

Introduction to Bayesian methods and their use in fusion data analysis

Geert Verdoolaege

Department of Applied Physics, Ghent University, Ghent, Belgium

IAEA Technical Meeting on Fusion Data Processing,
Validation and Analysis

Nov. 29, 2021

Overview

1. Classical probability and statistics
2. Principles of Bayesian probability theory
3. Applications
 - Classification
 - Regression analysis
 - Some other applications
4. Conclusions and references

1. Classical probability and statistics
2. Principles of Bayesian probability theory
3. Applications
 - Classification
 - Regression analysis
 - Some other applications
4. Conclusions and references

Frequentist probability

- Probability = frequency
- Straightforward:
 - Number of 1s in 60 dice throws ≈ 10 : $p = 1/6$
 - Probability of plasma disruption $p \approx N_{\text{disr}}/N_{\text{tot}}$
- Less straightforward:
 - Probability of fusion electricity by 2050?
 - Probability of mass of Saturn $90 m_A \leq m_S < 100 m_A$?

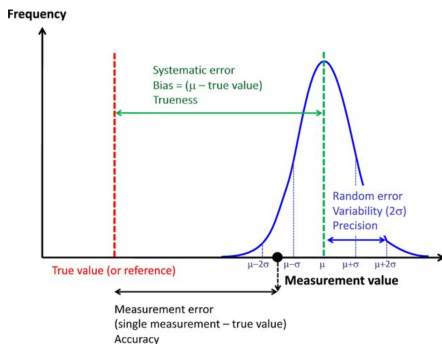
Flavors of uncertainty

Aleatoric/statistical/random uncertainty

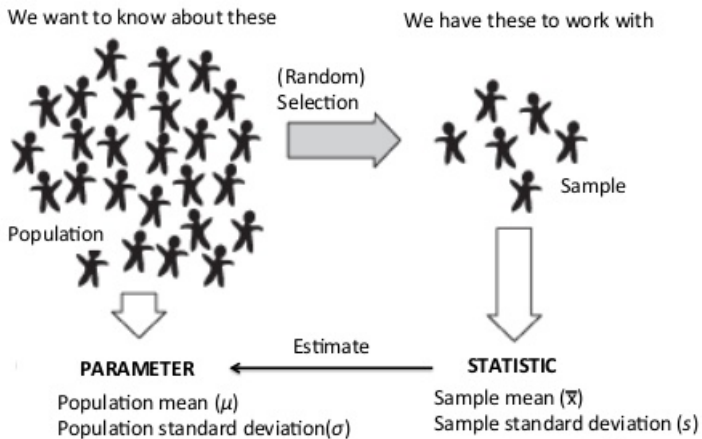
- Leads to different outcomes in multiple experimental trials
- Can be reduced by repeating measurement

Epistemic/systematic uncertainty

- 'Fixed' but unknown ('bias')
- Cannot be reduced through repeated measurement



Populations vs. sample



- E.g. weight w of Belgian men: unknown but *fixed* for every individual
- Average weight μ_w in population?
- **Random variable** W
- Sample: W_1, W_2, \dots, W_n
- Average weight: **statistic** (estimator) \bar{W}
- Central limit theorem:

$$W \sim p(W|\mu_w, \sigma_w) \Rightarrow \bar{W} \sim \mathcal{N}(\mu_w, \sigma_w / \sqrt{n})$$

Maximum likelihood parameter estimation

- *Maximum likelihood* (ML) principle:

$$\begin{aligned}\hat{\mu}_w &= \arg \max_{\mu_w \in \mathbb{R}^+} p(W_1, \dots, W_n | \mu_w, \sigma_w) \\ &\approx \arg \max_{\mu_w \in \mathbb{R}^+} \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_w} \exp \left[-\frac{(W_i - \mu_w)^2}{2\sigma_w^2} \right] \\ &= \arg \max_{\mu_w \in \mathbb{R}^+} \frac{1}{\sqrt{2\pi}\sigma_w} \exp \left[-\sum_{i=1}^n \frac{(W_i - \mu_w)^2}{2\sigma_w^2} \right]\end{aligned}$$

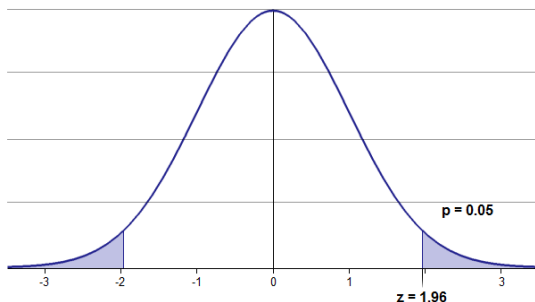
- ML estimator (known σ_w):

$$\hat{\mu}_w = \bar{W} = \frac{1}{n} \sum_{i=1}^n W_i$$

Frequentist hypothesis testing

- Weight of Dutch men compared to Belgian men (populations)
- Observed sample averages $\bar{W}_{NL}, \bar{W}_{BE}$
- **Null hypothesis** $H_0: \mu_{w,NL} = \mu_{w,BE}$
- Test statistic:

$$\frac{\bar{W}_{NL} - \bar{W}_{BE}}{\sigma_{\bar{W}_{NL} - \bar{W}_{BE}}} \sim \mathcal{N}(0, 1) \quad (\text{under } H_0)$$



1. Classical probability and statistics
2. Principles of Bayesian probability theory
3. Applications
 - Classification
 - Regression analysis
 - Some other applications
4. Conclusions and references

Probability theory: quantifying uncertainty

- Every piece of information has uncertainty
- Uncertainty = lack of information
- Observation may reduce uncertainty
- Probability (distribution) *quantifies* uncertainty



Example: physical sciences

- Measurement of physical quantity
- Origin of stochasticity:
 - Apparatus
 - Microscopic fluctuations
- Systematic uncertainty is assigned a probability distribution
- E.g. coin tossing, voltage measurement, probability of hypothesis vs. another, ...
- Bayesian: no 'random' variables



What is probability?

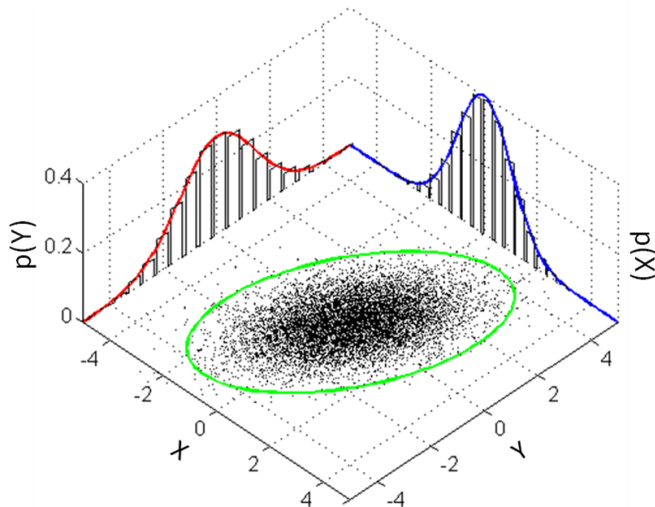
- Objective Bayesian view
- Probability = real number $\in [0, 1]$
- Always conditioned on known information

- Notation:

$$p(A|B) \quad \text{or} \quad p(A|I)$$

- Extension of logic: measure of degree to which *B implies A*
- Degree of plausibility, but subject to consistency
- Same information \Rightarrow same probabilities
- **Probability distribution:** outcome \rightarrow probability

Joint, marginal and conditional distributions



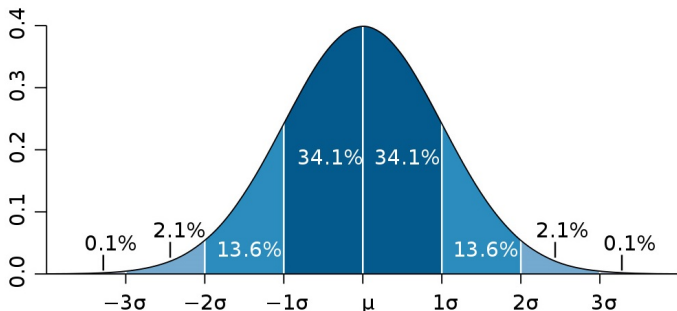
$$p(x, y), p(x), p(y), p(x|y), p(y|x)$$

Example: normal distribution

- Normal/Gaussian *probability density function* (PDF):

$$p(x|\mu, \sigma, I) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right]$$

- Probability $x_1 \leq x < x_1 + dx$
- Inverse problem: μ, σ given x ?



Updating information states

Bayes' theorem

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

\mathbf{x} = data vector
 $\boldsymbol{\theta}$ = vector of model parameters
 I = implicit knowledge

- **Likelihood**: misfit between model and data
- **Prior** distribution: 'expert' or diffuse knowledge
- **Evidence**:

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

- **Posterior** distribution

Updating information states

Bayes' theorem

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

\mathbf{x} = data vector
 $\boldsymbol{\theta}$ = vector of model parameters
 I = implicit knowledge

- **Likelihood**: misfit between model and data
- **Prior** distribution: 'expert' or diffuse knowledge
- **Evidence**:

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

- **Posterior** distribution

Updating information states

Bayes' theorem

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

\mathbf{x} = data vector
 $\boldsymbol{\theta}$ = vector of model parameters
 I = implicit knowledge

- **Likelihood**: misfit between model and data
- **Prior** distribution: 'expert' or diffuse knowledge
- **Evidence**:

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

- **Posterior** distribution

Updating information states

Bayes' theorem

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

\mathbf{x} = data vector
 $\boldsymbol{\theta}$ = vector of model parameters
 I = implicit knowledge

- **Likelihood**: misfit between model and data
- **Prior** distribution: 'expert' or diffuse knowledge
- **Evidence**:

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

- **Posterior** distribution

Mean of a normal distribution: uniform prior

- n measurements $x_i \rightarrow \mathbf{x}$
- Independent and identically distributed x_i :

$$p(\mathbf{x}|\mu, \sigma, I) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x_i - \mu)^2}{2\sigma^2}\right]$$

- Bayes' rule:

$$p(\mu, \sigma|\mathbf{x}, I) \propto p(\mathbf{x}|\mu, \sigma, I)p(\mu, \sigma|I)$$

- Suppose $\sigma \equiv \sigma_e \rightarrow$ delta function
- Assume $\mu \in [\mu_{\min}, \mu_{\max}] \rightarrow$ uniform prior:

$$p(\mu|I) = \begin{cases} \frac{1}{\mu_{\max} - \mu_{\min}}, & \text{if } \mu \in [\mu_{\min}, \mu_{\max}] \\ 0, & \text{otherwise} \end{cases}$$

- Let $\mu_{\min} \rightarrow -\infty, \mu_{\max} \rightarrow +\infty \Rightarrow$ *improper prior*
- Ensure proper posterior

Posterior for μ

- Posterior:

$$p(\mu|\mathbf{x}, I) \propto \exp \left[-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma_e^2} \right]$$

- Define

$$\bar{x} \equiv \frac{1}{n} \sum_{i=1}^n x_i, \quad \overline{(\Delta x)^2} \equiv \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

- Adding and subtracting $2n\bar{x}^2$ ('completing the square'),

$$p(\mu|\mathbf{x}, I) \propto \exp \left\{ -\frac{1}{2\sigma_e^2/n} \left[(\mu - \bar{x})^2 + \overline{(\Delta x)^2} \right] \right\}$$

- Retaining dependence on μ ,

$$p(\mu|\mathbf{x}, I) \propto \exp \left[-\frac{(\mu - \bar{x})^2}{2\sigma_e^2/n} \right]$$

- $\mu \sim \mathcal{N}(\bar{x}, \sigma_e^2/n)$

Mean of a normal distribution: normal prior

- Normal prior: $\mu \sim \mathcal{N}(\mu_0, \tau^2)$
- Posterior:

$$p(\mu | \mathbf{x}, I) \propto \exp \left[-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma_e^2} \right] \times \exp \left[-\frac{(\mu - \mu_0)^2}{2\tau^2} \right]$$

- Expanding and completing the square,

$$\mu \sim \mathcal{N}(\mu_n, \sigma_n^2),$$

where

$$\mu_n \equiv \sigma_n^2 \left(\frac{n}{\sigma_e^2} \bar{x} + \frac{1}{\tau^2} \mu_0 \right) \quad \text{and} \quad \sigma_n^2 \equiv \left(\frac{n}{\sigma_e^2} + \frac{1}{\tau^2} \right)^{-1}$$

- μ_n is weighted average of μ_0 and \bar{x}

Unknown mean and standard deviation

- Repeated measurements \rightarrow information on σ
- Scale variable $\sigma \rightarrow$ *Jeffreys' scale prior*:

$$p(\sigma|I) \propto \frac{1}{\sigma}, \quad \sigma \in]0, +\infty[$$

- Posterior:

$$p(\mu, \sigma | \mathbf{x}, I) \propto \frac{1}{\sigma^n} \exp \left[-\frac{(\mu - \bar{x})^2 + \overline{(\Delta x)^2}}{2\sigma^2/n} \right] \times \frac{1}{\sigma}$$

Marginal posterior for μ (1)

- **Marginalization** = integrating out a (nuisance) parameter:

$$\begin{aligned} p(\mu|x, I) &= \int_0^{+\infty} p(\mu, \sigma|x, I) d\sigma \\ &\propto \int_0^{+\infty} \frac{1}{2} \left[\frac{(\mu - \bar{x})^2 + \overline{(\Delta x)^2}}{2/n} \right]^{-\frac{n}{2}} s^{\frac{n}{2}-1} e^{-s} ds \\ &= \frac{1}{2} \Gamma\left(\frac{n}{2}\right) \left[\frac{(\mu - \bar{x})^2 + \overline{(\Delta x)^2}}{2/n} \right]^{-\frac{n}{2}}, \end{aligned}$$

where

$$s \equiv \frac{(\mu - \bar{x})^2 + \overline{(\Delta x)^2}}{2\sigma^2/n}$$

Marginal posterior for μ (2)

- After normalization:

$$p(\mu|\mathbf{x}, I) = \frac{\Gamma\left(\frac{n}{2}\right)}{\sqrt{\pi(\Delta x)^2} \Gamma\left(\frac{n-1}{2}\right)} \left[1 + \frac{(\mu - \bar{x})^2}{(\Delta x)^2} \right]^{-\frac{n}{2}}$$

- Changing variables,

$$t \equiv \frac{(\mu - \bar{x})^2}{\sqrt{(\Delta x)^2/(n-1)}}, \quad \text{with} \quad p(t|\mathbf{x}, I) dt \equiv p(\mu|\mathbf{x}, I) d\mu,$$

$$p(t|\mathbf{x}, I) = \frac{\Gamma\left(\frac{n}{2}\right)}{\sqrt{(n-1)\pi} \Gamma\left(\frac{n-1}{2}\right)} \left[1 + \frac{t^2}{n-1} \right]^{-\frac{n}{2}}$$

- Student's t -distribution with parameter $\nu = n - 1$
- If $n \gg 1$,

$$p(\mu|\mathbf{x}, I) \longrightarrow \frac{1}{\sqrt{2\pi(\Delta x)^2/n}} \exp\left[-\frac{(\mu - \bar{x})^2}{2(\Delta x)^2/n}\right]$$

Marginal posterior for σ

- Marginalization of μ :

$$\begin{aligned} p(\sigma|\mathbf{x}, I) &= \int_{-\infty}^{+\infty} p(\mu, \sigma|\mathbf{x}, I) \, d\mu \\ &\propto \frac{1}{\sigma^n} \exp \left[-\frac{\overline{(\Delta x)^2}}{2\sigma^2/n} \right] \end{aligned}$$

- Setting $X \equiv n\overline{(\Delta x)^2}/\sigma^2$,

$$p(X|\mathbf{x}, I) = \frac{1}{2^{\frac{k}{2}} \Gamma\left(\frac{k}{2}\right)} X^{\frac{k}{2}-1} e^{-\frac{X}{2}}, \quad k \equiv n - 1$$

- χ^2 distribution with parameter k

The Laplace approximation (1)

- Laplace (saddle point) approximation of distributions around the mode (= maximum)
- E.g. marginal for μ :

$$p(\mu|x, I) = \frac{\Gamma\left(\frac{n}{2}\right)}{\sqrt{\pi(\Delta x)^2} \Gamma\left(\frac{n-1}{2}\right)} \left[1 + \frac{(\mu - \bar{x})^2}{(\Delta x)^2}\right]^{-\frac{n}{2}}$$

- Taylor expansion around mode:

$$\begin{aligned}\ln[p(\mu|x, I)] &\approx \ln[p(\bar{x}|x, I)] + \frac{1}{2} \frac{d^2(\ln p)}{d\mu^2} \Big|_{\mu=\bar{x}} (\mu - \bar{x})^2 \\ &= \ln \left[\Gamma\left(\frac{n}{2}\right) \right] - \ln \left[\Gamma\left(\frac{n-1}{2}\right) \right] \\ &\quad - \frac{1}{2} \ln \left[\pi(\Delta x)^2 \right] - \frac{n}{2(\Delta x)^2} (\mu - \bar{x})^2\end{aligned}$$

The Laplace approximation (2)

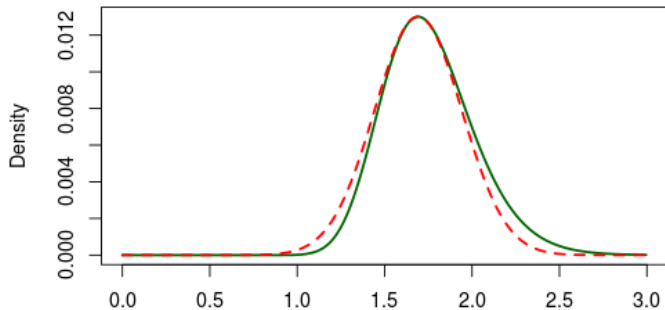
- On the original scale:

$$p(\mu|x, I) \approx \frac{\Gamma\left(\frac{n}{2}\right)}{\Gamma\left(\frac{n-1}{2}\right)} \frac{1}{\sqrt{\pi(\Delta x)^2}} \exp\left[-\frac{(\mu - \bar{x})^2}{2(\Delta x)^2/n}\right]$$

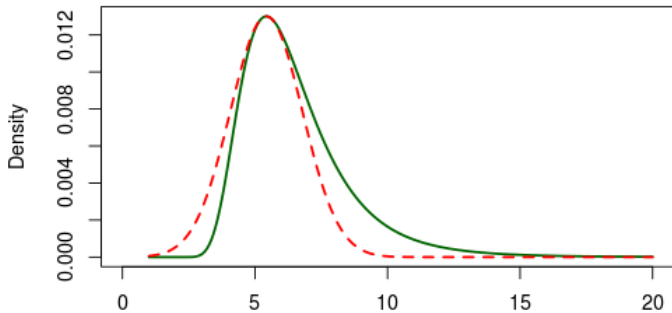
- Standard deviation $\sigma_L \rightarrow$ curvature of $\ln p$:

$$\sigma_L = \left[- \frac{d^2(\ln p)}{d\mu^2} \Big|_{\mu=\bar{x}} \right]^{-1/2}$$

Laplace approximation: example 1



Laplace approximation: example 2



Multivariate Laplace approximation

- For $\boldsymbol{\theta} = [\theta_1, \dots, \theta_p]^t$,

$$p(\boldsymbol{\theta}|\boldsymbol{\theta}_0, I) \propto \exp \left[\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^t [\nabla \nabla (\ln p)]_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} (\boldsymbol{\theta} - \boldsymbol{\theta}_0) \right]$$

- $\nabla \nabla (\ln p)$: Hessian matrix, where

$$\Sigma_L = - \{ [\nabla \nabla (\ln p)]_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \}^{-1}$$

Model comparison (hypothesis testing)

- Let $\{H_i\}$ be complete set of hypotheses
- Data D to support or reject hypotheses
- Bayes' rule:

$$p(H_i|D,I) = \frac{p(D|H_i,I)p(H_i|I)}{p(D|I)}, \quad p(D|I) = \sum_i p(D|H_i,I)p(H_i|I)$$

- Assume single hypothesis H and complement \bar{H}
- *Odds ratio* o :

$$o \equiv \frac{p(H|D,I)}{p(\bar{H}|D,I)} = \underbrace{\frac{p(D|H,I)}{p(D|\bar{H},I)}}_{\text{Bayes factor}} \underbrace{\frac{p(H|I)}{p(\bar{H}|I)}}_{\text{Prior odds}}$$

- $p(D|H,I) =$ *model evidence*

Testing a Gaussian mean (1)

- E.g. n measurements x_i of a quantity x
- Assume normal distribution with known variance σ^2
- Question: are the data compatible with mean $\mu = \mu_0$?
 - Yes: H
 - No: \bar{H}

- Under H :

$$p(\bar{x}|H, I) = C \exp \left[-\frac{1}{2\sigma^2/n} (\bar{x} - \mu_0)^2 \right]$$

- Under \bar{H} :

$$p(\bar{x}|\bar{H}, I) = \int p(\bar{x}|\mu, \sigma, I) p(\mu|\bar{H}, I) d\mu \quad (1)$$

Testing a Gaussian mean (2)

- Assume bounds μ_{\min} and μ_{\max} :

$$p(\mu|\bar{H}, I) = \frac{1}{|\mu_{\max} - \mu_{\min}|} \mathbf{1}(\mu_{\min} \leq \mu \leq \mu_{\max})$$

- Then (1) becomes

$$p(\bar{x}|\bar{H}, I) = \frac{C}{|\mu_{\max} - \mu_{\min}|} \int_{\mu_{\min}}^{\mu_{\max}} \exp\left[-\frac{1}{2\sigma^2/n}(\bar{x} - \mu)^2\right] d\mu$$

- Assume wide prior interval:

$$|\bar{x} - \mu_{\min}|, |\bar{x} - \mu_{\max}| \gg \text{SE},$$

$$\text{SE} = \text{standard error} \equiv \frac{\sigma}{\sqrt{n}}$$

Testing a Gaussian mean (3)

- Then

$$\begin{aligned} p(\bar{x}|\bar{H}, I) &\approx \frac{C}{|\mu_{\max} - \mu_{\min}|} \int_{-\infty}^{+\infty} \exp\left[-\frac{1}{2\sigma^2/n}(\bar{x} - \mu)^2\right] d\mu \\ &= \frac{C \text{ SE } \sqrt{2\pi}}{|\mu_{\max} - \mu_{\min}|} \end{aligned}$$

- Bayes factor BF:

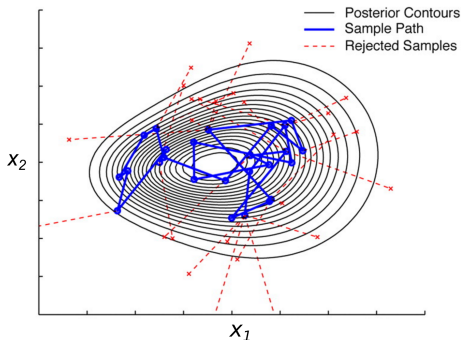
$$\text{BF} = \frac{p(\bar{x}|H, I)}{p(\bar{x}|\bar{H}, I)} = \frac{|\mu_{\max} - \mu_{\min}|}{\text{SE}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2},$$
$$z \equiv \frac{|\bar{x} - \mu_0|}{\text{SE}}$$

- Cf. frequentist hypothesis test

Sampling from distributions



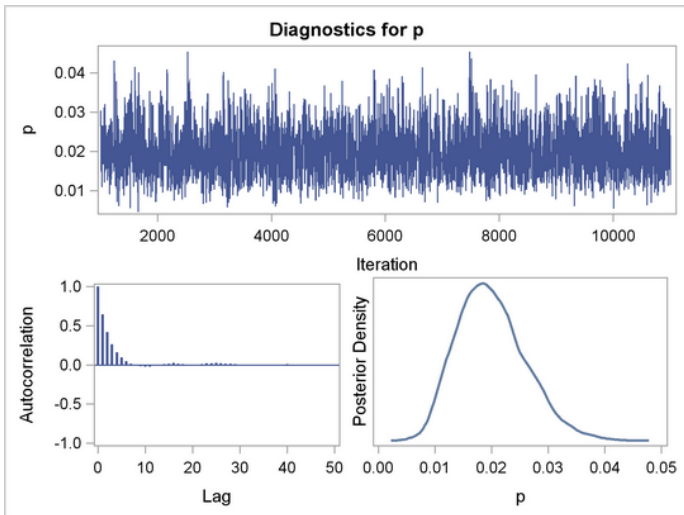
- Need samples from *target distribution* $p(\theta|I)$
- *Markov chain Monte Carlo* sampling
- Sample from *proposal distribution*



- Calculate Monte Carlo averages, e.g.

$$\bar{\theta}_j = \frac{1}{n} \sum_{i=1}^n \theta_j^{(t_c+i)}, \quad \overline{(\Delta\theta_j)^2} = \frac{1}{n} \sum_{i=1}^n \left(\theta_j^{(t_c+i)} - \bar{\theta}_j \right)^2$$

MCMC sampling

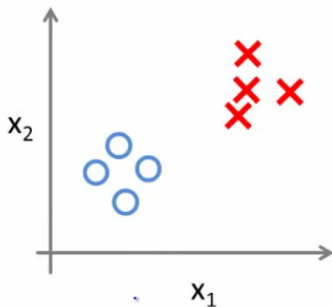


1. Classical probability and statistics
2. Principles of Bayesian probability theory
3. Applications
 - Classification
 - Regression analysis
 - Some other applications
4. Conclusions and references

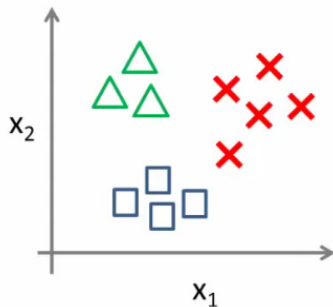
1. Classical probability and statistics
2. Principles of Bayesian probability theory
3. Applications
 - Classification
 - Regression analysis
 - Some other applications
4. Conclusions and references

Classification or clustering

Binary classification:



Multi-class classification:



Simple Bayesian classification

- M clusters of in total n data points \mathbf{x}_i in P -dimensional space
- Known class labels ω_j ($j = 1, \dots, M$) of \mathbf{x}_i
- Bayes' rule for new point \mathbf{x} :

$$p(\omega_j|\mathbf{x}, I) = \frac{p(\mathbf{x}|\omega_j, I)p(\omega_j|I)}{p(\mathbf{x}|I)}$$

- **Maximum a posteriori** (MAP) classification rule for \mathbf{x} :

$$\text{Assign } \mathbf{x} \text{ to } \omega_i = \arg \max_{\omega_j} p(\omega_j|\mathbf{x}, I) = \arg \max_{\omega_j} p(\mathbf{x}|\omega_j, I)p(\omega_j|I)$$

Examples of priors and likelihoods

- Examples of prior probabilities (indifference):

- $p(\omega_i|I) = p(\omega_j|I), \forall i, j$
- Count class membership:

$$p(\omega_i|I) \equiv \frac{n_i}{n}, \quad i = 1, \dots, M$$

- Examples of likelihoods:

- *Naive Bayesian classifier:*

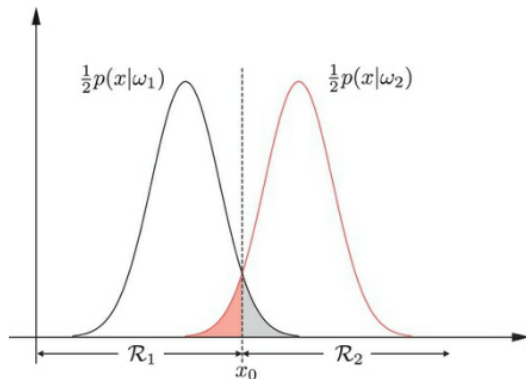
$$p(\mathbf{x}|\omega_i) = \prod_{k=1}^P p(x_k|\omega_i), \quad i = 1, \dots, M$$

- Multivariate Gaussian:

$$p(\mathbf{x}|\omega_j, I) = \frac{1}{(2\pi)^{p/2} |\Sigma_j|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^t \Sigma_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) \right]$$

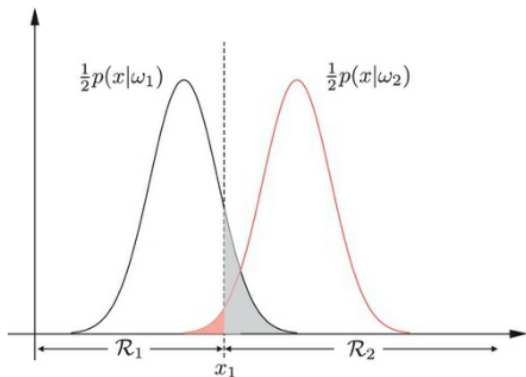
Optimality of Bayesian classifier

Bayesian classifier minimizes probability of misclassification



Optimality of Bayesian classifier

Bayesian classifier minimizes probability of misclassification



MAP classification: decision surfaces

- MAP: maximize w.r.t. ω_j :

$$\ln p(\omega_j|\mathbf{x}, I) = \ln p(\mathbf{x}|\omega_j, I) + \ln p(\omega_j|I)$$

- Define ($M = 2$):

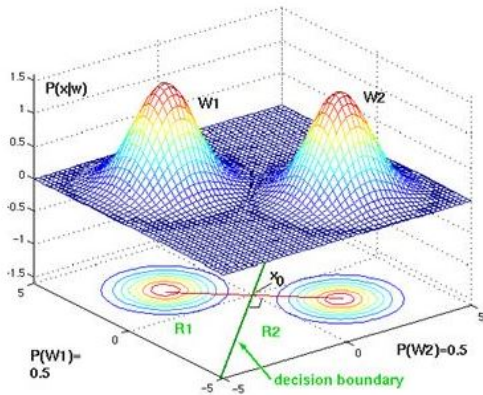
$$g(\mathbf{x}) \equiv \ln p(\omega_1|\mathbf{x}, I) - \ln p(\omega_2|\mathbf{x}, I)$$

- For normal likelihood:

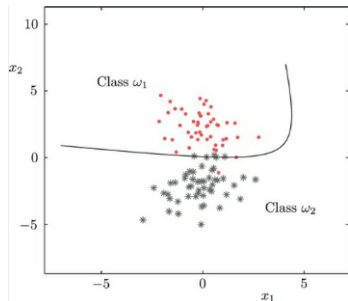
$$g(\mathbf{x}) = \underbrace{\frac{1}{2} \left(\mathbf{x}^t \Sigma_2^{-1} \mathbf{x} - \mathbf{x}^t \Sigma_1^{-1} \mathbf{x} \right)}_{\text{Quadratic}} + \underbrace{\mu_1^t \Sigma_1^{-1} \mathbf{x} - \mu_2^t \Sigma_2^{-1} \mathbf{x}}_{\text{Linear}} \\ - \underbrace{\frac{1}{2} \mu_1^t \Sigma_1^{-1} \mu_1 + \frac{1}{2} \mu_2^t \Sigma_2^{-1} \mu_2 + \frac{1}{2} \ln \frac{|\Sigma_2|}{|\Sigma_1|}}_{\text{Constant}} + \ln \frac{p(\omega_1|I)}{p(\omega_2|I)}$$

- $g(\mathbf{x})$ separates classes: *decision hypersurface*

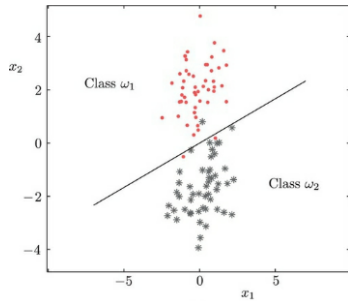
Discriminant analysis



QDA and LDA

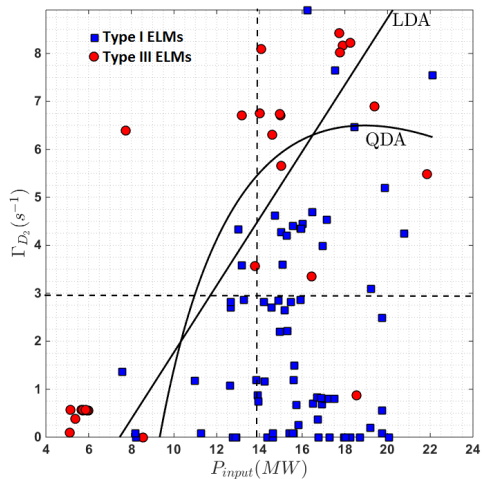


Quadratic discriminant analysis
(QDA)



Linear discriminant analysis
(LDA): $\Sigma_1 = \Sigma_2$

ELM type classification

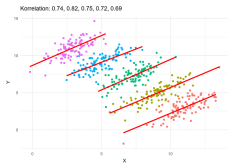
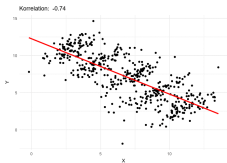
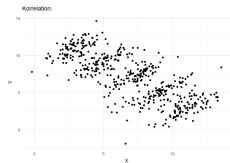
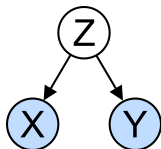


$$P_{input} - 1.41\Gamma_{D_2} = 7.47$$

1. Classical probability and statistics
2. Principles of Bayesian probability theory
3. Applications
 - Classification
 - Regression analysis
 - Some other applications
4. Conclusions and references

Uncertainties in regression analysis

- Measurement uncertainty
- Percentage errors from database
- Model uncertainty:
 - Power law
 - Missing variables
 - Confounding variables
- Predictor correlations (e.g. $I_p \propto B_t$)
- Heterogeneity: multi-machine database
- Simpson's paradox:



Multilinear regression and simple least squares

- Regression model (Gauss-Markov):

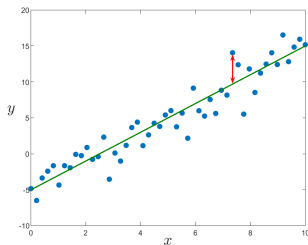
$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_p x_p + \epsilon \quad \text{Often loglinear!}$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2), \sigma \text{ known}$$

- Take n measurements:

$$\mathbf{y} \equiv [y_1, \dots, y_n]^t,$$

$$\mathbf{X} \equiv \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}$$



- Ordinary least squares (OLS):

$$\alpha_{\text{OLS}} = \arg \min_{\alpha} \left[(\mathbf{y} - \mathbf{X}\alpha)^t (\mathbf{y} - \mathbf{X}\alpha) \right] = \underbrace{(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t}_{\text{Moore-Penrose pseudoinverse}} \mathbf{y}$$

Maximum likelihood solution

- Likelihood:

$$p(y|x, \alpha, \sigma, I) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2\sigma^2} \left(y - \alpha_0 - \sum_{j=1}^p \alpha_j x_j \right)^2 \right],$$
$$\alpha \equiv [\alpha_0, \alpha_p^t]^t, \quad \alpha_p \equiv [\alpha_1, \dots, \alpha_p]^t$$

- Conditional independence:

$$p(\mathbf{y}|X, \alpha, \sigma, I) = (2\pi)^{-n/2} \sigma^{-n} \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - X\alpha)^t (\mathbf{y} - X\alpha) \right]$$

- ML solution:

$$0 = \nabla_{\alpha} (\mathbf{y} - X\alpha)^t (\mathbf{y} - X\alpha) = -2X^t \mathbf{y} + 2X^t X \alpha$$
$$\Rightarrow \alpha_{\text{ML}} = (X^t X)^{-1} X^t \mathbf{y} = \alpha_{\text{OLS}}$$

MAP solution and posterior

- Uniform priors on α_j (not the most uninformative!):

$$p(\boldsymbol{\alpha}|\mathbf{y}, X, \sigma, I) \propto \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - X\boldsymbol{\alpha})^t (\mathbf{y} - X\boldsymbol{\alpha}) \right]$$

- Due to linearity and Gaussianity: $\boldsymbol{\alpha}_{\text{MAP}} = \boldsymbol{\alpha}_{\text{ML}} = \boldsymbol{\alpha}_{\text{LS}}$
- Taylor expansion (exact!):

$$\begin{aligned} (\mathbf{y} - X\boldsymbol{\alpha})^t (\mathbf{y} - X\boldsymbol{\alpha}) &= (\mathbf{y} - X\boldsymbol{\alpha}_{\text{MAP}})^t (\mathbf{y} - X\boldsymbol{\alpha}_{\text{MAP}}) \\ &\quad + \frac{1}{2} (\boldsymbol{\alpha} - \boldsymbol{\alpha}_{\text{MAP}})^t 2X^t X (\boldsymbol{\alpha} - \boldsymbol{\alpha}_{\text{MAP}}) \end{aligned}$$

- Posterior distribution:

$$\begin{aligned} p(\boldsymbol{\alpha}|\mathbf{y}, X, \sigma, I) &= (2\pi)^{-n/2} |\Sigma|^{-1/2} \\ &\quad \times \exp \left[-\frac{1}{2} (\boldsymbol{\alpha} - \boldsymbol{\alpha}_{\text{MAP}})^t \Sigma^{-1} (\boldsymbol{\alpha} - \boldsymbol{\alpha}_{\text{MAP}}) \right], \\ \Sigma &\equiv \sigma^2 (X^t X)^{-1} \end{aligned}$$

Posterior predictive distribution (1)

- New predictions by the model?

- *Posterior predictive distribution:*

$$\begin{aligned} p(y_{\text{new}} | \mathbf{x}_{\text{new}}, \mathbf{y}, X, \sigma, I) &= \int_{\mathbb{R}^{p+1}} p(y_{\text{new}}, \boldsymbol{\beta} | \mathbf{x}_{\text{new}}, \mathbf{y}, X, \sigma, I) \, d\boldsymbol{\beta} \\ &= \int_{\mathbb{R}^{p+1}} p(y_{\text{new}} | \mathbf{x}_{\text{new}}, \boldsymbol{\beta}, I) p(\boldsymbol{\beta} | \mathbf{y}, X, \sigma, I) \, d\boldsymbol{\beta} \end{aligned}$$

- But

$$p(y_{\text{new}} | \mathbf{x}_{\text{new}}, \boldsymbol{\beta}, I) = \delta(y_{\text{new}} - \boldsymbol{\beta}^t \mathbf{x}_{\text{new}})$$

- Fix $\beta_0 = y_{\text{new}} - \beta_1 x_{\text{new},1} - \dots - \beta_p x_{\text{new},p}$

Posterior predictive distribution (2)

- Marginalize over β_p with flat priors:

$$p(y_{\text{new}} | \mathbf{x}_{\text{new}}, \mathbf{y}, X, \sigma, I)$$

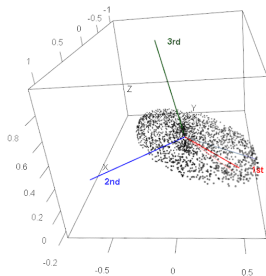
$$\propto \sigma^{-n} \int_{\mathbb{R}^p} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n \left[y_i - y_{\text{new}} + \sum_{j=1}^p \beta_j (x_{\text{new},j} - x_{ij}) \right]^2 \right\} d\beta_p$$

- After (quite some) algebra, one finds simply

$$y_{\text{new,MAP}} = \sum_{j=1}^p x_{\text{new},j} \beta_{\text{MAP},j} + \beta_{\text{MAP},0}$$

- Simpler derivation based on properties of \mathbb{E} and Var
- General posterior more complicated!

Multicollinearity



Detection

- Correlation matrices
- Belsley collinearity diagnostics

Remediation

- Eliminate predictor variables
- Principal component regression
- Regularization: ridge, lasso, elastic net, ...

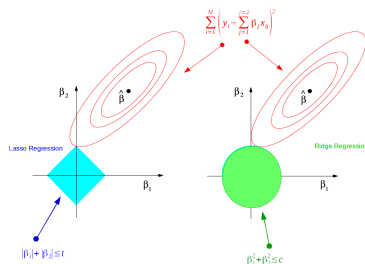
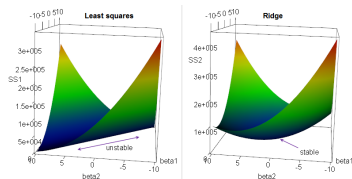
Ridge regression and lasso

- Ridge regression (Tikhonov regularization) or zero-mean normal prior:

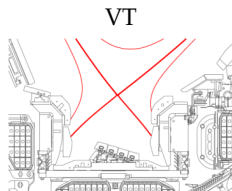
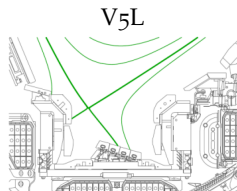
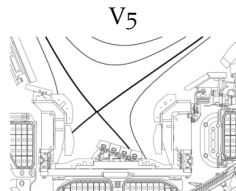
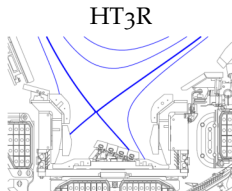
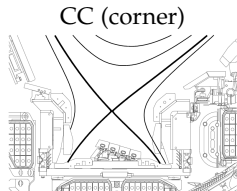
$$\alpha_{\text{ridge}} = \arg \min_{\alpha} \left[(\mathbf{y} - X\alpha)^t (\mathbf{y} - X\alpha) + \lambda \sum_{j=0}^p \alpha_j^2 \right]$$

- Lasso or zero-mean Laplace prior:

$$\alpha_{\text{lasso}} = \arg \min_{\alpha} \left[(\mathbf{y} - X\alpha)^t (\mathbf{y} - X\alpha) + \lambda \sum_{j=0}^p |\alpha_j| \right]$$



Divertor configurations at JET



Categorical variables

- Loglinear model:

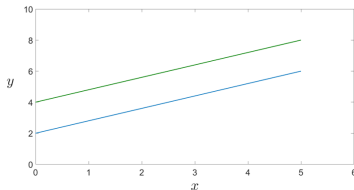
$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ \vdots \\ y_{n-1} \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & x_{B_t,1} & x_{n_e,1} \\ 1 & 0 & 0 & 0 & 0 & x_{B_t,2} & x_{n_e,2} \\ 0 & 1 & 0 & 0 & 0 & x_{B_t,3} & x_{n_e,3} \\ 0 & 1 & 0 & 0 & 0 & x_{B_t,4} & x_{n_e,4} \\ 0 & 1 & 0 & 0 & 0 & x_{B_t,5} & x_{n_e,5} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & 1 & x_{B_t,n-1} & x_{n_e,n-1} \\ 0 & 0 & 0 & 0 & 1 & x_{B_t,n} & x_{n_e,n} \end{bmatrix} \begin{bmatrix} \alpha_{0,CC} \\ \alpha_{0,HT3R} \\ \alpha_{0,V5} \\ \alpha_{0,V5L} \\ \alpha_{0,VT} \\ \alpha_B \\ \alpha_n \end{bmatrix}$$

$x_{0,CC}$ $x_{0,HT3R}$ $x_{0,V5}$ $x_{0,V5L}$ $x_{0,VT}$ x_{B_t} x_{n_e}

} Intercepts
 } Slopes

- Statistical model:

$$y = X\alpha + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$$



Global confinement scaling

- Power law + log-transform:

$$\tau_{E,\text{th}} = e^{\alpha_0} I_p^{\alpha_I} B_t^{\alpha_B} \bar{n}_e^{\alpha_n} P_{l,\text{th}}^{\alpha_P} R_{\text{geo}}^{\alpha_R} (1 + \delta)^{\alpha_\delta} \kappa_a^{\alpha_\kappa} e^{\alpha_\epsilon} M_{\text{eff}}^{\alpha_M}$$

$$\eta = \ln \tau_{E,\text{th}}, \quad \zeta_1 = \ln I_p, \quad \dots, \quad \zeta_p = \ln M_{\text{eff}}$$

- Errors in all variables:

$$\eta = \alpha_0 + \sum_{j=1}^p \alpha_j \zeta_j$$

$$y = \eta + \epsilon_y, \quad x_1 = \zeta_1 + \epsilon_{x_1}, \quad \dots, \quad x_p = \zeta_p + \epsilon_{x_p}$$

$$\epsilon_y \sim \mathcal{N}(0, \sigma_y^2), \quad \epsilon_{x_1} \sim \mathcal{N}(0, \sigma_{x_1}^2), \quad \dots, \quad \epsilon_{x_p} \sim \mathcal{N}(0, \sigma_{x_p}^2)$$

$$\sigma_{\text{mod}}^2 = \sigma_y^2 + \sum_{j=1}^p \alpha_j^2 \sigma_{x_j}^2$$

Robust Bayesian regression

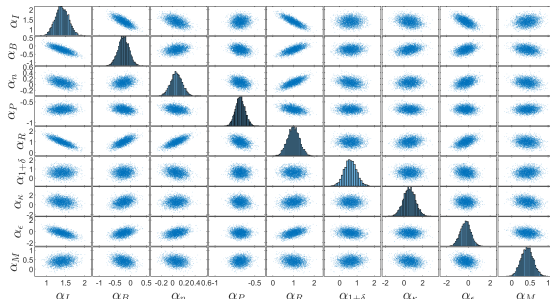
- Robust Bayesian regression **RBAYES**

$$p(\{y_{i_k,k}\}, \{x_{i_k,j,k}\} | \{\alpha_0, \alpha_j\}, \{\gamma_k\})$$

$$= \prod_k \prod_{i_k} \frac{1}{\sqrt{2\pi\gamma_k^2\sigma_{\text{mod},i_k,k}^2}} \exp \left[-\frac{1}{2} \frac{(y_{i_k,k} - \eta_{i_k,k})^2}{\gamma_k^2\sigma_{\text{mod},i_k,k}^2} \right]$$

1 for each device

- Sensitivity analysis \rightarrow practical error bars:



Geodesic least squares

- Geodesic least squares: **GLS**

$$\prod_k \prod_{i_k} \frac{1}{\sqrt{2\pi\sigma_{\text{tot},i_k,k}^2}} \exp \left[-\frac{1}{2} \frac{(y_{i_k,k} - \eta_{i_k,k})^2}{\sigma_{\text{mod},i_k,k}^2} \right]$$

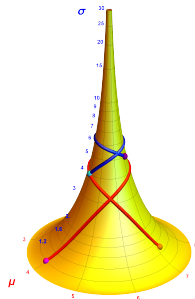
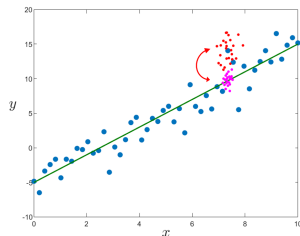


Rao geodesic distance (GD)

$$\frac{1}{\sqrt{2\pi} \sigma_{\text{obs}}} \exp \left[-\frac{1}{2} \frac{(y - y_i)^2}{\sigma_{\text{obs}}^2} \right]$$

G. Verdoolaege *et al.*, Nucl. Fusion, **55**, 113019, 2015

G. Verdoolaege *et al.*, Entropy, **17**, 4602–4626, 2015



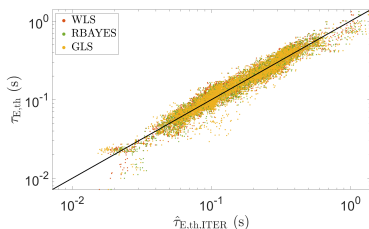
Multi-machine engineering scaling

STD5-IL ELMy H-mode (error bars from Bayesian analysis)

Engineering scaling *ITPA20-IL*

$$\tau_{E,th} = (0.067 \pm 0.060) I_p^{1.29 \pm 0.17} B_t^{-0.13 \pm 0.17} \bar{n}_e^{0.15 \pm 0.10} P_{1,th}^{-0.644 \pm 0.060} R_{geo}^{1.19 \pm 0.29} \\ \times (1 + \delta)^{0.56 \pm 0.35} \kappa_a^{0.67 \pm 0.65} M_{eff}^{0.30 \pm 0.17} \rightarrow \mathbf{H_{20}}$$

$$\hat{\tau}_{E,th,ITER} = 2.79 \pm 0.44 \text{ s}$$



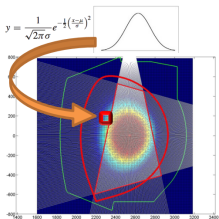
G. Verdoolaege *et al.*, Nucl. Fusion, **61**, 076006, 2021

G. Verdoolaege *et al.*, 27th IAEA Fusion Energy Conference, EX/P7-1, Gandhinagar, India, 2018

S. Kaye *et al.*, 60th Annual Meeting of the APS Division of Plasma Physics, TP11.00104, Portland, OR, USA, 2018

1. Classical probability and statistics
2. Principles of Bayesian probability theory
3. Applications
 - Classification
 - Regression analysis
 - **Some other applications**
4. Conclusions and references

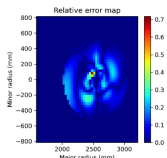
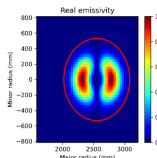
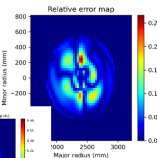
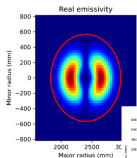
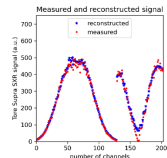
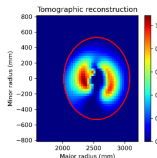
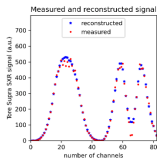
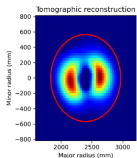
Gaussian process tomography



$$\Sigma = \begin{pmatrix} k(\vec{r}_1, \vec{r}_1) & \cdots & k(\vec{r}_1, \vec{r}_n) \\ \vdots & \ddots & \vdots \\ k(\vec{r}_n, \vec{r}_1) & \cdots & k(\vec{r}_n, \vec{r}_n) \end{pmatrix}, \quad k(\vec{r}_i, \vec{r}_j) = \sigma_f^2 \exp\left[-\left(\frac{\|\vec{r}_i - \vec{r}_j\|^2}{2\sigma_l^2}\right)\right]$$

J. Svensson, EFDA-ET-PR(11)24, 2011
 D. Li et al., Rev. Sci. Instrum. **84**, 083506, 2013
 T. Wang et al., Rev. Sci. Instrum. **89**, 063505, 2018

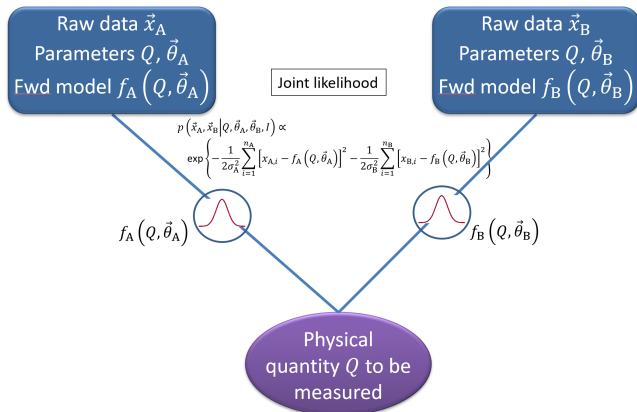
SXR tomography



Results by H. Wu, UGent

Information integration

- Data fusion / sensor fusion / integrated data analysis (IDA)



$$\underbrace{p(Q, \theta_A, \theta_B | x_A, x_B, I)}_{\text{Posterior}} \propto \underbrace{p(x_A, x_B | Q, \theta_A, \theta_B, I)}_{\text{likelihood}} \underbrace{p(Q, \theta_A, \theta_B | I)}_{\text{Prior}}$$

See IDA session on Friday!

1. Classical probability and statistics
2. Principles of Bayesian probability theory
3. Applications
 - Classification
 - Regression analysis
 - Some other applications
4. Conclusions and references

Conclusions

- Frequentist vs. Bayesian methods: interpretation of probability
- Bayesian probability: extension of logic to situations with uncertainty
- Posterior probability of parameters or hypotheses
- Numerical approach in general
- Underlies or explains many machine learning techniques

References

- D.S. Sivia and J. Skilling, *Data Analysis: A Bayesian Tutorial*, 2nd edition, Oxford University Press, 2006
- W. von der Linden, V. Dose and U. von Toussaint, *Bayesian Probability Theory: Applications in the Physical Sciences*, Cambridge University Press, 2014
- P. Gregory, *Data Analysis: A Bayesian Tutorial*, 2nd edition, Oxford University Press, 2006
- S. Theodoridis, *Machine Learning: A Bayesian and Optimization Perspective*, Academic Press (Elsevier), 2015
- E.T. Jaynes (G.L. Bretthorst, ed.), *Probability Theory: The Logic of Science*, Cambridge University Press, 2003
- S.B. McGrayne, *The theory that would not die: how Bayes' rule cracked the enigma code, hunted down Russian submarines, and emerged triumphant from two centuries of controversy*, Yale University Press, 2011