# **Gradient Boosting**

Gradient Boosting is a **supervised learning** algorithm used for **classification and regression**. It is an **ensemble technique** that builds multiple weak learners (usually decision trees) in a sequential manner, where each new tree **corrects the errors** of the previous ones using gradient descent.

### How Gradient Boosting Works

- 1. Train an initial weak learner (usually a decision tree).
- 2. Calculate residuals (errors) between actual and predicted values.
- 3. Fit a new weak learner to predict the residuals.
- 4. **Update the model** by adding this new learner to minimize error.
- 5. **Repeat** steps 2–4 for a fixed number of iterations or until errors are minimized.

#### Popular Variants of Gradient Boosting

- XGBoost (Extreme Gradient Boosting) → Faster and more regularized version.
- 2. **LightGBM (Light Gradient Boosting Machine)** → Optimized for large datasets.
- CatBoost (Categorical Boosting) → Designed for handling categorical features efficiently.

#### **6** When to Use Gradient Boosting?

- When you need **high accuracy** and are okay with longer training times.
- When working with structured/tabular data (like Kaggle competitions).
- When **feature importance** is useful for interpretability.

Gradient Boosting 1

## Day-8 guadient Boosting

learning algorithm that combines multiple week learners to exect a final Model: It can be used for both elassification and seguession purplems.

Greatent Boosting :

ideo! . Puedictor fr at stage to encure loss L(Fn(x), y)

· Train that to approximate negative guadrent:

$$h_{k+1}(x) \approx -\frac{\partial L(f_{k}(x), y)}{\partial f_{k}(x)}$$

· Update predictor by odding a multiple nk11 of hk11:

7. Consider squared loss

$$L(f_h(x_n), y_n) = \xi (f_h(x_n) - y_n)^2$$

· Negative guadient contresponds to residual or

$$-\frac{\partial L(f_{k}(x_{n}),y_{n})}{\partial f_{k}(x_{n})}=y_{n}-f_{k}(x_{n})=\delta_{n}$$

· Tuain hase learner that
with residual dataset { (xn, 5n) xn}

. Base learner to K+1 can be only non-lines - puededor (often a small drusson Tue)

=> Psuedo code: GBA:

· Initialize purdictor with a constant c:
fo(xn): arg mile &nh(c, yn)

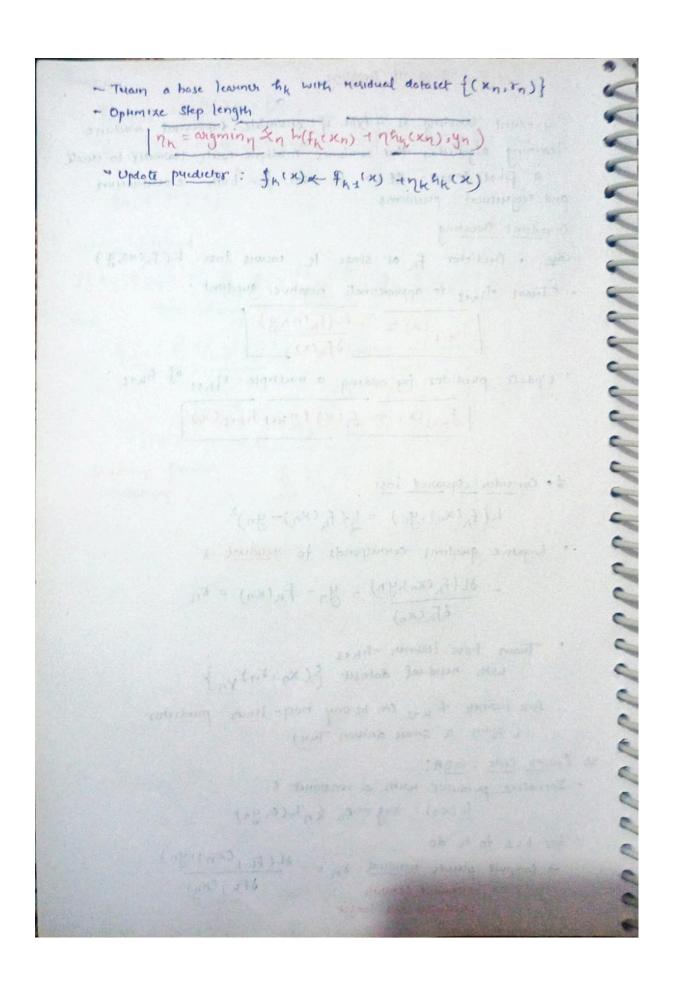
for k=1 to k do

-> compute period mesidual on = - dl (fk-1(2n), yn)

difference between

dfk-1(xh)

peredicton and target



Gradient Boosting 3