



DAY 24 — XGBoost (Extreme Gradient Boosting)

| *Boosting done right*

1 Why XGBoost Exists ?

Gradient Boosting is powerful, but:

- ✗ Slow
- ✗ Overfits easily
- ✗ Hard to tune
- ✗ Not optimized for large data

XGBoost was built to fix all of this

2 What is XGBoost?

| **XGBoost is an optimized, regularized, scalable version of Gradient Boosting.**

Think of it as:

- Gradient Boosting
 - math optimization
 - regularization
 - engineering excellence

3 What makes XGBoost Special?

🔥 a) Regularization

XGBoost **penalizes complex trees.**

It adds:

- L1 (alpha)
- L2 (lambda)

This prevents overfitting better than standard GB.

🔥 b) Second-Order Optimization

Instead of using only gradients:

- Uses **gradient + hessian**
- Learns faster
- More accurate splits

This is advanced math under the hood.

🔥 c) Pruning (Smart trees)

XGBoost:

- Grows tree fully
- Prunes branches that don't help

Result:

- Smaller
 - Better Trees
-

🔥 d) Handling Missing Values (AMAZING FEATURE)

XGBoost:

- Learns where missing values should go
 - No need to impute manually in many cases
-

🔥 e) Speed & Parallelization

- Parallel tree building
- Cache-aware
- Efficient memory use

This is why XGBoost dominates competitions.

4 XGBoost vs Random Forest vs Gradient Boosting

Feature	RF	GB	XGBoost
Parallel	✓	✗	✓

Regularization	✗	✗	✓
Speed	Medium	Slow	Fast
Accuracy	Good	Very Good	Excellent
Tuning	Easy	Hard	Medium

5 Important XGBoost Parameters

learning_rate (eta)

Same as Gradient Boosting.

Lower → safer → slower

max_depth

Tree complexity.

Small (3-6) is best.

n_estimators

Number of trees.

More trees → more power → slower.

subsample

Fraction of rows used per tree.

Prevents overfitting.

colsample_bytree

Fraction of features used per tree.

Adds randomness.

reg_alpha & reg_lambda

L1 & L2 regularization.

controls tree complexity.

6 Classification vs Regression in XGBoost

Task	Objective
Classification	binary:logistic

7 When Should you use XGBoost?

- Tabular data
- Medium-large datasets
- Need best accuracy
- Very small datasets
- Simple linear problems