



Módulo V - Aprendizaje Supervisado

Clase 16: Validación cruzada





¿Ponemos a grabar el
taller?

OBJETIVOS de hoy



- REPASO de temáticas anteriores

- Validación Cruzada



**REPASEMOS retomando clases
anteriores, ¿Dudas?**



VALIDACIÓN CRUZADA

Validación Cruzada: K-Fold

- La validación cruzada es un procedimiento de **remuestreo** que se utiliza para evaluar modelos de aprendizaje automático en una muestra de datos limitada.
- El procedimiento tiene un único parámetro llamado **k** que se refiere al número de grupos en que se dividirá una muestra de datos dada.



Validación Cruzada: K-Fold

Iteration 1	Test	Train	Train	Train	Train
Iteration 2	Train	Test	Train	Train	Train
Iteration 3	Train	Train	Test	Train	Train
Iteration 4	Train	Train	Train	Test	Train
Iteration 5	Train	Train	Train	Train	Test

Validación Cruzada: K-Fold

Pasos para realizar la validación cruzada:

- 1) Desordenar los datos
- 2) Separar en K folds del mismo tamaño
- 3) Para cada fold que separamos:
 - Elegir un fold como Test set, y los K-1 folds restantes como Train set.
 - Entrenar y evaluar el modelo.
 - Guardar el resultado de la evaluación y descartar el modelo.
- 4) Obtener una medida de performance del modelo como el promedio de las K evaluaciones obtenidas en (3). También es una buena práctica incluir una medida de la varianza de las métricas obtenidas.



Validación Cruzada: K-Fold

- Permite tener una estimación más realista de la performance del modelo.
- Generalmente resulta en una estimación menos sesgada de la habilidad del modelo para predecir.
- El valor de K no es trivial. Una mala elección puede resultar una mala evaluación del desempeño.

- Ejemplo para $k=5$:

Entrenamiento Validación

Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	→ Modelo 1	→ Performance 1	} Promedio
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	→ Modelo 2	→ Performance 2	
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	→ Modelo 3	→ Performance 3	
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	→ Modelo 4	→ Performance 4	
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	→ Modelo 5	→ Performance 5	





Descanso

Nos vemos en 10 minutos



**¿Cómo encuentro la óptima
combinación de atributos +
algoritmos + hiper
parámetros ?**

Optimización

Un **MODELO** es un *pipeline* que se conforma de:

- Distintos atributos (selección y transformación de atributos)
- Distintos algoritmos (árboles, KNN, Linear Regression, etc)
- Distintos **hiper parámetros** de cada algoritmo.



Optimización

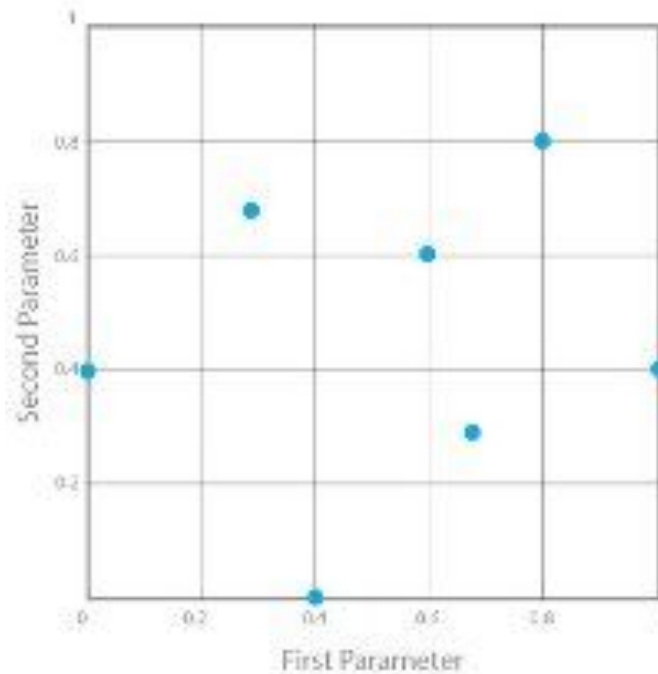
- 1) Definir el espacio de hiper parámetros
- 2) Explorar el espacio de búsqueda, definiendo un modelo para cada combinación posible de hiper parámetros. Para cada uno, entrenarlo y evaluar su desempeño.
- 3) Elegir la combinación con mejor desempeño, y entrenar el **único modelo** definitivo usando todos los datos.



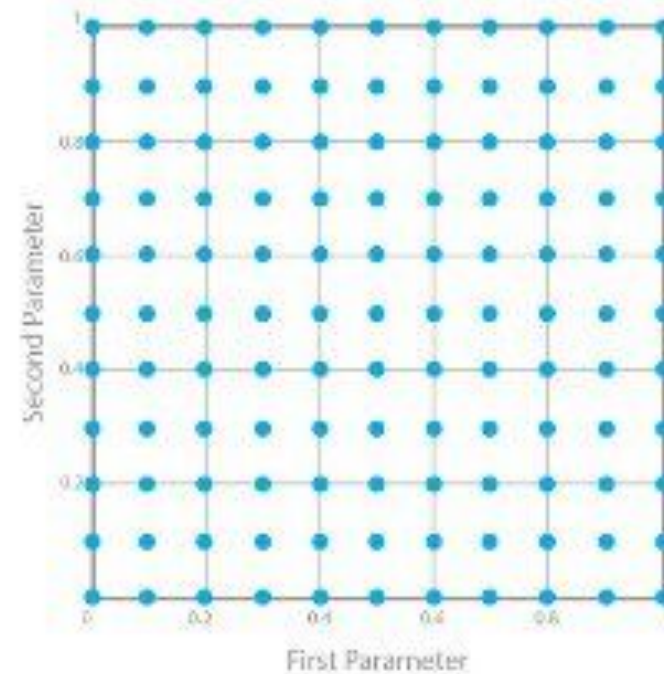
Optimización: Gridsearch

Técnica para buscar la mejor combinación de atributos + algoritmos + hiperparámetros de manera tal que el modelo sea el óptimo.

Manual Search



Grid Search



Validación Cruzada y Optimización: GridSearchCV

```
from sklearn.model_selection import GridSearchCV
```

```
parameters_optimize = {  
    'max_features': ['auto', 'sqrt', 'log2', None],  
    'max_depth': [2, 3, 4],  
    'criterion': ['gini', 'entropy'],  
    'bootstrap': [True, False],  
    'n_estimators': [2, 5, 10, 15, 20]  
}
```

```
random_forest_hyp = RandomForestClassifier()  
  
random_forest_search = GridSearchCV(random_forest_hyp,  
                                     cv = 20,  
                                     param_grid = parameters_optimize,  
                                     n_jobs = 3)  
random_forest_search.fit(X_train_imp, Y_train);
```

```
print('The best parameters after GridSearchCV', random_forest_search.best_params_)
```

The best parameters after GridSearchCV {'bootstrap': True, 'criterion': 'gini', 'max_depth': 4, 'max_features': 'auto', 'n_estimators': 20}



Repasamos en Kahoot



¿Dudas?



Descanso

Nos vemos en 10 minutos



Sección práctica:

**Métricas y overfitting con la
Notebook 18**

Trabajamos en salas de zoom

Métricas y overfitting

Trabajaremos con la Notebook 18

En los grupos establecidos,
aprendemos cómo evaluar los
modelos y evitar el overfitting
optimizando los
hiperparametros



50 minutos de actividad





REPASAMOS

**Revisamos los conceptos y código
trabajados en la notebook 18**



¿Dudas?

FUNDACIÓN
YPF

¡Muchas gracias!

