

Machine Learning Workshop

Mentor: Nicolas Käenzig

Email: nkaenzig@gmail.com

Workshop Repository: <https://github.com/nkaenzig/ml-workshop>

Contenido

Modulo 1

- Introducción ML
- Python crashcourse

Modulo 2

- Análisis de datos
- Preprocesamiento de datos
- Ejemplo ML

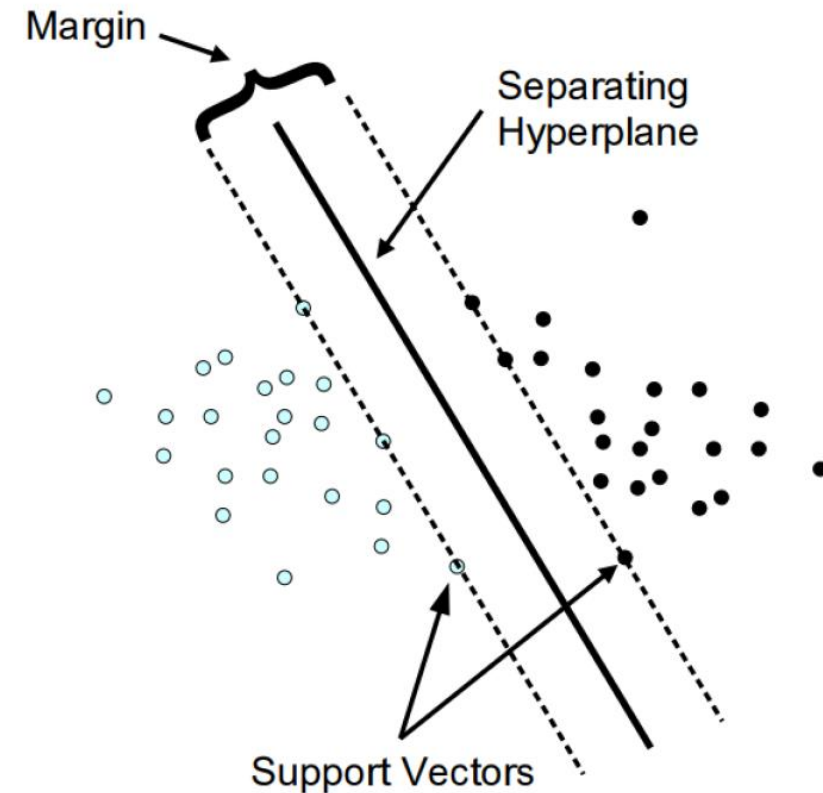
Modulo 3

- Modelos de ML
- Técnicas de evaluación
- Ejemplos ML

Machine Learning – Modelos

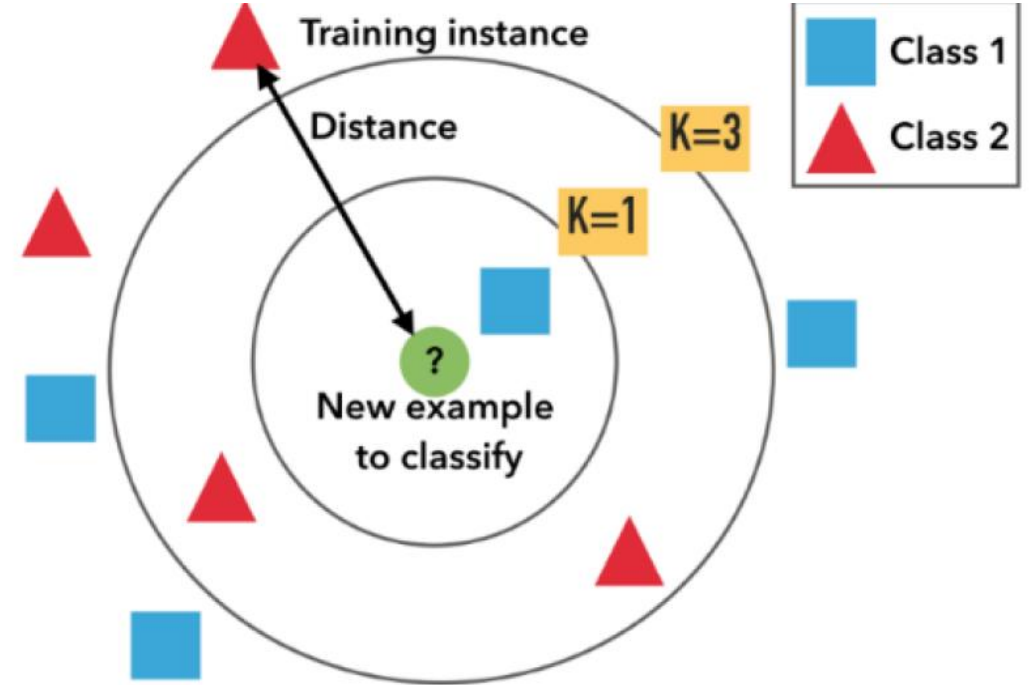
SVM (Support Vector Machine)

- Supervised Classification
- Scikit-Learn:
 - `sklearn.svm.SVC()`
- Time Complexity (Training):
 - Linear SVM: $O(n)$
 - Non-Linear SVM: $O(n^2)$ - $O(n^3)$



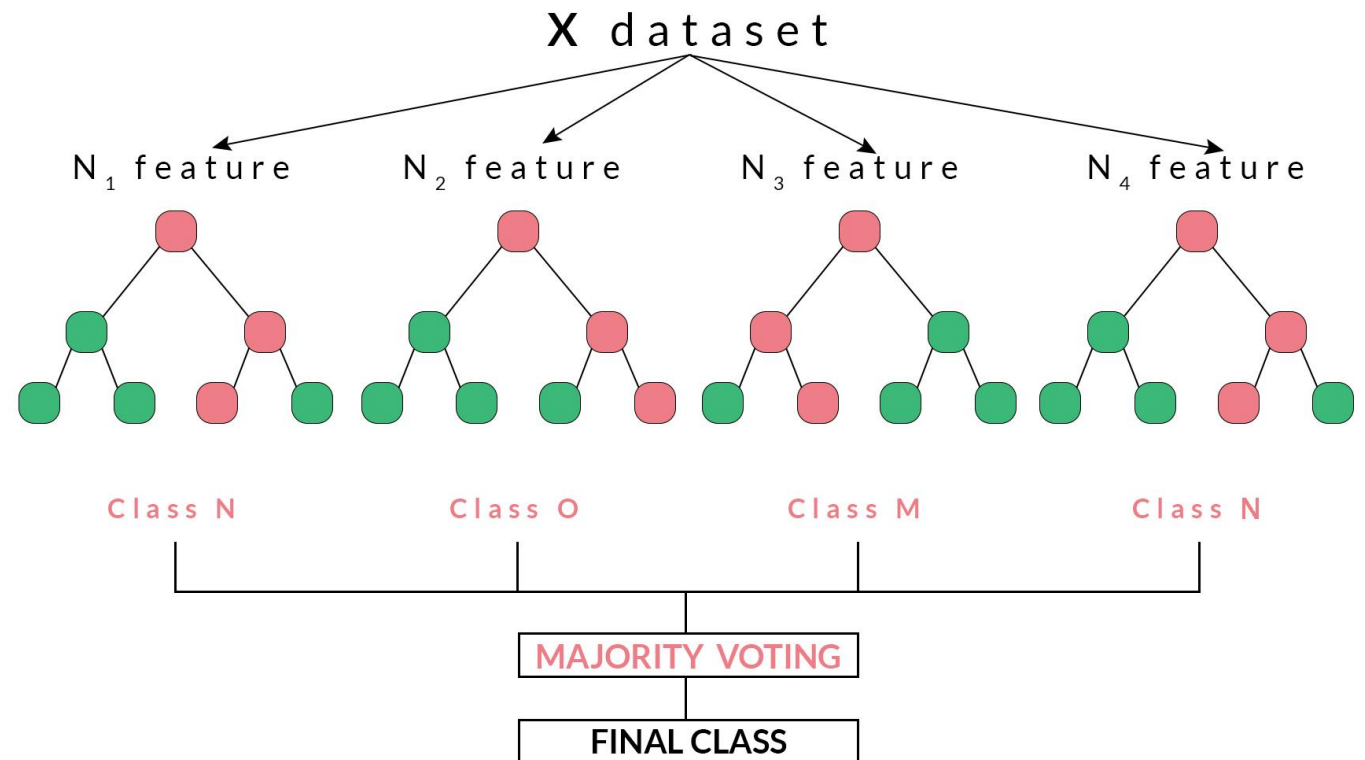
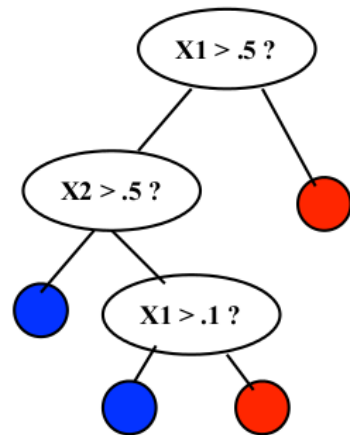
K-NN (K-Nearest Neighbors)

- Supervised Classification
- Scikit-Learn:
 - `sklearn.neighbors.KNeighborsClassifier()`
- Time Complexity (Training):
 - $O(n)$



Random Forest

- Supervised Classification
- Scikit-Learn:
 - `sklearn.ensemble.RandomForestClassifier`
- Time Complexity (Training):
 - $O(n \log(n))$



Ridge Regression

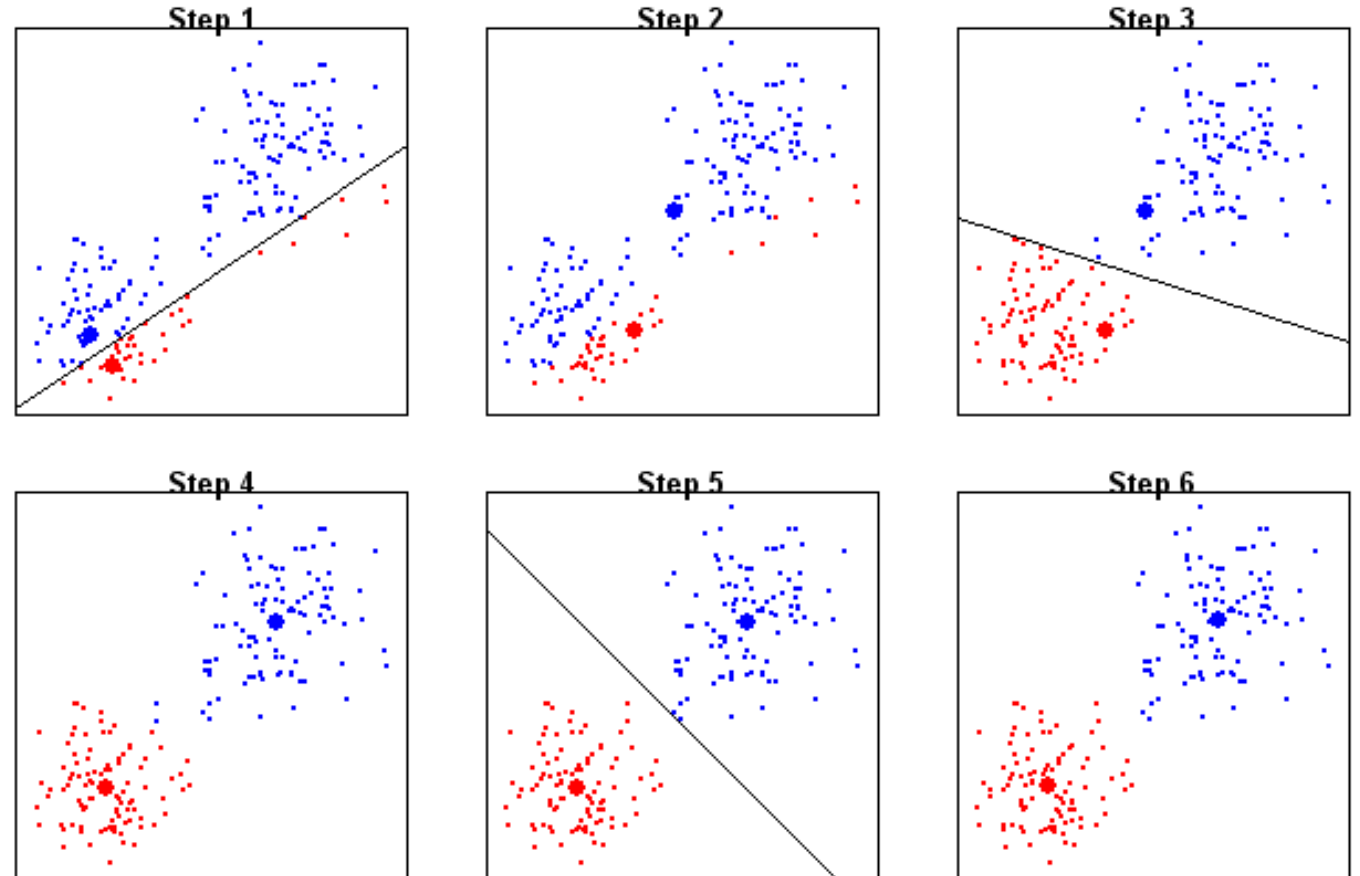
- Supervised Regression
- Scikit-Learn:
 - `sklearn.linear_model.Ridge`
- Time Complexity (Training):
 - **$O(n)$**

$$f(x, \theta) = \theta_0 + \theta_1 x$$

$$\theta^* = \underset{\theta}{\operatorname{argmin}} |f(x, \theta) - y|^2 + \alpha ||\theta||_2^2$$

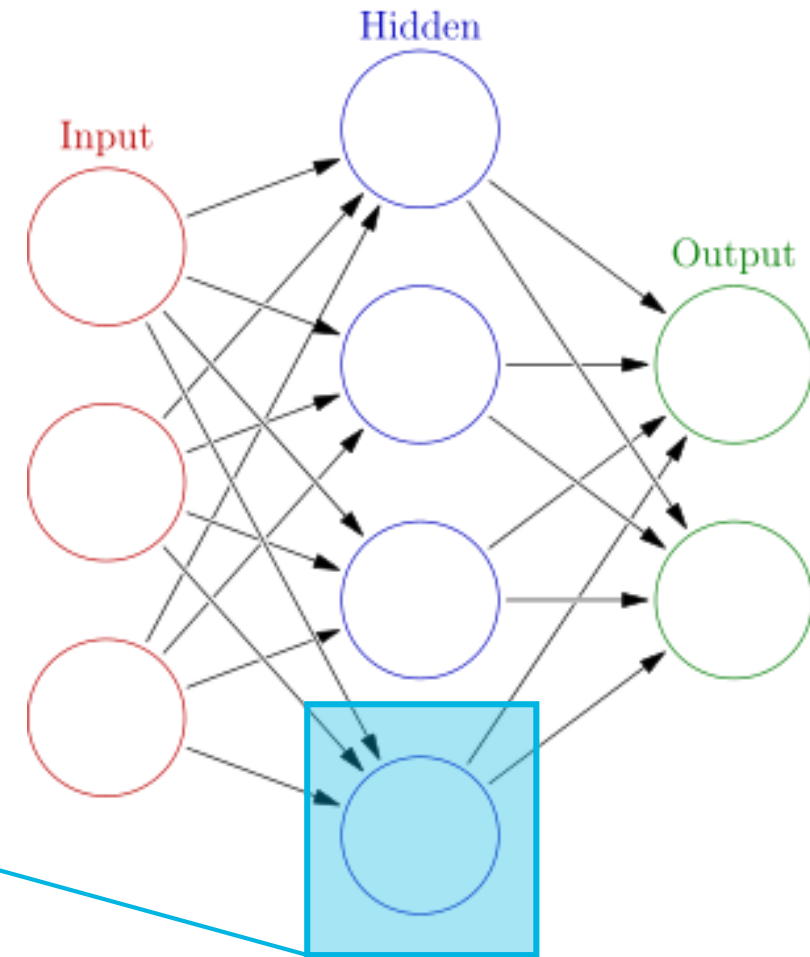
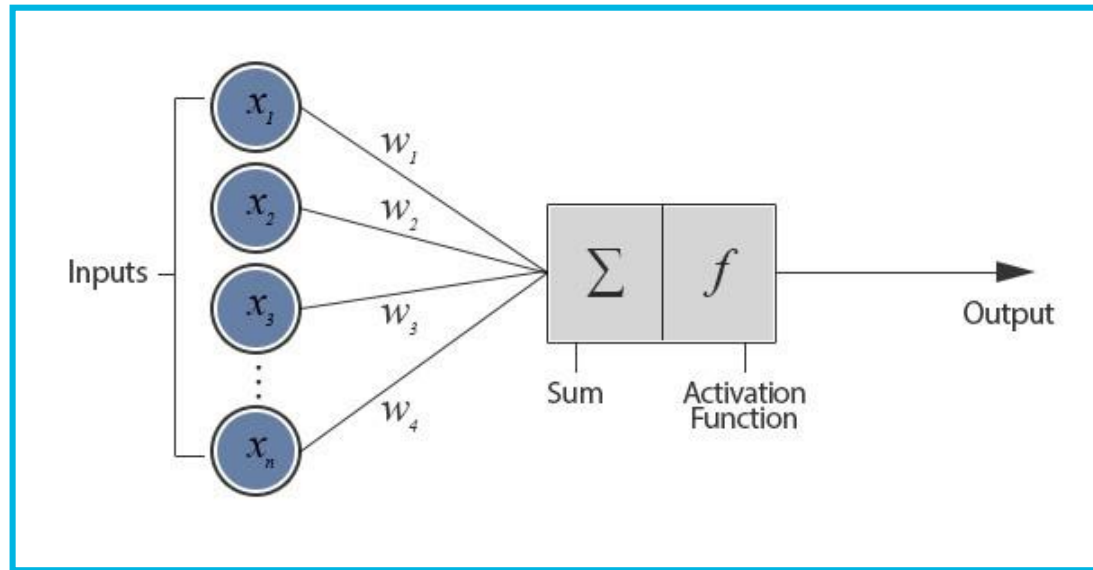
K-Means

- Unsupervised Clustering
- Scikit-Learn:
 - `sklearn.cluster.KMeans()`
- Time Complexity (Training):
 - $O(n)$



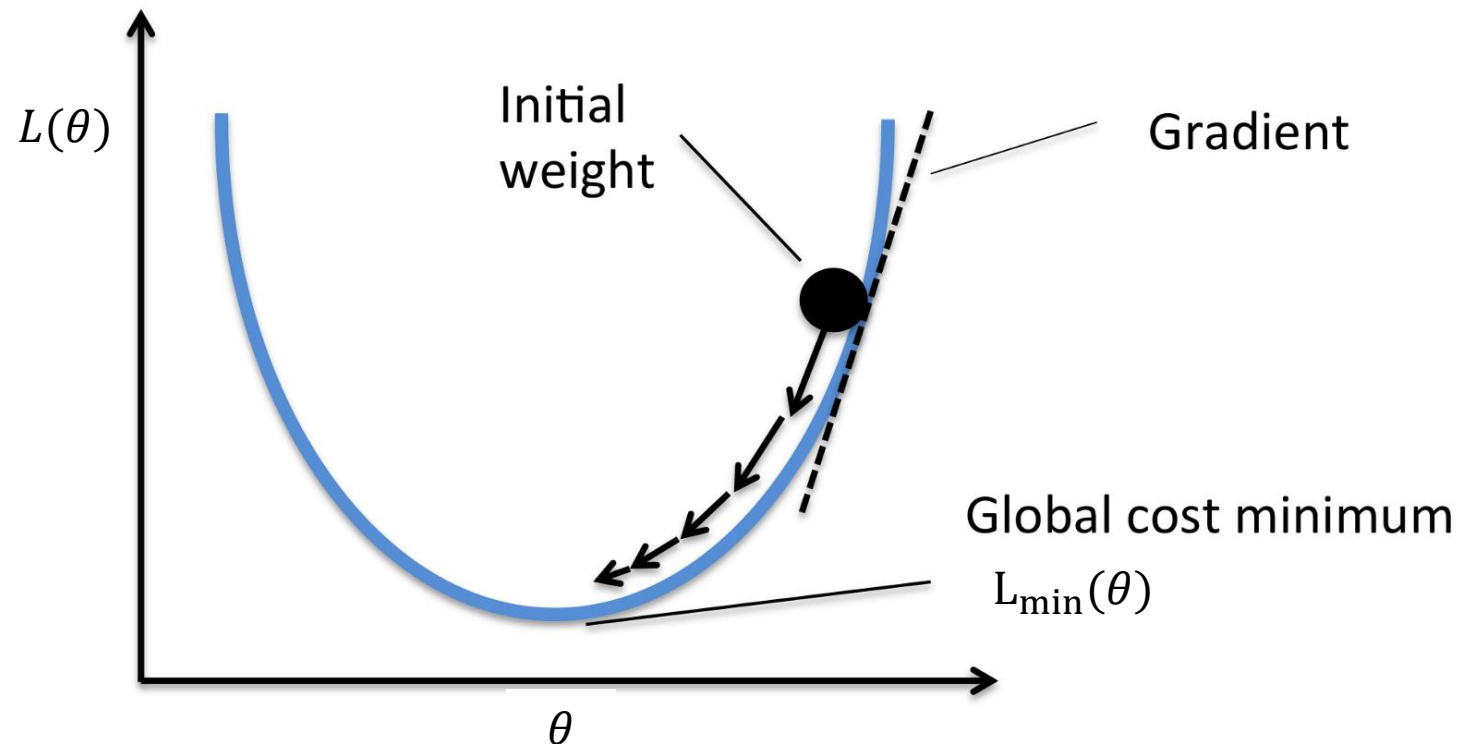
Artificial Neural Networks

Neurona artificial



Como funciona el Training?

- **Objetivo:** $\theta^* = \operatorname{argmin}_{\theta} L(x, y, \theta)$
- **Gradient Descent**



Algoritmos / Modelos

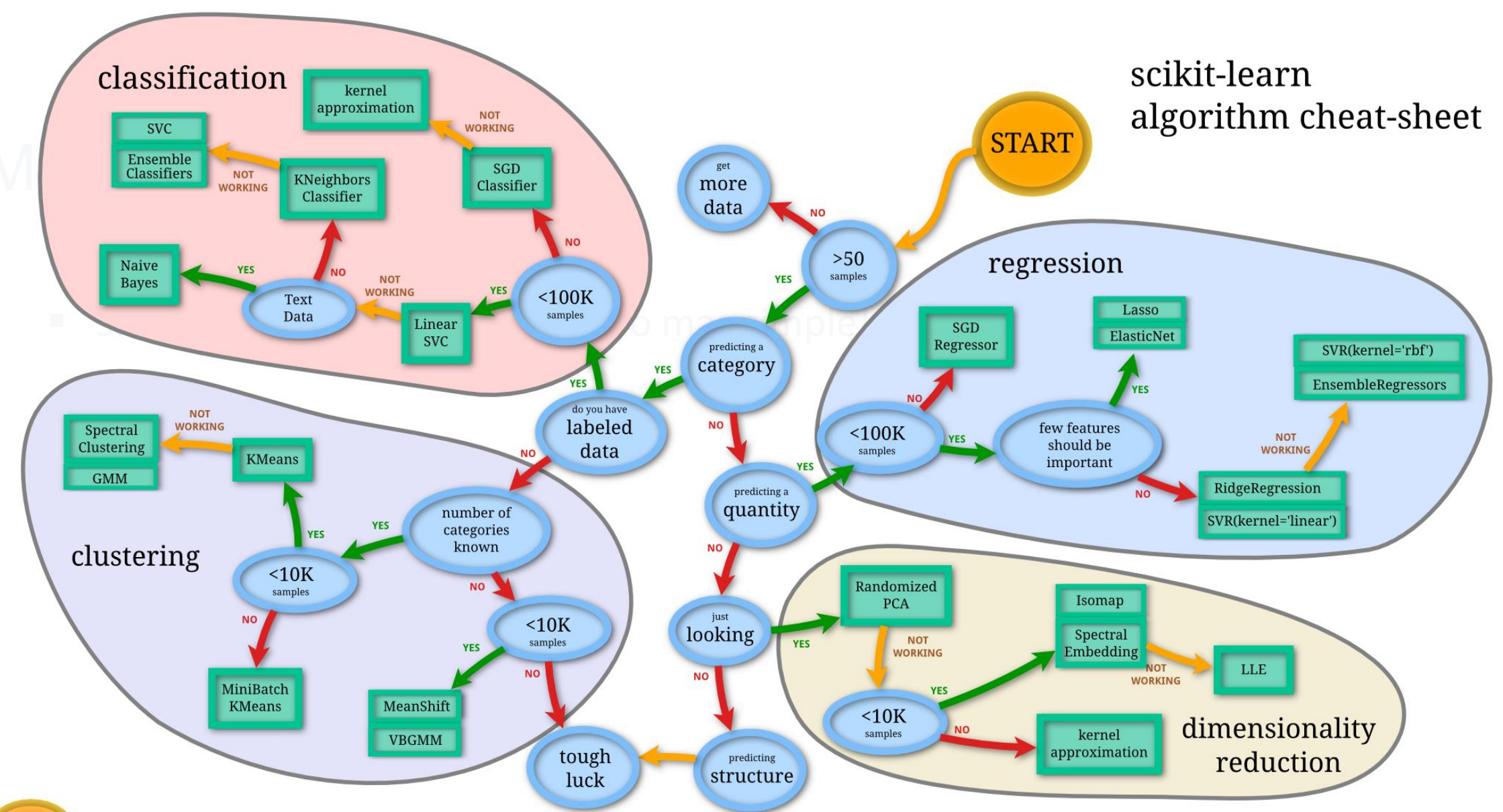
- Los modelos y algoritmos de optimización ya están implementados!
- Machine Learning:



- Deep Learning



scikit-learn algorithm cheat-sheet



Criterios para seleccionar un modelo

- Consejo: Siempre empieza con el modelo mas simple / fácil para usar
- **Preguntas:**
 - El modelo requiere standardization?
 - El modelo requiere one-hot encoding para los features categóricos?
 - El modelo requiere quitar columnas redundantes?
 - Que es la complejidad del modelo?
 - El modelo es lineal o no lineal?
 - Es difícil encontrar buenos parámetros para el modelo?
 - El modelo funciona con “imbalanced” datasets?

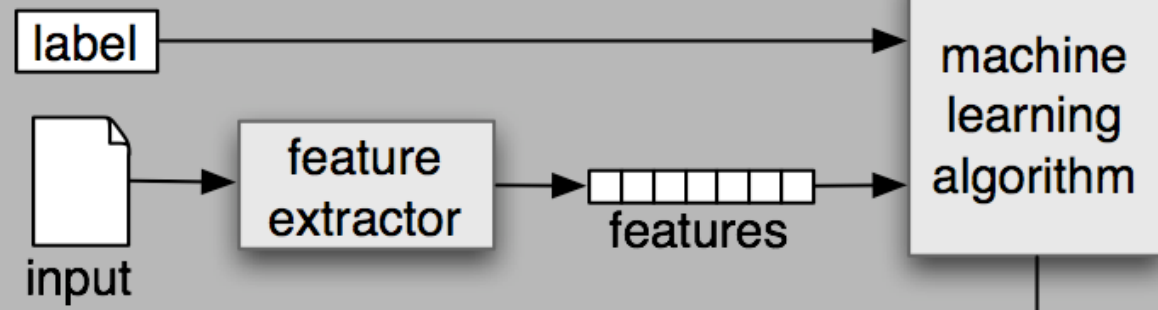
Criterios para seleccionar un modelo

Model	Standardization	Class Balancing	One-Hot Encoding	Non-Linear	Complexity
LinearSVM	yes	yes	yes	no	$O(n)$
SVM('rbf')	yes	yes	yes	yes	$O(n^2) - O(n^3)$
K-NN	yes	yes	yes	yes	$O(n)$
Random Forest	no	no	no	yes	$O(n \log(n))$

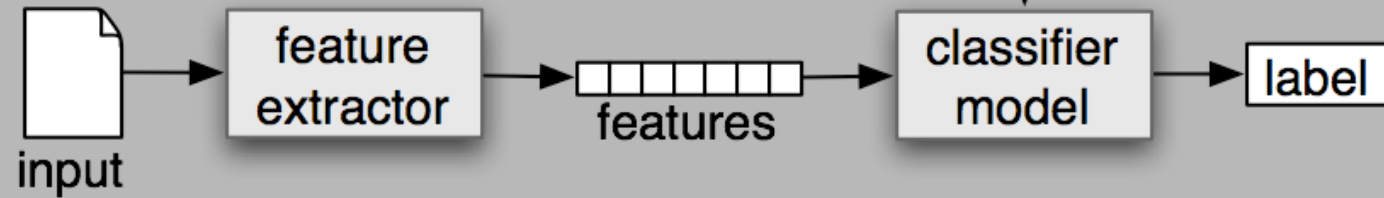
- Siempre empieza con el modelo mas simple / fácil para usar
→ Empezar con Random Forest siempre es una buena idea!

Machine Learning – Evaluación

(a) Training



(b) Prediction



Como medir la calidad de las predicciones?

- Clasificación:

$$\textit{Accuracy} = \frac{\# \textit{Correct predictions}}{\# \textit{Predictions}}$$

$$\textit{Precision} = \frac{\# \textit{True positives}}{\# \textit{True positives} + \# \textit{False Positives}}$$

$$\textit{Recall} = \frac{\# \textit{True positives}}{\# \textit{True positives} + \# \textit{False Negatives}}$$

Problemas con accuracy

- Dataset:
 - 1000 samples de pacientes **sin** cancer (N)
 - 5 samples de pacientes **con** cancer (P)
- Modelo solamente diagnostica 1 de los 5 pacientes con cancer

$$Accuracy = \frac{\# \text{ Correct predictions}}{\# \text{ Predictions}} = ?$$

$$Precision = \frac{\# \text{ True positives}}{\# \text{ True positives} + \# \text{ False Positives}} = ?$$

$$Recall = \frac{\# \text{ True positives}}{\# \text{ True positives} + \# \text{ False Negatives}} = ?$$

Problemas con accuracy

- Dataset:
 - 1000 samples de pacientes **sin** cancer (N)
 - 5 samples de pacientes **con** cancer (P)
- Modelo solamente diagnostica 1 de los 5 pacientes con cancer

$$Accuracy = \frac{\# \text{ Correct predictions}}{\# \text{ Predictions}} = \frac{1001}{1005} \% = 99.6\%$$

$$Precision = \frac{\# \text{ True positives}}{\# \text{ True positives} + \# \text{ False Positives}} = \frac{1}{1} \% = 100\%$$

$$Recall = \frac{\# \text{ True positives}}{\# \text{ True positives} + \# \text{ False Negatives}} = \frac{1}{5} \% = 20\%$$

Como medir la calidad de las predicciones?

- **Regresión:**

- Mean Absolute Error

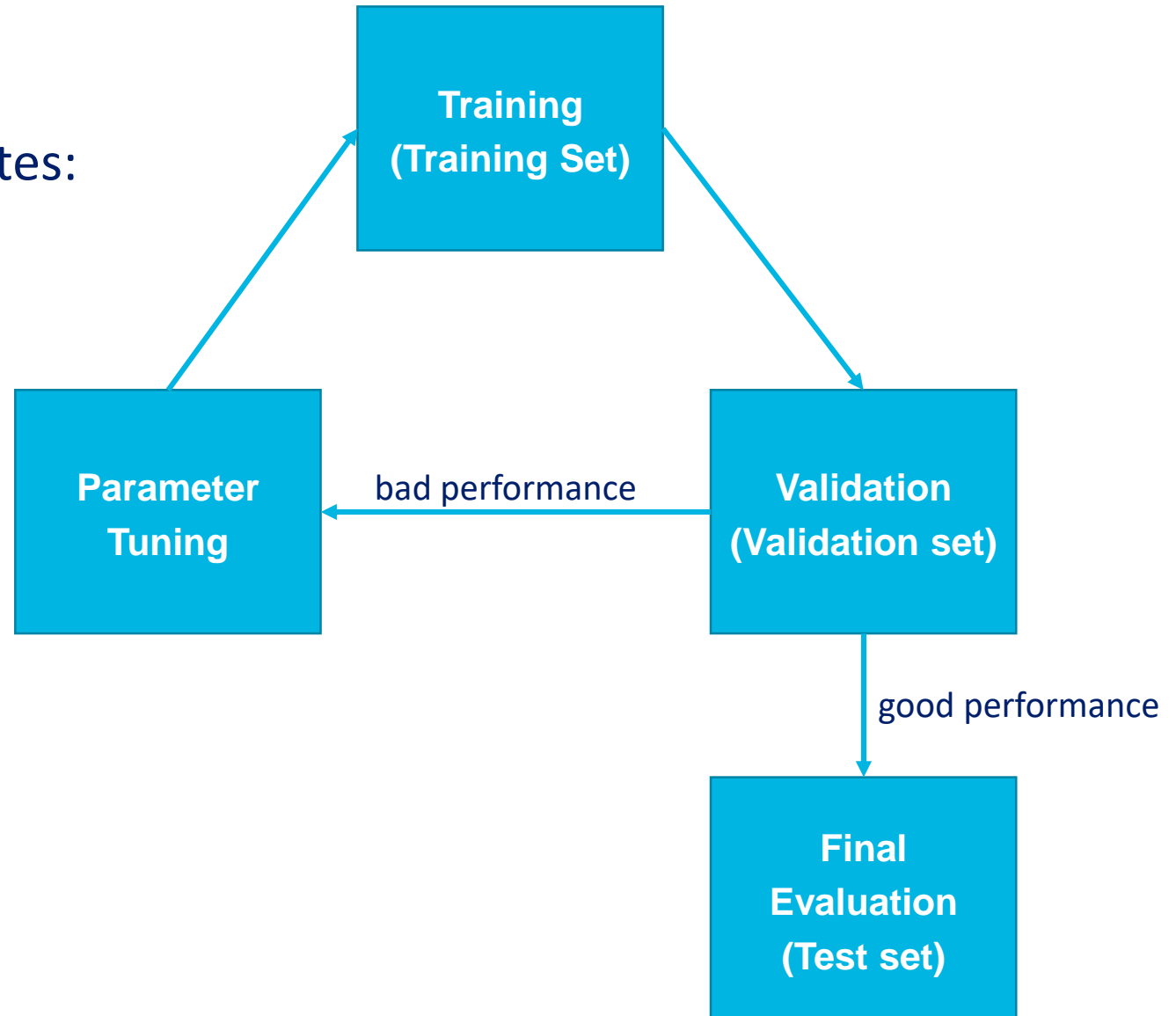
$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$$

- Root Mean Squared Error

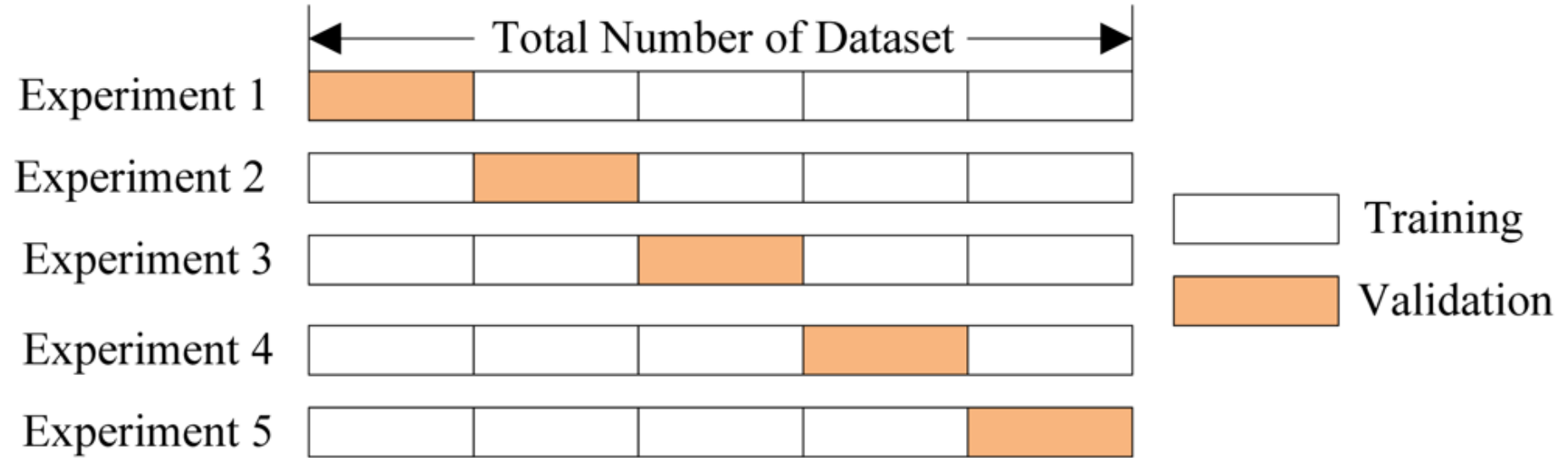
$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}$$

Datasets

- División de los datos en 3 partes:
 - Train set (70%)
 - Validation set (20%)
 - Test set (10%)



Cross-Validation



Sklearn

```
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score

# Data loading & preprocessing
# ...
# ...

# Define the model
model = SVC(kernel='linear');

# Train the model
model.fit(x_train, y_train)

# Make predictions
y_predicted = model.predict(x_test, y_test)

# Evaluate
accuracy_score(y_test, y_predicted)
```