



Universidad San Francisco de Quito

**Maestría en Ciencia de datos
Fundamentos de Ciencia de Datos**

Carla Parra

Abril 2025

CRISP-DM-Phase 1: Understanding of business (I)



- **Objetivo general:** El presente proyecto tiene como objetivo principal desarrollar un modelo predictivo que permita estimar el porcentaje de mujeres que ocupan cargos gerenciales en el sector privado, a partir de diversos indicadores socioeconómicos y laborales
- **Objetivos específicos:**
 1. **Preprocesar y transformar los datos.**
 2. **Evaluar la relación entre indicadores financieros, de educación y laborales** para predecir la participación femenina en cargos gerenciales en el sector privado en LATAM.
 3. **Construir y comparar modelos de regresión multivariada**, como Lasso y Random Forest para identificar los mejores predictores y evaluar su capacidad explicativa y predictiva.
 4. **Implementar técnicas de reducción de dimensionalidad (PCA)** para optimizar el desempeño del modelo y facilitar la interpretación de los resultados.



CRISP-DM-Phase 1: Understanding of business (II)



- **Pregunta de investigación:**

¿Qué indicadores socioeconómicos y laborales que se asocian con una mayor representación femenina en roles de liderazgo en el sector privado?

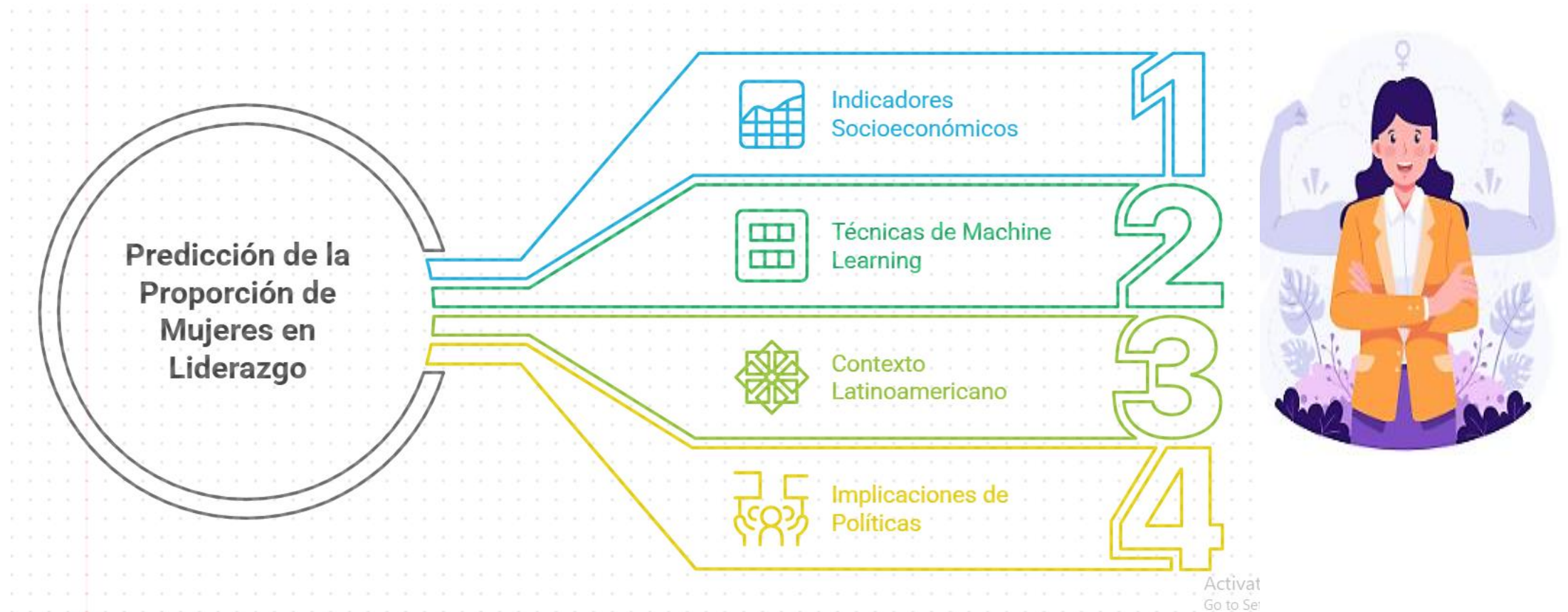
- **Hipótesis principal:**

La participación de mujeres en cargos gerenciales en el sector privado está significativamente influenciada por indicadores de inclusión financiera, nivel educativo y condiciones laborales del país.

- **Hipótesis secundaria:**

Modelos de regresión regularizada como Lasso pueden superar a modelos basados en árboles (como Random Forest) en la predicción de este indicador cuando las relaciones entre variables son principalmente lineales.

CRISP-DM-Phase 1: Understanding of business -Alcance (III)

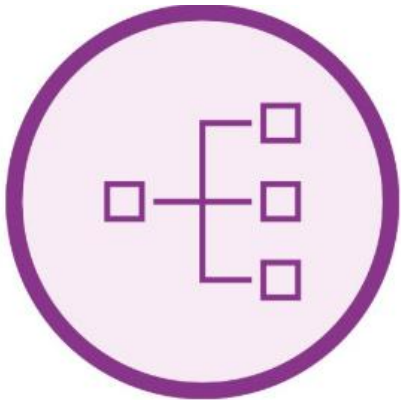


CRISP-DM-Phase 1: Understanding of business - Limitaciones (IV)

Pros	VS	Contras
 Rendimiento mejorado		 Sin inferencia causal
 Perspectivas basadas en datos		 Interpretabilidad reducida
		 Excluye variables culturales/legal
		 Ignora diferencias institucionales

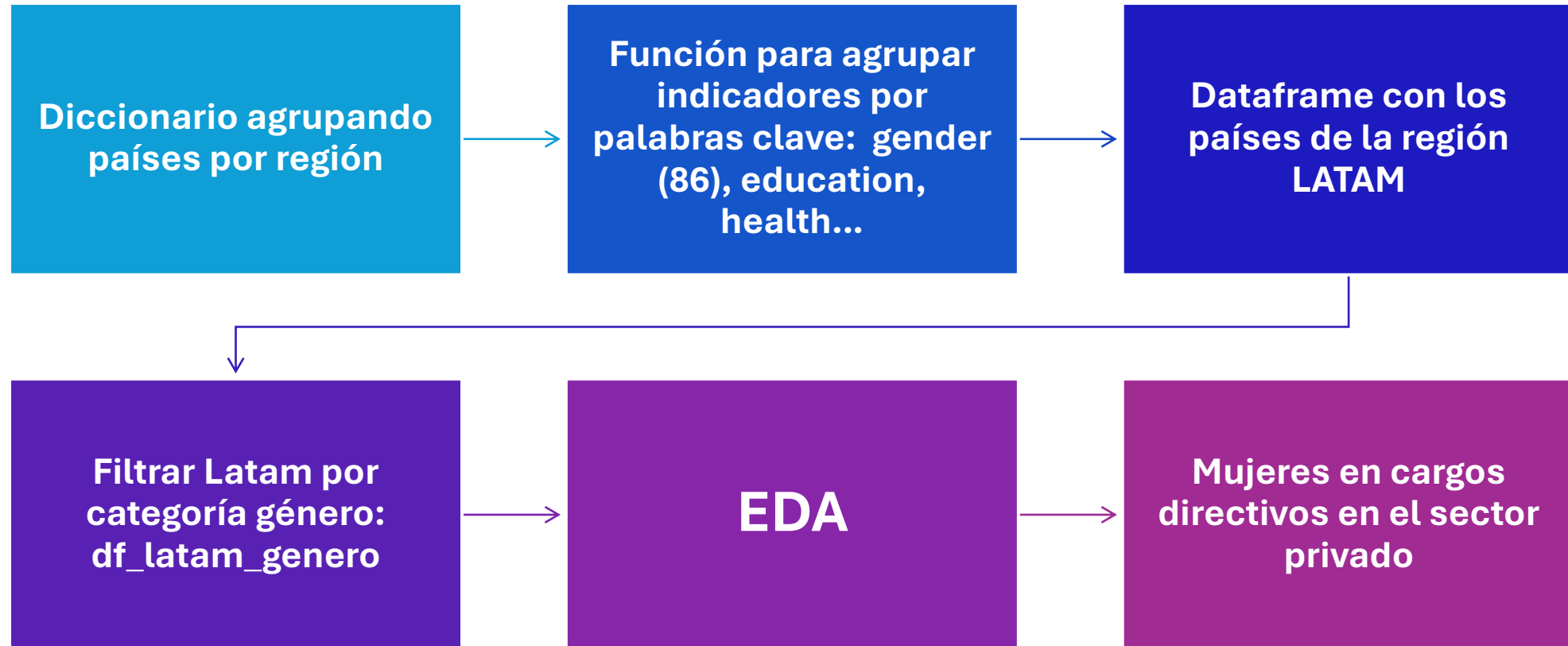


Phase 2: Understanding of data

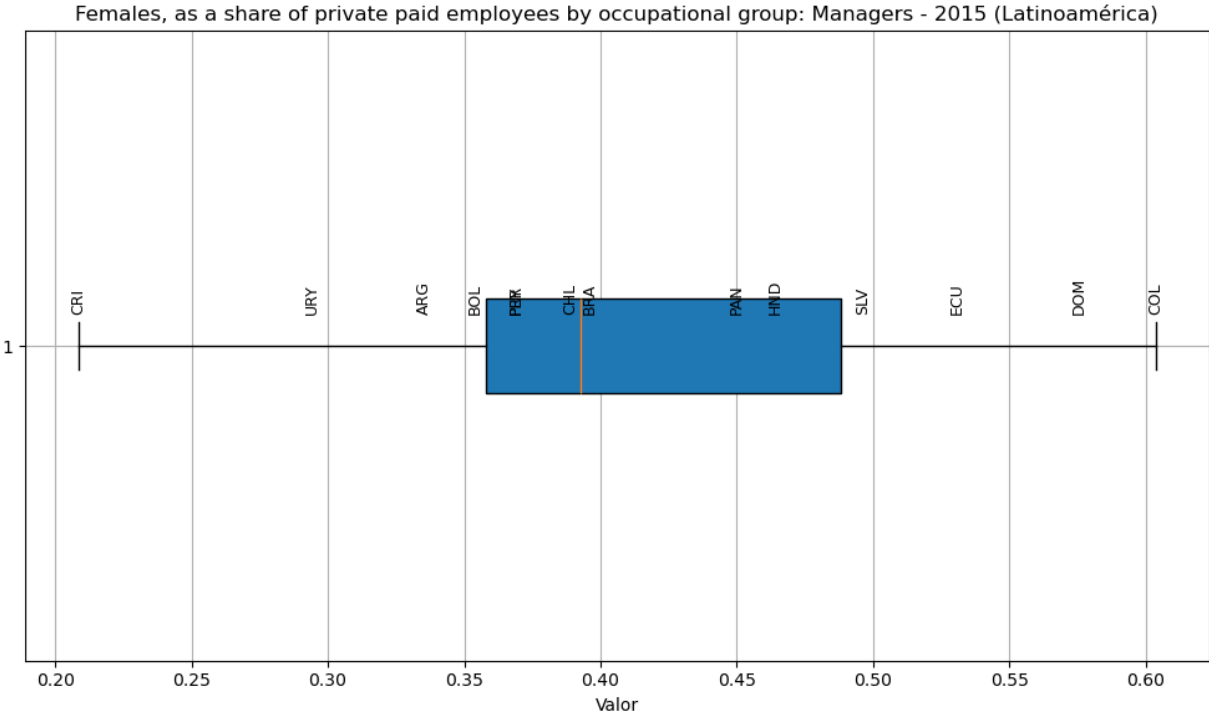


- Fuente principal: **Worldwide Bureaucracy Indicators (WWBI)**
(Enlace: <https://datacatalog.worldbank.org/search/dataset/0038132>)
- Incluye 302 indicadores para 202 economías, contruidos a partir de encuestas de hogares, fuerza laboral y datos administrativos. Estos abarcan características demográficas, brechas salariales, equidad de género, y la estructura del gasto público en salarios.
- **Filas (observaciones):** 61,004
- **Columnas (variables):** 27
- **Tipo de datos:**
 - **4 columnas categóricas:**
 - Country Name, Country Code, Indicator Name, Indicator Code
 - **23 columnas numéricas (años del 2000 al 2022)**

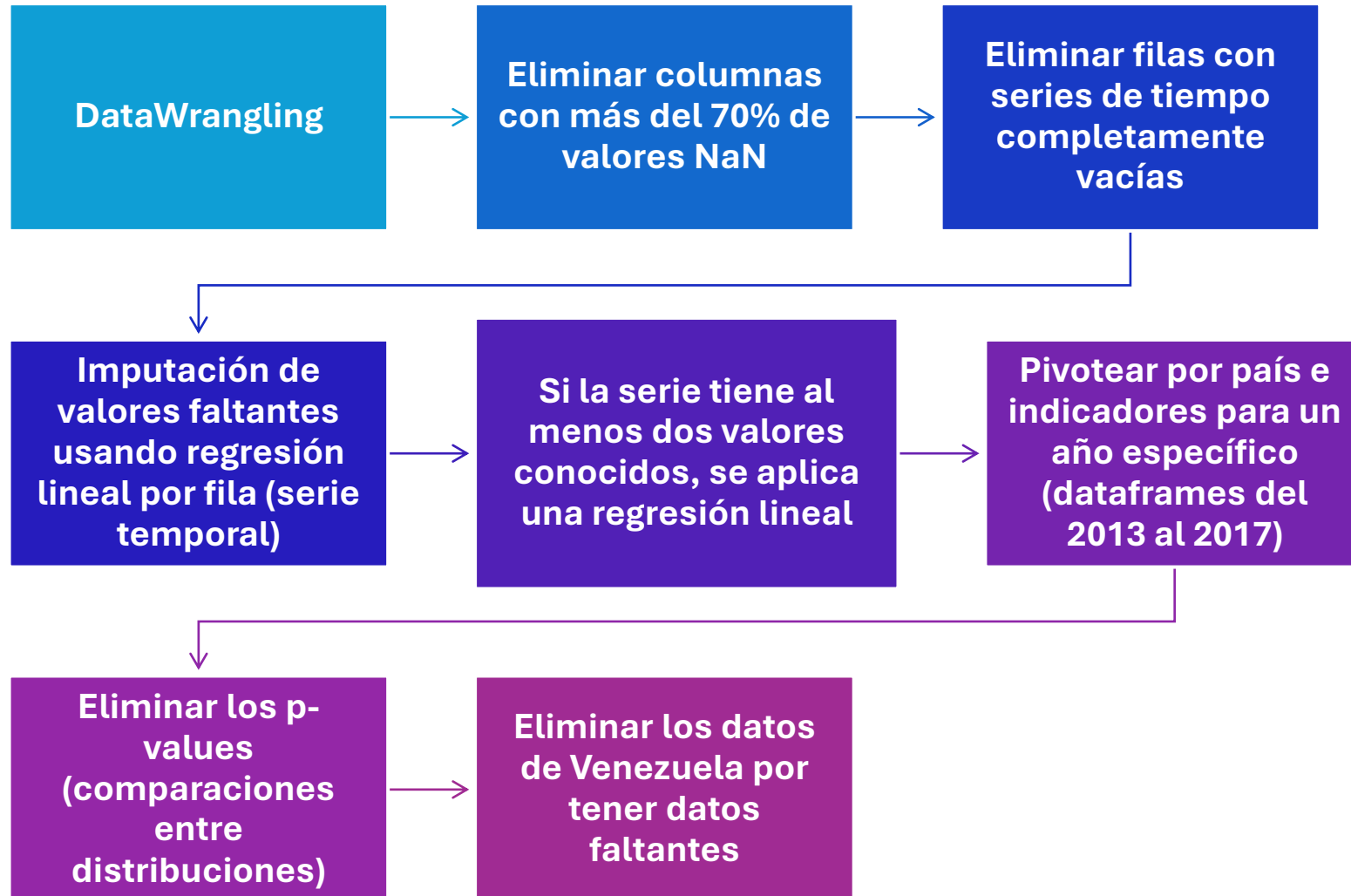
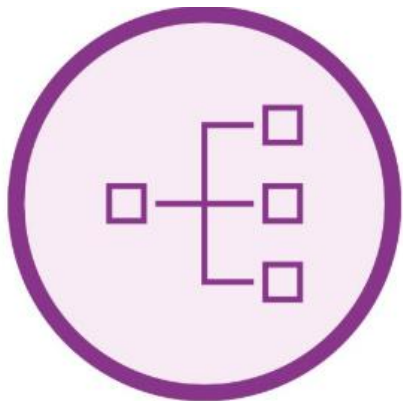
Phase 2: Understanding of data



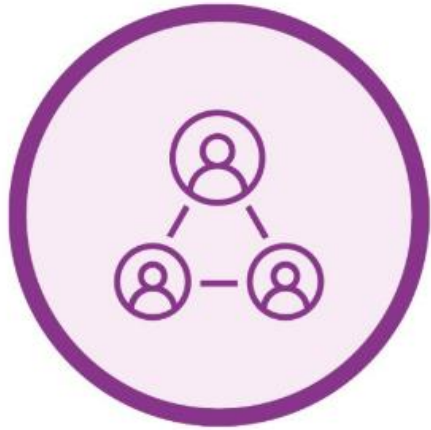
Phase 2: Understanding of data



Phase 3: Data preparation



Phase 4: Modeling



Split Dataset

Training and test sets created



Feature Selection

SelectKbest
Irrelevant features removed
Polynomial feature (no linear relations)



Standardization

Data normalized using z-score
PCA



Model Regression

Lasso or Random Forest applied



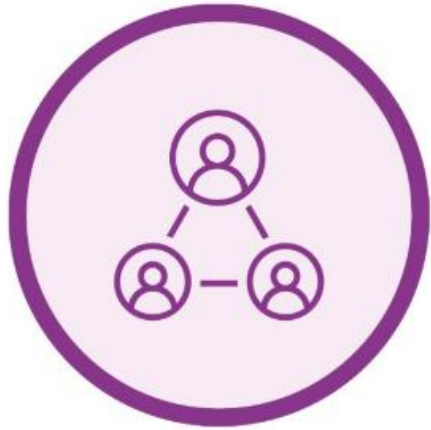
Evaluate Configurations

Performance across setups compared



Phase 4: Modeling

Which pipeline configuration yields the best model performance?



Include Polynomial Features

Expands feature space, potentially capturing non-linear relationships



Include PCA

Reduces dimensionality, potentially improving computational efficiency and reducing overfitting

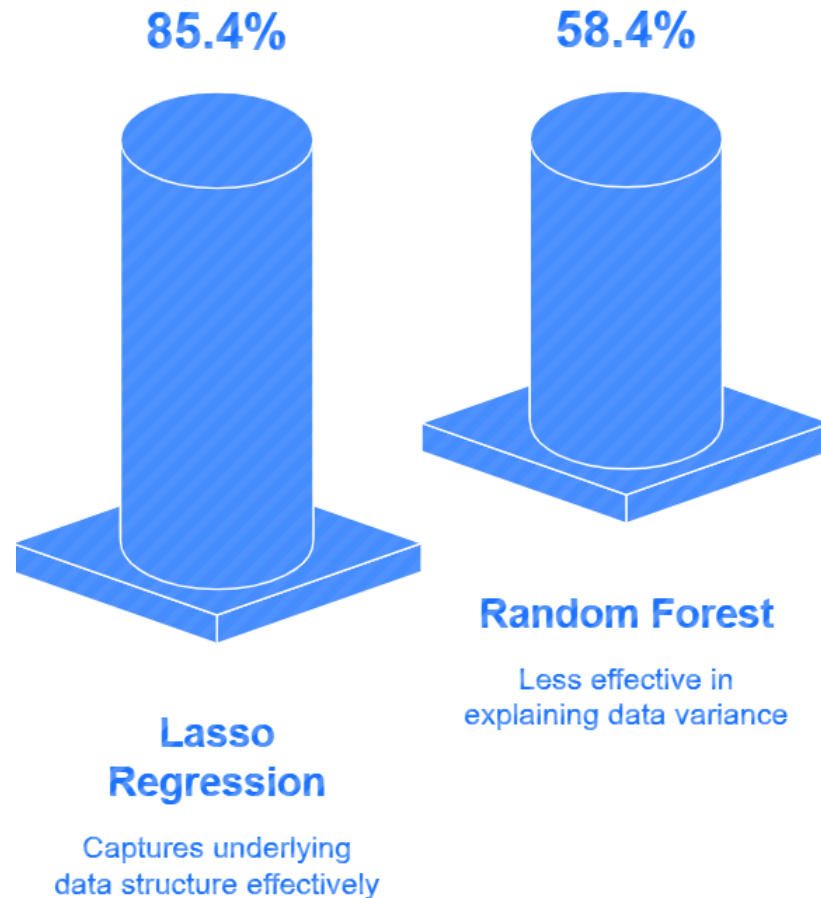


Exclude Both

Maintains original feature space, ensuring simplicity and effectiveness with proper selection

Phase 5: Evaluation

Comparison of R^2 Values for Lasso and Random Forest



Which model provides more reliable predictions?



Lasso Regression

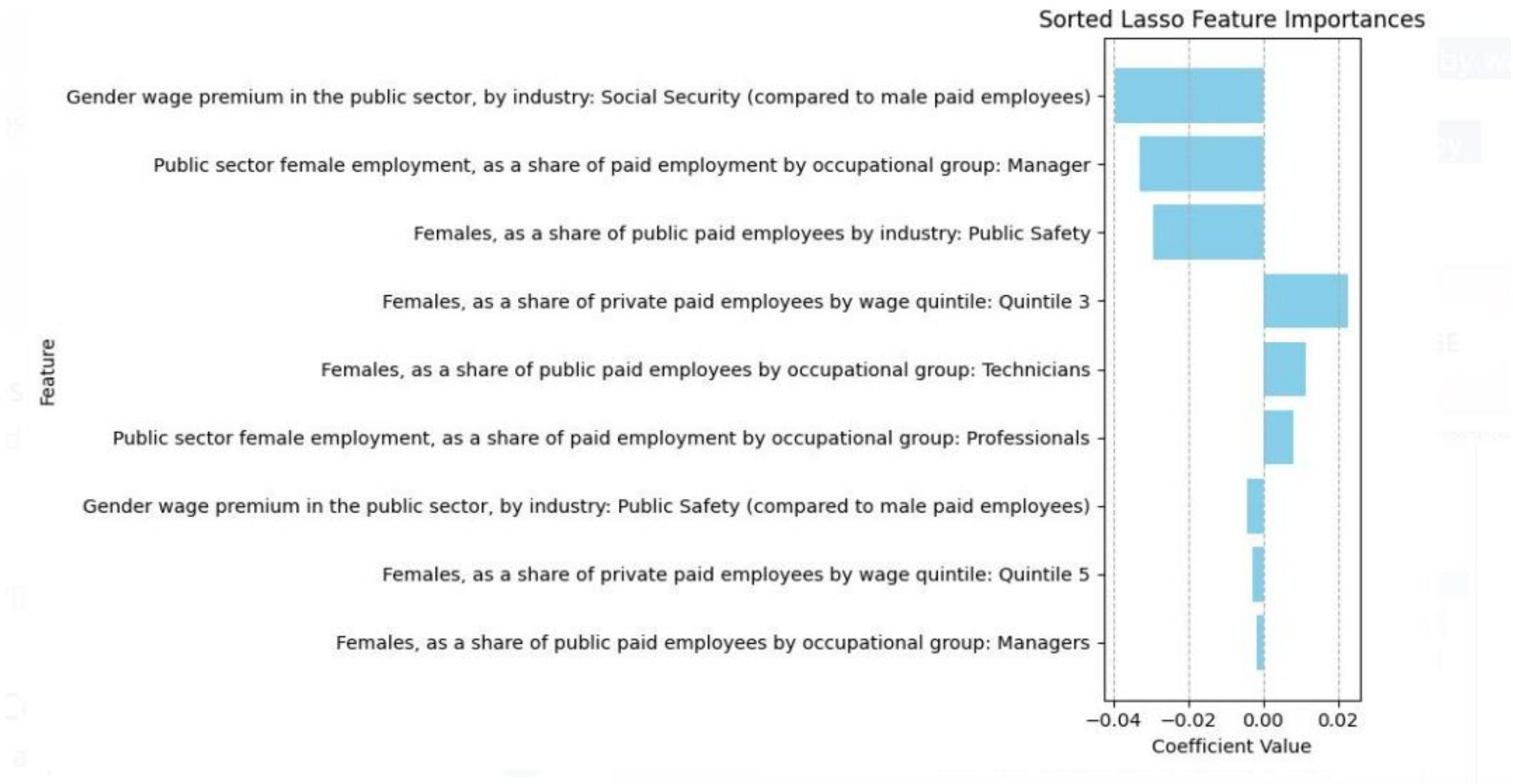
Achieves lower MSE and MAE, indicating higher accuracy



Random Forest

Higher MSE and MAE, suggesting less precision

Resultados:



Conclusiones:

1. Las brechas salariales en el sector público afectan al sector privado

→ Donde los hombres ganan mucho más en el sector público, hay menos mujeres como gerentes en el sector privado.

Brecha salarial de género en el sector público: Seguridad Social (-0.039)

Brecha salarial de género en el sector público: Seguridad Pública (-0.004)

2. El sector público atrae a las líderes mujeres

→ Más mujeres en cargos directivos públicos suele significar menos en el sector privado— puede haber competencia entre sectores.

Mujeres en el quintil 3 de salarios en el sector privado (+0.023)

3. El sector público atrae a las líderes mujeres

→ Más mujeres en cargos directivos públicos suele significar menos en el sector privado— puede haber competencia entre sectores.

Empleo femenino en el sector público como porcentaje de gerentes (-0.033)

Mujeres en cargos de gerencia en el sector público (-0.0015)

Phase 6: Deployment

Esquema de Despliegue del Modelo de Predicción de Liderazgo Femenino

