



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: Statistical Machine Learning

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

$$\frac{dy}{d\mathbf{x}} = \left[\frac{dy}{dx_1}, \dots, \frac{dy}{dx_n} \right] \quad \frac{d\mathbf{y}}{d\mathbf{x}} = \begin{bmatrix} \frac{dy_1}{dx_1} & \cdots & \frac{dy_1}{dx_n} \\ \vdots & \vdots & \vdots \\ \frac{dy_m}{dx_1} & \cdots & \frac{dy_m}{dx_n} \end{bmatrix}$$

- Some more matrix calculus

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

- 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:

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

works, if J has full column rank

Right Pseudo Inverse:

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

works, if J has full row rank



Statistics Refresher: Meet some old friends...

- **Gaussian Distribution**

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

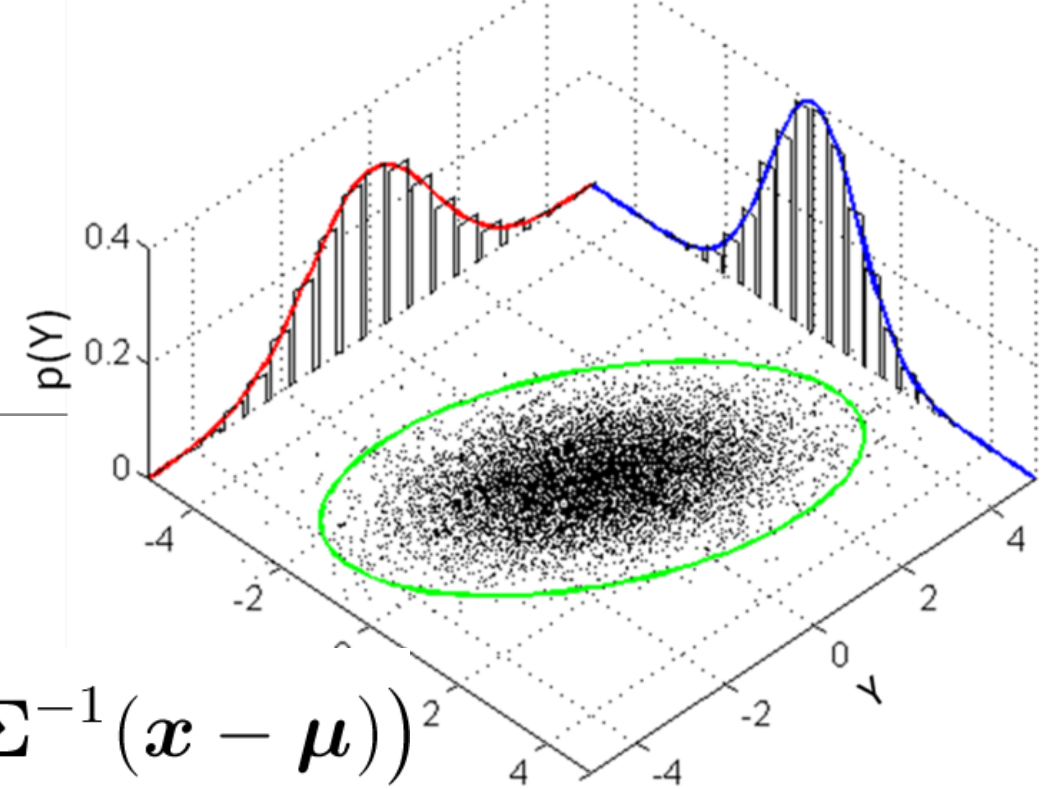
- Covariance matrix $\boldsymbol{\Sigma}$ captures **linear correlation**

- Product: **Gaussian stays Gaussian**

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

- **Mean is also the mode**

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



Statistics Refresher: Meet some old friends...



- Joint from **Marginal** and **Conditional**

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

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

- **Marginal** and **Conditional** Gaussian from Joint

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

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



Statistics Refresher: Meet some old friends...



- **Bayes Theorem** for Gaussians

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

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

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

Damped Pseudo Inverse

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



May I introduce you? The good old logarithm



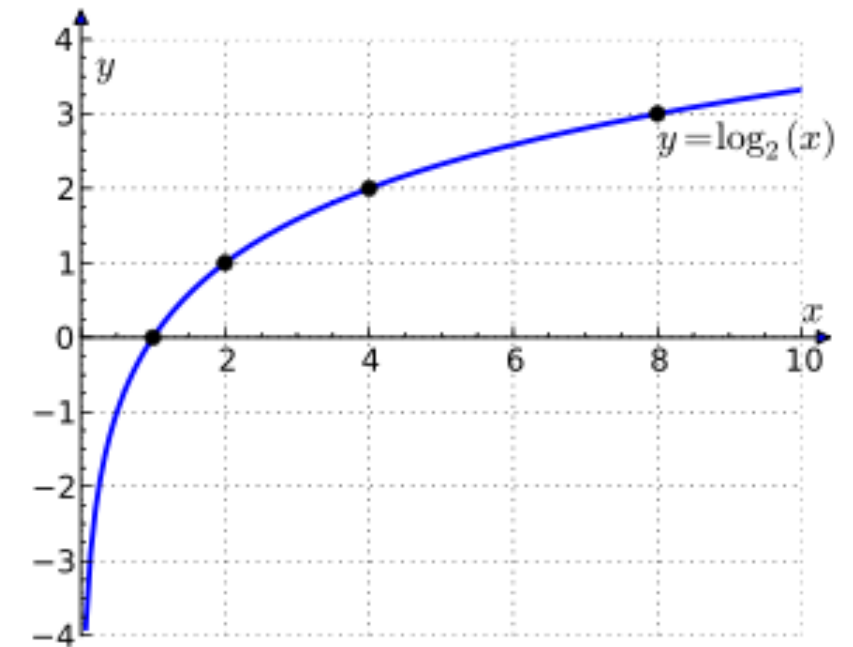
It's monoton $a > b \Rightarrow \log(a) > \log(b)$

... 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 (2.5×10^{15}) of information.

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

Why Machine Learning?

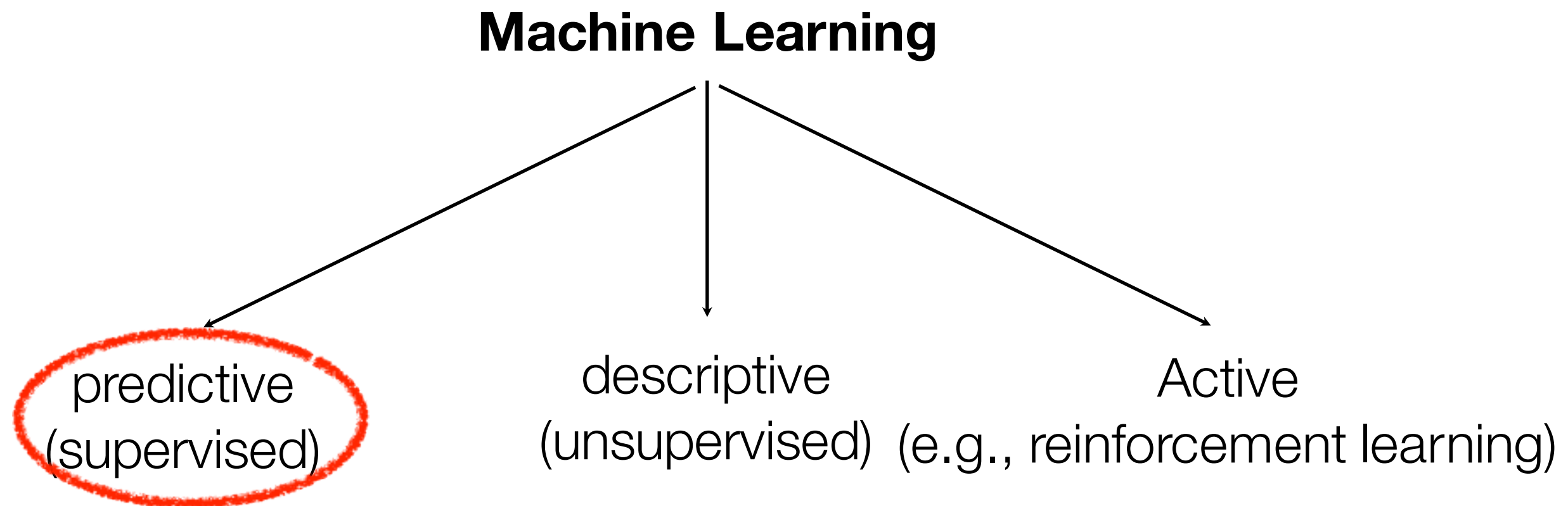


I keep saying the sexy job in the next ten years will be statisticians and machine learners. People think I'm joking, but who would've guessed that computer engineers would've been the sexy job of the 1990s? The ability to take data — to be able to understand it, to process it, to extract value from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

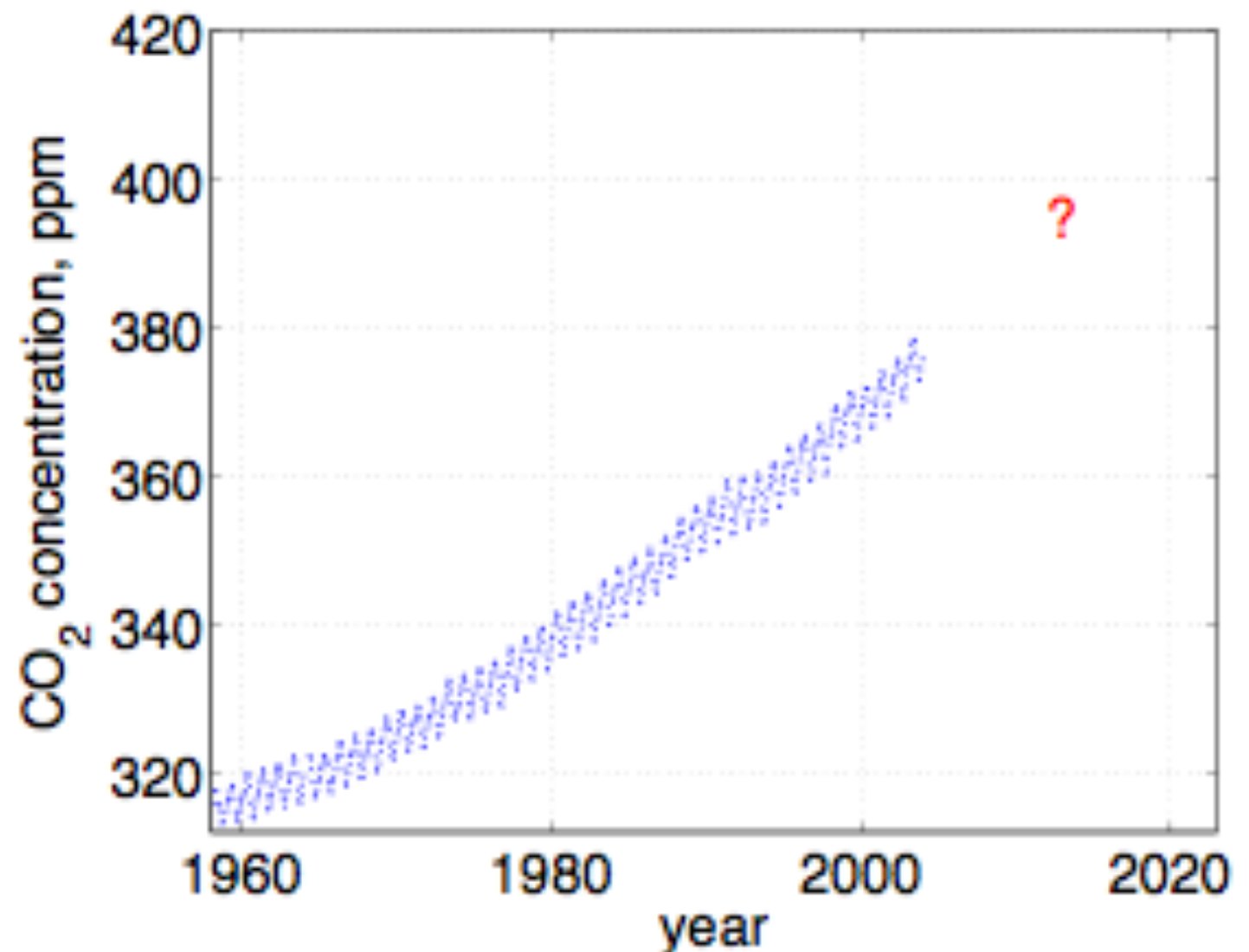
Types of Machine Learning



Prediction Problem (=Supervised Learning)



What will be the CO₂ concentration in the future?



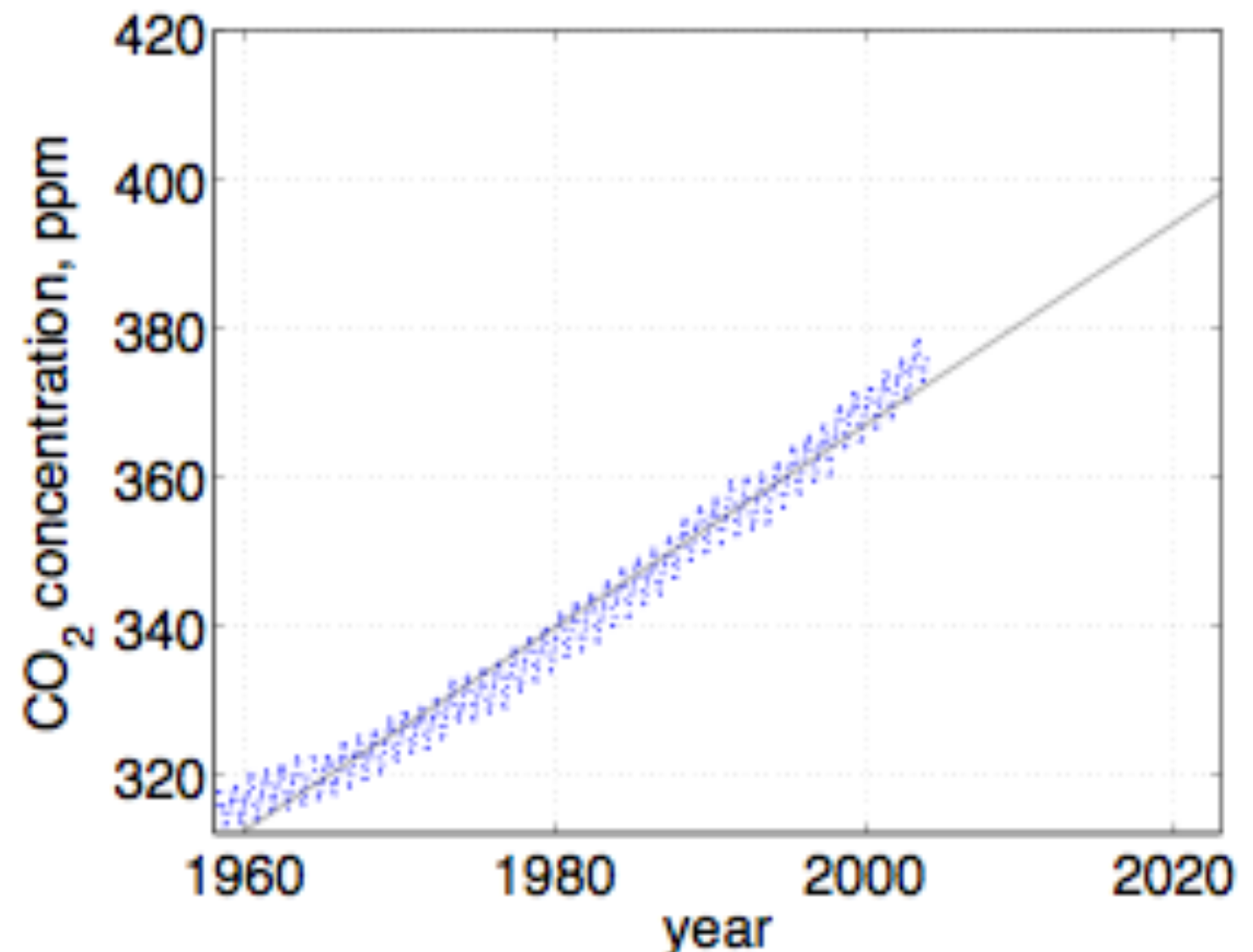
Prediction Problem (=Supervised Learning)



What will be the CO₂ concentration in the future?

Different prediction models possible

➔ Linear



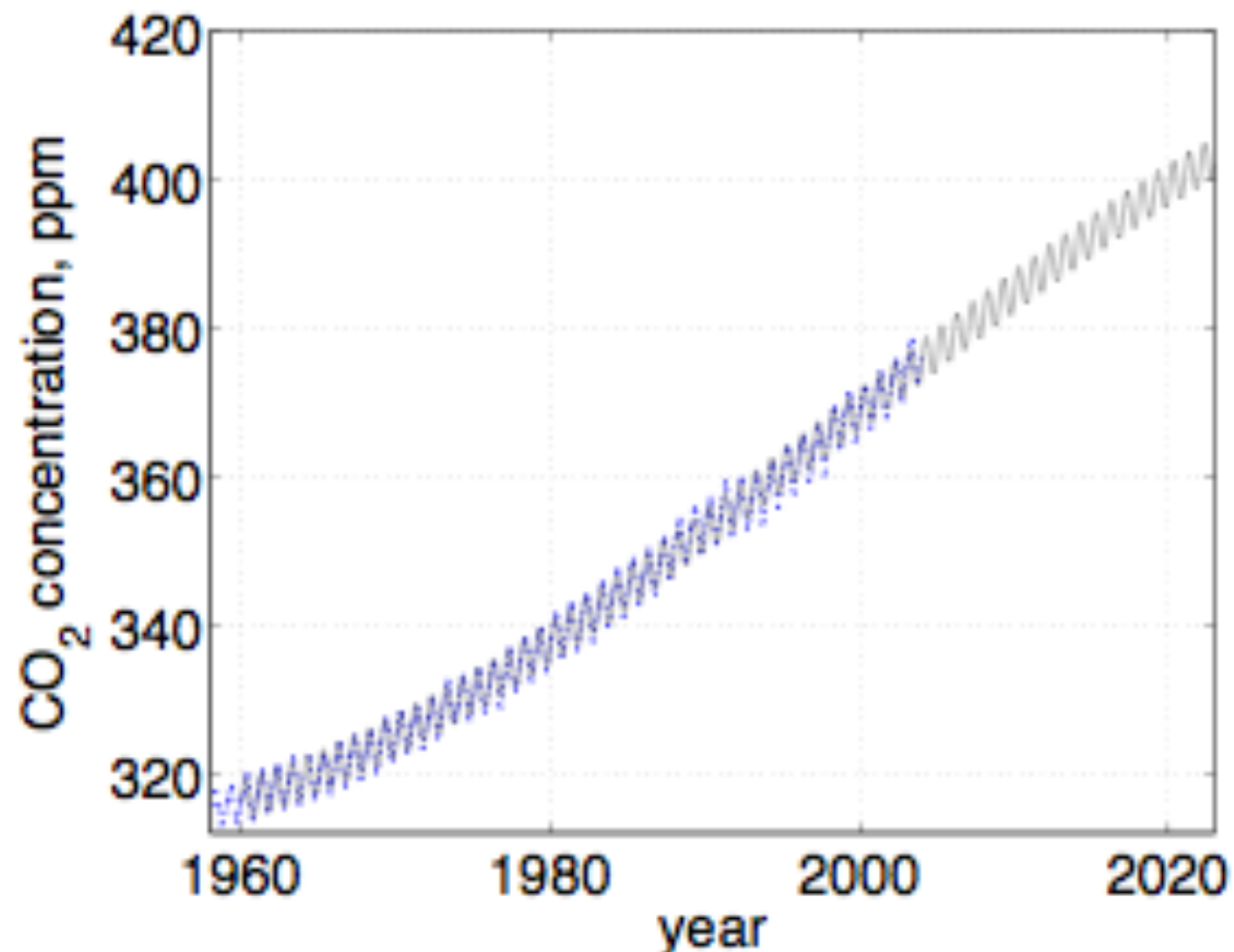
Prediction Problem (=Supervised Learning)



What will be the CO₂ 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

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

$x \dots$ inputs, $y \dots$ output / target

Two most prominent examples are:

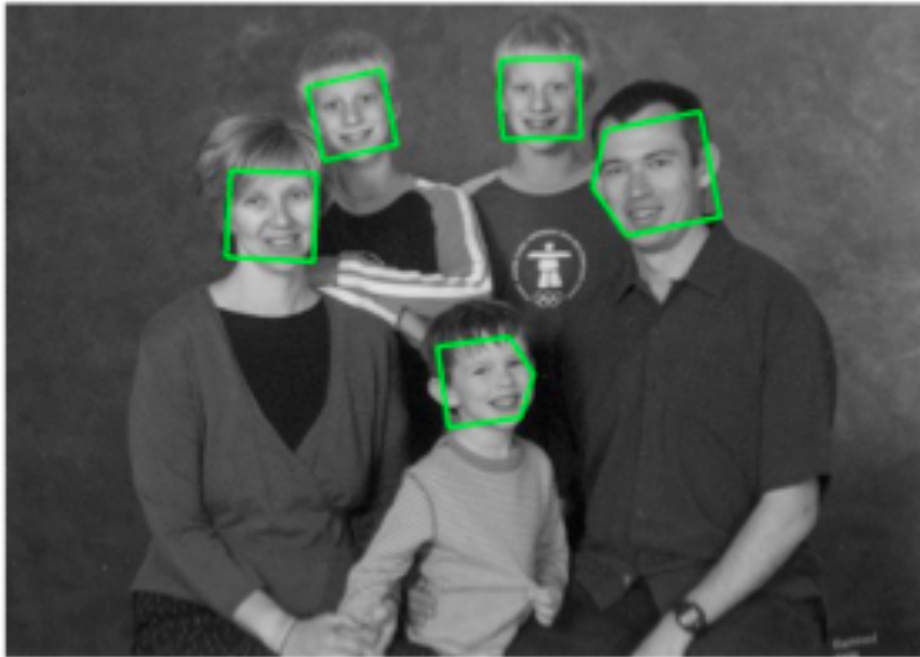
1. *Classification*: Discrete outputs or labels. $y_q \in \{0, 1\}$.

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

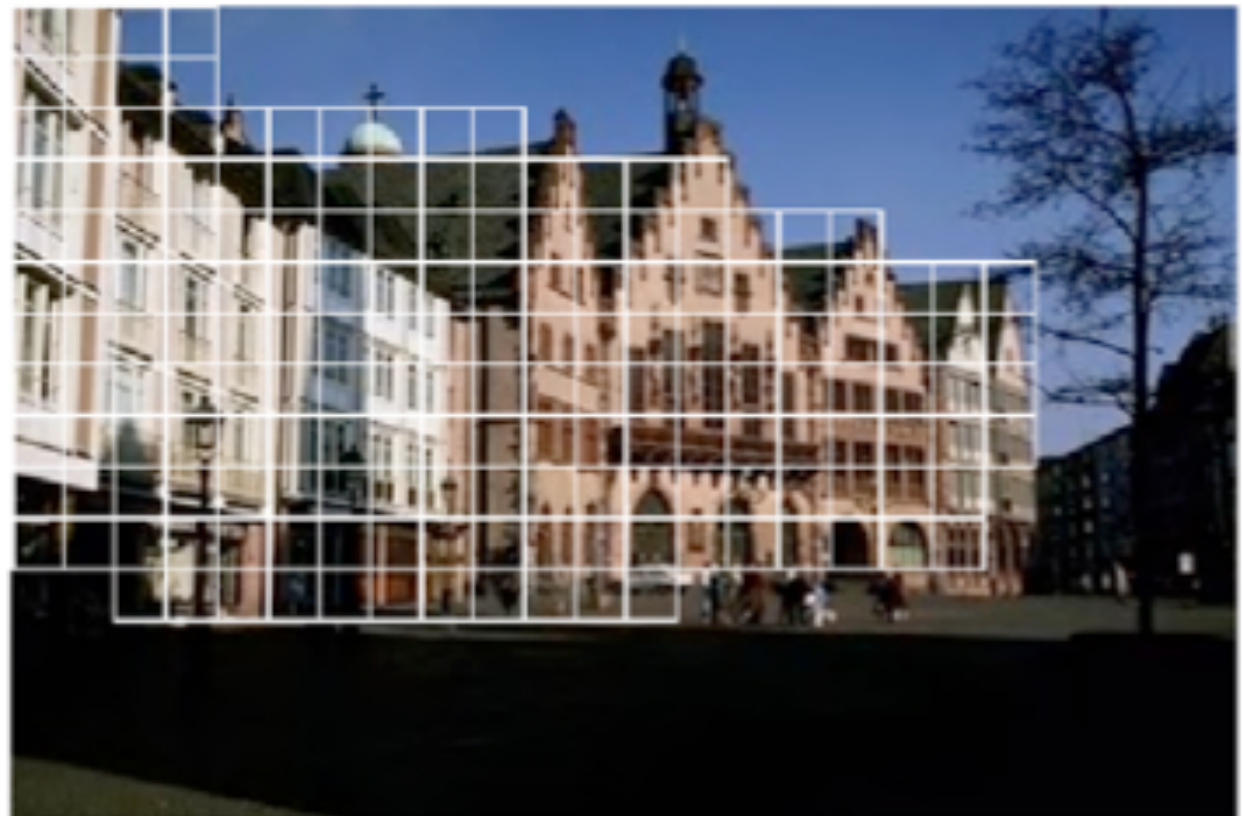
2. *Regression*: Continuous outputs or labels. $y_q \in \mathbb{R}$.

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

Examples of Classification



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

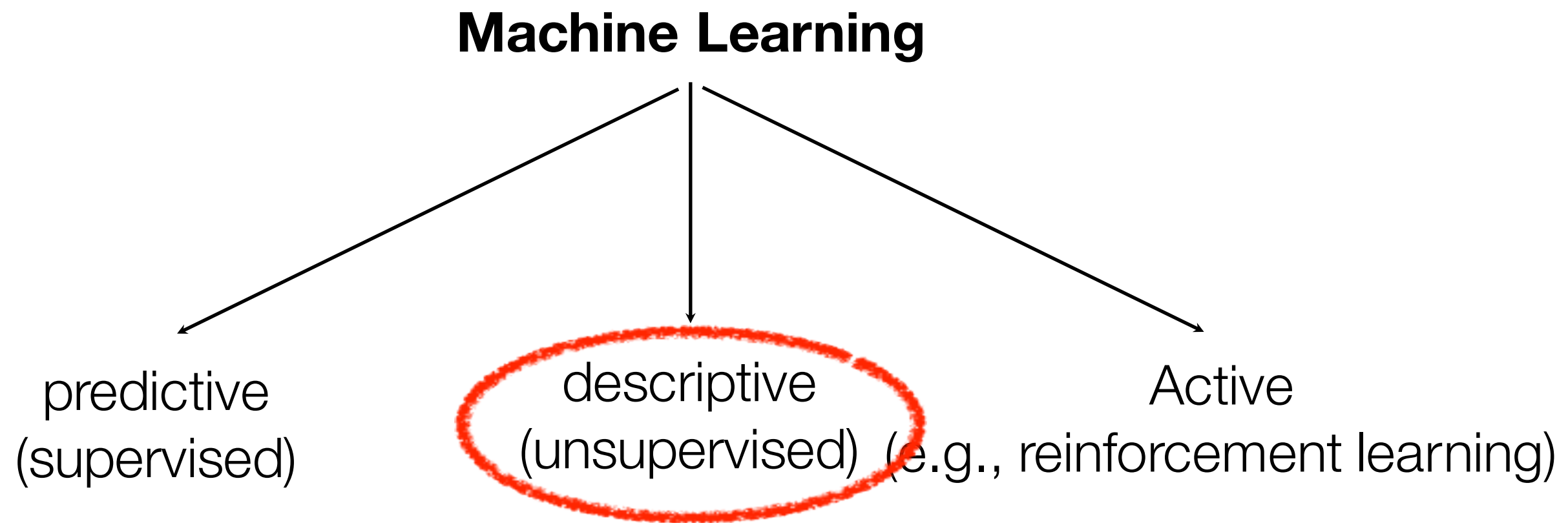


Examples of Regression



- 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



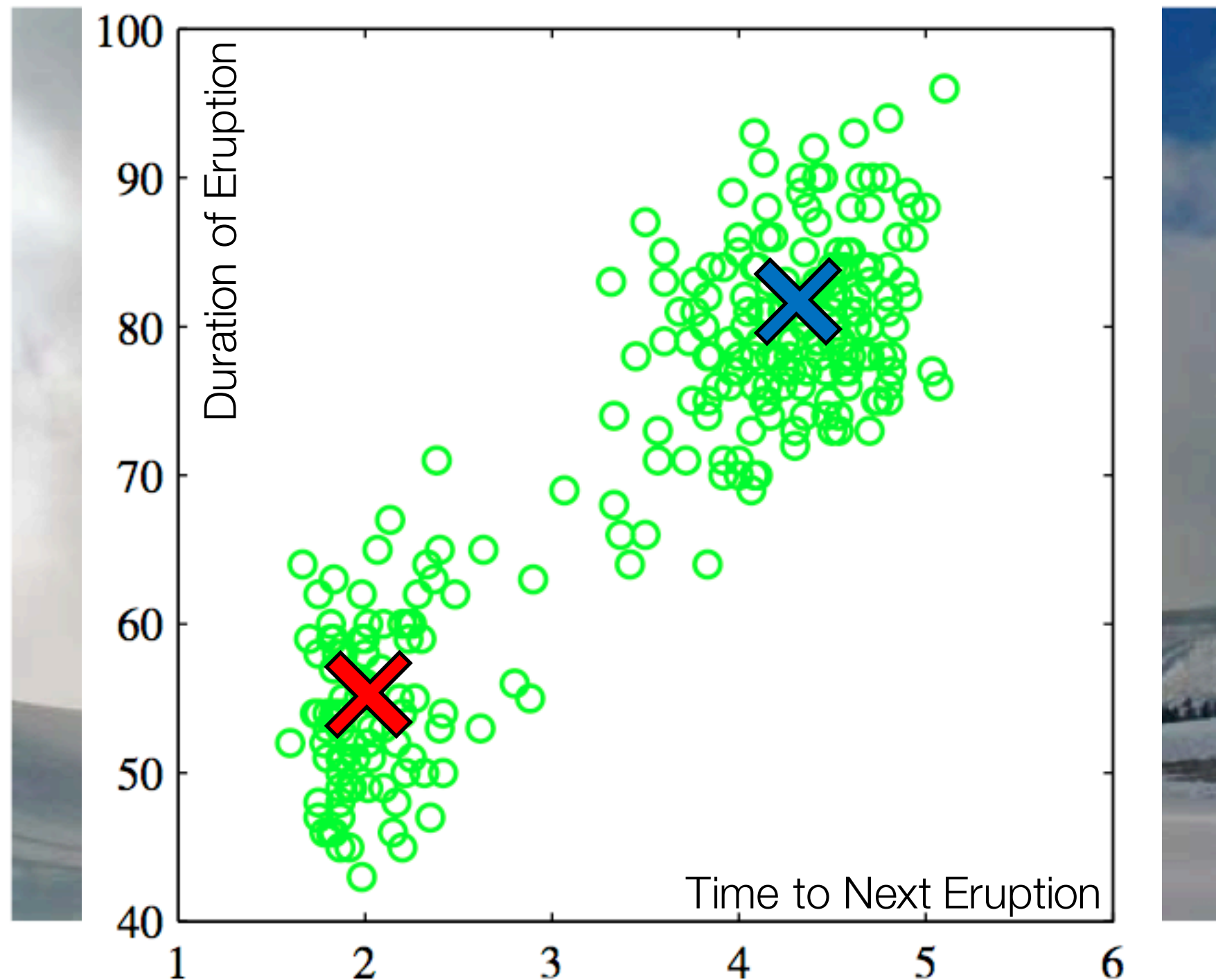
In descriptive problems, we have

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

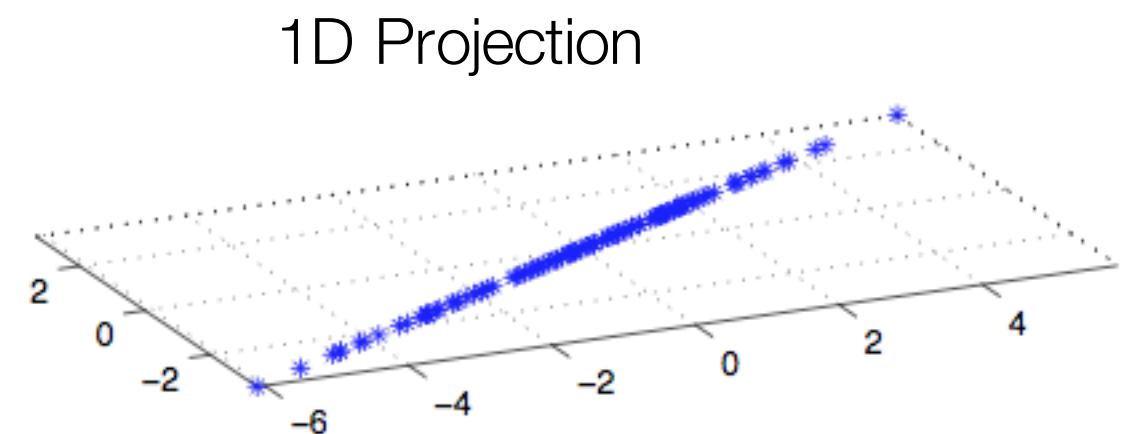
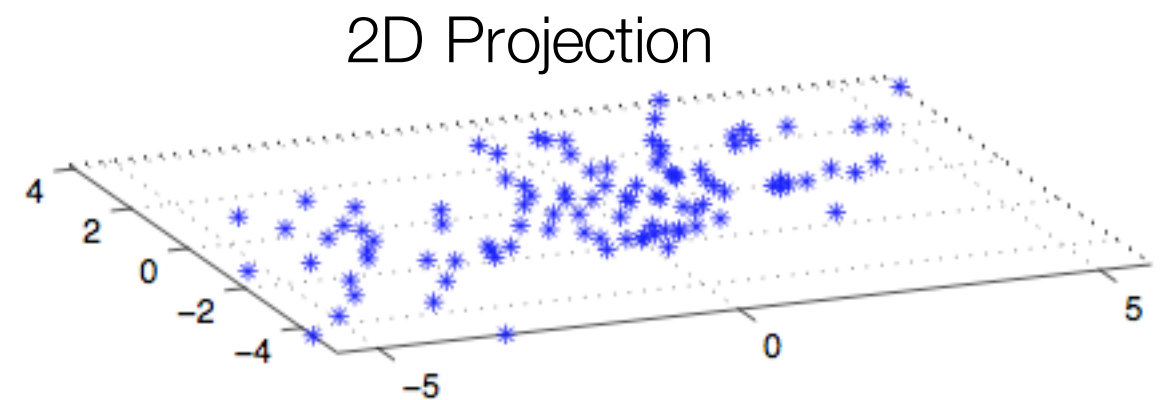
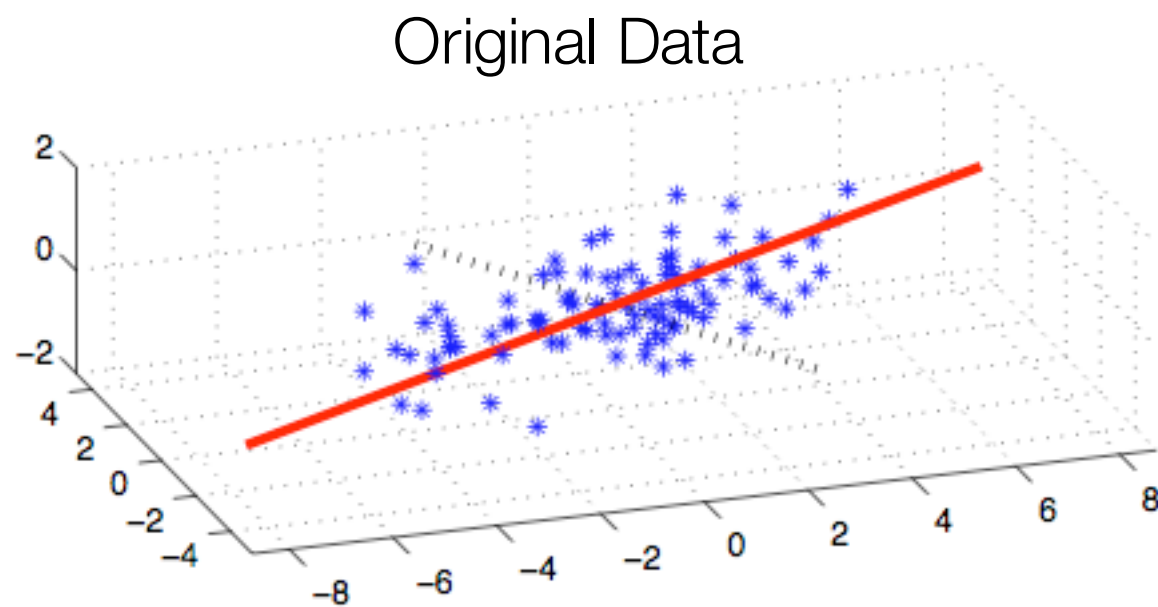
Three prominent examples are:

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...

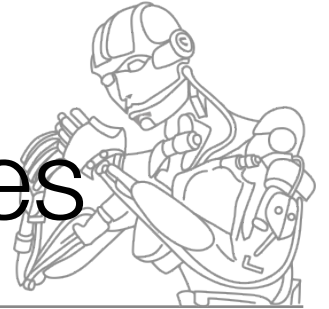
Old Faithful



Dimensionality Reduction



Dimensionality Reduction Example: Eigenfaces



How many faces do you need to characterize these?

Example: Eigenfaces



mean



principal basis 1



principal basis 2



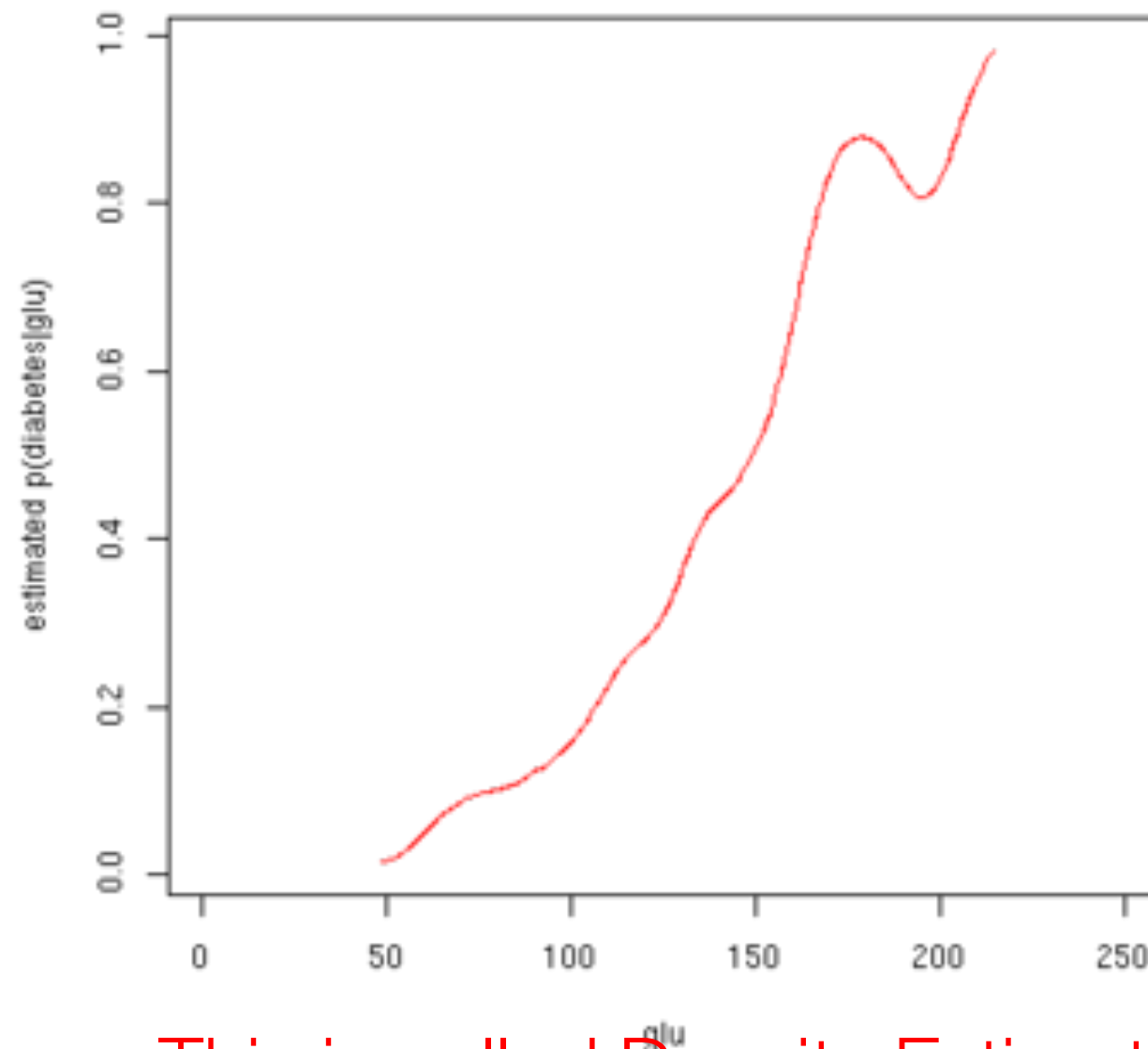
principal basis 3



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



Estimate relative occurrence of a data point



This is called Density Estimation!

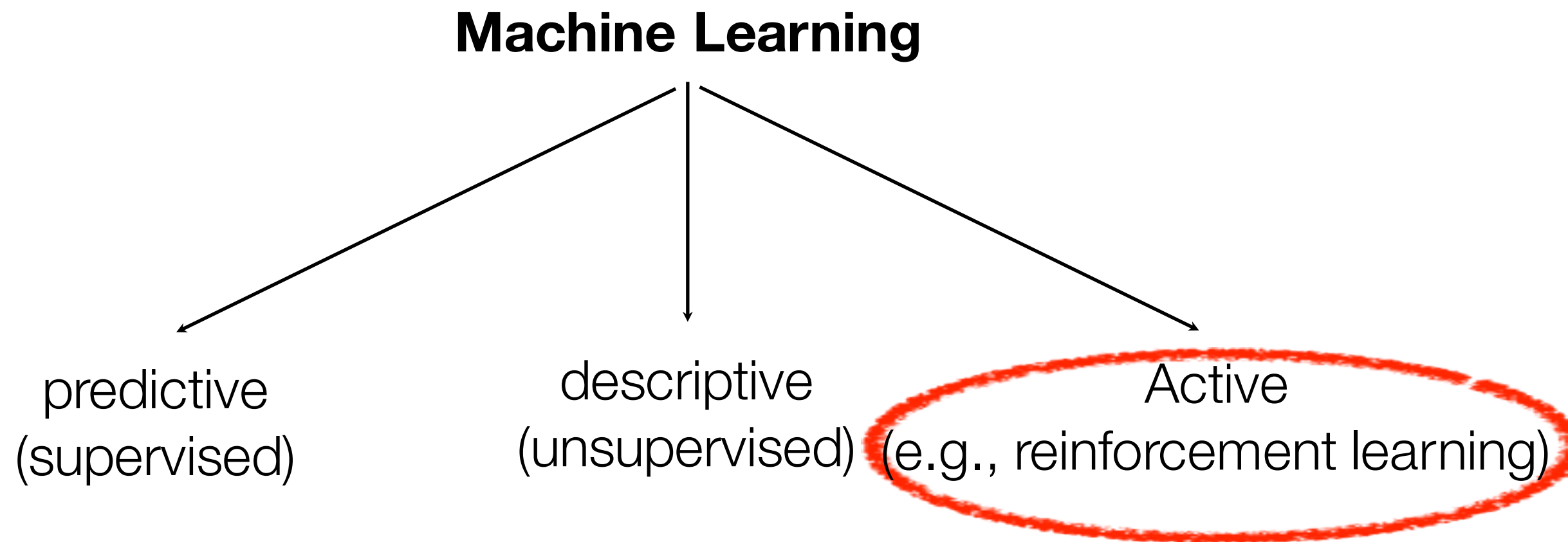
The bigger picture...



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?

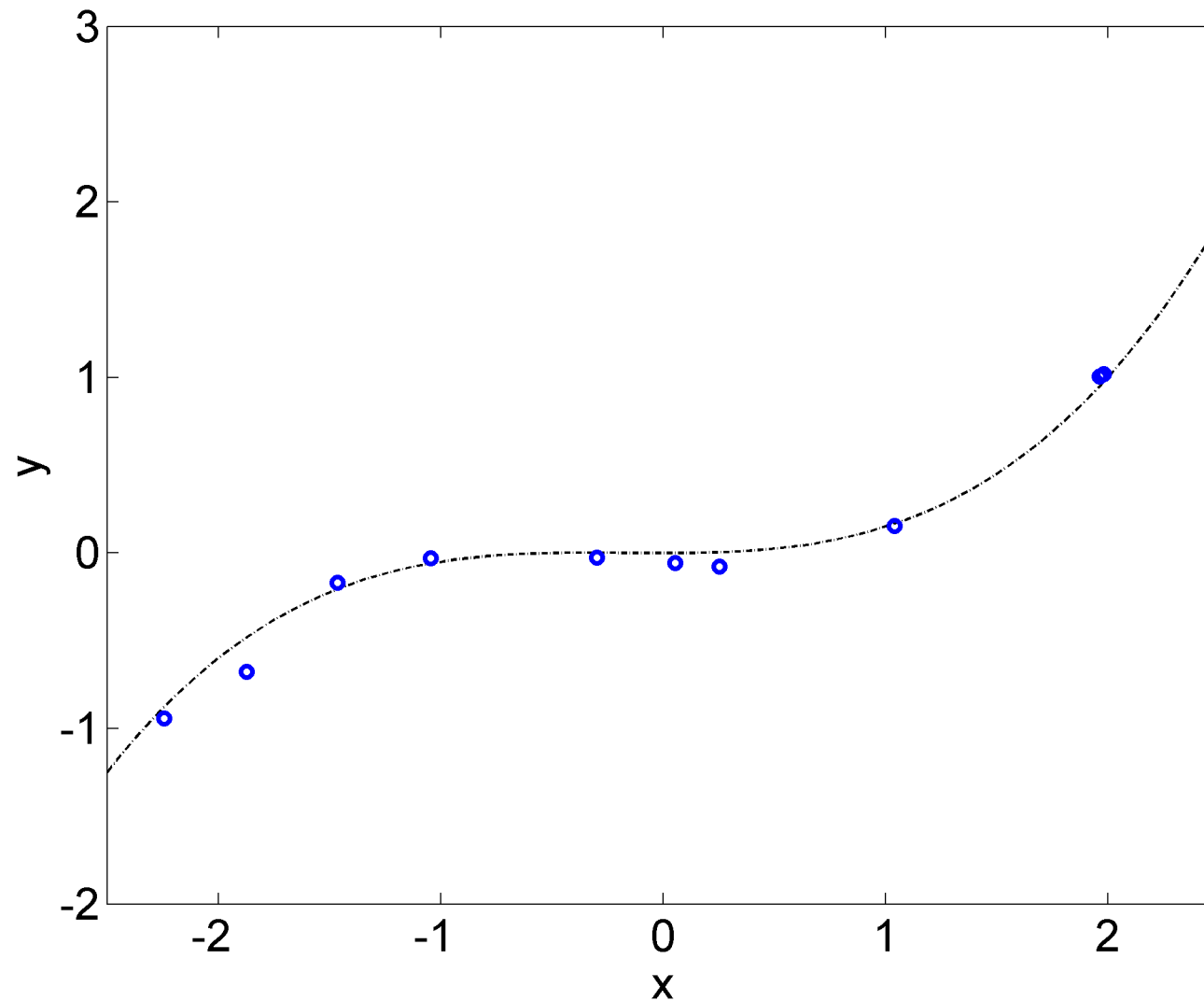
Important Questions



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 in simple: **linear in Features**

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

Important Questions



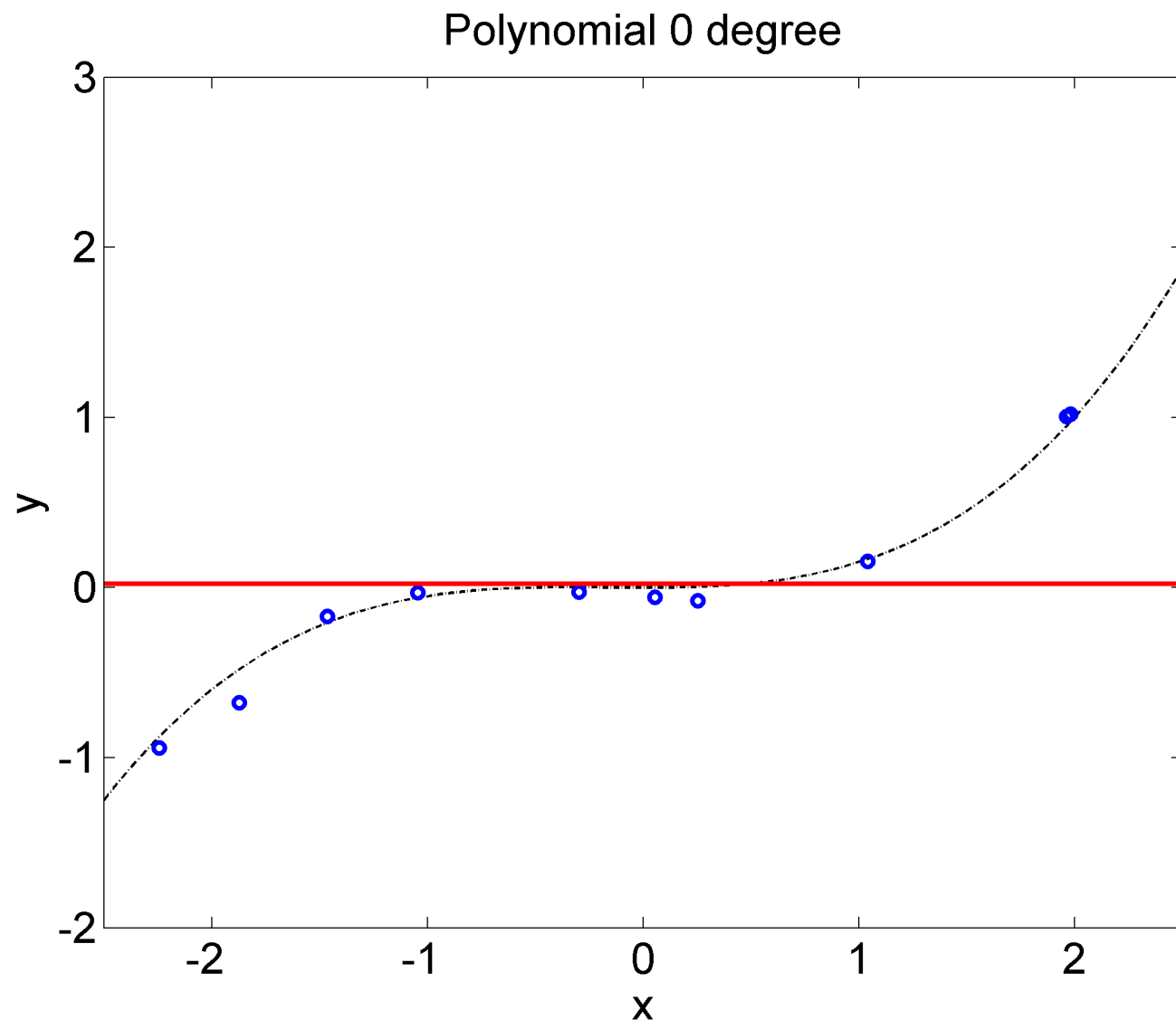
How does the data look like?

- 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 ...



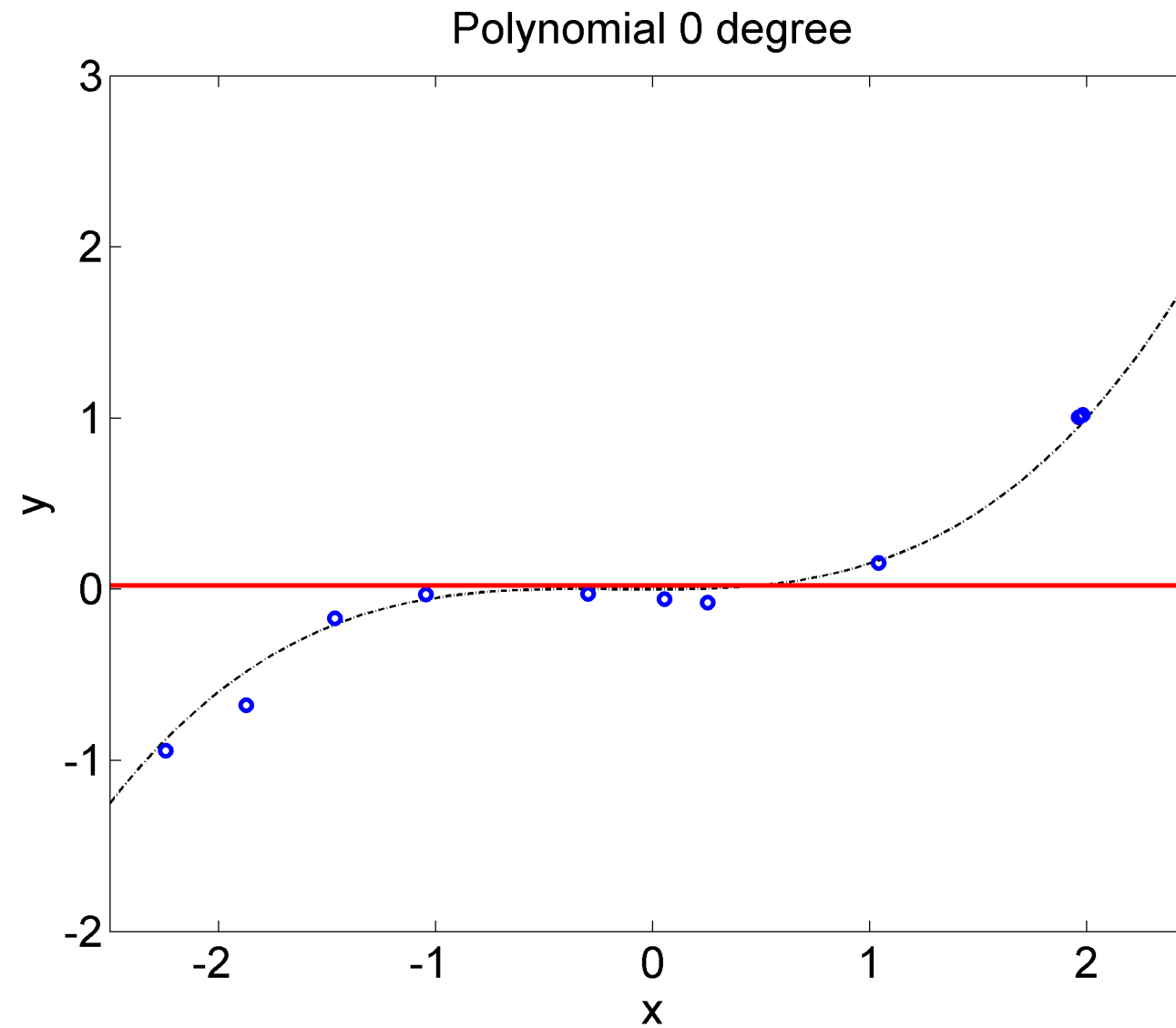
We need to answer:

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

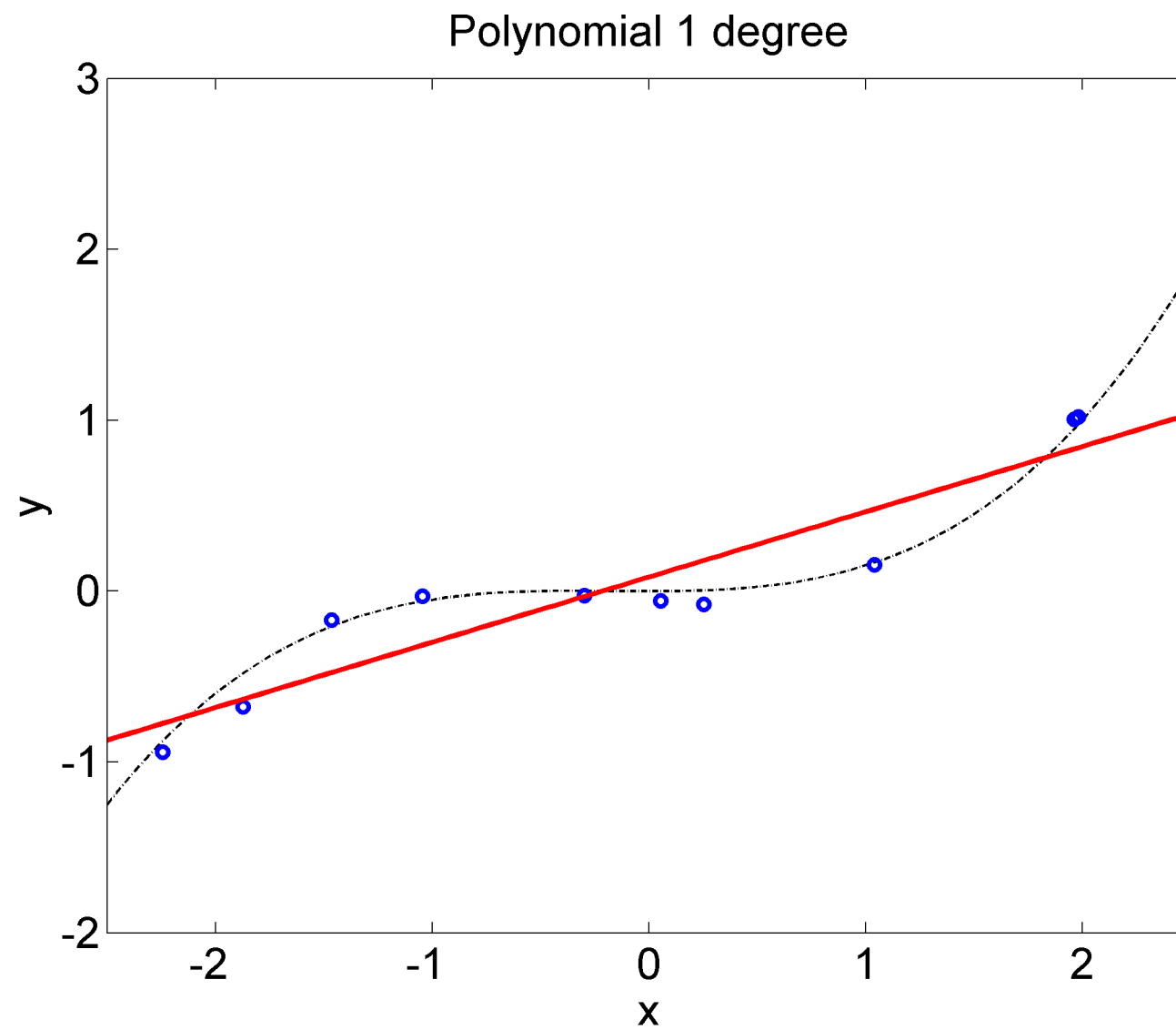
We assume a model class: polynomials of degree n

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

Fitting an Easy Model: $n=0$



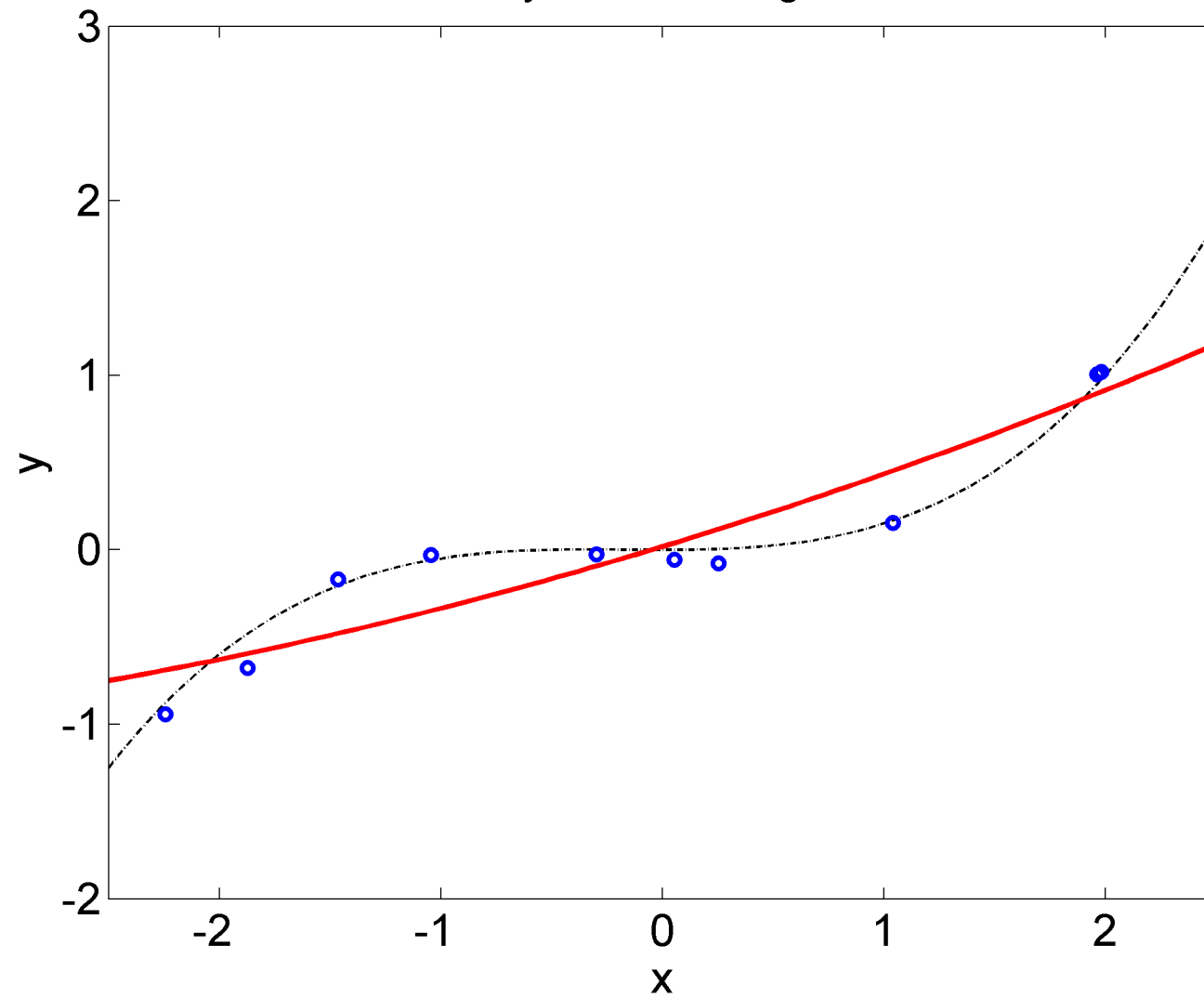
Add a Feature: $n=1$



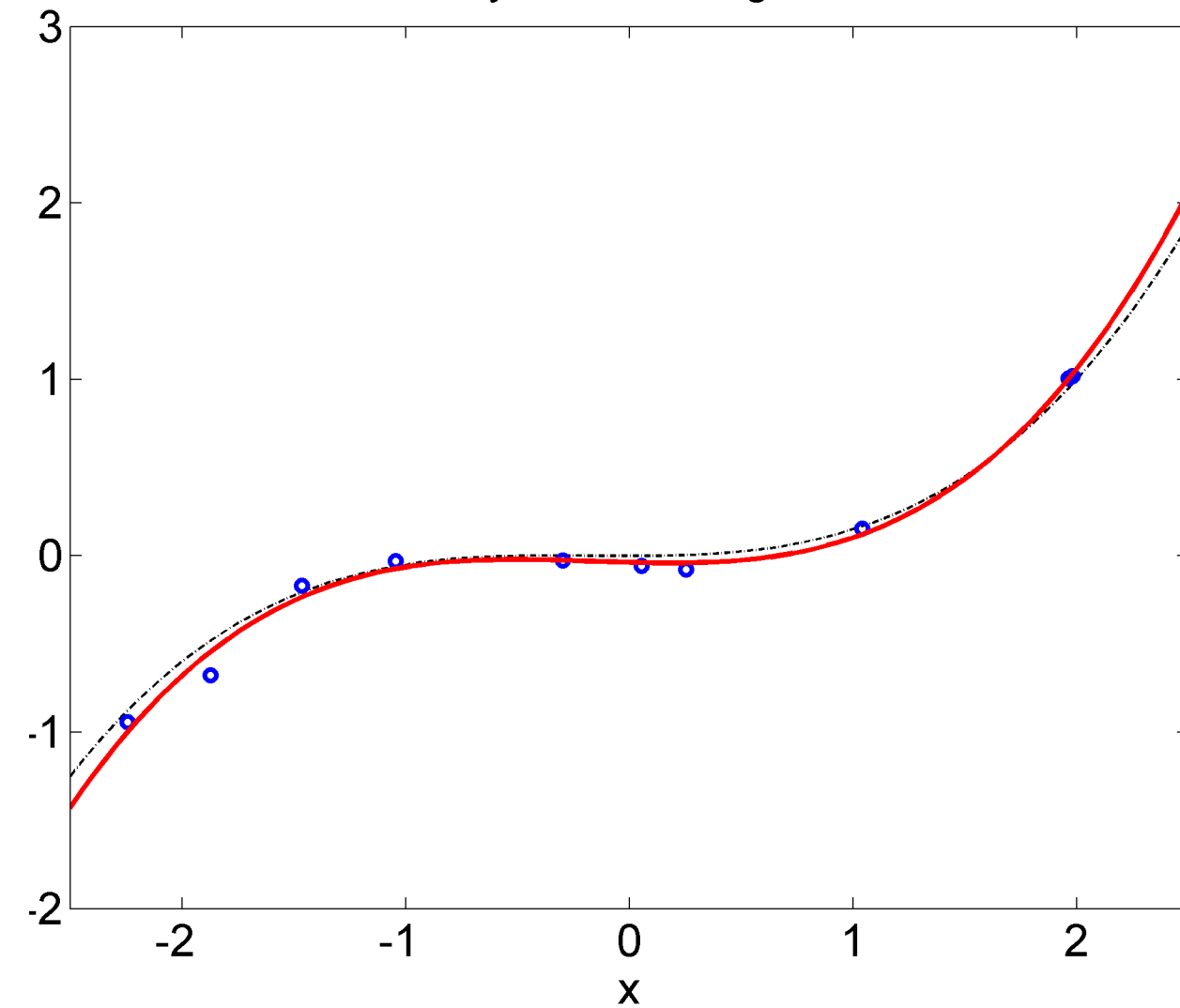
More features...



Polynomial 2 degree



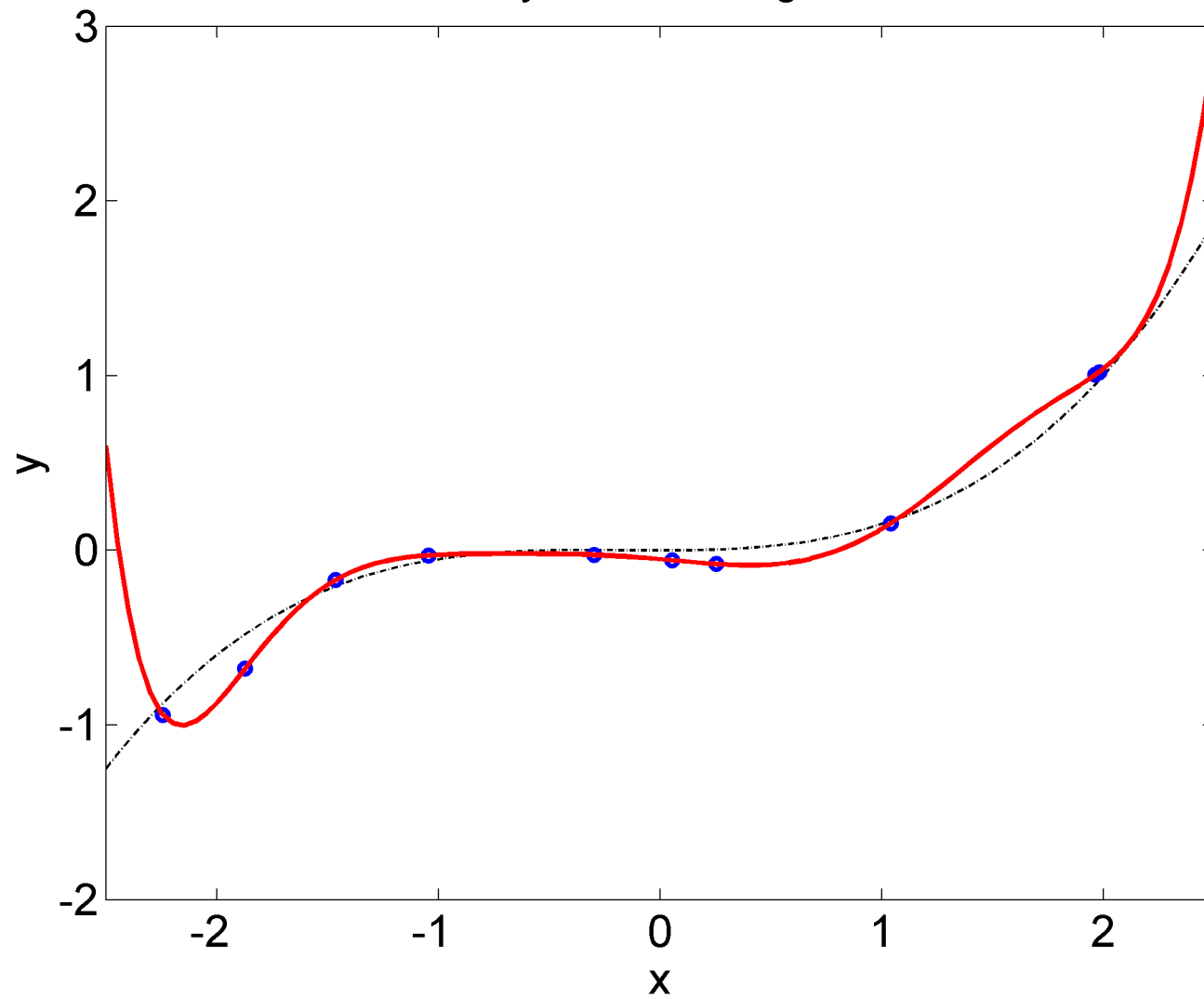
Polynomial 3 degree



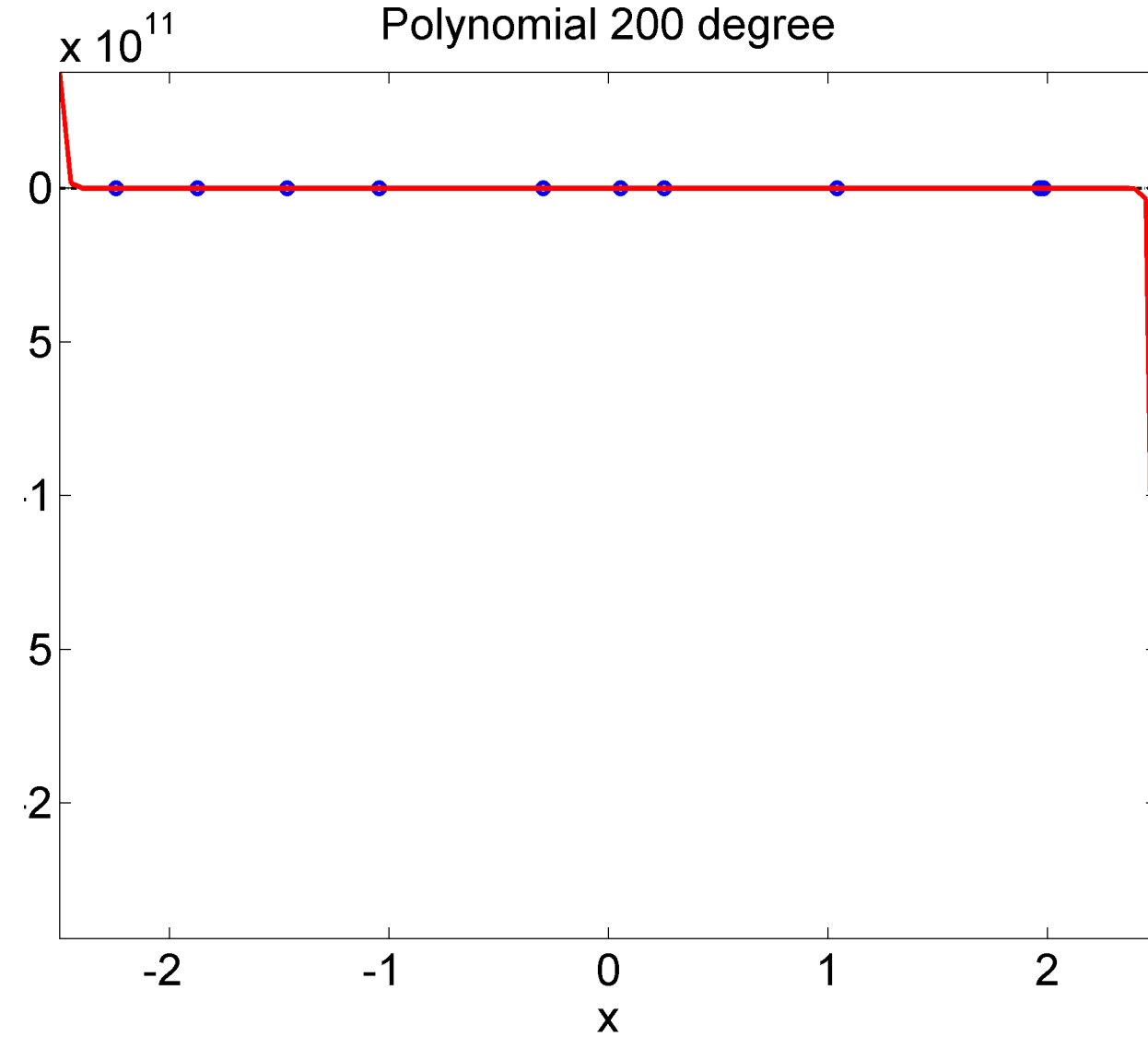
More features...



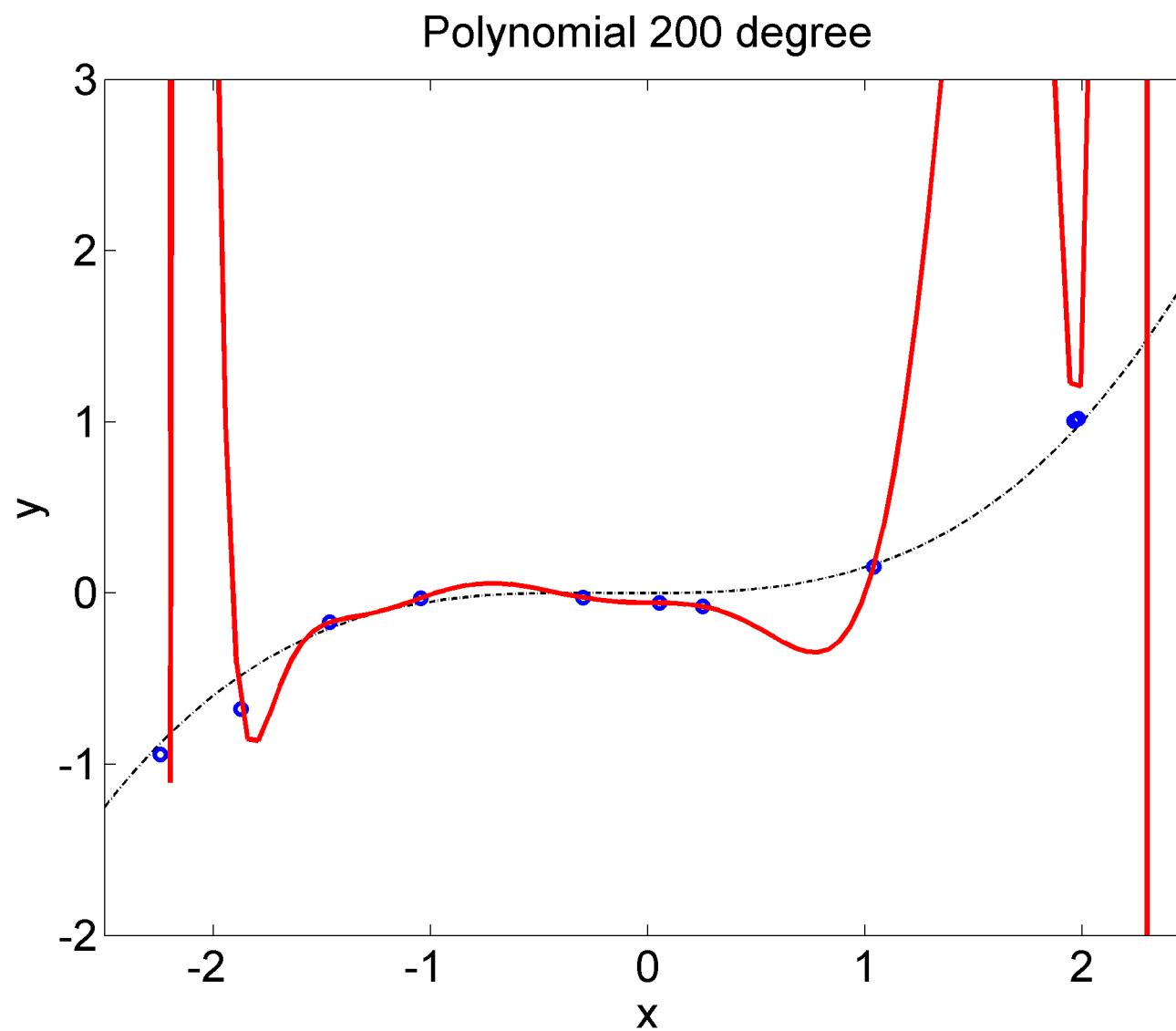
Polynomial 9 degree



Polynomial 200 degree



More features: $n=200$ (zoomed in)



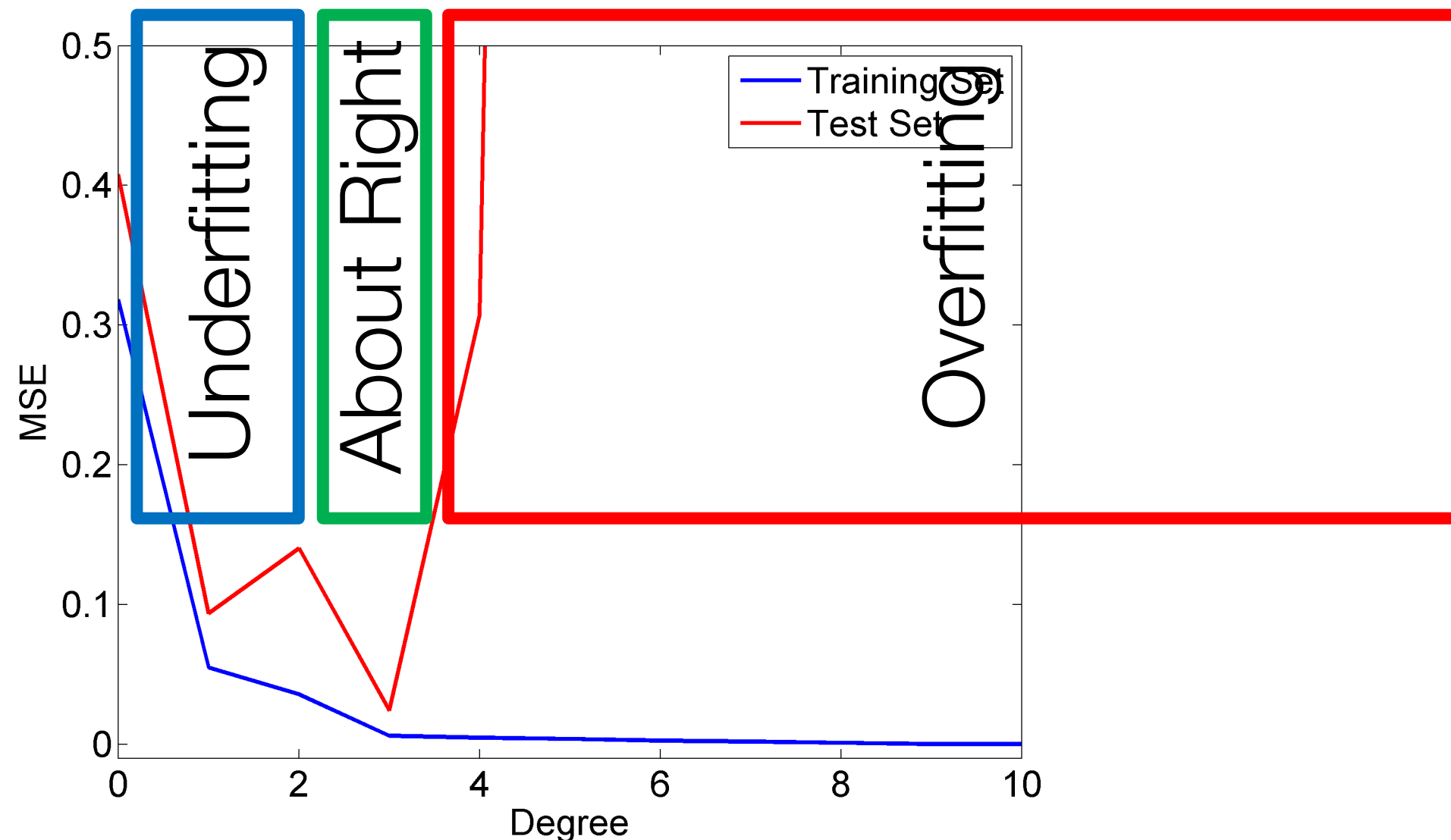
overfitting and numerical problems

Prominent example of overfitting...



Is there a tank in the picture?

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



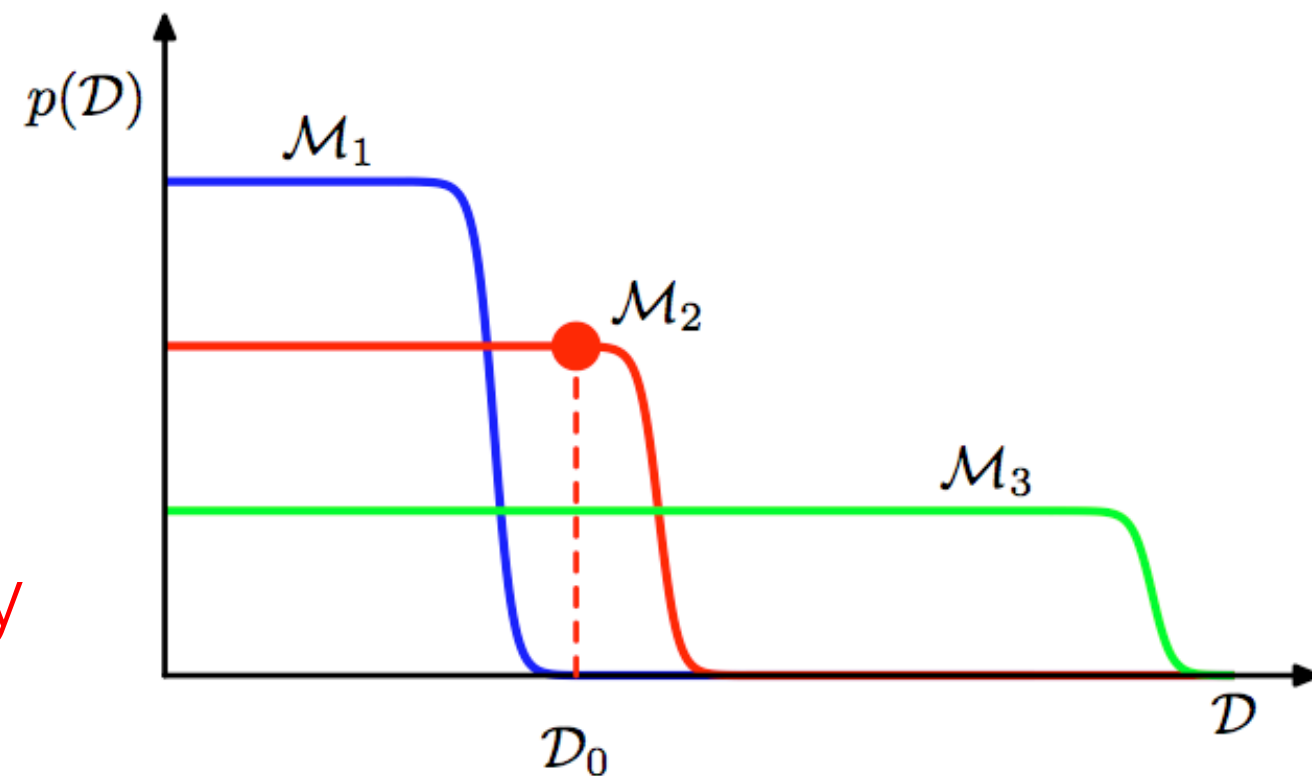
Model Selection: How can we...

- 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



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

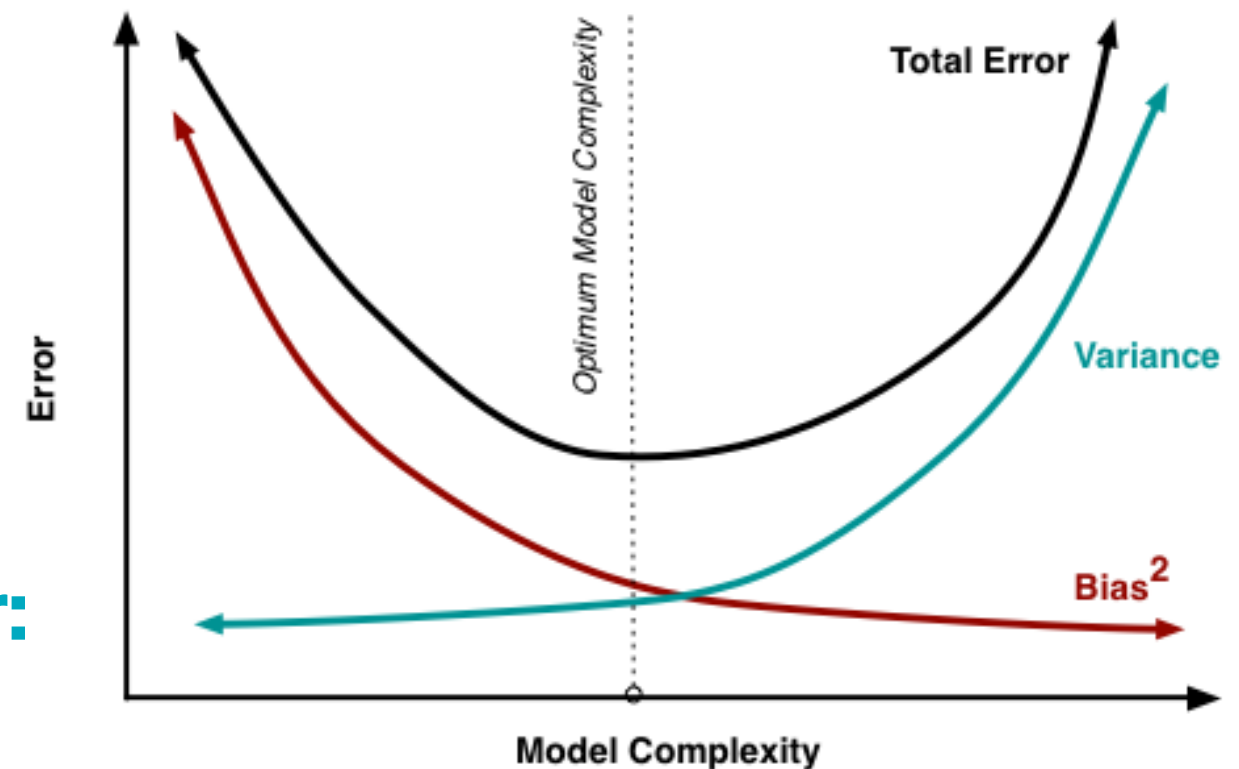
Typically, you can not minimize both!

Bias / Structure Error:

Error because our model can not do better

Variance / Approximation Error:

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

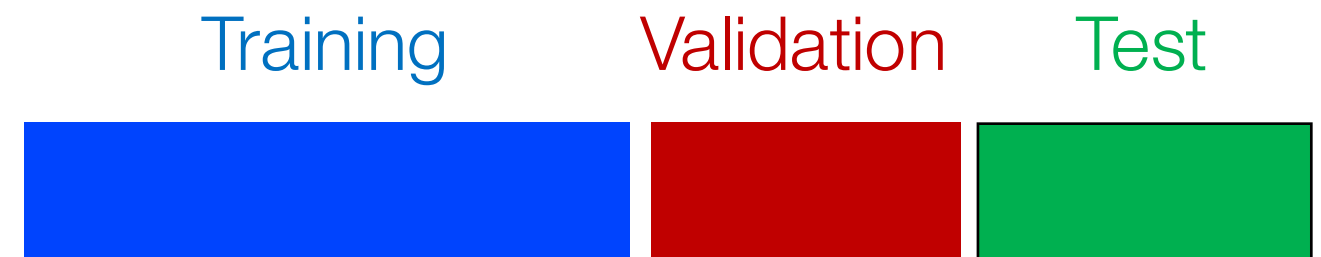


How do choose the model?



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!**

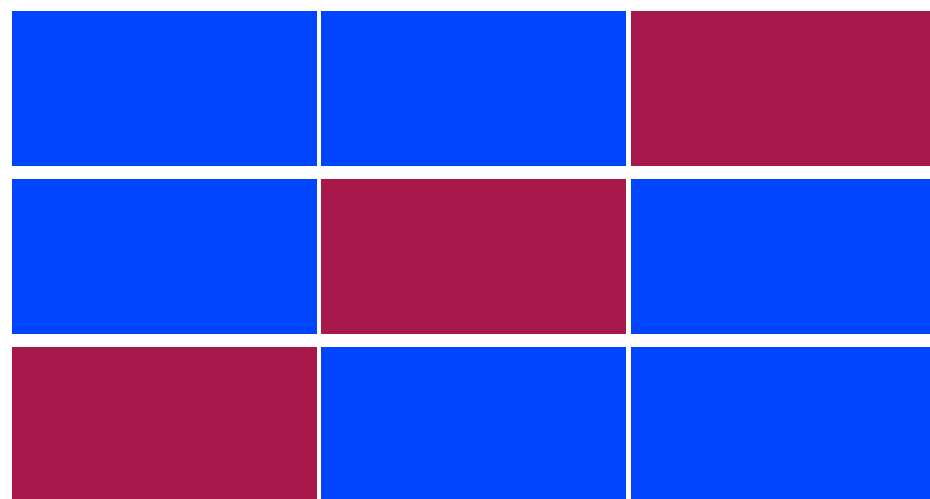
Model Selection: K-fold Cross Validation



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!!**



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

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

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

49 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 θ ?



Gauss

Let's find parameters through a cost function!



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

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

$$\begin{aligned} \Phi &= \begin{bmatrix} \phi(\mathbf{x}_1), & \phi(\mathbf{x}_2), & \phi(\mathbf{x}_3), & \dots, & \phi(\mathbf{x}_n) \end{bmatrix}^T, \\ \mathbf{Y} &= \begin{bmatrix} y_1, & y_2, & y_3, & \dots, & y_n \end{bmatrix}^T. \end{aligned}$$

we can rewrite it as

$$J = \frac{1}{2} (\mathbf{Y} - \Phi \boldsymbol{\theta})^T (\mathbf{Y} - \Phi \boldsymbol{\theta}).$$

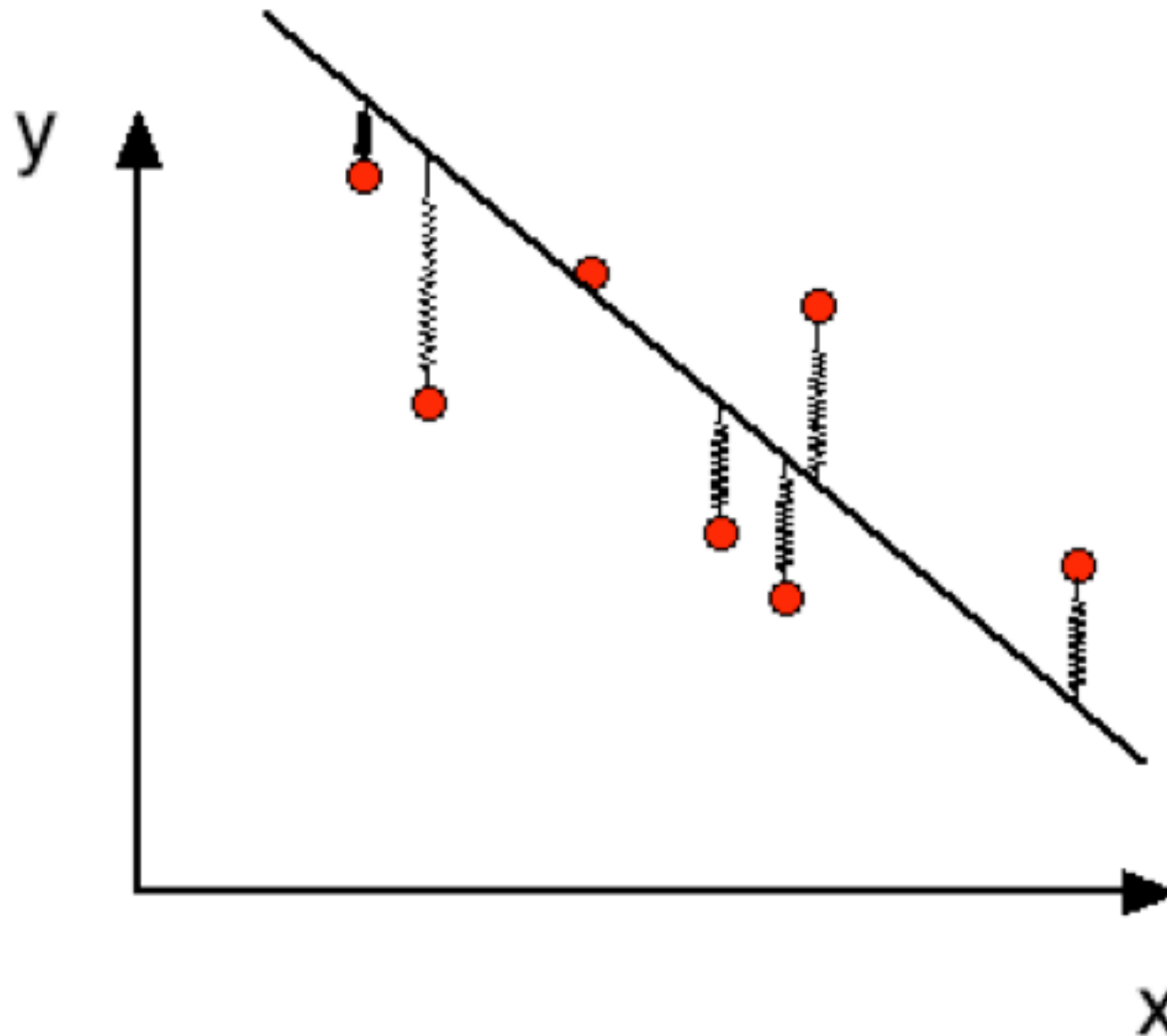
Scalar Product

and solve it

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

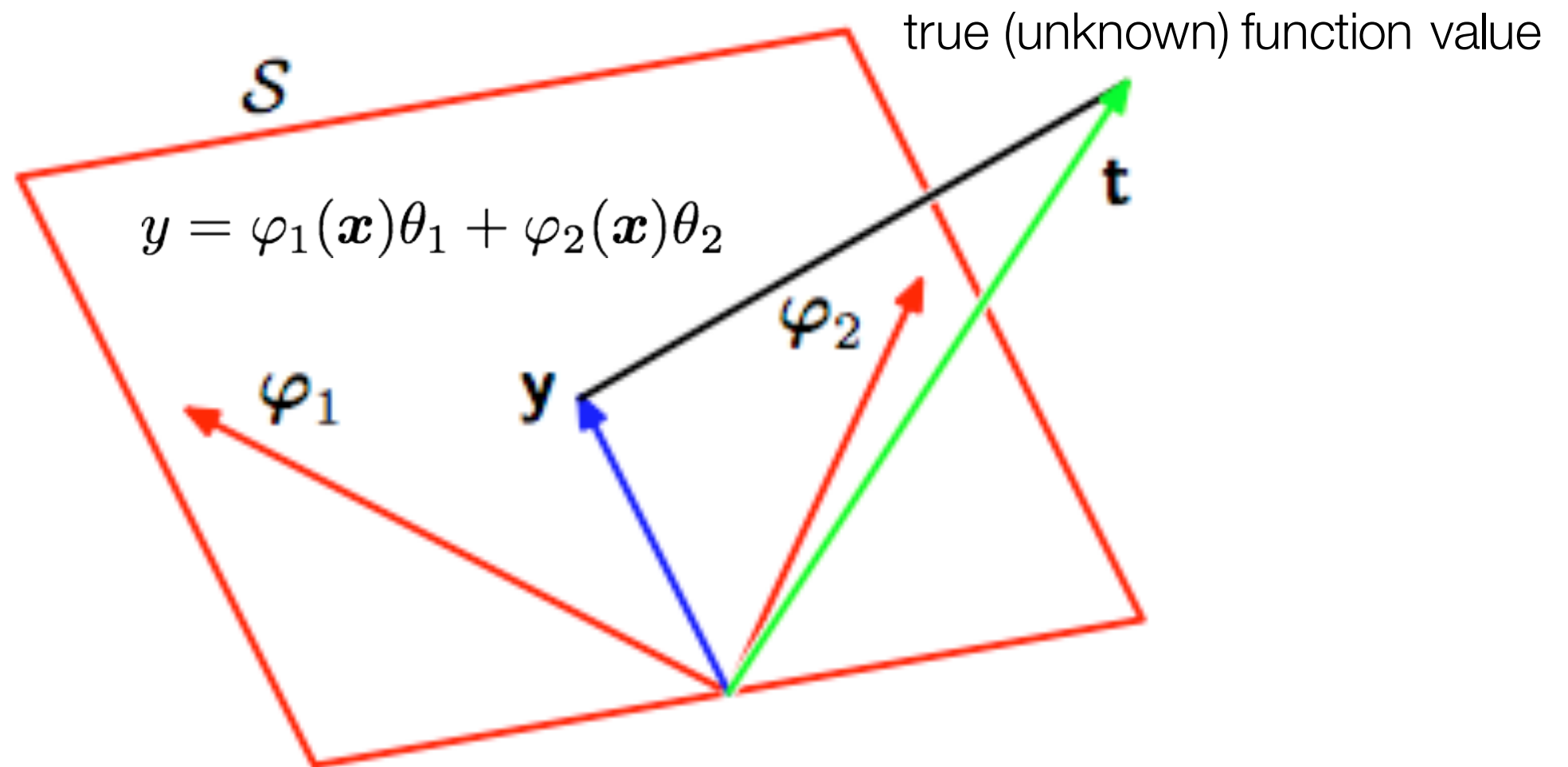
Least Squares solution contains left pseudo-inverse

Physical Interpretation



Energy of springs \sim squared lengths
➡ minimize energy of system

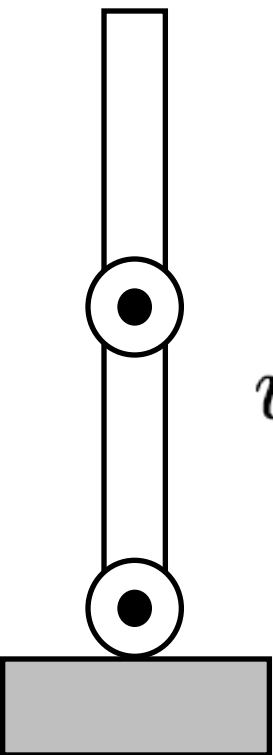
Geometric Interpretation



Minimize projection error \Rightarrow orthogonal projection

Robotics Example: Rigid-Body Dynamics

Known Features



$$\begin{aligned}
 u_1 = & [m_1 l_{g1}^2 + J_1 + m_2(l_1^2 + l_{g2}^2 + 2l_1 l_{g2} \cos \theta_2) + J_2] \ddot{\theta}_1 \\
 & + [m_2(l_{g2}^2 + l_1 l_2 \cos \theta_2) + J_2] \ddot{\theta}_2 && \text{Inertial Forces} \\
 & - 2m_2 l_1 l_{g2} \dot{\theta}_1 \dot{\theta}_2 \sin \theta_2 && \text{Coriolis Forces} \\
 & - 2m_2 l_1 l_{g2} \dot{\theta}_1^2 \sin \theta_2 && \text{Centripetal Forces} \\
 & + m_1 g l_{g1} \cos \theta_1 + m_2 g (l_1 \cos \theta_1 + l_{g2} \cos(\theta_1 + \theta_2)) \\
 u_2 = & [m_2(l_{g2}^2 + l_1 l_{g2} \cos \theta_2) + J_2] \ddot{\theta}_1 && \text{Gravity} \\
 & + (m_2 l_{g2}^2 + J_2) \ddot{\theta}_2 && \text{Inertial Forces} \\
 & - m_2 l_1 l_{g2} \dot{\theta}_1^2 \sin \theta_2 && \text{Centripetal Forces} \\
 & + m_2 g l_{g2} \cos(\theta_1 + \theta_2) && \text{Gravity}
 \end{aligned}$$

Robotics Example: Rigid-Body Dynamics



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

We can rewrite it as

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

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

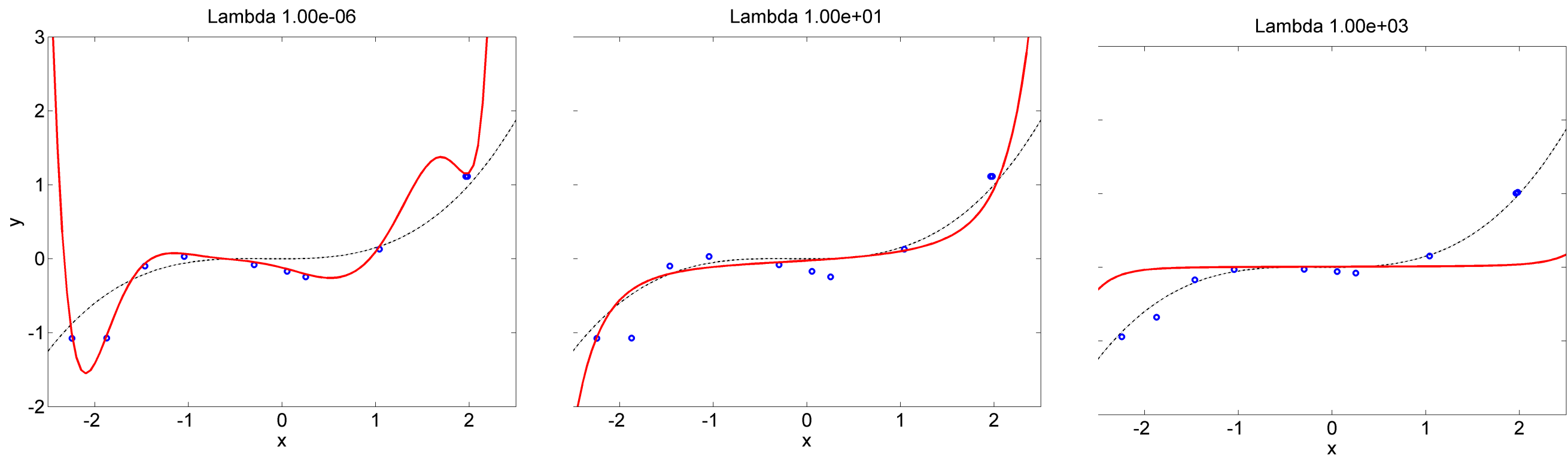
$$J_{\text{RR}} = (\mathbf{y} - \Phi\boldsymbol{\theta})^T (\mathbf{y} - \Phi\boldsymbol{\theta}) + \boldsymbol{\theta}^T \mathbf{W} \boldsymbol{\theta}$$

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

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

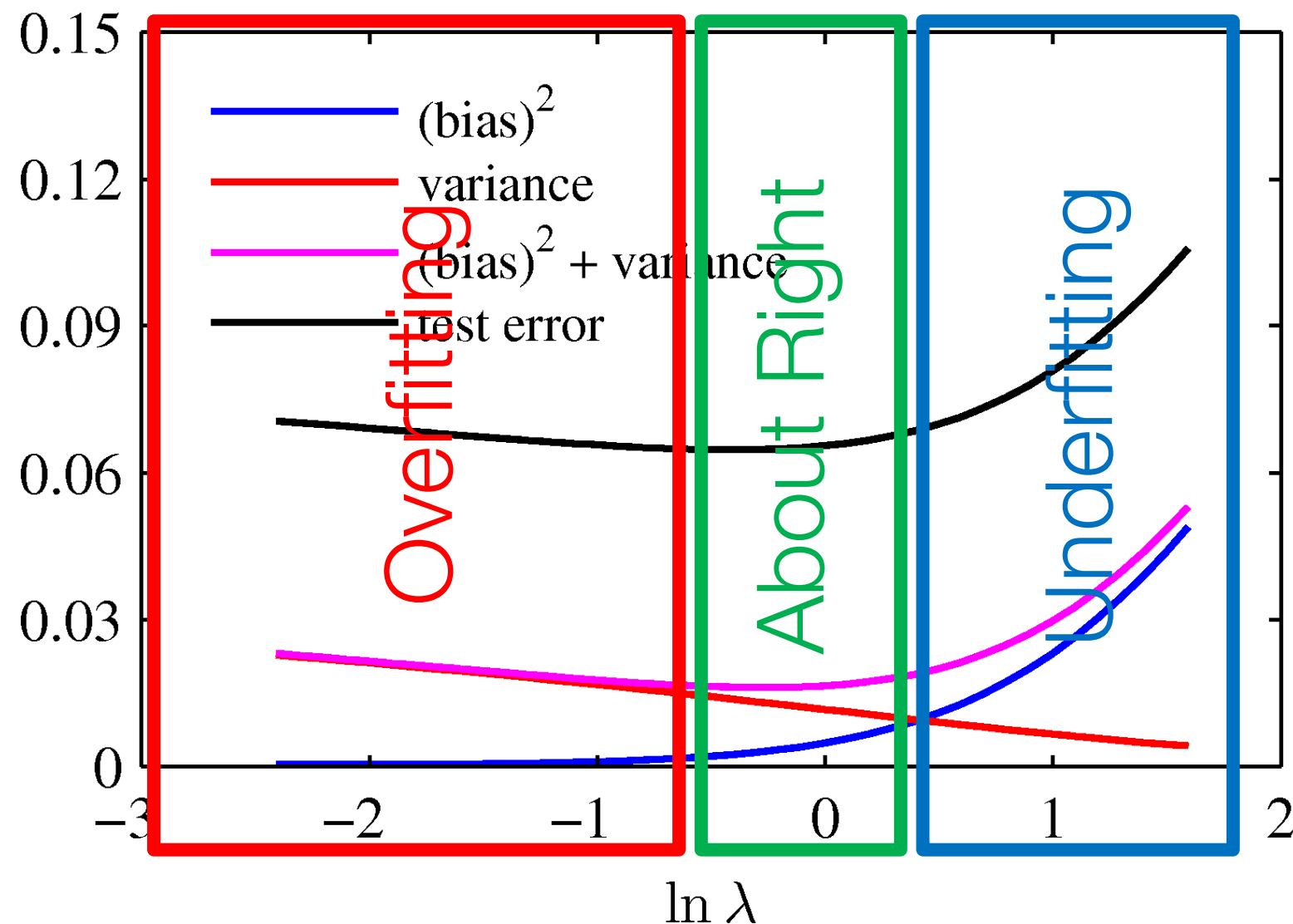
Numerically, this is much **more stable!** Even with **redundant features.**

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 θ ?



Fisher

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 $p(\mathbf{y}|\mathbf{X}, \theta)$.
- **We can obtain good parameters by maximizing likelihood of the outcome!**

Maximum-Likelihood (ML) estimate



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

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

That's hard!

Do the 'log-trick':

$$\begin{aligned} J_{\text{ML}} &= \arg \max_{\theta} \log p(\mathbf{y}|\mathbf{X}, \theta) = \arg \max_{\theta} \sum_{i=1}^N \log p(y_i|\mathbf{x}_i, \theta) \\ &= \arg \min_{\theta} \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_{\theta}(\mathbf{x}_i))^2 \end{aligned}$$

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?

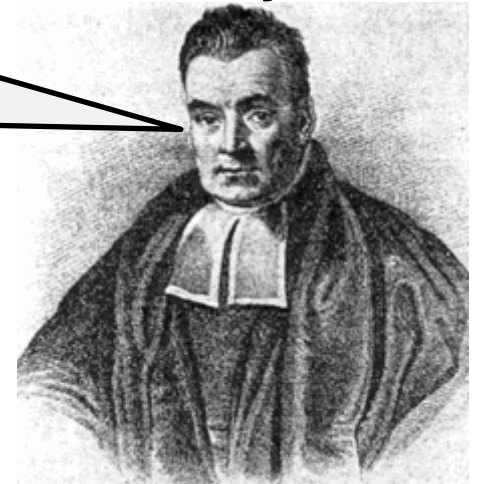
How to find the parameters θ ?



Bayesian:

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

Bayes



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

Diagram labels for the equation:

- likelihood: points to $p(y|X, \theta)$
- prior: points to $p(\theta)$
- evidence: points to $p(X|y)$
- posterior: points to $p(\theta|X, y)$

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**:

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

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

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

Ridge Regression is equivalent to MAP estimate with Gaussian prior

Predictions with the Model



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

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

$$p(y_* | \underset{\substack{\nearrow \\ \text{pred. function value}}}{x_*}, \underset{\substack{\nearrow \\ \text{test input}}}{\theta^*}) = \mathcal{N}(y_* | \mu(x_*), \sigma^2(x_*))$$

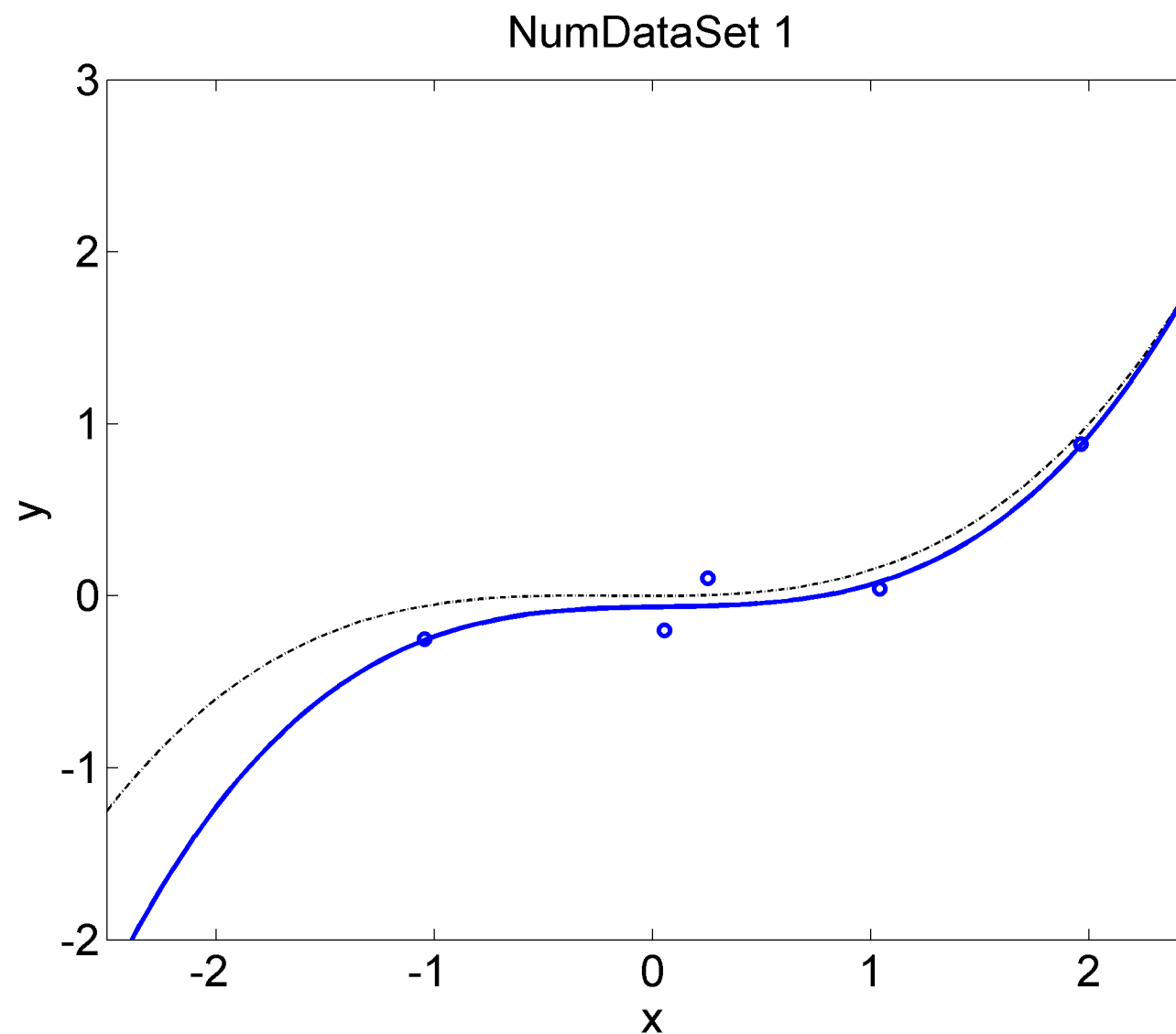
Predictive mean: $\mu(x_*) = \phi^T(x_*)\theta^*$

Predictive variance: $\sigma^2(x_*) = \sigma^2$

Comparing different data sets...



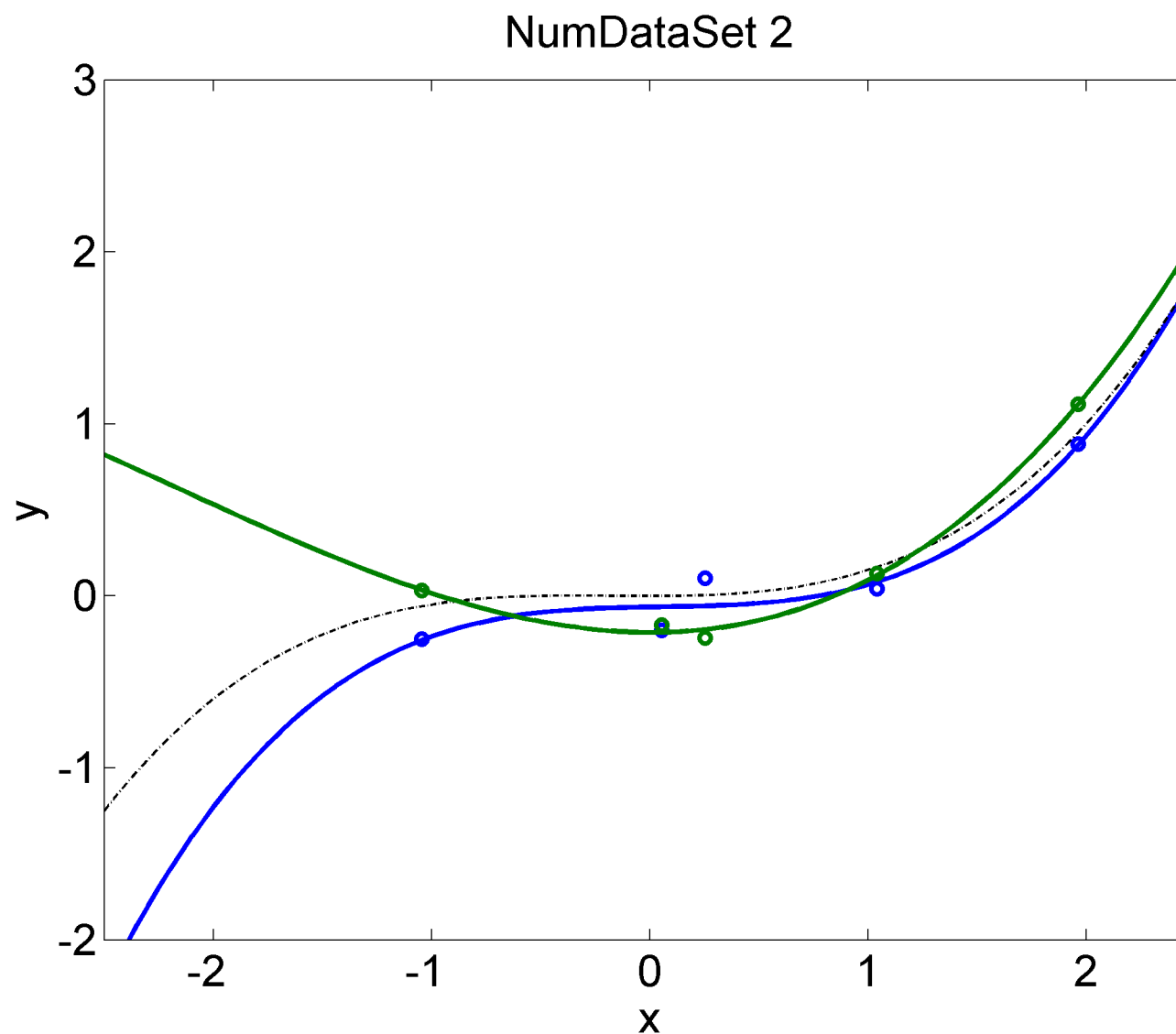
... with same input data, but different output values (due to noise):



Comparing different data sets...



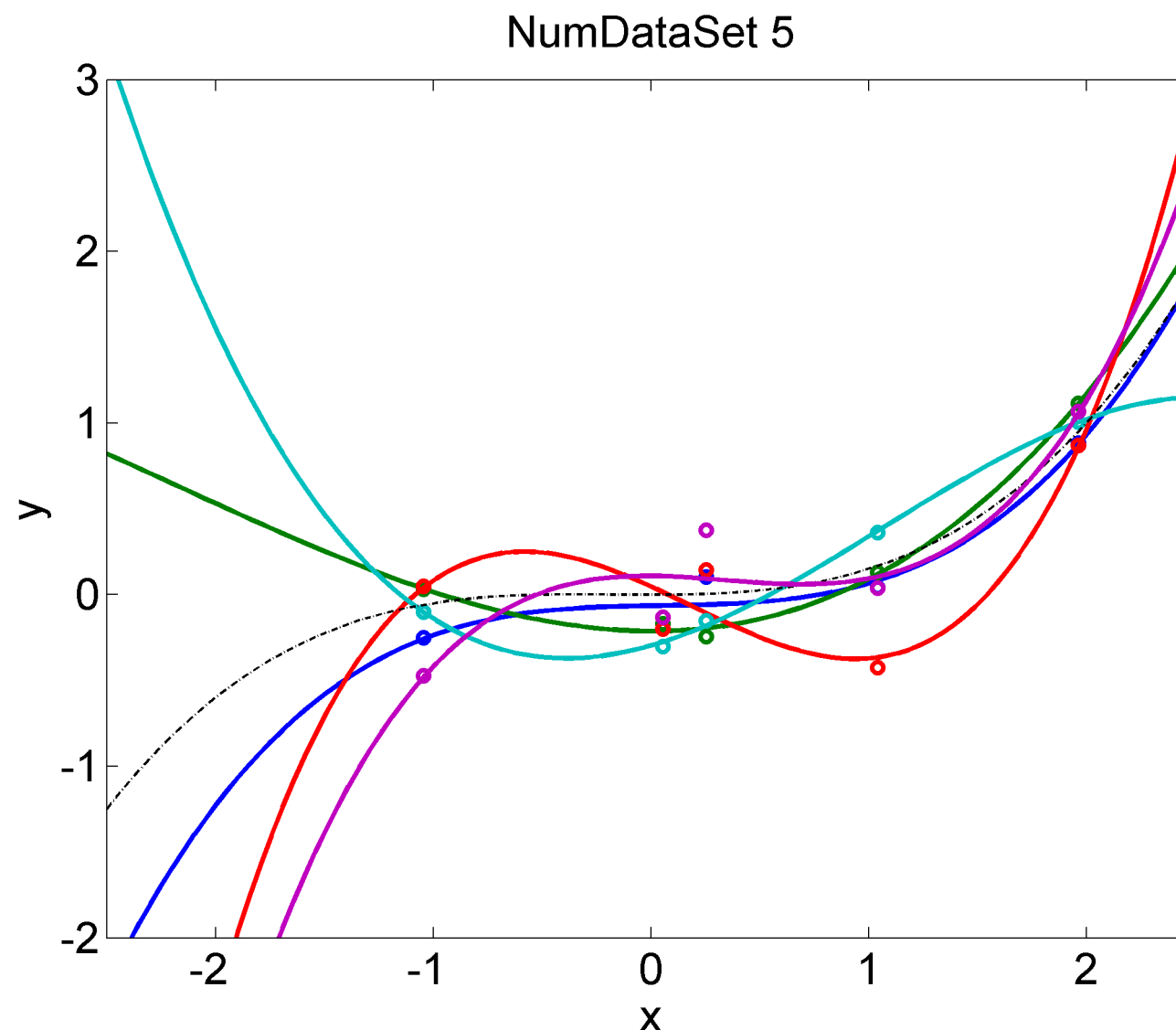
... with same input data, but different output values (due to noise):



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

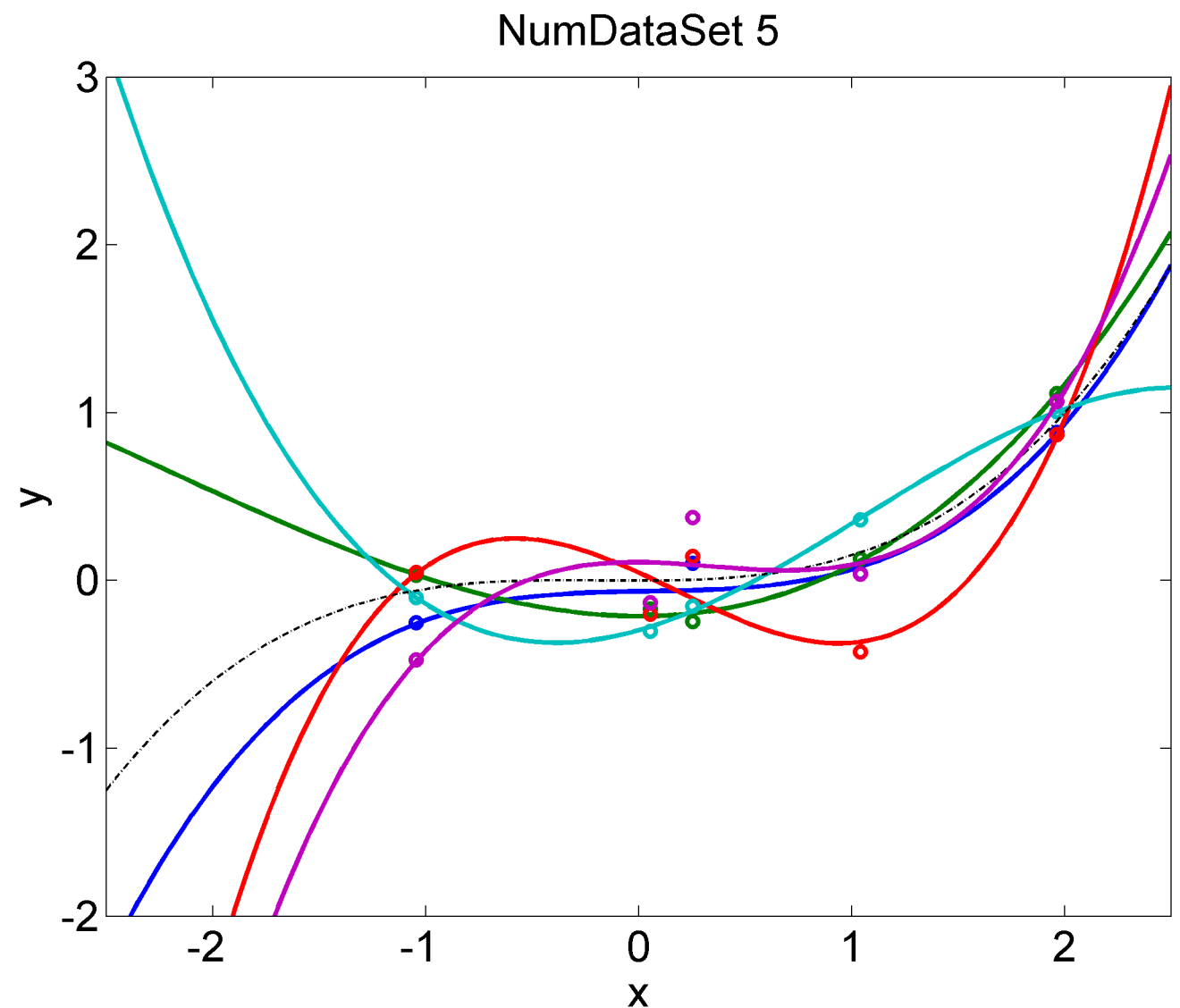
Comparing different data sets...



Can we also estimate our uncertainty in θ ?

Compute probability of θ given data

➔ $p(\theta | X, y)$



How to get the posterior?



Bayes Theorem for Gaussians

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

For our model:

Prior over parameters:

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

Data Likelihood

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

Posterior over parameters:

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

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

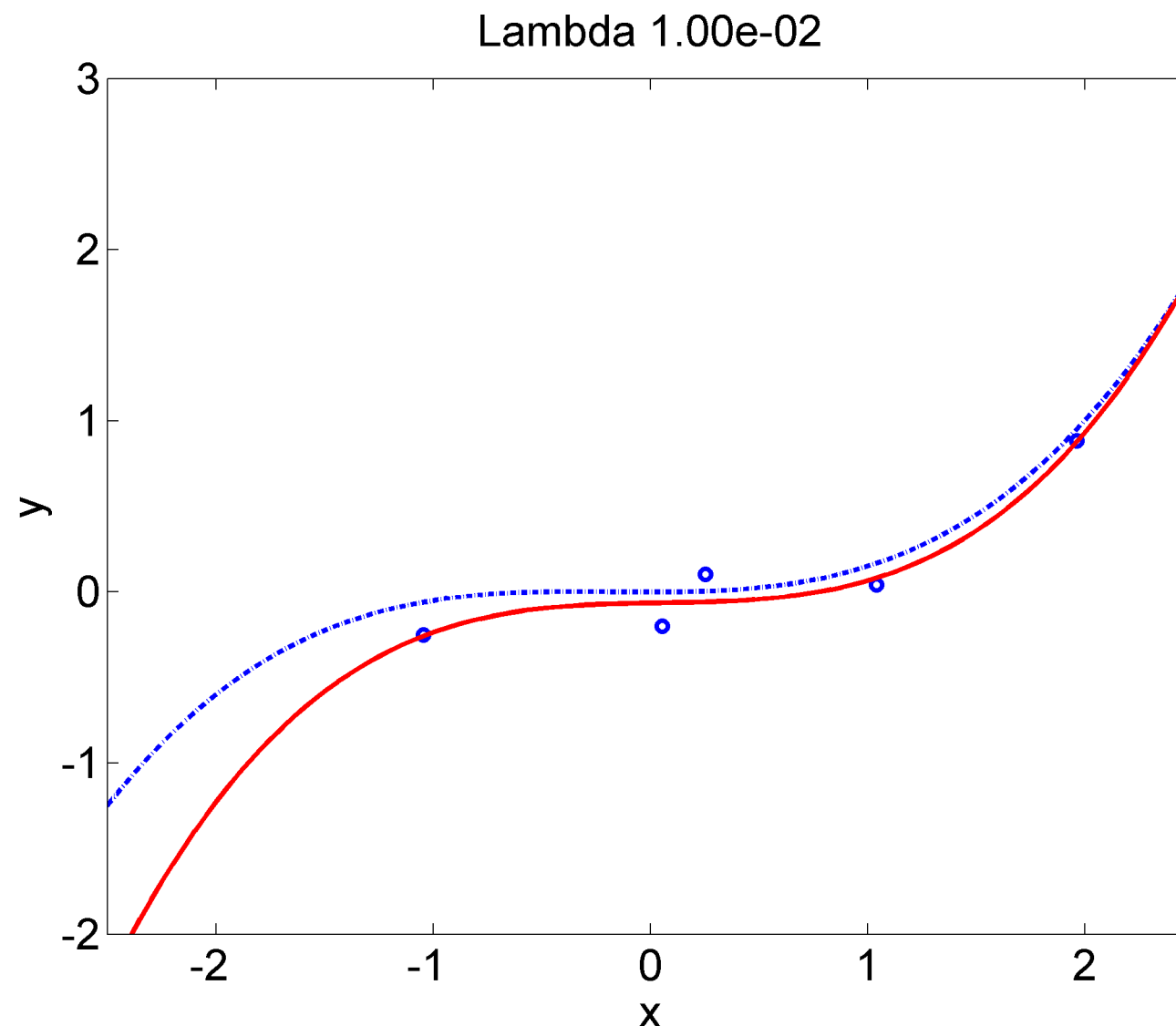
$$\boldsymbol{\mu}_N = \Sigma_N \Phi^T \mathbf{y}$$



What to do with the posterior?

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

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

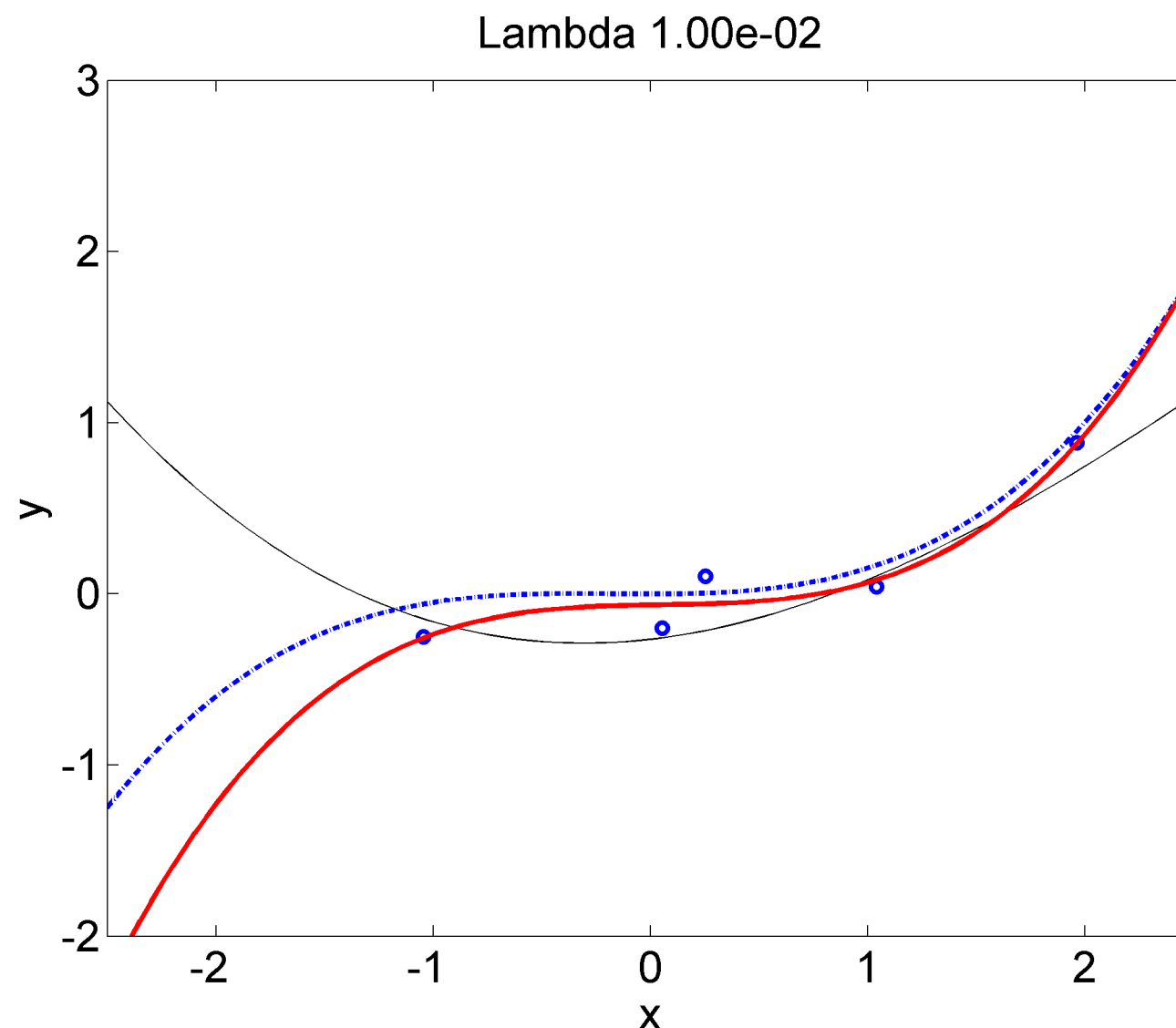




What to do with the posterior?

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

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

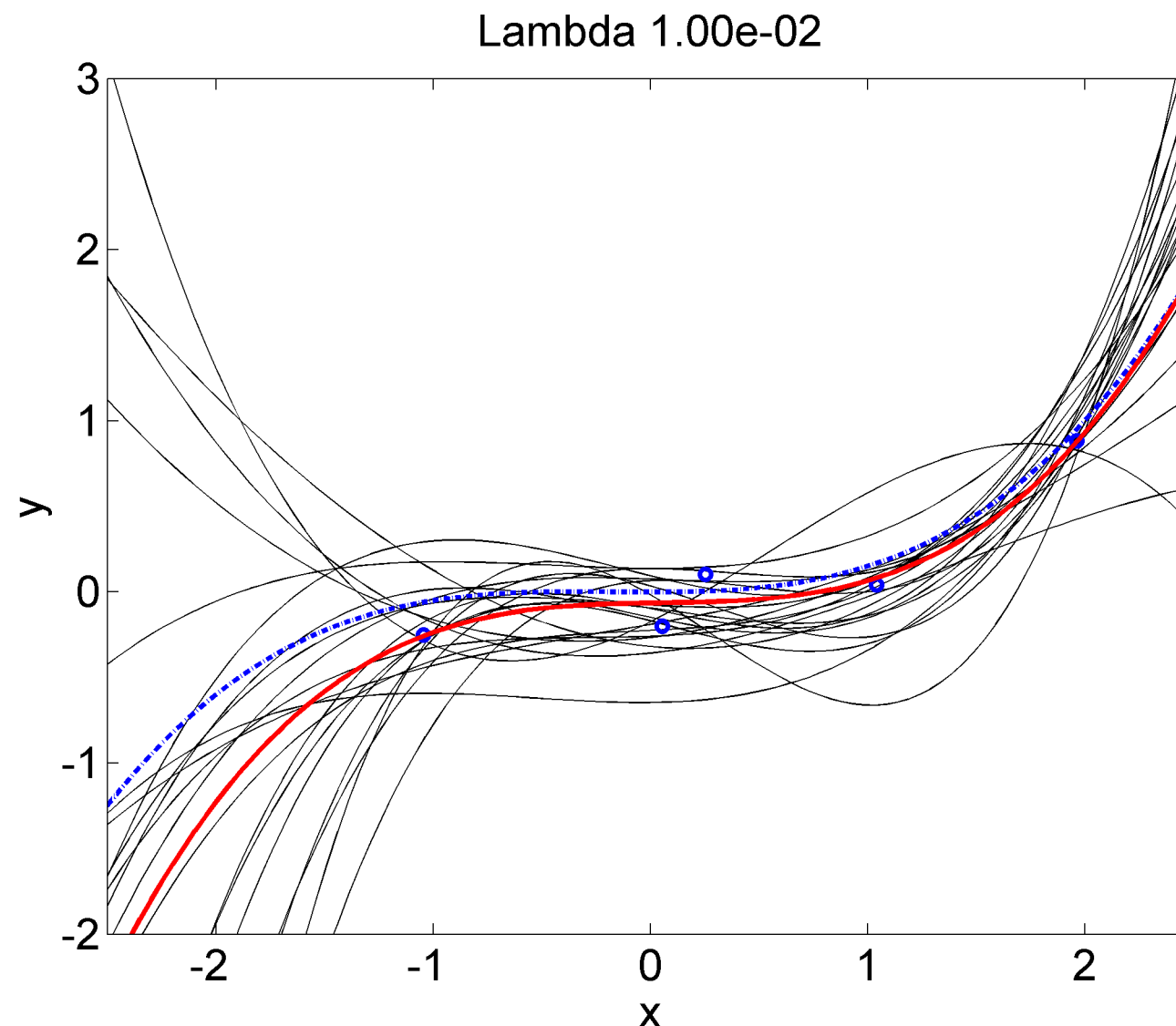




What to do with the posterior?

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 θ ?



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**

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

pred. function value test input training data likelihood parameter posterior

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

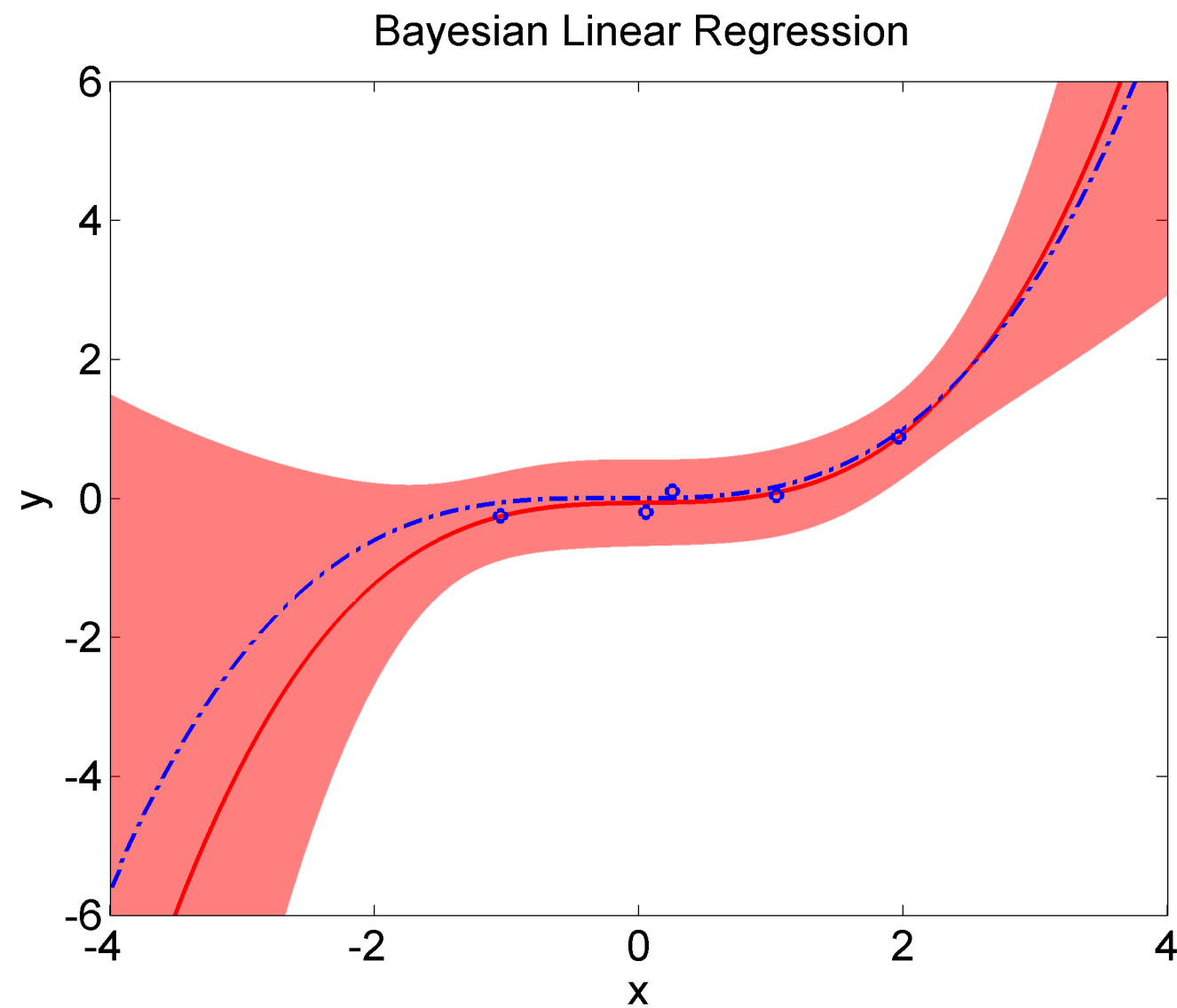
$$p(y_* | \mathbf{y}, \mathbf{X}, \mathbf{x}_*) = \mathcal{N}(\mathbf{y}_* | \mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

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

$$\sigma^2(\mathbf{x}_*) = \sigma^2 (1 + \boldsymbol{\phi}^T(\mathbf{x}_*) \boldsymbol{\Sigma}_N \boldsymbol{\phi}(\mathbf{x}_*))$$

State Dependent Variance!

Integrating out the parameters



Variance depends on the information in the data!

Quick Summary



- Models that are linear in the parameters: $y = \phi(x)^T \theta$
- **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