

FOSI: Hybrid First and Second Order Optimization

HADAR SIVAN



MOSHE GABEL



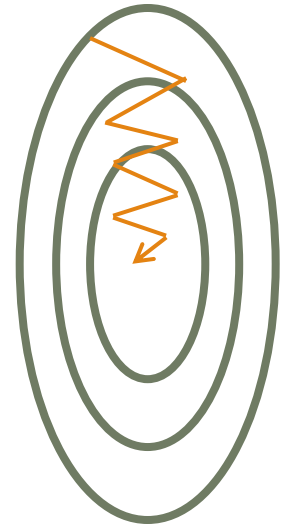
ASSAF SCHUSTER



Background

- $\min_{\Theta} f(\Theta)$
- Iterative optimizer:

$$\Theta = \Theta - \alpha \nabla f(\Theta)$$

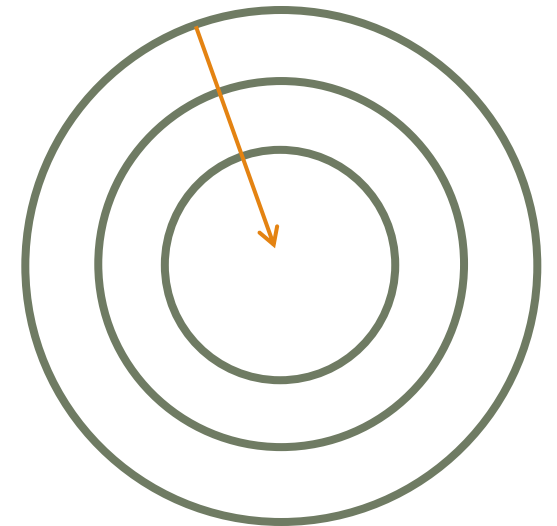


Background

➤ $\min_{\Theta} f(\Theta)$

➤ Iterative optimizer:

$$\Theta = \Theta - \alpha \overset{\text{preconditioner}}{\mathbf{P}^{-1}} \nabla f(\Theta)$$



Background

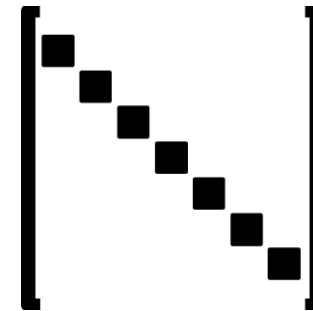
➤ $\min_{\Theta} f(\Theta)$

➤ Iterative optimizer:

$$\Theta = \Theta - \alpha P^{-1} \nabla f(\Theta)$$

- Adam
- AdamW
- RMSProp
- ...

$$P^{-1} \sim \text{diag}(H)^{-1}$$



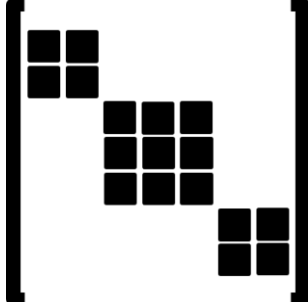
Background

- $\min_{\Theta} f(\Theta)$
- Iterative optimizer:

$$\Theta = \Theta - \alpha P^{-1} \nabla f(\Theta)$$

- K-FAC
- Shampoo
- K-BFGS
- ...

↓ ↓

$$P^{-1} \sim \text{block_diag}(H)^{-1}$$


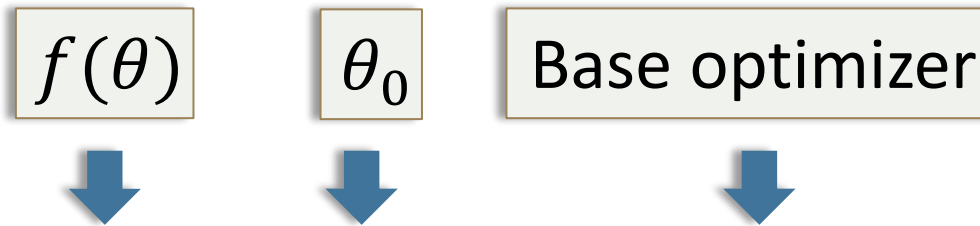
The diagram shows a large square matrix enclosed in square brackets. The matrix is block diagonal, with several smaller square blocks of varying sizes arranged along the main diagonal. The blocks are represented by small black squares. Two red arrows point downwards from the text above to the first and second blocks of the matrix.

FOSI

Improves performance of first-order optimizers by using two key ideas:

1. Split the problem to two parts with different optimizers
 2. Estimate P^{-1} directly, avoid inversion of P
- ✓ Drop-in replacement, no extra tuning, faster wall-time

FOSI: Overview



FOSI – repeat until convergence:

1. Gradient
2. Second-order info
3. Update θ_t



Notations

At iteration t :

$$H = \nabla^2 f(\theta_t)$$

eigenvalues

eigenvectors

$$\left[\lambda_1 > \dots > \lambda_k > \lambda_{k+1} > \dots > \lambda_{n-l} > \lambda_{n-l+1} > \dots > \lambda_n \right]$$

$$\begin{bmatrix} | & | & | & | & | \\ v_1 & \dots & v_k & v_{k+1} & \dots & v_{n-l} & v_{n-l+1} & \dots & v_n \\ | & | & | & | & | \end{bmatrix}$$

1. Gradient
2. Second-order info
3. Update θ_t




Notations

At iteration t :

$$H = \nabla^2 f(\theta_t)$$

eigenvalues

eigenvectors

1. Gradient
2. Second-order info 
3. Update θ_t

$$\left[\underbrace{\lambda_{k+1} > \dots > \lambda_{n-\ell}}_k \right]$$

$$\left[\begin{array}{c} | \quad | \\ v_{k+1} \quad \dots \quad v_{n-\ell} \\ | \quad | \end{array} \right]$$

$$\left[\lambda_1 \quad \dots \quad \lambda_k \quad \lambda_{n-\ell+1} \quad \dots \quad \lambda_n \right]$$

$$\left[\begin{array}{c} | \quad | \quad | \quad | \\ v_1 \quad \dots \quad v_k \quad v_{n-\ell+1} \quad \dots \quad v_n \\ | \quad | \quad | \quad | \end{array} \right]$$

Update Step

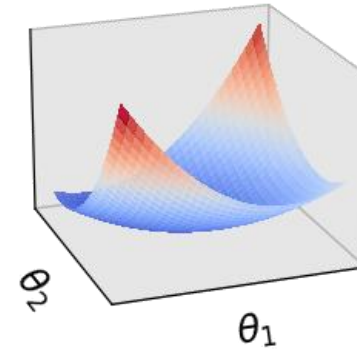
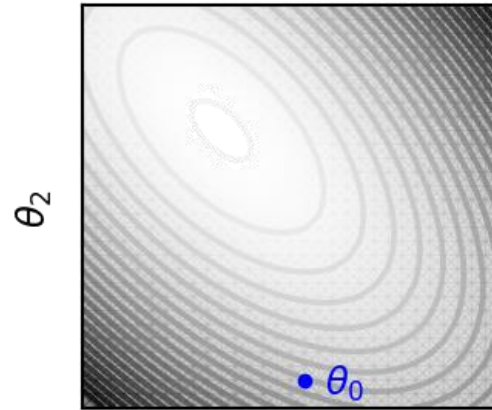
- a. Split f to f_1 and f_2
- b. Newton's step on f_1
- c. Base opt step on f_2

- 1. Gradient
- 2. Second-order info
- 3. Update θ_t

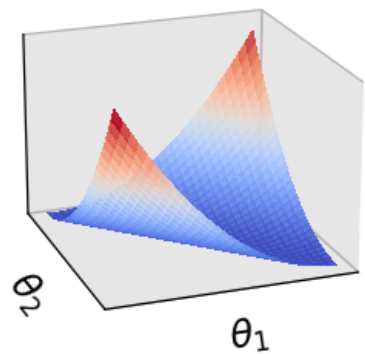


Splitting f ($k = 1, \ell = 0$)

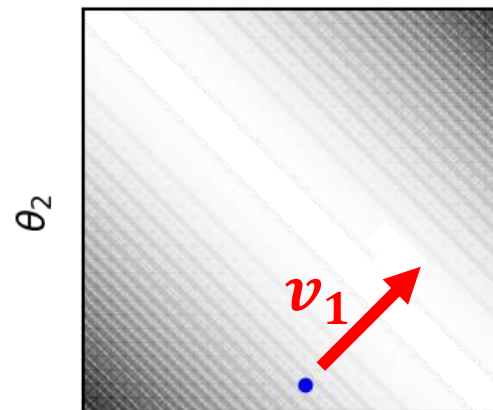
$$f(\theta) = \frac{1}{2} \theta^T \begin{pmatrix} 1 & -1 \\ \sqrt{2} & \sqrt{2} \\ 1 & 1 \\ \sqrt{2} & \sqrt{2} \end{pmatrix} \begin{pmatrix} 4 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & -1 \\ \sqrt{2} & \sqrt{2} \\ 1 & 1 \\ \sqrt{2} & \sqrt{2} \end{pmatrix}^T \theta$$



1. Gradient
2. Second-order info
3. Update θ_t
 - a. Split f to f_1 and f_2
 - b. Newton's step on f_1
 - c. Base opt step on f_2

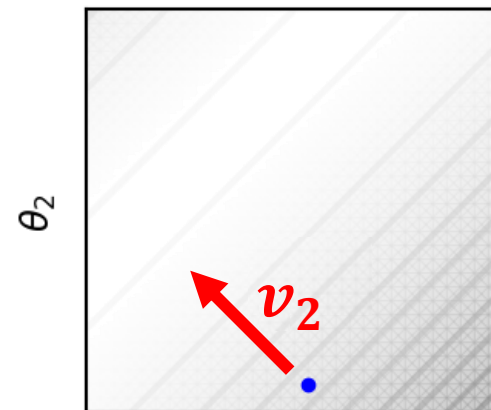


$f_1(\theta)$



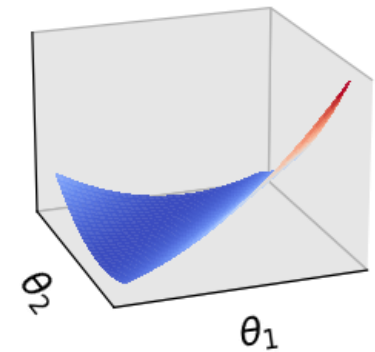
θ_1

$f_2(\theta)$



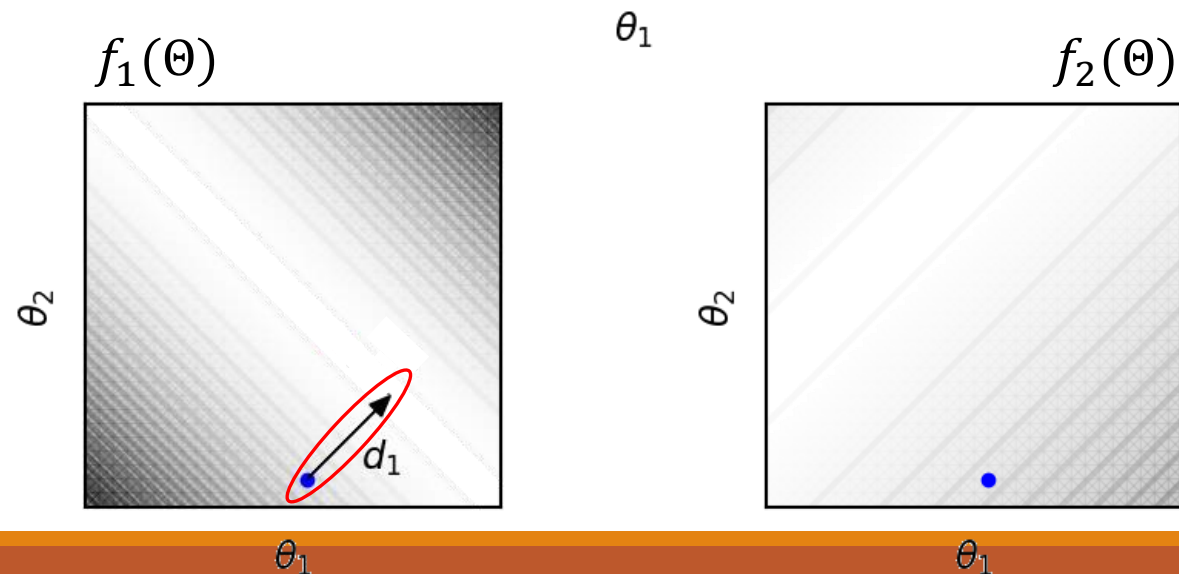
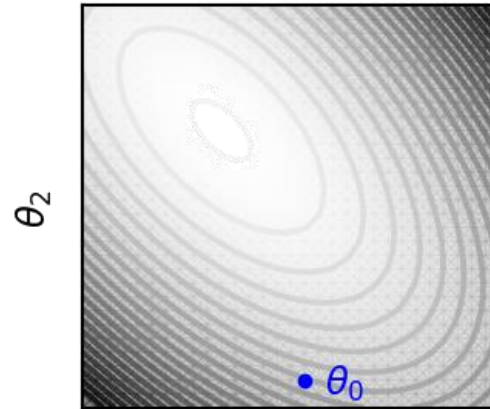
θ_1


θ_1



Minimize f_1

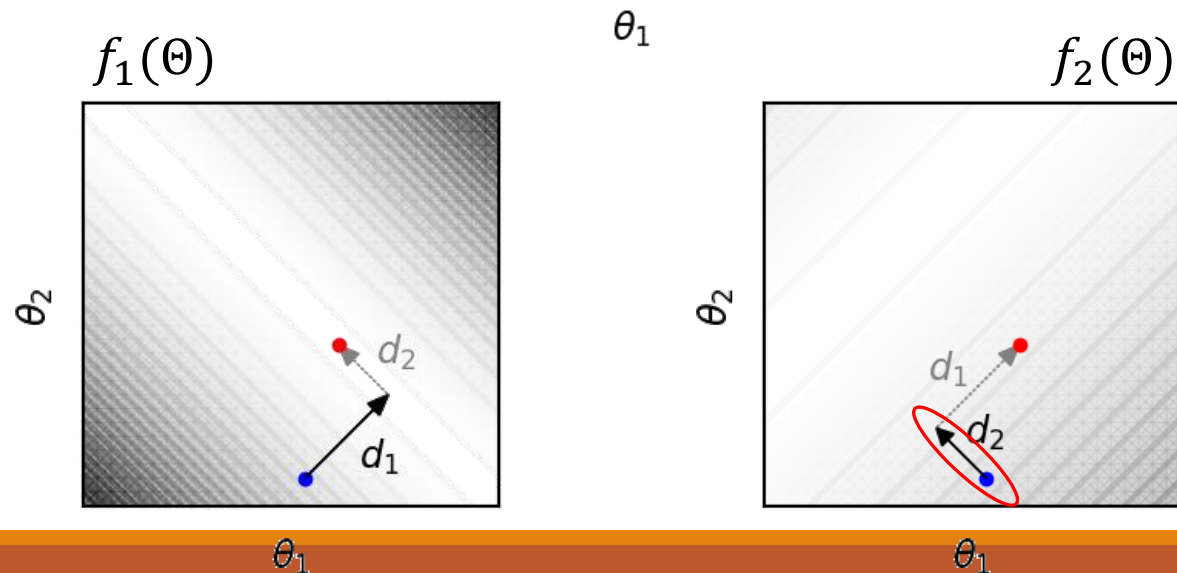
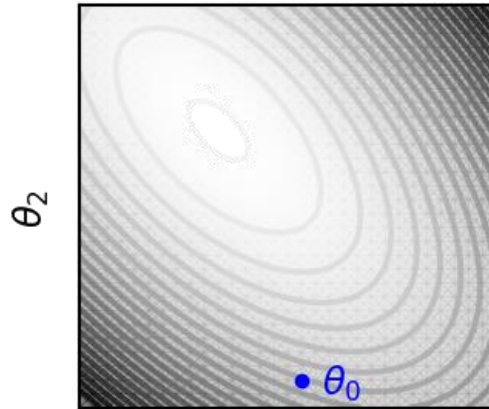
- ✓ d_1 brings f_1 to minimum in a single step
- ✓ Trivial matrix inversion



1. Gradient
2. Second-order info
3. Update θ_t
 - a. Split f to f_1 and f_2
 - b. Newton's step on f_1 
 - c. Base opt step on f_2

Minimize f_2

- ✓ d_1 brings f_1 to minimum in a single step
- ✓ Trivial matrix inversion
- ✓ d_1 orthogonal to d_2
- ✓ No mutual impact

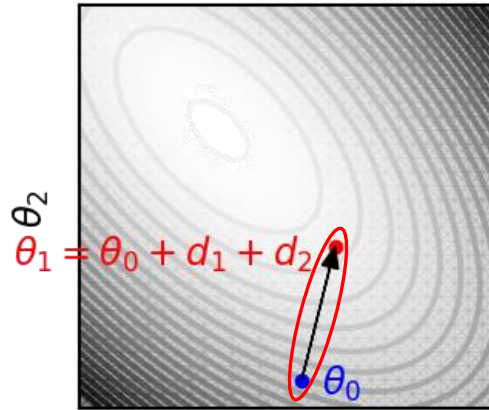


1. Gradient
2. Second-order info
3. Update θ_t
 - a. Split f to f_1 and f_2
 - b. Newton's step on f_1
 - c. Base opt step on f_2

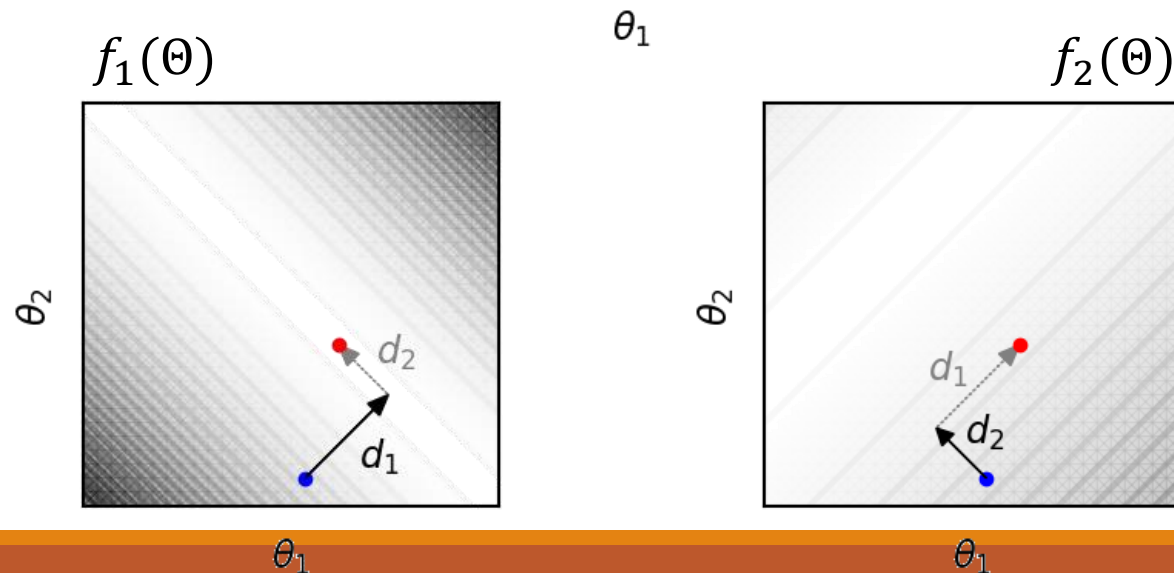


FOSI's Update Step

- ✓ d_1 brings f_1 to minimum in a single step
- ✓ Trivial matrix inversion
- ✓ d_1 orthogonal to d_2
- ✓ No mutual impact

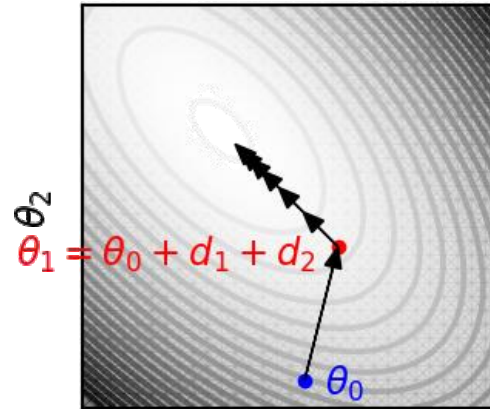


1. Gradient
2. Second-order info
3. Update θ_t
 - a. Split f to f_1 and f_2
 - b. Newton's step on f_1
 - c. Base opt step on f_2

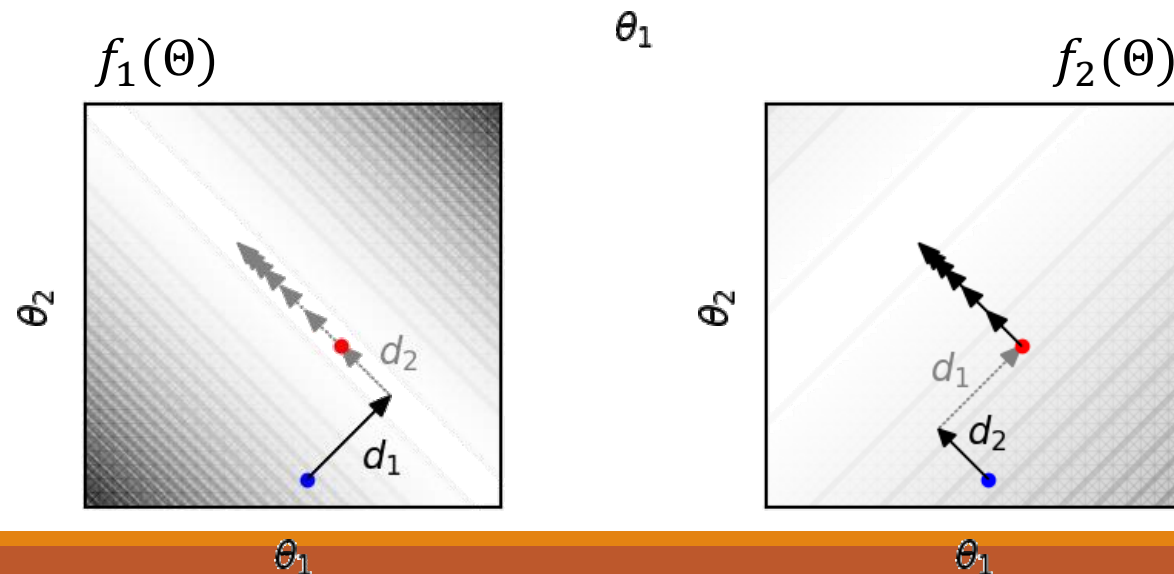


FOSI's Update Step

- ✓ d_1 brings f_1 to minimum in a single step
- ✓ Trivial matrix inversion
- ✓ d_1 orthogonal to d_2
- ✓ No mutual impact



1. Gradient
2. Second-order info
3. Update θ_t
 - a. Split f to f_1 and f_2
 - b. Newton's step on f_1
 - c. Base opt step on f_2



FOSI Algorithm

Inputs: $f(\theta)$, θ_0 , Base optimizer

Repeat until convergence:

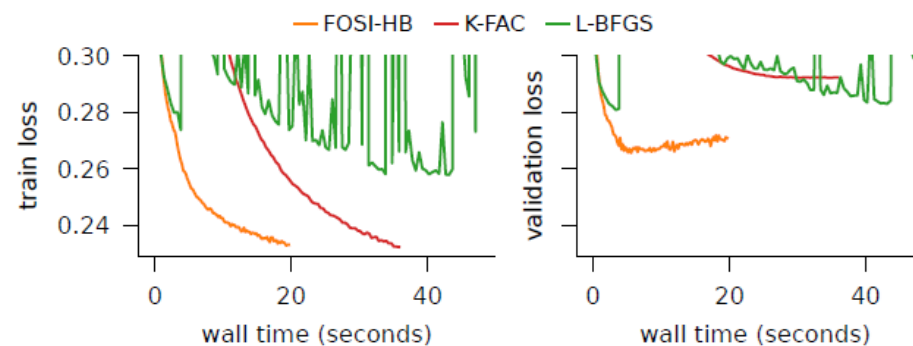
1. Gradient
2. Second-order info
3. Update θ_t
 - a. Split f to f_1 and f_2
 - b. Newton's step on f_1
 - c. Base opt step on f_2

Experiments and Results

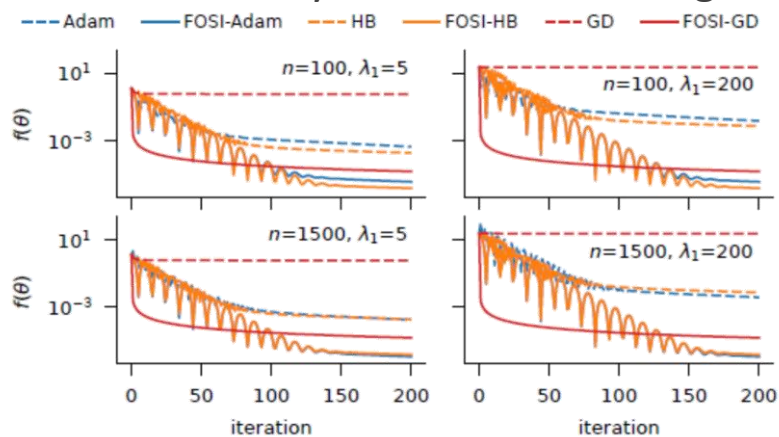
Comparison to base first-order optimizers

Task	HB	FOSI-HB	Adam	FOSI-Adam
AC	3822	1850 (40.4%)	5042	3911 (28.9%)
LM	269	207 (1.71)	270	219 (1.76)
AE	354	267 (52.46)	375	313 (51.26)
TL	93	53 (79.1%)	68	33 (79.0%)
LR	16	8 (92.8%)	12	18 (92.8%)

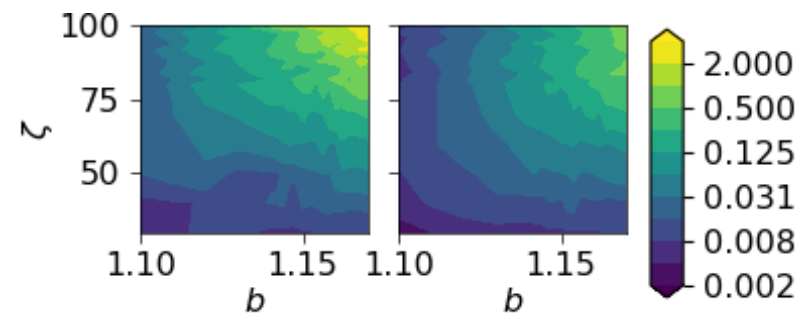
Comparison to second-order methods



Dimensionality and ill-conditioning



Ill-conditioning and diagonally dominance



Summary

- FOSI: hybrid-meta optimizer
 - ✓ **Meta-algorithm:** Applicable to any first order optimizer.
 - ✓ **Mathematical guarantees:** Improves effective condition number.
- Achieves the same loss as the base optimizer in 75% of the wall time.
- Open source: <https://github.com/hsivan/fosi>
 - ✓ Compatible with **JAX** (Optax) and **PyTorch** (TorchOpt)

