#### Universität Potsdam Institut für Informatik

Lehrstuhl Maschinelles Lernen



# **Bayesian Learning**

Tobias Scheffer, Niels Landwehr

#### **Remember: Normal Distribution**

- Distribution over  $x \in \mathbb{R}$ .
- Density function with parameters  $\mu \in \mathbb{R}$  (mean) and  $\sigma^2 \in \mathbb{R}$  (variance).



#### **Remember: Multivariate Normal Distribution**

- Distribution over vectors  $\mathbf{x} \in \mathbb{R}^{D}$ .
- Density function with parameters  $\mu \in \mathbb{R}^{D}, \Sigma \in \mathbb{R}^{D \times D}$ .

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

• Example *D*=2: density, sample from distribution



#### Overview

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

#### **Statistics & Machine Learning**

- Machine learning: tightly related to inductive statistics.
- Two areas in Statistics:
  - Descriptive Statistics: description and examination of the properties of data.

Mean values

Difference between Populations

 Inductive Statistics: What conclusions can be drawn from data about the underlying reality?

#### Model building

Variances

Explanations for observations

Relationships and patterns in the data

#### Frequentist vs. Bayesian Probabilities

- Frequentist probabilities
  - Describe the possibility of an occurrence of an intrinsically stochastic event (e.g., a coin toss).
  - Defined as limits of relative frequencies of possible outcomes in a repeatable experiment

*"If one throws a fair coin 1000 times, it will land on heads about 500 times"* 

"In 1 gram of Potassium-40, around 260,000 nuclei decay per second"

#### Frequentist vs. Bayesian Probabilities

- Bayesian "subjective" probabilities
  - Here, the reason for uncertainty is attributed to a lack of information.
  - How likely is it that suspect X killed the victim?
  - New Information (e.g., finger prints) can change these subjective probabilities.
- Bayesian view is more important in machine learning
- Frequentist and Bayesian perspectives are mathematically equivalent; in Bayesian statistics, probabilities are just used to model different things (lack of information).

#### **Bayesian Statistics**

- **1702-1761**
- "An essay towards solving a problem in the doctrine of chances", published in 1764.
- The work of Bayes laid the foundations for inductive Statistics.



 "Bayesian Probabilities" offer an important perspective on uncertainty and probability.

## **Bayesian Probability in Machine Learning**

- Model building: find an explanation for observations.
- What is the "most likely" model? Trade-off between
  - Prior knowledge (a priori distribution over models),
  - Evidence (data, observations).
- Bayesian Perspective:
  - Evidence (data) changes the "subjective" probability for models (explanation),
  - A posteriori model probability, MAP hypothesis.



Nature randomly determines  $\theta^*$  according to  $P(\theta)$ 

- "Nature" conducts random experiment, determines  $\theta^*$ .
- System with parameter  $\theta^*$  generates observations  $\mathbf{y} = f_{\theta^*}(\mathbf{X})$ .
- Bayesian inference inverts this process:
  - Given these assumptions about how y is generated,
  - Given the observed values of y for input matrix X,
  - What is the most likely true value of  $\theta$ ?
  - What is the most likely value  $y^*$  for a new input  $x^*$ ?



Maximum-likelihood (ML) model:

• 
$$\theta_{ML} = \arg \max_{\theta} P(\mathbf{y} | \mathbf{X}, \theta).$$

Maximum-a-positeriori (MAP) model:

• 
$$\theta_{MAP} = \arg \max_{\theta} P(\theta | \mathbf{y}, \mathbf{X})$$

A posteriori ("posterior") distribution

Likelihood

Maximum-likelihood (ML) model:

• 
$$\theta_{ML} = \arg \max_{\theta} P(\mathbf{y} | \mathbf{X}, \theta).$$

Maximum-a-positeriori (MAP) model:

• 
$$\theta_{MAP} = \arg \max_{\theta} P(\theta | \mathbf{y}, \mathbf{X}) = \arg \max_{\theta} \frac{P(\mathbf{y} | \mathbf{X}, \theta) P(\theta)}{P(\mathbf{y} | \mathbf{X})}$$
  
=  $\arg \max_{\theta} P(\mathbf{y} | \mathbf{X}, \theta) P(\theta)$   
Posterior  $\propto$  likelihood x prior

Maximum-likelihood (ML) model:

• 
$$\theta_{ML} = \arg \max_{\theta} P(\mathbf{y} | \mathbf{X}, \theta).$$

Maximum-a-positeriori (MAP) model:

• 
$$\theta_{MAP} = \arg \max_{\theta} P(\theta | \mathbf{y}, \mathbf{X}) = \arg \max_{\theta} \frac{P(\mathbf{y} | \mathbf{X}, \theta) P(\theta)}{P(\mathbf{y} | \mathbf{X})}$$
  
=  $\arg \max_{\theta} P(\mathbf{y} | \mathbf{X}, \theta) P(\theta)$ 

Most likely value y\* for new input x\* (Bayes-optimal decision):

• 
$$\mathbf{y}^* = \arg \max_{y} P(y | \mathbf{x}^*, \mathbf{y}, \mathbf{X})$$

Maximum-likelihood (ML) model:

• 
$$\theta_{ML} = \arg \max_{\theta} P(\mathbf{y} | \mathbf{X}, \theta).$$

Maximum-a-positeriori (MAP) model:

• 
$$\theta_{MAP} = \arg \max_{\theta} P(\theta | \mathbf{y}, \mathbf{X}) = \arg \max_{\theta} \frac{P(\mathbf{y} | \mathbf{X}, \theta) P(\theta)}{P(\mathbf{y} | \mathbf{X})}$$
  
=  $\arg \max_{\theta} P(\mathbf{y} | \mathbf{X}, \theta) P(\theta)$ 

Most likely value y\* for new input x\* (Bayes-optimal decision):

• 
$$\mathbf{y}^* = \arg \max_{y} P(y | \mathbf{x}^*, \mathbf{y}, \mathbf{X})$$
  
•  $P(y^* | \mathbf{x}^*, \mathbf{y}, \mathbf{X}) = \int P(y^*, \theta | \mathbf{x}^*, \mathbf{y}, \mathbf{X}) d\theta$   
Predictive  
distribution  
• Bayesian model averaging". Often computationally  
infeasible, but has a closed-form solution in some cases.

#### **Linear Regression Models**

• Training data:

• 
$$\mathbf{X} = \begin{pmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nm} \end{pmatrix}$$
  
•  $\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \stackrel{\frown \otimes}{[\odot]}$ 

Model

• 
$$f_{\theta} : X \to Y$$
  
 $\bigotimes \mapsto \bigotimes$   
•  $f_{\theta}(\mathbf{x}) = \mathbf{x}^{\mathrm{T}} \boldsymbol{\theta}$ 

#### **Overview**

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

#### **Probabilistic Linear Regression**

• Assumption 1: Nature generates parameter  $\theta^*$  of a linear function  $f_{\theta^*}(\mathbf{x}) = \mathbf{x}^T \theta^*$  according to  $p(\theta)$ .



#### **Probabilistic Linear Regression**

- Assumption 1: Nature generates parameter  $\theta^*$  of a linear function  $f_{\theta^*}(\mathbf{x}) = \mathbf{x}^T \theta^*$  according to  $p(\theta)$ .
- Assumption 2: Given inputs **X**, nature generates outputs **y**:

• 
$$y_i = f_{\theta^*}(\mathbf{x}_i) + \epsilon_i$$
 with  $\epsilon_i \sim N(\epsilon | 0, \sigma^2)$ .

• 
$$p(y_i | \mathbf{x}_i, \mathbf{\theta}^*) = N(y_i | \mathbf{x}_i^{\mathrm{T}} \mathbf{\theta}^*, \sigma^2)$$



In reality, we have y, X and want to make inferences about  $\theta$ .

Maximum-likelihood (ML) model:

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^{n} N(y_i | \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta}, \sigma^2)$$
Training instances are independent

Maximum-likelihood (ML) model:

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^{n} N(y_i | \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta}, \sigma^2)$$
$$= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} \left(y_i - \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta}\right)^2\right\}$$

Maximum-likelihood (ML) model:

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^{n} N(y_i | \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta}, \sigma^2)$$
$$= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} \left(y_i - \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta}\right)^2\right\}$$

n

- Log is a monononic transformation:
  - $\arg \max_{\boldsymbol{\theta}} P(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \log P(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$
- Also constant terms (constant in  $\theta$ ) can be dropped.

Maximum-likelihood (ML) model:

n

• Log is a monononic transformation:

•  $\arg \max_{\theta} P(\mathbf{y}|\mathbf{X}, \mathbf{\theta}) = \arg \max_{\theta} \log P(\mathbf{y}|\mathbf{X}, \mathbf{\theta})$ 

• Also constant terms (constant in  $\theta$ ) can be dropped

Maximum-likelihood (ML) model:

$$\boldsymbol{\theta}_{ML} = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} (y_i - \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta})^2$$

- Known as least-squares method in statistics.
- Setting the derivative to zero gives closed-form solution:  $\mathbf{\theta}_{ML} = (\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$
- Inversion of X<sup>T</sup>X is numerically unstable.

- The maximum-likelihood model is only based on the data, it is independent of any prior knowledge or domain assumptions.
- Calculating the maximum-likelihood model is numerically unstable.
- Regularized least squares works much better in practice.

#### **Overview**

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

Maximum-a-positeriori (MAP) model:

• 
$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} P(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) = \arg \max_{\boldsymbol{\theta}} \frac{P(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) P(\boldsymbol{\theta})}{P(\mathbf{y} | \mathbf{X})}$$

Maximum-a-positeriori (MAP) model:

• 
$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} P(\boldsymbol{\theta} | \boldsymbol{y}, \boldsymbol{X}) = \arg \max_{\boldsymbol{\theta}} \frac{P(\boldsymbol{y} | \boldsymbol{X}, \boldsymbol{\theta}) P(\boldsymbol{\theta})}{P(\boldsymbol{y} | \boldsymbol{X})}$$
  

$$= \arg \max_{\boldsymbol{\theta}} P(\boldsymbol{y} | \boldsymbol{X}, \boldsymbol{\theta}) P(\boldsymbol{\theta}) \qquad \text{Training instances} \\ \arg \max_{\boldsymbol{\theta}} \log (P(\boldsymbol{y} | \boldsymbol{X}, \boldsymbol{\theta}) P(\boldsymbol{\theta})) \qquad \text{Training instances} \\ = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log P(y_i | \boldsymbol{x}_i, \boldsymbol{\theta}) + \log P(\boldsymbol{\theta}) \\ = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log N(y_i | \boldsymbol{x}_i^T \boldsymbol{\theta}, \sigma^2) + \log P(\boldsymbol{\theta}) \\ = \cdots$$

#### MAP for Linear Regression: Prior

- Nature generates parameter  $\theta^*$  of linear model  $f_{\theta^*}(\mathbf{x}) = \mathbf{x}^T \theta^*$  according to  $p(\theta)$ .
- For convenience, assume  $p(\mathbf{\theta}) = N(\mathbf{\theta}|\mathbf{0}, \sigma_p^2 \mathbf{I})$ .

$$p(\mathbf{\theta}) = \mathcal{N}(\mathbf{\theta} | \mathbf{0}, \sigma_p^2 \mathbf{I})$$
$$= \frac{1}{2\pi^{m/2} \sigma_p^m} \exp\left(-\frac{1}{2\sigma_p^2} |\mathbf{\theta}|^2\right)$$

 $\sigma_p^2 \in \mathbb{R}$  controls strength of prior



Maximum-a-positeriori (MAP) model:

• 
$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log N(y_i | \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta}, \sigma^2) + \log P(\boldsymbol{\theta})$$
  

$$= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} (y_i - \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta})^2$$

$$+ \log \frac{1}{2\pi^{\frac{m}{2}} \sigma_p} - \frac{1}{2\sigma_p^2} \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\theta}$$
All terms that are constant in  $\boldsymbol{\theta}$  can be dropped.

Maximum-a-positeriori (MAP) model:

• 
$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log N(y_i | \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta}) + \log P(\boldsymbol{\theta})$$
  

$$= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} (y_i - \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta})^2$$

$$+ \log \frac{1}{2\pi^{\frac{m}{2}} \sigma_p} - \frac{1}{2\sigma_p^2} \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\theta}$$

$$= \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \frac{1}{2\sigma^2} (y_i - \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta})^2 + \frac{1}{2\sigma_p^2} \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\theta}$$

$$= \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{n} (y_i - \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta})^2 - \frac{\sigma^2}{\sigma_p^2} \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\theta}$$

$$\ell_2 \text{-regularized linear regression with}$$

squared loss (ridge regression).

Maximum-a-positeriori (MAP) model:

• 
$$\boldsymbol{\theta}_{MAP} = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} (y_i - \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta})^2 - \frac{\sigma^2}{\sigma_p^2} \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\theta}$$

- Same optimization criterion as ridge regression.
- Analytic solution (see lecture on ridge regression):

• 
$$\boldsymbol{\Theta}_{MAP} = \left(\mathbf{X}^{\mathrm{T}}\mathbf{X} + \frac{\sigma^2}{\sigma_p^2}\mathbf{I}\right)^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$

#### **Overview**

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

#### **Posterior for Linear Regression**

Posterior distribution of θ given y, X:

• 
$$P(\boldsymbol{\theta}|\mathbf{y},\mathbf{X}) = \frac{P(\mathbf{y}|\mathbf{X},\theta)P(\theta)}{P(\mathbf{y}|\mathbf{X})} = \frac{1}{Z}P(\mathbf{y}|\mathbf{X},\boldsymbol{\theta})P(\boldsymbol{\theta})$$

#### **Posterior for Linear Regression**

Posterior distribution of θ given y, X:



#### **Posterior for Linear Regression**

• Posterior distribution of  $\theta$  given y, X:

• 
$$P(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}) = \frac{P(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})P(\boldsymbol{\theta})}{P(\mathbf{y}|\mathbf{X})} = \frac{1}{Z}P(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})P(\boldsymbol{\theta})$$
  
$$= \frac{1}{Z}N(\mathbf{y}|\mathbf{X}^{\mathrm{T}}\boldsymbol{\theta}, \sigma^{2}\mathbf{I})N(\boldsymbol{\theta}|\mathbf{0}, \sigma_{p}^{2}\mathbf{I})$$
$$= N(\boldsymbol{\theta}|\overline{\boldsymbol{\theta}}, \mathbf{A}^{-1})$$

• With 
$$\overline{\mathbf{\Theta}} = \left(\mathbf{X}^{\mathrm{T}}\mathbf{X} + \frac{\sigma^{2}}{\sigma_{p}^{2}}\mathbf{I}\right)^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$
  
• And  $\mathbf{A}^{-1} = \sigma^{-2}\mathbf{X}^{\mathrm{T}}\mathbf{X} + \sigma_{p}^{-2}\mathbf{I}$ . Mean value

Mean value of posterior:  $\theta_{MAP}$
#### **Example MAP solution regression**

• Training data:

$$\mathbf{x}_1 = \begin{pmatrix} 2\\3\\0 \end{pmatrix}, \qquad \mathbf{x}_2 = \begin{pmatrix} 4\\3\\2 \end{pmatrix}, \qquad \mathbf{x}_3 = \begin{pmatrix} 0\\1\\2 \end{pmatrix},$$

 $y_1 = 2$   $y_2 = 3$   $y_3 = 4$ 

Matrix notation (adding constant attribute):

$$\mathbf{X} = \begin{pmatrix} 1 & 2 & 3 & 0 \\ 1 & 4 & 3 & 2 \\ 1 & 0 & 1 & 2 \end{pmatrix} \qquad \qquad \mathbf{y} = \begin{pmatrix} 2 \\ 3 \\ 4 \end{pmatrix}$$

#### **Example MAP solution regression**

- Choose
  - Variance of prior:  $\sigma_p = 1$
  - Noise parameter:  $\sigma = 0.5$

• Compute: 
$$\overline{\mathbf{\theta}} = (\mathbf{X}^{\mathrm{T}}\mathbf{X} + \frac{\sigma^{2}}{\sigma_{p}^{2}}\mathbf{I})^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$
  
 $\overline{\mathbf{\theta}} = \begin{pmatrix} 1 & 2 & 3 & 0 \\ 1 & 4 & 3 & 2 \\ 1 & 0 & 1 & 2 \end{pmatrix}^{\mathrm{T}} \begin{pmatrix} 1 & 2 & 3 & 0 \\ 1 & 4 & 3 & 2 \\ 1 & 0 & 1 & 2 \end{pmatrix}^{\mathrm{T}} + 0.25 \cdot \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}^{-1} \begin{pmatrix} 1 & 2 & 3 & 0 \\ 1 & 4 & 3 & 2 \\ 1 & 0 & 1 & 2 \end{pmatrix}^{\mathrm{T}} \begin{pmatrix} 2 \\ 3 \\ 4 \end{pmatrix}$   
 $\approx \begin{pmatrix} 0.7975 \\ -0.5598 \\ 0.7543 \\ 1.1217 \end{pmatrix}$ 

# **Example MAP solution regression**

• Predictions of model  $\overline{\theta}$  on the training data:

$$\hat{\mathbf{y}} = \mathbf{X}\overline{\mathbf{\theta}} = \begin{pmatrix} 1 & 2 & 3 & 0 \\ 1 & 4 & 3 & 2 \\ 1 & 0 & 1 & 2 \end{pmatrix} \begin{pmatrix} 0.7975 \\ -0.5598 \\ 0.7543 \\ 1.1217 \end{pmatrix} = \begin{pmatrix} 1.9408 \\ 3.0646 \\ 3.7952 \end{pmatrix}$$

# **Posterior and Regularized Loss Function**

MAP model:



# **Sequential Learning**

- Training examples arrive sequentially.
- Each training example  $(\mathbf{x}_i, y_i)$  changes prior  $p_{i-1}(\boldsymbol{\theta})$  into posterior  $p_{i-1}(\boldsymbol{\theta}|y_i, \mathbf{x}_i)$  which becomes the new prior  $p_i(\boldsymbol{\theta})$

• 
$$P(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}) = \frac{1}{Z} P_0(\boldsymbol{\theta}) P(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$$
  
 $= \frac{1}{Z} P_0(\boldsymbol{\theta}) \prod_{i=1}^n P(y_i|\mathbf{x}_i, \boldsymbol{\theta})$   
 $= \frac{1}{Z} \underbrace{\frac{P_0(\boldsymbol{\theta}) P(y_1|\mathbf{x}_1, \boldsymbol{\theta})}{P_1(\boldsymbol{\theta})} P(y_2|\mathbf{x}_2, \boldsymbol{\theta}) P(y_3|\mathbf{x}_3, \boldsymbol{\theta}) \dots P(y_n|\mathbf{x}_n, \boldsymbol{\theta})}_{P_3(\boldsymbol{\theta})}$ 

# **Example: Sequential Learning** [from Chris Bishop, Pattern Recognition and Machine Learning]

 $f_{\theta}(x) = \theta_0 + \theta_1 x$  (one-dimensional regression)

Sequential update:  $P_0(\boldsymbol{\theta})$ 



 $f_{\theta}(x) = \theta_0 + \theta_1 x$  (one-dimensional regression)

Sequential update:  $P_0(\mathbf{\theta})$ 



 $f_{\theta}(x) = \theta_0 + \theta_1 x$  (one-dimensional regression)

Sequential update:  $P_1(\mathbf{\theta}) \propto P_0(\mathbf{\theta}) P(y_1|x_1, \mathbf{\theta})$ 



 $f_{\theta}(x) = \theta_0 + \theta_1 x$  (one-dimensional regression)

Sequential update:  $P_2(\boldsymbol{\theta}) \propto P_1(\boldsymbol{\theta}) P(y_2|x_2, \boldsymbol{\theta})$ 



 $f_{\theta}(x) = \theta_0 + \theta_1 x$  (one-dimensional regression)

Sequential update:  $P_n(\mathbf{\theta}) \propto P_{n-1}(\mathbf{\theta})P(y_n|x_n, \mathbf{\theta})$ 



# **Learning and Prediction**

- So far, we have always separated *learning* from *prediction*.
- Learning:

• 
$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \widehat{R}(\mathbf{y}, \mathbf{X}, \boldsymbol{\theta}) + \Omega(\boldsymbol{\theta})$$

Prediction:

•  $y^* = f_{\mathbf{\theta}^*}(\mathbf{x}^*)$ 

For instance, in MAP linear regression, learning is

• 
$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} P(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}).$$

And prediction is

• 
$$y^* = \mathbf{\Theta}_{MAP}^T \mathbf{x}^*$$
.

# **Learning and Prediction**

- So far, we have always separated *learning* from *prediction*.
- And there are good reasons to do this:
  - Learning can require processing massive amounts of training data which can take a long time.
  - Predictions may have to be made in real time.
- However, sometimes, when relatively few data are available and an accurate prediction is worth waiting for, one can directly search for the best possible prediction:

• 
$$\mathbf{y}^* = \arg \max_{y} P(y | \mathbf{x}^*, \mathbf{y}, \mathbf{X})$$

• Most likely  $y^*$  for new input  $x^*$  given training data y, X.

# **Overview**

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

# **Bayes-Optimal Prediction**

Bayes-optimal decision: most likely value  $\mathbf{y}^*$  for new input  $\mathbf{x}^*$ 

• 
$$y^* = \arg \max_{y} P(y | \mathbf{x}^*, \mathbf{y}, \mathbf{X})$$

Predictive distribution for new input x\* given training data:



- Bayes-optimal decision is made by a weighted sum over all values of the model parameters.
- In general, there is no single model θ\* in the model space that always makes the Bayes-optimal decision.

# **Bayes-Optimal Prediction**

- Bayes-optimal decision is made by a weighted sum over all model parameters:
  - $P(y|\mathbf{x}^*, \mathbf{y}, \mathbf{X}) = \int P(y|\mathbf{x}^*, \mathbf{\theta}) P(\mathbf{\theta}|\mathbf{y}, \mathbf{X}) d\mathbf{\theta}$
- The prediction of the MAP model is only the prediction made by the single most likely model  $\theta_{MAP}$ .
  - Predictive distribution  $P(y|\mathbf{x}^*, \boldsymbol{\theta}_{MAP})$ .
  - Most likely prediction  $f_{\theta_{MAP}}(\mathbf{x}^*) = \arg \max_{y} P(y|\mathbf{x}^*, \boldsymbol{\theta}_{MAP}).$
- The MAP model  $\theta_{MAP}$  is an approximaton of this weighted sum by its element with the highest weight.

# **Bayes-Optimal Prediction**

- Bayes-optimal decision is made by a weighted sum over all model parameters:
  - $P(y|\mathbf{x}^*, \mathbf{y}, \mathbf{X}) = \int P(y|\mathbf{x}^*, \mathbf{\theta}) P(\mathbf{\theta}|\mathbf{y}, \mathbf{X}) d\mathbf{\theta}$
- Integration over the space of all model parameters is not generally possible.
- In some cases, there is a closed-form solution.
- In other cases, approximate numerical integration may be possible.

#### **Predictive Distribution for Linear Regression**

Predictive distribution for linear regression

• 
$$P(y|\mathbf{x}^*, \mathbf{y}, \mathbf{X}) = \int P(y|\mathbf{x}^*, \mathbf{\theta}) P(\mathbf{\theta}|\mathbf{y}, \mathbf{X}) d\mathbf{\theta}$$
  
=  $\int N(y|\mathbf{x}^*, \mathbf{\theta}) N(\mathbf{\theta}|\overline{\mathbf{\theta}}, \mathbf{A}^{-1}) d\mathbf{\theta}$   
=  $N(y|\overline{\mathbf{\theta}}^T \mathbf{x}^*, \sigma^2 + {\mathbf{x}^*}^T \mathbf{A}^{-1} \mathbf{x}^*)$ 

• With 
$$\overline{\mathbf{\Theta}} = \left(\mathbf{X}^{\mathrm{T}}\mathbf{X} + \frac{\sigma^2}{\sigma_p^2}\mathbf{I}\right)^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$

• And 
$$\mathbf{A}^{-1} = \sigma^{-2} \mathbf{X}^{\mathrm{T}} \mathbf{X} + \sigma_p^{-2} \mathbf{I}$$
.

Bayes-optimal prediction:

• 
$$y^* = \arg \max_{y} P(y | \mathbf{x}^*, \mathbf{y}, \mathbf{X}) = \overline{\mathbf{\theta}}^{\mathrm{T}} \mathbf{x}^*$$

#### **Predictive Distribution: Confidence Band**

Bayesian regression not only yields prediction  $y^* = \mathbf{x}^T \overline{\mathbf{\theta}}$ , but a distribution over y und therefore a confidence band.



# **Overview**

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

# Linear Classification

Training data:



- **Decision function:** 
  - $f_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{x}^{\mathrm{T}} \boldsymbol{\theta}$

• 
$$P(y|\mathbf{x}, \boldsymbol{\theta}) = \sigma(\mathbf{x}^{\mathrm{T}}\boldsymbol{\theta})$$

• 
$$y_{\theta}: \mathfrak{Q} \mapsto \mathfrak{B}$$

#### Linear Classification: Predictive Distribution

- For binary classification,  $y \in \{-1, +1\}$
- Predictive distribution given parameters θ of linear model:

• 
$$P(y = +1 | \mathbf{x}, \boldsymbol{\theta}) = \sigma(\mathbf{x}^{\mathrm{T}} \boldsymbol{\theta}) = \frac{1}{1 + \mathrm{e}^{-\mathrm{x}^{\mathrm{T}} \boldsymbol{\theta}}}$$

• 
$$P(y = -1 | \mathbf{x}, \mathbf{\theta}) = 1 - P(y = +1 | \mathbf{x}, \mathbf{\theta})$$

• Sigmoid function maps  $[-\infty, +\infty] \rightarrow [0,1]$ .



# **Logistic Regression**

- For binary classification,  $y \in \{-1, +1\}$
- Predictive distribution given parameters θ of linear model:

• 
$$P(y = +1 | \mathbf{x}, \mathbf{\theta}) = \sigma(\mathbf{x}^{\mathrm{T}} \mathbf{\theta}) = \frac{1}{1 + e^{-x^{\mathrm{T}} \mathbf{\theta}}}$$
  
•  $P(x = -1 | \mathbf{x}, \mathbf{\theta}) = 1$ 

• 
$$P(y = -1 | \mathbf{x}, \mathbf{\theta}) = 1 - \frac{1}{1 + e^{-\mathbf{x}^{T}\mathbf{\theta}}} = \frac{1}{1 + e^{\mathbf{x}^{T}\mathbf{\theta}}}$$

• Written jointly for both classes:

• 
$$P(y|\mathbf{x}, \boldsymbol{\theta}) = \sigma(y\mathbf{x}^{\mathrm{T}}\boldsymbol{\theta}) = \frac{1}{1 + \mathrm{e}^{-y\mathbf{x}^{\mathrm{T}}\boldsymbol{\theta}}}$$

Classification function:

• 
$$y_{\theta}(\mathbf{x}) = \arg \max_{y} P(y|\mathbf{x}, \theta)$$

 Called "logistic regression" even though it is a classification model.

# **Logistic Regression**

• Decision boundary:  $P(y = +1 | \mathbf{x}, \mathbf{\theta}) = P(y = -1 | \mathbf{x}, \mathbf{\theta}) = 0.5$ .

$$0.5 = \sigma(\mathbf{x}^{\mathrm{T}}\mathbf{\theta}) = \frac{1}{1 + \mathrm{e}^{-\mathbf{x}^{\mathrm{T}}\mathbf{\theta}}}$$
$$\Leftrightarrow 1 = \mathrm{e}^{-\mathbf{x}^{\mathrm{T}}\mathbf{\theta}}$$
$$\Leftrightarrow 0 = \mathbf{x}^{\mathrm{T}}\mathbf{\theta}$$

Decision boundary is a hyperplane in input space.



# **Logistic Regression**

- For multi-class classification:  $\boldsymbol{\Theta} = \begin{pmatrix} \boldsymbol{\Theta}_1 \\ \vdots \\ \boldsymbol{\Theta}_2 \end{pmatrix}$
- Generalize sigmoid function to softmax function:



• 
$$y_{\theta}(\mathbf{x}) = \arg \max_{y} P(y|\mathbf{x}, \theta)$$

 Called multi-class "logistic regression" even though it is a classification model.

# **Overview**

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

# Logistic Regression: ML Model

Maximum-likelihood model:

●

$$\begin{aligned} \mathbf{\theta}_{ML} &= \arg \max_{\mathbf{\theta}} P(\mathbf{y} | \mathbf{X}, \mathbf{\theta}) \\ &= \arg \max_{\mathbf{\theta}} \prod_{\substack{i=1 \\ n}}^{n} \frac{1}{1 + e^{-y_i \mathbf{x}_i^{\mathrm{T}} \mathbf{\theta}}} \\ &= \arg \min_{\mathbf{\theta}} \sum_{\substack{i=1 \\ n}}^{n} - \log \frac{1}{1 + e^{-y_i \mathbf{x}_i^{\mathrm{T}} \mathbf{\theta}}} \\ &= \arg \min_{\mathbf{\theta}} \sum_{\substack{i=1 \\ i=1}}^{n} \log \left(1 + e^{-y_i \mathbf{x}_i^{\mathrm{T}} \mathbf{\theta}}\right) \end{aligned}$$

 No analytic solution; numeric optimization, for instance, using (stochastic) gradient descent.

#### Logistic Regression: ML Model

Maximum-likelihood model:

• 
$$\boldsymbol{\theta}_{ML} = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log \left( 1 + e^{-y_i \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta}} \right)$$

• Gradient:

$$\begin{aligned} \bullet \quad & \frac{\partial}{\partial \theta} \sum_{i=1}^{n} \log \left( 1 + e^{-y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta} \right) \\ &= \sum_{i=1}^{n} \frac{\partial}{\partial \left( 1 + e^{-y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta} \right)} \log \left( 1 + e^{-y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta} \right) \frac{\partial}{\partial \left( -y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta \right)} \left( 1 + e^{-y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta} \right) \\ &= \sum_{i=1}^{n} \frac{1}{1 + e^{-y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta}} e^{-y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta} \left( -y_{i} \mathbf{x}_{i}^{\mathrm{T}} \right) = \sum_{i=1}^{n} -y_{i} \mathbf{x}_{i}^{\mathrm{T}} \frac{e^{-y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta}}{1 + e^{-y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta}} \\ &= \sum_{i=1}^{n} -y_{i} \mathbf{x}_{i}^{\mathrm{T}} \frac{1}{1 + e^{y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta}} \\ &= \sum_{i=1}^{n} y_{i} \mathbf{x}_{i} \left( 1 - \sigma (y_{i} \mathbf{x}_{i}^{\mathrm{T}} \theta) \right) \end{aligned}$$

# **Overview**

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

# Logistic Regression: MAP Model

• Maximum-a-posteriori model with prior  $P(\mathbf{\theta}) = N(\mathbf{\theta}|\mathbf{0}, \sigma^2 \mathbf{I})$ :

• 
$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) P(\boldsymbol{\theta})$$
  

$$= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^{n} \frac{1}{1 + e^{-y_i \mathbf{x}_i^T \boldsymbol{\theta}}} N(\boldsymbol{\theta}|\mathbf{0}, \sigma^2 \mathbf{I})$$

$$= \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{n} -\log \frac{1}{1 + e^{-y_i \mathbf{x}_i^T \boldsymbol{\theta}}} -\log N(\boldsymbol{\theta}|\mathbf{0}, \sigma^2 \mathbf{I})$$

$$= \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log \left(1 + e^{-y_i \mathbf{x}_i^T \boldsymbol{\theta}}\right) + \frac{1}{2\sigma^2} \boldsymbol{\theta}^T \boldsymbol{\theta}$$

 No analytic solution; numeric optimization, for instance, using (stochastic) gradient descent.

# Logistic Regression: MAP Model

• Maximum-a-posteriori model with prior  $P(\mathbf{\theta}) = N(\mathbf{\theta}|\mathbf{0}, \sigma^2 \mathbf{I})$ :

• 
$$\boldsymbol{\theta}_{MAP} = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log \left( 1 + e^{-y_i \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\theta}} \right) + \frac{1}{2\sigma_p^2} \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\theta}$$

• Gradient:

• 
$$\frac{\partial}{\partial \theta} \left( \sum_{i=1}^{n} \log \left( 1 + e^{-y_i \mathbf{x}_i^{\mathrm{T}} \mathbf{\theta}} \right) + \frac{1}{2\sigma_p^2} \mathbf{\theta}^{\mathrm{T}} \mathbf{\theta} \right)$$
  
=  $y_i \mathbf{x}_i \left( 1 - \sigma \left( y_i \mathbf{x}_i^{\mathrm{T}} \mathbf{\theta} \right) \right) + \frac{1}{2\sigma_p^2} \mathbf{\theta}$ 

# **Overview**

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

# **Bayes-Optimal Prediction for Classification**

Predictive distribution given the data

• 
$$P(y|\mathbf{x}^*, \mathbf{y}, \mathbf{X}) = \int P(y|\boldsymbol{\theta}, \mathbf{x}^*) P(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}) d\boldsymbol{\theta}$$
  
=  $\int \frac{1}{1 + e^{-y\mathbf{x}^*^T\boldsymbol{\theta}}} N(\boldsymbol{\theta}|\mathbf{0}, \sigma^2 \mathbf{I}) d\boldsymbol{\theta}$ 

- No closed-form solution for logistic regression.
- Possible to approximate by sampling from the posterior.
- Standard approximation: use only MAP model instead of integrating over model space.

# **Overview**

- Basic concepts of Bayesian learning
- Linear regression:
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Posterior distribution over models,
  - Bayesian prediction, predictive distribution,
- Linear classification (logistic regression):
  - Predictive distribution,
  - Maximum-likelihood model,
  - Maximum-a-posteriori model,
  - Bayesian Prediction.
- Nonlinear models: Gaussian processes.

# **Nonlinear Regression**

• Limitation of model discussed so far: only linear dependency between x and  $f_{\theta}(x)$ .



Nonlinear model



Now: nonlinear models.

# **Feature Mappings and Kernels**

- Use mapping  $\phi$  to embed instances  $\mathbf{x} \in X$  in higherdimensional feature space.
- Find linear model in higher-dimensional space, corresponds to non-linear model in input space X.
- Representer theorem:
  - Model  $f_{\theta^*}(\mathbf{x}) = {\theta^*}^T \phi(\mathbf{x})$
  - Has a representation  $f_{\alpha^*}(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* \underbrace{\phi(\mathbf{x}_i)^T \phi(\mathbf{x})}_{=k(\mathbf{x}_i,\mathbf{x})}$
- Feature mapping φ(x) does not have to be computed; only kernel function k(x<sub>i</sub>, x) is evaluated.
- Feature mapping  $\phi(\mathbf{x})$  can therefore be high- or even infinitedimensional.

#### Generalized Linear Regression (Finite-Dimensional Case)

- Assumption 1: Nature generates parameter  $\theta^*$  of a linear function  $f_{\theta^*}(\mathbf{x}) = \phi(\mathbf{x})^T \theta^*$  according to  $p(\theta) = N(\theta | \mathbf{0}, \sigma_p^2 \mathbf{I})$ .
- Assumption 2: Inputs are **X** with feature representation  $\Phi$ ; line *i* of  $\Phi$  contains row vector  $\phi(\mathbf{x}_i)^{\mathrm{T}}$ . Nature generates outputs **y**:

• 
$$y_i = f_{\theta^*}(\mathbf{x}_i) + \epsilon_i$$
 with  $\epsilon_i \sim N(\epsilon | 0, \sigma^2)$ .

• 
$$p(y_i | \mathbf{x}_i, \mathbf{\theta}^*) = N(y_i | \boldsymbol{\phi}(\mathbf{x}_i)^{\mathrm{T}} \mathbf{\theta}^*, \sigma^2)$$


## **Generalized Linear Regression**

- Generalized linear model:
  - $f_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \boldsymbol{\theta}$
  - $\mathbf{y} = \mathbf{\Phi}\mathbf{\theta}$
- Parameter  $\boldsymbol{\theta}$  governed by normal distribution  $N(\boldsymbol{\theta}|\boldsymbol{0}, \sigma_{\mathbf{p}}^2 \mathbf{I})$ .
- Therefore output vector y is also normally distributed.
  - Mean value  $E[\mathbf{y}] = E[\mathbf{\Phi}\mathbf{\theta}] = \mathbf{\Phi}E[\mathbf{\theta}] = \mathbf{0}$ .
  - Covariance  $E[\mathbf{y}\mathbf{y}^{\mathrm{T}}] = \mathbf{\Phi}E[\mathbf{\theta}\mathbf{\theta}^{\mathrm{T}}]\mathbf{\Phi}^{\mathrm{T}} = \sigma_{p}^{2}\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}} = \sigma_{p}^{2}\mathbf{K}.$

• With 
$$K_{ij} = \phi(\mathbf{x}_i)^{\mathrm{T}} \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j).$$

• Remember: 
$$\boldsymbol{\Phi}\boldsymbol{\Phi}^{\mathrm{T}} = \begin{bmatrix} - & \phi(\mathbf{x}_{1}) & - \\ & \dots & \\ - & \phi(\mathbf{x}_{n}) & - \end{bmatrix} \begin{bmatrix} | & & | \\ \phi(\mathbf{x}_{1}) & \vdots & \phi(\mathbf{x}_{n}) \\ | & & | \end{bmatrix} = \mathbf{K}$$

## Generalized Linear Regression (General Case)

- Data generation assumptions:
  - Given inputs **X**, nature generates target values  $\bar{\mathbf{y}} \sim N(\bar{\mathbf{y}}|0, \sigma_p^2 \mathbf{K})$ .
  - Then, nature generates observations  $y_i = \overline{\mathbf{y}}_i + \epsilon_i$  with noise  $\epsilon_i \sim N(\epsilon | 0, \sigma^2)$ .
- Bayesdian inference: determine predictive distribution  $P(\mathbf{y}^*|\mathbf{x}^*, \mathbf{y}, \mathbf{X})$  for new test instance.

## **Reminder: Linear Regression**

Bayes-optimal prediction:

• 
$$y^* = \arg \max_{y} P(y | \mathbf{x}^*, \mathbf{y}, \mathbf{X}) = \overline{\mathbf{\Theta}}^{\mathrm{T}} \mathbf{x}^*$$
  
• With  $\overline{\mathbf{\Theta}} = \left( \mathbf{X}^{\mathrm{T}} \mathbf{X} + \frac{\sigma^2}{\sigma_p^2} \mathbf{I} \right)^{-1} \mathbf{X}^{\mathrm{T}} \mathbf{y}.$ 

• Number of parameters  $\theta_i$  = number of attributes in **x**.

#### **Generalized Linear Regression**

• Mean value of predictive distribution  $P(\mathbf{y}^* | \mathbf{x}^*, \mathbf{y}, \mathbf{X})$  has the form

• 
$$y^* = \sum_{i=1}^n \bar{\alpha}_i k(\mathbf{x}_i, \mathbf{x}^*)$$

• With 
$$\overline{\alpha} = \left( \Phi \Phi^{\mathrm{T}} + \frac{\sigma^2}{\sigma_p^2} \mathbf{I} \right)^{-1} \mathbf{y} = \left( \mathbf{K} + \frac{\sigma^2}{\sigma_p^2} \mathbf{I} \right)^{-1} \mathbf{y}.$$

• Number of parameters  $\alpha_i$  = number of training instances.

## **Example nonlinear regression**

- Example for nonlinear regression
  - Generating nonlinear data by

$$y = \sin(2\pi x) + \varepsilon$$
  $\varepsilon \sim \mathcal{N}(\varepsilon \mid 0, \sigma^2), x \in [0, 1]$ 

Nonlinear kernel

• 
$$k(x, x') = \exp(-\theta |x - x'|)$$

How does the predictive distribution and the posterior over models look like?

#### Predictive distribution



#### Models sampled from posterior



79

# Summary

- Linear regression:
  - Maximum-likelihood model  $\arg \max_{\theta} P(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})$ ,
  - Maximum-a-posteriori model  $\arg \max_{\boldsymbol{\rho}} P(\boldsymbol{\theta} | \boldsymbol{y}, \boldsymbol{X})$ ,
  - Posterior distribution over models  $P(\theta|\mathbf{y}, \mathbf{X})$ ,
  - Bayesian prediction, predictive distribution arg max<sub>y</sub>  $P(\mathbf{y}^* | \mathbf{x}^*, \mathbf{y}, \mathbf{X})$ .
- Linear classification (logistic regression):
  - Predictive distribution  $P(\mathbf{y}^*|\mathbf{x}^*, \mathbf{\theta})$ ,
  - Maximum-likelihood model  $\arg \max_{\mathbf{x}} P(\mathbf{y}|\mathbf{X}, \mathbf{\theta})$ ,
  - Maximum-a-posteriori model  $\arg \max_{\boldsymbol{\rho}} P(\boldsymbol{\theta} | \boldsymbol{y}, \boldsymbol{X})$ ,
  - Bayesian Prediction  $\arg \max_{y} P(\mathbf{y}|\mathbf{x}^*, \mathbf{y}, \mathbf{X}).$
- Nonlinear models: Gaussian processes.