

# Modelagem Robusta de Observações Ausentes em Dados Sensoriados de Apiários

Daniel de Amaral da Silva<sup>1</sup>, Antonio Rafael Braga<sup>3</sup>,  
Juvêncio S. Nobre<sup>2</sup>, Danielo G. Gomes<sup>1</sup>

<sup>1</sup> Programa de Pós-Graduação em Engenharia de Teleinformática,  
Grupo de Redes de Computadores, Engenharia de Software e Sistemas (GREat)  
Centro de Tecnologia, Universidade Federal do Ceará,  
Fortaleza - CE, CEP 60.455-970.

<sup>2</sup>Departamento de Estatística e Matemática Aplicada (DEMA), Centro de Ciências,  
Universidade Federal do Ceará, Fortaleza - CE, CEP 60.440-900.

<sup>3</sup>Sistemas de Informação, Campus Quixadá,  
Universidade Federal do Ceará, Quixadá - CE, CEP 63.902-580.

{danielamaral}@alu.ufc.br

{rafaelbraga, juvencio, danielo}@ufc.br

**Abstract.** *Remote sensing of apiaries reduces the need for manual, unnecessary and invasive inspections of hives while providing beekeepers with advance warning of problems in loco. Temperature monitoring is one of several priorities when interested in diagnosing the health and well-being of a colony, since bees are strictly careful about maintaining the microclimate (climate inside the hive). However, eventual failures in the sensors can contaminate the data with anomalous values or simply cause "holes" in the information. Here we propose a statistical model capable of dealing with outliers local problems or missing data in a multivariate time series of temperature data from a grid with 36 sensors installed in a bee hive, with a metric  $R^2$  greater than 87% on 35 of the 36 sensors.*

**Resumo.** *O sensoriamento remoto de apiários reduz a necessidade de inspeções manuais, desnecessárias e invasivas nas colmeias ao mesmo tempo que propicia aos apicultores, com antecedência, alertas de problemas in loco. O monitoramento da temperatura é uma das diversas prioridades quando se tem interesse em diagnosticar o estado de saúde e bem-estar de uma colônia, posto que as abelhas têm um rigoroso cuidado quanto à sua manutenção do microclima (clima interno à colmeia). Entretanto, eventuais falhas nos sensores podem contaminar os dados com valores anômalos ou simplesmente causar "buracos" nas informações. Neste artigo, propomos um modelo estatístico capaz de contornar problemas locais de outliers ou de dados ausentes em uma série temporal multivariada de dados de temperatura provenientes de uma grade com 36 sensores instalados em uma colmeia de abelhas, com métrica  $R^2$  superior a 87% em 35 dos 36 sensores.*

## 1. Introdução

Muitos insetos controlam seu ambiente restringindo algumas condições, como umidade relativa e temperatura [Jones et al. 2005]. Por exemplo, as abelhas precisam manter o ambiente das crias com temperatura entre 32°C e 36 °C para seu desenvolvimento padrão [Becher et al. 2010]. No entanto, as colônias gastam muita energia para manter a temperatura em um intervalo controlado [Jones et al. 2004], o que sugere a importância de se monitorar a temperatura em toda a colmeia na busca de um bom diagnóstico da saúde e do bem-estar da colônia.

Normalmente, as empresas que fornecem sensores para colmeias<sup>1</sup> oferecem um sensor de temperatura por unidade monitorada. Apesar dos dados de temperatura capturados serem úteis, é possível haver erros de medição devido a problemas no sensor e por somente avaliar a temperatura em apenas um local da colmeia, geralmente em pontos centrais. O sensoriamento baseado em grade de sensores pode corrigir estes problemas e fornecer padrões que um único sensor não conseguiria reconhecer [Becher and Moritz 2009]. Entretanto, a maioria dos sistemas de sensoriamento sofrem falhas sistemáticas ou acidentais, como quebra de sensores, erros de comunicação, dentre outros. Tais problemas podem gerar valores anômalos (*outliers*) ou simplesmente não coletar os valores sensorizados em uma dada janela de tempo, gerando valores faltantes.

Tendo em vista essa problemática, nossa proposta utiliza uma série temporal multivariada de dados de temperatura, provenientes de uma grade com 36 sensores, instalados em uma colmeia. Nossa principal contribuição visa ao fornecimento de um modelo capaz de contornar problemas locais de *outliers* ou valores faltantes por meio da correlação global entre os sensores na grade, estimando janelas de temperatura com base no comportamento de outros sensores. Para essa tarefa foi escolhido o modelo de processo gaussiano com múltiplas saídas, composto de 36 processos gaussianos, para modelar o comportamento individual de cada sensor, e a modelagem intrínseca de correção regionalização, do inglês, *Intrinsic Model of Correogionalization* (ICM) [Bonilla et al. 2008] para modelar a correlação entre os sensores. A escolha desses modelos foi guiada pela necessidade de modelos estatísticos flexíveis para monitoramento e previsão de séries temporais multivariadas, que possam reconhecer padrões implícitos na série temporal multivariada de temperatura intra e entre sensores.

## 2. Processos Gaussianos Variacionais de Múltiplas Saídas

Processos gaussianos variacionais são alternativas não paramétricas em que consideram-se variáveis latentes de um processo gaussiano usual na família variacional. Vamos revisar o que são modelos variacionais, processos gaussianos, processos gaussianos variacionais, e por fim, apresentar a solução de correção regionalização entre processos gaussianos.

### 2.1. Modelos Variacionais

Seja  $p(\mathbf{z}|\mathbf{x})$  a distribuição a posteriori sobre um conjunto de  $d$  variáveis latentes<sup>2</sup>  $\mathbf{z} = \{z_1, z_2, \dots, z_d\}$ , dado uma amostra  $\mathbf{x}$ . A ideia da abordagem de inferência variacional é aproximar a distribuição  $p(\mathbf{z}|\mathbf{x})$ , inicialmente desconhecida, através de uma família conhecida de distribuições  $q(\mathbf{z}; \boldsymbol{\lambda})$  parametrizada por  $\boldsymbol{\lambda}$ , pela minimização da divergência

---

<sup>1</sup><https://ColonyMonitoring.com>

<sup>2</sup>Variáveis não diretamente observadas e sim inferidas a partir de outras variáveis observadas.

$\text{KL}(q(\mathbf{z}; \boldsymbol{\lambda}) || p(\mathbf{z}|\mathbf{x}))$ . Essa abordagem equivale a maximizar o limite inferior da evidência, em inglês, *Evidence Lower Bound* (ELBO), dado por:

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\lambda})}[\log p(\mathbf{z}|\mathbf{x})] - \text{KL}(q(\mathbf{z}; \boldsymbol{\lambda}) || p(\mathbf{z}|\mathbf{x})). \quad (1)$$

A ELBO (1) depende de  $q(\mathbf{z})$ , chamada de distribuição variacional. Dessa forma, a escolha da forma e parametrização da distribuição tem relevante impacto nos resultados da otimização. Um fator atrativo dos modelos variacionais, é que este converte o problema de inferência em um problema de otimização.

## 2.2. Processos Gaussianos

Seja um conjunto de dados  $\mathcal{D} = \{(\mathbf{x}_n^\top, \mathbf{y}_n^\top)\}_{n=1}^m$  de  $m$  pares de vetores de entradas e saídas, em que cada entrada  $\mathbf{x}_n$  é composta de  $c$  covariáveis pareadas com uma saída multivariada  $\mathbf{y}_n \in \mathbb{R}^d$ . Nosso principal objetivo é mapear uma função sobre as entradas  $\mathbf{x}_n$  de modo que  $\mathbf{y}_n = f(\mathbf{x}_n)$  [Bishop 2006]. Considere uma função  $f : \mathbb{R}^c \rightarrow \mathbb{R}^d$ , inicialmente desconhecida, decomposta da forma  $f = \{f_1, f_2, \dots, f_d\}$ , em que cada função  $f_i : \mathbb{R}^c \rightarrow \mathbb{R}, \forall i \in \{1, \dots, d\}$ . Um Processo gaussiano (GP) para regressão modela a forma funcional  $f$  por uma distribuição a priori, dada por:

$$p(f) = \prod_{i=1}^d \mathcal{GP}(f_i; \boldsymbol{\mu}_x, \mathbf{K}_{xx}),$$

em que  $\boldsymbol{\mu}_x$  é o vetor de médias, usualmente um vetor de zeros  $\boldsymbol{\mu} = \mathbf{0}$ , e  $\mathbf{K}_{xx'}$  denota a função de autocovariância ou *kernel*  $k(x, x')$  sobre os pares de entradas  $x$  e  $x'$ . Nesse trabalho, utilizamos uma função de kernel RBF (*Radial Basis Function*), dada por:

$$k(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2} \sum_{j=1}^d w_d^2(x_j, x'_j)\right). \quad (2)$$

A parametrização em (2) também é chamada de *kernel* ARD (*Automatic Relevance Determination*) que penaliza cada dimensão  $d$  [Murphy 2013]. Em alguns casos, no procedimento de inferência, algumas dimensões podem ter pesos  $w_d$  nulos, levando à redução de dimensionalidade automática.

A distribuição condicional do processo gaussiano dado um conjunto de pares de entrada e saída é dada por:

$$p(f|\mathcal{D}) = \prod_{i=1}^d \mathcal{GP}(f_i; \mathbf{K}_{\xi x} \mathbf{K}_{xx}^{-1} \mathbf{y}_i, \mathbf{K}_{\xi\xi} - \mathbf{K}_{\xi x} \mathbf{K}_{xx}^{-1} \mathbf{K}_{\xi x}^\top), \quad (3)$$

em que  $\mathbf{K}_{\xi x}$  representa a função de auto-covariância/*kernel* de uma entrada  $\xi$  sobre todas as entradas de  $\mathcal{D}$ , e  $\mathbf{y}_i$  representa a saída da  $i$ -ésima dimensão.

## 2.3. Processos Gaussianos Variacionais

Um processo gaussiano variacional é um modelo variacional bayesiano não-paramétrico que admite diferentes e arbitrarias estruturas de aproximação da distribuição preditiva (a

posteriori) dos dados [Tran et al. 2016]. A ideia básica por trás de um processo gaussiano variacional, em inglês, *Variational Gaussian Process* (VGP) é gerar as variáveis latentes  $\mathbf{z}$  através de entradas latentes, avaliando-as em um conjunto de funções candidatas e as usando como parâmetros para uma distribuição de *mean-field*.

O processo de geração das variáveis latentes  $\mathbf{z}$  [Tran et al. 2016], é dado pelas etapas:

1. Gere uma entrada latente  $\xi \in \mathbb{R}^c : \xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ;
2. Avalie a amostra obtida em  $f$ , condicionada em  $\mathcal{D} : f \sim \prod_{i=1}^d \mathcal{GP}(\mathbf{0}, \mathbf{K}_{\xi\xi})|\mathcal{D}$ ;
3. Aproxime as amostras da posteriori  $\mathbf{z} \in \text{suporte}(p) : \mathbf{z} = \{z_1, z_2, \dots, z_d\} \sim \prod_{i=1}^d q(f_i(\xi))$ .

Dessa forma a distribuição variacional  $q_{VGP}$ , parametrizada por  $\theta$  e pelo conjunto de dados variacional, é dada por:

$$q_{VGP}(\mathbf{z}; \theta, \mathcal{D}) = \iint \left[ \prod_{i=1}^d q(z_i | f_i(\xi)) \right] \left[ \prod_{i=1}^d \mathcal{GP}(f_i; \mathbf{0}, \mathbf{K}_{\xi\xi}) | \mathcal{D} \right] \mathcal{N}(\xi; \mathbf{0}, \mathbf{I}) df d\xi.$$

Diferentemente do modelo de processo gaussiano usual, o VGP pode capturar a correlação entre as variáveis latentes. Esse comportamento é devido à avaliação de todos os  $d$  processos gaussianos independentes na mesma entrada latente  $\xi$ . Dessa forma, induzindo correlação entre as saídas, os parâmetros da distribuição de *mean-field* e consequentemente entre variáveis latentes.

## 2.4. Corregionalização

Nesse trabalho, modelamos  $f$  como um processo gaussiano corregionalizado que assume um matriz de kernel da forma:

$$\text{cov}(f_i(\mathbf{x}), f_j(\mathbf{x}')) = k(\mathbf{x}, \mathbf{x}') \cdot b_{ij},$$

em que,  $b_{ij}$  é a entrada  $(i, j)$  de uma matriz  $\mathbf{B}$  positiva definida de dimensão  $P \times P$ , essa forma é conhecida como modelo de corregionalização intrínseca (ICM) [Bonilla et al. 2008]. Dessa forma a covariância entre  $f_i(\mathbf{x})$  e  $f_j(\mathbf{x}')$ , é dada pela relação de covariância entre as entradas  $\mathbf{x}$  e  $\mathbf{x}'$ , e a covariância entre as funções latentes  $f_i$  e  $f_j$ .

A inferência para o modelo de corregionalização pode ser realizada pelos procedimentos padrões ou variacionais vistos anteriormente, por exemplo, a média da distribuição a posteriori de um novo ponto  $x^*$  para a dimensão  $d$ , pode ser adaptada utilizando o método de inferência usual [Bishop 2006], dada por:

$$\bar{f}_d(x^*) = (b_d \otimes \mathbf{k}^*)^\top \Sigma^{-1} \mathbf{y}, \quad \Sigma^{-1} = \mathbf{B} \otimes \mathbf{K}_{xx} + \mathbf{D} \otimes \mathbf{I},$$

em que  $\otimes$  denota o produto de Kronecker,  $b_d$  seleciona a  $d$ -ésima coluna de  $\mathbf{B}$ ,  $\mathbf{k}^*$  é o vetor de covariâncias entre o novo padrão e os dados já observados,  $\mathbf{K}_{xx}$  é a matriz de covariância entre todos os pontos observados (exceto o novo ponto), e  $\mathbf{D}$  é uma matriz diagonal  $d \times d$ , que contém os valores dos ruídos  $\sigma_d^2$  de cada dimensão na coordenada  $(d, d)$ .

### 3. Material e Métodos

#### 3.1. Ambiente Experimental

Os experimentos foram executados no ambiente Python de uma máquina virtual da EC2 da *Amazon Web Services (AWS)*<sup>3</sup>, com 32 GB RAM, processador Xeon 2.3 GHz e HDD de 21 GB. O *framework* de modelagem escolhido foi o GPflow<sup>4</sup>.

#### 3.2. Conjunto de Dados e Delineamento Estratégico

Neste artigo, os dados foram obtidos de uma grade de 36 sensores de temperatura posicionada no topo de uma colmeia [Linton et al. 2020a], organizados como uma matriz de  $4 \times 9$ , monitorada no período de aproximadamente um ano. O sensoriamento ocorreu durante o período de 10/11/2019 a 30/12/2020 em Washington (EUA). Nas Figuras 1a e 1b são mostrados o topo da colmeia com e sem a grade de sensores.



(a) Visão da parte superior da colmeia



(b) Visão da grade de sensores ( $4 \times 9$ ) instalada

C0	C1	C2	CD0	CD1	CD2	D0	D1	D2
C3	C4	C5	CD3	CD4	CD5	D3	D4	D5
A0	A1	A2	AB0	AB1	AB2	B0	B1	B2
A3	A4	A5	AB3	AB4	AB5	B3	B4	B5

(c) Diagrama de localização e *tags* dos 36 sensores.

Figura 1. (a) *Cluster* de inverno; (b) topo da colmeia com a grade de sensores de temperatura; (c) disposição e *tags* dos sensores na grade. Fonte: [Linton et al. 2020a]

Cada sensor possui uma identificação (*tag*). Na Figura 1c são ilustradas as disposições dos 36 sensores com suas respectivas *tags* de identificação e localizações na

<sup>3</sup><https://aws.amazon.com/pt/ec2/>

<sup>4</sup><https://github.com/GPflow/GPflow>

grade, no mesmo posicionamento da Figura 1b. Várias combinações de sensores foram testadas em diferentes intervalos de medição, e.g. 18, 48 e 60 minutos, e apenas metade dos 36 sensores foram monitorados em períodos regulares de tempo. Devido à alta incidência de valores faltantes no início do experimento, aproximadamente 41% somente no mês de março/2020, por falhas no armazenamento de dados no cartão SD e problemas na bateria de alguns sensores [Linton et al. 2020a, Linton et al. 2020b], optamos por utilizar somente os meses subsequentes ao mês de março, já que grande parte dos sensores teve problemas de captura neste período provocando grandes períodos de observações ausentes (buracos). Assim, para utilizarmos todos os sensores e em uma escala de tempo única, optamos por efetuar um *down sample* diário nos dados, i.e, o valor médio dos dados em um dia de cada sensor.

A Figura 2 ilustra as medições diárias dos 36 sensores no novo período mencionado, i.e. 27/03/2020 a 30/12/202, em que a variável Dia é dada pela contagem de dias passados da data de referência, que no caso é 27/03/2020. Contudo, entre as datas de estudo, vários períodos ainda contêm grandes *gaps* de dados faltantes, como por exemplo, o período de aproximadamente um mês entre os dias 70 e 100, e nos últimos meses do ano entre os dias 240 e 260. Entretanto, apesar de termos janelas consideráveis de valores faltantes, o conjunto de dados possui dados de sensores disponíveis nesses períodos, i.e, em sua maior parte a não coleta de dados ocorreu em um subconjunto dos sensores. Ademais, nota-se altas correlações entre sensores, o que é esperado já que temos uma grade com sensores medindo a mesma variável (temperatura).

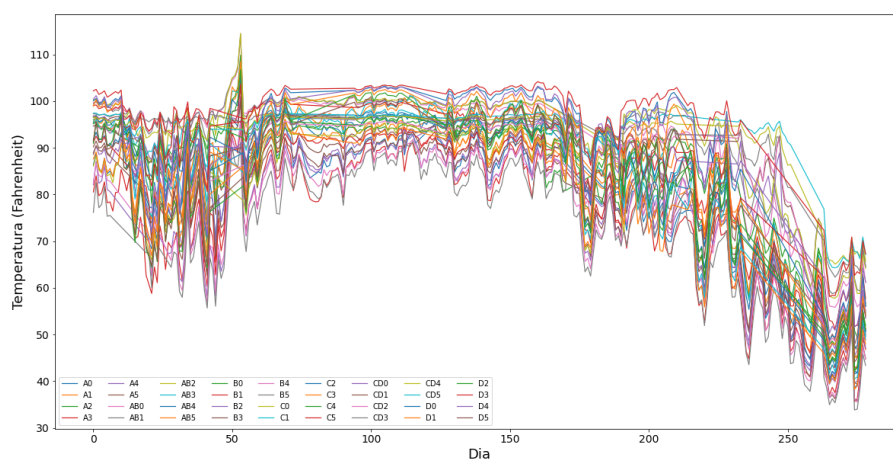


Figura 2. Medições de temperatura da grade de sensores ilustrada na Figura 1.

### 3.3. Métricas de Avaliação

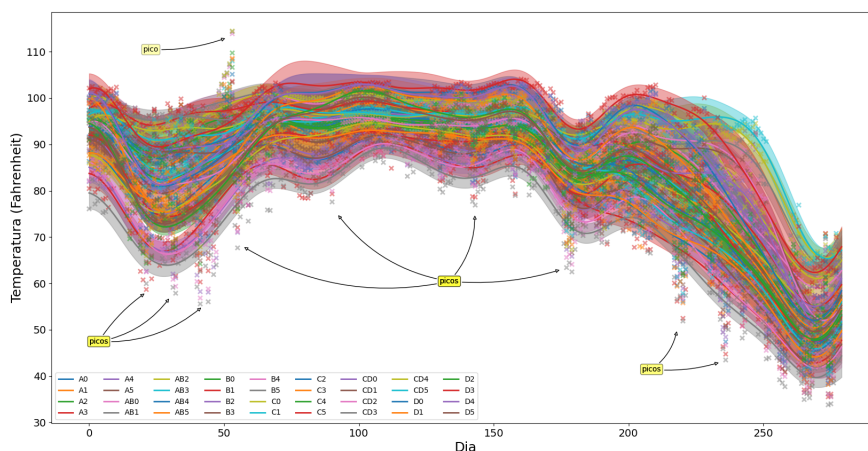
As métricas utilizadas no experimento foram:  $R^2$ , RMSE e MAE. Estas medidas de performance são comuns em experimentos de regressão univariada [Botchkarev 2019], e serão empregadas para um maior detalhamento da modelagem de cada sensor. Na Tabela 1 são mostradas as equações de obtenção das métricas utilizadas.

**Tabela 1. Equação e Descrição das Métricas de Desempenho do Modelo.**

Métrica	Equação	Descrição
$R^2$	$\frac{\sum_i (\hat{y}_i - \bar{y})^2}{\sum_i (y_i - \bar{y})^2}$	O coeficiente r-quadrado explica o grau em que as variáveis independentes explicam a variação da variável dependente. Por exemplo, um valor de $R^2$ de 0,90 indica que as variáveis de entrada explicam 90% da variabilidade da variável de saída
RMSE	$\sqrt{\frac{1}{n} \sum_i (y_i - \hat{y}_i)^2}$	Os valores da métrica RMSE são interpretados na mesma escala da variável dependente, e por ser uma métrica que avalia o erro, quanto menor seu valor melhor é o modelo na tarefa de predição
MAE	$\frac{1}{n} \sum_i  y_i - \hat{y}_i $	A métrica MAE também possui semelhante interpretação comparada a métrica RMSE, diferindo somente pela forma de obtenção dos desvios

#### 4. Resultados e Discussão

A série temporal multivariada de temperatura proveniente da grade de sensores [Linton et al. 2020a] foi modelada utilizando um processo gaussiano corregeionalizado, i.e, múltiplos modelos de processos gaussianos correlacionados. Para termos uma visão geral do comportamento do modelo aplicado aos dados sensorizados, plotamos a média de cada distribuição dos sensores e seus respectivos intervalos de credibilidade ( $\pm 2\sigma_{pred}$ ). A Figura 3 ilustra o ajuste do modelo corregeionalizado aos dados da grade de sensores.



**Figura 3. Distribuição preditiva do modelo de processos gaussianos corregeionalizados, com 200 iterações.**

Através da análise da Figura 3 podemos perceber picos que saltam as bandas de credibilidade do modelo, esses pontos anômalos podem ser ocasionados devido ao processo de obtenção dos dados diários (média horária), falha no sensor ou pela ação de uma fonte externa de variação esporádica, como por exemplo a inspeção *in loco* do apicultor. Como possível solução desses pontos, está a remoção dos mesmos se ultrapassarem o intervalo de credibilidade estabelecido de  $\pm 2\sigma_{pred}$  (para uma abordagem mais conservadora pode ser considerado  $\pm 3\sigma_{pred}$ ). Apesar de vários problemas detectados nos dados, o modelo conseguiu se ajustar satisfatoriamente aos dados, mesmo com poucas iterações. Na Tabela 2, são mostradas as métricas de desempenho na modelagem proposta.

**Tabela 2. Métricas do modelo de rocessos Gaussianos Corregionalizados aplicado aos dados da grade de sensores.**

Sensor	$R^2$	RMSE	MAE	Sensor	$R^2$	RMSE	MAE
A0	0.9033	4.5070	3.3867	C0	0.6473	4.7169	3.4465
A1	0.9315	3.8843	2.6604	C1	0.8846	4.8949	3.5782
A2	0.9394	3.3312	2.2592	C2	0.9104	4.0769	2.8963
A3	0.8845	4.9384	3.8158	C3	0.8931	4.4642	3.4011
A4	0.9073	4.4860	3.3454	C4	0.9181	3.8184	2.6784
A5	0.9212	4.1205	2.9887	C5	0.9470	2.6886	1.8049
AB0	0.9554	2.6280	1.6516	CD0	0.9312	0.2137	0.1382
AB1	0.9525	2.5146	1.6032	CD1	0.9325	0.1959	0.1232
AB2	0.9392	2.8471	1.8099	CD2	0.9180	0.2213	0.1463
AB3	0.9294	3.9037	2.7701	CD3	0.9712	0.0965	0.0611
AB4	0.9341	3.8824	2.6789	CD4	0.9754	0.0818	0.0544
AB5	0.9198	4.1692	2.9057	CD5	0.9543	0.1139	0.0766
B0	0.9157	3.5640	2.3126	D0	0.9121	4.0531	2.8148
B1	0.9015	4.2815	2.8616	D1	0.8925	4.6877	3.3749
B2	0.8925	4.6416	3.3364	D2	0.8797	5.0579	3.8093
B3	0.9104	4.4909	3.1697	D3	0.9271	2.8888	1.8656
B4	0.8974	4.7246	3.4007	D4	0.9042	3.9821	2.7928
B5	0.8918	4.8431	3.5891	D5	0.8880	4.5697	3.4238

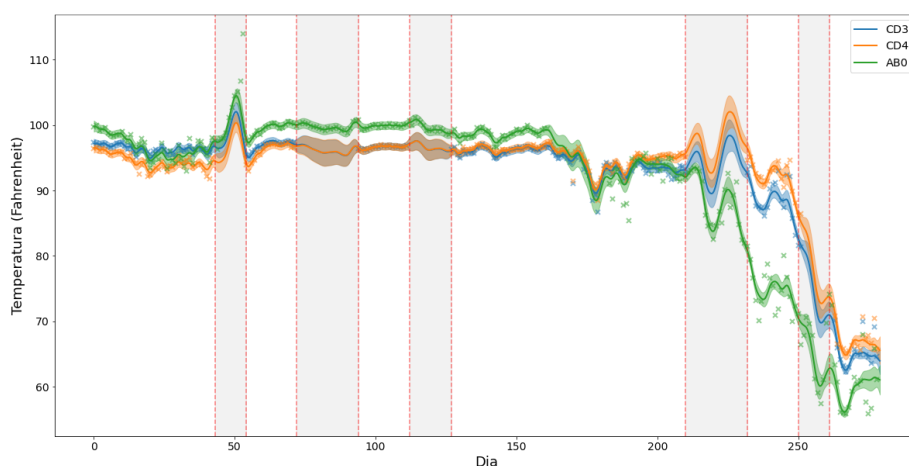
Notamos na Tabela 2 a métrica de  $R^2$  acima de 87% para todos os sensores com exceção de C0. As métricas de RMSE e MAE indicam que o modelo erra em média de 3 a 4 graus de temperatura do real, com exceção dos sensores centrais (AB e CD) em que o modelo proposto atingiu resultados muito bons. Este resultado preditivo para os sensores não-centrais sugere uma necessidade de mais iterações e/ou uma escolha de processos que tolerem o comportamento de alta variabilidade dos dados (e.g. Processos t-Student).

Sobre os valores faltantes, avaliamos o perfil da distribuição a posteriori (preditiva) dos sensores centrais CD3-CD5 e AB0-AB5. A Figura 4 mostra a distribuição preditiva dos sensores centrais CD3, CD4 e AB0, na qual pode-se notar alguns picos (rajadas) nos intervalos de credibilidade, com intervalos de cobertura médios de aproximadamente  $5^\circ F$ .

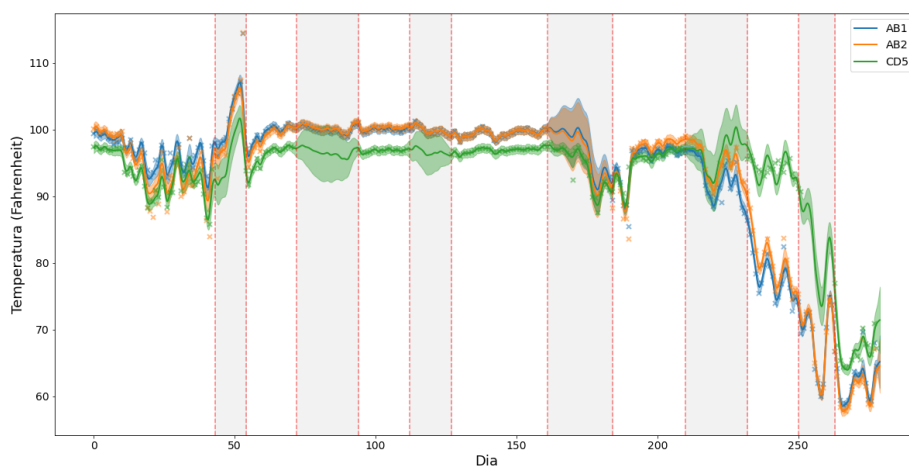
A ocorrência destas rajadas pode ser explicada pela ausência de informação (ver áreas em cinza) em dois dos três sensores. O comportamento modelado, por exemplo, entre os dias 70 e 90, mostra a capacidade do modelo inferir o comportamento faltante dos sensores com base nos dados observados no sensor AB0 (verde). Comportamento semelhante pode ser notado nos dias finais (210 a 230), nos quais novamente os sensores CD3 e CD4 são modelados pelas correlações com o sensor AB0.

Podemos perceber na Figura 5 novamente picos nos intervalos de credibilidade, com intervalos de cobertura médios de aproximadamente  $7^\circ F$ . Isto é devido à ausência de informação (ver áreas em cinza) em um dos três sensores. O comportamento modelado, por exemplo, entre os dias 70 e 90, mostra a capacidade do modelo em inferir o comportamento faltante dos sensores com base nos dados observados dos sensores AB1 e AB2 (azul e laranja). Comportamento semelhante pode ser notado nos dias finais (210 a 260), nos quais o modelo infere cerca de 30 dias do comportamento faltante do sensor CD5.





**Figura 4. Distribuição preditiva dos sensores centrais CD3, CD4 e AB0, com 5000 iterações. Em destaque, períodos com observações ausentes.**



**Figura 5. Distribuição preditiva dos sensores centrais AB1, AB2 e CD5 (5000 iterações).**

## 5. Conclusão

Nesse artigo, um modelo de processos gaussianos correionalizados foi ajustado em uma série temporal diária de temperatura multivariada de uma colmeia de abelhas através de uma grade de 36 sensores.

Como principal contribuição deste artigo, sugerimos a aplicação de modelos de processos gaussianos correionalizados como uma alternativa flexível e robusta para modelagem de observações faltantes em dados sensoriados no contexto de apicultura de precisão, com métricas de  $R^2$  superiores a 87% em 35 dos 36 sensores. Apesar de as métricas apontarem erros de 3 a 4 graus de temperatura, podemos concluir que os intervalos de credibilidade do modelo podem ser utilizados como limiares de tolerância para medições e que modelos de processos gaussianos correionalizados nos dão ideia da incerteza a respeito de uma predição pontual através da distribuição posteriori ou preditiva.

Na qualidade de possíveis perspectivas de trabalhos futuros, podemos sugerir (i) o uso de processos gaussianos esparsos para reduzir a complexidade computacional; (ii) o

uso de funções de média diferentes de 0 para modelar a tendência dos dados e solucionar a questão da restrição de média nula do processo gaussiano padrão e (iii) adição de variáveis climáticas provenientes de estações meteorológicas.

## Agradecimentos

O presente artigo foi realizado com apoio da Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES) - Código de Financiamento 001. Danielo G. Gomes agradece o suporte financeiro do Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) através dos processos 432585/2016-8 e 310317/2019-3.

## Referências

- Becher, M. A., Hildenbrandt, H., Hemelrijk, C. K., and Moritz, R. F. (2010). Brood temperature, task division and colony survival in honeybees: A model. *Ecological Modelling*, 221(5):769–776.
- Becher, M. A. and Moritz, R. F. (2009). A new device for continuous temperature measurement in brood cells of honeybees (*apis mellifera*). *Apidologie*, 40(5):577–584.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg.
- Bonilla, E. V., Chai, K., and Williams, C. (2008). Multi-task gaussian process prediction. In Platt, J., Koller, D., Singer, Y., and Roweis, S., editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc.
- Botchkarev, A. (2019). A new typology design of performance metrics to measure errors in machine learning regression algorithms. *Interdisciplinary Journal of Information, Knowledge, and Management*, 14:045–076.
- Jones, J. C., Helliwell, P., Beekman, M., Maleszka, R., and Oldroyd, B. (2005). The effects of rearing temperature on developmental stability and learning and memory in the honey bee, *apis mellifera*. *Journal of Comparative Physiology A*, 191:1121–1129.
- Jones, J. C., Myerscough, M. R., Graham, S., and Oldroyd, B. P. (2004). Honey bee nest thermoregulation: Diversity promotes stability. *Science*, 305(5682):402–404.
- Linton, F., Stumme, A., Padula, B., Ifshin, G., and Behrmann, G. (2020a). Monitoring honey bee colony activities with a temperature sensor grid. *American Bee Journal*. Parts 1 through 3.
- Linton, F., Stumme, A., Padula, B., Ifshin, G., and Behrmann, G. (2020b). Monitoring honey bee colony activities with a temperature sensor grid. Parts 2 through 3.
- Murphy, K. P. (2013). *Machine learning : a probabilistic perspective*. MIT Press, Cambridge, Mass. [u.a.].
- Tran, D., Ranganath, R., and Blei, D. M. (2016). The variational gaussian process.