



Geodesic least squares regression with applications to astrophysical data

G. Verdoolaege

Department of Applied Physics, Ghent University, B-9000 Ghent, Belgium
Laboratory for Plasma Physics, Royal Military Academy, B-1000 Brussels, Belgium

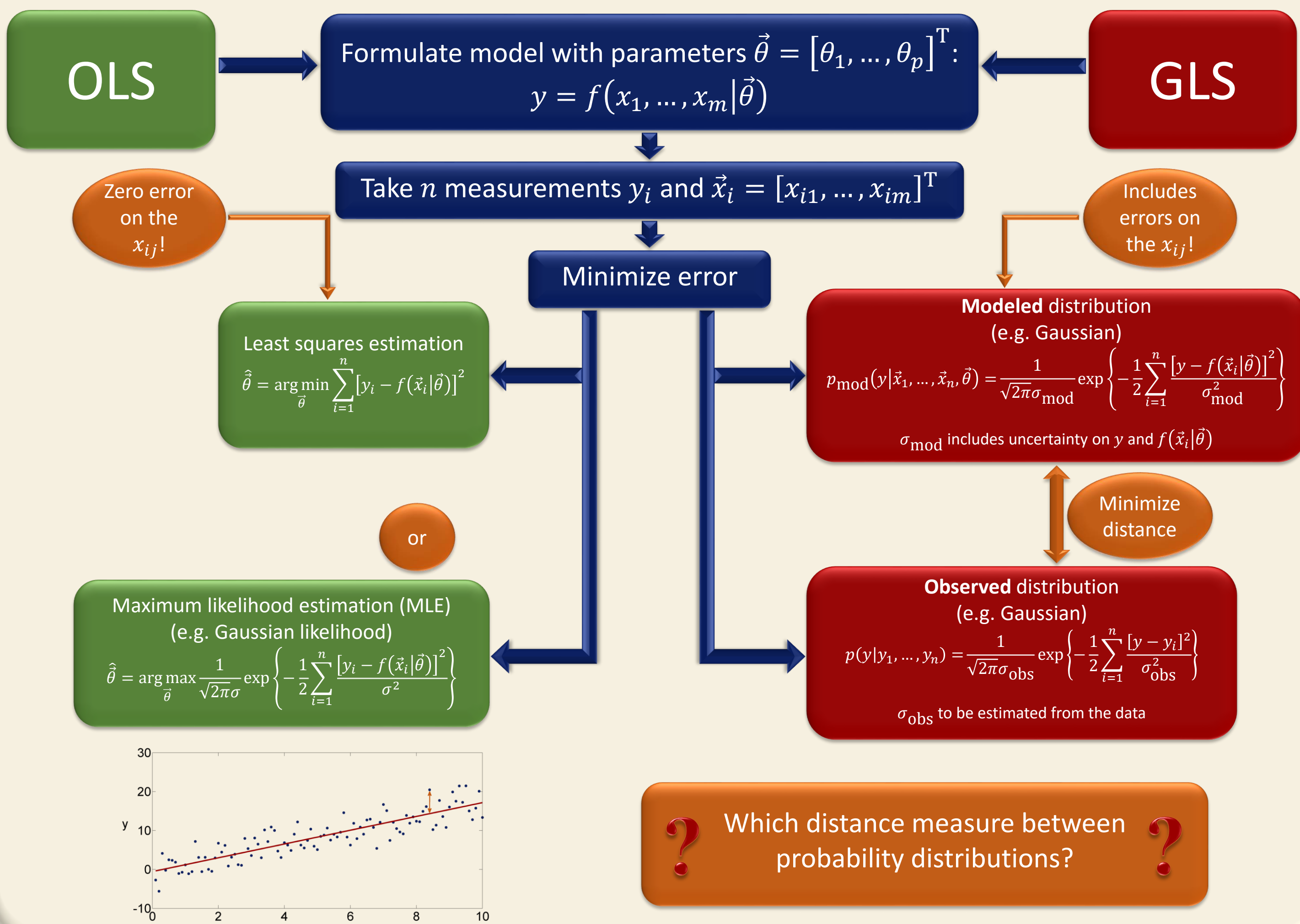
Abstract

Regression analysis on astrophysical data is a relatively challenging activity and it has long been recognized that many of the assumptions of ordinary least squares regression (OLS) are not valid in several applications in the field. Accordingly, various techniques from the domains of frequentist statistics and Bayesian probability theory have been proposed to address the shortcomings of OLS. We here present a new regression method called '**geodesic least squares regression**' (GLS), which has recently been developed and applied in the field of magnetic confinement fusion (MCF) [1,2]. The main difference with standard techniques is that, for the case of a single response variable, the distribution of the response conditional on a value suggested by the regression model ('modeled distribution'), need not be the same as its distribution conditional on an actual measurement ('observed distribution'). Then, instead of minimizing the difference between modeled and observed values of the response variable, which is the goal of standard OLS, GLS aims at minimizing the distance between the modeled and observed distributions. To this end, we use the Rao geodesic distance (GD) on the probabilistic manifold corresponding to the distribution in the regression model, equipped with the Fisher information metric. The method can handle errors in all variables, is robust against data outliers and uncertainty in the regression model, and can be used with arbitrary distribution models and regression functions. After introducing GLS and demonstrating its advantages on a synthetic data set, we show results of fitting MCF scaling laws as well as the baryonic Tully-Fisher relation in astronomy.

Motivation

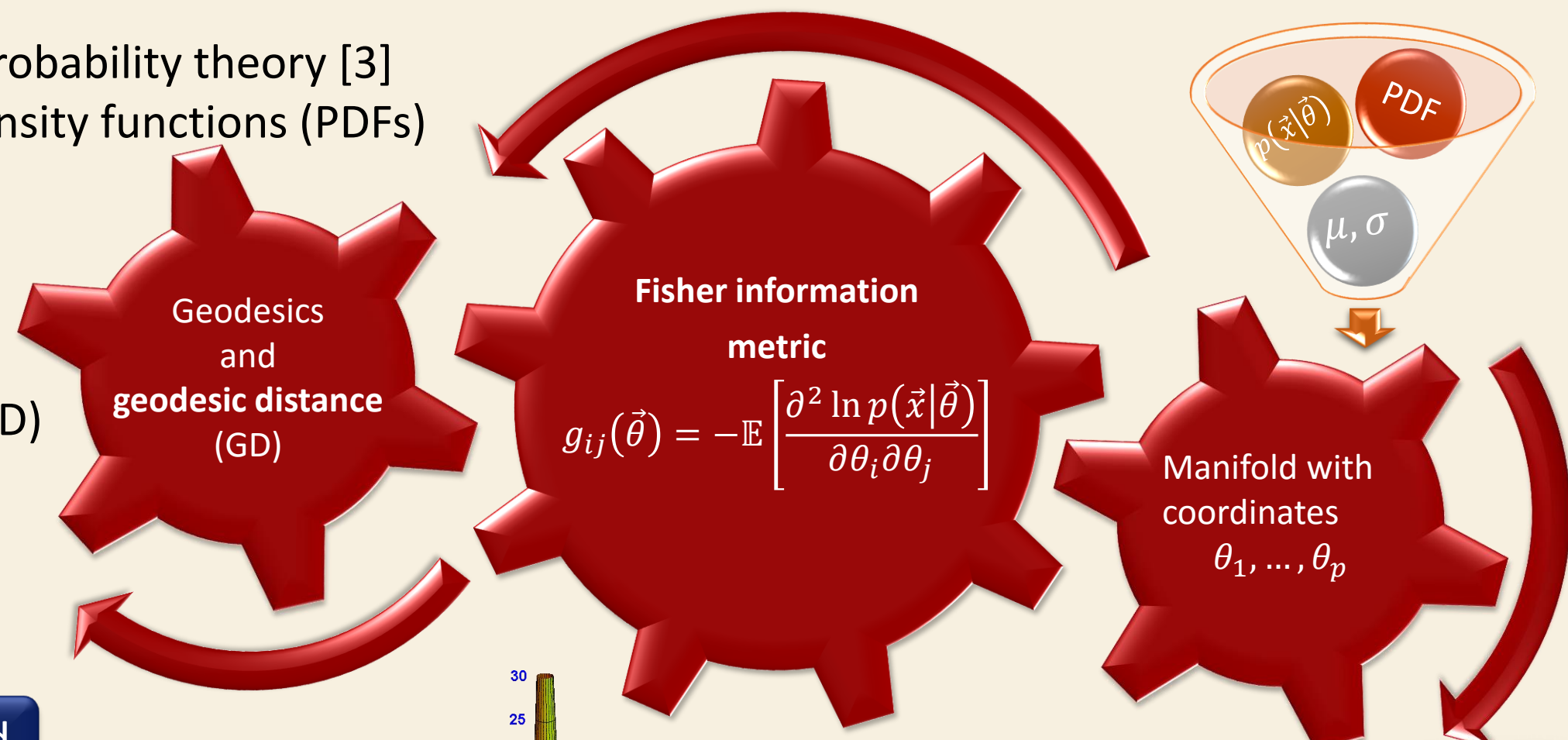
- In many areas of science, regression analysis is used:
 - As an aid to build and validate **theoretical models** from data and to find **parametric dependencies**
 - As a statistical tool to formulate **scaling laws** for the purpose of **extrapolation**
- Ordinary least squares regression (OLS)** is the workhorse
- Often, multiple assumptions underlying OLS are not fulfilled
- There may be various reasons:
 - Considerable measurement uncertainty: statistical and systematic
 - Uncertainty on response (dependent, y) and predictor (independent, x_j) variables
 - Model uncertainty: linear, power law, semi-empirical, ...
- Power law:** $y = b_0 x_1^{b_1} x_2^{b_2} \dots x_m^{b_m}$
- Heterogeneous data and error bars, correlations, non-Gaussian probability distributions
- Atypical observations (outliers)
- Near-collinearity of predictor variables
- Data transformations, e.g. $\log y = \log b_0 + b_1 \log x_1 + \dots + b_m \log x_m$
- Inferior regression analysis counteracts other efforts!**
- A **flexible, robust** and **user-friendly** regression tool is needed

Geodesic least squares regression (GLS)



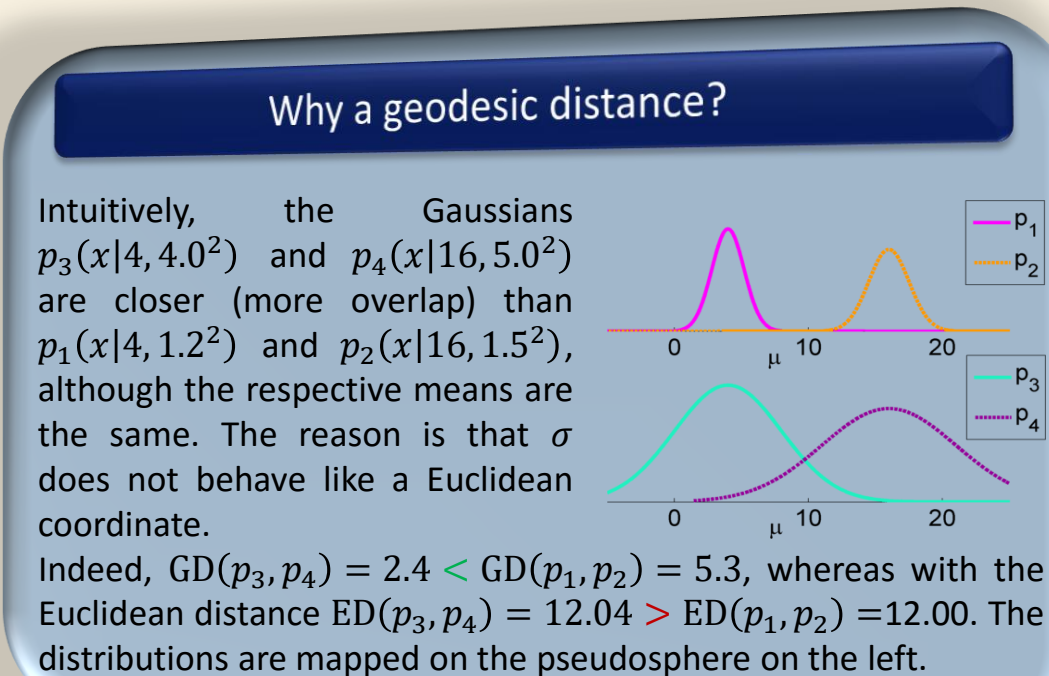
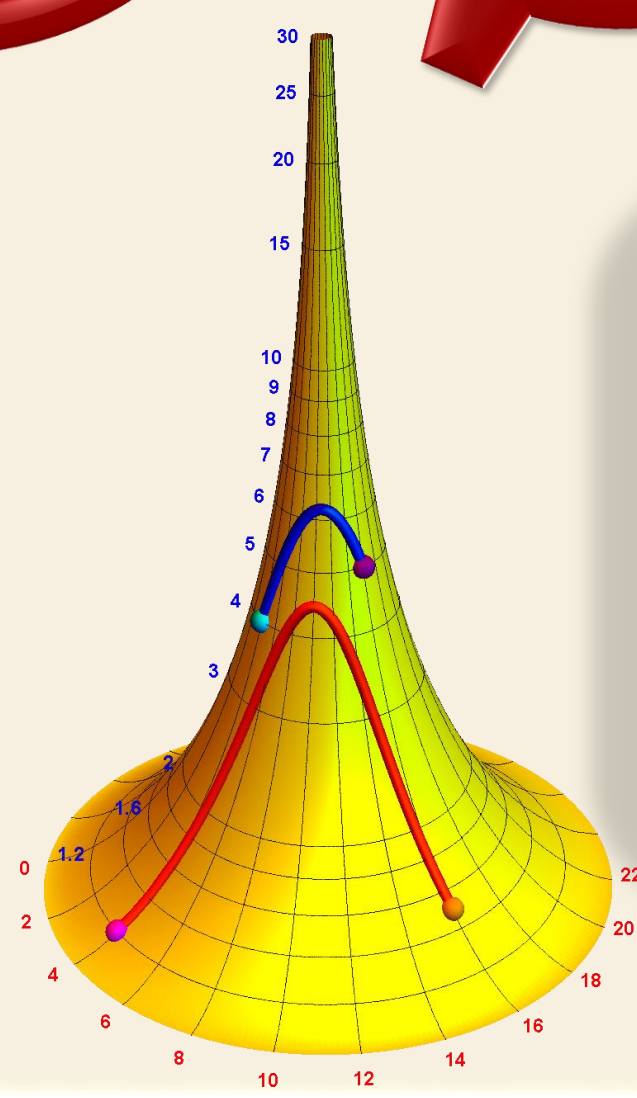
Information geometry

- Geometric approach to probability theory [3]
- A family of probability density functions (PDFs) forms a metric space, or **manifold**
- Fisher information** is the metric tensor
- Rao geodesic distance** (GD) is the shortest distance between points (PDFs)



Example: Gaussian manifold

- $p_1(x | \mu_1, \sigma_1) \leftrightarrow p_2(x | \mu_2, \sigma_2)$
- $GD(p_1, p_2) = 2\sqrt{2} \tanh^{-1} \delta$, $\delta = \frac{[(\mu_1 - \mu_2)^2 + 2(\sigma_1 - \sigma_2)^2]}{(\mu_1 - \mu_2)^2 + 2(\sigma_1 + \sigma_2)^2}$
- The **pseudosphere** (tractroid) is a model for the manifold of univariate Gaussian distributions, respecting the true geometry (μ in red, σ in blue)



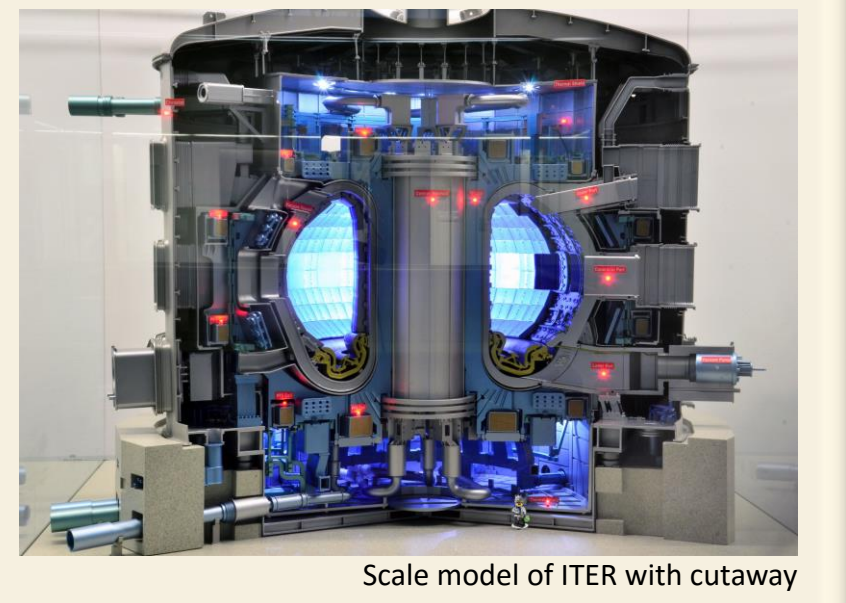
Scaling in magnetic confinement fusion

H-mode power threshold

- Controlled nuclear fusion: clean, safe, limitless energy
- Magnetic confinement fusion: tokamaks (ITER), stellarators, ...
- High confinement mode (H-mode): threshold P_{thr} on input power
- Scaling with classic power law:

$$P_{\text{thr}} = \beta_0 \bar{n}_e^{\beta_1} B_t^{\beta_2} S^{\beta_3}, \quad \begin{cases} \bar{n}_e = \text{average plasma density (10}^{19} \text{ m}^{-3}) \\ B_t = \text{toroidal magnetic field (T)} \\ S = \text{plasma surface area (m}^2) \end{cases}$$

- ITPA H-mode threshold database** [4]: 645 measurements from 7 tokamaks
- Logarithmic variables assumed Gaussian: single standard deviation = relative error from database



1. Synthetic data: outliers

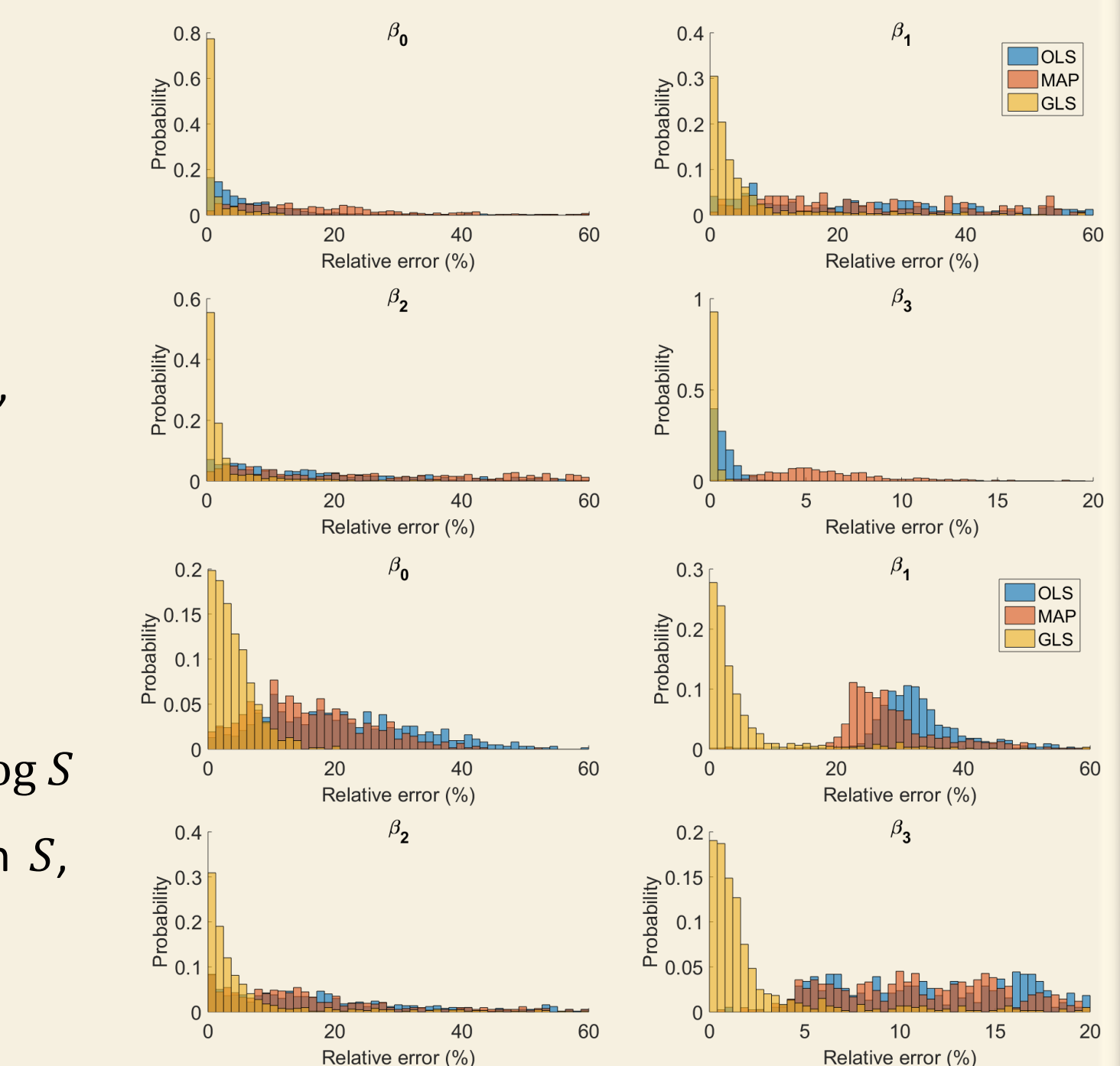
- Linear regression: $\eta = \beta_0 + \beta_1 \bar{n}_e + \beta_2 B_t + \beta_3 S$
- Artificial data sets: $\begin{cases} \beta_0 = 1, 1.1, \dots, 20 \\ \beta_1, \beta_2, \beta_3 = 0.1, 0.2, \dots, 2 \end{cases}$
- Gaussian noise: 4% on \bar{n}_e , 1% on B_t , 3% on S , 15% on P_{thr}

2. Synthetic data: logarithmic space

- Logarithmic transformation: $\eta = \log \beta_0 + \beta_1 \log \bar{n}_e + \beta_2 \log B_t + \beta_3 \log S$
- Gaussian noise: 20% on \bar{n}_e , 5% on B_t , 15% on S , 15% on P_{thr}

3. Real data: comparison of loglinear with nonlinear

Loglinear						Nonlinear					
Method	β_0	β_1	β_2	β_3	$P_{\text{thr},0.5}$ (MW)	Method	β_0	β_1	β_2	β_3	$P_{\text{thr},0.5}$ (MW)
GLS	0.043 ± 0.004	0.66 ± 0.07	0.80 ± 0.06	0.95 ± 0.03	48 \pm 5	GLS	0.040 ± 0.004	0.72 ± 0.07	0.75 ± 0.08	0.98 ± 0.03	52 \pm 4
OLS	0.051 ± 0.006	0.49 ± 0.07	0.87 ± 0.06	0.84 ± 0.04	38 \pm 4	OLS	0.027 ± 0.008	0.77 ± 0.09	1.0 ± 0.1	1.04 ± 0.07	70 \pm 20
MAP	0.045 ± 0.005	0.57 ± 0.08	0.87 ± 0.07	0.90 ± 0.04	46 \pm 5	MAP	0.046 ± 0.004	0.64 ± 0.07	0.79 ± 0.08	0.93 ± 0.03	44 \pm 4



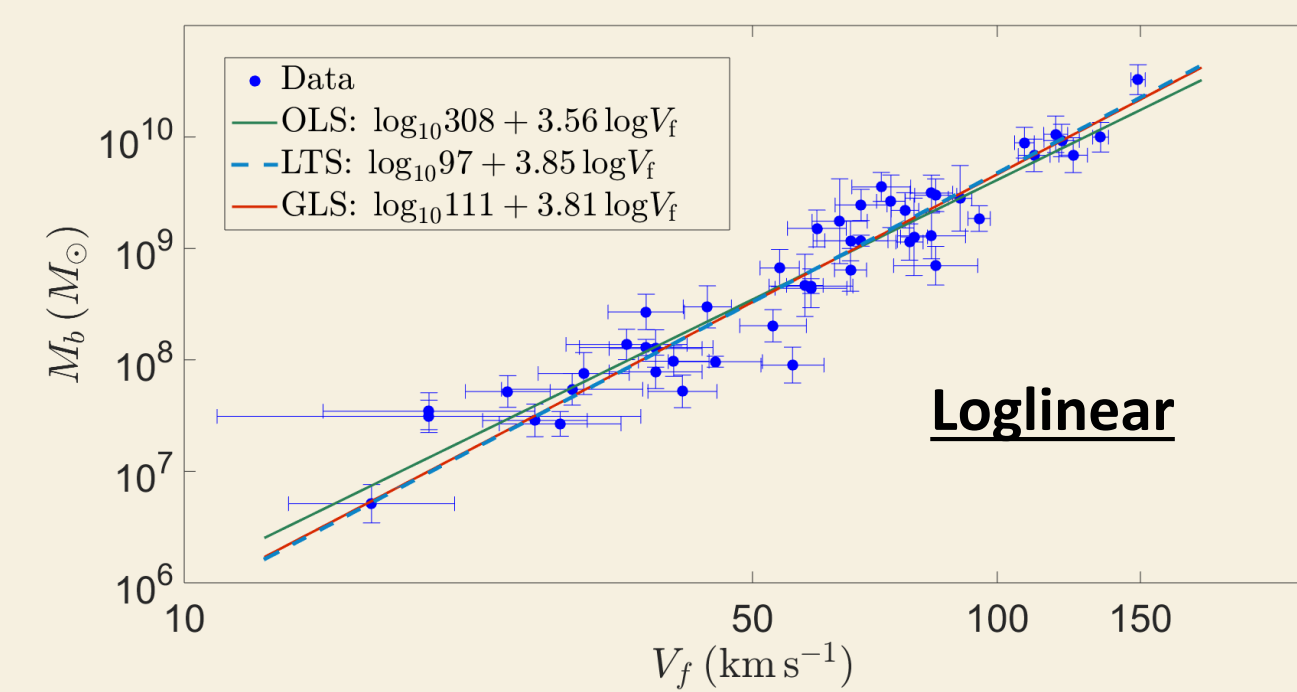
Baryonic Tully-Fisher relation in astronomy

- Relation between rotational velocity and baryonic mass of galaxies
- Various purposes:
 - Distance indicator
 - Constraints on galaxy formation models
 - Test for alternatives to Λ CDM (e.g. MOND) via slope and scatter

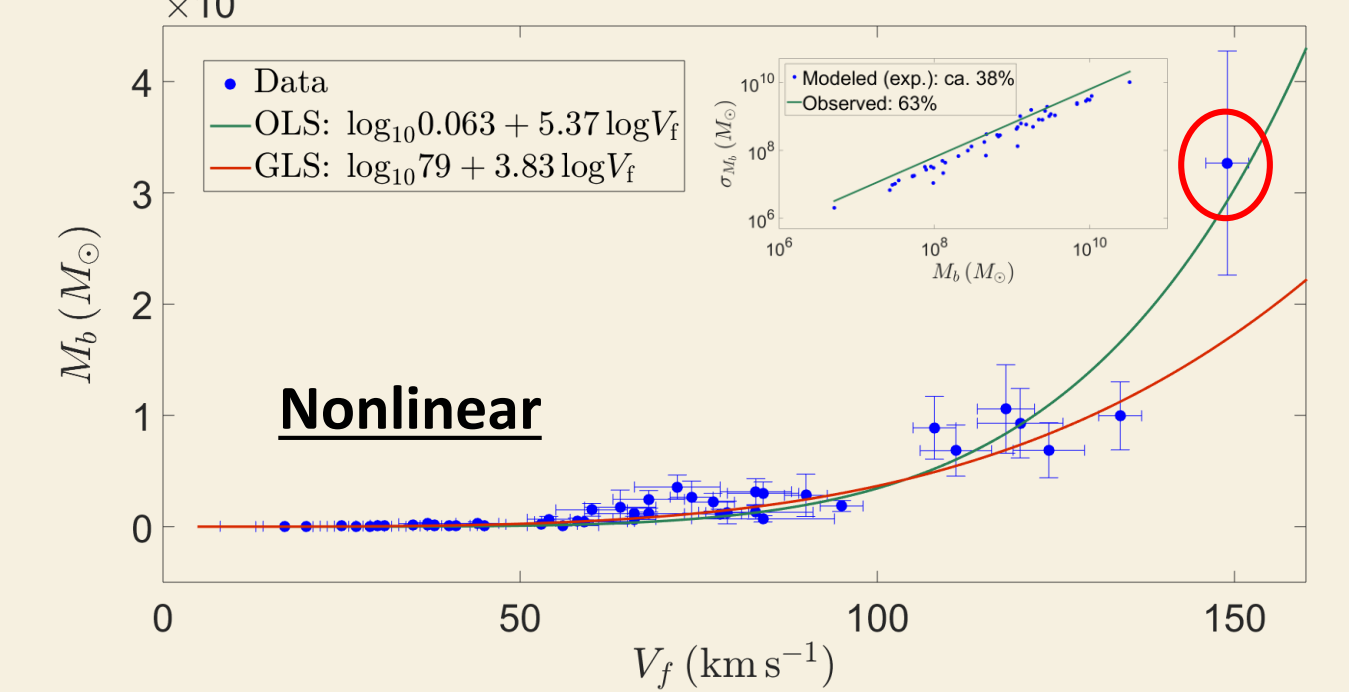


$$M_b = \beta_0 V_f^{\beta_1} \quad \begin{cases} V_f = \text{rotational velocity in flat part of rotation curve (km s}^{-1}) \\ M_b = M_* + M_{\text{gas}} = \text{total baryonic mass (M}_\odot) \\ \beta_0, \beta_1 = \text{constants} \end{cases}$$

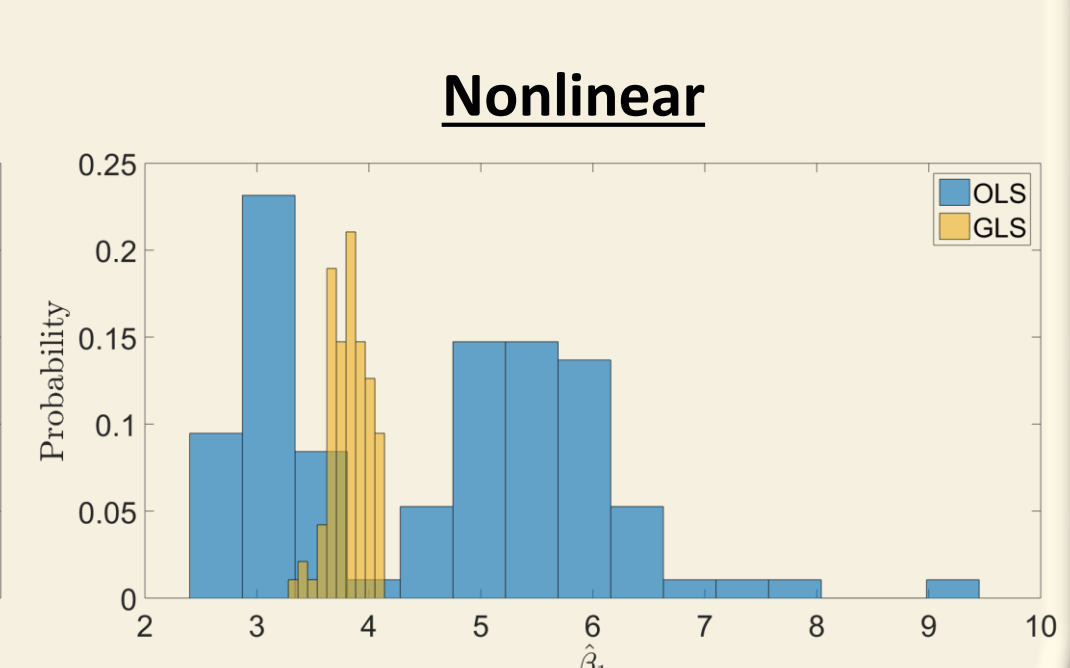
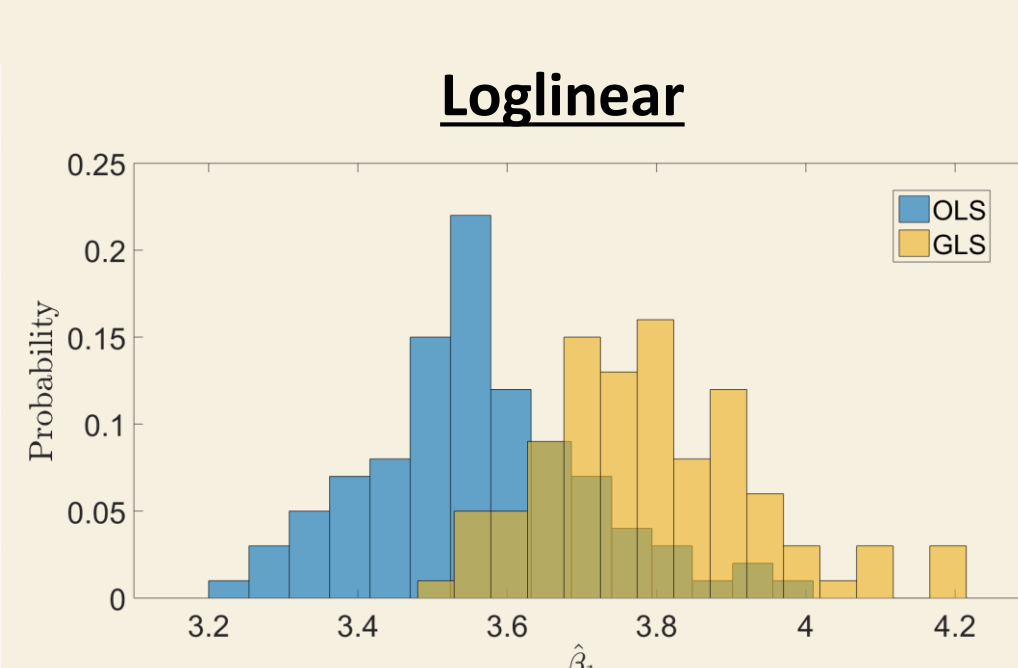
- Data from McGaugh [5] (gas-rich galaxies)



- Compare with least trimmed squares (LTS) [6,7]



100 bootstraps, 95% conf. interv.			Loglinear			Nonlinear		
Method	β_0	β_1	Method	β_0	β_1	Method	β_0	β_1
OLS loglin.	370 ± 400	3.56 ± 0.29	OLS loglin.	140 ± 150	3.80 ± 0.28	OLS loglin.	140 ± 150	3.80 ± 0.28
OLS nonlin.	2×10^3 $\pm 8 \times 10^3$	4.6 ± 2.7	OLS nonlin.	110 ± 220	3.82 ± 0.34	OLS nonlin.	110 ± 220	3.82 ± 0.34
GLS loglin.	140 ± 150	3.80 ± 0.28	GLS loglin.	110 ± 220	3.82 ± 0.34	GLS loglin.	110 ± 220	3.82 ± 0.34
GLS nonlin.	110 ± 220	3.82 ± 0.34	GLS nonlin.	110 ± 220	3.82 ± 0.34	GLS nonlin.	110 ± 220	3.82 ± 0.34



Conclusion

- Geodesic least squares** regression is **flexible** and **robust**
- GLS is **simple** but **powerful** due to strong mathematical foundations
- GLS offers **unified solution** to various regression problems
- Probability distributions **more informative** for regression
- Loglinear regression can be biased w.r.t. nonlinear analysis
- Future development: more accurate **error bars** on GLS estimates and predictions
- GLS will be implemented in a **public software package**

References

- G. Verdoolaege, Entropy **17**, 4602, 2015.
- G. Verdoolaege *et al.*, Nucl. Fusion **55**, 113019, 2015.
- S. Amari and H. Nagaoka, *Methods of Information Geometry*, AMS, New York, 2000.
- J.A. Snipes *et al.*, Fusion Energy 2002 (Proc. 19th Int. Conf. Lyon), IAEA, Vienna, CT/P-04.
- S. McGaugh, Astron. J. **143**, 40, 2012.
- M. Cappellari *et al.*, Mon. Not. R. Astron. Soc. **432**, 1709, 2013.
- P.J. Rousseeuw and A.M. Leroy, *Robust Regression and Outlier Detection*, Wiley, New York, 1987.