# Unified Analyses of Populations of Sources
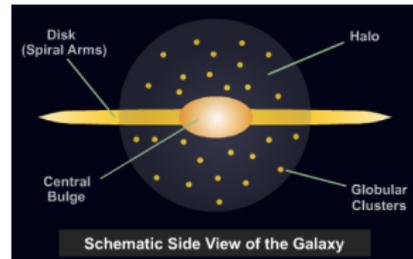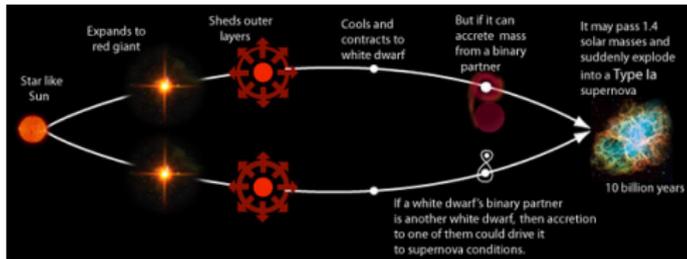## Advantages of "Shrinkage Estimates" in Astronomy

### David A. van Dyk

Statistics Section, Imperial College London

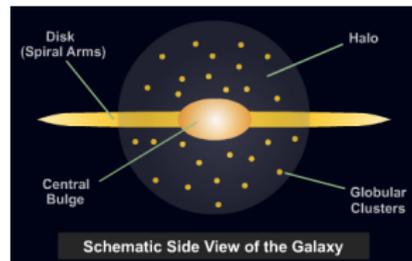University of Birmingham
April 2015

Imperial College
London

## Populations of Sources
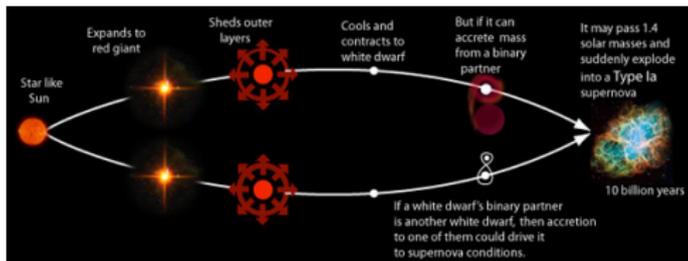


Estimating a property of each object in a population:

1. Intrinsic (absolute) magnitudes of Type Ia Super Novae.
   Or more simply: apparent magnitudes.
2. The ages of White Dwarfs in the galactic halo.
   Or more simply: ages of WDs in the galaxy.
3. Measured distance to Large Magellanic Cloud
   * With different methods, each with their own systematics.

**Imperial College**
London

# Estimating Source Characteristics



**Typical Strategy:** Estimate the magnitude, distance, or age for each source in a separate data analysis.

**Another Possibility:** Preform unified analysis, modeling dist'n of magnitudes, distances, or ages among sources.

- Relative advantages depends on *quality of individual estimates* and *degree of homogeneity* in population.
- Discuss from Frequntist and Bayesian perspectives.

Imperial College
London

**All Roads Lead to Rome**
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## Outline

**Imperial College**
London

David A. van Dyk    Unified Analyses of Populations of Sources

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## Outline

1. **All Roads Lead to Rome**
   - Frequentist origins of shrinkage estimates
   - Bayesian hierarchical models

2. Example 1: Using SNIa to Fit Cosmological Models
   - Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

3. Example 2: Ages of White Dwarfs in the Galactic Halo
   - Joint work with Ted von Hippel & Shijing Si

**Imperial College**
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## The Sample Mean

Suppose we wish to estimate a parameter, $\theta$, from repeated measurement or a single source:

$$y_i \stackrel{\text{indep}}{\sim} \mathsf{N}(\theta, \sigma^2) \quad \text{for} \quad i = 1, \ldots, n$$

Eg: calibrating detector from *n* measures of known source.

An obvious estimator:

$$\hat{\theta}^{\text{naive}} = \frac{1}{n} \sum_{i=1}^{n} y_i$$

*What is not to like about the arithmetic average?*

Imperial College
London

David A. van Dyk    Unified Analyses of Populations of Sources

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## Frequency Evaluation of an Estimator

- How far off is the estimator?

$$(\hat{\theta} - \theta)^2$$

- How far off do we expect it to be?

$$\mathrm{MSE}(\hat{\theta}|\theta) = \mathrm{E}\left[(\hat{\theta} - \theta)^2 \mid \theta\right] = \int \left(\hat{\theta}(y) - \theta\right)^2 f(y \mid \theta) dy$$

- This quantity is called the Mean Square Error of $\hat{\theta}$.
- An estimator is said to be inadmissible if there is an estimator that is uniformly better in terms of $\mathrm{MSE}$:

$$\mathrm{MSE}(\hat{\theta}|\theta) < \mathrm{MSE}(\hat{\theta}^{\mathrm{naive}}|\theta) \text{ for all } \theta.$$

**Imperial College London**

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## Inadmissibility of the Sample Mean

Suppose we wish to estimate more than one parameter:

$$y_{ij} \stackrel{\text{indep}}{\sim} N(\theta_j, \sigma^2) \text{ for } i = 1, \ldots, n \text{ and } j = 1, \ldots, G$$

The obvious estimator:

$$\hat{\theta}_j^{\text{naive}} = \frac{1}{n} \sum_{i=1}^{n} y_{ij} \quad \text{is inadmissible if } G \geqslant 3.$$

The James-Stein Estimator dominates $\theta^{\text{naive}}$:

$$\hat{\theta}_j^{\text{JS}} = \left(1 - \omega^{\text{JS}}\right) \hat{\theta}_j^{\text{naive}} + \omega^{\text{JS}} \nu \text{ for any } \nu$$
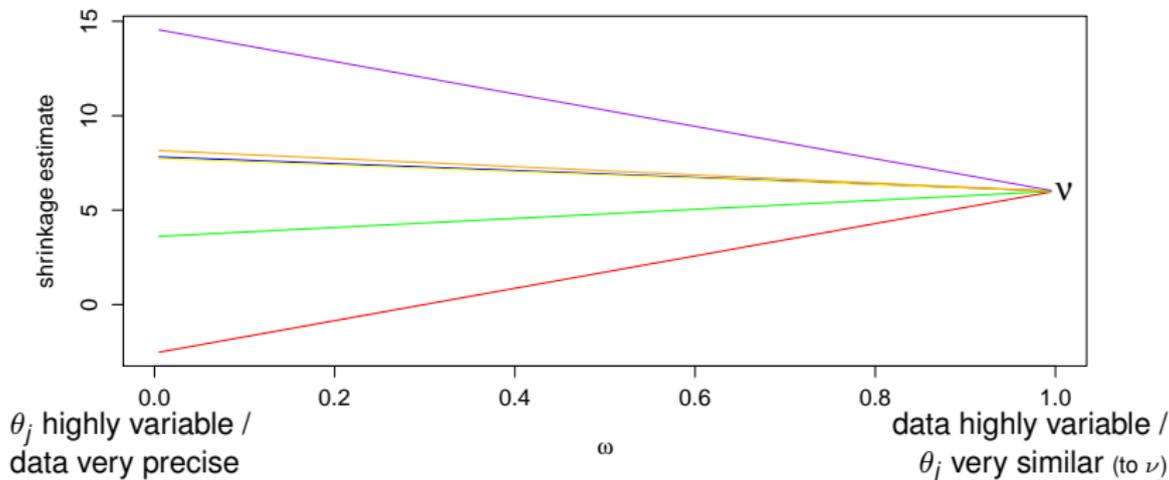
with $\omega^{\text{JS}} \approx \dfrac{\sigma^2/n}{\sigma^2/n + \tau_\nu^2}$ and $\tau_\nu^2 = \mathrm{E}[(\theta_j - \nu)^2]$.

**Imperial College** London

Specifically, $\omega^{\text{JS}} = (G - 2)\sigma^2 / n \sum_{j=1}^{G} (\hat{\theta}_j^{\text{naive}} - \nu)^2$.

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

# Shrinkage Estimators

**James-Stein Estimator** is a shrinkage estimator:

$$\hat{\theta}_j^{\text{JS}} = \left(1 - \omega^{\text{JS}}\right)\hat{\theta}_j^{\text{naive}} + \omega^{\text{JS}}\nu$$



$\theta_j$ highly variable /
data very precise

data highly variable /
$\theta_j$ very similar (to $\nu$)

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## To Whence To Shrink?

James-Stein Estimators

- Dominate the sample average for *any choice* of $\nu$.
- Shrinkage is mild and $\hat{\theta}^{\mathrm{JS}} \approx \hat{\theta}^{\mathrm{naive}}$ for most $\nu$.
- Can we choose $\nu$ to maximize shrinkage?

$$\hat{\theta}_j^{\mathrm{JS}} = \left(1 - \omega^{\mathrm{JS}}\right) \hat{\theta}_j^{\mathrm{naive}} + \omega^{\mathrm{JS}} \nu$$

with $\omega^{\mathrm{JS}} \approx \dfrac{\sigma^2/n}{\sigma^2/n + \tau_\nu^2}$ and $\tau_\nu^2 = \mathrm{E}[(\theta_j - \nu)^2]$.

- Minimize $\tau_\nu^2$.
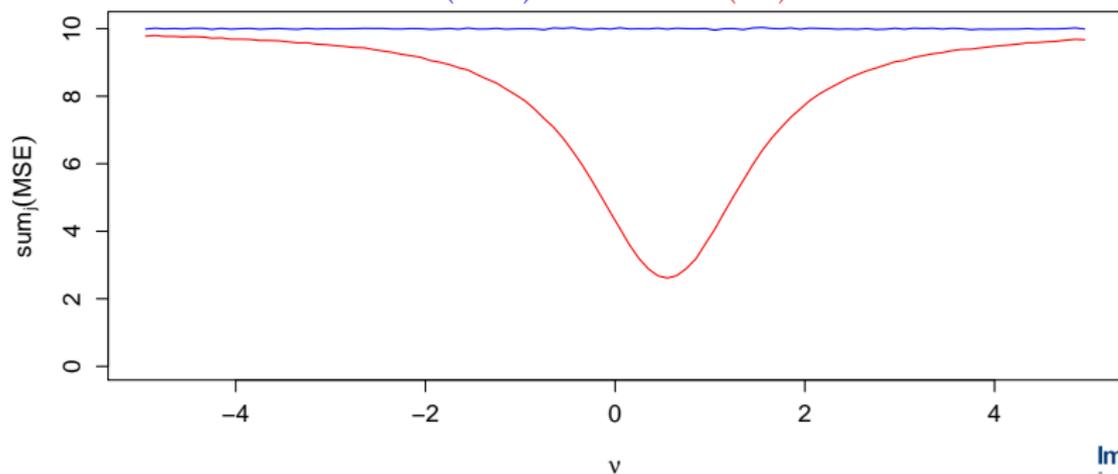
*The optimal choice of $\nu$ is the average of the $\theta_j$.*

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## Illustration

Suppose:

- $y_j \sim N(\theta_j, 1)$ for $j = 1, \ldots, 10$
- $\theta_j$ are evenly distributed on [0,1]

$\mathrm{MSE}(\hat{\theta}^{\mathrm{naive}})$ versus $\mathrm{MSE}(\hat{\theta}^{\mathrm{JS}})$:

Imperial College
London

**All Roads Lead to Rome**
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
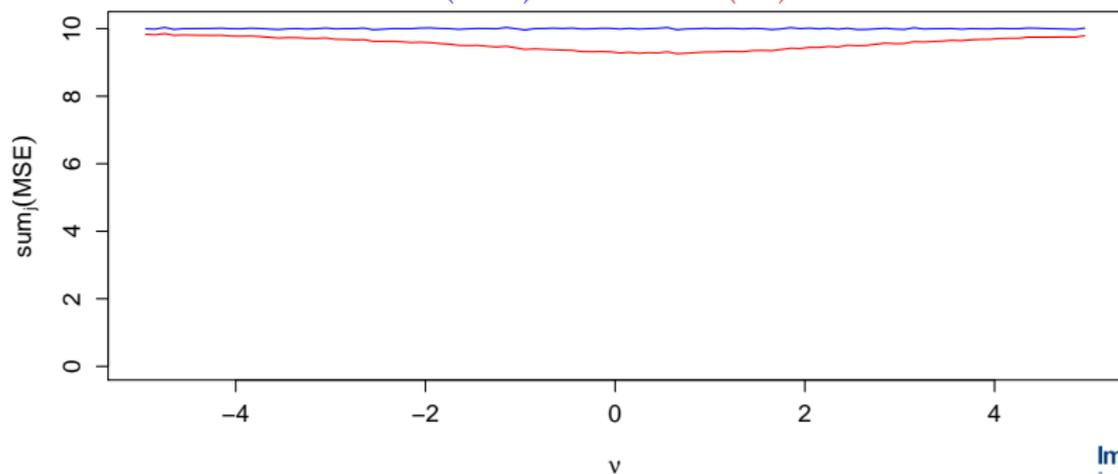Bayesian hierarchical models

## Illustration

Suppose:

- $y_j \sim N(\theta_j, 1)$ for $j = 1, \ldots, 10$
- $\theta_j$ are evenly distributed on **[-4,5]**

$\text{MSE}(\hat{\theta}^{\text{naive}})$ versus $\text{MSE}(\hat{\theta}^{\text{JS}})$:

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## Intuition

1. It you are estimating more than two parameters, it is always better to use shrinkage estimators.

2. Optimally shrink toward average of the parameters.

3. Most gain when the naive (non-shrinkage) estimators
   - $\star$ are noisy ($\sigma^2$ is large)
   - $\star$ are similar ($\tau^2$ is small)

4. Bayesian versus Frequentist:
   - $\star$ From frequentist point of view this is somewhat problematic.
   - $\star$ From a Bayesian point of view this is an opportunity!

5. James-Stein is a milestone in statistical thinking.
   - $\star$ Results viewed as paradoxical and counterintuitive.
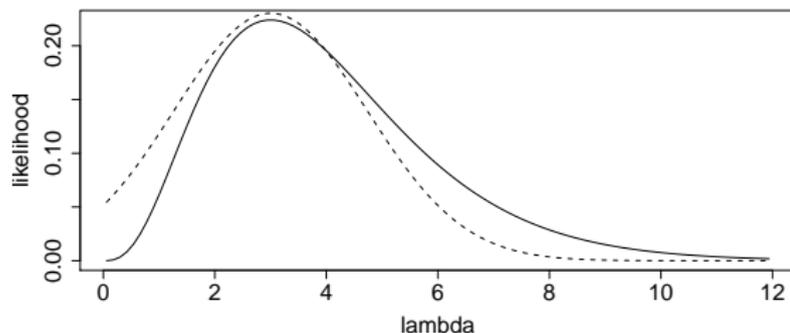   - $\star$ James and Stein are geniuses.

**Imperial College
London**

**All Roads Lead to Rome**
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
**Bayesian hierarchical models**

## Outline

**Imperial College**
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## Bayesian Statistical Analyses: Likelihood

<u>Likelihood Functions:</u> The distribution of the data given the model parameters. E.g., $Y \sim \text{Poisson}(\lambda_S)$:

$$\text{likelihood}(\lambda_S) = e^{-\lambda_S} \lambda_S^Y / Y!$$

<u>Maximum Likelihood Estimation:</u> Suppose $Y = 3$



*The likelihood and its normal approximation.*

*Can estimate $\lambda_S$ and its error bars.*

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
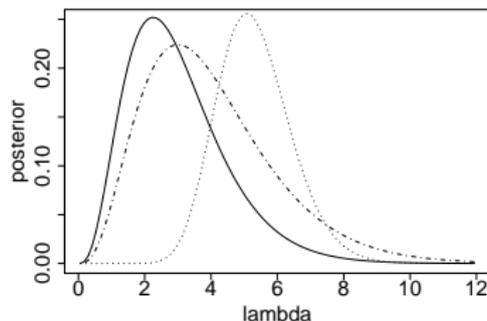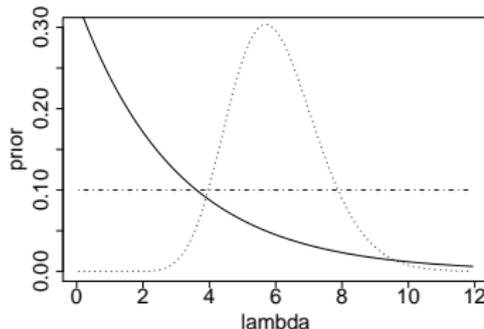Bayesian hierarchical models

## Bayesian Analyses: Prior and Posterior Dist'ns

Prior Distribution: Knowledge obtained *prior* to current data.

Bayes Theorem and Posterior Distribution:

$$\text{posterior}(\lambda \mid Y) \propto \text{likelihood}(\lambda; Y) \times \text{prior}(\lambda)$$

Combine past and current information:



*Bayesian analyses rely on probability theory* Imperial College London

**All Roads Lead to Rome**
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
**Bayesian hierarchical models**

## Bayesian Perspective

Back to shrinkage....

Frequentists tend to avoid quantities like:

1. $E(\theta_j)$ and $\text{Var}(\theta_j)$
2. $E\left[(\theta_j - \nu)^2\right]$

From a Bayesian point of view it is quite natural to consider

1. the *prior* distribution of a parameter or
2. the *common distribution of a group of parameters*.

*Models that are formulated in terms of the latter are*
*Hierarchical Models.*

**Imperial College**
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## A Simple Bayesian Hierarchical Model

Suppose

$$y_{ij}|\theta_j \overset{\text{indep}}{\sim} N(\theta_j, \sigma^2) \text{ for } i = 1, \ldots, n \text{ and } j = 1, \ldots, G$$

with

$$\theta_j \overset{\text{indep}}{\sim} N(\mu, \tau^2).$$

Let $\phi = (\sigma^2, \tau^2, \mu)$

$$E(\theta_j \mid Y, \phi) = (1 - \omega^{\text{HB}})\hat{\theta}^{\text{naive}} + \omega^{\text{HB}}\mu \text{ with } \omega^{\text{HB}} = \frac{\sigma^2/n}{\sigma^2/n + \tau^2}.$$

The Bayesian perspective

- automatically picks the best $\nu$,
- provides model-based estimates of $\phi$,
- requires priors be specified for $\sigma^2, \tau^2$, and $\mu$.

**Imperial College
London**

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

## Color Correction Parameter for SNIa Lightcurves

SNIa light curves vary systematically across color bands.

- Measure how peaked the color distribution is.
- Details in the next section!!
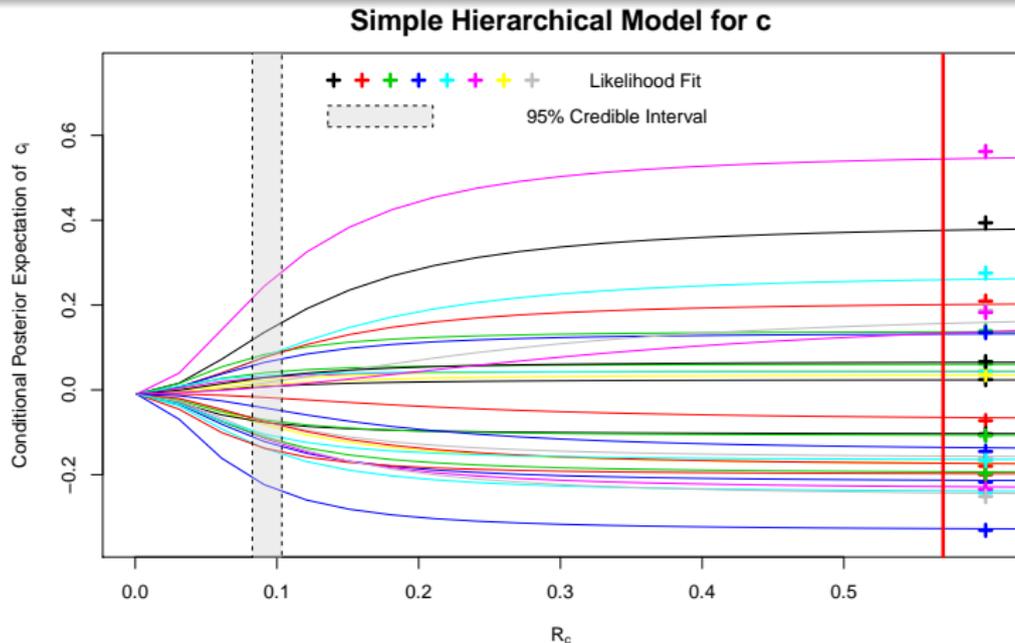- A hierarchical model:

$$\hat{c}_j | c_j \stackrel{\text{indep}}{\sim} N(c_j, \sigma_j^2) \text{ for } j = 1, \ldots, 288$$

with

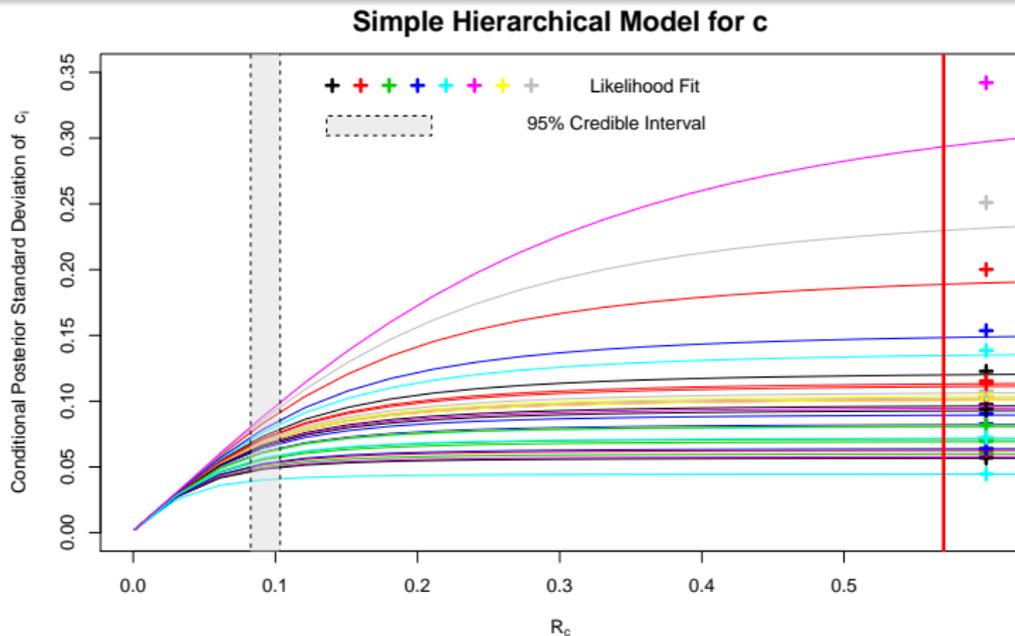$$c_j \stackrel{\text{indep}}{\sim} N(c_0, R_c^2) \text{ and } p(c_0, R_c) \propto 1.$$

- The measurement variances, $\sigma_j^2$ are assumed known.
- We could estimate each $c_j$ via $\hat{c}_j \pm \sigma_j$, or...

Imperial College
London

**All Roads Lead to Rome**
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
**Bayesian hierarchical models**

# Shrinkage of the Fitted Treatment Effects



**Simple Hierarchical Model for c**

*Pooling may dramatically change fits.*

**Imperial College**
London

**All Roads Lead to Rome**
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

# Standard Deviation of the Fitted Treatment Effects



**Simple Hierarchical Model for c**

*Borrowing strength for more precise estimates.*

**Imperial College**
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

# The Bayesian Perspective

**Advantages of Bayesian Perspective:**

- The advantage of James-Stein estimation is automatic.
  James-Stein had to find the estimator!

- Bayesians have a method to generate estimators.
  Even frequentists like this!

- General principle is easily tailored to any problem.

- Specification of level two model *may* not be critical.
  James-Stein derived same estimator using only moments.

**Cautions:**

- Results can depend on prior distributions for parameters
  that reside deep within the model, and far from the data.

**Imperial College**
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Frequentist origins of shrinkage estimates
Bayesian hierarchical models

# The Choice of Prior Distribution

Suppose

$$y_{ij}|\theta_j \overset{\text{indep}}{\sim} N(\theta_j, \sigma^2) \text{ for } i = 1, \ldots, n \text{ and } j = 1, \ldots, G$$

with

$$\theta_j \overset{\text{indep}}{\sim} N(\mu, \tau^2).$$

- Std non-informative prior for normal variance: $p(\sigma^2) \propto 1/\sigma^2$.
- Using this prior for the level-two variance,

$$p(\tau^2) \propto 1/\tau^2$$

leads to an improper posterior distribution:

$$p(\tau^2|y) \propto p(\tau^2) \sqrt{\frac{\text{Var}(\mu|y, \tau)}{(\sigma^2 + \tau^2)^G}} \exp\left\{ \sum_{j=1}^{G} -\frac{(\bar{y}_{.j} - \text{E}(\mu|y, \tau^2))^2}{2(\sigma^2 + \tau^2)} \right\}$$

**Imperial College
London**

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## Outline

**Imperial College London**

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## Outline

**Imperial College**
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## Type Ia Supernovae as Standardizable Candles

If mass surpasses "Chandrasekhar threshold" of $1.44 M_\odot$...



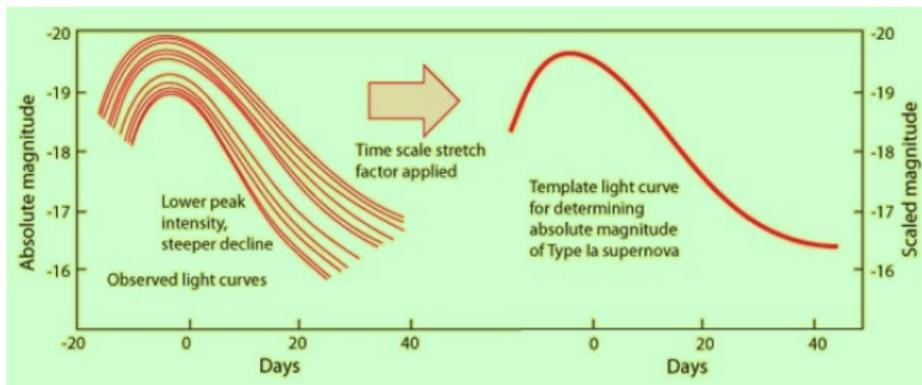Image Credit: http://hyperphysics.phy-astr.gsu.edu/hbase/astro/snovcn.html

Due to their common "flashpoint", SN1a have similar absolute magnitudes:

$$M_j \sim \mathrm{N}(M_0, \sigma_{\mathrm{int}}^2).$$

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## Predicting Absolute Magnitude

SN1a absolute magnitudes are correlated with characteristics of the explosion / light curve:

- $x_j$: rescale light curve to match mean template
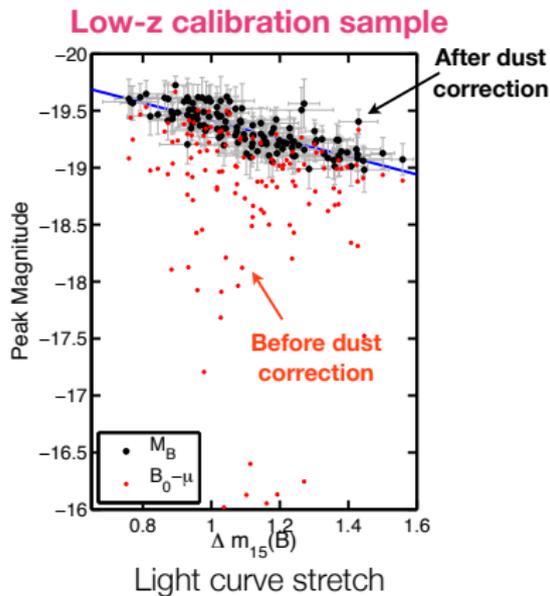- $c_j$: describes how flux depends on color (spectrum)



Credit: http://hyperphysics.phy-astr.gsu.edu/hbase/astro/snovcn.html

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## Phillips Corrections

- Recall: $M_j \sim N(M_0, \sigma_{int}^2)$.

- Regression Model:
  $$M_j = -\alpha x_j + \beta c_j + M_j^\epsilon,$$
  with $M_j^\epsilon \sim N(M_0, \sigma_\epsilon^2)$.

- $\sigma_\epsilon^2 \leqslant \sigma_{int}^2$

- Including $x_i$ and $c_i$ reduces variance and increases precision of estimates.

**Low-z calibration sample**



Mandel et al (2011)

Light curve stretch

*Brighter SNIa are slower decliners over time.*

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## Distance Modulus in an Expanding Universe

Apparent mag depends on absolute mag & distance modulus:

$$m_{Bj} = \mu_j + M_j = \mu_j + M_j^{\epsilon} - \alpha x_j + \beta c_j$$

Relationship between $\mu_i$ and $z_i$

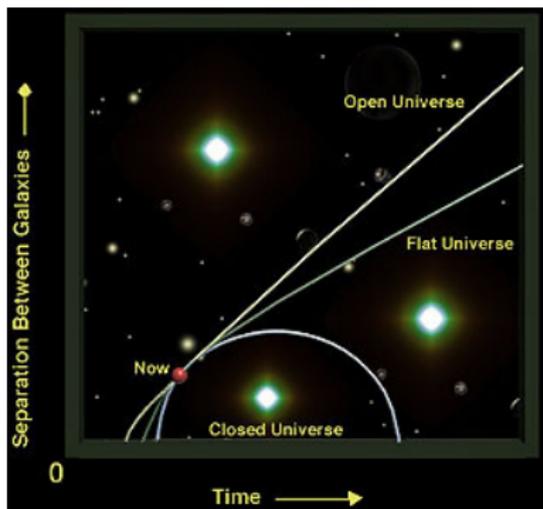- For nearby objects,

  $$z_j = \text{velocity}/c$$
  $$\text{velocity} = H_0 \text{ distance}.$$

  (Correcting for peculiar/local velocities.)

- For distant objects, involves expansion history of Universe:

  $$\mu_j = g(z_j, \Omega_\Lambda, \Omega_M, H_0)$$
  $$= 5 \log_{10}(\text{distance[Mpc]}) + 25$$

- We use peak B band magnitudes.



http://skyserver.sdss.org/dr1/en/astro/universe/universe.asp

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

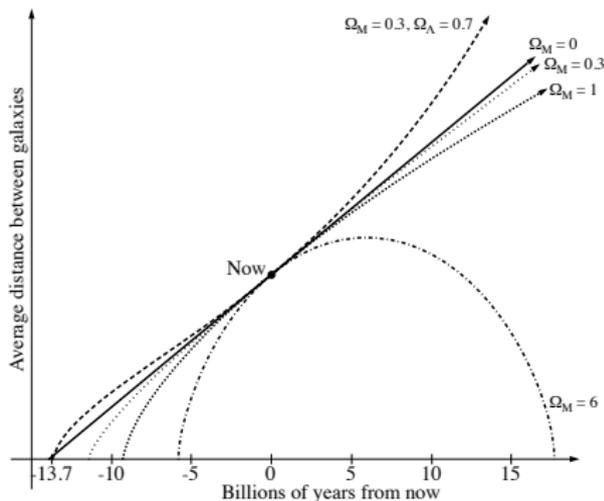## Accelerating Expansion of the Universe

- 2011 Physics Nobel Prize: discovery that expansion rate is increasing.

- Dark Energy is the principle theorized explanation of accelerated expansion.

- $\Omega_\Lambda$: density of dark energy (describes acceleration).

- $\Omega_M$: total matter.

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## A Hierarchical Model

**Level 1:** $c_j$, $x_j$, and $m_{Bj}$ are observed with error.

$$\begin{pmatrix} \hat{c}_j \\ \hat{x}_j \\ \hat{m}_{Bj} \end{pmatrix} \sim N \left\{ \begin{pmatrix} c_j \\ x_i \\ m_{Bj} \end{pmatrix}, \ \hat{C}_j \right\}$$

with $m_{Bj} = \mu_j + M_j^\epsilon - \alpha x_j + \beta c_j$ and $\mu_j = g(z_j, \Omega_\Lambda, \Omega_M, H_0)$

**Level 2:**

1. $c_j \sim N(c_0, R_c^2)$
2. $x_j \sim N(x_0, R_x^2)$
3. $M_j^\epsilon \sim N(M_0, \sigma_\epsilon^2)$

**Level 3:** Priors on $\alpha$, $\beta$, $\Omega_\Lambda$, $\Omega_M$, $H_0$, $c_0$, $R_c^2$, $x_0$, $R_x^2$ $M_0$, $\sigma_\epsilon^2$

**Imperial College London**

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

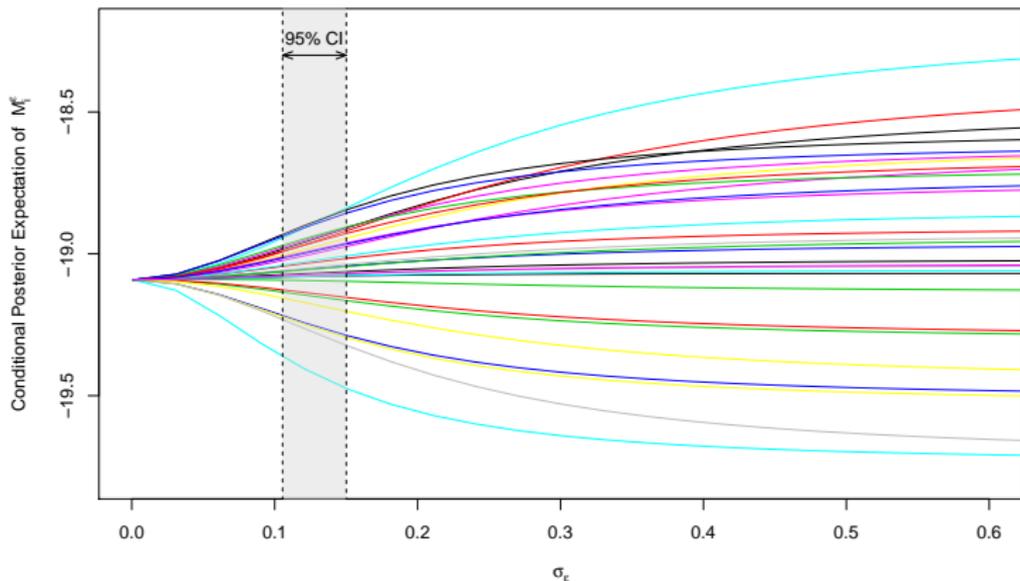Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## Other Model Features

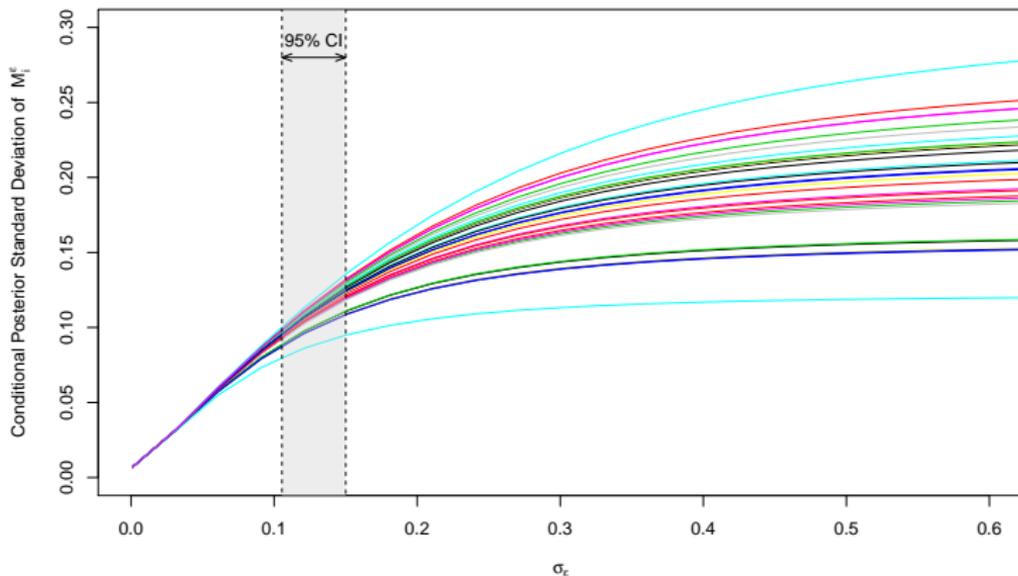Results are based on an SDSS (2009) sample of 288 SNIa.

In our full analysis, we also

1. account for systematic errors that have the effect of correlating observation across supernovae,

2. allow the mean and variance of $M_i^\epsilon$ to differ for galaxies with stellar masses above or below $10^{10}$ solar masses,

3. include a model component that adjusts for selection effects, and

4. use a larger JLA sample[1] of 740 SNIa observed with SDSS, HST, and SNLS.

**Imperial College London**

[1]Betoule, et al., 2014, arXiv:1401.4064v1

David A. van Dyk    Unified Analyses of Populations of Sources

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

# Shrinkage Estimates in Hierarchical Model

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

# Shrinkage Errors in Hierarchical Model

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## Fitting Absolute Magnitudes Without Shrinkage

Under the model, absolute magnitudes are given by

$$M_j^\epsilon = m_{Bj} - \mu_j + \alpha x_j - \beta c_j \text{ with } \mu_i = g(z_j, \Omega_\Lambda, \Omega_M, H_0)$$

Setting

1. $\alpha, \beta, \Omega_\Lambda$, and $\Omega_M$ to their minimum $\chi^2$ estimates,
2. $H_0 = 72 km/s/Mpc$, and
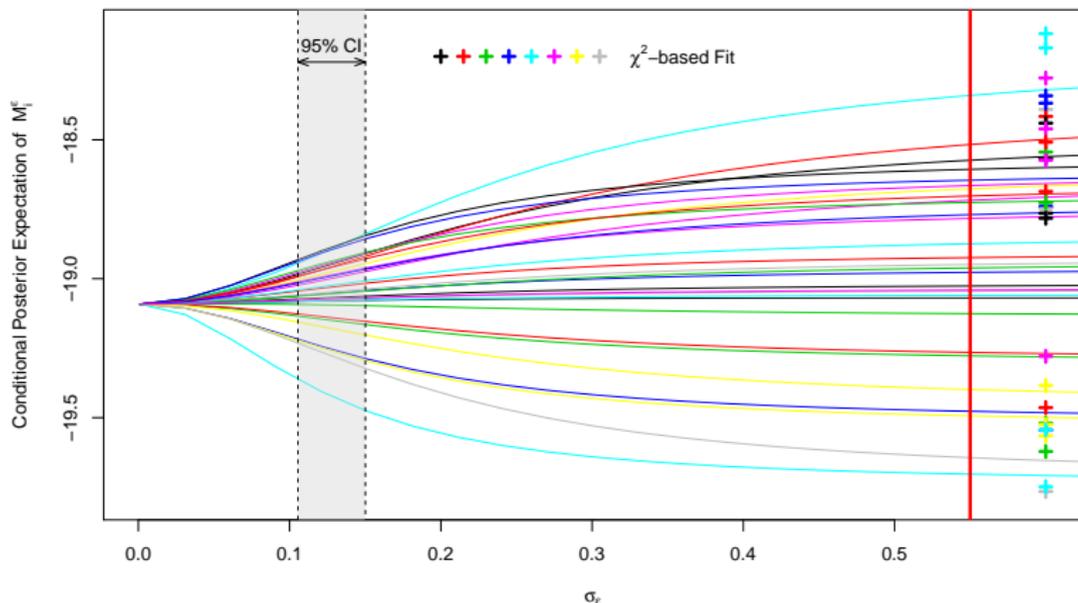3. $m_{Bj}, x_j$, and $c_j$ to their observed values

we have

$$\hat{M}_j^\epsilon = \hat{m}_{Bi} - g(\hat{z}_j, \hat{\Omega}_\Lambda, \hat{\Omega}_M, \hat{H}_0) + \hat{\alpha}\hat{x}_j - \hat{\beta}\hat{c}_j$$
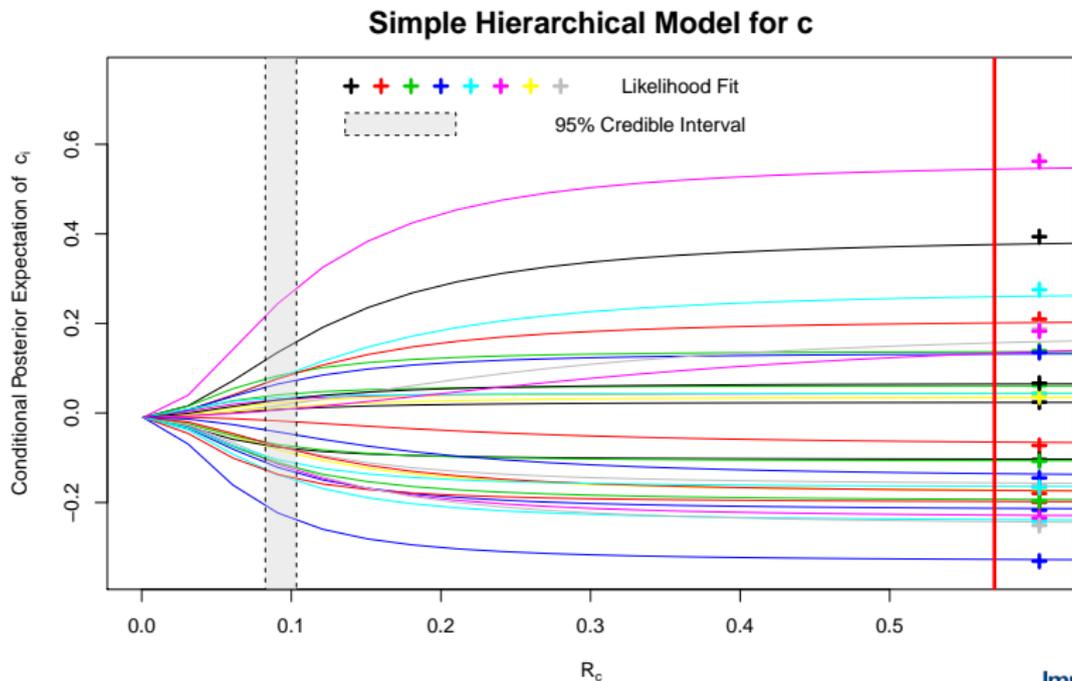
with error

$$\approx \sqrt{\text{Var}(\hat{m}_{Bj}) + \hat{\alpha}^2 \text{Var}(\hat{x}_j) + \hat{\beta}^2 \text{Var}(\hat{c}_j)}$$
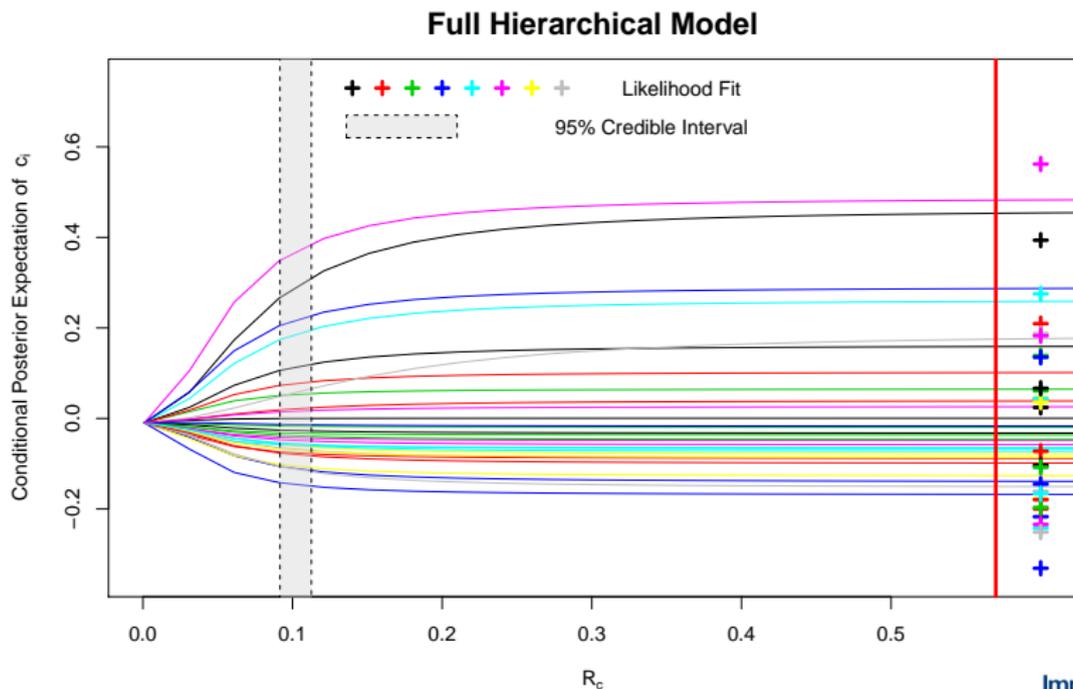
**Imperial College London**
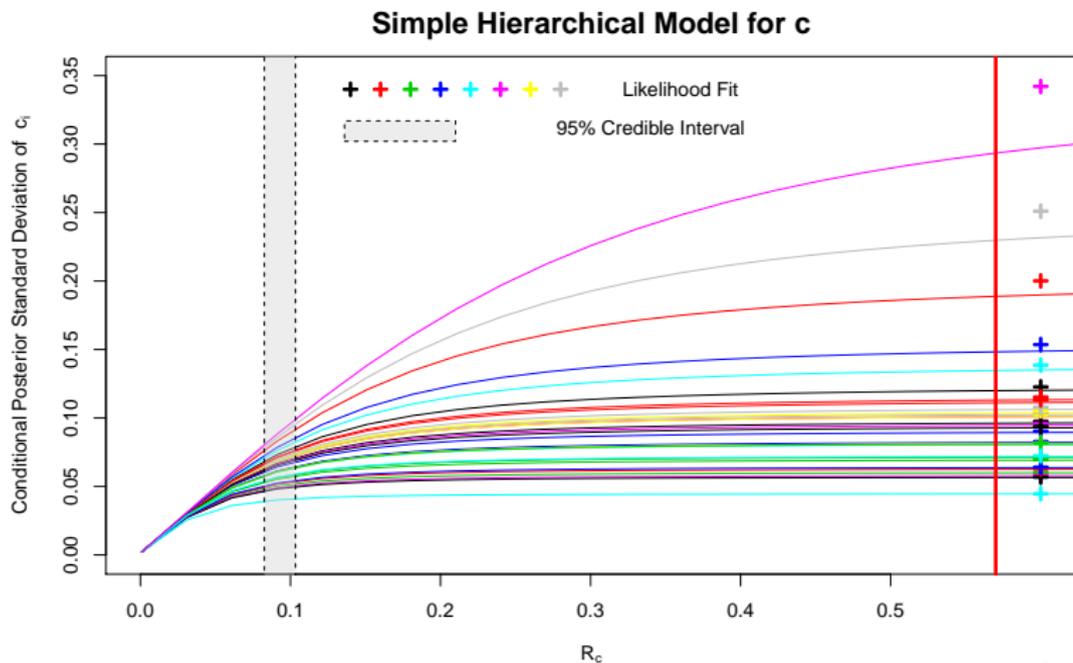
All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

# Comparing the Estimates

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

## Comparing the Estimates



*Offset estimates even without shrinkage.*

# Fitting a simple hierarchical model for $c_i$



Simple Hierarchical Model for c

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

# Additional shrinkage due to regression



Full Hierarchical Model

Imperial College
London

# Errors under simple hierarchical model for $c_i$



Simple Hierarchical Model for c

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

# Reduced errors due to regression

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint with Roberto Trotta, Xiyun Jiao, & Hikmatali Shariff

# Comparing the Estimates of $c_i$ and $x_i$

**Imperial College**
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

## Outline

**Imperial College London**

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

## Outline

**Imperial College London**

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Visitors from the Galactic Halo

- Age of galactic *halo* or *disk* can be estimated with their older stars.

- Halo stars pass through the galactic disk as they orbit the central bulge.



- Kilic et al. (ApJ, 2010) identified three nearby old halo white dwarfs in the SDSS; we have a sample of five.

*We would like to model the white dwarf colors to estimate their age and the age of galactic halo.*

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

## Fitting Dist'n of Stellar Ages in Galactic Halo

We observe seven photometric magnitudes for each WD:

$$(X_{1j}, \ldots X_{7j}) \sim \mathrm{MVN}\Big( G(\theta_j), V \Big)$$

where $\theta_j = (\log_{10}(\mathrm{age}_j), \mathrm{distance}_j, \mathrm{mass}_j)$ and

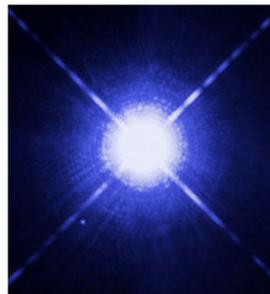$$\log_{10}(\mathrm{age}_j) \sim \mathrm{N}(\mu, \tau^2).$$

- If the WD are a representative sample, $\mu$ and $\tau^2$ are the population mean and variance for galactic halo.
- Even if sample is not representative, hierarchical model produces estimators with better statistical properties.

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

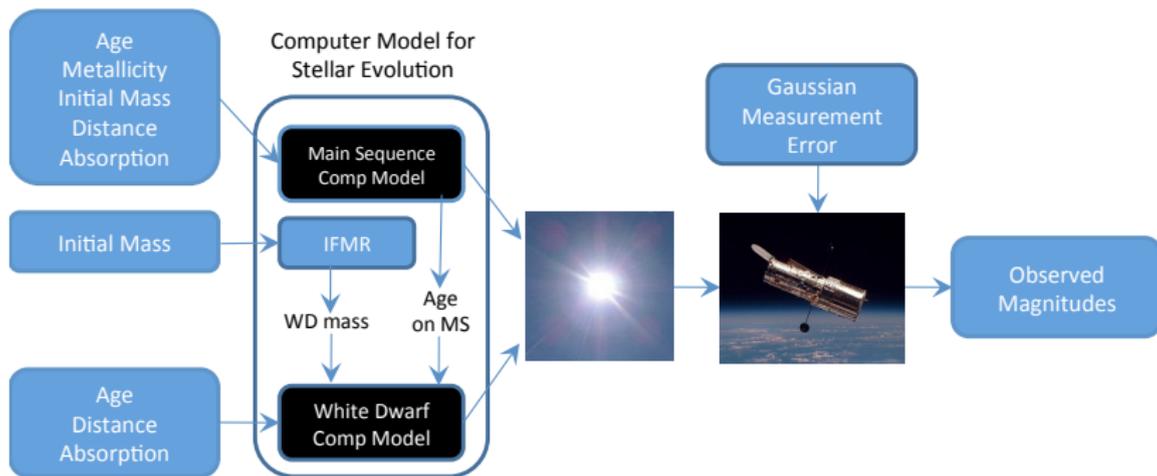# Computer Model for Main Sequence (& RG) Evolution



- Computer model predicts how the emergent and apparent spectra evolve as a function of input parameters.
- We observe photometric magnitudes, the apparent luminosity in each of several wide wavelength bands.

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

## White Dwarfs Physics



- White dwarf spectra are not predicted from MS/RG models
- Different physical processes require different models:
  1. Computer Model for White Dwarf Cooling
  2. Computer Model for White Dwarf Atmosphere
  3. Initial Final Mass Relationship (IFMR)

**Imperial College London**

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

*A parametric model for the IFMR forms a bridge between the computer models.*

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Complex Posterior Distributions

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

## Complex Posterior Distributions

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Fitting Each WD Individually (Kilic's sample)

*Posterior distributions exhibit similar structure
and similar fitted parameter values.*

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

## Fitting the Population Distribution of Halo WDs

**Model:**

$$(X_{1j}, \ldots X_{7j}) \sim \mathrm{MVN}\Big( G(\theta_j), V \Big)$$

where $\theta_j = (\log_{10}(\mathrm{age}_j), \ \mathrm{distance}_j, \ \mathrm{mass}_j)$ and

$$\log_{10}(\mathrm{age}_j) \sim \mathrm{N}(\mu, \tau^2).$$

**Maximum a posterior estimates:**

- $\hat{\mu} = 10.065$ (11.6 gigayears)
- $\widehat{\log_{10} \tau} = \log_{10}(0.053)$
- 95% range: $(9.1, 14.8)$ gigayears

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si
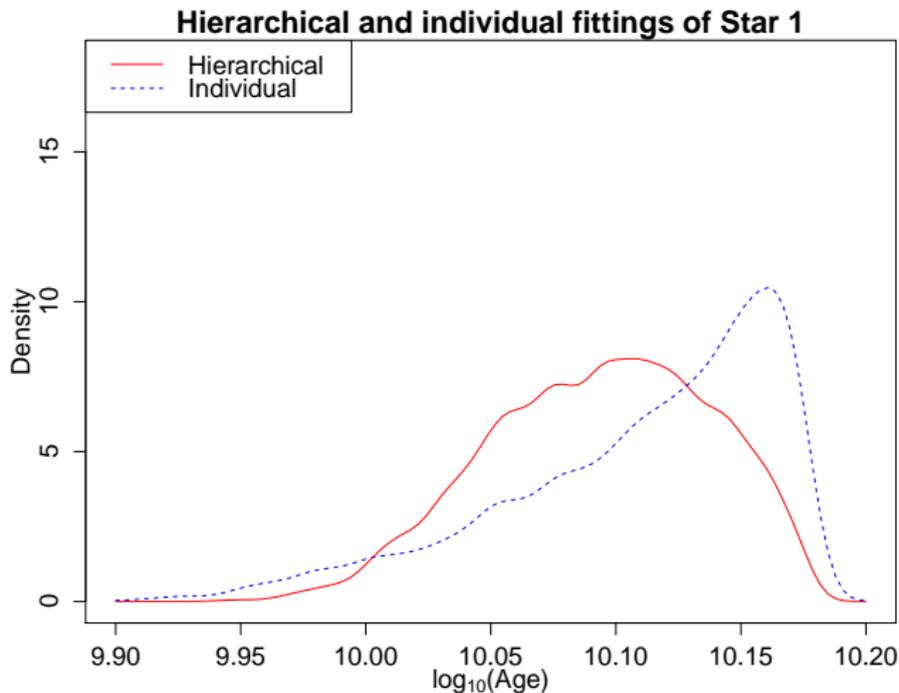
# Suppose $\mathrm{sd}(\log_{10}(\mathrm{age})) = 0.009$

Individual Fit

Hierarchical Fit

*Effect of Shrinkage for one halo WD.*

*Here we exaggerate the shrinkage by using $\tau = 0.009 < \hat{\tau}$.*

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Shrinkage in the Posterior of Age (with fitted $\tau$)



**Hierarchical and individual fittings of Star 1**

Legend:
- Hierarchical (red, solid)
- Individual (blue, dotted)

Y-axis: Density
X-axis: $\log_{10}(\text{Age})$

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Shrinkage in the Posterior of Age (with fitted $\tau$)



**Hierarchical and individual fittings of Star 5**

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Shrinkage in the Posterior of Age (with fitted $\tau$)



**Hierarchical and individual fittings of all Stars**

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Sensitivity of Results to $\mathrm{Var}(\log_{10}(\mathrm{age}))$



Ages of WDs from the Galactic Halo

Imperial College London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Sensitivity of Results to $\mathrm{Var}(\log_{10}(\mathrm{age}))$



*Gaia will provide photometric magnitudes for hundreds of galactic halo WDs.*

**Imperial College**
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

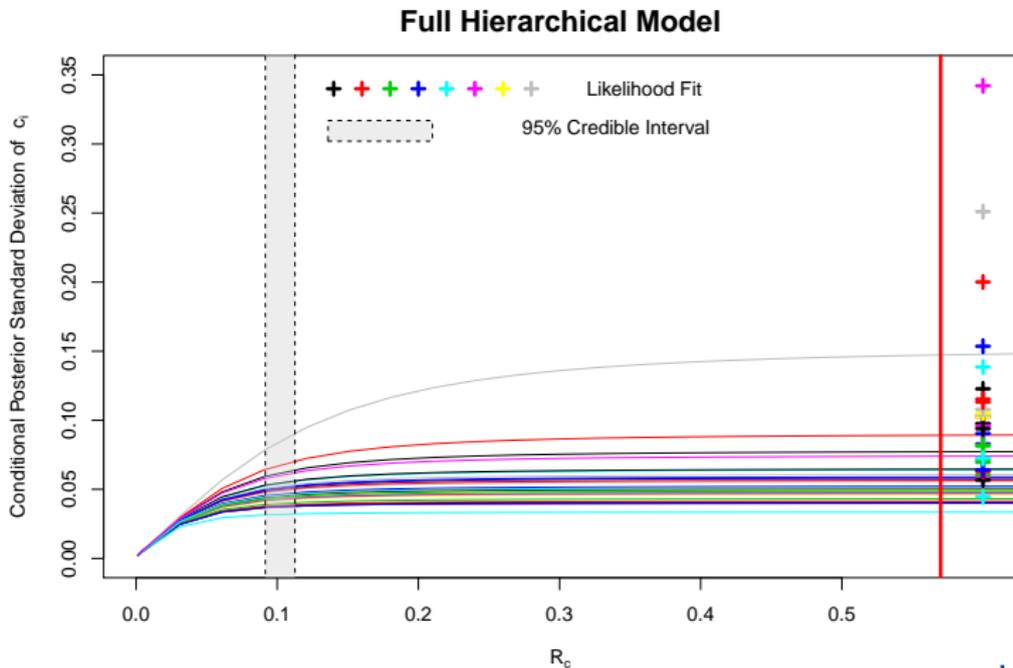Joint work with Ted von Hippel & Shijing Si

## Discussion

- Estimation of groups of parameters describing populations of sources is not uncommon in astronomy.
- These parameters may or may not be of primary interest.
- Modeling the distribution of object-specific parameters can dramatically reduce both error bars and MSE ...
- ... especially with noisy observations of similar objects.
- Shrinkage estimators are able to "borrow strength".
- May be little cost of freeing object-specific parameters (e.g., metallicity or distance of stars in a cluster).

*Don't throw away half of your toolkit!!*
*(Bayesian and Frequency methods)*

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Shrinkage Estimates of $c_i$ in Hierarchical Model



**Full Hierarchical Model**

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Shrinkage Errors of $c_i$ in Hierarchical Model



**Full Hierarchical Model**

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Shrinkage Estimates of $x_i$ in Hierarchical Model

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Shrinkage Errors of $x_i$ in Hierarchical Model

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

## A Non-Astronomical Example

The Educational Testing Service studied the effects of coaching programs on SAT-V scores in eight US high schools:[2]

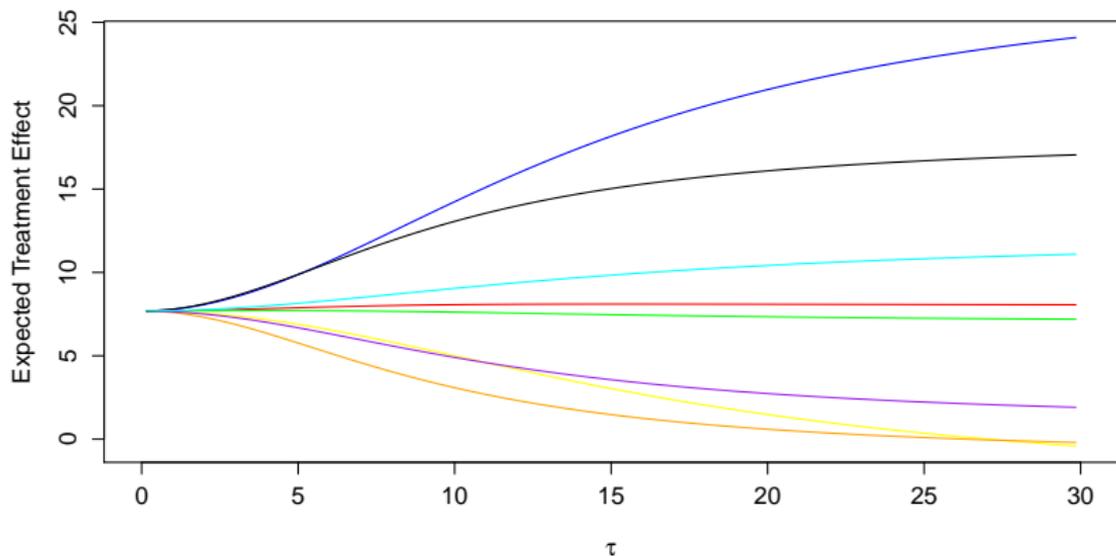$$y_j | \theta_j \overset{\text{indep}}{\sim} N(\theta_j, \sigma_j^2) \text{ for } j = 1, \ldots, 8$$

with

$$\theta_j \overset{\text{indep}}{\sim} N(\mu, \tau^2) \text{ and } p(\mu, \tau) \propto 1.$$

The $y_j$ are estimated treatment effects

- based on preliminary analyses
- adjust for PSAT (V and M) scores
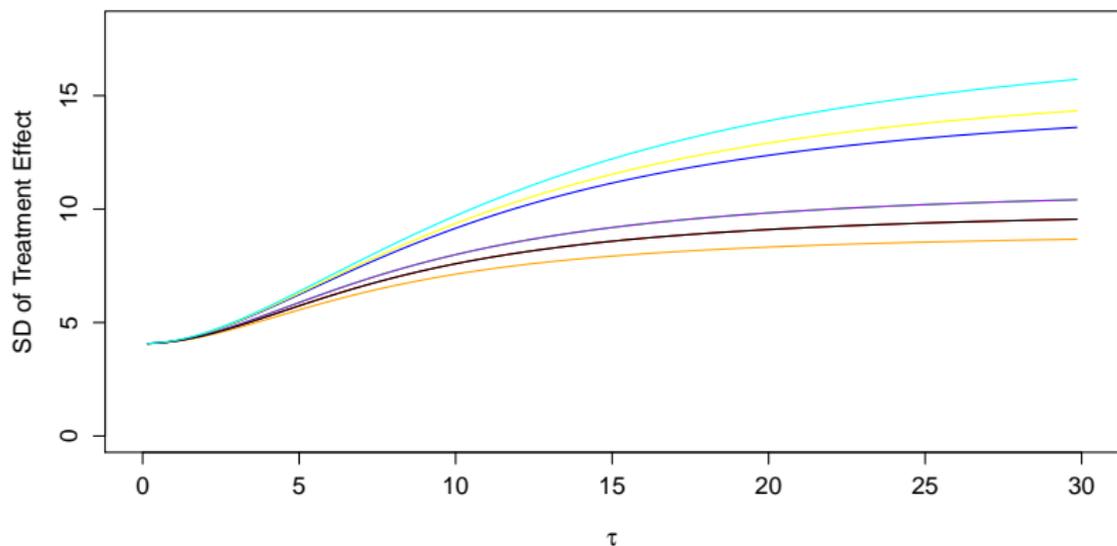- standard errors on estimated treatment effects are regarded as known

**Imperial College London**

[2]From Gelman et al. (2013), Bayesian Data Analysis, 3rd Edition, §5.5.

David A. van Dyk      Unified Analyses of Populations of Sources

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Shrinkage of the Fitted Treatment Effects



*Pooling may dramatically change fitted effects.*

Imperial College
London

All Roads Lead to Rome
Example 1: Using SNIa to Fit Cosmological Models
Example 2: Ages of White Dwarfs in the Galactic Halo

Joint work with Ted von Hippel & Shijing Si

# Standard Deviation of the Fitted Treatment Effects



*Pooling results in more precise estimates.*

Imperial College
London

## Fitting the Standard Deviation of the Treatment Effects



*Fitted $\tau$ determines the degree of pooling.*

Imperial College
London