

Calibrating SN Ia Standard Cepheids

function of
distance

$$M = m - M \quad \text{abs. mag}$$

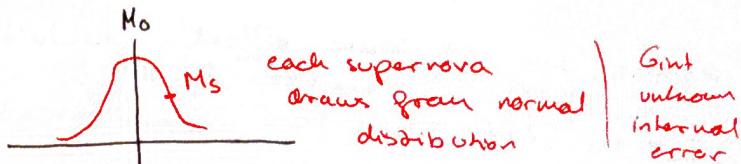
suppose have already calibrated cepheids

$s = 1, \dots, N$ measurements, $\hat{\mu}_s$, $\sigma_{\text{cepheid}, s}$ from cepheids

\hat{m}_s , $\sigma_{m, s}$ from SN Ia

$$\Rightarrow \hat{M}_s = \hat{m}_s - \hat{\mu}_s \Rightarrow \sigma_{\text{err}, s}^2 = \sigma_{m, s}^2 + \sigma_{\text{cepheid}, s}^2 \quad \text{assume } G_s \text{ known}$$

(assume $\hat{\mu}_s$, \hat{m}_s independent)



- assume $M_s \sim N(M_0, G_{\text{int}}^2)$

[could also say $M_s = M_0 + \varepsilon_s^M$, and $\varepsilon_s^M \sim N(0, G_{\text{int}}^2)$]

so the errors in meas normally dist. around true value

- Calibrators: SN $\hat{m}_s \sim N(m_s, \sigma_{m, s}^2)$ true values LATENT VARIABLES

Cepheid $\hat{\mu}_s \sim N(\mu_s, \sigma_{\text{cepheid}, s}^2)$

[could also write $\hat{m}_s = m_s + \varepsilon_{m, s}$ $\hat{\mu}_s = \mu_s + \varepsilon_{\mu, s}$] This is how you'd generate data by code
where $\varepsilon_{m, s} \sim N(0, \sigma_{m, s}^2)$ $\varepsilon_{\mu, s} \sim N(0, \sigma_{\text{cepheid}, s}^2)$

- Latent variable: $M_s = M_0 - \mu_s$, use estimator $\hat{M}_s = \hat{m}_s - \hat{\mu}_s$

$$\Rightarrow \hat{M}_s \sim N(M_s, \sigma_{\text{err}, s}^2 = \sigma_{m, s}^2 + \sigma_{\text{cepheid}, s}^2)$$

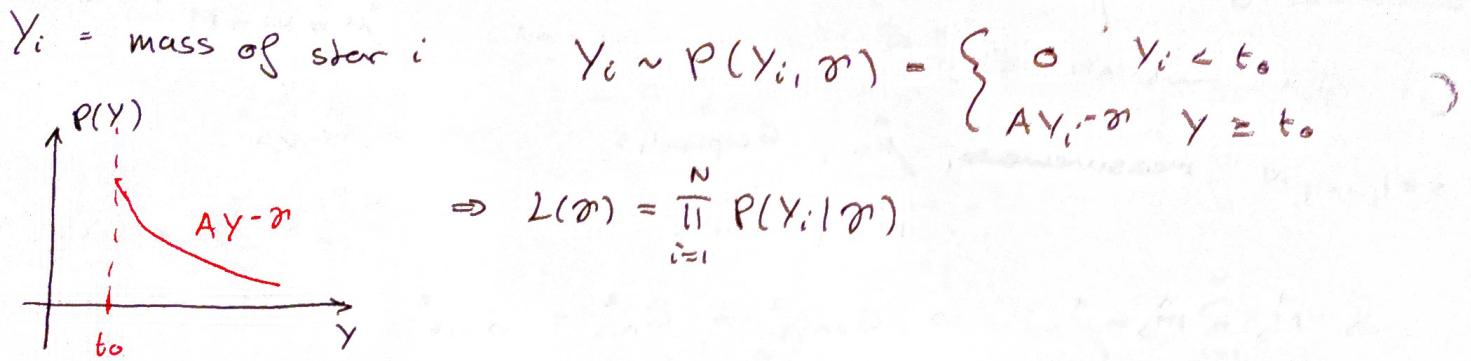
assume $\hat{\mu}_s$, $\sigma_{\text{cepheid}, s}$ known, measurement error and $\varepsilon_{m, s} \perp \varepsilon_{\mu, s}$ indep

$$P(\{\hat{M}_s\} | M_0, G_{\text{int}}^2) = \prod_{s=1}^N P(\hat{M}_s | M_0, G_{\text{int}}^2) \quad \hat{M}_s \sim N(M_0, G_{\text{int}}^2 + \sigma_{\text{err}, s}^2)$$

$$= \prod_{s=1}^N N(\hat{M}_s | M_0, G_{\text{int}}^2 + \sigma_{\text{err}, s}^2)$$

can't actually take analytic derivatives bc won't be able to solve for G_{int} \Rightarrow need to solve numerically

Inference of Stellar Mass Function



- suppose due to some effect couldn't see stars above some mass
- SELECTION EFFECT in this case truncation at t_1

→ introduce indicator $I_i = \begin{cases} 0 & \text{otherwise} \\ 1 & \text{if } Y_i \text{ observable} \end{cases}$

$$\Rightarrow P(I_i = 1 | Y_i, t_1) = \begin{cases} 1 & Y \leq t_1 \\ 0 & Y > t_1 \end{cases}$$

- Modified likelihood function:

$$P(X_i^{\text{obs}} | I_i = 1, \pi) = \frac{P(I_i = 1, Y_i | \pi)}{\int dY_i P(I_i = 1, Y_i | \pi)} = \frac{P(I_i = 1 | Y_i) P(Y_i | \pi)}{\int dY_i P(I_i = 1 | Y_i) P(Y_i | \pi)}$$

see example sheet for solution, this LHAD gives true value
while naive LHAD overestimates π

QUANTIFYING UNCERTAINTY

- Suppose $\vec{X} = (X_1, \dots, X_N) \stackrel{\text{iid}}{\sim} P(x)$

$$\vec{x}^{\text{obs}} = (x_1^{\text{obs}}, \dots, x_N^{\text{obs}}), \hat{G}(\cdot | \vec{x}^{\text{obs}}) \pm \sqrt{\text{Var}(\hat{G})}$$

- For case $X_i \stackrel{\text{iid}}{\sim} N(\mu, 1)$ $\vec{x}^{\text{obs}} = (-0.639, -0.9313, 0.1577, -0.8813)$

$$\bar{x}^{\text{obs}} = -0.5735 \quad G_x = Y_4 \Rightarrow G_{\bar{x}} = 0.5, \bar{x} \sim N(\mu, G^2/N)$$

$$\Rightarrow \bar{x} \pm \bar{G}_{\bar{x}} = [-1.0735, -0.0735]$$

This is a 68% confidence interval

- what does a confidence interval mean?

from a probability perspective (frequentist) only random variables have probabilities, so $P(\mu \in [L, U]) = m$ doesn't make sense because freq. takes μ as given, not a rv.

⇒ two estimators! $[L(\bar{X}), U(\bar{X})]$ functions of rv's → rvs

$$= C(\bar{X})$$

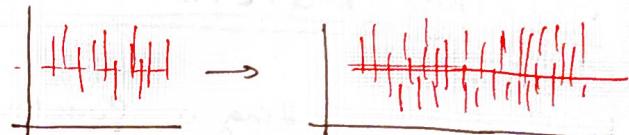
$$\begin{aligned} P(C(\bar{X}) \text{ contains } \mu) &\geq (1-\alpha) \\ &\geq \beta \end{aligned}$$

is an β percent confidence interval

the interval is random!

e.g.: run an experiment lots of times:

as $n \rightarrow \infty$ the probability of μ being in confidence interval goes to above



- More complicated cases:

$\bar{X} \sim p_G(\bar{x})$, $\hat{g}(\bar{X})$ not a simple function $\text{Var}(\hat{g}(\bar{X}))$ is intractable

- What could you do? above method worked bc gaussian, here not

simulate: $\vec{x}_{\text{sim},1}, \dots, \vec{x}_{\text{sim},M} \Rightarrow \text{Var}(\{\hat{g}_1, \dots, \hat{g}_M\}) = \frac{1}{M-1} \sum_{m=1}^M (\hat{g}_m - \bar{\hat{g}})^2$

but a lot of times can't do / too hard / model incorrect

* Bootstrap

suppose have $\vec{X} = (x_1, \dots, x_5)$

suppose real data that observe $\vec{x}^{\text{obs}} = (3, 8, 2, 4, 5)$

→ want to estimate skewness $\text{skew}(x_i) = \frac{\mathbb{E}[(x_i - \mu)^3]}{(G^2)^{3/2}}$ third moment

$$\begin{aligned} \text{sample skewness} &\Leftrightarrow \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3 \\ &= \frac{(\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2})^3}{\left(\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \right)^3} \end{aligned}$$

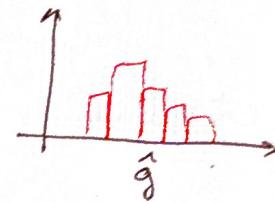
⇒ create bootstrapped datasets just drawing from \vec{X}^{obs}

$$\vec{X}_{\text{boot},1} = (2, 5, 4, 4, 4) \quad \text{sampling with replacement}$$

$$\vec{X}_{\text{boot},2} = (2, 4, 2, 8, 8)$$

: pretend these are equivalent to datasets that observed since you could have observed those numbers

$$\Rightarrow \hat{g}_1 = \hat{g}(\vec{X}_{\text{boot},1}), \hat{g}_2 = \hat{g}(\vec{X}_{\text{boot},2}), \text{ etc}$$



FITTING DISTRIBUTIONS

REGRESSION: fitting a function $E(y|x) = g(x; \theta)$
for mean relation between y and x

- Basic approaches: OLS, GLS, WLS, ML

real problems need
more complicated modelling

Linear regression fitting function is linear

- cepheids → lightcurve is time series

want to find relation of period-luminosity

$$Y = \underline{\beta} X + \underline{\epsilon} \quad E(\epsilon_i) = 0 \quad \text{Var}(\epsilon_i) = G \text{ is known}$$

X matrix

$$\text{minimise w.r.t } \underline{\beta} : \text{RSS} = (Y - \underline{\beta})^T (Y - \underline{\beta})$$

↑
residual sum of squares

$$\text{take derivative w.r.t } \underline{\beta} \Rightarrow \hat{\beta}_{\text{OLS}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T Y$$

If all assumptions are true can show $E(\hat{\beta}_{\text{OLS}}) = \underline{\beta}$ so unbiased

$$\hat{\sigma}^2 = \frac{1}{N-k} (Y - \underline{\beta})^T (Y - \underline{\beta})$$

Weighted Least Squares - $\hat{\beta}$ minimization

for heteroscedastic errors

$$\hat{Y} = \hat{\Xi}/\hat{\beta} + \hat{\varepsilon}$$

\uparrow
 $N \times k$ matrix

$$\mathbb{E}(\varepsilon_i) = 0$$

 $i=1, \dots, N$

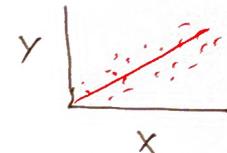
$$X^2 = \sum_{i=1}^N \frac{(y_i - \beta_0 - \sum \beta_j x_{ij})^2}{G_i^{-2}}$$

$X^2 \sim \chi_{N-k}^2$ χ^2 random variable with $N-k$ degrees of freedom

$$\mathbb{E}(X^2) = N-k \Rightarrow \frac{X^2}{N-k} \approx 1 \text{ for large } N-k$$

i.e. if have enough data

- e.g.: even for 0 measurement error there could be intrinsic dispersion



- Special cases: correlated error $\text{Var}(\varepsilon) = \text{Cov}[\varepsilon, \varepsilon^\top] = \underline{W}$ known

$$\Rightarrow \text{RSS} = (\hat{Y} - \hat{\Xi}/\hat{\beta})^\top \underline{W}^{-1} (\hat{Y} - \hat{\Xi}/\hat{\beta})$$

can think in terms of ML! assuming Gaussian with covariance matrix \underline{W}

$$L(\hat{\beta}) = P(Y | \hat{\Xi}, \hat{\beta}) = N(Y | \hat{\Xi}/\hat{\beta}, \underline{W})$$

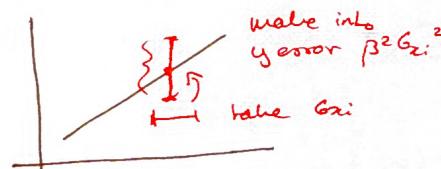
$$\uparrow$$

if $Y \sim N(\hat{\Xi}/\hat{\beta}, \underline{W})$

maximise
⇒ get previous results

- Example: TITEXY estimator idea

$$\text{physicalist } X^2 = \sum \frac{(y_i - \alpha - \beta x_i)^2}{G_{y,i}^{-2}}$$



Tremaine et al. 2002 add measurement error

G^2 to denominator as well

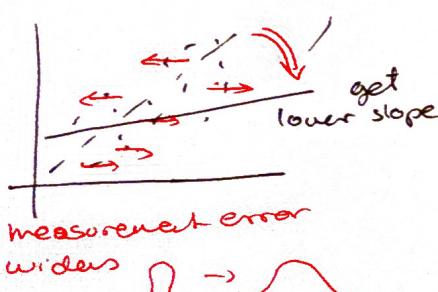
$$\frac{1}{G^2 + G_{y,i}^{-2} + \beta^2 G_{x,i}^{-2}}$$

Bad
bc could make G^2 huge to min.

Hely et al. 2017

⇒ OLS always underestimates
TITEXY always overestimates

But MLE is unbiased



Probabilistic Generative Modelling

Start with conceptual framework: FORWARD MODEL

can introduce latent variables & along the way integrate

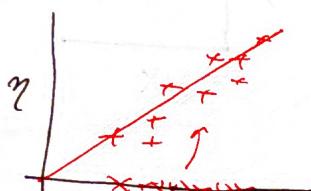
⇒ get sampling distribution $P(D|G) = \int P(D|\alpha) P(\alpha|G) d\alpha$

could integrate out "things that didn't see"

measured model

⇒ using observed D , draw inference from $L(\theta) = P(D|G)$

- quasars:



① generate slope

and map ξ to them

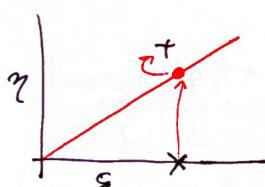
② add errors and get y

with model

③ generate 100 draws

$$\xi \sim N(\mu, \Sigma^2)$$

for one:

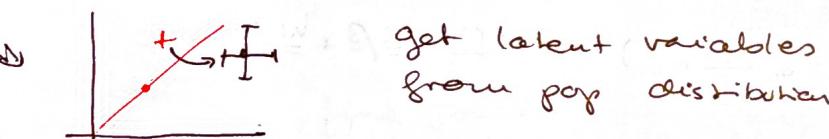


$$\xi \sim N(\mu, \Sigma^2)$$

$$\varepsilon \sim N(0, \sigma^2)$$

$$P(y_i | \xi_i, \alpha, \beta, G, \mu, \Sigma)$$

$$= P(y_i | \xi_i, \alpha, \beta, G) \cdot P(\xi_i | \mu, \Sigma)$$



get latent variables
from pop distribution

Now get observed data: $P((x_i, y_i) | \gamma_i, \xi_i) = N((\bar{x}_i, \bar{y}_i) | (\gamma_i, \xi_i), \Sigma)$

$$\Sigma = \begin{pmatrix} \sigma_{yy}^2 & G \\ 0 & \sigma_{xx}^2 \end{pmatrix}$$

$$\underline{\Theta} = (\alpha, \beta, G^2) \quad \underline{\Psi} = (\mu, \Sigma) \quad (\xi_i, \gamma_i) \quad (x_i, y_i)$$

reg. param.

indep variable
"hyperparam"

latent
(true)
variables

data

$$\textcircled{1} \quad \xi \sim N(\mu, \Sigma^2)$$

$$\textcircled{2} \quad y_i | \xi_i$$

⇒ formulate likelihood function
by integrating out latent variables