

Foundations of Machine Learning  
African Master's of Machine Intelligence



**AIMS** | African Institute for  
Mathematical Sciences  
RWANDA


**Imperial College  
London**

# Linear Regression

**Marc Deisenroth**

Quantum Leap Africa  
African Institute for Mathematical  
Sciences, Rwanda

Department of Computing  
Imperial College London

 @mpd37  
mdeisenroth@aimsammi.org

October 11, 2018

# Reference

## Mathematics for Machine Learning:

<https://mml-book.com>

### Chapter 9

# Overview

## Problem Setting

### Parameter Estimation

- Maximum Likelihood

- Maximum A Posteriori Estimation

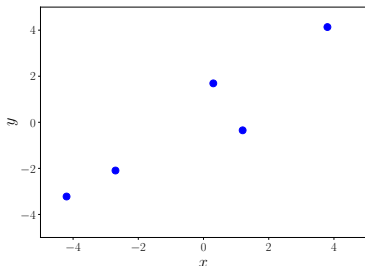
### Gaussian Identities

### Bayesian Linear Regression

# Regression Problems

## Regression (curve fitting)

Given inputs  $x$  and corresponding observations  $y$  find a function  $f$  that models the relationship between  $x$  and  $y$ .

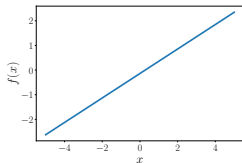


- ▶ Typically parametrize the function  $f$  with parameters  $\theta$
- ▶ Linear regression: Consider functions  $f$  that are **linear in the parameters**

# Linear Regression Functions

- ▶ Straight lines

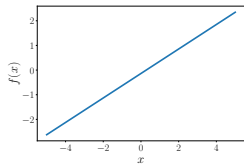
$$y = f(x, \boldsymbol{\theta}) = \theta_0 + \theta_1 x = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$



# Linear Regression Functions

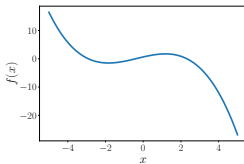
- ▶ Straight lines

$$y = f(x, \boldsymbol{\theta}) = \theta_0 + \theta_1 x = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$



- ▶ Polynomials

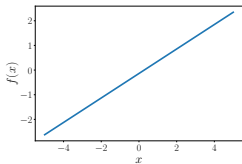
$$y = f(x, \boldsymbol{\theta}) = \sum_{m=0}^M \theta_m x^m = \begin{bmatrix} \theta_0 & \cdots & \theta_M \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ x^M \end{bmatrix}$$



# Linear Regression Functions

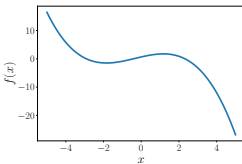
- ▶ Straight lines

$$y = f(x, \theta) = \theta_0 + \theta_1 x = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$



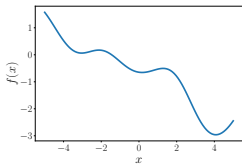
- ▶ Polynomials

$$y = f(x, \theta) = \sum_{m=0}^M \theta_m x^m = \begin{bmatrix} \theta_0 & \cdots & \theta_M \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ x^M \end{bmatrix}$$



- ▶ Radial basis function networks

$$y = f(x, \theta) = \sum_{m=1}^M \theta_m \exp\left(-\frac{1}{2}(x - \mu_m)^2\right)$$



# Linear Regression Model and Setting

$$y = \mathbf{x}^\top \boldsymbol{\theta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- ▶ Given a training set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  we seek optimal parameters  $\boldsymbol{\theta}^*$



# Linear Regression Model and Setting

$$y = \mathbf{x}^\top \boldsymbol{\theta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- ▶ Given a training set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  we seek optimal parameters  $\boldsymbol{\theta}^*$ 
  - ▶▶ **Maximum Likelihood Estimation**
  - ▶▶ **Maximum a Posteriori Estimation**

# Overview

Problem Setting

Parameter Estimation

- Maximum Likelihood

- Maximum A Posteriori Estimation

Gaussian Identities

Bayesian Linear Regression

# Maximum Likelihood

- ▶ Define  $\mathbf{X} = [x_1, \dots, x_N]^\top \in \mathbb{R}^{N \times D}$  and  $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$
- ▶ Find parameters  $\boldsymbol{\theta}^*$  that maximize the **likelihood**

# Maximum Likelihood

- ▶ Define  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$  and  $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$
- ▶ Find parameters  $\boldsymbol{\theta}^*$  that maximize the **likelihood**

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2)$$

# Maximum Likelihood

- ▶ Define  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$  and  $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$
- ▶ Find parameters  $\boldsymbol{\theta}^*$  that maximize the **likelihood**

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2)$$

- ▶ Log-transformation ► **Maximize the log likelihood**

# Maximum Likelihood

- ▶ Define  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$  and  $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$
- ▶ Find parameters  $\boldsymbol{\theta}^*$  that maximize the **likelihood**

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2)$$

- ▶ Log-transformation **▶ Maximize the log likelihood**

$$\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2),$$
$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

## Maximum Likelihood (2)

With

$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

we get

$$\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

## Maximum Likelihood (2)

With

$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

we get

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) &= \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const} \\ &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \text{const} \end{aligned}$$



## Maximum Likelihood (2)

With

$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

we get

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) &= \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const} \\ &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \text{const} \\ &= -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 + \text{const} \end{aligned}$$

## Maximum Likelihood (2)

With

$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

we get

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) &= \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const} \\ &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \text{const} \\ &= -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 + \text{const} \end{aligned}$$

- ▶ Computing the gradient with respect to  $\boldsymbol{\theta}$  and setting it to  $\mathbf{0}$  gives the **maximum likelihood estimator** (least-squares estimator)

$$\boldsymbol{\theta}^{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

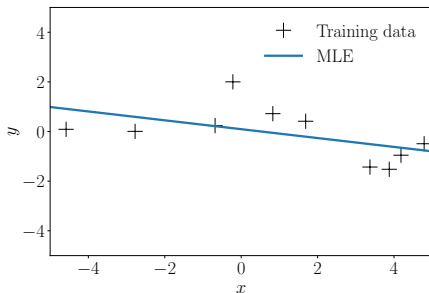
# Making Predictions

$$y = \mathbf{x}^\top \boldsymbol{\theta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Given an arbitrary input  $\mathbf{x}_*$ , we can predict the corresponding observation  $y_*$  using the maximum likelihood parameter:

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

## Example 1: Linear Functions



$$y = \theta_0 + \theta_1 x + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- At any query point  $x_*$  we obtain the mean prediction as

$$\mathbb{E}[y_* | \theta^{\text{ML}}, x_*] = \theta_0^{\text{ML}} + \theta_1^{\text{ML}} x_*$$

# Nonlinear Functions

$$y = \boldsymbol{\phi}(x)^\top \boldsymbol{\theta} + \epsilon = \sum_{m=0}^M \theta_m x^m + \epsilon$$

- ▶ Polynomial regression with features

$$\boldsymbol{\phi}(x) = [1, x, x^2, \dots, x^M]^\top$$

- ▶ Maximum likelihood estimator:

# Nonlinear Functions

$$y = \boldsymbol{\phi}(x)^\top \boldsymbol{\theta} + \epsilon = \sum_{m=0}^M \theta_m x^m + \epsilon$$

- ▶ Polynomial regression with features

$$\boldsymbol{\phi}(x) = [1, x, x^2, \dots, x^M]^\top$$

- ▶ Maximum likelihood estimator:

$$\boldsymbol{\theta}^{\text{ML}} = (\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{y}$$

## Example 2: Polynomial Regression

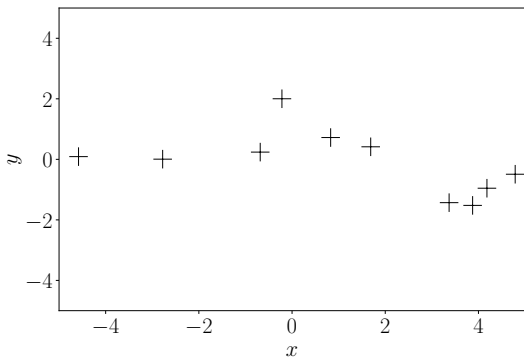


Figure: Training data

## Example 2: Polynomial Regression

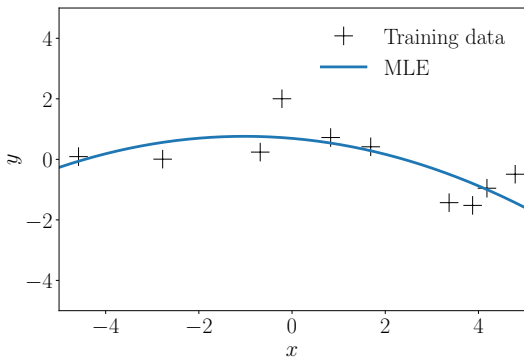


Figure: 2nd-order polynomial



## Example 2: Polynomial Regression

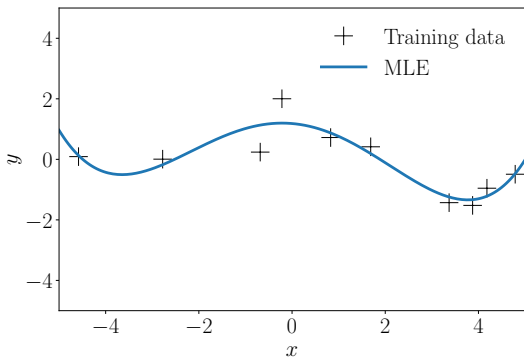


Figure: 4th-order polynomial

## Example 2: Polynomial Regression

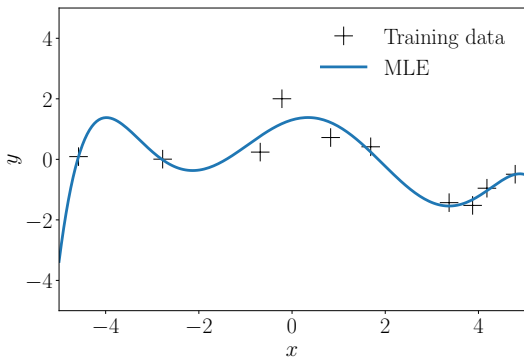


Figure: 6th-order polynomial

## Example 2: Polynomial Regression

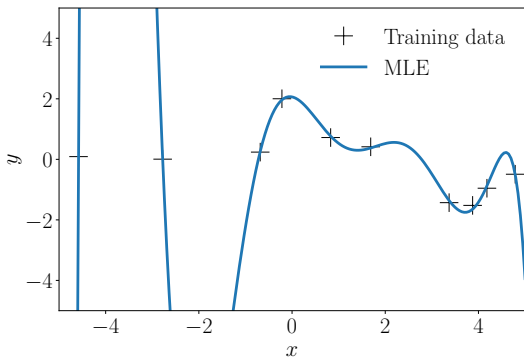


Figure: 8th-order polynomial

## Example 2: Polynomial Regression

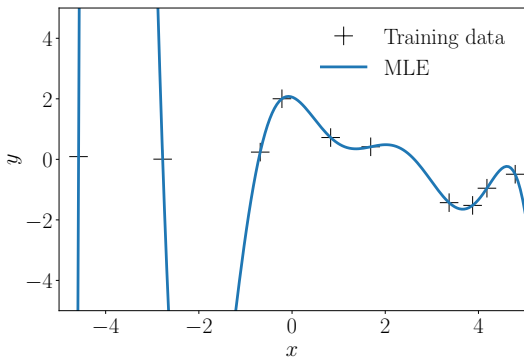
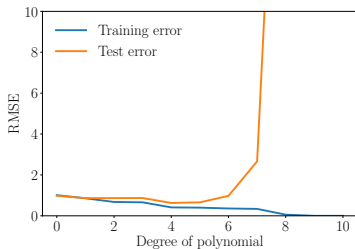


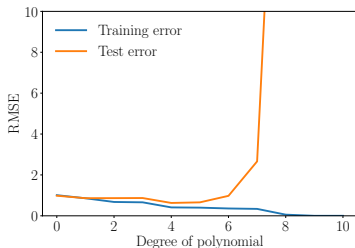
Figure: 10th-order polynomial

# Overfitting



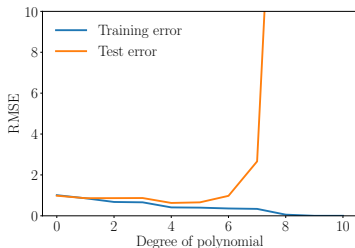
- ▶ Training error decreases with higher flexibility of the model

# Overfitting



- ▶ Training error decreases with higher flexibility of the model
- ▶ We are not so much interested in the training error, but in the **generalization error**: How well does the model perform when we predict at previously unseen input locations?

# Overfitting



- ▶ Training error decreases with higher flexibility of the model
- ▶ We are not so much interested in the training error, but in the **generalization error**: How well does the model perform when we predict at previously unseen input locations?
- ▶ Maximum likelihood often runs into **overfitting** problems, i.e., we exploit the flexibility of the model to fit to the noise in the data

# MAP Estimation

- ▶ **Observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values



# MAP Estimation

- ▶ **Observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values
- ▶ Mitigate the effect of overfitting by **placing a prior distribution  $p(\theta)$  on the parameters**
  - ▶▶ Penalize extreme values that are implausible under that prior

# MAP Estimation

- ▶ **Observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values
- ▶ Mitigate the effect of overfitting by placing a prior distribution  $p(\theta)$  on the parameters
  - ▶▶ Penalize extreme values that are implausible under that prior
- ▶ Choose  $\theta^*$  as the parameter that maximizes the (log) parameter posterior

$$\log p(\theta|X, \mathbf{y}) = \underbrace{\log p(\mathbf{y}|X, \theta)}_{\text{log-likelihood}} + \underbrace{\log p(\theta)}_{\text{log-prior}} + \text{const}$$

# MAP Estimation

- ▶ **Observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values
- ▶ Mitigate the effect of overfitting by **placing a prior distribution  $p(\theta)$  on the parameters**
  - ▶▶ Penalize extreme values that are implausible under that prior
- ▶ Choose  $\theta^*$  as the parameter that **maximizes the (log) parameter posterior**

$$\log p(\theta|X, \mathbf{y}) = \underbrace{\log p(\mathbf{y}|X, \theta)}_{\text{log-likelihood}} + \underbrace{\log p(\theta)}_{\text{log-prior}} + \text{const}$$

- ▶ Log-prior induces a direct penalty on the parameters

# MAP Estimation

- ▶ **Observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values
- ▶ Mitigate the effect of overfitting by **placing a prior distribution  $p(\theta)$  on the parameters**
  - ▶▶ Penalize extreme values that are implausible under that prior
- ▶ Choose  $\theta^*$  as the parameter that **maximizes the (log) parameter posterior**

$$\log p(\theta|X, \mathbf{y}) = \underbrace{\log p(\mathbf{y}|X, \theta)}_{\text{log-likelihood}} + \underbrace{\log p(\theta)}_{\text{log-prior}} + \text{const}$$

- ▶ Log-prior induces a direct penalty on the parameters
- ▶ **Maximum a posteriori estimate** (regularized least squares)

## MAP Estimation (2)

- ▶ Gaussian parameter prior  $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \alpha^2 \mathbf{I})$
- ▶ Log-posterior distribution:

$$\begin{aligned}\log p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) - \frac{1}{2\alpha^2} \boldsymbol{\theta}^\top \boldsymbol{\theta} + \text{const} \\ &= -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 - \frac{1}{2\alpha^2} \|\boldsymbol{\theta}\|^2 + \text{const}\end{aligned}$$

## MAP Estimation (2)

- ▶ Gaussian parameter prior  $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \alpha^2 \mathbf{I})$
- ▶ Log-posterior distribution:

$$\begin{aligned}\log p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) - \frac{1}{2\alpha^2} \boldsymbol{\theta}^\top \boldsymbol{\theta} + \text{const} \\ &= -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 - \frac{1}{2\alpha^2} \|\boldsymbol{\theta}\|^2 + \text{const}\end{aligned}$$

- ▶ Compute gradient with respect to  $\boldsymbol{\theta}$ , set it to  $\mathbf{0}$ 
  - ▶▶ **Maximum a posteriori estimate:**

$$\boldsymbol{\theta}^{\text{MAP}} = (\mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{\alpha^2} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

## Example: Polynomial Regression

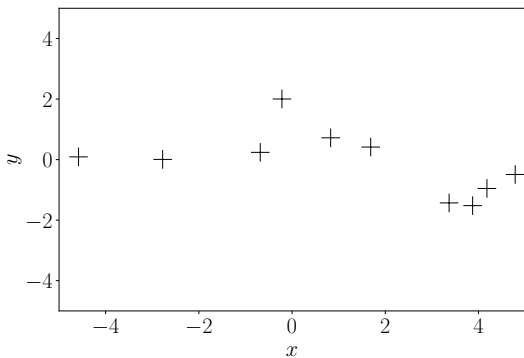


Figure: Training data

Mean prediction:

$$\mathbb{E}[y_* | \mathbf{x}_*, \boldsymbol{\theta}_{\text{MAP}}^*] = \boldsymbol{\phi}(\mathbf{x}_*)^\top \boldsymbol{\theta}_{\text{MAP}}^*$$

# Example: Polynomial Regression

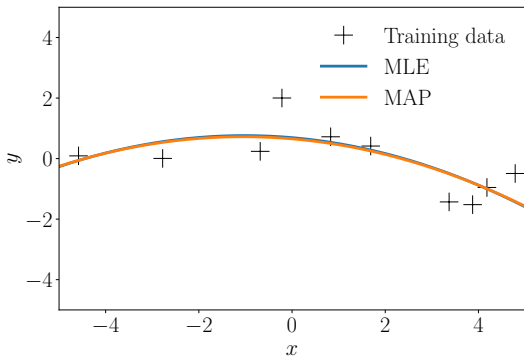


Figure: 2nd-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | x_*, \theta_{\text{MAP}}^*] = \boldsymbol{\phi}(x_*)^\top \boldsymbol{\theta}_{\text{MAP}}^*$$



# Example: Polynomial Regression

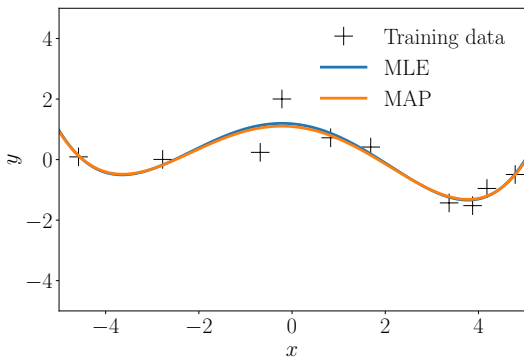


Figure: 4th-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | x_*, \theta_{\text{MAP}}^*] = \phi(x_*)^\top \theta_{\text{MAP}}^*$$

# Example: Polynomial Regression

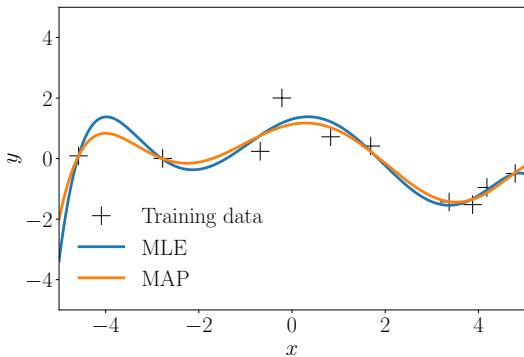


Figure: 6th-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | \mathbf{x}_*, \boldsymbol{\theta}_{\text{MAP}}^*] = \boldsymbol{\phi}(\mathbf{x}_*)^\top \boldsymbol{\theta}_{\text{MAP}}^*$$

# Example: Polynomial Regression

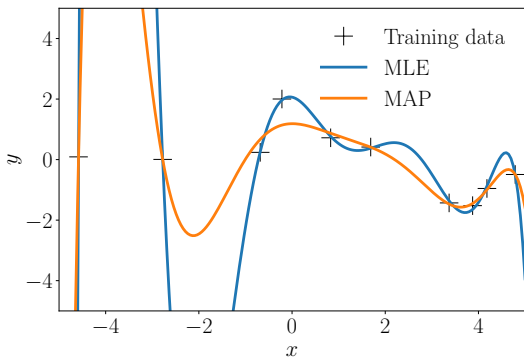


Figure: 8th-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | \mathbf{x}_*, \boldsymbol{\theta}_{\text{MAP}}^*] = \boldsymbol{\phi}(\mathbf{x}_*)^\top \boldsymbol{\theta}_{\text{MAP}}^*$$

# Example: Polynomial Regression

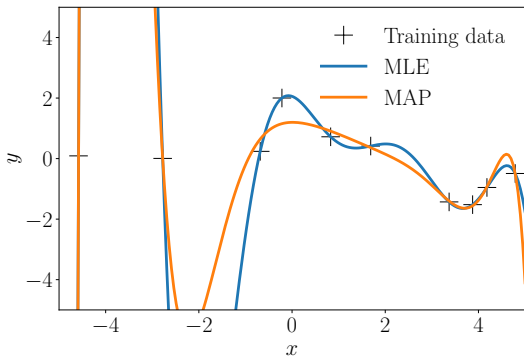
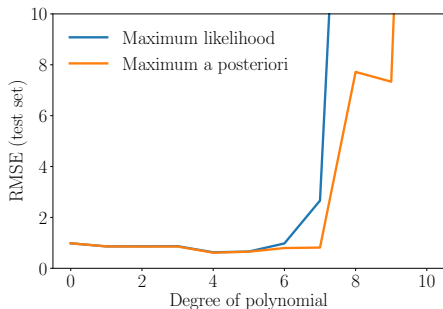


Figure: 10th-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | \mathbf{x}_*, \boldsymbol{\theta}_{\text{MAP}}^*] = \boldsymbol{\phi}(\mathbf{x}_*)^\top \boldsymbol{\theta}_{\text{MAP}}^*$$

# Generalization Error



- ▶ Maximum likelihood estimation “delays” the problem of overfitting
- ▶ It does not provide a general solution
- ▶▶ Need a more principled solution

# Overview

Problem Setting

Parameter Estimation

- Maximum Likelihood

- Maximum A Posteriori Estimation

Gaussian Identities

Bayesian Linear Regression

► Joint Gaussian distribution

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix} \right)$$

- ▶ Joint Gaussian distribution

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix} \right)$$

- ▶ Marginal:

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\ &= \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx}) \end{aligned}$$



- ▶ Joint Gaussian distribution

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix} \right)$$

- ▶ Marginal:

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\ &= \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx}) \end{aligned}$$

- ▶ Conditional:

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y}) \\ \boldsymbol{\mu}_{x|y} &= \boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} (\mathbf{y} - \boldsymbol{\mu}_y) \\ \boldsymbol{\Sigma}_{x|y} &= \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} \end{aligned}$$

# Linear Transformation of Gaussian Random Variables

If  $\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $\mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{b}$  then

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$$

## Product of Two Gaussians

$\mathbf{x} \in \mathbb{R}^D$ . Then:

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

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})$$

$$Z = (2\pi)^{-\frac{D}{2}} |\mathbf{A} + \mathbf{B}| \exp\left(-\frac{1}{2}(\mathbf{a} - \mathbf{b})^\top (\mathbf{A} + \mathbf{B})^{-1}(\mathbf{a} - \mathbf{b})\right)$$

## Product of Two Gaussians

$\mathbf{x} \in \mathbb{R}^D$ . Then:

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

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})$$

$$Z = (2\pi)^{-\frac{D}{2}} |\mathbf{A} + \mathbf{B}| \exp\left(-\frac{1}{2}(\mathbf{a} - \mathbf{b})^\top (\mathbf{A} + \mathbf{B})^{-1} (\mathbf{a} - \mathbf{b})\right)$$

- ▶ Product of two Gaussians is an unnormalized Gaussian

## Product of Two Gaussians

$\mathbf{x} \in \mathbb{R}^D$ . Then:

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

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})$$

$$Z = (2\pi)^{-\frac{D}{2}} |\mathbf{A} + \mathbf{B}| \exp\left(-\frac{1}{2}(\mathbf{a} - \mathbf{b})^\top (\mathbf{A} + \mathbf{B})^{-1} (\mathbf{a} - \mathbf{b})\right)$$

- ▶ Product of two Gaussians is an unnormalized Gaussian
- ▶ The “un-normalizer”  $Z$  has a Gaussian functional form:

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

Note: This is not a distribution (no random variables)

## Example: Marginalization of a Product

$$p_1(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mathbf{a}, \mathbf{A})$$

$$p_2(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mathbf{b}, \mathbf{B})$$

Then

$$\int p_1(\mathbf{x})p_2(\mathbf{x})d\mathbf{x} =$$

## Example: Marginalization of a Product

$$p_1(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mathbf{a}, \mathbf{A})$$

$$p_2(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mathbf{b}, \mathbf{B})$$

Then

$$\int p_1(\mathbf{x})p_2(\mathbf{x})d\mathbf{x} = Z = \mathcal{N}(\mathbf{a} \mid \mathbf{b}, \mathbf{A} + \mathbf{B})$$

## Example: Marginalization of a Product

$$p_1(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{a}, \mathbf{A})$$

$$p_2(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{b}, \mathbf{B})$$

Then

$$\int p_1(\mathbf{x})p_2(\mathbf{x})d\mathbf{x} = Z = \mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B})$$

$\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{Q})$ . Then

$$p(\mathbf{z}) = \int p(\mathbf{z}|\mathbf{x})p(\mathbf{x})d\mathbf{x}$$



## Example: Marginalization of a Product

$$p_1(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{a}, \mathbf{A})$$

$$p_2(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{b}, \mathbf{B})$$

Then

$$\int p_1(\mathbf{x})p_2(\mathbf{x})d\mathbf{x} = Z = \mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B})$$

$\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{Q})$ . Then

$$\begin{aligned} p(\mathbf{z}) &= \int p(\mathbf{z}|\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= \int \mathcal{N}(\mathbf{z} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{Q})\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})d\mathbf{x} \\ &= ?? \quad \blacktriangleright \text{later} \end{aligned}$$

# Overview

Problem Setting

Parameter Estimation

- Maximum Likelihood

- Maximum A Posteriori Estimation

Gaussian Identities

Bayesian Linear Regression

# Bayesian Linear Regression

$$y = \boldsymbol{\phi}(\mathbf{x})^\top \boldsymbol{\theta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- ▶ Avoid overfitting by not fitting any parameters:
  - ▶▶ Integrate parameters out instead of optimizing them

# Bayesian Linear Regression

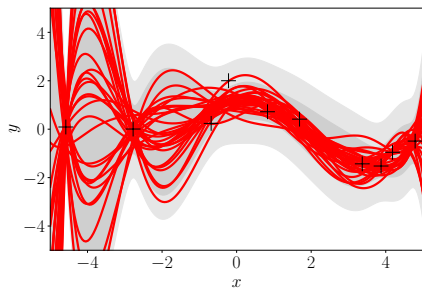
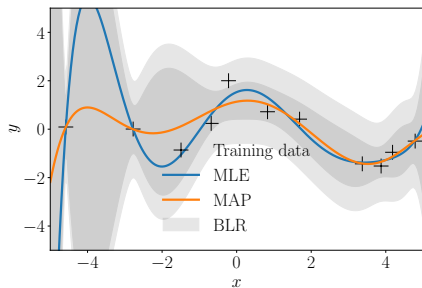
$$y = \boldsymbol{\phi}(\mathbf{x})^\top \boldsymbol{\theta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- ▶ Avoid overfitting by not fitting any parameters:
  - ▶ Integrate parameters out instead of optimizing them
- ▶ Use a full parameter distribution  $p(\boldsymbol{\theta})$  (and not a single point estimate  $\boldsymbol{\theta}^*$ ) when making predictions:

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

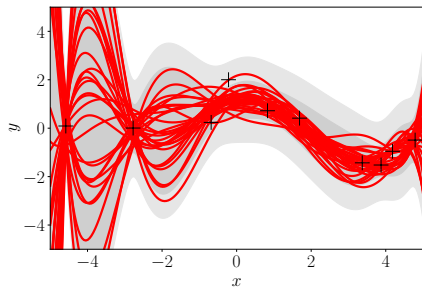
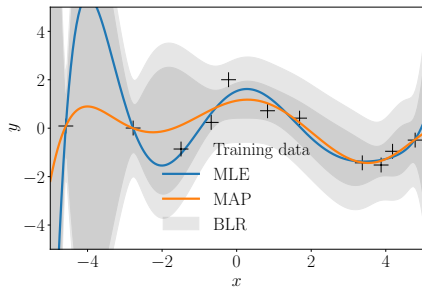
- ▶ Prediction no longer depends on  $\boldsymbol{\theta}$
- ▶ Predictive distribution reflects the uncertainty about the “correct” parameter setting

# Example



- ▶ Light-gray: uncertainty due to noise (same as in MLE/MAP)
- ▶ Dark-gray: uncertainty due to parameter uncertainty

# Example



- ▶ Light-gray: uncertainty due to noise (same as in MLE/MAP)
- ▶ Dark-gray: uncertainty due to parameter uncertainty
- ▶ Right: Plausible functions under the parameter distribution (every single parameter setting describes one function)

# Model for Bayesian Linear Regression

Prior  $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0),$

Likelihood  $p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y | \boldsymbol{\phi}^\top(\mathbf{x})\boldsymbol{\theta}, \sigma^2)$

- ▶ Parameter  $\boldsymbol{\theta}$  becomes a latent (random) variable
- ▶ Prior distribution induces a **distribution over plausible functions**
- ▶ Choose a conjugate Gaussian prior
  - ▶ Closed-form computations
  - ▶ Gaussian posterior

# Parameter Posterior and Predictions

- ▶ Prior  $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$  is Gaussian ► posterior is Gaussian:  
►► **Derive this**

$$\begin{aligned}p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}) &= \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N) \\ \mathbf{S}_N &= (\mathbf{S}_0^{-1} + \sigma^{-2}\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \\ \mathbf{m}_N &= \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \sigma^{-2}\boldsymbol{\Phi}^\top \mathbf{y})\end{aligned}$$



## Parameter Posterior and Predictions

- ▶ Prior  $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$  is Gaussian  $\blacktriangleright$  posterior is Gaussian:

$$\begin{aligned}p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}) &= \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N) \\ \mathbf{S}_N &= (\mathbf{S}_0^{-1} + \sigma^{-2}\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \\ \mathbf{m}_N &= \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \sigma^{-2}\boldsymbol{\Phi}^\top \mathbf{y})\end{aligned}$$

- ▶ Mean  $\mathbf{m}_N$  identical to MAP estimate

## Parameter Posterior and Predictions

- ▶ Prior  $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$  is Gaussian  $\blacktriangleright$  posterior is Gaussian:

$$\begin{aligned}p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}) &= \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N) \\ \mathbf{S}_N &= (\mathbf{S}_0^{-1} + \sigma^{-2}\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \\ \mathbf{m}_N &= \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \sigma^{-2}\boldsymbol{\Phi}^\top \mathbf{y})\end{aligned}$$

- ▶ Mean  $\mathbf{m}_N$  identical to MAP estimate
- ▶ Assume a Gaussian distribution  $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N)$ . Then

$$p(y|\mathbf{x}) = \mathcal{N}(y | \boldsymbol{\Phi}^\top(\mathbf{x})\mathbf{m}_N, \boldsymbol{\Phi}^\top(\mathbf{x})\mathbf{S}_N\boldsymbol{\Phi}(\mathbf{x}) + \sigma^2)$$

# Parameter Posterior and Predictions

- ▶ Prior  $p(\theta) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$  is Gaussian  $\blacktriangleright$  posterior is Gaussian:

$$\begin{aligned}p(\theta|\mathbf{X}, \mathbf{y}) &= \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N) \\ \mathbf{S}_N &= (\mathbf{S}_0^{-1} + \sigma^{-2}\mathbf{\Phi}^\top\mathbf{\Phi})^{-1} \\ \mathbf{m}_N &= \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \sigma^{-2}\mathbf{\Phi}^\top\mathbf{y})\end{aligned}$$

- ▶ Mean  $\mathbf{m}_N$  identical to MAP estimate
- ▶ Assume a Gaussian distribution  $p(\theta) = \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N)$ . Then

$$p(y|\mathbf{x}) = \mathcal{N}(y | \mathbf{\Phi}^\top(\mathbf{x})\mathbf{m}_N, \mathbf{\Phi}^\top(\mathbf{x})\mathbf{S}_N\mathbf{\Phi}(\mathbf{x}) + \sigma^2)$$

- ▶  $\mathbf{\Phi}^\top(\mathbf{x})\mathbf{S}_N\mathbf{\Phi}(\mathbf{x})$ : Contribution to uncertainty due to parameter distribution

**More details**  $\blacktriangleright$  <https://mml-book.com>, **Chapter 9**

# Marginal Likelihood

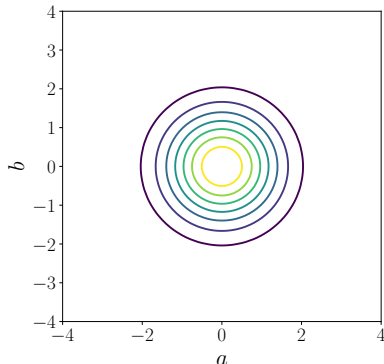
- ▶ Marginal likelihood can be computed analytically.
- ▶ With  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} = \mathcal{N}(\mathbf{y} | \boldsymbol{\Phi}\boldsymbol{\mu}, \boldsymbol{\Phi}\boldsymbol{\Sigma}\boldsymbol{\Phi}^\top + \sigma^2\mathbf{I})$$

# Distribution over Functions

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



# Sampling from the Prior over Functions

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$
$$f_i(x) = a_i + b_i x, \quad [a_i, b_i] \sim p(a, b)$$

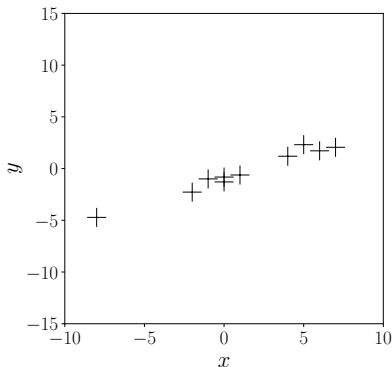
# Sampling from the Posterior over Functions

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$\mathbf{X} = [x_1, \dots, x_N]$ ,  $\mathbf{y} = [y_1, \dots, y_N]$  Training inputs/targets



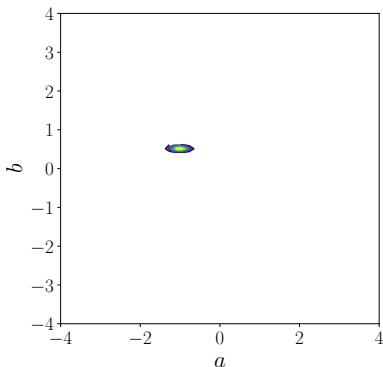
# Sampling from the Posterior over Functions

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$p(a, b | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N) \quad \text{Posterior}$$





# Sampling from the Posterior over Functions

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$[a_i, b_i] \sim p(a, b | \mathbf{X}, \mathbf{y})$$

$$f_i = a_i + b_i x$$

# Fitting Nonlinear Functions

- ▶ Fit nonlinear functions using (Bayesian) linear regression:  
Linear combination of nonlinear features

# Fitting Nonlinear Functions

- ▶ Fit nonlinear functions using (Bayesian) linear regression:  
Linear combination of nonlinear features
- ▶ Example: Radial-basis-function (RBF) network

$$f(\mathbf{x}) = \sum_{i=1}^n \theta_i \phi_i(\mathbf{x}), \quad \theta_i \sim \mathcal{N}(0, \sigma_p^2)$$

# Fitting Nonlinear Functions

- ▶ Fit nonlinear functions using (Bayesian) linear regression:  
Linear combination of nonlinear features
- ▶ Example: Radial-basis-function (RBF) network

$$f(\mathbf{x}) = \sum_{i=1}^n \theta_i \phi_i(\mathbf{x}), \quad \theta_i \sim \mathcal{N}(0, \sigma_p^2)$$

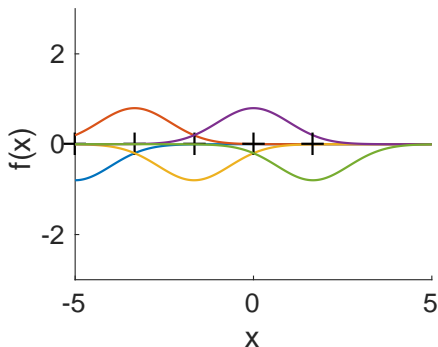
where

$$\phi_i(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^\top (\mathbf{x} - \boldsymbol{\mu}_i)\right)$$

for given “centers”  $\boldsymbol{\mu}_i$

# Illustration: Fitting a Radial Basis Function Network

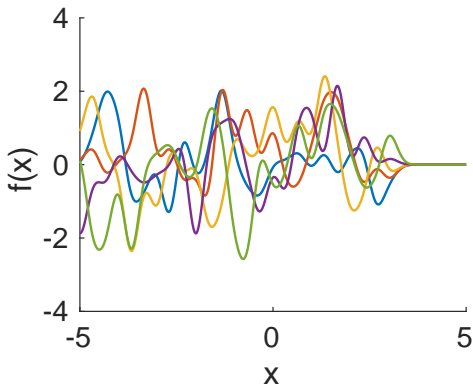
$$\phi_i(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^\top(\mathbf{x} - \boldsymbol{\mu}_i)\right)$$



- Place Gaussian-shaped basis functions  $\phi_i$  at 25 input locations  $\mu_i$ , linearly spaced in the interval  $[-5, 3]$

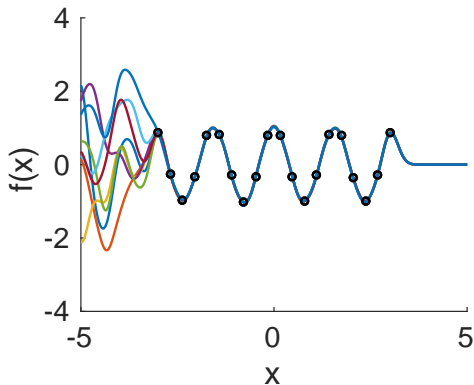
## Samples from the RBF Prior

$$f(\mathbf{x}) = \sum_{i=1}^n \theta_i \phi_i(\mathbf{x}), \quad p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

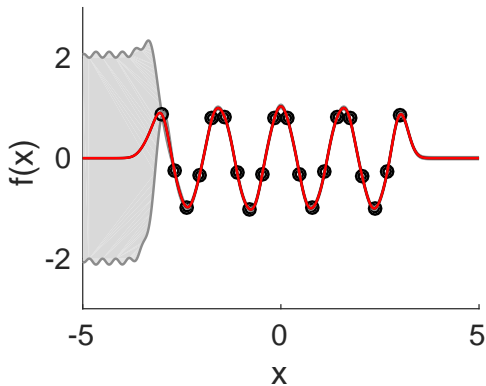


## Samples from the RBF Posterior

$$f(x) = \sum_{i=1}^n \theta_i \phi_i(x), \quad p(\theta | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N)$$

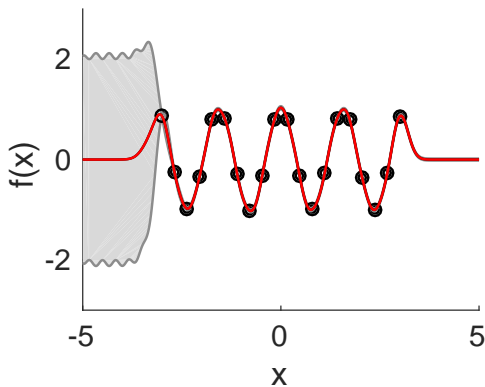


# RBF Posterior





## Limitations



- ▶ Feature engineering (what basis functions to use?)
- ▶ Finite number of features:
  - ▶ Above: Without basis functions on the right, we cannot express any variability of the function
  - ▶ Ideally: Add more (infinitely many) basis functions

# Approach

- ▶ Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
  - ▶▶ Place a prior on functions
  - ▶▶ Make assumptions on the distribution of functions

# Approach

- ▶ Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
  - ▶▶ Place a prior on functions
  - ▶▶ Make assumptions on the distribution of functions
- ▶ Intuition: function = infinitely long vector of function values
  - ▶▶ Make assumptions on the distribution of function values

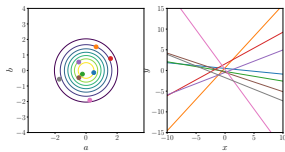
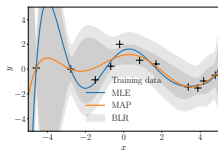
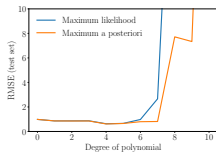
# Approach

- ▶ Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
  - ▶▶ Place a prior on functions
  - ▶▶ Make assumptions on the distribution of functions
- ▶ Intuition: function = infinitely long vector of function values
  - ▶▶ Make assumptions on the distribution of function values

# Approach

- ▶ Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
  - ▶▶ Place a prior on functions
  - ▶▶ Make assumptions on the distribution of functions
- ▶ Intuition: function = infinitely long vector of function values
  - ▶▶ Make assumptions on the distribution of function values
- ▶▶ **Gaussian process**

# Summary



- ▶ Regression = curve fitting
- ▶ Linear regression = linear in the parameters
- ▶ Parameter estimation via maximum likelihood and MAP estimation can lead to **overfitting**
- ▶ **Bayesian linear regression** addresses this issue, but may not be analytically tractable
- ▶ Predictive uncertainty in Bayesian linear regression explicitly depends on uncertainty of parameters
- ▶ Distribution over parameters induces distribution over functions