

# Labor Income Prediction

## Predicting Income

Sany León, Andrés Suárez, and Juan Rueda

2026-02-21



# Research Question

Can prediction-based income models reliably identify individuals with anomalous earnings patterns, and do large prediction errors reflect misreporting or structural model limitations?

# Results

	LOOCV RMSE		Test RMSE	
	With outliers	Without outliers	With outliers	Without outliers
Modelo 1	0.662	0.662	0.846	0.846
Modelo 2	0.676	0.676	0.855	0.855
Modelo 3	0.399	0.399	0.666	0.666
Modelo 4	0.391	0.391	0.655	0.655
Modelo 5	0.390	0.390	0.655	0.655
Modelo 6	0.389	0.389	0.652	0.652
Modelo 7	0.390	0.390	0.653	0.653
Modelo 8	0.388	0.388	0.650	0.650
Modelo 9	0.377	0.377	0.648	0.648

Train Obs. (Outliers-No outliers): 10334-10334

Dos observaciones fueron eliminadas en LOOCV para evitar leverage de 1 en el test set

Model 9 achieves the lowest validation RMSE (0.648). Notably, its LOOCV RMSE (0.377) is considerably smaller, indicating that LOOCV may understate the true out-of-sample prediction error relative to the validation-sample benchmark.

# Results

## Modelos

$$(1) \quad \log(\text{ingresos}_i) = \beta_0 + \beta_1 \text{Edad}_i + \beta_2 \text{Edad}_i^2 + u_i$$

$$(2) \quad \ln(\text{salariorio})_{i,f} = \beta_0 + \beta_1 \text{Mujer}_{i,f} + u_i$$

$$(3) \quad \ln(\text{salariorio})_{i,f} = \beta_0 + \beta_1 \text{Mujer}_{i,f} + \beta_2 \text{Edad}_{i,f} + \beta_3 \text{Edad}_{i,f}^2 \\ + \beta_4 \text{NivelEduc}_{i,f} + \beta_5 \text{Oficio}_{i,f} + \beta_6 \text{Relab}_{i,f} \\ + \beta_7 \text{TamFirma}_{i,f} + u_i$$

$$(4) \quad \ln(\text{salariorio})_{i,f} = \beta_0 + \beta_1 \text{Mujer}_{i,f} + \beta_2 \text{Edad}_{i,f} + \beta_3 \text{Edad}_{i,f}^2 \\ + \beta_4 \text{NivelEduc}_{i,f} + \beta_5 \text{Oficio}_{i,f} + \beta_6 \text{Relab}_{i,f} \\ + \beta_7 \text{TamFirma}_{i,f} + \beta_8 \text{NumMen}_f + \beta_9 \text{NumMay}_f + u_i$$

$$(5) \quad \log(y_i) = \beta_0 + \beta_1 \text{Sexo}_i + \beta_2 \text{Edad}_i + \beta_3 \text{Edad}_i^2 + \beta_4 \text{Edad}_i^3 \\ + \beta_5 \text{NivelEduc}_i + \beta_6 \text{Oficio}_i + \beta_7 \text{Relab}_i \\ + \beta_8 \text{TamFirma}_i + \beta_9 \text{Formalidad}_i + u_i$$

# Results

## Modelos

$$(6) \quad \log(y_i) = \beta_0 + \beta_1 \textit{Sexo}_i + \beta_2 \textit{Edad}_i + \beta_3 \textit{Edad}_i^2 \\ + \beta_4 (\textit{NivelEduc}_i \times \textit{Formalidad}_i) \\ + \beta_5 \textit{Oficio}_i + \beta_6 \textit{Relab}_i + \beta_7 \textit{TamFirma}_i + u_i$$

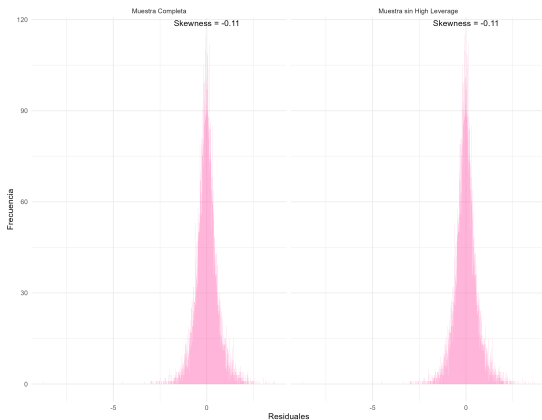
$$(7) \quad \log(y_i) = \beta_0 + \beta_1 \textit{Sexo}_i + \beta_2 \textit{Edad}_i + \beta_3 \textit{Edad}_i^2 \\ + \beta_4 \textit{NivelEduc}_i + \beta_5 \textit{Oficio}_i \\ + \beta_6 (\textit{Relab}_i \times \textit{TamFirma}_i) + \beta_7 \textit{Formalidad}_i + u_i$$

$$(8) \quad \log(y_i) = \beta_0 + \beta_1 \textit{Sexo}_i + \beta_2 \textit{Edad}_i + \beta_3 \textit{Edad}_i^2 \\ + \beta_4 \textit{Oficio}_i + \beta_5 \textit{Formalidad}_i \\ + \beta_6 (\textit{NivelEduc}_i \times \textit{Formalidad}_i) \\ + \beta_7 (\textit{Relab}_i \times \textit{TamFirma}_i) + u_i$$

$$(9) \quad \log(y_i) = \beta_0 + \beta_1 \textit{Sexo}_i + \beta_2 \textit{Edad}_i + \beta_3 \textit{Edad}_i^2 \\ + \beta_4 \textit{Oficio}_i + \beta_5 \textit{Formalidad}_i \\ + \beta_6 (\textit{NivelEduc}_i \times \textit{Formalidad}_i \times \textit{Edad}_i) \\ + \beta_7 (\textit{Relab}_i \times \textit{TamFirma}_i \times \textit{Edad}_i) + u_i$$

# Results

Figure 1: Distribution of the Prediction Residuals of Model from Equation (9)



# Discussion

- ▶ **Model performance:** Model 9 achieves the lowest validation RMSE (0.648), showing that including higher-order interactions improves predictive accuracy.
- ▶ **LOOCV vs. test error:** LOOCV RMSE values are substantially lower than validation RMSE (e.g., 0.377 vs. 0.648 for Model 9), indicating that LOOCV may underestimate true out-of-sample prediction error.
- ▶ **Outlier robustness:** Removing outliers does not affect RMSE, suggesting the models are robust to extreme observations and that large prediction errors likely reflect structural model limitations rather than anomalous data.