

SNM: Stochastic Newton Method for Optimization of Discrete Choice Models

Gael Lederrey

PhD student @ TRANSP-OR, EPFL

<http://glederrey.ch>



glederrey



gael.lederrey@epfl.ch

with

Pr. Virginie Lurkin (EPFL-TUE)

Pr. Michel Bierlaire (EPFL)

7th of November 2018



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE



Outline

1. Motivation
2. State-of-the-art
3. Optimization of DCMs
4. Future work
5. Conclusion



Motivation

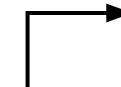
DCMs in a nutshell (1/2)



- Discrete Choice Theory (McFadden, 2000)
- Economics, Transportation, Health, ...
- Softwares: Biogeme, Larch, etc.



Interested? Have a look at **our MOOC on edX!** (Search “Discrete Choice”)
<https://courses.edx.org/courses/course-v1:EPFLx+DiscreteChoiceX+3T2017/course/>

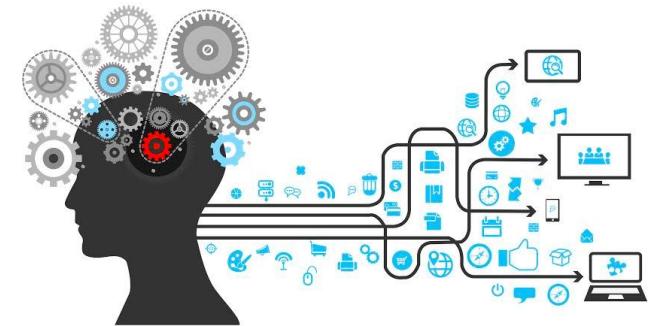


DCMs in a nutshell (2/2)

- Many success stories until now...

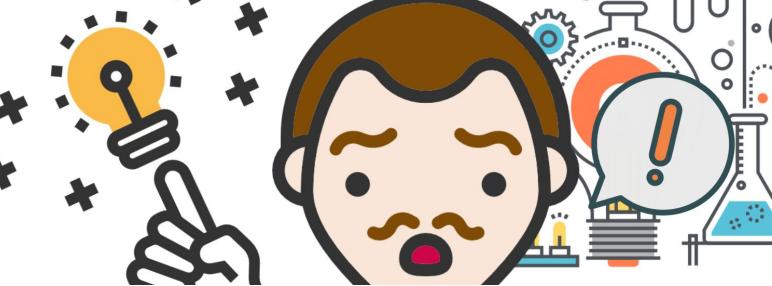


Solution: Machine Learning



- ML is **THE** field dealing with a lot of data!

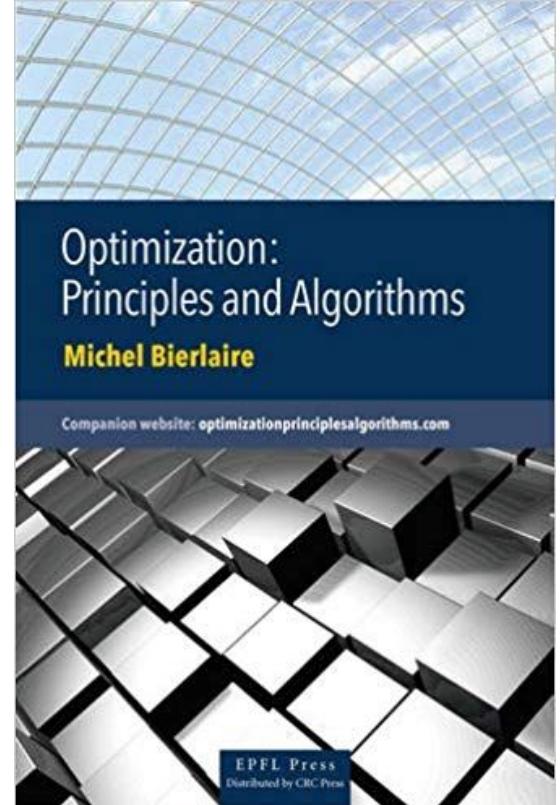
What can we do for DCMs?



1. Faster optimization of DCMs.

State-of-the-art

Optimization of DCM and ML



First-Order optimization - the ancestors

GD

(Cauchy, 1847)

Specificities:
Gradient computed on
all the data

Update step:

$$\theta = \theta - \alpha \cdot \nabla_{\theta} f(\theta; x)$$

Where

θ : Parameters

α : Step size

f : Function, $f \in C^1(\mathbb{R}^n)$

x : Data, $x \in \mathbb{R}^n$

SGD

(???, 1940's)

Specificities:
Gradient computed on
only one data

Update step:

$$\theta = \theta - \alpha \cdot \nabla_{\theta} f(\theta; x_i)$$

Where

θ : Parameters

α : Step size

f : Function, $f \in C^1(\mathbb{R}^n)$

x : Data, $x \in \mathbb{R}^n$

mbSGD

(???, 1940's)

Specificities:
Gradient computed on
a batch of data

Update step:

$$\theta = \theta - \alpha \cdot \nabla_{\theta} f(\theta; x_{\sigma(k)})$$

Where

θ : Parameters

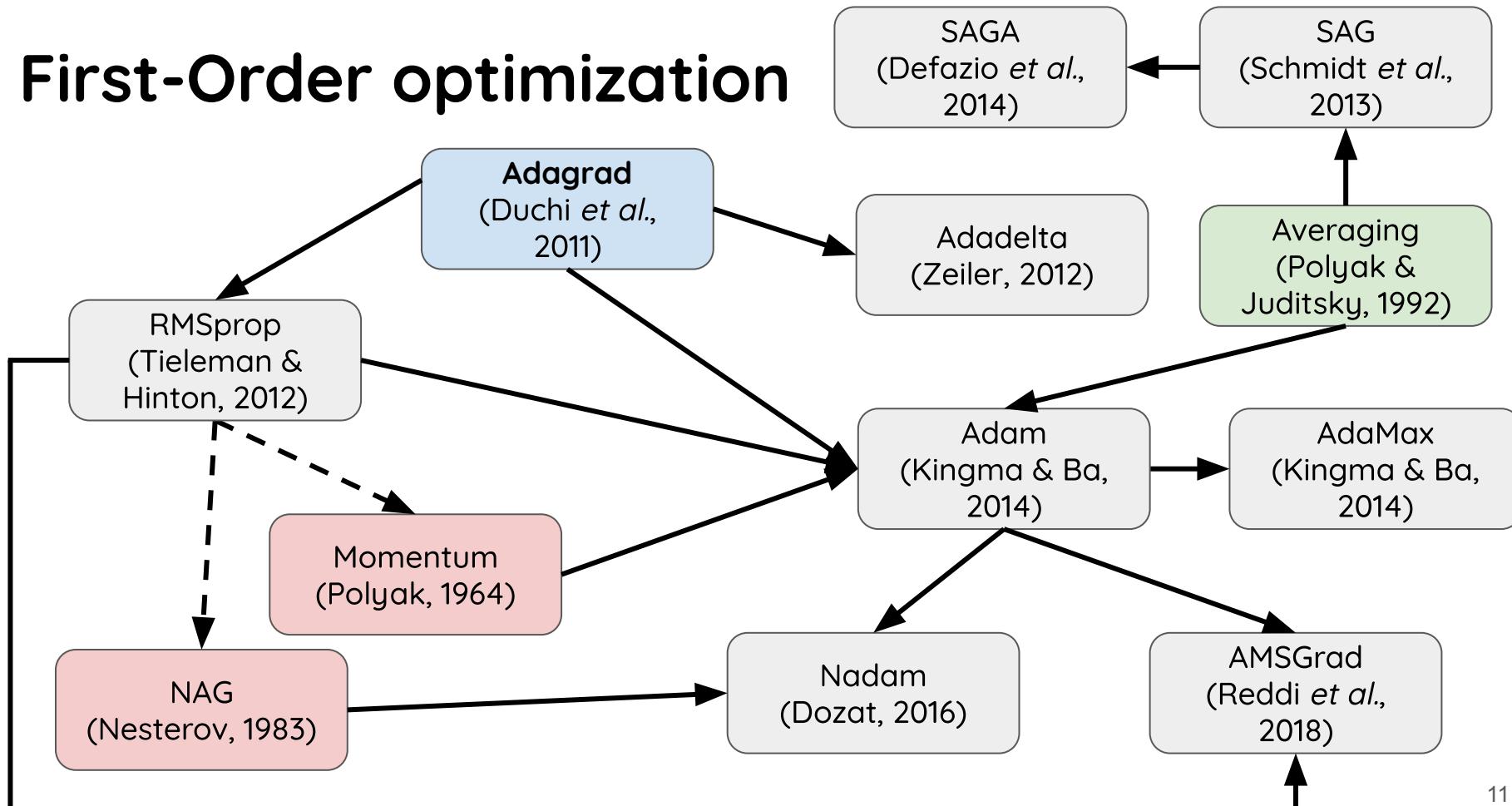
α : Step size

f : Function, $f \in C^1(\mathbb{R}^n)$

x : Data, $x \in \mathbb{R}^n$

$\sigma(k)$: Choice of k indices

First-Order optimization



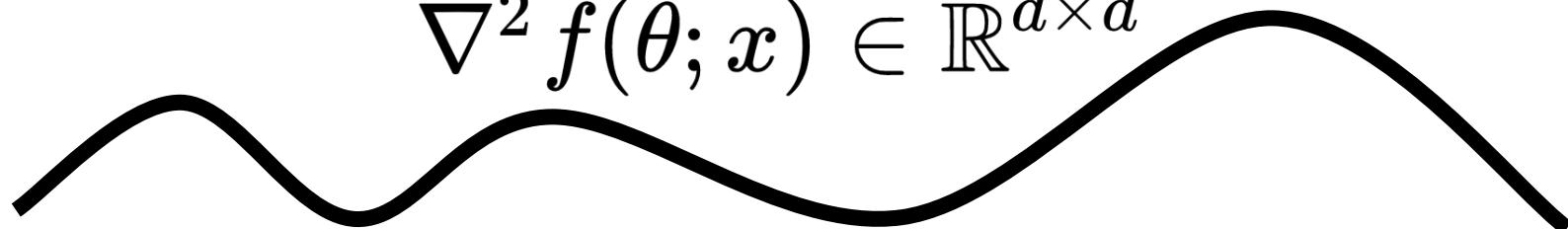
First-Order vs Second-Order

- Gradient is pretty cheap to compute

$$\nabla_{\theta} f(\theta; x) \in \mathbb{R}^d$$

- Computation of Hessian is difficult/impossible

$$\nabla^2 f(\theta; x) \in \mathbb{R}^{d \times d}$$



- Recently, more work on quasi-Newton methods.
Gower *et al.* (2018), Martens (2010), Bordes *et al.* (2009-2010)

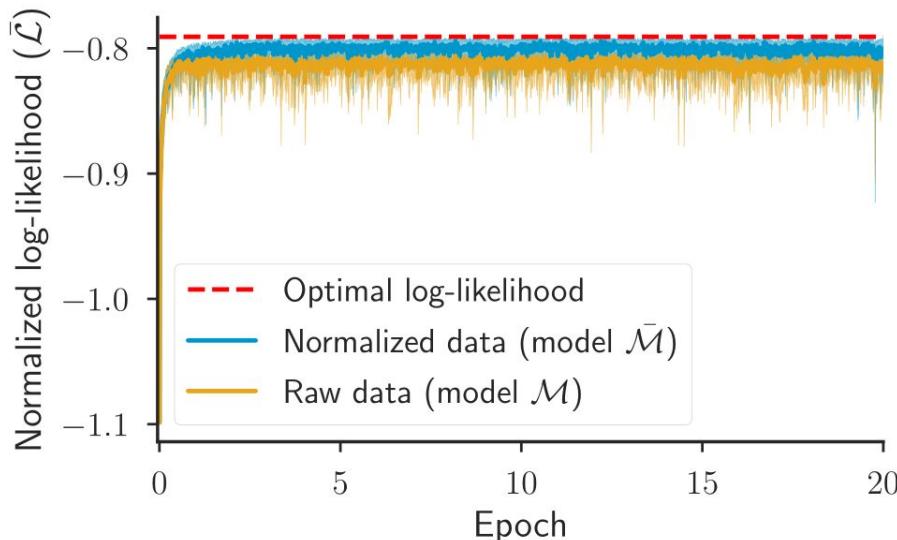
Optimization of DCMs

Simple MNL

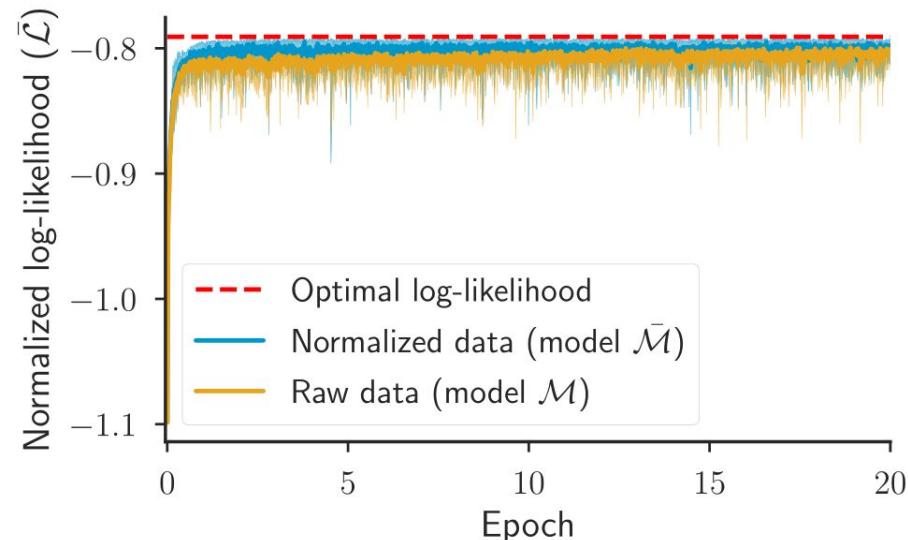
- Logit Model with 10 parameters (1 fixed)
- Normalized and raw observations (#9,036)
- **Normalized LL: -0.7908**
(Biogeme: ~10 epochs)



First-step - First-order



SGD

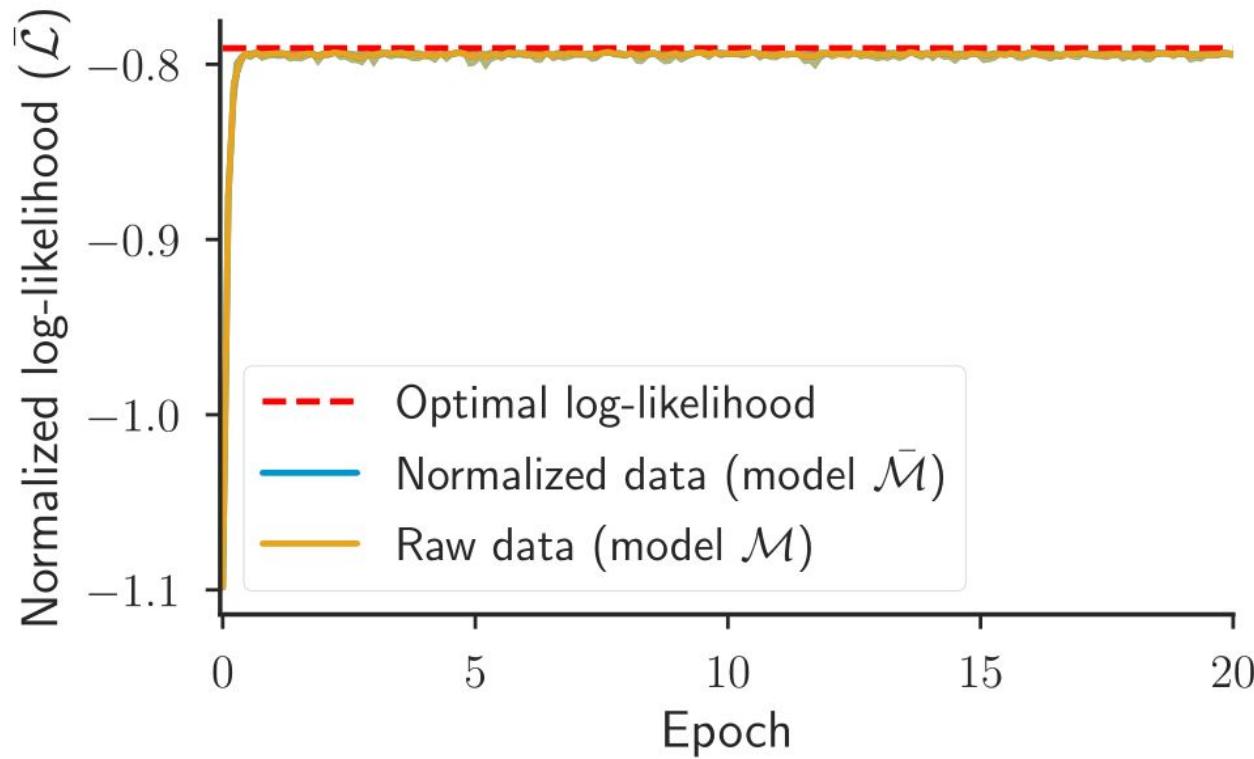


Adagrad

Why not Second-Order?

- Small number of parameters! (\neq Deep Learning)

SNM: The good, [...]



SNM: The bad, [...]

- Shift between Newton step and Gradient step
=> Conjugate Gradient

SNM: and the ugly!

- Still not able to converge!

~~Future~~ work
Current

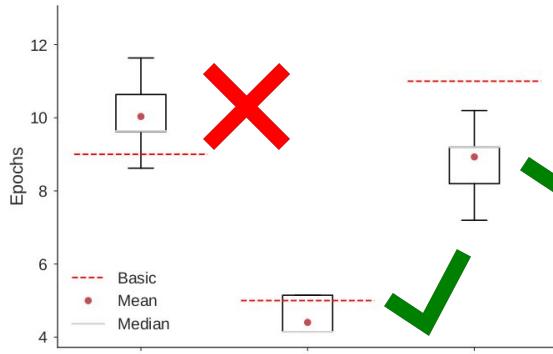
Adaptive Batch Size (ABS)



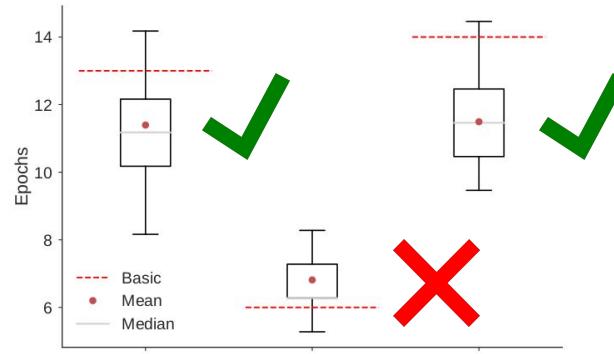
coming
soon!

Better: 10/12

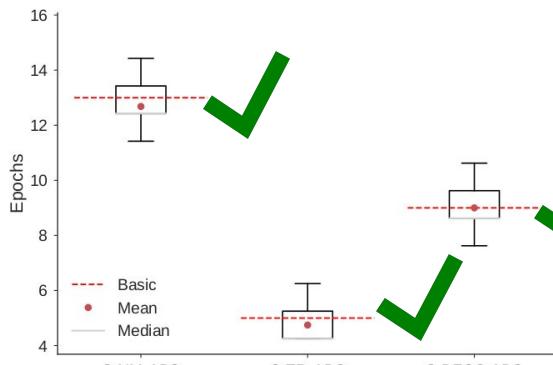
Results



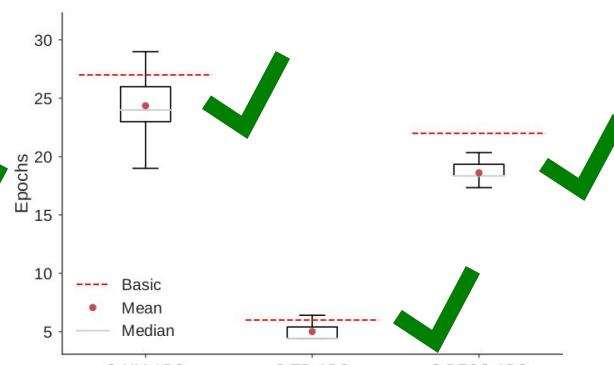
(a) MNL-SM



(b) NLM-SM



(c) LogReg-BS



(d) LogReg-BS-Full

Conclusion

Conclusion

- Estimation of DCMs has to be improved!
- Optimization is more tricky than expected!
- But it's working better everyday!

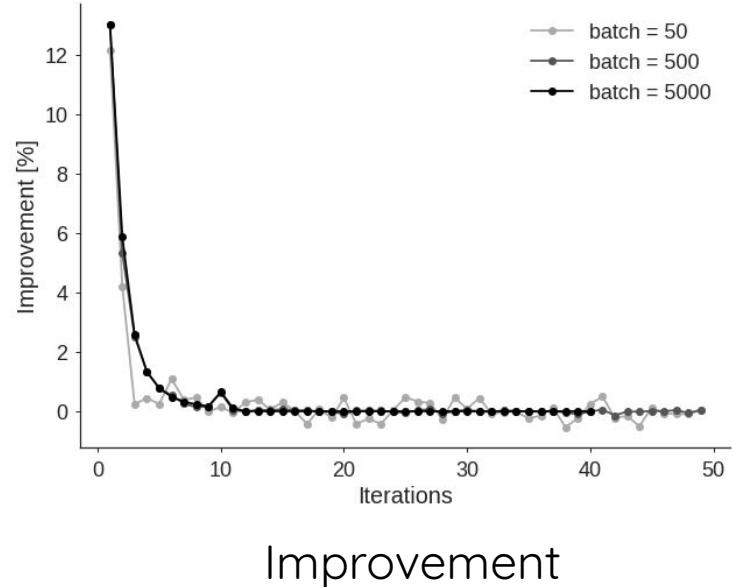
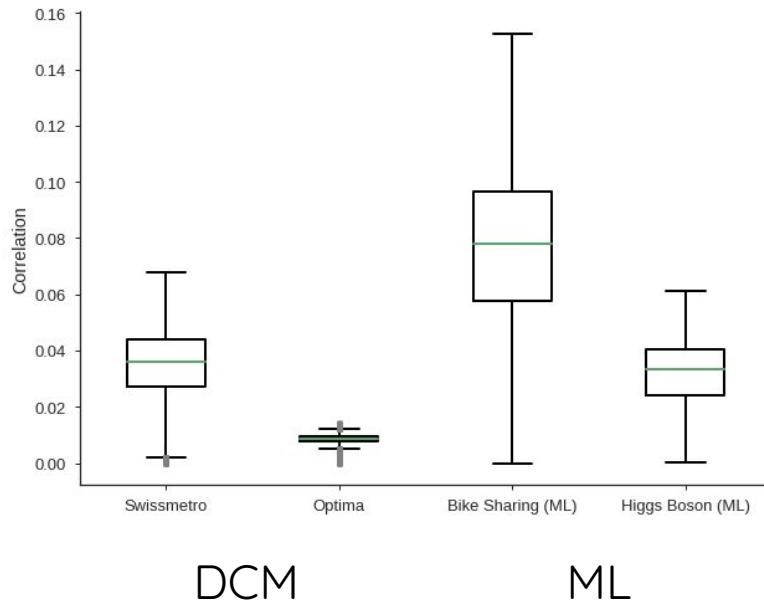
Mahalo



Back up
Slides

SNM not converging

- Data is not redundant enough!



SNM

Algorithm 1 Stochastic Newton Method (SNM)

Input: Initial parameter value (β_0), data (\mathcal{D}), function (f), gradient (∇f), Hessian ($\nabla^2 f$), #epochs (N_{ep}), batch size (N_{batch})

Output: Epochs (E), parameters (β), function values (f_v)

```
1: function SNM
2:    $(N_{\mathcal{D}}, M) = |\mathcal{D}|$                                 ▷ Number of samples and parameters
3:    $N_{iter} \leftarrow \lceil N_{ep} N_{\mathcal{D}} / N_{batch} \rceil$           ▷ Number of iterations
4:   Initialize  $E$ ,  $\beta$  and  $f_v$ . Set  $\beta[0] \leftarrow \beta_0$ 
5:   for  $i = 0 \dots N_{iter}$  do
6:      $E[i] \leftarrow i \cdot N_{batch} / N_{\mathcal{D}}$                       ▷ Store the epoch
7:      $f_v[i] \leftarrow f(\beta[i])$                                      ▷ Store the function value
8:      $idx \leftarrow N_{batch}$  indices from  $\mathcal{U}(0, N_{\mathcal{D}})$  without replacement
9:      $grad \leftarrow \nabla f_{idx}(\beta[i])$                            ▷ Gradient on the samples from  $idx$ 
10:     $hess \leftarrow \nabla^2 f_{idx}(\beta[i])$                          ▷ Hessian on the samples from  $idx$ 
11:    if  $hess$  is negative definite then
12:      Solve  $hess \cdot step = -grad$  to get  $step$                   ▷ Newton step
13:    else
14:       $step \leftarrow grad$                                          ▷ Gradient step
15:       $\alpha \leftarrow$  Backtracking Line Search with  $step$  on the subset of data with indices from  $idx$ 
16:       $\beta[i + 1] \leftarrow \beta[i] + \alpha \cdot step$ 
17:       $E[n_{iter}] \leftarrow N_{iter} \cdot N_{batch} / N_{\mathcal{D}}$ 
18:       $f_v[N_{iter}] \leftarrow f(\theta[N_{iter}])$ 
19:    return  $E$ ,  $\beta$  and  $f_v$ 
```
