

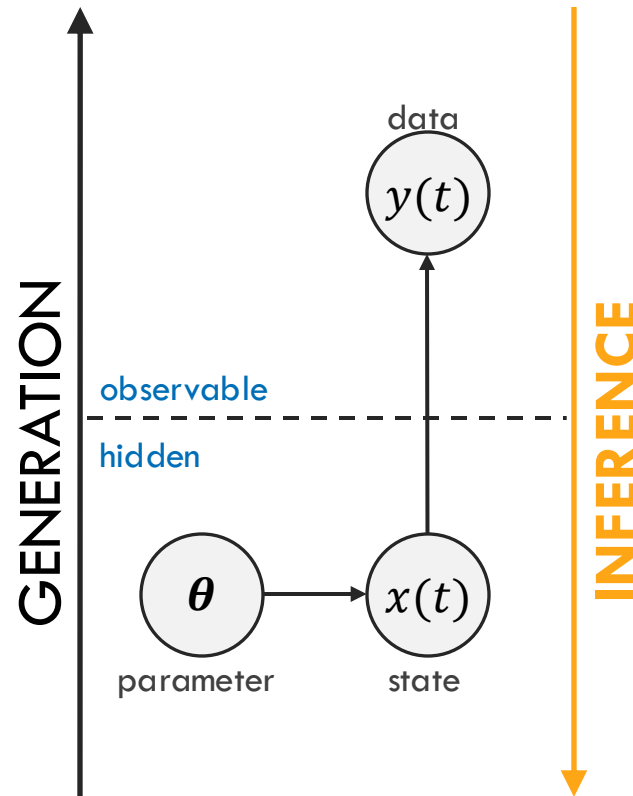
Fitting a Model: Maximum Likelihood Estimation (MLE)

Florian M. Schönleitner

Recap: generative modeling

Last talk:

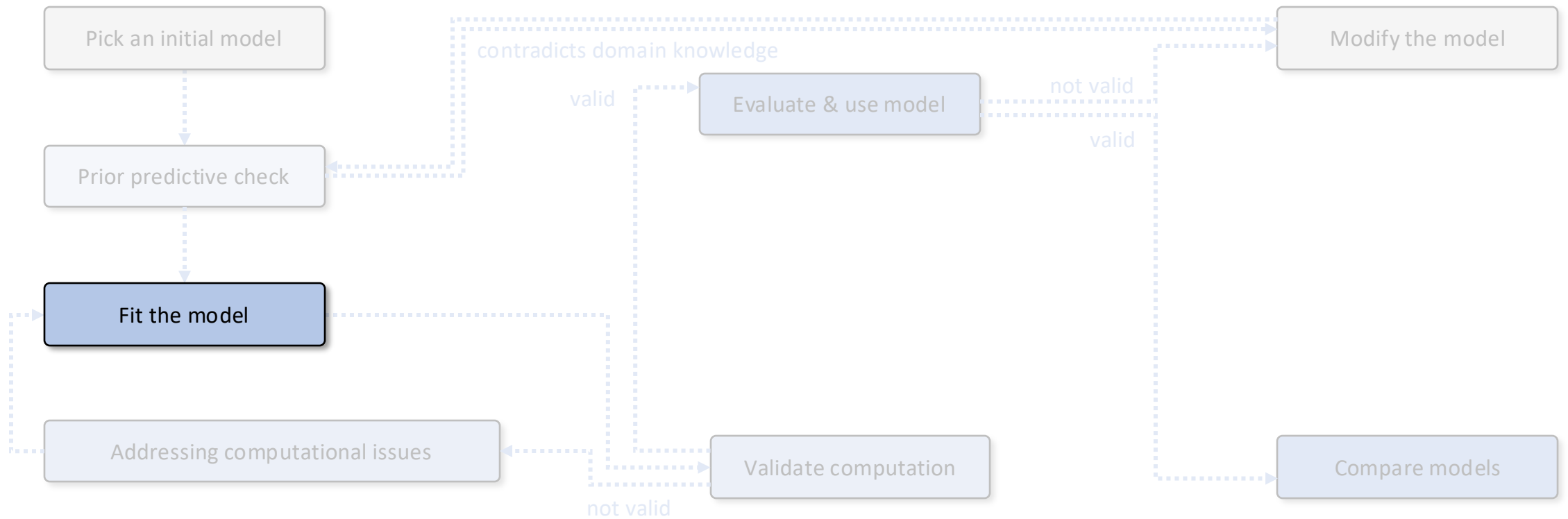
- ✓ Building a model
- ✓ Simulating data



This talk:

- ? Fitting the model to observed data

Recap: Bayesian workflow



MLE: maximum likelihood estimator

Principle:

Find the parameters θ for which the **acquired data Y is most likely** under the model m .

$$\theta_{MLE} = \arg \max_{\theta \in \Theta} \underbrace{p(Y \mid \theta, m)}_{\text{Likelihood}}$$

where

$$p(Y \mid \theta, m) = p(y_{1...T} \mid \theta, m)$$

m	model
Θ	parameter space
θ	model parameters
θ_{MLE}	mle estimate of θ
Y	observed dataset
y_t	single observation
T	number of trials



Example: slot machines

Understand how people learn to maximize their rewards in a case where the most rewarding choice is initially unknown.



vs.



$$p(\text{money} | \text{Slot machine 1}) = 0.8$$

$$p(\text{money} | \text{Slot machine 2}) = 0.2$$

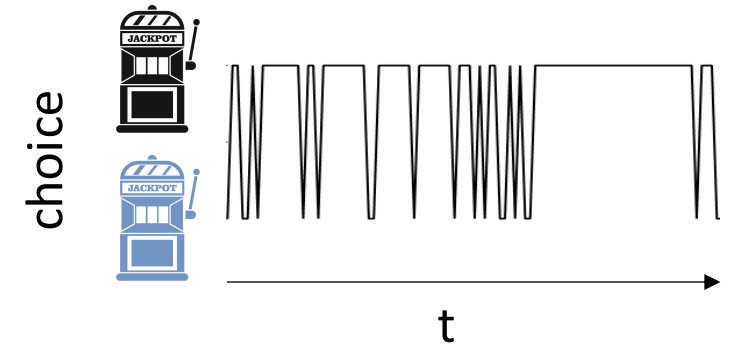
Experiment

Observations

Dataset:

Choice y_t in each trial t

$$Y = (y_1, \dots, y_T)$$



Specifying the likelihood function

Model 1

Random choice

$$p_t^1 = b$$

$$p_t^0 = 1 - b$$

$$0 \leq b \leq 1$$

$$\theta = \{b\}$$

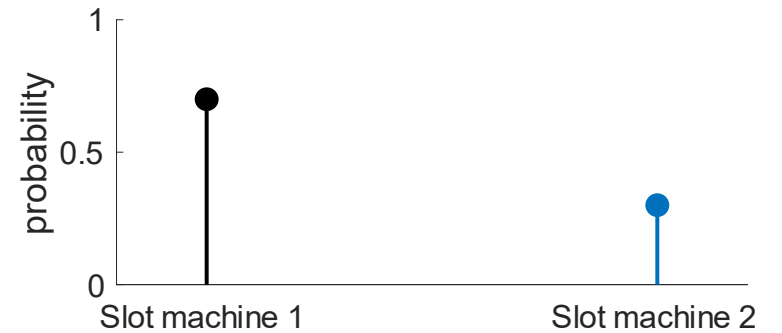
For a single trial t :

$$Y \stackrel{iid}{=} \{y_1, \dots, y_T\}$$

For all trials $1 \dots T$:

$$p(y_t | \theta, m) = \theta^{y_t}(1 - \theta)^{(1-y_t)} \\ = \textit{Bernoulli}$$

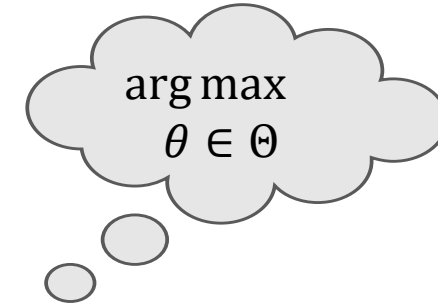
$$p(Y | \theta, m) = p(y_{1 \dots T} | \theta, m) = \prod_{t=1}^T \theta^{y_t}(1 - \theta)^{(1-y_t)}$$



Maximising the likelihood function

Likelihood function

$$p(Y | \theta, m) = \prod_{t=1}^T p(y_t | \theta, m)$$



Analytical solution

Is the likelihood tractable?

Is the likelihood differentiable?

→ Solve $\frac{d}{d\theta} p(Y | \theta, m) \stackrel{!}{=} 0$ and find maximum

Numerical solution

Use numerical optimization routines available in different software (MATLAB, Python, Julia, etc.)

→ Implement $p(Y | \theta, m)$ and find the maximum

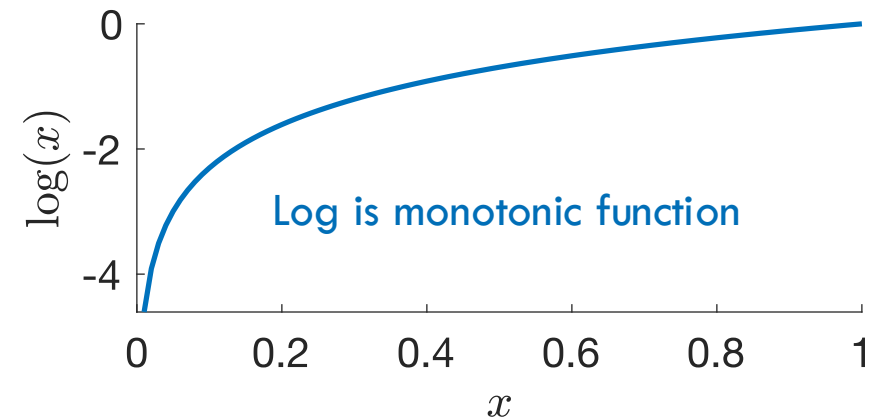
Maximising the likelihood function

Analytic solution

$$\begin{aligned} p(Y \mid \theta, m) &= \prod_{t=1}^T \theta^{y_t} (1 - \theta)^{(1-y_t)} \\ \downarrow 1. \\ \log p(Y \mid \theta, m) &= \log \prod_{t=1}^T \theta^{y_t} (1 - \theta)^{(1-y_t)} \\ &= \sum_{t=1}^T \log \theta^{y_t} (1 - \theta)^{(1-y_t)} \end{aligned}$$

1. Change product to sum by log-transformation:

$$\log \left(\prod_t x_t \right) = \sum_t \log x_t$$



Maximising the **likelihood function**

Analytic solution

$$\begin{aligned} p(Y \mid \theta, m) &= \prod_{t=1}^T \theta^{y_t} (1 - \theta)^{(1-y_t)} \\ \text{1.} \downarrow \\ \log p(Y \mid \theta, m) &= \log \prod_{t=1}^T \theta^{y_t} (1 - \theta)^{(1-y_t)} \\ &= \sum_{t=1}^T \log \theta^{y_t} (1 - \theta)^{(1-y_t)} \\ \text{2.} \downarrow \\ &= \sum_{t=1}^T y_t \log(\theta) + (1 - y_t) \log(1 - \theta) \end{aligned}$$

1. Change product to sum by log-transformation:

$$\log \left(\prod_t x_t \right) = \sum_t \log x_t$$

2. Logarithm of a power:
 $\log(x^a) = a \log x$

Maximising the **likelihood function**

Analytic solution

$$\frac{d}{d\theta} \sum_{t=1}^T y_t \log(\theta) + (1 - y_t) \log(1 - \theta) \stackrel{!}{=} 0$$

$$\left(\frac{d}{d\theta} \log(\theta) \right) \left(\sum_{t=1}^T y_t \right) + \left(\frac{d}{d\theta} \log(1 - \theta) \right) \left(\sum_{t=1}^T 1 - y_t \right) \stackrel{!}{=} 0$$

$$\frac{1}{\theta(1 - \theta)} \left(\sum_{t=1}^T y_t - \theta \sum_{t=1}^T y_t - \theta T + \theta \sum_{t=1}^T y_t \right) \stackrel{!}{=} 0$$

$$\sum_{t=1}^T y_t - \theta T \stackrel{!}{=} 0$$

MLE estimate

$$\theta_{MLE} = \frac{1}{T} \sum_{t=1}^T y_t$$



MLE estimate is arithmetic mean of data!

Advantages of maximum likelihood estimation

1. **Easy to compute:** simple optimisation problem

Advantages of maximum likelihood estimation

1. **Easy to compute:** simple optimisation problem
2. **Interpretable:** often θ_{MLE} is intuitively interpretable wrt. model parameters (see MLE of random choice model)

Advantages of maximum likelihood estimation

1. **Easy to compute:** simple optimisation problem
2. **Interpretable:** often θ_{MLE} is intuitively interpretable wrt. model parameters (see MLE of random choice model)
3. **Asymptotic properties:** consistency (true parameter value recovered) and efficiency (lowest possible parameter variance)

Advantages of maximum likelihood estimation

1. **Easy to compute:** simple optimisation problem
2. **Interpretable:** often θ_{MLE} is intuitively interpretable wrt. model parameters (see MLE of random choice model)
3. **Asymptotic properties:** consistency (true parameter value recovered) and efficiency (lowest possible parameter variance)
4. **Invariant to reparameterization:** if $\theta_{MLE} = \text{MLE}(\theta)$ for $\theta \in \Theta$, then $g(\theta_{MLE}) = \text{MLE}(g(\theta))$ for $g: \mathbb{R} \rightarrow \mathbb{R}$

Limitations of maximum likelihood estimation

1. **Point estimate:** θ_{MLE} is a point estimate \rightarrow no representation of uncertainty

Limitations of maximum likelihood estimation

1. **Point estimate:** θ_{MLE} is a point estimate \rightarrow no representation of uncertainty
2. **Existence & uniqueness:** The MLE might not be unique or even non existent, depending on properties of the likelihood function and parameter space

Limitations of maximum likelihood estimation

1. **Point estimate:** θ_{MLE} is a point estimate \rightarrow no representation of uncertainty
2. **Existence & uniqueness:** The MLE might not be unique or even non existent, depending on properties of the likelihood function and parameter space
3. **Overfitting:** MLE is limited to a finite set of observed datapoints

Limitations of maximum likelihood estimation

3. **Overfitting:** MLE is limited to observed data but does not (explicitly) take into account prior information

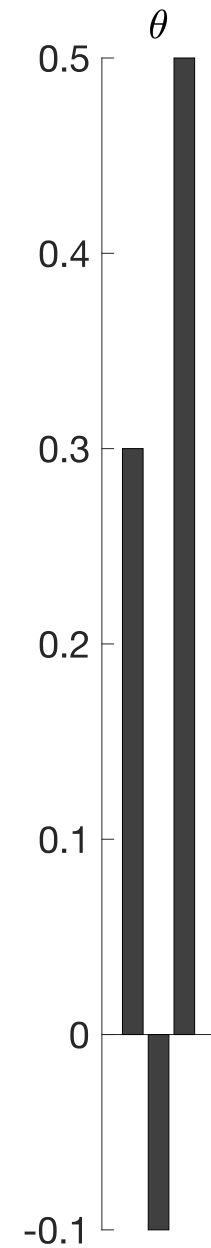
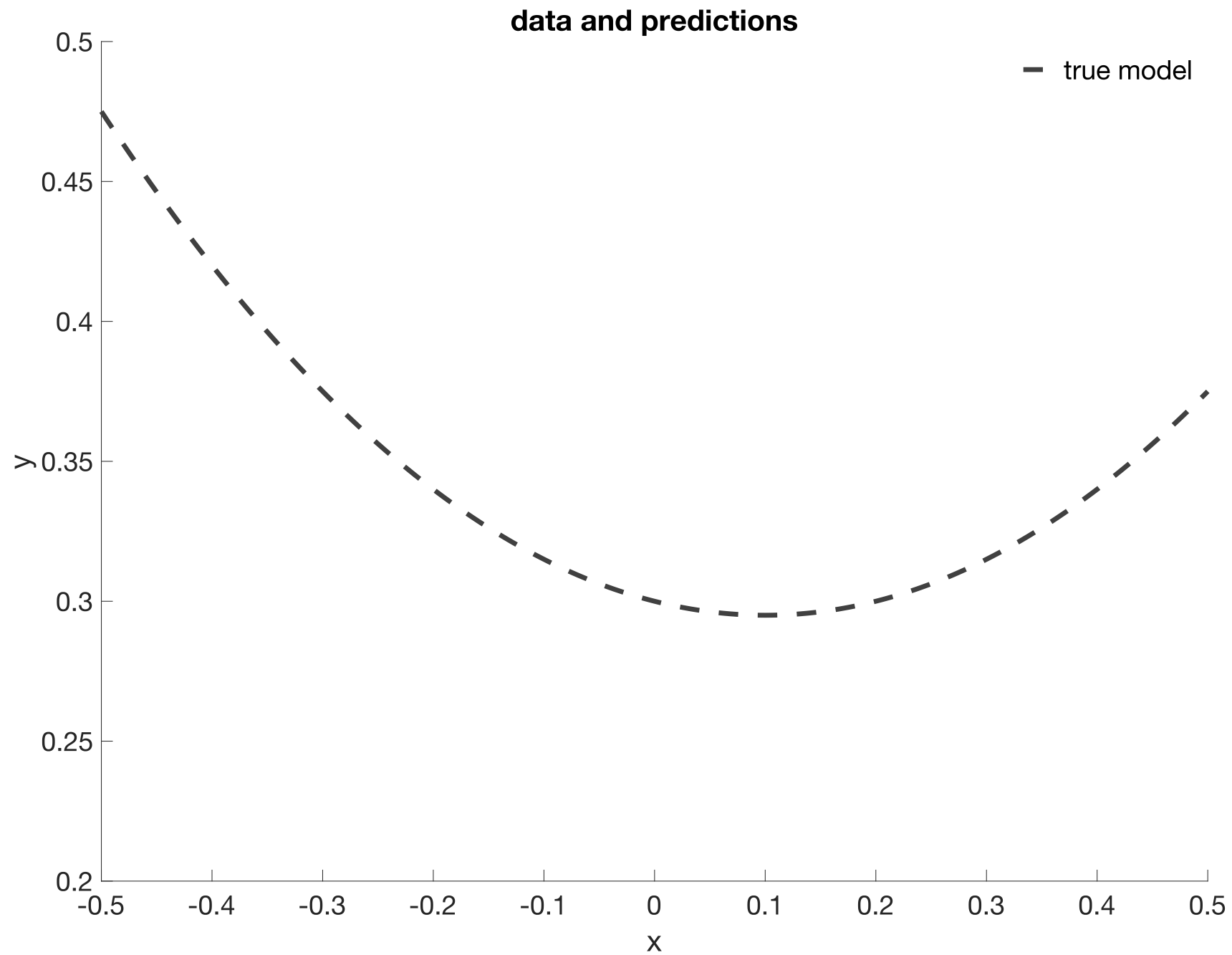
Example

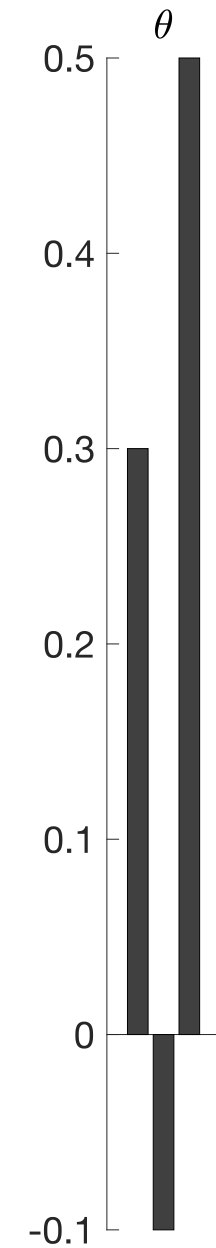
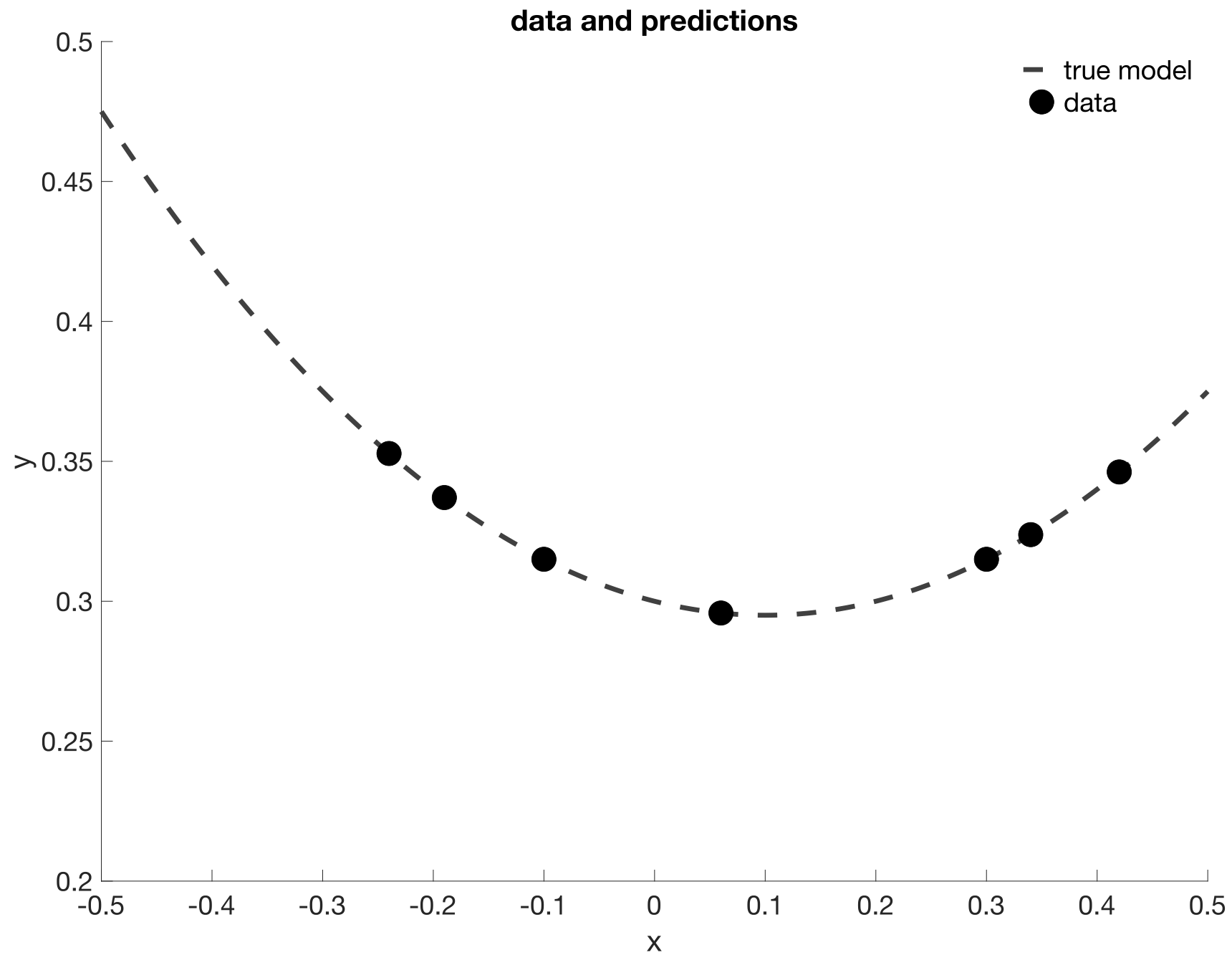
Polynomial model of order P with Gaussian noise

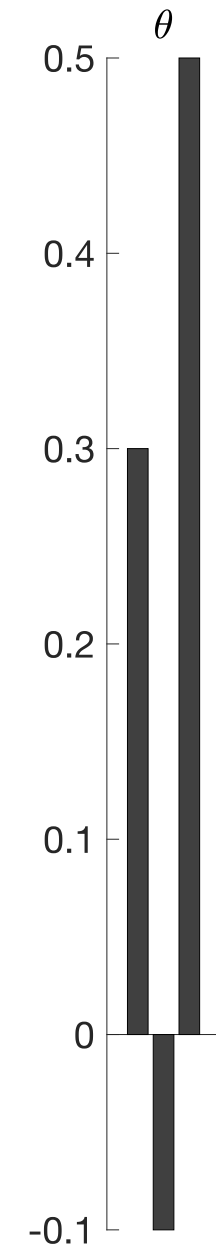
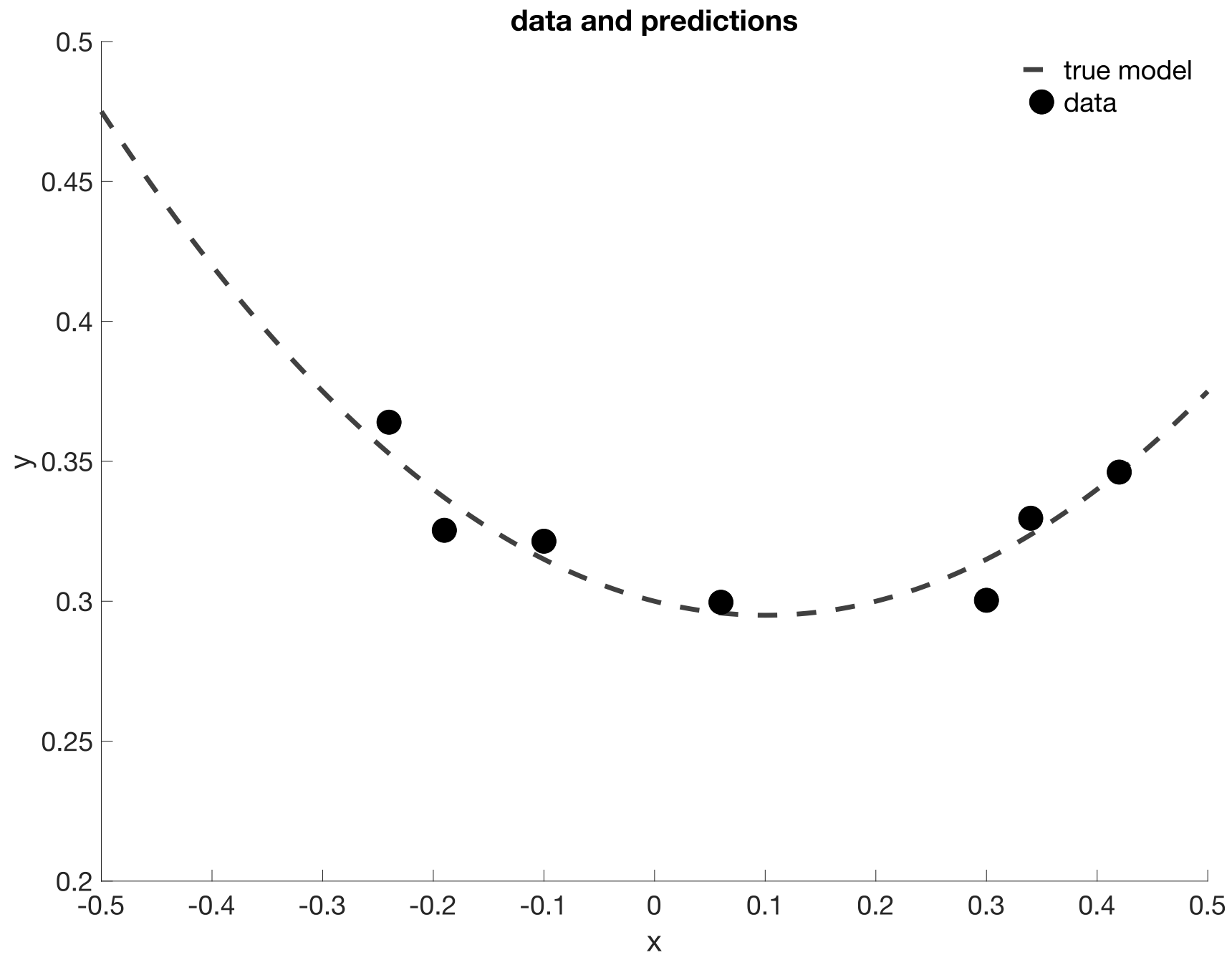
$$y = \theta_0 + \theta_1 x + \dots + \theta_P x^P + \epsilon = \mathbf{x}\boldsymbol{\theta} + \epsilon$$

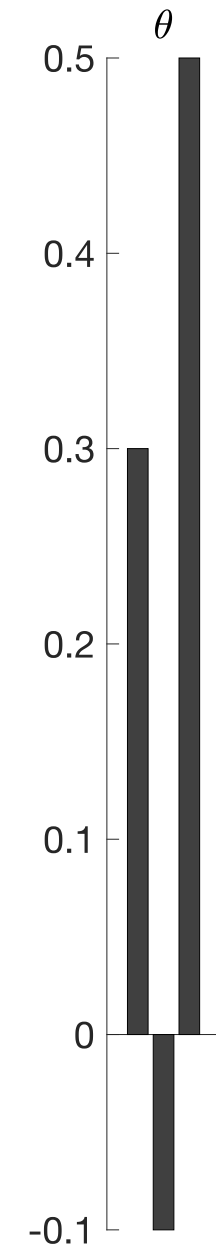
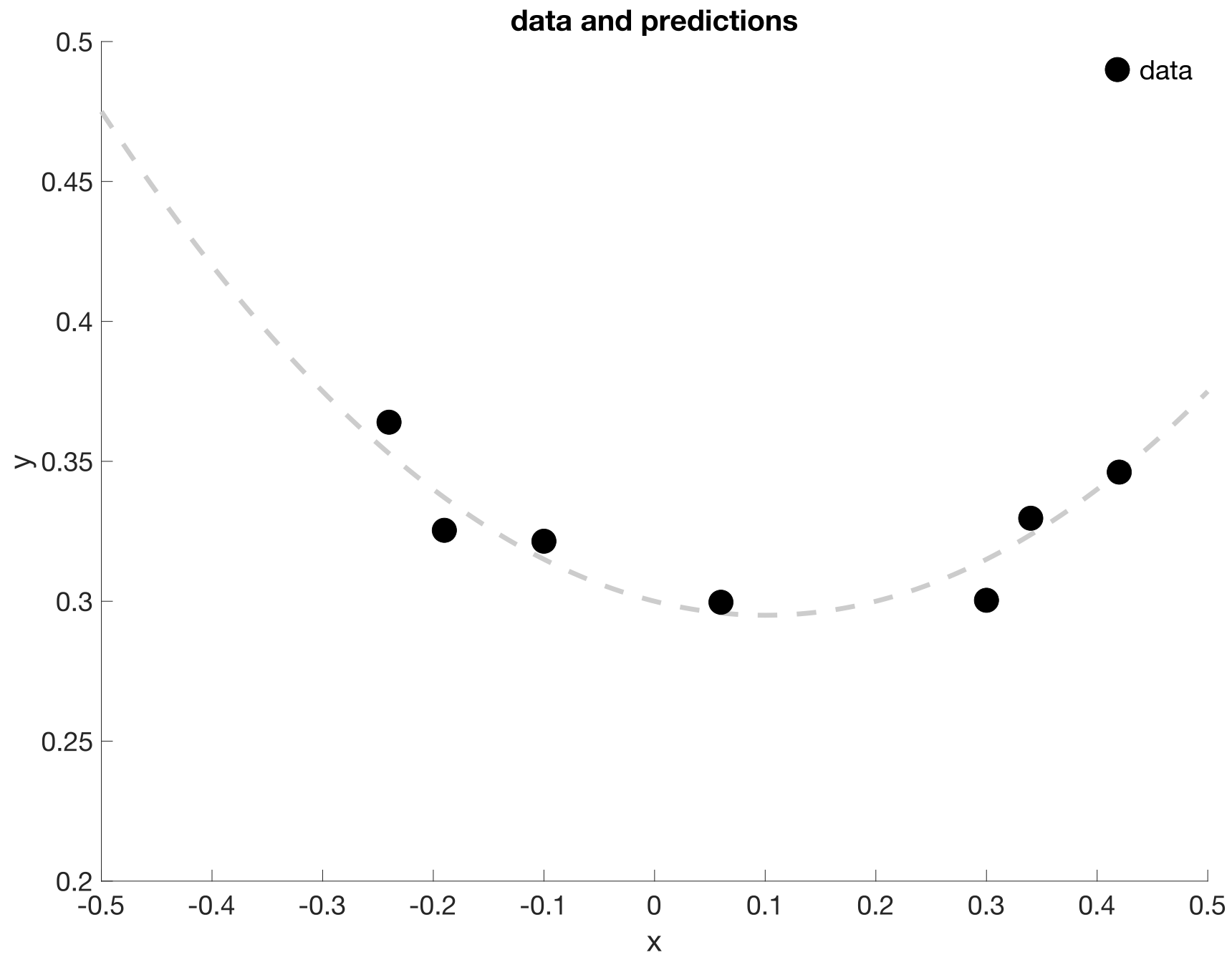
where

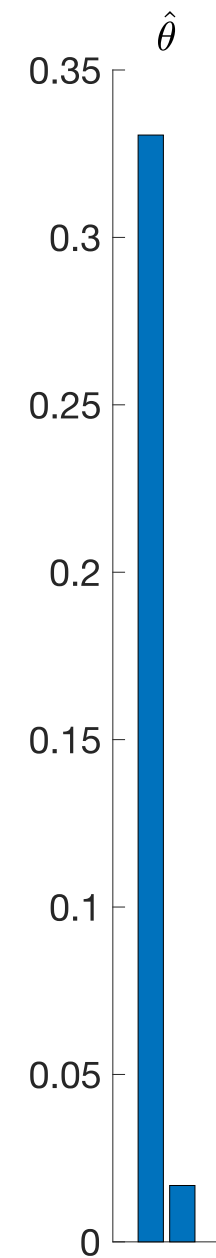
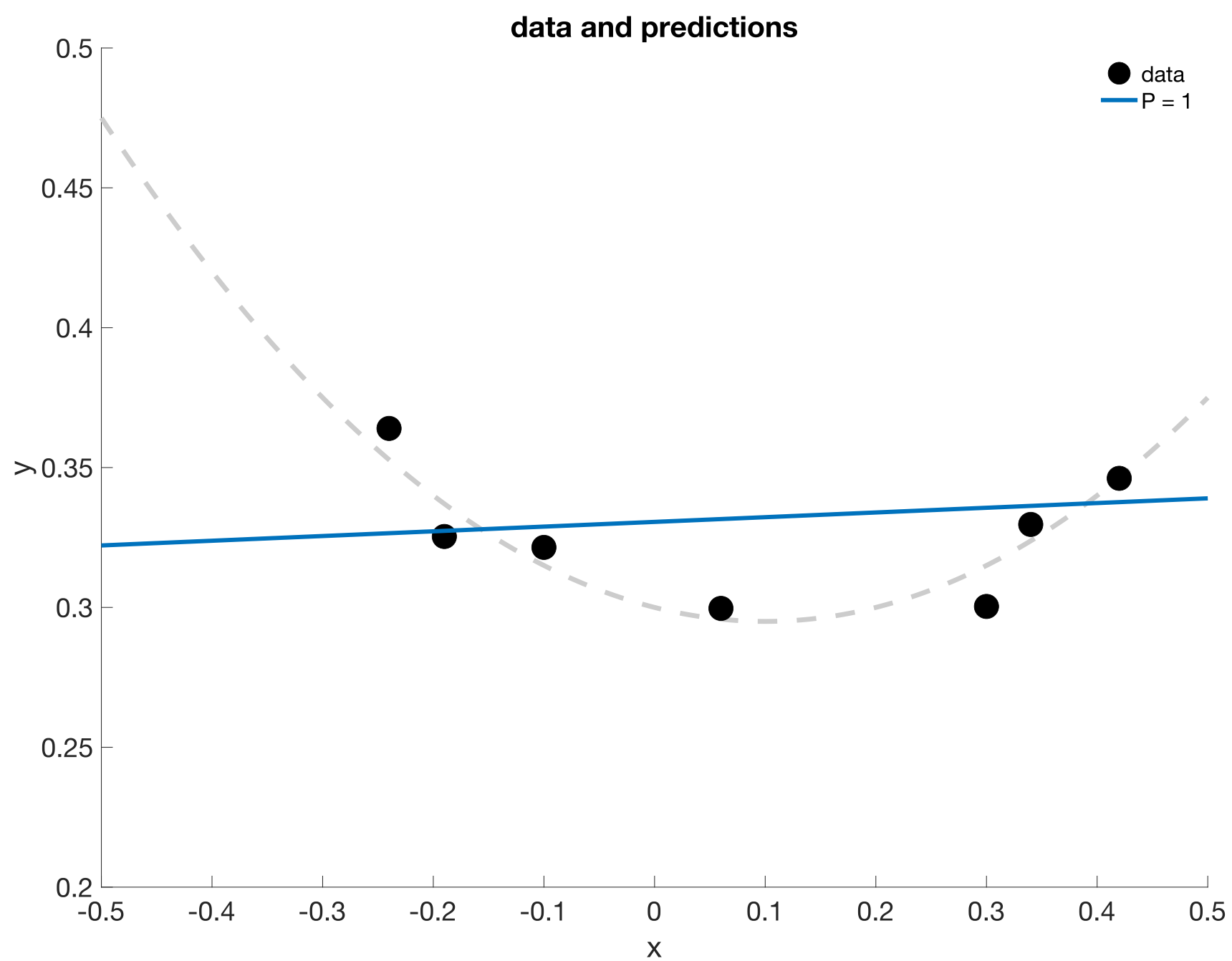
$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

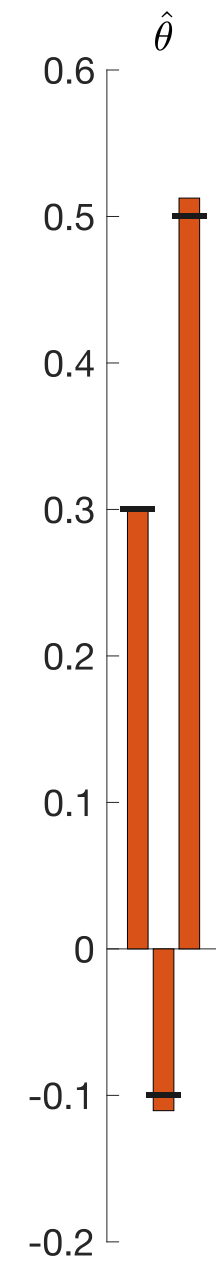
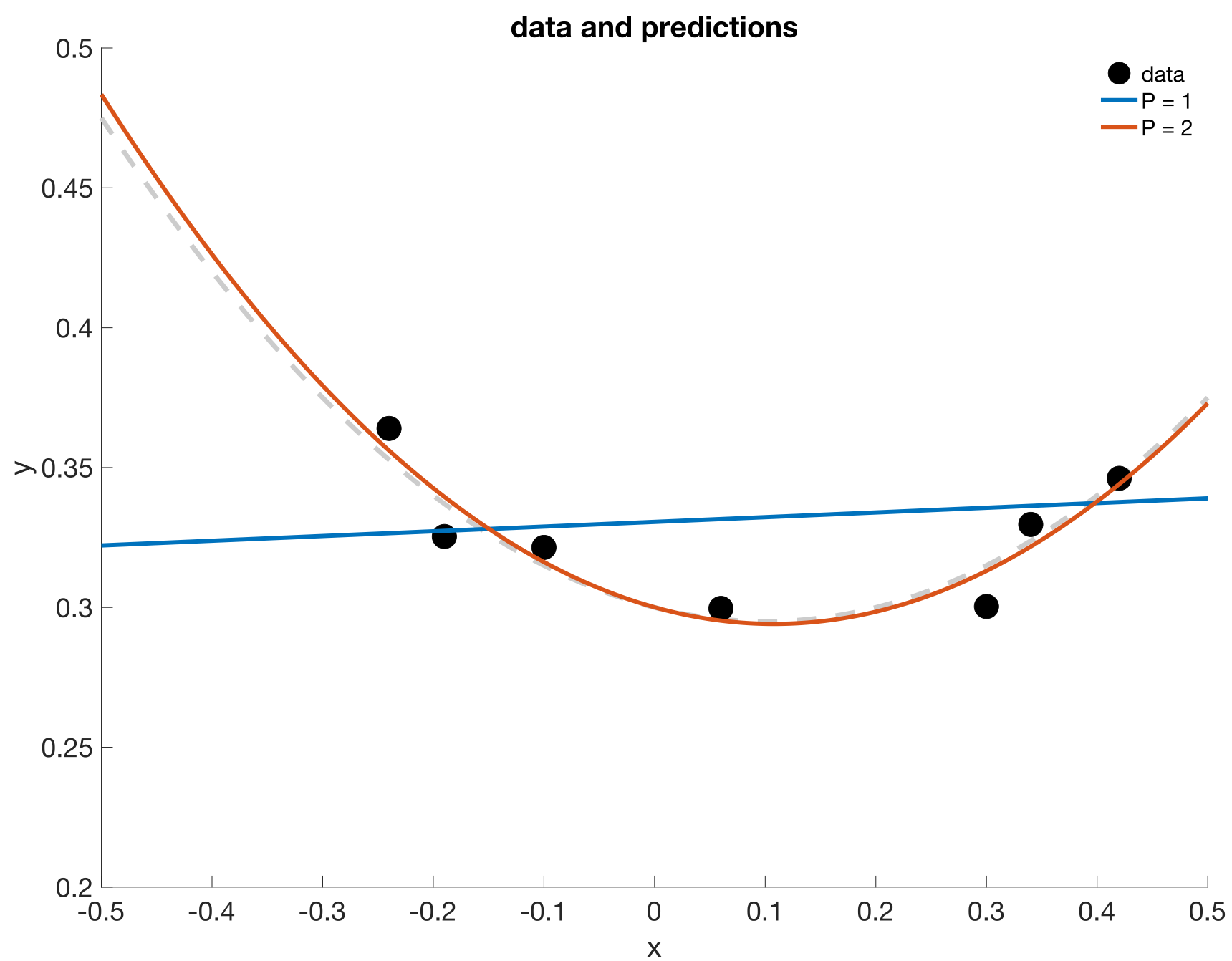


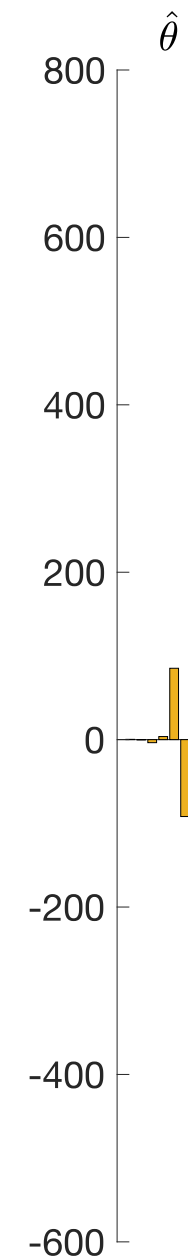
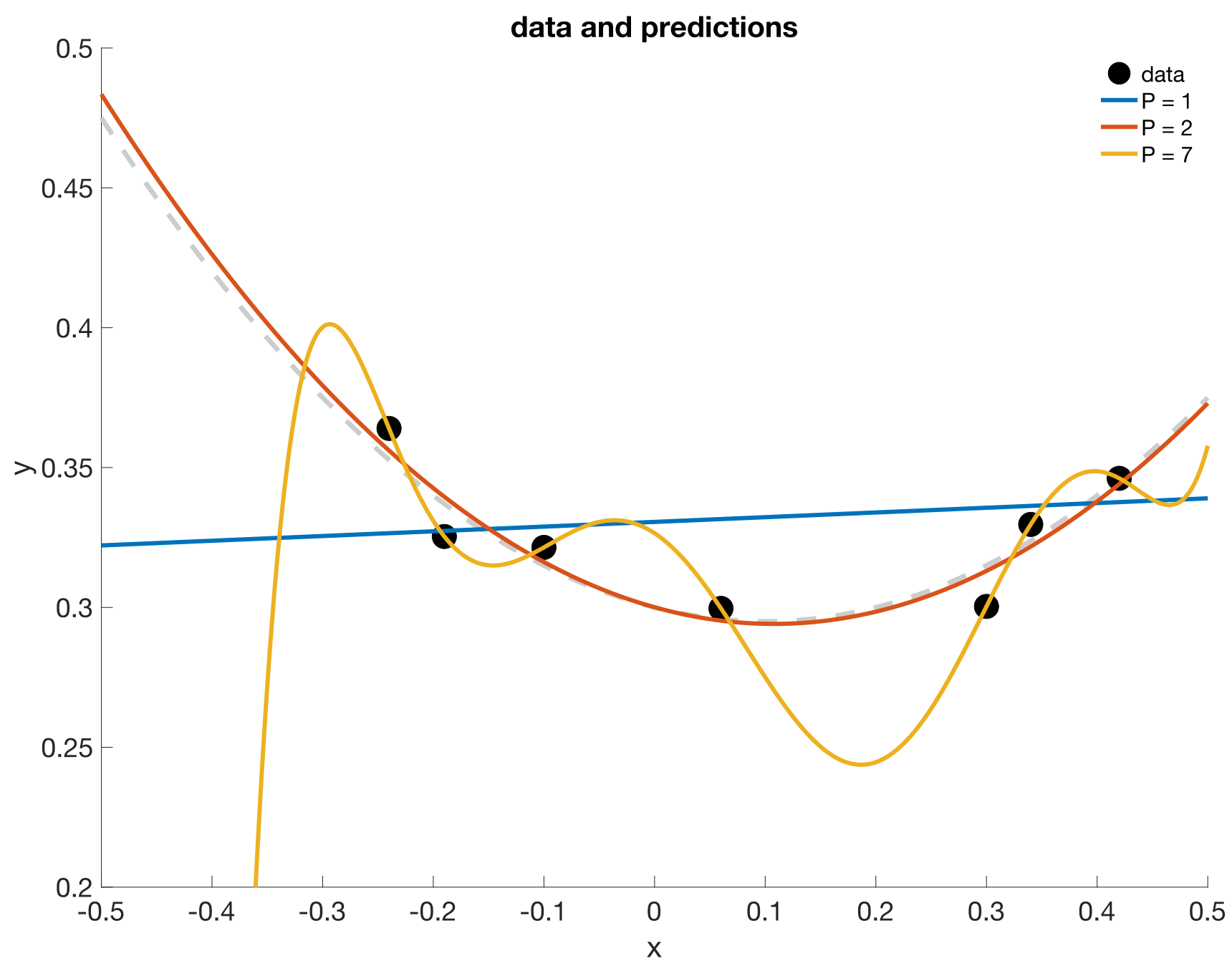


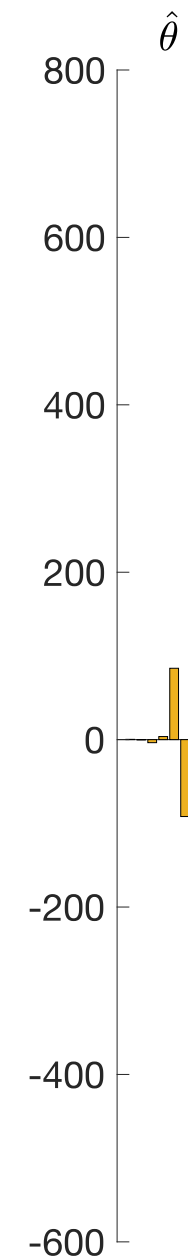
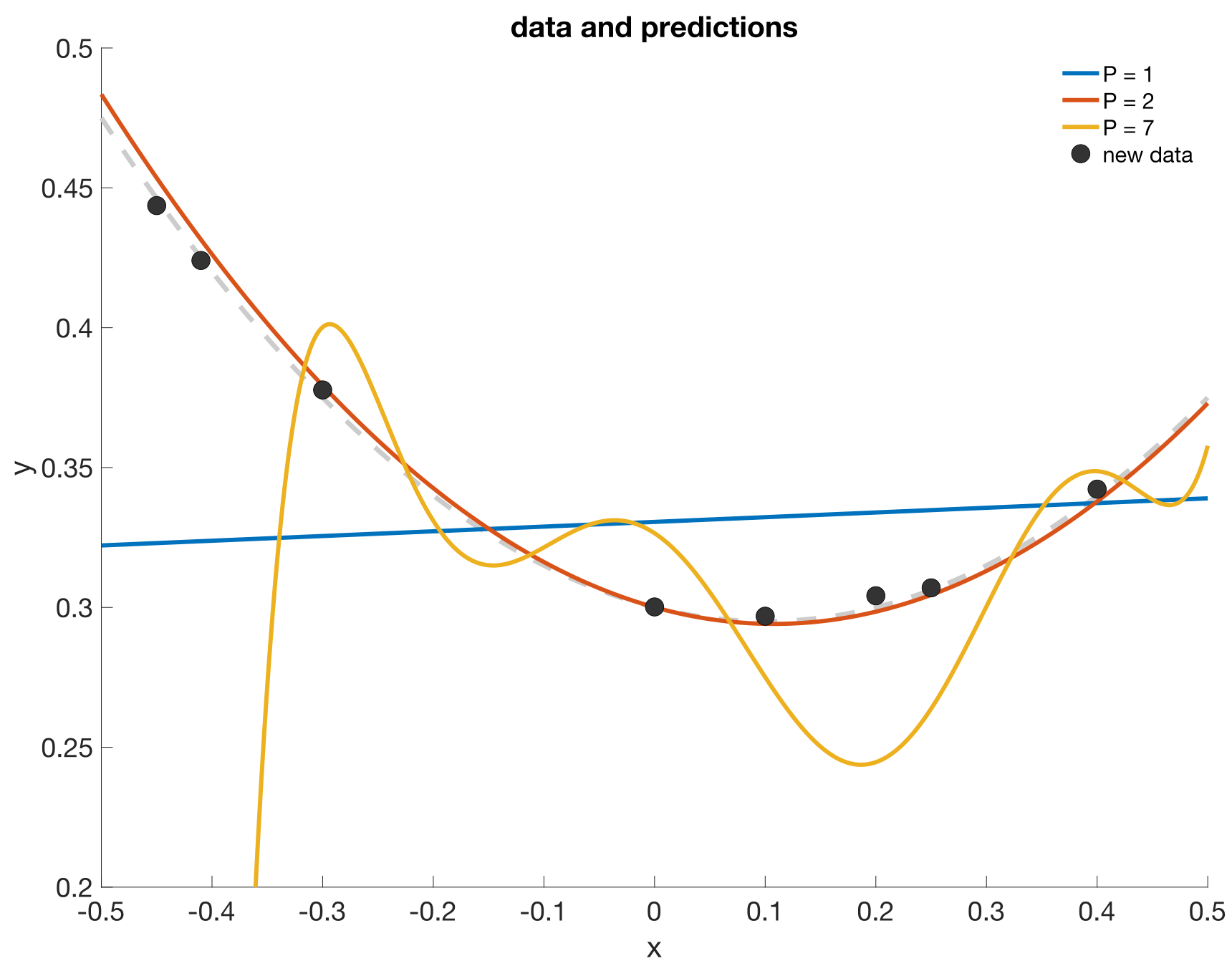













Maximum likelihood vs. full Bayesian inference

$$p(\theta \mid y, m) = \frac{p(y \mid \theta, m)p(\theta, m)}{p(y \mid m)}$$


Maximum-a-posteriori (**MAP**) estimation:

- Point estimate of the posterior
- Under a flat prior MAP=MLE

Variational Bayes (**VB**), sampling-based (**MCMC**) techniques

- Full posterior densities

Acknowledgement

Special thanks to my TNU colleagues



Stefan Frässle



Katharina
Wellstein



Jakob Heinzle



Klaas Enno
Stephan

QUESTIONS?



Schoenleitner@biomed.ee.ethz.ch



<https://github.com/computational-psychiatry-course/cpc2024>



<https://www.linkedin.com/in/florian-schönleitner-34201210a/>