

Astrostatistics and Astroinformatics

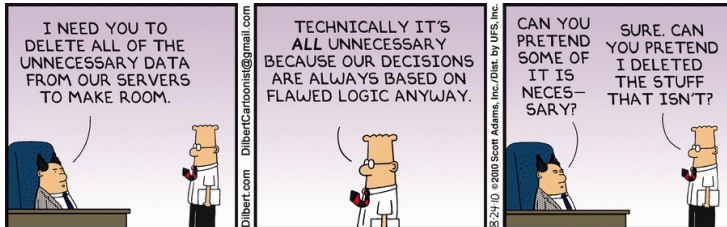
Big-Data in Astronomy and Astrophysics

Jason McEwen

www.jasonmcewen.org

[@jasonmcewen](https://twitter.com/jasonmcewen)

*Mullard Space Science Laboratory (MSSL)
University College London (UCL)*



UK Dark Energy Strategy 2020
Royal Astronomical Society, London, January 2016

Outline

- 1 Big-data in astronomy and astrophysics
- 2 Illustrative analyses
 - Planck
 - Euclid
 - LSST
 - SKA
- 3 Concluding remarks

What is big-data?

A. Gandomi, M. Haider / International Journal of Information Management 35 (2015) 137–144

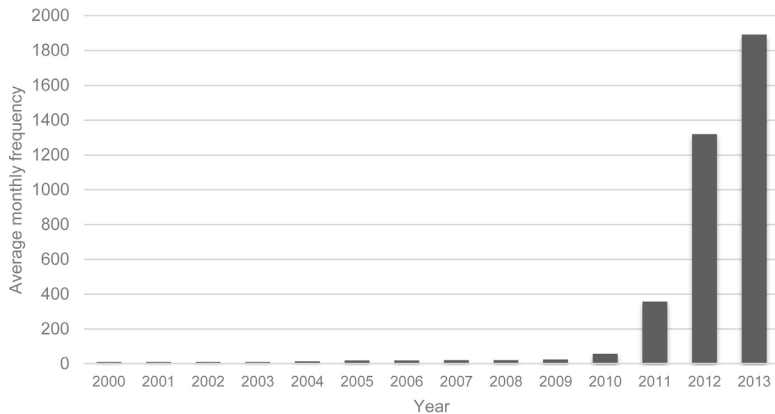


Fig. 1. Frequency distribution of documents containing the term “big data” in ProQuest Research Library.

What is big-data?

The nVs (originally 3Vs, then 6Vs, then 10Vs, ...):

- 1 **Volume**: many bytes (e.g. typically peta, exabytes)
- 2 **Variety**: structural heterogeneity (e.g. sub-populations, variety of sources)
- 3 **Velocity**: rate of generation and analysis
- 4 **Veracity**: unreliability in sources
- 5 **Variability**: variation in data flow rate
- 6 **Value**: low value density
- 7 ...

Typically (but not exclusively) characterised by:

- High-dimensional datum (**wide**)
- Massive number of datum (**deep**)

What is big-data?

The nVs (originally 3Vs, then 6Vs, then 10Vs, ...):

- 1 **Volume**: many bytes (e.g. typically peta, exabytes)
- 2 **Variety**: structural heterogeneity (e.g. sub-populations, variety of sources)
- 3 **Velocity**: rate of generation and analysis
- 4 **Veracity**: unreliability in sources
- 5 **Variability**: variation in data flow rate
- 6 **Value**: low value density
- 7 ...

Typically (but not exclusively) characterised by:

- High-dimensional datum (**wide**)
- Massive number of datum (**deep**)

What is big-data?

The nVs (originally 3Vs, then 6Vs, then 10Vs, ...):

- 1 **Volume**: many bytes (e.g. typically peta, exabytes)
- 2 **Variety**: structural heterogeneity (e.g. sub-populations, variety of sources)
- 3 **Velocity**: rate of generation and analysis
- 4 **Veracity**: unreliability in sources
- 5 **Variability**: variation in data flow rate
- 6 **Value**: low value density
- 7 ...

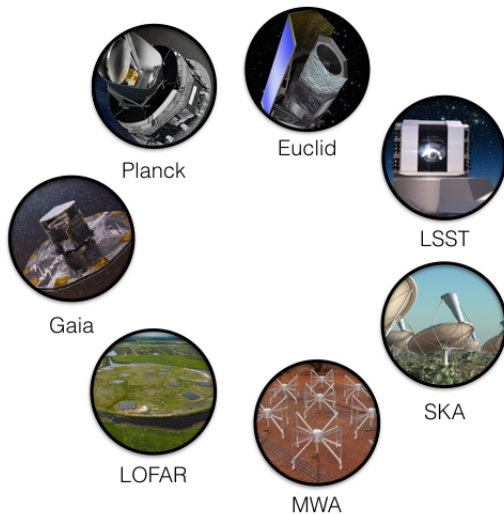
Typically (but not exclusively) characterised by:

- High-dimensional datum (**wide**)
- Massive number of datum (**deep**)

What is big-data in astronomy and astrophysics?

- Big **machines** (e.g. physical hardware, experiments)
- Big **theory**
- Big **simulations**
- Big **parameter space**
- Big **algorithms**
- Big **collaborations**
- Big **engagement** (e.g. outreach, industry)

What is big-data in astronomy and astrophysics?



Wide and deep observations (in addition to wide and deep data)

Challenges of big-data

A. Gandomi, M. Haider / *International Journal of Information Management* 35 (2015) 137–144

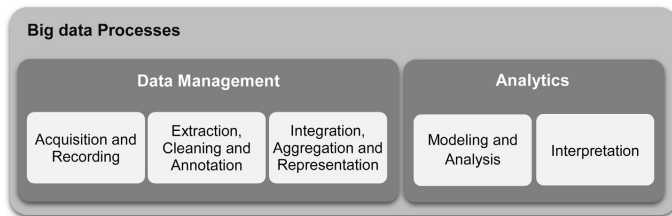


Fig. 3. Processes for extracting insights from big data.

Computational challenges:

- Data **too big** (to hold in memory)
- Access and analysis **too slow** (unfeasible)
- **Too much power/energy** required

Challenges of big-data

A. Gandomi, M. Haider / International Journal of Information Management 35 (2015) 137–144

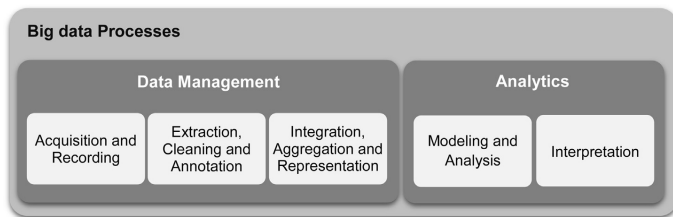


Fig. 3. Processes for extracting insights from big data.

Computational challenges:

- Data **too big** (to hold in memory)
- Access and analysis **too slow** (unfeasible)
- **Too much power/energy** required

Challenges of big-data

Analysis challenges (Fan et al. 2014):

- 1 **Heterogeneity**, e.g. sub-populations, different data sources, tension between data
- 2 Error accumulation, e.g. high-dimensional parameter spaces, bias
- 3 Spurious correlations, e.g. correlation vs causation, data dredging
- 4 Incident endogeneity, e.g. chance correlation between signal of interest and error

Challenges of big-data

Analysis challenges (Fan et al. 2014):

- 1 **Heterogeneity**, e.g. sub-populations, different data sources, tension between data
- 2 **Error accumulation**, e.g. high-dimensional parameter spaces, bias
- 3 **Spurious correlations**, e.g. correlation vs causation, data dredging
- 4 **Incident endogeneity**, e.g. chance correlation between signal of interest and error

Challenges of big-data

Analysis challenges (Fan et al. 2014):

- 1 **Heterogeneity**, e.g. sub-populations, different data sources, tension between data
- 2 **Error accumulation**, e.g. high-dimensional parameter spaces, bias
- 3 **Spurious correlations**, e.g. correlation vs causation, data dredging
- 4 **Incident endogeneity**, e.g. chance correlation between signal of interest and error

Challenges of big-data

Analysis challenges (Fan et al. 2014):

- 1 **Heterogeneity**, e.g. sub-populations, different data sources, tension between data
- 2 **Error accumulation**, e.g. high-dimensional parameter spaces, bias
- 3 **Spurious correlations**, e.g. correlation vs causation, data dredging
- 4 **Incident endogeneity**, e.g. chance correlation between signal of interest and error

Analysing big-data

Generic approaches to analysing big-data (Wang et al. 2015):

- Subsample
- Divide-and-conquer
- Stream processing

Additional approaches in astronomy and astrophysics:

- Exploit structure (geometry, symmetry, physics)
- Modelling:
 - Model-based consolidatory science
 - Model-agnostic exploratory science
- Approximation
- ...

Analysing big-data

Generic approaches to analysing big-data (Wang et al. 2015):

- Subsample
- Divide-and-conquer
- Stream processing

Additional approaches in astronomy and astrophysics:

- Exploit structure (geometry, symmetry, physics)
- Modelling:
 - Model-based consolidatory science
 - Model-agnostic exploratory science
- Approximation
- ...

Analysing big-data

Examples of specific methods:

- Bayesian analysis
- MCMC sampling
- Hierarchical probabilistic (Bayesian) models
- Variable selection
- Experimental design
- Machine learning
- Optimisation
- Wavelets
- Sparsity
- Compressed sensing
- ...

⇒ **Astrostatistics and Astroinformatics**

Outline

1 Big-data in astronomy and astrophysics

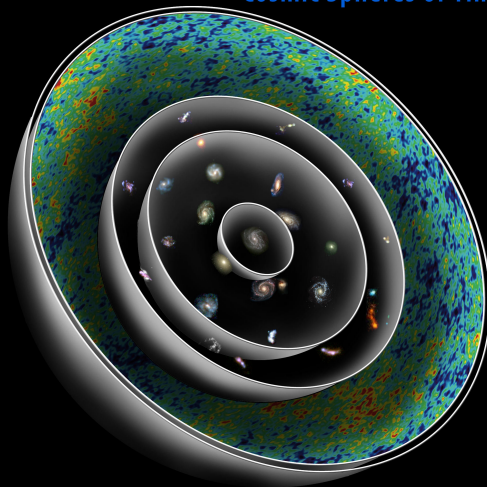
2 Illustrative analyses

- Planck
- Euclid
- LSST
- SKA

3 Concluding remarks

Observations made on the celestial sphere

Cosmic Spheres of Time



© 2006 Abrams and Primack, Inc.

Scale-discretised wavelets on the sphere

Transforms

- **Spin scale-discretised wavelet transform** is given by the projection onto each wavelet (Wiaux, McEwen *et al.* 2008, McEwen *et al.* 2013, McEwen *et al.* 2015):

$$W^s \Psi^j(\rho) = \underbrace{\langle sf, \mathcal{R}_{\rho s} \Psi^j \rangle}_{\text{projection}} = \int_{\mathbb{S}^2} d\Omega(\omega) sf(\omega) (\mathcal{R}_{\rho s} \Psi^j)^*(\omega) .$$

- Original function may be recovered exactly in practice from wavelet coefficients:

$$sf(\omega) = \underbrace{\sum_{j=0}^J}_{\text{finite sum}} \underbrace{\int_{\text{SO}(3)} d\rho W^s \Psi^j(\rho) (\mathcal{R}_{\rho s} \Psi^j)(\omega)}_{\text{wavelet contribution}} .$$

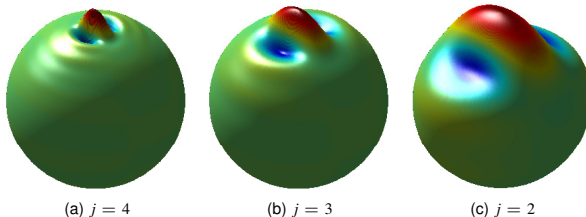


Figure: Scale-discretised wavelets on the sphere.

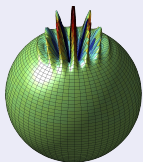
Scale-discretised wavelets on the sphere

Fast algorithms and codes

- **Fast algorithms essential** (McEwen, Leistedt *et al.* 2015, Leistedt, McEwen *et al.* 2013, McEwen *et al.* 2013, Leistedt McEwen *et al.* 2007, Wiaux, McEwen & Vielva 2007, Wiaux *et al.* 2005, Wandelt & Gorski 2001, Risbo 1996)

FastCSWT code

<http://www.fastcswt.org>



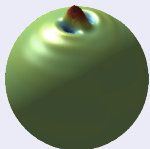
Fast directional continuous spherical wavelet transform algorithms

McEwen *et al.* (2007)

- Fortran
- Supports directional and steerable wavelets

S2DW code

<http://www.s2dw.org>



Exact reconstruction with directional wavelets on the sphere

Wiaux, McEwen, Vanderghyest, Blanc (2008)

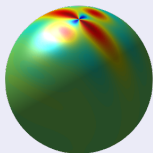
- Fortran
- Parallelised
- Supports directional and steerable wavelets
- Supports inversion

Scale-discretised wavelets on the sphere

Fast algorithms and codes

S2LET code

<http://www.s2let.org>



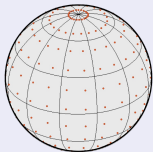
S2LET: Fast wavelet analysis on the sphere

McEwen, Leistedt, Büttner, Peiris & Wiaux (2015), Leistedt, McEwen, *et al.* (2012)

- C, Matlab, IDL, Python
- Supports directional and steerable wavelets, ridgelets and curvelets
- Supports inversion
- Supports spin
- Faster algorithms

SO3 code

<http://www.sothree.org>



SO3: Fast Wigner transforms on the rotation group

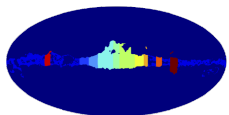
McEwen, Büttner, Leistedt, Peiris & Wiaux (2015)

- C, Matlab, Python
- Fast and exact Fourier transforms on the rotation group $SO(3)$

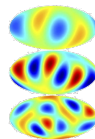
Planck component separation

SILC

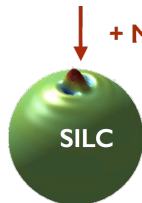
- **SILC**: Blind Planck component separation via Scale-discretised, directional wavelet Internal Linear Combination (Rogers, Peiris, Leistedt, McEwen & Pontzen 2016)



Spatial
WMAP Collab. (2003)

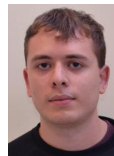


Harmonic
Tegmark et al. (2003)



+ Morphological

NILC: Delabrouille et al. (2009)
SILC: Rogers et al. (2016)
Wang et al. (2015)

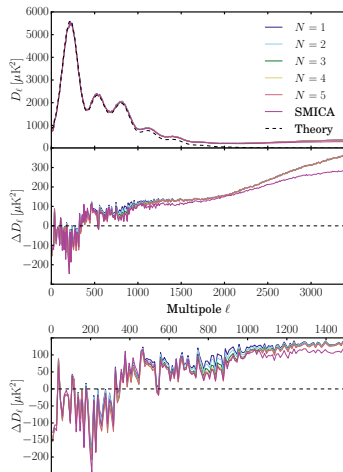
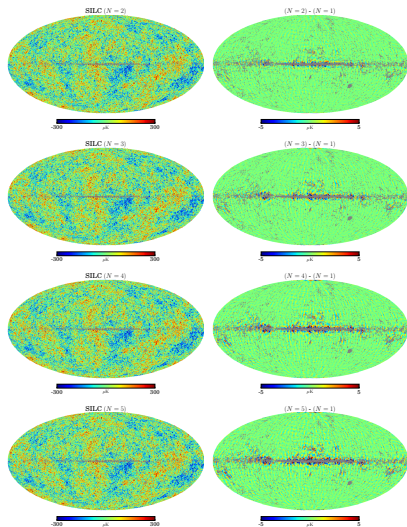


Keir Rogers

Planck component separation

SILC

- SILC (R1) maps available for download: <http://www.silc-cmb.org>



E/B separation

Exploiting scale-discretised wavelets

- E/B separation with spin directional wavelets for CMB polarisation and cosmic shear (Leistedt, McEwen, Büttner & Peiris, in prep.)



Boris Leistedt

Algorithm to recover E/B signals using scale-discretised wavelets

- 1 Compute spin wavelet transform of $Q + iU$:

$$(Q + iU)(\omega) \xrightarrow[\text{S2LET}]{\text{Spin wavelet transform}} W_{Q+iU}^{2\Psi^j}(\rho)$$

- 2 Account for mask in **harmonic and spatial** domains simultaneously:

$$W_{Q+iU}^{2\Psi^j}(\rho) \xrightarrow{\text{Mitigate mask}} \widehat{W}_{Q+iU}^{2\Psi^j}(\rho)$$

- 3 Construct E/B maps:

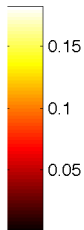
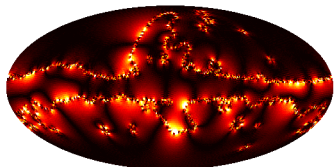
$$(a) \quad W_{\tilde{E}}^0(\rho) = -\text{Re} \left[\widehat{W}_{Q+iU}^{2\Psi^j}(\rho) \right] \xrightarrow[\text{S2LET}]{\text{Inverse scalar wavelet transform}} \tilde{E}(\omega)$$

$$(b) \quad W_{\tilde{B}}^0(\rho) = -\text{Im} \left[\widehat{W}_{Q+iU}^{2\Psi^j}(\rho) \right] \xrightarrow[\text{S2LET}]{\text{Inverse scalar wavelet transform}} \tilde{B}(\omega)$$

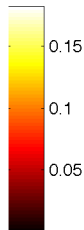
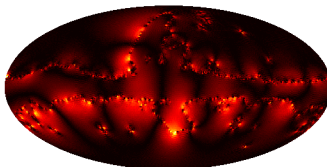
E/B separation

Preliminary results

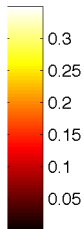
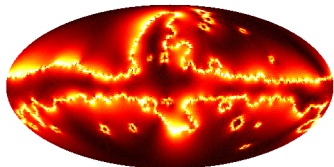
Mean of B maps
reconstructed using harmonics



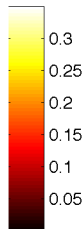
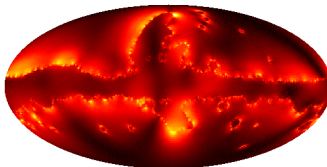
Mean of B maps
reconstructed using wavelets



Std dev of B maps
reconstructed using harmonics



Std dev of B maps
reconstructed using wavelets



Outline

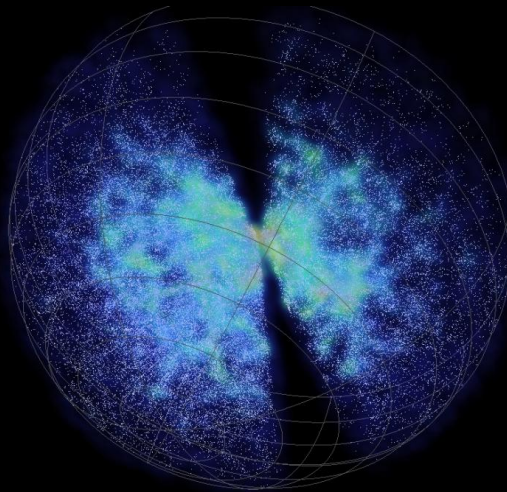
1 Big-data in astronomy and astrophysics

2 Illustrative analyses

- Planck
- **Euclid**
- LSST
- SKA

3 Concluding remarks

LSS on the 3D ball



Credit: SDSS

Fourier-LAGuerre wavelets (flaglets) on the ball

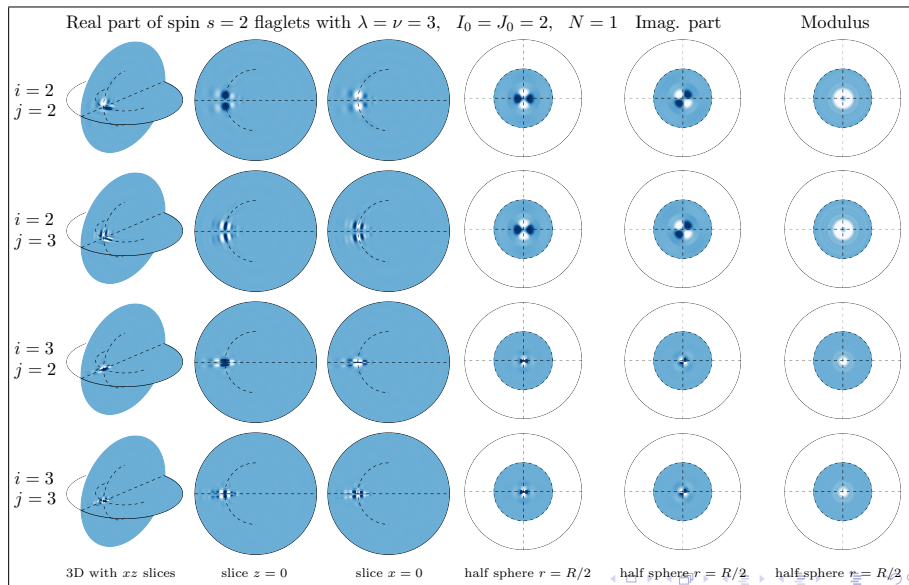
- Fourier-Laguerre wavelet (flaglet) transform is given by the projection onto each wavelet (Leistedt & McEwen 2012):

$$W^s \Psi^{jj'}(r, \rho) = \underbrace{\langle sf, \mathcal{T}_{(r, \rho)} s \Psi^{jj'} \rangle}_{\text{projection}} = \int_{\mathbb{B}^3} d^3 \mathbf{r} sf(\mathbf{r}) (\mathcal{T}_{(r, \rho)} s \Psi^{jj'})^*(\mathbf{r}).$$

- Original function may be recovered exactly in practice from wavelet coefficients:

$$sf(\mathbf{r}) = \underbrace{\sum_{j j'}}_{\text{finite sum}} \underbrace{\int_{\text{SO}(3)} d\varrho(\rho) \int_{\mathbb{R}^+} d\mathbf{r} W^s \Psi^{jj'}(r, \rho) (\mathcal{T}_{(r, \rho)} s \Psi^{jj'})^*(\mathbf{r})}_{\text{wavelet contribution}}.$$

Fourier-LAGuerre wavelets (flaglets) on the ball



3D weak lensing

- 3D weak lensing with spin wavelets on the ball (Leistedt, McEwen, Kitching, Peiris 2015).
- Wavelet transform of 3D cosmic shear:

$$W_{2\gamma}^{2\Psi^{ij}}(\mathbf{n}, r) = (2\gamma \odot {}_2\Psi^{ij})(\mathbf{n}, r)$$

- Wavelet covariance:

$$C^{ij, i'j'}(\mathbf{n}, \mathbf{n}', r, r') = \langle W_{2\gamma}^{2\Psi^{ij}}(\mathbf{n}, r) W_{2\gamma}^{2\Psi^{i'j'}} * (\mathbf{n}', r') \rangle$$

compute from data

- Theory wavelet covariance:

$$C^{ij, i'j'}(\mathbf{n} \cdot \mathbf{n}', r, r') = \frac{2}{\pi} \sum_{\ell} \frac{(N_{\ell, 2})^2}{4} \int_{\mathbb{R}^+} dk k^2 \int_{\mathbb{R}^+} dk' k'^2 C_{\ell}^{\phi\phi}(k, k') P_{\ell}(\mathbf{n} \cdot \mathbf{n}') {}_2\mathcal{H}_{\ell}^{ij}(k, r) {}_2\mathcal{H}_{\ell}^{i'j'} * (k', r')$$

compute from theory

- Simultaneous spatial and scale representation (can handle **complicated sky coverage** and **filter unreliable harmonic modes**).

Outline

1 Big-data in astronomy and astrophysics

2 Illustrative analyses

- Planck
- Euclid
- **LSST**
- SKA

3 Concluding remarks

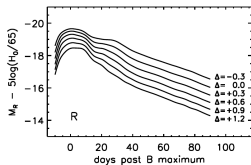
Photometric supernova classification

Machine learning

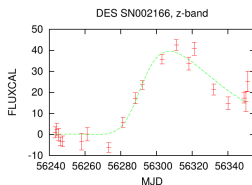
- Photometric supernova classification by machine learning (Lochner, McEwen, Peiris & Lahav, in prep.)
- Go beyond single techniques to study classes.
- Understand physical requirements (e.g. representative training, redshift).



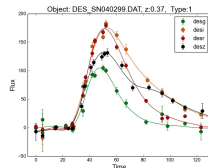
Michelle
Lochner



(a) Templates



(b) Generic parameterisations

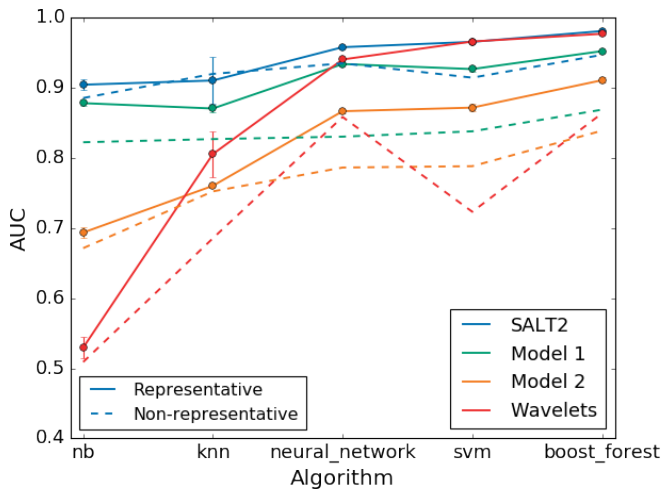


(c) Wavelets (non-parametric)

Figure: Feature selection classes (in order of increasing model independence)

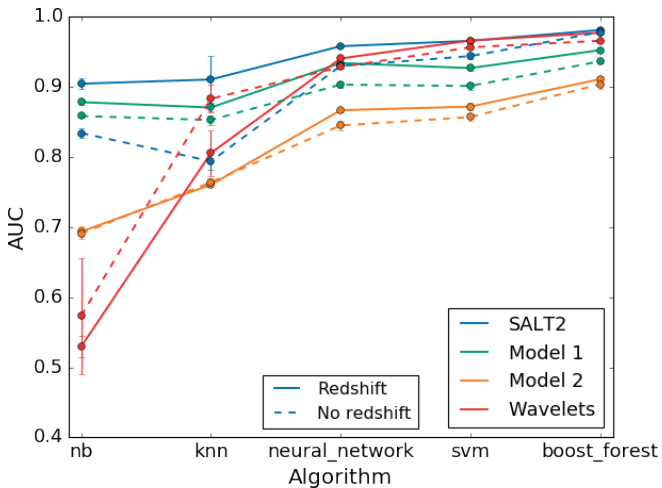
Photometric supernova classification

Importance of representative training data



Photometric supernova classification

Importance of redshift



Outline

1 Big-data in astronomy and astrophysics

2 Illustrative analyses

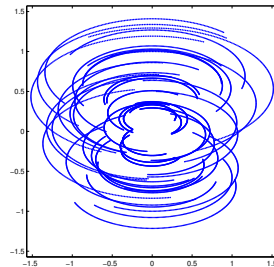
- Planck
- Euclid
- LSST
- SKA

3 Concluding remarks

Radio interferometric telescopes acquire “Fourier” measurements



“Fourier”
Measurements



Compressive sensing

- Developed by Candes *et al.* 2006 and Donoho 2006 (and others).
- Although many underlying ideas around for a long time.
- Exploits the **sparsity** of natural signals.
- Active area of research with many new developments.

SARA for radio interferometric imaging

Algorithm

- Sparsity averaging reweighted analysis (**SARA**) for RI imaging (Carrillo, McEwen & Wiaux 2012)
- Consider a dictionary composed of a **concatenation of orthonormal bases**, i.e.

$$\Psi = \frac{1}{\sqrt{q}} [\Psi_1, \Psi_2, \dots, \Psi_q],$$

thus $\Psi \in \mathbb{R}^{N \times D}$ with $D = qN$.

- We consider the following bases: **Dirac** (i.e. pixel basis); **Haar wavelets** (promotes gradient sparsity); **Daubechies wavelet bases two to eight**.
 \Rightarrow concatenation of 9 bases
- Promote average sparsity by solving the **reweighted ℓ_1 analysis problem**:

$$\min_{\bar{x} \in \mathbb{R}^N} \|W\Psi^T \bar{x}\|_1 \quad \text{subject to} \quad \|y - \Phi \bar{x}\|_2 \leq \epsilon \quad \text{and} \quad \bar{x} \geq 0,$$

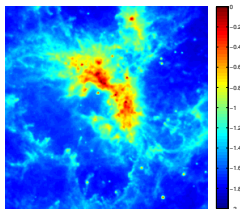
SARA

where $W \in \mathbb{R}^{D \times D}$ is a diagonal matrix with positive weights.

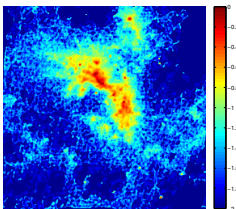
- Solve a sequence of reweighted ℓ_1 problems using the solution of the previous problem as the inverse weights \rightarrow **approximate the ℓ_0 problem**.

SARA for radio interferometric imaging

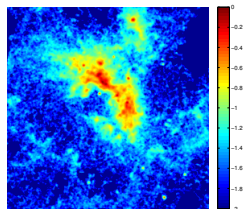
Results on simulations



(a) Original



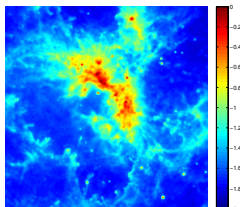
(b) "CLEAN" (SNR=16.67 dB)



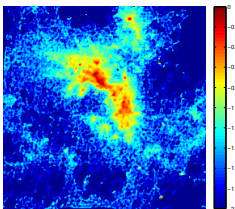
(c) "MS-CLEAN" (SNR=17.87 dB)

SARA for radio interferometric imaging

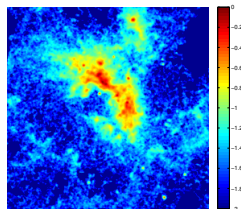
Results on simulations



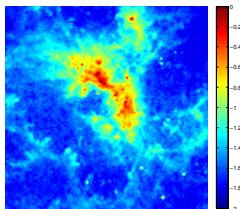
(a) Original



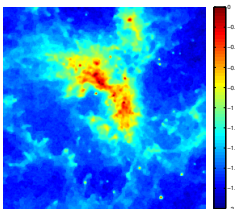
(b) "CLEAN" (SNR=16.67 dB)



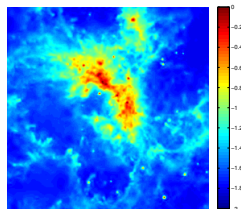
(c) "MS-CLEAN" (SNR=17.87 dB)



(d) BPDb8 (SNR=24.53 dB)



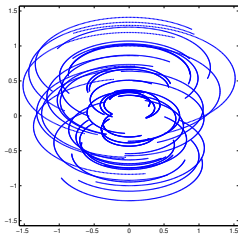
(e) TV (SNR=26.47 dB)



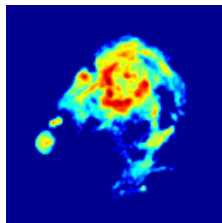
(f) SARA (SNR=29.08 dB)

Supporting continuous visibilities

Results on simulations



(a) Coverage

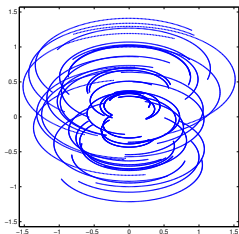


(b) M31 (ground truth)

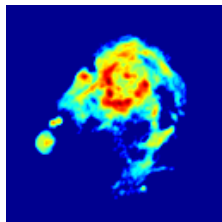
Figure: Reconstructed images from continuous visibilities.

Supporting continuous visibilities

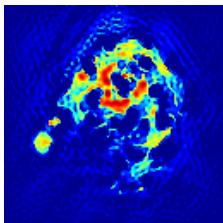
Results on simulations



(a) Coverage



(b) M31 (ground truth)

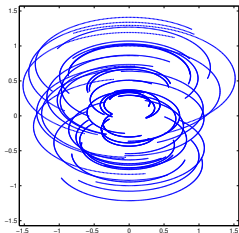


(c) "CLEAN" (SNR= 8.2dB)

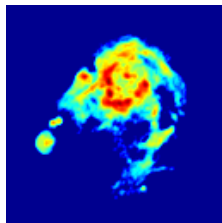
Figure: Reconstructed images from continuous visibilities.

Supporting continuous visibilities

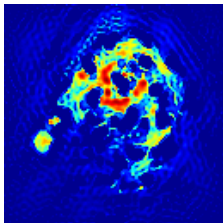
Results on simulations



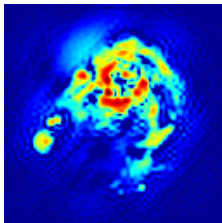
(a) Coverage



(b) M31 (ground truth)



(c) "CLEAN" (SNR= 8.2dB)

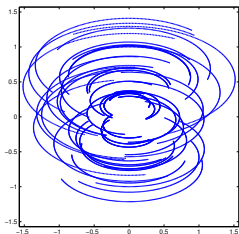


(d) "MS-CLEAN" (SNR= 11.1dB)

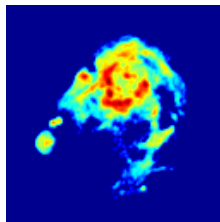
Figure: Reconstructed images from continuous visibilities.

Supporting continuous visibilities

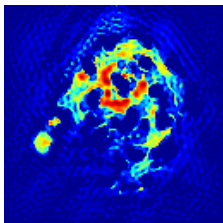
Results on simulations



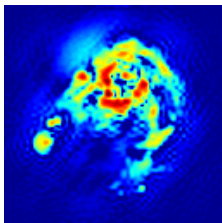
(a) Coverage



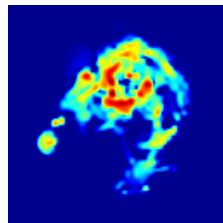
(b) M31 (ground truth)



(c) "CLEAN" (SNR= 8.2dB)



(d) "MS-CLEAN" (SNR= 11.1dB)



(e) SARA (SNR= 13.4dB)

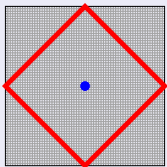
Figure: Reconstructed images from continuous visibilities.

Distributed algorithms and codes

- Distributed storage and computation (Onose *et al.* 2016) by **divide-and-conquer** and **sub-sampling** techniques

SOPT code

<http://basp-group.github.io/sopt/>



Sparse OPTimisation

Carrillo, McEwen, Wiaux

SOPT is an open-source code that provides functionality to perform sparse optimisation using state-of-the-art convex optimisation algorithms.

PURIFY code

<http://basp-group.github.io/purify/>



Next-generation radio interferometric imaging

Carrillo, McEwen, Wiaux

PURIFY is an open-source code that provides functionality to perform radio interferometric imaging, leveraging recent developments in the field of compressive sensing and convex optimisation.

Concluding remarks

- Increasingly **inter-disciplinary**, drawing on statistics, applied mathematics, computer science, information engineering, . . .
- Increasingly **intra-disciplinary** (e.g. Planck, Euclid, LSST, SKA, . . .)
- Many **methodological synergies**

Concluding remarks

How can we exploit synergies?

- 1 Open (unencumbered) data and open code
- 2 Develop **best practices** (e.g. code development, general codes, reproducible/replicable research, blinded analysis)
- 3 Explore **HPC synergies** (e.g. Dirac, Archer, Hartree, Google, Amazon, ...)
- 4 Develop appropriate **career progression** routes
- 5 Go beyond individual techniques to understand properties of **classes of approach**
- 6 Develop **common language**
- 7 Promote inter- and intra-disciplinary **collaboration and communication**, e.g. Alan Turing Institute (ATI), workshops (e.g. BASP conference), Hackathons, ...
- 8 ...

Concluding remarks

How can we exploit synergies?

- 1 Open (unencumbered) data and open code
- 2 Develop best practices (e.g. code development, general codes, reproducible/replicable research, blinded analysis)
- 3 Explore HPC synergies (e.g. Dirac, Archer, Hartree, Google, Amazon, ...)
- 4 Develop appropriate career progression routes
- 5 Go beyond individual techniques to understand properties of classes of approach
- 6 Develop common language
- 7 Promote inter- and intra-disciplinary collaboration and communication, e.g. Alan Turing Institute (ATI), workshops (e.g. BASP conference), Hackathons, ...
- 8 ...

Concluding remarks

How can we exploit synergies?

- 1 Open (unencumbered) data and open code
- 2 Develop best practices (e.g. code development, general codes, reproducible/replicable research, blinded analysis)
- 3 Explore HPC synergies (e.g. Dirac, Archer, Hartree, Google, Amazon, ...)
- 4 Develop appropriate career progression routes
- 5 Go beyond individual techniques to understand properties of classes of approach
- 6 Develop common language
- 7 Promote inter- and intra-disciplinary collaboration and communication, e.g. Alan Turing Institute (ATI), workshops (e.g. BASP conference), Hackathons, ...
- 8 ...

Concluding remarks

How can we exploit synergies?

- 1 Open (unencumbered) data and open code
- 2 Develop **best practices** (e.g. code development, general codes, reproducible/replicable research, blinded analysis)
- 3 Explore **HPC** synergies (e.g. Dirac, Archer, Hartree, Google, Amazon, ...)
- 4 Develop appropriate **career progression** routes
- 5 Go beyond individual techniques to understand properties of **classes of approach**
- 6 Develop **common language**
- 7 Promote inter- and intra-disciplinary **collaboration and communication**, e.g. Alan Turing Institute (ATI), workshops (e.g. BASP conference), Hackathons, ...
- 8 ...

Concluding remarks

How can we exploit synergies?

- 1 Open (unencumbered) data and open code
- 2 Develop best practices (e.g. code development, general codes, reproducible/replicable research, blinded analysis)
- 3 Explore HPC synergies (e.g. Dirac, Archer, Hartree, Google, Amazon, ...)
- 4 Develop appropriate career progression routes
- 5 Go beyond individual techniques to understand properties of classes of approach
- 6 Develop common language
- 7 Promote inter- and intra-disciplinary collaboration and communication, e.g. Alan Turing Institute (ATI), workshops (e.g. BASP conference), Hackathons, ...
- 8 ...

Concluding remarks

How can we exploit synergies?

- 1 Open (unencumbered) data and open code
- 2 Develop best practices (e.g. code development, general codes, reproducible/replicable research, blinded analysis)
- 3 Explore HPC synergies (e.g. Dirac, Archer, Hartree, Google, Amazon, ...)
- 4 Develop appropriate career progression routes
- 5 Go beyond individual techniques to understand properties of classes of approach
- 6 Develop common language
- 7 Promote inter- and intra-disciplinary collaboration and communication, e.g. Alan Turing Institute (ATI), workshops (e.g. BASP conference), Hackathons, ...
- 8 ...

Concluding remarks

How can we exploit synergies?

- 1 Open (unencumbered) data and open code
- 2 Develop best practices (e.g. code development, general codes, reproducible/replicable research, blinded analysis)
- 3 Explore HPC synergies (e.g. Dirac, Archer, Hartree, Google, Amazon, ...)
- 4 Develop appropriate career progression routes
- 5 Go beyond individual techniques to understand properties of classes of approach
- 6 Develop common language
- 7 Promote inter- and intra-disciplinary collaboration and communication, e.g. Alan Turing Institute (ATI), workshops (e.g. BASP conference), Hackathons, ...
- 8 ...