XXX CANARY ISLANDS WINTER SCHOOL, NOV 5-9, 2018.

DiRAC

Challenges in Image Processing

Turning Petabytes of Data into Terabytes of Information

Prof. Mario Juric DIRAC Institute | eScience Institute | UW Astronomy

With thanks to Andy Connolly, Robert Lupton, Ian Sullivan, David Reiss, Jake vanderPlas + LSST Project and Collaborations

DATA INTENSIVE RESEARCH IN ASTROPHYSICS AND COSMOLOGY COLLEGE OF ARTS & SCIENCES I UNIVERSITY of WASHINGTON

Context



Were used to performing analysis in catalog space, using objects and observations. But how those catalogs are created is a mystery to most astronomers.

Yet science need push us towards:

- Measurements of subtle (10-6) effects
- Ever higher precision (<1% photometry, ~mas astrometry)
- Ever more data (and thus the need for automated operation)

At these precision levels it becomes critical to understand our instruments and the measurement process. Otherwise, the instrumental noise (bias) may be misinterpreted for a signal.

This lecture is a high-level walk-through to give you a feel for how we turn images to catalogs, and the associated challenges today and in the future.

An analog universe (1950s)





Visual inspection: Lick survey, 1246 plates 1.6 million "pixels" 10x10 arcmin pixels

Processing rate Analysis 1947 – 1954 Published 1967

Map of galaxy counts down to V<19

A digitized universe (1980's)





Microdensitometers: 10¹² pixel sky 0.5 arcsec pixels

Processing rate 4 hours a plate 200 Mhz Pentium-pro

Automatic Plate Measuring facility in Cambridge, UK.

A digital universe (2000s)





Mosaic Cameras 10¹¹ pixels sky 0.45 arcsec pixels

Processing rate 4 MB/s 250K lines of code

The Large Synoptic Survey Telescope

30 m diameter dome

1.2 m diameter

atmospheric telescope Control room and heat producing equipment (lower level)

1,380 m² service and

maintenance facility



Stray light and Wind Screen 350 ton telescope

A project to build the biggest optical astronomical survey in existence, bringing data to topics ranging from the Solar System to Dark Energy.



Calibration Screen

LSST camera: A 3.2 Gigapixel camera







Modular design: 3200 Megapix = 189 x16 Megapix CCD 9 CCDs share electronics: raft (21=camera) 100 µm deep depletion devices (10 µm pixels)





Sensors at scale: 6.5x10⁹ pixels sky 0.2 arcsec pixels

Processing rate 170 MB/s 1.1M lines of code

Credit: John Peterson

A Big Data Universe



1 chip (4kx4k, 18 bits/pixel), 0.5% of the full image.

Expect ~2000 exposures per night, 300 nights a year, for 10 years.

Roughly 4 PB of raw imaging data per year.



What do astronomers care about?

- What's on the image?
 - Stars (point sources)
 - Galaxies (extended objects)
- Where is it?
 - Relatively (in pixels, to ~few hundredths of a pixel)
 - Absolutely (coordinates on the sky)
- How bright is it?
- Is it changing in time?
- Is it moving?
- What is its shape?
 - Of a particular object
 - Statistically, for a class of objects



Understanding Astronomical Images





GOODS-South Field, deep VLT stack



How do astronomical images come to be?

> Credit: John Peterson (Purdue) and the PhoSim Team

Optics

+Tracking

+Diffraction

+Detector Misalignments & Perturbations









+Mirror Misalignments Perturbations, & Micro-roughness +Detector

+High Altitude Atmosphere

+Mid Altitude Atmosphere +Low Altitude Atmosphere

+Pixelization

+Saturation & Blooming

Astronomical Images (approximated)



$I(x) = \phi(x) \otimes S(x) + \varepsilon(x)$ Observed PSF Truth Noise

Point Spread Function Estimation (Today)





Sampling, modeling, interpolation



Challenge: How could we do better?

Utilize all available information.

The PSF is a combination of deterministic and stochastic elements.

Deterministic: Camera, Optics Stochastic: Atmospheric perturbations



Directly measure the deterministic elements. Constrain the stochastic elements (e.g., using the knowledge of the power spectrum of fluctuations in the atmosphere).

e.g., Jee & Tyson 2011



Life is not so Simple: Tree Rings



Plazas, Bernstein & Sheldon 2014

Magnier et al. 2017



https://www.youtube.com/watch?v=jh2z-g7GJxE

Charge Diffusion

In LSST sensors, one sees tree rings variations at at a ~percent level.

Varying dopant density in silicon boules creates parasitic lateral E fields.

These DO NOT behave as QE variations. Naïve flat-fielding makes the problem WORSE by ~2x.

two pixels, 10 x 10 micron each Edite 100 micron depleted Channel stop detector volume implants confining 4 phase parallel potential clock lines in y direction, from parallel clocks confining potential in x direction, due to channel stops

Nonlinearities: the Brighter-Fatter Effect

Most of today's devices (DES, HSC, LSST, GPC1) are thick. The photon converting at the top has a long way to go to reach the bottom

- Tree rings (and related effects)
- "Brighter-fatter" effect

As the potential wells fill up with electrons, the bias voltage drops making it easier for electrons to be diverted to neighboring pixels.

Correlates the values of neighboring pixels; results in an intensity-dependent PSF.





Measurement = Modeling

Deconvolve?



$I(x) = \phi(x) \otimes S(x) + \varepsilon(x)$ Observed PSF Truth Noise

If we can estimate the PSF, can we simply deconvolve?

No. (at least not easily.) Deconvolution without regularization amplifies noise. Errors in the estimate of the PSF are further amplified in deconvolved images.

Forward Modeling

$$I(x) = \phi(x) \otimes S(x) + \varepsilon(x)$$

observed PSF Truth Noise

 $S(x) = \sum M_i(x, y)$

Object characterization Models:

- Stars: Point Source model
- Galaxies: Double Exponential models



Modeling object properties

Object characterization (models):

- Stars: Point Source model
- Galaxies: Double Exponential models





7757,301,1,74,187,6,8,12783867556709,26,627245975921,17,37402,17,92875,0.02894481,0.02568013 7757,301,1,74,188,3,8,12732322524192,26,6251199416623,20,1466,21,35297,0,3003744,0,3302762 4288.301.1.39.682.3.24.5161170422305.-1.16579446393527.22.97032.24.3259.0.2672399.0.5240437 4288.301.1.39.683.3.24.5179406515354.-1.1792069022485.22.62052.25.09109.0.1850479.0.6585805 4288.301.1.39.684.6.24.5189463293148.-1.15915086108891.21.4247.23.04125.0.06608655.0.1968172 4136,301,1,61,935,6,36,4715922759092,-1.06093938828308,22.71782,23.14112,0.158014,0.1799687 4136,301,1,61,936,3,36,4717583013136,-1,1378448207726,22,81683,23,88123,0,1742272,0,3260605 4136.301.1.61.937.3.36.4717582434391.-1.13784497192974.22.81147.23.87586.0.1734457.0.3247895 4288.301.1.40.311.3.24.6839203338022.-1.23631696217547.21.2002.21.67521.0.05694564.0.06316777 4288.301.1.40.312.3.24.6840602692246.-1.21784918362007.20.30287.21.04976.0.02972161.0.04032911 4288, 301, 1, 40, 313, 6, 24, 6840216690377, -1, 08292772289886, 24, 92263, 25, 72778, 0, 68427, 0, 5471938 5598,301,1,61,792,3,351.787950113407,6.14573538435867,22.43574,23.83793,0.1125768,0.2867745 5598,301,1,61,793,3,351,787950113434,6,14573538316393,22,43573,23,77753,0,1174541,0,2741905 5598.301.1.61.794.6.351.787349107439.6.14481612145222.24.6701.24.8507.0.4675894.0.4849904 2699.301.1.48.527.3.12.0760019016408.-3.32677418219699.22.18116.23.27577.0.126546.0.2270369 2699.301.1.48.528.6.12.0770027529666.-3.32913243320258.22.12757.23.79366.0.1215217.0.3403472 2699, 301, 1, 48, 529, 3, 12, 0832728187538, -3, 52539818226738, 22, 29741, 23, 32008, 0, 1377919, 0, 2253349 94,301,1,38,279,6,340,524768659138,-0.843090883870374,20,75028,21,13888,0,04839022,0,04644739 94,301,1,38,280,6,340,525793656628,-0,965210498356983,24,23321,26,07633,0,8058653,0,7702702 94.301.1.38.281.6.340.53257887691.-1.02365035629542.21.10608.21.15248.0.06605724.0.04672106 4288.301.1.76.766.3.30.0899167300738.-1.25189466355601.22.57054.22.91144.0.1807321.0.1984752 4288,301,1,76,767,3,30.0899962733195,-1.14314301954175,21.99419,23.78533,0.1094794,0.4162939 4288,301,1,76,768,6,30,0899156247035,-1,19236549812758,22,64196,22,91776,0,1927483,0,2022232 7937.301.1.84.354.3.5.59132226470183.26.6120638856974.22.63659.23.76799.0.2950225.0.5986399 7937.301.1.84.355.3.5.58984744314193.26.6157353187021.21.27009.22.13992.0.08808753.0.1490507 7937.301.1.84.356.6.5.59007089167375.26.5703384153732.24.52393.26.33231.1.060107.0.6528299 3996, 301, 1, 81, 615, 6, 216, 223187662076, 11, 9472 1,231,051,050 rows (SDSS DR10, PhotoObjAll table) 3996, 301, 1, 81, 616, 3, 216, 23835672109, 12, 05342 3996.301.1.81.617.3.216.219379232008.11.9171 ~500 columns 6354,301,1,29,2997,6,332.302798261878,41.217 <u>6354,301,1,29,2998,6,332.355200585822,41.244/00131300/,20.70133,22.35575,0.03523923,0.03734075</u>



Cataloging the Sky...





What we're doing is decomposing and modeling the sky in a way that makes physical sense.

... Compressing the Sky





What we're doing is decomposing and modeling the sky in a way that makes physical sense.

But you may also think of this as developing a very astronomyspecific lossy compression technique.

Challenges: Incomplete model space







C 🚽 🕨 D



Beyond a Single Image

Beyond a Single Image





Beyond a Single Image





Going Deep: Coaddition

Changes in Time: Image Differencing



Image Differencing

Why Difference?





Why Difference?





Why Difference?





First image

Second image

Difference image

Align, Subtract, Profit!

Right?...









Alard-Lupton Algorithm





 $I(x) = \phi(x) \otimes S(x) + \varepsilon(x)$ PSF Truth Noise $D(x) = I_1(x) - \kappa(x) \otimes I_2(x)$

Image Image Difference

How to match image quality



Difference imaging requires solving for the mapping kernel



$$k(x) = \sum_{i} a_{i} B_{i}(x) = \sum_{n,p,q} a_{i} e^{-\frac{u^{2}+v^{2}}{2\sigma_{n}^{2}}} u^{p} v^{q}$$

$$\frac{I_1 - \sum_i a_i B_i(x) \otimes I_2}{\sigma}$$

Alard and Lupton 97...

Works well but...

$D(x) = I_1(x) - \kappa(x) \otimes I_2(x)$



Assumes template PSF is smaller than image PSF

Deconvolution otherwise -> BAD!

Assumes the "template" has no noise

- Derived from a previous set of observations
- Recent work has shown that for Gaussian, heteroschedastic noise we can take a Fourier Transform, and compute the log-likelihood and prewhiten the images

$$\hat{D}(k) = (I_1(k) - \kappa(k)I_2(k)) \sqrt{\frac{\sigma_1^2 + \sigma_2^2}{\sigma_1^2 + \kappa^2(k)\sigma_2^2}}$$

Zackay et al 2016

This is running on ZTF!

Real world is more complicated...

Pan-STARRS1 Systematic False Detection Gallery



In real-world applications, differencing two images is never perfect. In fact, it's far from perfect!

Typically, the observed difference images are littered with artifacts, <u>false positives</u>, that make astronomers sad and unhappy.

Typically, we see 100:1 to 10:1 false-to-true detection ratios (that is not a typo!).



Real world is more complicated...

Pan-STARRS1 Systematic False Detection Gallery



Where do these things come from?

- Image misalignments
- Imperfect knowledge
 of PSF variation along
 the image
- CCD defects

....

- Readout electronics
 artifacts
- Optical ghosts and glints



Machine Learning to the Rescue!

Above: Figure 1, Goldstein al. (2015), AJ, 150, 82

Given a sample of real detections and false detections, teach the computer to recognize the difference between the two and judge assign a "score", τ , to each (τ =1 -> real, τ =0 -> false).



Fig. 7.— 5-fold cross-validated receiver operating characteristics of the best-performing classifier from $\S3.5$. Six visually indistinguishable curves are plotted: one translucent curve for each round of cross-validation, and one opaque curve representing the mean. Points on the mean ROC corresponding to different class discrimination boundaries τ are labeled. $\tau = 0.5$ was adopted in DES-SN.

2 2		No ML	ML ($\tau = 0.5$)	ML / No ML
1	N_c^{a}	$100,\!450$	$7,\!489$	0.075
	$N_A/N_{NA}\rangle^{\rm b}$	13	0.34	0.027
ε	F ^C	1.0	0.990	0.990

	autoScan	candidate-level	emciency	ior lake	Sive la.
Fake			v		
TURC					

SN	Bad Column	Sub	Sub	Sub

Above: Figure 1, Goldstein al. (2015), AJ, 150, 82

se detections, teach the computer to recognize the between the two and judge assign a "score", τ , to each (τ =1 -> real, τ =0 -> false).



TABLE 4 CAN ON REPROCESSED DES Y1 TRANSIENT CANDIDATE					
	No ML	ML ($\tau = 0.5$)	ML / No ML		
$N_c^{\mathbf{a}}$	$100,\!450$	$7,\!489$	0.075		
$\langle N_A/N_{NA}\rangle^{\rm b}$	13	0.34	0.027		
$\epsilon_F{}^{\mathrm{c}}$	1.0	0.990	0.990		

^aTotal number of science candidates discovered.

^bAverage ratio of artifact to non-artifact detections in human scanning pool.

^cautoScan candidate-level efficiency for fake SNe Ia.

Bad Column Sub Sub	Sub

Above: Figure 1, Goldstein al. (2015), AJ, 150, 82

Machine Learning to the Rescue!

Given a sample of real detections and false detections, teach the computer to recognize the difference between the two and judge assign a "score", τ , to each (τ =1 -> real, τ = 0 -> false).

Option #2: Understand the root cause(s)





Example: Astrometric misalignment introduces dipoles in the images, misalignment of >2% of a pixels will dominate the number of false positives

Major source is Differential Chromatic Refraction

Atmosphere refracts (shifts) a source more for blue light than red (even for light measured through the same filter)

A bandpass is not a delta function.





How differential refraction works





A bandpass is not infinitely narrow





A few more details...





"Template Image" (often a *coadd*) "Science Image"

Image Difference

A solution



Instead of having just one template, infer the template as a function of wavelength based on many observations at different airmasses.



=> "red" object

=> "blue" object

A solution



Instead of having just one template, infer the template as a function of wavelength based on many observations at different airmasses!

Effectively, given N exposures at different airmassess, infer the images $I_T(\lambda=blue)$, $I_T(\lambda=center)$, $I_T(\lambda=red)$ for the template(s).

Then, for each science image, we can construct the template exactly for that airmass, before doing the differencing.

DCR-corrected Templates: lan Sullivan et al.





Assume a model for the image, y, that is made up of series of images each of a different wavelength (ie a "hyperspectral" cube)

Byproduct: measuring the intra-band spectrum!





Figure 15 Example input spectrum for a type F star with surface temperature ~7130K (solid blue). The flux measured in each sub-band is marked with a with red '+', and the average values of the simulated spectrum across each subfilter is marked with a blue 'x' for comparison.

This is another example of utilizing all available information and understanding of image generation processes to extract additional information.



Figure 14 Source measurements in three sub-bands of the DCR sky model are converted to RGB values and used to fill the footprints of detected sources. The combined full-band model is displayed behind the footprint overlay.



Going Deeper: Coaddition (a.k.a. "astronomer's HDR")

Why co-add?



Pros:

See fainter objects!

Computationally

inexpensive

9

SDSS Southern Coadd



Left: Annis et al. 2011

FIG. 2.— Comparison between single pass (left) and coadd (right) images in *r*-band for run 206, camcol 3, field 505, RA=15, Dec=0. Images are shown with the same scale, contrast and stretch. The single pass counterpart (run 5800, camcol 3, field 505) is one out of 28 images used in the coaddition of this particular image. This example illustrates the fact that a large number of objects below the detection threshold of each image can be well detected and measured in the coadd.



Why co-add?





FIG. 2.— Comparison between single pass (left) and coadd (right) images in r-band for run 206, camcol 3, field 505, RA=15, Dec=0. Images are shown with the same scale, contrast and stretch. The single pass counterpart (run 5800, camcol 3, field 505) is one out of 28 images used in the coaddition of this particular image. This example illustrates the fact that a large number of objects below the detection threshold of each image can be well detected and measured in the coadd.

Pros:

See fainter objects! Computationally inexpensive

<u>Cons</u>: - Complicated PSF - Correlated noise - Loss of motion and time variability information - Loss of information.

Left: Annis et al. 2011



Better than Coaddition: Multi-Epoch Fitting (MultiFit)

A very simple idea: instead of co-adding pixels of individual observations, and then fitting the model to the result, why not fit the model directly to each individual observation?



MultiFit (Simultaneous Multi-Epoch Fitting)

Opportunity: Recovering motion from the noise



Lang (2009):

Fit moving source models to suspected moving stars in SDSS Stripe 82 survey.

Individual exposures: objects are undetected or marginally detected

Moving point-source and galaxy models are indistinguishable on the coadd



Downsides



Computationally extremely intensive. Scales with the number of epochs (so ~100-1000x more computationally expensive for modern surveys like LSST).

Worse, we're greedy: the physics we're trying to study is so sensitive to biases that ML estimators are not enough; we want posteriors PDFs for parameters of each observed galaxy! (20-200x more output storage!)

But computing is getting cheaper. LSST is building a ~2 PFLOP machine to do this (cca ~2025). Still cheaper than building a bigger telescope (and/or launching it into space)!

Putting it all together





Figure 2: Illustration of the conceptual design of LSST science pipelines for imaging processing.

A Common Theme: Understanding -> Information





Much of the information "loss" comes not from the physical instrument, but the data processing method.

By improving the processing (algorithms, software), we can do a factors of few better than <u>with the exact same dataset</u>.





The "next gen" survey may be just software!





SDSS DR12

Relative photometric calibration accuracy (Padmanabhan et al.

Summary



- Measurement used to be simple: "CCDs are linear", "pixels are independent photon buckets",
- We're in an era of high precision astronomy; this requires we drop these simplifications
 - Understand the physics of the instrument
 - Properly perform inference (measurement)
- Thinking about these problems can reveal interesting new opportunities.
 - Many instrumental effects are information preserving (or revealing!)
 - It is our poor measurement techniques that erase information.
 - We now have enough computing power to do things right.