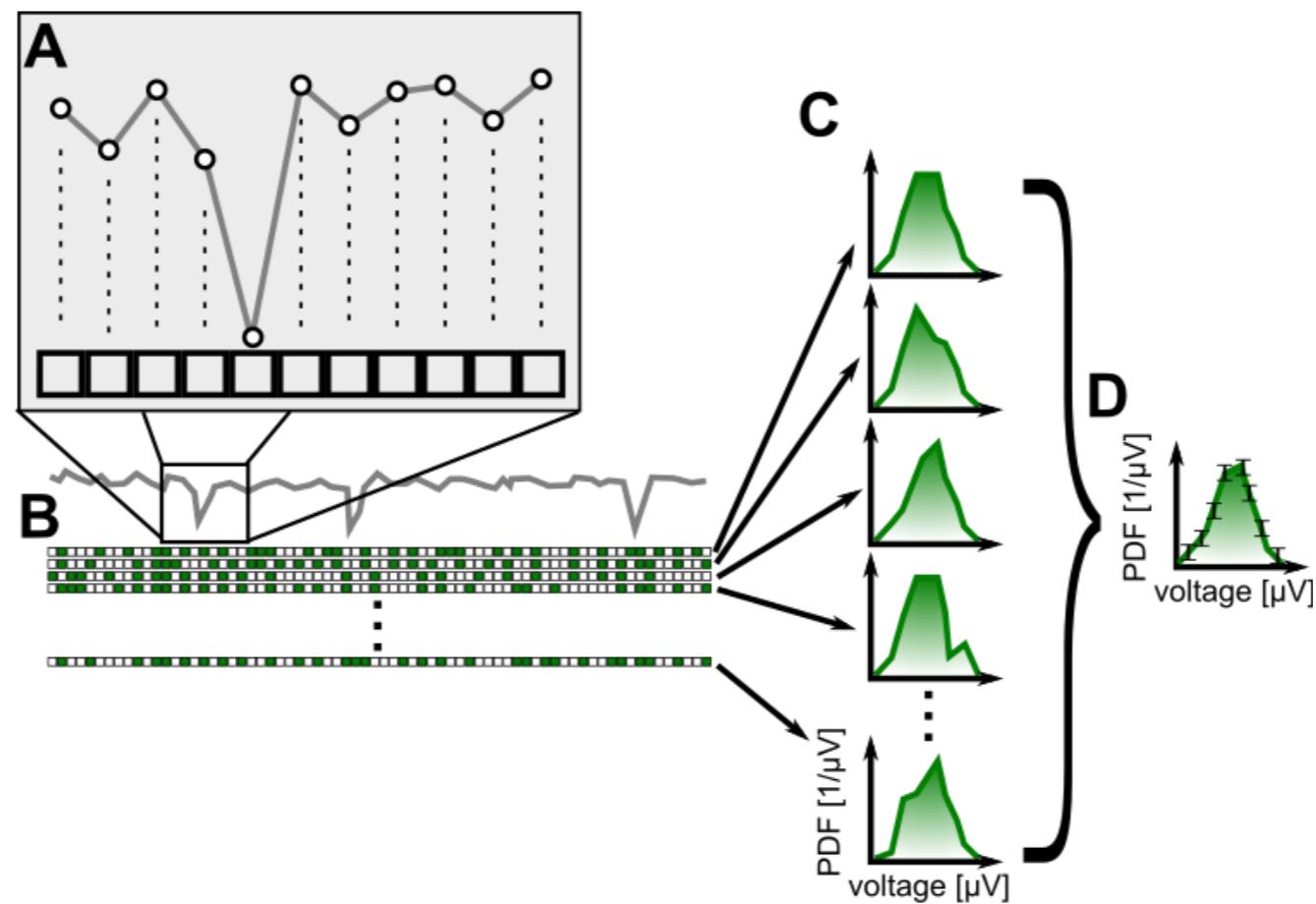


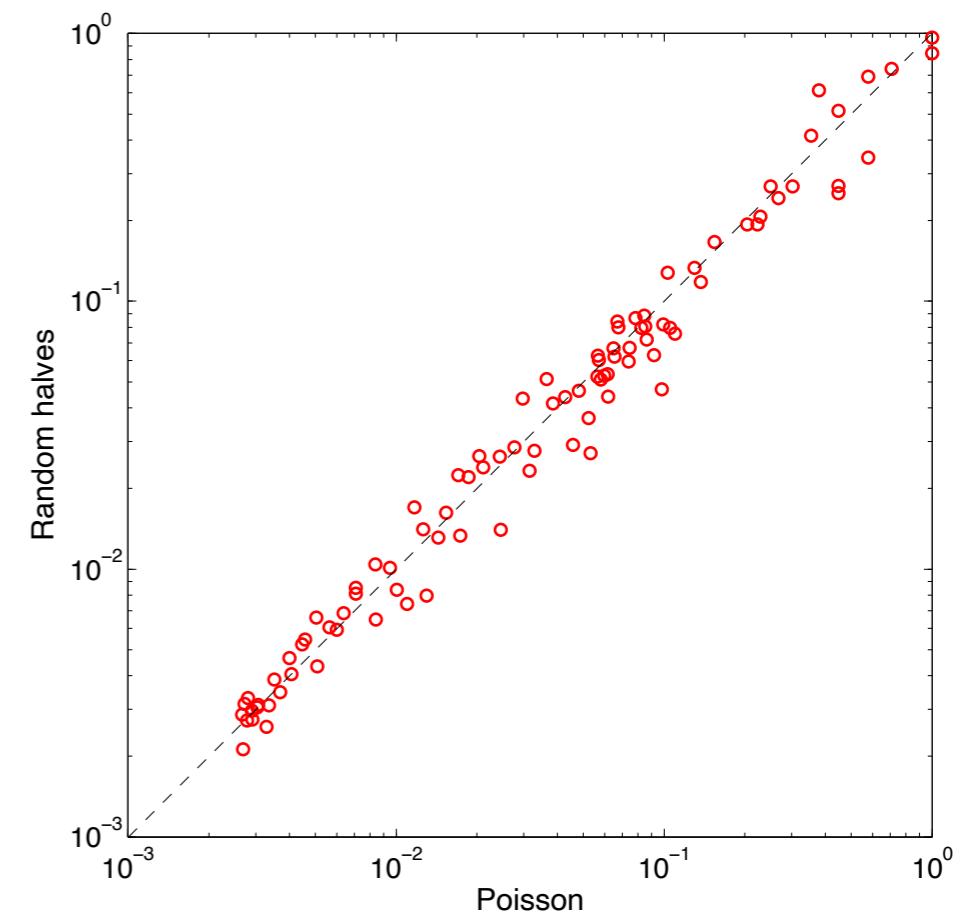
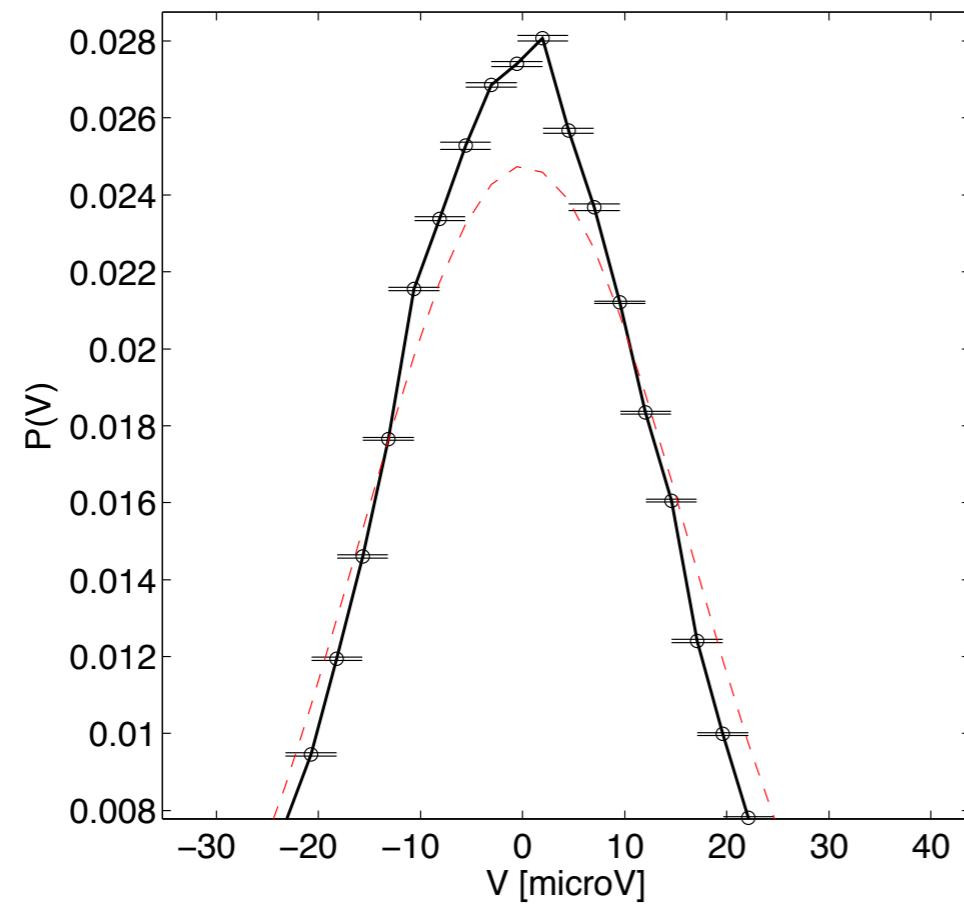
First resampling attempt at error bar estimation

- That's all good, but in normalizing the histogram, we lost the information about errors implicit in raw counts...
- Rough estimate: evaluate the histograms over random halves of the samples and take std over halves (over the same x-bins!)



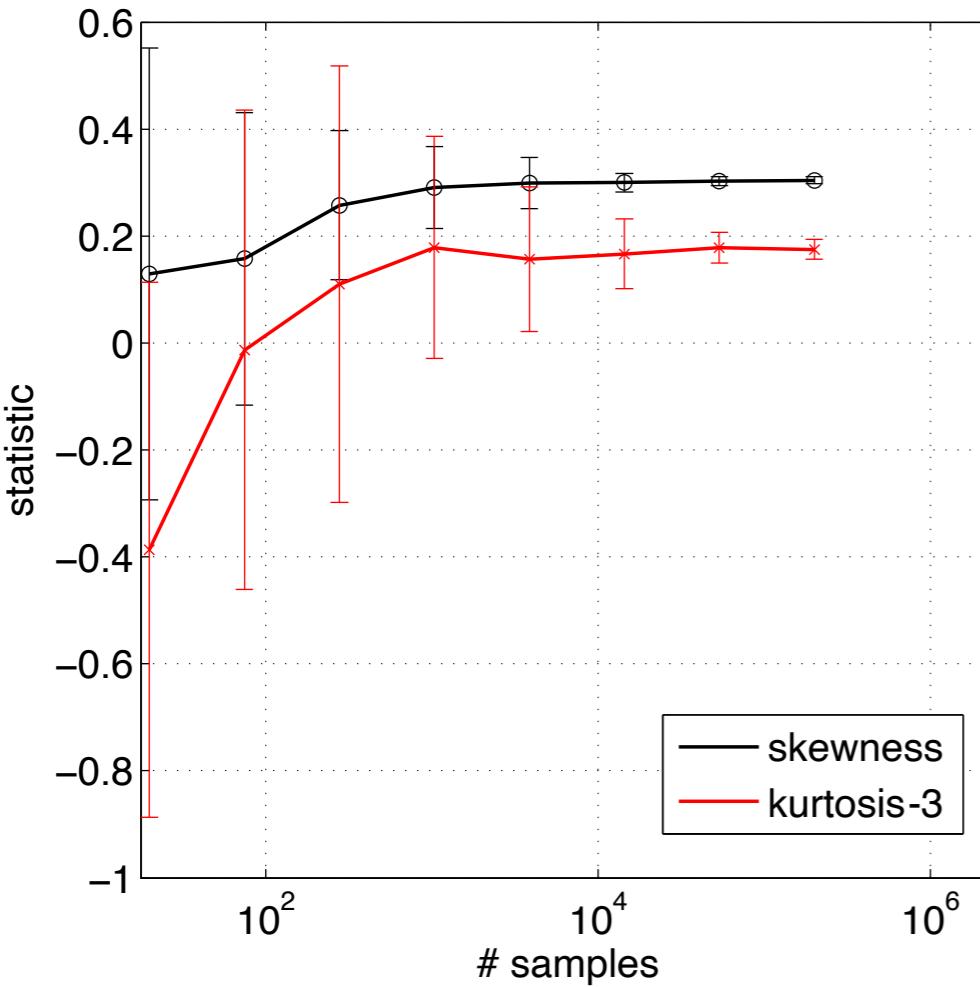
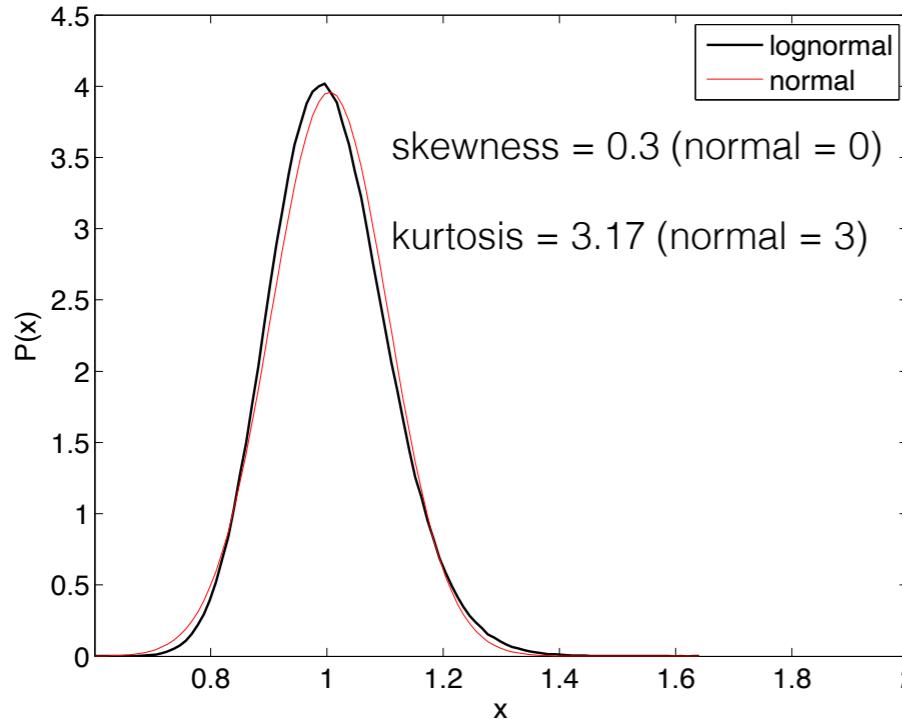
First resampling attempt at error bar estimation

- Looks like the PDF is indeed significantly different from a Gaussian!
- What do you expect for the error bar as a function of the counts? Counting / Poisson errors: $\text{Var}[N_i] = N_i$
- Let's say the error bar is $\sim \sqrt{\text{counts}}$, does this match with the rough bootstrap estimate?



Do you think this procedure for error bar estimation is actually correct for our data? (why Poisson? Why halves? Why independent errors? IID assumption?)

- Higher-order statistics can be very sensitive to sampling and outliers

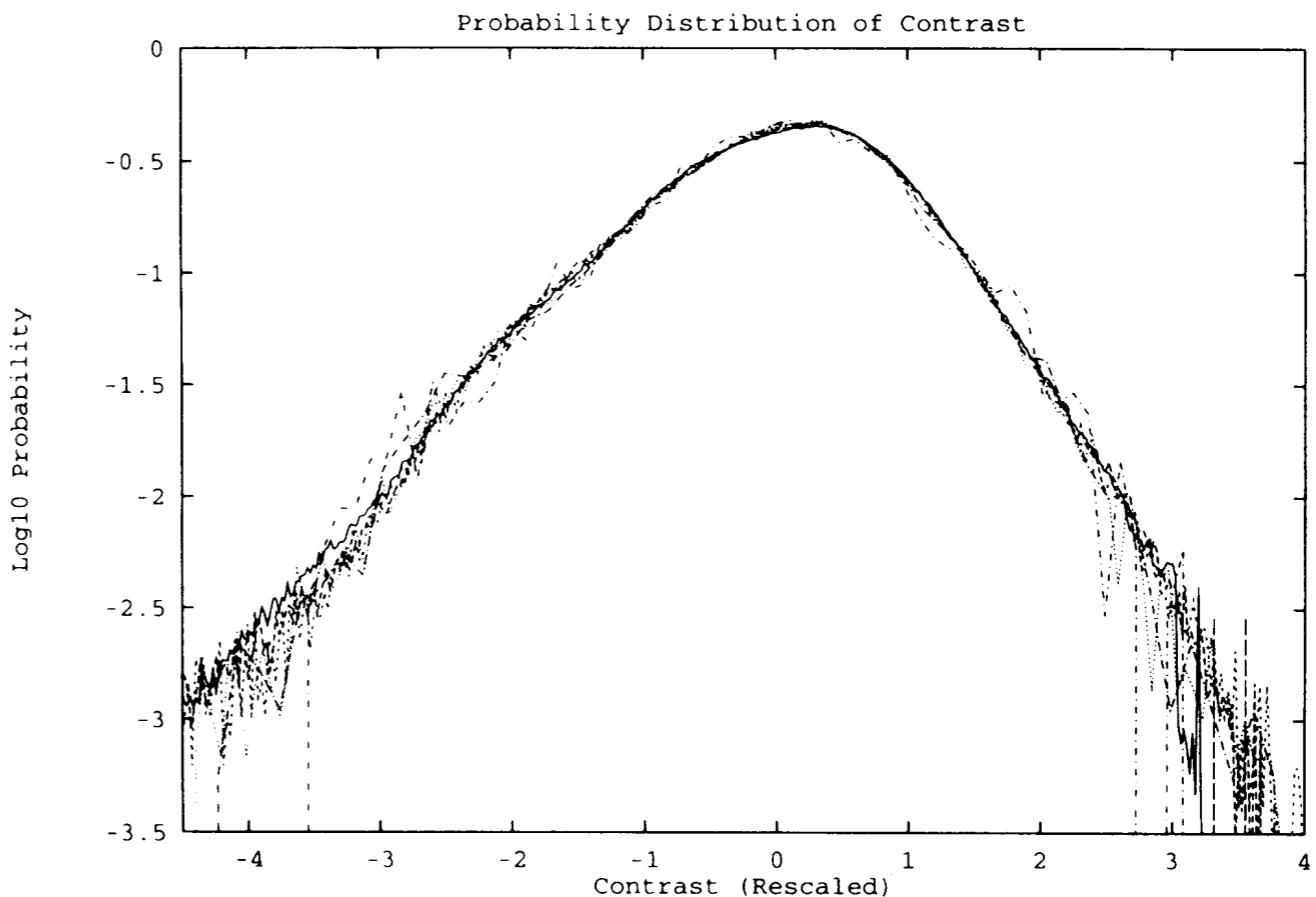
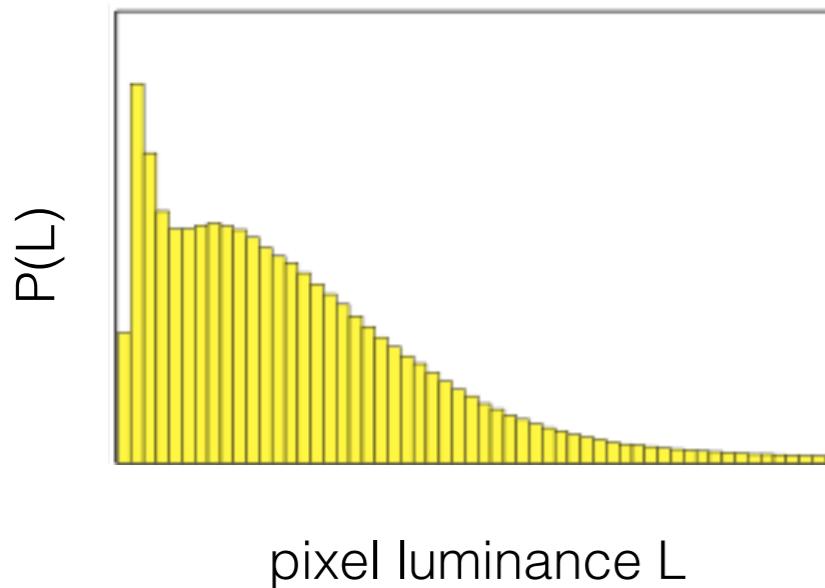


Quantile statistics are much more *robust*
 (decide whether to use moments or quantiles
 based on the amount of data + field)

... or try nonlinearly transforming the data?

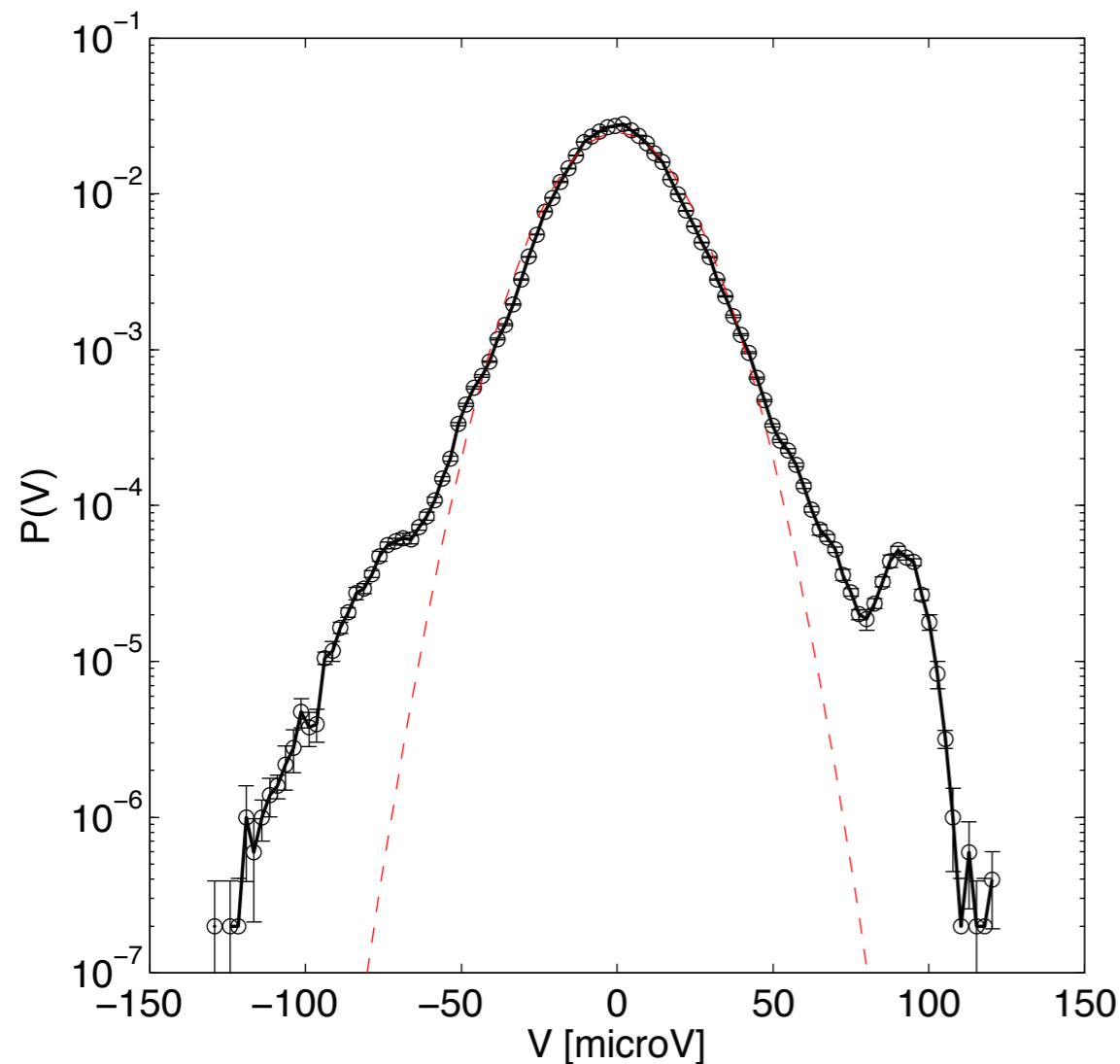
Excursion: data representation

- If you measure x , do you need to look at the histogram of x ? Why (not) $f(x)$?
- e.g. distribution of luminance over natural images (Ruderman et al, PRL 73 (1994)) which is very broad...



- if you define $\mathbf{C} = \log(\mathbf{I} / \mathbf{I}_0)$ where \mathbf{I}_0 is the mean over each image: **(i)** the data has a simpler (bi-exponential) form, better described by central moments (mean, std, skew, kurtosis); and **(ii)** the data collapses
- log-units for the measured variable sometimes make more sense in the context of a particular theory (e.g., log concentration instead of concentration), so it makes sense to reparametrize
- copula: transforming the marginals into the desired (simple) form to study correlations

- Lets circle back to error bars in real data...

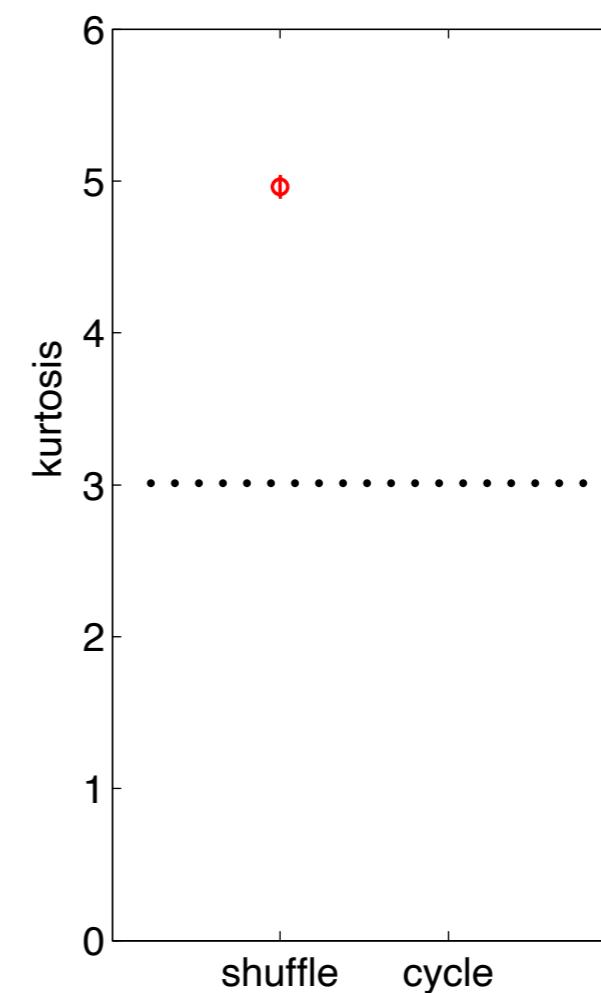


- A **huge difference in estimated error bar!**
- Non-IID samples due to temporal correlation?
- Estimate the effect: $\text{var}(\text{cycle})/\text{var}(\text{shuffle}) \sim 100$, 2M samples $\sim 20k$ IID samples?
Still not completely OK... depends on choice of data fractions etc.

skewness = 0.11 (normal = 0)

kurtosis = 4.97 (normal = 3)

- To see if this is significantly different from Gaussian, let's estimate over 1/20 (or 1/10) "fractions" of the data and get error bars
- What do we mean by fractions? Random shuffle draws, or contiguous blocks (e.g., by cycle-permuting data)?



Bootstrap re/sub-sampling for error estimation

When there is NO correlation in the data (IID samples)

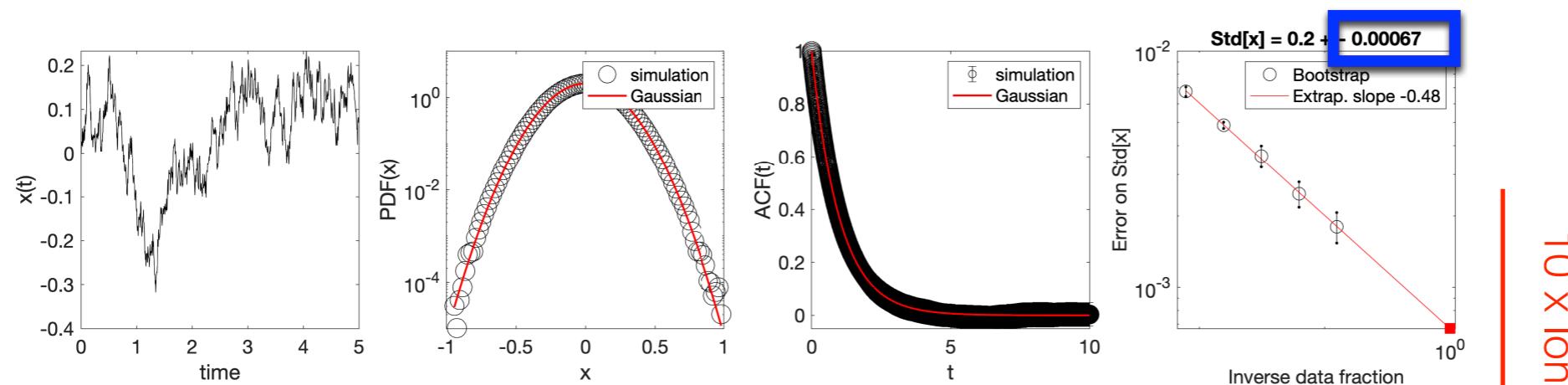
- **Basic idea:** given limited data, resampling from empirical distribution approximates sampling from a true (but unknown) underlying distribution
1. Take the empirical distribution based on N samples
 2. Resample (with replacement) / noise from the same distribution many times N samples
 3. Evaluate the statistic of interest and its error bar as std over resampling draws

When there are serial, short range correlations in the data

- **Basic idea:** subsample without replacement from *contiguous subsets* of data to preserve the correlation (assumptions?)
1. Split the data into blocks of different sizes (many times)
 2. Evaluate the statistic on each block and get its std across blocks
 3. Extrapolate the std in the statistic from different subset sizes to the whole data size

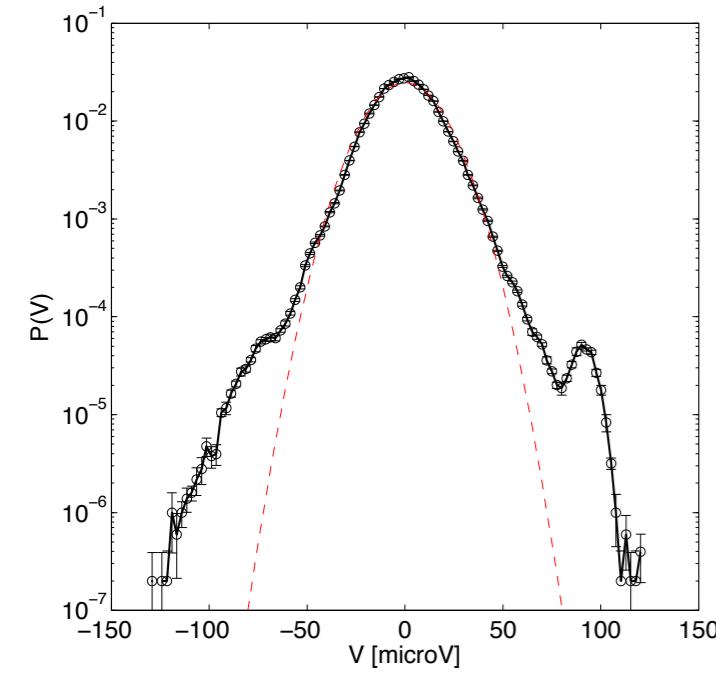
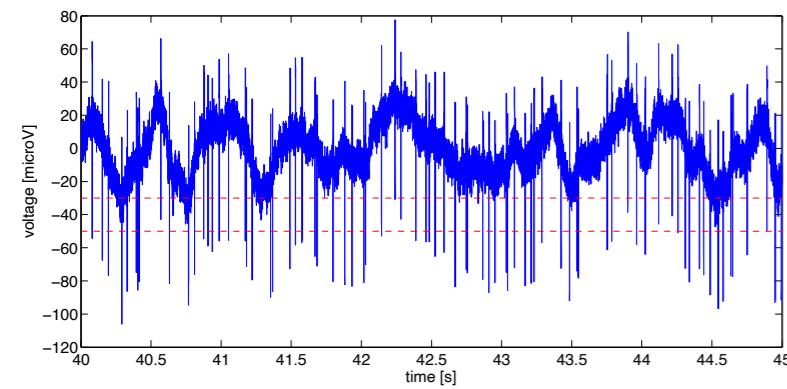
Need stationary data (second “I” in “IID”): statistics do not change with time!

Bootstrap on synthetic data - Ornstein-Uhlenbeck

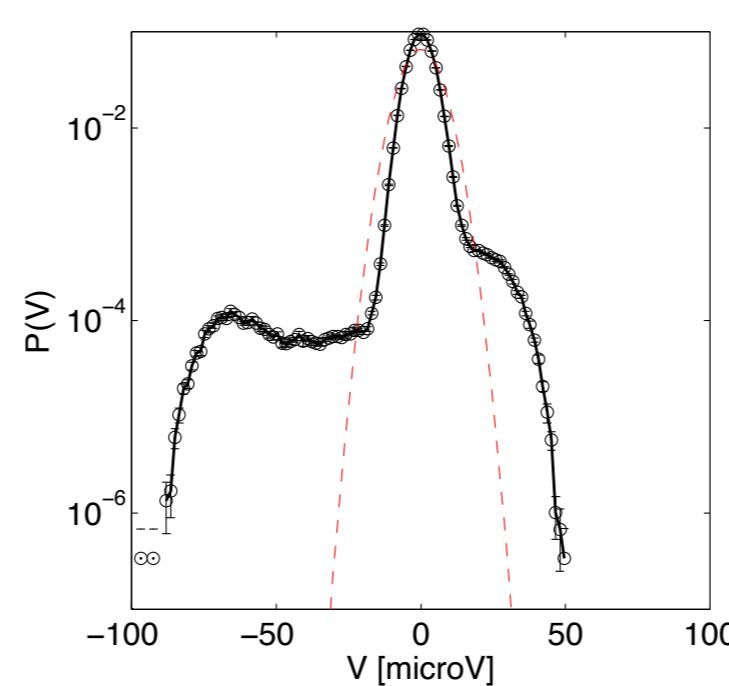
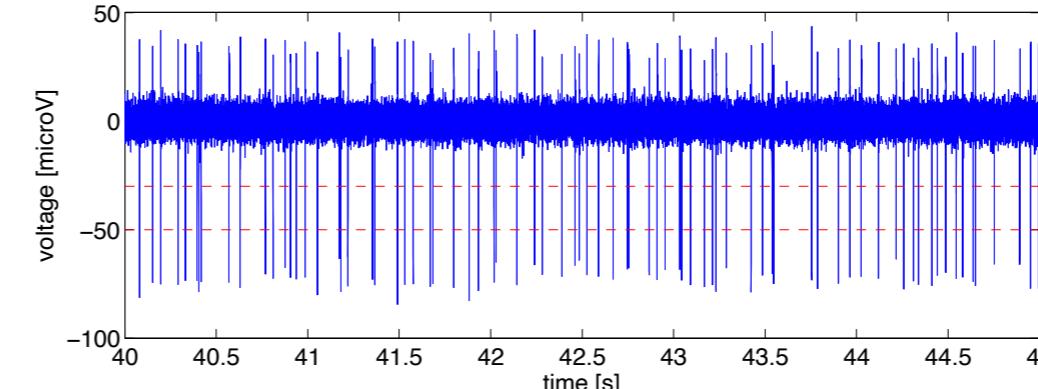


10 x longer correlation

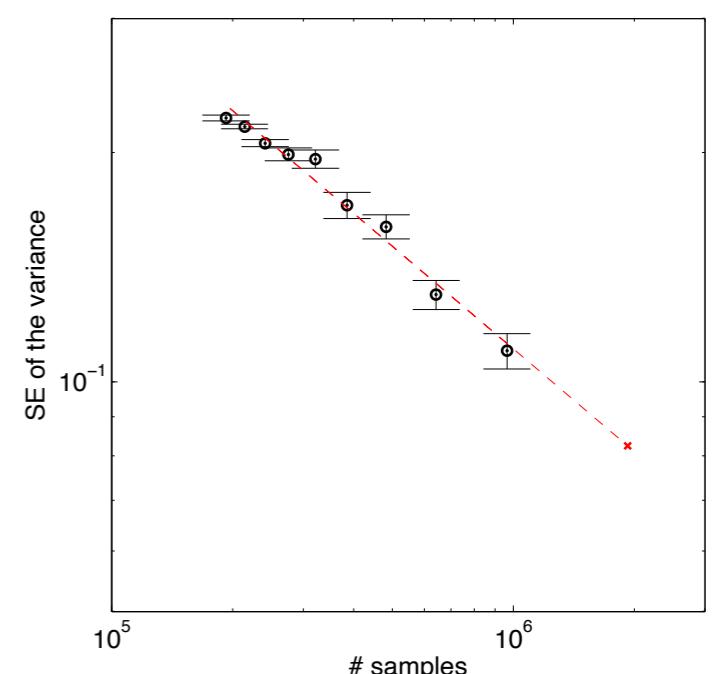
10 x shorter correlation



Original distribution



Baseline subtracted



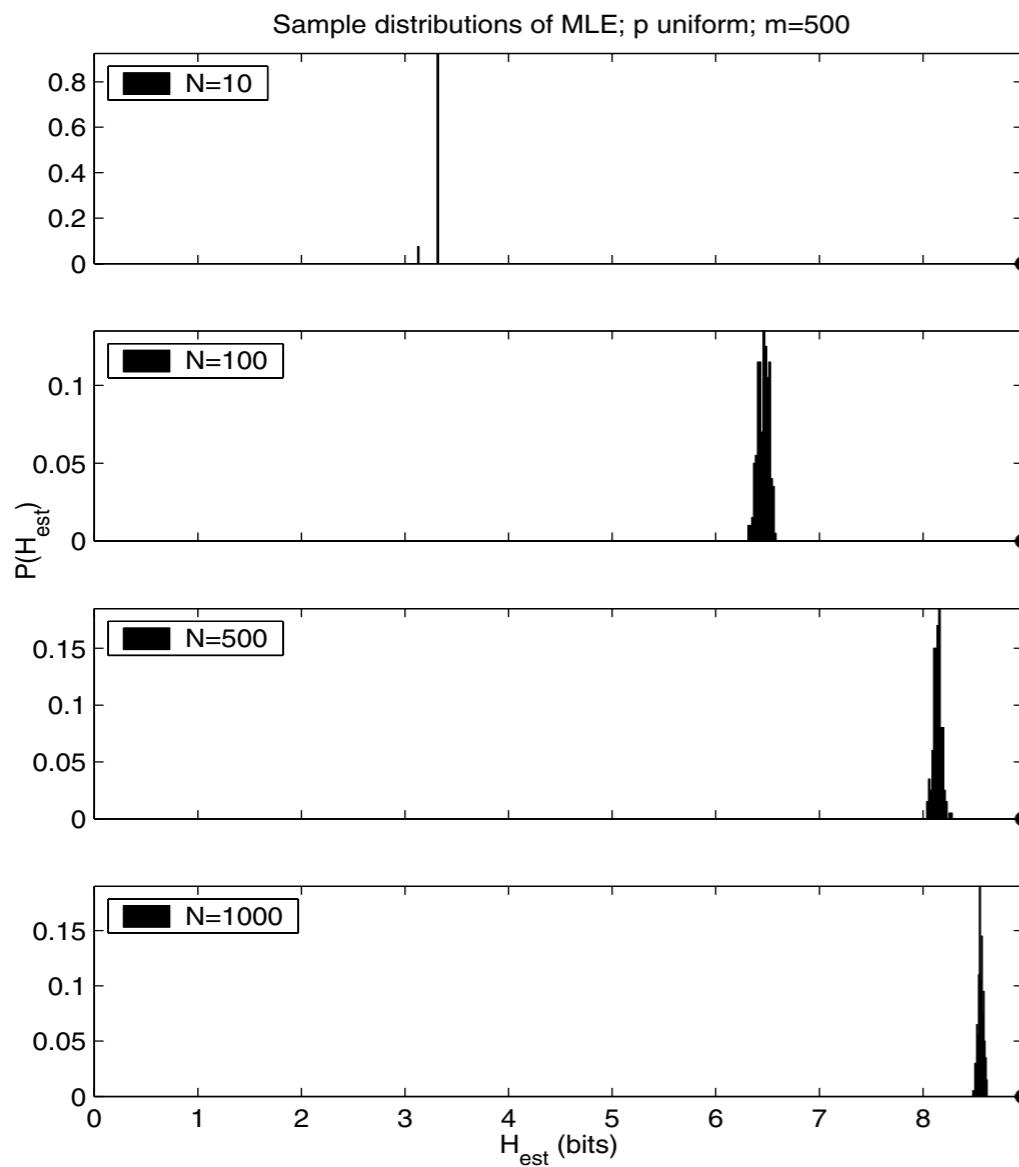
Let's use bootstrapping to estimate the variance
+- SE of the spike-subtracted signal...

$$\text{Variance} = 16.6 \pm 0.09 \text{ } \mu\text{V}^2$$

(+- 0.01 w/o correlations)

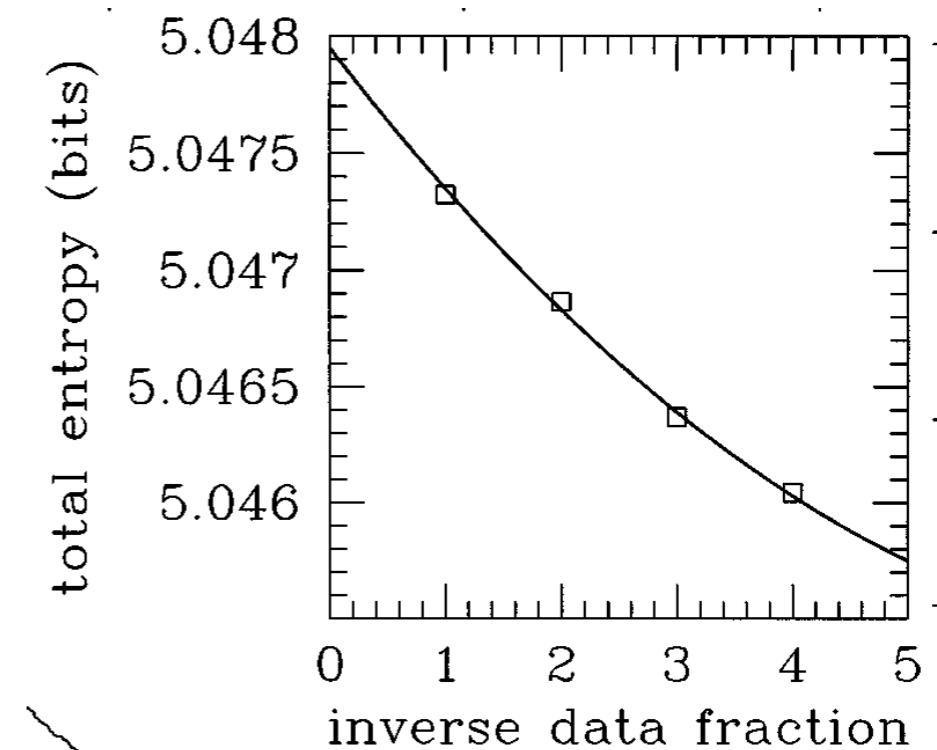
Excursion: using bootstrap for de-biasing

- **Task:** estimate the entropy of a discrete distribution given a finite number of samples
- **Problem:** maximum likelihood estimate (“naive” estimate / histogram statistic) is very biased at finite N (bias dominates over variance unless $N \gg m$)



$$S[p] = - \sum_i p_i \log_2 p_i$$

$$S(N) = S_{\text{true}} + S_1/N + S_2/N^2$$



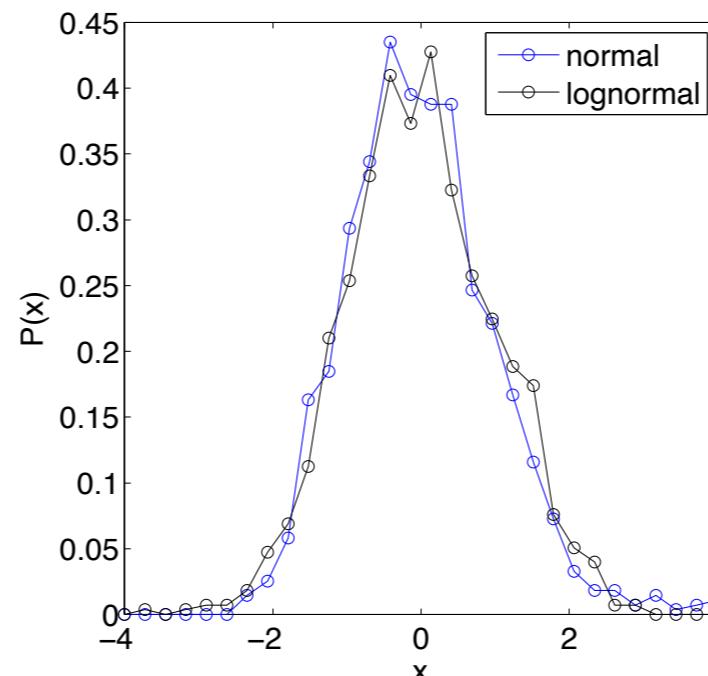
Strong SP et al, *Phys Rev Lett* 80 (1998)

Subsample data and extrapolate the bias to infinite data limit.

Resampling methods summary

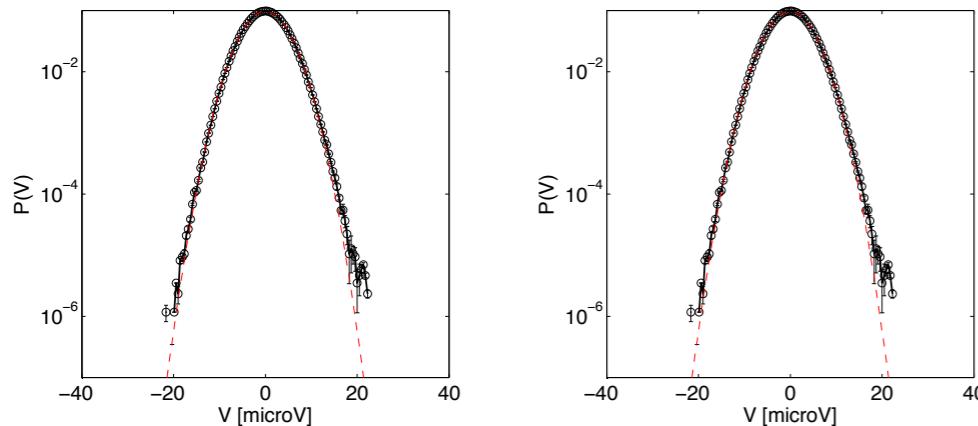
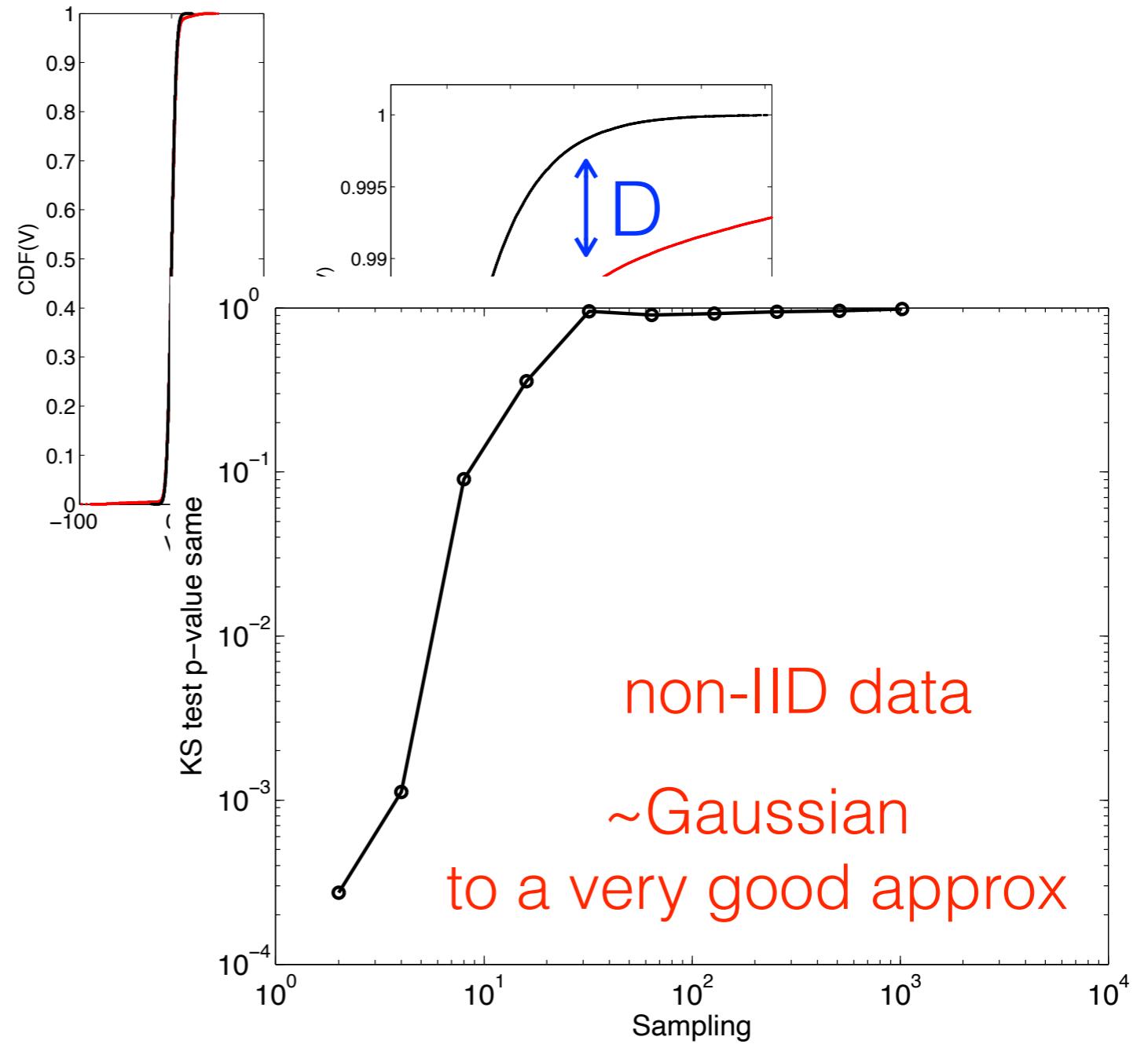
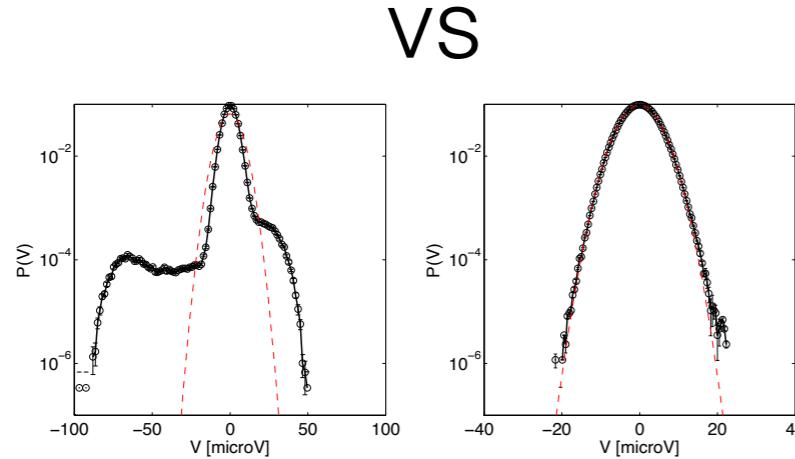
- Bootstrap resampling for IID data
- Bootstrap subsampling for correlated data
 1. Compute the statistic f on all leave-one-out subsamples of the data
- Jackknife
 2. $\sigma_f^2 = \frac{N-1}{N} \sum_{t=1}^N (f_{\mathcal{D} \setminus t} - f_{\mathcal{D}})^2$
- Debiasing estimators
- Creating null distributions to test for statistical significance (permutation tests)

CPU time is
cheap, permute
away (carefully)!



Are two distributions equal? Continuous data

- nonparametric Kolmogorov-Smirnov test for 1D distributions, reparametrization invariant
- $P(D > \text{observed} \mid \text{distributions same}) = \text{universal function of } \sqrt{N} \text{ and } D$
- not equally sensitive at all values, but generalizations exist



Are two distributions equal? Binned data

- Suppose we are given integer counts N_i in a binned histogram. Let's compare to expected number of counts (a model, or a reference distribution), n_i .
- Since the expected errors are $\sim \text{sqrt}(n_i)$ if the data is IID, we can define the chi2 statistic

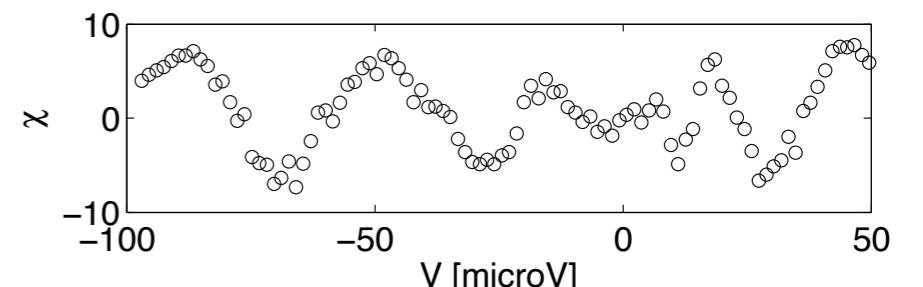
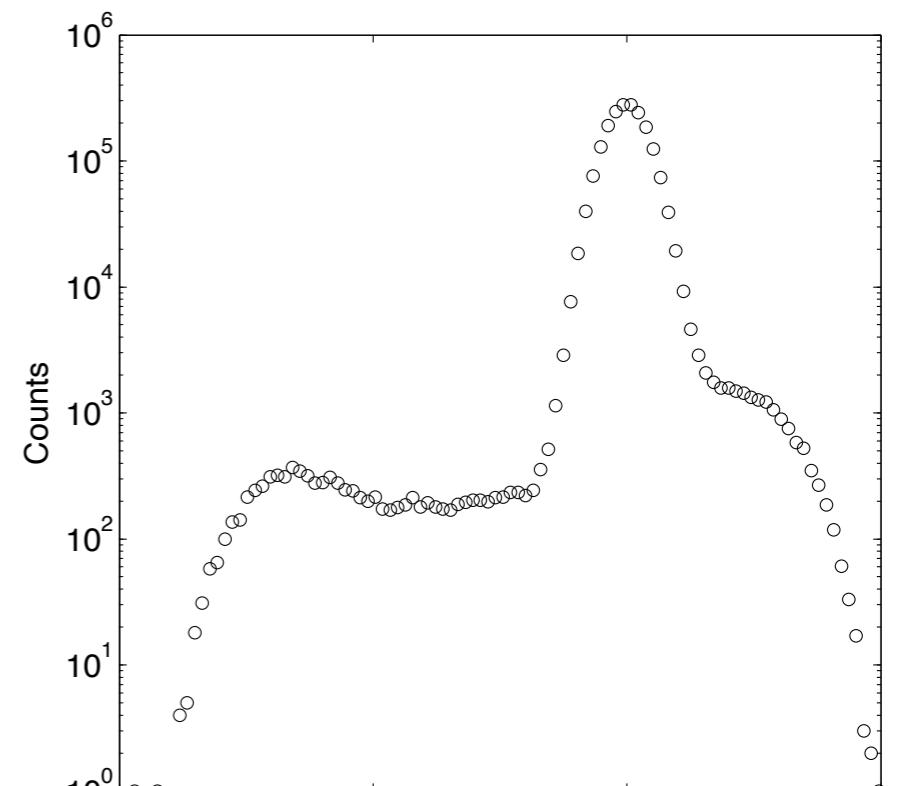
$$\chi^2 = \sum_i \frac{(N_i - n_i)^2}{n_i}$$

- Significance of the deviation is given by looking up getting the p-value from the chi2 distribution with k degrees of freedom

$$P(\chi^2 | k) = \frac{1}{2^{k/2} \Gamma(k/2)} \chi^{k-2} e^{-\chi^2/2}$$

We can use this as a error metric to fit distribution models to histograms, e.g., using nonlinear least-squares

Fitting histograms vs fitting Gaussian mixture models?



Excursion: excess events at LHC and dead salmons

One can view this also as set of hypothesis tests for the localized deviation...

but be careful about significance in multiple comparisons!

