Machine Learning Learning Machine

Machine Learning Basics

Yuri Mitrankov

Stony Brook University

December 05, 2024





Standard Model of ML running:

- **1**Inheriting ML Code
- 2Modify
- **3**Create your training sample
- 4Get total nonsense
- **5**Google for answers
- 6Eureka moment
- 7Try it

Repeat from step 5 until succeed





ML running basics











Validation Set (20%)

Testing Set (20%)

To prevent overfitting

Each ML model has multiple parameters: Run with various sets then find the best

To determine accuracy

4







Validation Set (20%)

Testing Set (20%)

To prevent overfitting

Each ML model has multiple parameters: Run with various sets then find the best **FOR YOU**

To determine accuracy



Classification in ML







Binary classification



12/06/2024

MLLM



Multi-class classification



12/06/2024

MLLM

https://www.kaggle.com/competitions/titanic/overview

12/06/2024

MLLM

Example

12/06/2024

11

Quality metrics

Quality metrics

Regression

Choose params smartly

Choose loss func and solver

Regression

 $\omega_1, \omega_2, \dots, \omega_n \to \omega_1^{-1}, \omega_2^2, \dots, \omega_n$

Just right

The optimal learning rate swiftly reaches the minimum point

Too large of a learning rate causes drastic updates which lead to divergent behaviors

 $J(\theta)$

 θ

Regression

12/06/2024

Transforming loss function by adding regularization term

12/06/2024

MLLM

- 1. Robust to overfitting
- 2. Easy to use
- 3. Parallelizable
- 4. Classification and regression

- 1. Harder to interpret
- 2. Slower and memory intense
- 3. Less accurate for small dataset

n_estimators – Number of trees; increase for better performance but higher computation. **max_depth** – Maximum depth of trees; controls overfitting (deeper trees capture more complexity). **min_samples_split** – Minimum samples to split a node; higher values prevent overfitting. **min** samples leaf – Minimum samples in a leaf; larger values create simpler models. max_features – Controls how many features the model looks at when making each split; choosing fewer features makes the model simpler and faster, while more features make it more accurate but slower. **bootstrap** – Whether to sample data with replacement; affects diversity of trees. **criterion** – Function to measure quality of a split (e.g., "gini" for classification, "mse" for regression).

21

Gradient Boosting

Gradient Boosting

Next is fixing

previous issues

Pros:

- 1. High accuracy
- 2. Handle complex tasks
- 3. Robust to outlier and -9999
- 4. Classification and regression Cons:
 - 1. Hard to interpret
 - 2. Slower training
 - 3. Sensitive to overfitting

n_estimators – Number of boosting stages; too many can overfit, too few underfit. **learning_rate –** Shrinks contribution of each tree; lower values -> more trees \rightarrow more overfit or precision. **max_depth** – Limits tree depth; helps control overfitting. **subsample** – Fraction of data used for fitting each tree; reduces variance. **colsample_bytree –** Fraction of features used per tree; improves generalization. Or one can use max_features **min_samples_split** – Minimum samples to split a node; prevents overfitting. **loss** – Specifies the loss function to minimize (e.g., "squared_error" for regression, "log_loss" for classification). criterion – Measures the quality of a split within each tree (e.g., "friedman_mse" for regression).

Gradient Boosting

Outline

I asked ChatGPT to summarize ML in HEP:

Data Reconstruction

Anomaly Detection

Data Reduction and Selection

Hyperparameter Tuning

Experimental Designs

Outline

I asked ChatGPT to summarize ML in HEP:

Data Reconstruction

Anomaly Detection

Data Reduction and Selection

Hyperparameter Tuning

Experimental Designs

Gradient descent

1

