

Modern Analysis Techniques for Large Data Sets

Miguel F. Morales
Bryna Hazelton

Graduate elective

- You get out what you put in
 - Homework (first ~7 weeks; 40%; 100%; + options)
 - Final project & presentation based on your research
 - No exams
- Fill in quiz describing what you'd like to get out of the course

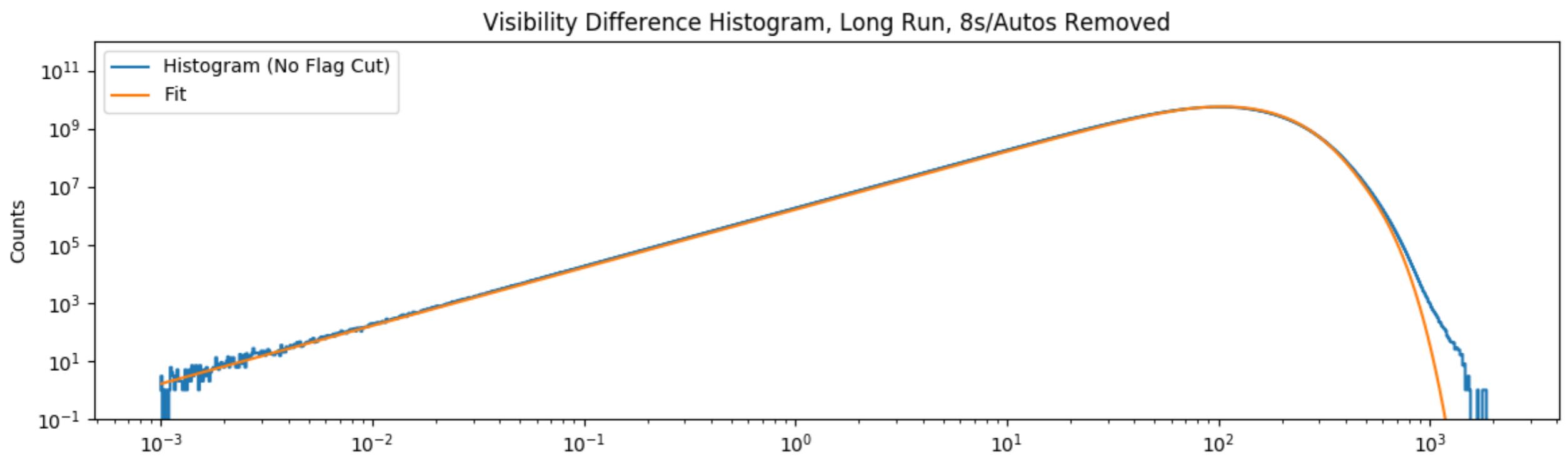
Key topics

- Practical statistics
 - Turning research questions into statistical math
 - Systematics, non-Gaussianity, covariance, avoiding common mistakes
- Data visualization
 - High quality visualizations (data density, perception, accessibility)
 - Types of visualizations (research, monitoring, presentation)
- Collaborative data analysis
 - Using GitHub to your advantage (branching, provenance)
 - Collaborative development
 - Peer reviewed code
- Advanced practices
 - Statistically rigorous visualization
 - Jackknife tests, testing thickets, testing below the noise
 - Question driven data rampages

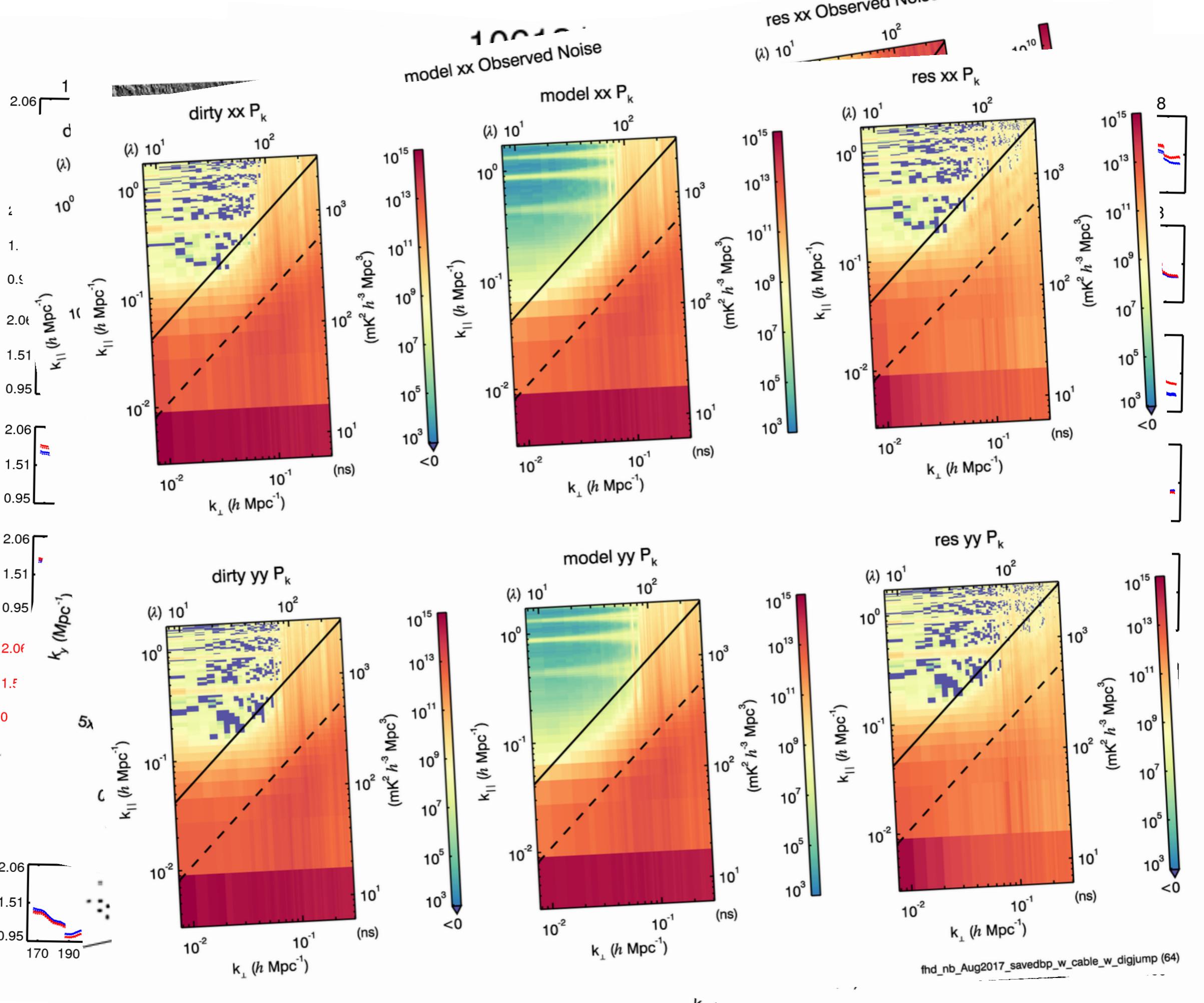
Course structure

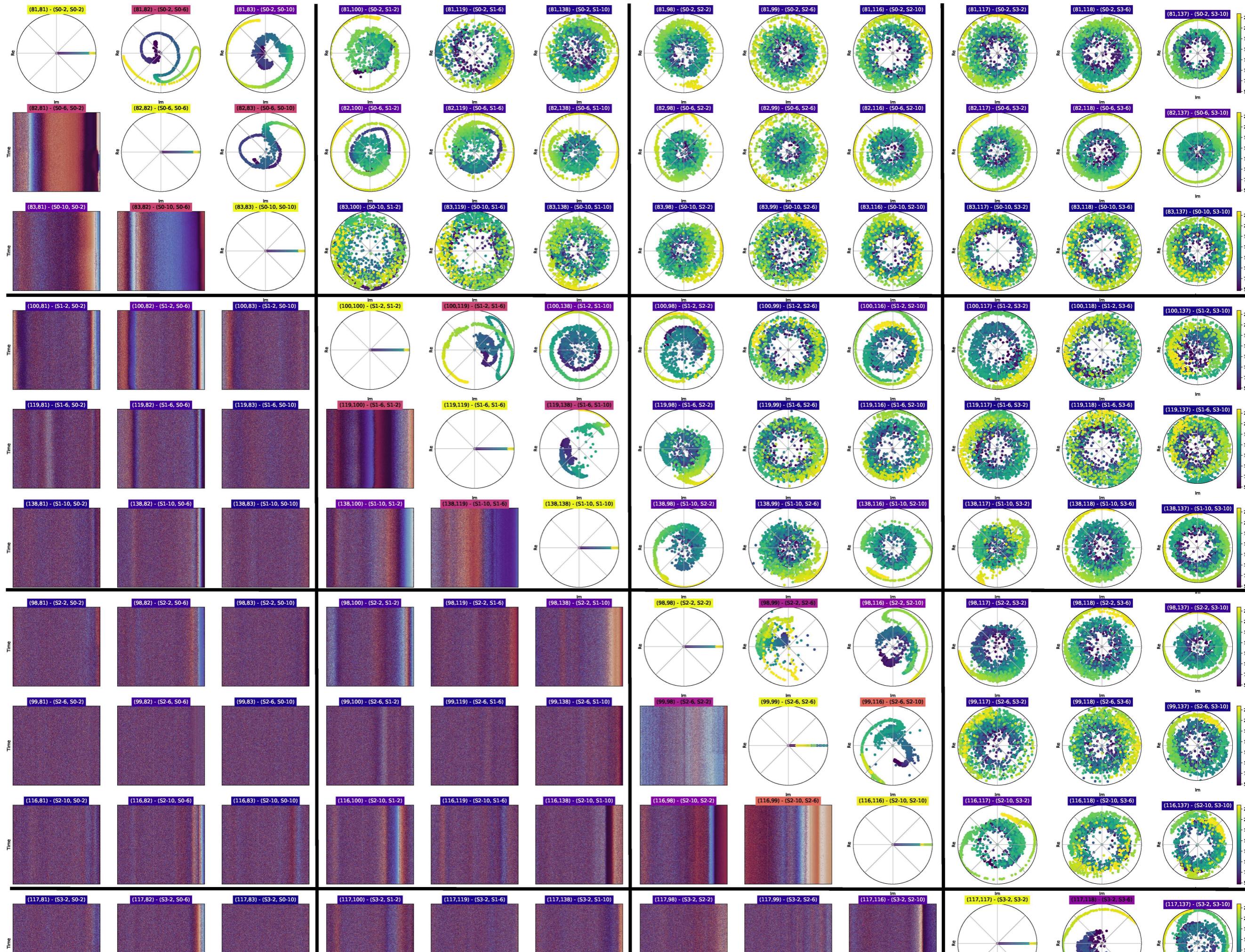
- Building blocks (stats; visualization; git & collaborative code)
- Putting it all together

Statistics



100 billion data points





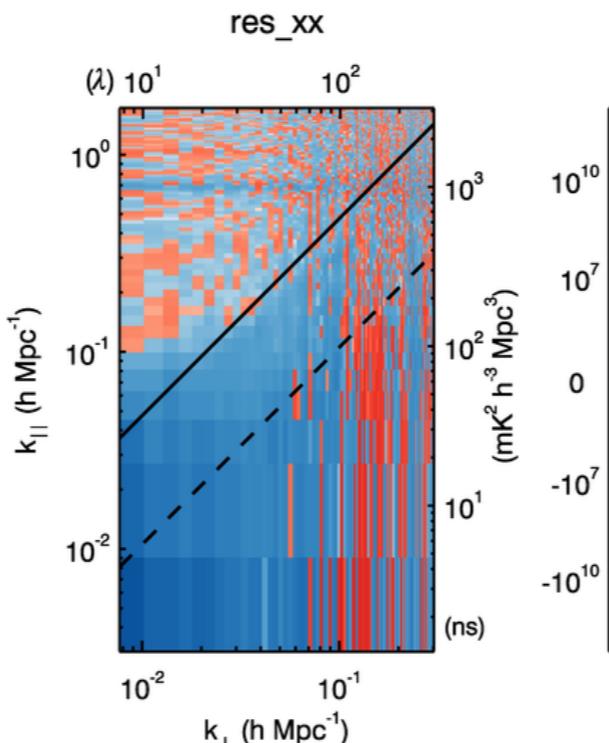
git & analysis traceability



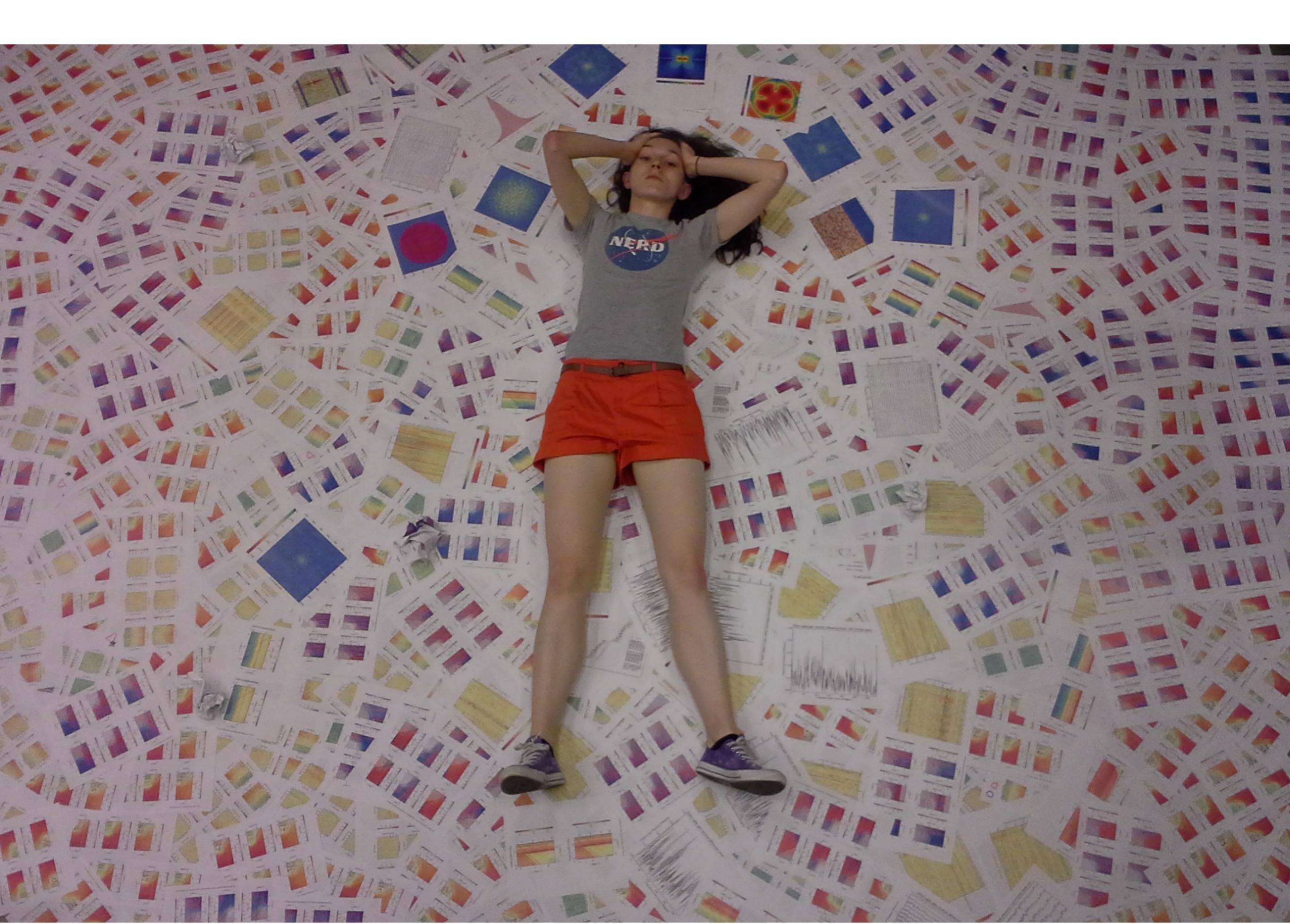
Data unit tests

2 fhd_core/fhd_struct_init_antenna.pro View

		@@ -86,7 +86,7 @@ dec_use=dec_arr[valid_i]
86	86	
87	87	;NOTE: Eq2Hor REQUIRES Jdate to have the same number of elements as RA and Dec for precession!!
88	88	;;NOTE: The NEW Eq2Hor REQUIRES Jdate to be a scalar! They created a new bug when they fixed the old one
89	89	-Eq2Hor,ra_use,dec_use,Jdate,alt_arr1,az_arr1,lat=obs.lat,lon=obs.lon,alt=obs.alt,precess=1
		+Eq2Hor,ra_use,dec_use,Jdate,alt_arr1,az_arr1,lat=obs.lat,lon=obs.lon,alt=obs.alt,precess=1,/nutate
90	90	za_arr=fltarr(psf_image_dim,psf_image_dim)+90. & za_arr[valid_i]=90.-alt_arr1
91	91	az_arr=fltarr(psf_image_dim,psf_image_dim) & az_arr[valid_i]=az_arr1
92	92	





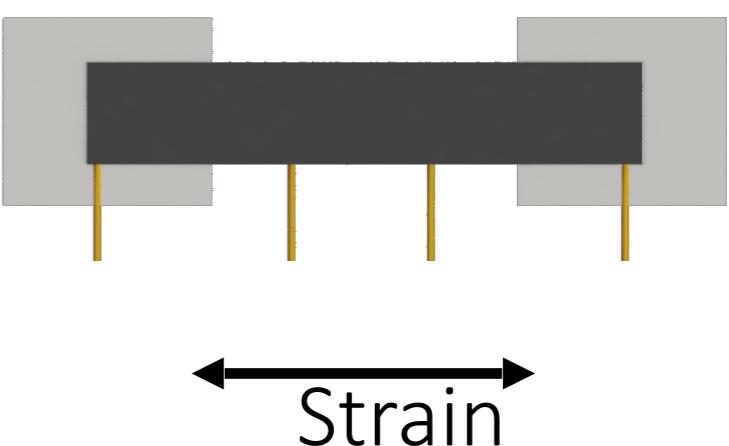
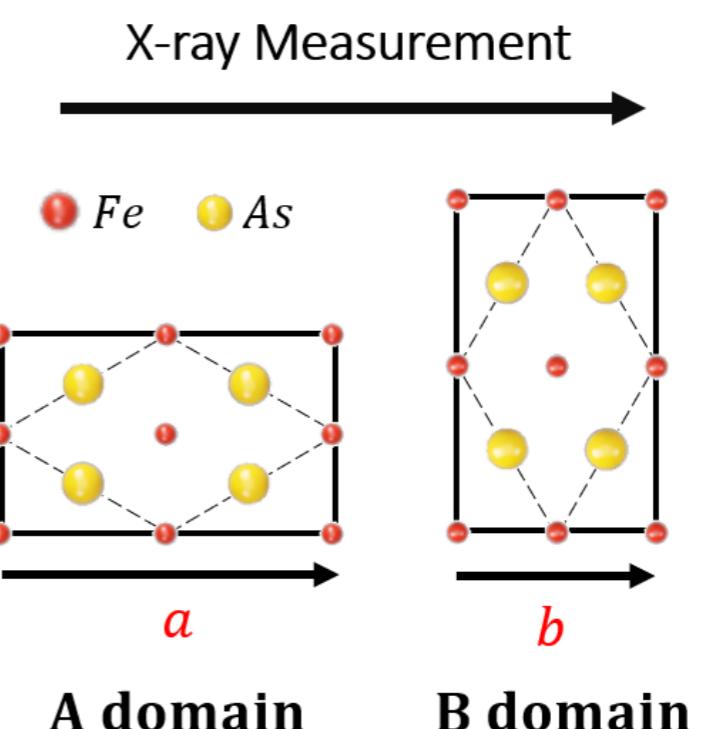


BaFe_2As_2 – Iron Based High Temperature Superconductor

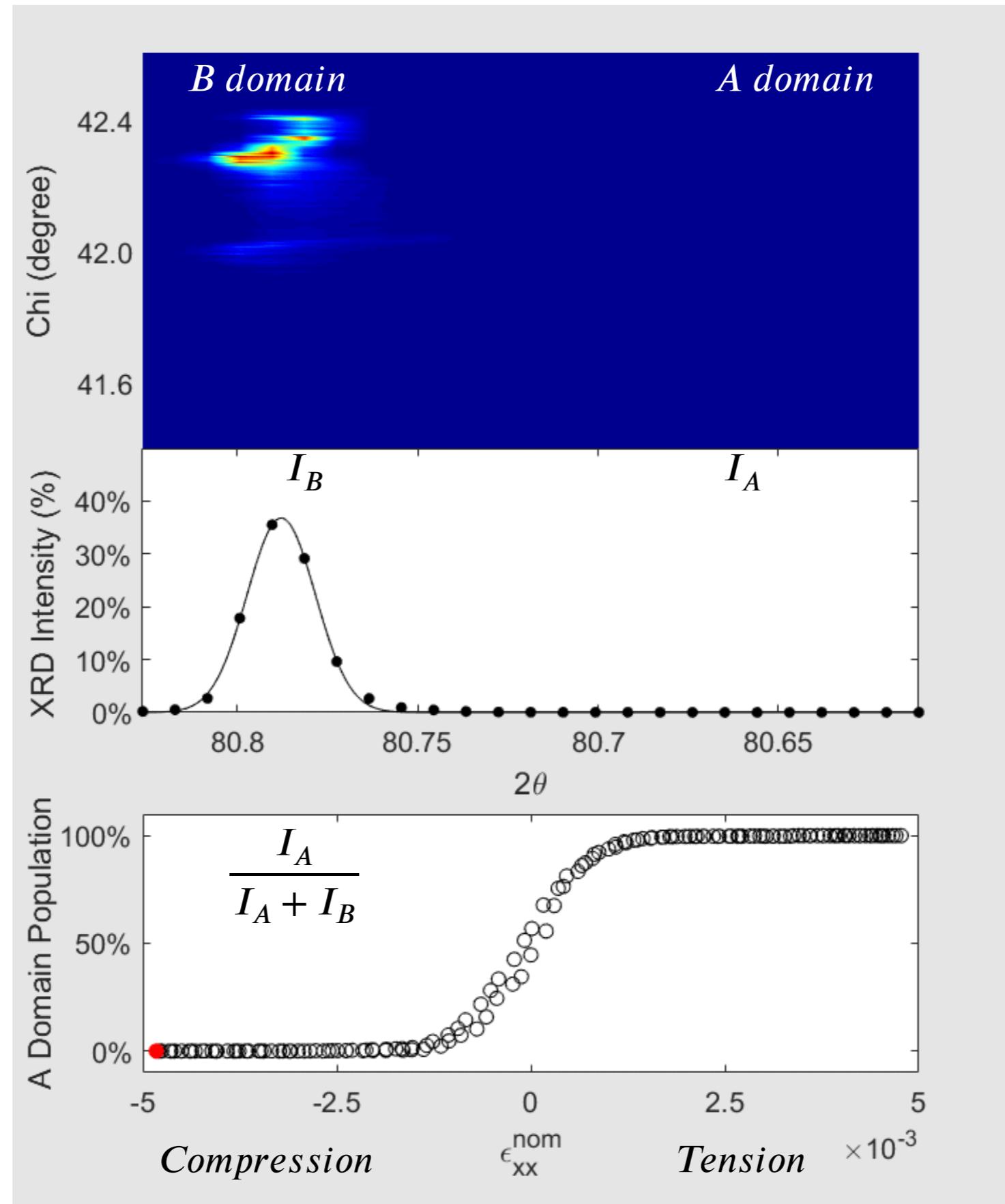


Shua Sanchez, Prof. Jiun-Haw Chu's Quantum Materials group

- Project goal: apply strains and x-ray to precisely detwin a sample and strain-tune superconductivity.
- This crystal has 2 structural domains (A and B) where iron atoms form rectangular lattice.
- Under zero stress, the A and B domains have the same total volume (domain population). Applying tension detwins the crystal to turn B domains into A, and compression to A to B.
- We combined x-rays to measure the a and b lattice constants directly while applying strain.
- The video shows 162 strain states sequentially and the intensity of the x-ray diffraction on the area detector



- (Top plot) the intensity position on the detector gives the length values of a and b which change with strain
- (Middle plot) the intensity is summed vertically and fit to 2 Gaussians
- (Bottom plot) The relative intensity $\frac{I_A}{I_A + I_B}$ gives the relative A domain population which change vs strain.



Interesting result!
Lattice constants freeze in place during detwinning!

Implies that the domain pinning is much softer than the crystal lattice

Can smoothly detwin the sample from B to A and back

Introductions

Thinking about statistics

What is significance anyway?

What is the question?

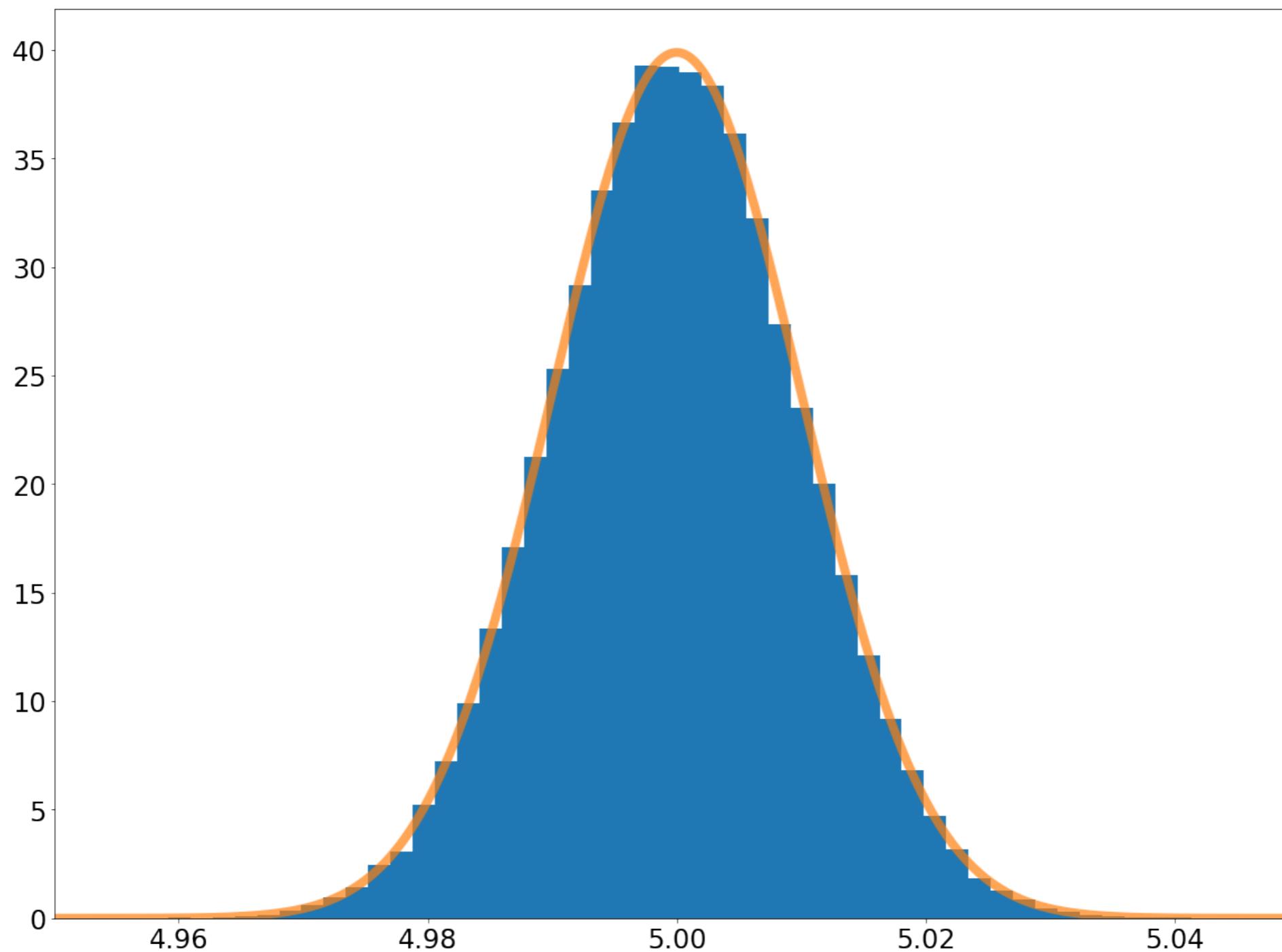
- Must clearly & precisely state the question

Example (null hypothesis):

If there is no signal; what is the probability that the background produces a signal that is equally or more signal-like than what I observed?

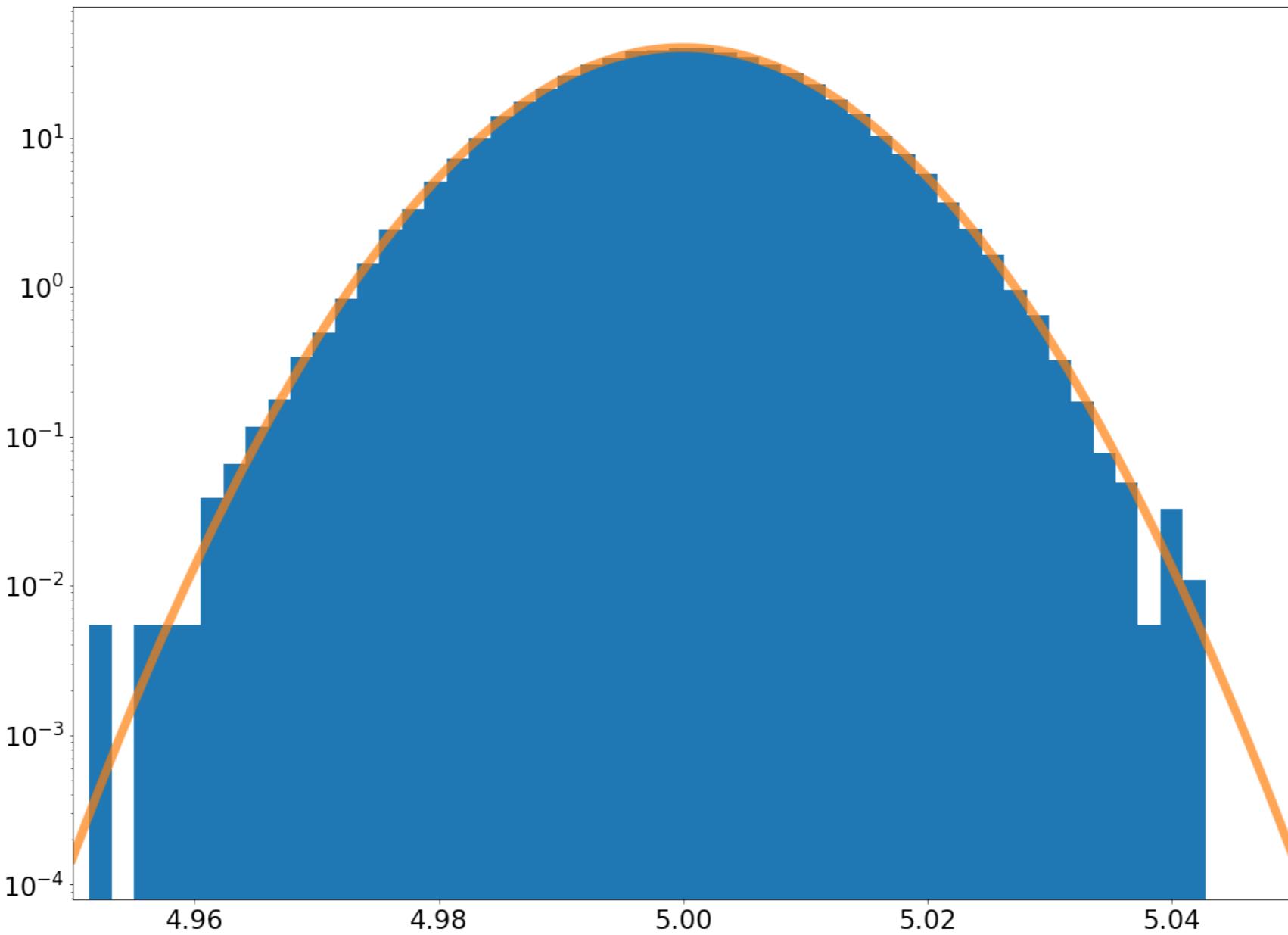
$$d = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

```
data = stats.norm.rvs(loc = 5., scale = 0.01, size = 100000)
```

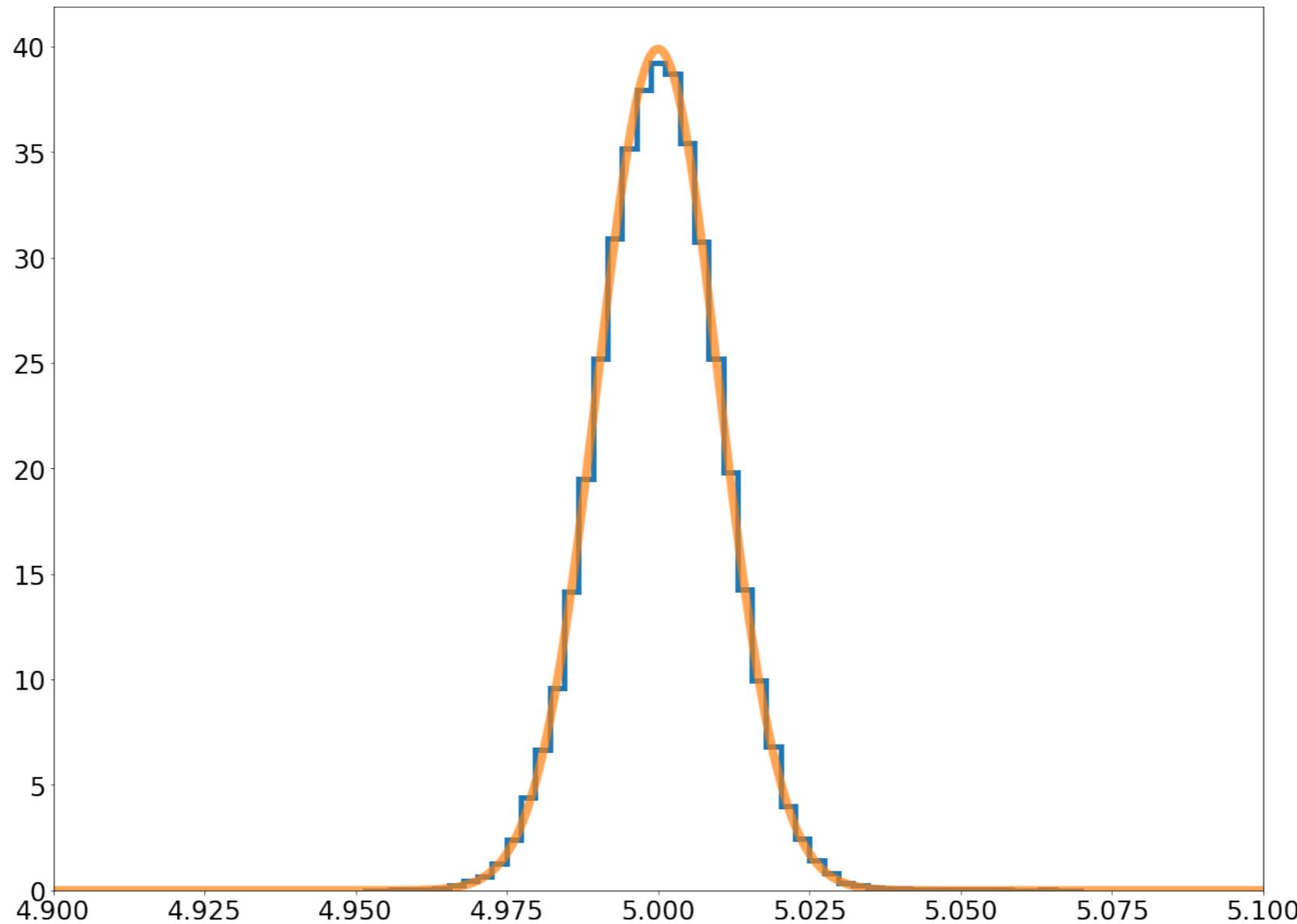


```
fig, ax = plt.subplots(1, 1)
ax.hist(data,50, density=True)
plt.yscale('log')
plt.tick_params(labelsize = 24)
plt.xlim([4.95,5.05])
x = np.linspace(4.95,5.05,1000)
ax.plot(x,stats.norm.pdf(x,loc = 5., scale = 0.01), linewidth = 8, alpha = 0.7)
plt.show()
```

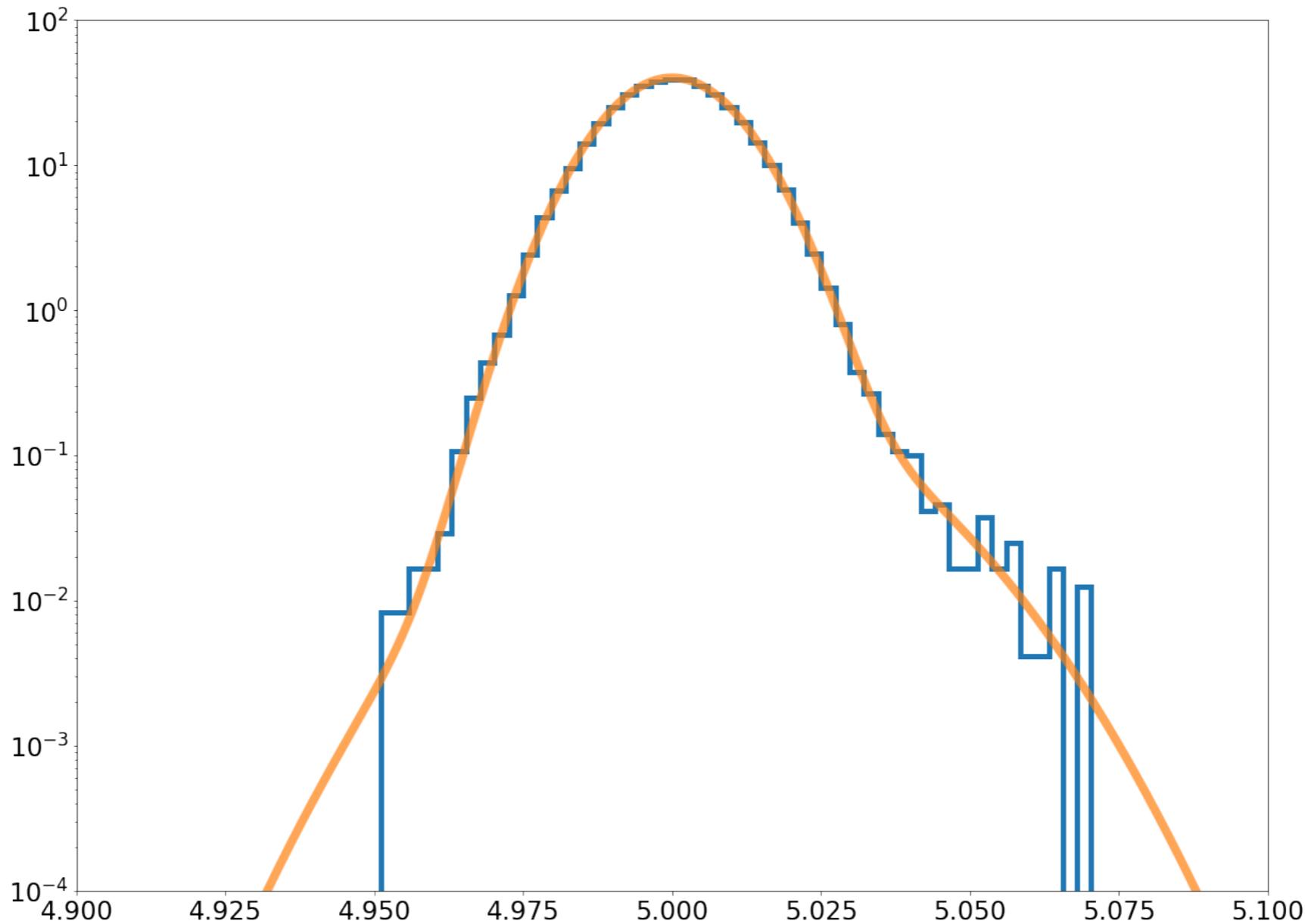
$$\log d = \frac{-(x - \mu)^2}{2\sigma^2} + c$$



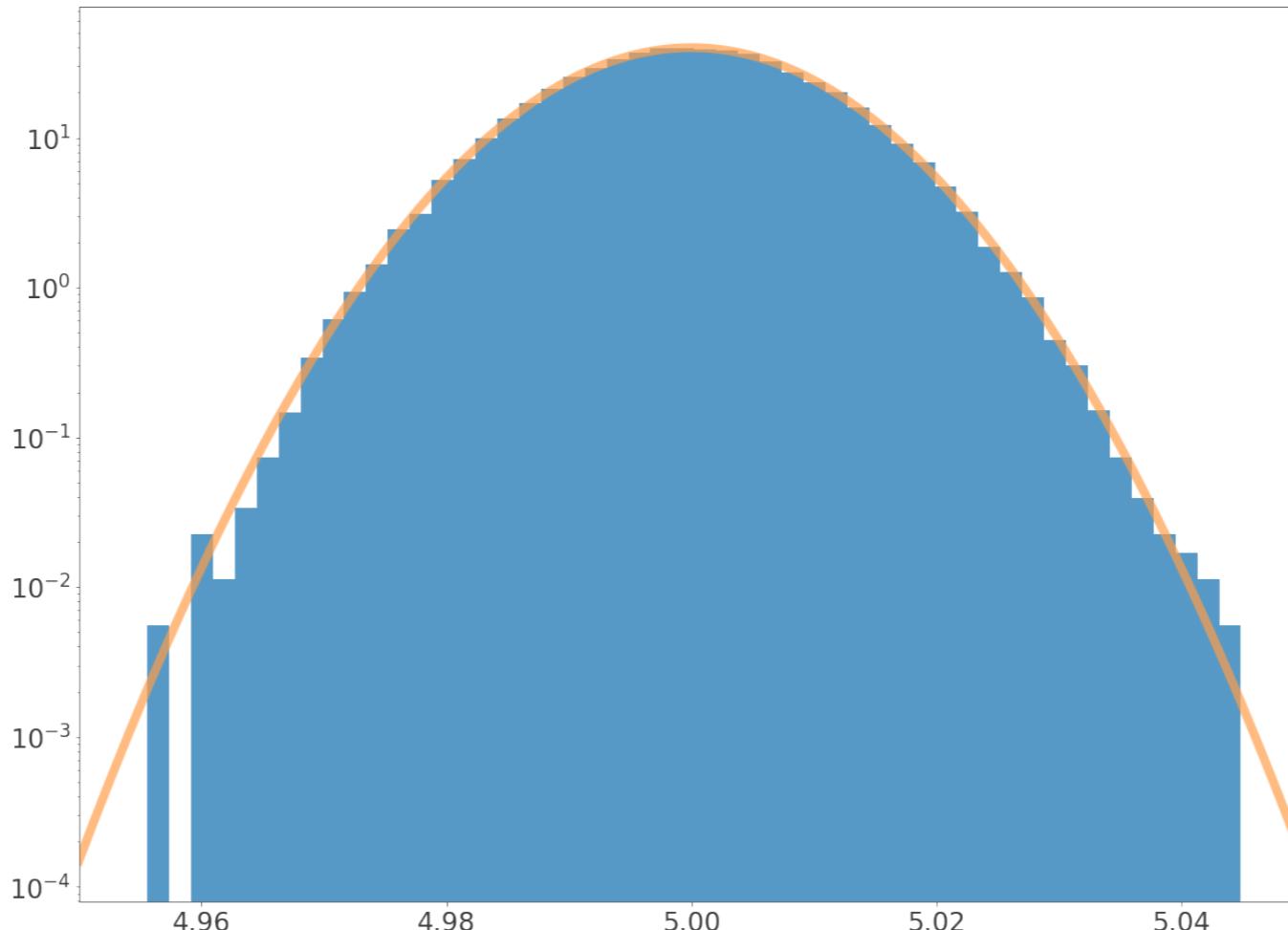
Never assume Gaussian statistics



Never assume Gaussian statistics



If there is no signal; what is the probability that the background produces a signal that is equally or more signal-like than what I observed?



$$\text{'Probability'} = \int_a^{\infty} \text{pdf}(x)dx$$

in this case!

In physics, $X\sigma$ is shorthand for a probability

- 5σ means: the probability of signal-free data giving a measurement that is equally or more signal like than your observation is less than 2.87×10^{-7} (or 1 in 3.5 million)

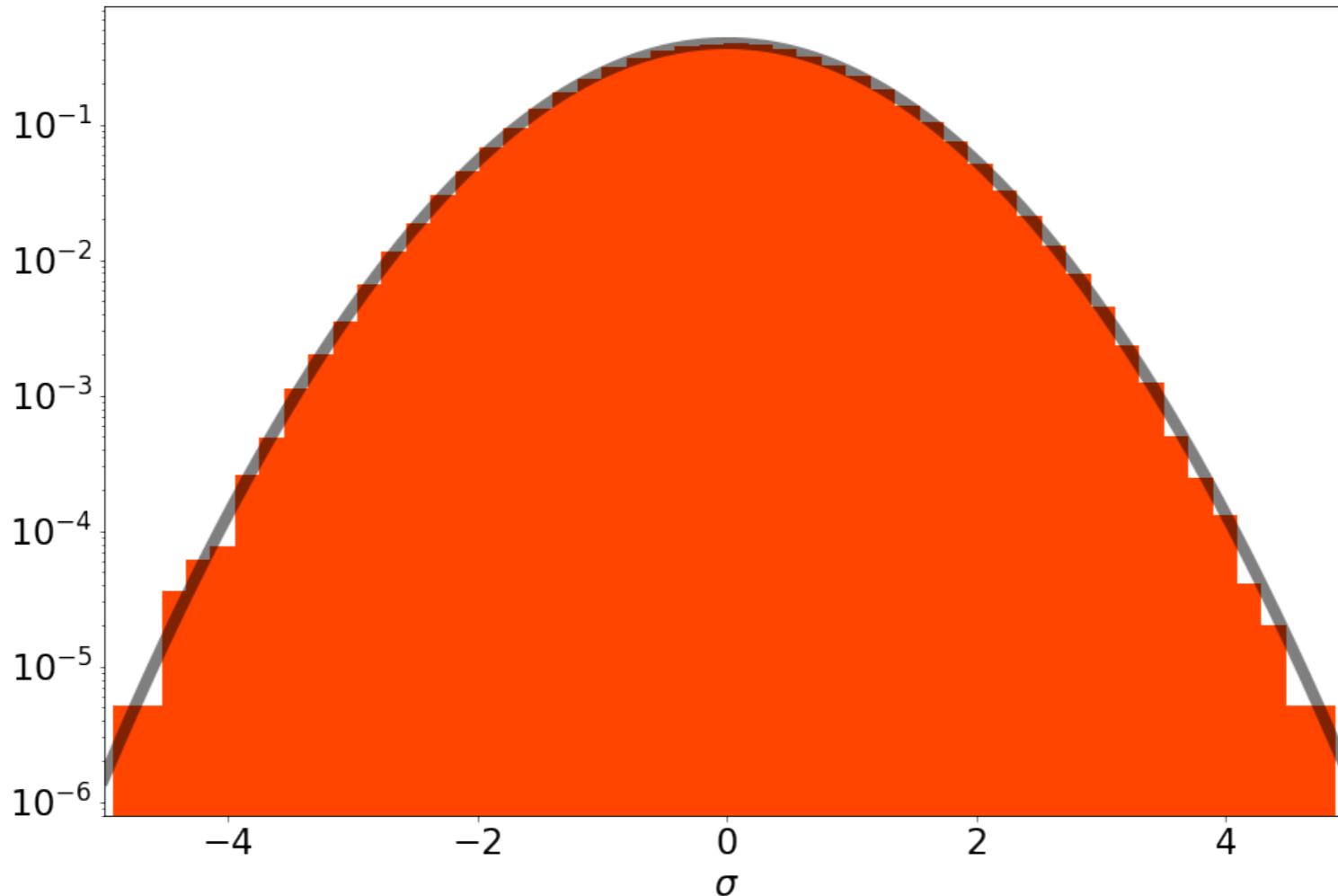
Common mistakes

- $X\sigma$ does not imply Gaussian distributed data
- $X\sigma$ is not $X\sigma$ away from the mean
- $X\sigma$ does not mean your question is one-sided

Best interpretation of $X\sigma$ (null hypothesis case)

- The probability of the background giving me a data point that looks as or more signal-like than the reading I have is the same probability as if my data was Gaussian and I was $X\sigma$ away from the mean

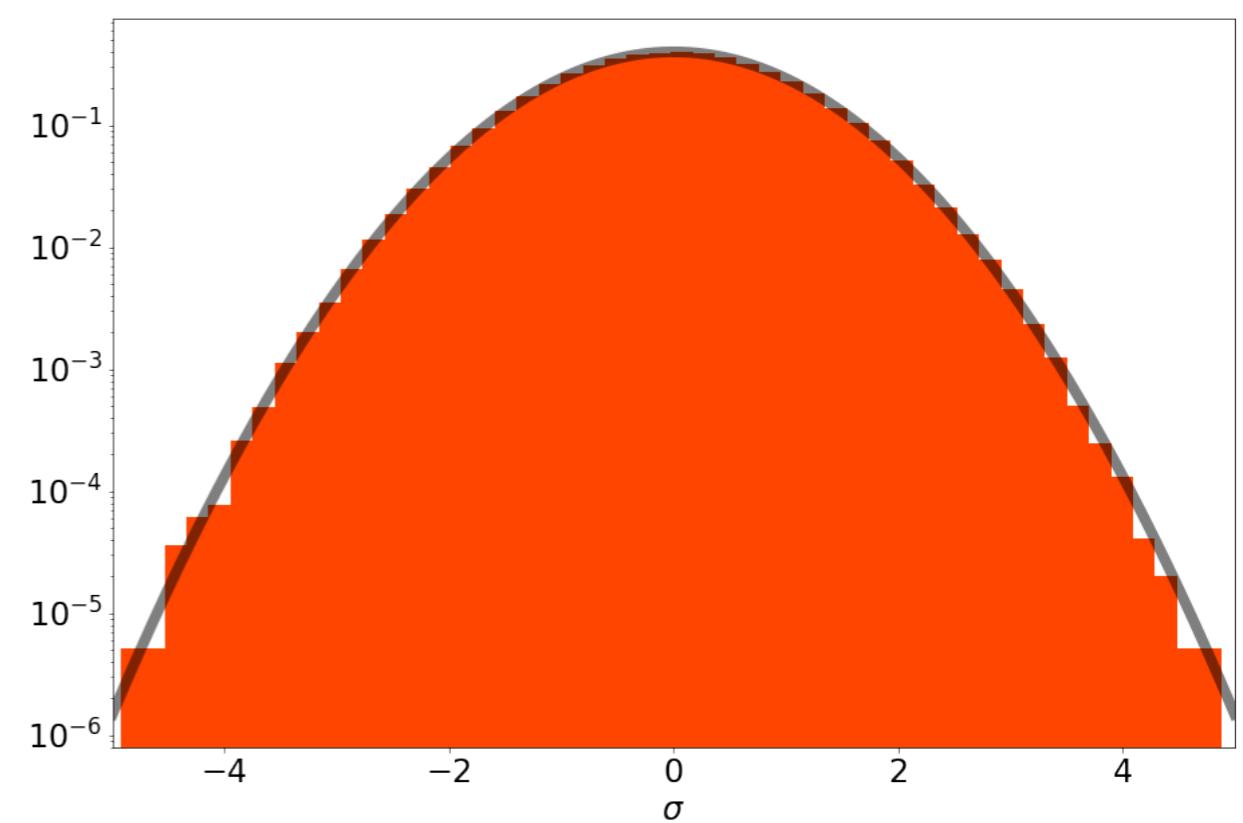
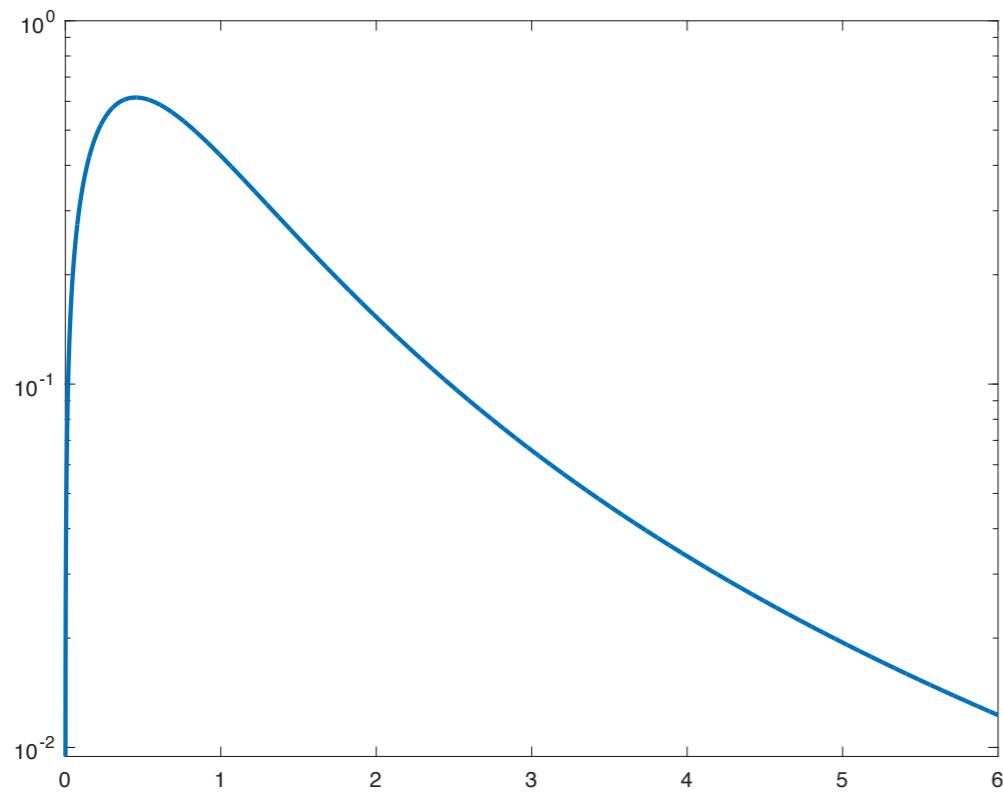
If there is no signal; what is the probability that the background produces a signal that is equally or more signal-like than what I observed?



$$\text{Probability } X\sigma = \int_{X\sigma}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}} dx = \frac{1}{2} \operatorname{erfc}\left(\frac{X}{\sqrt{2}}\right)$$

Best interpretation of $X\sigma$ (null hypothesis case)

- The probability of the background giving me a data point that looks as or more signal-like than the reading I have is the same probability as if my data was Gaussian and I was $X\sigma$ away from the mean



Key statistical steps

- Clearly state the question (& turn into math)
- Determine the background distribution
- Integrate background to find probability
- Convert probability into equivalent sigma

$$\text{'Probability'} = \int_a^{\infty} \text{pdf}(x)dx$$

