**Part-A**

Assume least squares objective,

$$f(w) = \frac{1}{2n} \| Xw - y \|_2^2, \text{ gradient } g(w) = \nabla f(w) = \frac{1}{n} X^T (Xw - y), \text{ Hessian } H = \frac{1}{n} X^T X$$

→ As we iterate $w_k$, take a gradient descent $p_k$ ( for vanilla GD, $p_k = -g_k$), stepsize $t_k > 0$

## (i) Exact line search (closed form for least squares)

for quadratic $f$, with PSD $H$, the exact minimizer along the line $w_k + t p_k$ is:

$$ t_k^* = \arg\min_{t \geq 0} \phi(t) = f(w_k + t p_k) = \frac{-g_k^T p_k}{p_k^T H p_k} $$

Special Case: (steepest descent)

$$ p_k = -g_k $$

$$ t_k^* = \frac{g_k^T g_k}{g_k^T H g_k} $$

Algorithm:- ① Compute $g_k = \nabla f(w_k)$ → ② Set $p_k = -g_k$ → ③ Compute $t_k = \frac{-g_k^T p_k}{p_k^T H p_k}$ → ④ Update $w_{k+1} = w_k + t_k p_k$

## (ii) $\alpha-\beta$ Backtracking line search (Armijo rule)

Parameters : choose $\alpha \in [0, 0.5)$, $\beta \in [0, 1]$  (typical $\alpha = 10^{-4}$, $\beta = 0.5$)

Goal: find smallest $m \in \{0, 1, 2, \dots\}$ s.t Armijo decreases holds.

Algorithm: - ① Compute $g_k = \nabla f(w_k)$ → ② Initialise $t \leftarrow 1$ → ③ while $f(w_k + t p_k) > f(w_k) + \alpha t g_k^T p_k$, set $t \leftarrow \beta t$.

④ Set $t_k \leftarrow t$ and update $w_{k+1} = w_k + t_k p_k$

## (iii) Ternary search along a line (for unimodel $\phi(t) = f(w_k + t p_k)$)

Use when $\phi(t)$ is unimodal on an interval $[0, \tau]$, if $\tau$ is unknown first bracket it.

Bracketing (doubling)

① Set $t_0 \leftarrow 0$, $t_1 \leftarrow \tau > 0$ → ② while $\phi(t_1) < \phi(t_0)$: set $t_0 \leftarrow t_1$, $t_1 \leftarrow 2t_1$ → ③ Now $[0, t_1]$ brackets the min$^m$.

Ternary search on $[L, R]$ (with $L=0$, $R=\tau$):

Repeat until $R - L < \epsilon$, ①

- $m_1 = L + \frac{R-L}{3}$, $m_2 = R - \frac{R-L}{3}$
- if $\phi(m_1) < \phi(m_2)$, set $R \leftarrow m_2$ else $L \leftarrow m_2$

→ ② Return $t_k \in [L, R]$ → ③ Update $w_{k+1} = w_k + t_k p_k$

**Part-B**

we use the kaggle house-prices — Advanced Regression techniques train set (train-csv). To keep a pure least-square objective, only numerical predictors are used Missingvalues are imputed with column means, features are standardized and a bias column is added.

We fit a linear regression model via gradient descent $(p = -\nabla f)$ on

$$f(w) = \frac{1}{2n} \|Xw - y\|_2^2 \qquad \nabla f(w) = \frac{1}{n} X^T(Xw - y).$$

1. $\alpha-\beta$ backtracking $(\alpha = 10^{-4}, \beta = 1, t_{init} = 1)$.

2. Ternary search along the line w.r.t.p with bracketing by doubling (start $T=1$, growth $=2$) and $\epsilon = 10^{-7}$ (tolerance)

Stopping rule :→ $|f_k - f_{k-1}| < 10^{-10}$ or 5000 iterations.

Results :→  • Backtracking  462 iterations,  0.2154s, final $f = 5.7550 \times 10^8$,
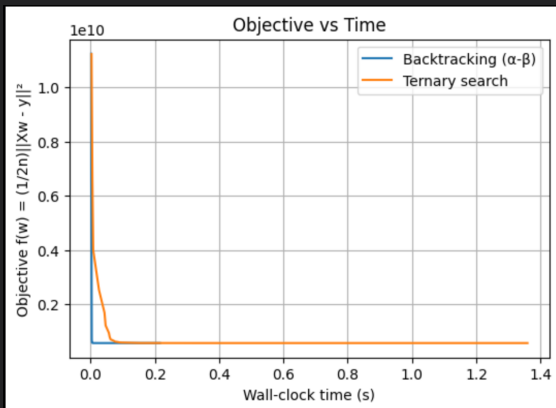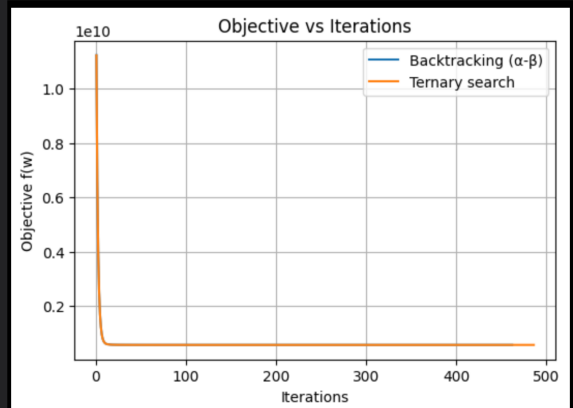           • Ternary  486 iterations,  1.3594s, final $f = 5.7550 \times 10^8$

Graphs:



fig 1. Objective $f(w)$ vs Time (s)



fig 2: Objective $f(w)$ vs iterations