

# Why Line-Search When You Can Plane-Search?

ISMP 2024, Montreal

---

Betty Shea, Mark Schmidt  
University of British Columbia



## We could be using better step sizes

For many common machine learning (ML) problems, these step size strategies all have the same asymptotic cost

- fixed step size,
- line search,
- plane search

## Polyak's Heavy Ball Method (PHB)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \nabla f(\mathbf{x}_t) + \beta_t (\mathbf{x}_t - \mathbf{x}_{t-1})$$

- Usually  $\alpha_t = \alpha$  and  $\beta_t = \beta$

## Polyak's Heavy Ball Method (PHB)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \nabla f(\mathbf{x}_t) + \beta_t (\mathbf{x}_t - \mathbf{x}_{t-1})$$

- Usually  $\alpha_t = \alpha$  and  $\beta_t = \beta$
- But you can do a **plane search** for the same cost

$$\alpha_t, \beta_t = \arg \min_{\alpha, \beta} f(\mathbf{x}_t - \alpha \nabla f(\mathbf{x}_t) + \beta (\mathbf{x}_t - \mathbf{x}_{t-1}))$$

## Polyak's Heavy Ball Method (PHB)

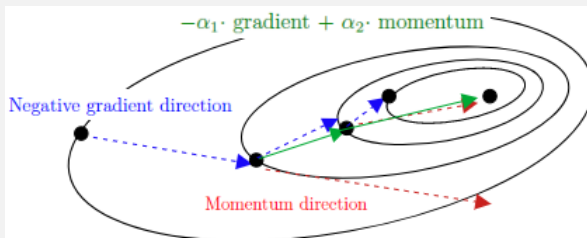
$$\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_t \nabla f(\mathbf{x}_t) + \beta_t (\mathbf{x}_t - \mathbf{x}_{t-1})$$

- Usually  $\alpha_t = \alpha$  and  $\beta_t = \beta$
- But you can do a **plane search** for the same cost

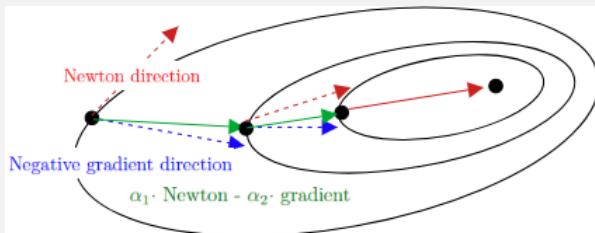
$$\alpha_t, \beta_t = \arg \min_{\alpha, \beta} f(\mathbf{x}_t - \alpha \nabla f(\mathbf{x}_t) + \beta (\mathbf{x}_t - \mathbf{x}_{t-1}))$$

- Potentially more progress at every iteration
- No hyperparameter tuning

## Adaptable gradient and momentum learning rates



### Hybrid first- and second- order methods



## Subproblem to solve at every iteration $t$

Step sizes and search directions are

$$\alpha_t = \begin{bmatrix} \alpha_{t,1} \\ \vdots \\ \alpha_{t,k} \end{bmatrix}, P_t = \begin{bmatrix} | & & | \\ p_{t,1} & \cdots & p_{t,k} \\ | & & | \end{bmatrix}$$



## Subproblem to solve at every iteration $t$

Step sizes and search directions are

$$\alpha_t = \begin{bmatrix} \alpha_{t,1} \\ \vdots \\ \alpha_{t,k} \end{bmatrix}, P_t = \begin{bmatrix} | & & | \\ p_{t,1} & \cdots & p_{t,k} \\ | & & | \end{bmatrix}$$

The update becomes

$$\mathbf{x}_{t+1} = \mathbf{x}_t + P_t \alpha_t$$

## Subproblem to solve at every iteration $t$

Step sizes and search directions are

$$\alpha_t = \begin{bmatrix} \alpha_{t,1} \\ \vdots \\ \alpha_{t,k} \end{bmatrix}, P_t = \begin{bmatrix} | & & | \\ p_{t,1} & \cdots & p_{t,k} \\ | & & | \end{bmatrix}$$

The update becomes

$$\mathbf{x}_{t+1} = \mathbf{x}_t + P_t \alpha_t$$

Plane search solves this problem

$$\alpha_t \in \arg \min_{\alpha} f(\mathbf{x}_t + P_t \alpha_t)$$

## Matrix-vector multiplications are bottlenecks

Linear composition problems (LCPs)

$$f(\mathbf{x}) = g(A\mathbf{x})$$

Examples:

- least squares:  $g(A\mathbf{x}) = \frac{1}{2} \|A\mathbf{x} - \mathbf{y}\|^2$
- logistic regression:  $g(A\mathbf{x}) = \log(1 + \exp(-\mathbf{y} \circ A\mathbf{x}))$

\*Narkiss and Zibulevsky (2005) SESOP

## Cost of calculating function value

$$f(\mathbf{x}_{t+1}) = g(A(\mathbf{x}_t + P_t \alpha_t)) = g(\mathbf{v}_t + \tilde{P}_t \alpha_t)$$

where

$$\mathbf{v}_t = A\mathbf{x}_t \in \mathbb{R}^m \text{ and } \tilde{P}_t = AP_t \in \mathbb{R}^{m \times k}$$

Matrix-vector multiplications:

1.  $\mathbf{v}_t = A\mathbf{x}_t$
2.  $\tilde{P}_{t,1} = A^\top \mathbf{p}_{t,1}$
3. ...
4.  $\tilde{P}_{t,k} = A^\top \mathbf{p}_{t,k}$

## Subproblem at iteration $t$

Store  $\mathbf{v}_t$  and  $\tilde{\mathbf{P}}_t$ . Plane search is

$$\alpha_t \in \arg \min_{\alpha} f(\mathbf{x}_{t+1}) = g(\mathbf{v}_t - \tilde{\mathbf{P}}_t \alpha) \quad (1)$$

- evaluating (1) requires no matrix-vector multiplication
- does require  $k$  vector-scalar multiplications
- significant gain if  $A$  is big and  $k$  is small
- buy 1  $\alpha$  get all the other  $\alpha$ s for free

## Subproblem at iteration $t$

PHB is a special case of  $k$ -direction plane search

$$\alpha_t = \begin{bmatrix} \alpha_{t,1} \\ \alpha_{t,2} \end{bmatrix}, \quad P_t = \begin{bmatrix} -\nabla f(x_t) & (x_t - x_{t-1}) \end{bmatrix}$$

$$\alpha_t \in \arg \min_{\alpha} f(\mathbf{x}_{t+1}) = g(\mathbf{v}_t + \alpha_{t,1} \tilde{\mathbf{p}}_{t,1} + \alpha_{t,2}(\mathbf{v}_t - \mathbf{v}_{t-1})) \quad (2)$$

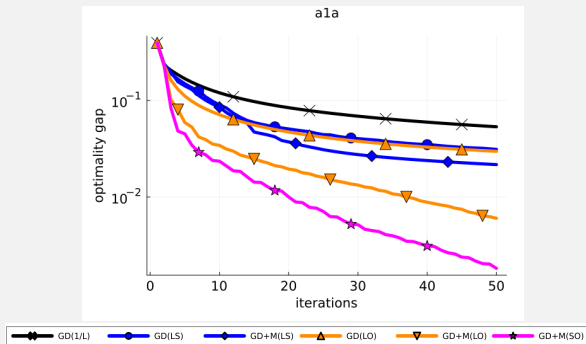
## Matrix-vector multiplications

Matrix-vector multiplications

1.  $\tilde{\mathbf{p}}_{t,1} = -A^T \nabla f(\mathbf{x}_t)$ .
2.  $\mathbf{v}_t = A\mathbf{x}_t$ . (Use  $\mathbf{v}_{t-1}$  from previous iteration.)

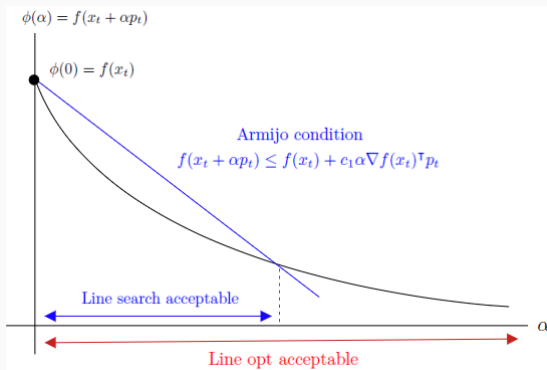
Same as gradient descent with fixed step size and no momentum!

## Fixed versus varying momentum weight



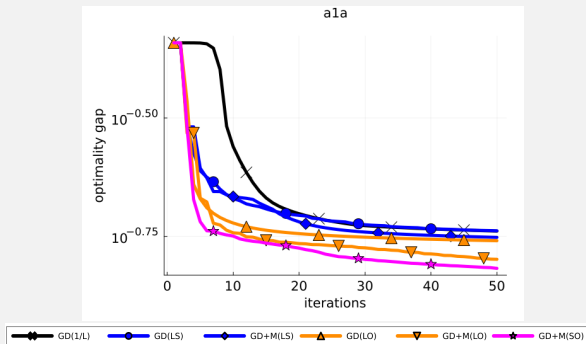


Armijo line search may rule out good step sizes

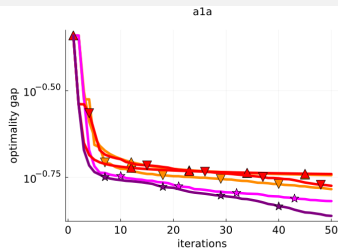
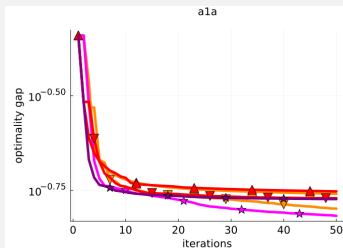


This affects Wolfe conditions too.

## No per-layer step sizes



## Per-layer step sizes without and with regularization



## Miscellaneous things we saw from the experiments

- large  $k$  (e.g.  $k > 2$ ) does not appear to help much
- 3-term plane search implementation of Nesterov's acceleration
- can improve on Adam, quasi-Newton methods, etc.

## Software development: optimizers with built in plane search

- **Alyssa Zhang**,  
NSERC Undergraduate Student  
Research Award



## Questions

- Does this work in the stochastic case?
- Does this work for modern architectures?



(or many) free step sizes

## Accelerated gradient methods

