

Statistics and Inference in Astrophysics

Bayesian and frequentist inference

Probability theory

- We cannot directly measure/observe what we are interested in (think Ω , or “the formation of the Milky Way”)
- Connection between models and data is often statistical, and data has noise
- Need theory to express uncertain knowledge and to update it

Two definitions of “probability”

- Great schism between two definitions of probability:
 - Frequentist: Long-run relative frequency of occurrence of an event in repeated experiments.
E.g., $P(\text{heads}) = 0.5$ bc half of coin-tosses of ideal coin result in heads
 - Bayesian: Real-valued measure of the plausibility of a proposition, closely follows intuitive reasoning.
E.g., $P(\text{it will rain in 10 minutes}|\text{cloudy}) = 0.5$.

Likelihood

- The likelihood is a function both used in frequentist and Bayesian inference
- Essentially encodes how the data are produced by the model (e.g., straight line, intrinsic flux) and observing procedure (e.g., noise)
- Once model is fixed and observing procedure is known, *no freedom*
- Many desirable properties

Likelihood

- Abstract:

$L = p(\text{data} \mid \text{model}, \text{observing procedure}, \text{other necessary knowledge})$

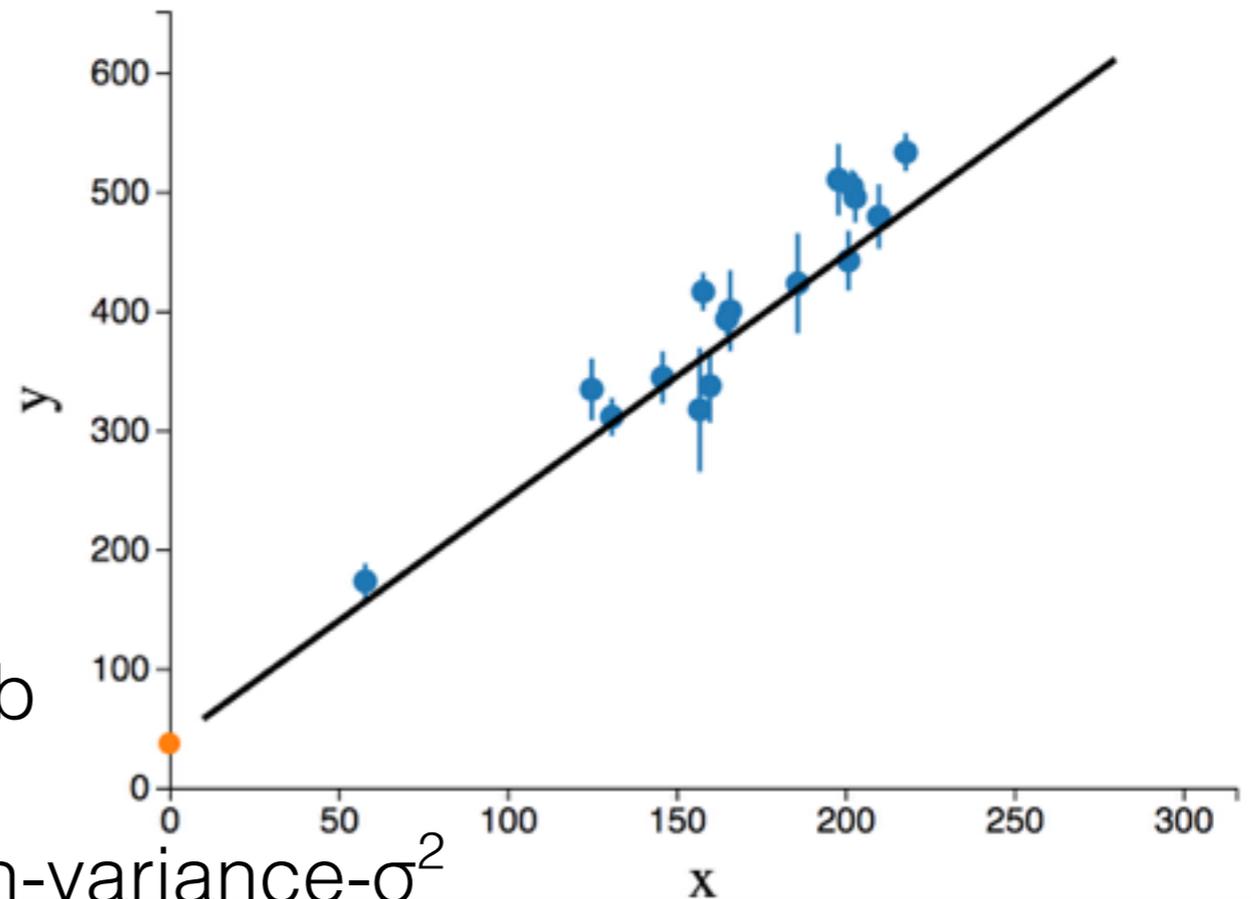
- Example: Straight line fit

- Given x : model $\longrightarrow y_{\text{true}} = mx + b$

- $y_{\text{obs}} = y_{\text{true}} + \text{Gaussian-noise-with-variance-}\sigma^2$

- $L = p(y_{\text{obs}} \mid \text{model}, x, \sigma) = p(y_{\text{obs}} \mid m, b, x, \sigma)$
 $= p(y_{\text{obs}} \mid y_{\text{true}} = mx + b, \sigma)$
 $= N(y_{\text{obs}} \mid y_{\text{true}} = mx + b, \sigma^2)$

- Or $-2 \ln L = (y_{\text{obs}} - [mx + b])^2 / \sigma^2 = \chi^2$



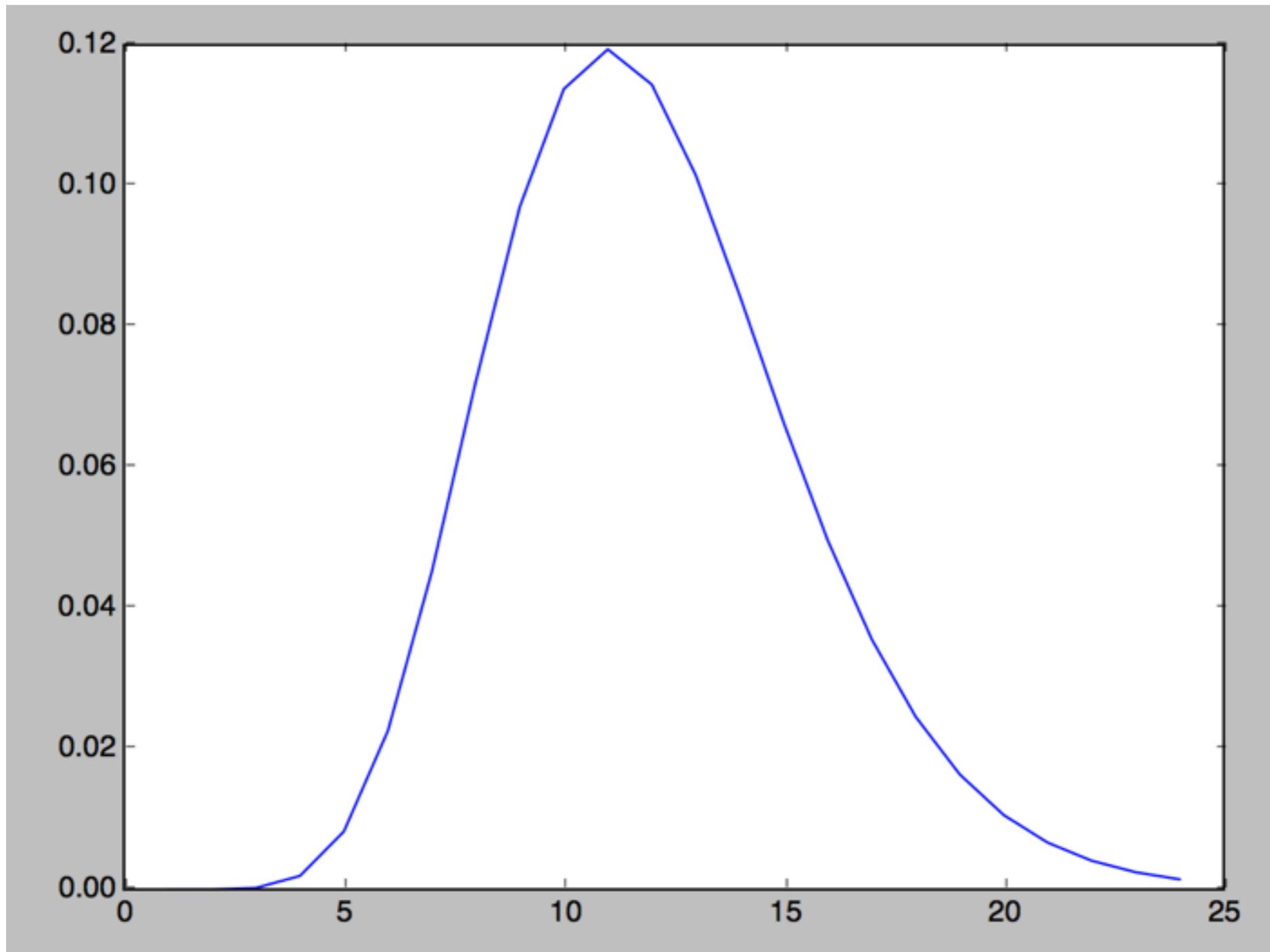
Likelihood

- Abstract:

$L = p(\text{data} \mid \text{model, observing procedure, other necessary knowledge})$

- Example: data = 11 photons, observed with dark noise equivalent to 1 photon
- $p(11 \text{ photons} \mid \text{model}=9 \text{ photons, dark}=1 \text{ photon})$
 $= \text{Poisson}(11 \mid \text{mean} = 9+1, \text{variance} = 9+1)$
- = 0.11373639611012128

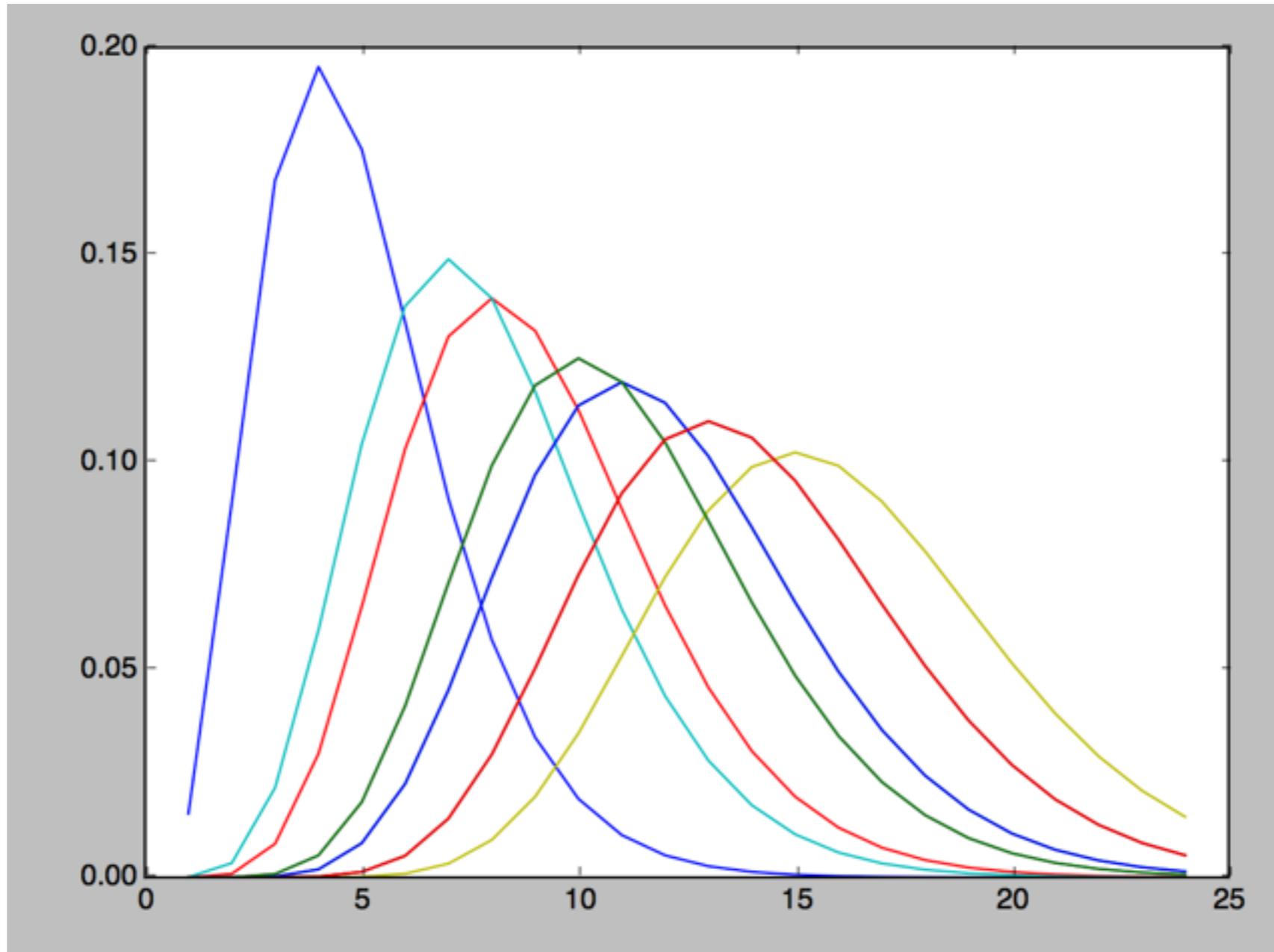
- $p(11 \text{ photons} \mid \text{model} = x-1 \text{ photons, dark} = 1 \text{ photon})$



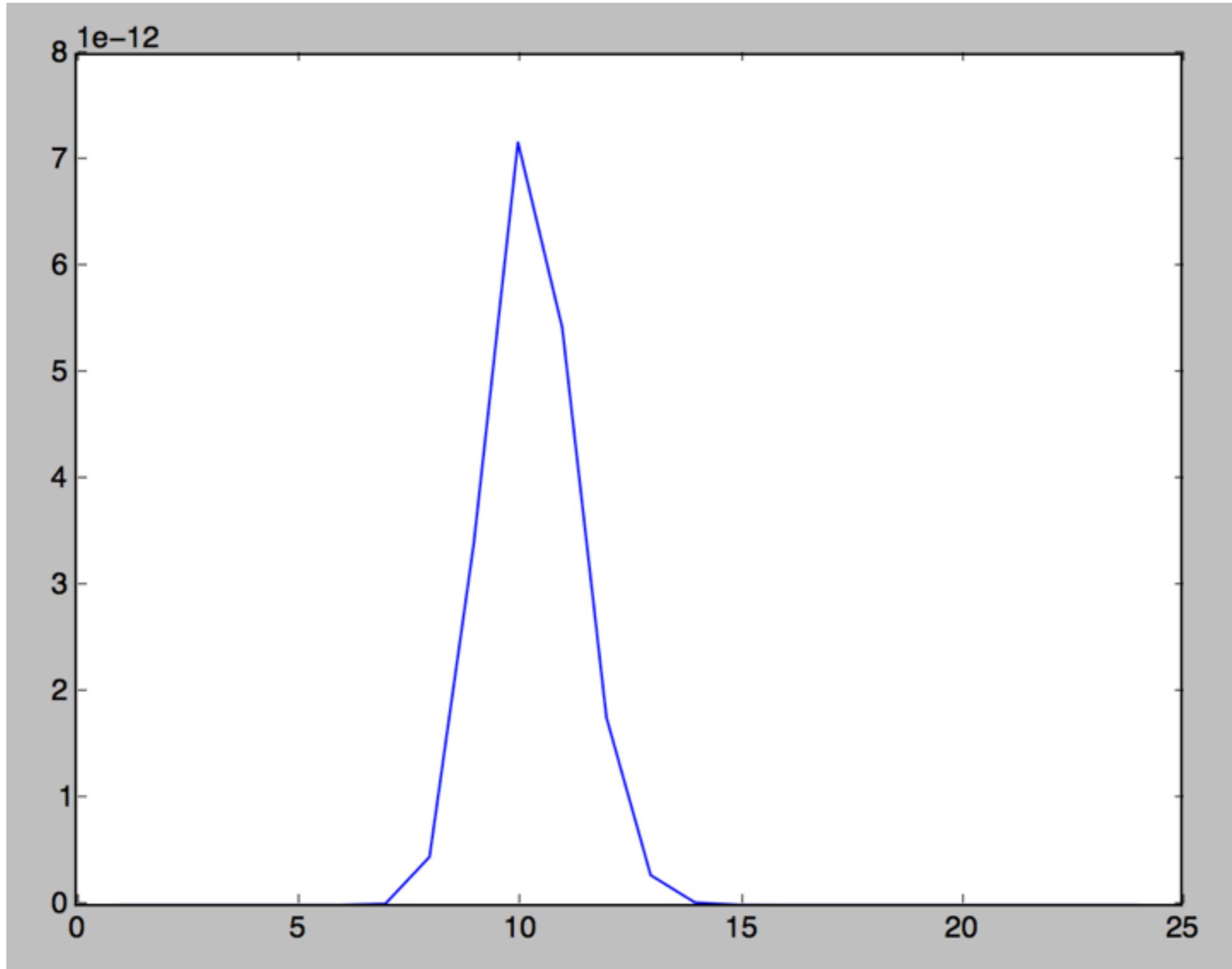
Likelihood

- For multiple data points:
- Suppose I observe the source 10 times, get {4, 11, 8, 7, 10, 15, 13, 11, 10, 13}
- Assume average model flux = 9 photons
- $L = \text{Poi}(4|10) \times \text{Poi}(11|10) \times \text{Poi}(8|10) \times \text{Poi}(7|10) \times \text{Poi}(10|10) \times \text{Poi}(15|10) \times \text{Poi}(13|10) \times \text{Poi}(11|10) \times \text{Poi}(10|10) \times \text{Poi}(13|10)$
- = 7.1695477633905203e-12
- Typically use $\ln L$!!

All individual likelihoods $Poi(obs|x)$



Product



Likelihood

- Assuming multiple measurements are independent, multiply together individual likelihoods:

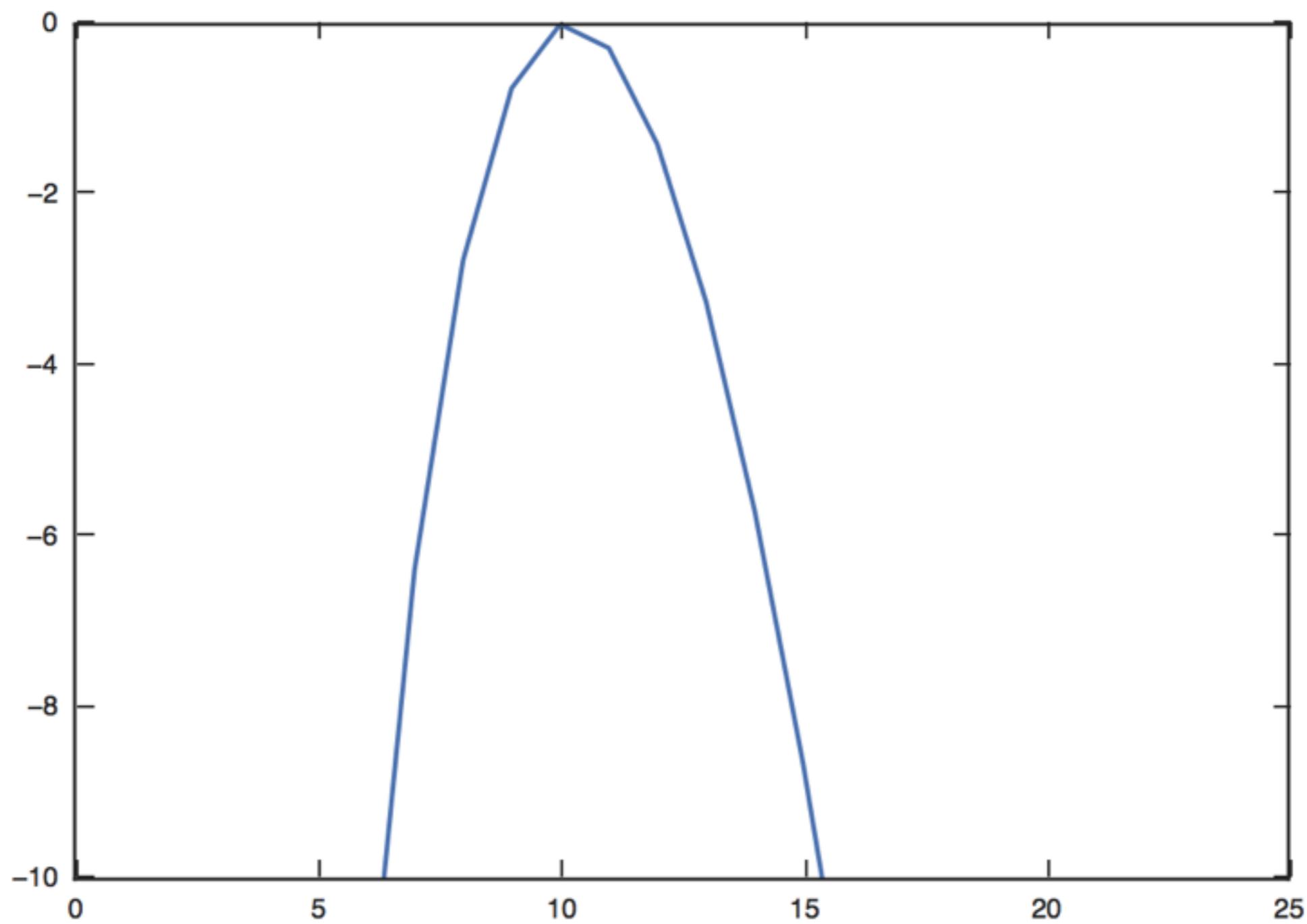
$$L = p(\text{data}_1|\text{model}) \times p(\text{data}_2|\text{model}) \times \dots \times p(\text{data}_N|\text{model})$$

- L completely determined by model and observing:
 - Photometry: intrinsic flux + dark noise + read noise \rightarrow Poisson / Gaussian for large counts (more than ~ 100)
 - Measurements of constant A with Gaussian noise $s \rightarrow$ Gaussian with mean= A , noise= s
 - Model: Velocity distribution with mean A and velocity dispersion $s \rightarrow$ Gaussian with mean= A , noise= s

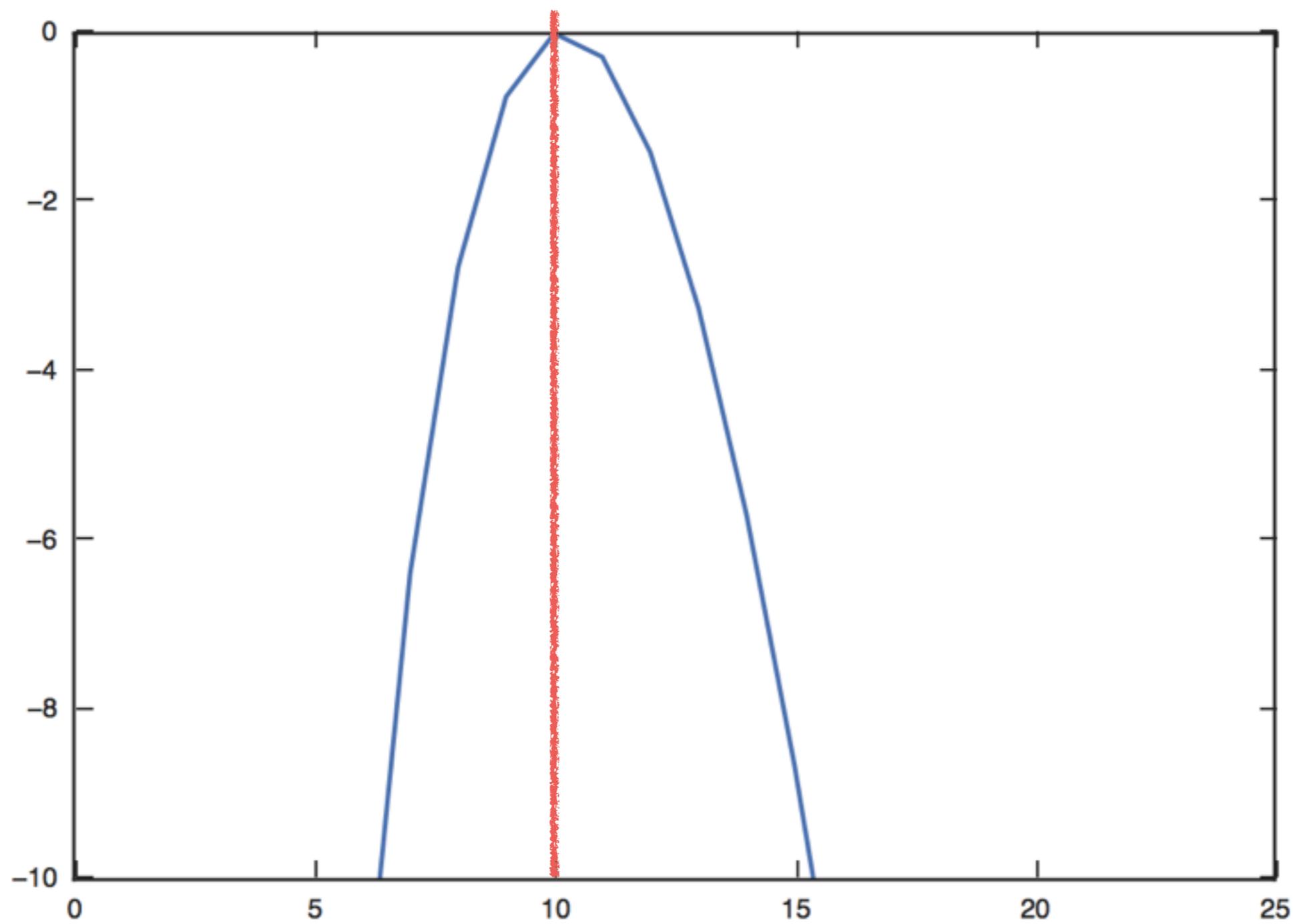
Maximum likelihood Estimator (MLE)

- Fit parameters by finding the maximum of the likelihood
- Likelihood = probability of data given model →> makes sense to maximize this!

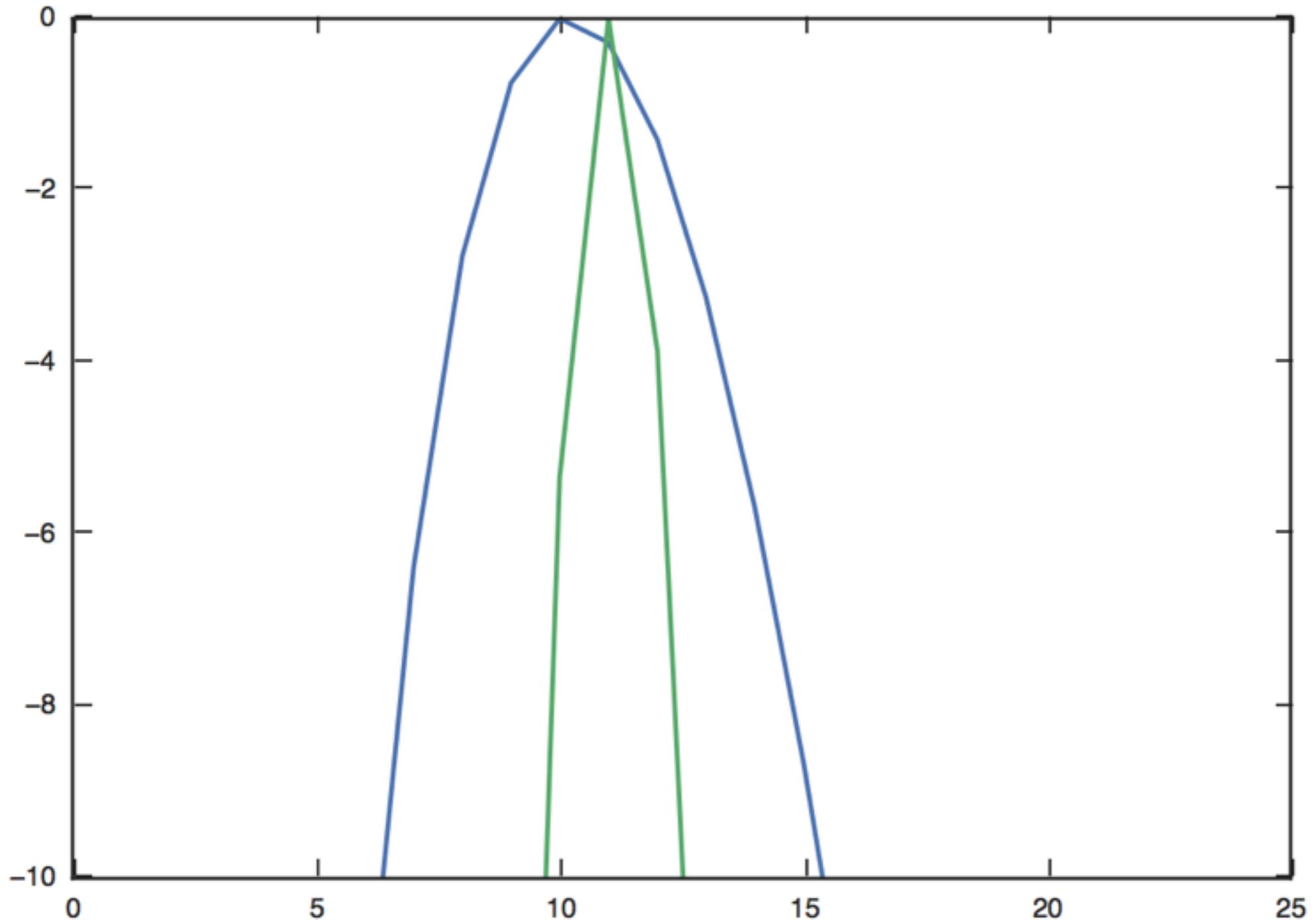
Sum In L



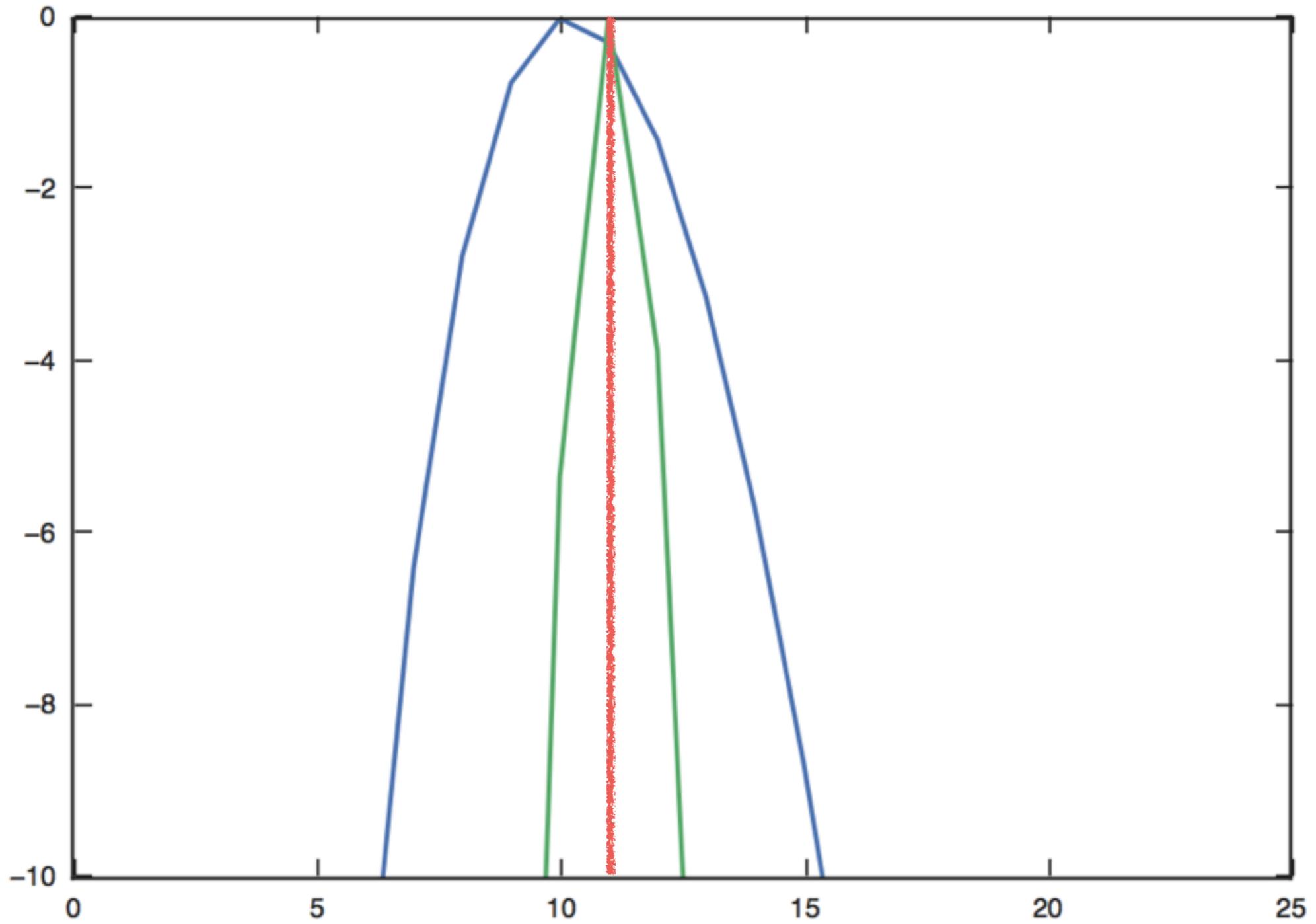
Sum In L



Sum In L, 100 observations



Sum In L, 100 observations



Desirable properties of maximum likelihood

- Units: $1/\text{data}$ \rightarrow maximum doesn't change when changing parametrization of model! (*functional invariance*)
- *Consistent*: approaches true value with probability 1 when N goes to infinity (\sim asymptotically unbiased)
- *Asymptotically normal*: Estimator becomes true value \pm Gaussian error
- *Asymptotically efficient*: Saturates Cramer-Rao bound when data goes to infinity (cannot get better estimate)

Example: Gaussian

- Have N measurements x_i with error σ , model = m

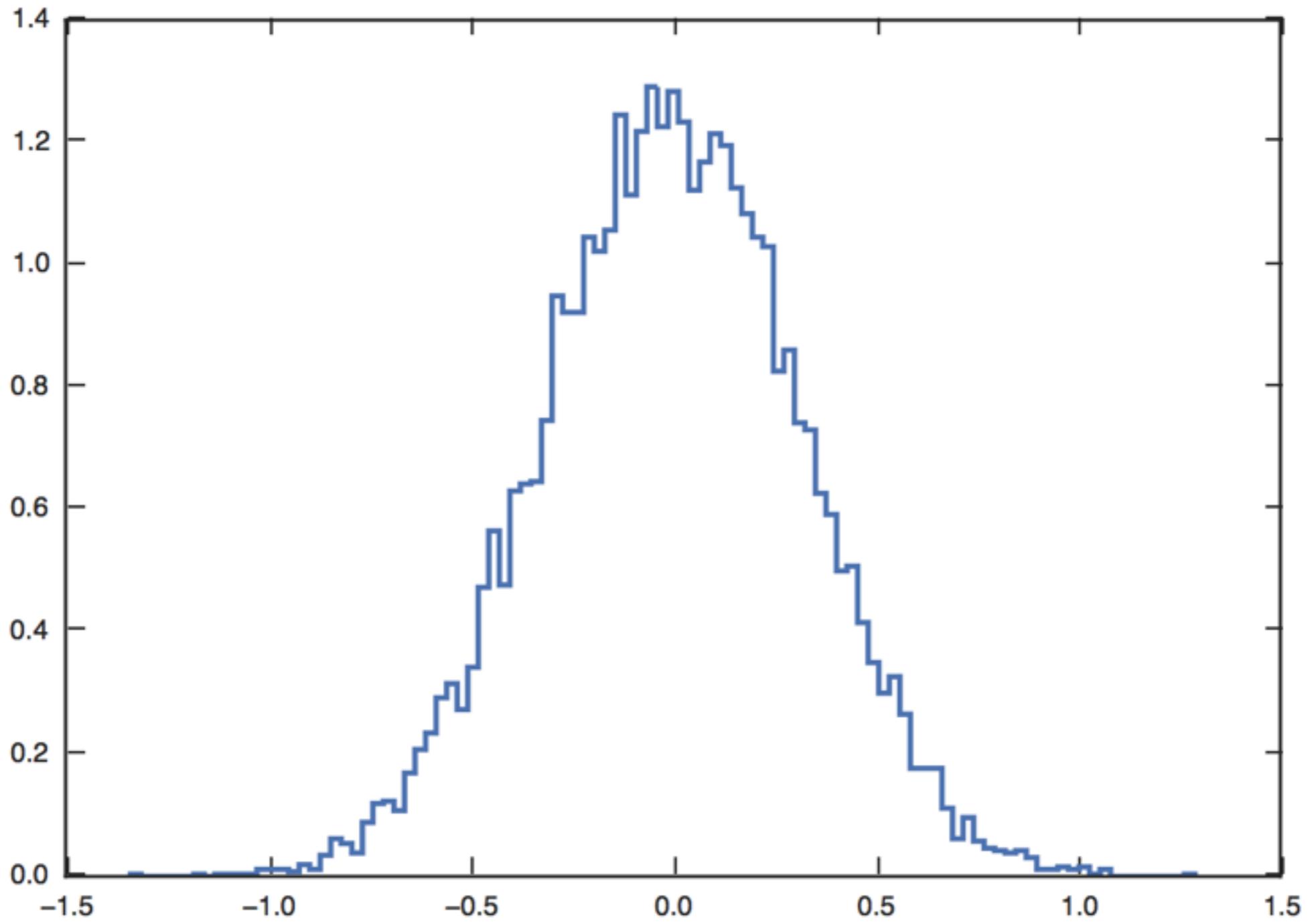
- $$\mathcal{L} = \prod_i p(x_i|m, \sigma) = \prod_i \mathcal{N}(x_i|m, \sigma^2)$$

- $$\ln \mathcal{L} = -\frac{1}{2} \sum_i \frac{(x_i - m)^2}{\sigma^2} + \text{constant}$$

- $$\frac{d \ln \mathcal{L}}{dm} = \sum_i \frac{(x_i - m)}{\sigma^2} = 0$$

- $$\sum_i x_i = Nm \rightarrow m = \langle x_i \rangle$$

- Unbiased!



Mean = -0.0037968773546516459

Example: Gaussian variance

- Have N measurements x_i with mean m , draw from Gaussian with variance V
- Mean is the same!

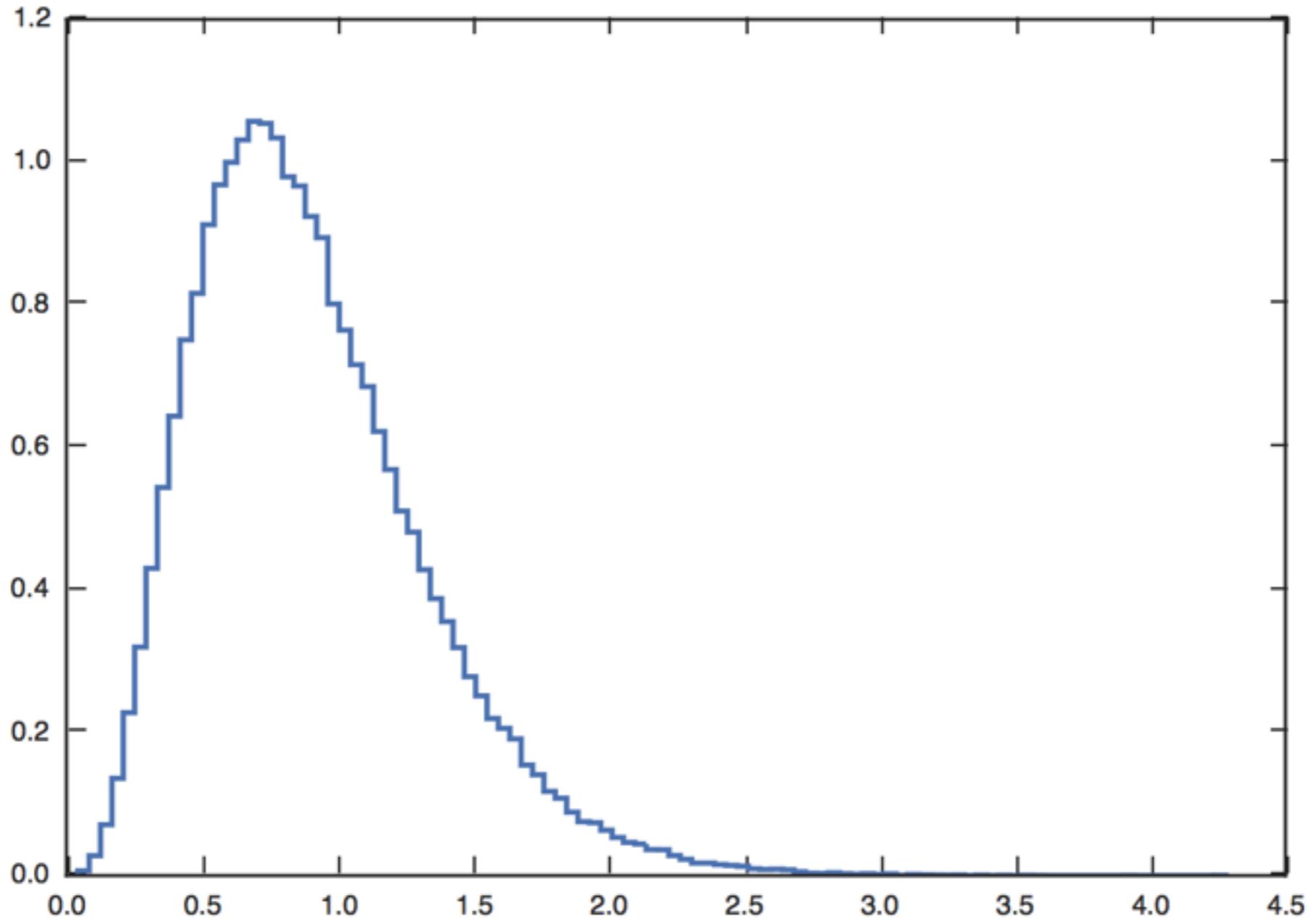
$$\mathcal{L} = \prod_i p(x_i|m, V) = \prod_i \mathcal{N}(x_i|m, V)$$

$$\ln \mathcal{L} = -\frac{1}{2} \sum_i \left[\frac{(x_i - m)^2}{V} + \ln V + \text{constant} \right]$$

$$\frac{d \ln \mathcal{L}}{dV} = \frac{1}{2} \sum_i \left[\frac{(x_i - m)^2}{V^2} - \frac{1}{V} \right] = 0$$

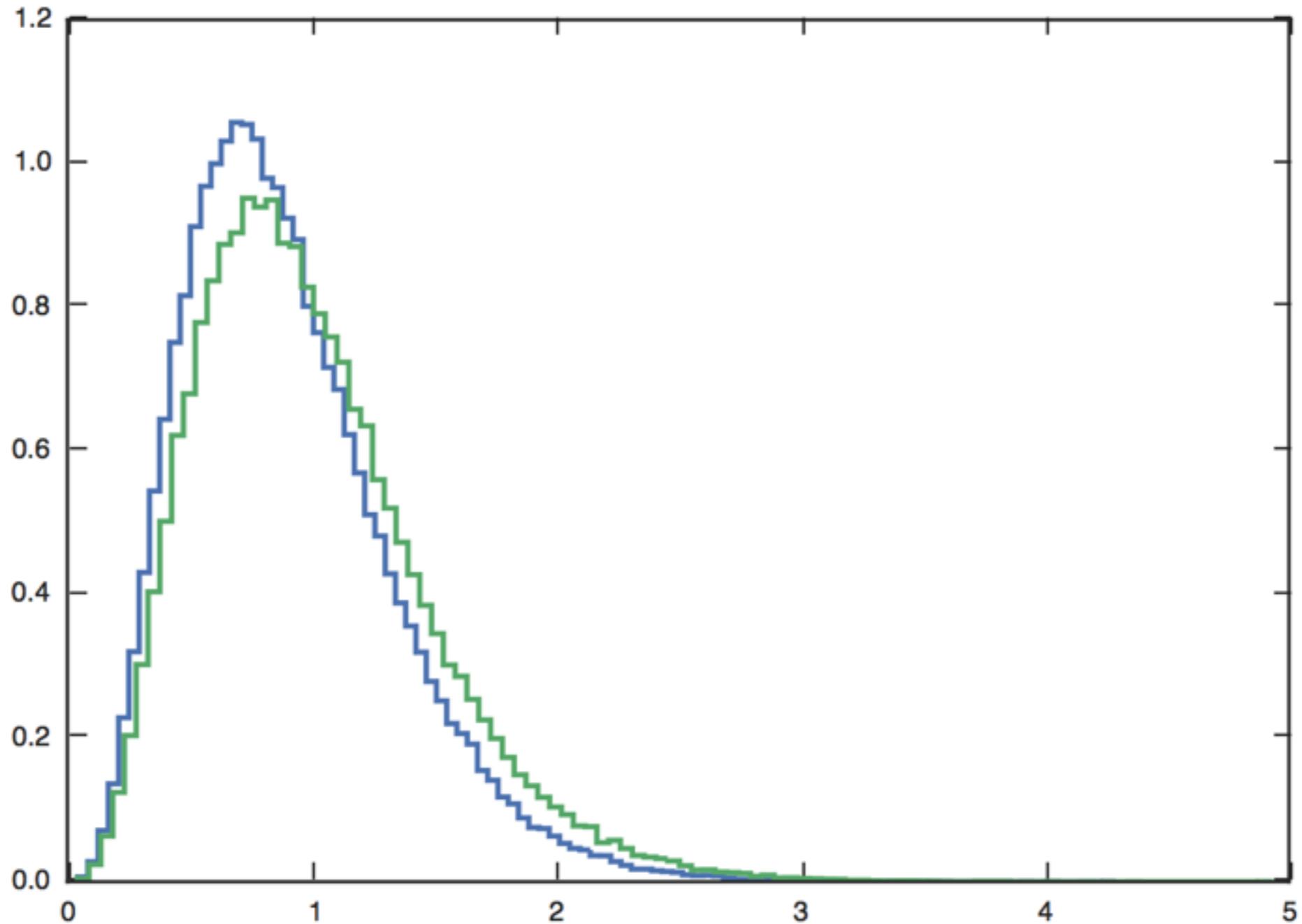
$$V = \frac{1}{N} \sum_i (x_i - m)^2$$

- Biased! (Unbiased has $1/[N-1]$)



mean=0.90158043895813211

Bessell correction: only N-1 constraints,
because 1 used for mean



Mean=0.99903451943557275

Confidence intervals

- Without Bayes, the likelihood on its own is *not* a probability distribution for the estimator
- Can derive *confidence intervals*: 95-percent confidence interval contains the true value 95% of the time
- Typically need to simulate data to figure this out; analytic results for some distributions
- Asymptotic normality: when N becomes large, difference between estimate and true value is Gaussian with variance

$$V_{ij} = -1/(d^2 \ln L / d \text{ model}_1 d \text{ model}_2) \text{ evaluated at MLE}$$

Example: Gaussian

- Have N measurements x_i with error σ , model = m

$$\mathcal{L} = \prod_i p(x_i | m, \sigma) = \prod_i \mathcal{N}(x_i | m, \sigma^2)$$

$$\ln \mathcal{L} = -\frac{1}{2} \sum_i \frac{(x_i - m)^2}{\sigma^2} + \text{constant}$$

$$\frac{d \ln \mathcal{L}}{dm} = \sum_i \frac{(x_i - m)}{\sigma^2}$$

$$\frac{d^2 \ln \mathcal{L}}{dm^2} = -\sum_i \frac{1}{\sigma^2} = -\frac{N}{\sigma^2}$$

- Uncertainty on m : $\frac{\sigma}{\sqrt{N}}$

Bayesian probability theory

- Bayesian probability theory follows from three axioms:
 - Degrees of plausibility are represented by real numbers
 - Qualitative consistency with common sense (e.g., $p(A|C) \uparrow$ then $p(\text{not } A|C) \downarrow$; small increases in plausibility lead to small increases in the real number representing it)
 - Consistency (internal, use of all information, indifference)

Bayesian probability theory

- Three axioms lead to probability calculus similar to deductive logic (see Chapters 1 & 2 of Jaynes' Probability Theory: The Logic of Science)
 - $P(A \cup B | C) = P(A | C) + P(B | C) - P(A \cap B | C)$
 - $P(A \cap B | C) = P(A | B \cap C) \times P(B | C)$
 - $P(A | B \cap C) = P(B | A \cap C) \times P(A | C) / P(B | C)$

Inference using Bayes's theorem

- Bayesian probability theory allows you to compute $p(\text{model} \mid \text{data})$
- Bayes's theorem:

$$p(\text{model} \mid \text{data}) = \frac{p(\text{data} \mid \text{model}) \times p(\text{model})}{p(\text{data})}$$

or

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

- Posterior probability distribution can be directly interpreted as probability of the model (parameters)

Posterior probabilities

- The fact that $p(\text{model}|\text{data})$ is a probability distribution has advantages and disadvantages:
 - **Bad:** $p(\text{model}|\text{data})$ is not functionally independent: changing the parametrization of the model will change $p(\text{model}|\text{data})$ \rightarrow maximum-a-posteriori estimate, mean, etc. depend on parametrization
 - **Good:** Can directly derive *credibility intervals* from $p(\text{model}|\text{data})$
 - **Good:** Can marginalize over nuisance parameters: $p(\text{model}|\text{data}) = \int d \text{nuisance } p(\text{model}, \text{nuisance}|\text{data})$
 - **Good:** Can carry full $p(\text{model}|\text{data})$ forward to ‘new data’
 $p(\text{model} | \text{new data}, \text{data}) = p(\text{new data} | \text{model}) p(\text{model}|\text{data}) / p(\text{new data})$
- All good things come at the cost of introducing the *prior* $p(\text{model})$, which many people find hard to stomach...

A word on priors

- Any application of Bayes's theorem requires priors, often considered a disadvantage
- As the name implies, these typically encode one's prior knowledge of the model (parameters) under investigation
- Long literature on "uninformative priors": rules of thumb:
 - Unitless parameter: flat prior over reasonable range
 - Parameter with units: flat prior on $\ln(\text{parameter})$; puts equal weight on different orders of magnitude
 - However, if you know the order of magnitude, a flat linear prior might be more appropriate
 - If prior matters much, then your data is not that informative!
- Use freedom in specifying the prior to your advantage (hierarchical modeling)

“Uninformative” priors

- One is typically expected to use “non-informative priors”: priors that do not strongly constrain the posterior
- Note: choosing the model is often a very strong prior!
- For example: unitless parameter A : 1, 1.5, 2.5, 3.3, ... no reason to prefer any $\rightarrow p(A) = \text{constant}$ (improper!)
- Scale parameter V (has units): prior shouldn't depend on units \rightarrow should be invariant under re-scaling

$$p(V) dV = p_W(W=sV) d(sV) = p(W=sV) d(sV) \rightarrow$$

$$p(V) \sim 1/V$$

Example: Gaussian variance

- Have N measurements x_i with mean m , draw from Gaussian with variance V
- Prior on the mean: constant, prior on the variance $\sim 1/\text{variance} = 1/V$
- Mean is the same as MLE

$$\mathcal{L} = \prod_i \mathcal{N}(x_i | m, V) \rightarrow p(V | x_i) \propto \mathcal{L} \frac{1}{V}$$

$$\ln p(V | x_i) = -\frac{1}{2} \sum_i \frac{(x_i - m)^2}{V} - \frac{N}{2} \ln V - \ln V + \text{constant}$$

$$\frac{d \ln p(V | x_i)}{dV} = \frac{1}{2} \sum_i \frac{(x_i - m)^2}{V^2} - \frac{N + 2}{2V}$$

$$V = \frac{1}{N + 2} \sum_i (x_i - m)^2$$

- Biased! (Unbiased has $1/[N-1]$)

So far,
uniform for unitless parameters ,
 $1/\text{param}$ for unit-full parameters
has served me pretty well...

Advanced approaches to determining priors

- Jeffreys prior:
prior \sim square-root (determinant Fisher Information)

$$\text{Fisher information} = E[-(d^2 \ln L / d \text{ model}^2)]$$

- Invariant under change of variables (good!)

Advanced approaches to determining priors

- Conjugate priors: For computational ease, useful to get $p(\text{model}|\text{data})$ that has the same form as $p(\text{model})$

- So want

$$p(\text{model}|\text{data}) \sim p(\text{data}|\text{model}) p(\text{model})$$

to have the same form as $p(\text{model}) \rightarrow p(\text{model})$ set by likelihood

- For example, mean of a Gaussian: conjugate prior on mean is Gaussian
- Useful if you want an *informative* prior, but want to be able to, e.g., compute the maximum of the posterior probability analytically

Advanced approaches to determining priors

- Maximum entropy: If you want as uninformative prior as possible, but have some constraints (information)
- Maximize entropy = $-\sum_i p_i \ln[p_i]$ (or integral generalization) under certain constraints (Lagrange multipliers and all that)

Okay, you have a prior and the likelihood, now what do you do with the posterior probability distribution?

What to do with PDFs

- Bayes's theorem:

$$p(\text{model} \mid \text{data}) = \frac{p(\text{data} \mid \text{model}) \times p(\text{model})}{p(\text{data})}$$

- *Some* people would claim that you need to publish $p(\text{model} \mid \text{data})$ somehow
- Practically, need *summaries*
- Single-point summaries: MAP (maximum-a-posteriori value), mean, median, ...
- Width: variance? Some range of quantiles, like 68% around single-point
- Latter: Start at (max, mean, median, ...) and integrate outward at constant p until you have 68% of the area; works in multi-D
- Multi-modal PDFs: Sorry! Do something sensible.

Bayesian inference recap

- Likelihood: $p(\text{data}|\text{model})$, comes from underlying (physical/empirical) model + observing procedure (noise, PSF, ...)
- Pick reasonable prior: uninformative or based on previous results
- compute posterior PDF \sim likelihood \times prior: Can use grid for low-dim, sampling methods for higher dim (next week)
- Compute summaries of PDF to list in tables, abstracts, press releases

Bayesians vs. frequentists

- Like most of such battles, there is very little actually at stake; at high SNR, all good (unbiased, efficient) methods return the same answer
- Bayes's theorem proven to be optimal way to do inference; so will get best results by using it!
- Likelihood-based frequentist methods often *very* similar to corresponding Bayesian method
- Bayesian inference has more freedom than frequentist inference: can open up the prior to modeling (empirical Bayes, hierarchical modeling)
- Difficult to do marginalization in frequentist approach → difficult to integrate over lack of knowledge