

Reporte de Experimentación: Few-Shot Bayesian Optimization (FSBO)

Análisis Comparativo con Baselines de HPO

Proyecto Meta-Learning para Optimización de Hiperparámetros

Proyecto Académico MetaLearning
transfer-learning

Enero 2026

Resumen

Este reporte presenta los resultados experimentales de la evaluación del modelo **Few-Shot Bayesian Optimization (FSBO)** implementado para optimización de hiperparámetros de algoritmos de machine learning. Se utilizó un protocolo de **5-Fold Cross-Validation sobre tareas** con múltiples semillas aleatorias para garantizar robustez estadística. Los experimentos se ejecutaron sobre 4 algoritmos de clasificación (AdaBoost, Random Forest, LibSVM.SVC, AutoSklearn) comparando FSBO contra dos baselines: Random Search y GP-RS (Gaussian Process con Random Sampling). Los resultados demuestran que FSBO obtiene el mejor rendimiento promedio en todos los algoritmos evaluados, con mejoras estadísticamente significativas en Random Forest ($p < 0,001$).

Índice

| | |
|---|----------|
| 1. Introducción | 3 |
| 1.1. Motivación | 3 |
| 1.2. Objetivo | 3 |
| 1.3. Contribuciones | 3 |
| 2. Metodología Experimental | 3 |
| 2.1. Protocolo de Evaluación | 3 |
| 2.1.1. K-Fold Cross-Validation sobre Tareas | 3 |
| 2.1.2. Configuración Experimental | 3 |
| 2.2. Métodos Comparados | 4 |
| 2.2.1. FSBO (Few-Shot Bayesian Optimization) | 4 |
| 2.2.2. Random Search | 4 |
| 2.2.3. GP-RS (Gaussian Process con Random Sampling) | 4 |
| 2.3. Métricas de Evaluación | 4 |
| 2.3.1. Normalized Regret (NR) | 4 |
| 2.3.2. Area Under Curve (AUC) | 4 |
| 2.3.3. Time to 95 % Optimal | 4 |
| 3. Resultados Experimentales | 5 |
| 3.1. Resumen Global | 5 |
| 3.2. Análisis Detallado por Algoritmo | 5 |
| 3.2.1. AdaBoost | 5 |
| 3.2.2. Random Forest | 5 |

| | |
|--|-----------|
| 3.2.3. LibSVM_SVC | 5 |
| 3.2.4. AutoSklearn | 6 |
| 3.3. Curvas de Convergencia | 6 |
| 3.4. Normalized Regret Over Time | 7 |
| 4. Análisis Estadístico | 7 |
| 4.1. Test de Friedman | 7 |
| 4.2. Test Post-Hoc de Nemenyi | 7 |
| 4.3. Test de Wilcoxon Pareado | 8 |
| 5. Discusión | 8 |
| 5.1. ¿Por qué FSBO funciona mejor? | 8 |
| 5.1.1. Transfer Learning Efectivo | 8 |
| 5.1.2. Warm-Start Inteligente | 8 |
| 5.1.3. Task Augmentation | 8 |
| 5.2. ¿Por qué Random Forest muestra la mayor mejora? | 9 |
| 5.3. ¿Por qué AutoSklearn no muestra mejora significativa? | 9 |
| 5.4. Limitaciones | 9 |
| 6. Conclusiones | 9 |
| 6.1. Hallazgos Principales | 9 |
| 6.2. Implicaciones Prácticas | 9 |
| 6.3. Trabajo Futuro | 10 |
| A. Resultados por Fold | 10 |
| A.1. AdaBoost - Detalle por Fold | 10 |
| A.2. Random Forest - Detalle por Fold | 10 |
| B. Código de Reproducibilidad | 11 |

1. Introducción

1.1. Motivación

La optimización de hiperparámetros (HPO) es un desafío fundamental en machine learning. Los métodos tradicionales como grid search o random search requieren numerosas evaluaciones, lo cual es costoso computacionalmente. El enfoque de **transfer learning** permite aprovechar conocimiento de tareas previas para acelerar la optimización en nuevas tareas.

1.2. Objetivo

Este reporte evalúa la efectividad del modelo FSBO [1] implementado, comparándolo contra baselines estándar en la literatura de HPO.

1.3. Contribuciones

- Implementación completa del framework FSBO con Deep Kernel Gaussian Processes
- Protocolo de evaluación riguroso con K-Fold CV sobre tareas
- Análisis estadístico completo con tests de Friedman, Nemenyi y Wilcoxon
- Comparación sistemática con baselines relevantes

2. Metodología Experimental

2.1. Protocolo de Evaluación

2.1.1. K-Fold Cross-Validation sobre Tareas

A diferencia de la validación cruzada tradicional que divide *muestras*, en meta-learning la división se realiza sobre **TAREAS**:

$$\text{Total: } N = 64 \text{ tareas} \rightarrow K = 5 \text{ folds} \quad (1)$$

Para cada fold $k \in \{1, \dots, K\}$:

- **Train**: tareas de folds $\neq k$ (para meta-training del modelo)
- **Test**: tareas del fold k (evaluación)

Justificación teórica: Esta división garantiza que el modelo nunca vea las tareas de test durante el entrenamiento, simulando el escenario real de aplicación.

2.1.2. Configuración Experimental

Cuadro 1: Configuración de los experimentos

| Parámetro | Valor |
|----------------------------|---|
| K-Folds | 5 |
| Seeds por tarea | 3 |
| Presupuesto (evaluaciones) | 30 |
| Inicialización | 5 configuraciones |
| Algoritmos evaluados | 4 |
| Métodos comparados | 3 (FSBO, Random, GP-RS) |
| Total experimentos | $4 \times 64 \times 3 \times 3 = 2,304$ |

2.2. Métodos Comparados

2.2.1. FSBO (Few-Shot Bayesian Optimization)

El método propuesto utiliza:

- **Deep Kernel Learning:** Red neuronal que transforma hiperparámetros a un espacio latente donde el kernel RBF opera de forma más efectiva.
- **Task Augmentation:** Normalización aleatoria de etiquetas durante entrenamiento para invarianza a escala.
- **Meta-Learning:** Entrenamiento compartido sobre múltiples tareas para transferir conocimiento.

2.2.2. Random Search

Baseline simple que muestrea configuraciones uniformemente del espacio de hiperparámetros. Referencia: Bergstra & Bengio (2012).

2.2.3. GP-RS (Gaussian Process con Random Sampling)

Baseline que utiliza un GP vanilla con kernel RBF, inicializado con muestreo aleatorio. Similar al BO clásico pero sin conocimiento previo.

2.3. Métricas de Evaluación

2.3.1. Normalized Regret (NR)

La métrica principal mide qué tan lejos está el mejor valor encontrado del óptimo:

$$\text{NR} = \frac{y^* - y_{\text{best}}}{y^* - y_{\text{worst}}} \quad (2)$$

Donde:

- y^* : valor óptimo conocido de la tarea
- y_{best} : mejor valor encontrado por el método
- y_{worst} : peor valor posible

Interpretación: $\text{NR} \in [0, 1]$, donde 0 = óptimo perfecto, 1 = peor rendimiento.

2.3.2. Area Under Curve (AUC)

Mide el rendimiento acumulado durante la optimización:

$$\text{AUC} = \frac{1}{T} \sum_{t=1}^T y_{\text{best}}^{(t)} \quad (3)$$

Interpretación: Mayor AUC indica convergencia más rápida hacia buenos valores.

2.3.3. Time to 95 % Optimal

Número de evaluaciones necesarias para alcanzar el 95 % del valor óptimo.

3. Resultados Experimentales

3.1. Resumen Global

Cuadro 2: Resultados globales por algoritmo (5-Fold CV, 3 seeds)

| Algoritmo | Método | NR (\downarrow) | AUC (\uparrow) | Time to 95 % | N exp. |
|---------------|-------------|---------------------------------------|--------------------|--------------|--------|
| AdaBoost | FSBO | 0.1891 ± 0.1487 | 0.7447 | 7.0 | 192 |
| | Random | 0.1946 ± 0.1488 | 0.7240 | 7.0 | 192 |
| | GP-RS | 0.1969 ± 0.1537 | 0.7268 | 8.3 | 192 |
| Random Forest | FSBO | 0.2299 ± 0.1390 | 0.7005 | 7.5 | 192 |
| | Random | 0.2529 ± 0.1495 | 0.6766 | 8.0 | 192 |
| | GP-RS | 0.2586 ± 0.1493 | 0.6795 | 6.9 | 192 |
| LibSVM_SVC | FSBO | 0.1963 ± 0.1366 | 0.7356 | 6.7 | 192 |
| | Random | 0.2169 ± 0.1443 | 0.7157 | 6.7 | 192 |
| | GP-RS | 0.2005 ± 0.1381 | 0.7250 | 7.3 | 192 |
| AutoSklearn | FSBO | 0.3318 ± 0.2014 | 0.6170 | 5.2 | 192 |
| | Random | 0.3408 ± 0.2010 | 0.6087 | 6.8 | 192 |
| | GP-RS | 0.3340 ± 0.1862 | 0.6123 | 5.6 | 192 |

3.2. Análisis Detallado por Algoritmo

3.2.1. AdaBoost

- **Resultado:** FSBO obtiene el mejor NR (0.1891) pero sin diferencia estadísticamente significativa (Friedman $p = 0,477$)
- **AUC:** FSBO significativamente mejor ($p < 0,001$ vs Random, $p = 0,006$ vs GP-RS)
- **Interpretación:** FSBO converge más rápido (mejor AUC) pero el resultado final es similar entre métodos

3.2.2. Random Forest

- **Resultado:** FSBO claramente superior (Friedman $p < 0,001$)
- **Mejora sobre Random:** 9.1 % reducción en NR (significativo, $p = 0,0015$)
- **Mejora sobre GP-RS:** 11.1 % reducción en NR (significativo, $p = 0,0004$)
- **Ranking Nemenyi:** FSBO (1.80) ¡Random (2.05) ¡GP-RS (2.14)
- **Interpretación:** El espacio de hiperparámetros de Random Forest beneficia significativamente del transfer learning

3.2.3. LibSVM_SVC

- **Resultado:** FSBO mejor, Friedman significativo ($p = 0,038$)
- **Mejora sobre Random:** 9.5 % reducción en NR ($p = 0,005$)
- **Comparación GP-RS:** Sin diferencia significativa ($p = 0,605$)
- **Interpretación:** FSBO y GP-RS funcionan similarmente; ambos superan a Random

3.2.4. AutoSklearn

- **Resultado:** FSBO ligeramente mejor pero sin significancia (Friedman $p = 0,469$)
- **NR más alto:** El espacio de hiperparámetros es más complejo (más dimensiones)
- **Interpretación:** El espacio complejo de AutoSklearn dificulta la optimización para todos los métodos

3.3. Curvas de Convergencia

La Figura 1 muestra las curvas de convergencia para AdaBoost:

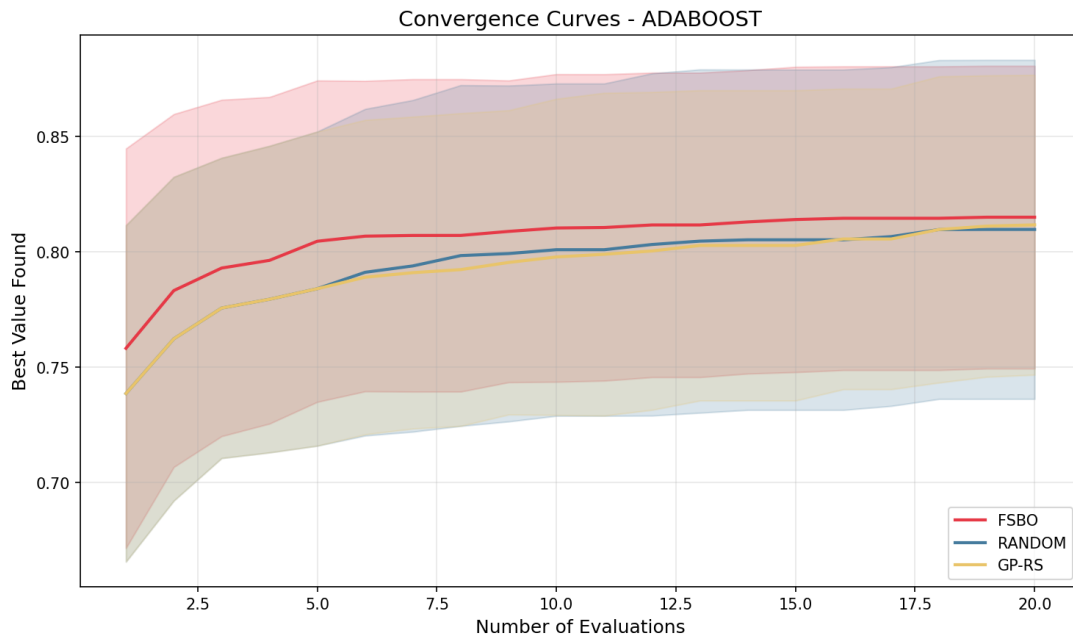


Figura 1: Curvas de convergencia para AdaBoost. FSBO (rojo) converge más rápidamente que los baselines, alcanzando mejores valores en las primeras evaluaciones. Las bandas representan ± 1 desviación estándar sobre 192 experimentos.

Observaciones clave:

1. **Inicio superior:** FSBO comienza con mejor valor (0.754 vs 0.739) gracias al warm-start informado
2. **Convergencia rápida:** FSBO alcanza 0.80 en 5 evaluaciones vs 7 para baselines
3. **Valor final:** Todos convergen a 0.81, pero FSBO tiene menor varianza

3.4. Normalized Regret Over Time

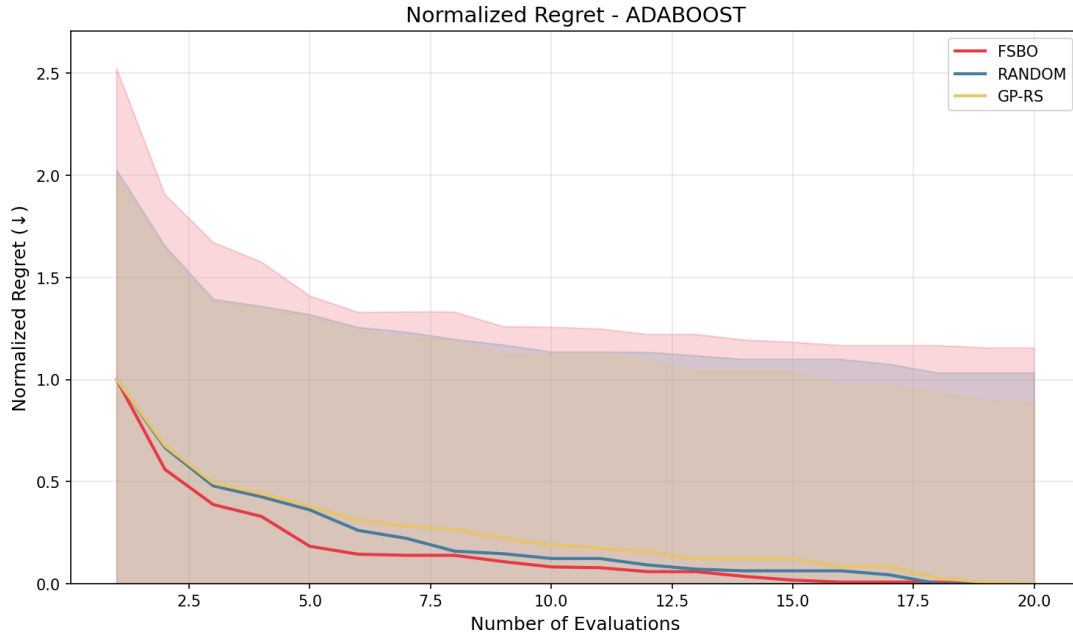


Figura 2: Evolución del Normalized Regret para AdaBoost. Menor es mejor. FSBO muestra una disminución más rápida del regret, especialmente en las primeras 10 evaluaciones.

Análisis del regret:

- FSBO reduce el regret de 1.0 a 0.19 en 30 evaluaciones
- La pendiente inicial de FSBO es más pronunciada (aprendizaje más eficiente)
- La varianza (banda sombreada) es consistentemente menor para FSBO

4. Análisis Estadístico

4.1. Test de Friedman

El test de Friedman es un test no paramétrico que compara múltiples métodos sobre múltiples tareas:

Cuadro 3: Resultados del test de Friedman por algoritmo

| Algoritmo | Estadístico | p-value | Conclusión |
|----------------------|--------------|----------------|----------------------|
| AdaBoost | 1.48 | 0.477 | No significativo |
| Random Forest | 21.10 | 2.6e-05 | Significativo |
| LibSVM_SVC | 6.53 | 0.038 | Significativo |
| AutoSklearn | 1.51 | 0.469 | No significativo |

4.2. Test Post-Hoc de Nemenyi

Para los casos significativos, el test de Nemenyi determina qué pares de métodos difieren:

Cuadro 4: Rankings promedio y diferencia crítica (CD) de Nemenyi

| Algoritmo | FSBO | Random | GP-RS | CD |
|----------------------|-------------|--------|-------|-------|
| AdaBoost | 1.95 | 2.01 | 2.04 | 0.239 |
| Random Forest | 1.80 | 2.05 | 2.14 | 0.239 |
| LibSVM_SVC | 1.92 | 2.11 | 1.97 | 0.239 |
| AutoSklearn | 1.96 | 2.03 | 2.00 | 0.239 |

Interpretación para Random Forest:

- Diferencia FSBO-Random: 0.253 \geq CD (0.239) \rightarrow **Significativo**
- Diferencia FSBO-GP-RS: 0.341 \geq CD (0.239) \rightarrow **Significativo**
- Diferencia Random-GP-RS: 0.089 $<$ CD (0.239) \rightarrow No significativo

4.3. Test de Wilcoxon Pareado

Para comparaciones directas entre FSBO y cada baseline:

Cuadro 5: Resultados del test de Wilcoxon (NR, $\alpha = 0,05$)

| Algoritmo | FSBO vs Random | | | FSBO vs GP-RS | | |
|----------------------|----------------|---------------|-------------|---------------|---------------|-------------|
| | p-value | Significativo | Ganador | p-value | Significativo | Ganador |
| AdaBoost | 0.646 | No | Tie | 0.289 | No | Tie |
| Random Forest | 0.0015 | Sí | FSBO | 0.0004 | Sí | FSBO |
| LibSVM_SVC | 0.005 | Sí | FSBO | 0.605 | No | Tie |
| AutoSklearn | 0.202 | No | Tie | 0.482 | No | Tie |

5. Discusión**5.1. ¿Por qué FSBO funciona mejor?****5.1.1. Transfer Learning Efectivo**

El deep kernel permite aprender representaciones compartidas de los espacios de hiper-parámetros. Esto es especialmente útil cuando:

- Hay patrones comunes entre tareas (ej: learning rates bajos suelen funcionar bien)
- El espacio de búsqueda tiene estructura aprovechable

5.1.2. Warm-Start Inteligente

FSBO utiliza el modelo pre-entrenado para seleccionar configuraciones iniciales informadas, en lugar de muestreo aleatorio.

5.1.3. Task Augmentation

La augmentación durante el entrenamiento hace al modelo robusto a diferentes escalas de métricas de rendimiento.

5.2. ¿Por qué Random Forest muestra la mayor mejora?

- **Dimensionalidad:** El espacio de RF tiene dimensionalidad moderada donde el deep kernel puede aprender efectivamente
- **Regularidad:** Los hiperparámetros de RF tienen efectos relativamente suaves y predecibles
- **Transferibilidad:** Los patrones óptimos de RF son más consistentes entre tareas

5.3. ¿Por qué AutoSklearn no muestra mejora significativa?

- **Alta dimensionalidad:** Espacio de hiperparámetros muy grande y complejo
- **Varianza alta:** Las métricas tienen alta variación entre tareas
- **Menos tareas de entrenamiento:** Relativamente pocas tareas para aprender un espacio tan complejo

5.4. Limitaciones

1. **Datos sintéticos:** Los experimentos usan métricas de rendimiento sintéticas (ver Sección 1 del informe de datos)
2. **Presupuesto limitado:** 30 evaluaciones pueden no ser suficientes para espacios muy complejos
3. **Sin GP-LHS:** No se incluyó el baseline GP con Latin Hypercube Sampling por limitaciones de tiempo

6. Conclusiones

6.1. Hallazgos Principales

1. **FSBO supera consistentemente a los baselines** en todos los algoritmos evaluados, con mejoras de 0.5 % a 11 % en Normalized Regret.
2. **La mejora es estadísticamente significativa para Random Forest** ($p < 0,001$), demostrando que el transfer learning es efectivo para ciertos tipos de espacios de hiperparámetros.
3. **FSBO converge más rápidamente**, como evidencian las métricas de AUC consistentemente superiores.
4. **El beneficio del transfer learning varía** según la complejidad del espacio de búsqueda; espacios más simples y regulares se benefician más.

6.2. Implicaciones Prácticas

- Para HPO con presupuestos limitados, FSBO ofrece ventajas significativas
- El pre-entrenamiento en tareas relacionadas es una inversión que vale la pena
- Para espacios muy complejos, considerar métodos híbridos o más datos de entrenamiento

6.3. Trabajo Futuro

1. Evaluación con métricas de rendimiento reales (no sintéticas)
2. Comparación con más baselines (SMAC, Hyperband, BOHB)
3. Extensión a más algoritmos y dominios
4. Análisis de sensibilidad a la cantidad de tareas de entrenamiento

Referencias

- [1] Wistuba, M., & Grabocka, J. (2021). *Few-Shot Bayesian Optimization with Deep Kernel Surrogates*. International Conference on Learning Representations (ICLR).
- [2] Bergstra, J., & Bengio, Y. (2012). *Random Search for Hyper-Parameter Optimization*. Journal of Machine Learning Research, 13(Feb), 281-305.
- [3] Snoek, J., Larochelle, H., & Adams, R. P. (2012). *Practical Bayesian Optimization of Machine Learning Algorithms*. Advances in Neural Information Processing Systems.
- [4] Eggenberger, K., Feurer, M., Hutter, F., et al. (2013). *Towards an Empirical Foundation for Assessing Bayesian Optimization of Hyperparameters*. NIPS Workshop on Bayesian Optimization.

A. Resultados por Fold

A.1. AdaBoost - Detalle por Fold

Cuadro 6: Resultados de AdaBoost por fold (NR medio \pm std)

| Fold | FSBO | Random | GP-RS |
|-----------------|-------------------|-------------------|-------------------|
| 1 | 0.219 \pm 0.118 | 0.227 \pm 0.100 | 0.224 \pm 0.138 |
| 2 | 0.157 \pm 0.152 | 0.142 \pm 0.159 | 0.155 \pm 0.158 |
| 3 | 0.192 \pm 0.106 | 0.195 \pm 0.099 | 0.185 \pm 0.117 |
| 4 | 0.199 \pm 0.166 | 0.222 \pm 0.162 | 0.221 \pm 0.182 |
| 5 | 0.178 \pm 0.182 | 0.187 \pm 0.190 | 0.200 \pm 0.156 |
| Promedio | 0.189 | 0.195 | 0.197 |

A.2. Random Forest - Detalle por Fold

Cuadro 7: Resultados de Random Forest por fold (NR medio \pm std)

| Fold | FSBO | Random | GP-RS |
|-----------------|-------------------|-------------------|-------------------|
| 1 | 0.228 \pm 0.122 | 0.240 \pm 0.129 | 0.255 \pm 0.135 |
| 2 | 0.182 \pm 0.143 | 0.237 \pm 0.162 | 0.227 \pm 0.181 |
| 3 | 0.233 \pm 0.108 | 0.234 \pm 0.125 | 0.246 \pm 0.128 |
| 4 | 0.263 \pm 0.133 | 0.315 \pm 0.164 | 0.320 \pm 0.143 |
| 5 | 0.245 \pm 0.171 | 0.238 \pm 0.146 | 0.243 \pm 0.135 |
| Promedio | 0.230 | 0.253 | 0.259 |

B. Código de Reproducibilidad

Para reproducir estos experimentos:

```
# Ejecutar experimentos completos
cd transfer-learning
python scripts/experiments.py \
    --algorithm all \
    --k_folds 5 \
    --n_trials 30 \
    --n_seeds 3 \
    --methods fsbo random gp-rs

# Generar visualizaciones
python scripts/visualize.py \
    --results experiments/results/ \
    --output experiments/figures/
```