

# Model Criticism in Latent Space

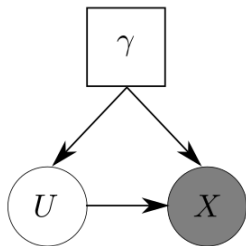
Sohan Seth

University of Edinburgh, School of Informatics, Edinburgh, EH8 9AB, UK

September 10, 2019



THE UNIVERSITY *of* EDINBURGH



- $X$  denotes observed variables
- $U$  denotes unknown variables
- $\gamma$  denotes known variables

Estimate the mean  $\mu$  of  $n$  observations

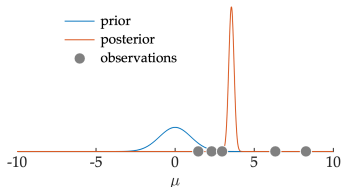
$$\mu \sim \mathcal{N}(\mu_0, \sigma_0^2) \quad (\text{prior})$$

$$X_i | \mu \sim \mathcal{N}(\mu, \sigma^2) \quad (\text{likelihood})$$

- $X = \{X_1, X_2, \dots, X_n\}$
- $U = \{\mu\}$
- $\gamma = \{\mu_0, \sigma_0^2, \sigma^2\}$

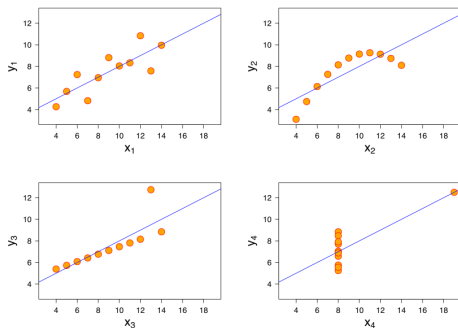
## Bayes' rule

$$p_{U|X}(u | x^{\text{obs}}, \gamma) = \frac{p_{X|U}(x^{\text{obs}} | u, \gamma) p_U(u | \gamma)}{p_X(x^{\text{obs}} | \gamma)}$$



# Motivation

- Statistical models are approximation of complex natural processes
- “all models are wrong but some are useful” [Box and Draper, 1987, p. 424]
- Is the simplification meaningful?
- Are the assumptions we make reasonable?
- Knowing the limitations can guide us to build a better model
- *Model criticism* is the process of assessing the limitations of a model



[SOURCE]

Figure: Anscombe's Quartet [Francis Anscombe 1973]

# Model Criticism in Observation Space

“if the model fits, then replicated data [ $X^{\text{rep}}$ ] generated under the model should look similar to observed data [in terms of discrepancy measure  $D$ ]” [Gelman et al., 2004, p. 165]

$$D(x, u) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

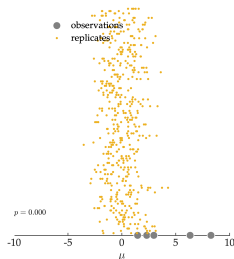
# Model Criticism in Observation Space

“if the model fits, then replicated data  $[X^{\text{rep}}]$  generated under the model should look similar to observed data [in terms of discrepancy measure  $D$ ]” [Gelman et al., 2004, p. 165]

- *prior predictive p-value* [Box, 1980]

$$p_{\text{prior}} = \Pr(D(X^{\text{rep}}, U) > D(x^{\text{obs}}, U)) \text{ where } X^{\text{rep}}, U \sim P(X, U) \quad (1)$$

$$D(x, u) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$



# Model Criticism in Observation Space

“if the model fits, then replicated data  $[X^{\text{rep}}]$  generated under the model should look similar to observed data [in terms of discrepancy measure  $D$ ]” [Gelman et al., 2004, p. 165]

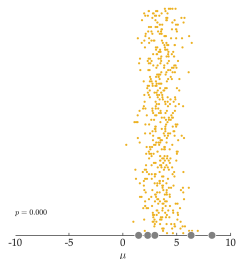
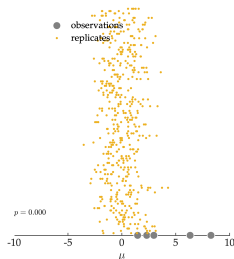
- *prior predictive p-value* [Box, 1980]

$$p_{\text{prior}} = \Pr(D(X^{\text{rep}}, U) > D(x^{\text{obs}}, U)) \text{ where } X^{\text{rep}}, U \sim P(X, U) \quad (1)$$

- *posterior predictive p-value* [Rubin, 1984]

$$p_{\text{post}} = \Pr(D(X^{\text{rep}}, U) > D(x^{\text{obs}}, U) | x^{\text{obs}}) \text{ where } X^{\text{rep}}, U \sim P(X, U | x^{\text{obs}}) \quad (2)$$

$$D(x, u) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$



# Model Criticism in Latent Space

Posterior predictive check requires [Johnson, 2007]

- ① generating replicate observations
- ② crafting an appropriate discrepancy measure
- ③ approximating the null distribution, and
- ④ “double use” of data

# Model Criticism in Latent Space

Posterior predictive check requires [Johnson, 2007]

- ① generating replicate observations
- ② crafting an appropriate discrepancy measure
- ③ approximating the null distribution, and
- ④ “double use” of data

If the model fits, then posterior inferences should match the prior assumptions.

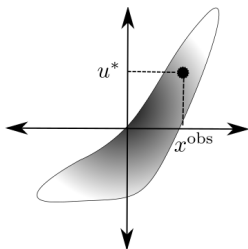


# Model Criticism in Latent Space

Posterior predictive check requires [Johnson, 2007]

- ① generating replicate observations
- ② crafting an appropriate discrepancy measure
- ③ approximating the null distribution, and
- ④ “double use” of data

If the model fits, then posterior inferences should match the prior assumptions.



$$x^{\text{obs}} \sim P(X) \text{ and } u^* | x^{\text{obs}} \sim P_{U|X}(u | x^{\text{obs}}) \Rightarrow (u^*, x^{\text{obs}}) \sim P_{U,X}(u, x) \Rightarrow u^* \sim P(U)$$

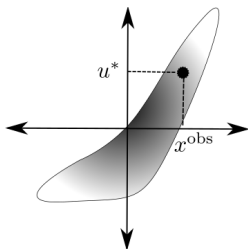
If  $x^{\text{obs}}$  is a sample from  $P(X | \gamma)$ , then a sample  $u^*$  from  $P(U | x^{\text{obs}}, \gamma)$  will be a draw from  $P(U | \gamma)$ .

# Model Criticism in Latent Space

Posterior predictive check requires [Johnson, 2007]

- ① generating replicate observations
- ② crafting an appropriate discrepancy measure
- ③ approximating the null distribution, and
- ④ “double use” of data

If the model fits, then posterior inferences should match the prior assumptions.



$$x^{\text{obs}} \sim P(X) \text{ and } u^* | x^{\text{obs}} \sim P_{U|X}(u | x^{\text{obs}}) \Rightarrow (u^*, x^{\text{obs}}) \sim P_{U,X}(u, x) \Rightarrow u^* \sim P(U)$$

If  $x^{\text{obs}}$  is a sample from  $P(X | \gamma)$ , then a sample  $u^*$  from  $P(U | x^{\text{obs}}, \gamma)$  will be a draw from  $P(U | \gamma)$ .

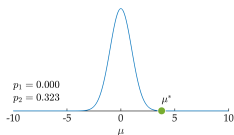
$u_1^*, \dots, u_m^* \not\sim P(U | \gamma)$ , i.e.,  $m$  posterior samples are not independent samples from the prior

# Aggregated Posterior Check

**Require:** Observed data  $x^{\text{obs}}$

**Require:** Bayesian model  $P(X | U, \gamma)P(U | \gamma)$  with latent variables  $U$

- 1: Generate a posterior sample  $u^*$  from  $P(U | x^{\text{obs}}, \gamma)$
- 2:
- 3: Compare  $u^*$  sample with corresponding prior distribution
- 4: **return** p-value of the test



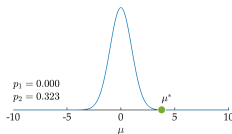
$$(1) \mu^* \sim \mathcal{N}(\mu_0, \sigma_0^2) \text{ and } (2) X_1, \dots, X_n \sim \mathcal{N}(\mu^*, \sigma^2) \quad (3)$$

# Aggregated Posterior Check

**Require:** Observed data  $x^{\text{obs}}$

**Require:** Bayesian model  $P(X | U, \gamma)P(U | \gamma)$  with latent variables  $U$

- 1: Generate a posterior sample  $u^*$  from  $P(U | x^{\text{obs}}, \gamma)$
- 2: Generate aggregated posterior sample
- 3: Compare aggregated posterior sample with corresponding prior distribution
- 4: **return** p-value of the test



$$(1) \mu^* \sim \mathcal{N}(\mu_0, \sigma_0^2) \text{ and } (2) X_1, \dots, X_n \sim \mathcal{N}(\mu^*, \sigma^2) \quad (3)$$

- Often  $U$  is a collection of variables, i.e.,  $U = (U_1, \dots, U_K)$ , and  $P(U | \gamma) = \prod_{k=1}^K P_{u_k}(U_k | \gamma)$
- Instead of testing if  $(u_1^*, \dots, u_K^*)$  is a sample from  $P(U | \gamma)$ , test if the *aggregated* variables  $\{u_1^*, \dots, u_K^*\}$  is independent and identical draws from  $P_{u_k}(\cdot | \gamma)$ .

# 1. Probabilistic Matrix Factorization

$$\begin{array}{c} n \\ \boxed{\mathbf{X}} \\ m \end{array} = \begin{array}{c} K \\ \boxed{\Theta} \end{array} \times \begin{array}{c} n \\ \boxed{\mathbf{Z}} \end{array}$$

For  $i = 1, \dots, n$

$$\mathbf{z}_i \sim \text{LatentDist} \mid \tau_z \quad (4)$$

$$\mathbf{x}_i \sim \mathcal{N}(\Theta \mathbf{z}_i + \mathbf{b}, \tau^{-1} \mathbf{I}) \quad (5)$$

Given  $\mathbf{Z}^* = [\mathbf{z}_1^*, \dots, \mathbf{z}_n^*]$  and  $\tau_z^*$ , (e.g.,  $\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \tau_z^{-1} \mathbf{I})$ )

$$\begin{aligned} \{z_{ki}^*\} &\sim P(z \mid \tau_z^*) && \text{(univariate)} \\ \{(z_{k_1 i}^*, z_{k_2 i}^*) : k_1 \neq k_2\} &\sim P(z_1, z_2 \mid \tau_z^*) && \text{(bivariate)} \end{aligned}$$

# 1. Image Patches [Zoran and Weiss, 2012]

$n = 50,000$ ,  $8 \times 8$  image patches, i.e.,  $m = 64$  and we consider  $k = 16$

$$\tau_z \sim \text{Gamma}(\alpha, \beta), z \sim \mathcal{N}(0, \tau_z^{-1}), \quad (6)$$

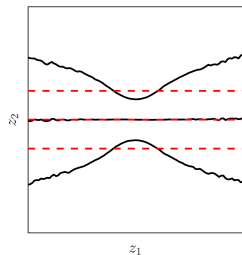
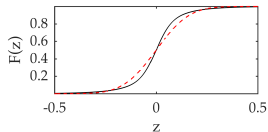


Figure: Dotted line is the prior distribution and straight line is the aggregated posterior

# 1. Image Patches [Zoran and Weiss, 2012]

$n = 50,000$ ,  $8 \times 8$  image patches, i.e.,  $m = 64$  and we consider  $k = 16$

$$\tau_z \sim \text{Gamma}(\alpha, \beta), z \sim \text{Laplace}(0, \tau_z). \quad (6)$$

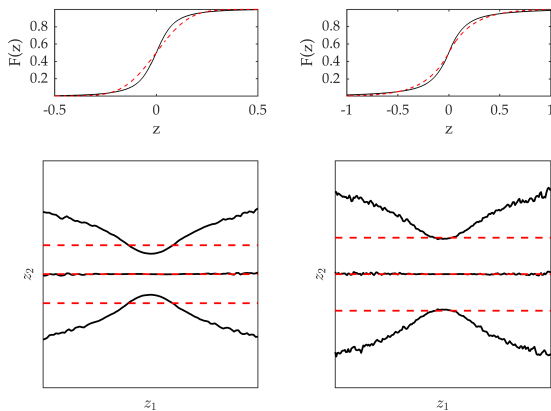


Figure: Dotted line is the prior distribution and straight line is the aggregated posterior

# 1. Image Patches [Zoran and Weiss, 2012]

$n = 50,000$ ,  $8 \times 8$  image patches, i.e.,  $m = 64$  and we consider  $k = 16$

$$\pi \sim \text{Dir}(\mathbf{1}), \tau_m \sim \text{Gamma}(\alpha, \beta), \mathbf{z} \sim \sum_{m=1}^8 \pi_m \mathcal{N}(\mathbf{0}, \tau_m^{-1} \mathbf{I}) \quad (6)$$

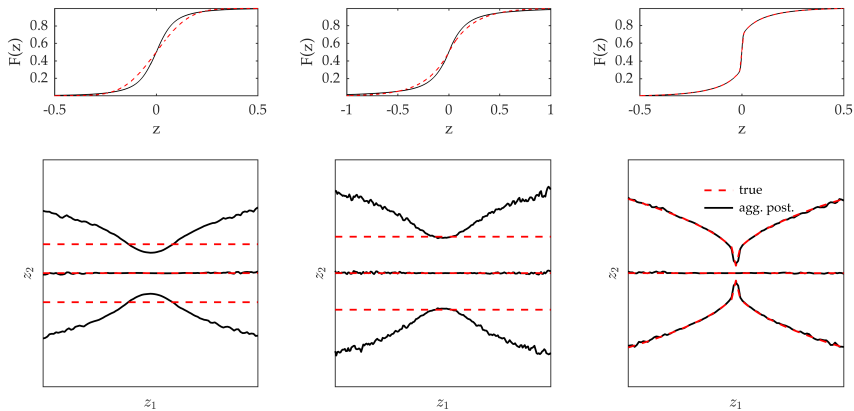
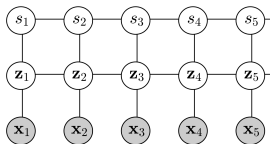


Figure: Dotted line is the prior distribution and straight line is the aggregated posterior



## 2. Linear Dynamical Systems



$$s_1 = 1, \mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (7)$$

$$s_t \sim \text{Cat}(\boldsymbol{\pi}^{(s_{t-1})}) \quad \forall t = 2, \dots, n \quad (8)$$

$$\mathbf{z}_t \sim \mathbf{A}^{(s_t)} \mathbf{z}_{t-1} + \boldsymbol{\epsilon}_t, \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{(s_t)^{-1}}) \quad \forall t = 2, \dots, n \quad (9)$$

$$\mathbf{x}_t \sim \mathbf{B} \mathbf{z}_t + \boldsymbol{\psi}_t, \boldsymbol{\psi}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R}^{-1}) \quad \forall t = 1, \dots, n \quad (10)$$

Standardized residuals follow  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  distribution

- *standardized* latent residuals

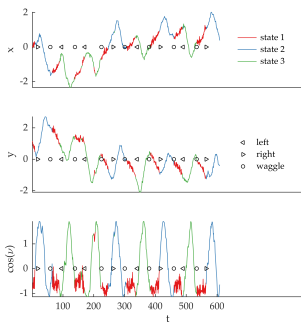
$$\tilde{\mathbf{z}}_t = (\mathbf{Q}^{(s_t^*)})^{0.5} (\mathbf{z}_t^* - \mathbf{A}^{(s_t^*)} \mathbf{z}_{t-1}^*) \quad \forall t = 2, \dots, n, \quad (11)$$

- *standardized* observation residuals (or innovations)

$$\tilde{\mathbf{x}}_t = (\mathbf{R}^*)^{0.5} (\mathbf{x}_t^{\text{obs}} - \mathbf{B}^* \mathbf{z}_t^*) \quad \forall t = 2, \dots, n. \quad (12)$$

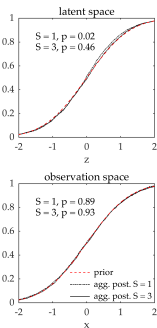
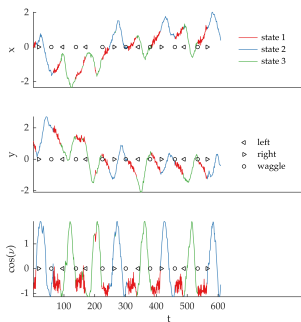
## 2. Honey Bee

- Four measurements of  $(x, y)$  coordinate and cosine and sine of head angle ( $v$ )
- Three distinct dynamical regimes, namely, left turn, right turn and waggle



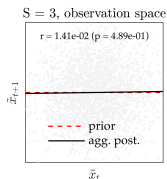
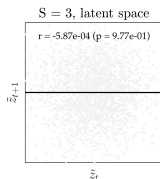
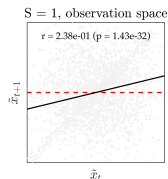
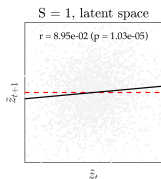
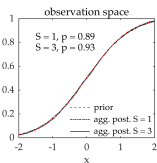
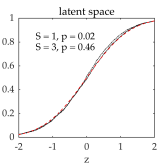
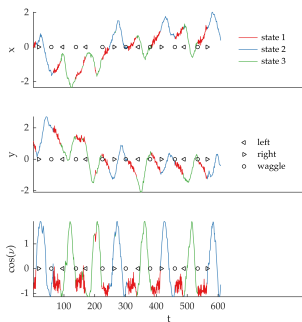
## 2. Honey Bee

- Four measurements of  $(x, y)$  coordinate and cosine and sine of head angle ( $v$ )
- Three distinct dynamical regimes, namely, left turn, right turn and waggle



## 2. Honey Bee

- Four measurements of  $(x, y)$  coordinate and cosine and sine of head angle ( $v$ )
- Three distinct dynamical regimes, namely, left turn, right turn and waggle

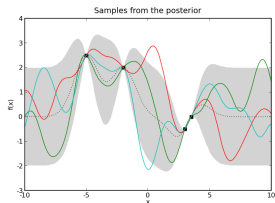
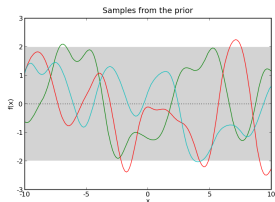


### 3. Gaussian Process Regression

$$\theta, \zeta, \tau \sim p(\theta) p(\zeta) p(\tau), \quad (13)$$

$$f(x) \sim \mathcal{GP}(m(x | \theta), \kappa(x, x' | \zeta)), \quad (14)$$

$$y_i \sim \mathcal{N}(f(x_i), \tau^{-1}) \quad \forall i = 1, \dots, n, \quad (15)$$



[SOURCE]

### 3. Gaussian Process Regression

$$\theta, \zeta, \tau \sim p(\theta) p(\zeta) p(\tau), \quad (13)$$

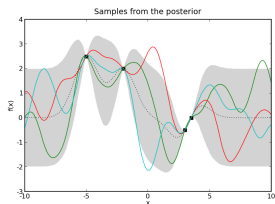
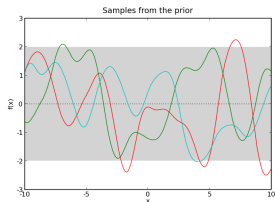
$$f(x) \sim \mathcal{GP}(m(x | \theta), \kappa(x, x' | \zeta)), \quad (14)$$

$$y_i \sim \mathcal{N}(f(x_i), \tau^{-1}) \quad \forall i = 1, \dots, n, \quad (15)$$

$$\mathbf{m} = (m(\mathbf{x}_1 | \theta), \dots, m(\mathbf{x}_n | \theta))^\top \quad (16)$$

$$\mathbf{K}_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j | \zeta) + \tau^{-1} \delta(\mathbf{x}_i, \mathbf{x}_j) \quad (17)$$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{m}, \mathbf{K}) \quad (18)$$



[SOURCE]

### 3. Gaussian Process Regression

$$\vartheta, \zeta, \tau \sim p(\vartheta) p(\zeta) p(\tau), \quad (13)$$

$$f(x) \sim \mathcal{GP}(m(x | \vartheta), \kappa(x, x' | \zeta)), \quad (14)$$

$$y_i \sim \mathcal{N}(f(x_i), \tau^{-1}) \quad \forall i = 1, \dots, n, \quad (15)$$

$$\mathbf{m} = (m(\mathbf{x}_1 | \vartheta), \dots, m(\mathbf{x}_n | \vartheta))^T \quad (16)$$

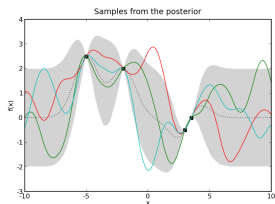
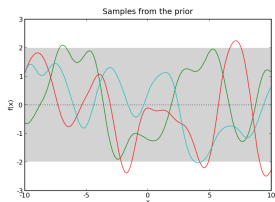
$$\mathbf{K}_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j | \zeta) + \tau^{-1} \delta(\mathbf{x}_i, \mathbf{x}_j) \quad (17)$$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{m}, \mathbf{K}) \quad (18)$$

$$\mathbf{K} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \quad (19)$$

$$\mathbf{c} = \mathbf{U}^T (\mathbf{y} - \mathbf{m}) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}) \quad (20)$$

$$\mathbf{z} = \mathbf{\Lambda}^{-1/2} \mathbf{U}^T (\mathbf{y} - \mathbf{m}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (21)$$



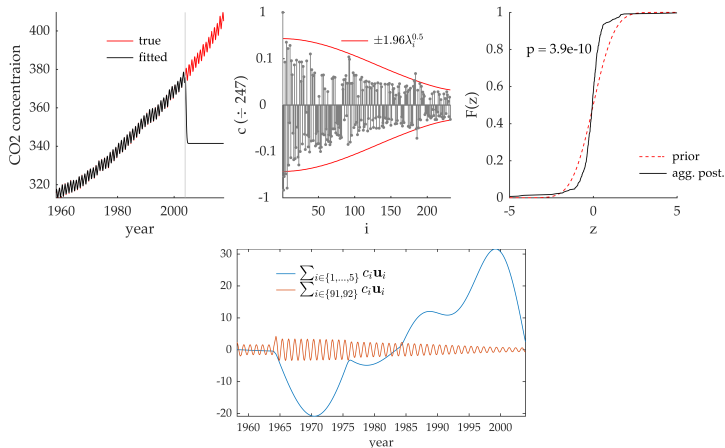
[SOURCE]

Given  $\nu^*, \zeta^*, \tau^*$

$$\{z_1^*, \dots, z_n^*\} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (22)$$

### 3. CO2 Emmision

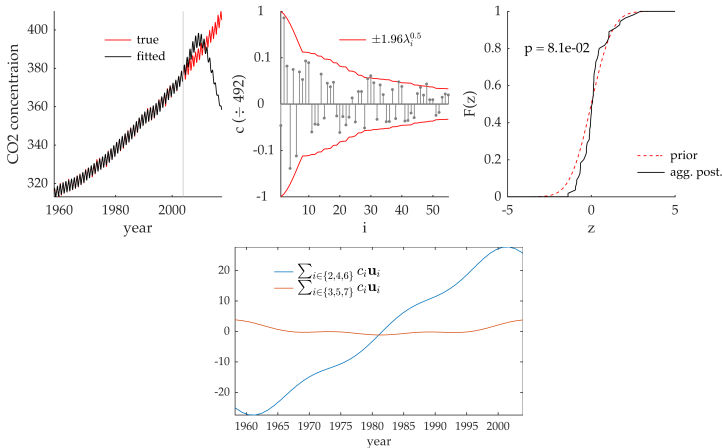
$$\kappa_{\text{se}}(x, x' | \zeta) = \sigma_f^2 \exp\left(-\frac{(x - x')^2}{2l^2}\right) \quad (23)$$





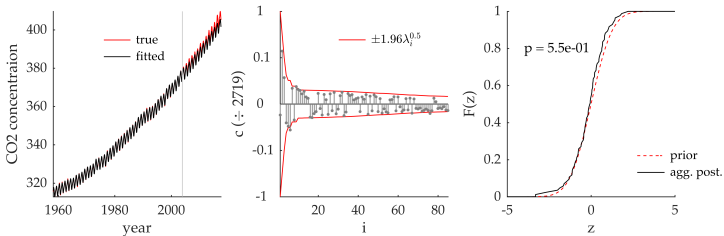
### 3. CO<sub>2</sub> Emmission

$$\kappa_{\text{pe}}(x, x' | \zeta) = \sigma_f^2 \exp\left(-\frac{2 \sin^2(\pi(x - x')/p)}{l_p^2}\right) \exp\left(-\frac{(x - x')^2}{2l_d^2}\right) \quad (23)$$



### 3. CO2 Emmision

$$\kappa_{se}(x, x' | \zeta_s) + \kappa_{se}(x, x' | \zeta_I) + \kappa_{pe}(x, x' | \zeta) \quad (23)$$



- Model criticism goodness-of-fit and graphical illustration for understanding a limitations of the model with the hope that a better model can be found
- The term *model criticism* is preferred over *model validation* and *model checking* since the former has a more active tone of looking to discover problems, while the latter may seem a more passive activity that does not expect to uncover any problems O'Hagan [2003, p423].
- Model criticism is contrasted with *model comparison* in that model criticism assesses a single model, while model comparison deals with at least two models to decide which model is a better fit.
- Model comparison can be applied to compare the original and the extended model after model criticism and extension [O'Hagan, 2003, p. 2].
- Aggregated Posterior Check complements Posterior Predictive Check by criticising the latent space rather than the observation space, and has been used in the literature in different forms Meulders et al. [1998], Buccigrossi and Simoncelli [1999], O'Hagan [2003], Tang et al. [2012]

- Model criticism goodness-of-fit and graphical illustration for understanding a limitations of the model with the hope that a better model can be found
- The term *model criticism* is preferred over *model validation* and *model checking* since the former has a more active tone of looking to discover problems, while the latter may seem a more passive activity that does not expect to uncover any problems O'Hagan [2003, p423].
- Model criticism is contrasted with *model comparison* in that model criticism assesses a single model, while model comparison deals with at least two models to decide which model is a better fit.
- Model comparison can be applied to compare the original and the extended model after model criticism and extension [O'Hagan, 2003, p. 2].
- Aggregated Posterior Check complements Posterior Predictive Check by criticising the latent space rather than the observation space, and has been used in the literature in different forms Meulders et al. [1998], Buccigrossi and Simoncelli [1999], O'Hagan [2003], Tang et al. [2012]

Sohan Seth, Iain Murray, and Christopher K. I. Williams.  
Model Criticism in Latent Space.  
*Bayesian Analysis*, 14(3):703–725, 2019.  
<https://projecteuclid.org/euclid.ba/1560240024>.

## References

- G. E. P. Box and N. R. Draper. *Empirical Model-Building and Response Surfaces*. Wiley, 1987.
- G. E.P Box. Sampling and Bayes' Inference in Scientific Modelling and Robustness. *Journal of the Royal Statistical Society*, 143(4):383–430, 1980.
- R. P. Buccirossi and E. P. Simoncelli. Image Compression via Joint Statistical Characterization in the Wavelet Domain. *IEEE Transactions on Signal Processing.*, 8(12):1688–1701, 1999.
- A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*. Chapman and Hall, London, 2004. Second edition.
- V. E. Johnson. Bayesian model assessment using pivotal quantities. *Bayesian Anal.*, 2(4):719–733, 12 2007. doi: 10.1214/07-BA229.
- M. Meulders, A. Gelman, I. Van Mechelen, and P. De Boeck. Generalizing the Probability Matrix Decomposition Model: an Example of Bayesian Model Checking and Model Expansion. In Hox, J. J. and de Leeuw, E. D., editor, *Assumptions, Robustness and Estimation Methods in Multivariate Modeling*. TT-Publikaties, Amsterdam, 1998.
- A. O'Hagan. HSSS Model Criticism. In P. J. Green, N. L. Hjort, and S. Richardson, editors, *Highly Structured Stochastic Systems*, pages 422–444. Oxford University Press, 2003.
- D. B. Rubin. Bayesianly Justifiable and Relevant Frequency Calculations for the Applied Statistician. *Annals of Statistics*, 12:1151–1172, 1984.
- Y. Tang, R. Salakhutdinov, and G. E. Hinton. Deep Mixtures of Factor Analysers. In *Proceedings of the 29th International Conference on Machine Learning*, 2012.
- D. Zoran and Y. Weiss. Natural Images, Gaussian Mixtures and Dead Leaves. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1736–1744. Curran Associates, Inc., 2012.

Thank you for your patience!