poses a problem: estimated regression coefficients can be attenuated (i.e., it doesn’t just increase standard errors, but can drive the coefficient down).

### A.2.1 Binary outcome models

The outcome variable  $y|X$  is binary, i.e., it follows a Bernoulli distribution:

$$y|X \sim \text{Ber}(\pi) \equiv \begin{cases} 1 & \text{with probability } \pi \\ 0 & \text{with probability } 1 - \pi \end{cases}$$

The conditional mean  $\mathbb{E}[y|X]$  is equal to the conditional probability  $\pi \equiv P(y=1|X)$ .<sup>31</sup> A regression model is therefore formed by expressing  $\pi$  as a function of  $X$  and  $\beta$ ;<sup>32</sup> and we look for a link function  $g()$  that maps the  $[0,1]$  interval to the real line.

$$y_i|X_i \sim \text{Ber}(\pi_i), \quad \pi_i = \mathbb{E}[y_i|X_i] = g^{-1}(X_i, \beta)$$

### Models

- **Linear probability model**

$$y_i|X_i \sim \mathcal{F}(\pi_i), \quad \pi_i = X_i' \beta$$

This model is probably the first one that comes to mind. It is not appropriate, as the identity link is not a CDF, it will not constrain the predicted values to be in  $[0,1]$ , since the predictor  $X_i' \beta$  can take any real value. Yet, it is still frequently preferred to Logit or Probit, on grounds that it is computationally simpler, the estimated marginal effects are easier to interpret, and are usually very similar anyway, especially with a large sample size.

However, [Horrace and Oaxaca \(2006\)](#) show that in almost all circumstances, the LPM yields biased and, most importantly, *inconsistent* estimates. I.e., the LPM gives the wrong answer, with almost certainty, even with an infinitely large sample: “consistency seems to be an exceedingly rare occurrence as one would have to accept extraordinary restrictions on the joint distribution of the regressors. Therefore, OLS is frequently a biased estimator and almost always an inconsistent estimator of the LPM.”

- **Logit model = Logistic regression model**

$$y_i|X_i \sim \text{Ber}(\pi_i)$$

$$\pi_i = \text{logit}^{-1}(X_i' \beta) \equiv \frac{e^{X_i' \beta}}{1 + e^{X_i' \beta}} \iff \text{logit}(\pi_i) = X_i' \beta$$

We choose as link function  $g()$  the logit  $\text{logit}(\cdot) \equiv \ln\left(\frac{\cdot}{1-\cdot}\right)$  (i.e., we choose as  $g^{-1}()$  the CDF of the logistic distribution:  $\text{logit}^{-1}()$ ), which maps  $[0,1]$  to  $[-\infty, \infty]$ . We transform the probability outcome using this logit or “log-odds” transformation. As this new outcome need not be in  $[0,1]$ , we can model it as a *linear* function of the covariates.

Interpretation of each coefficient  $\hat{\beta}_k$  (keeping all the other predictors fixed):

- logit scale  $[-\infty, \infty]$       “a 1-unit difference in  $x$  corresponds to a  $\hat{\beta}_k$ -unit difference in  $\text{log-odds}(y=1)$ ”
- odds<sup>33</sup> scale  $[0, \infty]$       “a 1-unit difference in  $x$  corresponds to a  $e^{\hat{\beta}_k}$  multiplicative difference in  $\text{odds}(y=1)$ ”
- probability scale  $[0,1]$       “a 1-unit difference in  $x$  corresponds to a  $\frac{\hat{\beta}_k}{4}$ -unit maximum<sup>34</sup> difference in  $P(y=1)$ ”

<sup>31</sup> As  $\mathbb{E}[y|X] = 1 \times P(y=1|X) + 0 \times P(y=0|X) = P(y=1|X)$ .

<sup>32</sup> The function  $g^{-1}()$  should be a *cumulative distribution function*, to ensure that  $0 \leq \pi_i \leq 1$ .

<sup>33</sup> The odds of success are defined as the ratio of the probability of success  $\pi$  over the probability of failure. Here, where “success” is  $y=1$ , the odds of  $y=1$  are  $\frac{\pi}{1-\pi}$  to 1.

<sup>34</sup>  $\Delta$  The logistic function  $\text{logit}^{-1}()$  is nonlinear, so the expected difference in  $P(y=1)$  from a given difference in  $x$  is not a constant along  $x$ . We must choose where to evaluate changes, if we want to interpret them on the probability scale. The slope of the logistic regression curve is steepest at its halfway point ( $\text{logit}^{-1}() = 0.5$ ) and is  $\beta/4$ . I.e., the largest change in  $\pi$  from a 1-unit change in  $x$  is  $\beta/4$ .

- **Probit model**

$$y_i|X_i \sim \text{Ber}(\pi_i)$$

$$\pi_i = \text{probit}^{-1}(X_i'\beta) \equiv \int_{-\infty}^{X_i'\beta} \phi(t)dt \iff \text{probit}(\pi_i) = X_i'\beta$$

We choose as link function  $g()$  the probit, which is the quantile function of the standard normal distribution (i.e., we choose as  $g^{-1}()$  the CDF of the normal distribution:  $\text{probit}^{-1}()$ ), which maps  $[0, 1]$  to  $[-\infty, \infty]$ .

We cannot interpret each coefficient  $\hat{\beta}_k$  directly, we need to *compute* the marginal effects.

*Note: As a rule of thumb, probit regression coefficients are roughly equal to logistic regression coefficients divided by 1.6.*

**Estimation** by Maximum Likelihood, as the distribution of the data  $y|X$  must be the Bernoulli. The conditional density of each observation is:  $f(y_i|X_i) = \pi_i^{y_i}(1 - \pi_i)^{(1-y_i)}$ . Given independence over  $i$ , the (log-)likelihood of the data is then the (log-)likelihood for  $n$  independent Bernoulli observations:

$$\begin{aligned} \hat{\theta}_{\text{ML}} &= \underset{\theta}{\text{argmax}} \log \mathcal{L}(y|X, \theta) = \underset{\theta}{\text{argmax}} \log \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i} \\ &= \underset{\theta}{\text{argmax}} \sum_{i=1}^n y_i \ln(\pi_i) + (1 - y_i) \ln(1 - \pi_i) \\ &= \underset{\beta}{\text{argmax}} \sum_{i=1}^n y_i \ln(g^{-1}(X_i'\beta)) + (1 - y_i) \ln(1 - g^{-1}(X_i'\beta)) \end{aligned}$$

$\triangle$  *Don't fit logistic models for binary outcomes when the underlying continuous variable is available. For inference or prediction, it is much more efficient to model the underlying continuous variable and then map it back to the probability of the discrete outcome. Ex:*

- *basketball game: model the expected score differential, and then map it to  $P(\text{winning})$ .*
- *elections: predict vote differential and then map that to  $P(\text{winning})$ .*
- *health: model change in blood pressure, and then convert it to the binary disease state  $P(\text{hypertension})$ .*

## A.2.2 Count data models

$y_i \in \{0, 1, 2, \dots\}$ : number of occurrences of an event. *Ex: number of children in a household, number of doctor visits per year.*

- **Poisson regression model**

The Poisson distribution  $\text{Pois}(\lambda)$  models the number of events occurring in a fixed interval (of time or space), when these events occur at random, independently in time, with the constant mean rate  $\lambda$ . Its probability mass function is therefore  $P(y|\lambda) = \frac{e^{-\lambda}\lambda^y}{y!}$ , which further implies  $\mathbb{E}[y] = \mathbb{V}[y] = \lambda$ . The general Poisson regression model is:

$$y_i|X_i \sim \text{Pois}(\lambda_i), \quad \lambda_i = g^{-1}(X_i'\beta)$$

A common choice of link function  $g()$  is  $\ln()$ . The Poisson regression model is therefore fitted as a log-linear regression with Poisson error distribution:  $y_i|X_i \sim \text{Pois}(e^{X_i'\beta})$ .

Estimation by Maximum Likelihood:

$$\begin{aligned}
\hat{\beta}_{\text{ML}} &= \operatorname{argmax}_{\beta} \log \mathcal{L}(y|X, \beta) = \operatorname{argmax}_{\beta} \log \prod_{i=1}^n P(y_i|X_i, \beta) \\
&= \operatorname{argmax}_{\beta} \sum_{i=1}^n \log \frac{e^{-e^{X_i' \beta}} (e^{X_i' \beta})^{y_i}}{y_i!} \\
&= \operatorname{argmax}_{\beta} \sum_{i=1}^n \left[ -e^{X_i' \beta} + y_i (X_i' \beta) - \ln(y_i!) \right]
\end{aligned}$$

A limitation of the Poisson model is that it implies equi-dispersion:  $\mathbb{V}[y_i|x_i] = \mathbb{E}[y_i|x_i]$ , whereas we often see overdispersion in the data (ex: a few traders will do many trades, many traders will do a few). To accomodate overdispersion, some softwares (e.g., R) have packages that permit an “adjusted” Poisson regression, or we can turn to the negative binomial distribution.

- **Negative binomial model**

The negative binomial distribution  $\text{NB}(p, r)$  models the number of successes in a sequence of iid Bernoulli( $p$ ) trials before  $r$  failures occur. Its probability mass function is therefore  $P(y|p, r) = \binom{y+r-1}{y} p^y (1-p)^r = \frac{\Gamma(y+r)}{y! \Gamma(r)} p^y (1-p)^r$ .<sup>35</sup> The NB distribution converges to Poisson as  $r \rightarrow \infty$ , but has larger variance (overdispersion) for small  $r$ , s.t.  $r$  is called the “reciprocal dispersion” parameter.

Using as link function  $g()$  the usual logarithmic transformation  $\ln()$ , the NB regression model is:

$$y_i|X_i \sim \text{NB}(\mu_i, r), \quad \mu_i = e^{X_i' \beta}$$

Estimation by Maximum Likelihood:

$$\begin{aligned}
\hat{\theta}_{\text{ML}} &\equiv \operatorname{argmax}_{\theta} \log \mathcal{L}(y|X, \theta) = \operatorname{argmax}_{\theta} \log \prod_{i=1}^n P(y_i|X_i, \theta) \\
&= \operatorname{argmax}_{\theta} \sum_{i=1}^n \log \left( \frac{\Gamma(y_i + r)}{y_i! \Gamma(r)} \left( \frac{\mu_i}{r + \mu_i} \right)^{y_i} \left( \frac{r}{r + \mu_i} \right)^r \right) \\
&= \operatorname{argmax}_{\beta, r} \sum_{i=1}^n \log \left( \frac{\Gamma(y_i + r)}{y_i! \Gamma(r)} \left( \frac{e^{X_i' \beta}}{r + e^{X_i' \beta}} \right)^{y_i} \left( \frac{r}{r + e^{X_i' \beta}} \right)^r \right)
\end{aligned}$$

Interpretation of each coefficient  $\hat{\beta}_k$  (keeping all the other predictors fixed):

- log scale  $[-\infty, \infty]$       “a 1-unit difference in  $x$  corresponds to a  $\hat{\beta}_k$ -unit difference in  $\log(\mu)$ .”
- incidence rate ratio scale  $[0, \infty]$       “a 1-unit difference in  $x$  corresponds to a  $e^{\hat{\beta}_k}$  multiplicative difference in  $\mu$ .”

### A.3 Nonparametric models

**Splines** A spline is a piecewise polynomial function with additional constraints:

- **piecewise polynomial:** we divide  $X$  into intervals, fit each interval with a separate polynomial of low-degree  $m$ . There is seldom any good reason to go beyond a cubic spline ( $m = 3$ ).

<sup>35</sup>The negative binomial distribution is typically specified in terms of the parameters  $(p, r)$ . In a regression framework, it is more intuitive to specify it in terms of its mean  $\mu = \frac{pr}{1-p}$  and  $r$ . We rewrite the probability mass function as  $P(y|\mu, r) = \frac{\Gamma(y+r)}{y! \Gamma(r)} \left( \frac{\mu}{r+\mu} \right)^y \left( \frac{r}{r+\mu} \right)^r$ .



- additional constraints: we impose that the  $m - 1$  first derivatives are continuous at the knots (equal values on both sides), so that the total fit is smooth.  $\implies$  For  $K$  knots, there are  $K + 4$  total degrees of freedom.

As the behavior of polynomials tends to be erratic near the boundaries, we can use a “natural” or “restricted” spline which further imposes that the polynomial fits beyond the boundary knots are linear. I.e., for a natural cubic spline, it sets the cubic and quadratic parts there to 0. The total degrees of freedom are reduced from  $K + 4$  to  $K$ .

Restricted cubic splines are an easy way of including an explanatory variable in a smooth non-linear way in a wide variety of models. With low-degree polynomials, we don’t observe high oscillations of the curve around the data.

In practice: restricted cubic splines are just a transformation of an explanatory variable. This transformed variable can then be entered in any regression command (fixest, logit, glm...)

Interpreting the results:

- graph the adjusted predictions: the predicted outcome against the spline variable. If there are other covariates, show the predicted outcome for an observation with typical values on the other covariates.
- graph the marginal effects, i.e., the spline’s first derivative. = How much does the predicted outcome change for a unit change in the explanatory variable?

## A.4 Multilevel models

**Context** There is some hierarchical structure in our data: individuals belong to groups.<sup>36</sup> Hence observations are not independent, we expect some group effects. And we may be interested in *both* within-group and between-group relationships/effects/variation.

**Structure** We consider a simple 2-level model. The lower level is a regression model with variables at the individual-level, and some of its coefficients (the intercept, and eventually slopes) are allowed to vary by group and are *modeled* i.e., they are given a probability model: this forms the higher group-level model, which can include predictors (group-level variables) or not.

Multilevel models are thereby generalizations of classical regression models, which may allow for varying coefficients but do not *model* them. The generalization proceeds along two dimensions; simple examples of models are given below:<sup>37</sup>

- (i) whether only intercepts or also slopes are modeled;
- (ii) whether group-level predictors are included.

$$\forall i = 1, \dots, N, \quad j = 1, \dots, J:$$

	varying intercept	varying intercept and slope
w/o group-level predictors	$y_i \sim \mathcal{F}_1(a_{j[i]} + \beta x_i, \sigma_y^2)$ $a_j \sim \mathcal{F}_2(\alpha, \sigma_a^2)$	$y_i \sim \mathcal{F}_1(a_{j[i]} + b_{j[i]}x_i, \sigma_y^2)$ $\begin{bmatrix} a_j \\ b_j \end{bmatrix} \sim \text{MVN} \left( \begin{bmatrix} \alpha \\ \beta \end{bmatrix}, \begin{bmatrix} \sigma_a^2 & \sigma_{ab} \\ \sigma_{ab} & \sigma_b^2 \end{bmatrix} \right)$
with group-level predictors	$y_i \sim \mathcal{F}_1(a_{j[i]} + \beta x_i, \sigma_y^2)$ $a_j \sim \mathcal{F}_2(\alpha_0 + \alpha_1 z_j, \sigma_a^2)$	$y_i \sim \mathcal{F}_1(a_{j[i]} + b_{j[i]}x_i, \sigma_y^2)$ $\begin{bmatrix} a_j \\ b_j \end{bmatrix} \sim \text{MVN} \left( \begin{bmatrix} \alpha_0 + \alpha_1 z_j \\ \beta_0 + \beta_1 z_j \end{bmatrix}, \begin{bmatrix} \sigma_a^2 & \sigma_{ab} \\ \sigma_{ab} & \sigma_b^2 \end{bmatrix} \right)$

### Terminology and relation to classical models

For simplicity, let us consider a normally-distributed outcome and a single regressor with a homogeneous effect, so group effects would vary the intercept only.<sup>38</sup>

There is a group structure in our data, so we expect group effects. The question is: How do we account for them?

The answer is: By making an assumption on how they are distributed. More specifically, by making an assumption on their variance  $\sigma_a^2$ , i.e., how much they are pooled toward their mean (pooling = sharing information). We can distinguish the two extremes, and the data-driven compromise from multilevel modeling:

- *Complete pooling*: we don't allow for group effects

<sup>36</sup>Example 1: observations of  $N$  students in  $J$  classrooms; each student  $i$  belongs to a group  $j[i]$ . Example 2: longitudinal data; each dated observation  $\{i, t\}$  belongs to a unit  $i$ . For simplicity, we consider in this section only such hierarchical data with two levels. However, results extend to more levels, that can be nested or not.

<sup>37</sup>For clarity, we use Latin letters for variables or varying parameters, and Greek letters for fixed parameters (as distinguished in a frequentist framework). The fixed parameters that inform varying parameters are also called the “hyperparameters” of the full model.

<sup>38</sup>With a normal distribution as  $\mathcal{F}_1$ , the multilevel model can also be written in a single line:  $y_i = a_{j[i]} + \beta x_i + e_i$ ,  $a_j \sim \mathcal{F}(\alpha, \sigma_a^2)$ ,  $e_i \sim \mathcal{N}(0, \sigma_y^2)$ .

I.e., we neither model them nor adjust for them. We fit the basic model:

$$\begin{aligned} y_i &= \alpha + \beta x_i + e_i, & e_i &\sim \mathcal{N}(0, \sigma_y^2) \\ \iff y_i &= a_{j[i]} + \beta x_i + e_i, & e_i &\sim \mathcal{N}(0, \sigma_y^2), a_j \sim \mathcal{F}(\alpha, 0) \end{aligned}$$

- *No pooling*: we allow for group effects, but don't model them

We include a dummy variable for each  $j$  (i.e., fit a separate intercept per group), or equivalently, we de-mean the data by group. This eliminates the group effects. We fit the “fixed effects” model:

$$\begin{aligned} y_i &= \alpha_{j[i]} + \beta x_i + e_i, & e_i &\sim \mathcal{N}(0, \sigma_y^2) \\ \iff y_i &= a_{j[i]} + \beta x_i + e_i, & e_i &\sim \mathcal{N}(0, \sigma_y^2), a_j \sim \mathcal{F}(\alpha, \infty) \\ \iff (y_i - \bar{y}_{j[i]}) &= \beta(x_i - \bar{x}_{j[i]}) + e_i, & e_i &\sim \mathcal{N}(0, \sigma_y^2) \end{aligned}$$

- *Partial pooling*: we allow for group effects, and model them

The group-level parameters  $a_j$  are partially pooled toward the mean  $\alpha$ , by an amount that depends on the sample size of each group and on  $\sigma_a$ ,<sup>39</sup> which is also estimated from the data. We fit the multilevel model:<sup>40</sup>

$$\begin{cases} y_i = a_{j[i]} + \beta x_i + e_i & e_i \sim \mathcal{N}(0, \sigma_y^2) \\ a_j \sim \mathcal{F}(\alpha, \sigma_a^2) \end{cases}$$

*Partial pooling* is a compromise between *complete pooling* (equivalent to setting  $\sigma_a \rightarrow 0$ ) and *no pooling* (setting  $\sigma_a \rightarrow \infty$ , which risks overstating the variation between groups, i.e., overfitting the data) (Gelman and Hill, 2006, ch. 12).

## Why use multilevel models?

In short: to (1) acknowledge and (2) analyze within-group *and* between-group variations.

- *To account for the dependence in our data (e.g., with time series, spatial correlation, networks...)*. Traditional regression techniques assume independent observations. If some dependence structure in our data isn't modeled, it will be left out in the error term, and the corresponding standard errors of regression coefficients will be underestimated (esp. that of higher-level regressors). A common way to deal with that in econometrics is to not model the group-dependence but to cluster standard errors by group, i.e., after estimation of the regression coefficients, to compute a new estimator of the error covariance matrix that adjusts for the dependence within groups. However, we cannot analyze this dependence. Instead, multilevel modeling models the dependence. A multilevel model is equivalent to a classical regression model with correlated *and modeled* errors.<sup>41</sup>

$$\begin{cases} y_i = \alpha_{j[i]} + X_i \beta + e_i, & e_i \sim \mathcal{N}(0, \sigma_y^2) \\ \alpha_j = \mu_\alpha + \eta_j, & \eta_j \sim \mathcal{N}(0, \sigma_\alpha^2) \end{cases} \iff y_i = \mu_\alpha + X_i \beta + e_i + \underbrace{\eta_{j[i]}}_{e_i^{\text{all}}}, \quad e^{\text{all}} \sim \mathcal{N}(0, \Sigma)$$

<sup>39</sup>Partial pooling is proportional to the variance, not the standard deviation.

<sup>40</sup>In the regression framework, multilevel regression models are a particular case of “random effect” models that pool information across groups, or “mixed effect” models when they include both a “random effect” component and a “fixed effect” component (e.g., a varying-intercept, fixed-slope model). This terminology is confusing, as different disciplines use the terms “fixed” and “random” effects to refer to different things (e.g., sometimes a slope or effect is called “fixed” if it does not vary by group, other times — notably in econometrics — it refers to a model where coefficients vary by group but where this variation is not estimated using a probability model). Gelman and Hill (2006, p. 245) suggests to avoid these terms and describe the model explicitly: for example, varying intercepts and constant slopes.

<sup>41</sup>The error  $e_i^{\text{all}}$  is the sum of an individual-level noise  $e_i$  and a group-level error  $\eta_{j[i]}$  which induces correlation in  $e^{\text{all}}$ . The covariance matrix  $\Sigma$  is parameterized in some way, and these parameters are estimated from the data.

- *To increase efficiency by pooling information.* By treating groups as a “random-effect” within the model, we can pool shared information about means across the groups. *Partially pooling* the varying coefficients will produce more efficient (less noisy) estimates of the  $J$  regression lines than by including group indicators, especially when the number of observations in some groups is small.
- *To model heterogeneity in the relationship to a covariate.* For example, in causal inference, we may be interested not just in an average treatment effect but also in how the effect varies across the population. With a multilevel model, we can model variation in the expected treatment effect, for example as a function of pre-treatment covariates  $X$ .
- *To generalize results to a population not well-represented in the sample.*
  - To do inference for the population of groups when our data are not random samples: one can generalize to a larger population using *multilevel regression and poststratification (MRP)*.
  - To estimate  $y$  for particular groups: notably to get reasonable estimates even for groups with small sample sizes (which is difficult with classical regression).
  - To predict a new observation in a new group: multilevel regression enables to quantify sources of variation, and hence to propagate the uncertainty about the new group into the uncertainty about the new individual in this group; this distinction isn’t provided in classical regression.

Note that when the number of groups  $J$  is small, it is difficult to estimate the between-group variation  $\sigma_a^2$  precisely. As this  $\sigma_a^2$  determines the amount of partial pooling, a bad estimation of its value results in pooling by a somewhat random amount. Hence multilevel modeling adds little to no-pooling models.

## Endogeneity concern, solved by REWB

Context: We want to estimate the causal effect  $\beta$  of  $x_i$  on  $y_i$ . But there are unobserved group (i.e., level-2) effects that are correlated with  $x_i$  and determine  $y_i$ . For example with panel data, these are time-constant factors. Hence the slope estimate in a simple regression of  $y_i$  on  $x_i$  would suffer from omitted variable bias.<sup>42</sup> How to tackle this bias on level 1 coefficients due to omitted variables at level 2?

We saw that the  $J$  unobserved group effects can be considered as either *fixed effects* (i.e., unrelated) or *random effects* (i.e., belonging to a shared distribution), with important implications and caveats:

- FE model = no pooling
  - Conceptualization:  $J$  group effects are unconnected to each other.
  - Method: We eliminate them — by either de-meaning by group or including a dummy variable for each group.
  - Results: The estimator of the within effect is not biased by between effects. But as all the group-level variance is accounted for, no group-level variable can be identified, we can’t say anything about relationships with group-level variables (we throw away a lot of information). We also don’t use information about one group to learn about another — they are deemed unrelated.
- RE model = partial pooling
  - Conceptualization:  $J$  group effects are random variables with a joint distribution.
  - Method: We model them with error, i.e., we estimate the parameters of that distribution.

---

<sup>42</sup>This endogeneity concern can arise with any form of nested data (e.g., Bafumi and Gelman (2007) consider a general multilevel structure, with observations at level 1, indexed by  $i$ , belonging to groups — level 2 — indexed by  $s[i]$ ). A very common case is that of longitudinal data ( $it$  = level 1,  $i$  = level 2). The concern is then of time-constant effects that are plausibly correlated with  $x_{it}$ .

\* Option 1: Simple RE model

Results: The random-effects estimator is biased for the within effect, as it is a weighted average of the within estimator and the between estimator.<sup>43</sup> The RE estimates of the group effects are similar to the fixed effects in a ‘dummy variable’ FE model, but are drawn closer towards their mean — unreliably estimated (groups with few observations) and more extreme values are shrunk the most.

\* Option 2: Within-Between RE (REWB) (Bell et al., 2019) or ‘correlated RE’ (CRE) (Wooldridge, 2013, 14.3) model:  $\bar{x}_{j[i]}$  is added as a regressor, to extract the problematic correlation from the composite error term. Results: By decomposing  $x$  into a within-group component and a between-group component, we recover an estimator of the within effect that is not biased by group-level confounders — as the FE model does — but also enables us to:

- test whether the B. and the W. effects are significantly different (Hausman test: tests whether the difference between the two coefficients in the REWB model is statistically different from 0);
- estimate the between relationship, as well as any relationship between the outcome and a group-level covariate (though we can’t interpret their coefficients as causal);
- estimate the level-2 variance and compare it to the level-1 variance.

△ With a *non-identity link function*, unbiasedness is guaranteed only if  $u_j$  is a linear function of  $\bar{x}_{j[i]}$ .<sup>44</sup> However, the available evidence (from simulations) suggests that the bias of the REBW method remains small in most situations (Bell et al., 2019). One can also include functions of  $\bar{x}_{j[i]}$  as regressors to characterize more flexible functional forms of the correlation.<sup>45</sup> If the estimates of the added coefficients aren’t significant and the estimate of  $\beta_W$  doesn’t change much, it suggests the linearity assumption is reasonable and bias should not be such an issue.

CCL: do FE only when really don’t care about between-group relationships + don’t worry about efficiency (have a very large sample size).

## Inference

• **Frequentist (point estimation)**

First the hyperparameters are estimated via Maximum Likelihood or Restricted Maximum Likelihood (REML), then inference is performed for the coefficients conditional on the estimated hyperparameters.

△ ML estimators are unbiased for the group-level mean parameters but downward-biased for the group-level variance parameters (especially when the number of groups is small), because the mean parameters are assumed to be known with certainty when estimating the variance parameters. Instead, REML accounts for the number of mean parameters estimated, losing 1 degree of freedom for each, and so produces unbiased estimators of variance parameters. However, for model comparison: likelihood ratio tests for REML require exactly the same fixed effects specification in both models. So, when comparing models with *different fixed effects* with an LR test, ML must be used.

<sup>43</sup>  $\beta_1 = \frac{w_W \beta_W + w_B \beta_B}{w_W + w_B}$ , where  $w_W \equiv 1/SE[\hat{\beta}_W]^2$  and  $w_B \equiv 1/SE[\hat{\beta}_B]^2$  are the precisions of the within estimate and the between estimate, respectively. As there are more data at level 1 (and therefore higher precision of the within estimate),  $\hat{\beta}_1$  will often tend towards the within estimate (Bell et al., 2019).

<sup>44</sup> I.e., to get an unbiased effect, we are trading one assumption of linearity for another? Standard model: assume identity link function. REBW model with non-identity link function: assume  $u_j$  is a linear function of  $\bar{x}_j$ .

<sup>45</sup> P. Allison suggests using polynomial functions of the means, i.e., including not only  $\bar{x}_i$  but also  $\bar{x}_i^2, \bar{x}_i^3$  as regressors, or other cluster-level functions of the  $x_{it}$ , such as their standard deviation: <https://statisticalhorizons.com/problems-with-the-hybrid-method/>

👍 computational: fast.

🗨️ A. Gelman: *“The usual non-Bayesian procedures are designed to work well asymptotically (in the case of hierarchical models, this is the limit as the number of groups approaches infinity). But as noted Bayesian J. M. Keynes could’ve said, asymptotically we’re all dead.”*

- **Bayesian**

All levels are fitted simultaneously. The hyperparameters are given a prior distribution, and we estimate their whole posterior distribution.

👍 accounts for all the uncertainty in the parameter estimates when predicting the varying intercepts and slopes, and their associated uncertainty.

🗨️ computational: slow. Markov chain Monte Carlo simulations are generally much slower than (RE)ML estimation.

FE	$y_{it} = \beta_1(x_{it} - \bar{x}_i) + u_i + e_{it}, \quad u_i \sim \mathcal{N}(u_0, \infty)$ $\iff (y_{it} - \bar{y}_i) = \beta_1(x_{it} - \bar{x}_i) + e_{it}$
RE	$\begin{cases} y_{it} = \beta_0 + \beta_1 x_{it} + \beta_3 z_i + u_i + e_{it} \\ u_i \sim \mathcal{N}(0, \sigma_u^2) \end{cases}$ $\iff \begin{cases} y_{it} = a_i + \beta_1 x_{it} + e_{it} \\ a_i \sim \mathcal{N}(\beta_0 + \beta_3 z_i, \sigma_a^2) \end{cases}$
<p>Group effect of <math>x_i</math>:</p> <ul style="list-style-type: none"> <li>homogeneous (varying intercepts, fixed slopes) <math display="block">\begin{cases} y_{it} = \beta_0 + \beta_1(x_{it} - \bar{x}_i) + \beta_2 \bar{x}_i + \beta_3 z_i + \nu_i + e_{it} \\ \nu_i \sim \mathcal{N}(0, \sigma_\nu^2) \end{cases}</math> <math display="block">\iff \begin{cases} y_{it} = a_i + \beta_1(x_{it} - \bar{x}_i) + e_{it} \\ a_i \sim \mathcal{N}(\beta_0 + \beta_2 \bar{x}_i + \beta_3 z_i, \sigma_a^2) \end{cases}</math> </li> <li>heterogeneous (varying intercepts and slopes) <ul style="list-style-type: none"> <li>group-level predictors for the intercepts only</li> </ul> </li> </ul>	
REWB	$\begin{cases} y_{it} = \beta_0 + \beta_1(x_{it} - \bar{x}_i) + \beta_2 \bar{x}_i + \beta_3 z_i + \nu_{i0} + \nu_{i1}(x_{it} - \bar{x}_i) + e_{it} \\ \begin{bmatrix} \nu_{i0} \\ \nu_{i1} \end{bmatrix} \sim \text{MVN} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{\nu_0}^2 & \sigma_{\nu_{01}} \\ \sigma_{\nu_{01}} & \sigma_{\nu_1}^2 \end{bmatrix} \right) \end{cases}$ $\iff \begin{cases} y_{it} = a_i + b_i(x_{it} - \bar{x}_i) + e_{it} \\ \begin{bmatrix} a_i \\ b_i \end{bmatrix} \sim \text{MVN} \left( \begin{bmatrix} \beta_0 + \beta_2 \bar{x}_i + \beta_3 z_i \\ \beta_1 \end{bmatrix}, \begin{bmatrix} \sigma_a^2 & \sigma_{ab} \\ \sigma_{ab} & \sigma_b^2 \end{bmatrix} \right) \end{cases}$ <ul style="list-style-type: none"> <li>group-level predictors for the intercepts &amp; slopes</li> </ul> $\begin{cases} y_{it} = a_i + b_i(x_{it} - \bar{x}_i) + e_{it} \\ \begin{bmatrix} a_i \\ b_i \end{bmatrix} \sim \text{MVN} \left( \begin{bmatrix} \beta_0 + \beta_2 \bar{x}_i + \beta_3 z_i \\ \gamma_0 + \gamma_2 \bar{x}_i + \gamma_3 z_i \end{bmatrix}, \begin{bmatrix} \sigma_a^2 & \sigma_{ab} \\ \sigma_{ab} & \sigma_b^2 \end{bmatrix} \right) \end{cases}$ $\iff y_{it} = [\beta_0 + \beta_2 \bar{x}_i + \beta_3 z_i + \nu_{i0}] + [\gamma_0 + \gamma_2 \bar{x}_i + \gamma_3 z_i + \nu_{i1}](x_{it} - \bar{x}_i) + e_{it}$ $\iff y_{it} = \beta_0 + \beta_2 \bar{x}_i + \beta_3 z_i + \nu_{i0} + \delta_0(x_{it} - \bar{x}_i) + \delta_2 \bar{x}_i(x_{it} - \bar{x}_i) + \delta_3 z_i(x_{it} - \bar{x}_i) + \nu_{i1}(x_{it} - \bar{x}_i) + e_{it}$