# シミュレーションデータベースと統計学 吉田直紀, 西道啓博, 大木平 (KAVLI IPMU)

- Cosmology with Subaru HSC survey
- A large number of cosmological N-body simulations
- Parameter space exploration and the Gaussian process

## Observational data and forward modeling

銀河周りの物質分布プロファイル (SDSSの観測データ) 9 (1095 halos)  $\chi^2 / dof = 75.03/77$ Average DM distribution around massive clusters  $a = 0.56^{+0.04}_{-0.05}$ 10<sup>2</sup> 10<sup>2</sup>  $\Delta \Sigma(\mathbf{R}) \; [hM_{\odot}/\text{comoving pc}^2 \;]$  $\Delta \Sigma(\mathbf{R}) \left[ h M_{\odot} / \text{comoving pc}^2 \right]$  $b = 1.48^{+0.09}_{-0.09}$  $\sigma_{logM} = 0.31^{+0.06}_{-0.07}$ 10<sup>1</sup> 10<sup>1</sup> 10<sup>0</sup> 10<sup>0</sup>  $10^{-1}$  $10^{0}$  $10^{1}$  $10^{-1}$  $10^{0}$ 10<sup>1</sup>  $R [h^{-1} comoving Mpc]$  $R [h^{-1} comoving Mpc]$  $40.0 <\!\!\lambda <\!\!55.0,\, 0.1 <\!\!z <\!\!0.33,\, P_{\rm cen}\!>\!\!0.9$  (609 halos)  $55.0 < \lambda < 140.0, 0.1 < z < 0.33, P_{cen} > 0.9 (382 halos)$ 10<sup>2</sup> 10<sup>2</sup>  $\Delta \Sigma(\mathbf{R}) \left[ h M_{\odot} / \text{comoving pc}^2 \right]$  $\Delta \Sigma(\mathbf{R}) \ [hM_{\odot}/\text{comoving pc}^2 \ ]$ 10<sup>1</sup>  $10^{1}$ line=N-body emulator 10<sup>0</sup> 10<sup>0</sup> data = SDSS measurements  $10^{-1}$  $10^{0}$  $10^{1}$  $10^{-1}$ 10<sup>0</sup> 10<sup>1</sup>  $R [h^{-1} comoving Mpc]$  $R [h^{-1} comoving Mpc]$ 

## **COSMOLOGICAL PARAMETER DEPENDENCE**



## **COSMOLOGICAL PARAMETER DEPENDENCE**



## **COSMOLOGICAL PARAMETER DEPENDENCE**



# **HSC GALAXY-GALAXY LENSING**





- Known lens objects
  - host halo properties well known
  - photo-z accuracy not important
  - can measure Σ as a function of scale R instead of angle on the sky



### **EFFICIENT SAMPLING IN MULTI DIMENSIONAL SPACE: LATIN HYPERCUBE**

- Each sample is the only one in each axisaligned hyperplane containing it
  - One can find many realizations of such design (ex. diagonal design)
  - Impose additional condition such as "the sum of the distances to the nearest design point is maximal" (maximin distance)



cosmological parameter 1



## **EFFICIENT SAMPLING IN MULTI DIMENSIONAL SPACE: LATIN HYPERCUBE**

### Simulation spec

- ✓ N of particles: 2048<sup>3</sup>
- ✓ Size: 5 billion lightyears
- ✓ 21 outputs per model
- ✓ 5 Tbyte data per model





## 84 sims are already available

#### parameter space sweep

- "sliced" LH design (Ba, Brenneman & Myers '15)
- generate >100 sample eventually
- maxi-min distance LH design for each 20 model set (e.g., red/blue points)

- two suites of runs
  - keep the initial random number seed (20 done)
  - different seeds (40)
    - red: emulator
    - blue: validation

## **SUPERCOMPUTING AND BIG DATA**

#### 10240 parallel earths



### **EFFICIENT SAMPLING IN MULTI DIMENSIONAL SPACE: LATIN HYPERCUBE**

#### fiducial model

- PLANCK15 (flat ACDM)
- 24 realizations done
- assess statistical error
- emulator accuracy check





# 119 sims are already available parameter space sweep

- "sliced" LH design (Ba, Brenneman & Myers '15)
- generate >200 sample eventually
- maxi-min distance LH design for each 20 model set (e.g., red/blue points)

- two suites of runs
  - keep the initial random number seed (20 done)
  - different seeds (40)
    - red: emulator
    - blue: validation

# **SIMULATION SPEC**

- ✓ N of particles: 2048<sup>3</sup>
- ✓ box size: 1h<sup>-1</sup>Gpc
  - resolve a  $10^{12}$  h<sup>-1</sup>M<sub>solar</sub> halo with ~100 particles
- ✓ 2nd-order Lagrangian PT initial condition @ z<sub>in</sub>=59

(vary slightly for different cosmologies to keep the RMS displacement about 25% of the inter-particle separation)

#### ✓ Tree-PM force by L-Gadget2 (w/ 4096<sup>3</sup> PM mesh)

### ✓ 21 outputs in $0 \le z \le 1.5$

(equispaced in linear growth factor)

- ✓ Data compression (256GB -> 48GB par snapshot)
  - ✓ positions -> displacement (16 bits par dimension; accuracy ~1h<sup>-1</sup>kpc)
  - velocity: discard after halo identification
  - ID: rearrange the order of particles by ID and then discard
  - ✓ already consuming ~200TB in half a year (~observational data)

#### SIMULATIONS



# **GAUSSIAN PROCESS**

- A machine-learning technique to do inference in function space
  - non-parametic Bayesian inference
  - nonlinear regression analysis
- Basic quantities f(x) ~ P [μ (x), k(x, x')]
  - mean function (cf. mean)
  - covariance function (cf. variance)
- covariance function is characterized by a simple function with several hyper parameters

ex. 
$$C(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) = \theta_1 \exp\left[-\frac{1}{2} \sum_{i=1}^{I} \frac{(x_i - x'_i)^2}{r_i^2}\right] + \theta_2.$$
"length scale" r: for each x:



- ✓ infer hyper parameters θ from training data (x<sub>i</sub>, t<sub>i</sub>)
  - Given a point x<sub>N+1</sub>, infer t<sub>N+1</sub> from θ and (x<sub>i</sub>, t<sub>i</sub>)

$$P(t_{N+1} | \mathbf{t}_N) \propto \exp\left[-\frac{1}{2} \left[\mathbf{t}_N \ t_{N+1}\right] \mathbf{C}_{N+1}^{-1} \begin{bmatrix} \mathbf{t}_N \\ t_{N+1} \end{bmatrix}\right]$$
$$\hat{t}_{N+1} = \mathbf{k}^T \mathbf{C}_{N+1}^{-1} \mathbf{t}_N$$

answer:

$$\hat{t}_{N+1} = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N \sigma_{\hat{t}_{N+1}}^2 = \kappa - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}$$

# **GAUSSIAN PROCESS**

- A machine-learning technique to do inference in function space
  - non-parametic Bayesian inference
  - nonlinear regression analysis
- Basic quantities f(x) ~ P [μ (x), k(x, x')]
  - mean function (cf. mean)
  - covariance function (cf. variance)
- covariance function is characterized by a simple function with several hyper parameters

ex. 
$$C(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}) = \theta_1 \exp\left[-\frac{1}{2} \sum_{i=1}^{I} \frac{(x_i - x'_i)^2}{r_i^2}\right] + \theta_2.$$
"length scale" r: for each x:



 infer hyper parameters θ from training data (x<sub>i</sub>, t<sub>i</sub>)

Given a point x<sub>N+1</sub>, infer t<sub>N+1</sub> from θ and (x<sub>i</sub>, t<sub>i</sub>)

$$P(t_{N+1} | \mathbf{t}_N) \propto \exp\left[-\frac{1}{2} \begin{bmatrix} \mathbf{t}_N \ t_{N+1} \end{bmatrix} \mathbf{C}_{N+1}^{-1} \begin{bmatrix} \mathbf{t}_N \\ t_{N+1} \end{bmatrix}\right]$$

answer:

$$\hat{t}_{N+1} = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t}_N$$

$$\sigma_{\hat{t}_{N+1}}^2 = \kappa - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}.$$

## OUR $\Delta \Sigma$ EMULATOR DEMONSTRATIONS



## **GAUSSIAN PROCESS ACCURACY**

# **training** with 20 models (red) **validation** with 20 other models (blue)





# **HMF EMULATOR PERFORMANCE**



# **HMF EMULATOR PERFORMANCE**



# **HMF EMULATOR PERFORMANCE**





# PLANCK 2015

cparam = np.array([[0.02225,0.1198,0.6844,3.094,0.9645,-1]]) set\_cosmo(cparam) give your cosmological params ~5s set\_redshift(z) and redshifts ~600ms; HMF GP called inside lognh = mh\_to\_logdens(Mmin) convert M\_min to n\_h ~50µs plt.loglog(Rplot,get\_dsigma(ascale, lognh, Rplot),lw=2,color='red')

Evaluate !! ~1ms

# SUMMARY + FUTURE

- Modeling the halo mass function and galaxygalaxy lensing signal
  - Latin hypercube design + fitting/GP/spline
  - handy emulator in python ready
  - accuracy test undergoing, aimed at 5% accuracy
- To come
  - ▶ scikit-learn → george
  - RSD emulator to combine g-g lensing and 3D clustering (needs bigger volume)
  - further extension under discussion
    - e.g., non-flat, w0-wa cosmologies

other dependence (time, scale, mass, ...)



(6++)-D cosmological parameter space

## Don't forget, however, the effect of baryons



# Need for accurate templates



Parameter bias with respect to naive estimate from dark matter only simulations

Osato, Shirasaki, NY, 2015