

Modeling Fusion Data in Probabilistic Metric Spaces for the Identification of Confinement Regimes and Scaling Laws

Geert Verdoolaege and Guido Van Oost

Department of Applied Physics, Ghent University, Ghent, Belgium

IAEA FEC 2012 San Diego, CA, October 10, 2012 Pattern recognition plays an important role and has great potential in fusion data analysis. However, a drawback is that individual measurements are usually represented as unstructured points in a Euclidean data space. We argue that a fundamentally probabilistic approach offers significant advantages. It allows representing the data in a non-Euclidean probabilistic space, wherein the patterns of interest are much more distinct, simply because they are based on more information. In this work, we address the identification of confinement regimes and the establishment of a scaling law for the energy confinement time, using data from the International Global H-mode Confinement Database. We propose a single-level and a Bayesian multilevel model for capturing the statistical data uncertainty. We then show that pattern recognition operations working in the associated probability space are considerably more powerful than their counterparts in a Euclidean data space. This opens up new possibilities for analyzing confinement data and for fusion data processing in general.

Pattern recognition opportunities

- Dimensionality reduction: data visualization
- Clustering/classification: grouping of data points
- Regression: (nonlinear) deterministic relation between variables

Objectives

- Physics from information: contribute to physics studies by extracting patterns, structure and relations from data
- 2 Contribute to plasma control through real-time data interpretation

Probability is fundamental

- Measurements are uncertain due to a lack of information:
 - Systematic: estimated through cross-validation
 - Stochastic: needs probability theory
- Traditional measurement: value + error bar
- Measurement = sample from underlying (non-Gaussian?) probability distribution
- Goal of measurement = probing the underlying distribution
- Probability density function (PDF) contains all information about measurement
- PDF is fundamental object resulting from measurement

Problem statement

- Pattern recognition usually operates in Euclidean data spaces with structureless data points
- This neglects additional information in the PDF!
- A huge potential remains unexplored: use additional information on the data to determine patterns:
 - Data probability distribution
 - Established theoretical models
 - Previous experiments

Challenge

To construct a probabilistic pattern recognition framework that exploits all available information for fast and efficient pattern recognition.

 \rightarrow Pattern recognition based on non-Euclidean geometry in probability spaces

Probability + geometry: a happy marriage

- Probabilistic manifold:
 - PDF = point on manifold
 - Coordinates = PDF parameters
- Distance between PDFs?
- Information geometry



Information geometry

- Riemannian differential geometry
- Fisher information = unique metric tensor:

Parametric probability model:
$$p\left(\overrightarrow{x}|\overrightarrow{\theta}\right) \Longrightarrow$$

 $g_{\mu\nu}\left(\overrightarrow{\theta}\right) = -\mathbb{E}\left[\frac{\partial^2}{\partial\theta^{\mu}\partial\theta^{\nu}}\ln p\left(\overrightarrow{x}|\overrightarrow{\theta}\right)\right], \quad \mu,\nu = 1\dots N$

$$\stackrel{
ightarrow}{ heta}=$$
 N-dimensional parameter vector

Line element:

$$\mathrm{d}s^2 = g_{\mu\nu}\mathrm{d}\theta^{\mu}\mathrm{d}\theta^{\nu}$$

- Minimum-length curve: geodesic
- Geodesic distance (GD)
- Natural and theoretically well motivated distance between PDFs

Univariate Gaussian distribution

• PDF:

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

$$\mathrm{d}s^2 = \frac{\mathrm{d}\mu^2}{\sigma^2} + 2\frac{\mathrm{d}\sigma^2}{\sigma^2}$$

• Hyperbolic geometry: Poincaré half-plane model

Poincaré half-plane

$p_1: \mu_1 = -4, \sigma_1 = 0.7; \quad p_2: \mu_2 = 3, \sigma_2 = 0.2$



ITPA Confinement Database

- ITPA Global H Mode Confinement Database (DB3)
 ITER H-Mode Database Working Group
 D.C. McDonald *et al.*, Nucl. Fusion **47**, pp. 147–174, 2007
 http://efdasql.ipp.mpg.de/hmodepublic
- ullet ~ 10⁴ entries from 19 tokamaks
- Approximate error estimates: limited information on PDF!
- Assume standard deviations → Gaussian PDFs (maximum entropy)
- Different machines \rightarrow different error estimates: difficult to handle using classic approach!

- Distinguish between L- and H-mode: 3845 L and 6207 H
- 8 global engineering variables: I_p, B_t, $\bar{\textit{n}}_{e},$ P_loss, R, a, M_{eff}, κ
- ullet Variables statistically independent ightarrow product of Gaussians

Note: this does not exclude the variables to be related through a deterministic relation!

Gaussian product manifold

- Plasma and machine variables: $x_{\lambda} \rightarrow \vec{x}$, $\lambda = 1, \dots, 8$
- Distribution parameters: μ_{λ} , σ_{λ}
- Gaussian product distribution:

$$p\left(\overrightarrow{x}|\mu_1,\ldots,\mu_8,\sigma_1,\ldots,\sigma_8\right) = \prod_{\lambda=1}^8 \mathcal{N}\left(x_\lambda|\mu_\lambda,\sigma_\lambda\right)$$

• Measurements A and B: $\overrightarrow{\mu}^{A,B} = \left(\mu_1^{A,B}, \dots, \mu_8^{A,B}\right), \ \overrightarrow{\sigma}^{A,B} = \left(\sigma_1^{A,B}, \dots, \sigma_8^{A,B}\right)$ • GD in closed form:

$$\operatorname{GD}\left(\stackrel{\rightarrow}{\mu}{}^{A}, \stackrel{\rightarrow}{\sigma}{}^{A}||\stackrel{\rightarrow}{\mu}{}^{B}, \stackrel{\rightarrow}{\sigma}{}^{B}\right) = \sqrt{2} \left[\sum_{\lambda=1}^{8} \ln^{2} \left(\frac{1+\delta_{\lambda}^{AB}}{1-\delta_{\lambda}^{AB}}\right)\right]^{1/2},$$
$$\delta_{\lambda}^{AB} = \left[\frac{\left(\mu_{\lambda}^{A}-\mu_{\lambda}^{B}\right)^{2}+2\left(\sigma_{\lambda}^{A}-\sigma_{\lambda}^{B}\right)^{2}}{\left(\mu_{\lambda}^{A}-\mu_{\lambda}^{B}\right)^{2}+2\left(\sigma_{\lambda}^{A}+\sigma_{\lambda}^{B}\right)^{2}}\right]^{1/2}$$

Bayesian multilevel model

- Introduce additional information on tokamak and database means and standard deviations for every variable
- \bullet At the minimum: from which tokamak were data obtained? \rightarrow Expected range of variables
- Bayesian hierarchical/multilevel model
- Conjugate prior distributions for means, maximum-likelihood estimates for standard deviations



Level	Model
1	$x_{\lambda,ijk} \sim \mathcal{N}(x_{\lambda,ijk} \mu_{\lambda,jk}, \sigma_{\lambda,jk})$
2	$\mu_{\lambda,jk} \sim \mathcal{N}(\mu_{\lambda,jk} \mu_{\lambda,k}, \sigma_{\lambda,k})$
3	$\mu_{\lambda,k} \sim \mathcal{N}(\mu_{\lambda,k} \mu_{\lambda,0}, \sigma_{\lambda,0})$
4	$\mu_{\lambda,0}\sim\mathcal{N}(\mu_{\lambda,0} \phi_{\lambda}, au_{\lambda})$
5	$\phi_{\lambda} \sim \mathcal{U}(-\infty, +\infty), \tau_{\lambda} = 0.1 \phi_{\lambda}$

Multilevel posterior distributions

• Bayes' rule:

$$egin{aligned} & p(\mu_{\lambda,jk},\mu_{\lambda,k},\mu_{\lambda,0},\phi_{\lambda}|x_{\lambda,ijk},orall j,k) \ & \sim \prod_{j,k} \mathcal{N}(x_{\lambda,ijk}|\mu_{\lambda,jk},\sigma_{\lambda,jk})\mathcal{N}(\mu_{\lambda,jk}|\mu_{\lambda,k},\sigma_{\lambda,k})\mathcal{N}(\mu_{\lambda,k}|\mu_{\lambda,0},\sigma_{\lambda,0}) \ & imes \mathcal{N}(\mu_{\lambda,0}|\phi_{\lambda}, au_{\lambda}) \end{aligned}$$

• Conditional posteriors for all parameters are also Gaussian, e.g.:

$$\begin{split} \mu_{\lambda,jk} | \mu_{\lambda,k}, \mu_{\lambda,0}, \phi_{\lambda}, \overrightarrow{x}_{\lambda} \sim \mathcal{N}(\hat{\mu}_{\lambda,jk}, \hat{\sigma}_{\lambda,jk}) \\ \hat{\mu}_{\lambda,jk} &= \frac{\frac{\mu_{\lambda,k}}{\sigma_{\lambda,k}^2} + \frac{n_{jk}}{\sigma_{\lambda,jk}^2} \overline{x}_{\lambda,jk}}{\frac{1}{\sigma_{\lambda,k}^2} + \frac{n_{jk}}{\sigma_{\lambda,jk}^2}} \qquad \qquad \hat{\sigma}_{\lambda,jk}^2 = \frac{1}{\frac{1}{\frac{1}{\sigma_{\lambda,k}^2} + \frac{n_{jk}}{\sigma_{\lambda,jk}^2}}} \end{split}$$

• Estimate parameters via (Gibbs) sampling

Step 1. Calculate all pairs of GDs

 \rightarrow proximity matrix $[D_{ii}]$ Step 2. Plot points arbitrarily in 2D Euclidean spaceMulti-
dimensional
scalingStep 3. Calculate Euclidean proximity matrix $[E_{ij}]$ Step 3

Step 4. Minimize $\sum_{i,j} (D_{ij} - E_{ij})^2$

Step 5. Plot final configuration

Confinement visualization



k-nearest neighbor classification

- Confinement mode identification
- Training: 1%, testing: 99%
- k = 1: nearest neighbor
- Correct classification rates (%)



Mode	Euclidean	Euclidean	GD	GD
	w/o errors	with errors	single-level	multilevel
L	89.7	91.2	91.9	93.2
H	89.1	90 5	93 3	94 3
п	09.1	90.5	95.5	94.5

Energy confinement scaling

• Energy confinement time $\tau_{\rm E}$:

$$\tau_{\rm E} = \beta_0 \ I_{\rm p}^{\beta_1} \ B_{\rm t}^{\beta_2} \ \bar{n}_{\rm e}^{\beta_3} \ P_{\rm loss}^{\beta_4} \ R^{\beta_5} \ \epsilon^{\beta_6} \ M_{\rm eff}^{\beta_7} \ \kappa^{\beta_8}$$

$$\Rightarrow \ln(\tau_{\rm E}) = \ln(\beta_0) + \beta_1 \ln(I_{\rm p}) + \beta_2 \ln(B_{\rm t}) + \beta_3 \ln(\bar{n}_{\rm e}) + \beta_4 \ln(P_{\rm loss}) \\ + \beta_5 \ln(R) + \beta_6 \ln(\epsilon) + \beta_7 \ln(M_{\rm eff}) + \beta_8 \ln(\kappa)$$

- Question: are all assumptions fulfilled? See e.g. D.C. Mc Donald *et al.*, PPCF **48**, pp. A439–A447, 2006
- Several methods:
 - Simple linear regression: ordinary least squares (OLS)
 - Errors-in-variables (EIV): total least squares
 - Geodesic regression (GR)
 - . . .

Ordinary least squares



Minimize SS_{err} = ∑_i(y_i − ŷ_i)²
Coefficient of determination:

$$\mathcal{R}^2 = 1 - rac{\mathrm{SS}_{\mathrm{err}}}{\mathrm{SS}_{\mathrm{tot}}} = 1 - rac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - ar{y})^2}$$

Errors-in-variables



Geodesic regression (GR)



- Minimize sum of squared GD
- Use geodesic centroid

Synthetic data

$$\begin{cases} y^* = 0.7 + 1.6x^* \\ x = x^* + \eta, \quad \eta \sim \mathcal{N}(0, 4) \\ y = y^* + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 8) \end{cases}$$

	Original	OLS	EIV	GR
β_0	0.70	7.21	8.41	0.28
β_1	1.60	0.45	0.23	1.61



DB3 standard data set

• Multiplicative Gaussian noise, both dependent and independent

 \bullet Logarithmic transformation \Longrightarrow additive non-Gaussian noise

• Approximate with Gaussian distributions

GR empirical and model distribution

Empirical:

$$\ln(\tau_{\rm E}) \sim \mathcal{N}(\mu_{\ln \tau_{\rm E}}, \sigma_{\ln \tau_{\rm E}})$$

Model:

$$\begin{split} \ln(\tau_{\rm E}^{*}) &= \ln(\beta_{0}) + \beta_{1} \ln(l_{\rm p}^{*}) + \beta_{2} \ln(B_{\rm t}^{*}) + \beta_{3} \ln(\bar{n}_{\rm e}^{*}) + \beta_{4} \ln(P_{\rm loss}^{*}) \\ &+ \beta_{5} \ln(R^{*}) + \beta_{6} \ln(\epsilon^{*}) + \beta_{7} \ln(M_{\rm eff}^{*}) + \beta_{8} \ln(\kappa^{*}) \\ \mu_{\ln\tau_{\rm E}^{*}} &= \ln(\beta_{0}) + \beta_{1} \mu_{\ln}_{l_{\rm p}^{*}} + \beta_{2} \mu_{\ln}_{R_{\rm t}^{*}} + \beta_{3} \mu_{\ln}_{n_{\rm e}^{*}} + \beta_{4} \mu_{\ln}_{P_{\rm loss}^{*}} + \beta_{5} \mu_{\ln}_{R^{*}} \\ &+ \beta_{6} \mu_{\ln\epsilon^{*}} + \beta_{7} \mu_{\ln}_{M_{\rm eff}^{*}} + \beta_{8} \mu_{\ln\kappa^{*}} \\ \sigma_{\ln\tau_{\rm E}^{*}}^{2} &= \beta_{1}^{2} \sigma_{\ln}^{2}_{l_{\rm f}} + \beta_{2}^{2} \sigma_{\ln}^{2}_{R_{\rm t}^{*}} + \beta_{3}^{2} \sigma_{\ln}^{2}_{l_{\rm n}} R^{*} + \beta_{4}^{2} \sigma_{\ln}^{2}_{P_{\rm loss}^{*}} + \beta_{5}^{2} \sigma_{\ln}^{2}_{R^{*}} \\ &+ \beta_{6}^{2} \sigma_{\ln\epsilon^{*}}^{2} + \beta_{7}^{2} \sigma_{\ln}^{2}_{R_{\rm eff}^{*}} + \beta_{8}^{2} \sigma_{\ln\kappa^{*}}^{2} \\ &+ \beta_{6}^{2} \sigma_{\ln\epsilon^{*}}^{2} + \beta_{7}^{2} \sigma_{\ln}^{2}_{R_{\rm eff}^{*}} + \beta_{8}^{2} \sigma_{\ln\kappa^{*}}^{2} \\ &+ \beta_{6}^{2} \sigma_{\ln\epsilon^{*}}^{2} + \beta_{7}^{2} \sigma_{\ln}^{2}_{R_{\rm eff}^{*}} + \beta_{8}^{2} \sigma_{\ln\kappa^{*}}^{2} \\ &+ \beta_{6}^{2} \sigma_{\ln\epsilon^{*}}^{2} + \beta_{7}^{2} \sigma_{\ln}^{2}_{R_{\rm eff}^{*}} + \beta_{8}^{2} \sigma_{\ln\kappa^{*}}^{2} \\ &+ \beta_{6}^{2} \sigma_{\ln\epsilon^{*}}^{2} + \beta_{7}^{2} \sigma_{\ln}^{2}_{R_{\rm eff}^{*}} + \beta_{8}^{2} \sigma_{\ln\kappa^{*}}^{2} \\ &+ \beta_{6}^{2} \sigma_{\ln\epsilon^{*}}^{2} + \beta_{7}^{2} \sigma_{\ln\kappa^{*}}^{2} \\ &+ \beta_{6}^{2} \sigma_{\ln\epsilon^{*}}^{2} \\ &+ \beta_{6}^{2} \sigma_{\ln\epsilon^{$$

Regression results



 R^2 values:

OLS	EIV	GR single	GR multi	OLS in GR single	EIV in GR single
0.94	0.97	0.71	0.78	0.44	0.52

• GR yields full probability distributions

• Precise error estimates are not required, but may improve estimates

Conclusion

- Huge potential for pattern recognition in fusion
 - Physics studies
 - Plasma control
- Probability distributions are maximally informative
- Full probability structure actively determines patterns
- Geodesic distance is natural
- Even approximate probabilities are useful
- Very flexible to other probability models

Probability is fundamental

Probability distributions contain useful information for pattern recognition. Any useful information can be incorporated, including trustable theoretical models!