

*Bayesian and Frequentist Issues
in Modern Inference*

Bradley Efron

Stanford University

Small Data

- **Classical statistics** Direct tests and estimates of individual parameters within well-defined models (MLE, Neyman–Pearson)
- **Not much:**
 - ▶ data-based model selection
 - ▶ Bayesian combination of related problems
- *Today* Methodology (not Philosophy)

Bayesian Inference

- Parameter: $\mu \in \Omega$
- Observed data: x
- Prior: $g(\mu)$
- Probability distributions: $\{f_\mu(x), \mu \in \Omega\}$
- Parameter of interest: $\theta = t(\mu)$

$$E\{\theta|x\} = \int_{\Omega} t(\mu)f_\mu(x)g(\mu) d\mu \Big/ \int_{\Omega} f_\mu(x)g(\mu) d\mu$$

Jeffreysonian Bayes Inference

“Uninformative Priors”

- What if we don't know prior g ?
- Jeffreys: $g(\mu) = |\mathcal{I}(\mu)|^{1/2}$ where $\mathcal{I}(\mu) = \text{cov} \{ \nabla_{\mu} \log f_{\mu}(x) \}$
(the Fisher information matrix)
- Can still use Bayes theorem but how accurate are the estimates?
- Frequentist variability of $E \{ t(\mu) | x \}$

General Accuracy Formula

- μ and $x \in \mathcal{R}^p$
- $V_\mu = \text{cov}_\mu(x)$
- $\alpha_x(\mu) = \nabla_x \log f_\mu(x) = \left(\dots, \frac{\partial \log f_\mu(x)}{\partial x_i}, \dots \right)^T$

Lemma

$E = E\{t(\mu)|x\}$ has gradient $\nabla_x E = \text{cov}\{t(\mu), \alpha_x(\mu)|x\}$.

Theorem

The delta-method standard deviation of E is

$$\text{sd}(E) = \left[\text{cov}\{t(\mu), \alpha_x(\mu)|x\}^T V_x \text{cov}\{t(\mu), \alpha_x(\mu)|x\} \right]^{1/2}.$$

Implementation

- Posterior sample from $\mu|x$

$\{\mu_1, \mu_2, \dots, \mu_B\}$ (MCMC?)

- Each μ_i gives $t_i = t(\mu_i)$ and

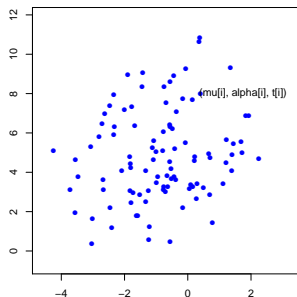
$$\alpha_i = \alpha_x(\mu_i)$$

- $\hat{E} = \sum t_i/B \doteq E\{t(\mu)|x\}$

- $\widehat{\text{cov}} = \sum_{i=1}^B (\alpha_i - \hat{\alpha})(t_i - \bar{t})/B$

- $\widehat{\text{sd}} = \left[\widehat{\text{cov}}^T V_x \widehat{\text{cov}} \right]^{1/2}$

- No additional sampling for $\widehat{\text{cov}}$



Diabetes Data

Efron et al. (2004), "LARS"

- $n = 442$ subjects
- $p = 10$ predictors: age, sex, bmi, glu, . . .
- Response: $y =$ disease progression at one year
- Model:
$$\underset{n \times 1}{\mathbf{y}} = \underset{n \times p}{\mathbf{X}} \underset{p \times 1}{\boldsymbol{\beta}} + \underset{n \times 1}{\mathbf{e}} \quad [\mathbf{e} \sim \mathcal{N}_n(\mathbf{0}, I)]$$

Bayesian Lasso

Park and Casella (2008)

- Model: $\mathbf{y} \sim \mathcal{N}_n(\mathbf{X}\beta, I)$
- Prior: $g(\beta) = e^{-\gamma L_1(\beta)}$ [$\gamma = 0.37$]
- Then posterior mode at Lasso $\hat{\beta}_\gamma$
- Subject 125: $\theta_{125} = \mathbf{x}_{125}^T \beta$
- How accurate are Bayes posterior inferences for θ_{125} ?

Bayesian Analysis

- MCMC: posterior sample $\{\beta_i \text{ for } i = 1, 2, \dots, 10,000\}$
- Gives $\{\theta_{125,i} = \mathbf{x}_{125}^T \beta_i, i = 1, 2, \dots, 10,000\}$

$$\theta_{125,i} \sim 0.248 \pm 0.072$$

- **General accuracy formula** frequentist sd **0.071** for $E = 0.248$

$$[\alpha_x(\mu) = \mathbf{X}^T \mathbf{X} \beta]$$

Posterior CDF for Subject 125

- $\text{cdf}_{\mathbf{y}}(c) = \Pr\{\theta_{125} \leq c | \mathbf{y}\}$

- $s_i = \begin{cases} 1 & \text{as } t_i \leq c \\ 0 & \text{as } t_i > c \end{cases}$

- $\widehat{\text{cdf}}_{\mathbf{y}}(c) = \sum_1^B s_i / B$

- For $c = 0.3$:

$$\widehat{\text{cdf}}_{\mathbf{y}}(c) = 0.762 \pm 0.304$$

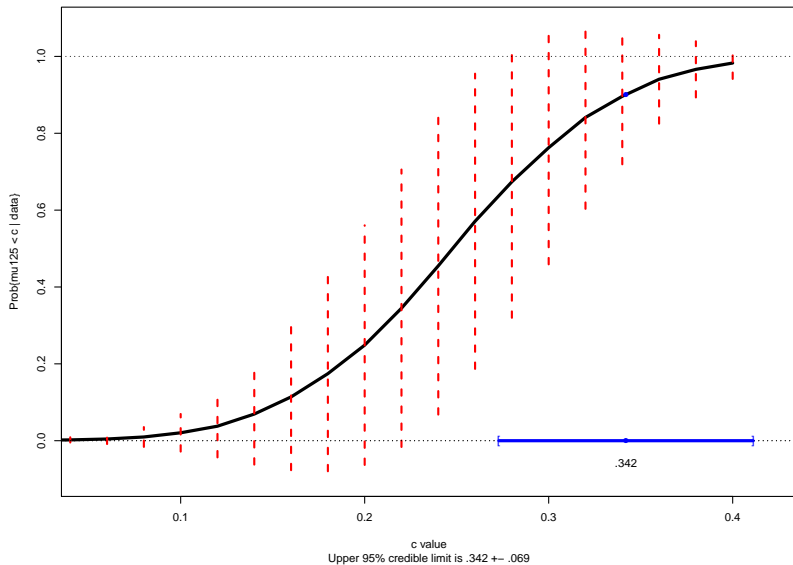


Bayes
estimate



frequentist
sd

Posterior cdf for mu125, Diabetes data, 10000 MCMC draws,
Prior $\exp\{-.37 \cdot L1(\beta)\}$; verts are \pm One Frequentist Standard Dev



Exponential Families

- $f_\alpha(\hat{\beta}) = e^{\alpha^T \hat{\beta} - \psi(\alpha)} f_0(\hat{\beta})$, with $\alpha, \hat{\beta}$ in \mathcal{R}^p
- Natural parameter α , sufficient statistic $\hat{\beta}$, expectation $\beta = E_\alpha\{\hat{\beta}\}$
[Poisson : $f_\mu(x) = e^{-\mu} \mu^x / x! : x = \hat{\beta}, \mu = \beta, \alpha = \log(\mu)$]
- **General accuracy formula** For $E = E\{t(\beta)|\hat{\beta}\}$,

$$\widehat{\text{sd}}(E) = \left\{ \text{cov}(t, \alpha | \hat{\beta})^T V_{\hat{\alpha}} \text{cov}(t, \alpha | \hat{\beta}) \right\}^{1/2}$$

with $V_{\hat{\alpha}} = \text{cov}_{\alpha=\hat{\alpha}}(\hat{\beta})$.

Better Frequentist Inferences

- For $E = E \{t(\beta) | \hat{\beta}\}$ in exfam $f_\alpha(\hat{\beta}) = e^{\alpha^T \hat{\beta} - \psi(\alpha)} f_0(\hat{\beta})$
- **Parametric bootstrap** $f_{\hat{\alpha}}(\cdot) \rightarrow [\hat{\beta}_1^*, \hat{\beta}_2^*, \dots, \hat{\beta}_j^*, \dots, \hat{\beta}_J^*]$
 $\rightarrow [\dots E_j^* = E \{t(\beta) | \hat{\beta}_j^*\} \dots] \rightarrow$ bootstrap conf int for E
- **Trouble** Need new MCMC sample for each $\hat{\beta}_j^*$
- **Shortcut** Reweight original MCMC sample (importance sampling)

Digression: *Posterior Exponential Family*

- $f_\alpha(\hat{\beta}) = e^{\alpha^T \hat{\beta} - \psi(\alpha)} f_0(\hat{\beta})$: natural param α , suff stat $\hat{\beta}$, “carrier” $f_0(\hat{\beta})$
- Posterior exponential family

$$g(\alpha | \text{suff stat } b) = e^{(b - \hat{\beta})^T \alpha - \phi(b)} g(\alpha | \hat{\beta})$$

- natural param b , suff stat α , carrier $g(\alpha | \hat{\beta})$
- **Importance sampling** Reweight the original MCMC realizations:

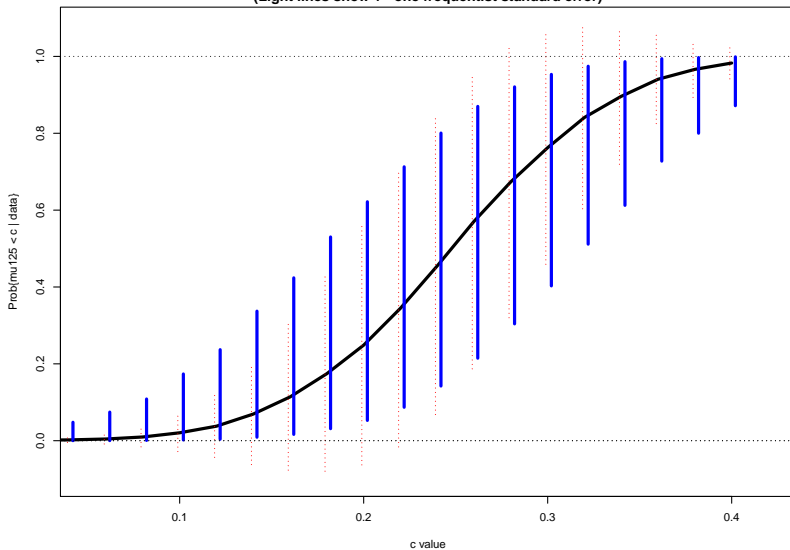
$$\hat{E}\{t_i | b\} = \frac{\sum_{i=1}^B t_i W_i(b)}{\sum_{i=1}^B W_i(b)} \quad \left[W_i(b) = e^{(b - \hat{\beta})^T \alpha_i} \right]$$

Bootstrap Intervals Without Bootstrapping

DiCiccio and Efron (1992)

- “**abc**” Investigate $\hat{E}\{t|b\}$ for b near $\hat{\beta}$
- Requires $p + 2$ numerical 2nd derivatives of \hat{E} function
- *Next:* Applied to posterior cdf for mu125

Heavy curve is posterior cdf for mu125, diabetes data.
Vertical bars frequentist central 68% abc confidence intervals.
(Light lines show \pm one frequentist standard error)



Estimation After Model Selection

- **Usually:**
 - (a) look at data
 - (b) choose model (linear, quad, cubic . . . ?)
 - (c) fit estimates using chosen model
 - (d) analyze as if pre-chosen
- *Today* Include model selection process in the analysis
- **Question** Effects on standard errors, confidence intervals, etc.?

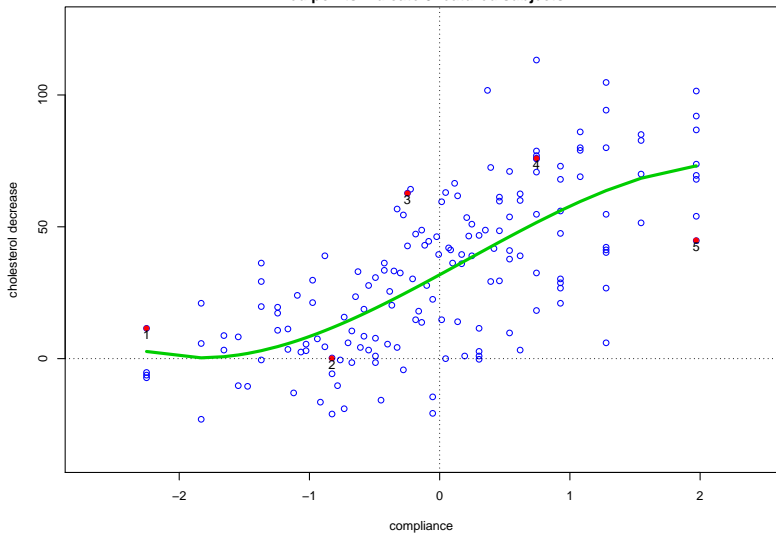
Cholesterol Data

- $n = 164$ men took Cholestyramine for ~ 7 years
- $x = \text{compliance measure}$ (adjusted: $x \sim \mathcal{N}(0, 1)$)
- $y = \text{cholesterol decrease}$
- Wish to estimate regression values

$$\mu_j = E\{y|x = x_j\} \quad \text{for } j = 1, 2, \dots, 164$$

- $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_{164})^T$

Cholesterol data, n=164 subjects: cholesterol decrease plotted versus adjusted compliance; Green curve is OLS cubic regression; Red points indicate 5 featured subjects



C_p Selection Criterion

- Regression model $\mathbf{y} = \mathbf{X} \beta + \mathbf{e} \quad [e_i \sim (0, \sigma^2)]$
 $n \times 1 \quad n \times m \quad m \times 1 \quad n \times 1$

- C_p criterion $\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2 + 2m\sigma^2$

$\hat{\beta}$ = OLS estimate, m = “degrees of freedom”

- Model selection From possible models X_1, X_2, X_3, \dots
choose the one minimizing C_p .
- Then use OLS estimate from chosen model.

C_p for Cholesterol Data

Model	df	$C_p - 80000$	(Boot %)
\mathcal{M}_1 (linear)	2	1132	(19%)
\mathcal{M}_2 (quad)	3	1412	(12%)
\mathcal{M}_3 (cubic)	4	667	(34%)
\mathcal{M}_4 (quartic)	5	1591	(8%)
\mathcal{M}_5 (quintic)	6	1811	(21%)
\mathcal{M}_6 (sextic)	7	2758	(6%)

($\sigma = 22$ from “full model” \mathcal{M}_6)

Nonparametric Bootstrap Analysis

- **data** = $\{(x_i, y_i), i = 1, 2, \dots, n = 164\}$ gave original estimate

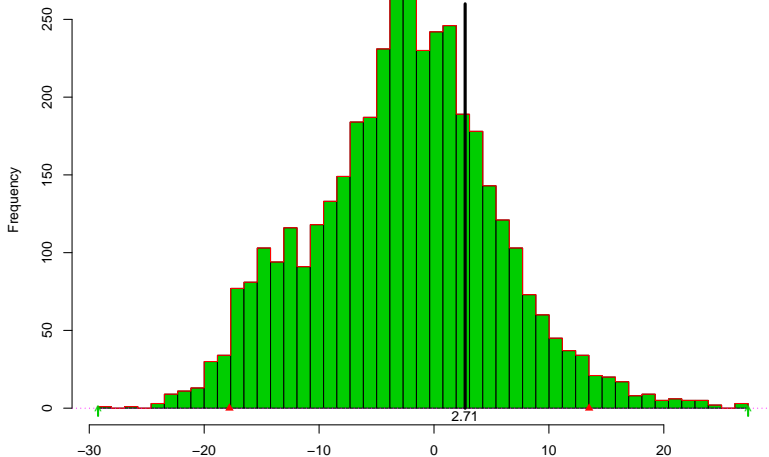
$$\hat{\mu} = X_3 \hat{\beta}_3$$

- Bootstrap data set: **data*** = $\{(x_j, y_j)^*, j = 1, 2, \dots, n\}$ where $(x_j, y_j)^*$ drawn randomly and with replacement from **data**:

$$\text{data}^* \xrightarrow{C_p} m^* \xrightarrow{\text{OLS}} \hat{\beta}_{m^*}^* \longrightarrow \hat{\mu}^* = X_{m^*} \hat{\beta}_{m^*}^*$$

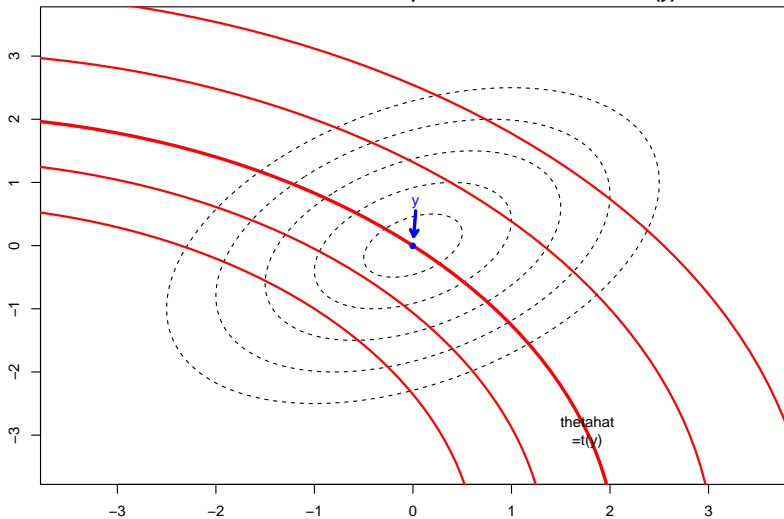
- I did this all $B = 4000$ times.

B=4000 nonparametric bootstrap replications for the model-selected regression estimate of Subject 1; boot (m,stdev)=(-2.63,8.02); 76% of the replications less than original estimate 2.71

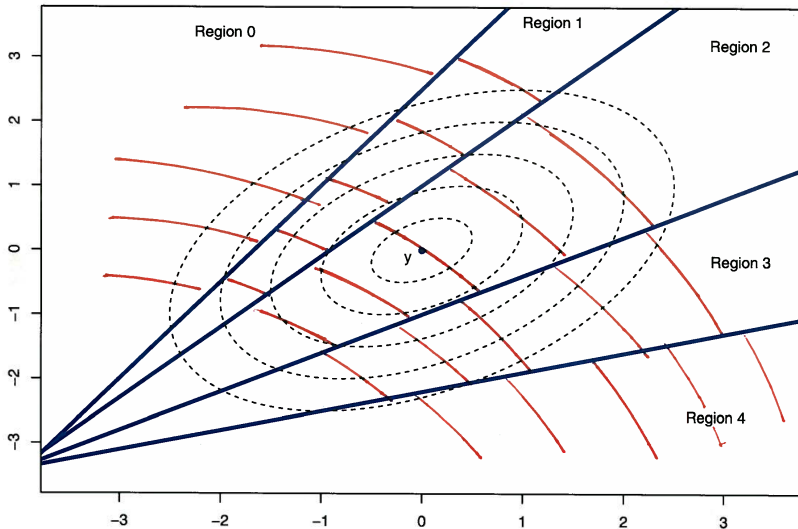


bootstrap estimates for subject 1
Red triangles are 2.5th and 97.5th boot percentiles

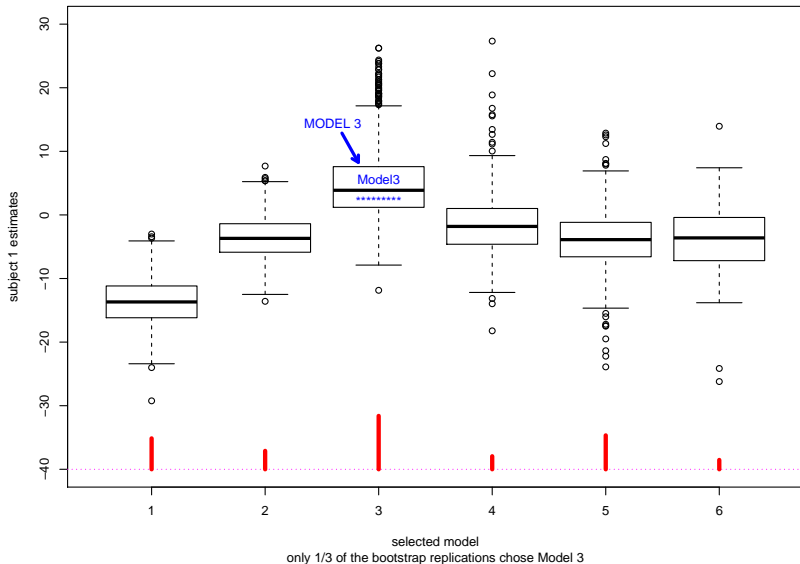
Smooth Estimation Model: ' y ' is observed data;
Ellipses indicate bootstrap distribution for ' y^* ' ;
Red curves level surfaces of equal estimation for $\hat{\theta} = t(y)$



Estimation with model selection: now the curves of equal estimation are discontinuous across the Model region boundaries



Boxplot of Cp boot estimates for Subject 1; B=4000 bootreps;
Red bars indicate selection proportions for Models 1-6



Bootstrap Smoothing

- *Idea* Replace original estimator $t(\mathbf{y})$ with bootstrap average

$$s(\mathbf{y}) = \sum_{i=1}^B t(\mathbf{y}_i^*) / B$$

- Model averaging
- Same as *bagging* (“bootstrap aggregation,” Breiman)
- Removes discontinuities, reduces variance
- Approximate confidence interval: $s(\mathbf{y}) \pm 1.96 \cdot \widetilde{\text{sd}}$

Accuracy Theorem

- **Notation** $s_0 = s(\mathbf{y})$, $t_i^* = t(\mathbf{y}_i^*)$, $i = 1, 2, \dots, B$
- $Y_{ij}^* = \#$ of times j th data point appears in i th boot sample
- $\text{cov}_j = \sum_{i=1}^B Y_{ij}^* \cdot (t_i^* - s_0) / B$ [covariance Y_{ij}^* with t_i^*]

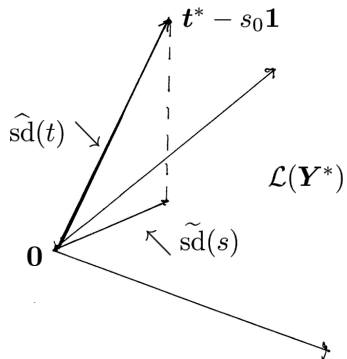
Theorem

The delta method standard deviation estimate for s_0 is

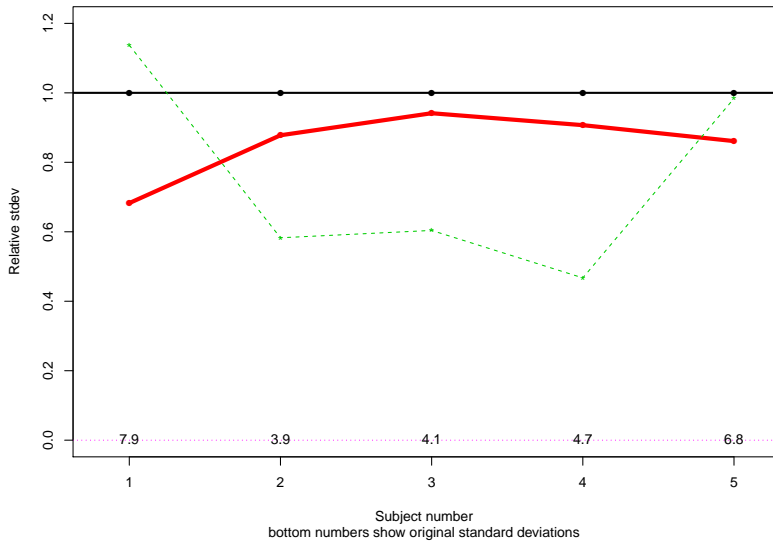
$$\widetilde{\text{sd}} = \left[\sum_{j=1}^n \text{cov}_j^2 \right]^{1/2},$$

always $\leq \left[\sum_{i=1}^B (t_i^* - s_0)^2 / B \right]^{1/2}$, the boot stdev for $t(\mathbf{y})$.

Projection Interpretation



Standard Deviation of smoothed estimate relative to original (Red)
for five subjects; green line is stdev Naive Cubic Model



Model Probability Estimates

- 34% of the 4000 bootreps chose the cubic model
- Poor man's Bayes posterior prob for "cubic"
- How accurate is that 34%?
- Apply accuracy theorem to indicator function for choosing "cubic"

Model	Boot %	\pm Standard Error
\mathcal{M}_1 (linear)	19%	± 24
\mathcal{M}_2 (quad)	12%	± 18
\mathcal{M}_3 (cubic)	34%	± 24
\mathcal{M}_4 (quartic)	8%	± 14
\mathcal{M}_5 (quintic)	21%	± 27
\mathcal{M}_6 (sextic)	6%	± 6