

# Convex and non-convex worlds in machine learning

**Anna Choromanska**

Courant Institute of Mathematical Sciences  
New York University

## Convex and non-convex worlds

**Machine learning and optimization** - many machine learning problems are formulated as minimization of some loss function on a training set of examples. Loss functions express the discrepancy between the predictions of the model being trained and the actual problem instances. Optimization algorithms can then minimize this loss. (*Wikipedia*)

# Convex and non-convex worlds

**Machine learning and optimization** - many machine learning problems are formulated as minimization of some loss function on a training set of examples. Loss functions express the discrepancy between the predictions of the model being trained and the actual problem instances. Optimization algorithms can then minimize this loss. (*Wikipedia*)

## Convex world

**local min = global min**

**strictly convex: unique min**

**efficient solvers**

**strong theoretical guarantees**

# Convex and non-convex worlds

**Machine learning and optimization** - many machine learning problems are formulated as minimization of some loss function on a training set of examples. Loss functions express the discrepancy between the predictions of the model being trained and the actual problem instances. Optimization algorithms can then minimize this loss. (*Wikipedia*)

## Convex world

local min = global min  
strictly convex: unique min  
efficient solvers  
strong theoretical guarantees

## Non-convex world

multiple local min  $\neq$  global min  
many solvers come from convex world  
weak theoretical guarantees if any

# Layout of the talk

**Optimization solvers:** generic optimization vs bound majorization  
**partition function-based objectives**

**Design (convex) solver: quadratic bound majorization**

[JC12, CKJ12, ACJK14]

## Layout of the talk

**Optimization solvers:** generic optimization vs bound majorization  
partition function-based objectives

**Design (convex) solver:** quadratic bound majorization  
[JC12, CKJ12, ACJK14]

**Challenging problems:** multi-class classification

**Design objective:** statistical and computational constraints  
online multi-class partition trees for logarithmic time predictions  
[CL14, CAL13, CCB15, CCJM15, CCJM13, CJKMM13, BCCL15, CM12]

# Layout of the talk

**Optimization solvers:** generic optimization vs bound majorization  
**partition function-based objectives**

**Design (convex) solver:** quadratic bound majorization  
 [JC12, CKJ12, ACJK14]

**Challenging problems:** multi-class classification

**Design objective:** statistical and computational constraints  
**online multi-class partition trees for logarithmic time predictions**  
 [CL14, CAL13, CCB15, CCJM15, CCJM13, CJKMM13, BCCL15, CM12]

**Non-convex problems:** deep learning  
**highly non-convex objective**

**Build understanding:** new theoretical results  
 [CHMCL15, CLB15, ZCL15]

# How to build good efficient convex solver?



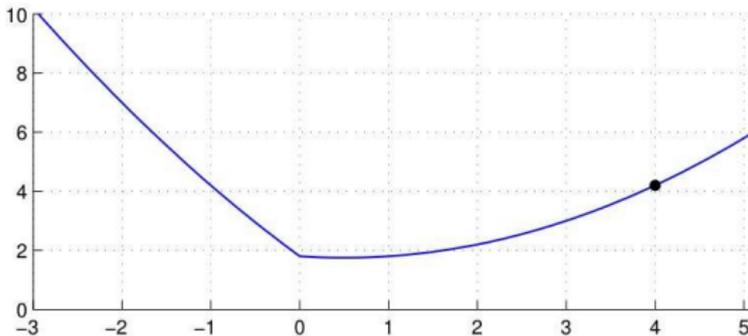
# Bound majorization

If cost function  $\theta^* = \arg \min_{\theta} C(\theta)$  has no closed form solution

Majorization uses a surrogate  $Q$  with closed form solution

Monotonically improves from initial  $\theta_0$

- Find bound  $Q(\theta, \theta_i) \geq C(\theta)$  where  $Q(\theta_i, \theta_i) = C(\theta_i)$
- Update  $\theta_{i+1} = \arg \min_{\theta} Q(\theta, \theta_i)$
- Repeat until converged































# Experiments - (Semi-)Stochastic Bound Majorization

- Computing the bound is  $\mathcal{O}(t)$   $\rightarrow$  intractable for large  $t$
- Semi-stochastic: compute bound on data mini-batches
  - convergence to a stationary point under weak assumptions (in particular convexity is not required)
  - **linear** convergence rate for logistic regression problem when batch size grows sufficiently fast

# Experiments - (Semi-)Stochastic Bound Majorization

- Computing the bound is  $\mathcal{O}(t) \rightarrow$  intractable for large  $t$
- Semi-stochastic: compute bound on data mini-batches
  - convergence to a stationary point under weak assumptions (in particular convexity is not required)
  - **linear** convergence rate for logistic regression problem when batch size grows sufficiently fast

## Theorem

For each iteration we have (for any  $\epsilon > 0$ )

$$J(\boldsymbol{\theta}_k) - J(\boldsymbol{\theta}^*) \leq \left(1 - \frac{\rho}{L}\right)^k [J(\boldsymbol{\theta}_0) - J(\boldsymbol{\theta}^*)] + \mathcal{O}(C_k),$$

with  $C_k = \max\{B_k, (1 - \frac{\rho}{L} + \epsilon)^k\}$  and  $B_k = \|\nabla J(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_k} - \mathbf{g}_{\mathcal{T}}^k\|^2$ , where  $\mathcal{T}$  is the mini-batch.



# How to design good objective function?

# eXtreme multi-class classification problem

## Problem setting

- classification with large number of classes
- data is accessed online

## Goal:

- good predictor with **logarithmic training and testing time**
- reduction to tree-structured binary classification
- top-down approach for class partitioning allowing gradient descent style optimization



# How do you learn the structure?

- Not all partitions are equally difficult, e.g.
  - if you do  $\{1, 7\}$  vs  $\{3, 8\}$ , the next problem is hard;
  - if you do  $\{1, 8\}$  vs  $\{3, 7\}$ , the next problem is easy;
  - if you do  $\{1, 3\}$  vs  $\{7, 8\}$ , the next problem is easy.

# How do you learn the structure?

- Not all partitions are equally difficult, e.g.
  - if you do  $\{1, 7\}$  vs  $\{3, 8\}$ , the next problem is hard;
  - if you do  $\{1, 8\}$  vs  $\{3, 7\}$ , the next problem is easy;
  - if you do  $\{1, 3\}$  vs  $\{7, 8\}$ , the next problem is easy.
- [BWG10]: **Better to confuse near leaves than near root.**  
Intuition: The root predictor tends to be overconstrained while the leafwards predictors are less constrained.

# How do you learn the structure?

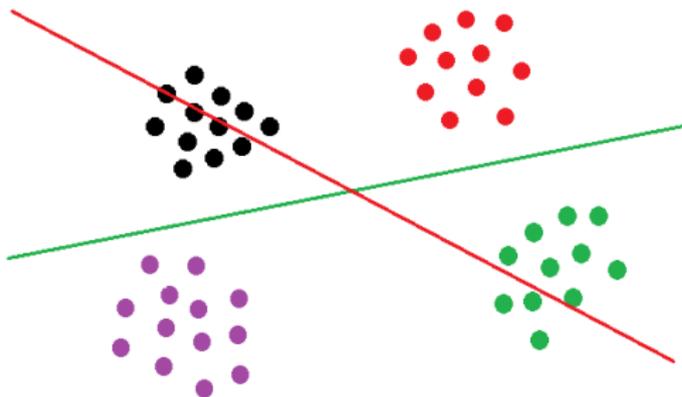
Our approach:

- **top-down** approach for class partitioning
- splitting criterion guaranteeing  
**balanced tree**  $\Rightarrow$  **logarithmic training and testing time**  
and  
**small classification error**





# Pure split and balanced split



- $k_r(x)$ : number of data points in the same class as  $x$  on the right side of the partitioning
- $k(x)$ : total number of data points in the same class as  $x$
- $n_r$ : number of data points on the right side of the partitioning
- $n$ : total number of data points

Measure of **balanceness**:  $\frac{n_r}{n}$



## Pure split and balanced split

- $k$ : number of classes
- $\mathcal{H}$ : hypothesis class (typically: linear classifiers)
- $\pi_y = \frac{|\mathcal{X}_y|}{n}$
- balance =  $Pr(h(x) > 0)$
- purity =  $\sum_{y=1}^k \pi_y \min(Pr(h(x) > 0|y), Pr(h(x) < 0|y))$



## Pure split and balanced split

- $k$ : number of classes
- $\mathcal{H}$ : hypothesis class (typically: linear classifiers)
- $\pi_y = \frac{|\mathcal{X}_y|}{n}$
- $\text{balance} = \Pr(h(x) > 0)$
- $\text{purity} = \sum_{y=1}^k \pi_y \min(\Pr(h(x) > 0|y), \Pr(h(x) < 0|y))$

### Definition (Balanced split)

The hypothesis  $h \in \mathcal{H}$  induces a balanced split iff

$$\exists_{c \in (0, 0.5]} c \leq \text{balance} \leq 1 - c.$$

### Definition (Pure split)

The hypothesis  $h \in \mathcal{H}$  induces a pure split iff

$$\exists_{\delta \in [0, 0.5)} \text{purity} \leq \delta.$$

# Objective function

$$\begin{aligned}
 J(h) &= 2 \sum_{y=1}^k \pi_y |P(h(x) > 0) - P(h(x) > 0|y)| \\
 &= 2\mathbb{E}_{x,y} [|P(h(x) > 0) - P(h(x) > 0|y)|]
 \end{aligned}$$

$J(h) \Rightarrow$  **Splitting criterion (objective function)**

Given a set of  $n$  examples each with one of  $k$  labels, find a **partitioner**  $h$  that maximizes the objective.

## Lemma

*For any hypothesis  $h : \mathcal{X} \mapsto \{-1, 1\}$ , the objective  $J(h)$  satisfies  $J(h) \in [0, 1]$ . Furthermore,  $h$  induces a maximally pure and balanced partition iff  $J(h) = 1$ .*





# What is the quality of obtained tree?

- In each node of the tree  $\mathcal{T}$  optimize the splitting criterion
- Apply recursively to construct a tree structure
- Measure the quality of the tree using entropy

$$G_{\mathcal{T}} = \sum_{I \in \text{leafs of } \mathcal{T}} w_I \sum_{y=1}^k \pi_{I,y} \ln \left( \frac{1}{\pi_{I,y}} \right)$$

Why?

Small entropy of leafs  $\Rightarrow$  pure leafs

**Goal:** maximizing the objective reduces the entropy

# What is the quality of obtained tree?

## Definition (Weak Hypothesis Assumption)

Let  $m$  denotes any node of the tree  $\mathcal{T}$ , and let  $\beta_m = P(h_m(x) > 0)$  and  $P_{m,i} = P(h_m(x) > 0|i)$ . Furthermore, let  $\gamma \in \mathbb{R}^+$  be such that for all  $m$ ,  $\gamma \in (0, \min(\beta_m, 1 - \beta_m)]$ . We say that the *weak hypothesis assumption* is satisfied when for any distribution  $\mathcal{P}$  over  $\mathcal{X}$  at each node  $m$  of the tree  $\mathcal{T}$  there exists a hypothesis  $h_m \in \mathcal{H}$  such that  $J(h_m)/2 = \sum_{i=1}^k \pi_{m,i} |P_{m,i} - \beta_m| \geq \gamma$ .

## Theorem

*Under the Weak Hypothesis Assumption, for any  $\epsilon \in [0, 1]$ , to obtain  $G_{\mathcal{T}} \leq \epsilon$  it suffices to make  $(\frac{1}{\epsilon})^{\frac{4(1-\gamma)^2 \ln k}{\gamma^2}}$  splits.*

# What is the quality of obtained tree?

## Definition (Weak Hypothesis Assumption)

Let  $m$  denotes any node of the tree  $\mathcal{T}$ , and let  $\beta_m = P(h_m(x) > 0)$  and  $P_{m,i} = P(h_m(x) > 0|i)$ . Furthermore, let  $\gamma \in \mathbb{R}^+$  be such that for all  $m$ ,  $\gamma \in (0, \min(\beta_m, 1 - \beta_m)]$ . We say that the *weak hypothesis assumption* is satisfied when for any distribution  $\mathcal{P}$  over  $\mathcal{X}$  at each node  $m$  of the tree  $\mathcal{T}$  there exists a hypothesis  $h_m \in \mathcal{H}$  such that  $J(h_m)/2 = \sum_{i=1}^k \pi_{m,i} |P_{m,i} - \beta_m| \geq \gamma$ .

## Theorem

Under the Weak Hypothesis Assumption, for any  $\epsilon \in [0, 1]$ , to obtain  $G_{\mathcal{T}} \leq \epsilon$  it suffices to make  $(\frac{1}{\epsilon})^{\frac{4(1-\gamma)^2 \ln k}{\gamma^2}}$  splits.

- Tree depth  $\approx \log \left[ \left( \frac{1}{\epsilon} \right)^{\frac{4(1-\gamma)^2 \ln k}{\gamma^2}} \right] = \mathcal{O}(\ln k) \Rightarrow$   
 $\Rightarrow$  **logarithmic training and testing time**

# LOMtree algorithm

- Recall the objective function we consider at every tree node:

$$J(h) = 2\mathbb{E}_y[|\mathbb{E}_x[\mathbb{1}(h(x) > 0)] - \mathbb{E}_x[\mathbb{1}(h(x) > 0|y)]|].$$

Problem: discrete optimization

Relaxation: drop the indicator operator and look at margins

# LOMtree algorithm

- Recall the objective function we consider at every tree node:

$$J(h) = 2\mathbb{E}_y[|\mathbb{E}_x[\mathbb{1}(h(x) > 0)] - \mathbb{E}_x[\mathbb{1}(h(x) > 0|y)]|].$$

Problem: discrete optimization

Relaxation: drop the indicator operator and look at margins

# LOMtree algorithm

- Recall the objective function we consider at every tree node:

$$J(h) = 2\mathbb{E}_y[|\mathbb{E}_x[\mathbb{1}(h(x) > 0)] - \mathbb{E}_x[\mathbb{1}(h(x) > 0|y)]|].$$

Problem: discrete optimization

Relaxation: drop the indicator operator and look at margins

- The objective function becomes

$$J(h) = 2\mathbb{E}_y[|\mathbb{E}_x[h(x)] - \mathbb{E}_x[h(x)|y]|]$$

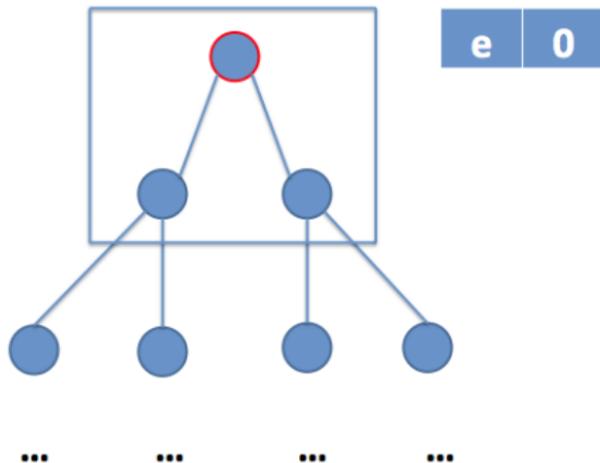
- Keep the online empirical estimates of these expectations.
- The sign of the difference of two expectations decides whether to send an example to the left or right child node.

## LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

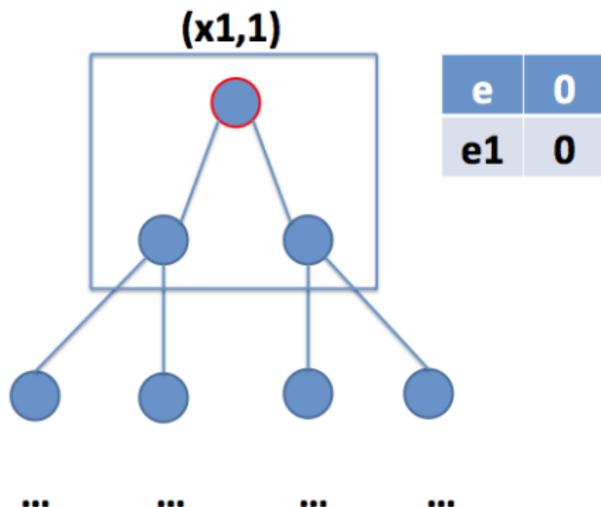


# LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

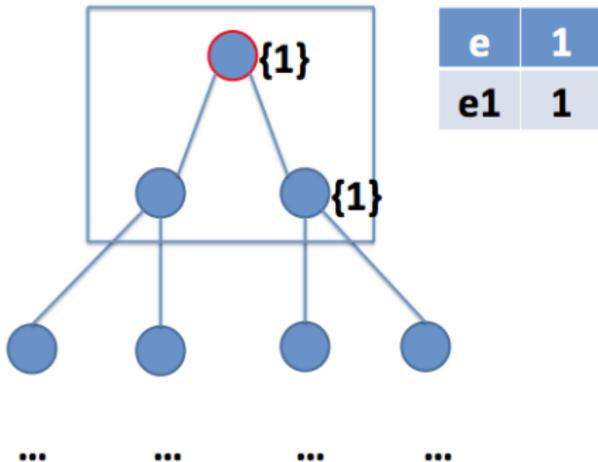


# LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

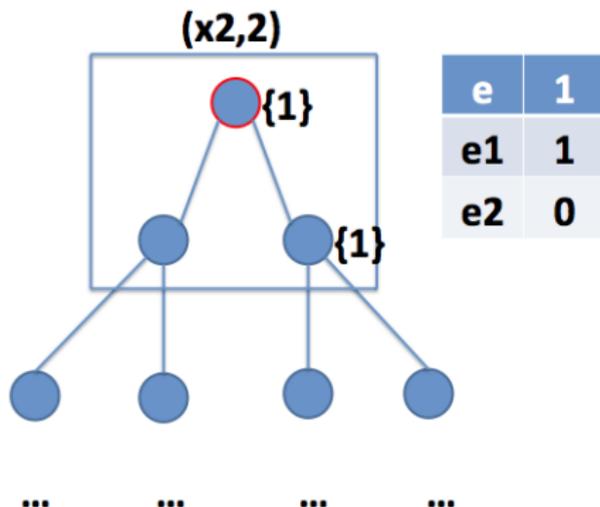


# LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

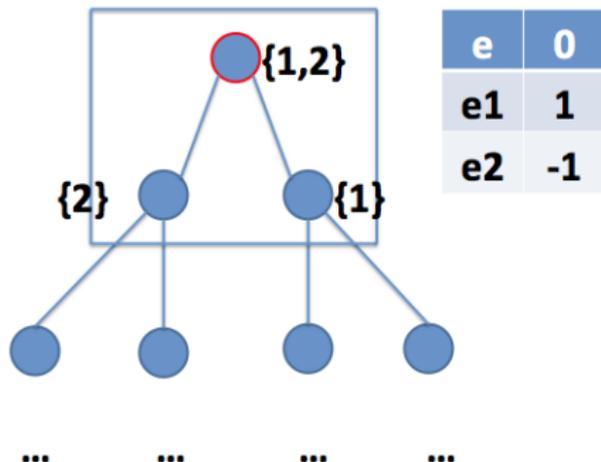


## LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

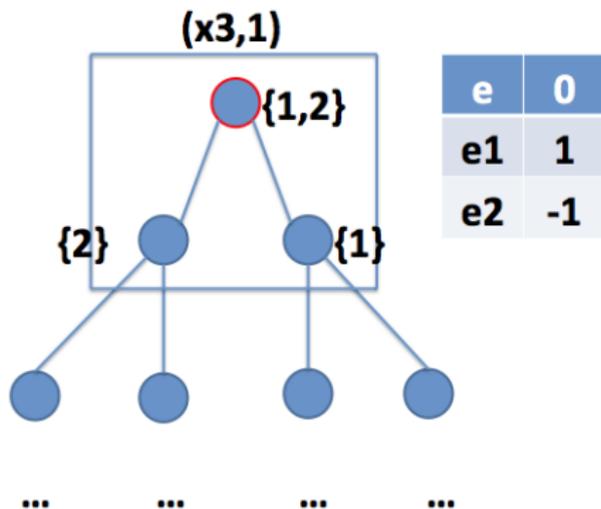


## LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

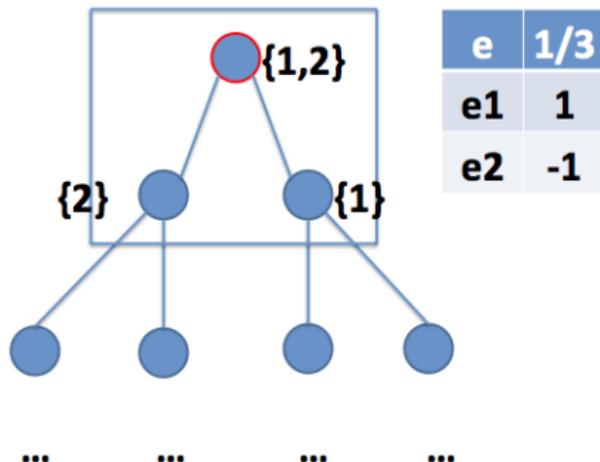


# LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

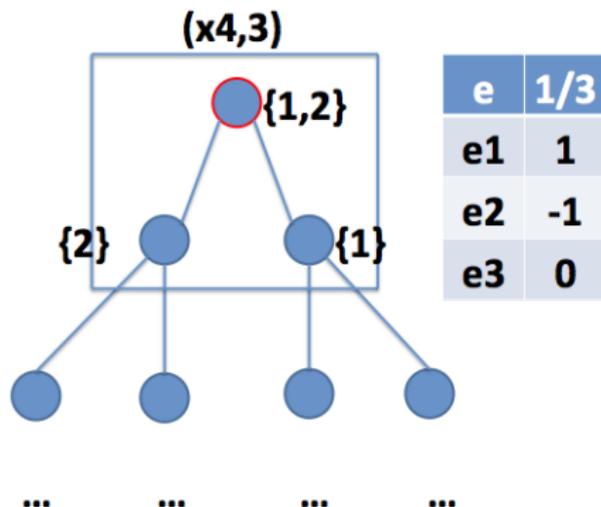


## LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

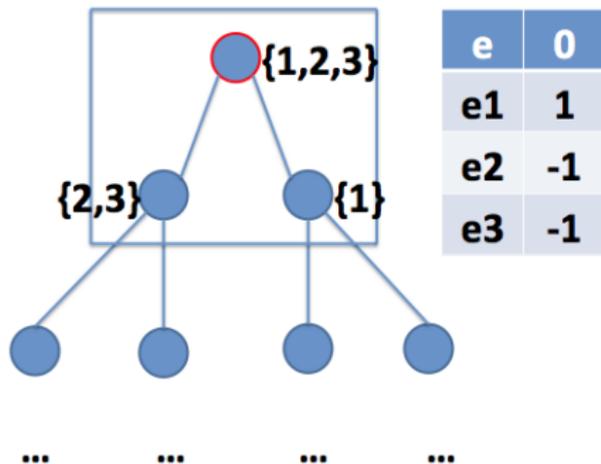


# LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

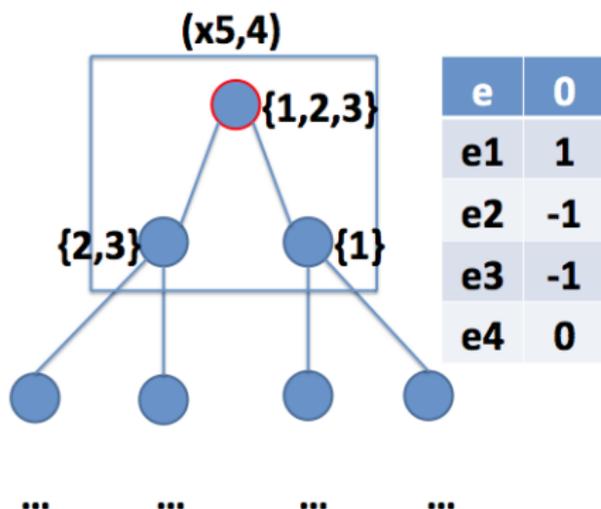


## LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

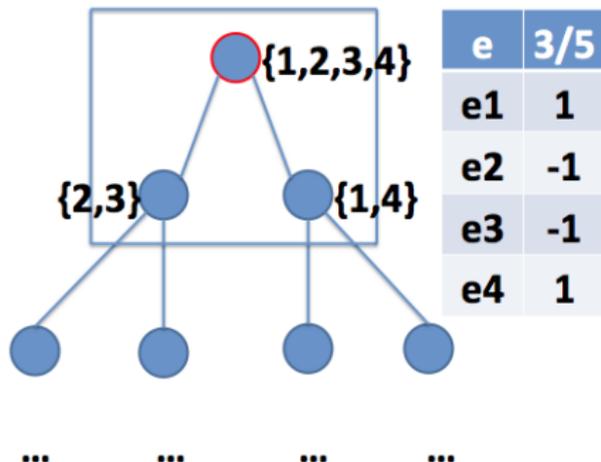


## LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$



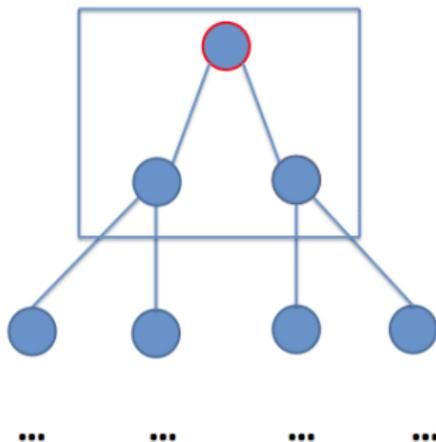
# LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

Apply recursively to  
construct a tree structure.



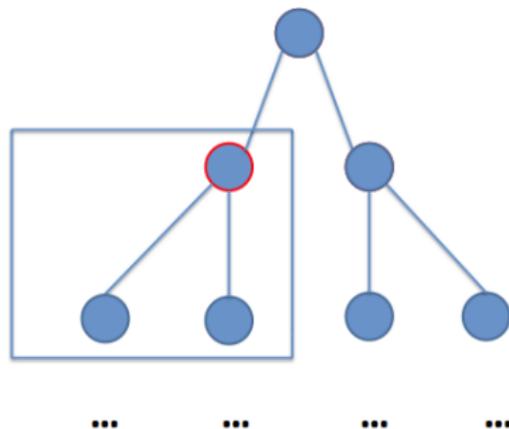
# LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

Apply recursively to  
construct a tree structure.



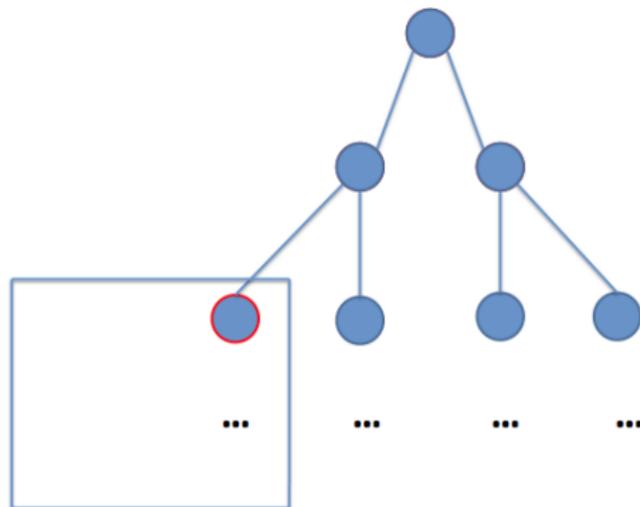
# LOMtree algorithm

Let  $e = 0$  and for all  $y$ ,  $e_y = 0$ ,  $n_y = 0$

For each example  $(x, y)$

- if  $e_y < e$  then  $b = -1$  else  $b = 1$
- Update  $w$  using  $(x, b)$
- $n_y \leftarrow n_y + 1$
- $e_y \leftarrow \frac{(n_y - 1)e_y}{n_y} + \frac{w \cdot x}{n_y}$
- $e \leftarrow \frac{(n - 1)e}{n} + \frac{w \cdot x}{n}$

Apply recursively to  
construct a tree structure.



# Experiments

Table : Training time on selected problems.

	Isolet	Sector	Aloi
LOMtree	<b>16.27s</b>	<b>12.77s</b>	<b>51.86s</b>
OAA	19.58s	18.37s	11m2.43s

Table : Per-example test time on all problems.

	Isolet	Sector	Aloi	ImNet	ODP
LOMtree	<b>0.14ms</b>	<b>0.13ms</b>	<b>0.06ms</b>	<b>0.52ms</b>	<b>0.26ms</b>
OAA	0.16 ms	0.24ms	0.33ms	0.21s	1.05s

Table : Test error (%) and confidence interval on all problems.

	LOMtree	Rtree	Filter tree
Isolet ( <b>26</b> )	<b>6.36</b> $\pm$ 1.71	16.92 $\pm$ 2.63	15.10 $\pm$ 2.51
Sector ( <b>105</b> )	16.19 $\pm$ 2.33	<b>15.77</b> $\pm$ 2.30	17.70 $\pm$ 2.41
Aloi ( <b>1000</b> )	<b>16.50</b> $\pm$ 0.70	83.74 $\pm$ 0.70	80.50 $\pm$ 0.75
ImNet ( <b>22K</b> )	<b>90.17</b> $\pm$ 0.05	96.99 $\pm$ 0.03	92.12 $\pm$ 0.04
ODP ( <b>105K</b> )	<b>93.46</b> $\pm$ 0.12	93.85 $\pm$ 0.12	93.76 $\pm$ 0.12



# How to understand non-convex optimization?

# Why non-convex optimization and deep learning?

State-of-the art results on number of problems:

- image recognition [KSH12, CMGS10]
- speech recognition [HDYDMJSVNSK12, GMH13]
- natural language processing [WCA14]
- video recognition [KTSLSF-F14, SZ14]



# Challenge

**Goal:** Understanding loss function in deep learning.

**Recent related works:** Choromanska et al., 2015, Goodfellow et al., 2015, Dauphin et al., 2014, Saxe et al., 2014.

# Challenge

**Goal:** Understanding loss function in deep learning.

**Recent related works:** Choromanska et al., 2015, Goodfellow et al., 2015, Dauphin et al., 2014, Saxe et al., 2014.

## Questions:

- Why the result of multiple experiments with multilayer networks consistently give very similar performance despite the presence of many local minima?

# Challenge

**Goal:** Understanding loss function in deep learning.

**Recent related works:** Choromanska et al., 2015, Goodfellow et al., 2015, Dauphin et al., 2014, Saxe et al., 2014.

## Questions:

- Why the result of multiple experiments with multilayer networks consistently give very similar performance despite the presence of many local minima?
- What is the role of saddle points in the optimization problem?

# Challenge

**Goal:** Understanding loss function in deep learning.

**Recent related works:** Choromanska et al., 2015, Goodfellow et al., 2015, Dauphin et al., 2014, Saxe et al., 2014.

## Questions:

- Why the result of multiple experiments with multilayer networks consistently give very similar performance despite the presence of many local minima?
- What is the role of saddle points in the optimization problem?
- Is the surface of the loss function of multilayer networks structured?

# Multilayer network and spin-glass model

**Can we use the spin-glass theory to explain the optimization paradigm with large multilayer networks?**

# Multilayer network and spin-glass model

Can we use the spin-glass theory to explain the optimization paradigm with large multilayer networks?

What assumptions need to be made?



# Loss function in deep learning and assumptions

Consider hinge loss

$$L(\mathbf{w}) = \max(0, 1 - Y_t Y),$$

where  $Y_t$  corresponds to the true data labeling ( $1/-1$ ), and  $\mathbf{w}$  denotes all network weights.

- max operator is often modeled as Bernoulli r.v. (0/1). Denote it as  $M$  and its expectation as  $\rho'$ . Therefore

$$L(\mathbf{w}) = M(1 - Y_t Y) = M + \frac{1}{\Lambda^{(H-1)/2}} \sum_{i=1}^{\Psi} Z_i l_i \prod_{k=1}^H w_i^{(k)}, \quad (1)$$

where  $Z_i = -Y_t X_i$ , and  $l_i = M A_i$  is a Bernoulli r.v. (0/1).

- assume  $l_1, l_2, \dots, l_{\Psi}$  are identically distributed (A1p)
- assume each  $X_i$  is a standard Gaussian r.v. (A2p)





# Loss function in deep learning and assumptions

- assume the independence of  $Z_{i_1, i_2, \dots, i_H}$  and  $I_{i_1, i_2, \dots, i_H}$  (A5u)

$$\mathbb{E}_{M, I_1, I_2, \dots, I_\Psi} [L(\mathbf{w})] = \rho' + \rho \underbrace{\frac{1}{\Lambda^{(H-1)/2}} \sum_{i_1, i_2, \dots, i_H=1}^{\Lambda} Z_{i_1, i_2, \dots, i_H} w_{i_1} w_{i_2} \dots w_{i_H}}_{\mathcal{L}_{\Lambda, H}}$$

- assume that  $Z$ 's are independent (A6u)
- impose spherical constraint (A7p)

$$\frac{1}{\Lambda} \sum_{i=1}^{\Lambda} w_i^2 = 1.$$

**We obtain the Hamiltonian of the spherical spin-glass model!!!**

Question: What happens when  $\Lambda \rightarrow \infty$ ?

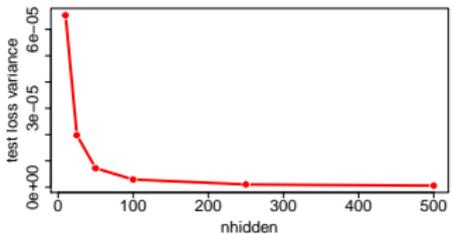
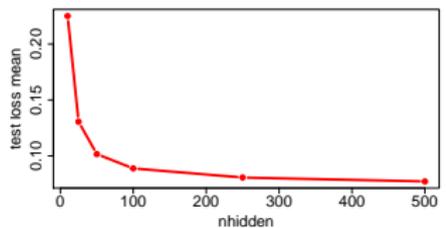
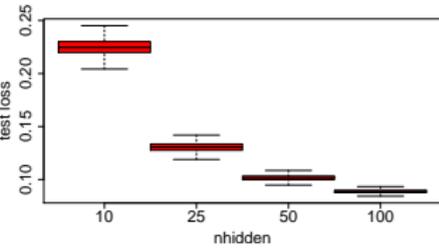
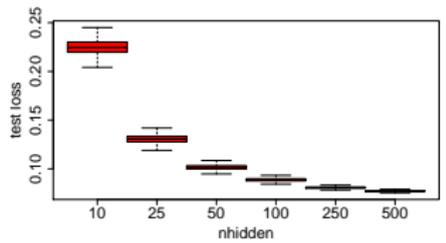
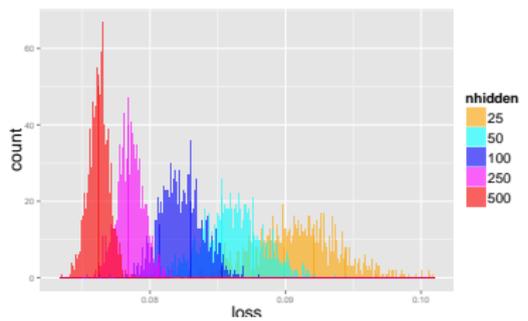
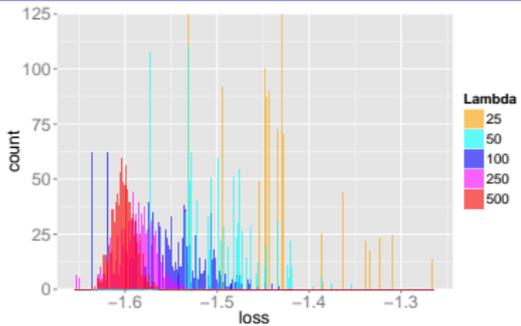






Deep networks versus spherical spin-glass models

# Comparison of deep network and spherical spin-glass model







# Spherical spin-glass versus deep network

## Conjecture (Deep learning)

*For large-size networks, most local minima are equivalent and yield similar performance on a test set.*

## Spherical spin-glass

Critical points form an ordered structure such that there exists an *energy barrier*  $\Lambda E_{-\infty}$  (a certain value of the Hamiltonian) below which with overwhelming probability one can find only low-index critical points, most of which are concentrated close to the barrier.

# Spherical spin-glass versus deep network

## Conjecture (Deep learning)

*The probability of finding a “bad” (high value) local minimum is non-zero for small-size networks and decreases quickly with network size.*

## Spherical spin-glass

Low-index critical points are 'geometrically' lying closer to the ground state than high-index critical points.



# Spherical spin-glass versus deep network

## Conjecture (Deep learning)

*Struggling to find the global minimum on the training set (as opposed to one of the many good local ones) is not useful in practice and may lead to overfitting.*

$n_1$	25	50	100	250	500
$\rho$	0.7616	0.6861	0.5983	0.5302	0.4081

**Table :** Pearson correlation between training and test loss for different numbers of hidden units of a network with one hidden layer. *MNIST* dataset.

## Spherical spin-glass

Recovering the ground state, i.e. global minimum, takes exponentially long time.

# Take-home message

- For large-size networks, most local minima are equivalent and yield similar performance on a test set.
- The probability of finding a “bad” (high value) local minimum is non-zero for small-size networks and decreases quickly with network size.
- Struggling to find the global minimum on the training set (as opposed to one of the many good local ones) is not useful in practice and may lead to overfitting.

# Open problem

**Can we establish a stronger connection between the loss function of the deep model and the spherical spin-glass model by dropping the unrealistic assumptions?**

# Convexity and non-convexity: challenges

## Convex and non-convex world

- Building solvers, i.e. bound majorization
  - New and tight **quadratic bound on the partition function**
  - **Linear convergence** of the batch/semi-stochastic variants
  - **Competitive/better** than **state-of-the-art** methods
  - Admits multiple extensions
- Designing problem-specific, i.e. multi-classification, objectives
  - **Logarithmic training and testing time**
  - **Reduction** from multi-class to binary classification
  - New **splitting criterion** with desirable properties
    - allows gradient descent style optimization
    - makes decision trees applicable to multi-class classification

## Non-convex world

- Understanding why non-convex approaches work
  - **Deep learning: state-of-the-art** in numerous problems
  - Possible connection between **spin-glass theory** and **deep learning**
  - Landscape is **highly non-convex** but most likely **structured**

# Acknowledgments

**Courant Institute of Mathematical Sciences:** Yann LeCun,  
G erard Ben Arous

**Columbia University:** Tony Jebara, Shih-Fu Chang

**George Washington University:** Claire Monteleoni

**NYU Polytechnic School of Engineering:** Mariusz Bojarski

**Microsoft Research:** John Langford, Alekh Agarwal

**Google Research:** Krzysztof Choromanski, Dimitri Kanevsky

**IBM T. J. Watson Research Center:** Aleksandr Aravkin

**ATT Shannon Research Laboratories:** Phyllis Weiss, Alice Chen

**LEAR team of Inria:** Zaid Harchaoui

**PostDocs:** Pablo Sprechmann

**PhD students:** Michael Mathieu, Mikael Bruce Henaff, Sixin  
Zhang, Ross Goroshin, Rahul Krishnan, Wojciech Zaremba,  
Hyungtae Kim