



Basics of Bayesian learning – Basically Bayes

Jan Larsen



Intelligent Signal Processing Group
Informatics and Mathematical Modelling
Technical University of Denmark

www.imm.dtu.dk
jl@imm.dtu.dk



Bayesian learning

T. Bayes “An Essay Towards Solving a Problem in the Doctrine of Chances, *Phil. Trans. Roy. Soc.*, **53**, 370–418, 1783

- The world is uncertain.....

inference: assign probabilities to hypothesis from specific data set

decision theory: choose between actions to minimize loss/risk

- Basic axiom systems for decision theory and inference leads to that rational analysis must corresponds to a Bayesian paradigm [Berger]
- You are probably already doing Bayes – even if you don't know it



Intention

- Crash course in Bayesian learning for those unfamiliar with this paradigm
- Experienced people hopefully gets new inspiration



Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



General Resources

- D. MacKay: *Information Theory, Inference and Learning Algorithms*, <http://www.inference.phy.cam.ac.uk/mackay/itprnn>
- J.O. Berger: *Statistical Decision Theory and Bayesian Analysis*, Springer-Verlag, 2nd edition, 1985.
- C.P. Robert: *The Bayesian Choice: A Decision-Theoretic Motivation*, Springer-Verlag, 1994.
- J.J.K. Ó Ruanaidh and W.J. Fitzgerald: *Numerical Bayesian Methods Applied to Signal Processing*, Springer-Verlag, 1996.



Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



Why Bayesian learning?

- principal framework which combines available uncertain knowledge
 - data, prior etc.
- Bayesian learning is optimal - if you are
- Bayesian learning is typically more robust to mis-specifications and small data sets
- classical learning schemes are special cases
- known to give better performance for most models
- offers model selection as an integrated part



Why Bayesian learning?

- predictions/forecast comes with errorbars (credible sets, highest posterior density credibility set)
- new approaches such as Variational Bayes, Expectation Propagation makes the Bayesian learning computational attractive



Applications of Bayesian learning

- clustering using mixture of Gaussians (Attias, Rasmussen)
- mixture of factor analyzers clustering and dimensionality reduction (Ghahramani+Beal)
- principal components analysis (Bishop)
- independent component analysis
(Højen-Sørensen+Whinter+Hansen, Lee, Attais, Valpola, Miskin+MacKay)
- state-space models, e.g., extended Kalman filters
(Ghahramani+Beal, de Freitas, Niranjana, Wan, Doucet, Gordon)
- time series modeling (Roberts+Penny, Quiñero+Girard+Larsen+Rasmussen)



Applications of Bayesian learning (cont.)

- mixture of experts (Ueda)
- hidden Markov models (MacKay)
- Bayesian networks, graphical models (Heckerman, Jordan, Ghahramani, Bishop, Spiegelhalter)



Software Tools

- VIBES (Bishop, Spiegelhalter, Winn) <http://vibes.sourceforge.net/>
- Matthew Beal <http://www.cse.buffalo.edu/faculty/mbeal/software.html>
- ICA toolbox (DTU) <http://isp.imm.dtu.dk/toolbox>
- Bayes Blocks (HUT) <http://www.cis.hut.fi/projects/bayes/software>
- Bayes Net toolbox (Kevin Murphy) <http://bnt.sourceforge.net>
- ReBEL : Recursive Bayesian Estimation Library (E. Wan)
<http://choosh.ece.ogi.edu/rebel>
- NetLab (Bishop) <http://www.ncrg.aston.ac.uk/netlab/index.php>



Basic ideas of Bayesian framework

- all variables have associated probability densities
- variables not required in the final estimate are integrated out



Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



The basic ingredients

- variables which we want to infer
- problems (unsupervised/supervised)
- data
- model
- prior
- predictive distribution through Bayes theorem
- loss and Bayes risk
- Bayes estimate



The basic ingredients - variables

Predict y from measurement x

- x multivariate input
- y multivariate output
- z multivariate hidden/latent variables

introduction of hidden variables often facilitate model specification



The basic ingredients

Problems

- unsupervised modeling if only x
- predictive modeling if x and y
 - y continuous is regression, e.g., time-series modeling
 - y discrete is classification
- state-space models, mixture models use continuous or discrete hidden variables



The basic ingredients - data

$$\mathcal{D} = \{\boldsymbol{x}_k, \boldsymbol{y}_k\}_{k=1}^N$$

Usually i.i.d. samples



The basic ingredients - models

Models

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

- $\boldsymbol{\theta}$ are model parameters usually not amenable for interpretation
- m index a particular model structure
- we consider usually flexible universal approximation model families
neural networks, Gaussian processes, mixture models

Example

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-[y - f(\mathbf{x}, \mathbf{w})]^2/2\sigma^2)$$

$$\boldsymbol{\theta} = (\sigma^2, \mathbf{w})$$



The basic ingredients - priors

$$p(\theta)$$

expresses the degree of belief

- probability is limit of frequency $\# \text{outcomes} / \# \text{total}$
- properties beliefs lead to same rules as for probabilities, hence using probability to measure belief

.... more on choice of prior later



The basic ingredients - Bayes theorem

Combining information using Bayes theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}, P(B) = \sum_A P(B|A)P(A)$$

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})}$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{prob of data}}$$

$$p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{k=1}^N p(\mathbf{y}_k|\mathbf{x}_k, \boldsymbol{\theta}, m)$$



The basic ingredients - predictive distribution

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}, m) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}, m) \cdot p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta}$$

is the result of Bayesian learning and provides a full conditional distribution for new inputs \mathbf{x}

Relation to classical learning

■ MAP: $p(\boldsymbol{\theta}|\mathcal{D}) = \delta(\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})$, $\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathcal{D})$

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}, m) = p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_{MAP}, m)$$

■ ML: no prior $\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} p(\mathcal{D}|\boldsymbol{\theta})$

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}, m) = p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_{ML}, m)$$



The basic ingredients - transductive learning

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}, m) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}, m) \cdot p(\boldsymbol{\theta}|\mathcal{D}, \mathbf{x}) d\boldsymbol{\theta}$$

Model is updated for every test input \mathbf{x}

$$p(\boldsymbol{\theta}|\mathcal{D}, \mathbf{x}) = \frac{p(\mathcal{D}, \mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D}, \mathbf{x})}$$

Requires a model for input distribution as well!



Predictive distribution example

Model

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

MAP (or ML)

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

Gaussian prior on $\boldsymbol{\theta}$

$$p(y|\mathbf{x}, \mathcal{D}) = \mathcal{N}(\mathbf{x}^\top \hat{\boldsymbol{\theta}}, \sigma^2 + \mathbf{x}^\top \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \mathbf{x})$$



The basic ingredients - loss function

$$L(\mathbf{y}, \hat{\mathbf{y}}) = L(\mathbf{y}(\mathbf{x}), \hat{\mathbf{y}}(\mathbf{x}|\mathcal{D}))$$

defines how close our estimate $\hat{\mathbf{y}}$ is from the truth \mathbf{y} .

Can formally be defined through axiomatic utility theory Berger

Examples

square loss for continuous variable

$$L(\mathbf{y}, \hat{\mathbf{y}}) = (\mathbf{y} - \hat{\mathbf{y}})^2$$

zero-one loss for classification $y, \hat{y} \in [1; C]$

$$L(y, \hat{y}) = \begin{cases} 0, & y = \hat{y} \\ 1, & y \neq \hat{y} \end{cases}$$



The basic ingredients - risk

Frequentist risk

$$R(\mathbf{y}, \hat{\mathbf{y}}) = \int L(\mathbf{y}, \hat{\mathbf{y}}(\mathcal{D})) \cdot p(\mathcal{D}) d\mathcal{D}$$

average over all possible data sets

Bayes risk

$$\rho(\hat{\mathbf{y}}|\mathcal{D}) = \int L(\mathbf{y}, \hat{\mathbf{y}}(\mathcal{D})) \cdot p(\mathbf{y}|\mathbf{x}, \mathcal{D}) d\mathbf{y}$$

average w.r.t. predictive distribution and **conditioned** on data



The basic ingredients - risk

Integrated frequency risk

$$\begin{aligned} r(\hat{\mathbf{y}}) &= \int L(\mathbf{y}, \hat{\mathbf{y}}(\mathcal{D})) \cdot p(\mathbf{y}|\mathbf{x}, \mathcal{D}) p(\mathcal{D}) d\mathbf{y} d\mathcal{D} \\ &= \int \rho(\hat{\mathbf{y}}|\mathcal{D}) \cdot p(\mathcal{D}) d\mathcal{D} \\ &= \int L(\mathbf{y}, \hat{\mathbf{y}}(\mathcal{D})) \cdot p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{y} \end{aligned}$$



Outline

- Why Bayesian learning?
- Basic ingredients
- **Bayes estimators**
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



Bayes estimator

$$\hat{\mathbf{y}}_B = \arg \min_{\hat{\mathbf{y}}} r(\hat{\mathbf{y}}) = \arg \min_{\hat{\mathbf{y}}} \rho(\hat{\mathbf{y}}|\mathcal{D})$$

Risk of Bayes estimator

$$r(\hat{\mathbf{y}}_B)$$



Bayes estimator in continuous case

Square loss

$$\hat{\mathbf{y}}_B = \arg \min_{\hat{\mathbf{y}}} \rho(\hat{\mathbf{y}}|\mathcal{D}) = \arg \min_{\hat{\mathbf{y}}} \int (\mathbf{y} - \hat{\mathbf{y}})^2 \cdot p(\mathbf{y}|\mathbf{x}, \mathcal{D}) d\mathbf{y}$$

$$\hat{\mathbf{y}}_B = \int \mathbf{y} \cdot p(\mathbf{y}|\mathbf{x}, \mathcal{D}) d\mathbf{y} = E_{pred}[\mathbf{y}]$$

Absolute loss

$$\hat{\mathbf{y}}_B = \arg \min_{\hat{\mathbf{y}}} \int |\mathbf{y} - \hat{\mathbf{y}}| \cdot p(\mathbf{y}|\mathbf{x}, \mathcal{D}) d\mathbf{y}$$

$\hat{\mathbf{y}}_B$ is the median of predictive distribution



Bayes estimator for classification

Loss matrix

Penalty of estimating class $\hat{y} \in [1; C]$ if the truth is class $y \in [1; C]$

$$L(y, \hat{y}) = \begin{cases} 0 & \text{if } \hat{y} = y \text{ (correct decision)} \\ l(y, \hat{y}) & \text{if } \hat{y} \neq y \in [1; C] \\ t & \text{if } \hat{y} = C + 1 \text{ (rejection)} \end{cases}$$

Zero-one loss with rejection

$$L(y, \hat{y}) = \begin{cases} 0 & \text{if } \hat{y} = y \text{ (correct decision)} \\ 1 & \text{if } \hat{y} \neq y \in [1; C] \\ t & \text{if } \hat{y} = C + 1 \text{ (rejection)} \end{cases}$$



Bayes decision rule - general loss

$$\rho(\hat{y}|\mathcal{D}) = \sum_y L(y, \hat{y}) \cdot p(y|\mathbf{x}, \mathcal{D})$$

For $\hat{y} = C + 1$ then $\sum_y L(y, \hat{y}) \cdot p(y|\mathbf{x}, \mathcal{D}) = t \sum_y p(y|\mathbf{x}, \mathcal{D}) = t$

$$\hat{y}_B = \begin{cases} k & \text{if } \min_{\hat{y} \leq C} \sum_y L(y, \hat{y}) \cdot p(y|\mathbf{x}, \mathcal{D}) < t \\ C + 1 & \text{otherwise} \end{cases}$$



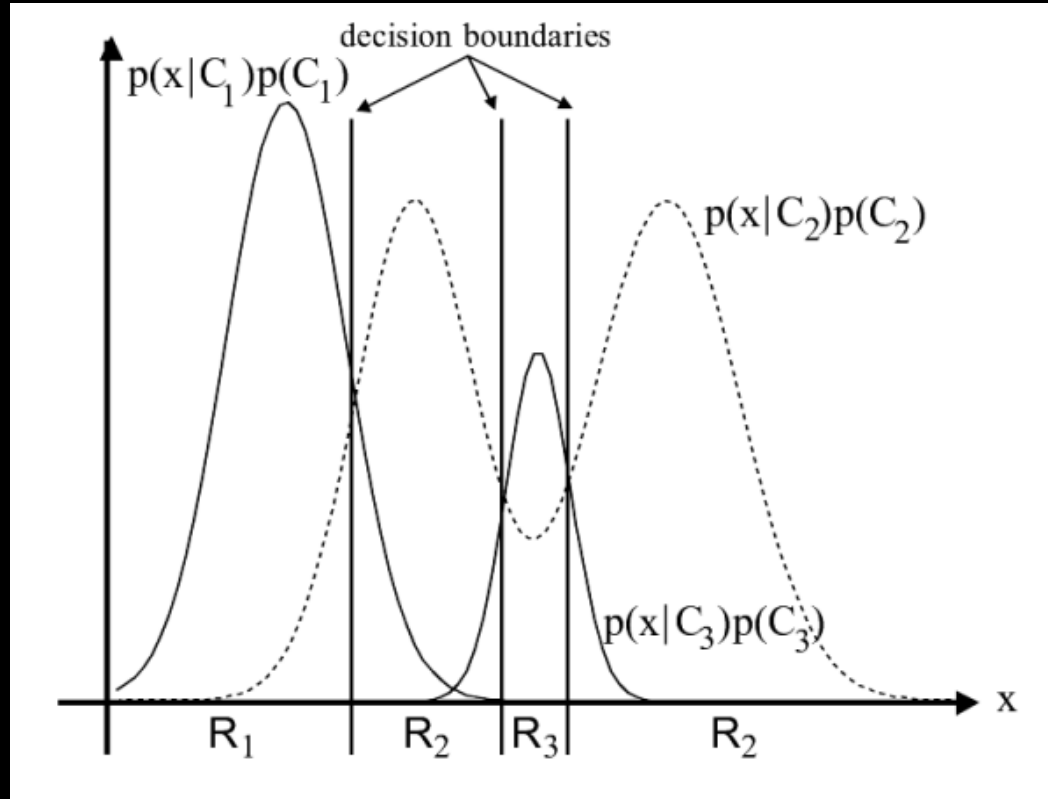
Bayes decision rule - zero-one loss

$$\sum_y L(y, \hat{y}) \cdot p(y|\mathbf{x}, \mathcal{D}) = \begin{cases} 1 - p(j|\mathbf{x}, \mathcal{D}), & \hat{y} = j \in [1; C] \\ t, & \hat{y} = C + 1 \end{cases}$$

$$\hat{y}_B = \begin{cases} \arg \max_{y \leq C} p(y|\mathbf{x}, \mathcal{D}), \text{ and } \text{prob} > 1 - t \\ C + 1 \text{ if all } p(y|\mathbf{x}, \mathcal{D}) \leq 1 - t \end{cases}$$

$p(y|\mathbf{x}, \mathcal{D}) \geq 1/C$ which means $1 - t > 1/C$ for rejection to occur

Bayes classifier



decision boundaries are specified by $p(y = i|\mathbf{x}, \mathcal{D}) = p(y = j|\mathbf{x}, \mathcal{D})$



Optimality through admissibility

Admissibility

$\hat{\mathbf{y}}_0$ is **inadmissible** if

$$\forall \mathbf{y}, R(\mathbf{y}, \hat{\mathbf{y}}_0) \geq R(\mathbf{y}, \hat{\mathbf{y}}_1), \exists \mathbf{y}_0, R(\mathbf{y}_0, \hat{\mathbf{y}}_0) > R(\mathbf{y}_0, \hat{\mathbf{y}}_1)$$

Generalized Bayes estimator is admissible under regularity conditions

- $p(\mathbf{y}|\mathbf{x}) = p(\mathbf{y}|\mathbf{x}, \mathcal{D})p(\mathcal{D})/p(\mathcal{D}|\mathbf{y}, \mathbf{x}) > 0$ for all data
- Bayes risk is finite (might fail for generalized - improper prior case)
- $R(\mathbf{y}, \hat{\mathbf{y}})$ is continuous in \mathbf{y}



Optimality through generalization error

Randomized estimators

use predictive distribution $\hat{p} = p(\mathbf{y}|\mathbf{x}, \mathcal{D})$ as random estimator rather than point estimate $\hat{\mathbf{y}}$

Kullback-Leibler information as average loss between distributions

$$\begin{aligned} KL(p|\hat{p}) &= \int p(\mathbf{y}|\mathbf{x}) \log \frac{p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}|\mathbf{x}, \mathcal{D})} d\mathbf{y} \\ &= E\{L(p, \hat{p})\} \end{aligned}$$

■ the loss is defined as

$$L(p, \hat{p}) = \log p(\mathbf{y}|\mathbf{x}) - \log p(\mathbf{y}|\mathbf{x}, \mathcal{D})$$

■ $KL \geq 0$ with 0 if and only if $p \equiv \hat{p}$



On the property of KL

Inequality

$$-\log \lambda \geq 1 - \lambda \text{ for } \lambda > 0 \quad \text{and} \quad -\log \lambda = 1 - \lambda \text{ for } \lambda = 1$$

Proof Define $\lambda(y) = \hat{p}(y)/p(y)$. That is $-\log \lambda(y) = \log p(y)/\hat{p}(y)$

$$\begin{aligned} KL &= \int p(y) [-\log \lambda(y)] dy \\ &\geq \int p(y) [1 - \lambda(y)] dy \\ &= \int p(y) \left[1 - \frac{\hat{p}(y)}{p(y)} \right] dy \\ &= 0 \end{aligned}$$



Optimality through generalization error

$$KL = \underbrace{G}_{\text{generalization error}} - \underbrace{H(p)}_{\text{entropy}}$$

Average generalization error

$$\Gamma = E_{\mathcal{D}} E_{\mathbf{x}} \{G\} = - \int \log p(\mathbf{y}|\mathbf{x}, \mathcal{D}) p(\mathcal{D}) p(\mathbf{x}) d\mathbf{y} d\mathbf{x} d\mathcal{D}$$

is integrated frequency risk (also averaged w.r.t $p(\mathbf{x})$) up to a constant

L.K. Hansen: “Bayesian Averaging is Well-Tempered,” NIPS99, 265–271, 2000
shows optimality in generalization error



Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework



More on priors

- **subjective priors:** consider relative likelihood of various parameters values
- **empirical priors:** obtained from past experience data
- **structural priors:**
 - independence of some parameters?
 - imposing functional smoothness
 - invoking constraints
- **convenience priors:**
 - nice functional form in order to make calculations simple
 - conjugate priors: posterior and prior have same shape

exponential family is important



More on priors

- **hierarchical:** $p(\theta) = \int p(\theta|\lambda)p(\lambda) d\lambda$
- **non-informative:** make the influence of the prior as small as possible
- **improper:** improper priors do not integrate to one. Leads to generalized Bayes estimator which typically also is admissible



Non-informative priors

- discrete parameter taking C values: $p(\theta) = 1/C$
- continuous parameter $p(\theta) = 1$ which is improper $\int p(\theta) d\theta = \infty$

Location parameter

invariance to choice of parameterization

$$\eta = \theta + c, \forall c$$

$$\int p(\eta) d\eta = \int p(\theta) d\theta$$
$$\int p(\eta - c) d\eta = \int p(\eta) d\eta$$

for all η . With $\eta = c$ then $p(c) = p(0)$ thus $p(\theta) = 1$



Non-informative priors

Scale parameter $\eta = c\theta, \forall c > 0$

$$\int p(\eta) d\eta = \int p(\theta) d\theta$$
$$\int p(\eta c^{-1}) c^{-1} d\eta = \int p(\eta) d\eta$$

thus with $\eta = c$ and $p(c) = c^{-1}p(1)$. Setting $p(1) = 1$ then $p(\theta) = \theta^{-1}$

Jeffrey's non-informative

$$p(\boldsymbol{\theta}) = \sqrt{\det \mathbf{I}(\boldsymbol{\theta})}$$

$$\mathbf{I}(\boldsymbol{\theta}) = -E[\partial^2 \log p(\mathbf{x}|\boldsymbol{\theta}) / \partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top]$$



Optimizing hyperparaters using evidence

$$p(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\boldsymbol{\lambda})$$

$\boldsymbol{\lambda}$ are hyperparameters

Evidence - marginal likelihood

$$p(\mathcal{D}) = p(\mathcal{D}|\boldsymbol{\lambda}) = \int p(\mathcal{D}|\boldsymbol{\theta}, \boldsymbol{\lambda})p(\boldsymbol{\theta}|\boldsymbol{\lambda}) d\boldsymbol{\theta}$$

$$\boldsymbol{\lambda}_{ML-II} = \arg \max_{\boldsymbol{\lambda}} p(\mathcal{D}|\boldsymbol{\lambda})$$



Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



Generalization is the ultimate frequentist objective

Assume that the underlying system is stationary (time-invariant).

How well are we doing on future data?

Generalization error

$$G(\mathcal{D}) = \int L(\mathbf{y}, \hat{\mathbf{y}}(\mathcal{D})) \cdot p(\mathbf{x}, \mathbf{y}) d\mathbf{x}d\mathbf{y}$$

- $\{\mathbf{x}; \mathbf{y}\}$ is a sample independent of all samples in the training set
- $L(\cdot)$ is any loss function



Generalization error decomposition

Average generalization error / averaged integrated risk

$$\Gamma = E_{\mathcal{D}} \{G(\mathcal{D})\}$$

Decomposition

$$\Gamma = \textit{Inherent Noise} + \textit{Bias} + \textit{Variance}$$

- *Inherent Noise* (minimal Bayes risk) can not be modeled
- *Bias* is due to an incomplete model
- *Variance* is due to a finite training set



A general B/V decomposition

Heskes, 1998

Properties

- bias only depends on true and average distributions
- variance is non-negative not a function of true distribution, and zero only if and only if distributions are equal
- mean-square error is a special case

Decomposition

$$\begin{aligned}\Gamma = E_{\mathcal{D}} \{G(\mathcal{D})\} &= H(p) + KL(p|\hat{p}) \\ &= H(p) + KL(p|\bar{p}) + E_{\mathcal{D}} \{KL(\bar{p}|\hat{p})\} \\ &= \textit{Inherent Noise} + \textit{Bias} + \textit{Variance}\end{aligned}$$

with average model

$$\bar{p} = Z^{-1} \cdot \exp(E_{\mathcal{D}} \{\log \hat{p}\})$$



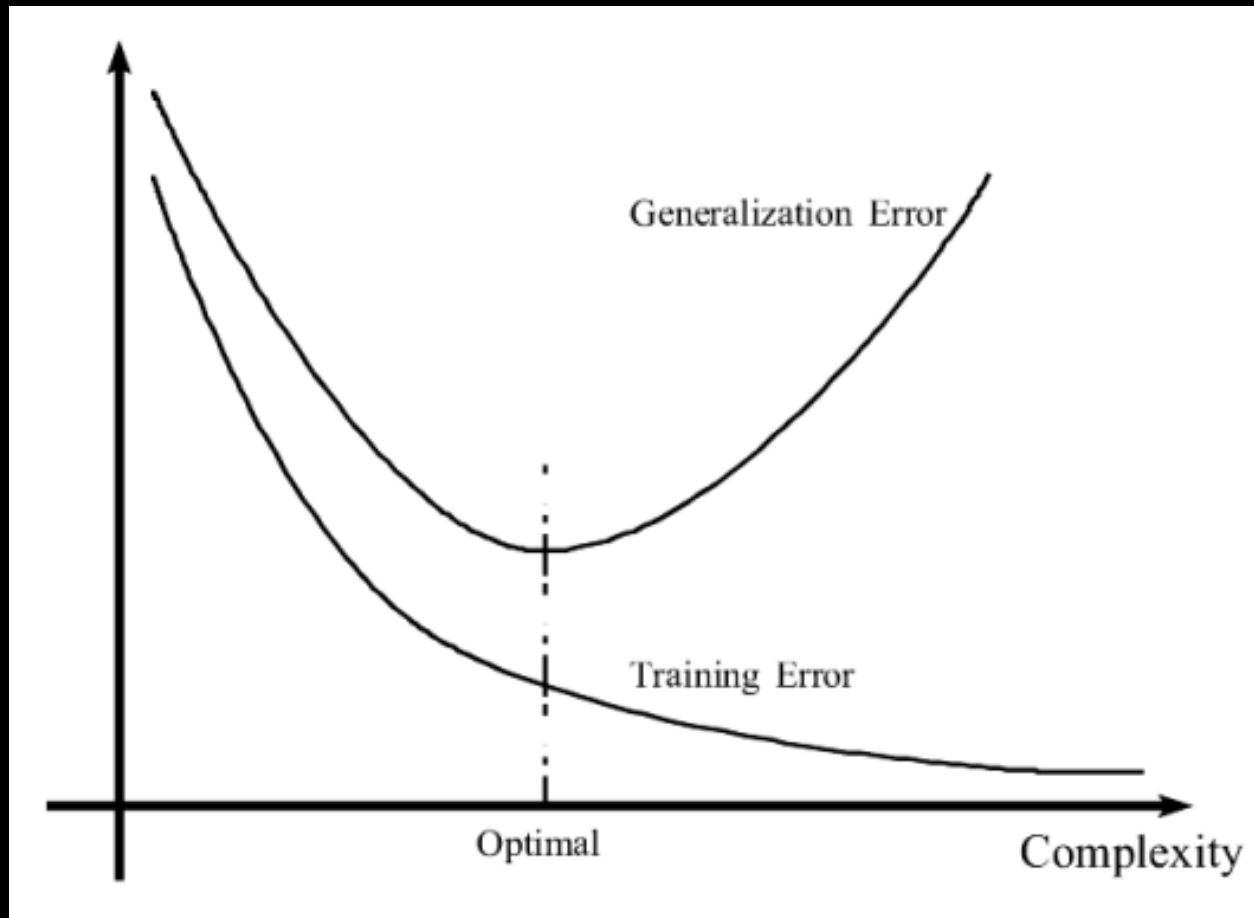
Bias/variance dilemma

The model should learn rather than memorize training data

Model complexity

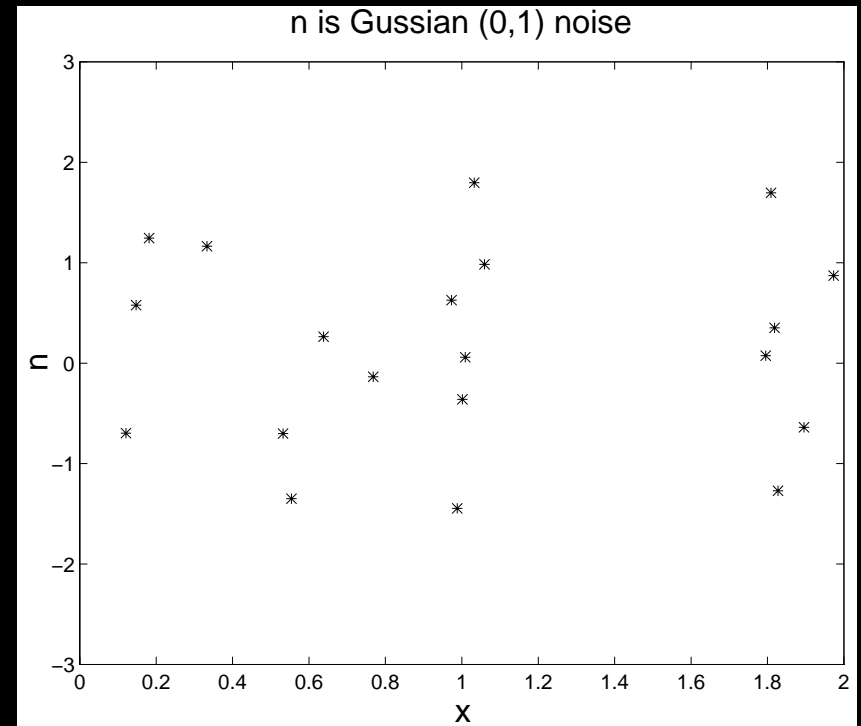
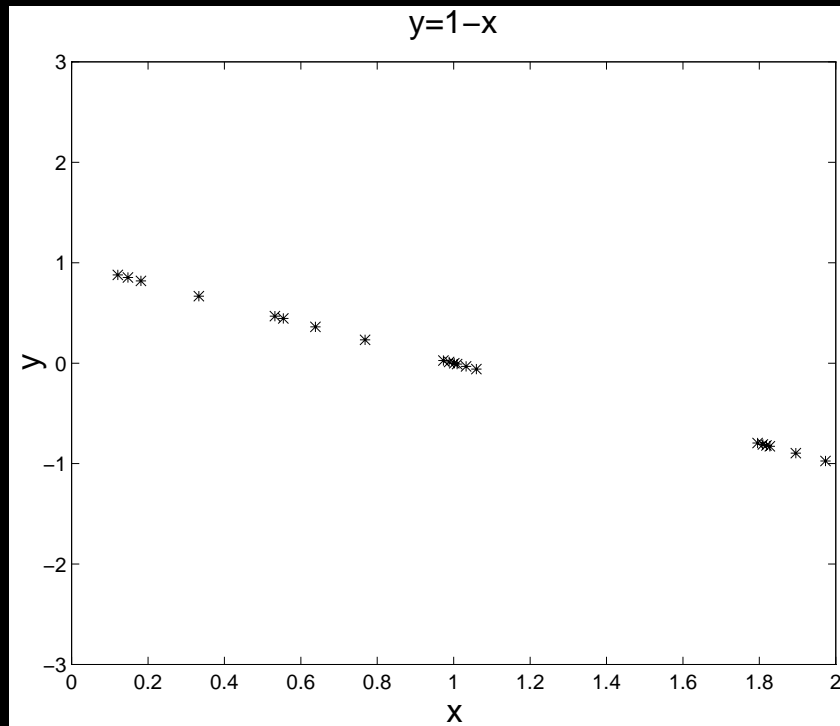
low	high
cannot fit	fits to noise

Bias/variance dilemma



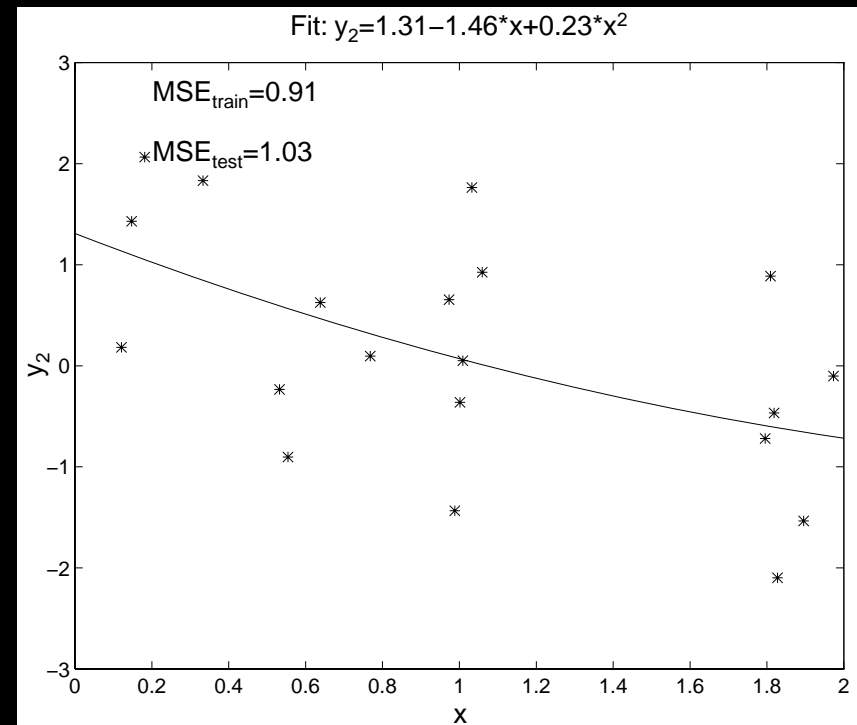
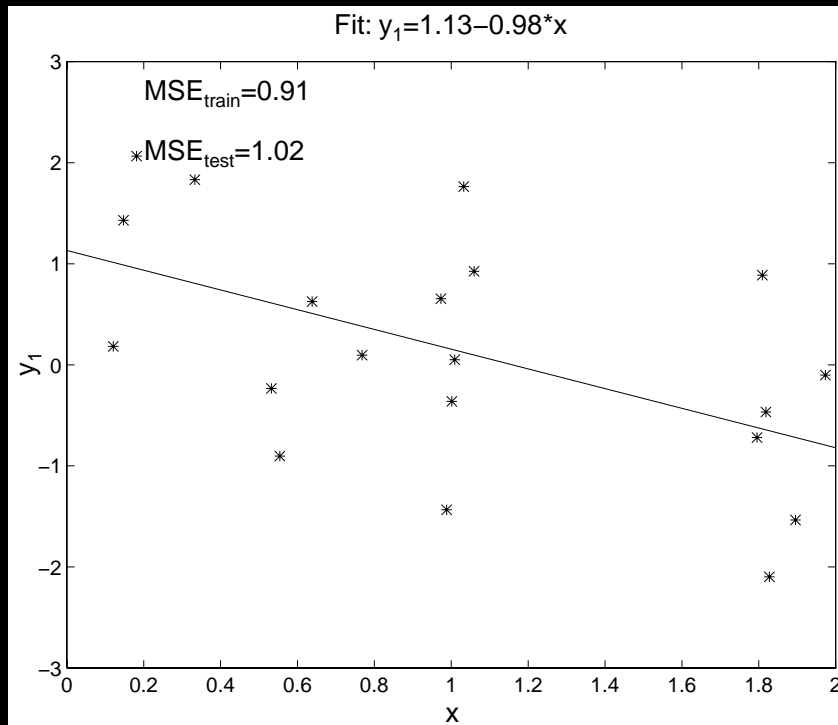


Overfitting



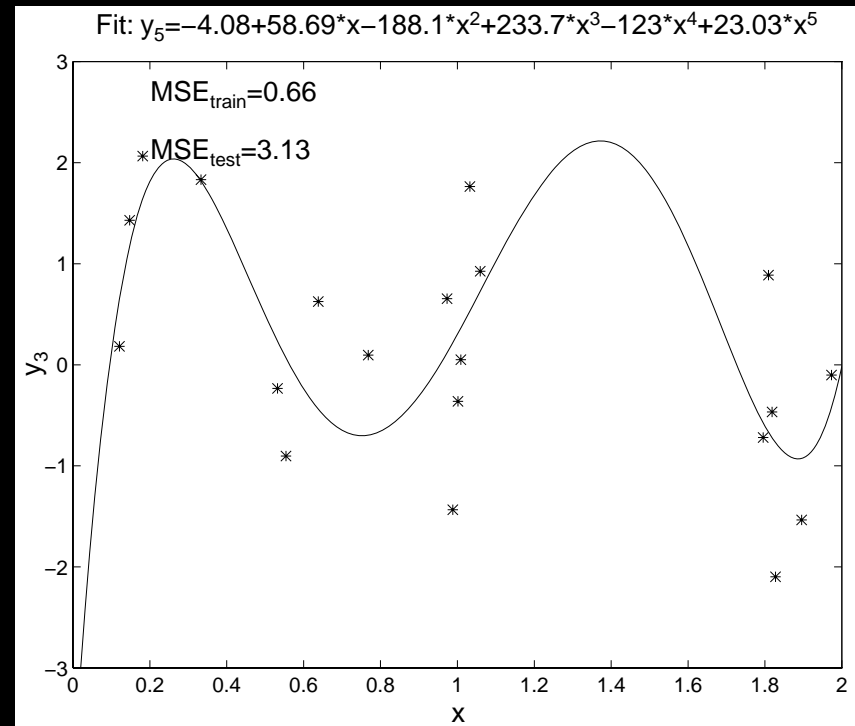
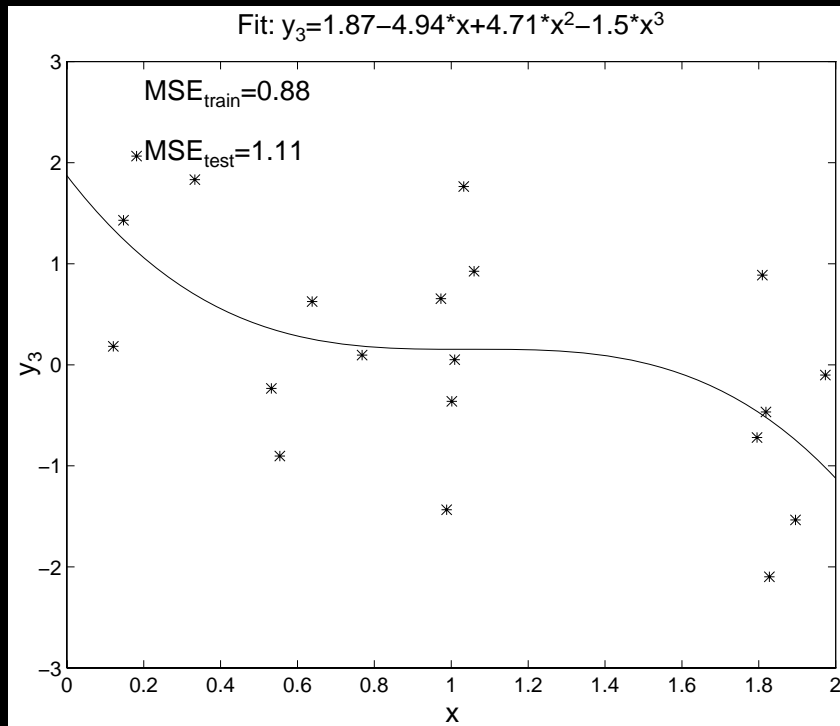


Overfitting



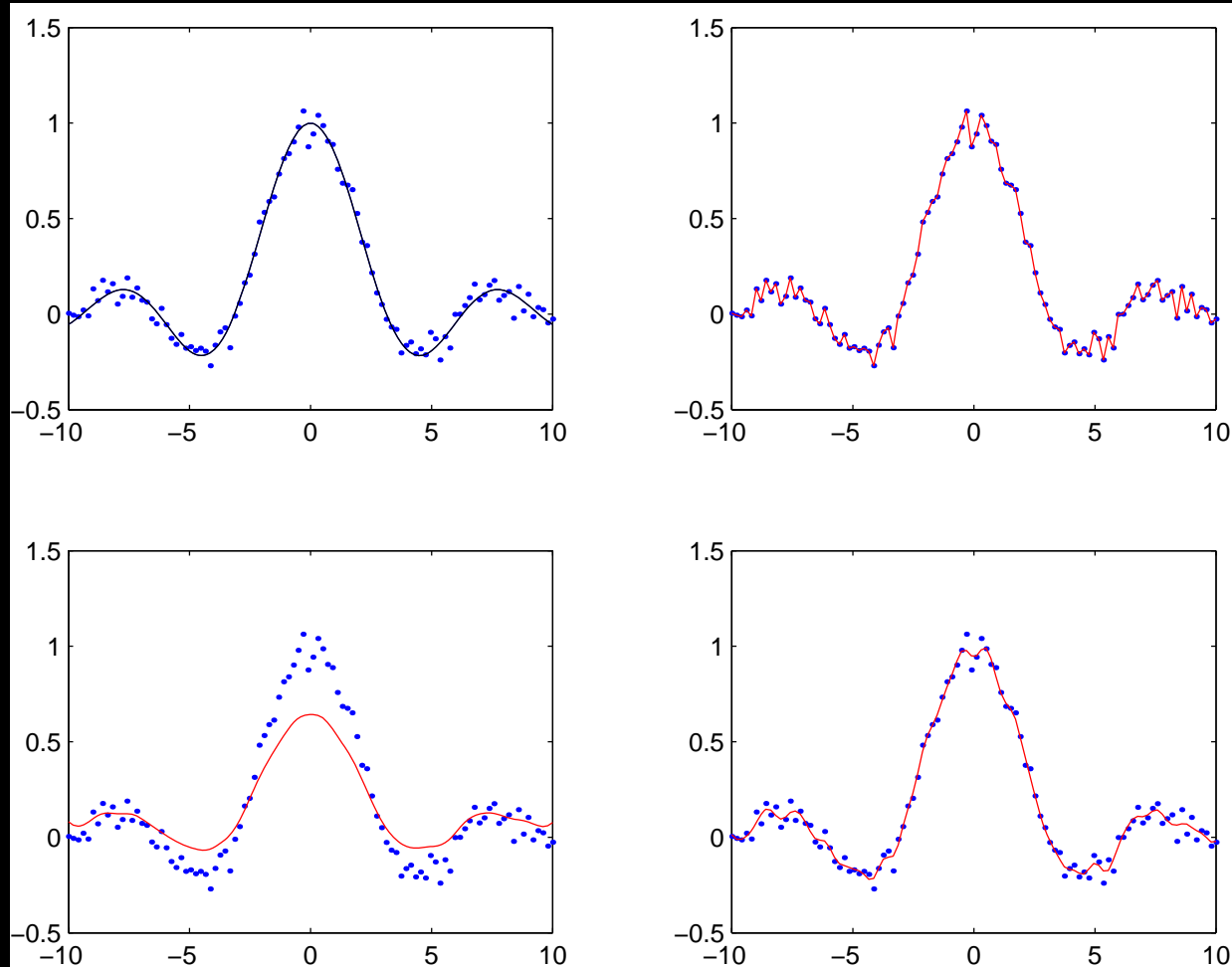


Overfitting





Overfitting in an RBF network





Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



Generalization estimation

Approaches

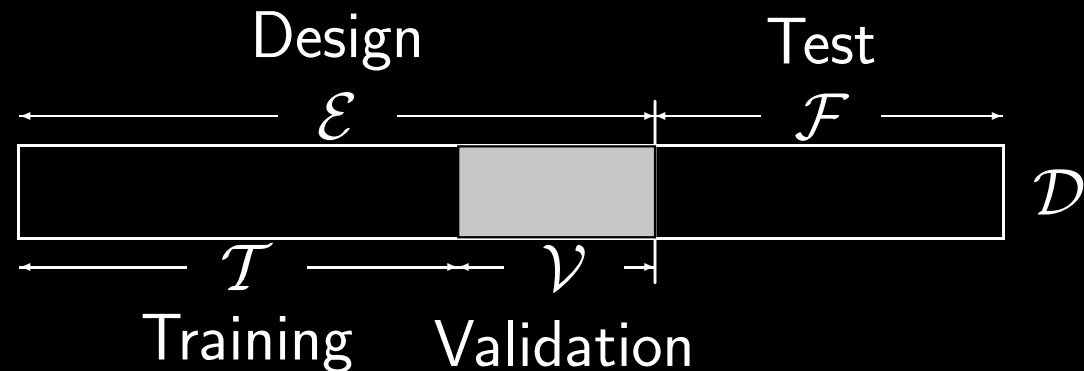
- Asymptotic theory leading to algebraic estimates of the **average generalization error**
- Resampling approaches (cross-validation, jackknife, and bootstrap) of the **generalization error** or **average generalization error**

Purposes

- Assessing the final quality and reliability of the model
- Model selection



Limited data is always a challenge



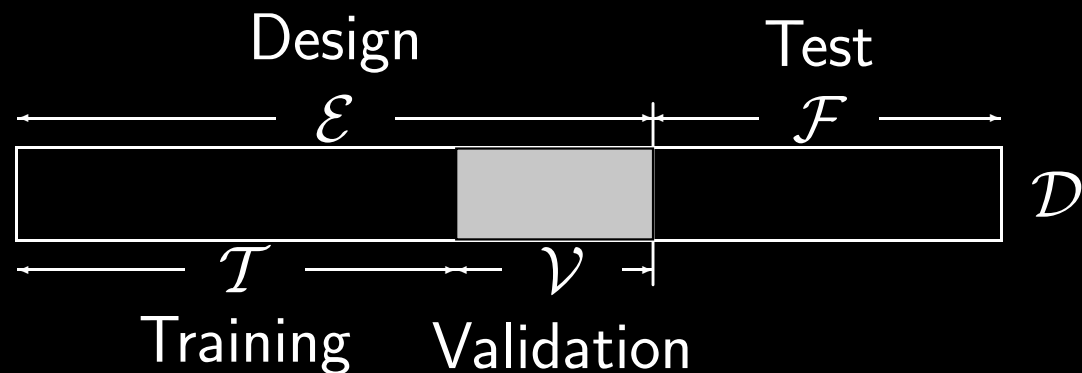
Design/Test Split

- Test set is exclusively used for final assessment of model designed from \mathcal{E}
- Objective is high generalization ability and reliable assessment

J. Larsen and C. Goutte: "On Optimal Data Split for Generalization Estimation and Model Selection," in Proceedings of the IEEE NNSP Workshop IX, pp. 225–234, 1999



Limited data is always a challenge



Training/Validation Split

- Model is trained on training set. Validation set is used to select optimal model or tune additional hyperparameters
- Objective is high generalization ability.



Hold-Out Cross-Validation

Data Splitting

- $\gamma = i/N$ $i = 1, 2, \dots, N - 1$ is the *split ratio*
- $N_{\mathcal{F}} = \gamma N$ for testing and $N_{\mathcal{E}} = (1 - \gamma)N$ for design

HO Estimate

$$\hat{G}_{\text{HO}} = N_{\mathcal{F}}^{-1} \sum_{k \in \mathcal{F}} -\log p(\mathbf{y}(k)|x(k), \mathcal{E}, m)$$

Property HO is an unbiased estimate of the generalization error for i.i.d. samples



Quality of the HO Estimator

$$\begin{aligned} MSE_{HO}(\gamma) &= E_{\mathcal{D}} \left\{ \left(\hat{G}_{HO} - G^* \right)^2 \right\} \\ &= \underbrace{E_{\mathcal{D}} \left\{ \left(\hat{G}_{HO} - G(\mathcal{D}) \right)^2 \right\}}_{\text{variance}} + \underbrace{E_{\mathcal{D}} \left\{ \left(G(\mathcal{D}) - G^* \right)^2 \right\}}_{\text{bias}} \end{aligned}$$

where G^* is the minimum achievable gen. error for the current model, i.e., infinitely data

Property

- Bias \downarrow as $\gamma \downarrow$
- Variance \uparrow as $\gamma \downarrow$



K -fold Cross-Validation

Procedure

- Split data into K disjoint subsets \mathcal{F}_j (approx. equal sizes)
- $K = \lfloor 1/\gamma \rfloor$, for $\gamma < 1/2$,
 $\gamma = \{1/N, 1/(N-1), \dots, 1/2, \dots, 1-1/N\}$
- Evaluate on each subset the model designed on the remaining data, $\mathcal{E}_j = \mathcal{D} \setminus \mathcal{F}_j$

Estimate

$$\hat{\Gamma}_{\text{KCV}} = \frac{1}{N} \sum_{j=1}^K \sum_{k \in \mathcal{F}_j} -\log p(\mathbf{y}(k)|x(k), \mathcal{E}_j, m)$$

Property Unbiased estimator the average generalization error, Γ , based on $N_{\mathcal{E}}$ data.



Quality of KCV Estimator

$$\begin{aligned} MSE_{KCV}(\gamma) &= E_{\mathcal{D}} \left\{ \left(\hat{\Gamma}_{KCV} - G^* \right)^2 \right\} \\ &= \underbrace{E_{\mathcal{D}} \left\{ \left(\hat{\Gamma}_{KCV} - \Gamma \right)^2 \right\}}_{\text{variance}} + \underbrace{E_{\mathcal{D}} \left\{ \left(\Gamma - G^* \right)^2 \right\}}_{\text{bias}} \end{aligned}$$



Results

- Analytical expression for opt. design/test splits for location parameter model. Tends to hold asymp. for other models:
 - HO: $\gamma_{opt} \rightarrow 1$, as $N \rightarrow \infty$
 - KCV: $\gamma_{opt} = 1/N$, (LOO)
- Model selection using KCV:
 - $\gamma_{opt} \rightarrow 1$, as $N \rightarrow \infty$
 - LOO seems to optimal when N is small both wrt. generalization error and probability of selecting correct model
- Model selection using HO:
 - Conflict between opt. gen. error and probability of selecting correct model



Algebraic Generalization Error Approach

Properties

- Asymptotic estimates valid for large training sets
- Various assumptions on model bias and example dependencies
- No data need to be set aside for validation



Estimator	Model	Est.
Exact (Hansen 93)	Unbiased lin. zero-mean Gaus data	ML
FPE (Akaike 69)	Unbiased, no prior	MSE
FPER (Larsen 94)	Unbiased with prior	Pen. MSE
AIC (Akaike 73)	Unbiased no prior	ML
AICc (Hurvich&Tsai 1989)	Unbiased no prior	ML
GEN (Larsen 92, 2000)	No restrictions, auto corr data	Pen. MSE/MAP
GPE (Moody 91)	nonlin with prior	MSE
NIC (Murata 94)	nonlin(NN),nested, i.i.d. data	MAP
TIC (Takeuchi 76)	nonlin i.i.d. data	MAP
GIC (Konishi&Kitagawa 96)	general i.i.d.	MAP,Bayes
DIC (Spiegelhalter et al. 2002)	general i.i.d.	Bayes



GEN

Major Assumptions

- Asymptotic validity, $o(1/N)$
- Applies to bias and regularized models
- Estimates average generalization error

MAP approach

$$C_{\mathcal{D}}(\boldsymbol{\theta}) = N^{-1} \sum_{k=1}^N \ell(\mathbf{y}(k) | \mathbf{x}(k), \boldsymbol{\theta}) + R(\boldsymbol{\theta}) = S_{\mathcal{D}}(\boldsymbol{\theta}) + R(\boldsymbol{\theta})$$



GEN Estimator

$$\Gamma_{GEN} = E_{\mathcal{D}}\{S_{\mathcal{D}}(\hat{\boldsymbol{\theta}})\} + \frac{m_{eff}}{N} \\ - \frac{A}{N} \cdot \frac{\partial R}{\partial \boldsymbol{\theta}^{\top}}(\boldsymbol{\theta}^*) \mathbf{J}^{-1}(\boldsymbol{\theta}^*) \frac{\partial R}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}^*)$$

- Optimal parameters: $\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} G(\boldsymbol{\theta})$ where the expected cost: $C(\boldsymbol{\theta}) = E_{\mathcal{D}}\{C_{\mathcal{D}}(\boldsymbol{\theta})\} = G(\boldsymbol{\theta}) + R(\boldsymbol{\theta})$
- For practical use an unbiased $o(1/N)$ estimator is obtained by, neglecting the expectation, replacing $\boldsymbol{\theta}^*$ by $\hat{\boldsymbol{\theta}}$, and \mathbf{J} by $\mathbf{J}_{\mathcal{D}}$



GEN Estimator

- Effective number of parameters

$$m_{eff} = \text{tr} \left[\mathbf{J}^{-1}(\boldsymbol{\theta}^*) \left(\mathbf{K}(0) + \sum_{n=1}^{\bar{M}} \frac{N-n}{N} (\mathbf{K}(n) + \mathbf{K}^\top(n)) \right) \right]$$

$$= \text{tr} \left[\mathbf{J}^{-1}(\boldsymbol{\theta}^*) \mathbf{L} \right]$$

where $\bar{M} = \min(M, N-1)$, M is the time dependence length
(for i.i.d. examples $M = 0$),

$$\mathbf{L} = \bar{M} + 1 - \bar{M}(\bar{M} + 1)/2N$$

- $\mathbf{K}(n) = E\{\partial \ell(k)/\partial \boldsymbol{\theta} \cdot \partial \ell(k+n)/\partial \boldsymbol{\theta}^\top\}$ with

$$\ell(k) \equiv \ell(\mathbf{y}(k)|\mathbf{x}(k), \boldsymbol{\theta}^*)$$

- $\mathbf{J}(\boldsymbol{\theta})$ is the Hessian matrix of the expected cost function $C(\boldsymbol{\theta})$



Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- **Bayesian model selection**
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



Bayesian model selection

Bayes optimal decision rule (under 0/1 loss function) leads to the optimal model

$$m_{opt} = \arg \max_m p(m|\mathcal{D})$$

$$p(m|\mathcal{D}) = \frac{p(\mathcal{D}|m)P(m)}{\sum_{m=1}^M p(\mathcal{D}|m)P(m)}$$

is the probability of the model given data

The Bayes model probability is only correct if the true model is among the hypothesis models!



Probabilistic Model Selection

Uniform Model Belief

In the case of equal model priors, i.e., $P(m) = 1/M$, the model selection concerns computing the **evidence** $p(\mathcal{D}|m)$

Evidence

$$p(\mathcal{D}|m) = \int p(\mathcal{D}, \boldsymbol{\theta}|m) d\boldsymbol{\theta} = \int p(\mathcal{D}|\boldsymbol{\theta}, m)p(\boldsymbol{\theta}|m) d\boldsymbol{\theta}$$

where

- $\boldsymbol{\theta}$ are model parameters
- $p(\mathcal{D}|\boldsymbol{\theta}, m)$ is the likelihood
- $p(\boldsymbol{\theta}|m)$ is the prior which is normally assumed vague and normalizable



Approximations

- Laplace approximation
- BIC approximation (large sample Laplace)
- variational Bayes (“EM like”)
- expectation propagation
- ensemble learning
- (annealed) importance sampling, particle filtering
- Gibbs sampling
- Markov chain Monte Carlo methods (Metropolis-Hastings, Parallel tempering (gets marginal likelihood), particle path filter)



Models with hidden variables

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}) = \int p(\mathbf{y}, \mathbf{z}|\mathbf{x}, \boldsymbol{\theta}, m) \cdot p(\boldsymbol{\theta}|\mathcal{D}) d\mathbf{z}d\boldsymbol{\theta}$$

In particular ensemble learning and Variational Bayes are useful

- approximate by integrating w.r.t. proposal distributions $q(\mathbf{z}, \boldsymbol{\theta})$
- ensemble learning uses simple functional forms often fully factorized $q(\mathbf{z}, \boldsymbol{\theta}) = \prod_i q(z_i) \prod_j q(\theta_j)$
- Variational Bayes uses functional forms as priors partly factorized $q(\mathbf{z}, \boldsymbol{\theta}) = q(\mathbf{z}) \prod_j q(\theta_j)$
- works for on-line models Ghahramani, Valpola



Variational Bayes learning

Key ingredients

- Use factorized approximate posterior distribution

$$q(\boldsymbol{\theta}, \mathbf{z}) = q(\boldsymbol{\theta}) \cdot q(\mathbf{z})$$

- Use Jensen's inequality to bound marginal likelihood (evidence)

Resources

- www.variational-bayes.org
- Z. Ghahramani, C. Bishop, G. Hinton, M.I. Jordan, D. MacKay, C. Rasmussen, R. Neal, M. Beal



Variational Bayes learning

$$\begin{aligned}\log p(\mathcal{D}|m) &= \log \int p(\mathcal{D}, \boldsymbol{\theta}, \mathbf{z}|m) d\mathbf{z}d\boldsymbol{\theta} \\ &= \log \int \frac{q(\boldsymbol{\theta})q(\mathbf{z})p(\mathcal{D}, \boldsymbol{\theta}, \mathbf{z}|m)}{q(\boldsymbol{\theta})q(\mathbf{z})} d\mathbf{z}d\boldsymbol{\theta} \\ &= \mathcal{F}_m(q(\boldsymbol{\theta}), q(\mathbf{z}), \mathcal{D}) + KL_{post}(q||p)\end{aligned}$$

Variational free energy \mathcal{F} bounds the evidence

$$\begin{aligned}\log p(\mathcal{D}|m) &\geq \mathcal{F}_m(q(\boldsymbol{\theta}), q(\mathbf{z}), \mathcal{D}) \\ &= \int q(\boldsymbol{\theta})q(\mathbf{z}) \log \frac{p(\mathcal{D}, \boldsymbol{\theta}, \mathbf{z}|m)}{q(\boldsymbol{\theta})q(\mathbf{z})} d\mathbf{z}d\boldsymbol{\theta}\end{aligned}$$



Variational Bayes learning

Kullback-Liebler between true and approximate posterior

$$KL_{post}(q|p) = \int q(\boldsymbol{\theta})q(\mathbf{z}) \log \frac{p(\boldsymbol{\theta}, \mathbf{z}|\mathcal{D}, m)}{q(\boldsymbol{\theta})q(\mathbf{z})} d\mathbf{z}d\boldsymbol{\theta}$$



Variational Bayes learning

Maximize $\mathcal{F}_m(q(\boldsymbol{\theta}), q(\mathbf{z}), \mathcal{D})$ w.r.t. $q(\boldsymbol{\theta})$ and $q(\mathbf{z})$ i.e., minimize $KL_{post}(q||p)$

“E-step” - estimate posterior over hidden variables

$$q^{(j+1)}(\mathbf{z}) \propto \exp \left[\int \log p(\mathcal{D}, \mathbf{z} | \boldsymbol{\theta}, m) q^{(j)}(\boldsymbol{\theta}) \right] d\boldsymbol{\theta}$$

“M-step” - estimate posterior over parameters

$$q^{(j+1)}(\boldsymbol{\theta}) \propto p(\boldsymbol{\theta}) \exp \left[\int \log p(\mathcal{D}, \mathbf{z} | \boldsymbol{\theta}, m) q^{(j+1)}(\mathbf{z}) \right] d\mathbf{z}$$

Reduces to classical EM when $q(\boldsymbol{\theta}) = \delta(\boldsymbol{\theta} - \boldsymbol{\theta}_0)$



Exponential distribution family

Complete likelihood

$$p(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta}, \mathbf{x}) = f(\mathbf{y}, \mathbf{x}, \mathbf{z}) g(\boldsymbol{\theta}) \exp \left[\boldsymbol{\phi}(\boldsymbol{\theta})^\top \mathbf{u}(\mathbf{y}, \mathbf{z}, \mathbf{x}) \right]$$

Conjugate prior

$$p(\boldsymbol{\theta} | \eta, \boldsymbol{\nu}) = h(\eta, \boldsymbol{\nu}) g(\boldsymbol{\theta})^\eta \exp \left[\boldsymbol{\phi}(\boldsymbol{\theta})^\top \boldsymbol{\nu} \right]$$

- many standard distribution belongs to exponential family
- also complete likelihood for many mixture models, classes of Markov models, etc.



Optimum under exponential distribution family

i.i.d. Data $\mathcal{D} = \{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^N$

Exact posterior as in regular EM but using averaged natural parameters

$$q(\mathbf{z}_k) \propto f(\mathbf{y}_k, \mathbf{x}_k, \mathbf{z}_k) \exp \left[\bar{\boldsymbol{\phi}}(\boldsymbol{\theta})^\top \mathbf{u}(\mathbf{y}_k, \mathbf{z}_k, \mathbf{x}_k) \right] = p(\mathbf{z}_k | \mathbf{y}_k, \mathbf{x}_k, \bar{\boldsymbol{\phi}}(\boldsymbol{\theta}))$$

$\bar{\boldsymbol{\phi}}(\boldsymbol{\theta}) = \langle \boldsymbol{\phi}(\boldsymbol{\theta}) \rangle_{q(\boldsymbol{\theta})}$ are natural parameters

$q(\boldsymbol{\theta})$ also exponential family conjugate

$$q(\boldsymbol{\theta}) = h(\tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}}) g(\boldsymbol{\theta})^{\tilde{\boldsymbol{\eta}}} \exp \left[\boldsymbol{\phi}(\boldsymbol{\theta})^\top \tilde{\boldsymbol{\nu}} \right]$$

$$\tilde{\boldsymbol{\eta}} = \boldsymbol{\eta} + N, \tilde{\boldsymbol{\nu}} = \boldsymbol{\nu} + \sum_{n=1}^N \bar{\mathbf{u}}(\mathbf{y}_k, \mathbf{x}_k),$$

$$\bar{\mathbf{u}}(\mathbf{y}_k, \mathbf{x}_k) = \langle \mathbf{u}(\mathbf{y}_k, \mathbf{x}_k) \rangle_{q(\mathbf{z})}$$



Expectation propagation

- Minka: focus on approximating marginals for each sample:
 $t_k(\boldsymbol{\theta}) = p(\boldsymbol{\theta})p(y_k|\mathbf{x}_k, \boldsymbol{\theta})$.
- Use $KL(p|q)$ not $KL(q|p)$ as in VB which typically under-estimates variability.
- Iterate for each sample k
 - Deletion: delete $t_k(\boldsymbol{\theta})$
 - Projection: update $\tilde{t}_k(\boldsymbol{\theta})$
 - Inclusion: update $q(\boldsymbol{\theta})$
- No proof of convergence.



Simple evidence approximation

Simpler than more involved methods like Laplace, variational Bayes and MCMC

Normalized log-posterior

$$C_{\mathcal{D}}(\boldsymbol{\theta}) = \frac{1}{N} (\log p(\mathcal{D}|\boldsymbol{\theta}, m) + \log p(\boldsymbol{\theta}|m))$$

and the maximum a posteriori (MAP) solution

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} C_{\mathcal{D}}(\boldsymbol{\theta})$$



Approximations

Gaussian MAP approximation

$$C_{\mathcal{D}}(\boldsymbol{\theta}) = C_{\mathcal{D}}(\hat{\boldsymbol{\theta}}) - \frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})\mathbf{J}_{\mathcal{D}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^{\top}$$

$$\mathbf{J}_{\mathcal{D}} = -\frac{\partial^2 C_{\mathcal{D}}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^{\top}} \bigg|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} = O(1)$$

Essential assumptions

- $\mathbf{J}_{\mathcal{D}}$ should be of full rank, hence $\mathbf{J}_{\mathcal{D}}^{-1}$ should exist
- $\mathbf{J}_{\mathcal{D}} = O(1)$ is usually fulfilled with N^{-1} normalization. Sinusoidal model is a counter example Stoica



Approximations

Laplace Approximation

$$\begin{aligned} p(\mathcal{D}|m) &= \int \exp(NC_{\mathcal{D}}(\boldsymbol{\theta})) d\boldsymbol{\theta} \\ &\approx \int \exp\left(NC_{\mathcal{D}}(\hat{\boldsymbol{\theta}}) - \frac{N}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})\mathbf{J}_{\mathcal{D}}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^{\top}\right) d\boldsymbol{\theta} \\ &\approx p(\mathcal{D}|\hat{\boldsymbol{\theta}}, m) \cdot p(\hat{\boldsymbol{\theta}}|m) \cdot \left(\frac{2\pi}{N}\right)^{\frac{\dim(\boldsymbol{\theta})}{2}} \cdot |\mathbf{J}|_{\mathcal{D}}^{-\frac{1}{2}} \end{aligned}$$



Bayesian Information Criterion

Since $\mathbf{J}_{\mathcal{D}} = O(1)$, the leading term for large N does not involve the often complicated Hessian, hence, the evidence is approximated as

$$p(\mathcal{D}|m) \approx p(\mathcal{D}|\hat{\boldsymbol{\theta}}, m) \cdot p(\hat{\boldsymbol{\theta}}|m) \cdot \left(\frac{2\pi}{N}\right)^{\frac{\dim(\boldsymbol{\theta})}{2}}$$

$$\log p(\mathcal{D}|m)/N \approx \text{BIC} = C_{\mathcal{D}}(\hat{\boldsymbol{\theta}}) + \dim(\hat{\boldsymbol{\theta}}) \cdot \log(N)/(2N)$$



Views on Bayesian model selection

- BIC gives a consistent model selection as $N \rightarrow \infty$ if the true model is among the candidates
- GEN/AIC consistently overfit for $N \rightarrow \infty$ has a smaller penalty $\dim(\boldsymbol{\theta})/N$ compared to $\dim(\boldsymbol{\theta}) \cdot \log(N)/2N$ in BIC

Connection between BIC and GEN/AIC

- In GEN/AIC $E_{\mathcal{D}} \{G(\mathcal{D})\}$ is approximated by a 2nd order Taylor
- In BIC, $E_{\mathcal{D}} \{\exp(G(\mathcal{D}))\}$ is approximated a 2nd order Taylor; hence log is performed



Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



Discussion of Bayesian Learning

Objectivity

We want to be as objective as possible, however without prior expectation nothing can be learned

- no-free-lunch theorems
- link to philosophical theories
- J. Friedman: “no methods dominates all others over all possible situations”



Discussion of Bayesian Learning

Prior knowledge

likelihood, loss function, model family, parameter priors

- G.E.P. Box (1976) and Stephen Strother: “all models are wrong – but some are useful”
- use flexible models with careful model optimization
- be as data-driven as possible, minimum non-informative prior assumptions
- use Bayes for formal incorporation of all available knowledge
- use careful model evaluation (generalization performance, robustness to changes in assumptions, sensitivity analysis)



Discussion of Bayesian Learning

Robustness

Slight changes in model assumptions should lead to slight changes in conclusions/decisions

- sensitivity analysis
- generalization error - test performance
- extensive cross-validation
- learning curves
- information conveyed by data and by prior - if they clash we want likelihood dominance.
- errorbars and fluctuations in predictive distributions



Discussion of Bayesian Learning

Robustness

- iterated modeling until desired performance/robustness is obtained
- trade-off between performance and robustness for specific limited data set



Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



Bayesian learning in RBF nets - relevance vector machine

RBF network

$$y = \sum_{j=1}^{n_H} \phi_j(\mathbf{x}) w_{ij}^O + \varepsilon = \boldsymbol{\phi}^\top(\mathbf{x}) \boldsymbol{\theta} + \varepsilon$$

- $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ and i.i.d
- $\phi_j(\mathbf{x}) = \exp(-\|\mathbf{x} - \mathbf{x}(j)\|^2 / 2\nu^2),$

Vector notation of training data

$$\mathbf{y} = \boldsymbol{\Phi}^\top \boldsymbol{\theta} + \boldsymbol{\varepsilon}$$



Getting the predictive distribution

Ingredients

Prior	$p(\boldsymbol{\theta} \mathbf{A}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}), \mathbf{A} = \text{diag}(\boldsymbol{\alpha}), \alpha_j$ is (inverse) individual weight decay or hyperparameter
Likelihood	$p(\mathcal{D} \boldsymbol{\theta}, \sigma^2, \nu^2) = \prod_{k=1}^N \mathcal{N}(y(k) - \boldsymbol{\phi}^\top(\mathbf{x}(k))\boldsymbol{\theta}, \sigma^2)$

Posterior weight distribution

$$p(\boldsymbol{\theta}|\mathcal{D}, \sigma^2, \mathbf{A}, \nu^2) = \frac{p(\mathcal{D}|\boldsymbol{\theta}, \sigma^2, \nu^2)p(\boldsymbol{\theta}|\mathbf{A})}{p(\mathcal{D}|\sigma^2, \mathbf{A}, \nu^2)} \sim \mathcal{N}(\hat{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})$$

$$\boldsymbol{\Sigma}_{\boldsymbol{\theta}} = (\sigma^{-2}\boldsymbol{\Phi}\boldsymbol{\Phi}^\top + \mathbf{A})^{-1}$$

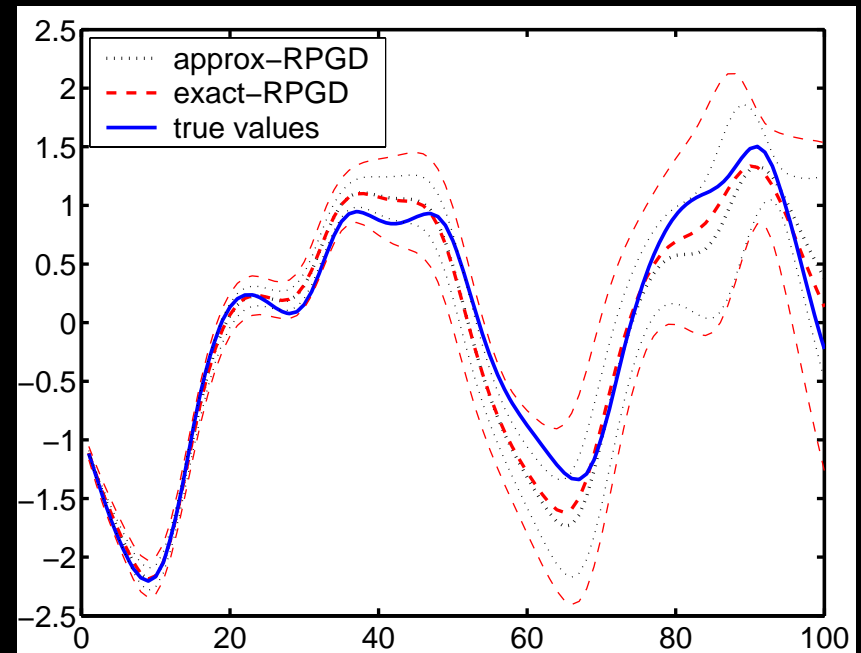
$$\hat{\boldsymbol{\theta}} = \sigma^{-2}\boldsymbol{\Sigma}_{\boldsymbol{\theta}}\boldsymbol{\Phi}\mathbf{y} = (\boldsymbol{\Phi}\boldsymbol{\Phi}^\top + \sigma^2\mathbf{A})^{-1}\boldsymbol{\Phi}\mathbf{y}$$



Getting the predictive distribution

$$p(y|\mathbf{x}, \mathcal{D}) = \int p(y|\mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}, \sigma^2, \mathbf{A}, \nu^2) d\boldsymbol{\theta} \\ \sim \mathcal{N}(\hat{y}, \sigma_y^2)$$

$$\hat{y} = \boldsymbol{\phi}^\top(\mathbf{x}) \hat{\boldsymbol{\theta}} \\ \sigma_y^2 = \sigma^2 + \boldsymbol{\phi}^\top(\mathbf{x}) \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \boldsymbol{\phi}(\mathbf{x})$$





Optimizing hyperparameters

σ^2 , \mathbf{A} , ν^2 are optimized by maximizing the **evidence** using EM and simple search

$$p(\mathcal{D}|\sigma^2, \mathbf{A}, \nu^2) = p(\mathcal{D}|\boldsymbol{\theta}, \sigma^2, \nu^2)p(\boldsymbol{\theta}|\mathbf{A})$$

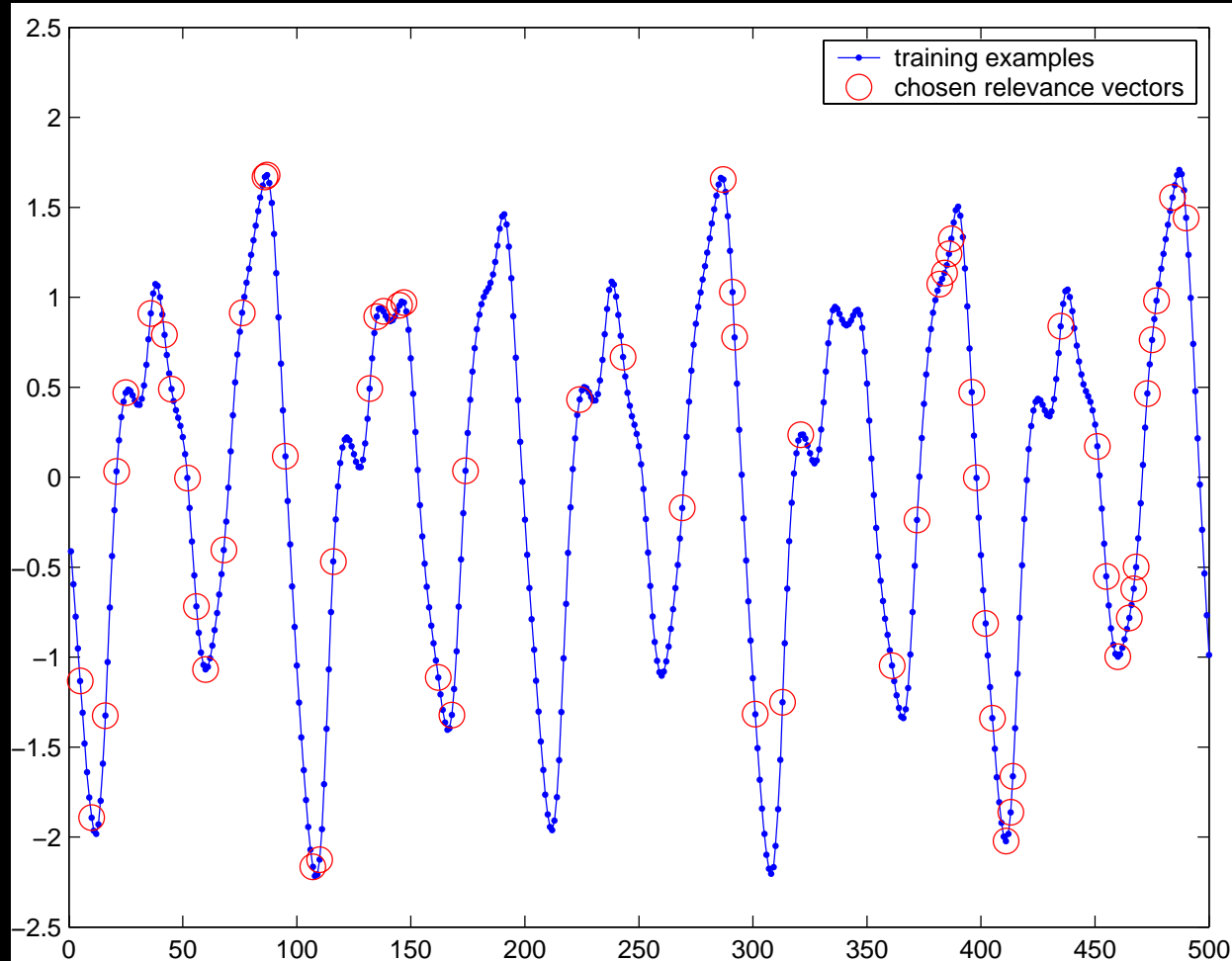
Literature:

Quiñonero-Candela, J., Girard, A., Larsen, J., Rasmussen, C. E.: “Propagation of Uncertainty in Bayesian Kernel Models - Application to Multiple-Step Ahead Forecasting,” ICASSP, vol. 2, pp. 701-704, 2003

Quiñonero-Candela, J., Hansen, L. K.: “Time Series Prediction Based on the Relevance Vector Machine with Adaptive Kernels,” ICASSP, pp. 985-988, 2002



Example





Outline

- Why Bayesian learning?
- Basic ingredients
- Bayes estimators
- More on selection of priors
- Generalization and bias/variance
- Generalization estimation
- Bayesian model selection
- Discussion of Bayesian framework
- Example of Bayesian learning: RVM
- Bayesian signal detection



Bayesian signal detection model

Literature:

Hansen, L. K., Nielsen, F. ., Larsen, J.; “Exploring fMRI Data for Periodic Signal Components,” Artificial Intelligence in Medicine, vol. 25, pp. 25–44, 2002

$$y(n) = \sum_{k=1}^K b_k \cdot x_k(n) + \epsilon(n)$$

$$\mathbf{y} = \hat{\mathbf{y}} + \boldsymbol{\epsilon} = \mathbf{X}\mathbf{b} + \boldsymbol{\epsilon}$$

- observed signal $\mathbf{y} = \{y(n)\}$ is a $N \times 1$ (data \mathcal{D}) vector
- $K = 2\kappa$ periodic basis functions $k \in [1; \kappa]$

$$x_{2k}(n) = \cos(k\omega_0 n), \quad x_{2k-1}(n) = \sin(k\omega_0 n)$$

$\mathbf{X} = \{x_k(n)\}$ is a $N \times K$ matrix



- b_k linear coefficients $\mathbf{b} = \{b_k\}$ is a $K \times 1$ vector
- noise: $\epsilon \sim \mathcal{N}(0, \sigma^2)$



Objective

Estimate unknown fundamental frequency ω_0 and the number of components K

Bayesian learning

- integrate out undesired model parameters $\theta = (\mathbf{b}, \sigma^2)$
- select model $\omega_0(m)$, $K(m)$, $m = [1; M]$
- $m = 0$ corresponds to only noise, i.e., $\mathbf{X} \equiv 0$.

$$p(\omega_0, K | \mathbf{y}) = \frac{p(\mathbf{y} | \omega_0, K) p(\omega_0, K)}{p(\mathbf{y})}$$



Marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\omega_0, K) &= \int p(\mathbf{y}|\mathbf{b}, \sigma^2, \omega_0, K) \cdot p(\mathbf{b}, \sigma^2) d\mathbf{b} d\sigma^2 \\ &= \int (2\pi\sigma^2)^{N/2} \exp\left(-\|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2/2\sigma^2\right) \cdot p(\mathbf{b}, \sigma^2) d\mathbf{b} d\sigma^2 \end{aligned}$$

Conjugate prior: normal-inverse-gamma

$$\begin{aligned} p(\mathbf{b}, \sigma^2|a, d, K, \mathbf{m}, \mathbf{V}) &= \frac{(a/2)^{d/2} \cdot (\sigma^2)^{-(d+K+2)/2}}{(2\pi)^{K/2} \cdot \det \mathbf{V}^{1/2} \cdot \Gamma(d/2)} \\ &\quad \exp\left(-(\mathbf{b} - \mathbf{m})^\top (2\sigma^2 \mathbf{V})^{-1} (\mathbf{b} - \mathbf{m}) - a/2\sigma^2\right) \end{aligned}$$



Priors

$$p(\mathbf{b}|a, d, K, \mathbf{m}, \mathbf{V}) = \int p(\mathbf{b}, \sigma^2|a, d, K, \mathbf{m}, \mathbf{V}) d\sigma^2 \\ \sim \mathcal{T}(\mathbf{m}, a\mathbf{V}/(d-2))$$

$$p(\sigma^2|a, d) = \int p(\mathbf{b}, \sigma^2|a, d, K, \mathbf{m}, \mathbf{V}) d\mathbf{b} \sim \mathcal{IG}(a/(d-2))$$

- mean of noise variance $a/(d-2) = \hat{\sigma}_y^2 = \mathbf{y}^\top \mathbf{y}/N$
- $d = 3$ is smallest value for which prior is finite, hence “weak”
- $\mathbf{m} = \mathbf{0}$ for no prior assumption of mean amplitude of periodic components



Priors

- $\mathbf{V} = v\mathbf{I}$ for simplicity

$$\begin{aligned} E_{\text{prior}}[\mathbf{y}^{\top} \mathbf{y}] / N &= \text{Tr}[\mathbf{X} \mathbf{X}^{\top} E_{\text{prior}}[\mathbf{b} \mathbf{b}^{\top}]] / N \\ &= v \cdot a / (d - 2) \cdot \text{Tr}[\mathbf{X} \mathbf{X}^{\top}] / N \end{aligned}$$

If $E_{\text{prior}}[\mathbf{y}^{\top} \mathbf{y}] / N = \hat{\sigma}_y^2 = a / (d - 2)$ then $v = N \cdot \text{Tr}[\mathbf{X} \mathbf{X}^{\top}]^{-1}$



Marginal likelihood

posterior is NIG thus the marginal likelihood is its normalization integral

$$p(\mathbf{y}|\omega_0, K) = \left(\frac{\det \mathbf{V}_P \cdot a^d}{\det \mathbf{V} \cdot a_P^{d_P} \cdot \pi^N} \right)^{1/2} \frac{\Gamma(d_P/2)}{\Gamma(d/2)}$$

- $\mathbf{V}_P^{-1} = \mathbf{V}^{-1} + \mathbf{X}^\top \mathbf{X}$
- $\mathbf{m}_P = \mathbf{V}_P(\mathbf{V}^{-1}\mathbf{m} + \mathbf{X}^\top \mathbf{y})$
- $a_P = a + \mathbf{m}^\top \mathbf{V}^{-1}\mathbf{m} + \mathbf{y}^\top \mathbf{y} - \mathbf{m}_P^\top \mathbf{V}_P^{-1}\mathbf{m}_P$
- $d_P = d + N$



Model selection

uniform models prior $p(\omega_0(m), K(m)) = 1/(M + 1)$

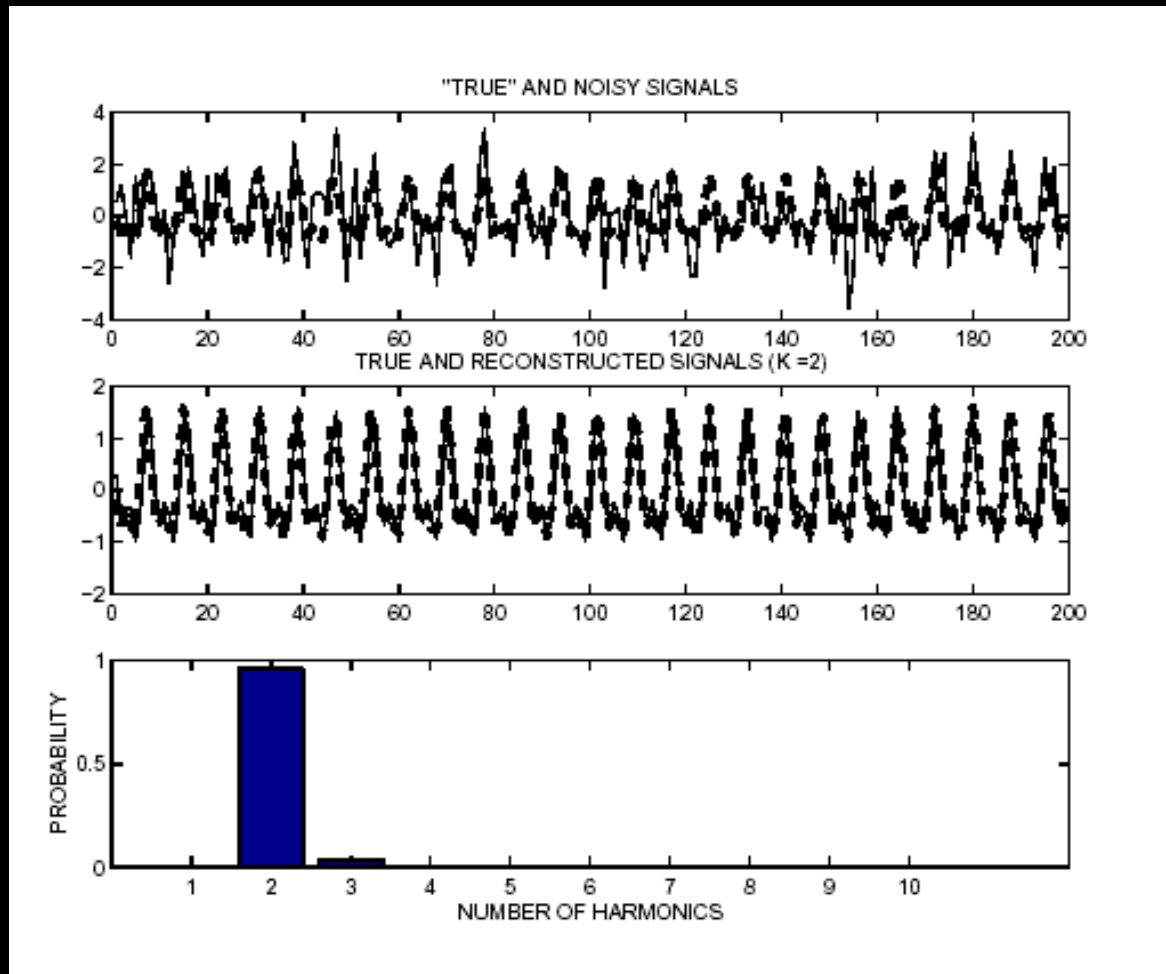
$$p(m|\mathbf{y}) = \frac{p(\mathbf{y}|m)}{\sum_{m=0}^M p(\mathbf{y}|m)}$$

with

- $p(\mathbf{y}|m) = p(\mathbf{y}|\omega(m), K(m)), m > 0$
- $p(\mathbf{y}|0) = p(\mathbf{y}|\omega, K)$ with $\mathbf{X} = 0$



Example





Suggestions for a roadmap

- More research on the evaluation of the learning process. Account for all variation – also data set variability
- Development of better and easy to communicate approximation schemes
- More research on online learning in a non-stationary switching dynamics settings
- Bayes does not tell you anything about the domain exterior to the model - hence, more focus on integrating the data representation, feature selection, and preprocessing steps
- Systems interact with other systems and humans – model the man in the loop, model irrationality



Wrap up

- Bayesian learning combines all available knowledge in principled way
- Ingredients: variable, data, model, prior, loss
- Bayes is optimal in admissibility/generalization sense
- Bayes framework is complete as it offers model selection and confidence
- Robustness needs to be tested, model mis-specifications can cause arbitrary errors