

Integrated Data Analysis of the DIII-D Density Profile

L. Stagner¹, W.W. Heidbrink¹, MA Van Zeeland², and L. Zeng²

¹University of California–Irvine, Irvine, California, USA

²General Atomics, P.O. Box 85608, San Diego, California 92186-5608, USA

Abstract

1 Introduction

Introstuff

2 Integrated Data Analysis

As fusion experiments expand, so does the number of diagnostics used to probe the underlying physics. With the increase in diagnostics the number of interrelations among the diagnostics increases accordingly. While each diagnostic can act independently it is often advantageous to combine the data sets to acquire the most accurate result that is consistent with all the information available. This causes difficulties when trying to combine multiple data sets that cover complementary areas of interest. Bayesian data analysis give a comprehensive, scalable, and automated framework for combining complementary diagnostics.

2.1 The Bayesian perspective

Bayesian statistics differs from the prevailing frequentist statistics that is taught in most introductory courses. The difference between the two views is in the definition of probability. In the Bayesian perspective, probability represented a degree of belief or plausibility that an event was true, based on the information available. This school of thought is opposed to frequentist view that probability is a long run relative frequency with which the event occurs, given infinitely many repeated experimental trials.[1]

While the frequentist view is seemingly more objective, it is limited in its straightforward application since some things cannot have long run relative frequencies and, as a result, cannot be analyzed by probability theory. For instance a constant, such as the mass of an object, cannot be analyzed by a straightforward application of the frequentist definition of probability as a long run relative frequency. To use the frequentist view you first have to relate the mass to the data through some function called the *statistic*. Since statistic is subject to random noise it becomes the random variable to which the rules of probability can then be applied. The issue with this approach is that there is no natural way of choosing the best statistic.[2] The founders of the field have a plethora of ad hoc rules and tests to determine which statistic should be used in any situation. These rules and tests hides much of the mechanics and assumptions inherent in the analysis and the consequences can often turn up at the worst possible moments.

The bayesian definition has no such difficulty since nearly anything can be thought of as a probability and the framework explicitly formulates all assumptions. That is not to say that it is without its issues. The frequentist view is often much easier to use because it hides most of the difficulties from the user. In bayesian statistics you have to be very careful with the assumptions that are made since it can drastically affect the results. Also, the calculation of optimal parameters is also a challenging numerical exercise in Bayesian statistics. It was partially for this reason that the frequentist view arose into prominence. It was only in the mid 20th century with the advent of computers and efficient sampling algorithms that Bayesian statistics has begun its resurgence. In recent years Bayesian analysis has become essential to many scientific fields such as cosmology and artificial intelligence.[1]

2.2 Mathematical Formulation

In 1946, Richard Cox began to formulate the rules for logical and consistent reasoning by considering how we might quantify our beliefs about the truth of an event or object. He concluded that the real numbers we attached to our beliefs followed the rules of probability:[2]

$$prob(X|I) + prob(\bar{X}|I) = 1 \quad (2.1)$$

and

$$prob(X, Y|I) = prob(X|Y, I) \times prob(Y|X, I) \quad (2.2)$$

Here \bar{X} denotes that X is false and the vertical bar ‘|’ means ‘given’. All the probabilities are conditional on I to denote the relevant background information available, since there is no such thing as an absolute probability.

Equation 2.1 and 2.2 are called the sum and product rules, respectively. The sum and product rule form the basic equations of probability theory. From these two rules we can derive the core equations used in bayesian statistics: Bayes’ theorem

$$prob(X|Y, I) = \frac{prob(Y|X, I) \times prob(X|I)}{prob(Y|I)} \quad (2.3)$$

and marginalization.

$$prob(X|I) = \int_{-\infty}^{\infty} (X, Y|I) dY \quad (2.4)$$

The power behind Bayes theorem can best be recognized with a change of notation.

$$\begin{aligned} &prob(hypothesis|\{data\}, I) = \\ &\frac{prob(\{data\}|hypothesis, I) \times prob(hypothesis|I)}{prob(\{data\}|I)} \end{aligned} \quad (2.5)$$

or as denoted in most parameter estimation problems:

$$\mathcal{P}(\vec{\alpha}|\vec{d}, I) = \frac{\mathcal{L}(\vec{d}|\vec{\alpha}, I)\pi(\vec{\alpha}|I)}{\mathcal{Z}(\vec{d}|I)} \quad (2.6)$$

where $\vec{\alpha}$ are the parameters to be inferred. The terms in the above equation have formal names. $prob(hypothesis|I)$ is called the *prior*. The prior

encodes the current state of knowledge about the parameters before we have analyzed any new data. When we have acquired new information the prior is modified by the first term, $\text{prob}(\{\text{data}\}|\text{hypothesis}, I)$ called the *likelihood*. The product of the likelihood and the prior give $\text{prob}(\text{hypothesis}|\{\text{data}\}, I)$, which is called the *posterior*. The posterior encodes our final state of knowledge after all the data has been incorporated. The maximum value of the posterior gives us the best estimate of the parameters we trying to infer. The last term, $\text{prob}(\{\text{data}\}|I)$, is called the *evidence* or *marginal likelihood*. In parameter estimation problems the evidence can be ignored since it is essentially a normalization constant. However, in model comparison problems it is of vital importance.

2.3 Choosing the Best Prior

It was noted earlier that one has to be extra careful with the assumptions that are made prior to incorporating data. Quite fittingly, these assumptions are encoded in the choice of prior. This process can be difficult and if often done incorrectly. In order to choose the best prior three main schools of thought for choosing priors have been developed.

The first school of thought is to choose a prior that expresses specific, definite information about a parameter. If, for instance, a parameter is known, through past experience or expert testimony, to have a value around 1 ± 0.2 with upper bound of 2. It would not be uncommon for someone to assign a prior that is a Gaussian with a mean of 1 and standard deviation of 0.2 that is zero for values greater than 2. Informative priors are both a major advantage and can also be a major pitfall. If our prior knowledge was incorrect and the upper bound is around 5, depending on the amount of data, we may have biased the posterior such that it doesn't reflect the truth. We must then be careful not to bias our results without proper justification. Informative priors also suffer from the fact that there is not a well defined procedure for choosing them.

The second school of thought is to be as non-informative as possible. This method should be used when we have no prior information about the parameter in question. The most common type of non-informative prior is called the Jeffreys' Prior. The key property of the Jeffreys' Prior is that it is invariant under re-parametrization. The Jeffreys' prior is defined as

$$\pi(\vec{\alpha}) \propto \sqrt{\det \mathcal{I}(\vec{\alpha})} \quad (2.7)$$

where \mathcal{I} is the Fisher information. One of the main problems with the Jeffreys' Prior is that it can be improper in the sense that it cannot be normalized to one (although there are ways around this). Improper priors should not be used since it can lead to paradoxes.[1] The Jeffreys' Prior is also sometimes impossible to calculate and is therefore of limited use.

The third school of thought is to choose the prior that has the largest Shannon information entropy that is still consistent with the given testable information. The idea is that entropy is a measure of “uninformativeness” so picking the most ‘uninformative’ prior that is still consistent with testable information would be the best choice. Mathematically, this is the constrained optimization problem:

$$Q = - \int \pi(x) \log\left(\frac{\pi(x)}{m(x)}\right) dx - \lambda C(x) \quad (2.8)$$

where $p(x)$ is the prior, $m(x)$ is the Lebesgue measure which ensures invariance under transformation, λ is the Lagrange multiplier, and $C(x)$ is the constraint function. For example, if we knew the mean of a parameter, using the principle of maximum entropy the best prior would be a Poisson distribution. Likewise, if we knew both the mean and the variance of the parameter the principle would yield a Gaussian. The principle of maximum entropy provides a balanced and systematic approach to problem of picking the right prior. It is for these reasons that we will endeavour to apply the principle of maximum entropy whenever possible.

2.4 Model Comparison

In a bayesian framework it is possible to compare two competing models to see which one best describes the observed measurements. Let \vec{d} be some observed data. A model M has parameters $\vec{\alpha}$, with prior $\pi(\vec{\alpha}|M)$ and likelihood $\mathcal{L}(\vec{d}|\vec{\alpha}, M)$. Applying marginalization, the model evidence is defined as:

$$\mathcal{Z}(\vec{d}|M) = \int \mathcal{L}(\vec{d}|\vec{\alpha}, M) \pi(\vec{\alpha}|M) d\vec{\alpha} \quad (2.9)$$

You will notice the above equation is of the same form as Eq. 2.6 with $M \rightarrow I$. To compare two models M_1 and M_2 one may compute the ratio of the models evidences, called the Bayes Factor[1]:

$$B(M_1, M_2) = \frac{\mathcal{Z}(\vec{d}|M_1)}{\mathcal{Z}(\vec{d}|M_2)} \quad (2.10)$$

Bayes factors greater than one favour model M_1 and values less than one favour model M_2 . Jefferys[3] gave the following qualitative interpretation of a Bayes factor:

$\mathbf{B}(M_1, M_2)$	Interpretation
$B(M_1, M_2) < .10$	Strong Evidence for M_2
$.10 < B(M_1, M_2) < .33$	Moderate Evidence for M_2
$.33 < B(M_1, M_2) < 1.0$	Weak Evidence for M_2
$1.0 < B(M_1, M_2) < 3.0$	Weak Evidence for M_1
$3.0 < B(M_1, M_2) < 10.$	Moderate Evidence for M_1
$B(M_1, M_2) > 10.$	Strong Evidence for M_1

2.5 Combining multiple diagnostics

One of the advantages of bayesian statistics is that it provides a systematic framework for combining data sets from multiple sources. This can be seen from Eq. 2.6. Consider just two data points D_1 and D_2 . Bayes theorem would yield:

$$prob(H|D_1, D_2, I) \propto prob(D_1, D_2|H, I) \times prob(H|I) \quad (2.11)$$

We can use Bayes theorem to express the posterior to be conditional on D_1 .

$$prob(H|D_1, D_2, I) \propto prob(D_2|H, D_1, I) \times prob(H|D_1, I) \quad (2.12)$$

This shows that the prior in Eq. 2.11 can be replaced by the posterior based on D_1 .

If D_1 and D_2 are independent then

$$prob(D_1|H, D_2, I) = prob(D_1|H, I) \quad (2.13)$$

and

$$prob(D_2|H, D_1, I) = prob(D_2|H, I) \quad (2.14)$$

Substituting Eq. 2.14 into Eq. 2.12 and applying Bayes theorem to the prior yields:

$$prob(H|D_1, D_2, I) \propto prob(D_2|H, I) \times prob(D_1|H, I) \times prob(H|I) \quad (2.15)$$

This result can be generalized for K data points as

$$prob(H|D, I) \propto \left(\prod_{k=0}^K prob(D_k|H, I) \right) \times prob(H|I) \quad (2.16)$$

This result forms the basis for integrated data analysis. So long as all the data is independent from each other, we can assign different likelihoods for each data point. This allows for the data to have different types of errors. For instance, we could combine a diagnostic that is subject to systematic error with a diagnostic whose error is distributed according to a Poisson or Gaussian distribution. Using orthodox methods this could not be done easily, but within a bayesian framework it flows naturally from the basic equations.

Looking more closely at Eq. 2.16 you will notice that each likelihood is dependent on H . If we put Eq. 2.16 into the notation of Eq. 2.6:

$$\mathcal{P}(\vec{\alpha}|D, I) \propto \left(\prod_{k=0}^K \mathcal{L}(D_k|\vec{\alpha}, I) \right) \times \pi(\vec{\alpha}|I) \quad (2.17)$$

We can easily see that integrated data analysis requires that we have a forward model for each diagnostic that depends on a set of common parameters $\vec{\alpha}$. The requirement of a forward model for each diagnostic forces one to formulate a model linking the physics underlying a measurement to the resulting measurement. Mathematically, for J number of diagnostics:

$$d_j = f_j(\vec{\beta}, \vec{\gamma}_j) \quad (2.18)$$

where d_j and f_j is the theoretical measurement and forward model for diagnostic j , respectively. The vector $\vec{\beta}$ are the common parameters that are shared between diagnostics. The vector $\vec{\gamma}_j$ are parameters that are unique to each diagnostic. These are also called *nuisance* parameters since they are of no particular interest but are needed for the forward model. These nuisance parameters mean that we need to modify the definition of the prior in Eq. 2.18 to

$$\pi(\vec{\alpha}|I) = \pi(\vec{\beta}|I) \times \left(\prod_{j=0}^J \pi(\vec{\gamma}_j|I) \right) \quad (2.19)$$

If there are no nuisance parameters for the diagnostics the above equation reduces to the prior of Eq. 2.17.

In the following sections we will infer the DIII-D density profile from multiple diagnostics. The analysis is based off van Milligen's work on inferring the density profile at the TJ-II stellarator.[4]

3 Forward modelling of density diagnostics

3.1 Functional form of the density profile

In following the analysis performed by van Milligen, we can parametrize the density profile as a Fourier-Bessel series:

$$n_e(\rho) = \sum_{k=0}^K \beta_k J_0(\lambda_k \rho / \rho_{max}) \quad (3.1)$$

where ρ is the magnetic flux, λ_k is the k^{th} zero of the Bessel function J_0 and β_k are the parameters to be inferred. In addition to and in contrast with previous works, we allow for ρ_{max} to be a free parameter that will also be inferred. By choosing the functional form to be a function of magnetic flux we will assume that the error introduced by the process of mapping to flux coordinates from machine coordinates is negligible.

This form of the density profile has several advantageous properties:

- The derivative and therefore the particle flux at $\rho = 0$ is zero
- The expansion goes to zero at $\rho = \rho_{max}$, which is constrained to be greater than one
- The expansion is capable of approximating any continuous and square-integrable function with arbitrary accuracy
- The expansion consists of orthogonal functions.

These properties of the density profile allow for a fast and stable maximization of the profile parameters.

3.2 Interferometry

DIII-D uses a fiber optic, heterodyne, two color interferometer to give a line averaged electron density.[5] The basic principle is that the light from a single color interferometer will experience a phase shift from the plasma as well as from vibrations of the associated optics. The phase shift for a given wavelength can be expressed as:

$$\phi = A\lambda + \frac{2\pi V}{\lambda} \quad (3.2)$$

where the first term ($A\lambda$) is the plasma contribution with A related to the line-integrated density n_e according to

$$A = k \int n_e dl = k \bar{n}_e L \quad (3.3)$$

where $k = 2.82 \times 10^{-15} m$ and the second term ($2\pi V/\lambda$) is the phase shift due to vibrations V in the direction of beam propagation. To extract the plasma density a second laser(hence the name two color) is made collinear with the first. This creates a system of equations which can be solved exactly.

$$\bar{n}_e L = \frac{\lambda_2}{k(\lambda_2^2 - \lambda_1^2)} \left(\phi_2 - \frac{\lambda_1 \phi_1}{\lambda_2} \right) \quad (3.4)$$

$$2\pi V = \frac{\lambda_1 \lambda_2^2}{(\lambda_2^2 - \lambda_1^2)} \left(\phi_1 - \frac{\lambda_1 \phi_2}{\lambda_2} \right) \quad (3.5)$$

This system of equations relates the measured phase to the line-averaged density and the vibration of the plasma. As seen from Eq. 3.3 the forward model for the interferometry data is:

$$d_0(\vec{\beta}) = \bar{n}_e = \frac{1}{L} \int n_e(l, \vec{\beta}) dl \quad (3.6)$$

3.3 Reflectometry

$$d_1(\vec{\beta}) = \quad (3.7)$$

3.4 Thomson scattering

The Thomson scattering diagnostics gives local measurements of the electron density and temperature. A full forward model of the Thomson scattering diagnostic would require a parametrized temperature profile. This would introduce un-necessary complexity. To avoid this, the density measurements and errors are calculated using standard data analysis techniques. Due to drifts in the calibration of the diagnostic, an additional scale parameter A_{TS} is introduced. The forward model for a Thomson scattering measurement at ρ_i is then:

$$d_3(\vec{\beta}, A_{TS}) = A_{TS} \times n_e(\rho_i, \vec{\beta}) \quad (3.8)$$

3.5 Beam emission spectroscopy

stuff

4 Likelihood and prior probabilities

4.1 Priors

The functional form of the density profile presents some difficulty when trying to form informative priors. While we often have intimate knowledge of the general shape of a density profile, that knowledge does not easily transfer to the free parameters. The difficulty of transferring information from data space (density) to parameter space encourages the use of pseudo– data points to realize physical constraints on the data space. One could think of the use of pseudo–points as a type of forward modelling of the prior. Since we are using pseudo–points in data space, parameter space priors must be uninformative to reduce unnecessary bias.

4.1.1 Parameter space priors

The coefficients of the Fourier-Bessel series, $\vec{\beta}$, are location parameters. As such, the least informative prior would be uniform. Taking into account that at the density at $\rho = 0$ is unlikely to exceed 10^{20} m^{-3} we limit the values to be between -10^{20} m^{-3} and 10^{20} m^{-3} . Mathematically, with β_k in units of 10^{19} m^{-3} , the prior is

$$\pi(\vec{\beta}|I) \propto \prod_{k=0}^K (\Theta(\beta_k + 10) - \Theta(\beta_k - 10)) \quad (4.1)$$

The density goes to zero outside the last closed flux surface ($\rho = 1$). By design, the functional form of the density profile goes to zero at ρ_{max} . It is therefore proper to assign a prior that is uniform between $\rho = 1.1$ and $\rho = 1.2$. Mathematically this is

$$\pi(\rho_{max}|I) \propto \Theta(\rho_{max} - 1.1) - \Theta(\rho_{max} - 1.2) \quad (4.2)$$

This prior replaces the data space prior that is van Millagen’s analysis that limits the

4.1.2 Data space priors

Data space priors, henceforth called likelihood priors, use psuedo-points and therefore resemble likelihood functions that are associated with measurements. To avoid confusion with likelihood functions \mathcal{L} , we will denote likelihood priors with the superscript p .

To keep the density profile positive on the region of $0 \leq \rho \leq 1$ we introduce a set of n_1 equally spaced psuedo-points in the region with a log likelihood prior of

$$\ln(\mathcal{L}_1^p) = \sum_{i=0}^{n_1} -\frac{(min(0, n_e(\rho_i))^2}{2\sigma_1^2} \quad (4.3)$$

We choose $n_1 = 20$ and $\sigma_1 = 0.0001$ with n_e in units of $10^{19} m^{-3}$. This yeilds a harsh penalty for negative densities but does not penalize positive densities.

In a departure from the analysis of van Milligen, we do not put any constraints on the first derivative of the density profile. Instead we put a weak constraint on the second derivative to avoid large oscillatory behavior. To do this we introduce a set of n_2 equally spaced psuedo-points in the region $0 \leq \rho \leq 0.8$ with a log likelihood prior of

$$\ln(\mathcal{L}_2^p) = \sum_{i=0}^{n_2} -\frac{(\partial_\rho^2 n_e(\rho_i))^2}{2\sigma_2^2} \quad (4.4)$$

with $n_2 = 20$ and $\sigma_2 = 10$ with n_e in units of $10^{19} m^{-3}$.

4.2 Likelihoods

stuff

4.2.1 Interferometry

stuff

4.2.2 Reflectometry

stuff

4.2.3 Thomson scattering

stuff

4.2.4 Beam emission spectroscopy

stuff

4.2.5 Combined likelihood probability

stuff

5 Total posterior probability

stuff

5.1 Exploring the posterior

stuff

6 Density profile reconstructions

stuff

6.1 H-mode reconstruction

stuff otherstuff

6.2 L-mode reconstruction

stuff

7 Representation of Error

stuff

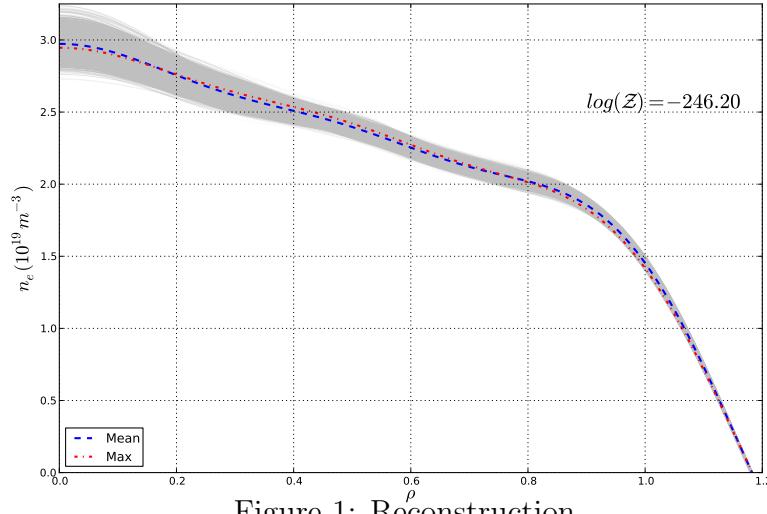


Figure 1: Reconstruction

8 Determining discrepant diagnostics

stuff

9 Conclusions and future work

stuff

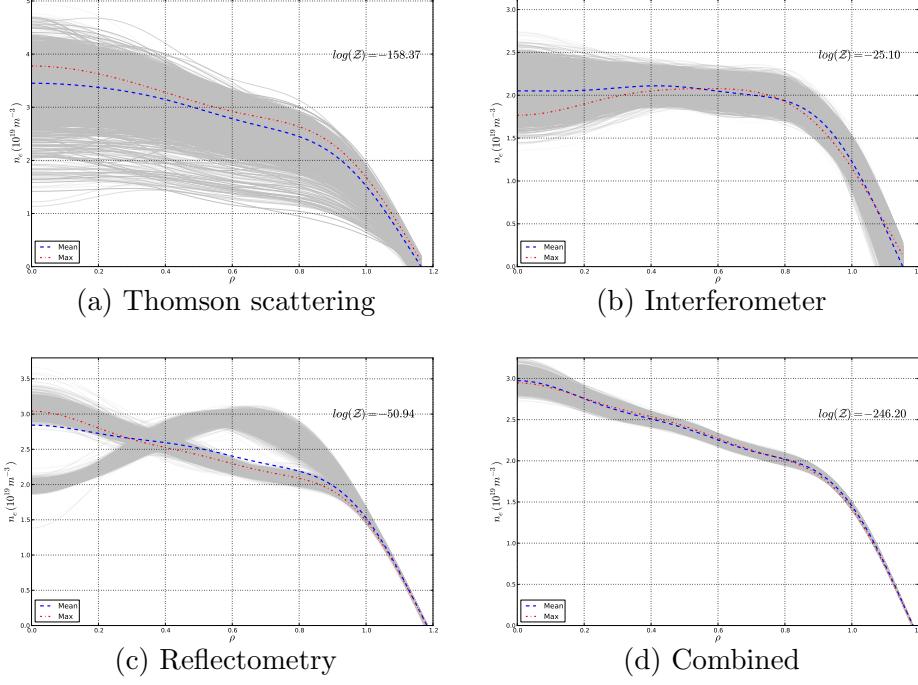


Figure 2: The l-o-n-g caption for all the subfigures (FirstFigure through FourthFigure) goes here.

References

- [1] Udo von Toussaint. Bayesian inference in physics. *Reviews of Modern Physics*, 83(3):943, 2011.
- [2] Devinderjit Sivia and John Skilling. *Data Analysis: A Bayesian Tutorial*. Oxford University Press, USA, 2006.
- [3] Harold Jeffreys. *Theory of probability*. Oxford University Press, 1998.
- [4] B Ph van Milligen, T Estrada, E Ascáibar, D Tafalla, D López-Bruna, A López Fraguas, JA Jiménez, I García-Cortés, A Dinklage, and R Fischer. Integrated data analysis at tj-ii: The density profile. *Review of Scientific Instruments*, 82(7):073503–073503, 2011.
- [5] MA Van Zeeland, RL Boivin, TN Carlstrom, T Deterly, and DK Finkenthal. Fiber optic two-color vibration compensated interferometer

for plasma density measurements. *Review of scientific instruments*, 77(10):10F325–10F325, 2006.