

# Universidad Externado

## Bayesian Hotel Pricing Lab

Winter School Montevideo

Montevideo, Uruguay



### Objective

This lab aims to build **Bayesian models** to understand hotel price dynamics across Europe. Students will explore whether prices follow a global trend or are influenced by city-specific or country-specific effects. The main goal is to detect **global vs. local pricing patterns** and learn how **hierarchical models** can uncover hidden structures in real-world data.

### Dataset

You will work with two datasets:

- `hotelseuropefeatures.csv`: general hotel info: `hotel_id`, `city`, `country`, `stars`, `rating`, `neighbourhood`, `accommodation_type`, etc.
- `hotelseuropeprice.csv`: pricing and availability: `price`, `offer`, `weekend`, `holiday`, `nnights`, `scarce_room`, etc.

These datasets must be merged using the `hotel_id` column.

### Project Structure

#### 1. Problem Definition

- Formulate a research question: *Do prices follow a global trend or do they differ by city/country?*

- Choose your grouping structure: **City** vs. **Country**
- Optional: restrict your dataset to one specific country if you prefer a more focused analysis.

## 2. Data Preparation

- Merge datasets
- Handle missing values
- Select and transform predictors (categorical/numerical)
- Choose target variable: **price**

## 3. Exploratory Analysis

Limit to a maximum of **three informative plots**, such as:

- Boxplots of **price** per city or country
- Correlations between numerical variables
- Scatterplots between predictors and target

## 4. Modeling Phase (Flexibility Encouraged )

You are **free to define the level of complexity** of your model based on your learning goals and data quality. Choose **at least two** of the following approaches:

### a. Global Model

- Bayesian linear regression using all data
- No group structure
- Example:

$$price_i \sim \mathcal{N}(\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots, \sigma)$$

### b. Local Models by City or Country

- Separate models per city or country
- Analyze how posterior distributions vary across locations

### c. Hierarchical Model

- Multilevel intercept or slope model by **city** or **country**
- Example:

$$\alpha_j \sim \mathcal{N}(\mu_\alpha, \sigma_\alpha) price_i \sim \mathcal{N}(\alpha_{j[i]} + \beta x_i, \sigma)$$

## 5. Posterior Analysis & Visualization

- Plot posterior distributions (with **arviz**)
- Show 90% HDIs for parameters
- Boxplots of intercepts per group (city/country)
- Visual comparison of global vs. local parameters

## 6. Discussion

- What model structure explains price best?
- Are there major differences in price dynamics across cities?
- Does a global model overlook important group-specific effects?
- How sensitive are your results to prior assumptions?

## Model Interpretation Guide

Model Type	Description	Expected Insights
Global Model	Fits a single linear regression across all observations without grouping.	Identifies general price drivers; may overlook local variations.
Local Models	Fits a separate model for each city or country independently.	Captures heterogeneity but lacks pooling; useful when city behavior is distinct.
Hierarchical Model	Models parameters (e.g., intercepts) as varying by city but sharing a common global prior.	Combines global and local structure; shrinks estimates in cities with few data points.

Table 1: Modeling options to explain hotel price dynamics.

## How to Compare and Select Models

- **Posterior predictive checks:** Are predicted values consistent with observed data?
- **WAIC/LOO:** Use information criteria for model comparison.
- **Interpretability:** Can you explain results clearly?
- **Uncertainty:** Look at HDIs and trace plots for stable estimates.
- **Sample size per group:** Prefer hierarchical models when some groups are small.

## Technical Requirements

- Python 3.10+
- Libraries: `pymc`, `arviz`, `pandas`, `matplotlib`, `seaborn`
- Use MCMC sampling via `pm.sample`
- Ensure good diagnostics:  $\hat{R} < 1.01$ , high ESS, trace plots

## Deliverables

Each group must submit:

- A well-documented `.ipynb` notebook or `.py` script
- Clean visualizations
- Model interpretation
- Code reproducibility

## Tip

You are **not required** to build the most complex model **build what you can justify** and what best answers your question. Quality of reasoning and clarity of interpretation matter more than complexity.

## Presentation Guidelines

Each group will present their work in a **10-minute oral presentation**. The presentation should include:

### Structure:

1. Problem motivation and research question
2. Dataset description and cleaning
3. Exploratory visualizations (max 3)
4. Modeling strategy and results
5. Interpretation of posteriors
6. Conclusions and reflections