Large-scale analyses in functional brain Imaging

Bertrand Thirion bertrand.thirion@inria.fr



12/02/2016

Outline

- The importance of using large(r) datasets in medical imaging
- Combining modalities to increase power
- Identify ensuing computational bottlenecks and algorithmic solutions



12/02/2016

Imaging in the low power regime

Low power → unreliable findings

Low positive predictive value: low likelihood that a statistically significant finding is actually true





Exaggerated estimate of the magnitude of discovered effects "*winner's curse*"

Medical imaging in the big data era?

- Neuroscientists reluctant to using big data : this is a burden
- Common belief : you should only focus on large effects
 → small datasets are good enough
- But...

Nature Reviews Neuroscience | AOP, published online 10 April 2013; doi:10.1038/nrn3475



Power failure: why small sample size undermines the reliability of neuroscience

Katherine S. Button^{1,2}, John P. A. Ioannidis³, Claire Mokrysz¹, Brian A. Nosek⁴, Jonathan Flint⁵, Emma S. J. Robinson⁶ and Marcus R. Munafò¹

Example

OASIS dataset, effet of age on gray matter density



n=10 n=20 n=30 n=50 n=70 n=100

Reproducibility by bootstrap: 7% 19% 32%

53%

66%

75%



Multivariate analysis



12/02/2016

Multivariate_analysis



Large-scale connectome analysis for disease prediction



871 subjects, 300-400GB of data

12/02/2016

Large scale connectome analysis for disease prediction

- Classification of ASD patients
 - Public ABIDE dataset [Di Martino et al. Mol. Psych. 2014]
 - Large heterogeneity: multi-centric study
 - Functional connectivity: weakly sensitive marker
- 68% correct classification (65% when generalizing across sites)

Large scale connectome analysis for disease prediction

What are the most important factors of the pipeline ?





Response: the choice of the atlas.

[Abraham et al. Subm to NIMG]

12/02/2016

Large scale connectome analysis for disease prediction

More subjects \rightarrow higher accuracy

Impact of the number of regions

Asymptote not yet reached at n=871





Leveraging multi-modal datasets



Combining images from Different modalities



12/02/2016

Trans-modal priors improve discriminative models

Diagnosis of Alzheimer disease (Alzheimer vs MCI) using resting-state fMRI and a PET prior



12/02/2016

Trans-modal priors improve discriminative models



[Rahim et al. MICCAI 2015]

12/02/2016

Using rest data to better model task data

- A wealth of "restingstate" data available
- Little cognitive data with proper annotations
- Idea: use unlabelled data to improve the learning of discriminative models



[Bzdok et al. NIPS 2015]

12/02/2016

Using rest data to better model task data

- A wealth of "restingstate" data available
- Little cognitive data with proper annotations
- Idea: use unlabelled data to improve the learning of discriminative models



[Bzdok et al. NIPS 2015]

12/02/2016

Using rest data to better model task data data SS-F LogReg

 The discriminative patterns for many task data is much cleaner thanks to the composite objective



[Bzdok et al. NIPS 2015]

12/02/2016

Technical aspects of big data analysis



Big data in medical imaging ?

HCP mailing list, Jan 19th, 2015

"Has anyone on the run group-wise analysis on the HCP resting state data, and if so what tools did you use?

I am having memory issues when running more than 10 subjects and I was wondering if anyone has a way of getting around the large memory requirements when concatenating in time."

Diagnosis and strategy

- Modern, multiple, datasets: 100 GB-TBs
 - Do not fit in memory. Need online learning
 - Memory access is the main bottleneck
 - Large p, large n
 - Better file formats would help, but compatibility with existing tool stacks.
- Agile approach: readable Python code, runs on your laptop
 - Online learning, SGD algorithms
 - compression

Caching & parallel comp.: joblib

Table Of Contents

Introduction

- Vision
- Main features
 User manual

Module reference

Next topic

Mby	ioblib	· pro	ioct	aoal	C
	unuu	. 010	ECL	uuai	5

Quick search

Enter search terms or a module, class or function name.

Go

Mailing list

Joblib: running Python functions as pipeline jobs

Introduction

Joblib is a set of tools to provide **lightweight pipelining in Python**. In particular, joblib offers:

- 1. transparent disk-caching of the output values and lazy re-evaluation (memoize pattern)
- 2. easy simple parallel computing
- 3. logging and tracing of the execution

Joblib is optimized to be **fast** and **robust** in particular on large data and has specific optimizations for *numpy* arrays. It is **BSD-licensed**.

Fast image compression

• Cost of Logistic regression $\propto pn^2$

+ memory consumption \propto pn

- Solution: reduce p by compression
 - Random projections [Johnson, Lindenstrauss, AMS 1984]
 - Explicit control on precision
 - Clustering: informed by data [Thirion et al. Stamlins 2015]



12/02/2016

Efficient discriminative models

Example: HCP dataset, multi-class prediction (k=18), clustering + fit time



- Time savings are more than linear wrt data volume
- Data compression also benefits to accuracy

(Sparse) Principal Components Analysis

• Identify structure in data with *dictionary learning / sparse principal components analysis*



- Memory cost: n.p → scales poorly with large n (large number of subjects): 2TB of data on HCP
 - Run this on computers with 8GB of RAM
 - Possible with stochastic (online) methods !

[Mensch et al. ISBI 2016, ICML 2016]



12/02/2016

Summary

- Large(r) sample sizes to make imagine findings more reliable
 - Even though less homogeneous data
- Combining multi-modal data to improve statistical analysis
- Improve software to support large data analysis
 - Reduce memory requirements: online methods
 - Leverage parallel computers



Acknowledgements

Parietal

- G. Varoquaux,
- P. Ciuciu
- A. Gramfort,
- O. Grisel,
- L. Estève,
- Y. Schwartz,
- F. Pedregosa,
- E. Dohmatob,
- A. H. Idrobo,
- V. Fritsch,
- M. Eickenberg,
- R. Bricquet
- A. Abraham
- D. Bzdok
- A. Frau

12/02/2016

N. Chauffert

Other collaborators (thanks for the data)

S. Dehaene E. Eger, R. Poldrack, K. Jimura, J. Haxby C. F. Gorgolevski JB. Poline

Brainpedia project Brainpedia project Microsoft Research - Inria JOINT CENTRE Human Brain Project