



Analysis of Frank-Wolfe and Pairwise Frank-Wolfe Algorithms on the LASSO Problem

Optimization for Data Science

Roya Ghamari - 2071969

July 2024

Introduction

Objective: Study the application of FW and PFW algorithms in minimizing the LASSO regression problem

Algorithms: Frank-Wolfe (FW), Pairwise Frank-Wolfe (PFW)

Focus on convergence behavior and optimization efficiency.

Datasets: Concrete Compressive Strength, Boston Housing

LASSO: Ideal for feature selection in high-dimensional data, explicitly recovering sparse solutions.

Constraints and Rules: ℓ_1 - Ball Constraint, Duality Gap, Armijo Rule

LASSO Problem

$$1) \min_{x \in \mathbb{R}^n} f(x) := \|Ax - b\|^2, \quad \text{s.t. } \|x\|_1 \leq \tau$$

$$2) T = \{(r_i, b_i) \in \mathbb{R}^n \times \mathbb{R} : i \in [1:m]\}$$

$$3) C = \{x \in \mathbb{R}^n : \|x\|_1 \leq \tau\} = \text{conv}\{\pm \tau e_i : i \in [1:n]\}$$

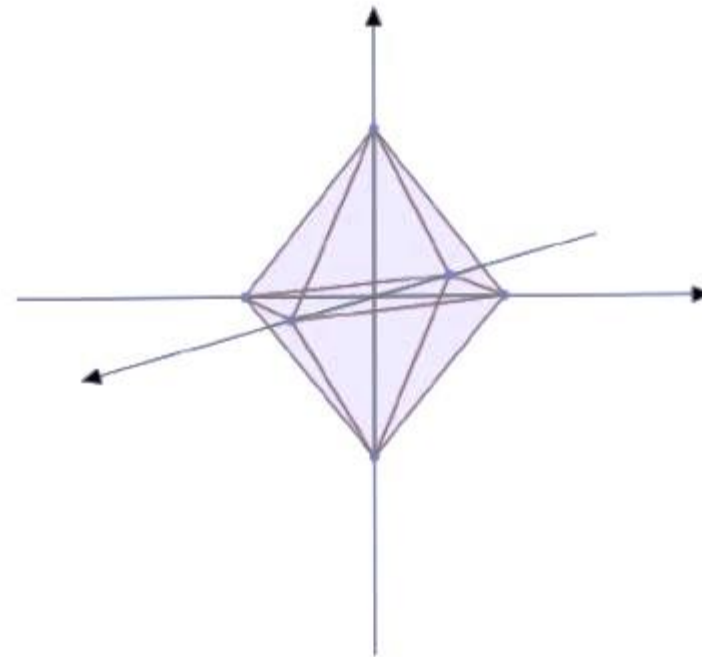
$$4) LMO_C(\nabla f(x_k)) = \text{sign}(-\nabla_{i_k} f(x_k)) \cdot \tau e_{i_k}, \quad i_k \in \arg \max_i |\nabla_i f(x_k)|$$

l_1 -ball constraint

1) $B_1(\tau) = \{x \in \mathbb{R}^n : \|x\|_1 \leq \tau\}$

2) $\|x\|_1 = \sum_{i=1}^n |x_i|$

l_1 -ball Constraint in 3D



Frank Wolfe Algorithm

Algorithm 1 Frank-Wolfe method

- 1 Choose a point $x_0 \in C$
 - 2 For $k = 0, \dots$
 - 3 If x_k satisfies some specific condition (duality gap), then STOP
 - 4 Compute $s_k \in LMO_C(\nabla f(x_k))$
 - 5 Set $d_k^{FW} = s_k - x_k$
 - 6 Set $x_{k+1} = x_k + \alpha_k d_k^{FW}$, with $\alpha_k \in (0,1]$ a suitably chosen stepsize
 - 7 End for
-

$$f(x) := \|Ax - b\|^2$$

$$LMO_C(\nabla f(x_k)) = \text{sign}(-\nabla_{i_k} f(x_k)) \cdot \tau e_{i_k}$$

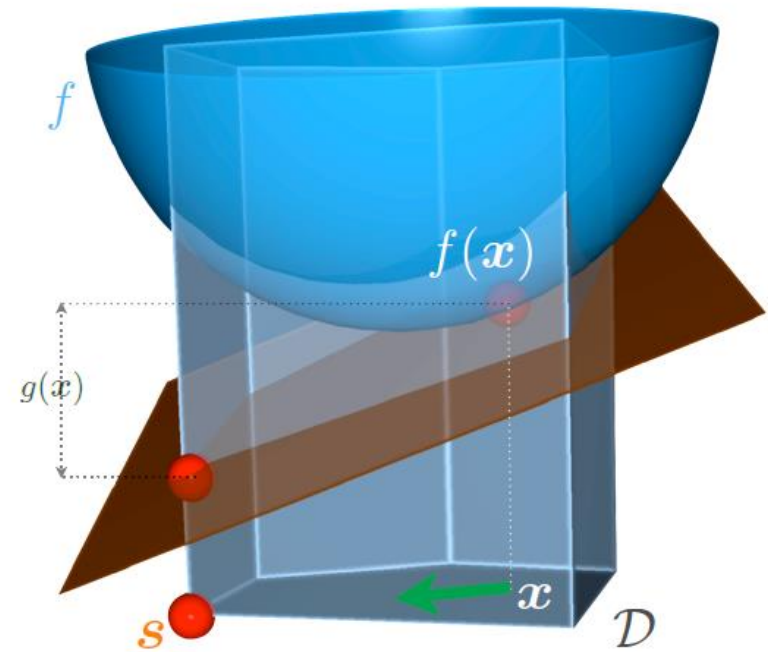
Pairwise Frank Wolfe Algorithm

Algorithm 2 Pairwise Frank-Wolfe method

```
1  Let  $x^{(0)} \in \mathcal{A}$  and  $\mathcal{S}^{(0)} := \{x^{(0)}\}$ 
2  For  $t = 0$  to  $T$  do
3      Let  $s_t := LMO_{\mathcal{A}}(\nabla f(x^t))$  and  $d_t^{FW} := s_t - x^{(t)}$ 
4      Let  $v_t \in \operatorname{argmax}_{v \in \mathcal{S}^{(t)}} \langle \nabla f(x^t), v \rangle$  and  $d_t^A = x^{(t)} - v_t$ 
5      If  $g_t^{FW} := \langle -\nabla f(x^t), d_t^{FW} \rangle \leq \epsilon$  then return  $x_t$ 
6      Else
7           $d_t = d_t^{PFW} = d_t^{FW} + d_t^A = s_t - v_t$ ,  $\gamma_{max} := |\alpha_{v_t}|$ 
8          Determine  $\gamma_t \in [0, \gamma_{max}]$  a suitably chosen stepsize
9          Update  $x^{(t+1)} := x^t + \gamma_t d_t$ 
10         Update  $\mathcal{S}^{(t+1)} := \{v \in \mathcal{A} \mid |\alpha_v^{(t+1)}| > 0\}$ 
11     End
12 End
```

The Duality Gap

- $g(x) = \max_{s \in \mathcal{D}} \nabla f(x)^T (x - s) = \max_{s \in \mathcal{D}} -\nabla f(x)^T (s - x)$
- $g(x) \geq f(x) - f(x^*)$



Duality gap shown in a step of frank-wolfe algorithm

Armijo Line Search

- Fix the parameters $\delta \in (0,1)$ and $\gamma \in (0,1/2)$
- A starting step size Δ_k
- The steps $\alpha = \delta^m \Delta_k$ with $m = \{0, 1, 2, \dots\}$ are tried until:

$$f(x_k + \alpha d_k) \leq f(x_k) + \gamma \alpha \nabla f(x_k)^T d_k$$

- Set $\alpha_k = \alpha$

Datasets

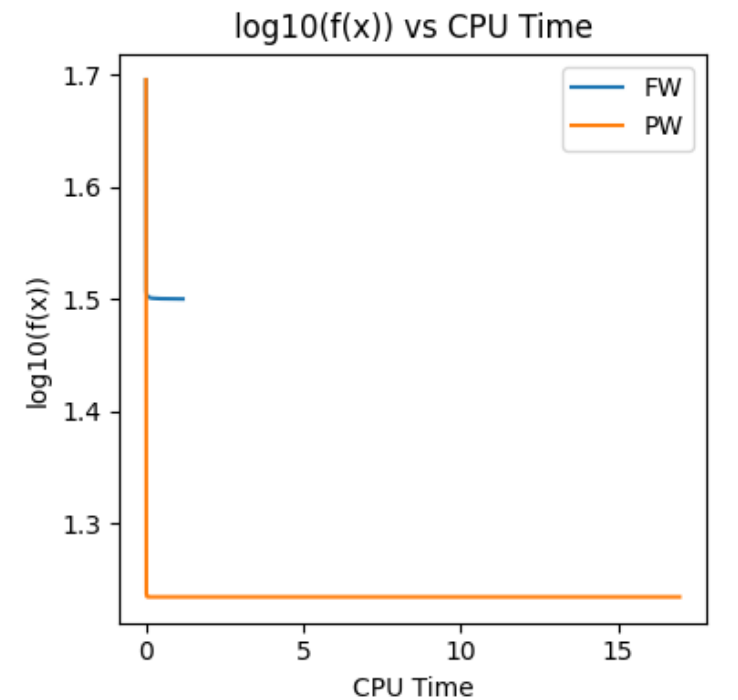
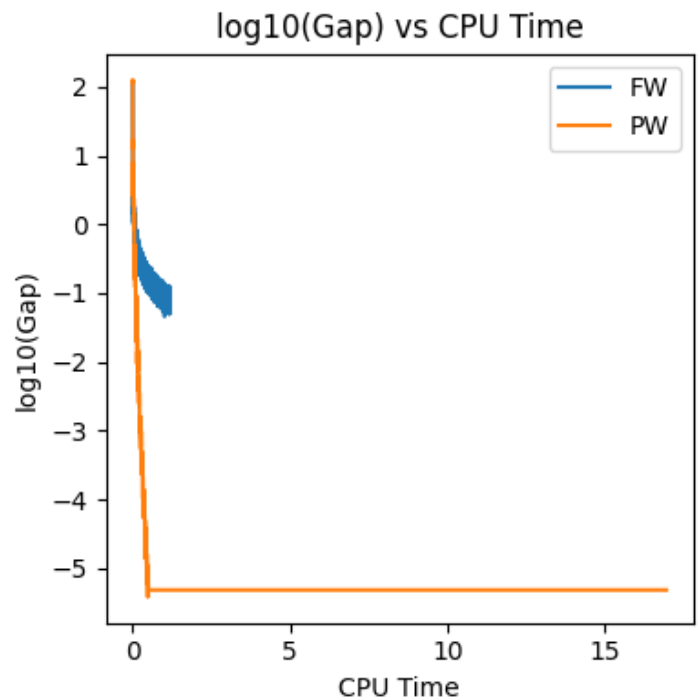
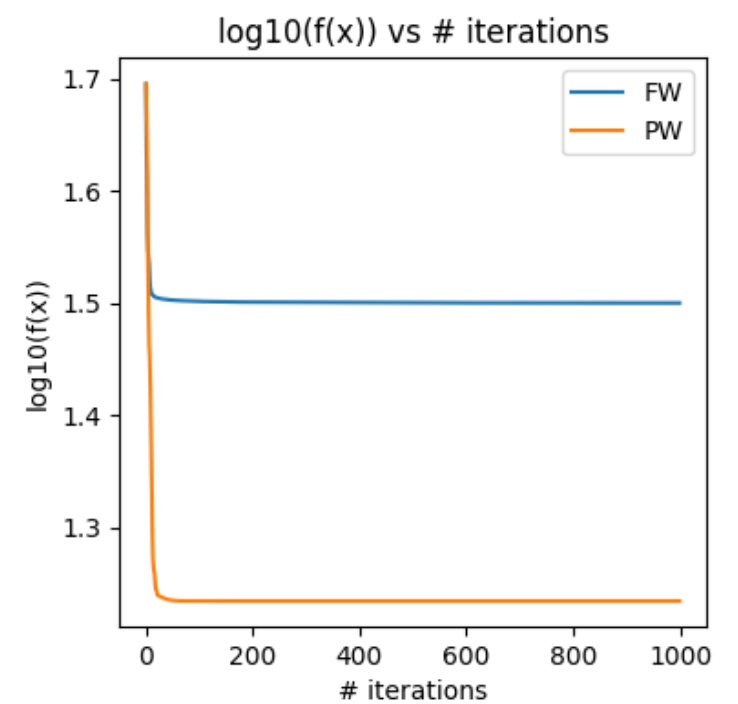
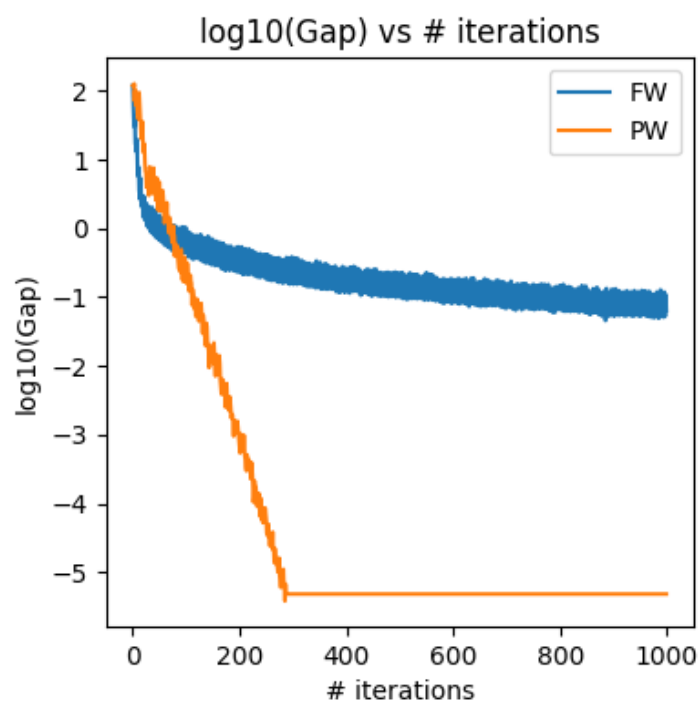
Concrete Strength Dataset (1030 instances)

- **Cement:** Amount of cement in the mixture
- **Blast Furnace Slag:** Amount of blast furnace slag in the mixture
- **Fly Ash:** Amount of fly ash in the mixture
- **Water:** Amount of water in the mixture
- **Superplasticizer:** Amount of superplasticizer in the mixture
- **Coarse Aggregate:** Amount of coarse aggregate in the mixture
- **Fine Aggregate:** Amount of fine aggregate in the mixture
- **Age:** Age of the concrete
- **Compressive Strength:** Compressive strength of the concrete

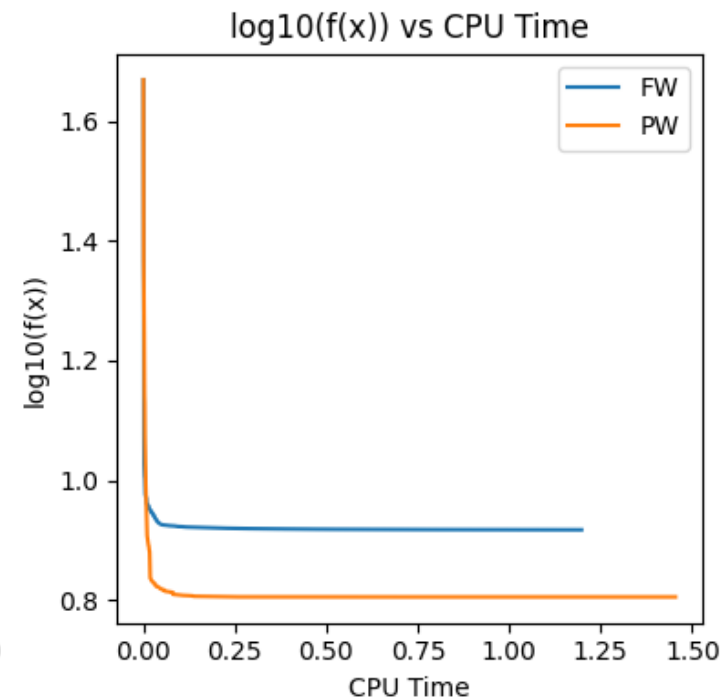
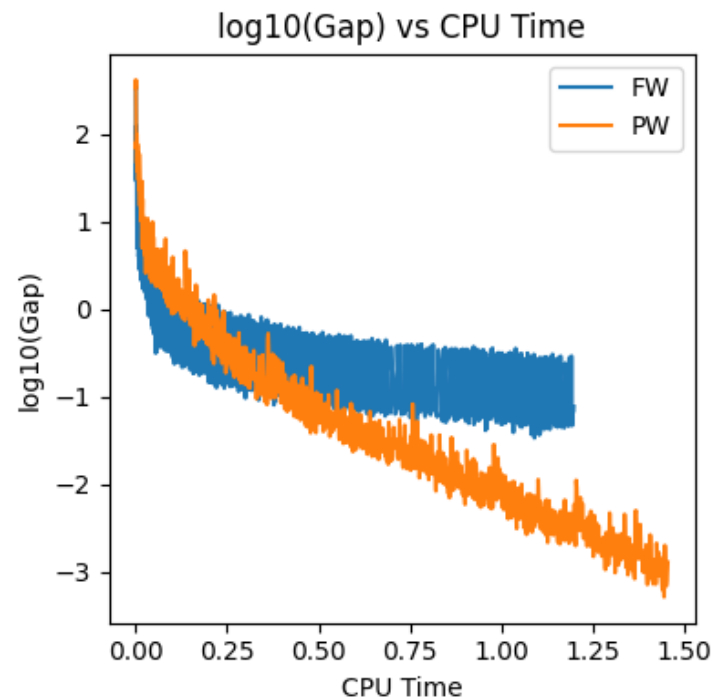
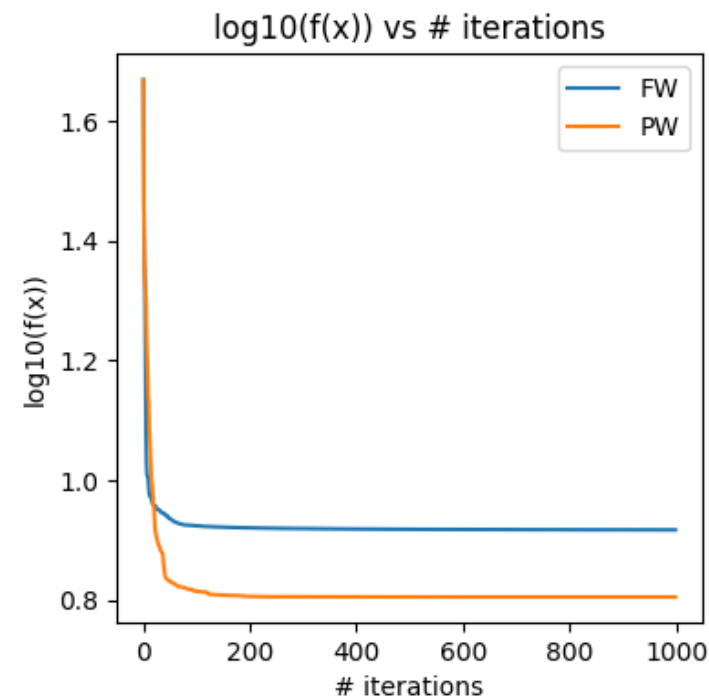
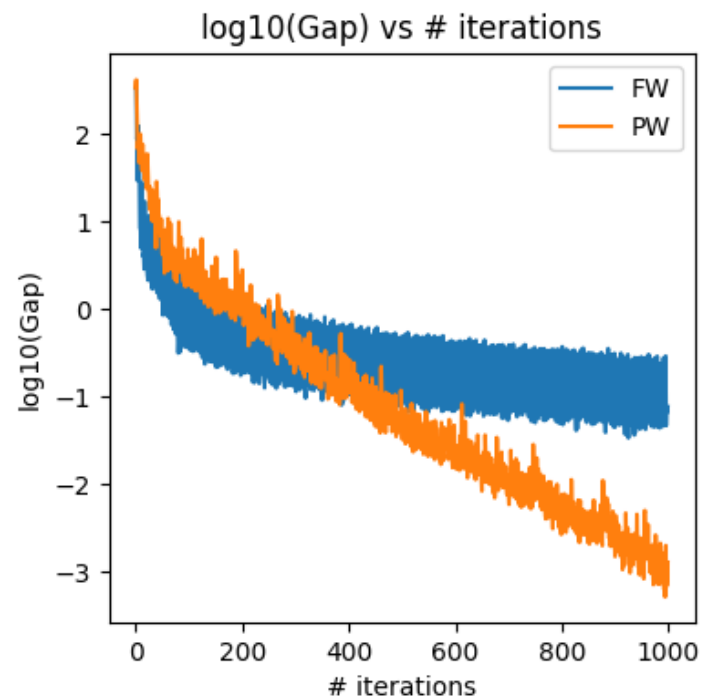
Boston Housing Dataset (506 instances)

- **CRIM:** Per capita crime rate by town
- **ZN:** Proportion of residential land zoned for lots over 25,000 sq.ft.
- **INDUS:** Proportion of non-retail business acres per town
- **CHAS:** Charles River dummy variable (1 if tract bounds river; 0 otherwise)
- **NOX:** Nitric oxides concentration (parts per 10 million)
- **RM:** Average number of rooms per dwelling
- **AGE:** Proportion of owner-occupied units built prior to 1940
- **DIS:** Weighted distances to five Boston employment centers
- **RAD:** Index of accessibility to radial highways
- **TAX:** Full-value property-tax rate per \$10,000
- **PTRATIO:** Pupil-teacher ratio by town
- **B:** $1000(B_k - 0.63)^2$ where B_k is the proportion of blacks by town
- **LSTAT:** Percentage of lower status of the population
- **MEDV:** Median value of owner-occupied homes in \$1000's

Results For Concrete Strength Dataset



Results For Boston Housing Dataset



Conclusion



FW Algorithm:

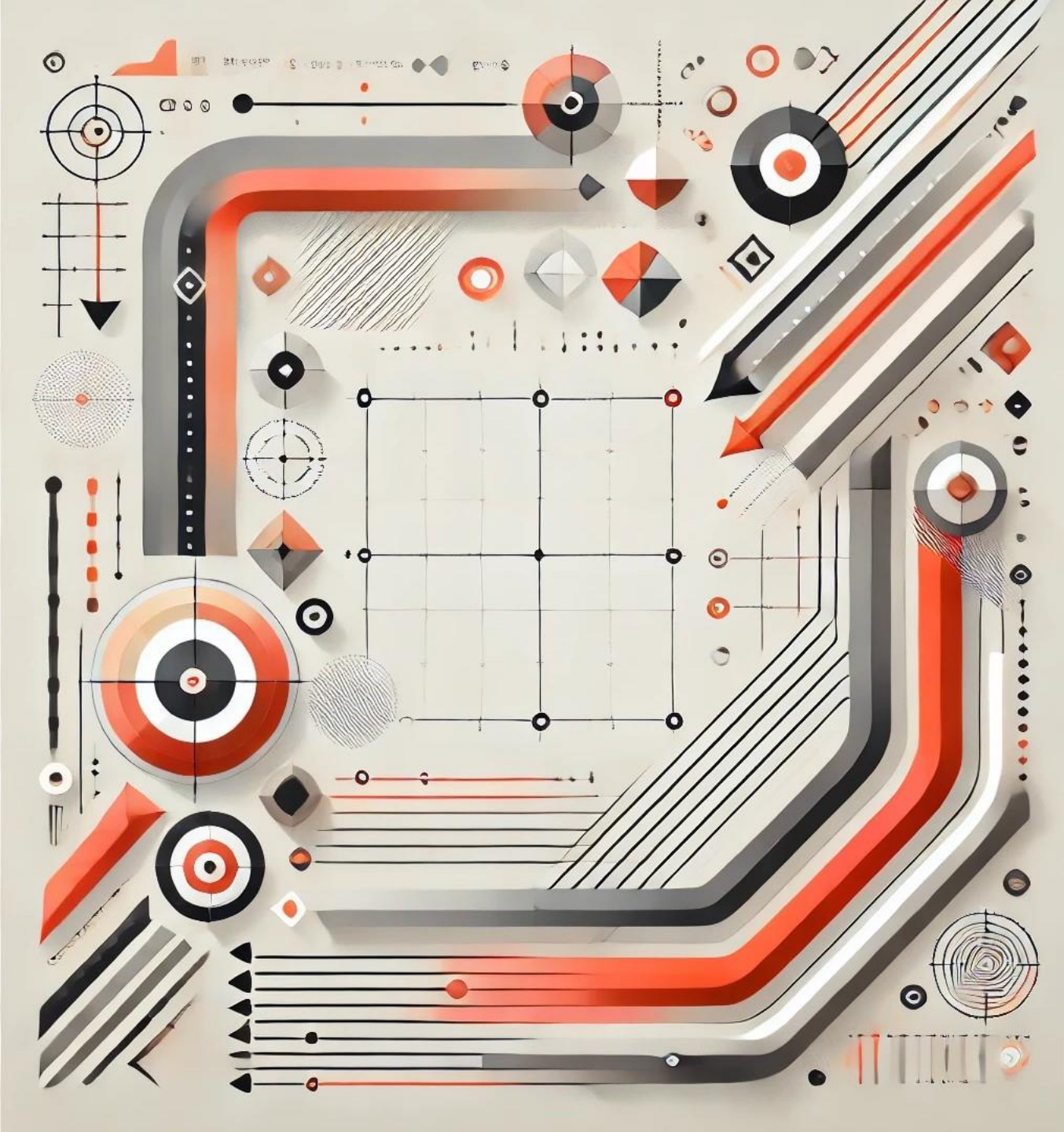
- Efficient in terms of CPU time
- Suitable for quick, significant improvements
- Tends to plateau early, resulting in suboptimal solutions in complex datasets

PFW Algorithm:

- Better at finding global optimal
- Higher computational cost
- Preferable for applications prioritizing the quality of the final solution

Datasets:

- Similar trends but less pronounced differences for Boston Housing
- Higher dimensionality and fewer instances made the optimization landscape more complex



Thank you
