

# Higher-order Factorization Machines



Mathieu Blondel

NTT Communication Science Laboratories  
Kyoto, Japan

Joint work with M. Ishihata, A. Fujino and  
N. Ueda

2016/9/28

# Regression analysis

---

- Variables

$y \in \mathbb{R}$ : target variable

$\mathbf{x} \in \mathbb{R}^d$ : explanatory variables (features)

- Training data

$$\mathbf{y} = [y_1, \dots, y_n]^T \in \mathbb{R}^n$$

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$$

- Goal

- Learn model parameters
- Compute prediction  $y$  for a new  $\mathbf{x}$

# Linear regression

---

- Model

$$\hat{y}_{LR}(\mathbf{x}; \mathbf{w}) := \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{j=1}^d w_j x_j$$

- Parameters

$\mathbf{w} \in \mathbb{R}^d$ : feature weights

- Pros and cons

- ☺  $O(d)$  predictions
- ☺ Learning  $\mathbf{w}$  can be cast as a convex optimization problem
- ☹ Does not use feature interactions

# Polynomial regression

---

- Model

$$\hat{y}_{PR}(\mathbf{x}; \mathbf{w}) := \langle \mathbf{w}, \mathbf{x} \rangle + \mathbf{x}^T \mathbf{W} \mathbf{x} = \langle \mathbf{w}, \mathbf{x} \rangle + \sum_{j,j'=1}^d w_{j,j'} x_j x_{j'}$$

- Parameters

$\mathbf{w} \in \mathbb{R}^d$ : feature weights

$\mathbf{W} \in \mathbb{R}^{d \times d}$ : weight matrix

- Pros and cons

☺ Learning  $\mathbf{w}$  and  $\mathbf{W}$  can be cast as a convex optimization problem

☹  $O(d^2)$  time and memory cost

# Kernel regression

---

- Model

$$\hat{y}_{KR}(\mathbf{x}; \boldsymbol{\alpha}) := \sum_{i=1}^n \alpha_i \mathcal{K}(\mathbf{x}_i, \mathbf{x})$$

- Parameters

$\boldsymbol{\alpha} \in \mathbb{R}^n$ : instance weights

- Pros and cons

- ☺ Can use non-linear kernels (RBF, polynomial, etc...)
- ☺ Learning  $\boldsymbol{\alpha}$  can be cast as a convex optimization problem
- ☹  $O(dn)$  predictions (linear dependence on training set size)

# Factorization Machines (FMs) (Rendle, ICDM 2010)

- Model

$$\hat{y}_{FM}(\mathbf{x}; \mathbf{w}, \mathbf{P}) := \langle \mathbf{w}, \mathbf{x} \rangle + \sum_{j' > j} \langle \bar{\mathbf{p}}_j, \bar{\mathbf{p}}_{j'} \rangle x_j x_{j'}$$

$j^{\text{th}}$  row of  $\mathbf{P}$

- Parameters

$\mathbf{w} \in \mathbb{R}^d$ : feature weights

$\mathbf{P} \in \mathbb{R}^{d \times k}$ : weight matrix

- Pros and cons

☺ Takes into account feature combinations

☺  $O(2dk)$  predictions (linear-time) instead of  $O(d^2)$

☹ Parameter estimation involves a non-convex optimization problem

# Application 1: recsys without features

- Formulate it as a matrix completion problem

	Movie 1	Movie 2	Movie 3	Movie 4
Alice	★★	?	★★★	?
Bob	★	?	★★	?
Charlie	★★	?	?	★★

- Matrix factorization: find  $\mathbf{U}$ ,  $\mathbf{V}$  that approximately reconstruct the rating matrix

$$\mathbf{R} \approx \mathbf{UV}^T$$

# Conversion to a regression problem

	Movie 1	Movie 2	Movie 3	Movie 4
Alice	**	?	***	?
Bob	*	?	**	?
Charlie	**	?	?	**

⇓ one-hot encoding

$$\underbrace{\begin{bmatrix} ** \\ *** \\ * \\ ** \\ ** \\ ** \end{bmatrix}}_y \quad \underbrace{\begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}}_x$$

Using this representation, FMs are equivalent to MF!



# Generalization ability of FMs

---

- The weight of  $x_j x_{j'}$  is  $\langle \bar{\mathbf{p}}_j, \bar{\mathbf{p}}_{j'} \rangle$  compared to  $w_{j,j'}$  for PR
- The same parameters  $\bar{\mathbf{p}}_j$  are shared for the weight of  $x_j x_{j'} \quad \forall j > j'$
- This increases the amount of data used to estimate  $\bar{\mathbf{p}}_j$  at the cost of introducing some bias (low-rank assumption)
- This allows to generalize to feature interactions that were not observed in the training set

## Application 2: recsys with features

Rating	User		Movie	
	Gender	Age	Genre	Director
★★	M	20-30	Adventure	S. Spielberg
★★★	F	0-10	Anime	H. Miyazaki
★	M	20-30	Drama	A. Kurosawa
⋮	⋮	⋮	⋮	⋮

- Interactions between **categorical variables**
  - Gender  $\times$  genre:  $\{M, F\} \times \{\text{Adventure, Anime, Drama, ...}\}$
  - Age  $\times$  director:  $\{0-10, 10-20, \dots\} \times \{S. Spielberg, H. Miyazaki, A. Kurosawa, \dots\}$
- In practice, the number of interactions can be **huge!**

# Conversion to regression

Rating	User		Movie	
	Gender	Age	Genre	Director
**	M	20-30	Adventure	S. Spielberg
***	F	0-10	Anime	H. Miyazaki
*	M	20-30	Drama	A. Kurosawa
⋮	⋮	⋮	⋮	⋮

⇓ one-hot encoding

$$\underbrace{\begin{bmatrix} ** \\ *** \\ * \\ \vdots \end{bmatrix}}_{\mathbf{y}} \quad \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & \dots \\ 0 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & \dots \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots \end{bmatrix}}_{\mathbf{X}}$$

very sparse  
binary data!

# FMs revisited (Blondel+, ICML 2016)

- ANOVA kernel of degree  $m = 2$  (Stitson+, 1997; Vapnik, 1998)

$$\mathcal{A}^2(\mathbf{p}, \mathbf{x}) := \sum_{j' > j} p_j x_j p_{j'} x_{j'}$$

- Then

$$\begin{aligned}\hat{y}_{FM}(\mathbf{x}; \mathbf{w}, \mathbf{P}) &= \langle \mathbf{w}, \mathbf{x} \rangle + \sum_{j' > j} \langle \bar{\mathbf{p}}_j, \bar{\mathbf{p}}_{j'} \rangle x_j x_{j'} \\ &= \langle \mathbf{w}, \mathbf{x} \rangle + \sum_{s=1}^k \mathcal{A}^2(\mathbf{p}_s, \mathbf{x})\end{aligned}$$

↑  $s^{\text{th}}$  column of  $\mathbf{P}$

# ANOVA kernel (arbitrary-order case)

---

- ANOVA kernel of degree  $2 \leq m \leq d$

$$\mathcal{A}^m(\mathbf{p}, \mathbf{x}) := \sum_{j_m > \dots > j_1} (p_{j_1} x_{j_1}) \dots (p_{j_m} x_{j_m})$$

↑ All possible  $m$ -combinations of  $\{1, \dots, d\}$

- Intuitively, the kernel uses all  $m$ -combinations of features without replacement:  $x_{j_1} \dots x_{j_m}$  for  $j_1 \neq \dots \neq j_m$
- Computing  $\mathcal{A}^m(\mathbf{p}, \mathbf{x})$  naively takes  $O(d^m)$  ☹

# Higher-order FMs (HOFMs)

---

- Model

$$\hat{y}_{HOFM}(\mathbf{x}; \mathbf{w}, \{\mathbf{P}^t\}_{t=2}^m) := \langle \mathbf{w}, \mathbf{x} \rangle + \sum_{t=2}^m \sum_{s=1}^k \mathcal{A}^t(\mathbf{p}_s^t, \mathbf{x})$$

- Parameters

$\mathbf{w} \in \mathbb{R}^d$ : feature weights

$\mathbf{P}^2, \dots, \mathbf{P}^m \in \mathbb{R}^{d \times k}$ : weight matrices

- Pros and cons

- ☺ Takes into account higher-order feature combinations
- ☺  $O(dkm^2)$  prediction cost using our proposed algorithms
- ☹ More complex than 2nd-order FMs

# Learning HOFMs (1/2)

---

- We use alternating minimization w.r.t.  $\mathbf{w}$ ,  $\mathbf{P}^2$ , ...,  $\mathbf{P}^m$
- Learning  $\mathbf{w}$  alone reduces to linear regression
- Learning  $\mathbf{P}^m$  can be cast as minimizing

$$F(\mathbf{P}) := \frac{1}{n} \sum_{i=1}^n \ell \left( y_i, \sum_{s=1}^k \mathcal{A}^m(\mathbf{p}_s, \mathbf{x}_i) + o_i \right) + \frac{\beta}{2} \|\mathbf{P}\|^2$$

where  $o_i$  is the contribution of degrees other than  $m$

# Learning HOFMs (2/2)

---

- Stochastic gradient update

$$\mathbf{p}_s \leftarrow \mathbf{p}_s - \eta \ell'(y_i, \hat{y}_i) \nabla \mathcal{A}^m(\mathbf{p}_s, \mathbf{x}_i) - \eta \beta \mathbf{p}_s$$

where  $\eta$  is a learning rate hyper-parameter and

$$\hat{y}_i := \sum_{s=1}^k \mathcal{A}^m(\mathbf{p}_s, \mathbf{x}_i) + o_i$$

- We propose  $O(dm)$  (linear time) DP algorithms for
  - Evaluating ANOVA kernel  $\mathcal{A}^m(\mathbf{p}, \mathbf{x}) \in \mathbb{R}$
  - Computing gradient  $\nabla \mathcal{A}^m(\mathbf{p}, \mathbf{x}) \in \mathbb{R}^d$



# Evaluating the ANOVA kernel (1/3)

---

- Recursion (Blondel+, ICML 2016)

$$\mathcal{A}^m(\mathbf{p}, \mathbf{x}) = \mathcal{A}^m(\mathbf{p}_{\neg j}, \mathbf{x}_{\neg j}) + p_j x_j \mathcal{A}^{m-1}(\mathbf{p}_{\neg j}, \mathbf{x}_{\neg j}) \quad \forall j$$

where  $\mathbf{p}_{\neg j}, \mathbf{x}_{\neg j} \in \mathbb{R}^{d-1}$  are vectors with the  $j^{\text{th}}$  element removed

- We can use this recursion to remove features until computing the kernel becomes trivial

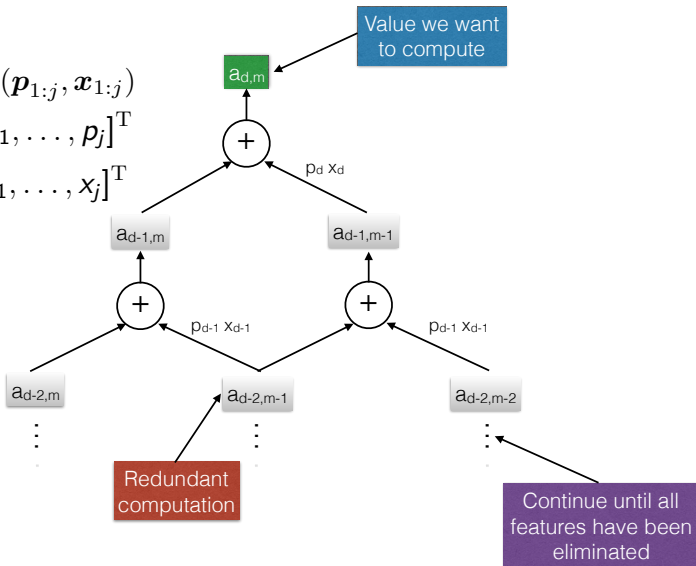
# Evaluating the ANOVA kernel (2/3)

Shortcut:

$$a_{j,t} := \mathcal{A}^t(\mathbf{p}_{1:j}, \mathbf{x}_{1:j})$$

$$\mathbf{p}_{1:j} := [p_1, \dots, p_j]^T$$

$$\mathbf{x}_{1:j} := [x_1, \dots, x_j]^T$$



# Evaluating the ANOVA kernel (3/3)

Ways to avoid redundant computations:

- Top-down approach with memory table
- **Bottom-up dynamic programming (DP)**

	$j = 0$	$j = 1$	$j = 2$	$\dots$	$j = d$
$t = 0$	1	1	1	1	1
$t = 1$	0	$\overset{\text{start}}{\rightarrow} a_{1,1}$	$a_{2,1}$	$\dots$	$a_{d,1}$
$t = 2$	0	0	$\overset{p_{2 \times 2}}{\rightarrow} a_{2,2}$	$\dots$	$a_{d,2}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$t = m$	0	0	0	$\dots$	$a_{d,m} \leftarrow \text{goal}$

---

**Algorithm 1** Evaluating  $\mathcal{A}^m(\mathbf{p}, \mathbf{x})$  in  $O(dm)$ 

---

**Input:**  $\mathbf{p} \in \mathbb{R}^d$ ,  $\mathbf{x} \in \mathbb{R}^d$

$a_{j,t} \leftarrow 0 \forall t \in \{1, \dots, m\}, j \in \{0, 1, \dots, d\}$

$a_{j,0} \leftarrow 1 \forall j \in \{0, 1, \dots, d\}$

**for**  $t := 1, \dots, m$  **do**

**for**  $j := t, \dots, d$  **do**

$a_{j,t} \leftarrow a_{j-1,t} + p_j x_j a_{j-1,t-1}$

**end for**

**end for**

**Output:**  $\mathcal{A}^m(\mathbf{p}, \mathbf{x}) = a_{d,m}$

---

# Backpropagation (chain rule)

---

Ex: compute derivatives of composite function  $f(g(h(\mathbf{p})))$

- Forward pass


$$a = h(\mathbf{p})$$

$$b = g(a)$$

$$c = f(b)$$

- Backward pass

Only the last part  
depends on  $j$

$$\frac{\partial c}{\partial p_j} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial p_j} = f'(b) g'(a) h'_j(\mathbf{p})$$


Can compute all derivatives in one pass!

# Gradient computation (1/2)

- We want to compute  $\nabla \mathcal{A}^m(\mathbf{p}, \mathbf{x}) = [\tilde{p}_1, \dots, \tilde{p}_d]^T$
- Using the chain rule, we have

$$\tilde{p}_j := \frac{\partial a_{d,m}}{\partial p_j} = \sum_{t=1}^m \underbrace{\frac{\partial a_{d,m}}{\partial a_{j,t}}}_{:= \tilde{a}_{j,t}} \underbrace{\frac{\partial a_{j,t}}{\partial p_j}}_{= a_{j-1,t-1} x_j} = \sum_{t=1}^m \tilde{a}_{j,t} a_{j-1,t-1} x_j$$

since  $p_j$  influences  $a_{j,t} \forall t \in [m]$

- $\tilde{a}_{j,t}$  can be computed recursively in reverse order

$$\tilde{a}_{j,t} = \tilde{a}_{j+1,t} + p_{j+1} x_{j+1} \tilde{a}_{j+1,t+1}$$

## Gradient computation (2/2)

goal	$j = 1$	$j = 2$	$\dots$	$j = d - 1$	$j = d$
$t = 1$	$\tilde{a}_{1,1}$	$\tilde{a}_{2,1}$	$\dots$	0	0
$t = 2$	0	$\tilde{a}_{2,2}$	$\dots$	0	0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\tilde{a}_{d-1,t-1}$	0
$t = m$	0	0	0	1	1 ← start
$t = m + 1$	0	0	0	0	0

---

**Algorithm 2** Computing  $\nabla \mathcal{A}^m(\mathbf{p}, \mathbf{x})$  in  $O(dm)$ 

---

**Input:**  $\mathbf{p} \in \mathbb{R}^d$ ,  $\mathbf{x} \in \mathbb{R}^d$ ,  $\{a_{j,t}\}_{j,t=0}^{d,m}$

$\tilde{a}_{j,t} \leftarrow 0 \forall t \in [m+1], j \in [d]$

$\tilde{a}_{d,m} \leftarrow 1$

**for**  $t := m, \dots, 1$  **do**

**for**  $j := d-1, \dots, t$  **do**

$\tilde{a}_{j,t} \leftarrow \tilde{a}_{j+1,t} + \tilde{a}_{j+1,t+1} p_{j+1} x_{j+1}$

**end for**

**end for**

$\tilde{p}_j := \sum_{t=1}^m \tilde{a}_{j,t} a_{j-1,t-1} x_j \forall j \in [d]$

**Output:**  $\nabla \mathcal{A}^m(\mathbf{p}, \mathbf{x}) = [\tilde{p}_1, \dots, \tilde{p}_d]^T$

---



# Summary so far

---

- HOFMs can be expressed using the ANOVA kernel  $\mathcal{A}^m$
- We proposed  $O(dm)$  time algorithms for computing  $\mathcal{A}^m(\mathbf{p}, \mathbf{x})$  and  $\nabla \mathcal{A}^m(\mathbf{p}, \mathbf{x})$
- The cost per epoch of stochastic gradient algorithms for learning  $\mathbf{P}^m$  is therefore  $O(dnkm)$
- The prediction cost is  $O(dkm^2)$

# Other contributions

---

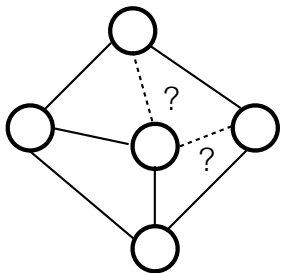
- Coordinate-descent algorithm for learning  $\mathbf{P}^m$  based on a different recursion
  - Cost per epoch is  $O(dnkm^2)$  😞
  - However, no learning rate to tune! 😊
- HOFMs with shared parameters:  $\mathbf{P}^2 = \dots = \mathbf{P}^m$ 
  - Total prediction cost is  $O(dkm)$  instead of  $O(dkm^2)$  😊
  - Corresponds to using new kernels derived from the ANOVA kernel

# Experiments

# Application to link prediction

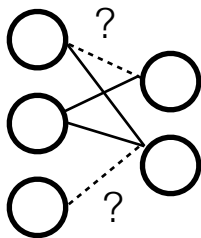
---

Goal: predict missing links between nodes in a graph



Graph:

- Co-author network
- Enzyme network



Bipartite graph:

- User-movie
- Gene-disease

# Application to link prediction

---

- We assume two sets of nodes  $A$  (e.g., users) and  $B$  (e.g., movies) of size  $n_A$  and  $n_B$
- Nodes in  $A$  are represented by feature vectors  $\mathbf{a}_i \in \mathbb{R}^{d_A}$
- Nodes in  $B$  are represented by feature vectors  $\mathbf{b}_j \in \mathbb{R}^{d_B}$
- We are given a matrix  $\mathbf{Y} \in \{-1, +1\}^{n_A \times n_B}$  such that  $y_{i,j} = +1$  if there is a link between  $\mathbf{a}_i$  and  $\mathbf{b}_j$
- Number of positive samples is  $n_+$

# Datasets

Dataset	$n_+$	Columns of $\mathbf{A}$	$n_A$	$d_A$	Columns of $\mathbf{B}$	$n_B$	$d_B$
NIPS	4,140	Authors	2,037	13,649			
Enzyme	2,994	Enzymes	668	325			
GD	3,954	Diseases	3,209	3,209	Genes	12,331	25,275
ML 100K	21,201	Users	943	49	Movies	1,682	29

Features:


- NIPS: word occurrence in author publications
- Enzyme: phylogenetic information, gene expression information and gene location information
- GD: MimMiner similarity scores (diseases) and HumanNet similarity scores (genes)
- ML 100K: age, gender, occupation, living area (users); release year, genre (movies)

# Models compared

---

Goal: predict if there is a link between  $\mathbf{a}_i$  and  $\mathbf{b}_j$

vector concatenation

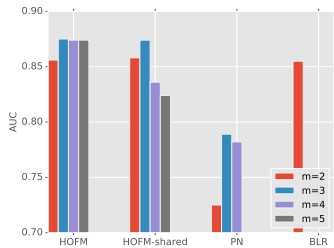
- HOFM:  $\hat{y}_{i,j} = \hat{y}_{HOFM}(\mathbf{a}_i \oplus \mathbf{b}_j; \mathbf{w}, \{\mathbf{P}^t\}_{t=2}^m)$   

- HOFM-shared: same but with  $\mathbf{P}^2 = \dots = \mathbf{P}^m$
- Polynomial network (PN): replace ANOVA kernel by polynomial kernel
- Bilinear regression (BLR):  $\hat{y}_{i,j} = \mathbf{a}_i \mathbf{U} \mathbf{V}^T \mathbf{b}_j$

# Experimental protocol

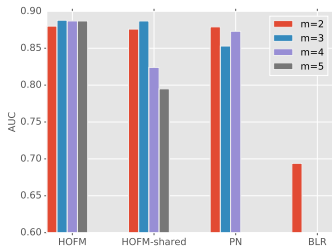
---

- We sample  $n_- = n_+$  negatives samples (missing edges are treated as negative samples)
- We use 50% for training and 50% for testing
- We use ROC-AUC (area under ROC curve) for evaluation
- $\beta$  tuned by CV,  $k$  fixed to 30
- $\mathbf{P}^2, \dots, \mathbf{P}^m$  initialized randomly
- $\ell$  is set to the squared loss

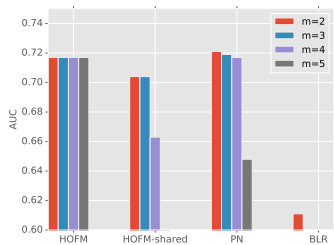




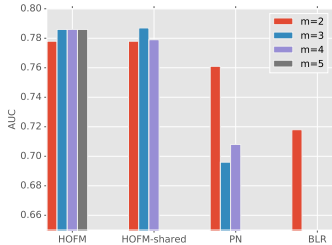
(a) NIPS



(b) Enzyme



(c) GD



(d) ML100K

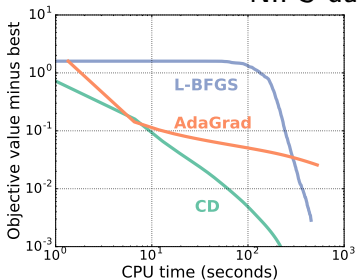
# Solver comparison

---

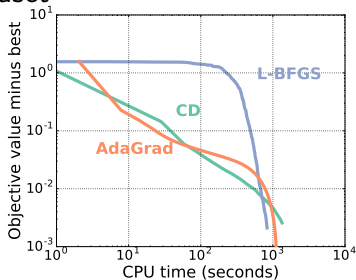
- Coordinate descent
- AdaGrad
- L-BFGS

AdaGrad and L-BFGS use the proposed DP algorithm to compute  $\nabla \mathcal{A}^m(\mathbf{p}, \mathbf{x})$

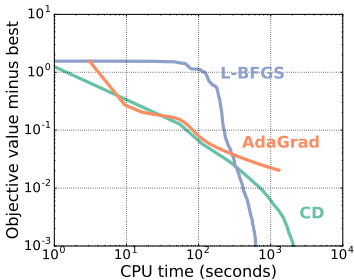
# NIPS dataset



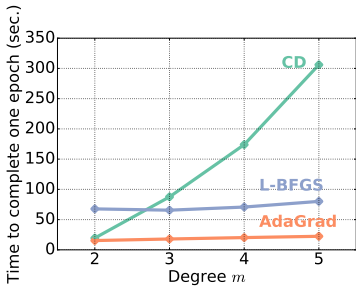
(a) Convergence when  $m = 2$



(b) Convergence when  $m = 3$



(c) Convergence when  $m = 4$



(d) Scalability w.r.t.  $m$