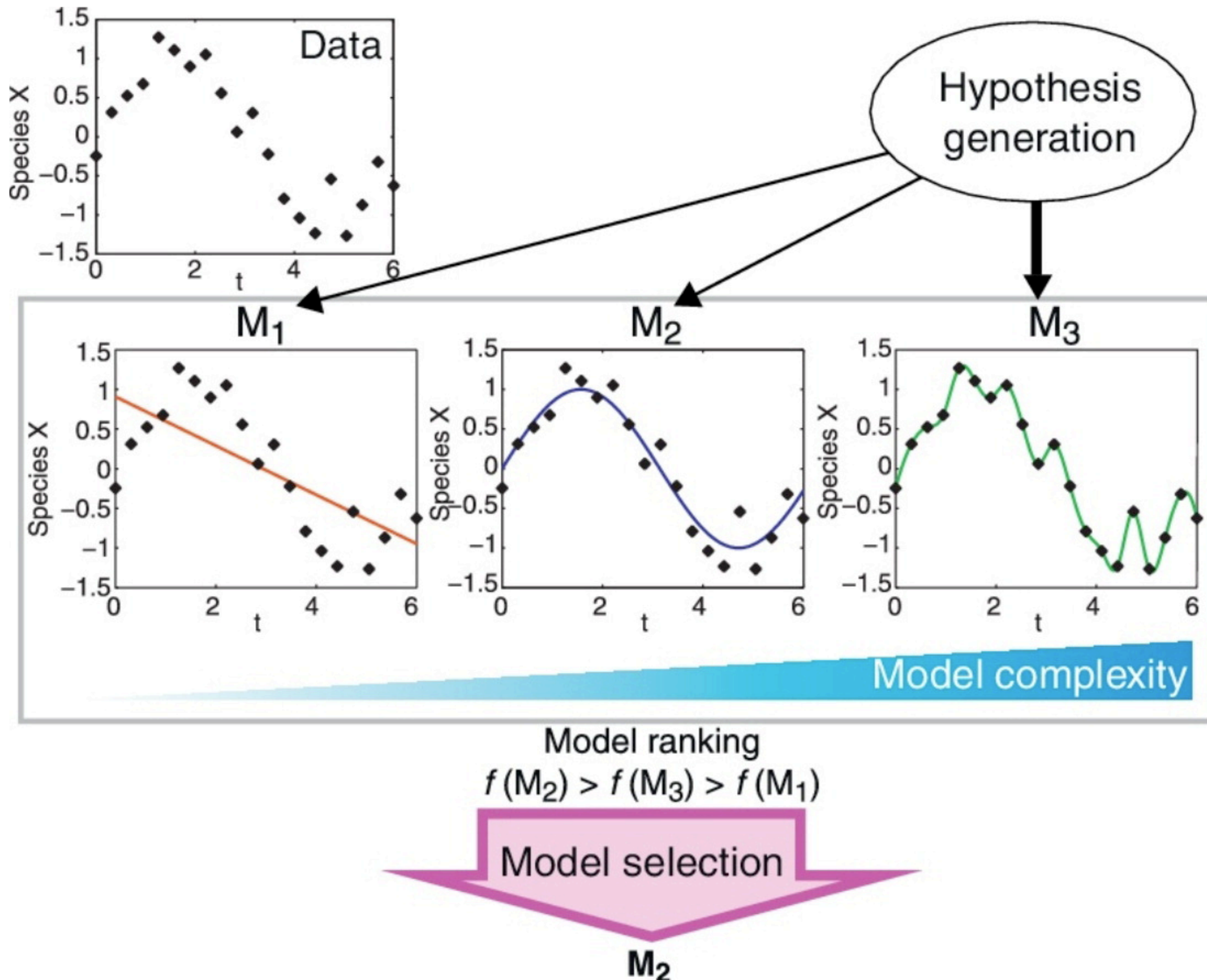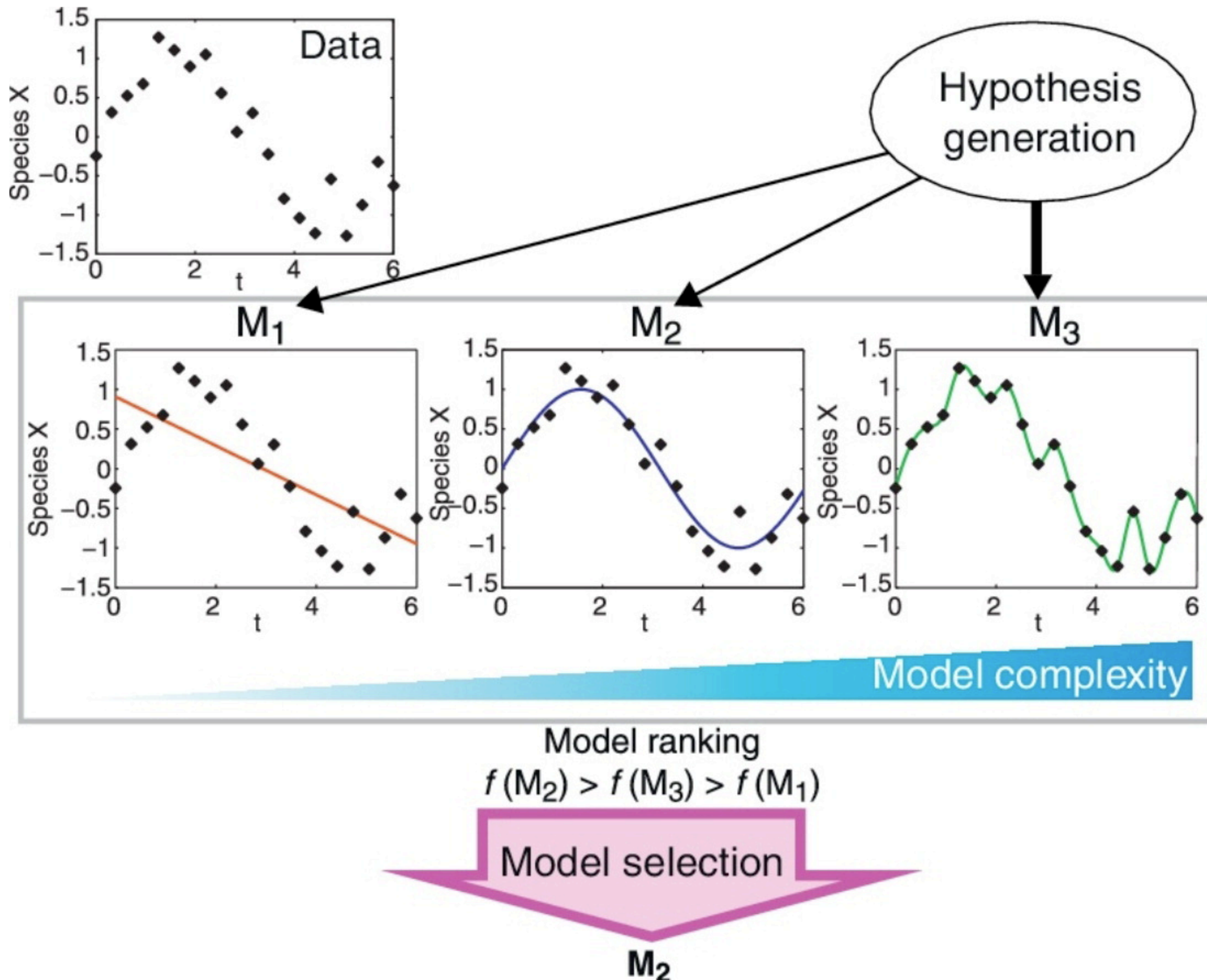# COSMOS

## Tutorial 3: Model comparison and robustness

**Wataru Toyokawa & Charley Wu**
**July 6th**

1. Generate hypotheses

2. Build models for each hypothesis

3. Fit models to data

4. Determine the best model

5. Interpretation

1. Generate hypotheses

2. Build models for each hypothesis

3. Fit models to data

4. **Determine the best model**

5. **Interpretation**

# Outline

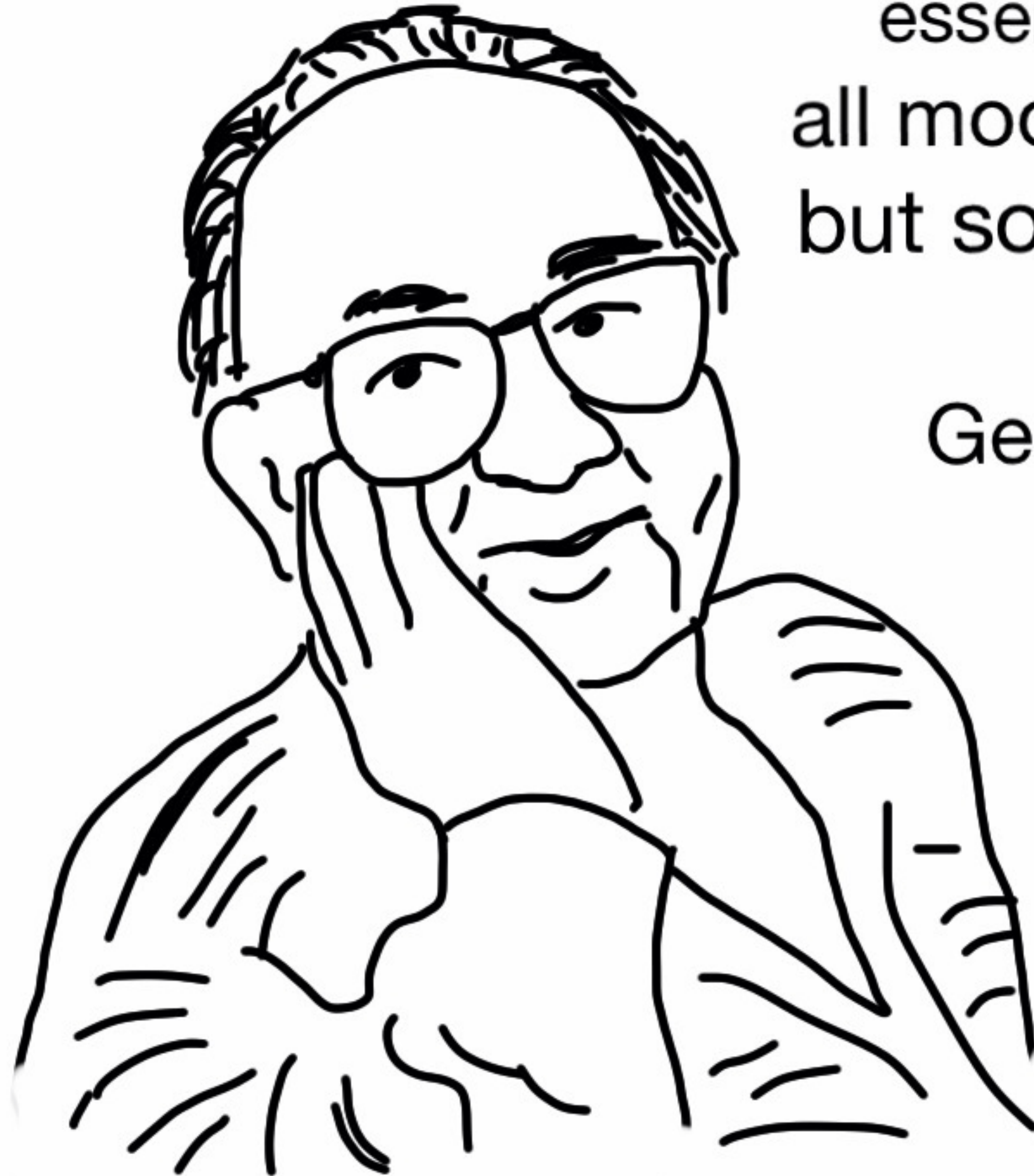Part 1. Model Comparison

Part 2. Robustness

# Part 1. Model Comparison

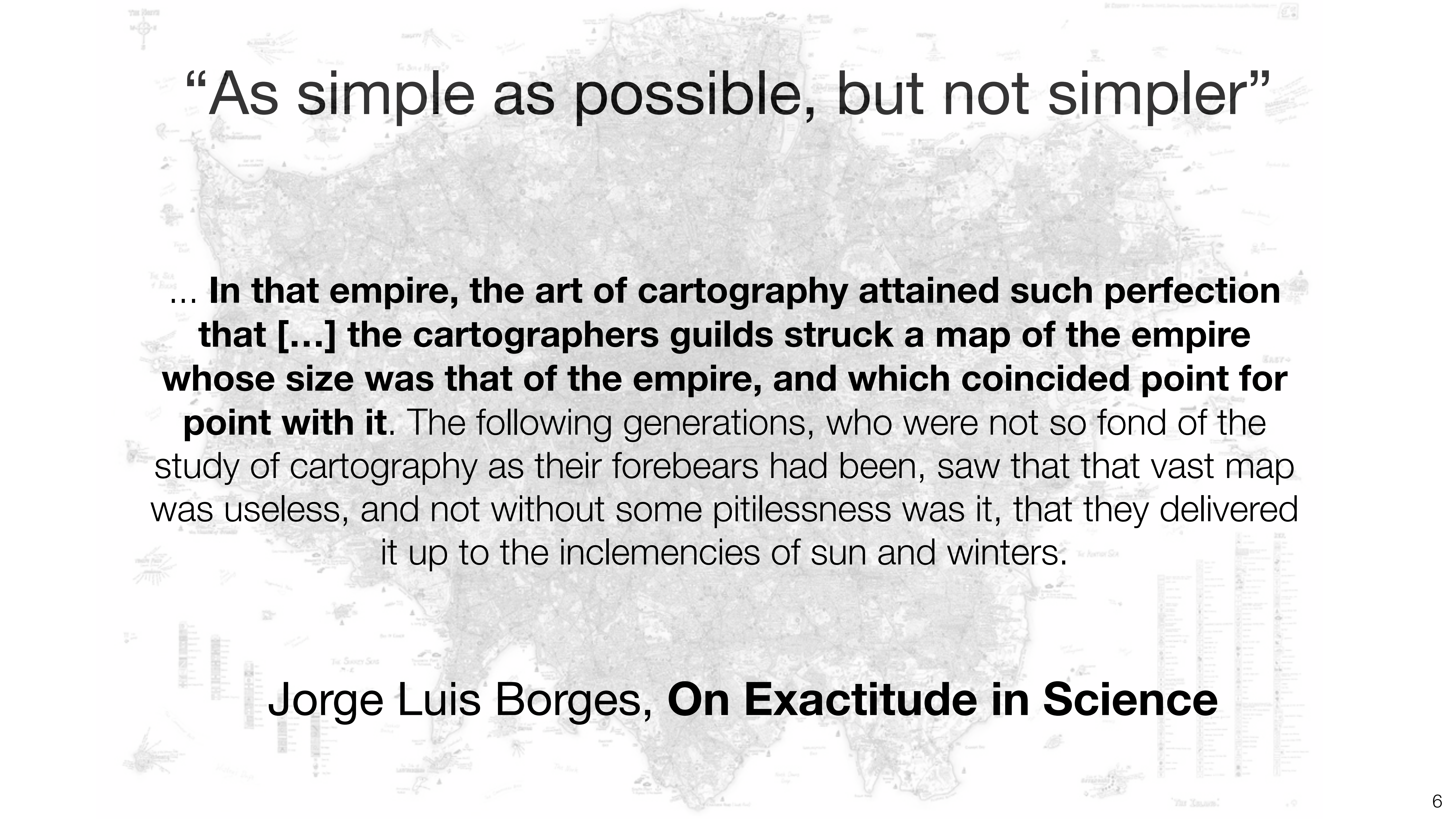# What makes a good model?

# What makes a good model?

essentially,
all models are wrong,
but some are useful

George E. P. Box

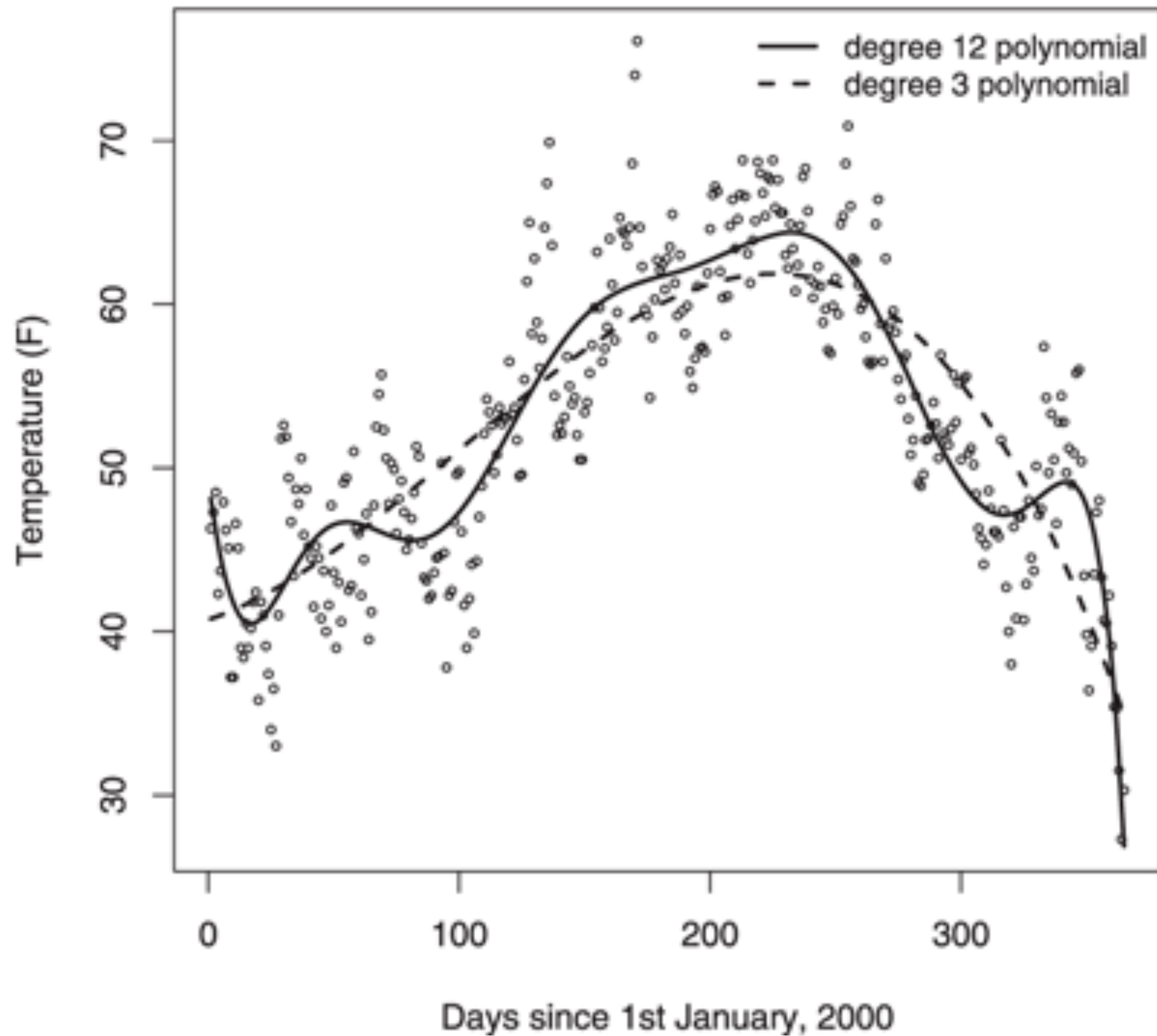# "As simple as possible, but not simpler"

# "As simple as possible, but not simpler"

… **In that empire, the art of cartography attained such perfection that […] the cartographers guilds struck a map of the empire whose size was that of the empire, and which coincided point for point with it**. The following generations, who were not so fond of the study of cartography as their forebears had been, saw that that vast map was useless, and not without some pitilessness was it, that they delivered it up to the inclemencies of sun and winters.
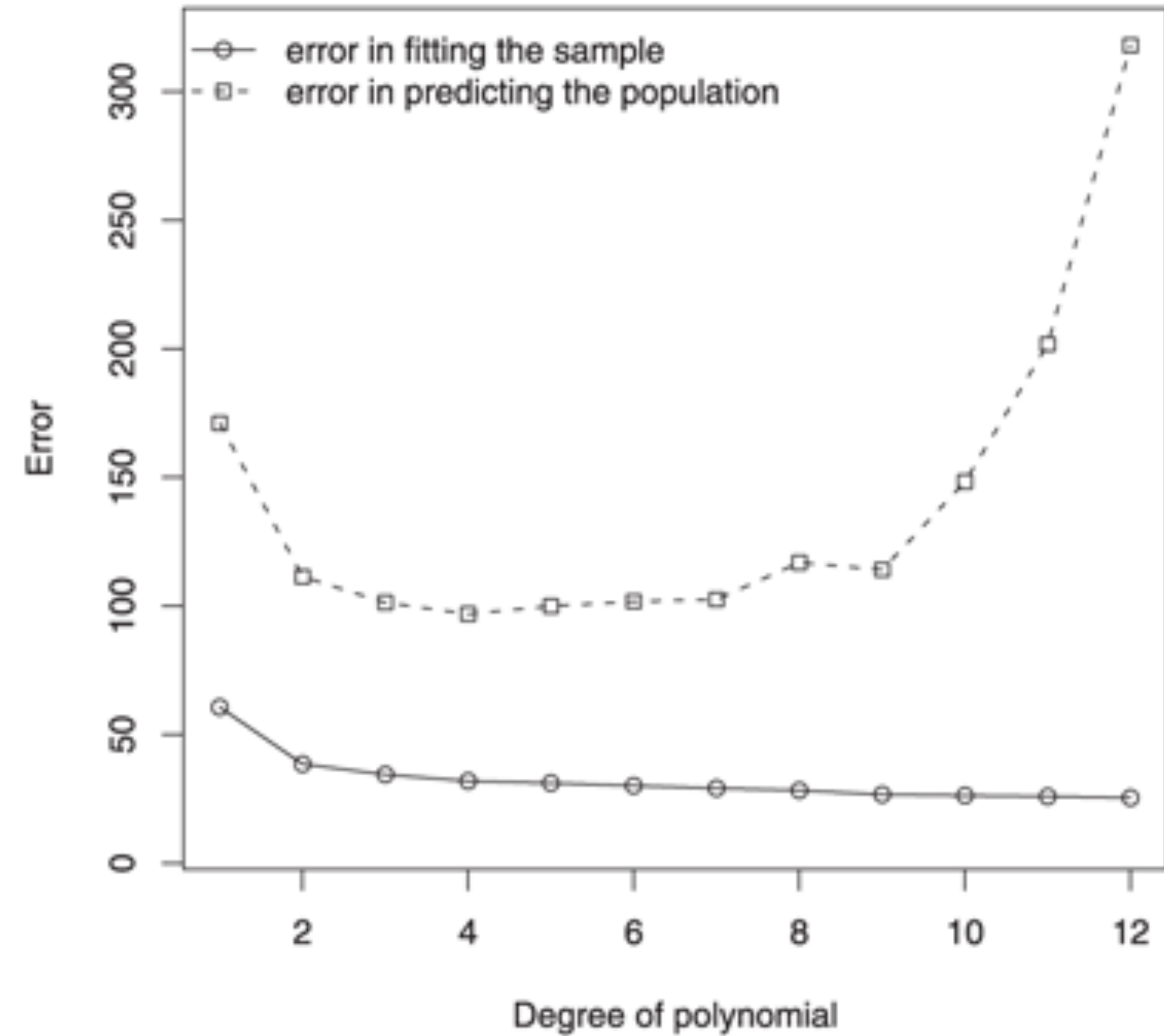
Jorge Luis Borges, **On Exactitude in Science**

# "As simple as possible, but not simpler"
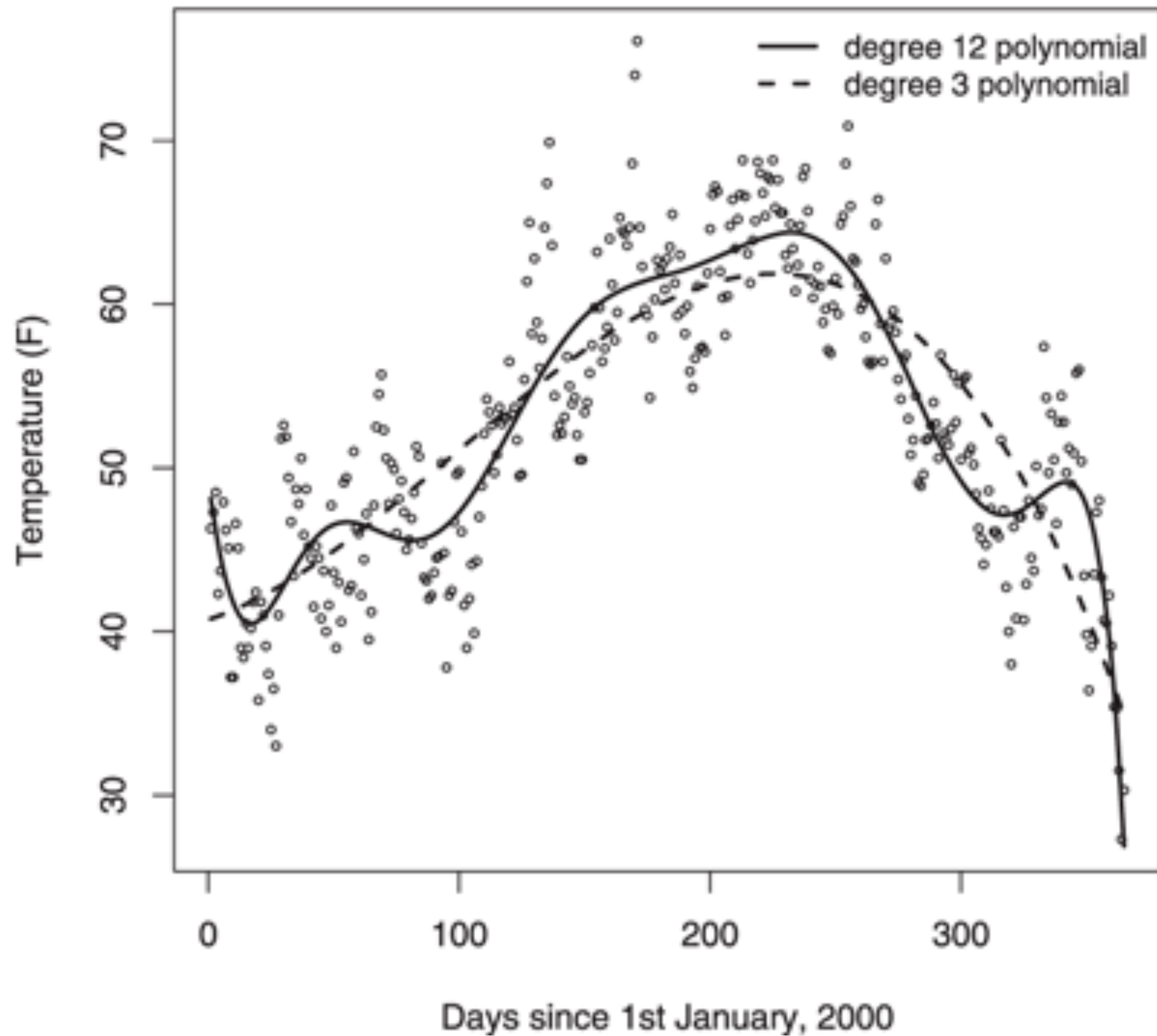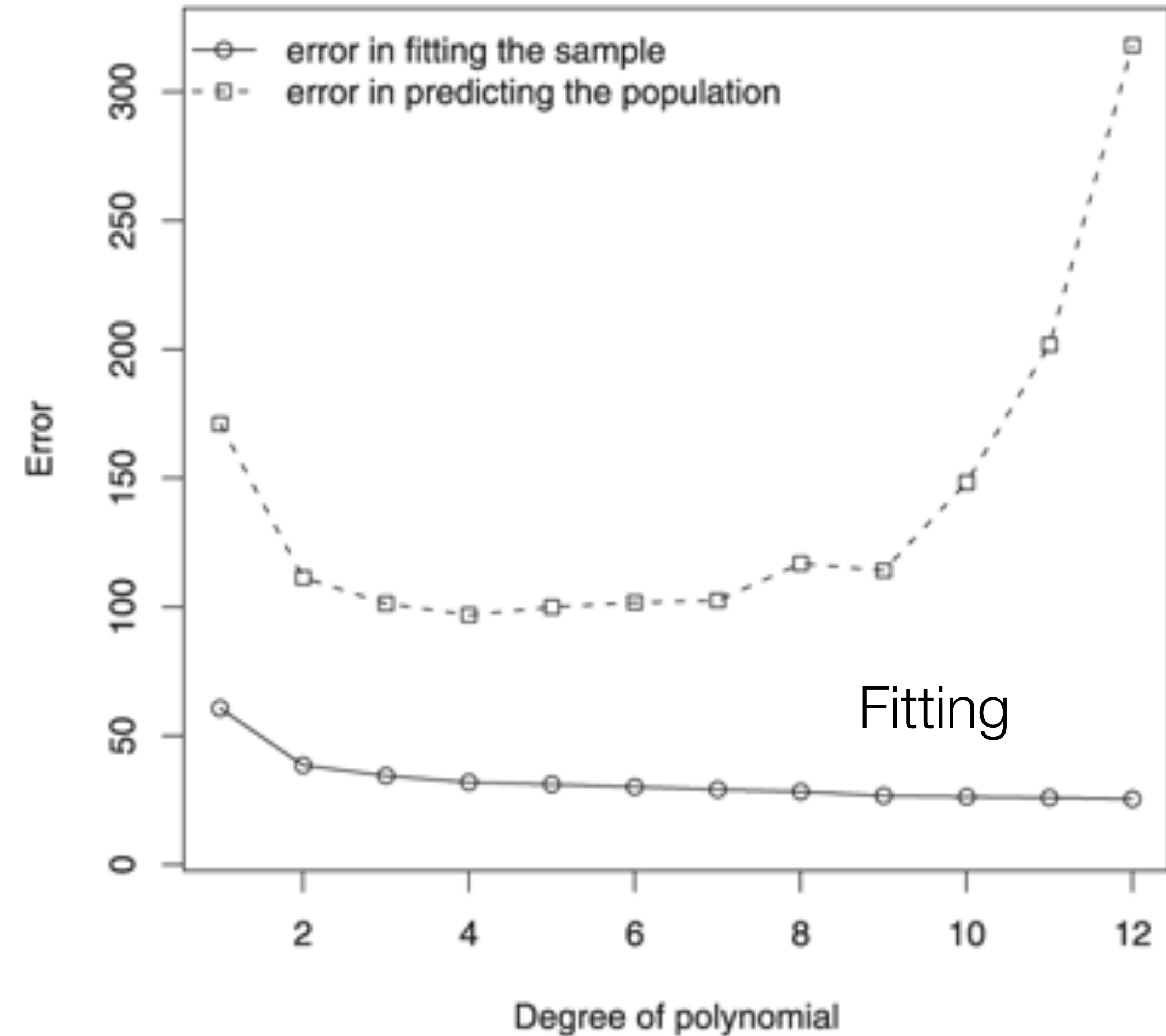


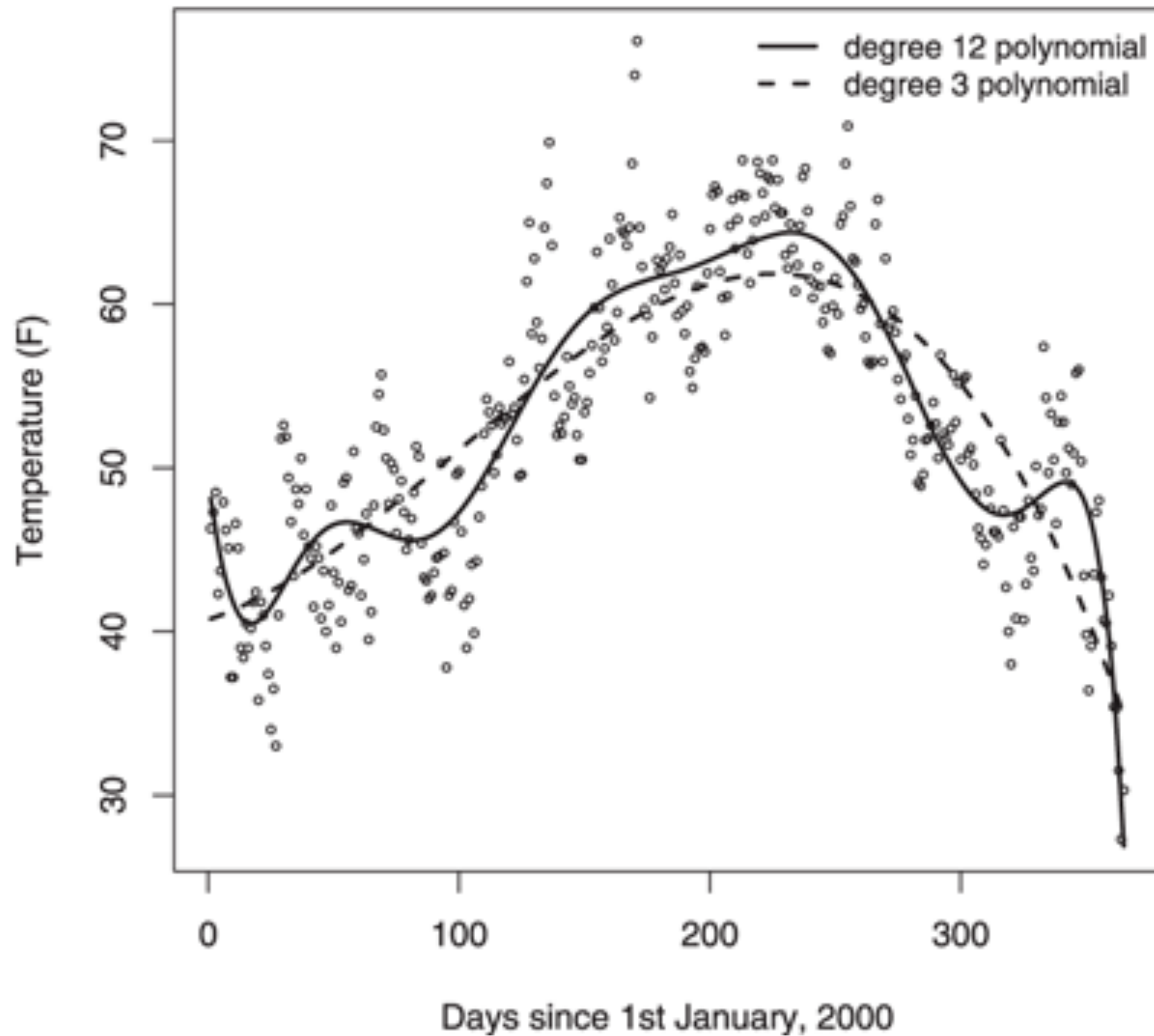London's daily temperature in 2000



Model performance for London 2000 temperatures

# "As simple as possible, but not simpler"



London's daily temperature in 2000



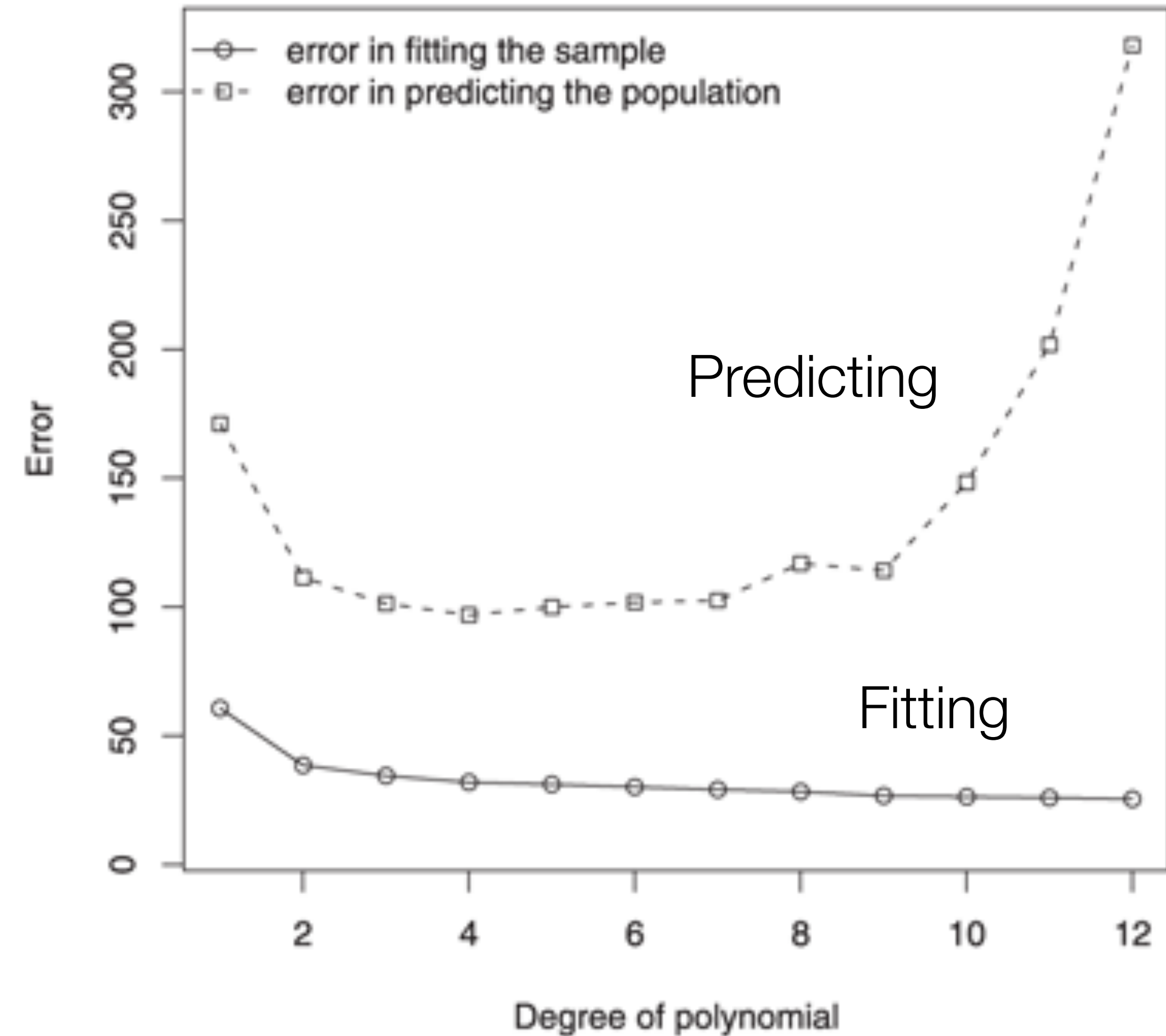Model performance for London 2000 temperatures

# "As simple as possible, but not simpler"



London's daily temperature in 2000



Model performance for London 2000 temperatures

# Goodness of Fit



Simplicity ←――――――――――――――――――――→ Fit

# Goodness of Fit Measures

| | **Maximum Likelihood** $P(D\,|\,m,\hat{\theta})$ | **Bayesian Model Selection** $\dfrac{P(D\,|\,m_1)}{P(D\,|\,m_2)}$ |
|---|---|---|
| **Penalizing for parameters** | Akaike's Information Criterion (AIC) | Bayesian Information Criterion (BIC) |
| **Prediction error/ Bayesian Occam's Razor** | Cross-validation loss | Model evidence using Markov Chain Monte Carlo (MCMC) |

# Goodness of Fit Measures

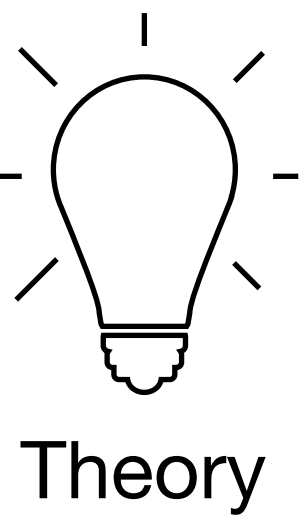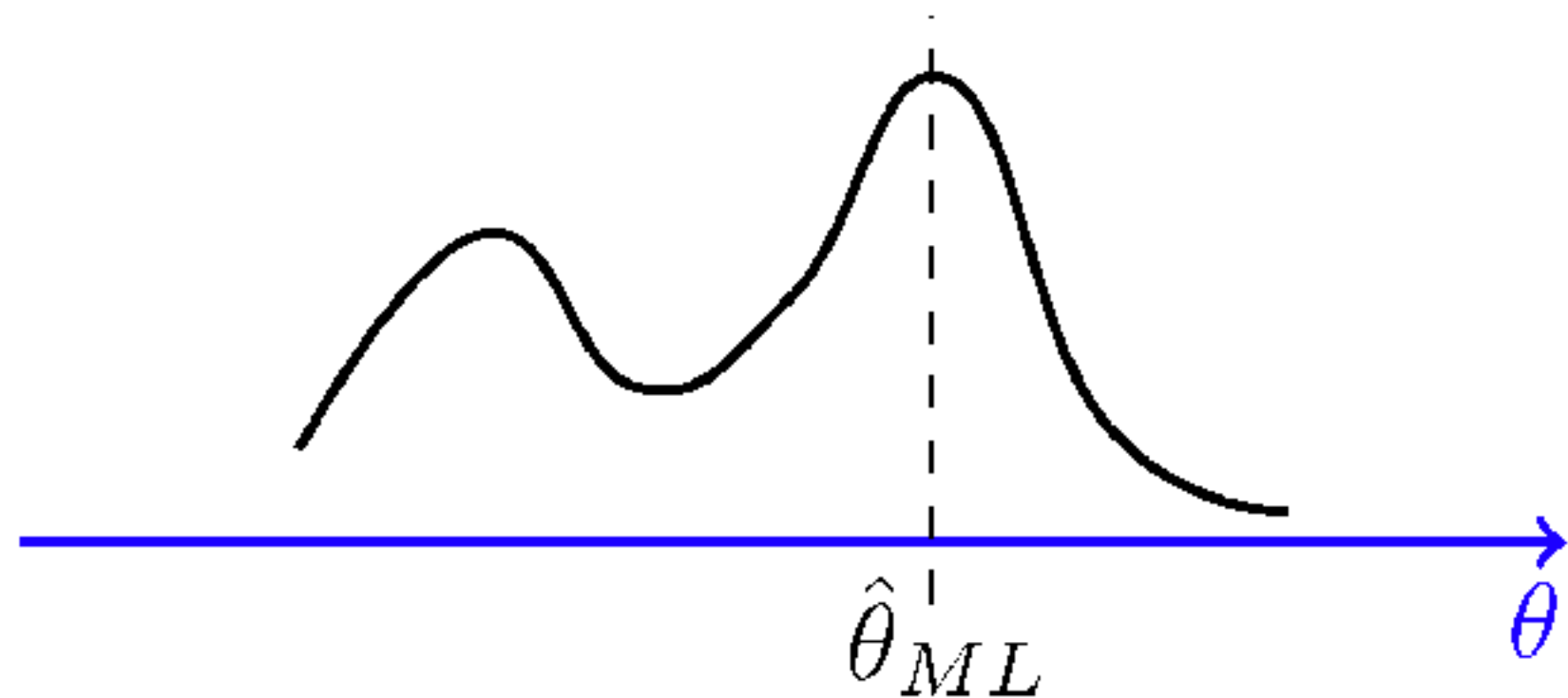| | **Maximum Likelihood** $P(D \mid m, \hat{\theta})$ | **Bayesian Model Selection** $\dfrac{P(D \mid m_1)}{P(D \mid m_2)}$ |
|---|---|---|
| **Penalizing for parameters** | Akaike's Information Criterion (AIC) | Bayesian Information Criterion (BIC) |
| **Prediction error/ Bayesian Occam's Razor** | Cross-validation loss | Model evidence using Markov Chain Monte Carlo (MCMC) |

Theory

Practice

# Maximum likelihood estimation (MLE)

VS.

# Bayesian model selection

- **Goal**: Quantify the goodness fit for a single set of parameter values $\hat{\theta}$ that provides the best fit to the data:

$$\arg \max_{\hat{\theta}} P(D \,|\, m, \hat{\theta})$$

- Overfitting is avoided by penalizing for the number of parameters (e.g., AIC) or using cross-validation to test predictive power



$\hat{\theta}_{ML}$   $\theta$

- **Goal**: quantify how well a given model $m$ captures the data using the *marginal likelihood*:

$$P(D \,|\, m) = \int P(D \,|\, m, \theta) P(\theta \,|\, m) d\theta$$

- This integrates over all possible parameter values, allowing for a natural penalization of more complex models (i.e., Bayesian Occam's Razor)

  - You don't only test the model at it's best, but also at it's worse

- Intractable in most settings, so approximated using BIC or through MCMC sampling

# Goodness of Fit Measures

| | **Maximum Likelihood** | **Bayesian Model Selection** |
|---|---|---|
| | $P(D \mid m, \hat{\theta})$ | $\dfrac{P(D \mid m_1)}{P(D \mid m_2)}$ |
| **Penalizing for parameters** | Akaike's Information Criterion (AIC) | Bayesian Information Criterion (BIC) |
| **Prediction error/ Bayesian Occam's Razor** | Cross-validation loss | Model evidence using Markov Chain Monte Carlo (MCMC) |

# Akaike's Information Criterion (AIC)

$$AIC = -2 \log P(D \,|\, \hat{\theta}) + 2k$$
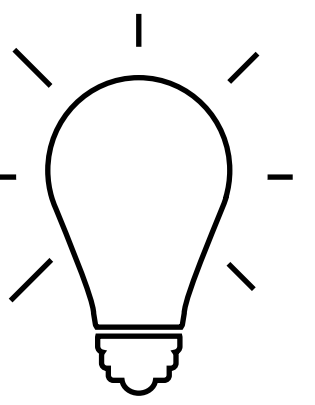
# Akaike's Information Criterion (AIC)

$$AIC = -\underbrace{2 \log P(D|\hat{\theta})}_{\text{Fit}} + 2k$$

1. Perform MLE and compute 2x the negative Log Likelihood (aka *deviance*)

# Akaike's Information Criterion (AIC)

$$AIC = -2 \log P(D | \hat{\theta}) + 2k$$

Fit        Complexity

1. Perform MLE and compute 2x the negative Log Likelihood (aka *deviance*)

2. Penalize by adding an additional loss that is 2x the number of parameters $k$
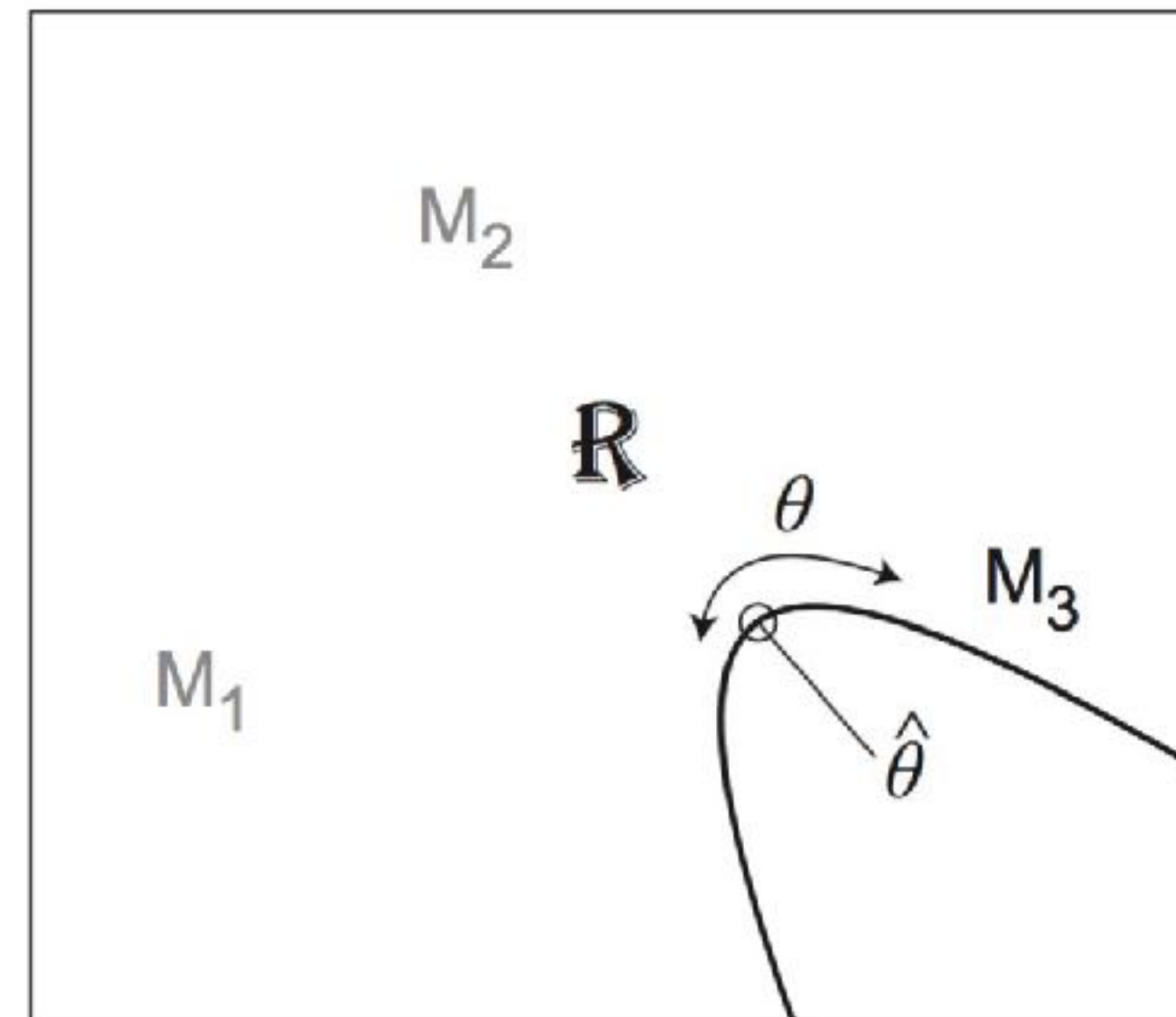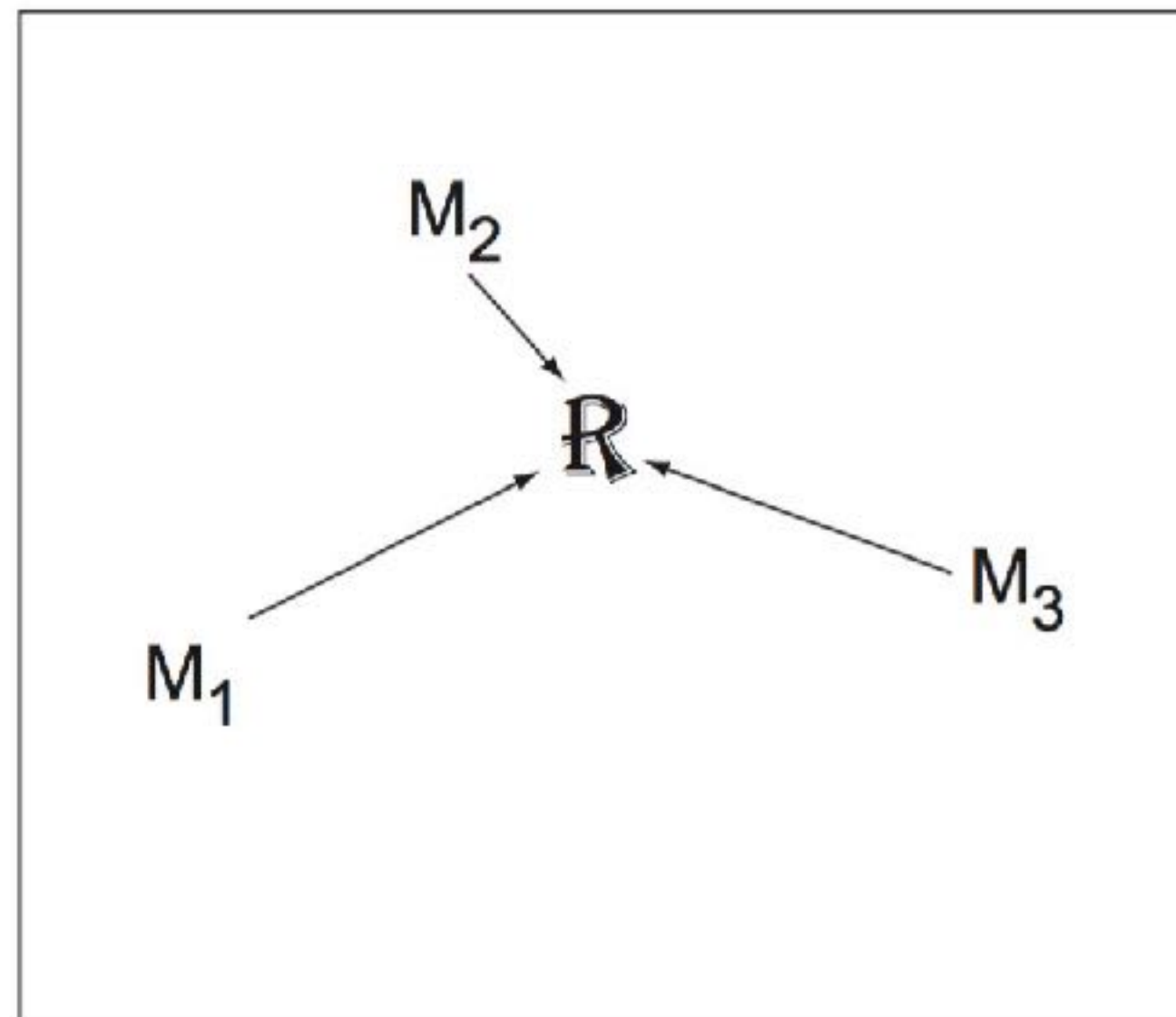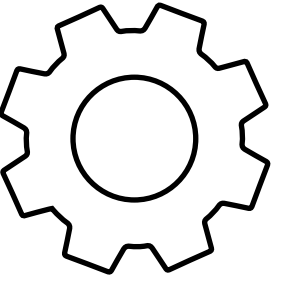
# Akaike's Information Criterion (AIC)

A measure of the relative information lost by a given model that is trying to capture some objective reality R(x)

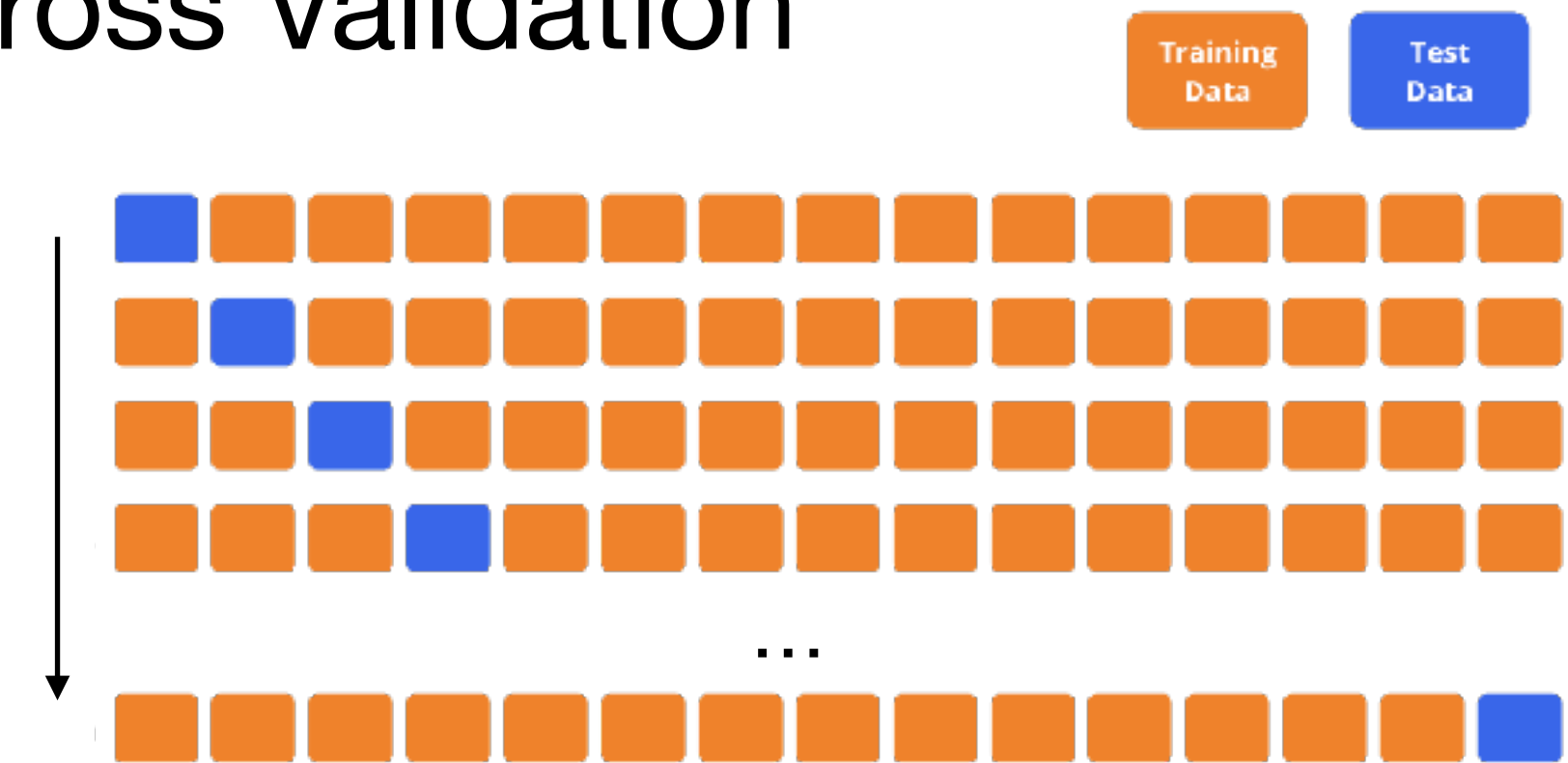$$KL = \int R(x) \log R(x) dx - \int R(x) \log P(x \mid \theta) dx$$

# Akaike's Information Criterion (AIC)

Asymptotically, AIC is equivalent to Leave-One-Out-Cross Validation (Stone, 1977)

- for linear regression and mixed-effects regression

- in the limit of infinite data

… yet for it's simplicity, AIC is commonly used for non-linear models and certainly always short of infinite data

In practice, AIC can be considered the most lax of the goodness of fit measures we introduce, and is more prone to preferring an overfit model

# Bayesian Information Criterion (BIC)

$$BIC = -2\log P(D\,|\,\hat{\theta}) + k\log n$$
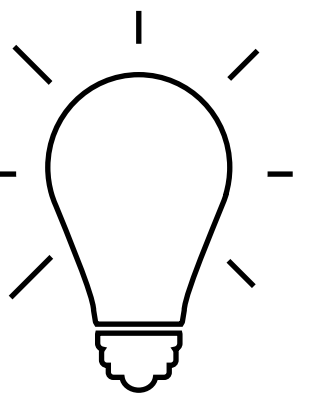
# Bayesian Information Criterion (BIC)

$$BIC = -\underbrace{2\log P(D \mid \hat{\theta})}_{\text{Fit}} + k \log n$$

1. Perform MLE and compute 2x the negative Log Likelihood (aka *deviance*)

# Bayesian Information Criterion (BIC)

$$BIC = \underbrace{-2\log P(D \mid \hat{\theta})}_{\text{Fit}} + \underbrace{k\log n}_{\text{Complexity}}$$

1. Perform MLE and compute 2x the negative Log Likelihood (aka *deviance*)

2. Penalize by adding an additional loss that is **the number of parameters $k$ times the log of the number of data points $n$**
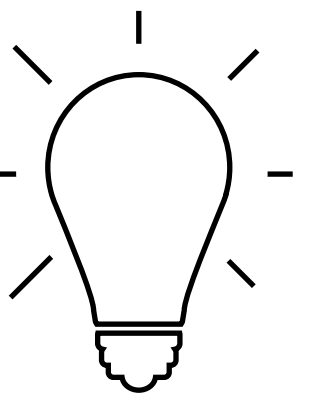
# Bayesian Information Criterion (BIC)

Bayesian model selection sometimes relies on Bayes Factors (BFs) to quantify the evidence of one model $m_1$ over another $m_2$

$$BF_{1,2} = \frac{P(D \mid m_1)}{P(D \mid m_2)}$$

- BF = 1; no evidence for either model

- BF >> 1; evidence for model 1

- BF << 1; evidence for model 2

BIC approximates the marginal likelihood using the MLE and by making some assumptions about the prior (Schwartz, 1975)
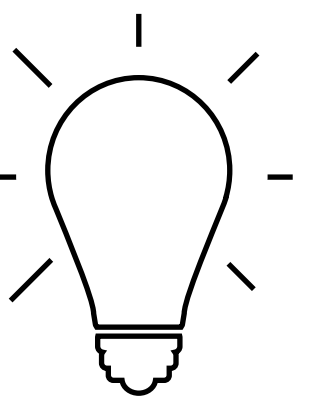
# Bayesian Information Criterion (BIC)

Bayesian model selection sometimes relies on Bayes Factors (BFs) to quantify the evidence of one model $m_1$ over another $m_2$

$$BF_{1,2} = \frac{P(D|m_1)}{P(D|m_2)}$$

- BF = 1; no evidence for either model

- BF >> 1; evidence for model 1

- BF << 1; evidence for model 2

$$P(D|m) = \int P(D|\theta, m)P(\theta|m)d\theta$$

BIC approximates the marginal likelihood using the MLE and by making some assumptions about the prior (Schwartz, 1975)

# Bayesian Information Criterion (BIC)

Bayesian model selection sometimes relies on Bayes Factors (BFs) to quantify the evidence of one model $m_1$ over another $m_2$
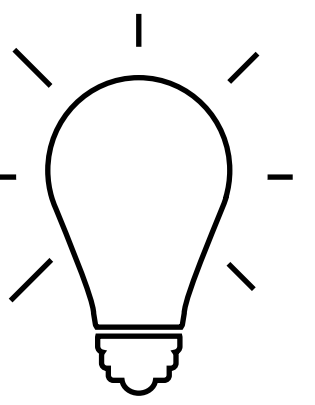
- BF = 1; no evidence for either model

- BF >> 1; evidence for model 1

- BF << 1; evidence for model 2

BIC approximates the marginal likelihood using the MLE and by making some assumptions about the prior (Schwartz, 1975)

$$BF_{1,2} = \frac{P(D|m_1)}{P(D|m_2)}$$

$$P(D|m) = \int P(D|\theta, m)P(\theta|m)d\theta$$

$$P(D|m) \approx BIC$$

# Bayesian Information Criterion (BIC)

Bayesian model selection sometimes relies on Bayes Factors (BFs) to quantify the evidence of one model $m_1$ over another $m_2$

- BF = 1; no evidence for either model

- BF >> 1; evidence for model 1

- BF << 1; evidence for model 2

BIC approximates the marginal likelihood using the MLE and by making some assumptions about the prior (Schwartz, 1975)

$$BF_{1,2} = \frac{P(D|m_1)}{P(D|m_2)}$$

$$P(D|m) = \int P(D|\theta, m)P(\theta|m)d\theta$$

$$P(D|m) \approx BIC$$

$$BF_{1,2} = \exp\left(-\frac{1}{2}(BIC_1 - BIC_2)\right)$$

# Bayesian Information Criterion (BIC)

Bayesian interpretation is not without controversy (see Lewandowsky & Farrell, 2010 for a discussion) and the assumptions are hardly ever met or even unpacked
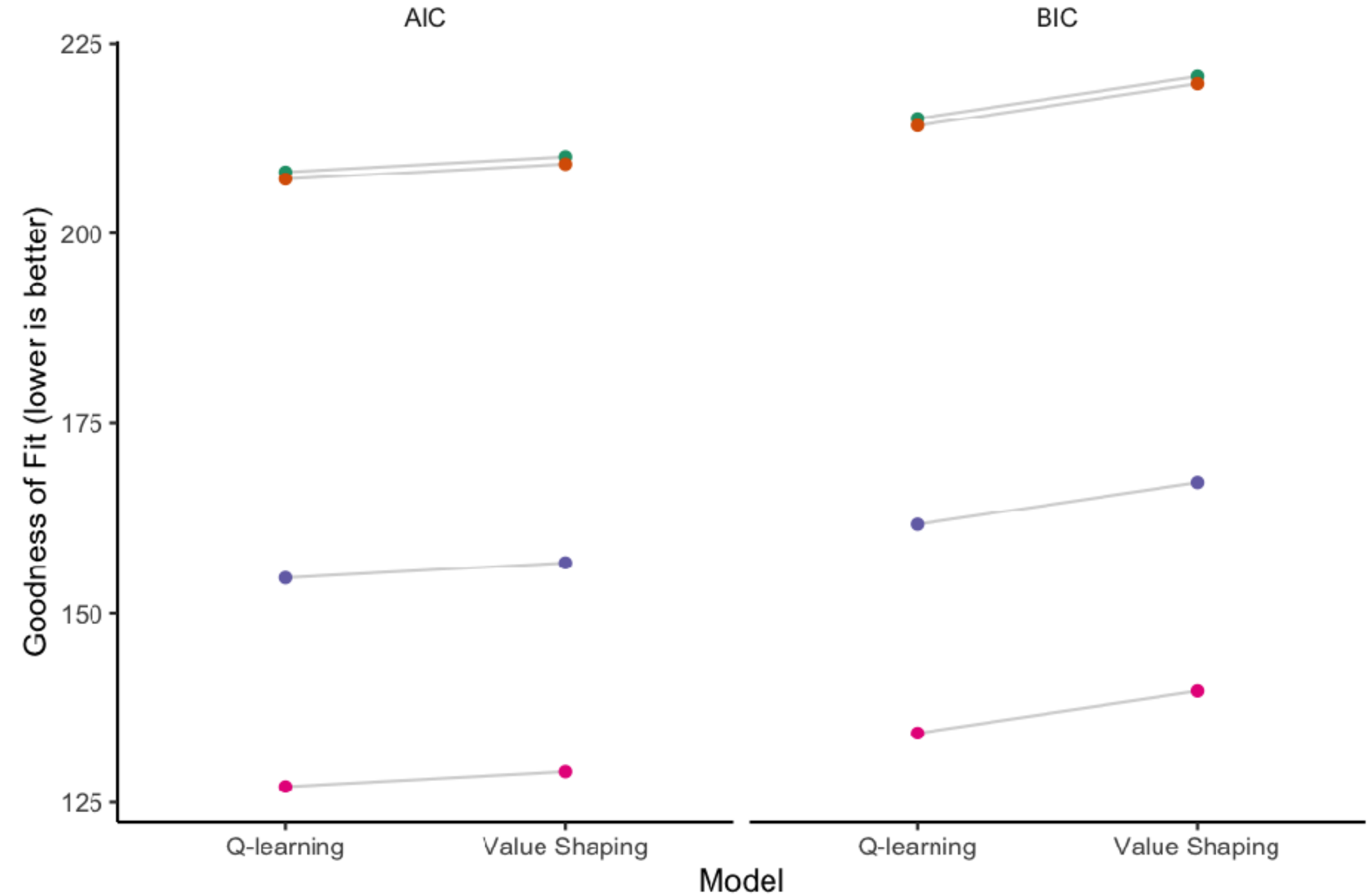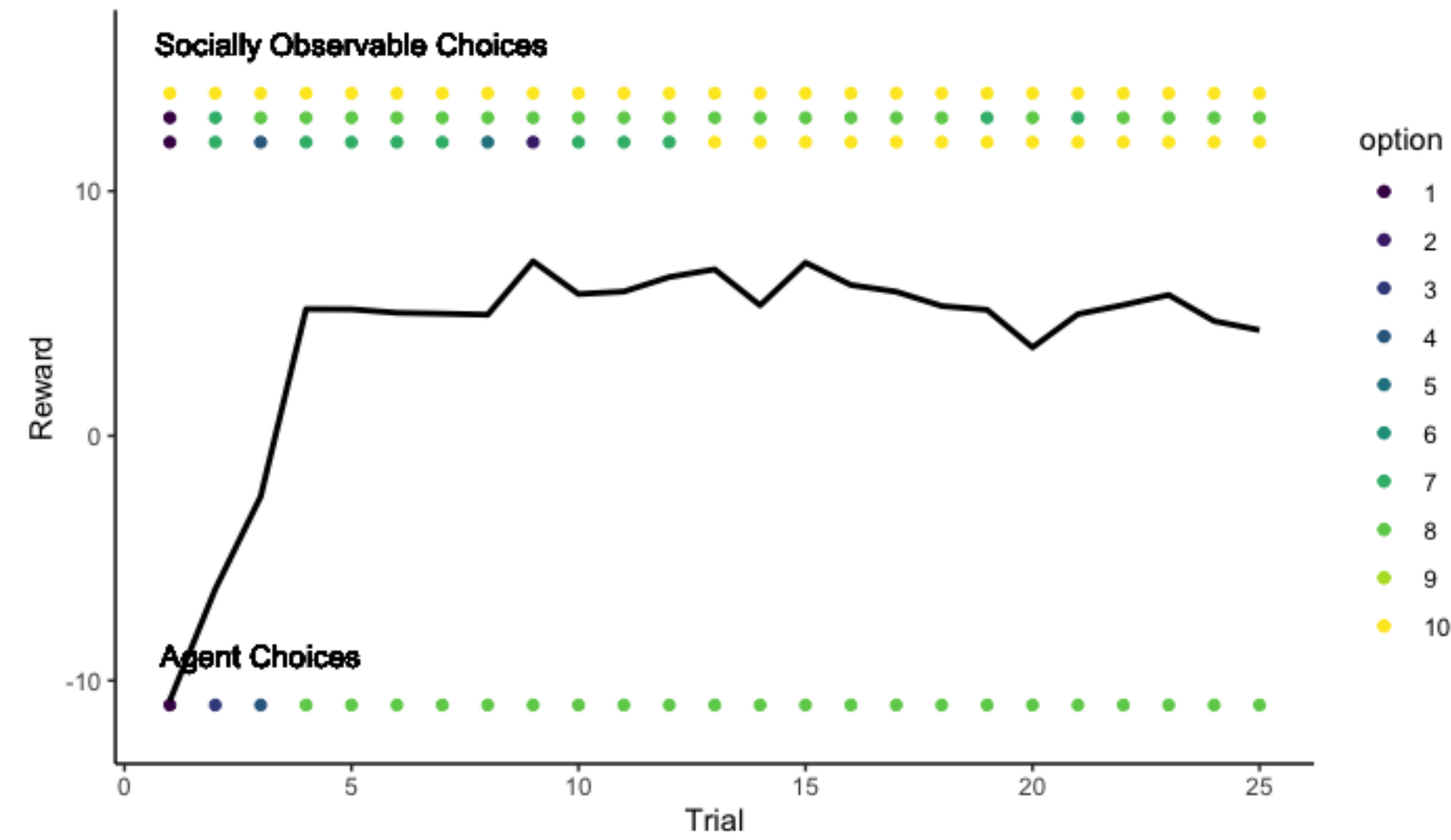
But in practice, BIC is generally a more strict approach to penalizing for complexity compared to AIC and is less likely to prefer an overfit model:

$\log(n) > 2$ when there are at least 8 data points

# AIC vs. BIC

## Simulated data from a Q-learning agent

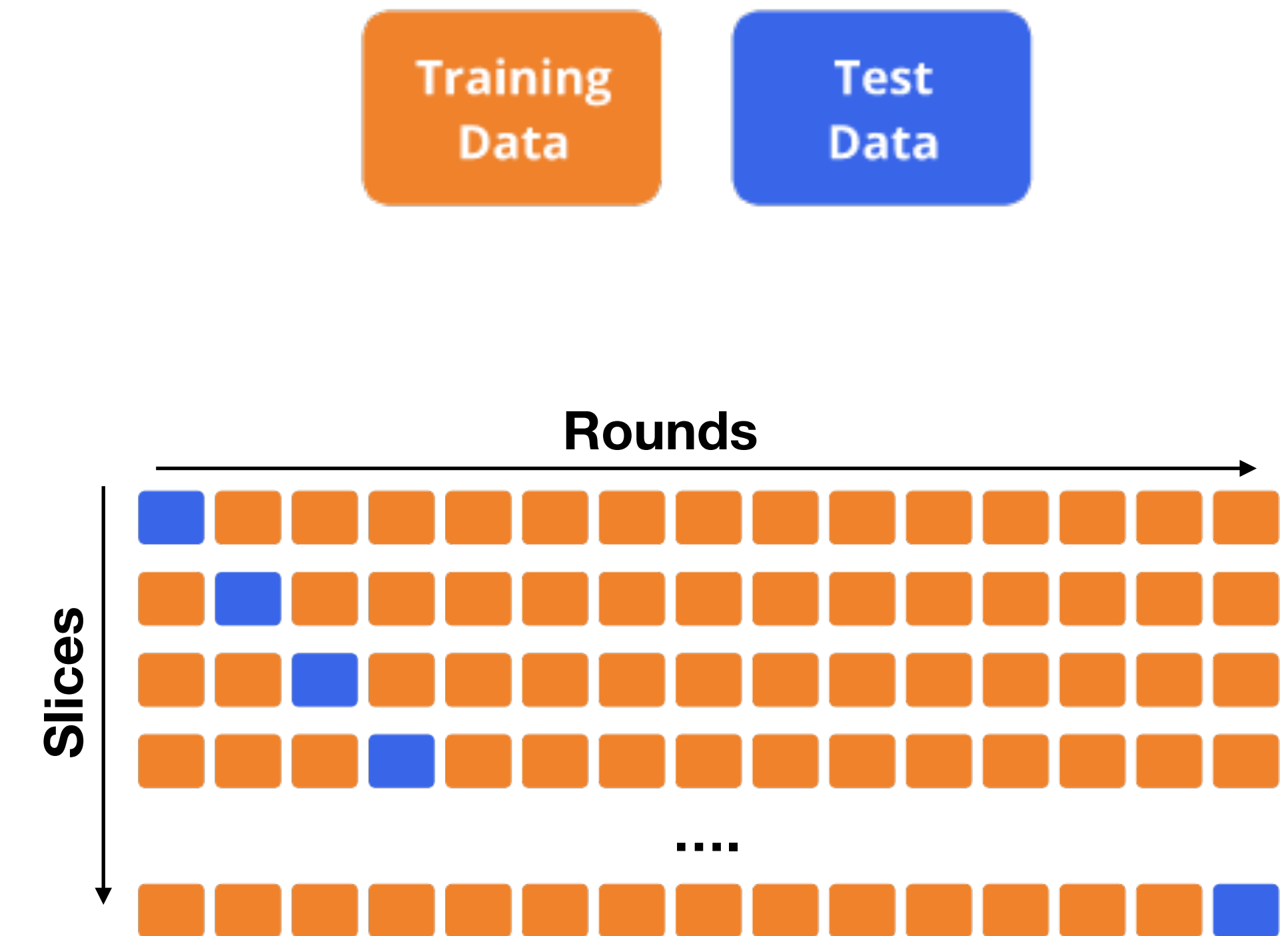## Q-learning vs. Value shaping

# Goodness of Fit Measures

| | **Maximum Likelihood** | **Bayesian Model Selection** |
|---|---|---|
| | $P(D \mid m, \hat{\theta})$ | $\dfrac{P(D \mid m_1)}{P(D \mid m_2)}$ |
| **Penalizing for parameters** | Akaike's Information Criterion (AIC) | Bayesian Information Criterion (BIC) |
| **Prediction error/ Bayesian Occam's Razor** | Cross-validation loss | Model evidence using Markov Chain Monte Carlo (MCMC) |

# Cross Validation

Rather than penalizing for complexity posthoc, we can actively test the predictive accuracy of a model through cross validation

1. Iteratively split the data into training and test sets

2. Estimate MLE on the training set, and then predict out-of-sample on the test set

3. Goodness of fit is the summed negative log likelihood of all out-of sample predictions:
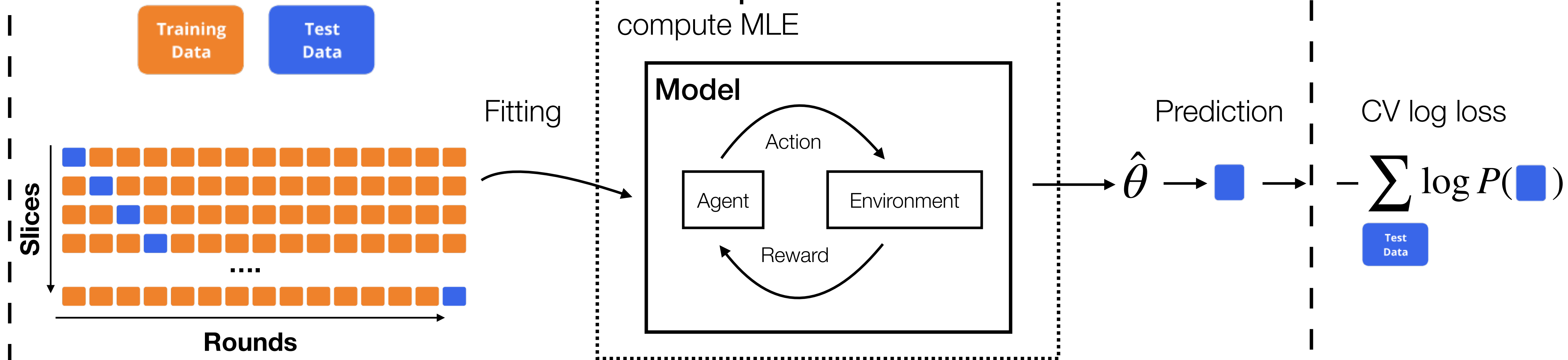
■ + ■ + ■ + … + ■

# Cross validation



**Outer loop:** Iterate over cross validation slices

Training Data

Test Data

Slices

Rounds

Fitting

**Inner loop:** Iterate over data and compute MLE

**Model**

Action

Agent

Environment

Reward

Prediction

$\hat{\theta}$

CV log loss

$-\sum \log P(\ \ )$

Test Data

21

# Variants of Cross validation

- **Leave-one-round-out cross validation**: Use the natural distinction between independent rounds or blocks in an experiment

- **k-Fold cross validation**: when there is no natural structure in the data, we can break it into *k* equally sized slices

- **Leave-one-out-cross validation**: most extreme case, where we iteratively leave a single data point out of the training set
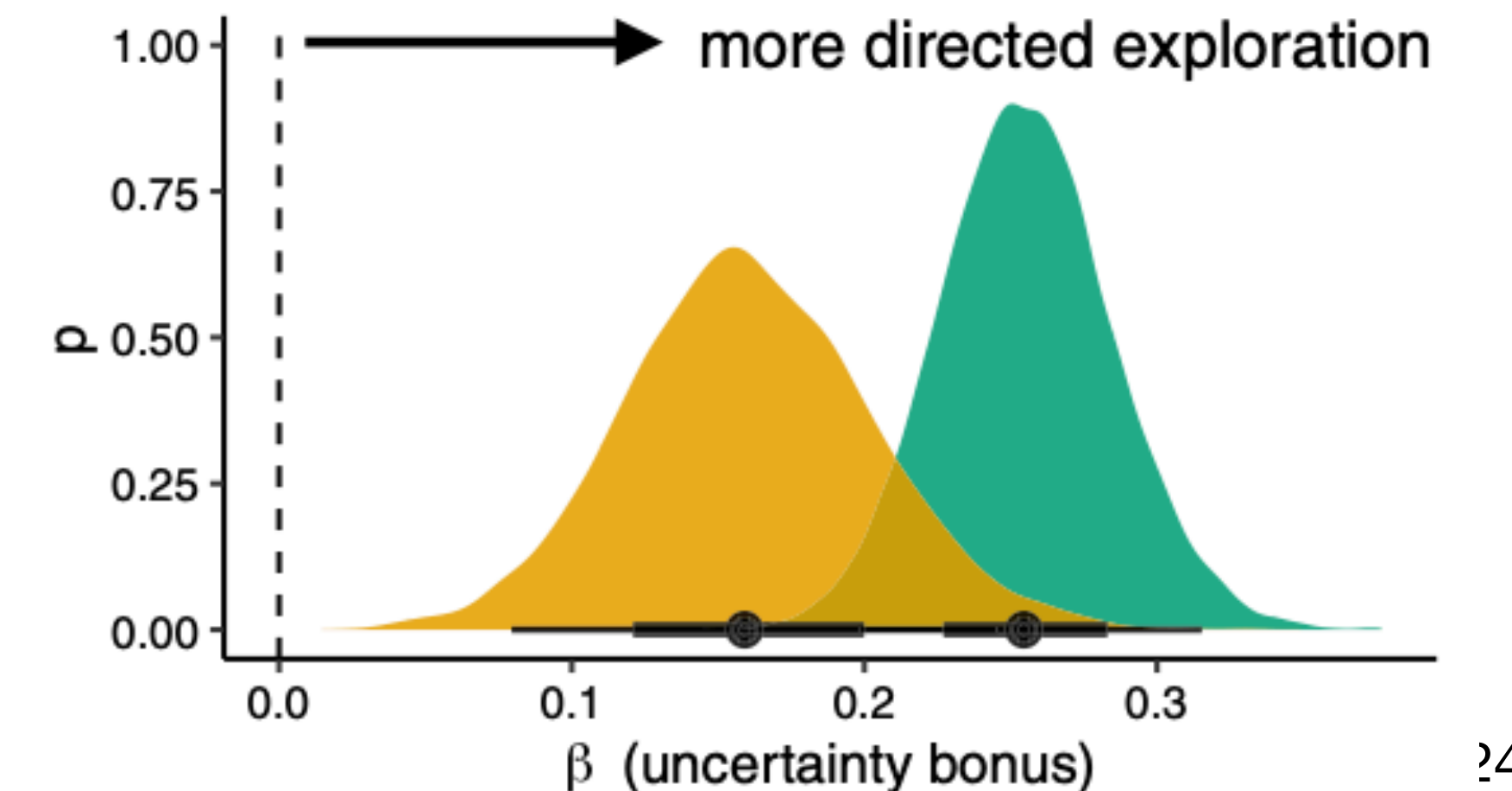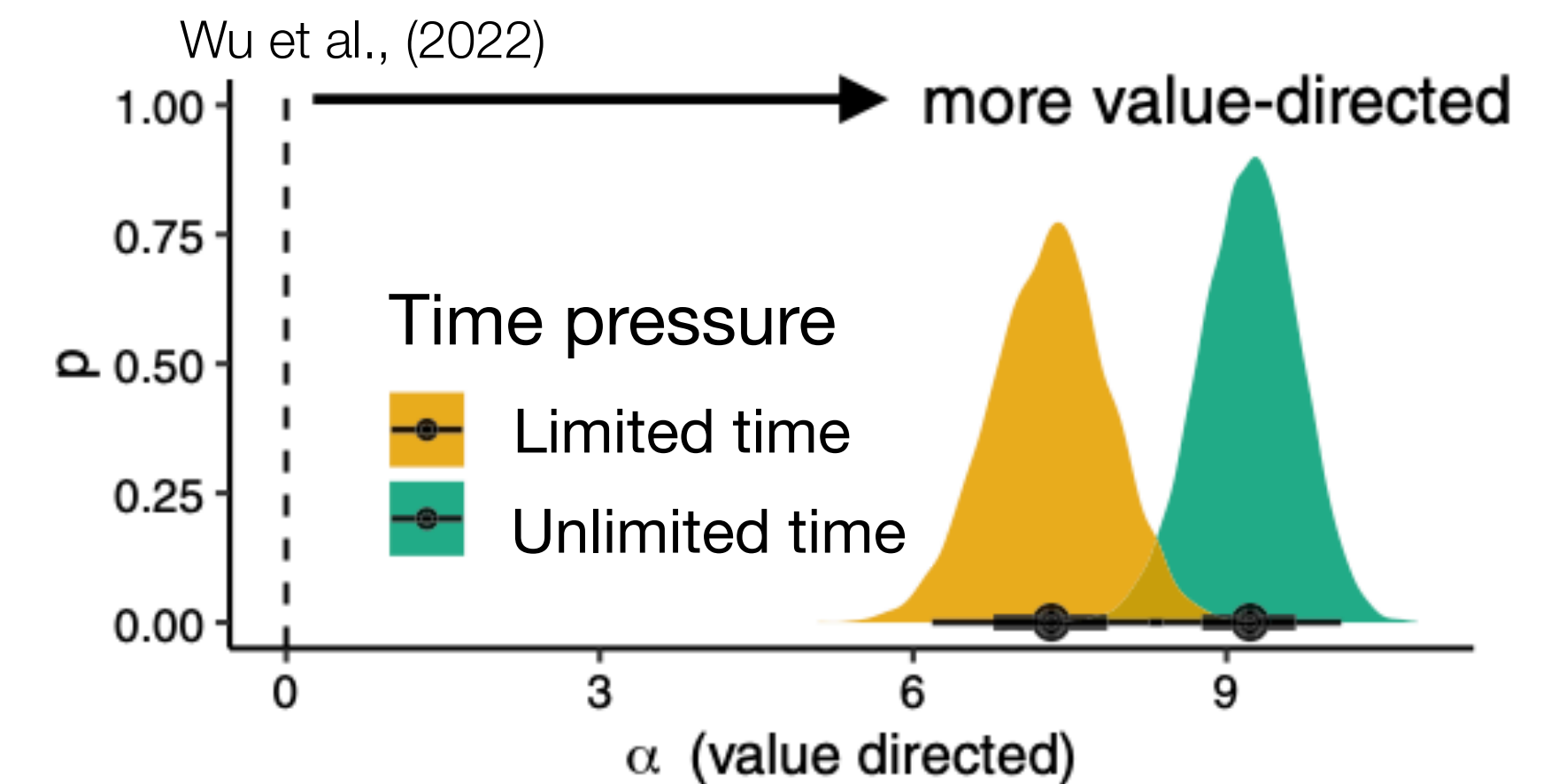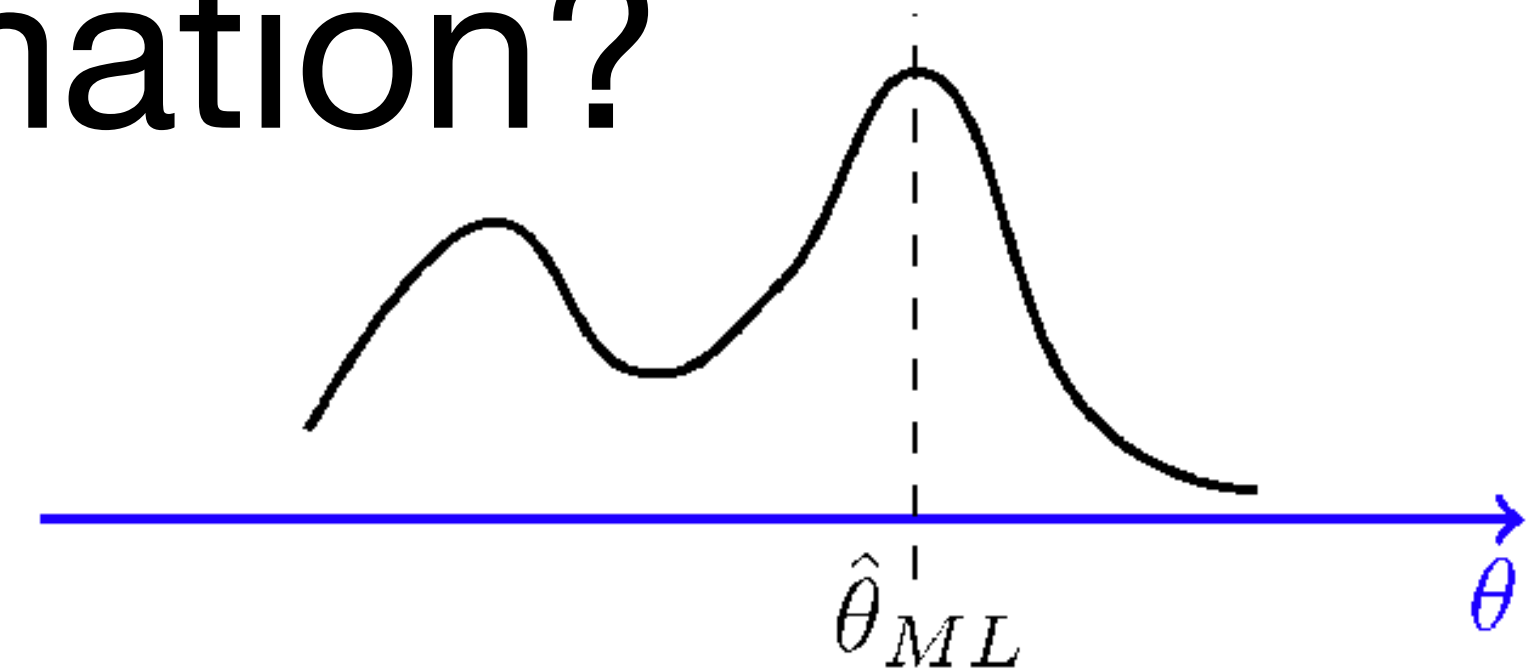
Let's get distributional!

23

# Why Bayesian model estimation?

1. Not just a point estimate, but an entire **probability distribution over parameters**

2. Rather than only assuming participants are independent samples, we can model **hierarchical relationships**

3. Naturally avoid overfitting through **Bayesian Occam's Razor**, since we evaluate the model across the entire range of parameters



$\hat{\theta}_{ML}$

$\theta$

Wu et al., (2022)

more value-directed

Time pressure

Limited time

Unlimited time

$\alpha$ (value directed)

more directed exploration

$\beta$ (uncertainty bonus)

$\alpha_1$  $\alpha_2$  $\alpha_3$  ~

A distribution of $\alpha$

# Posterior distribution over parameters

- Previously, we only used MLE to provide a point estimate of the best parameters $\hat{\theta}$

- Here, we want to estimate the full distribution of parameters suggested by the data and our choice of model:

$$P(\theta \,|\, D, m) \propto P(D \,|\, \theta, m) P(\theta, m)$$

- $P(\theta \,|\, D, m)$ is the **posterior** distribution, which we compute using Bayes' rule combining:

  - The **likelihood** $P(D \,|\, \theta, m)$ of the data given a specific model and set of parameters

  - A **prior** $P(\theta, m)$ over parameters, capturing our initial guess before we see the data

# Markov Chain Monte Carlo

# Markov Chain Monte Carlo

- **Problem**: We want to model a probability distribution that is difficult to compute analytically

# Markov Chain Monte Carlo

- **Problem**: We want to model a probability distribution that is difficult to compute analytically

- **Solution**: acquire random samples that approximate this distribution

# Markov Chain Monte Carlo

- **Problem**: We want to model a probability distribution that is difficult to compute analytically

- **Solution**: acquire random samples that approximate this distribution

- *Markov Chain*

  - sequential process, where each random sample is used as a stepping stone to generate the next sample

  - Special property: Markov Chain has as it's equilibrium distribution the target distribution we are trying to approximate

# Markov Chain Monte Carlo

- **Problem**: We want to model a probability distribution that is difficult to compute analytically

- **Solution**: acquire random samples that approximate this distribution

- *Markov Chain*

  - sequential process, where each random sample is used as a stepping stone to generate the next sample

  - Special property: Markov Chain has as it's equilibrium distribution the target distribution we are trying to approximate

- *Monte Carlo*

  - Law of large numbers —> enough randomly drawn samples will approximate the underlying distribution

# Metropolis-Hastings MCMC



**Psuedocode**

1. Sample $\theta^i$ from $P(\theta^i | \theta^{i-1})$

2. Compute likelihood of data given these parameters $P(D|\theta^i)$

3. Accept the sample with probability proportional to how much of an improvement $P(D|\theta^i)$ is over $P(D|\theta^{i-1})$

The final collection of samples approximates the posterior parameter estimate $P(\theta|D)$

Initial Sample ($\theta^0$)

Prior distribution $p(\theta)$

Posterior distribution $P(\theta|D)$

Lee, Sung, & Choi (2015)

Step 1:  $r(\theta_{new}, \theta_{t\text{-}1}) = \dfrac{\text{Posterior}(\theta_{new})}{\text{Posterior}(\theta_{t\text{-}1})} = \dfrac{\text{Beta}(1,1,0.306) \times \text{Binomial}(10,4,0.306)}{\text{Beta}(1,1,0.429) \times \text{Binomial}(10,4,0.429)} = 0.834$

Step 2:  Acceptance probability  $\alpha(\theta_{new}, \theta_{t\text{-}1}) = \min\{r(\theta_{new}, \theta_{t\text{-}1}), 1\} = \min\{0.834, 1\} = 0.834$

Step 3:  Draw $u \sim \text{Uniform}(0,1) = 0.617$

Step 4:  If  $u < \alpha(\theta_{new}, \theta_{t\text{-}1})$  $\rightarrow$  If  $0.617 < 0.834$    Then    $\theta_t = \theta_{new} = 0.306$
Otherwise  $\theta_t = \theta_{t\text{-}1} = 0.429$

Density — MCMC Iteration

$N(\theta_{t-1}, \sigma)$

Step 1:  $r(\theta_{new}, \theta_{t-1}) = \dfrac{Posterior(\theta_{new})}{Posterior(\theta_{t-1})} = \dfrac{Beta(1,1,0.306) \times Binomial(10,4,0.306)}{Beta(1,1,0.429) \times Binomial(10,4,0.429)} = 0.834$

Step 2:  Acceptance probability  $\alpha(\theta_{new}, \theta_{t-1}) = \min\{r(\theta_{new}, \theta_{t-1}), 1\} = \min\{0.834, 1\} = 0.834$

Step 3:  Draw $u \sim Uniform(0,1) = 0.617$

Step 4:  If  $u < \alpha(\theta_{new}, \theta_{t-1})$  →  If  $0.617 < 0.834$      Then      $\theta_t = \theta_{new} = 0.306$
Otherwise  $\theta_t = \theta_{t-1} = 0.429$

28

# MCMC Samplers

```
data {
  int<lower=0> N;                    // N >= 0
  int<lower=0,upper=1> y[N];         // y[n] in { 0, 1 }
}
parameters {
  real<lower=0,upper=1> theta;       // theta in [0, 1]
}
model {
  theta ~ beta(1,1);                 // prior
  y ~ bernoulli(theta);              // likelihood
}
```

```python
with pm.Model() as hierarchical_model_centered:
    # hyperpriors for group nodes
    mu_a = pm.Normal('mu_a', mu=0., sd=100**2)
    sigma_a = pm.HalfCauchy('sigma_a', 5)
    mu_b = pm.Normal('mu_b', mu=0., sd=100**2)
    sigma_b = pm.HalfCauchy('sigma_b', 5)

    # intercept about each county
    a = pm.Normal('a', mu=mu_a, sd=sigma_a, shape=n_counties)
    b = pm.Normal('b', mu=mu_b, sd=sigma_b, shape=n_counties)

    # error
    eps = pm.HalfCauchy('eps', 5)

    # regression
    radon_est = a[county_idx] + b[county_idx] * Dat.floor.values

    # likelihood
    radon_like = pm.Normal('radon_like', mu=radon_est, sd=eps, observed=Dat.log_radon)
```

# Posterior over parameters

# Bayesian model comparison

## Information Criteria

AIC – Akaike information criterion

DIC – Deviance Information Criterion

WAIC – Widely Applicable Information Criterion

(Watanabe–Akaike information criterion)

finding the model that has the highest out-of-sample predictive accuracy

approximation to LOO

WAIC: using entire posterior distribution

AIC/DIC: using point estimation

**Tutorial 4**

52

# Part 1 Summary

|  | **Maximum Likelihood** $P(D\,|\,m,\hat{\theta})$ | **Bayesian Model Selection** $\dfrac{P(D\,|\,m_1)}{P(D\,|\,m_2)}$ |
|---|---|---|
| **Penalizing for parameters** | Akaike's Information Criterion (AIC) | Bayesian Information Criterion (BIC) |
| **Prediction error/ Bayesian Occam's Razor** | Cross-validation loss | Model evidence using Markov Chain Monte Carlo (MCMC) |

# Notebook

https://cosmos-konstanz.github.io/notebooks/tutorial-3-model-comparisons.html#model-fitting-exercise

# Model fitting exercise

**meteor.csv**　　　　　**comet.csv**　　　　**Self-contained model-fitting code**

# Which model best explains each dataset?

# Part 2. Robustness

**(5 minute break)**

# Robustness checks

1. **Model recovery**

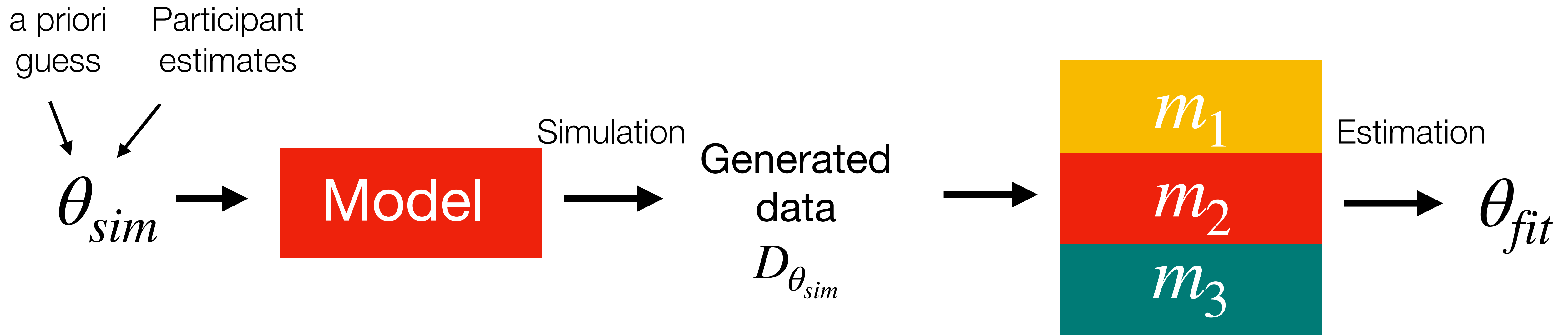   - Can the data actually differentiate between the models we are considering? Could there be model mimicry, where the wrong model can mistakenly win?

2. **Parameter recovery**

   - Are the parameters of the model capturing distinct phenomenon? Can changes in one parameter be acommodated by changes in another parameter (i.e., misspecification)?

3. **Simulated data**

   - Can the model generate realistic participant behavior? Is it capturing the mechanisms that matter for performance, rather than simply fitting the noise?

# Model recovery

a priori
guess

Participant
estimates

$$\theta_{sim} \longrightarrow \boxed{\text{Model}} \xrightarrow{\text{Simulation}} \begin{array}{c} \text{Generated} \\ \text{data} \\ D_{\theta_{sim}} \end{array} \longrightarrow \boxed{\begin{array}{c} m_1 \\ m_2 \\ m_3 \end{array}} \xrightarrow{\text{Estimation}} \theta_{fit}$$

1.  Use models to simulate data, parameterized with $\theta_{sim}$ either an *a priori* guess or from participant estimates

2.  Use the same model estimation procedure on the simulated data to estimate $\theta_{fit}$ for each model under consideration

3.  How often does the correct model provide the best fit?

# Model recovery

**Confusion matrix** p(fit|sim)



**Inversion matrix** p(sim|fit)



Which alternative models mimic a given simulation model?

If a given model wins a model competition, how likely is it to actually be the true generative model?

# Parameter Recovery

**Goal**: Determine if parameters are distinct and behaviorally specific

1. Use either participant parameter estimates or some prior guess to simulate data (x-axis)
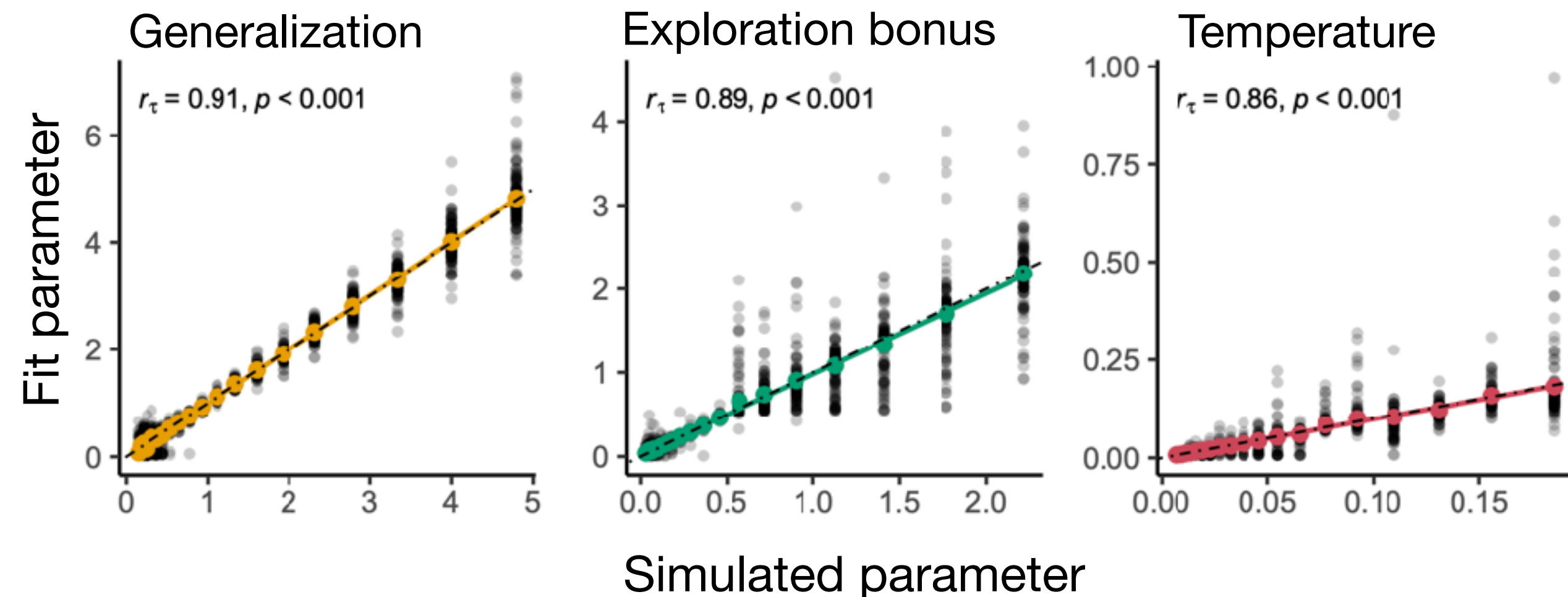
2. Run model fitting to estimate new parameters on simulated data (y-axis)

3. Do the fit parameters correspond to the simulated parameters?

[Bonus] Counterfactual parameter recovery: Systematically vary simulating parameters across a range of plausible values. Does the entire hypothesis space recover?



Wilson & Collins (*eLife* 2019)


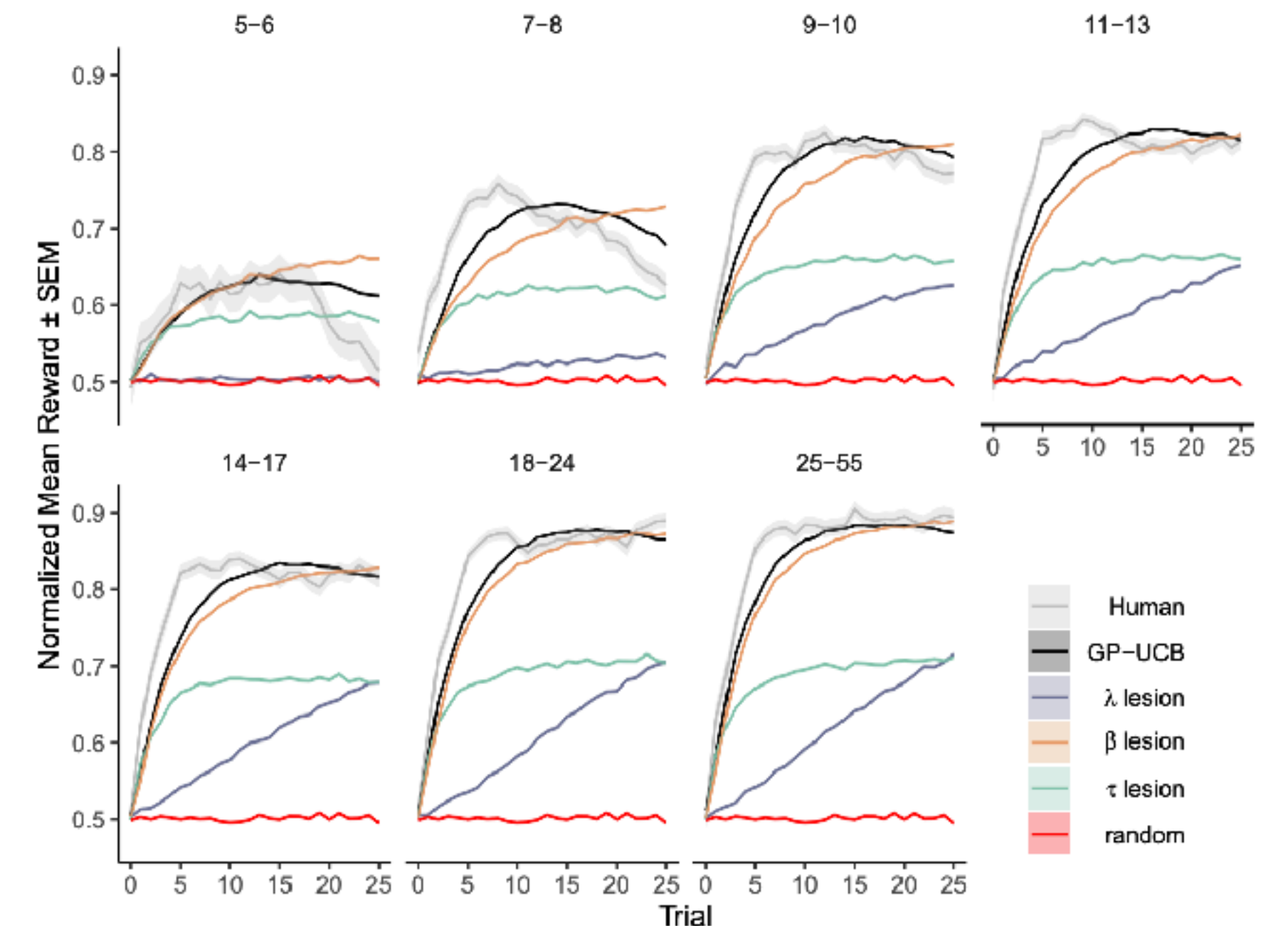
Giron*, Ciranka* et al., (2023)

# Simulated Performance

- Goodness of fits don't always tell the full story

- Sometimes you need to check that models can reproduce important patterns of human behavior

  - Can also be used to probe hidden components of the model, such as value representations

- Compare simulated model performance to human performance

# Simulated Performance

- Goodness of fits don't always tell the full story

- Sometimes you need to check that models can reproduce important patterns of human behavior

  - Can also be used to probe hidden components of the model, such as value representations

- Compare simulated model performance to human performance

  - Can the model replicate differences across experimental manipulations or from different populations
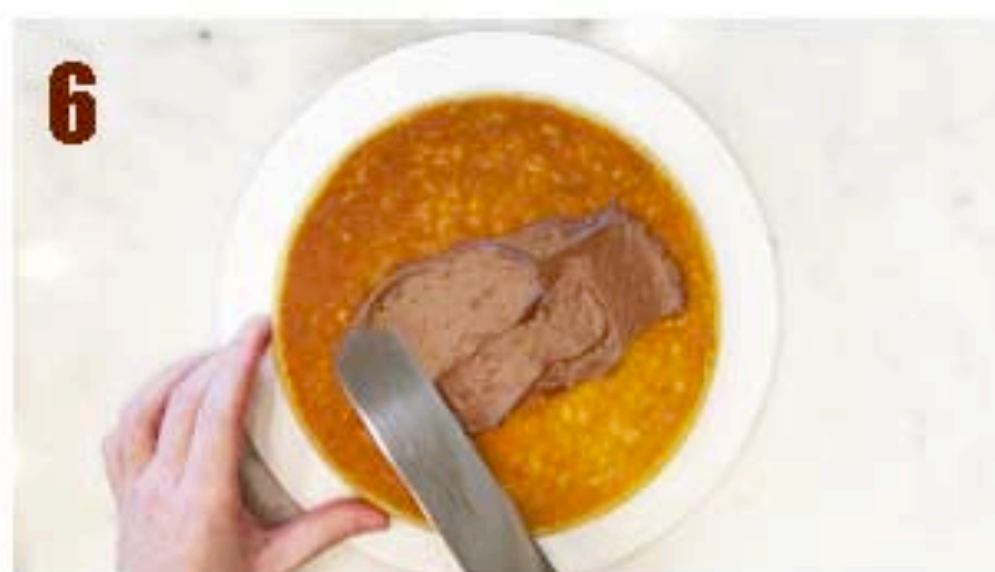




Giron*, Ciranka*, et al., (2022)

# General Recipe for Cognitive Modeling

# Not fixed, step by step instructions…

**Not fixed, step by step instructions…**



**… but an adaptive set of principles**

# General Recipe

1. What are your hypotheses? Turn them into models

2. How will you estimate the model parameters and perform model comparison?

3. Is your modeling framework robust? If not, rethink your task, the models, and/or your modeling framework.

4. [Collect data]

5. Analyze and interpret results

6. Test if recoverability still works with participant parameters

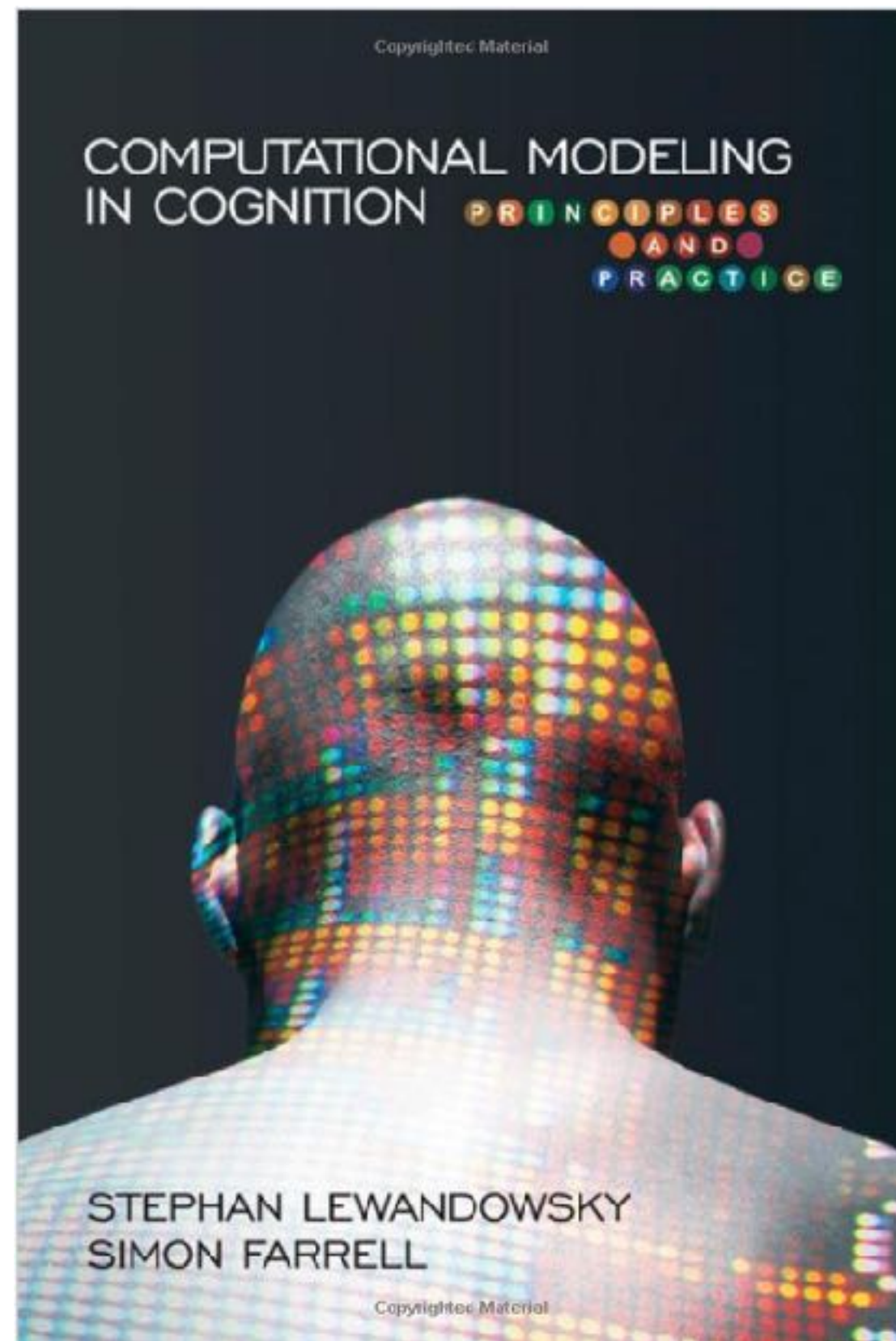# What can you justify?



Reviewers
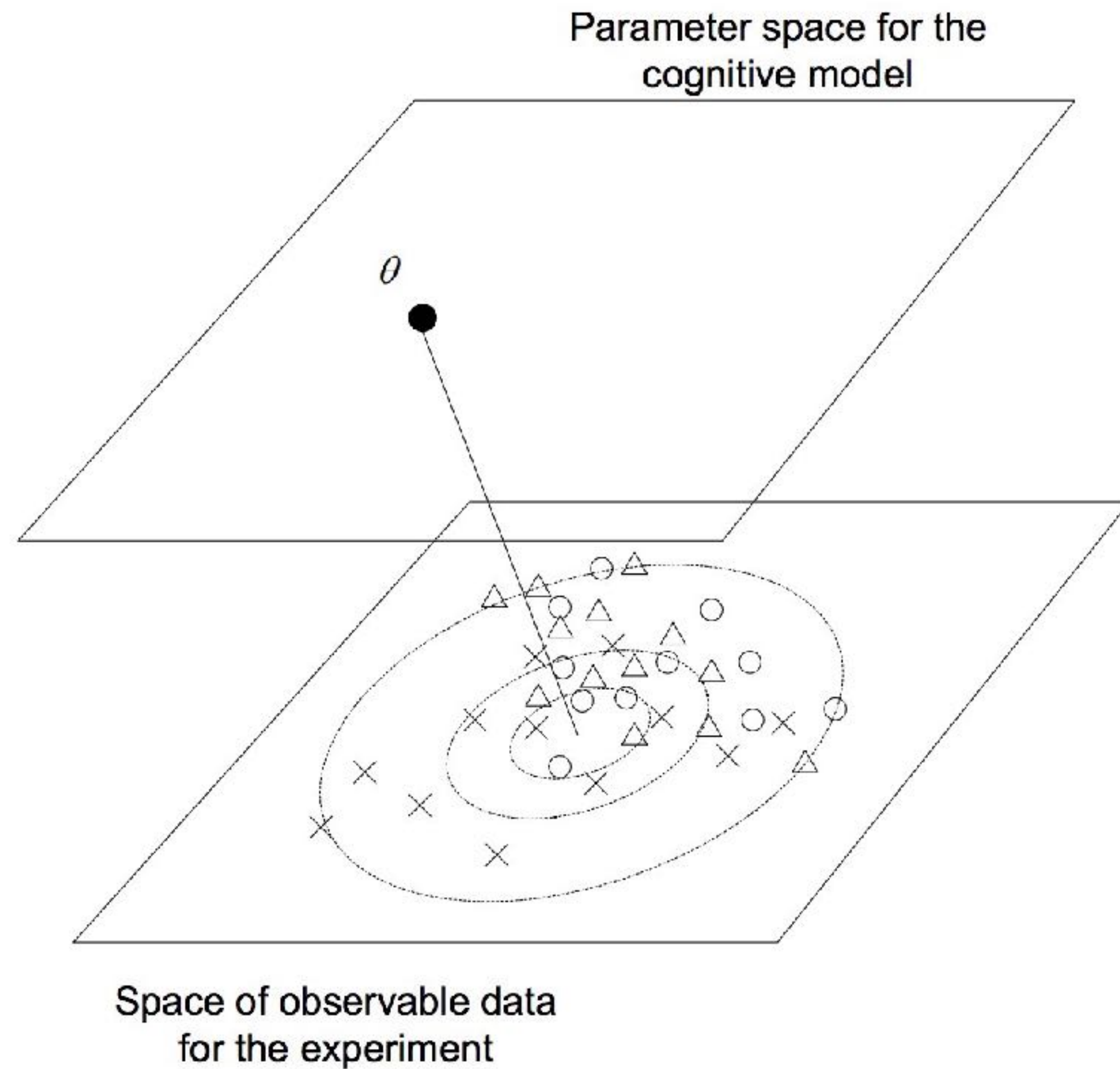
You

# Social Learning Specific Challenges

- Social learning strategies have **frequency-dependent fitness**.

  - Performance of both real and simulated agents, don't only depend on their model parameters, but also on the make-up of the group it is interacting with

  - Objective performance can only be demonstrated with evolutionary simulations

- We only covered **conformity biased** social learning strategies that treat all other individuals as the same

  - Much of social learning is *selective* in learning from successful or prestigious individuals

  - More we need models to account for selectivity biases, but without ballooning in complexity

- We only very briefly touched on **Theory of Mind**, where individuals infer the hidden mental states of others

  - Modeling ToM is very difficult, even using sample-based approximations

  - Even more so, due to infinite recursion of an agent reasoning about what other individuals think about themselves, *ad infinitum*

- Capturing sophistication of social learning may come with the trade-off of needing to simplify individual learning mechanisms
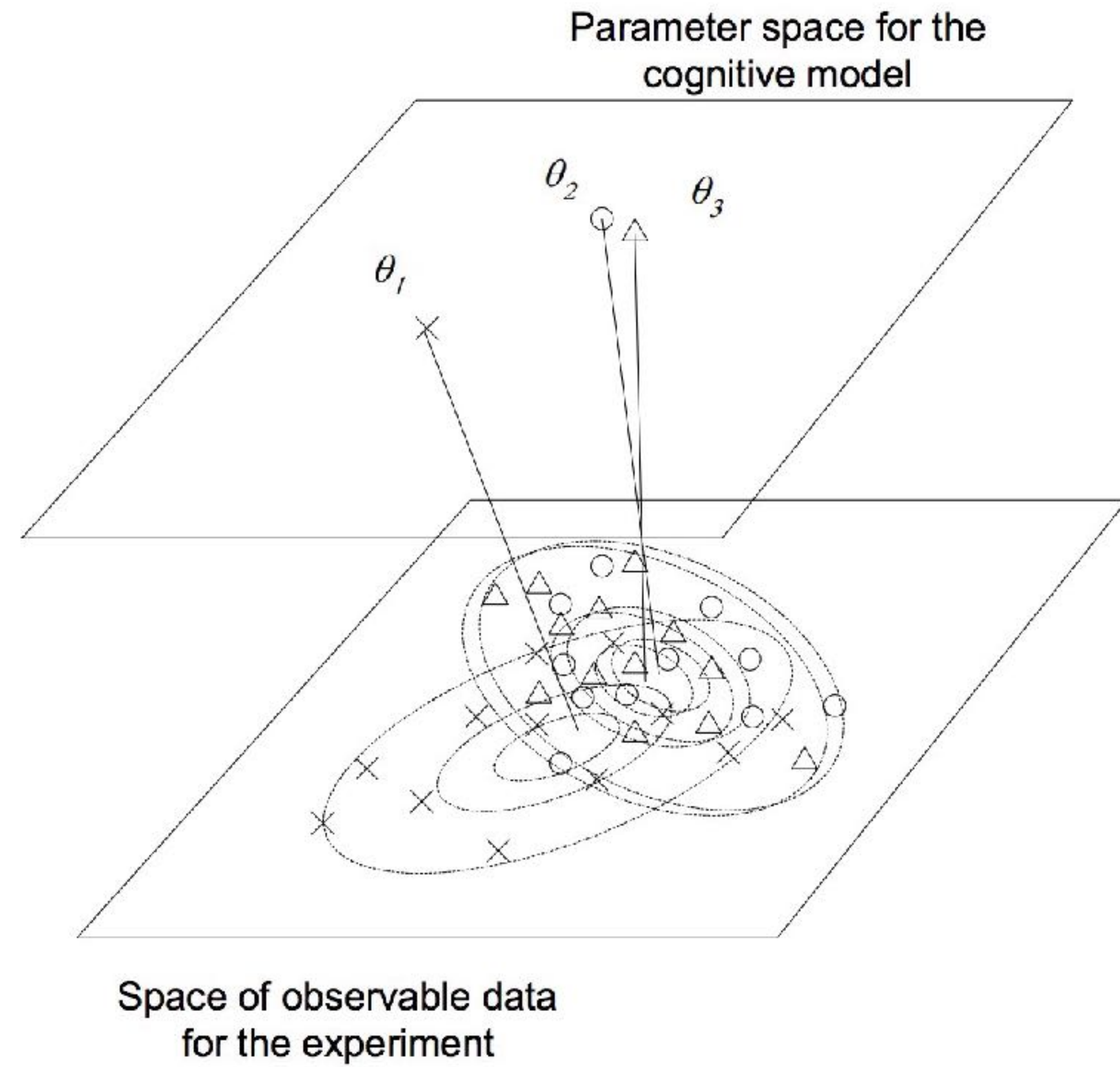
# Recommended Readings
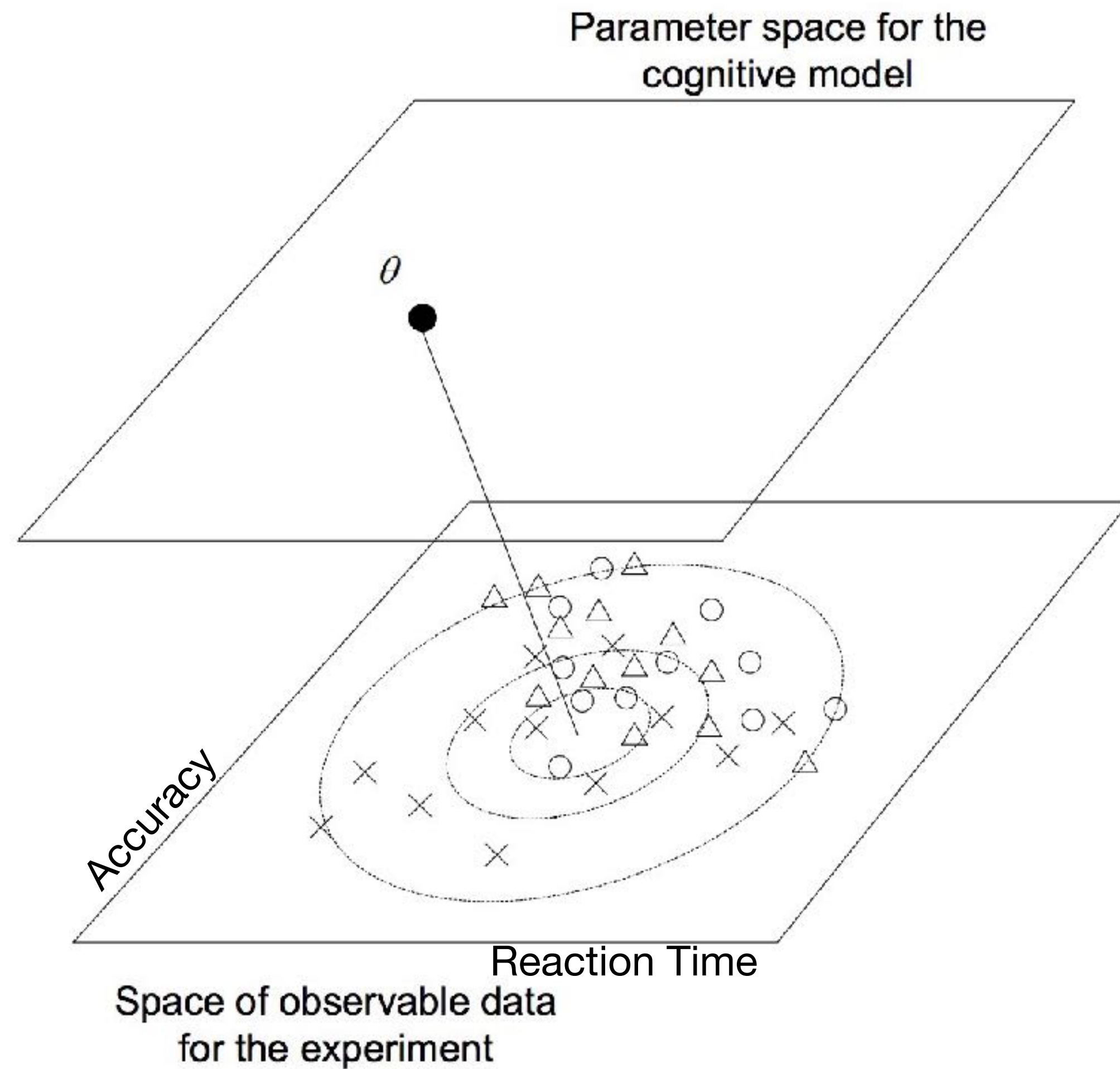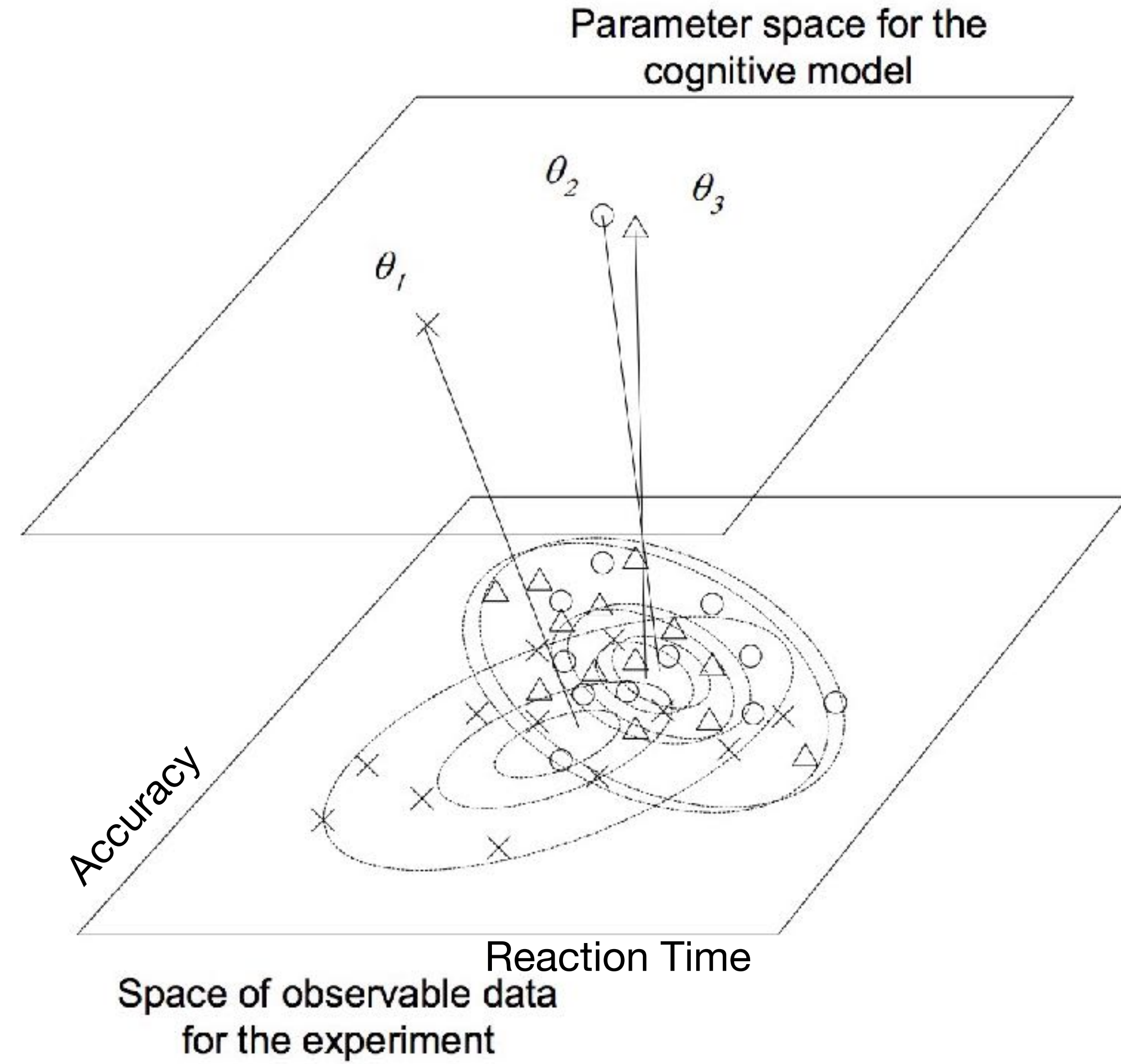
# Supplemental Slides

# Aggregate vs. Individual



Parameter space for the cognitive model

$\theta$

Space of observable data for the experiment

(a)

Parameter space for the cognitive model

$\theta_2$  $\theta_3$

$\theta_1$

Space of observable data for the experiment

(b)

Navarro, Griffiths, Steyvers, & Lee (MathPsych, 2006)   49

# Aggregate vs. Individual



Parameter space for the cognitive model

Accuracy

Reaction Time

Space of observable data for the experiment

(a)

Parameter space for the cognitive model

$\theta_2$ $\theta_3$

$\theta_1$

Accuracy

Reaction Time

Space of observable data for the experiment

(b)

# Discrete vs. Continuous data

# Discrete vs. Continuous data

Choices are discrete outcomes

# Discrete vs. Continuous data

Choices are discrete outcomes

Which flavour of ice-cream?
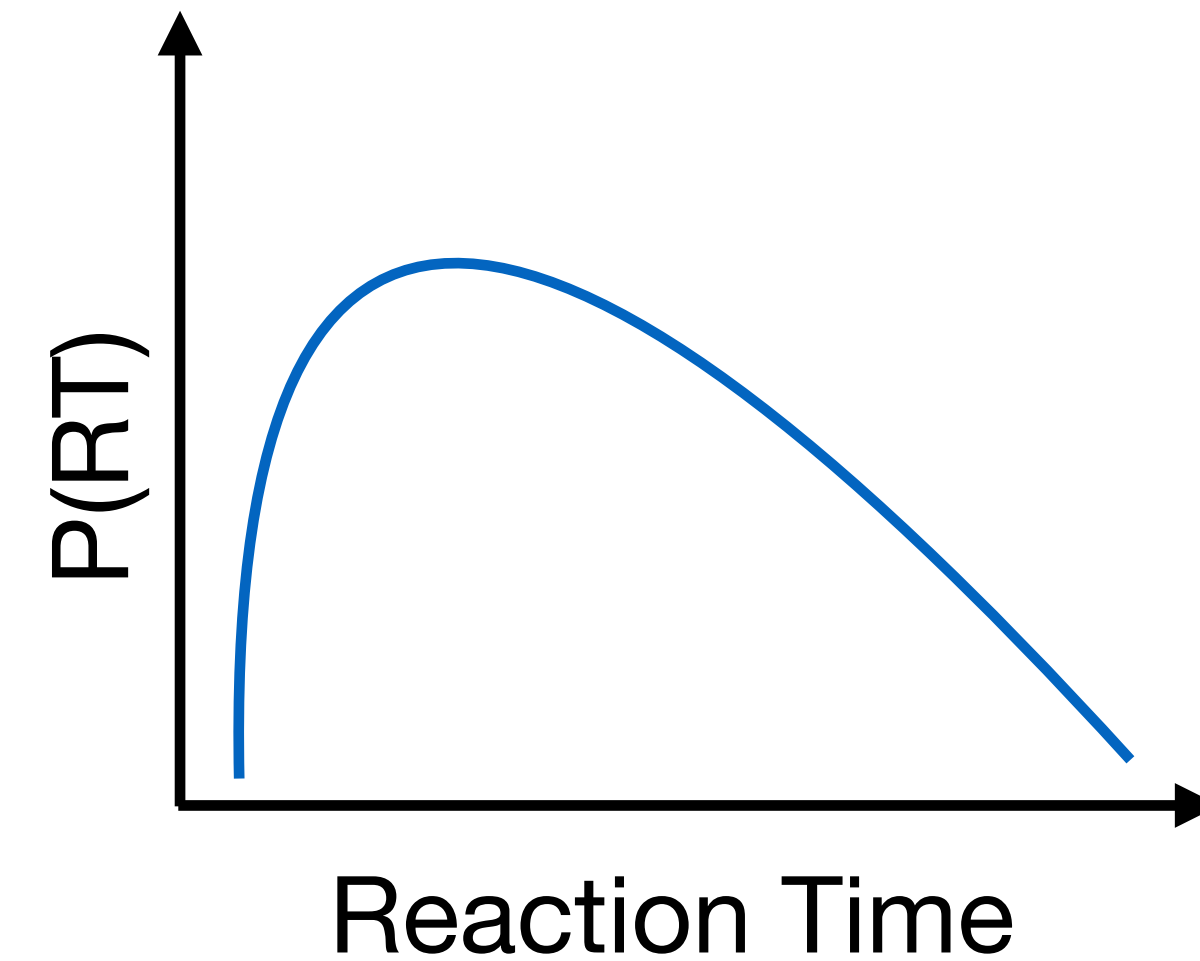
# Discrete vs. Continuous data

Choices are discrete outcomes

Judgments and reaction times are continuous measures

Which flavour of ice-cream?

# Discrete vs. Continuous data

Choices are discrete outcomes
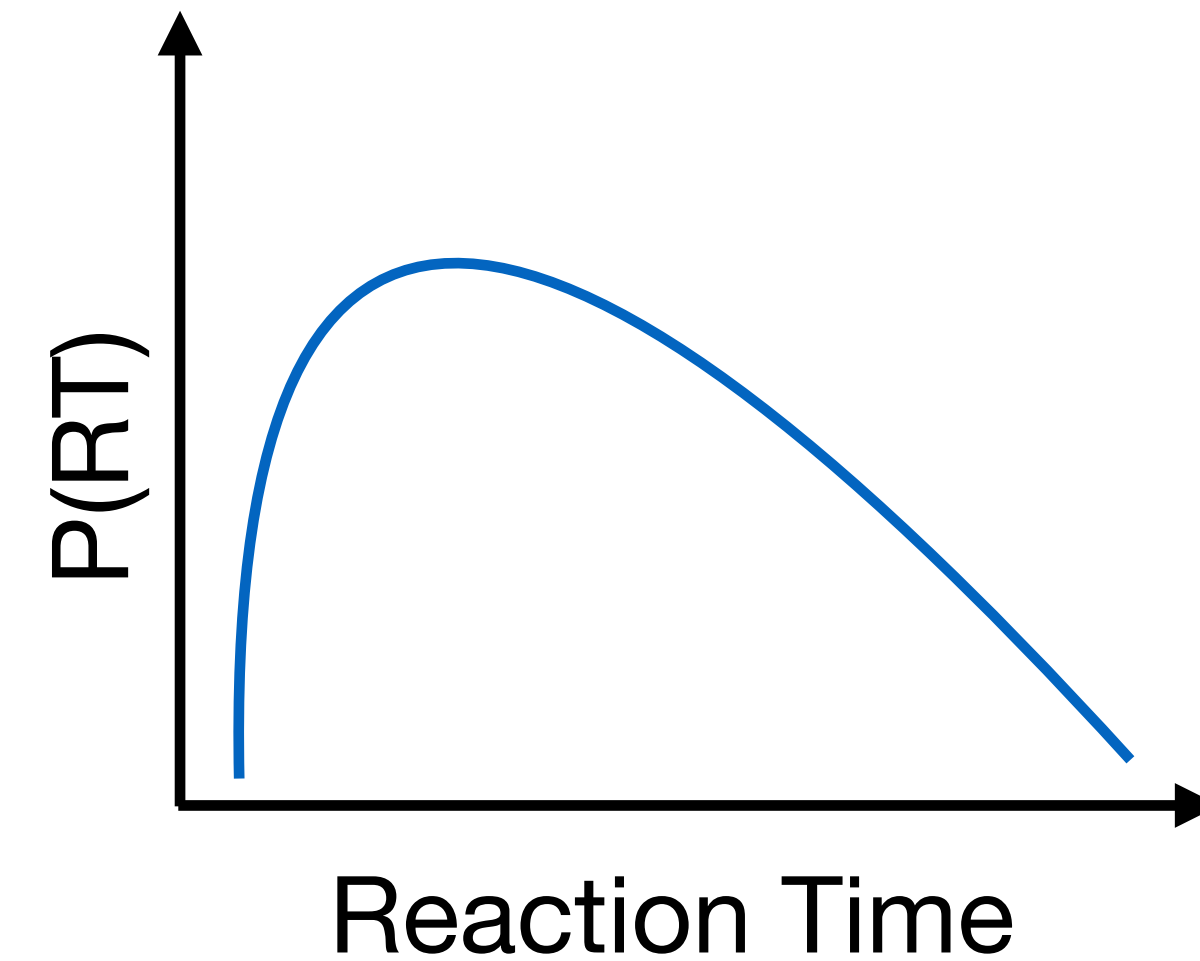
Judgments and reaction times are continuous measures

Which flavour of ice-cream?



How much do you like ice-cream?



Not at all

Extreme

# Discrete vs. Continuous data

Choices are discrete outcomes

Judgments and reaction times are continuous measures

Which flavour of ice-cream?

How much do you like ice-cream?

Not at all                    Extreme

P(RT)

Reaction Time

# Discrete vs. Continuous data

Choices are discrete outcomes

Which flavour of ice-cream?



Model predictions



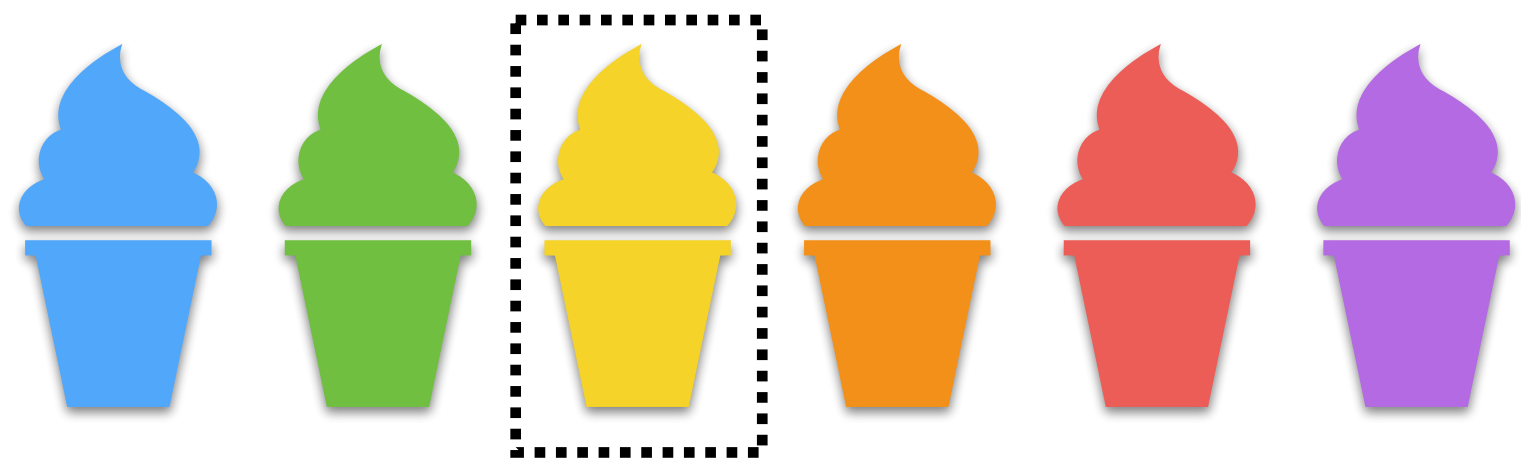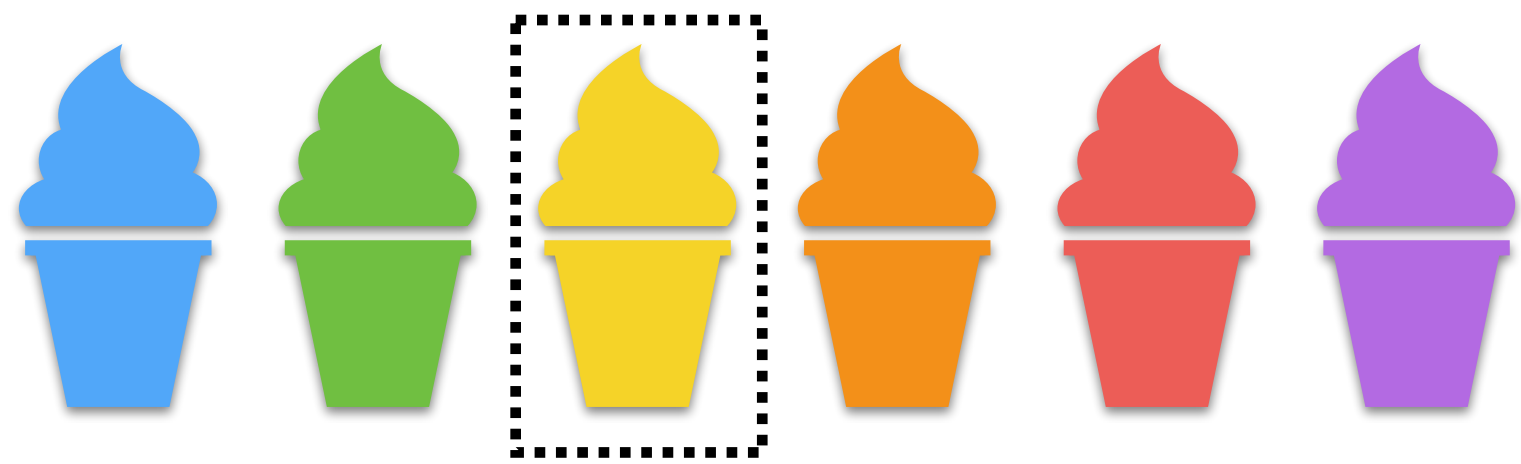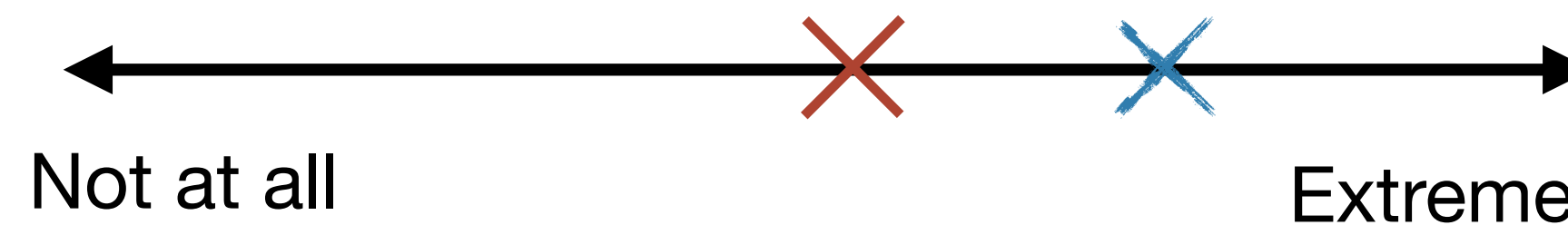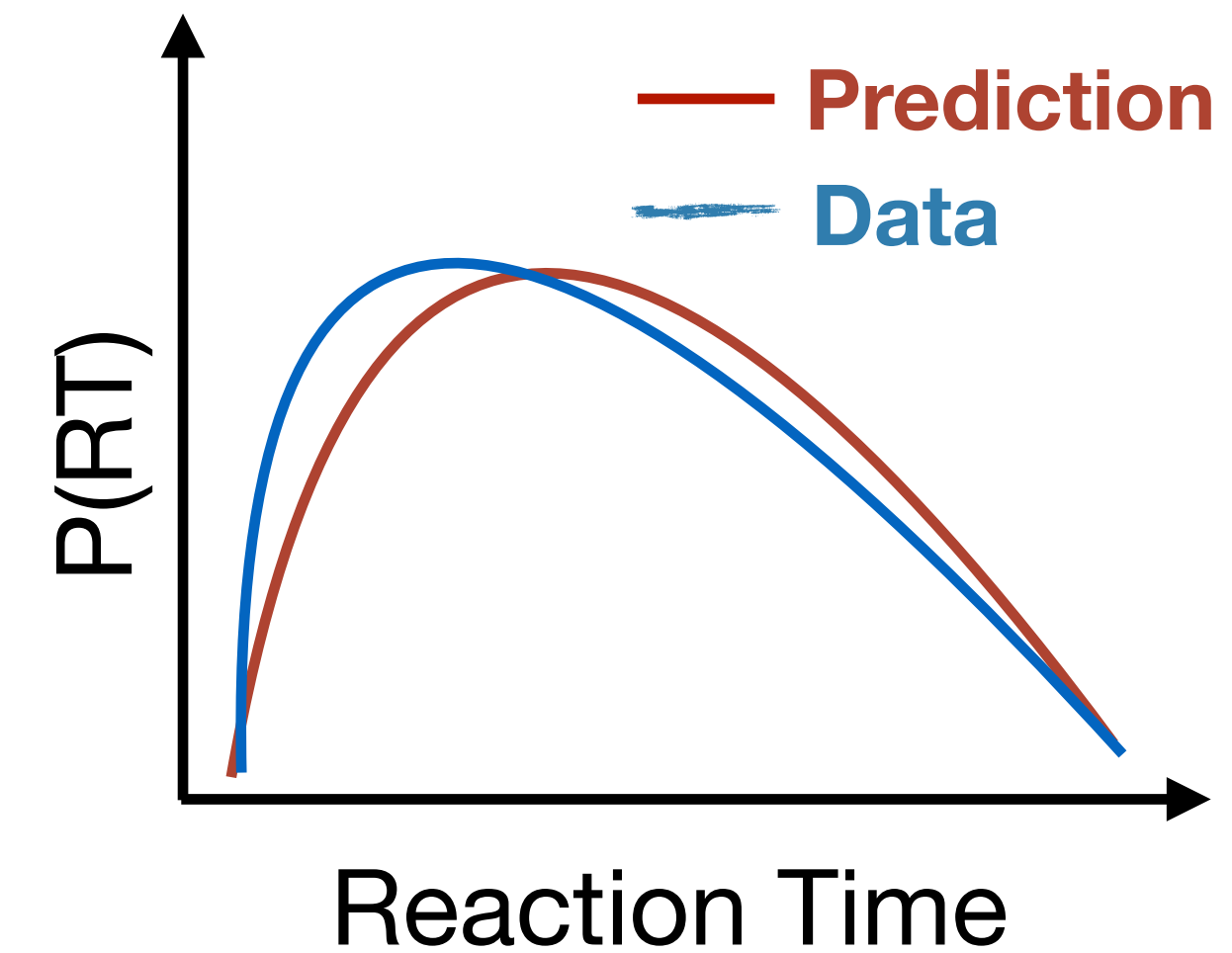Judgments and reaction times are continuous measures
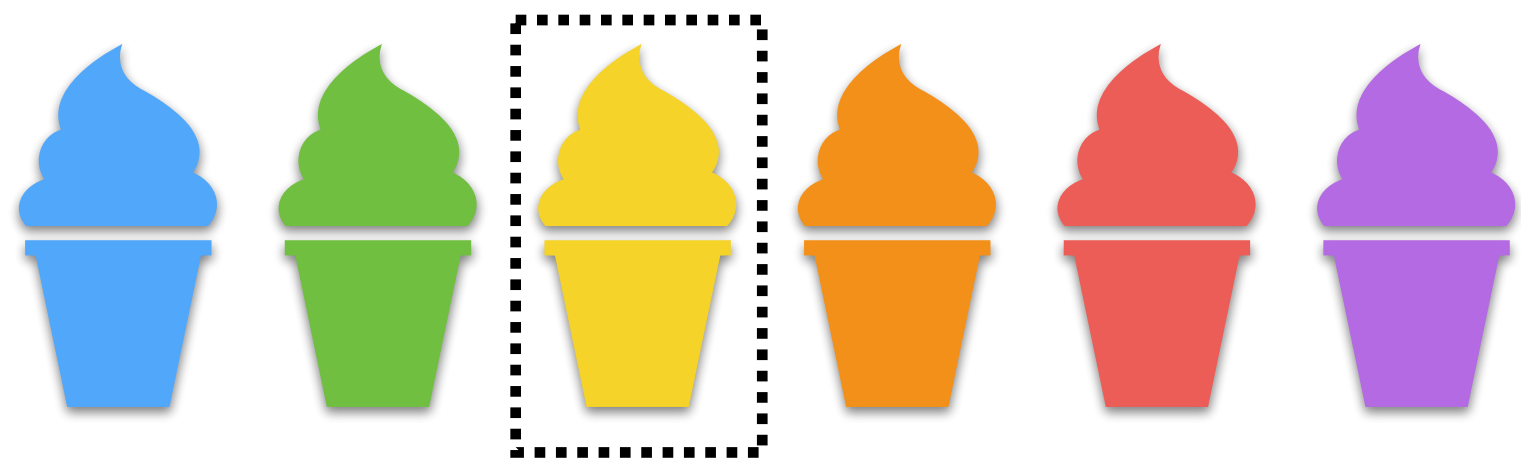
How much do you like ice-cream?

Not at all                    Extreme



50

# Discrete vs. Continuous data

Choices are discrete outcomes

Judgments and reaction times are continuous measures

Which flavour of ice-cream?



How much do you like ice-cream?

Not at all                    Extreme



P(RT)

Reaction Time

Model predictions



P(x)

Flavour

# Discrete vs. Continuous data

Choices are discrete outcomes

Which flavour of ice-cream?



Judgments and reaction times are continuous measures

How much do you like ice-cream?

Not at all    Extreme

⨯ **Prediction**    ⨯ **Data**



Model predictions



- Log P(yellow)

# Discrete vs. Continuous data

## Choices are discrete outcomes

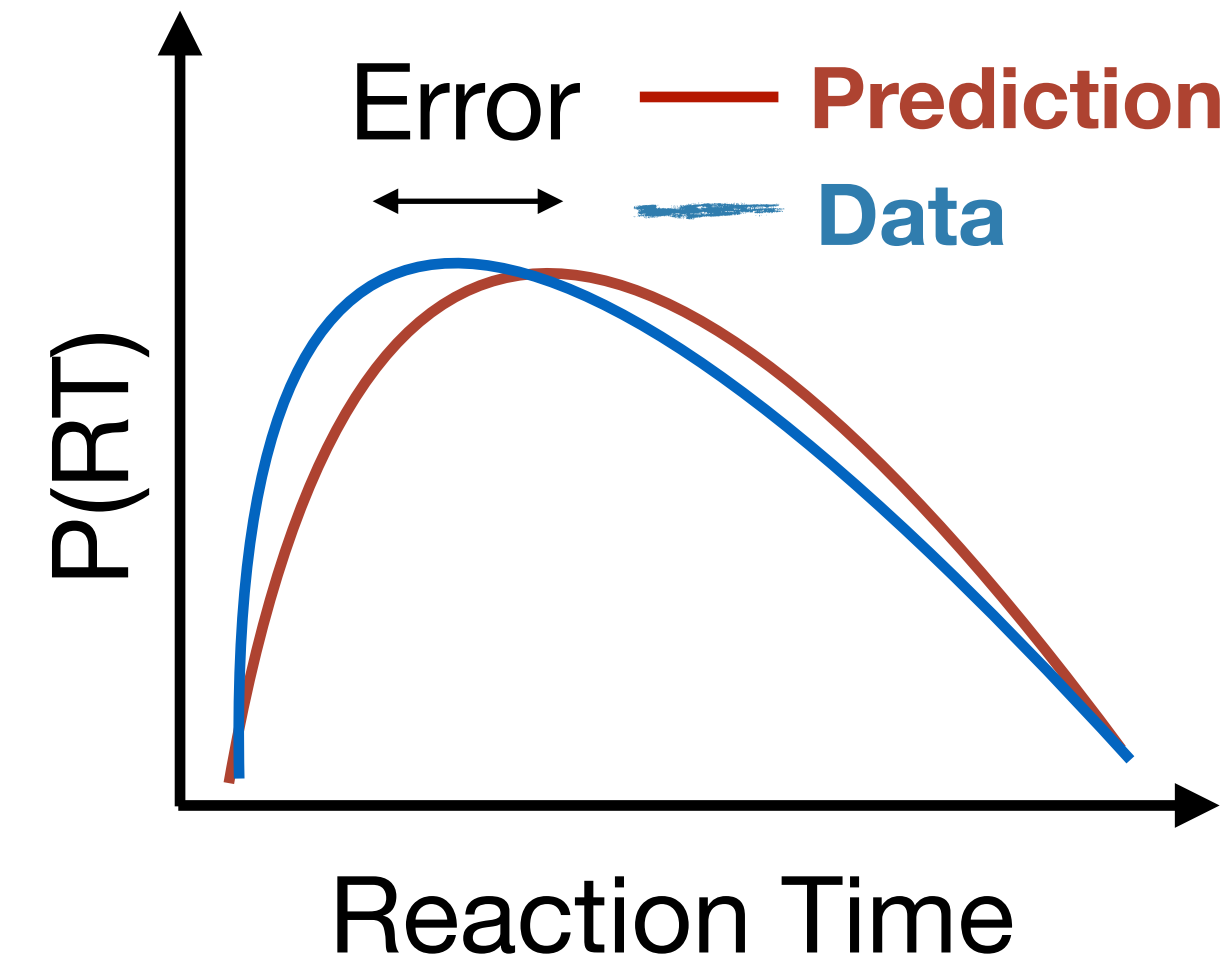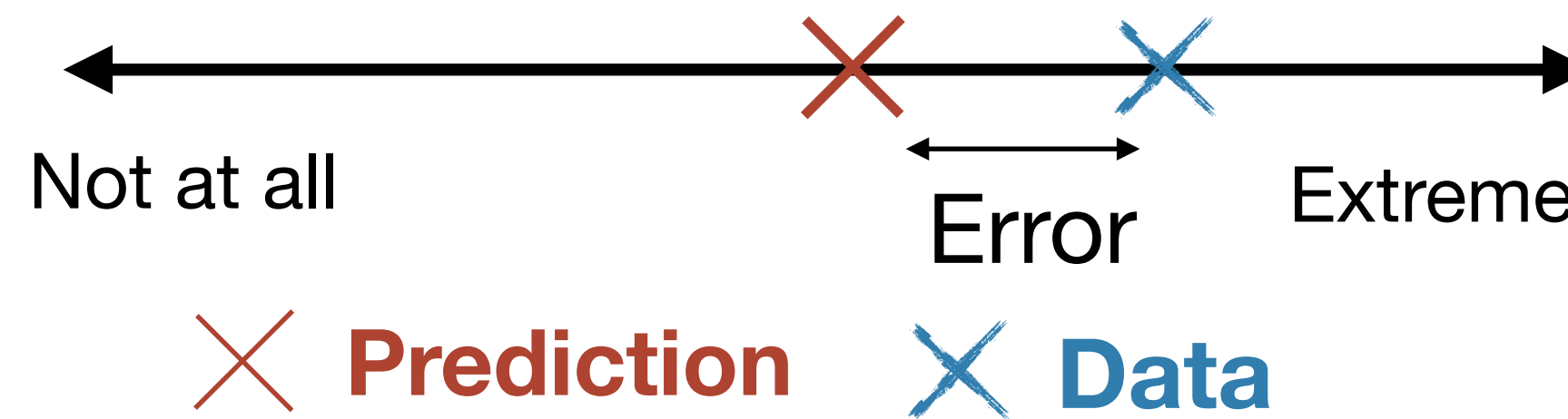### Which flavour of ice-cream?



### Model predictions



## Judgments and reaction times are continuous measures

### How much do you like ice-cream?

# Discrete vs. Continuous data

Choices are discrete outcomes

Which flavour of ice-cream?



Model predictions



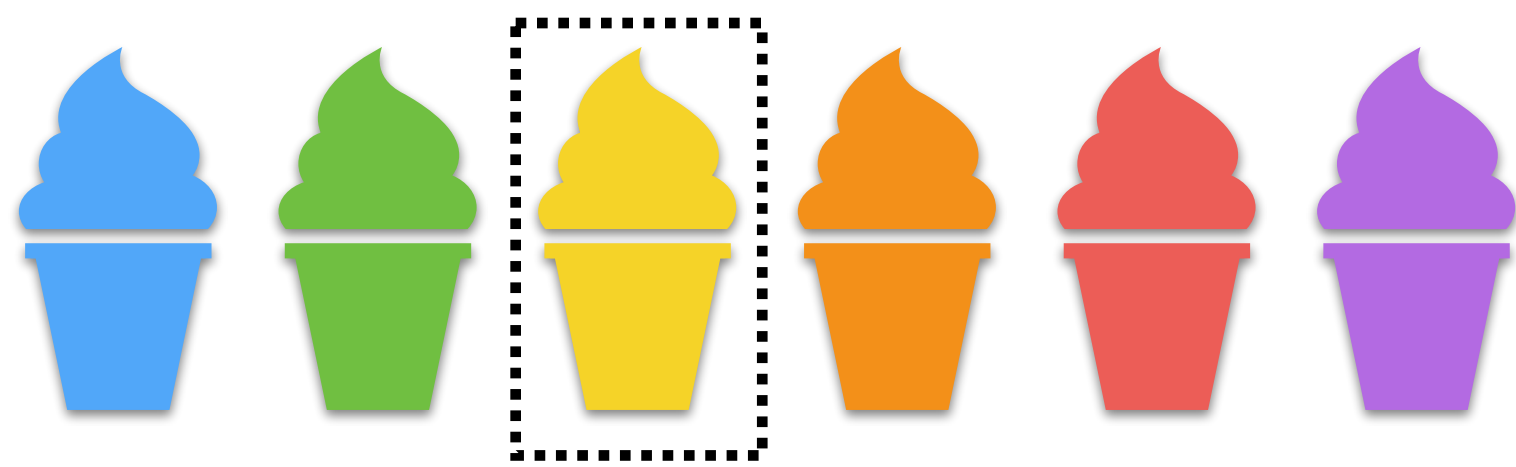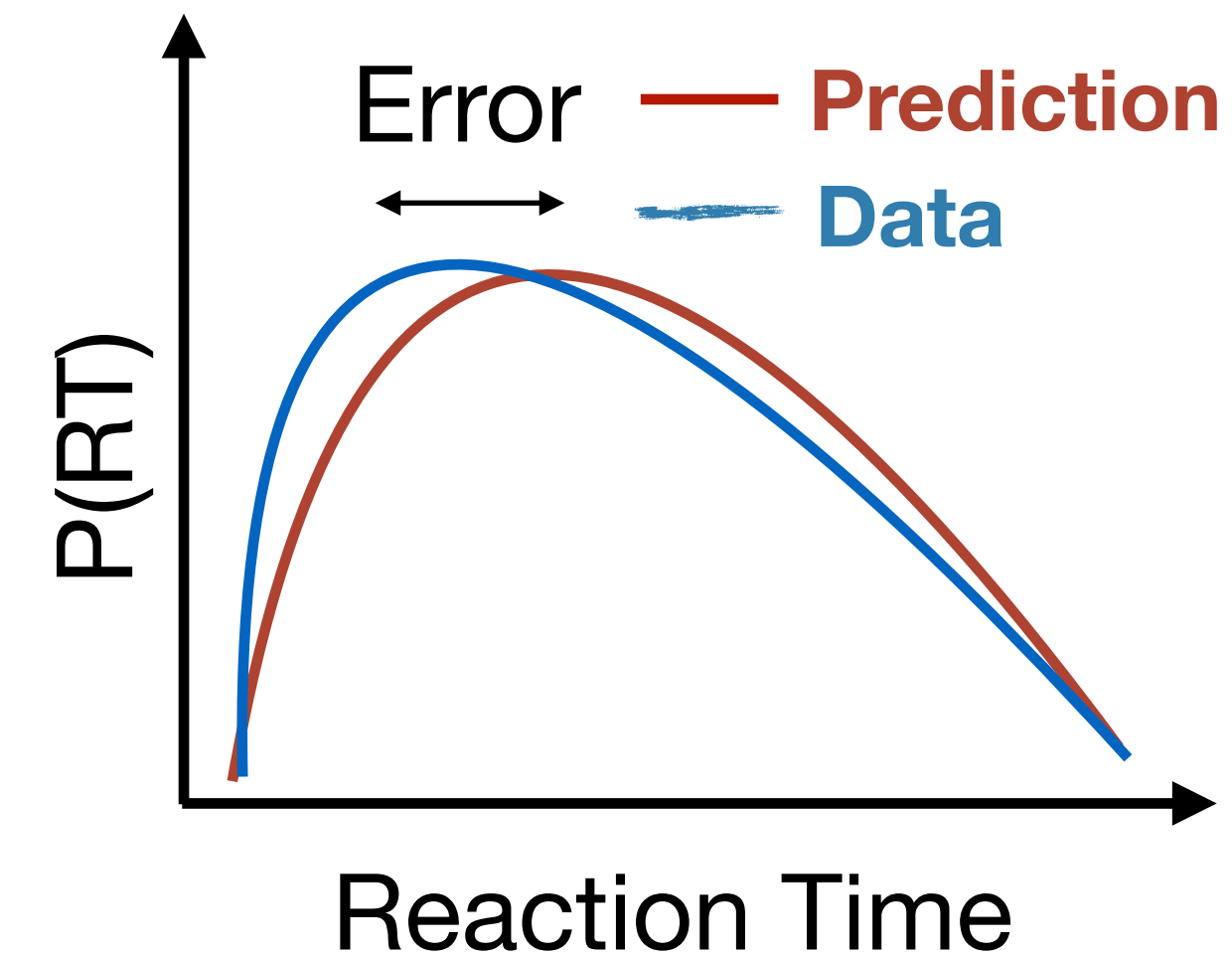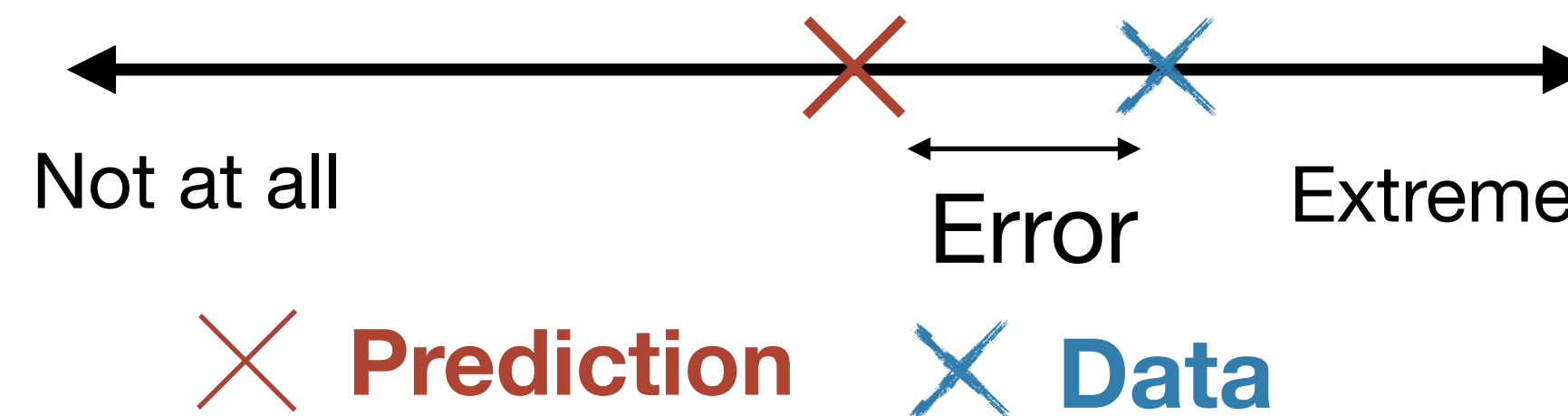Judgments and reaction times are continuous measures

How much do you like ice-cream?



Not at all        Error        Extreme

❌ **Prediction**    ❌ **Data**



· Maximizing likelihood is equivalent to:

  · minimizing Mean Squared Error (MSE)

  · minimizing KL-Divergence

· MSE and KL-Divergence can also be transformed into likelihoods