

CS 744: BISMARCK

Shivaram Venkataraman Fall 2019

ADMINISTRIVIA

- Assignment 2 out!
- Project groups extension -> dill oct 7th
- OH / Setup meeting by email

fill OCt 7th Oct 17th Introduction

COURSE PROJECT PROPOSAL

Propose topic, group (2 sentences) – Oct 7 Project Proposal (2 pages) – Oct 17

Introduction Related Work Timeline (with eval plan)

WRITING AN INTRODUCTION

- I-2 paras: what is the problem you are solving why is it important (need citations)
 I-2 paras: How other people solve and why they fall short

 La Target and ince : Familier work
 La Target and ince : Familier work
 La Target and ince : Familier work
- I para: Anticipated results or what experiments you will use

WRITING RELATED WORK

Group related work into two/three buckets (I-2 para per bucket)

Explain what the papers / projects do Why are they different / insufficient ML Francosiles:

Scheduling Franceworks



MACHINE LEARNING

Classification

IM GENET

GMail - Span filter

N

Recommendation

Maller

amazon

NETFLIX



CONVEX OPTIMIZATION

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^N f(w, z_i) + P(w)$$

What is convex ?
Linear Regression, Linear SVM
Kernel SVMs, Logistic Regression,

What is not convex ? Graph mining, Deep Learning

What is

GRADIENT DESCENT

$$w^{(k+1)} = w^{(k)} - \alpha_k \nabla f(w^{(k)}) \xrightarrow{\text{stable}} t^k$$

Initialize $w \xrightarrow{\text{random}} t^{\text{suitialization}} \xrightarrow{\text{the prime}} t^{\text{the iteration}}$
For many iterations:
Compute Gradient $\xrightarrow{\text{How}} t^{\text{range}} dt_k$ points?
Update model
End
IGD: All of your training deta (n)
SGD: Stochastic Gradied (B=128 or 512)

INCREMENTAL GRADIENT DESCENT

$$w^{(k+1)} = w^{(k)} - \alpha_k \nabla f_{\eta(k)}(w^{(k)})$$

Initialize w

For many iterations:

Pick one point

Compute Gradient

Update model

End

Analytics Task	Objective
Logistic Regression (LR)	$\sum_{i} \log(1 + exp(-y_i w^T x_i)) + \mu \ \vec{w}\ _1$
Classification (SVM)	$\int_{i} (1 - y_i w^T x_i)_+ + \mu \ \vec{w} \ _1$
Recommendation (LMF)	$\int_{(i,j)\in\Omega} (L_i^T R_j - M_{ij})^2 + \mu \ L, R\ _F^2$
Labeling (CRF) $[48]$	$\sum_{k} \left[\sum_{j} w_{j} F_{j}(y_{k}, x_{k}) - \log Z(x_{k}) \right]$
Kalman Filters	$\int_{t=1}^{T} Cw_t - f(y_t) _2^2 + w_t - Aw_{t-1} _2^2$
Portfolio Optimization	$p^T w + w^T \Sigma w$ s.t. $w \in \Delta$

BISMARCK ARCHITECTURE



BISMARCK: USER DEFINED AGGREGATE

Three steps:



BISMARCK: LOGISTIC REGRESSION

LR_Transition(ModelCoef *w, Example e) { ...
wx = Dot_Product(w, e.x);
sig = Sigmoid(-wx * e.y);
c = stepsize * e.y * sig;
Scale_And_Add(w, e.x, c); ... }



n(t_shugg + t_epoch) < n,t_epoch + t_shugg

Random sampling

- Sample without replacement
- Shuffle the data after each epoch

Shuffle once

- Avoids pathological ordering
- Much cheaper





RESERVOIR SAMPLING

On the k^{th} additional item s = random in [0, m + k)

> if s < mPut in slot s else Drop the item





PARALLEL GRADIENTS

Vf (wk

IGM

Shared Memory:

- Compute gradients in parallel
- Average their updates
- Or update in parallel
 - Locks?

- AIG : Afonic Increment Instructions 25 fast, implemented in hardware · Lock-free: Iteration times very fast, conflicts?

DISCUSSION

https://forms.gle/nFNEi2NZMNhZio1f7

How would an implementation of GD look in Spark? Try to sketch an implementation. What would be similar / different to Bismarck?

Shared memory Lyno broadcast LAIG2 & Luck free What are some ML scenarios where Bismarck architecture might prove to be limited?



NEXT STEPS

- Next class: Parameter Server
- Assignment 2 out!
- **Project Proposal**
 - Groups by Oct 7
 - 2 pager by Oct 17