

• ALL HW DUE WED in class

LAST TIME

LIKELIHOOD FUNCTION

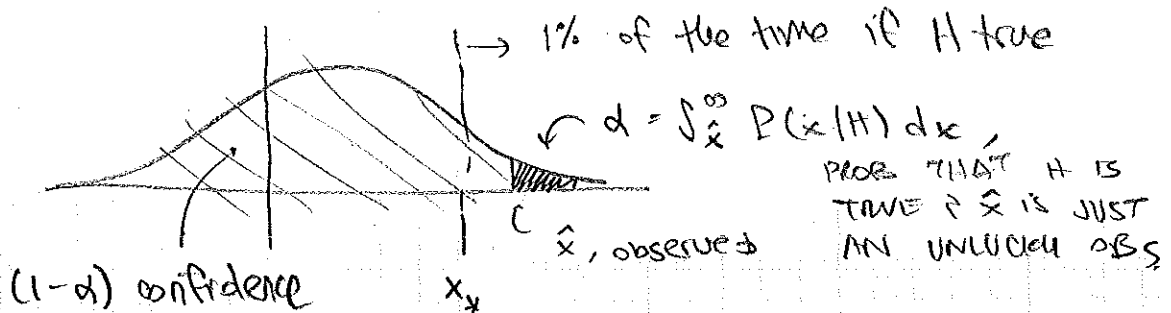
$$\mathcal{L}(H) = P(\hat{x}|H)$$

\hat{x} : OBSERVED DATA

eg H can dep on a parameter a

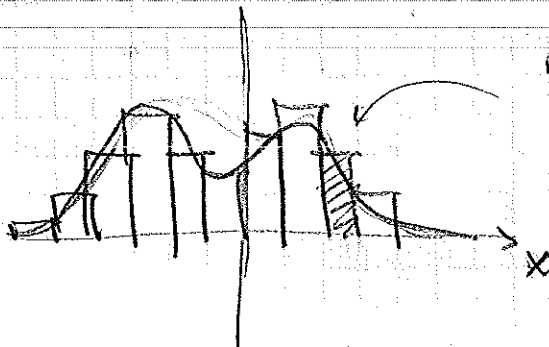
give p-values, α
confidence $(1-\alpha)$

if P is Gaussian
then σ 's



PRACTICAL MATTER: OBS DATA IS OFTEN BINNED

→ DISCRETIZED
just like when we talked
about discretizing function space



in k th BIN:

$$P(k|H) = \int_{\text{BIN } k} dx P(x|H)$$

PROBABILISTIC PROCESS!

eg. unstable particle has lifetime τ ... but how to measure?

$$P(t|\tau) dt = \frac{1}{\tau} e^{-t/\tau} dt \quad \leftarrow \text{FACT}$$

then take N particles & time their decays

$$\hookrightarrow \hat{t}_1, \hat{t}_2, \dots, \hat{t}_N = \{\hat{t}_i\}$$

$$P(\{\hat{t}_i\} | \tau) = \prod_{i=1}^N \left(\frac{1}{\tau} e^{-\hat{t}_i / \tau} \right)$$

LIKELIHOOD,
 $L(\tau)$

many terms! but can see

$$= \left(\frac{1}{\tau}\right)^N e^{-\sum_{i=1}^N \hat{t}_i / \tau}$$

easy to write

$$\mathcal{L}(\tau) \equiv \log L(\tau) = \sum_{i=1}^N \left[\log \frac{1}{\tau} - \frac{\hat{t}_i}{\tau} \right]$$

same info as L

SUM EASIER THAN PRODUCT
(esp computationally)

HOW TO ESTIMATE τ ? $\hat{\tau}$ IS THE VALUE THAT MAXIMIZES L
 τ ESTIMATOR

$$\leftrightarrow \text{MAX } \mathcal{L}$$

$$\leftrightarrow \frac{\partial}{\partial \tau} \mathcal{L} \Big|_{\hat{\tau}} = \sum_{i=1}^N \left(-\frac{1}{\tau} + \frac{\hat{t}_i}{\tau^2} \right) = \left[-\frac{N}{\tau} + \frac{\sum \hat{t}_i}{\tau^2} \right] = 0$$

$$\Rightarrow \boxed{\hat{\tau} = \frac{1}{N} \sum_{i=1}^N \hat{t}_i} \quad \leftarrow \text{no surprise, the average}$$

ACTUALLY, there's another way to approach this:

$$p(t | \tau) = \frac{1}{\tau} e^{-t/\tau} \mapsto p(t | \Gamma) = \Gamma e^{-t\Gamma}$$

$\Gamma = 1/\tau$: DECAY RATE.

EQUIVALENT INFO.
DIFFERENT CHOICE OF
PARAMETER.

Q: DO WE GET THE SAME ESTIMATOR

$$\hat{\Gamma} = 1/\hat{\tau} ?$$

$$\mathcal{L} = \sum_{i=1}^N [\log \tau - \tau \hat{t}_i]$$

$$\left. \frac{\partial \mathcal{L}}{\partial \tau} \right|_{\hat{\tau}} = \frac{N}{\hat{\tau}} - \sum_{i=1}^N \hat{t}_i = 0 \Rightarrow \hat{\tau} = \underbrace{\left[\frac{1}{N} \sum_{i=1}^N \hat{t}_i \right]^{-1}}_{\hat{\tau} \checkmark}$$

THIS DEMONSTRATES A GENERAL PROPERTY:

the maximum of the likelihood is the same, no matter what parameterization.

$$\uparrow \quad \hat{\tau} = \hat{\tau}^{-1} \quad \text{in the same way that } \tau = \tau^{-1}$$

from $\max \mathcal{L}$

BUT EVEN THOUGH THE ESTIMATORS $\hat{\tau} \neq \hat{\tau}^{-1}$ ARE CONSISTENT, THEY ARE NOT EQUALLY USEFUL!

eg. EXPECTATION OF THE MEASUREMENTS \hat{t}_i :

$$\begin{aligned} \langle \hat{t}_i \rangle &= \int \hat{t}_i \frac{1}{\tau} e^{-\hat{t}_i/\tau} d\hat{t}_i \\ &= \int \left(\frac{\hat{t}_i}{\tau} \right) e^{-(\hat{t}_i/\tau)} d\left(\frac{\hat{t}_i}{\tau} \right) \tau \\ &= \tau \end{aligned}$$

\Rightarrow EXPECTATION OF ESTIMATOR:

$$\langle \hat{\tau} \rangle = \left\langle \frac{1}{N} \sum_{i=1}^N \hat{t}_i \right\rangle = \frac{1}{N} \sum_{i=1}^N \langle \hat{t}_i \rangle = \tau \checkmark$$

AGREES $\forall N$.

on the other hand.

$$\langle \hat{\Gamma} \rangle = \left\langle \frac{N}{\sum \hat{t}_i} \right\rangle = N \left\langle \int_0^\infty dw \prod_{i=1}^N \left[e^{-w \hat{t}_i} \right] \right\rangle$$

TRICK: $\frac{1}{x} = \int_0^\infty dw e^{-wx}$
for $x > 0$

$$= N \int_0^\infty dw \left[\prod_{i=1}^N \int_0^\infty dt_i \left[e^{-w \hat{t}_i} \right] \Gamma e^{-\Gamma \hat{t}_i} \right]$$

(by def. of Γ)

$$= N \int_0^\infty dw \left(\frac{\Gamma}{\Gamma + w} \right)^N$$

$$= N \Gamma^N \int_\Gamma^\infty du u^{-N}$$

$$= N \Gamma^N \frac{1}{1-N} u^{1-N} \Big|_{u=\Gamma}^{u=\infty}$$

$$= N \Gamma^N \cdot \frac{-1}{1-N} (\Gamma^{1-N})$$

$$= \left[\frac{N}{N-1} \Gamma \right] \neq \Gamma !!$$

BIASED ESTIMATOR OF $\hat{\Gamma}$

LESSON: MAX OF LIKELIHOOD: DOESN'T MATTER WHAT PARAMETER.

BUT BIAS DOES CARE, @ LEAST FOR FINITE N.

ANOTHER EXAMPLE : GAUSSIANS

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

$$\mathcal{L} = \log \mathcal{L} = \sum_i^N \left[\log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{(\hat{x}_i - \mu)^2}{2\sigma^2} \right]$$

$$\uparrow$$

$$P(\{x_i\} | \mu, \sigma^2)$$

↑ measurements

$$\frac{\partial \mathcal{L}}{\partial \mu} \Big|_{\hat{\mu}, \hat{\sigma}^2} = + \frac{1}{\hat{\sigma}^2} \sum_i^N (\hat{x}_i - \hat{\mu}) = 0 \Rightarrow \boxed{\hat{\mu} = \frac{1}{N} \sum_i^N x_i}$$

unbiased

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \sigma^2} \Big|_{\hat{\mu}, \hat{\sigma}^2} &= \sum_i \left[\underbrace{\frac{1}{\sqrt{2\pi\hat{\sigma}^2}}}_{-\frac{1}{2} \frac{1}{\hat{\sigma}^2}} \cdot \underbrace{\left(-\frac{1}{2} (\hat{\sigma}^2)^{-3/2} \right)}_{-\frac{1}{2} \frac{1}{\hat{\sigma}^2}} + \frac{(\hat{x}_i - \hat{\mu})^2}{2(\hat{\sigma}^2)^2} \right] \\ &= \frac{-N}{2\hat{\sigma}^2} + \frac{\sum_i (\hat{x}_i - \hat{\mu})^2}{2(\hat{\sigma}^2)^2} = 0 \end{aligned}$$

$$\Rightarrow \boxed{\hat{\sigma}^2 = \frac{1}{N} \sum_i (\hat{x}_i - \hat{\mu})^2}$$

↑

BUT THIS IS BIASED : $\hat{\mu}$ IS BUILT OUT OF \hat{x}_i ,

so if we had a vector

$$\underline{\hat{x}} = (\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N)$$

then $\sum_i (\hat{x}_i - \hat{\mu})$ subtracts out fluctuations along $(1, 1, 1, \dots, 1)$

so only fluctuations in $(N-1)$ other dir.
are gaussian \rightarrow VARIANCE IN LENGTH OF VECTOR IS $(N-1)\sigma^2$

\hookrightarrow unbiased estimator : $S^2 = \frac{1}{N-1} \sum_i (x_i - \mu)^2$

A LOOK AT CONFIDENCE INTERVALS

eg measure \hat{x} , this is the estimator for x
eg by likelihood.

BUT THE TRUE VALUE x_t IS PROBABLY NOT \hat{x} .
WHAT WINDOW AROUND \hat{x} ARE YOU REASONABLY
SURE THAT x_t IS IN?

eg (68%) sure that $x_1 < x_t < x_2$

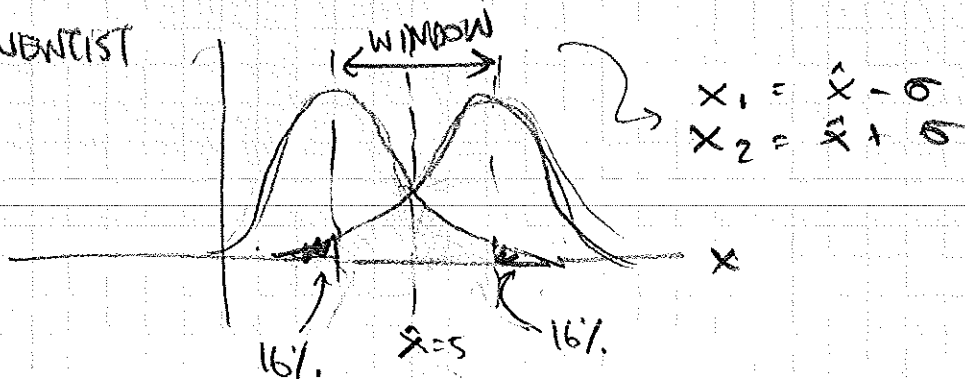
CONSIDER FIRST: GAUSSIAN

SUPPOSE $\sigma^2 = 1$
 $\hat{x} = 5$ ← measured

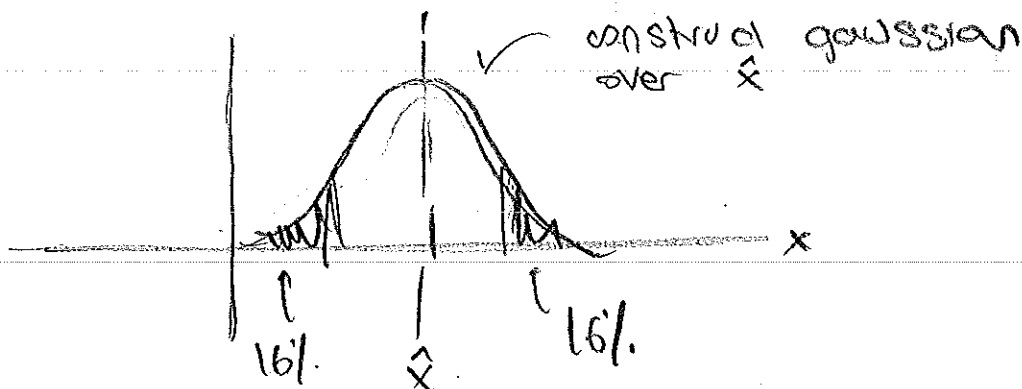
no matter what x_t is
68% of WINDOWS BY
AN ENSEMBLE OF
EXPTS WILL INCLUDE
 x_t .

① ONE WAY TO DO THIS $(100\% - 68\%) \frac{1}{2} = 16\%$

FREQUENTIST



USUALLY WE THINK OF THIS AS



this mental construction is justified by likelihood

$$L = p(\hat{x} | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\hat{x} - \mu)^2}{2\sigma^2}}$$

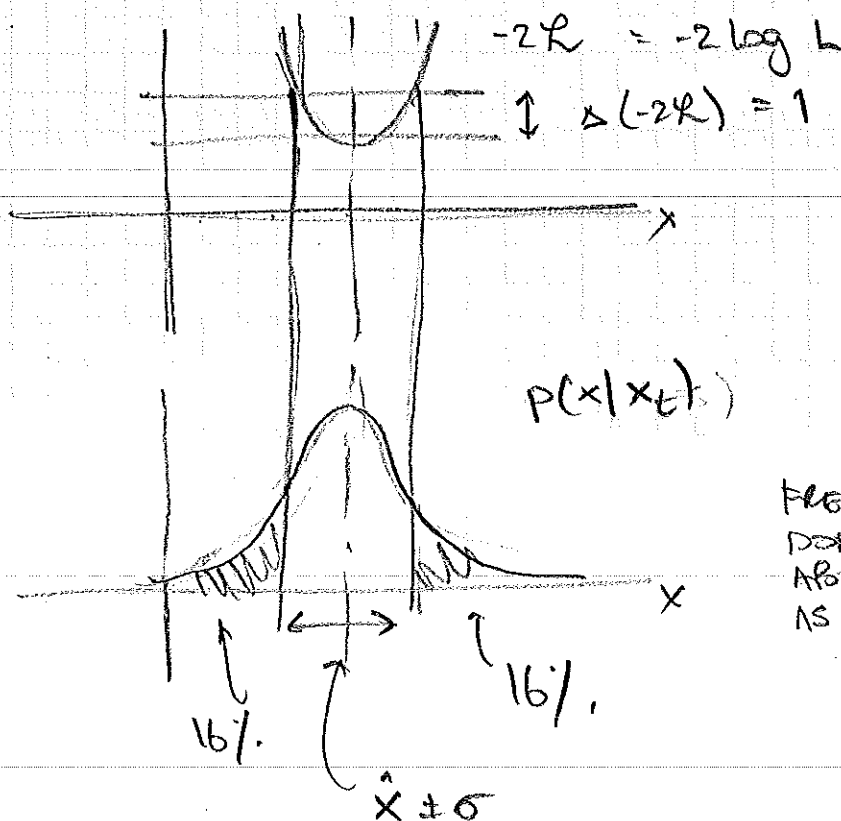
$$\mathcal{L} = \log L = \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{(\hat{x} - \mu)^2}{2\sigma^2} \leftarrow \text{ASSUME } \sigma \text{ KNOWN}$$

to construct conf. interval, take difference w/r + maximum likelihood

USEFUL TO WORK WITH

$$-2\mathcal{L} = \frac{1}{\sigma^2} (\hat{x} - \mu)^2 + (\text{indep of } \hat{x})$$

② then: when $-2\mathcal{L}$ changes by 1 unit, that sets boundary of window
 $\hookrightarrow \pm \sigma$ for GAUSSIAN $p(x | x_t)$



FREQUENTIST!
 DON'T THINK
 ABOUT $p(x | x_t)$
 AS A PDF!

so $-2 \log L$ arg supports "draw L " picture

\hookrightarrow BUT THAT TREATS L AS A PDF

THE PROBLEM: $L(x_t) \neq P(x_t | \hat{x})$

this is a pdf over x_t
shows up as $P(x_t | \hat{x}) dx_t$

nothing about this is a pdf

$L = P(\hat{x} | x_t) \leftarrow$ pdf w/rt \hat{x}
BUT NOT x_t
will NEVER see $P(\hat{x} | x_t) dx_t$

... at least as FREQUENTIST

AS A BAYESIAN:

this prior is a pdf in x_t

$$P(x_t | \hat{x}) = \frac{L(x_t) P(\hat{x})}{P(x_t)}$$

$P(\hat{x}) = 1$

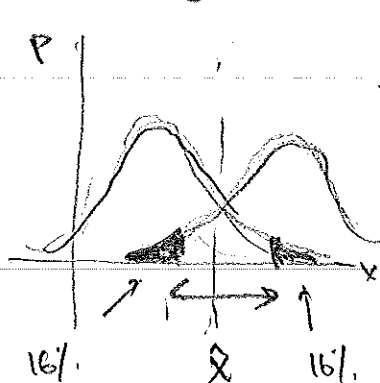
effectively \rightarrow

$P(\hat{x} | x_t) P(x_t) \leftarrow$ maybe uniform!
(subj to caveat: what parameterize?)

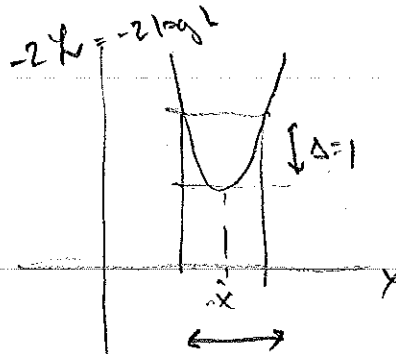
$$\int dx_t P(\hat{x} | x_t) P(x_t) \leftarrow \text{normalize}$$

only makes sense b/c of $P(x_t)$

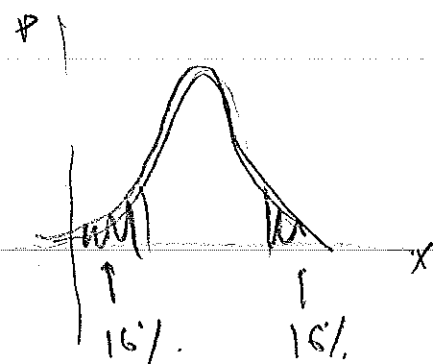
3 eqs of CONFIDENCE INTERVAL



FREQUENTIST



FREQ. w/
LIKELIHOOD



BAYES w/
FLAT PRIOR

WHY THIS MATTERS

things are more subtle when the pdf is not Gaussian.

eg. POISSON STATISTICS

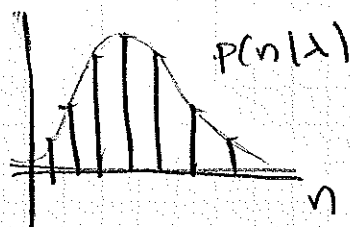
$$P(n|\lambda) = \frac{1}{n!} e^{-\lambda} \lambda^n$$

\uparrow
n events

\uparrow
 λ expected

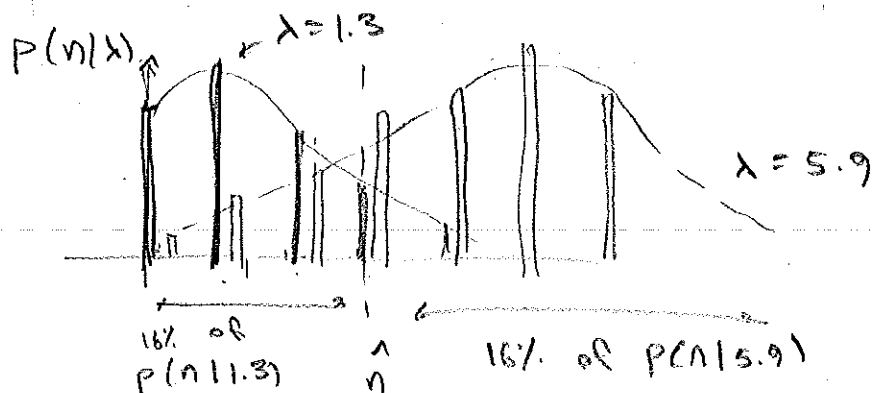
eg. rare decay \rightarrow PROBABILISTIC; "WOULD BE BINOMIAL" BUT I DON'T HAVE A SENSE OF # OF TRIALS.

OBS. DISCRETE AS A "DISTRIBUTION" IN n



SUPPOSE YOU MEASURE $\hat{n} = 3$ events.
WANT TO GIVE CONFIDENCE INTERVAL FOR λ .

① FREQUENTIST: 16% BELOW, 16% ABOVE



② LIKELIHOOD RATIO

$$L(\lambda) = \frac{\lambda^{\hat{n}}}{\hat{n}!} e^{-\lambda}$$

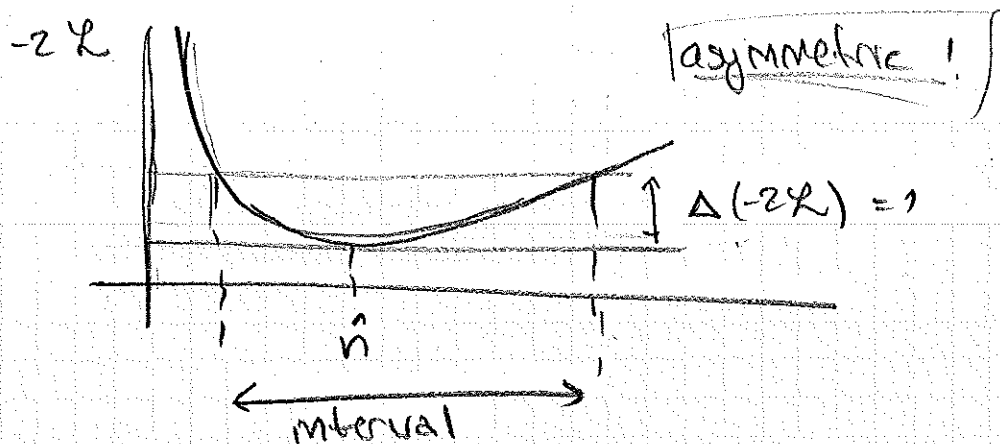
← continuous function of λ
still not a pdf

$$P(\hat{n}|\lambda)$$

← nothing about this allows integration over $d\lambda$

not Gaussian! no longer any sense that $\Delta(-2\log L) = \Delta(-2\chi^2) \Rightarrow 15$

BUT CAN STILL PLOT & use those guidelines



WHAT IF WE'RE BAYESIAN?

STICK TO FLAT PRIOR: $P(\lambda) = \text{const.}$

↳ non-flat changes things a lot — but let's see how even a flat prior changes the analysis

now we have:

$$P(\lambda|\hat{n}) = \frac{P(\hat{n}|\lambda) P(\lambda)}{\int d\lambda} \quad \leftarrow \text{normalize}$$

↑ this is a pdf

↑ $L(\lambda)$ val a pdf in λ

↑ this is a pdf in λ

WE'RE STILL LOOKING FOR 68% CONFIDENCE

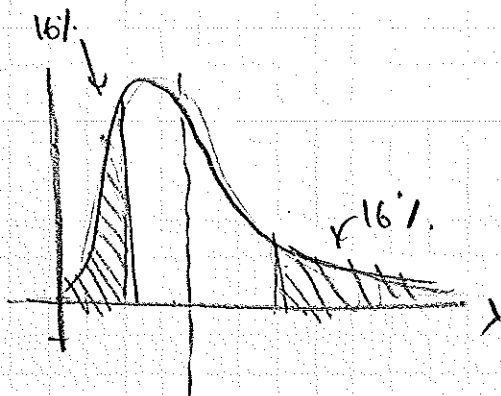
(some window around \hat{n} where true value λ is 68% likely to fall into

↑

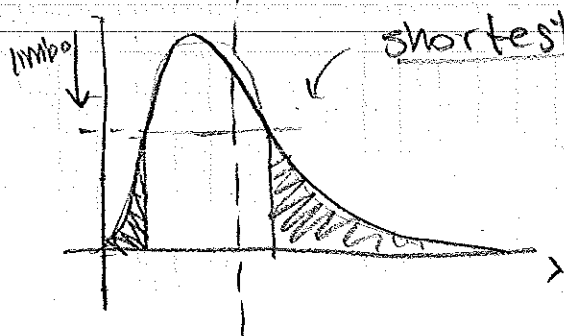
if we applied this procedure to many different labs, the true value λ will fall into the windows of 68% of them.

BUT $p(\lambda | \hat{n})$ is asymmetric... so we have a choice of how to do this

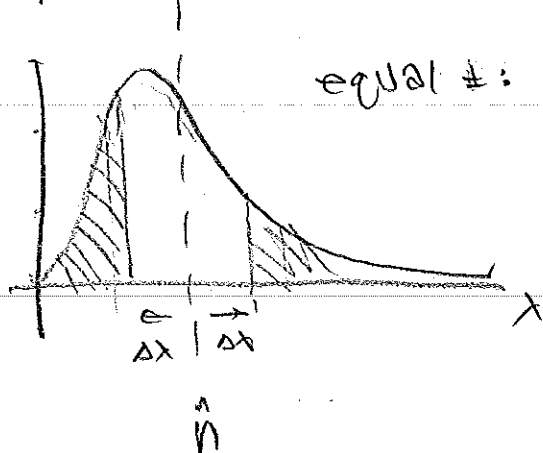
→ BOTTOM LINE: SAY EXACTLY WHAT YOU DID



central 16% to left
16% to right



shortest: $p(\lambda | \hat{n})$ outside interval should be less than $p(\lambda | \hat{n})$ inside.



equal ±: interval is $\hat{n} \pm \Delta\lambda$

this is just to illustrate some of the nuances when giving confidence intervals

→ to highlight the operational differences between FREQUENTIST & BAYESIAN.

seems kind of academic, until you have an experiment that says

"the best fit for the neutrino mass

is $m^2 = (-50 \pm 30) \text{ eV}^2$ "

↑

m^2 cannot be negative!!

the theory upon which these measurements were fit assume $m > 0$.

so HERE A BAYESIAN PRIOR

OF $P(m^2) = \Theta(m^2)$

MAY BE USEFUL.