



# Machine Learning 101a

---

Jan Peters  
Gerhard Neumann

# Purpose of this Lecture

---



- Statistics and Math Refresher
- Foundations of machine learning tools for robotics
- We focus on regression methods and general principles
  - Often needed in robotics
  - More on machine learning in general: Machine Learning
    - Statistical Approaches 1

# Content of this Lecture

---



- **Math and Statistics Refresher**
- What is Machine Learning?
- Model-Selection
- Linear Regression
  - Gauss' Approach
  - Gauss' Approach
  - Frequentist Approach
  - Bayesian Approach

# Statistics Refresher: Sweet memories from High School...

---



## What is a **random variable** $X$ ?

$X$  is a variable whose value  $x$  is subject to variations due to chance

## What is a **distribution** $p(X = x)$ ?

Describes the probability that the random variable will be equal to a certain value

## What is an **expectation**?

$$\mathbb{E}_{p(X)}[f(x)] = \int p(x)f(x)dx$$



# Statistics Refresher: Sweet memories from High School...

---



- **What is a **joint**, a **conditional** and a **marginal** distribution?**

$$p(x, y) = p(y|x)p(x)$$

joint = conditional  $\times$  marginal

- **What is **independence** of random variables?**

$$p(x, y) = p(x)p(y)$$

- **What does **marginalization** mean?**

$$p(x) = \int p(x, y) dy$$

- **And finally... what is **Bayes Theorem**?**

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$



# Math Refresher: Some more fancy math...



- From now on, matrices are your friends... derivatives too

$$\frac{dy}{d\mathbf{x}} = \left[ \frac{dy}{dx_1}, \dots, \frac{dy}{dx_n} \right]$$

$$\frac{d\mathbf{y}}{d\mathbf{x}} = \begin{bmatrix} \frac{dy_1}{dx_1} & \dots & \frac{dy_1}{dx_n} \\ \vdots & \vdots & \vdots \\ \frac{dy_m}{dx_1} & \dots & \frac{dy_m}{dx_n} \end{bmatrix}$$

- Some more matrix calculus

$$\frac{d\mathbf{a}^T \mathbf{x}}{d\mathbf{x}} = \frac{d\mathbf{x}^T \mathbf{a}}{d\mathbf{x}} = \mathbf{a}^T$$

$$\frac{d\mathbf{A}\mathbf{x}}{d\mathbf{x}} = \mathbf{A}$$

$$\frac{d\mathbf{x}^T \mathbf{A} \mathbf{x}}{d\mathbf{x}} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T)$$

- Need more ? [Wikipedia on Matrix Calculus](#) or [The Matrix Cookbook](#)

# Math Refresher: Inverse of matrices

---



**How can we invert a matrix that is not a square matrix?**

$$\mathbf{J} \in \mathbb{R}^{n \times m}$$

**Left-Pseudo Inverse:**

works, if  $\mathbf{J}$  has full column rank

$$\mathbf{J}^\dagger \mathbf{J} = \underbrace{(\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T}_{\text{left multiplied}} \mathbf{J} = \mathbf{I}_m$$

**Right Pseudo Inverse:**

works, if  $\mathbf{J}$  has full row rank

$$\mathbf{J} \mathbf{J}^\dagger = \mathbf{J} \underbrace{\mathbf{J}^T (\mathbf{J} \mathbf{J}^T)^{-1}}_{\text{right multiplied}} \mathbf{J}^T = \mathbf{I}_n$$



# Statistics Refresher: Meet some old friends...

- **Gaussian Distribution**

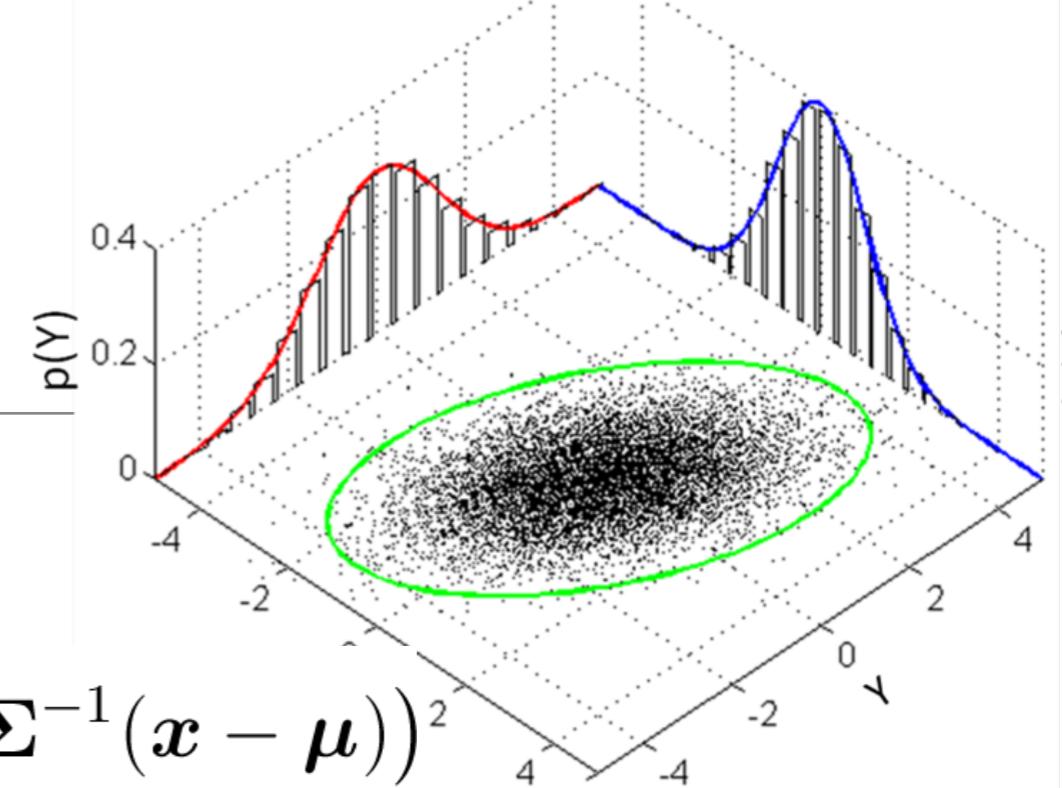
$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{|2\pi\Sigma|^{\frac{1}{2}}} \exp\left(-0.5(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

- Covariance matrix  $\Sigma$  captures **linear correlation**
- Product: **Gaussian stays Gaussian**

$$\mathcal{N}(x|a, A) \mathcal{N}(x|b, B) = \mathcal{N}(x| \dots)$$

- Mean is also the mode

$$\mu = \operatorname{argmax}_x \mathcal{N}(x|\mu, \Sigma)$$



# Statistics Refresher: Meet some old friends...



- Joint from Marginal and Conditional

$$\mathcal{N}(\mathbf{x}|\mathbf{a}, \mathbf{A}) \mathcal{N}(\mathbf{y}|\mathbf{b} + \mathbf{F}\mathbf{x}, \mathbf{B}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \mid \begin{bmatrix} \mathbf{a} \\ \mathbf{b} + \mathbf{F}\mathbf{a} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{F}^T \mathbf{A}^T \\ \mathbf{F}^T \mathbf{A}^T & \mathbf{B} + \mathbf{F}\mathbf{A}^T \mathbf{F}^T \end{bmatrix}\right)$$

$$p(\mathbf{x})p(\mathbf{y}|\mathbf{x}) \longrightarrow p(\mathbf{x}, \mathbf{y})$$

- Marginal and Conditional Gaussian from Joint

$$\mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \mid \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{B} \end{bmatrix}\right) = \mathcal{N}(\mathbf{x}|\mathbf{a}, \mathbf{A}) \mathcal{N}(\mathbf{y}|\mathbf{b} + \mathbf{C}^T \mathbf{A}^{-1}(\mathbf{x} - \mathbf{a}), \mathbf{B} - \mathbf{C}^T \mathbf{A}^{-1} \mathbf{C})$$

$$p(\mathbf{x}, \mathbf{y}) \longrightarrow p(\mathbf{x}) \qquad p(\mathbf{y}|\mathbf{x})$$



# Statistics Refresher: Meet some old friends...



- Bayes Theorem for Gaussians

$$p(\mathbf{x})p(\mathbf{y}|\mathbf{x}) \longrightarrow p(\mathbf{x}|\mathbf{y})$$

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{A}) &\longrightarrow p(\mathbf{y}) &= \mathcal{N}(\mathbf{y}|F\mathbf{a}, \sigma^2 \mathbf{I} + F\mathbf{A}F^T) \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|F\mathbf{x}, \sigma^2 \mathbf{I}) &p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\Sigma F^T \mathbf{y}, \sigma^2 \Sigma) \end{aligned}$$

$$\Sigma = (F^T F + \sigma^2 A^{-1})^{-1}$$

**Damped Pseudo Inverse**

- Not enough? Find more stuff [here](#) (credit to Marc Toussaint)



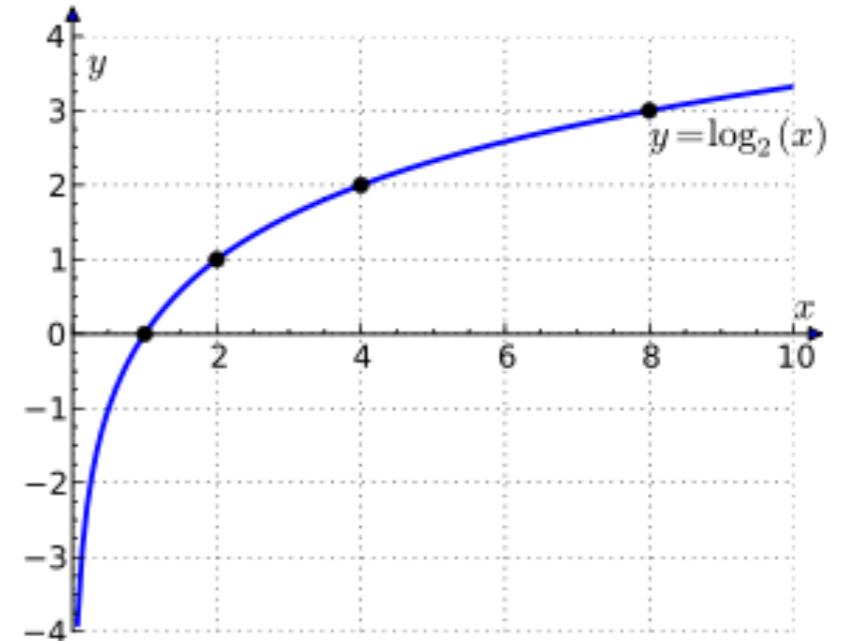
# May I introduce you? The good old logarithm



It's monoton  $a > b \Rightarrow \log(a) > \log(b)$

... but not boring, as:

- Product is easy...  $\log \prod_{i=1}^N a_i = \sum_{i=1}^N \log a_i$
- Division a piece of cake...  $\log \frac{1}{a} = -\log a$
- Exponents also...  $\log(a^b) = b \log a$



# Content of this Lecture

---



- Math and Statistics Refresher
- What is Machine Learning?
- Model-Selection
- Linear Regression
  - Gauss' Approach
  - Frequentist Approach
  - Bayesian Approach

# Why Machine Learning

---



*We are drowning in information and starving for knowledge.*

-John Naisbitt.

## Era of big data:

- In 2008 there are about 1 trillion web pages
- 20 hours of video are uploaded to YouTube every minute
- Walmart handles more than 1M transactions per hour and has databases containing more than 2.5 petabytes ( $2.5 \times 10^{15}$ ) of information.



**No human being can deal with the data avalanche!**

# Why Machine Learning?

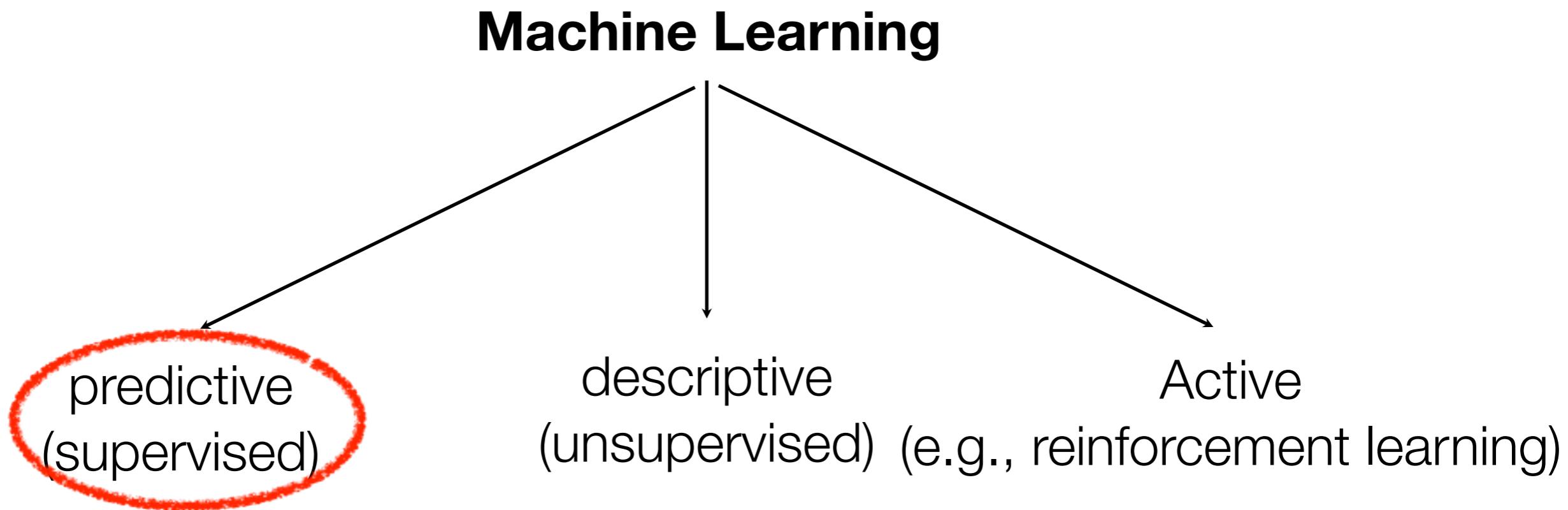
---



I keep saying the sexy job in the next ten years will be statisticians and machine learners. People think I'm joking, but who would've guessed that computer engineers would've been the sexy job of the 1990s? The ability to take data — to be able to understand it, to process it, to extract value from it, to visualize it, to communicate it — that's going to be a hugely important skill in the next decades.

**Hal Varian, 2009**  
*Chief Engineer of Google*

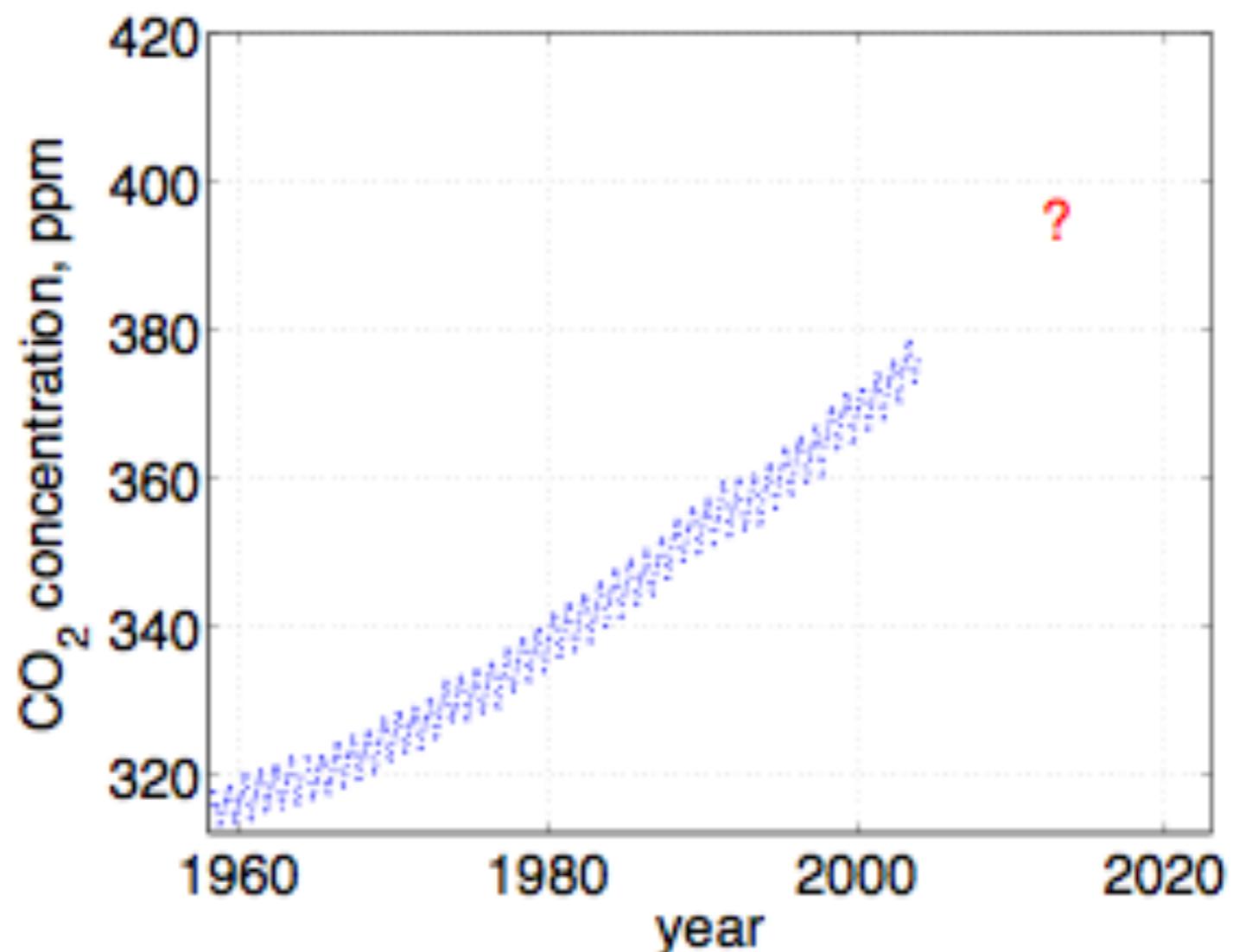
# Types of Machine Learning



# Prediction Problem (=Supervised Learning)



What will be the CO<sub>2</sub> concentration in the future?



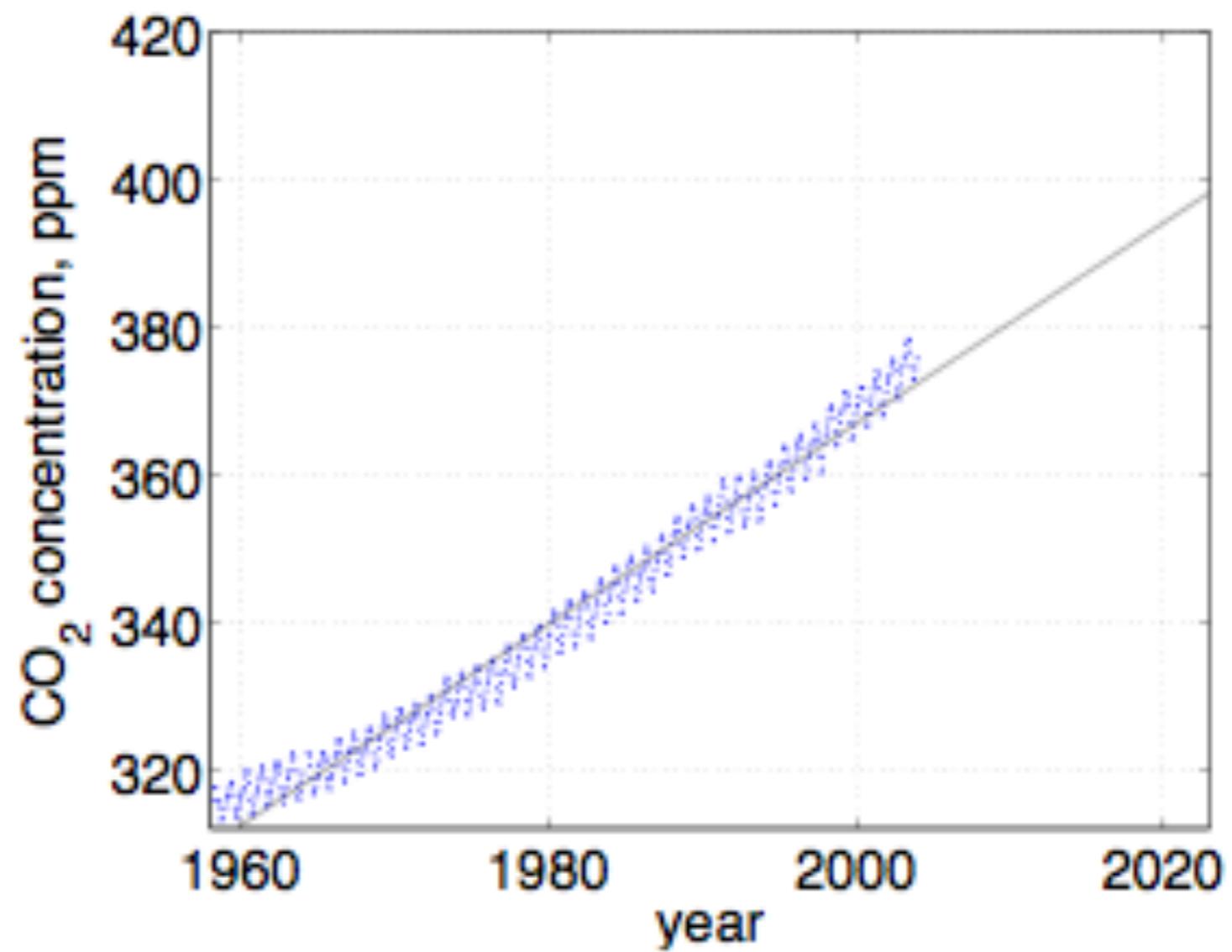
# Prediction Problem (=Supervised Learning)



What will be the CO<sub>2</sub> concentration in the future?

Different prediction models possible

→ Linear



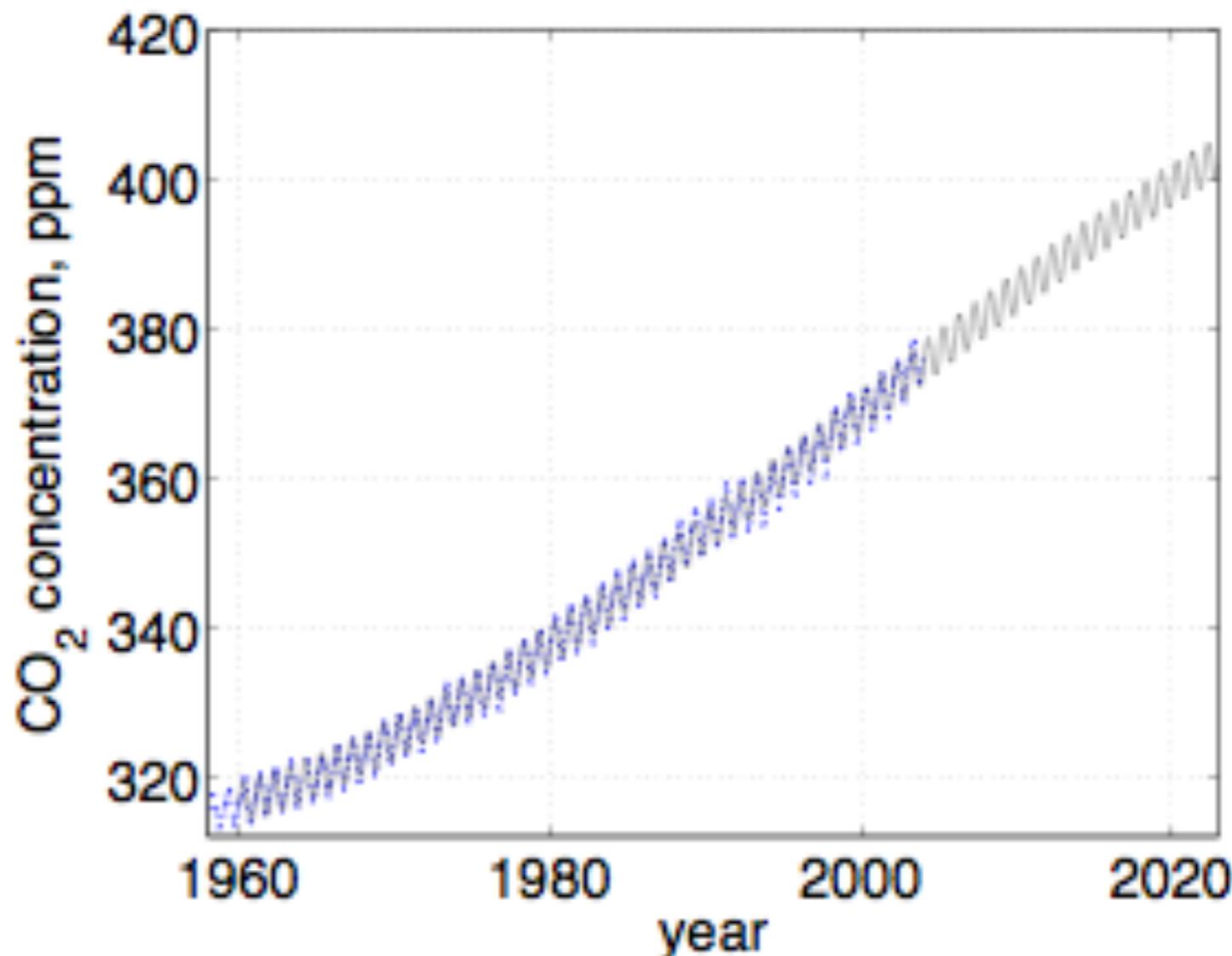
# Prediction Problem (=Supervised Learning)



What will be the CO<sub>2</sub> concentration in the future?

Different prediction models possible

- Linear
- Exponential with seasonal trends





# Formalization of Predictive Problems

---

**In predictive problems, we have the following data-set**

$$\mathcal{D} = \{(x_i, y_i) \mid i = 1, 2, 3, \dots, n\}$$

$x$  . . . inputs,  $y$  . . . output / target

**Two most prominent examples are:**

1. *Classification*: Discrete outputs or labels:  $y_q \in \{0, 1\}$ .

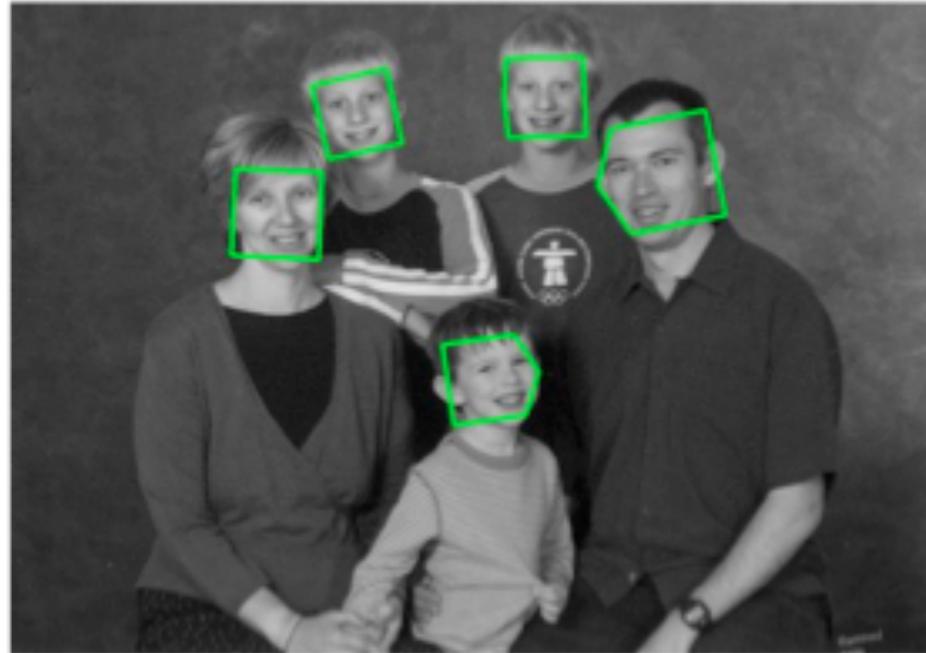
Most likely class:  $y_q = \operatorname{argmax}_y p(y | \mathbf{x} = \mathbf{x}_q)$ .

2. *Regression*: Continuous outputs or labels:  $y_q \in \mathbb{R}$ .

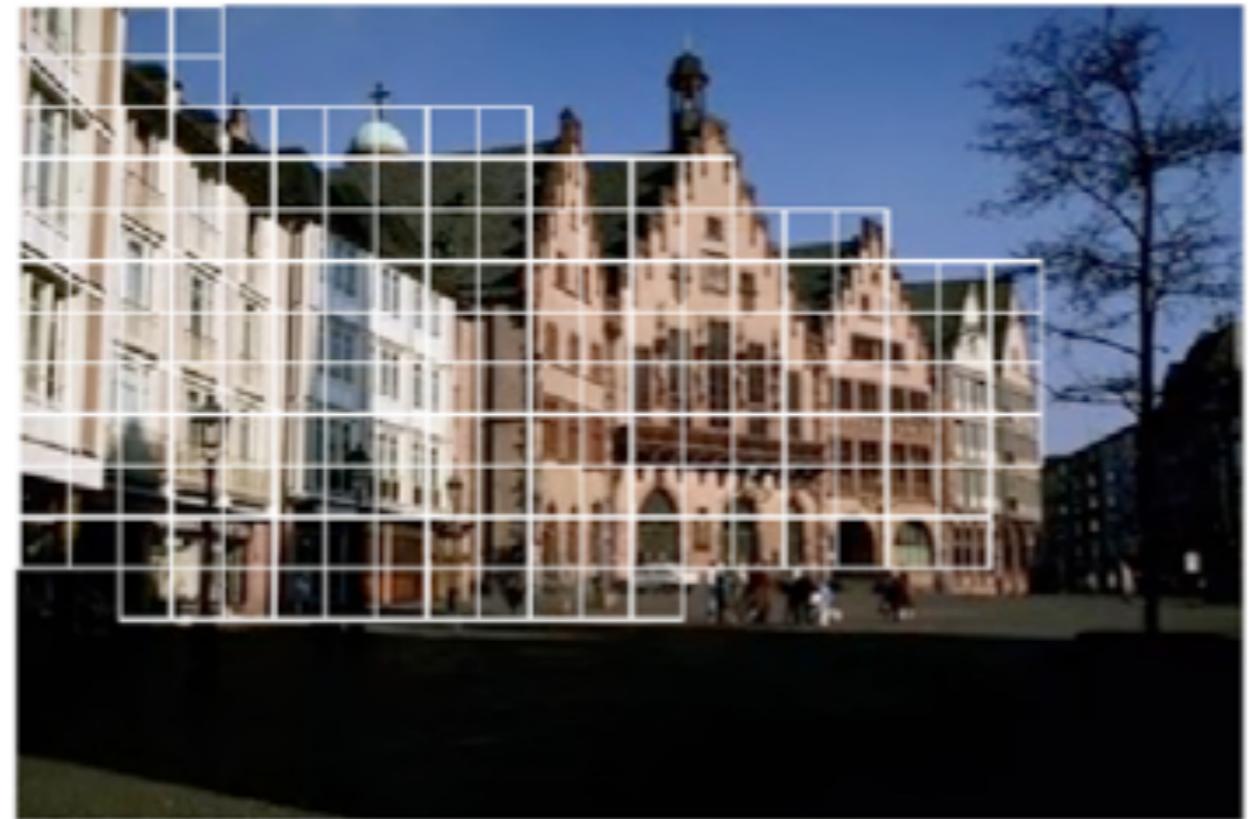
Expected output:  $y_q = E\{y | \mathbf{x} = \mathbf{x}_q\} = \int y p(y | \mathbf{x} = \mathbf{x}_q) dy$ .



# Examples of Classification



- Document classification, e.g., Spam Filtering
- Image classification: Classifying flowers, face detection, face recognition, handwriting recognition, ...



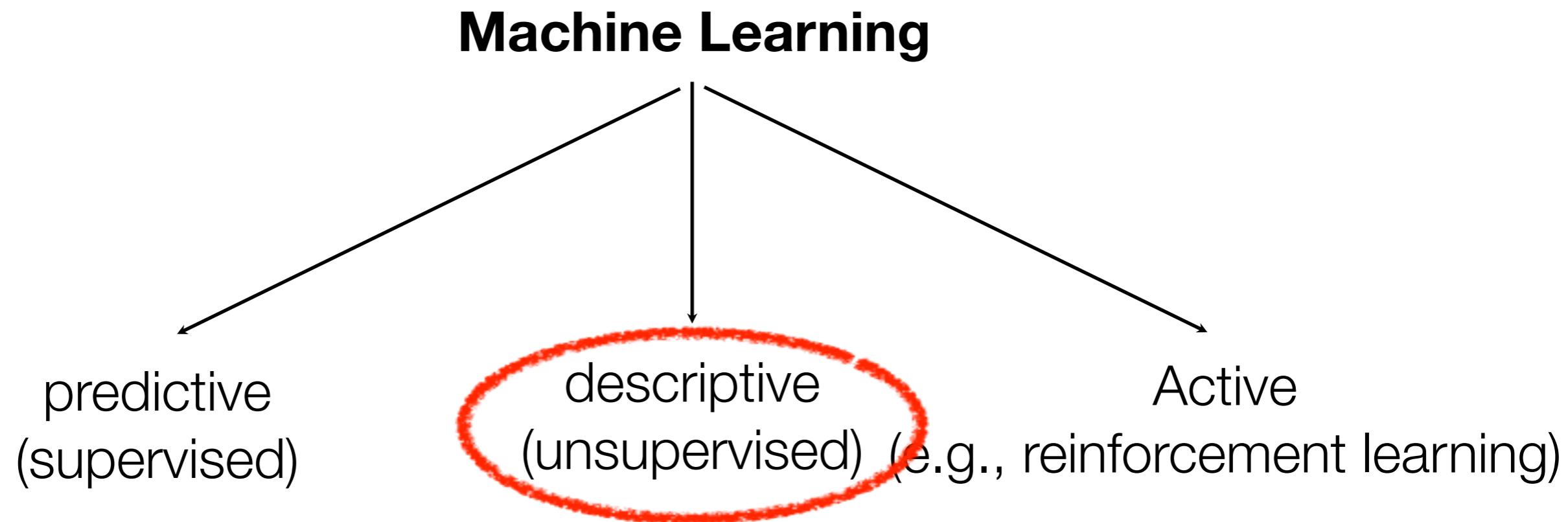


# Examples of Regression

---

- Predict tomorrow's stock market price given current market conditions and other possible side information.
  - Predict the amount of prostate specific antigen (PSA) in the body as a function of a number of different clinical measurements.
  - Predict the temperature at any location inside a building using weather data, time, door sensors, etc.
  - Predict the age of a viewer watching a given video on YouTube.
- Many problems in robotics can be addressed by regression!

# Types of Machine Learning



# Formalization of Descriptive Problems

---



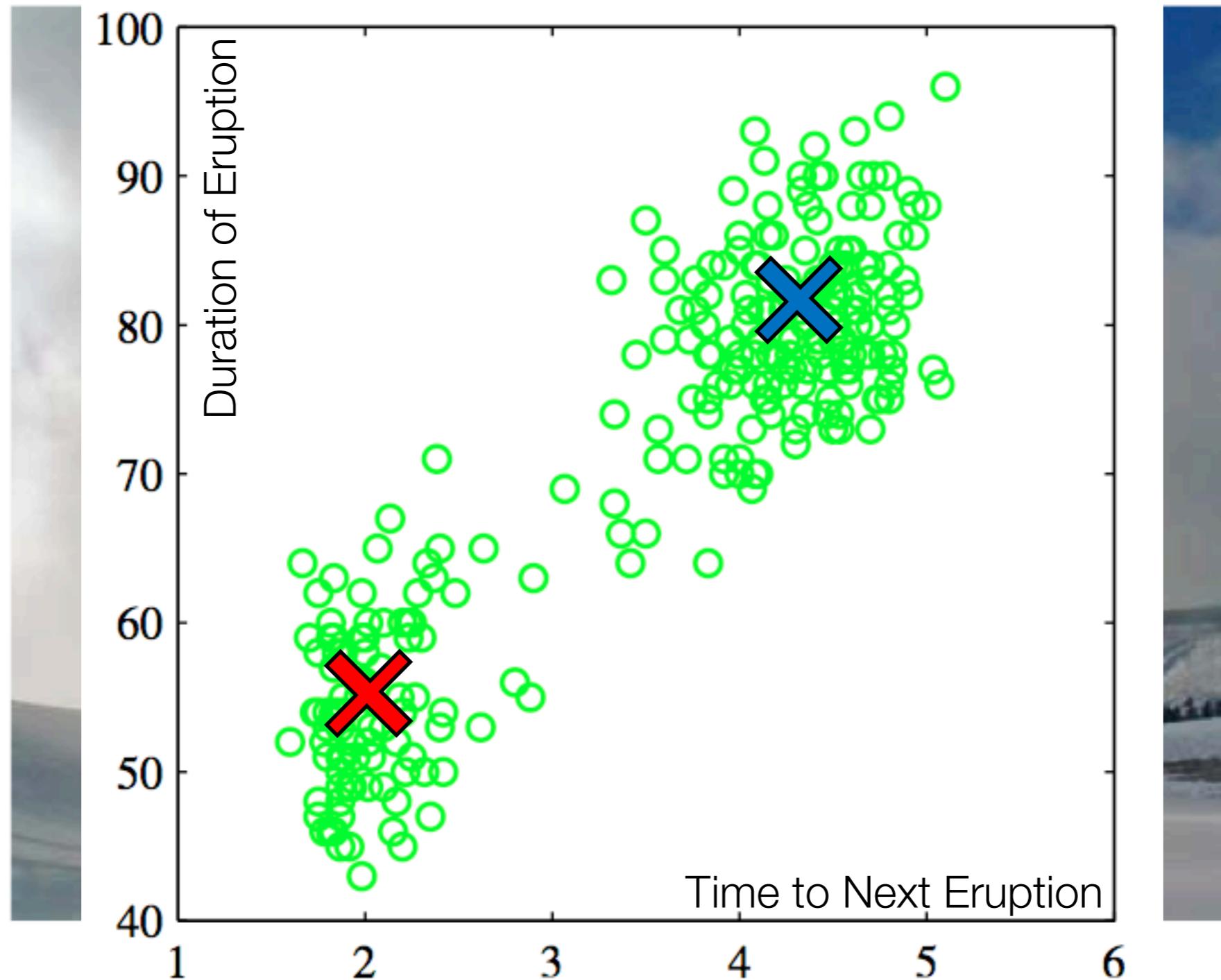
**In descriptive problems, we have**

$$\mathcal{D} = \{ \mathbf{x}_i \mid i = 1, 2, 3, \dots, n \}$$

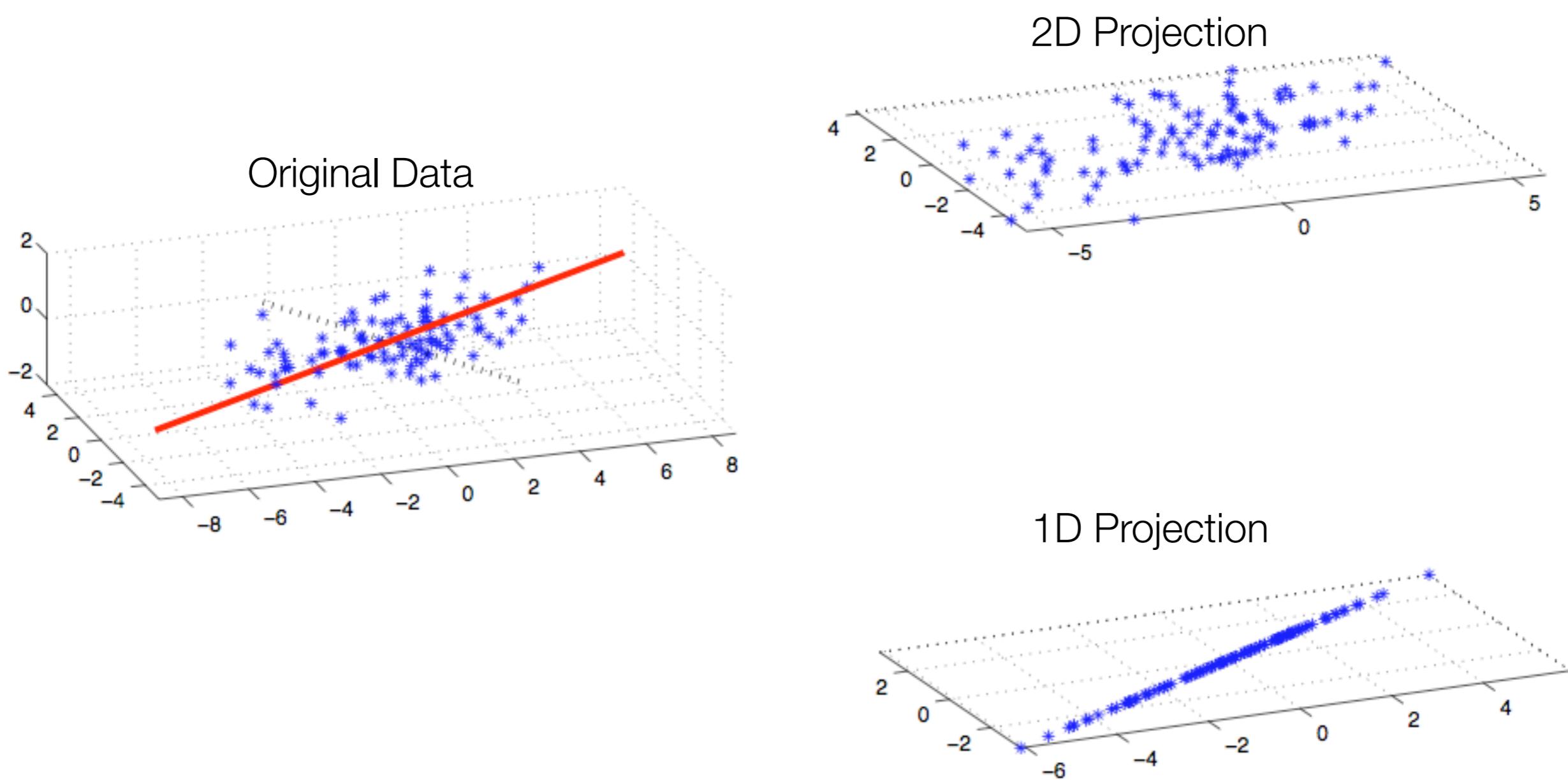
**Three prominent examples are:**

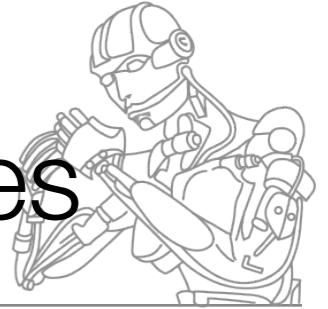
1. *Clustering*: Find groups of data which belong together.
2. *Dimensionality Reduction*: Find the latent dimension of your data.
3. *Density Estimation*: Find the probability of your data...

# Old Faithful



# Dimensionality Reduction





# Dimensionality Reduction Example: Eigenfaces



How many faces do you need to characterize these?

# Example: Eigenfaces



mean



principal basis 1



principal basis 2



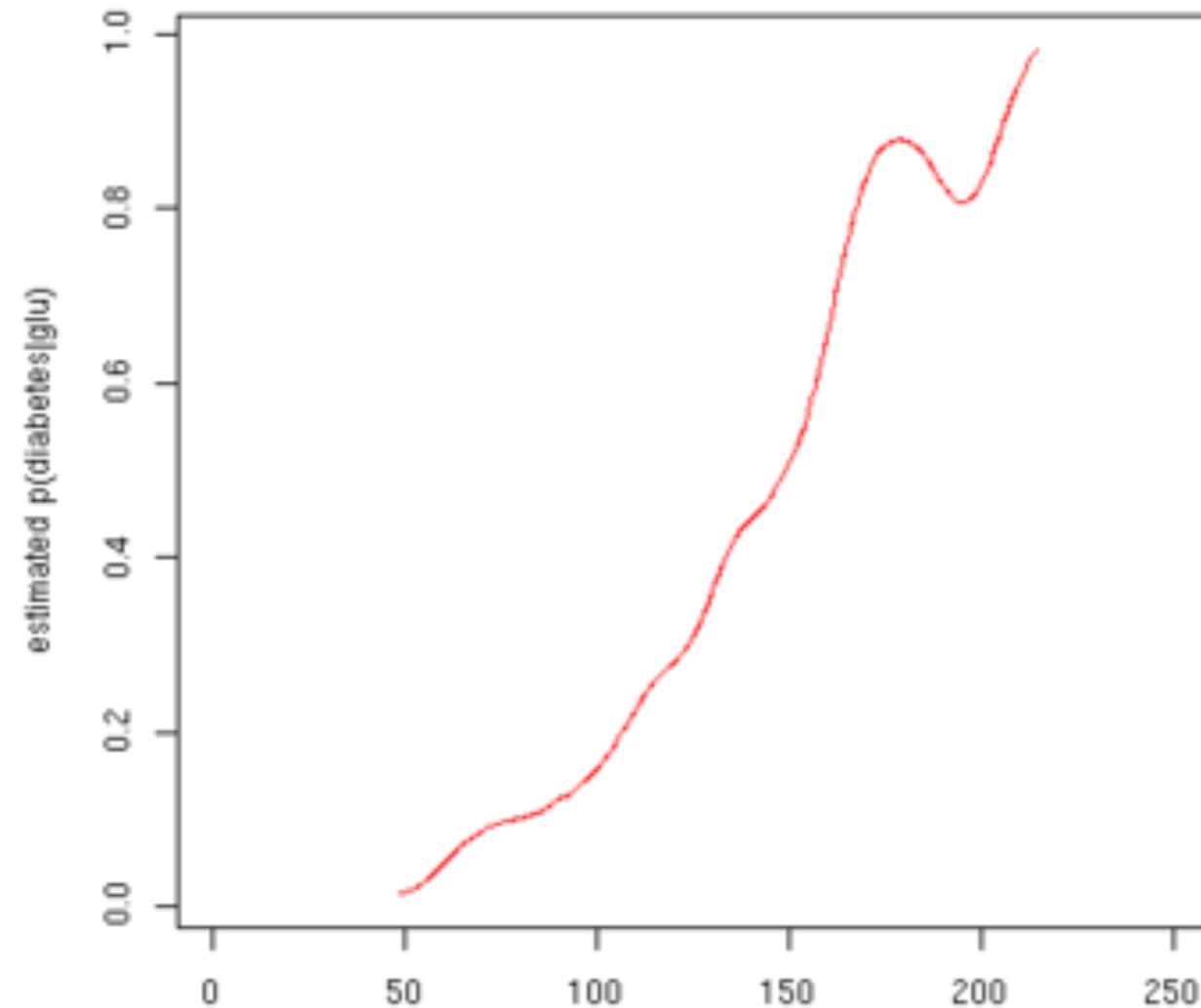
principal basis 3



Example: density of glu (plasma glucose concentration) for diabetes patients



Estimate relative occurrence of a data point



This is called Density Estimation!

# The bigger picture...

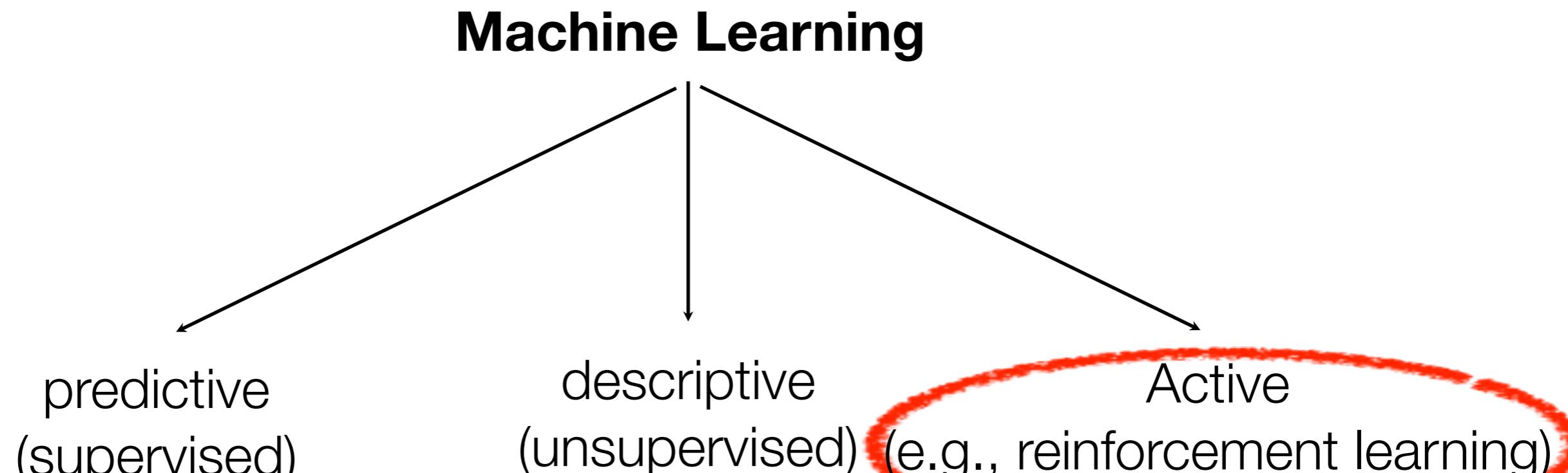
---



When we're learning to see, nobody's telling us what the right answers are – we just look. Every so often, your mother says 'that's a dog,' but that's very little information. You'd be lucky if you got a few bits of information — even one bit per second — that way. The brain's visual system has  $10^{14}$  neural connections. And you only live for  $10^9$  seconds. So it's no use learning one bit per second. You need more like  $10^5$  bits per second. And there's only one place you can get that much information: from the input itself.

— Geoffrey Hinton, 1996

# Types of Machine Learning



That will be the main topic of the lecture!

# How to attack a machine learning problem?

---



Machine learning problems essentially always are about two entities:

*(i) data model assumptions:*

- Understand your problem
- generate good features which make the problem easier
- determine the model class
- Pre-processing your data

*(ii) algorithms that can deal with (i):*

- Estimating the parameters of your model.

We are gonna do this for regression...



# Content of this Lecture

---

- Math and Statistics Refresher
- What is Machine Learning?
- **Model-Selection**
- Linear Regression
  - Gauss' Approach
  - Frequentist Approach
  - Bayesian Approach

# Important Questions

---



## How does the data look like?

- Are you really learning a function?
- What data types do our outputs have?
- Outliers: Are there “data points in China”?

## What is our model (relationship between inputs and outputs)?

- Do you have features?
- What type of noise / What distribution models our outputs?
- Number of parameters?
- Is your model sufficiently rich?
- Is it robust to overfitting?

# Important Questions

---



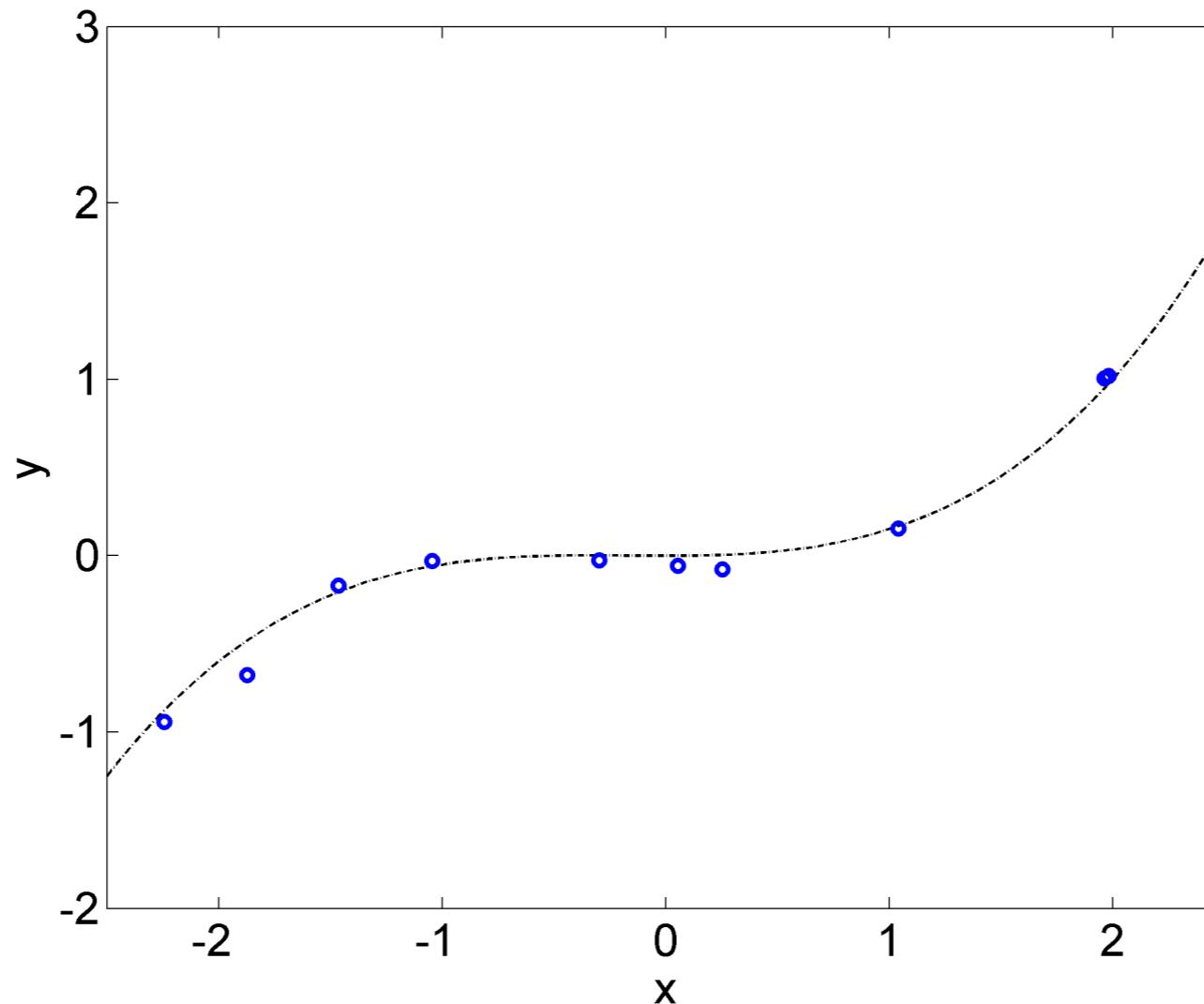
## Requirements for the solution

- accurate
- efficient to obtain (computation/memory)
- interpretable



# Example Problem: a data set

---



Task: Describe the outputs as a function of the inputs (**regression**)



# Model Assumptions: Noise + Features

---

Additive **Gaussian** Noise:

$$y = \mathbf{f}_\theta(\mathbf{x}) + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Equivalent Probabilistic Model

$$p(y|\mathbf{x}) = \mathcal{N}(y|f_\theta(\mathbf{x}), \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(y - f_\theta(\mathbf{x}))^2}{\sigma^2}\right)$$

Lets keep it simple: **linear in Features**

$$\mathbf{f}_\theta(\mathbf{x}) = \phi(\mathbf{x})^T \boldsymbol{\theta}$$



# Important Questions

---

## How does the data look like?

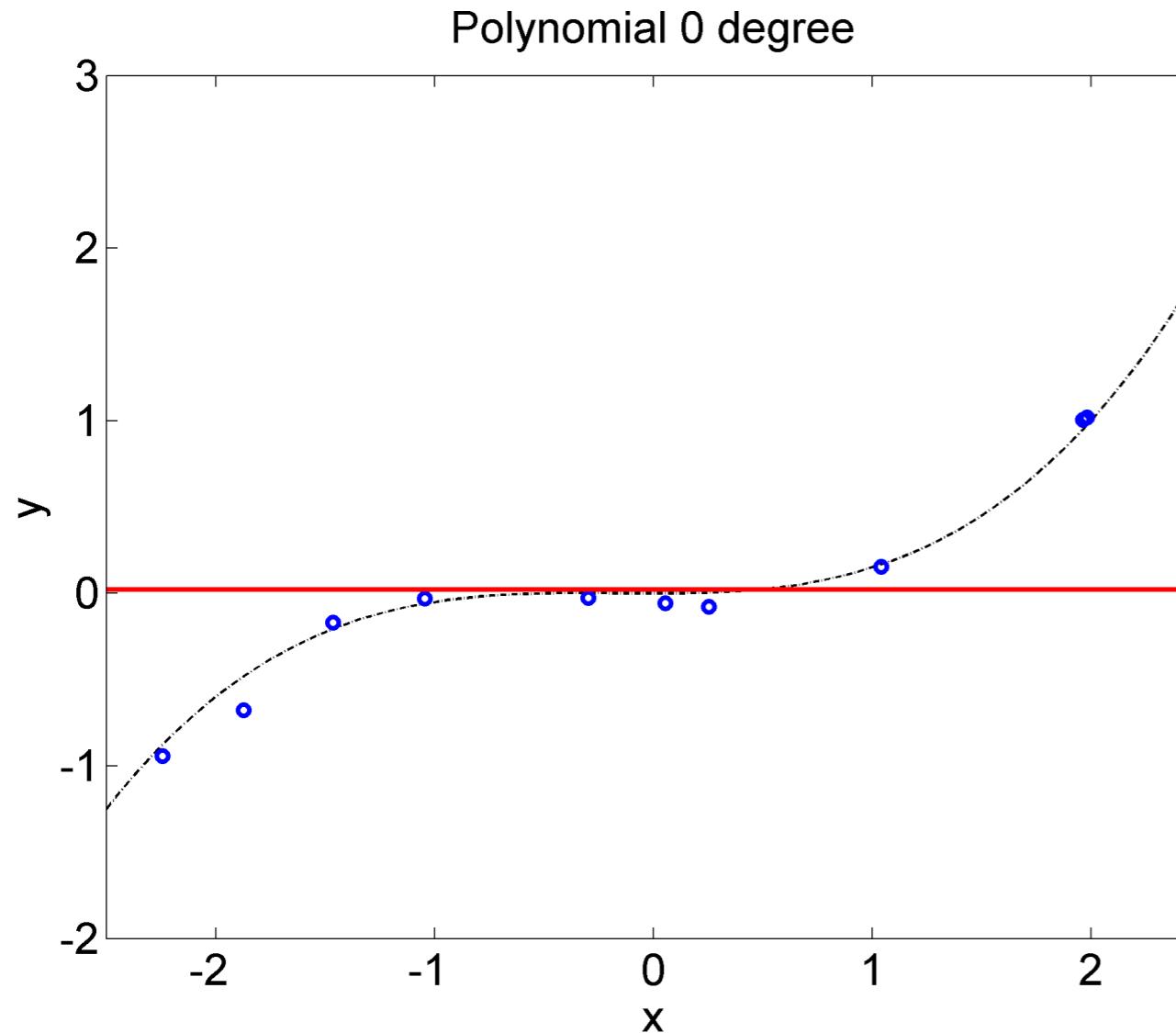
- What data types do our outputs have?  $y_q \in \mathbb{R}$
- Outliers: Are there “data points in China”? NO
- Are you really learning a function? YES

**What is our model?**  $y = \phi(\mathbf{x})^T \boldsymbol{\theta} + \epsilon$

- Do you have features?  $\phi(\mathbf{x})$
- What type of noise / What distribution models our outputs?  $\epsilon \sim \mathcal{N}(0, \sigma^2)$
- Number of parameters?
- Is your model sufficiently rich?
- Is it robust to overfitting?



# Let us fit our model ...



**We need to answer:**

- How many parameters?
- Is your model sufficiently rich?
- Is it robust to overfitting?

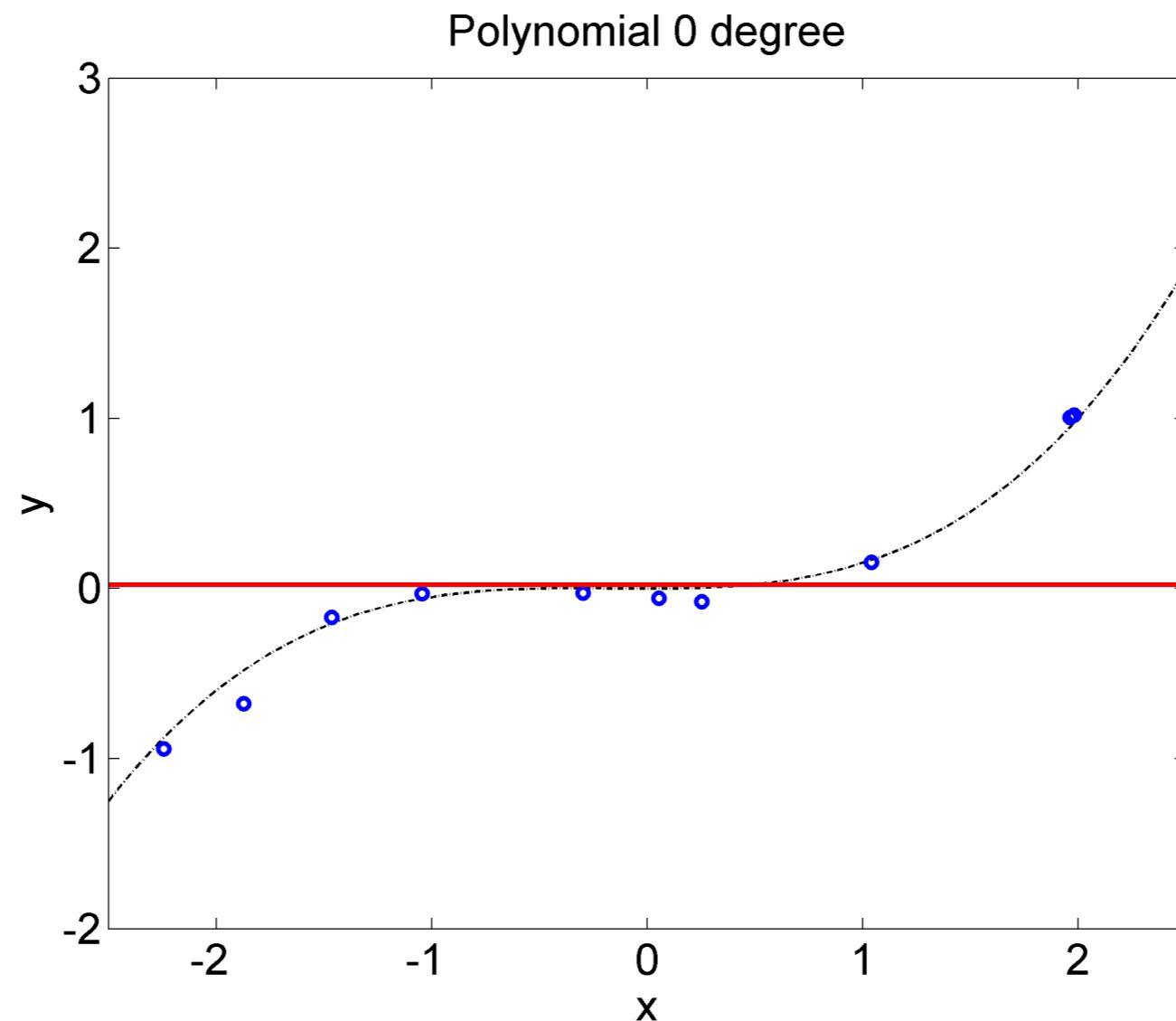
**We assume a model class:** polynomials of degree n

$$y = \phi(x)^T \theta + \epsilon = [1, x, x^2, x^3, \dots, x^n]^T \theta + \epsilon$$



# Fitting an Easy Model: n=0

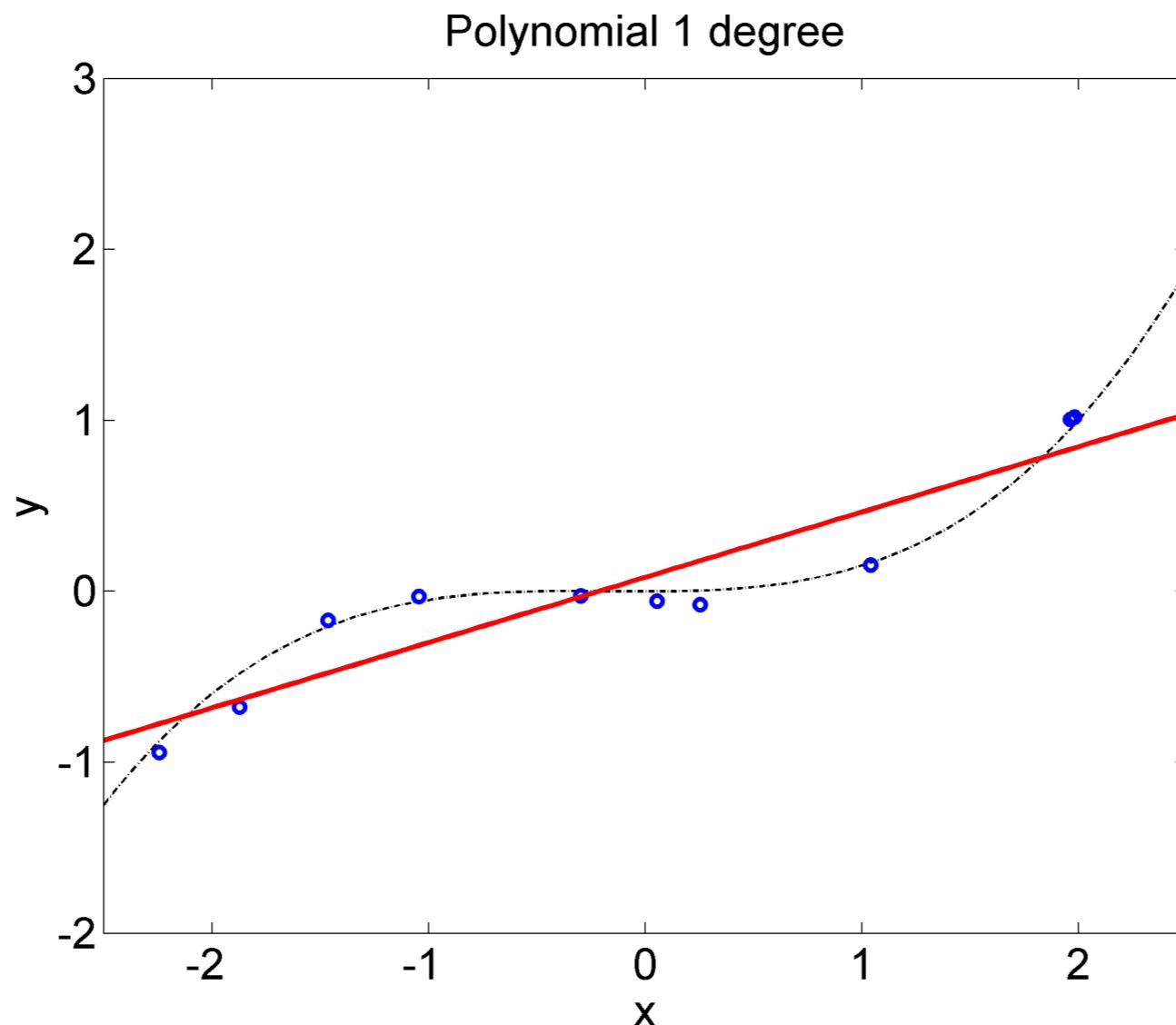
---





# Add a Feature: n=1

---

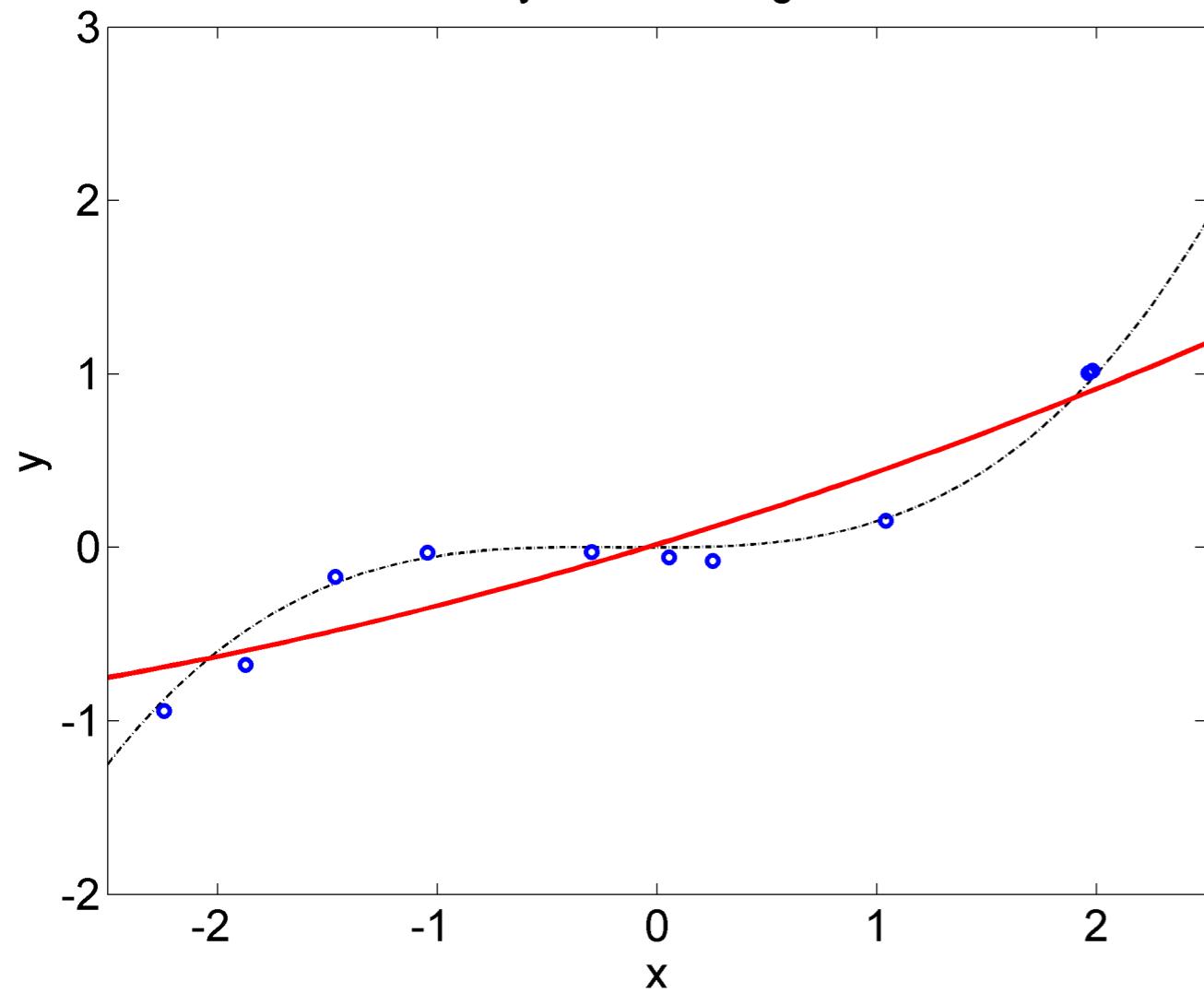


# More features...

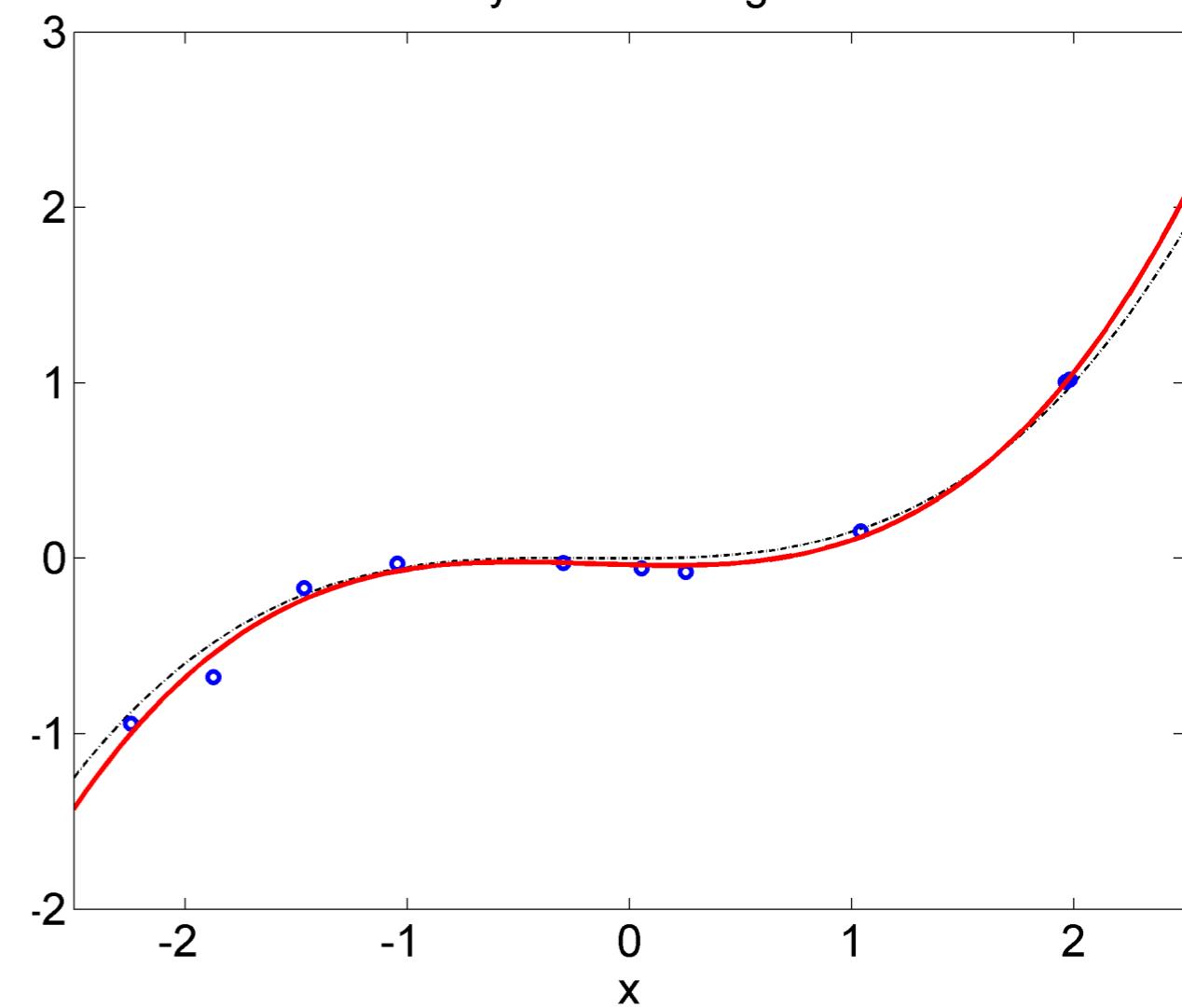
---



Polynomial 2 degree

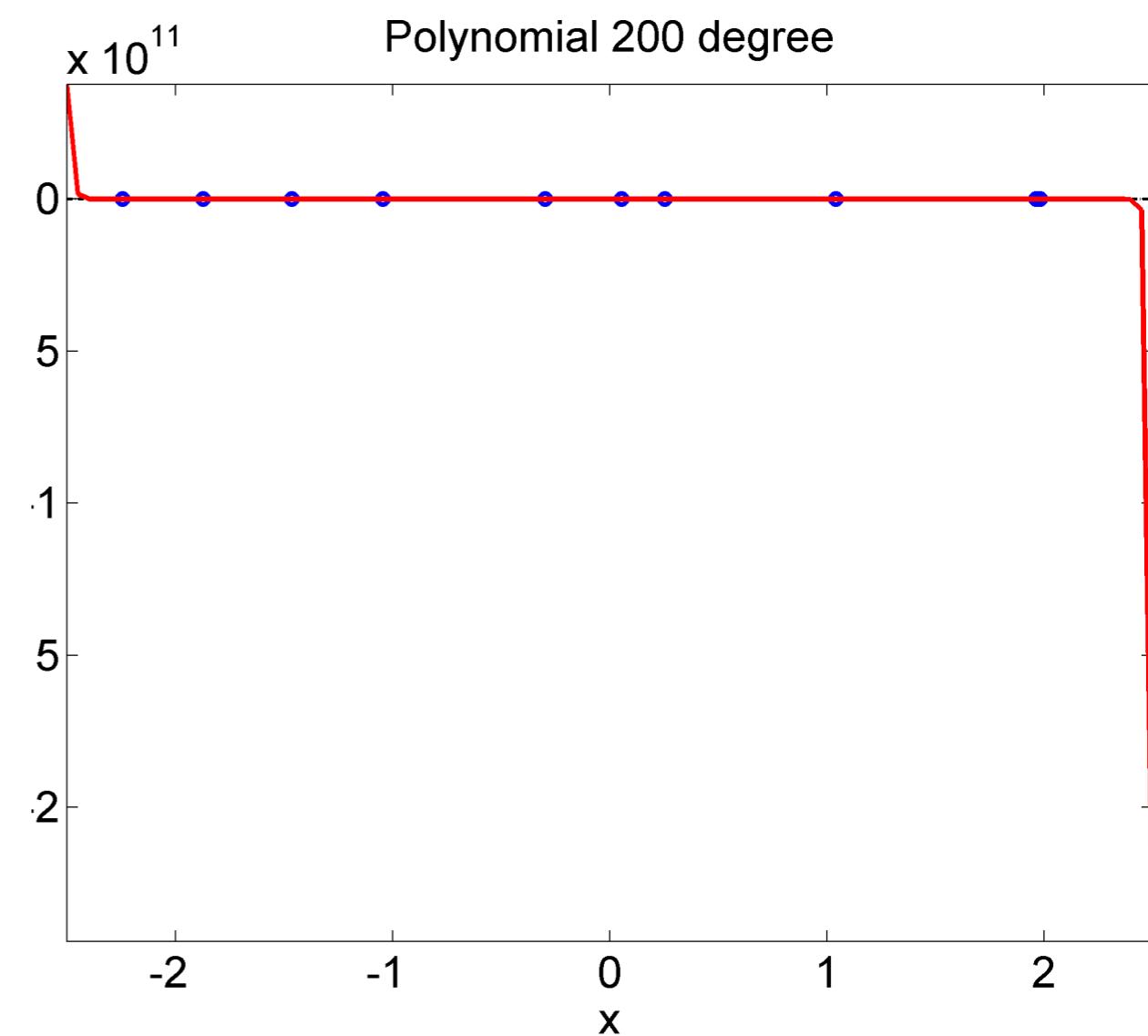
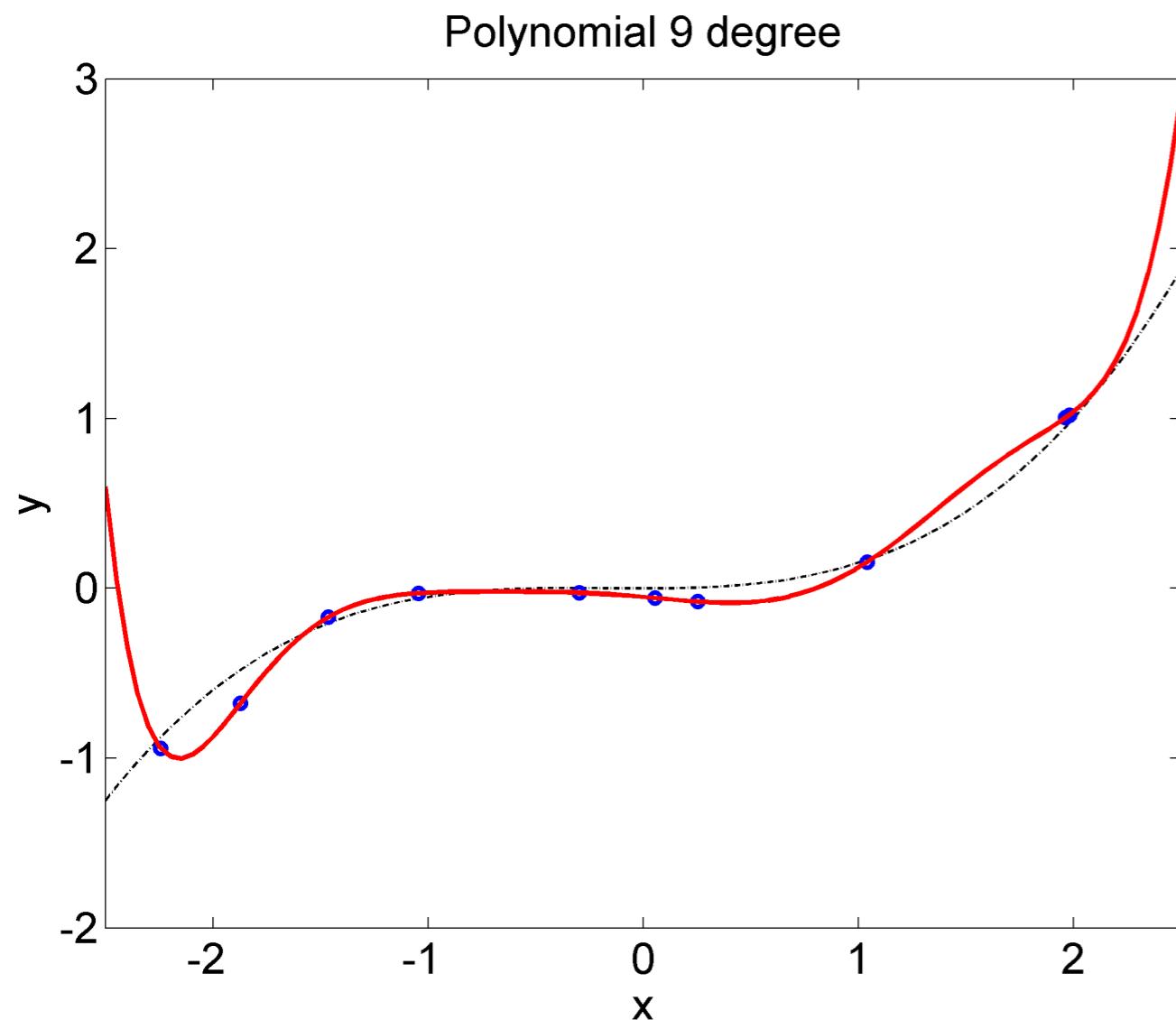


Polynomial 3 degree



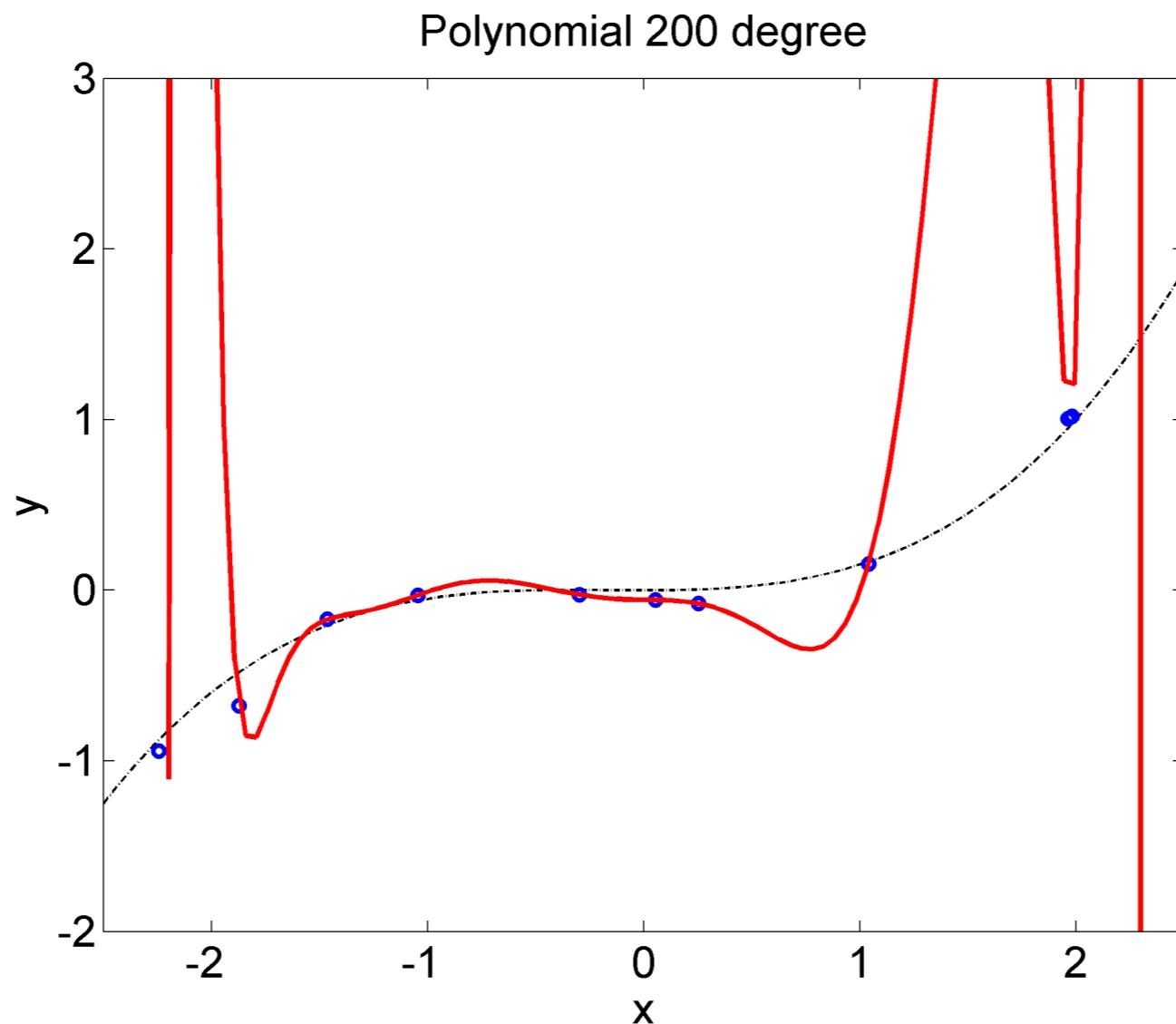
# More features...

---





# More features: n=200 (zoomed in)



overfitting and numerical problems



Prominent example of overfitting...

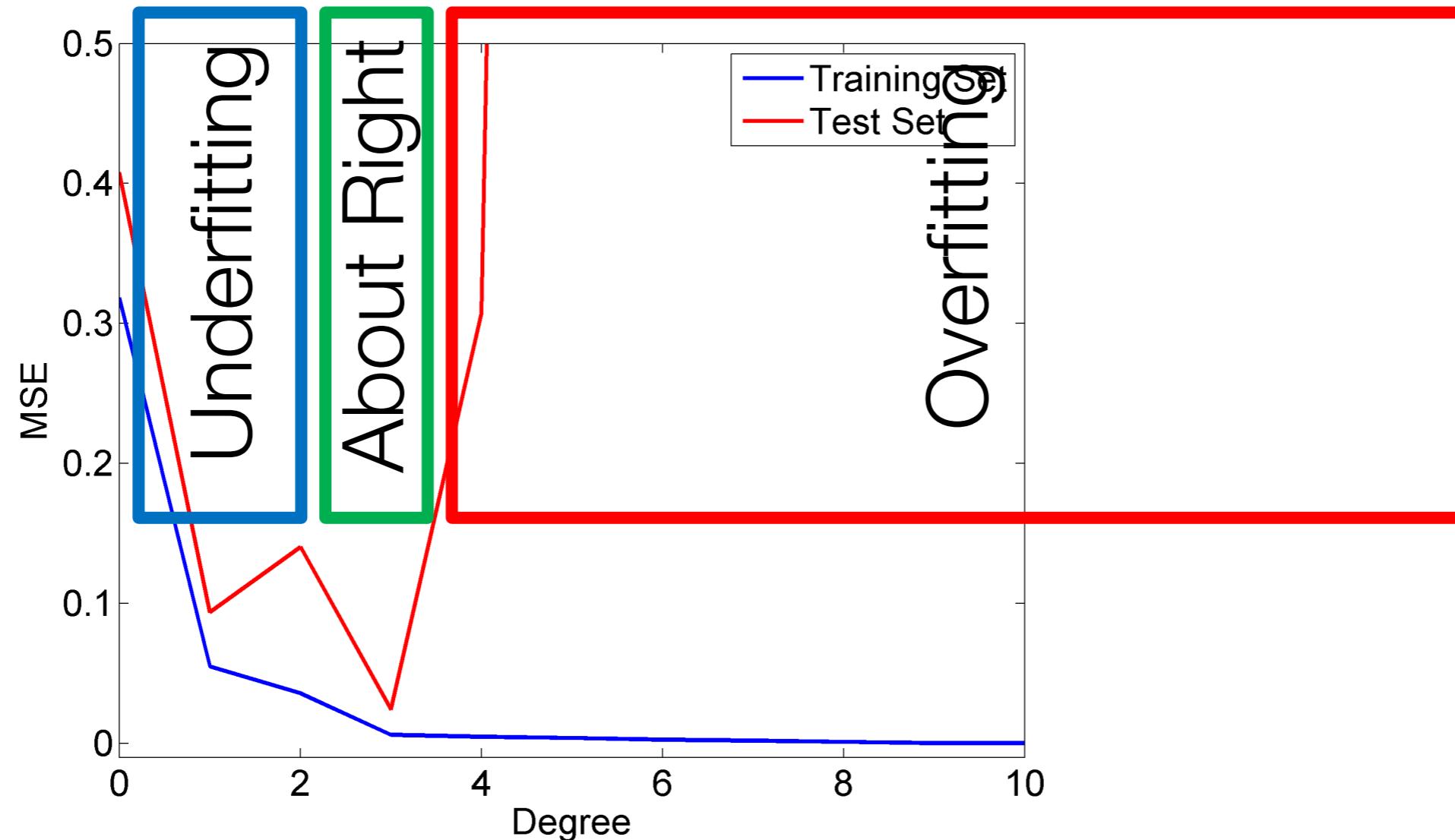


Is there a tank in the picture?

DARPA Neural Network Study (1988-89), AFCEA International Press



# Test Error vs Training Error



Does a small training error lead to a good model ??

**NO ! We need to do model selection**

# Occam's Razor and Model Selection



## Model Selection: How can we...

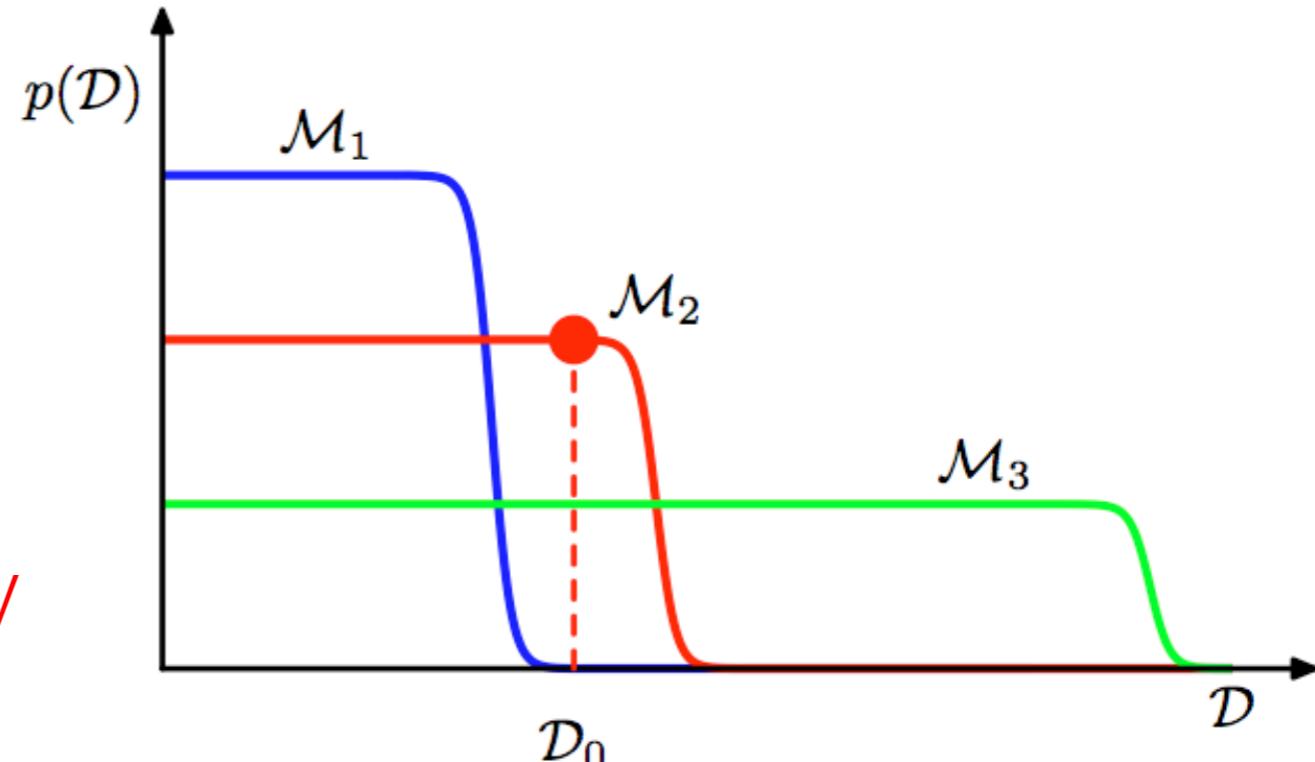
- choose number of features/parameters?
- choose type of features?
- prevent overfitting?

## Some insights:

Always choose the model  
that fits the data and has  
the **smallest model complexity**



called **occam's razor**



# Bias-Variance Tradeoff



$$\text{Expected Total Error} = \text{Bias}^2 + \text{Variance}$$

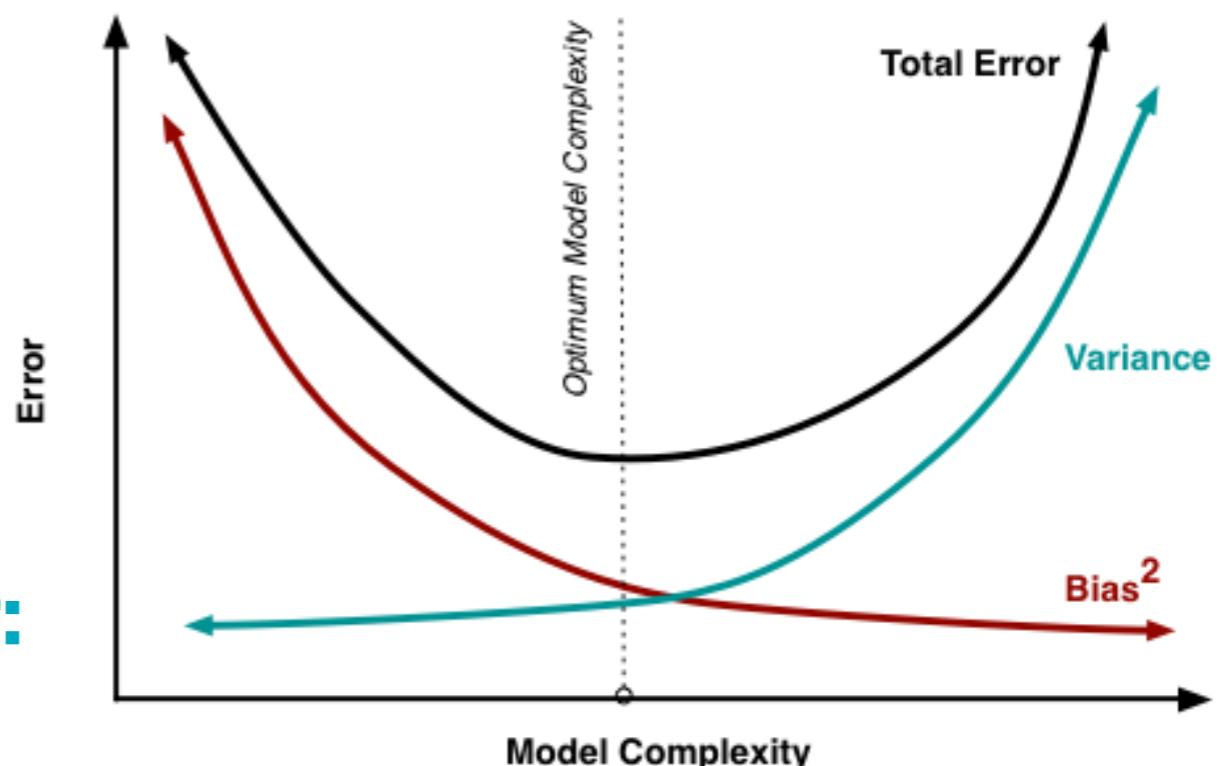
**Typically, you can not minimize both!**

## Bias / Structure Error:

Error because our model can not do better

## Variance / Approximation Error:

Error because we estimate parameters on a **limited data set**

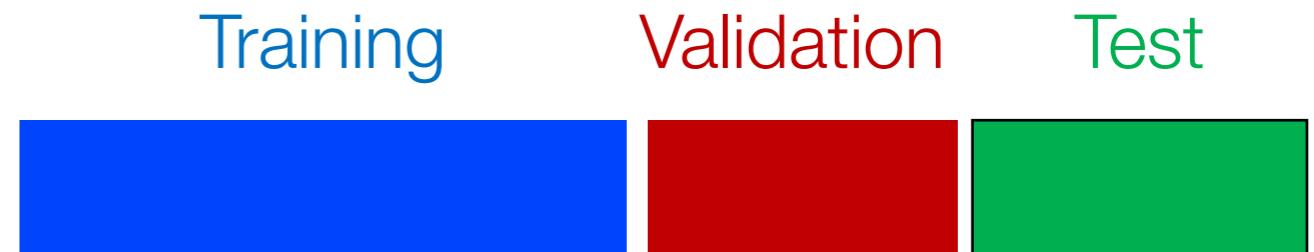




# How do choose the model?

**Goal:** Find a good model  $\mathcal{M}$  (e.g., good set of features)

**Split the dataset into:**



1. **Training Set:** Fit Parameters
2. **Validation Set:** Choose model class or single parameters
3. **Test Set:** Estimate prediction error of trained model

→ **Error needs to be estimated on independent set!**

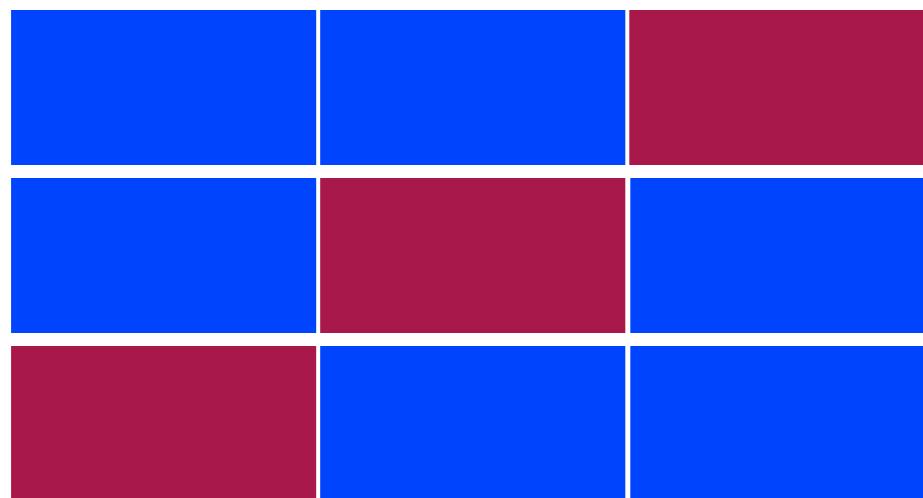
# Model Selection: K-fold Cross Validation



Partition data into K sets, use K-1 as **training set** and 1 set as **validation set**



For all possible ways of partitioning compute the **validation error**  $J$     ➔    **computationally expensive!!**



$$J(p_1, \mathcal{M}_i)$$

$$J(p_2, \mathcal{M}_i)$$

$$J(p_3, \mathcal{M}_i)$$

Choose model  $\mathcal{M}_i$  with **smallest average validation error**

# Content of this Lecture

---



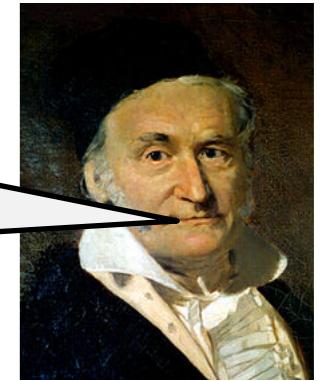
- Math and Statistics Refresher
- What is Machine Learning?
- Model-Selection
- Linear Regression
- Gauss' Approach
- Frequentist Approach
- Bayesian Approach

# How to find the parameters $\theta$ ?



Gauss

Let's find parameters through a cost function!



$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} \sum_{i=1}^N (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i)^2$$

Objective is defined by minimizing a **certain cost function**



# Gauss' view: Least Squares

The classical cost function is the one of **least-squares**

$$J = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{f}_\theta(\mathbf{x}_i))^2. \quad y_i = \phi(\mathbf{x}_i)^T \boldsymbol{\theta} + \epsilon$$

Using

$$\Phi = [\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \phi(\mathbf{x}_3), \dots, \phi(\mathbf{x}_n)]^T,$$
$$\mathbf{Y} = [y_1, y_2, y_3, \dots, y_n]^T.$$

we can rewrite it as

$$J = \frac{1}{2} (\mathbf{Y} - \Phi \boldsymbol{\theta})^T (\mathbf{Y} - \Phi \boldsymbol{\theta}) \rightarrow \text{Scalar Product}$$

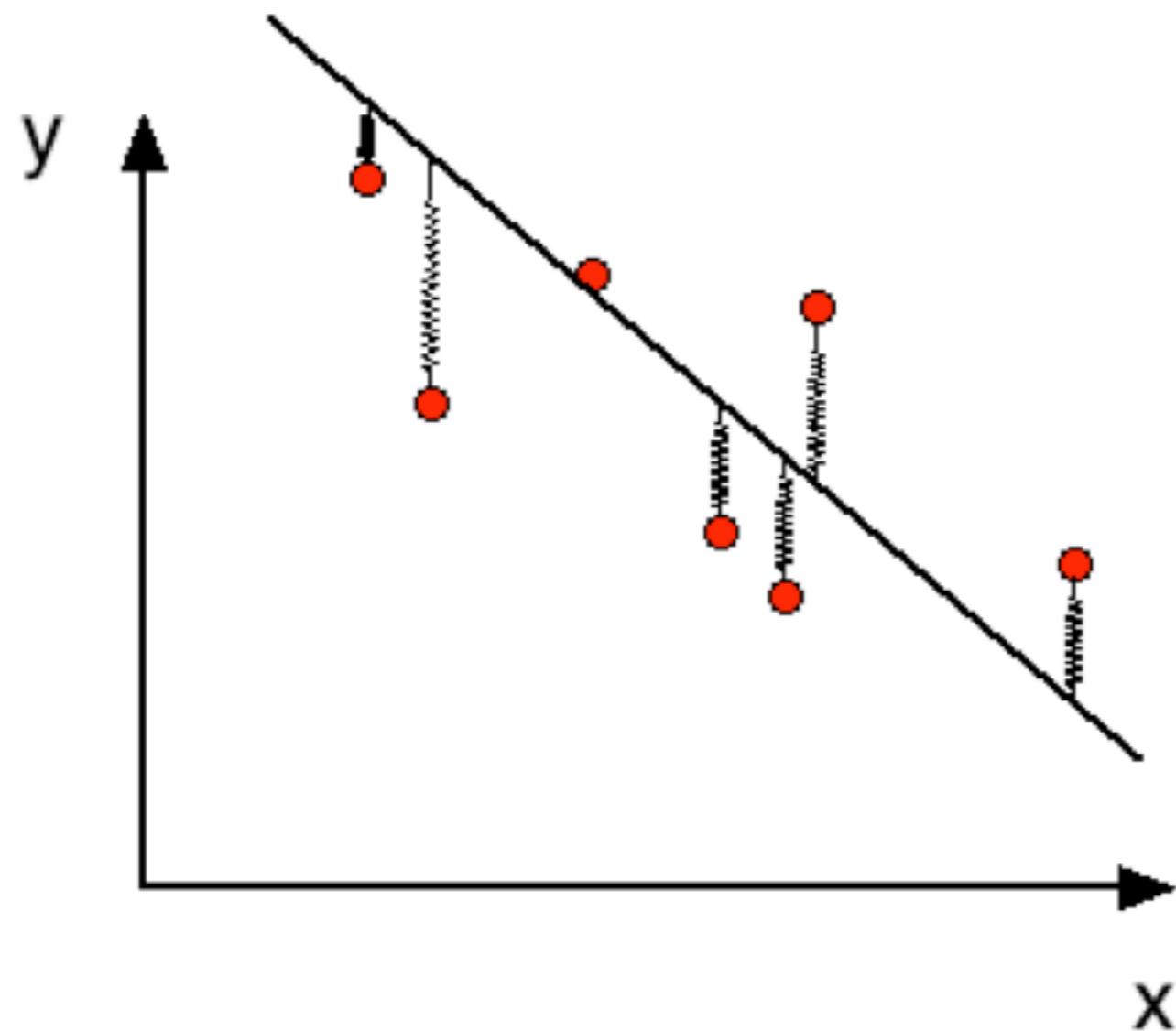
and solve it

$$\boldsymbol{\theta} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{Y}$$

Least Squares solution contains left pseudo-inverse



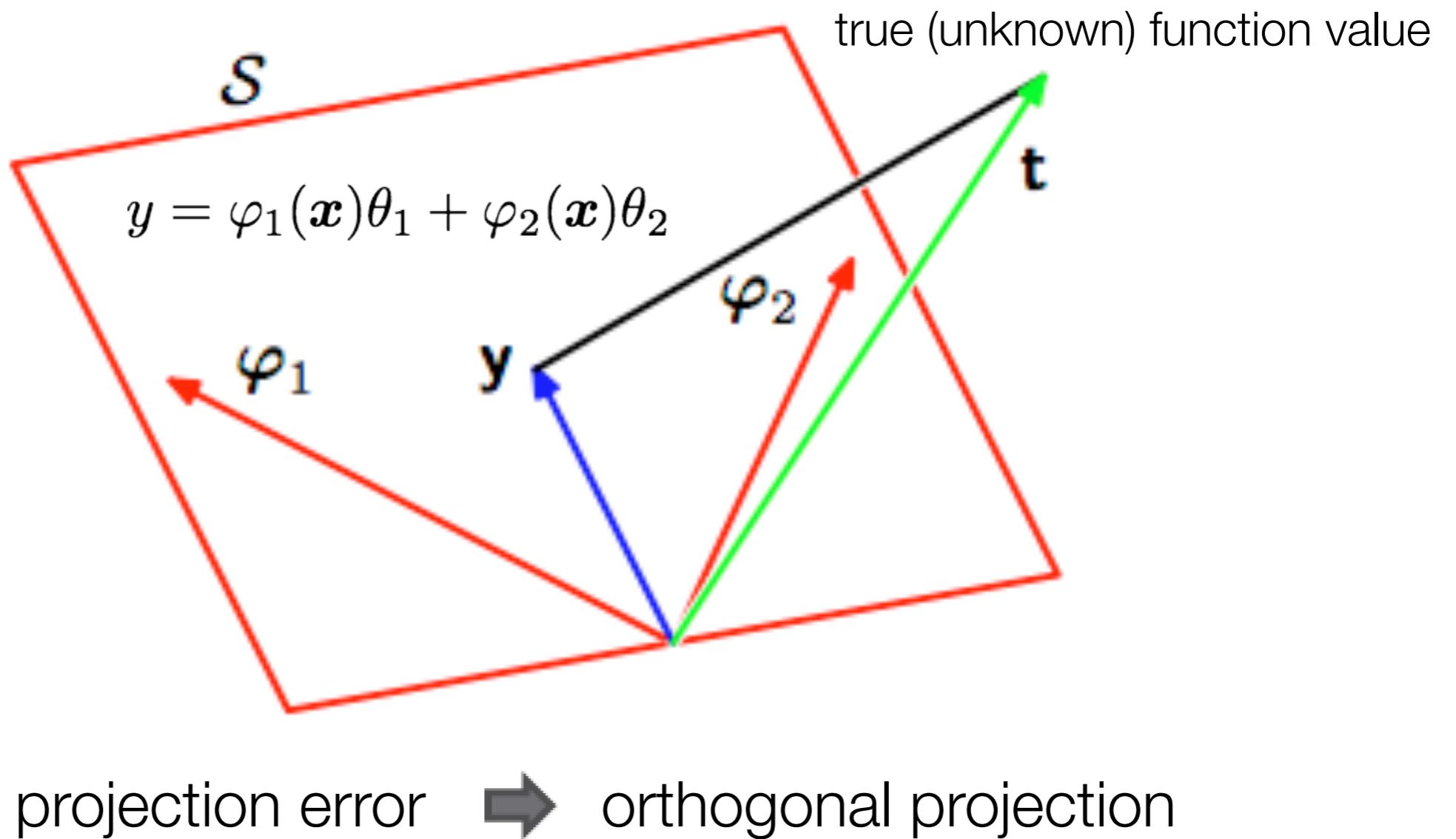
# Physical Interpretation



Energy of springs  $\sim$  squared lengths  
→ minimize energy of system



# Geometric Interpretation



# Robotics Example: Rigid-Body Dynamics

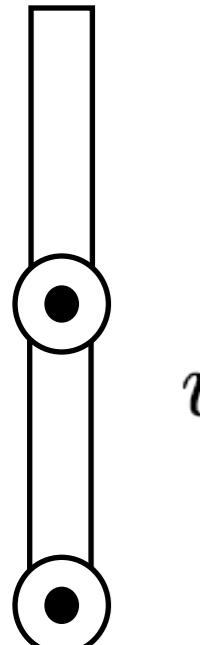
## Known Features



$$u_1 = [m_1 l_{g1}^2 + J_1 + m_2(l_1^2 + l_{g2}^2 + 2l_1 l_{g2} \cos \theta_2) + J_2] \ddot{\theta}_1$$

$$+ [m_2(l_{g2}^2 + l_1 l_2 \cos \theta_2) + J_2] \ddot{\theta}_2 \quad \text{Inertial Forces}$$

$$- 2m_2 l_1 l_{g2} \dot{\theta}_1 \dot{\theta}_2 \sin \theta_2 \quad \text{Coriolis Forces}$$



$$u_2 = [m_2(l_{g2}^2 + l_1 l_{g2} \cos \theta_2) + J_2] \ddot{\theta}_1$$

$$+ (m_2 l_{g2}^2 + J_2) \ddot{\theta}_2 \quad \text{Gravity}$$

$$- m_2 l_1 l_{g2} \dot{\theta}_1^2 \sin \theta_2 \quad \text{Inertial Forces}$$

$$+ m_2 g l_{g2} \cos(\theta_1 + \theta_2) \quad \text{Centripetal Forces}$$

$$+ m_2 g l_{g2} \cos(\theta_1 + \theta_2) \quad \text{Gravity}$$

# Robotics Example: Rigid-Body Dynamics



We realize that rigid body dynamics is **linear in the parameters**

We can rewrite it as

$$\mathbf{u} = \boldsymbol{\phi}(\theta, \dot{\theta}, \ddot{\theta})^T \boldsymbol{\psi}$$

accelerations, velocities, sin and cos terms

masses, lengths, inertia, ...

For finding the parameters we can apply even the first machine learning method that comes to mind: Least-Squares Regression

# Cost Function II: Ridge Regression

---



We punish the magnitude of the parameters

→ Controls model complexity

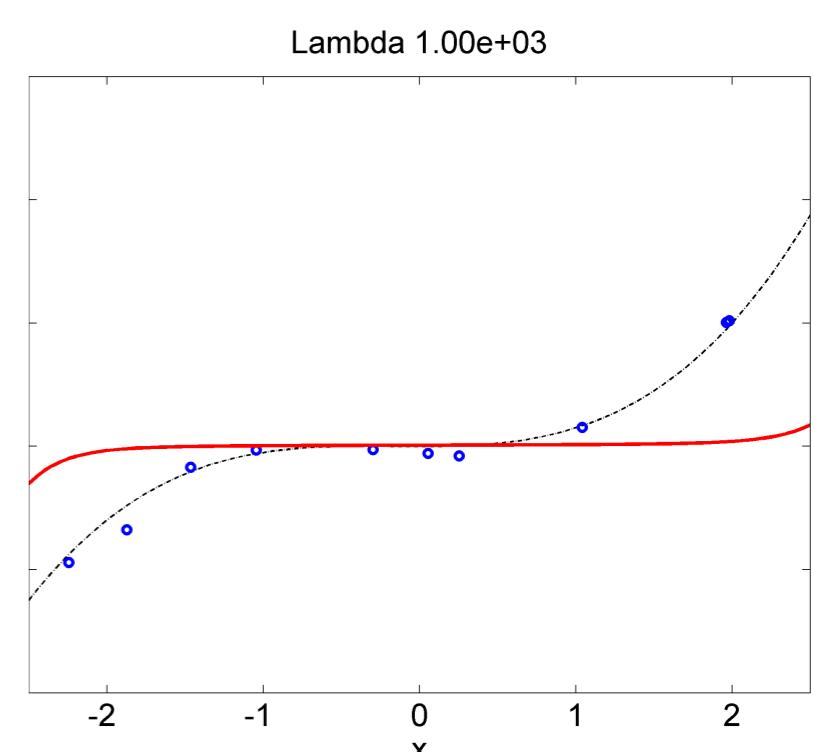
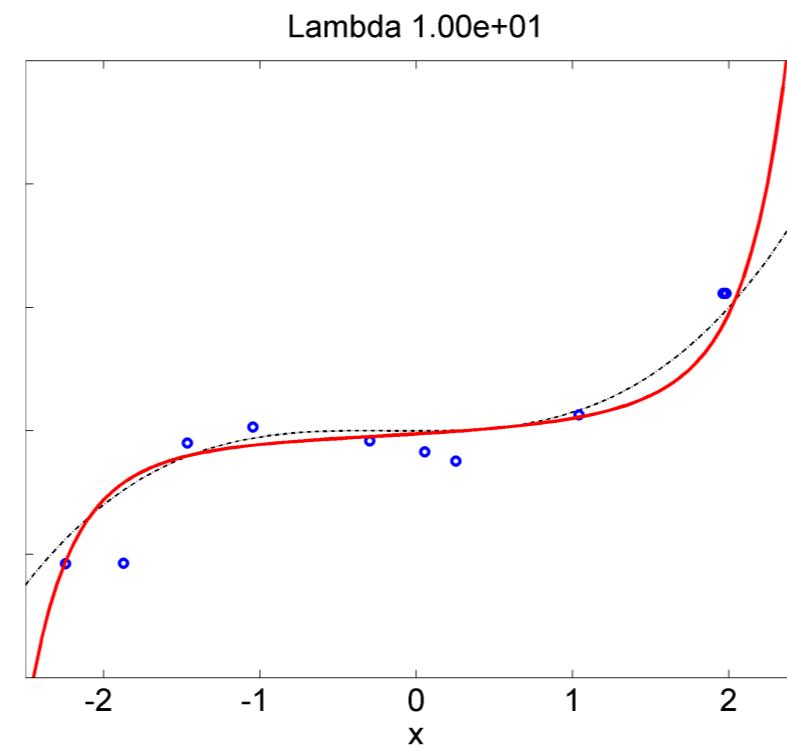
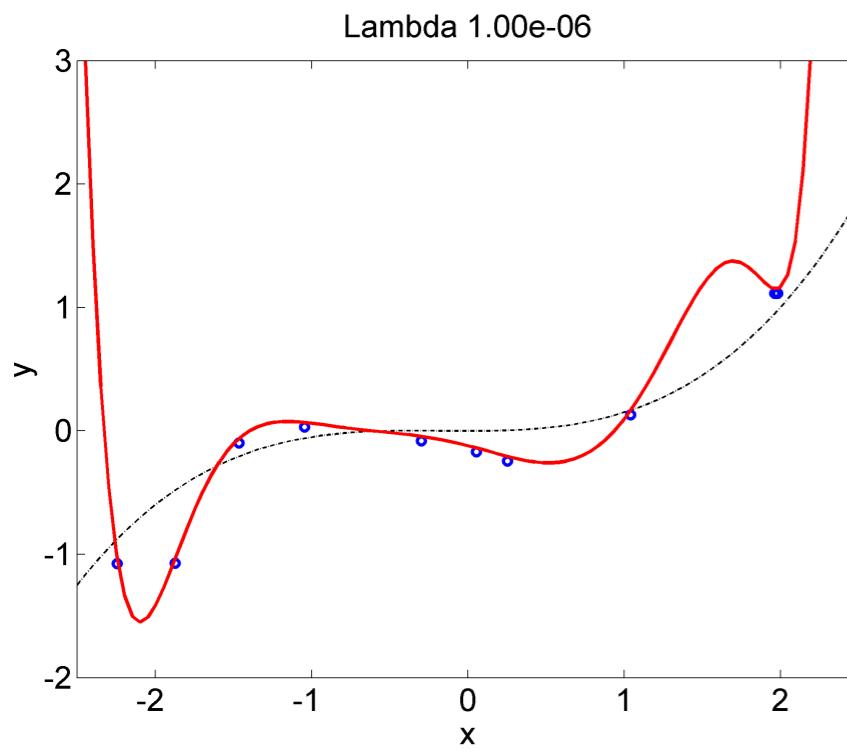
$$J_{\text{RR}} = (\mathbf{y} - \Phi\boldsymbol{\theta})^T(\mathbf{y} - \Phi\boldsymbol{\theta}) + \boldsymbol{\theta}^T \mathbf{W} \boldsymbol{\theta}$$

This yields **ridge regression**  $\boldsymbol{\theta} = (\Phi^T \Phi + \mathbf{W})^{-1} \Phi^T \mathbf{Y}$

with  $\mathbf{W} = \lambda \mathbf{I}$ , where  $\lambda$  is called ridge parameter. For features normalized by variance, typically  $\lambda \in [10^{-9}, \dots, 10^{-5}]$ .

Numerically, this is much **more stable!** Even with **redundant features.**

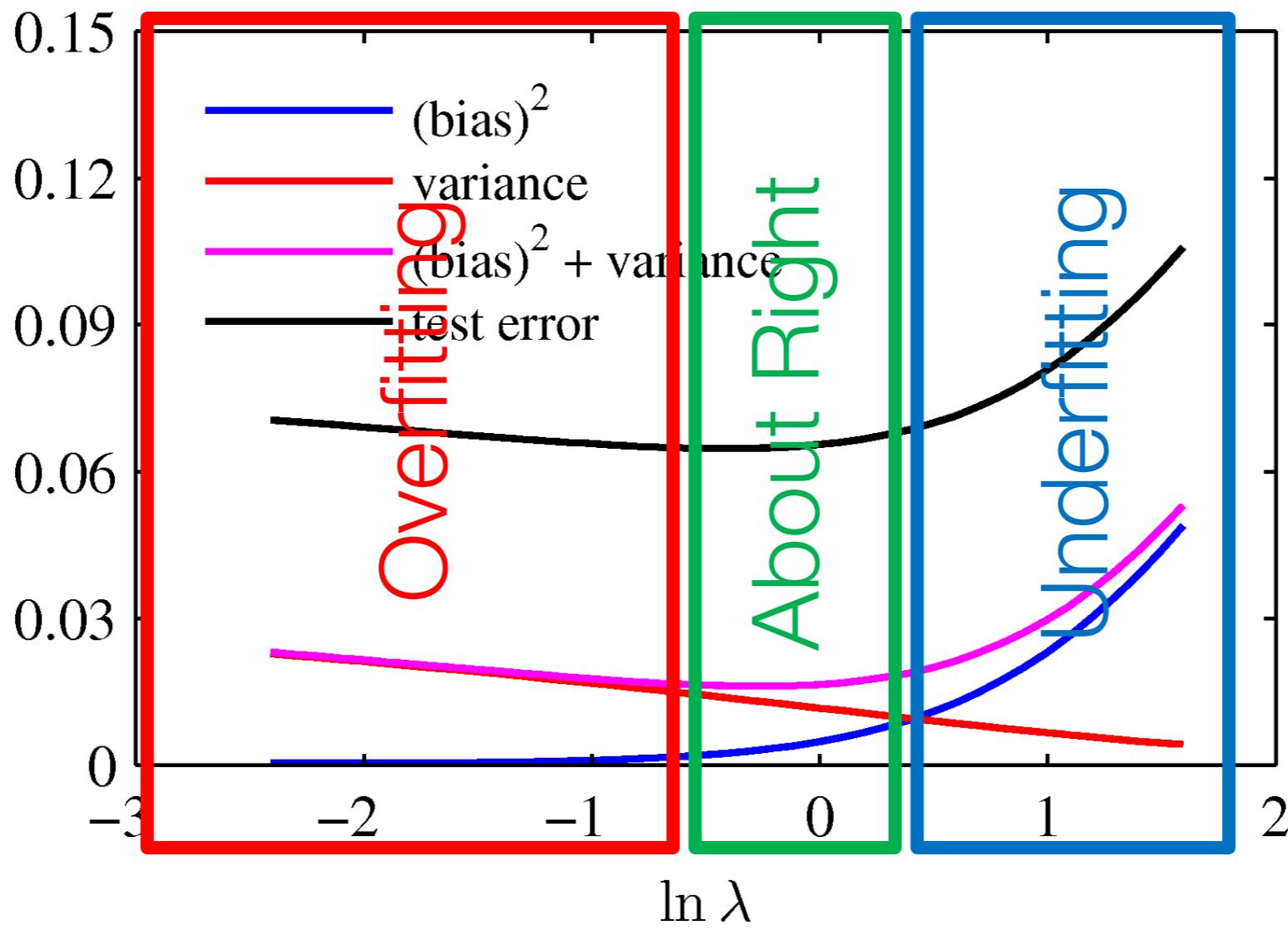
# Ridge regression: n=15



Influence of the regularization constant



# MAP: Back to the Overfitting Problem



We can also scale the model complexity with the regularization parameter!

**Smaller lambda** **higher model complexity**

# Content of this Lecture

---



- Math and Statistics Refresher
- What is Machine Learning?
- Model-Selection
- **Linear Regression**
  - Gauss' Approach
  - **Frequentist approach**
  - Bayesian Approach

# How to find the parameters $\theta$ ?



Frequentist:

Probabilities are frequencies of a repeated experiment.



- There are some true parameters of the experiment which we cannot observe.
- They reveal themselves by the frequency (i.e., likelihood) at which we can repeat the outcome of the experiment  $p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$ .
- **We can obtain good parameters by maximizing likelihood of the outcome!**



# Maximum-Likelihood (ML) estimate

We can maximize the **likelihood of the outcome**:

$$\arg \max_{\theta} p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \arg \max_{\theta} \prod_{i=1}^N p(y_i|\mathbf{x}_i, \boldsymbol{\theta})$$

**That's hard!**

Do the ‘log-trick’:

$$J_{\text{ML}} = \arg \max_{\theta} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \arg \max_{\theta} \sum_{i=1}^N \log p(y_i|\mathbf{x}_i, \boldsymbol{\theta})$$

$$= \arg \min_{\theta} \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_{\theta}(\mathbf{x}_i))^2$$

**That's easy!**

$$\rightarrow \theta_{\text{LS}}^* = \theta_{\text{ML}}^*$$

**Least Squares Solution is equivalent to  
ML solution with Gaussian noise!!**

# Content of this Lecture

---



- Math and Statistics Refresher
- What is Machine Learning?
- Model-Selection
- Linear Regression
  - Gauss' Approach
  - Frequentist Approach
  - **Bayesian Approach**



# Does this make sense?

---

- Maximizing  $p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$  basically means we only care about the accuracy of the reproduction of the outcomes.
- What if there is no fully true  $\boldsymbol{\theta}^*$ ?
- In this case, we rather need to study the probability of different  $\boldsymbol{\theta}$ .
- Thus, our quantity of interest is  $p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y})$ .
- But where how can we obtain this quantity?

# How to find the parameters $\theta$ ?



## Bayesian:

Parameters are just random variables. We can encode our subjective belief in the prior.

Bayes



$$p(\theta|X, y) = \frac{p(y|X, \theta)p(\theta)}{p(X|y)}$$

likelihood                                      prior  
  ↓  
    evidence  
posterior    ↑

**Intuition:** If you assign each parameter estimator a “probability of being right”, the average of these parameter estimators will be better than the single one



# Maximum a posteriori (MAP) estimate

Put a **prior** on our parameters

- E.g.,  $\boldsymbol{\theta}$  should be small:  $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{W}^{-1})$

Find parameters that **maximize the posterior**

$$p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y}|\mathbf{X})} \propto p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})$$

Do the ‘log trick’ again:

$$\begin{aligned} \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}) &= \arg \max_{\boldsymbol{\theta}} p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})p(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \\ &= \arg \min_{\boldsymbol{\theta}} -\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) =: J_{\text{MAP}} \end{aligned}$$



# Maximum a posteriori (MAP) estimate

The prior is just additive costs...

$$J_{\text{MAP}} = -\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) - \log p(\boldsymbol{\theta})$$

Lets put in **our Model:**

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_i p(y_i|\mathbf{x}_i, \boldsymbol{\theta}) = \prod_i \mathcal{N}(y_i|\boldsymbol{\theta}^T \phi(\mathbf{x}_i), \sigma^2) \quad p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{W}^{-1})$$

$$\begin{aligned} J_{\text{MAP}} &= \frac{1}{2} \sum_i \frac{(y_i - \boldsymbol{\theta}^T \phi(\mathbf{x}_i))^2}{\sigma^2} + \frac{1}{2} \boldsymbol{\theta}^T \mathbf{W} \boldsymbol{\theta} \\ &= \frac{1}{2\sigma^2} (\mathbf{y} - \Phi(\mathbf{X})\boldsymbol{\theta})^T (\mathbf{y} - \Phi(\mathbf{X})\boldsymbol{\theta}) + \frac{1}{2} \boldsymbol{\theta}^T \mathbf{W} \boldsymbol{\theta} = \frac{1}{\sigma^2} J_{\text{RR}} \end{aligned}$$



$$\boldsymbol{\theta}_{\text{RR}}^* = \boldsymbol{\theta}_{\text{MAP}}^*$$

**Ridge Regression is equivalent  
to MAP estimate with Gaussian prior**



# Predictions with the Model

---

We found an amazing parameter set  $\theta^*$  (e.g., ML, MAP)

Let's do **predictions!** parameter estimate

$$p(y_* | \mathbf{x}_*, \boldsymbol{\theta}^*) = \mathcal{N}(y_* | \mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

pred. function value                            test input

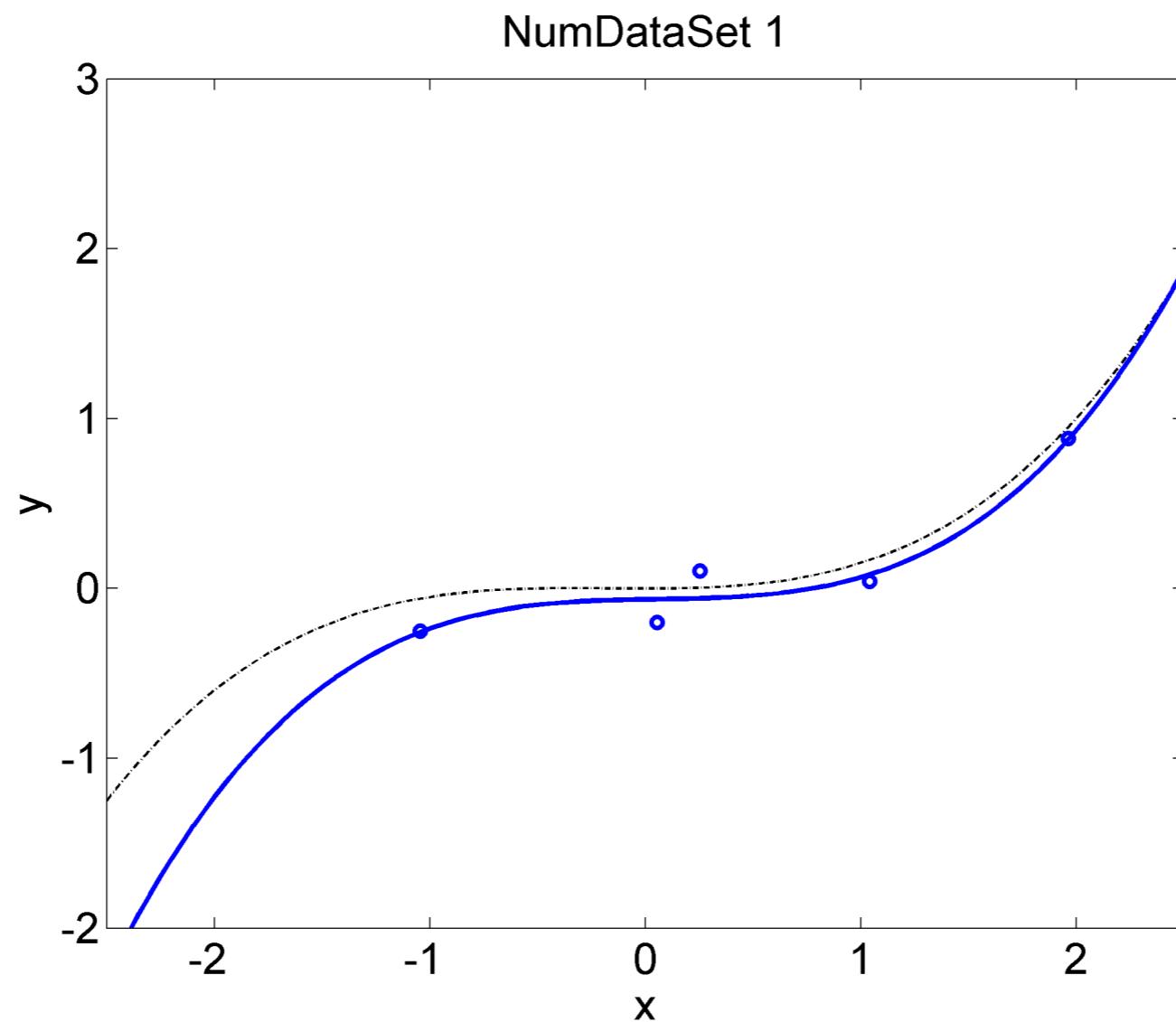
**Predictive mean:**  $\mu(\mathbf{x}_*) = \boldsymbol{\phi}^T(\mathbf{x}_*)\boldsymbol{\theta}^*$

**Predictive variance:**  $\sigma^2(\mathbf{x}_*) = \sigma^2$



# Comparing different data sets...

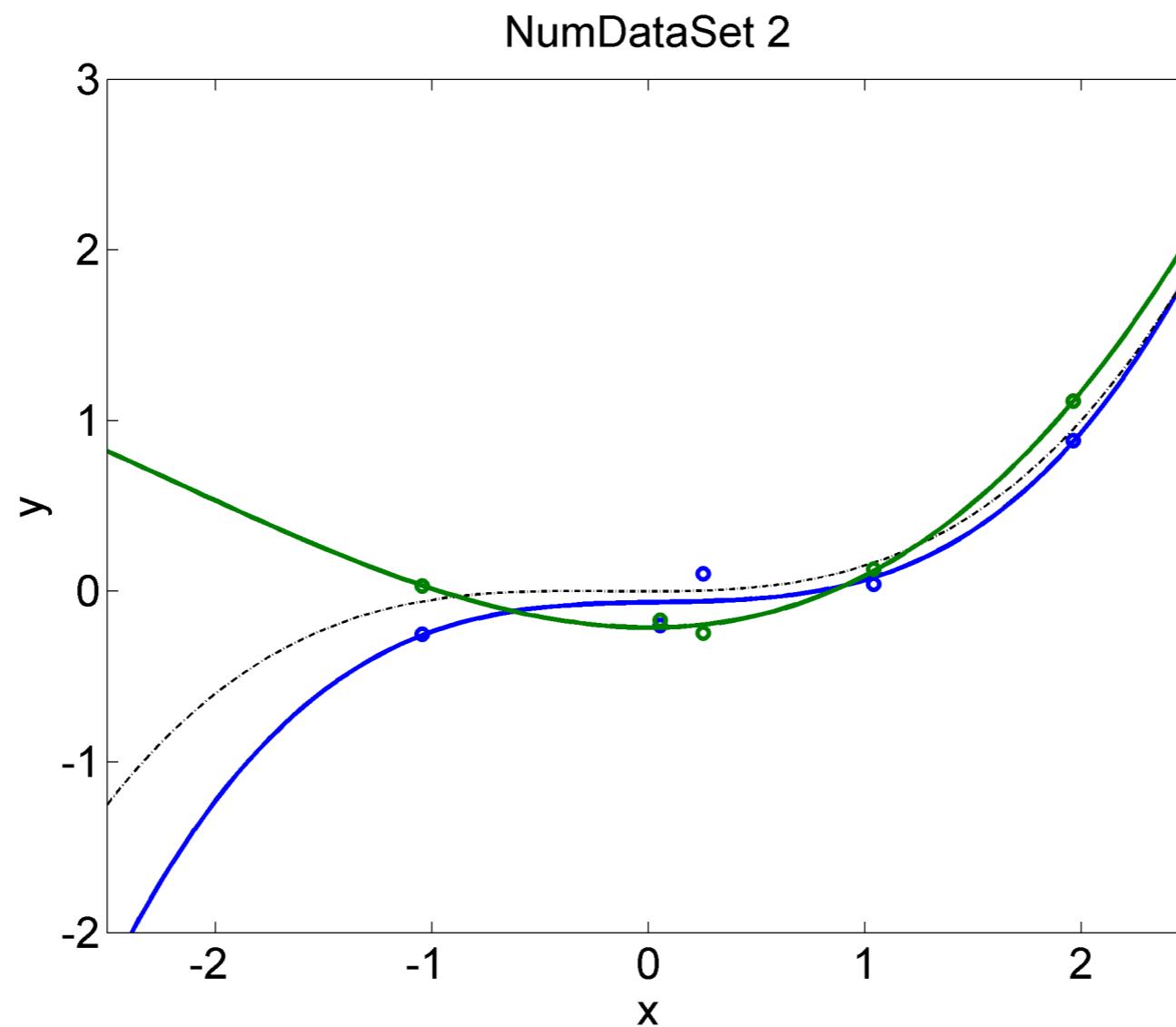
... with same input data, but different output values (due to noise):





# Comparing different data sets...

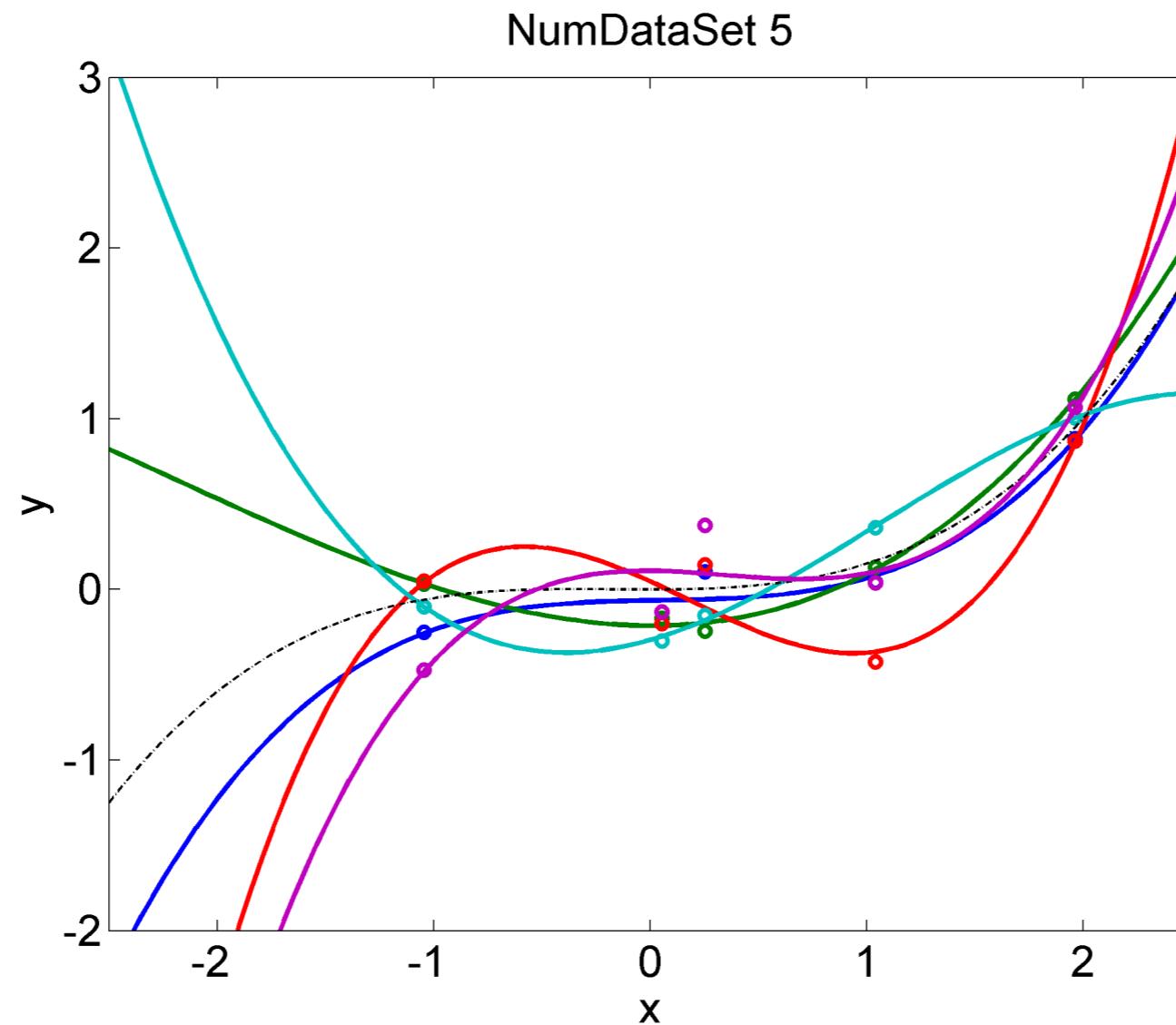
... with same input data, but different output values (due to noise):





# Comparing different data sets...

... with same input data, but different output values (due to noise):



Our parameter estimate is also noisy!

It depends on the noise in the data

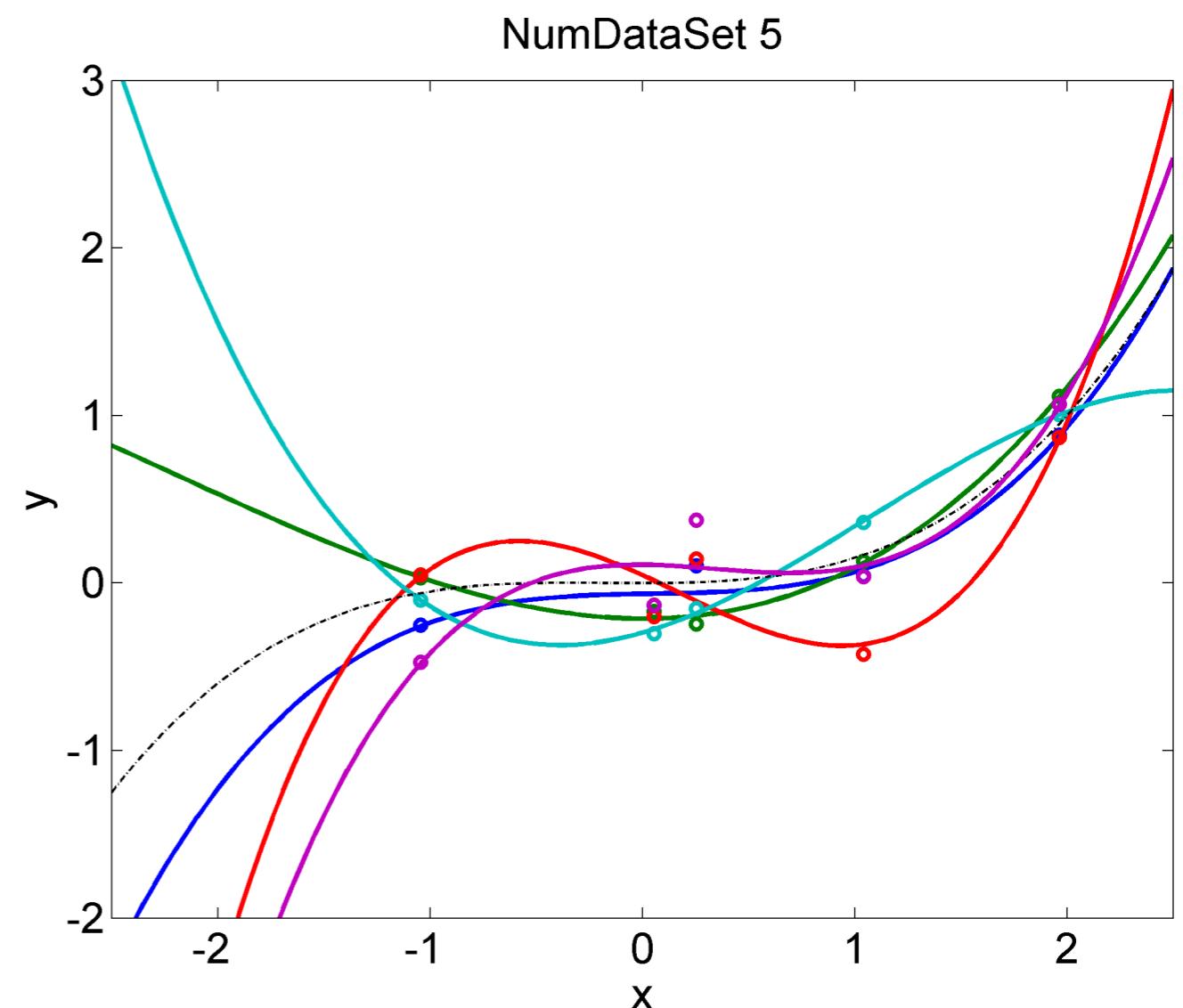


# Comparing different data sets...

Can we also estimate our uncertainty in  $\theta$  ?

Compute probability of  $\theta$  given data

$$\rightarrow p(\theta | X, y)$$





# How to get the posterior?

## Bayes Theorem for Gaussians

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{A}) \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|F\mathbf{x}, \sigma^2 \mathbf{I}) \end{aligned} \longrightarrow \begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\Sigma F^T \mathbf{y}, \sigma^2 \Sigma) \\ \Sigma &= (F^T F + \sigma^2 A^{-1})^{-1} \end{aligned}$$

For our model:

Prior over parameters:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \lambda^{-1} \mathbf{I})$$

Data Likelihood

$$p(\mathbf{y}|X, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}|\Phi\boldsymbol{\theta}, \sigma^2 \mathbf{I})$$

Posterior over parameters:

$$p(\boldsymbol{\theta}|\mathbf{y}, X) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}_N, \Sigma_N)$$

$$\Sigma_N = (\Phi^T \Phi + \sigma^2 \lambda \mathbf{I})^{-1}$$

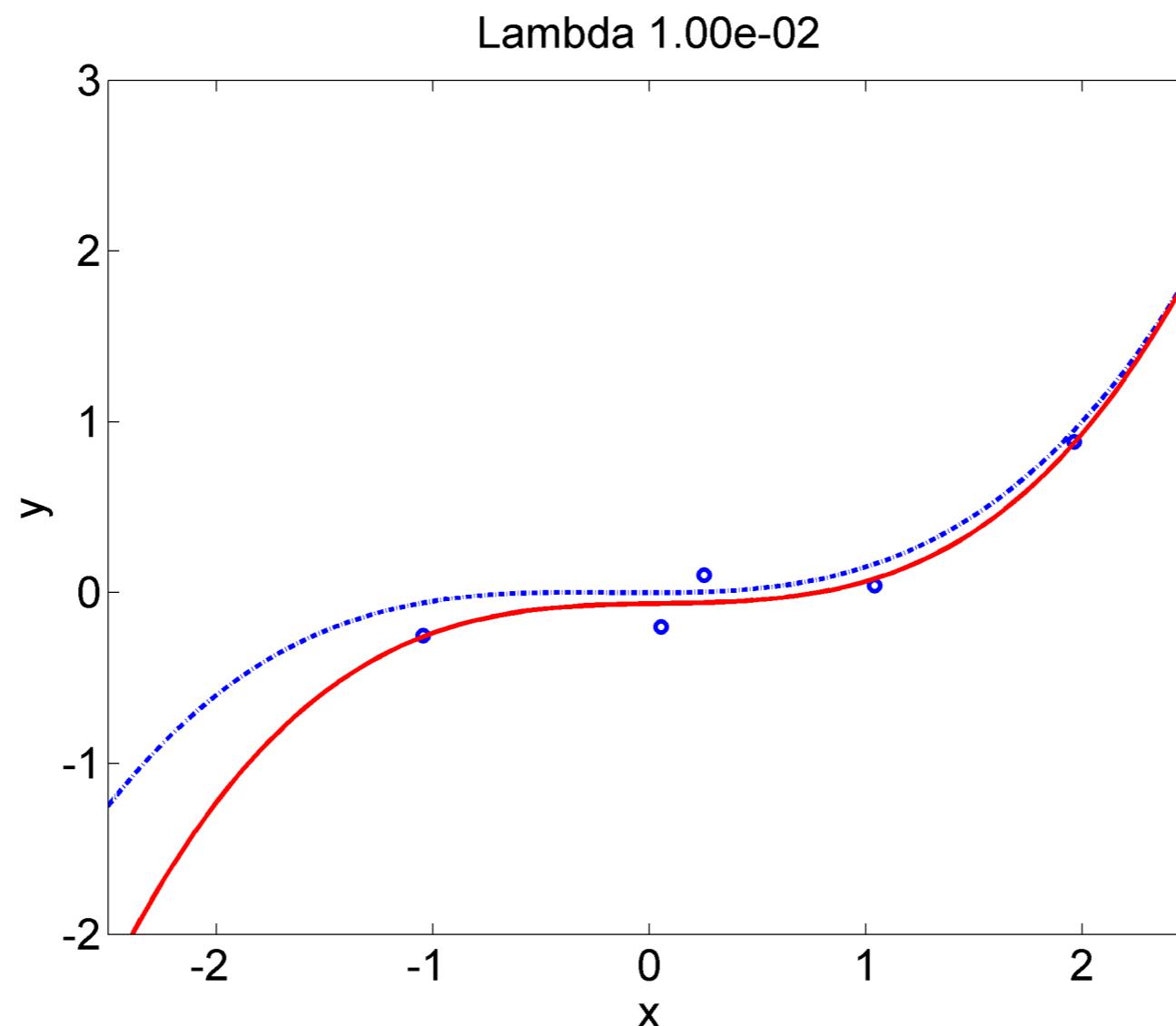
$$\boldsymbol{\mu}_N = \Sigma_N \Phi^T \mathbf{y}$$



# What to do with the posterior?

We could sample from it to **estimate uncertainty**

$$\theta_i \sim p(\theta|y, X) = \mathcal{N}(\theta|\mu_N, \Sigma_N)$$

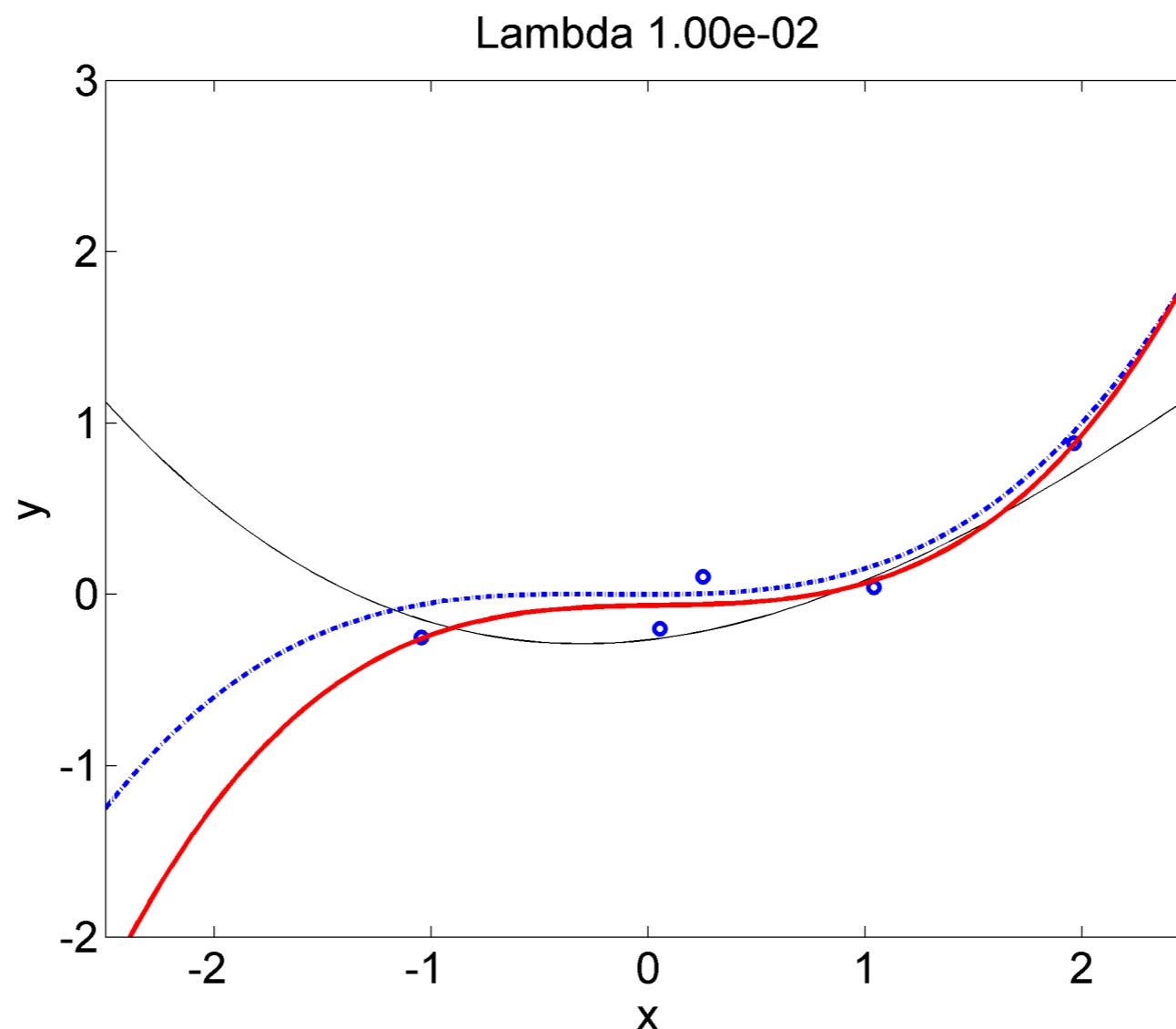




# What to do with the posterior?

We could sample from it to **estimate uncertainty**

$$\theta_i \sim p(\theta|y, X) = \mathcal{N}(\theta|\mu_N, \Sigma_N)$$

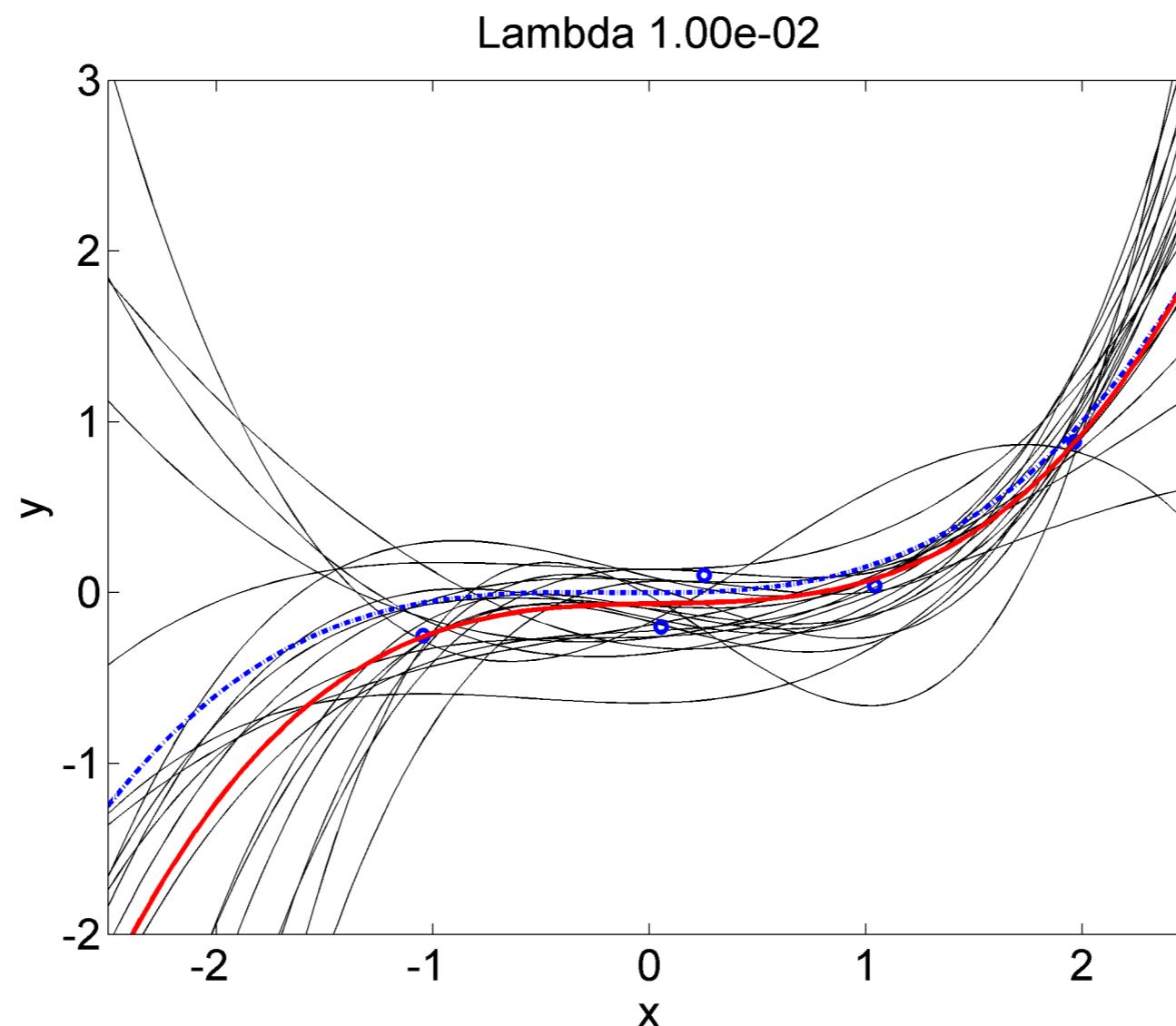




# What to do with the posterior?

We could sample from it to **estimate uncertainty**

$$\theta_i \sim p(\theta|y, X) = \mathcal{N}(\theta|\mu_N, \Sigma_N)$$



# Can we avoid the parameters $\theta$ ?



**Bayesian Fundamentalism:** We should not!

We don't care about parameters. We care  
about predictions!

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y})$$

# Full Bayesian Regression



We can also do that in closed form: integrate out all possible parameters

**Intuition:** If you assign each parameter estimator a “probability of being right”, the average of these parameter estimators will be better than the single one

# Full Bayesian Regression



We can also do that in closed form: integrate out all possible parameters

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* | \mathbf{x}_*, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) d\boldsymbol{\theta}$$

pred. function value    test input    training data

likelihood                              parameter posterior

Predictive Distribution is again a Gaussian

$$p(y_* | \mathbf{y}, \mathbf{X}, \mathbf{x}_*) = \mathcal{N}(\mathbf{y}_* | \mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

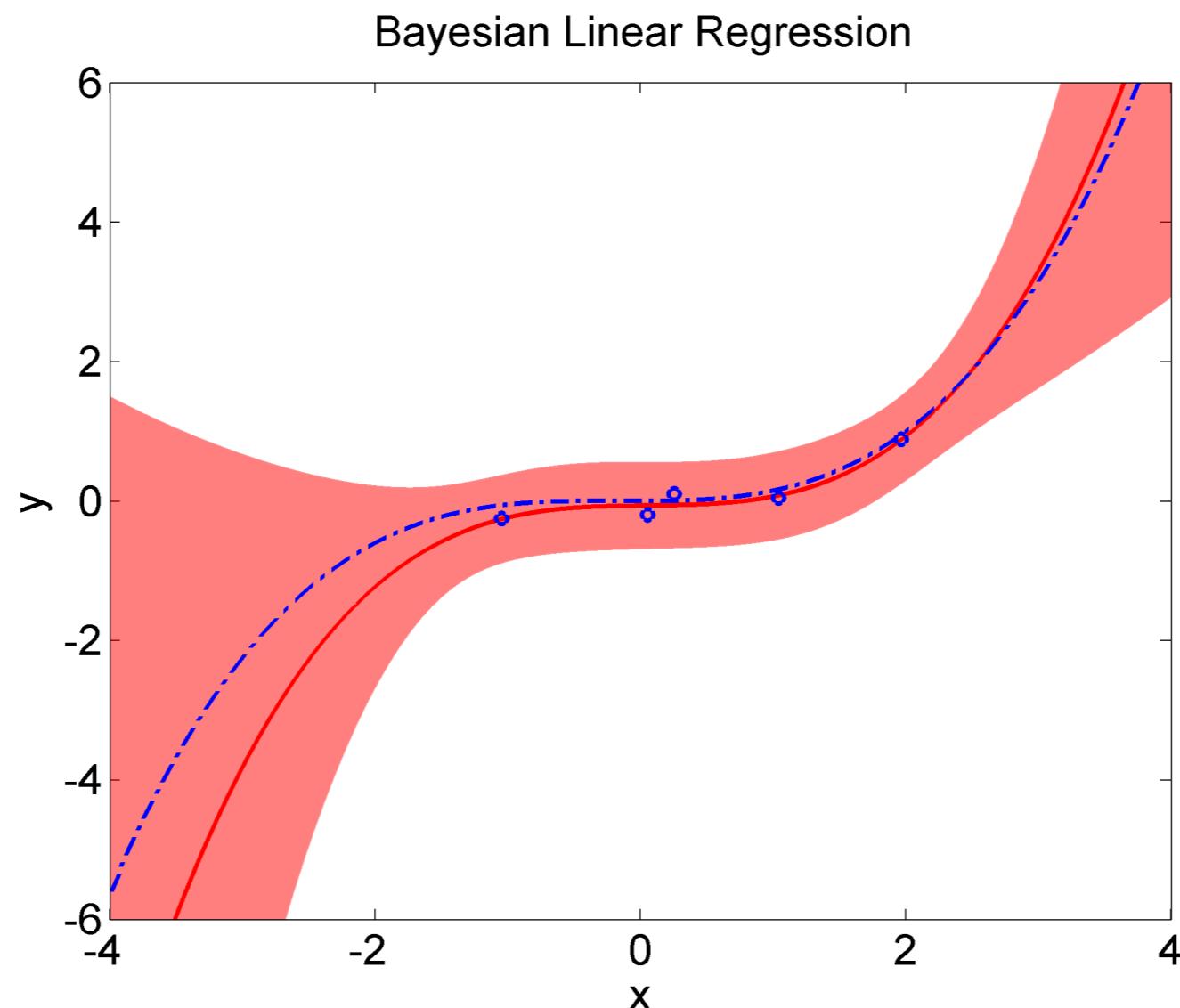
$$\mu(\mathbf{x}_*) = \boldsymbol{\phi}^T(\mathbf{x}_*)(\Phi^T\Phi + \sigma^2\lambda I)^{-1}\Phi^T\mathbf{y}$$

$$\sigma^2(x_*) = \sigma^2(1 + \phi^T(x_*)\Sigma_N\phi(x_*))$$

# State Dependent Variance!



# Integrating out the parameters



Variance depends on the information in the data!



# Quick Summary

- Models that are linear in the parameters:  $y = \phi(\mathbf{x})^T \boldsymbol{\theta}$
- **Overfitting** is bad
- **Model selection** (leave-one-out cross validation)
- **Parameter Estimation in Regression:** Frequentist vs. Bayesian
  - **Cost functions like Least Squares** go back to Gauss...
  - **Least Squares** ~ Maximum Likelihood estimation (ML; Frequentist)
  - **Ridge Regression** ~ Maximum a Posteriori estimation (MAP; Bayesian)
- **Full Bayesian Regression** integrates out the parameters when predicting
  - State dependent uncertainty



# Machine Learning 101b

---

Jan Peters  
Gerhard Neumann

# Purpose of this Lecture

---



- How can we define such features for general machine learning problems?
- Can we avoid or automate the feature specification?
- Familiarize you **with non-parametric models**

# Content of this Lecture

---



## Constructing Basis Functions

→ Radial Basis Function Networks

## Non-Parametric Approaches

→ Locally Weighted Regression

→ Kernel Methods



# What we did so far...

---

- Models that are linear in the parameters:  $y = \phi(\mathbf{x})^T \boldsymbol{\theta}$
- **Parameter Estimation in Regression**
  - **Least Squares** ~ Maximum Likelihood estimation (ML; Frequentist)
  - **Ridge Regression** ~ Maximum a Posteriori estimation (MAP; Bayesian)
  - **Full Bayesian Regression** integrates out the parameters when predicting
    - State dependent uncertainty

However, for most problems **good features are not easy to find**

# What to do when you don't know the features?



In most real applications, we know **some** good features.

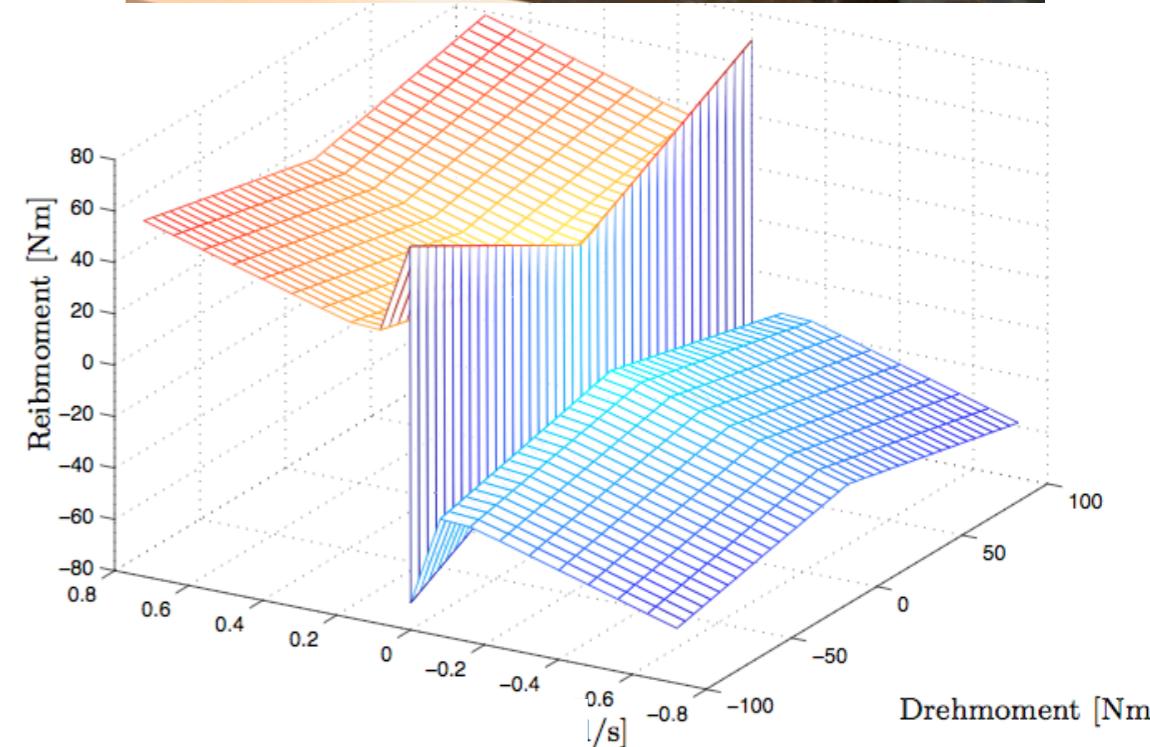
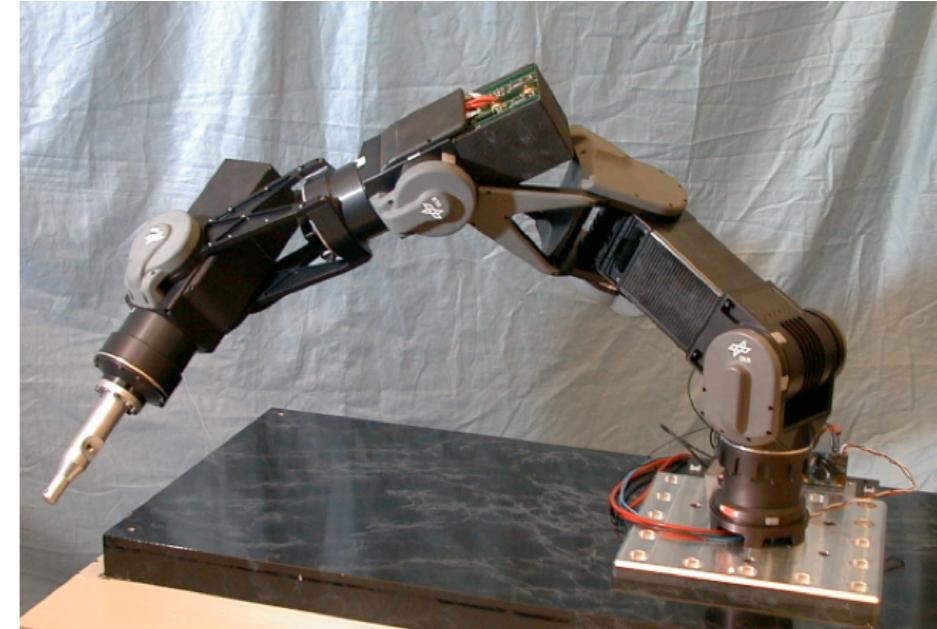
However, we almost certainly **don't know all** features we need.

**Example:** Rigid body dynamics

- Friction has no good features
- Unknown dynamics causes huge problems (requires more state variables).

There may also be way too many features!

Hand-crafted features are almost never enough...



# Can we avoid having to find good features?



**Yes, we can!**

We need to find machine learning approaches that **generate the features automatically** from data.

- **Type 1:** *Automatic Basis Function Construction* constructs basis functions from data.
- **Type 2:** *Non-Parametric Regression* look at data locally and interpolate with similar data.
- **Type 3:** *Kernel Regression* finds the features implicitly by going into *function space* using a *kernel*

# Type 1: Construct Basis Functions from Data

---



## Classical idea behind “neural networks”

- Multi-Layer Perceptrons (see Machine Learning: Statistical Approaches)
- Radial Basis Function Networks



# Radial Basis Function Network

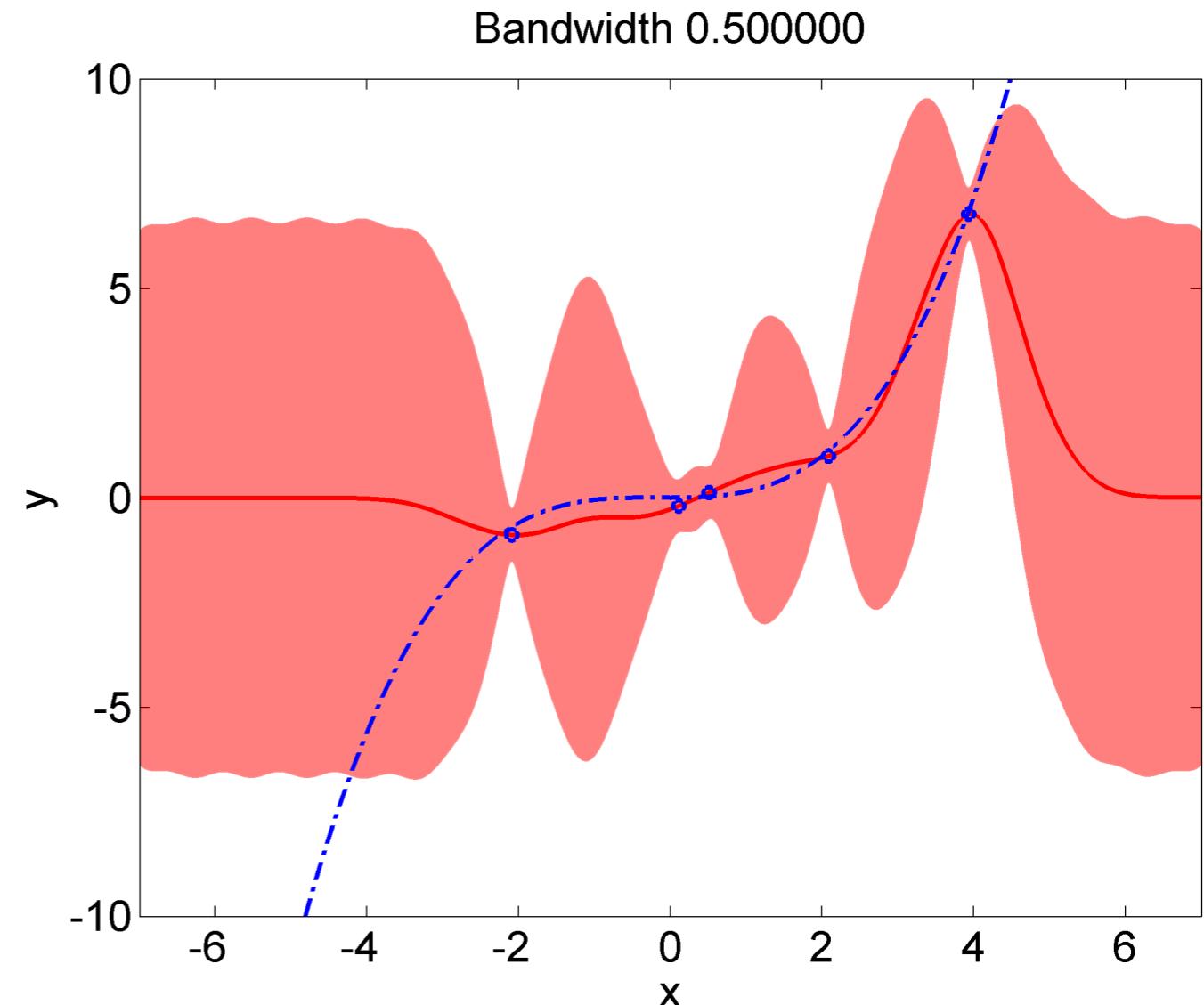
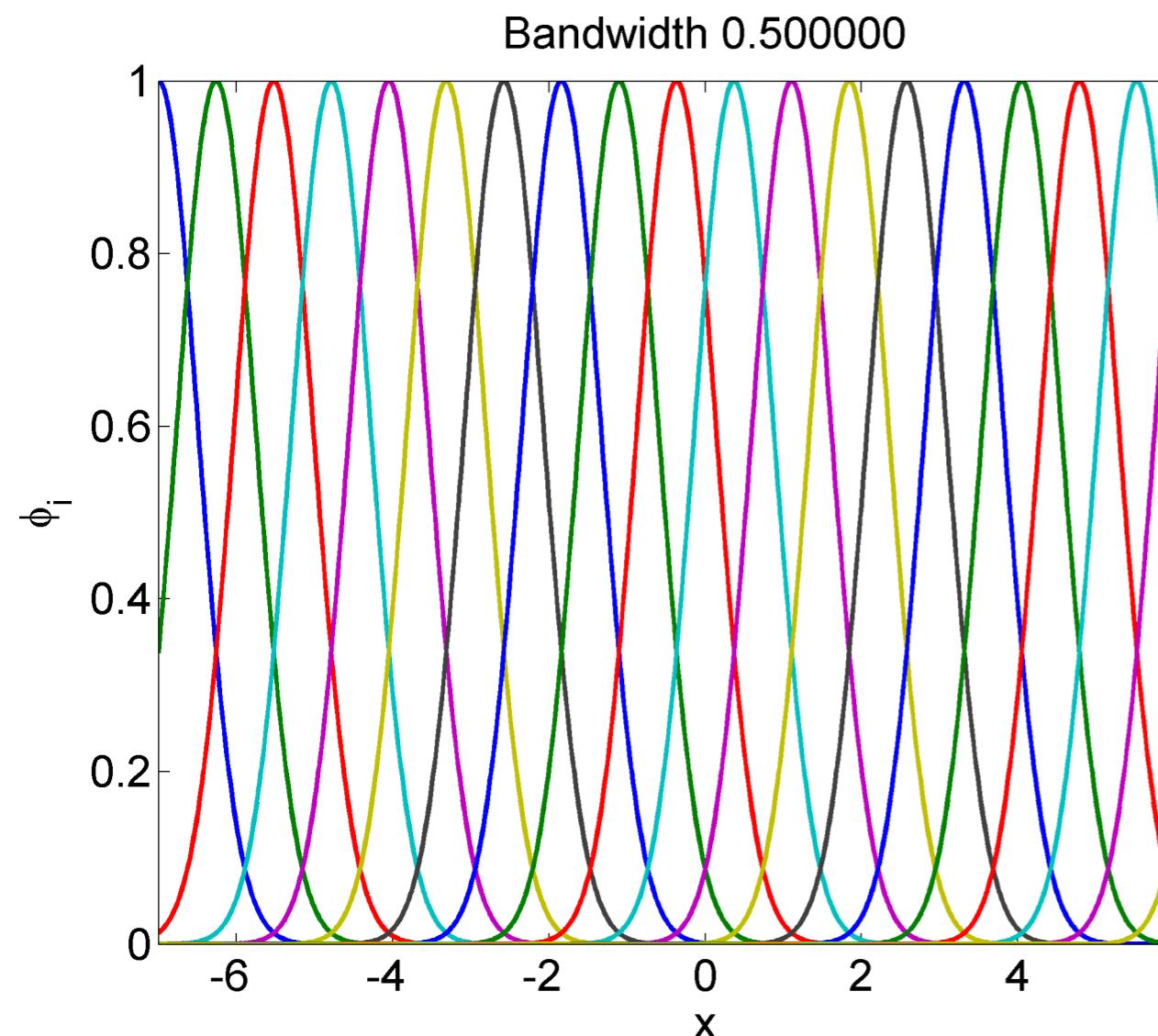
Assume a smoothness prior and obtain the cost function

$$J = \frac{1}{2} \sum_{i=1}^N \left[ (y_i - \mathbf{f}_\theta(\mathbf{x}_i))^2 + \left\| \frac{d^2}{dx^2} \mathbf{f}_\theta(\mathbf{x}_i) \right\|^2 \right]$$

This prior yields radial basis functions as features:

$$\begin{aligned} f_\theta(\mathbf{x}) &= \sum_i \theta_i \exp\left(\frac{-\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2l^2}\right) \\ &= \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\theta}, \quad \text{with } \phi_i(\mathbf{x}) = \exp\left(\frac{-\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2l^2}\right) \end{aligned}$$

# Example: Radial Basis Function Features



# Radial Basis Functions Hyperparameters



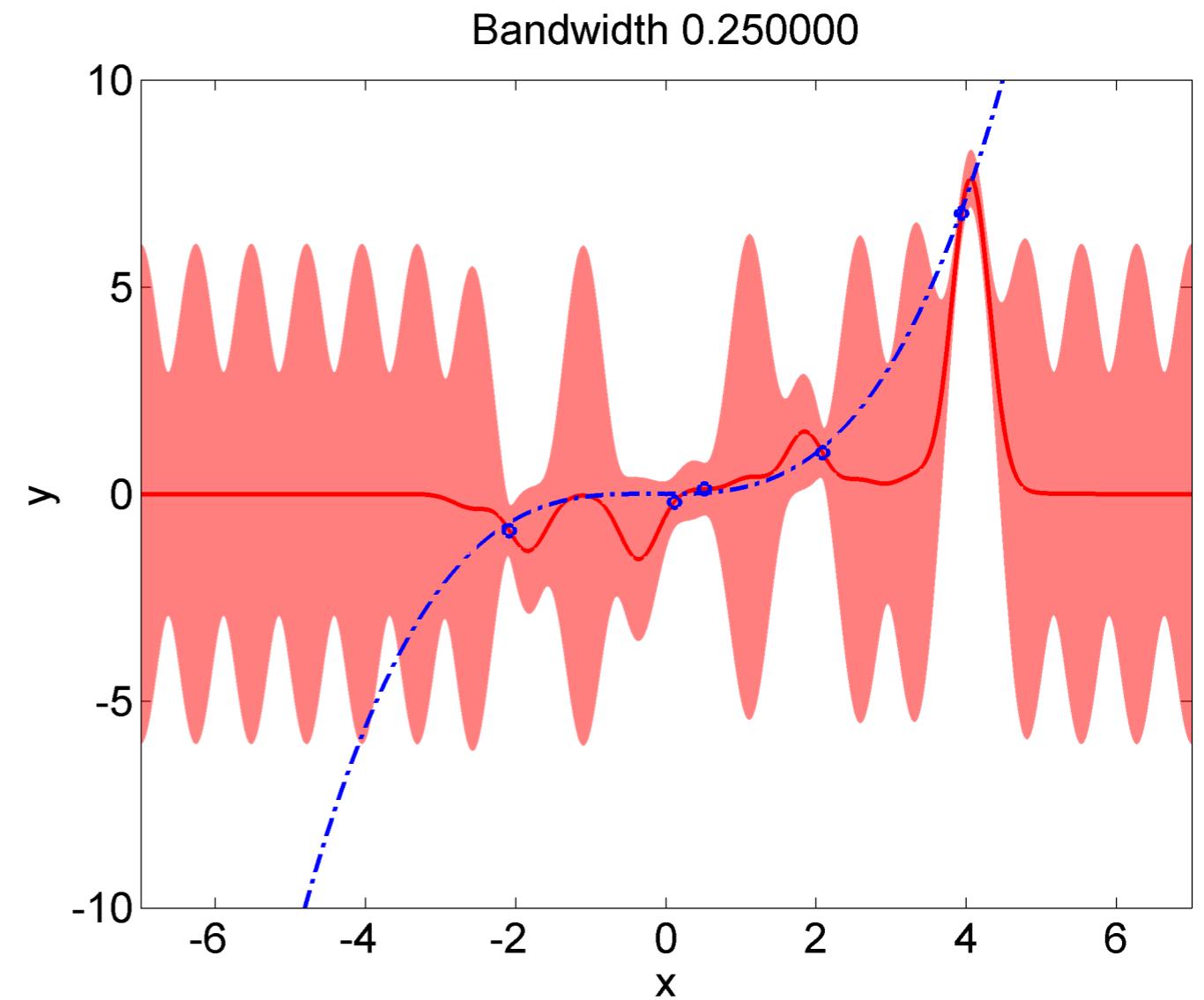
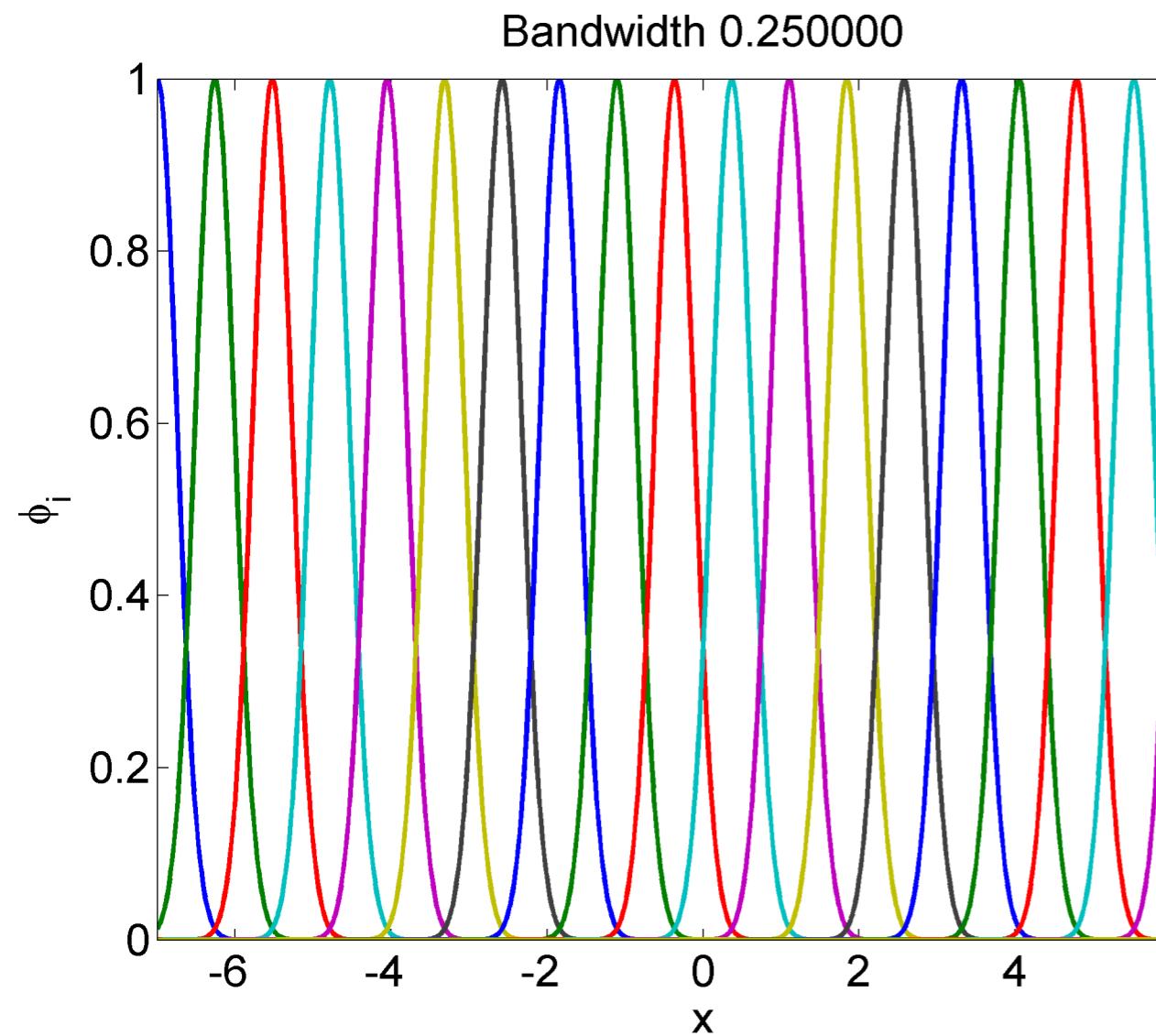
**Let's look again at**

$$\phi_i(\mathbf{x}) = \exp\left(\frac{-\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2l^2}\right)$$

- How do I find the width  $l$  of the basis functions or the centers  $\boldsymbol{\mu}_i$  ?
- Linear regression? Nope: not linear in  $l$  or  $\mu$  !
- We need to optimize this width on the training set
- We can do that by gradient descent: Write down a loss function, take the derivative w.r.t.  $l$ , and use an algorithm for non-convex optimization

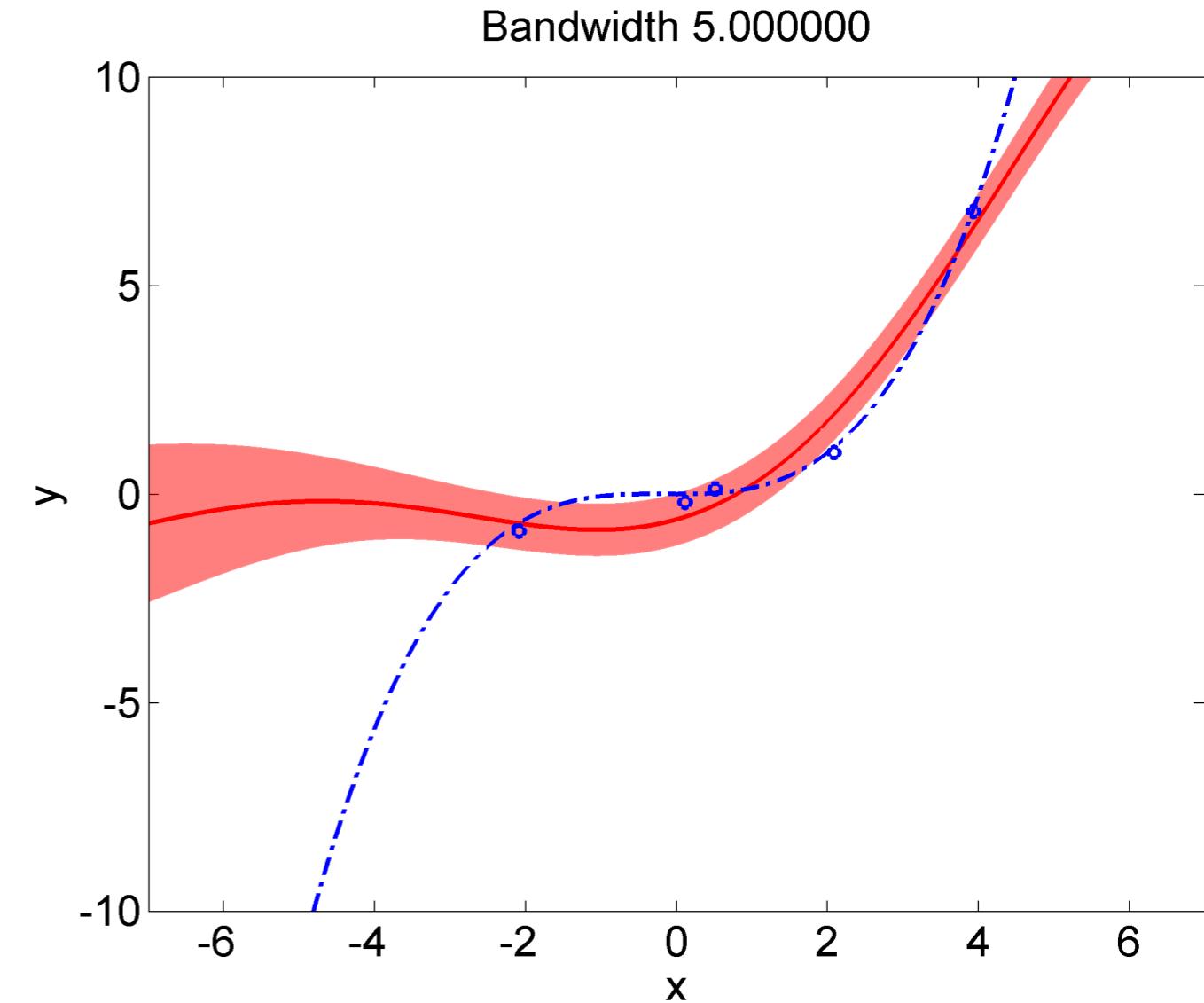
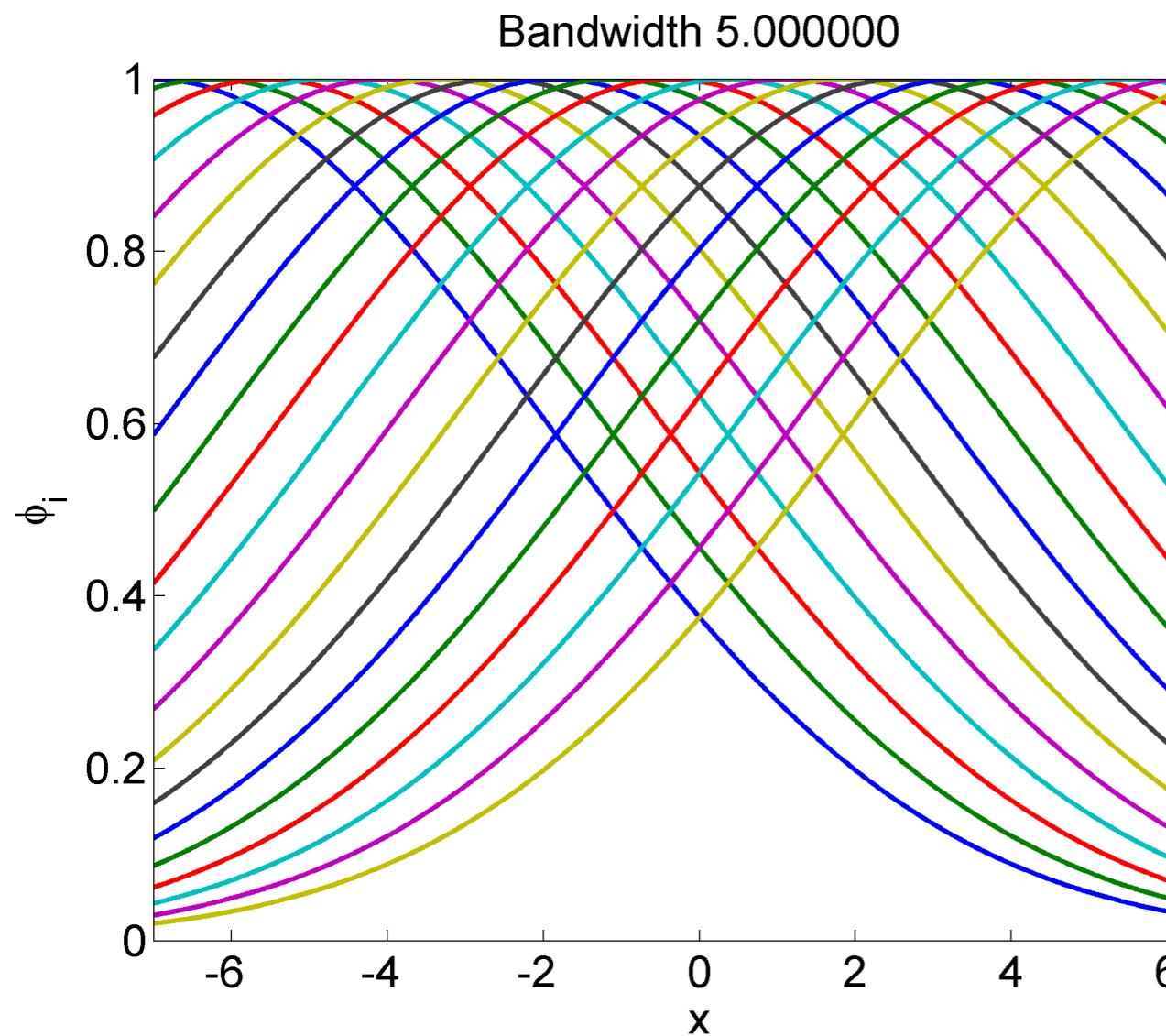


# Example: Bandwidth too small





# Example: Bandwidth too large



# Content of this Lecture

---



## Constructing Basis Functions

- Radial Basis Function Networks

## Non-Parametric Approaches

- Locally Weighted Regression
- Kernel Methods



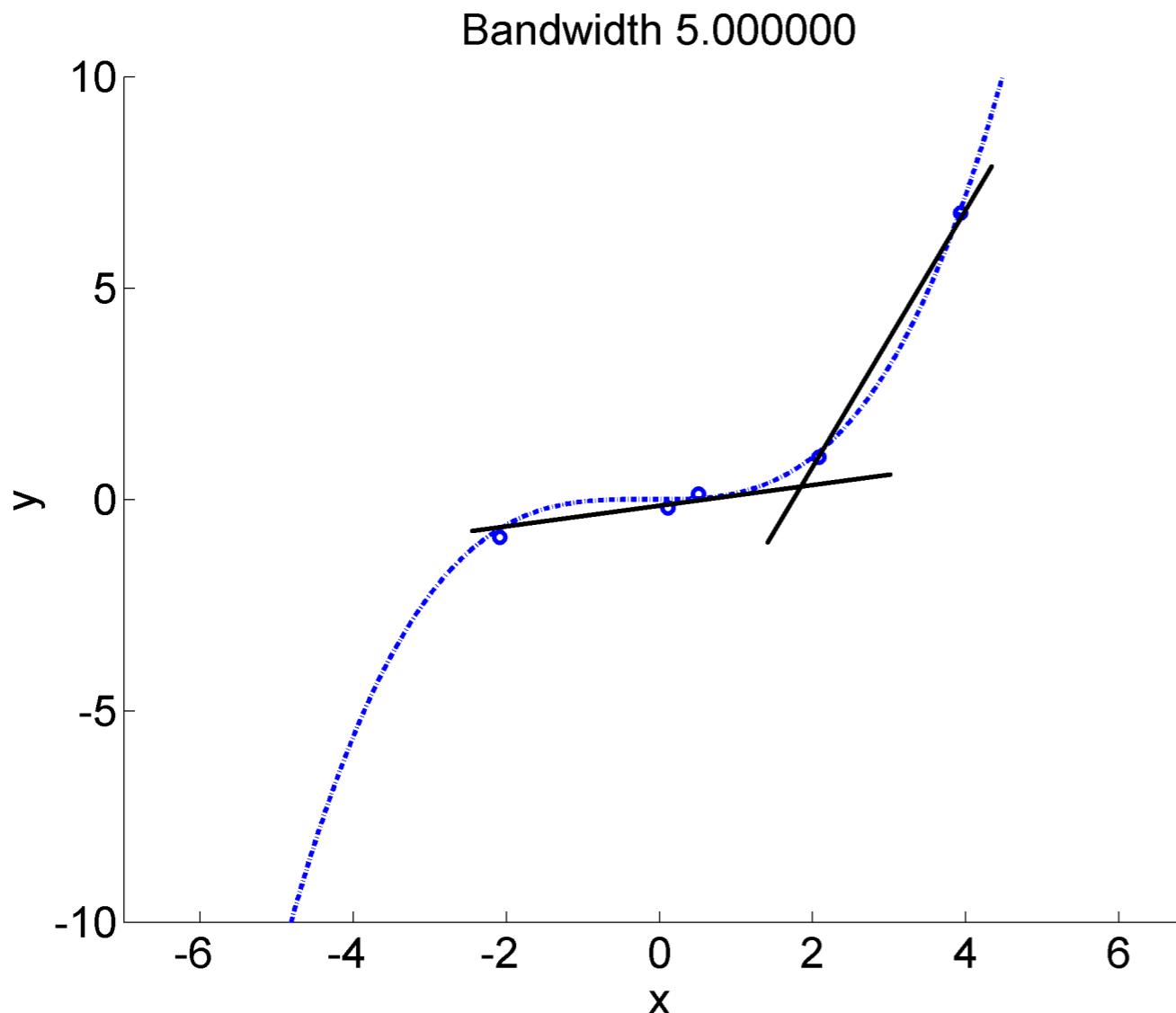
# Type 2: Non-Parametric Regression

---

- If you choose to have one feature/basis function per sample, you have a “**non-parametric method**” ➔ Don’t need to select the number of bases
- Non-parametric means
  - infinitely many parameters not no parameters
  - expressiveness of the model depends on the number of data points
  - No predetermined “parametric” form necessary
    - (e.g., “5th-degree polynomial”)
- One of them is **locally-weighted linear regression...**



# Example: Locally Linear Solutions



Locally all data is linear!

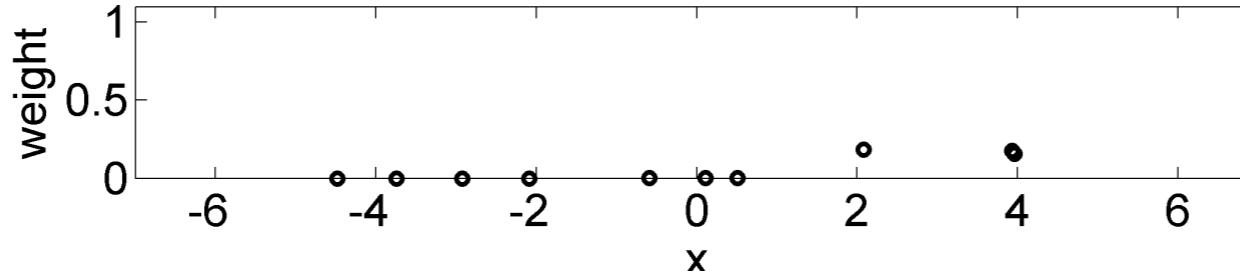
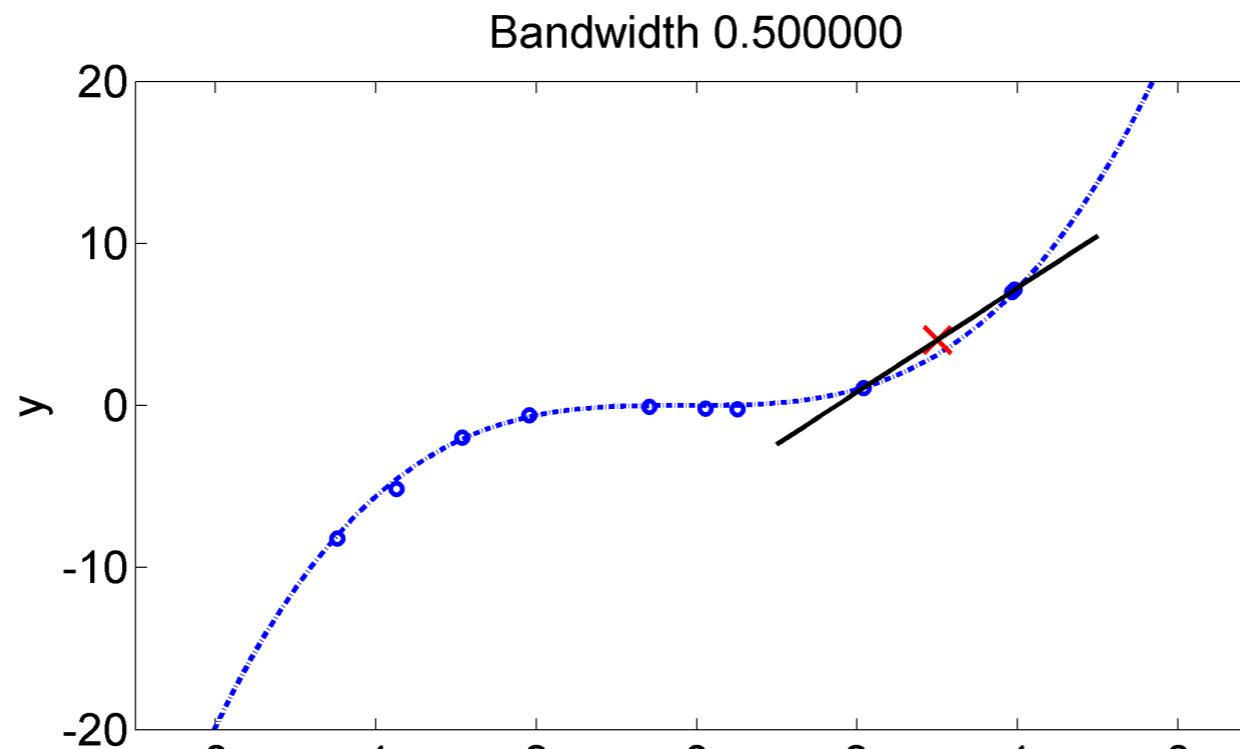


# Locally all data is linear ...

so why don't we **take the neighboring data points to predict the solution?**

- Use **higher importance** or weighting of **neighboring data points**
- For each query point  $x$ , weight training points  $x_i$  by

$$w_i(x) = \exp\left(-\frac{\|x - x_i\|^2}{2l^2}\right) \dots \text{squared exponential weighting}$$





# Weighted Linear Regression

Weighted cost function

$$J = \frac{1}{2} \sum_{i=1}^N w_i(\mathbf{x})(y_i - \mathbf{f}_\theta(\mathbf{x}_i))^2, \quad w_i(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2l^2}\right)$$

The function is linear in  $\mathbf{x}$

$$f_\theta(\mathbf{x}) = \boldsymbol{\theta}^T \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} = \boldsymbol{\theta}^T \tilde{\mathbf{x}}$$

In matrix form with  $W = \text{diag}(w_1, w_2, w_3, \dots, w_n)$  :

$$J = \frac{1}{2} (\tilde{\mathbf{X}} \boldsymbol{\theta} - \mathbf{y})^T W (\tilde{\mathbf{X}} \boldsymbol{\theta} - \mathbf{y})$$



# Weighted Linear Regression

**The solution to this problem:** weighted pseudo inverse

$$\theta = (\tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{W} \mathbf{y}$$

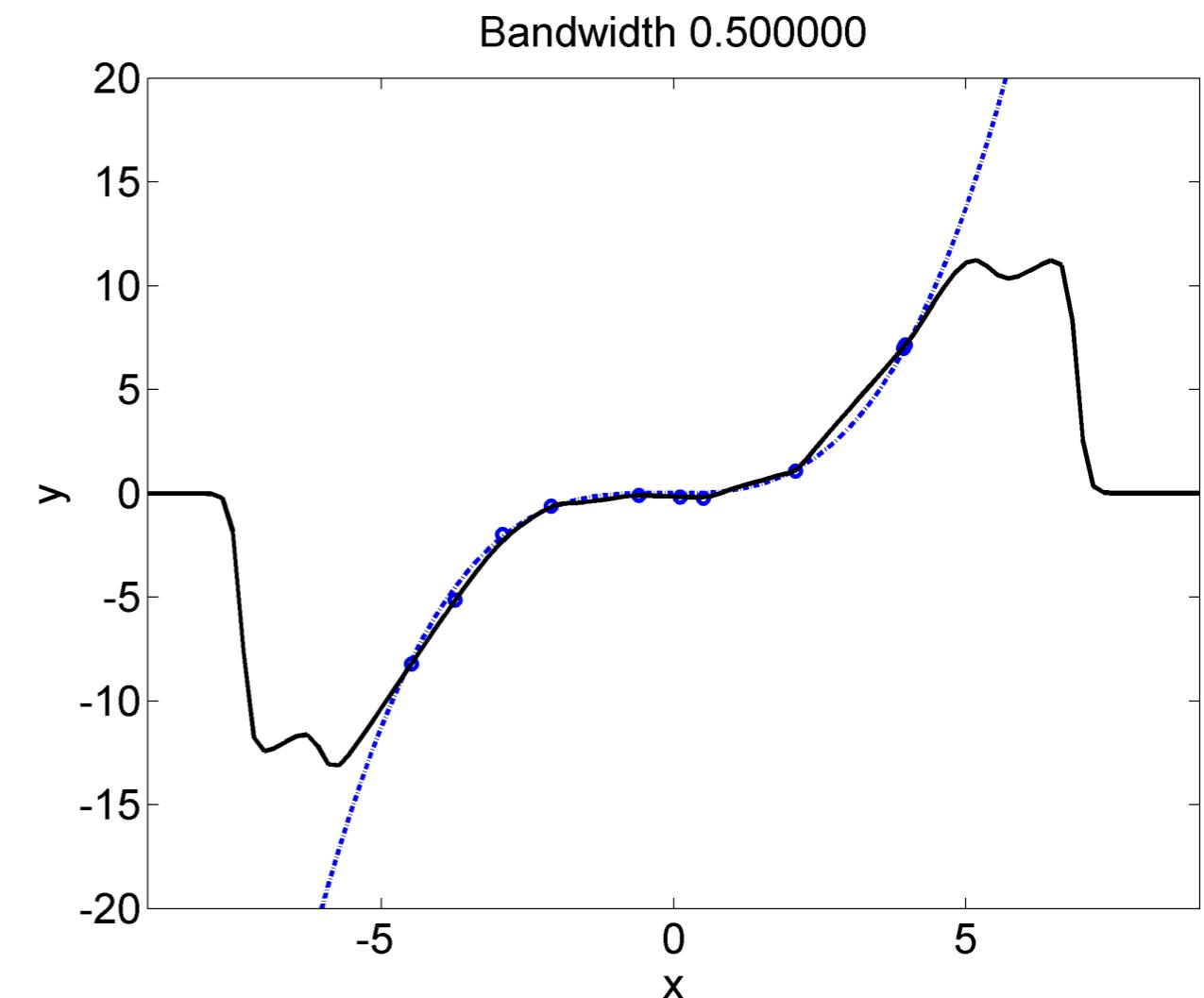
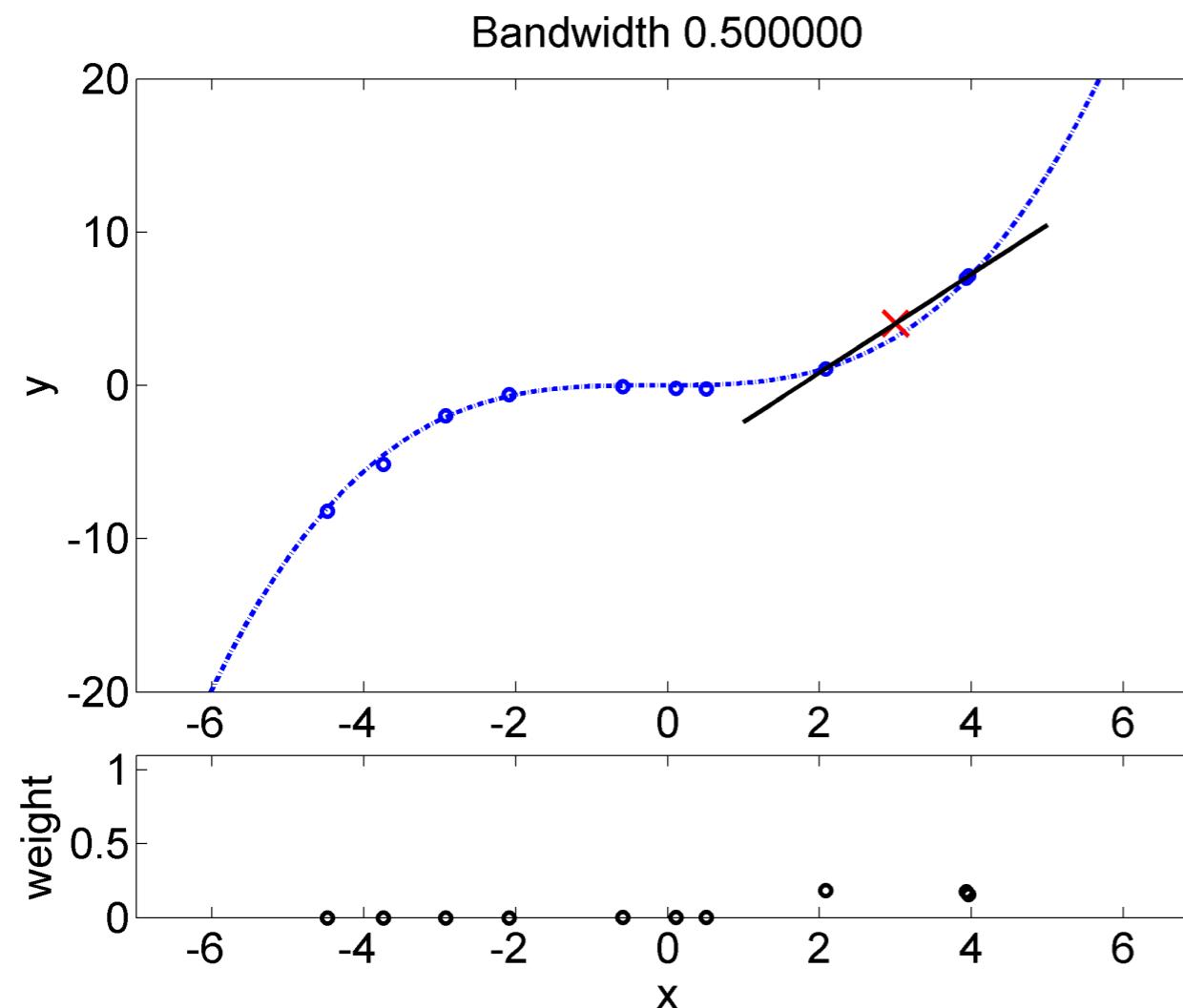
- $\mathbf{W}$  can be large - don't implement it like this...
- Dismiss data points with small weights / use bsxfun

## Local Ridge Regression:

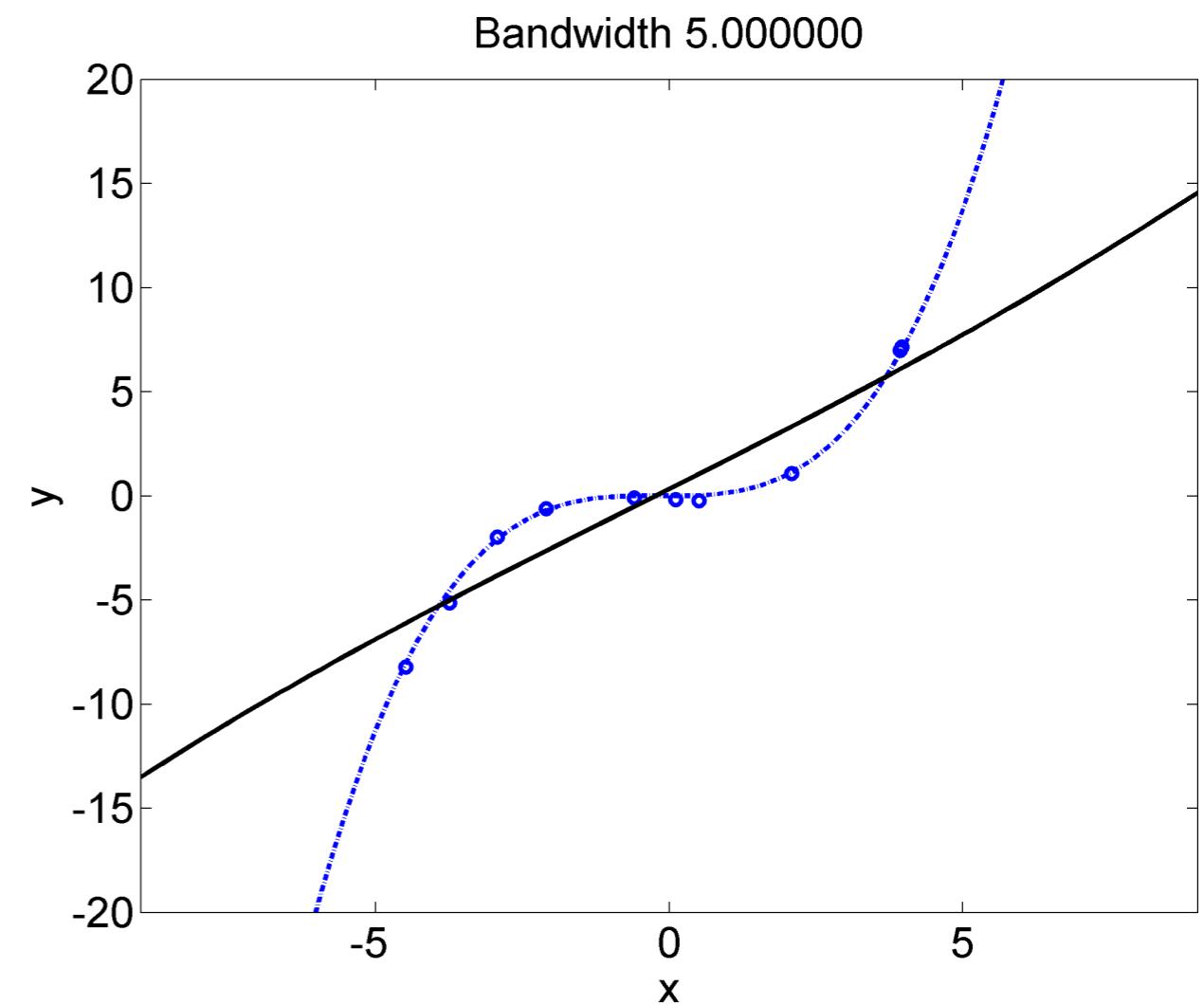
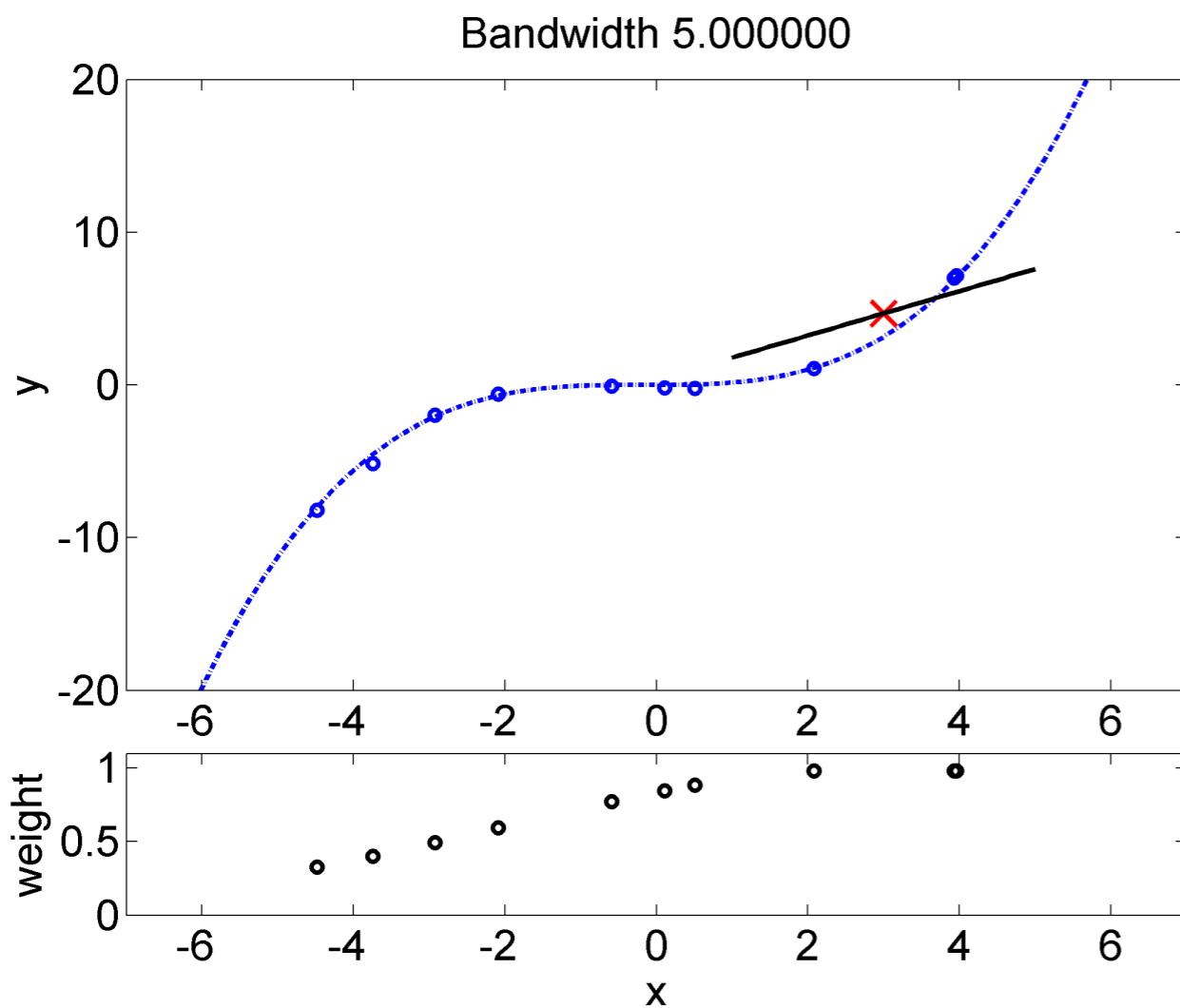
$$\theta = (\tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \mathbf{W} \mathbf{y}$$

Advantages: Fast(real-time capable), scales(lots of data), interpolates linearly(useful in control)  
Disadvantage: Tuning is not easy

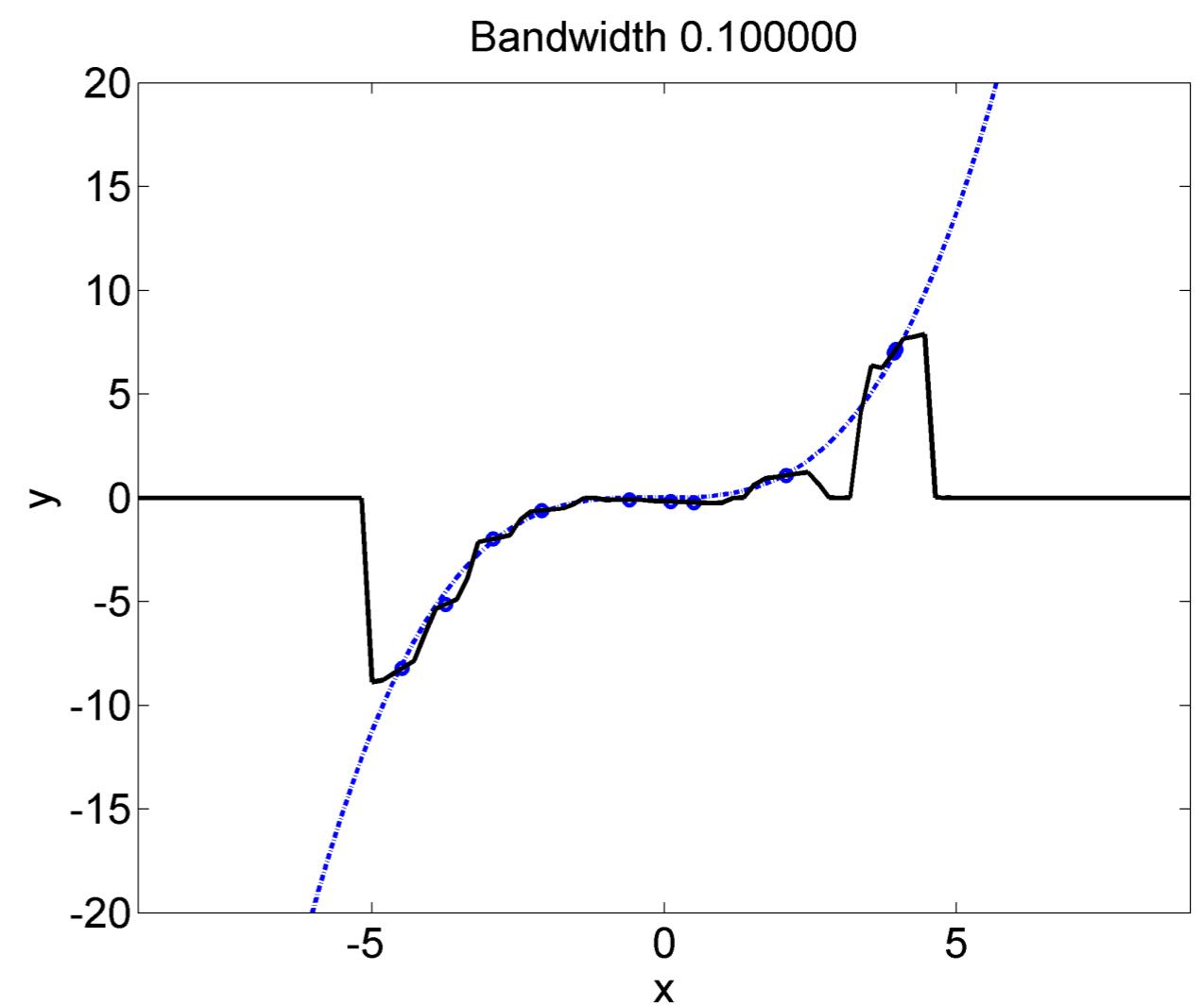
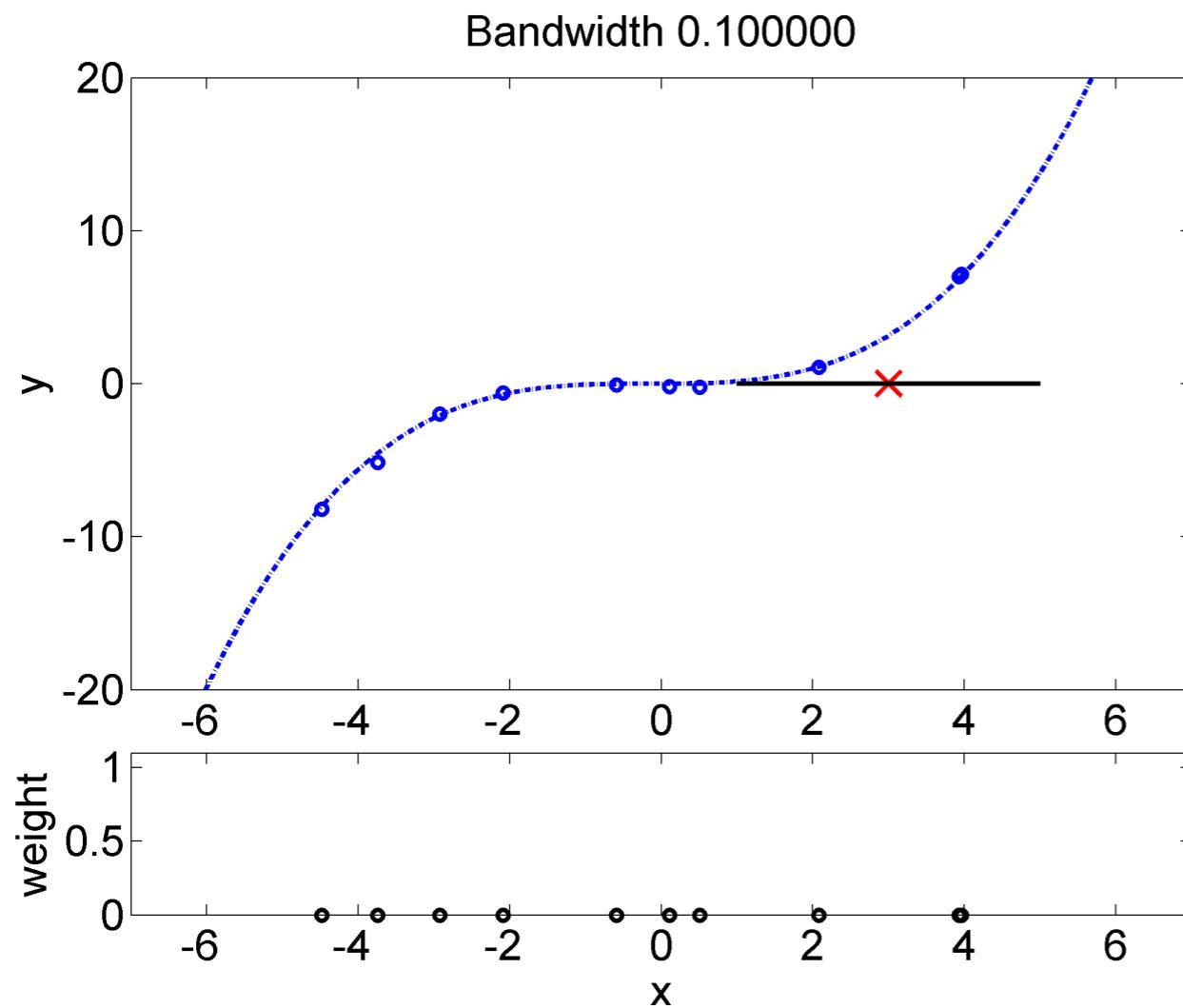
# Solution with Locally-Weighted Regression



# Solution with Locally-Weighted Regression



# Solution with Locally-Weighted Regression



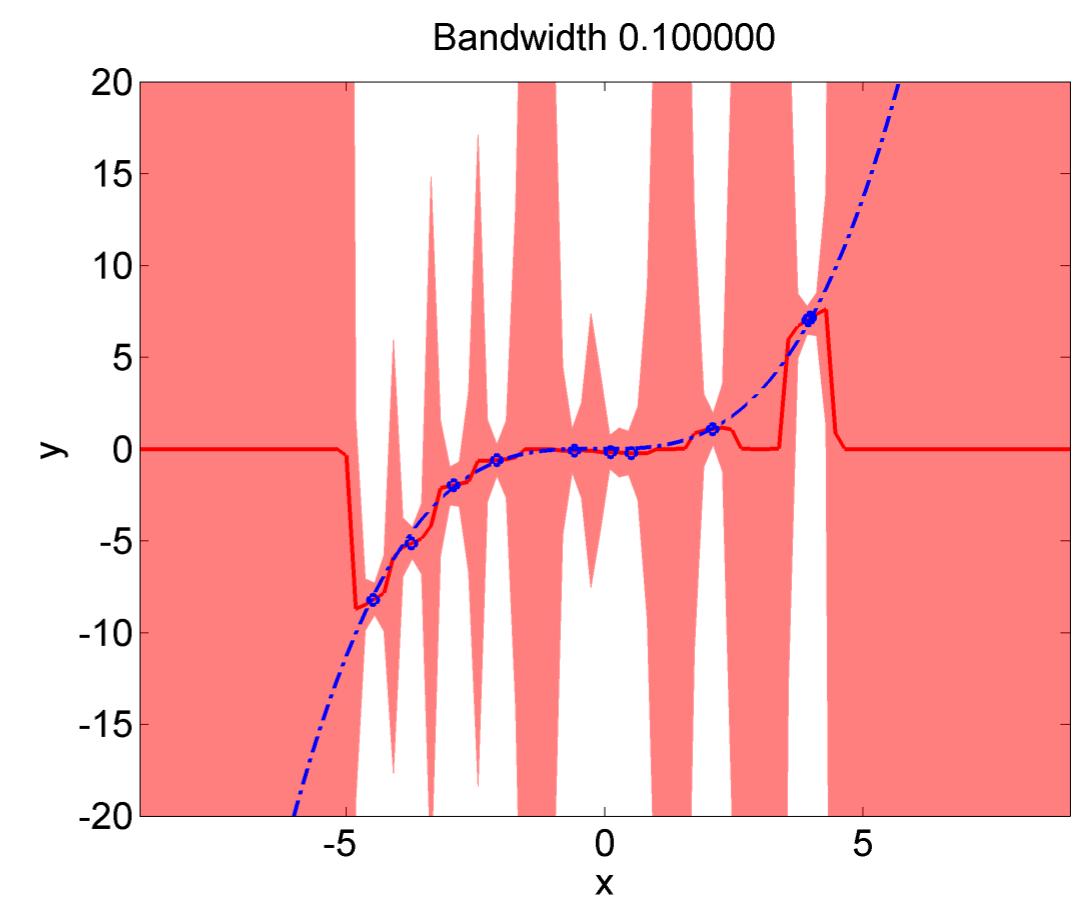
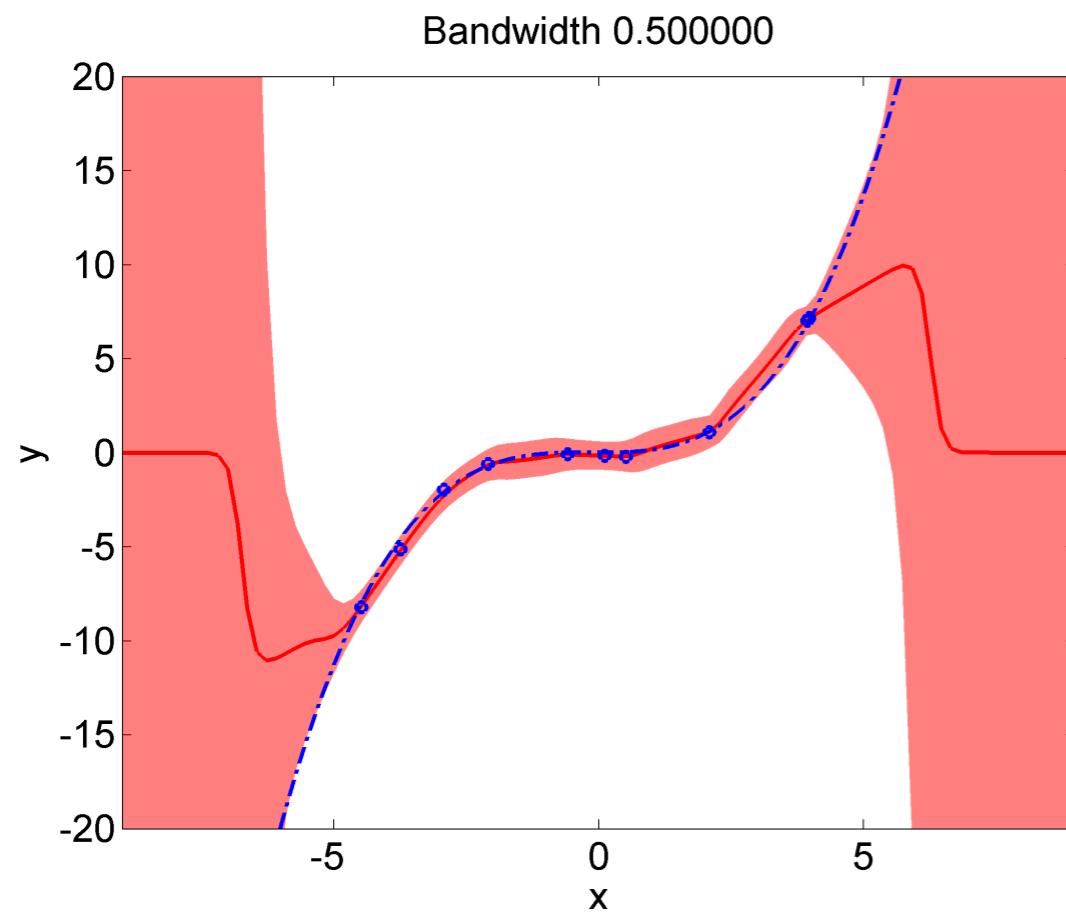
# Weighted Linear Regression



## Locally Weighted Bayesian Linear Regression

$$p(\theta|y, X, W) = \mathcal{N}(\theta|\mu_N, \Sigma_N)$$

$$\Sigma_N = (\tilde{X}^T W \tilde{X} + \sigma^2 \lambda I)^{-1} \quad \mu_N = \Sigma_N \tilde{X} W y$$



# Content of this Lecture

---



## Constructing Basis Functions

- Radial Basis Function Networks

## Non-Parametric Approaches

- Locally Weighted Regression
- Kernel Methods



# Type 3: Kernel Methods

**Kernel methods rely on the ‘kernel trick’**

- It is sufficient to evaluate **the scalar product between two samples in feature space**, called kernel

$$k(\mathbf{x}_1, \mathbf{x}_2) = \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2)$$

**Why is this useful?**

- Kernels are easier to design than features
- The feature space can be possibly infinite dimensional.
- We just need to be able to compute the scalar product



## Type 3: Kernel Methods

---

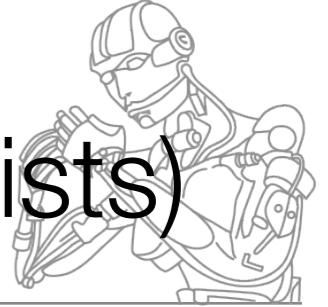
**Example:** One RBF feature at **every position c**

$$\begin{aligned} k(\mathbf{x}_1, \mathbf{x}_2) &= \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) = \int \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{c}\|^2}{2l^2}\right) \exp\left(-\frac{\|\mathbf{x}_2 - \mathbf{c}\|^2}{2l^2}\right) d\mathbf{c} \\ &= \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{4l^2}\right) \end{aligned}$$

→ Reduces to an RBF feature at **each sample**

General conditions for kernels

- **symmetric:**  $k(\mathbf{x}_1, \mathbf{x}_2) = k(\mathbf{x}_2, \mathbf{x}_1)$
- **positive definite...**



# Kernel Ridge Regression (Kernels for Frequentists)

Look at the predictions with the MAP/RR estimator (linear regression) again:

$$y(\mathbf{x}_*) = \phi(\mathbf{x}_*)^T \boldsymbol{\theta} = \phi(\mathbf{x}_*)^T (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y}$$

Even more general, the **Woodbury identity for matrix inversion** yields

$$(\Phi^T \Phi + \lambda \mathbf{I}_D)^{-1} \Phi^T = \Phi^T (\Phi \Phi^T + \lambda \mathbf{I}_N)^{-1}$$

This yields

$$y(\mathbf{x}_*) = \phi(\mathbf{x}_*)^T \underbrace{(\Phi^T \Phi + \lambda \mathbf{I}_D)^{-1} \Phi^T}_{D \times D} \mathbf{y}$$

$$= \phi(\mathbf{x}_*)^T \underbrace{\Phi^T (\Phi \Phi^T + \lambda \mathbf{I}_N)^{-1}}_{N \times N} \mathbf{y}$$

Equivalent solution to ridge regression

Why is this potentially useful?

# Kernel Ridge Regression (Kernels for Frequentists)



Let's say, we have an inner product for our features:

$$k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

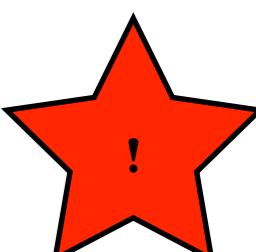
We can use this to rewrite

$$y = \phi(x)^T \Phi (\Phi \Phi^T + \lambda^{-1} I)^{-1} Y$$

Into

$$y = k(x)(K + \lambda^{-1} I)^{-1} Y$$

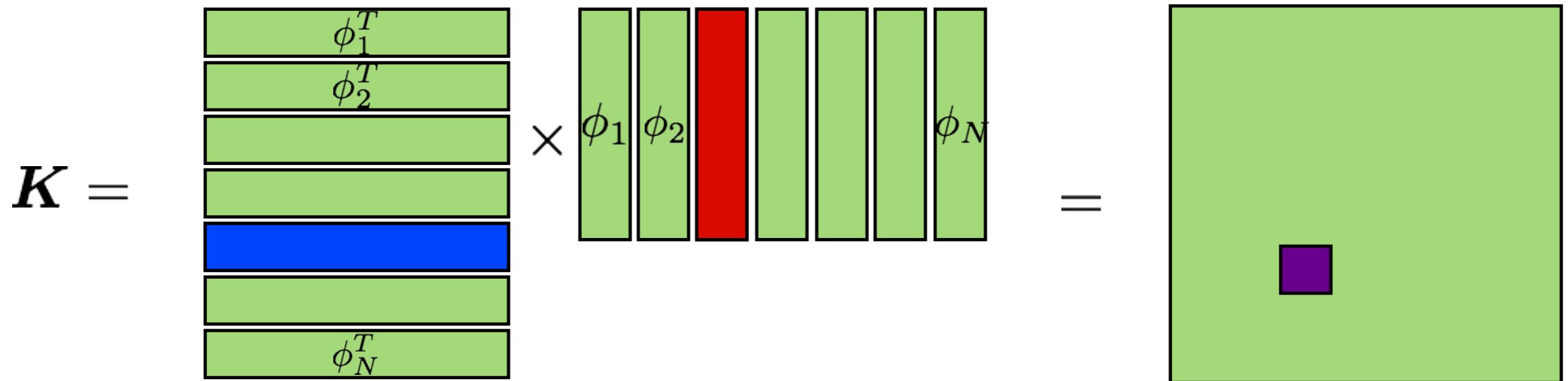
This a “Kernelization” of regression! **But why is this a good idea?**





# Kernels are scalar products in feature space!

$$K_{ij} = \lambda^{-1} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = k(\mathbf{x}_1, \mathbf{x}_2) \quad \text{scalar products in feature space}$$



Kernels can measure the similarity between data points in feature space without evaluating or explicitly knowing all features!

# Bayesian Linear Regression revisited



We have:

Data-Likelihood:  $p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{X}) = \mathcal{N}(\mathbf{y}|\Phi\boldsymbol{\theta}, \sigma^2 \mathbf{I})$

Prior:  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \lambda^{-1} \mathbf{I})$

If we integrate out the weights, we get

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int \mathcal{N}(\mathbf{y}|\Phi\boldsymbol{\theta}, \sigma^2 \mathbf{I}) \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \lambda^{-1} \mathbf{I}) d\boldsymbol{\theta} \\ &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \sigma^2 \mathbf{I} + \lambda^{-1} \Phi \Phi^T) \end{aligned}$$

Defines a multivariate Gaussian distribution over the samples

→ Samples are correlated as the marginalized weight vector is the same for each sample

# Bayesian Kernel Regression: Gaussian Processes (GPs)

---



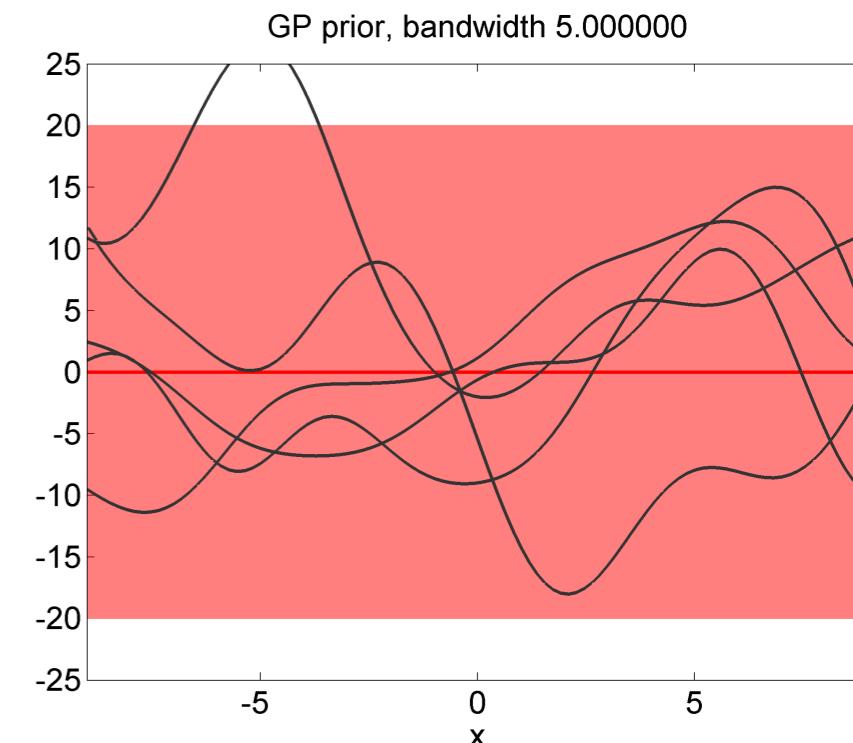
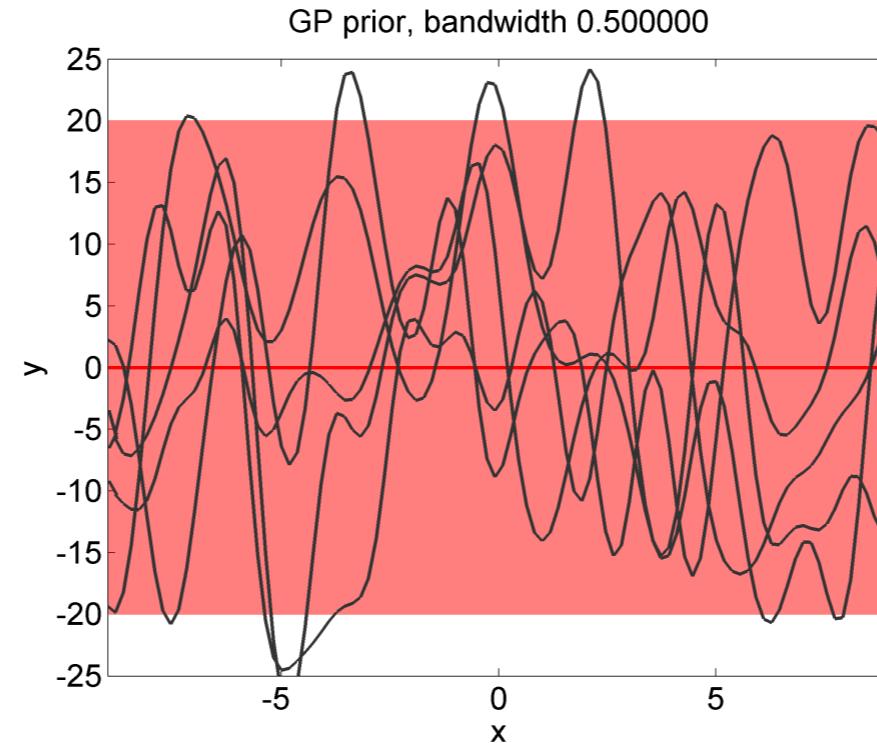
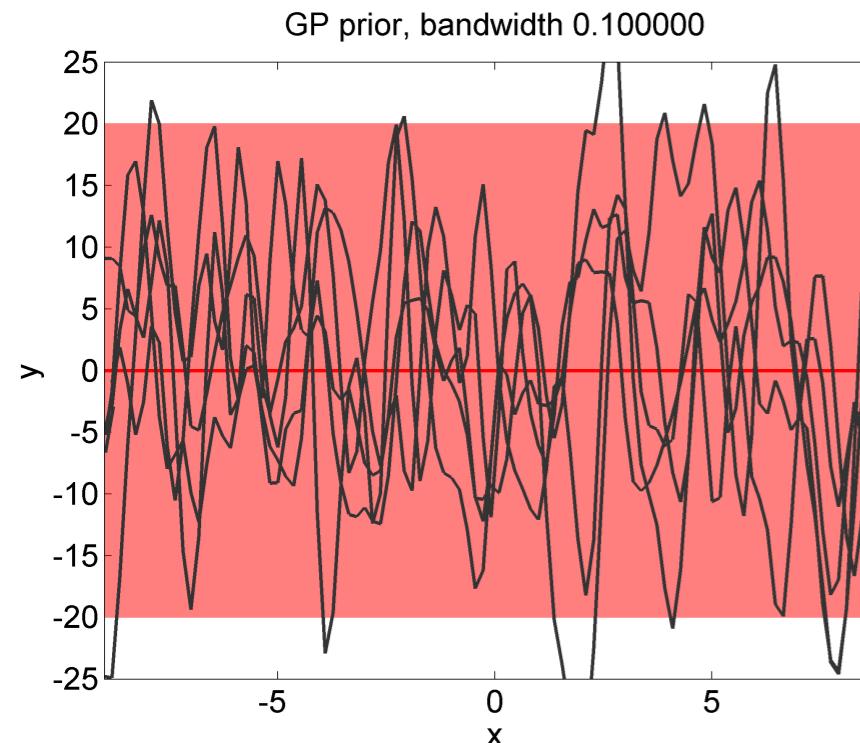
Replace the features in Bayesian Linear Regression by a Kernel and you obtain:

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{y}|0, \sigma^2 \mathbf{I} + \mathbf{K}) \text{ with } \mathbf{K} = \lambda^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^T$$

This method is called a Gaussian Process  $\mathcal{GP}(\mathbf{0}, \mathbf{K})$  with covariance function  $k$



# Sampling from the GP-Prior



The kernel bandwidth of the exponential kernel  
is a prior on the smoothness on the function

$$k(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{4l^2}\right)$$



# GP-Posterior

---

**Now we observe a data set given by  $\mathbf{y}$  and  $\mathbf{X}$  and we want to predict  $y^*$  for  $x^*$**

- We can write down the **GP prior for the concatenated data**

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{X} \\ \mathbf{x}_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K} & \mathbf{k}(\mathbf{X}, \mathbf{x}_*) \\ \mathbf{k}(\mathbf{x}_*, \mathbf{X}) & \mathbf{k}(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} + \sigma^2 \mathbf{I}\right)$$

- We get the **GP-posterior by Gaussian conditioning** (see refresher)

$$p(y_* | \mathbf{y}, \mathbf{X}, \mathbf{x}_*) = \mathcal{N}(y_* | \mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

# GP-Posterior



$$p(y_* | \mathbf{y}, \mathbf{X}, \mathbf{x}_*) = \mathcal{N}(\mathbf{y}_* | \mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

Same solution  
as in Kernel Ridge  
Regression

## Predictive mean

$$\mu(\mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

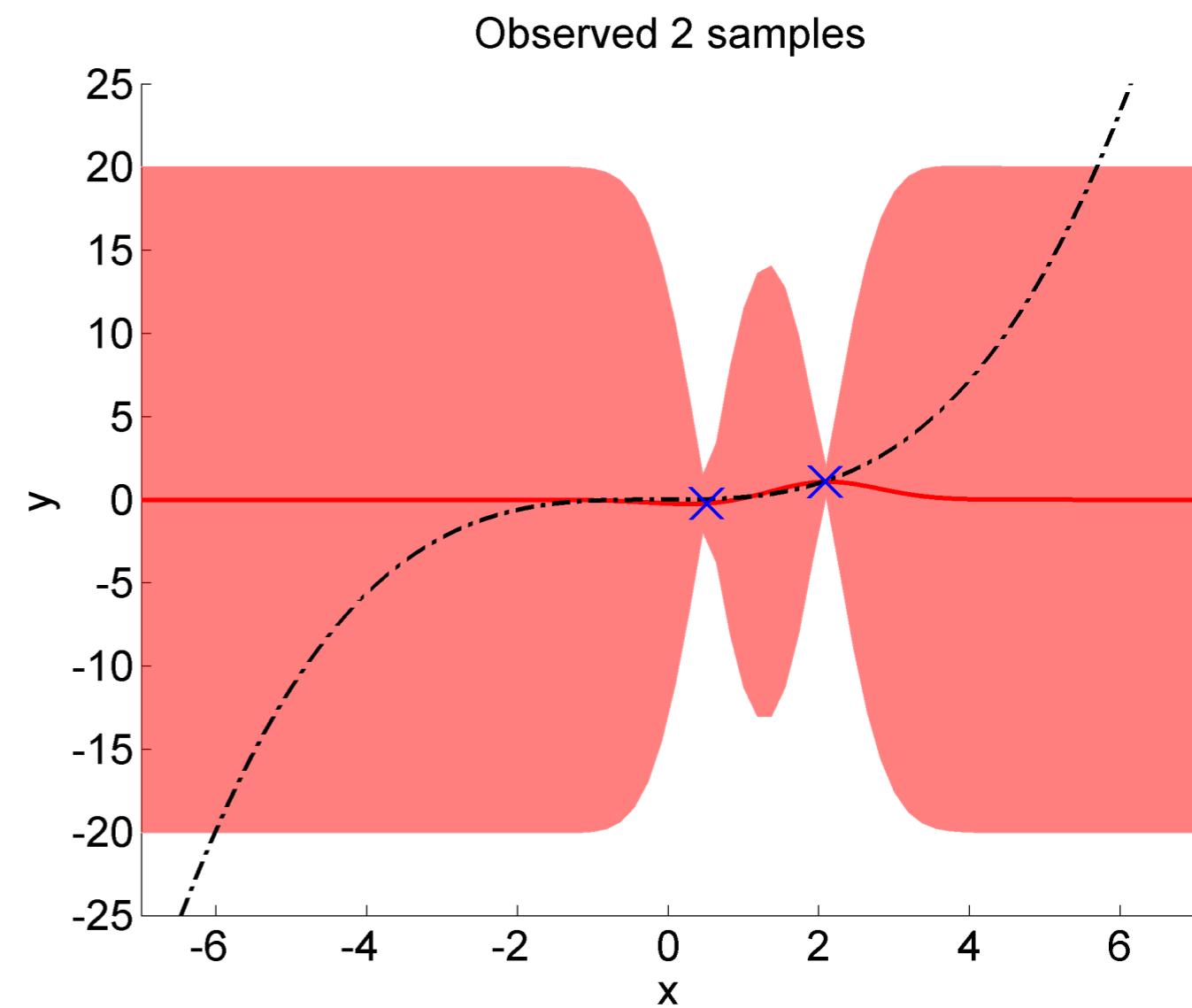
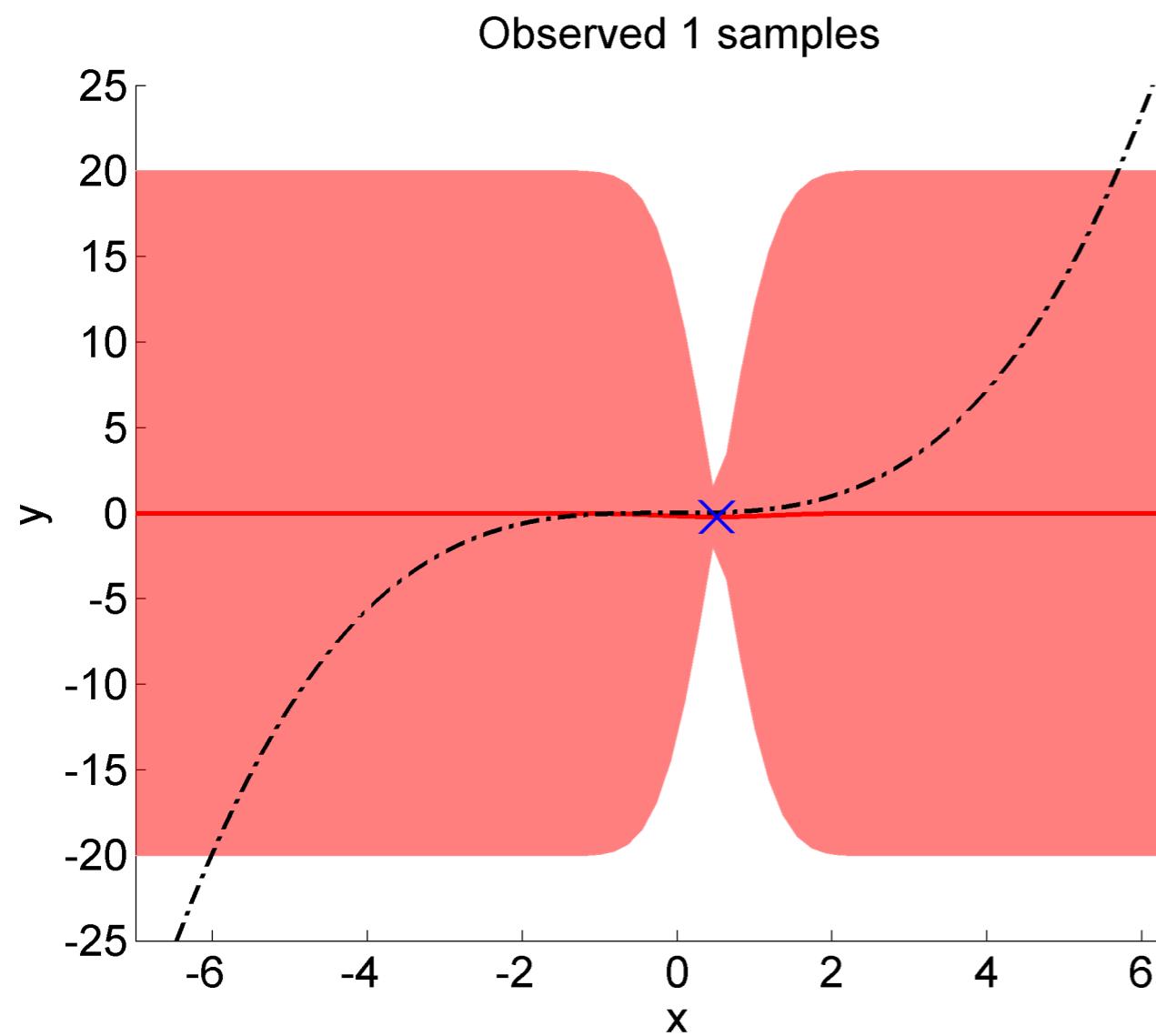
## Predictive variance

$$\begin{aligned} \sigma^2(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2 \\ &\quad - k(\mathbf{x}_*, \mathbf{X})(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}_*) \end{aligned}$$

GPs = Kernel Ridge Regression + Knowledge on your Uncertainty!

# GP-Posterior

---

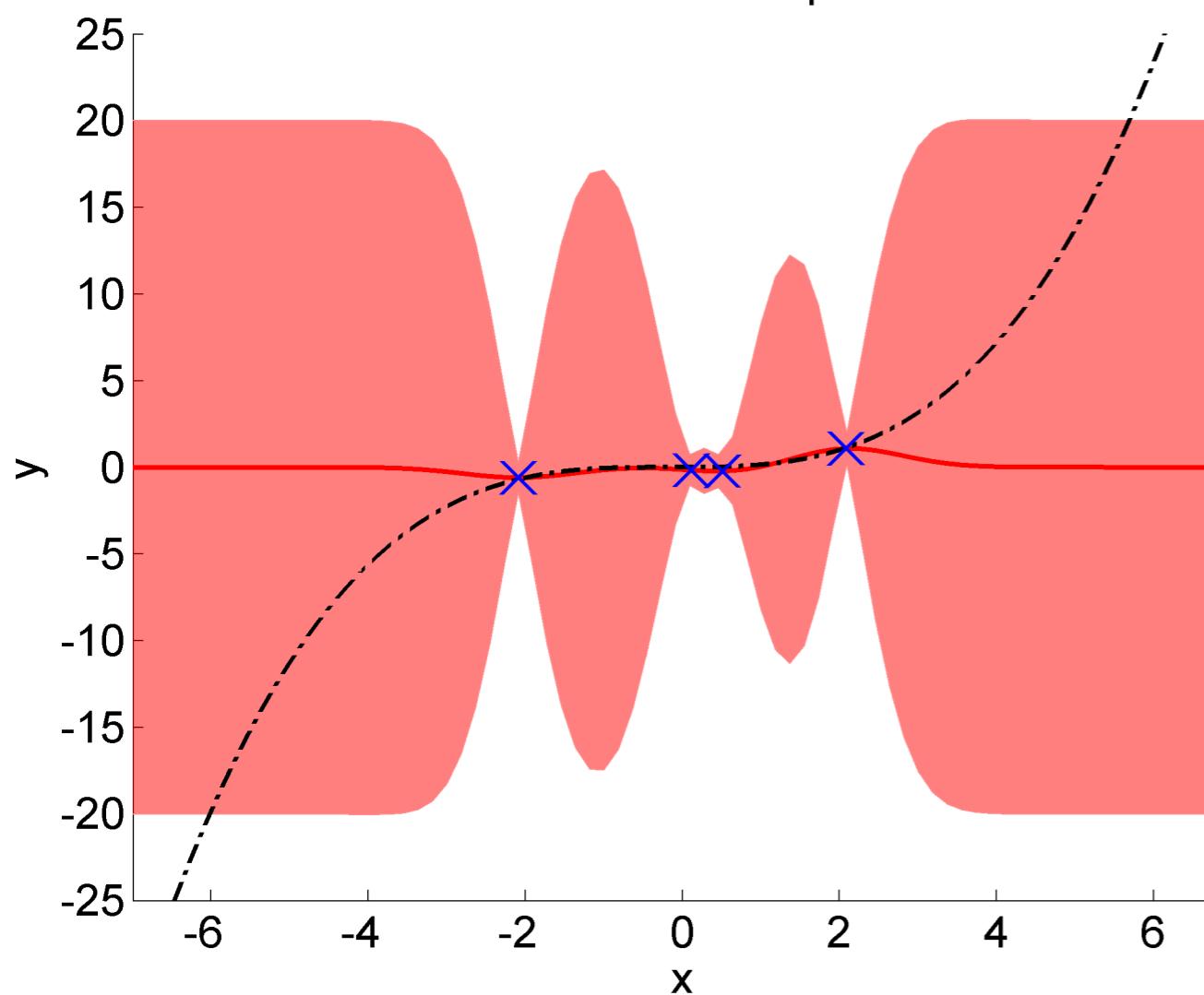


# GP-Posterior

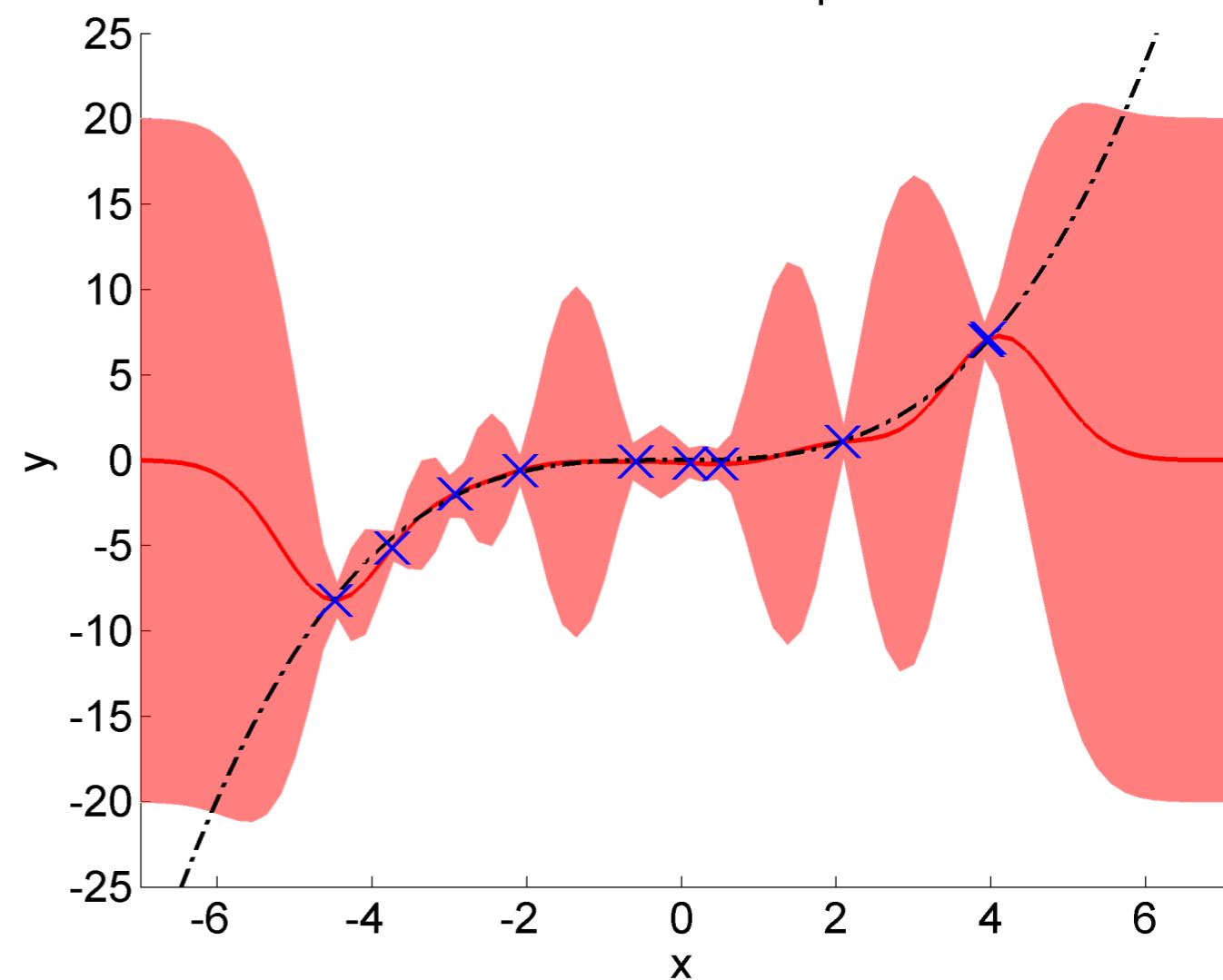
---



Observed 4 samples



Observed 10 samples



# Gaussian Processes

A key advantage over  
Kernel Ridge Regression



## Optimization of Hyper-Parameters (All learning becomes optimization)

- The parameters of the kernel are called **hyper-parameters**
- Cross validation or **maximization of marginal log-likelihood**

## GPs vs. Bayesian Linear Regression:

- GP is Kernel Ridge Regression with uncertainty! Same mean!
- GPs is **kernelized Bayesian Linear Regression**
- Kernels are often easier to use than features!

**GPs are currently the gold standard for regression!**





# Summary

---

- You should have a really good overview of machine learning by now.
- You should remember the following **regression methods**
  - Least-Squares Regression / Ridge Regression
  - Bayesian Regression
  - Radial-Basis Function Regression
  - Locally-Weighted Linear Regression
  - Kernel Ridge Regression and GPs
- You should know how to **choose the right method** for a regression problem