

Statistical Rethinking

Week 4: Ockham, Ulysses, and the Model

Richard McElreath

NARODOWY BANK POLSKI

KK 4859628

1000

TYSIĄC ZŁOTYCH

WARSZAWA, 1 CZERWCA 1982 r.

PREZES

J. Pawłowski

GŁÓWNY
SKARBNIK

M. Kowalski

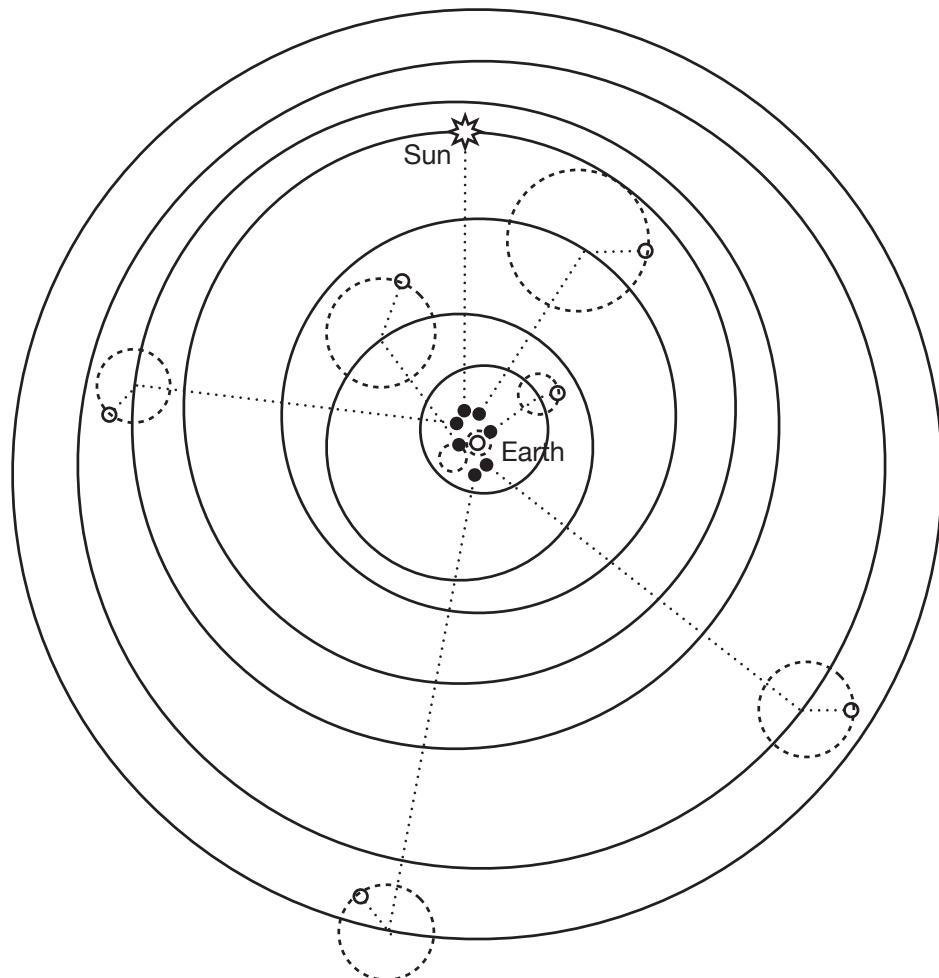
1000

KK 4859628

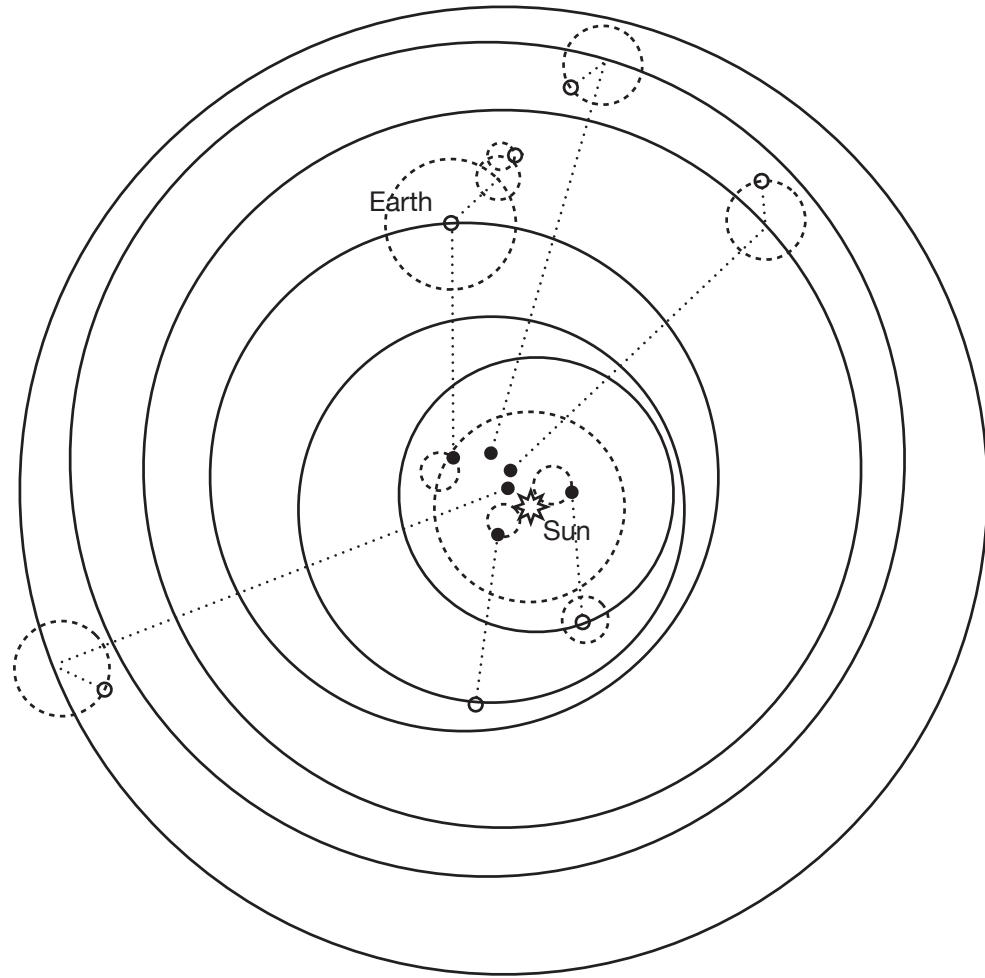
MIKOŁAJ
KOPERNIK

Mikołaj Kopernik (1473–1543)

Ptolemaic Model



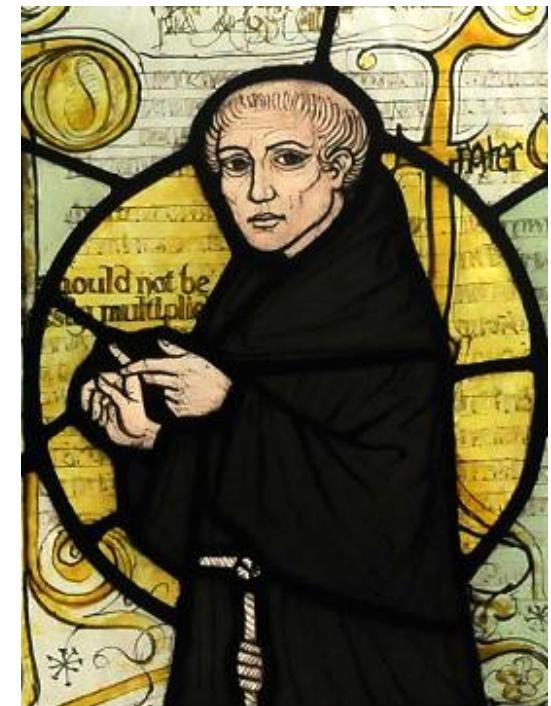
Copernican Model



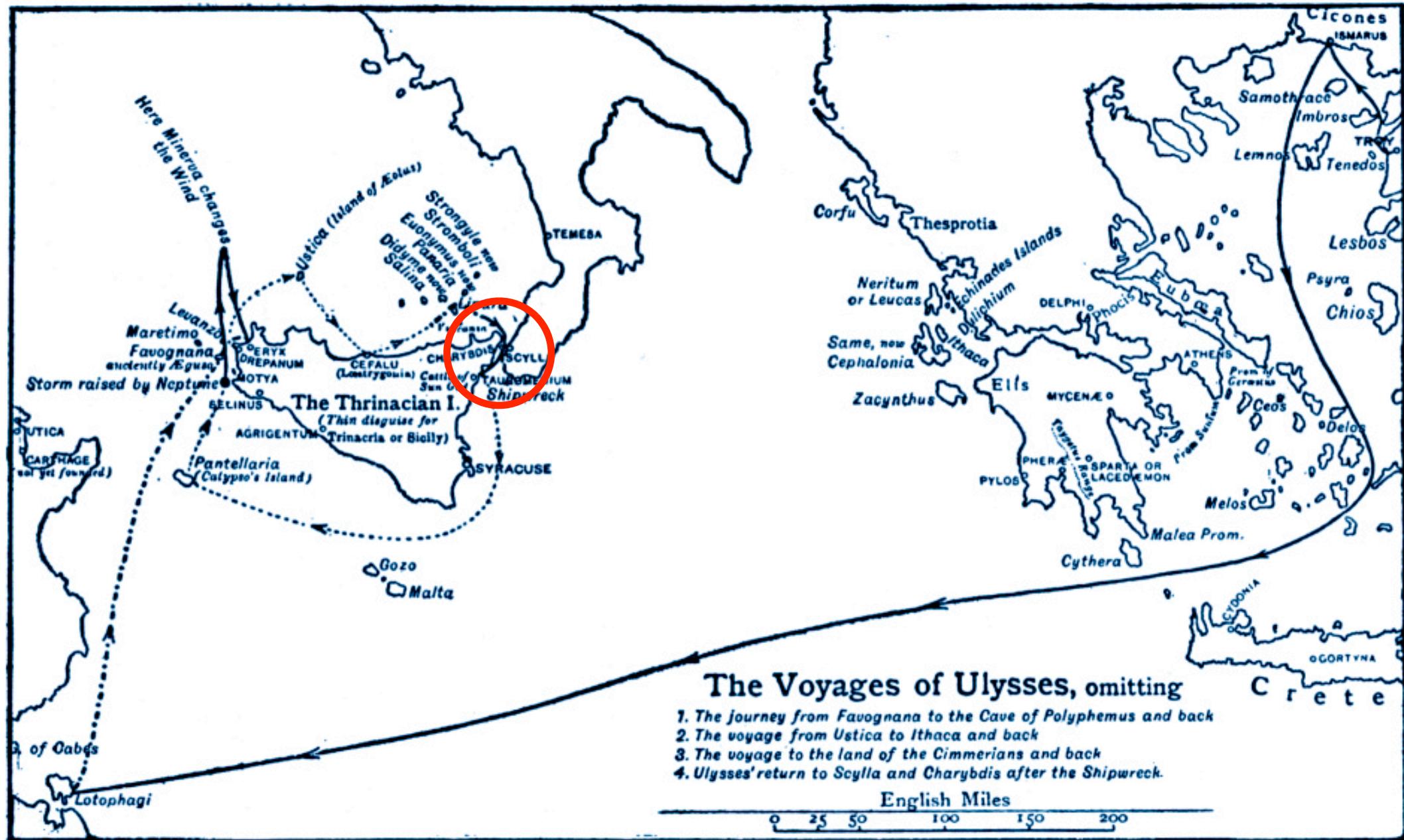
Ockham's Razor?

*Numquam ponenda est pluralitas
sine necessitate.*

(Plurality should never be posited without necessity.)



William of Ockham
(c.1288–c.1348)



Ulysses' Compass

- Two major hazards:
 1. Too simple
 2. Too complex



Stargazing

- *Stargazing*: Using asterisks ($p < 0.05$) to decide which variables improve prediction
- Arbitrary 5% is arbitrary

Coefficients:

	Estimate	Std. Error	z value	Pr(z)	
a	1.5699e+02	9.3802e-16	1.6736e+17	< 2.2e-16	***
b1	1.6540e-01	6.6628e-14	2.4825e+12	< 2.2e-16	***
b2	-4.7063e-02	3.2586e-13	-1.4443e+11	< 2.2e-16	***
b3	1.9168e-03	5.6805e-11	3.3743e+07	< 2.2e-16	***
b4	-1.4002e-05	6.6694e-11	-2.0994e+05	< 2.2e-16	***
b5	-4.7965e-07	4.7818e-08	-1.0031e+01	< 2.2e-16	***
b6	6.6002e-09	9.5819e-10	6.8882e+00	5.651e-12	***
tau	1.2132e-01	5.2829e-20	2.2965e+18	< 2.2e-16	***

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘					



Goals this week

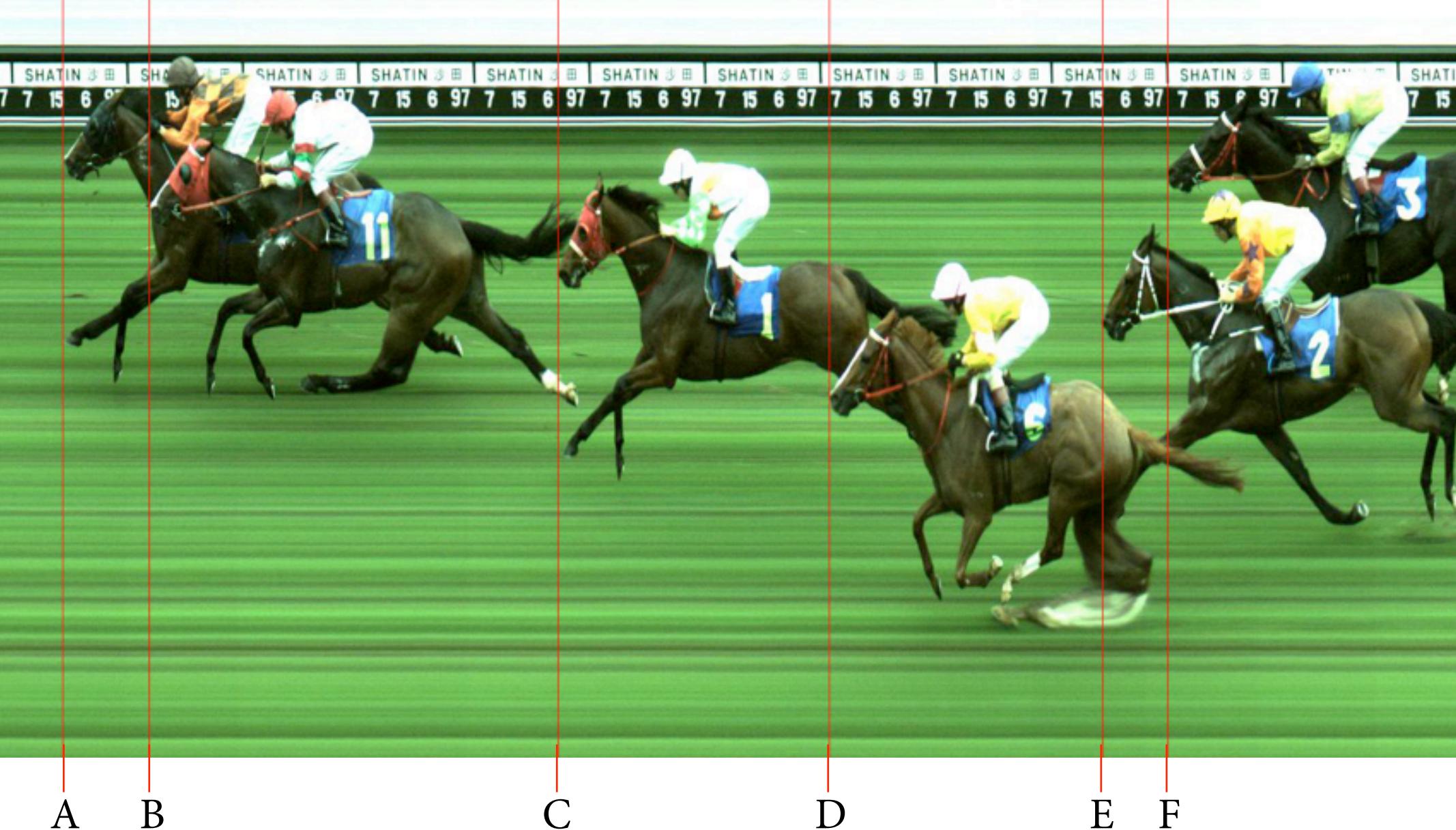
- Understand *overfitting* and *underfitting*
- Learn AIC/DIC/WAIC as ways to:
 - guard against overfitting and underfitting
 - explicitly compare models
- Introduce *regularizing* priors as complementary strategy
- Learn how to average *predictions* across models

AIC

WAIC

DIC





The Problem with Parameters

- *Underfitting*: Learning too little from the data. Too simple models both fit and predict poorly.
- *Overfitting*: Learning too much from the data. Complex models always fit better, but often predict worse.
- Need to find a model that navigates between underfitting and overfitting

The Problem with Parameters

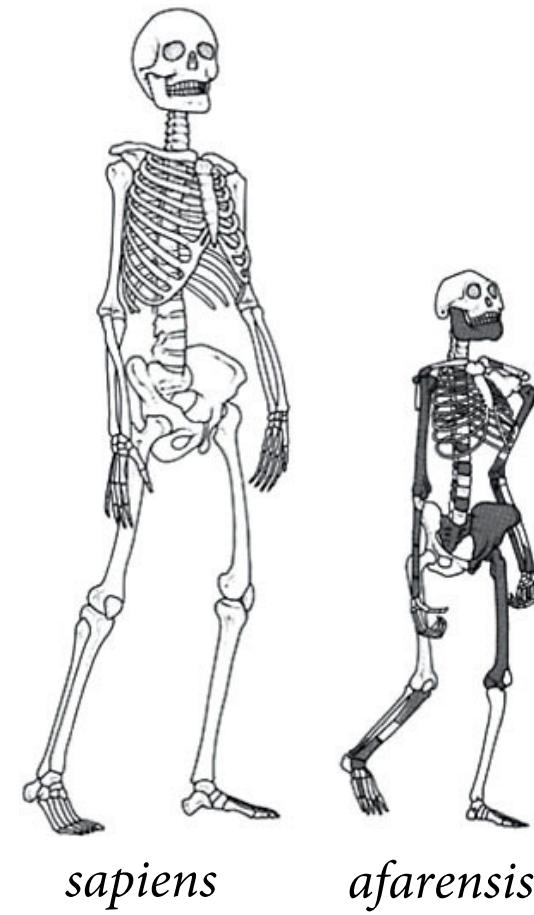
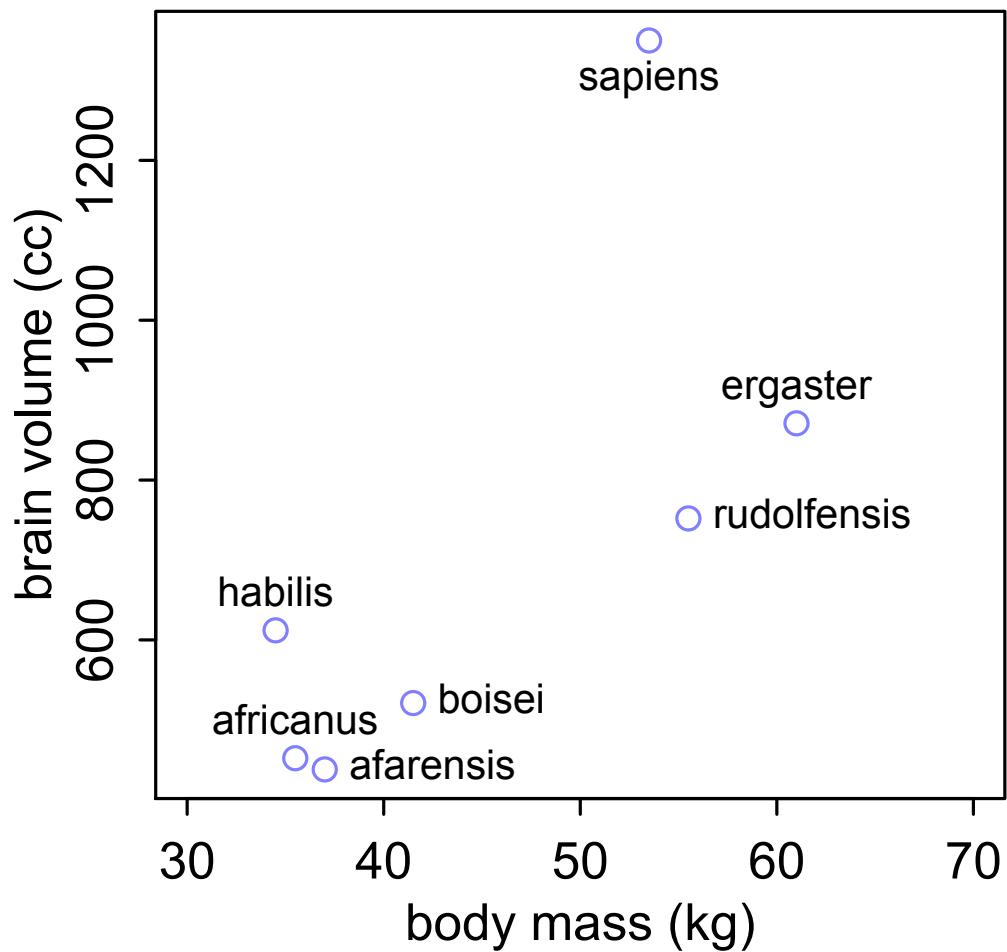


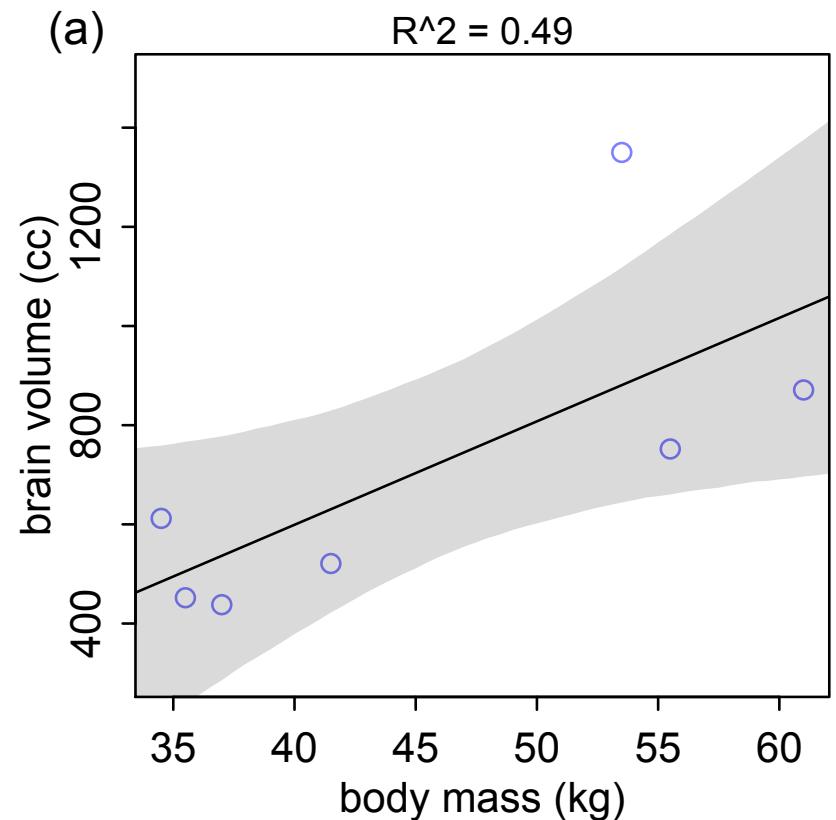
Figure 6.2

Hominin brains

- Simplest model:

$$v_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta_1 m_i$$

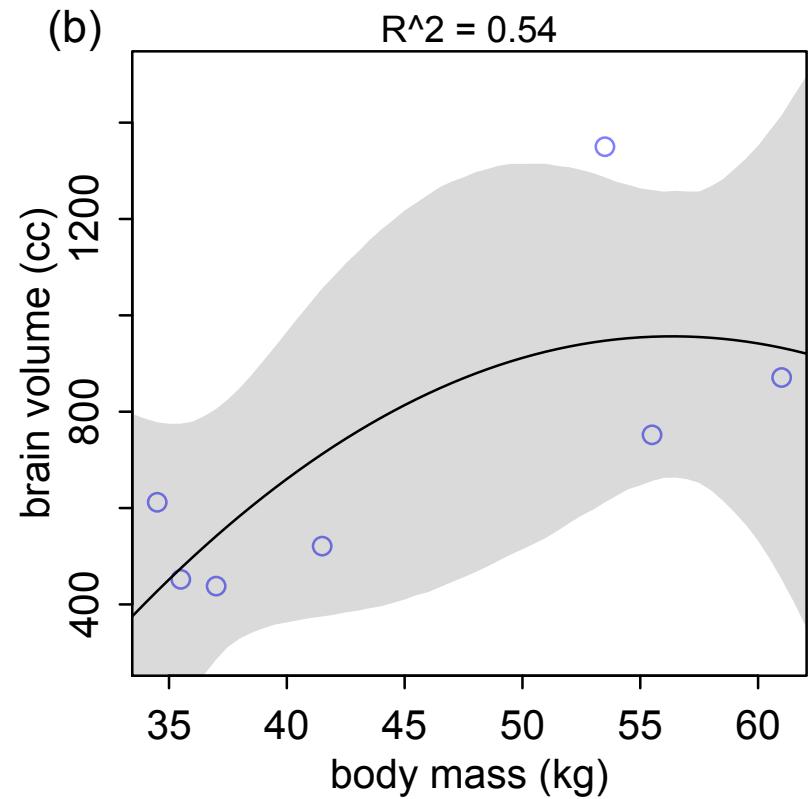


Hominin brains

- Why not parabola?

$$\nu_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta_1 m_i + \beta_2 m_i^2$$



Hominin brains

- Why not higher order polynomials?

$$\mu_i = \alpha + \beta_1 m_i + \beta_2 m_i^2 + \beta_3 m_i^3 + \beta_4 m_i^4 + \beta_5 m_i^5 + \beta_6 m_i^6$$

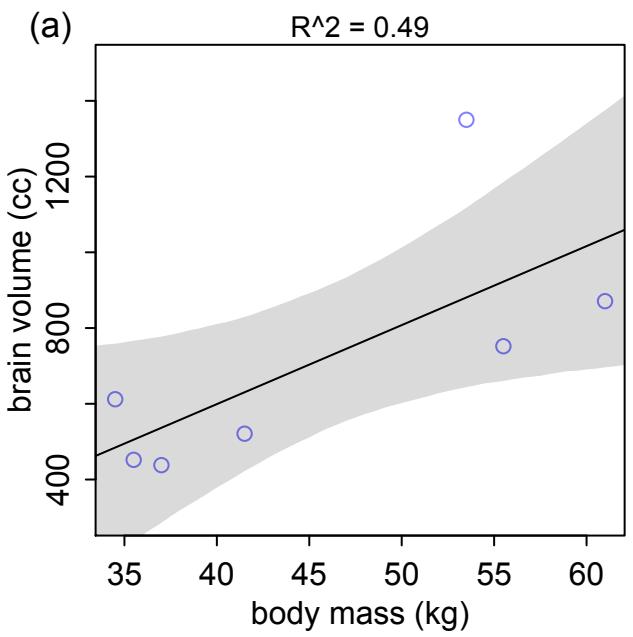


Figure 6.3

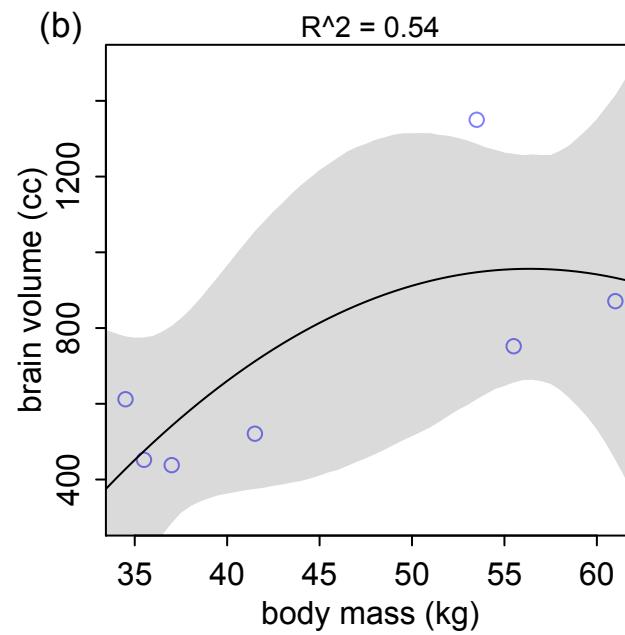
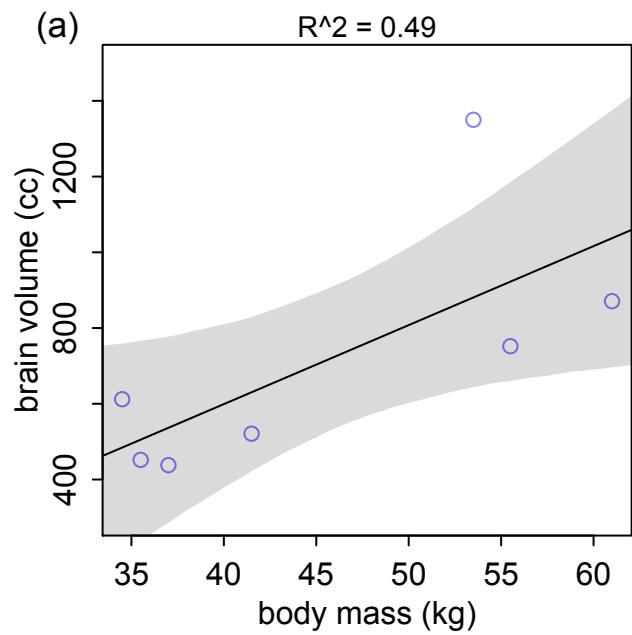


Figure 6.3

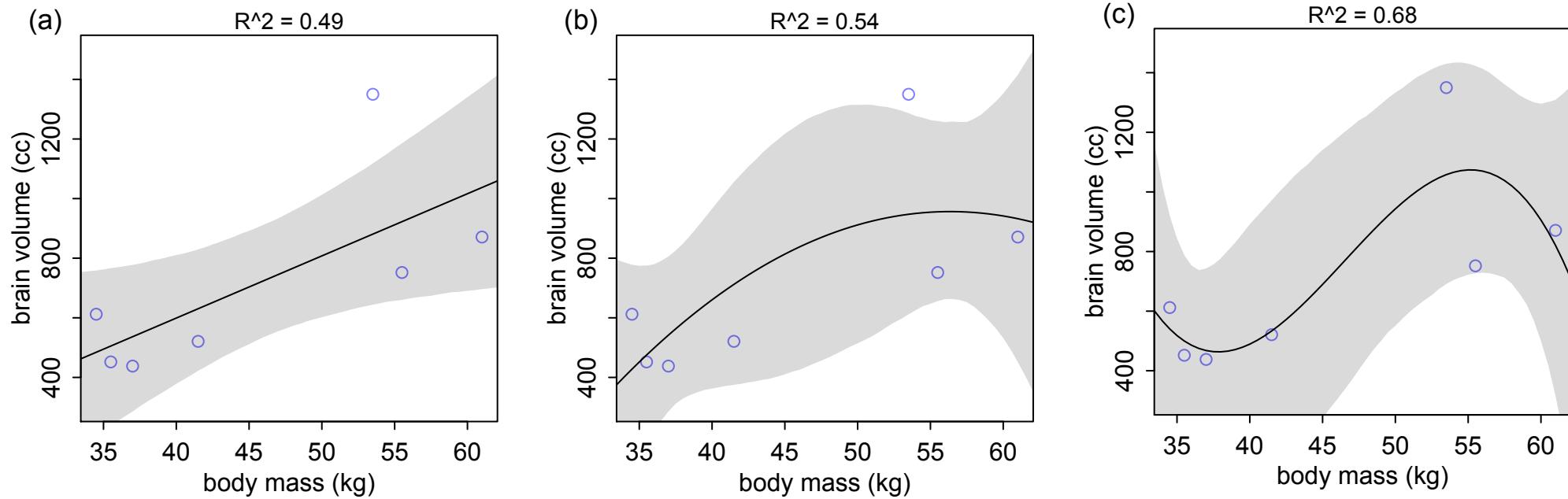


Figure 6.3

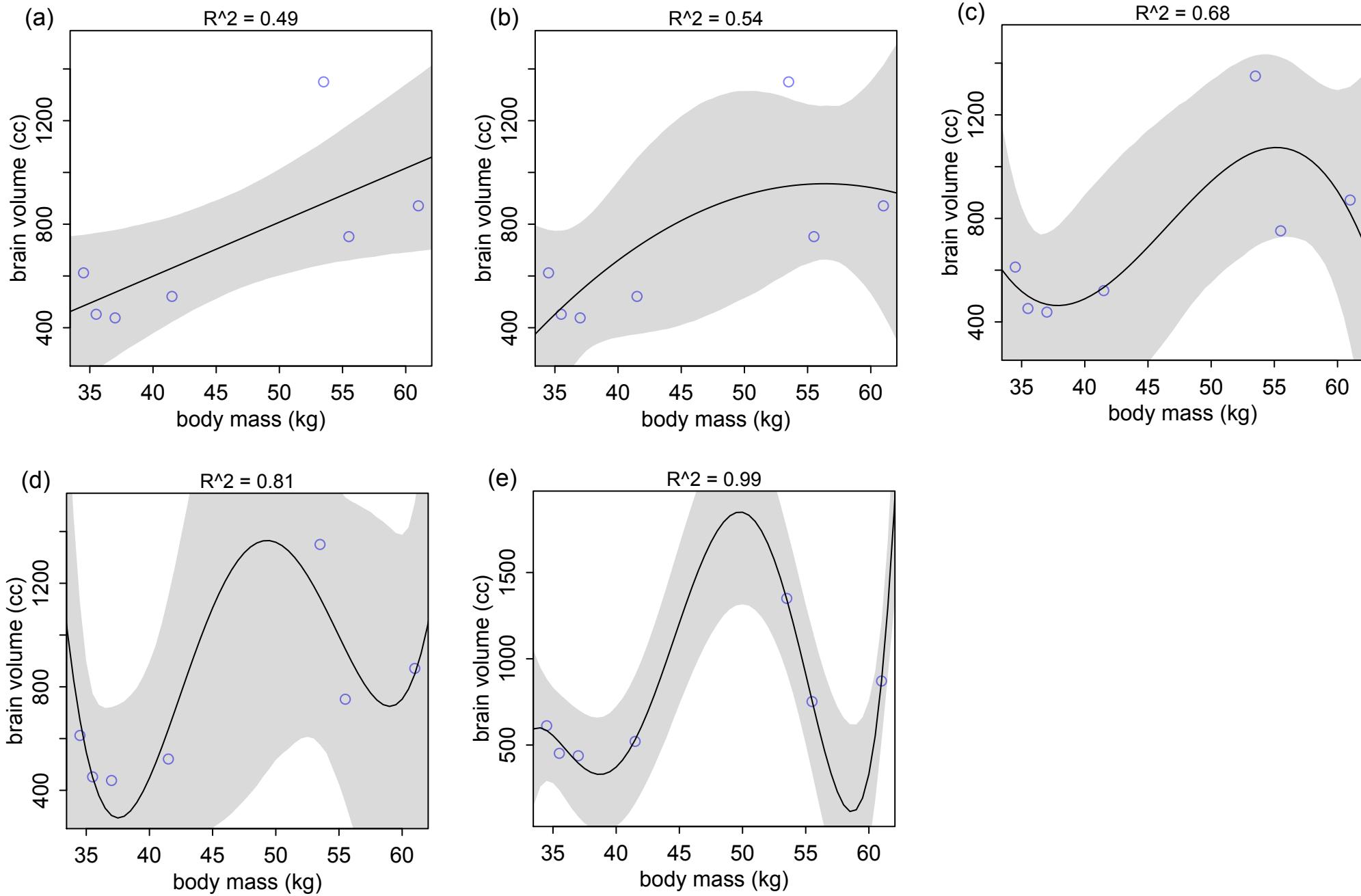


Figure 6.3

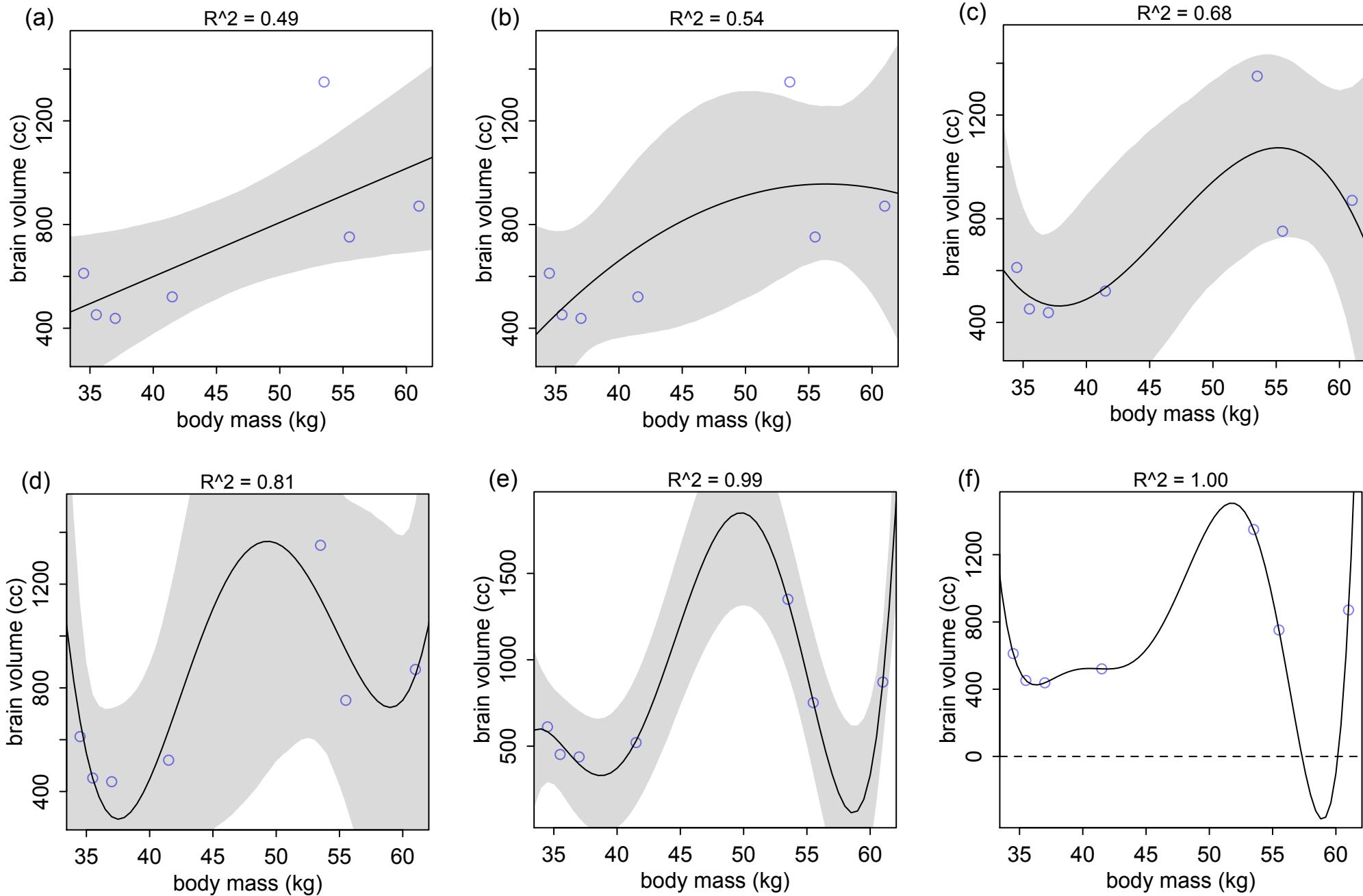


Figure 6.3

Underfitting
Insensitive to
exact data

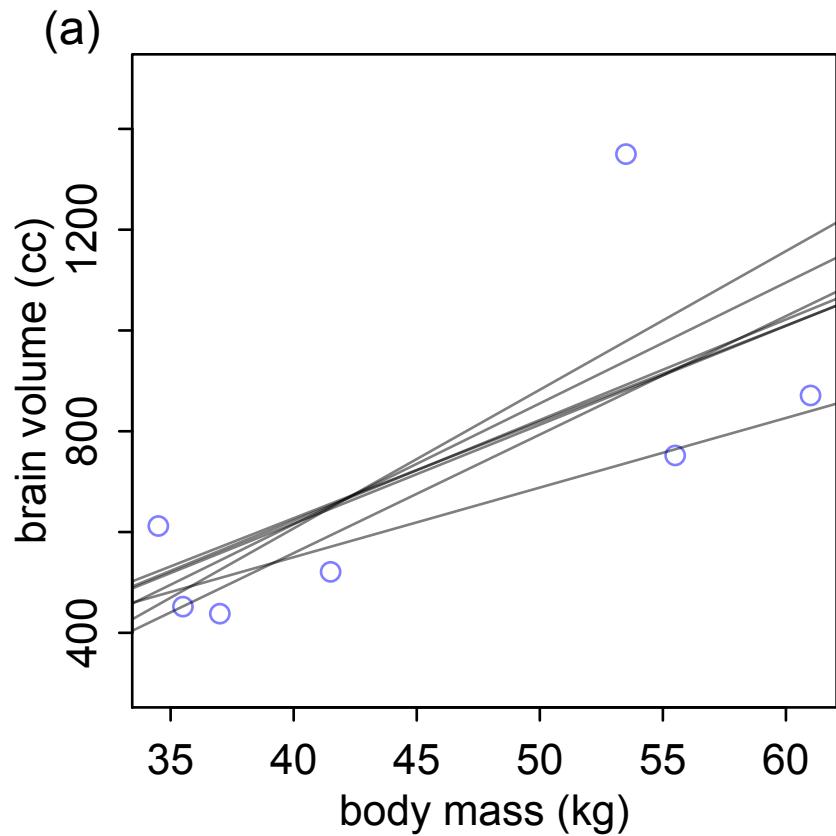
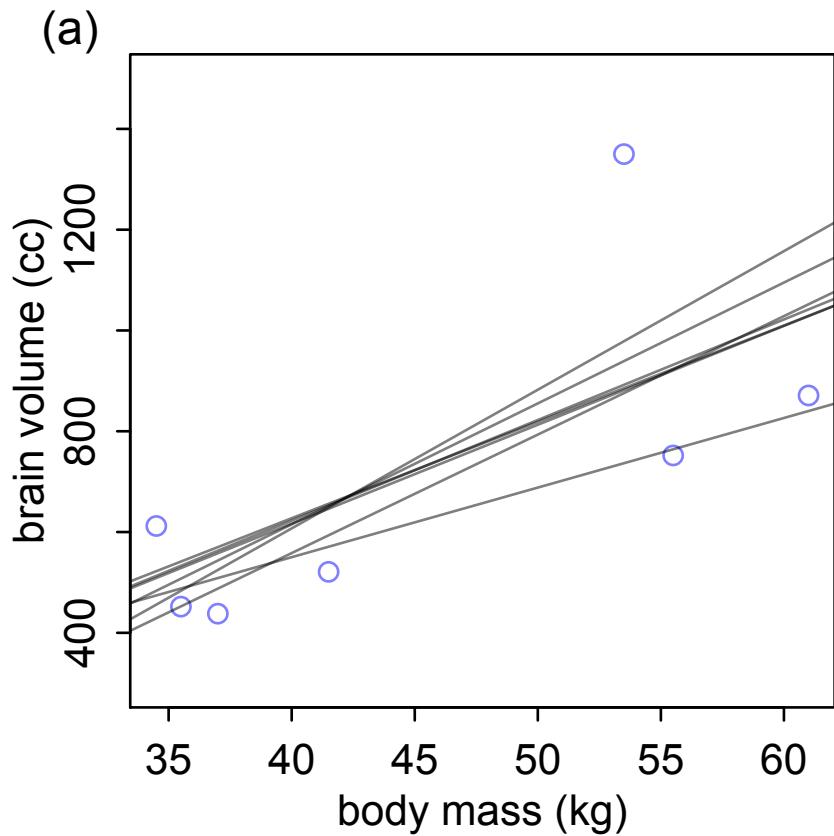


Figure 6.5

Underfitting
Insensitive to
exact data



Overfitting
Very sensitive to
exact data

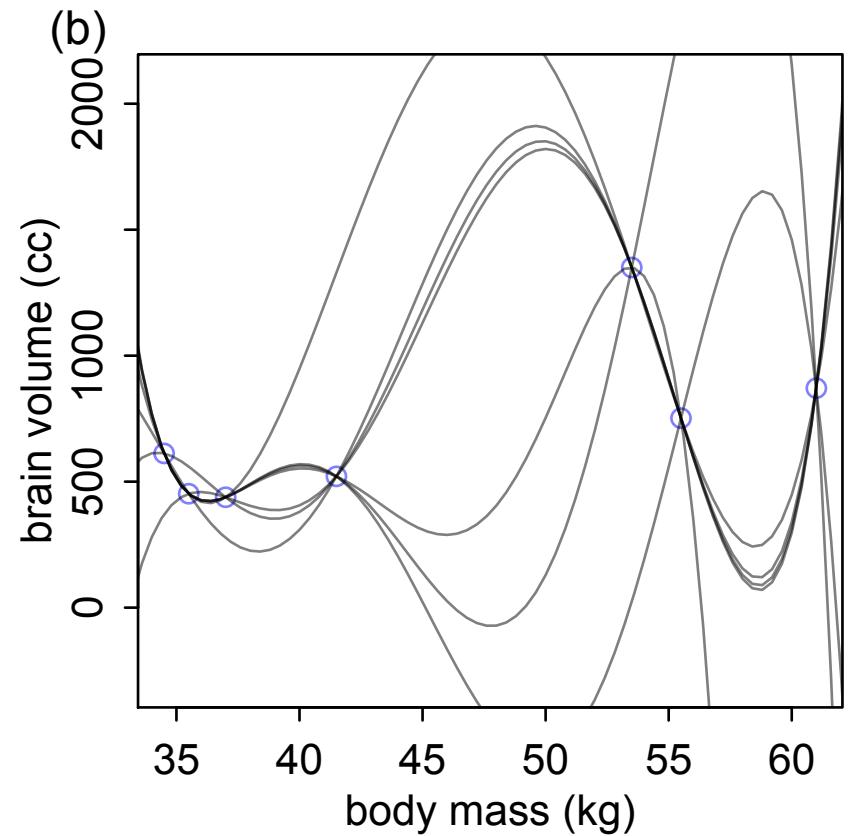
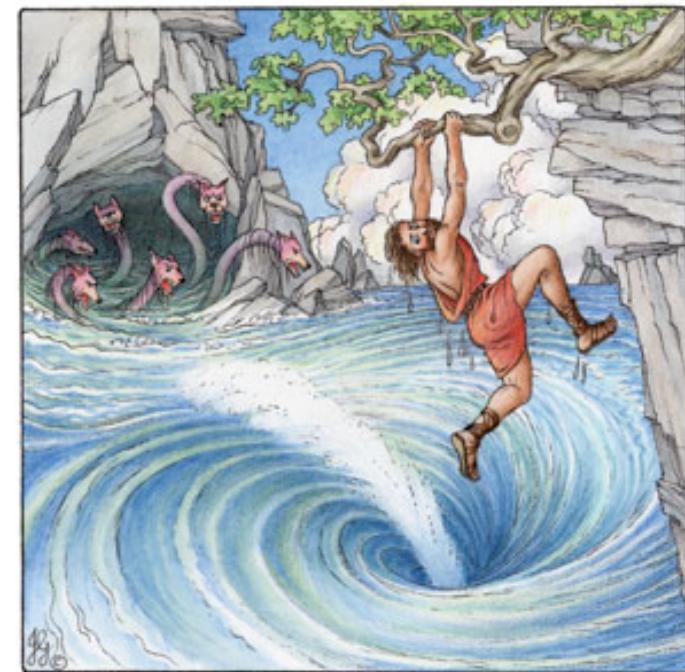


Figure 6.5

Importance of being *regular*

- Want the *regular* features of the sample
- Strategies
 - Cross-validation
 - Regularizing priors (penalized likelihood)
 - Information criteria
 - Science! (iterative group learning)
- Proper approach depends upon purpose



The road to AIC/DIC/WAIC

- What's a good target?
 - average correct?
 - average probability correct?
 - average log probability correct?
- How measure distance from the target?
- How can we estimate that distance?
- How can we adjust that estimate to account for overfitting?



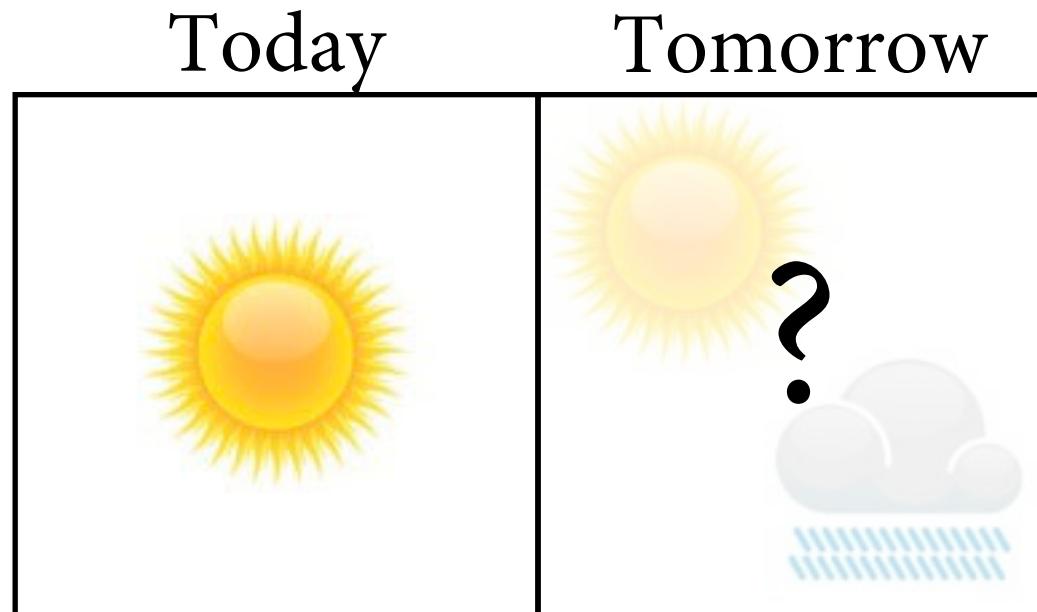
How far from truth?

- *Truth*: The real joint probability of events
 - *Truth* defines probability distribution
 - *Model* defines another
- Need a way to measure distance of a model from truth
- Distance needs to accommodate complexity of prediction task

Day	1	2	3	4	5	6	7	8	9	10
Observed	☁️	☁️	☁️	☀️	☀️	☀️	☀️	☀️	☀️	☀️
Points										
Current	-1	-1	-1	-0.6	-0.6	-0.6	-0.6	-0.6	-0.6	-0.6
Newcomer	-5	-5	-5	0	0	0	0	0	0	0

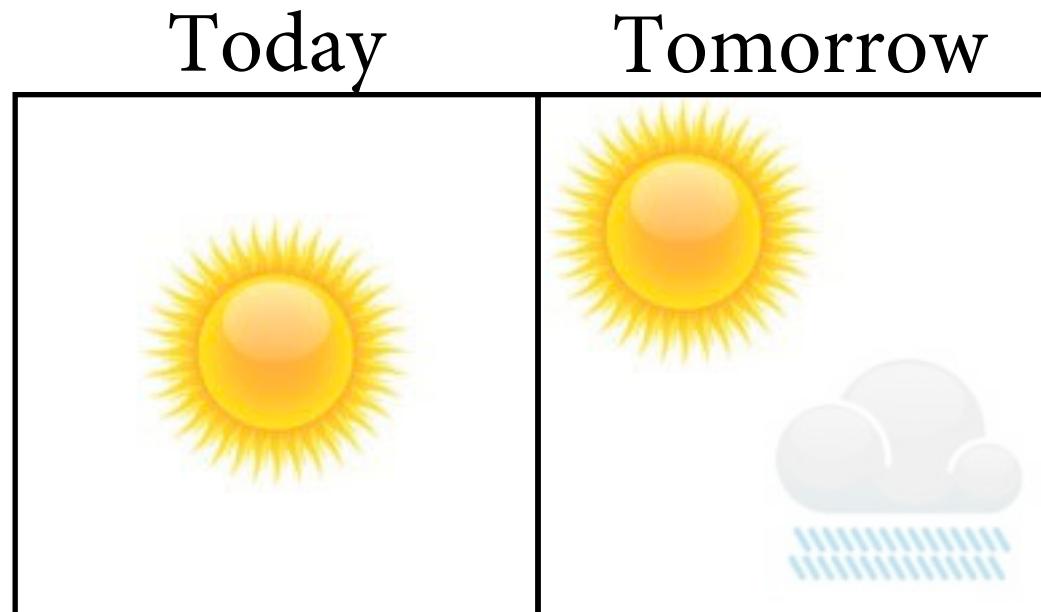
Information theory

- *Information:* Reduction in uncertainty caused by learning an outcome.



Information theory

- *Information*: Reduction in uncertainty caused by learning an outcome.

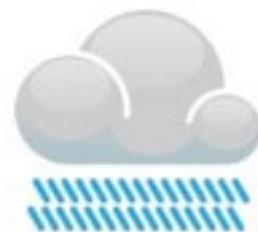


Today Tomorrow

Los Angeles



Glasgow



New York



Information theory

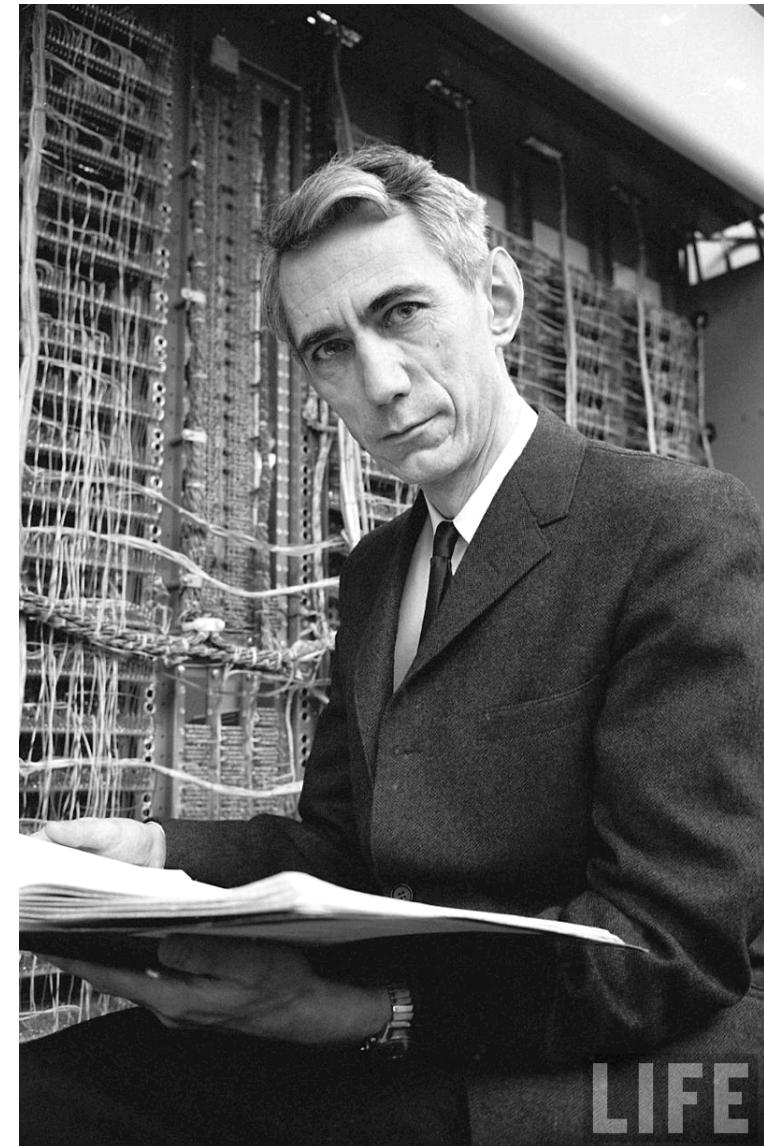
- *Information*: Reduction in uncertainty caused by learning an outcome.
- How to quantify *uncertainty*? Should be:
 1. Continuous
 2. Increasing with number of possible events
 3. Additive
- These criteria intuitive, but effectiveness is why we keep using them
 - Like Bayes: intuitive, but effectiveness is reason to use

Information entropy

- 1948, Claude Shannon derived *information entropy*:

$$H(p) = - \text{E} \log(p_i) = - \sum_{i=1}^n p_i \log(p_i)$$

Uncertainty in a probability distribution is average (minus) log-probability of an event.



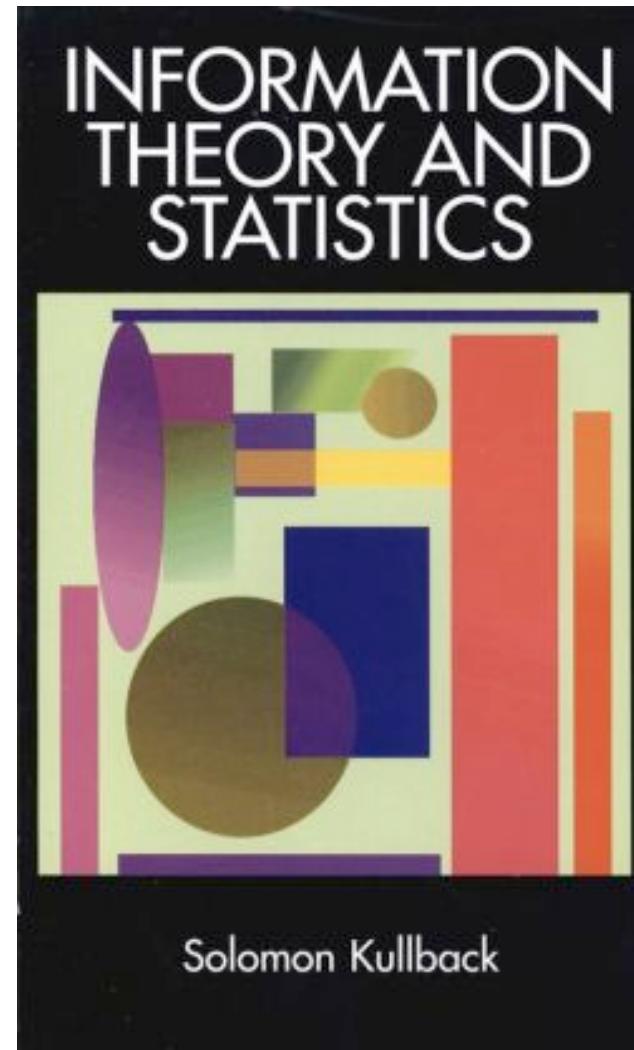
Shannon (1916–2001)

Entropy to accuracy

- Two probability distributions: p, q
- How accurate is q , for describing p ?
- Distance from q to p : *Divergence*

$$D_{\text{KL}}(p, q) = \sum_i p_i (\log(p_i) - \log(q_i))$$

Distance from q to p is the average difference in log-probability.



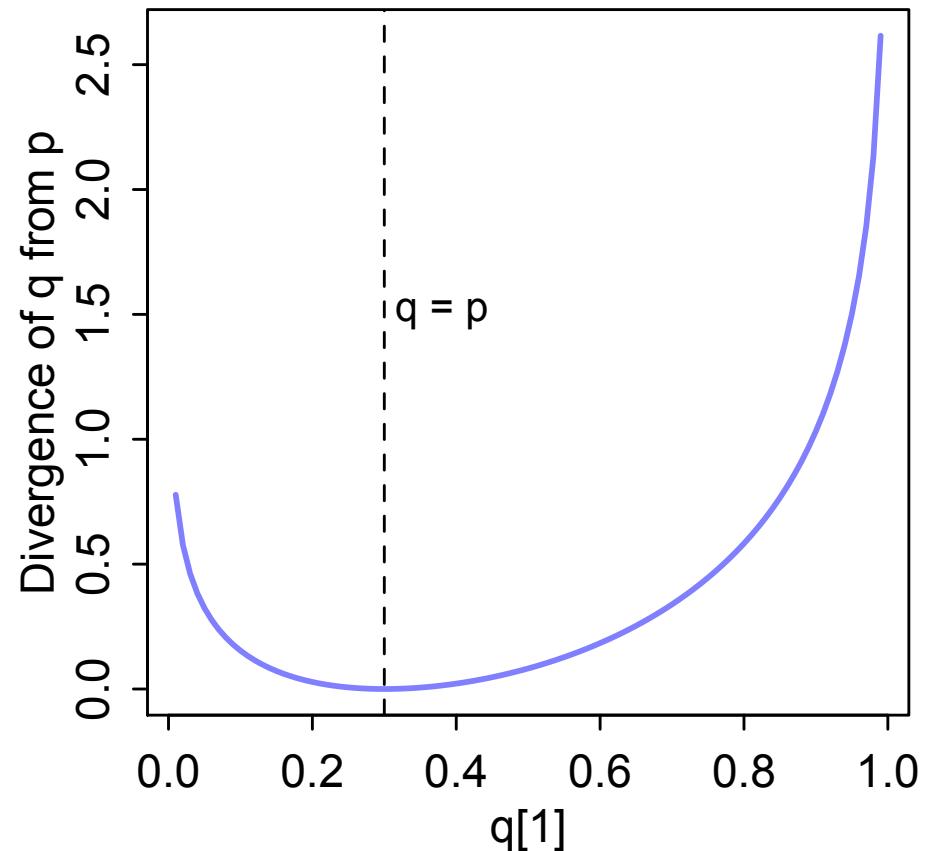
Entropy to accuracy

$$D_{\text{KL}}(p, q) = \sum_i p_i (\log(p_i) - \log(q_i))$$

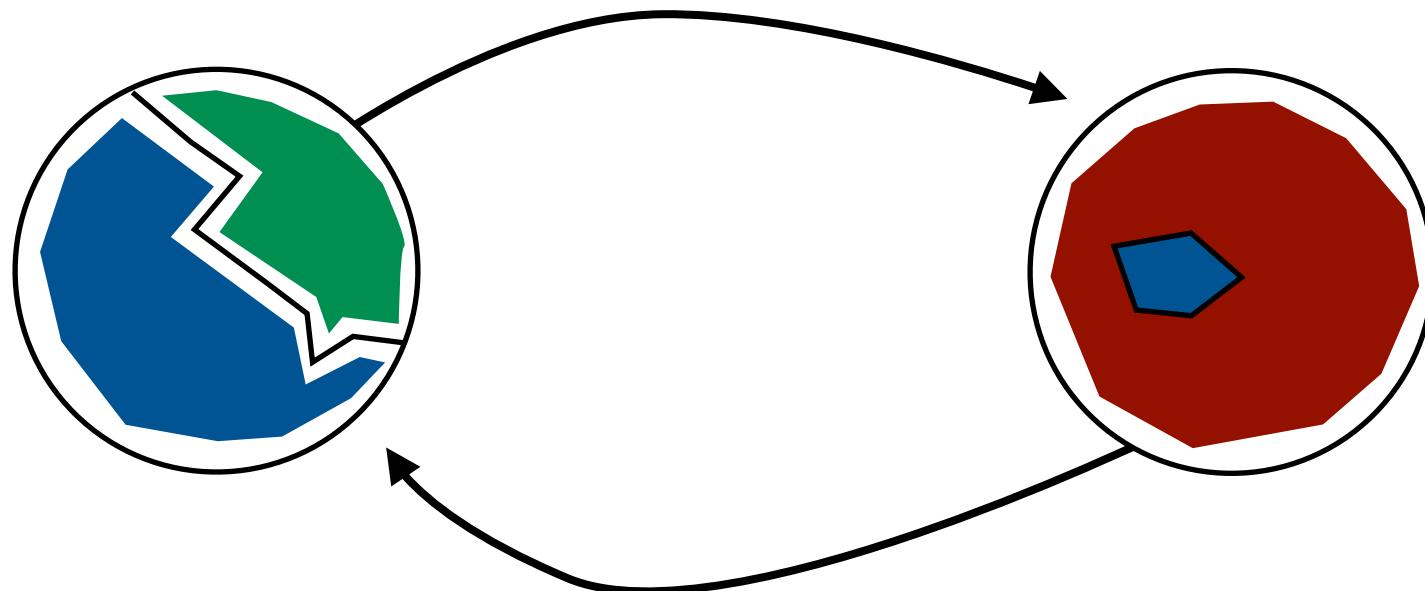
```
p <- c(0.3,0.7)

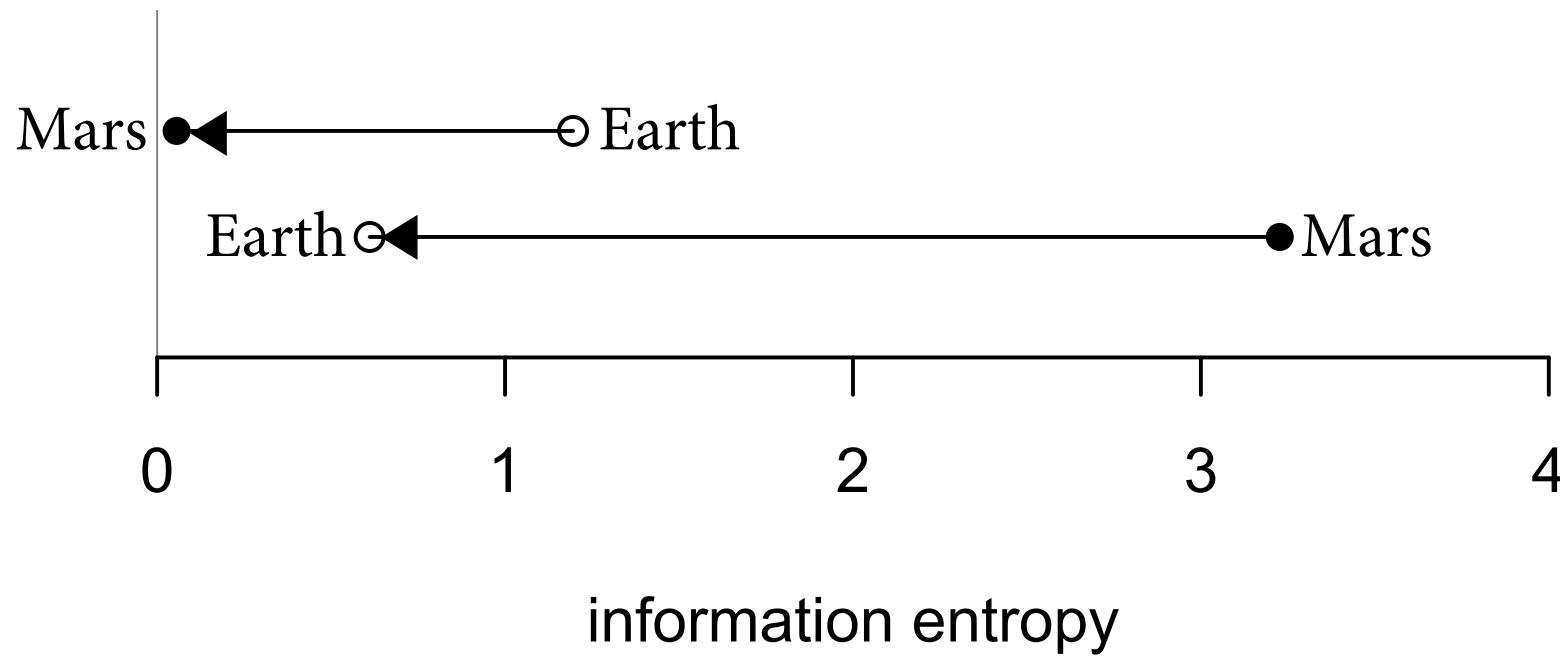
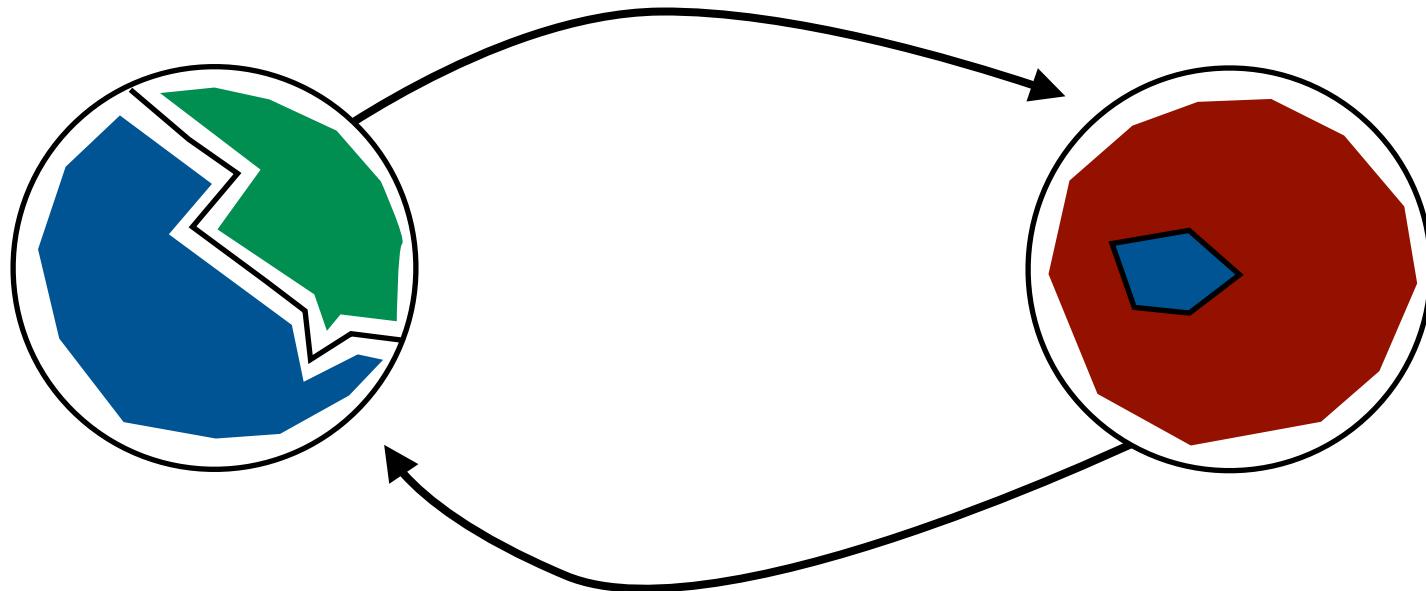
DKL <- function(p,q)
  sum(p*(log(p)-log(q)))

q1seq <- seq(from=0.01,to=0.99,by=0.01)
DKLseq <- sapply(q1seq,
  function(q1) DKL(p,c(q1,1-q1)) )
plot( q1seq , DKLseq )
```



Direction matters





Estimating divergence

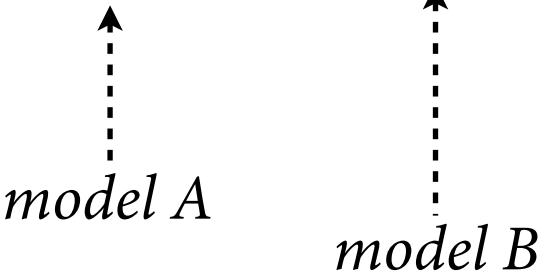
- How to estimate D_{KL} ?

$$D_{\text{KL}}(p, q) = \sum_i p_i (\log(p_i) - \log(q_i))$$


The diagram illustrates the components of the KL divergence formula. Two vertical dashed arrows point upwards from the labels "truth" and "model" to the corresponding variables p and q in the equation. The label "truth" is positioned to the left of the first arrow, and the label "model" is positioned to the right of the second arrow.

Estimating divergence

- Don't know p ! Don't need it. Focus on difference between two approximating models:

$$D_{\text{KL}}(p, q) - D_{\text{KL}}(p, r) = - \sum_i p_i (\log q_i - \log p_i) - \left(- \sum_i p_i (\log r_i - \log p_i) \right)$$


The diagram consists of two vertical dashed arrows. The left arrow points from the label "model A" to the term $\log q_i$ in the first sum. The right arrow points from the label "model B" to the term $\log r_i$ in the second sum.

Estimating divergence

- Don't know p ! Don't need it. Focus on difference between two approximating models:

$$D_{\text{KL}}(p, q) - D_{\text{KL}}(p, r) = - \sum_i p_i (\log q_i - \log p_i) - \left(- \sum_i p_i (\log r_i - \log p_i) \right)$$

\uparrow \uparrow
model A *model B*

$$= - \sum_i p_i (\log q_i - \log r_i) = -(\mathbb{E} \log q_i - \mathbb{E} \log r_i)$$

- log-probability scores (deviance e.g.) provide estimate of $\mathbb{E} \log q_i$

Deviance (classic estimate)

$$D(q) = -2 \sum_i \log(q_i)$$

- How *bad* the model is, not how *good*
- Compute it:
 - Compute log probability of each observation
 - Sum all of these log probabilities
 - Multiply by -2
- Common to use MAP estimates for probabilities, but can use entire posterior
 - Will do so later, when compute WAIC as estimate of deviance

The road to AIC/DIC/WAIC

- ✓ What's a good prediction?
- ✓ How far is the model from the target?
- ✓ How can we estimate that distance?
- How can we adjust that estimate to account for overfitting?



Everybody overfits

- A meta-model of forecasting:
 - Two samples: *training* and *testing*, size N
 - Fit model to *training* sample, get D_{train}
 - Use fit to *training* to compute D_{test}
 - Difference $D_{\text{test}} - D_{\text{train}}$ is overfitting

Everybody overfits

Data generating model:

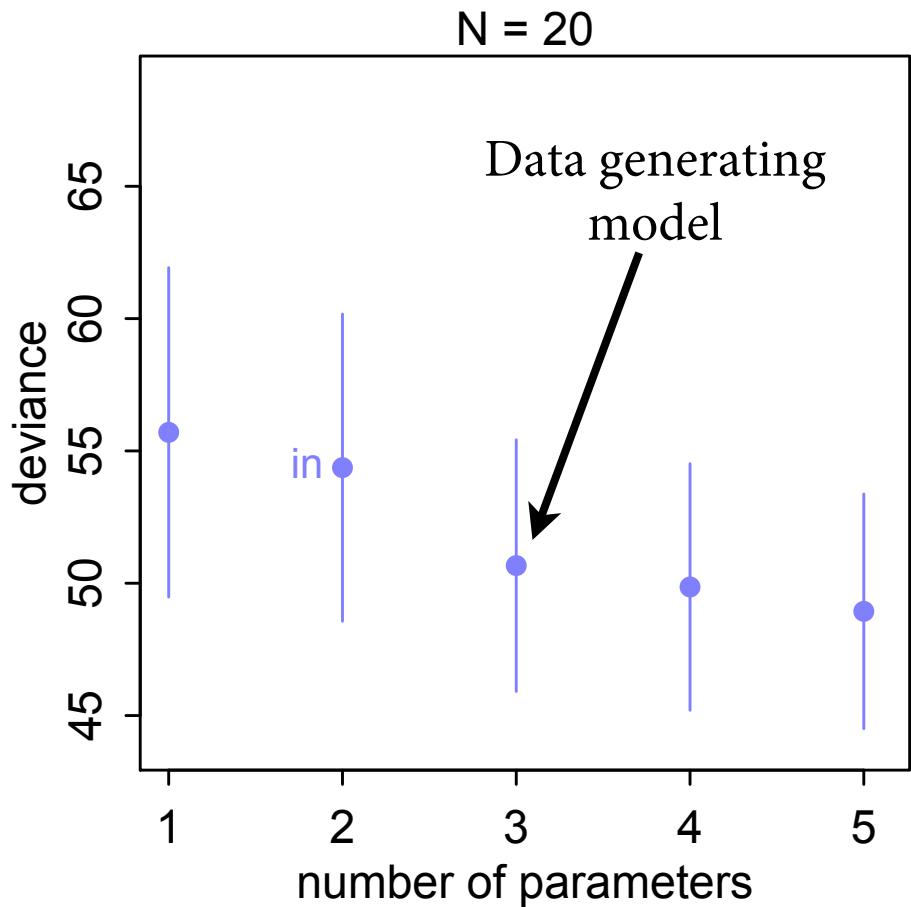
$$y_i \sim \text{Normal}(\mu_i, 1)$$
$$\mu_i = (0.15)x_{1,i} - (0.4)x_{2,i}$$

Models fit to data:

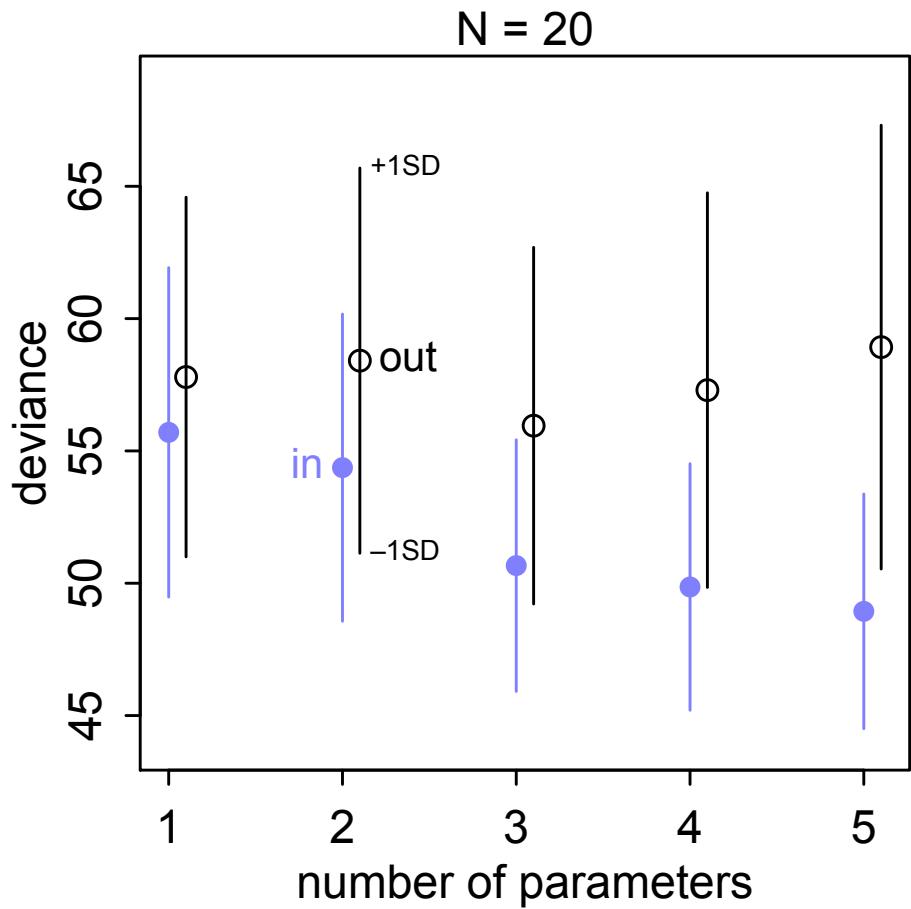
(flat priors)

$$\mu_i = \alpha$$
$$\mu_i = \alpha + \beta_1 x_{1,i}$$
$$\mu_i = \alpha + \beta_1 x_{1,i} + \beta_2 x_{2,i}$$
$$\mu_i = \alpha + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \beta_3 x_{3,i}$$
$$\mu_i = \alpha + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \beta_3 x_{3,i} + \beta_4 x_{4,i}$$

Everybody overfits



Everybody overfits



Everybody overfits

