



Universidade Federal da Bahia  
Instituto de Matemática

Programa de Pós-Graduação em Ciência da Computação

...

Ruivaldo Azevedo Lobão Neto

DISSERTAÇÃO DE MESTRADO

Salvador  
14 de Junho de 2019



RUIVALDO AZEVEDO LOBÃO NETO

...

Esta Dissertação de Mestrado foi apresentada ao Programa de Pós-Graduação em Ciência da Computação da Universidade Federal da Bahia, como requisito parcial para obtenção do grau de Mestre em Ciência da Computação.

Orientador: Ricardo Araújo Rios

Salvador  
14 de Junho de 2019



## RESUMO

A quantidade de informações produzidas por sistemas computacionais tem crescido de forma acentuada nas últimas décadas. Diante disso, pesquisadores passaram a aplicar técnicas de Aprendizado de Máquina para extrair informações úteis de grandes conjuntos de dados. Contudo, uma parcela expressiva dessas informações é produzida na forma de fluxos contínuos de dados, que são sequências constantes e potencialmente infinitas. Esses fluxos são, em sua maioria, não-estacionários, podendo sofrer variações na distribuição dos dados ou no contexto do processo gerador. Estas alterações são denominadas mudanças de conceito e podem impactar negativamente a acurácia da técnica de aprendizado aplicada. Para mitigar este problema, pesquisadores desenvolveram métodos especializados na detecção de mudanças de conceito. Entretanto, os métodos propostos apresentam limitações ao serem aplicados em cenários com fluxos contínuos, como a necessidade de rotulação por especialistas ou a incapacidade de atender às restrições de tempo de processamento e de uso dos recursos computacionais desses cenários. Este projeto de trabalho de mestrado propõe um novo método de detecção baseado em Redes de Função de Base Radial (RBF). As redes RBF realizam, em sua camada intermediária, um processo de ativação de neurônios que, implicitamente, produz agrupamentos das observações recebidas ao longo do tempo. O algoritmo proposto nesta pesquisa visa monitorar o centro ativo desse agrupamento e detecta mudanças de conceito quando este centro é alterado. O método proposto se diferencia dos trabalhos existentes na literatura por detectar mudanças em tempo de execução, de forma computacionalmente eficiente e independente de rótulos. Os resultados preliminares obtidos permitem demonstrar a viabilidade desta proposta de mestrado, apresentando um desempenho superior ou similar ao de algoritmos conhecidos como ADWIN, CUSUM e PageHinkley, em termos de acurácia média e tempo de processamento.

**Palavras-chave:** Aprendizado de Máquina, Fluxos Contínuos de Dados, Mudanças de Conceito, Redes de Função de Base Radial, Não supervisionado



## SUMÁRIO





## LISTA DE FIGURAS



## **LISTA DE TABELAS**



## LISTA DE ALGORITMOS

1	RBFDRIFTDETECTOR . . . . .	22
---	----------------------------	----



## **INTRODUÇÃO**

### **1.1 CONTEXTO E MOTIVAÇÃO**

Nos últimos anos, o volume de dados produzidos por sistemas computacionais tem crescido de forma acentuada. Esse crescimento foi favorecido por avanços tecnológicos recentes, como a pervasividade dos dispositivos móveis, a popularização das redes sociais e a expansão da internet das coisas (??). A dimensão desse aumento foi verificada por ??), os quais estimaram que, entre os anos de 2014 e 2020, a quantidade de informações produzidas anualmente irá aumentar de 4,4 zettabytes (trilhões de gigabytes) para 44 zettabytes.

Parte significativa dessas informações é produzida na forma de sequências ininterruptas e potencialmente infinitas (??). Na literatura, sequências com essas características são denominadas Fluxos Contínuos de Dados (FCDs) e estão presentes em diversos domínios de aplicação como, por exemplo, monitoramento do mercado financeiro (??), acompanhamento de tráfego rodoviário (??), gerenciamento de redes de telecomunicação (??), análise de sentimento em tempo real (??) e sistemas de prevenção e identificação de intrusos (??).

Para extrair informações úteis dessa grande quantidade de dados, pesquisadores têm aplicado técnicas da área de Aprendizado de Máquina (AM), a qual estuda algoritmos que melhoram seu desempenho conforme ganham experiência (??). Entretanto, as estratégias tradicionais de Aprendizado de Máquina têm aplicação limitada para contextos com fluxos contínuos de dados, pois nesses cenários os algoritmos devem atender a severas restrições de tempo de execução e de uso dos recursos computacionais (??).

Além dessas limitações, as técnicas de Aprendizado de Máquina, quando aplicadas em contextos com fluxos contínuos, também devem lidar com variações na distribuição dos dados ou no contexto do processo gerador. Essas alterações são denominadas Mudanças de Conceito (??) e a sua ocorrência pode impactar a acurácia da técnica de aprendizado aplicada.

Inicialmente, a atualização periódica de modelos obtidos a partir das técnicas de AM foi utilizada como estratégia para evitar a perda de acurácia causada por tais mudanças.

Contudo, esta solução é pouco sofisticada e computacionalmente custosa. Diante disso, pesquisadores propuseram técnicas de detecção de mudanças de conceito baseadas em monitoramento (??). Estes métodos identificam a ocorrência de mudanças com maior precisão, permitindo que o modelo de decisão seja atualizado somente quando necessário. Exemplos de algoritmos baseados nesta abordagem, incluem: DDM (??), EDDM (??), ADWIN (??), ECDD (??), PL (??), FCWM (??) e STEPDP (??).

Entretanto, as técnicas baseadas em monitoramento necessitam que o rótulo correto de cada exemplo esteja disponível. Em muitos cenários, o tempo ou o custo para obter esses rótulos é proibitivo (??). Consequentemente, foram desenvolvidos novos algoritmos independentes de rótulos. Nestes métodos, a detecção se baseia na identificação de exemplos que não se enquadram na estrutura conhecida dos dados (??). Essa análise é implementada com base em técnicas de agrupamento, detecção de *outliers* e medidas de dissimilaridade (??). Os seguintes algoritmos são exemplos desta abordagem: OLINDDA (??), MINAS (??), ECSMiner (??) e GC3 (??).

Todavia, segundo ??, as técnicas de detecção de mudanças de conceito presentes na literatura ainda apresentam limitações ao serem aplicadas em cenários com fluxos contínuos de dados. Os algoritmos dependentes de rótulo se tornam inviáveis, por causa do custo e do tempo necessário para obter os rótulos corretos. Enquanto as técnicas independentes têm dificuldade em atender as severas restrições de tempo de execução e de uso dos recursos computacionais desses cenários.

Visando resolver essas limitações, este projeto de mestrado discute uma abordagem baseada em Redes de Função de Base Radial (RBF) para detecção de mudanças de conceito em fluxos contínuos de dados.

## 1.2 HIPÓTESE E OBJETIVO

Com base nas observações citadas anteriormente, a seguinte hipótese foi formulada:

*“A aplicação de Redes de Função de Base Radial em fluxos contínuos de dados permite a detecção de mudanças de conceito em tempo de execução, de forma computacionalmente eficiente e independente de rótulos.”*

Assim, o objetivo deste trabalho de mestrado será a validação desta hipótese. Para atingir este objetivo, será desenvolvido um método para detecção de mudanças de conceito baseado em Redes de Função de Base Radial. A técnica proposta será validada através de comparações com o estado da arte. Os dados utilizados durante a validação serão divididos em dois conjuntos. Um conjunto formado por dados sintéticos, que permitirão uma análise detalhada da abordagem, uma vez que as características e os comportamentos dos fluxos serão conhecidos. O outro conjunto será composto por dados obtidos a partir de sistemas computacionais utilizados na indústria, visando apresentar uma aplicação prática para a solução proposta.

O restante deste projeto está organizado conforme a seguinte estrutura: O **Capítulo ??** apresenta uma revisão bibliográfica dos principais conceitos utilizados como, por exemplo, Fluxos Contínuos de Dados e Aprendizado de Máquina, Mudança de Conceito e Redes



de Função de Base Radial; No **Capítulo ??** o plano de pesquisa é detalhado, identificando a metodologia que será aplicada na pesquisa e o cronograma de atividades. Por fim, o **Capítulo ??** apresenta um conjunto de experimentos preliminares e a análise dos resultados obtidos.



## **REVISÃO BIBLIOGRÁFICA**

### **2.1 CONSIDERAÇÕES INICIAIS**

Este capítulo apresenta uma discussão geral sobre os principais conceitos utilizados neste projeto. Inicialmente, será abordada a relação entre fluxos contínuos de dados e técnicas de Aprendizado de Máquina. Em seguida, o fenômeno mudança de conceito e os métodos de detecção serão discutidos. Posteriormente, as Redes de Função de Base Radial serão detalhadas. Por fim, serão apresentados os trabalhos relacionados encontrados na literatura.

### **2.2 FLUXOS CONTÍNUOS DE DADOS E APRENDIZADO DE MÁQUINA**

Fluxos Contínuos de Dados (FCDs) podem ser definidos como sequências ininterruptas e potencialmente infinitas de eventos (??). Nestes fluxos, os eventos ocorrem em alta frequência, sendo necessário processá-los em tempo real. Além disso, por serem de tamanho potencialmente ilimitado, não é possível armazená-los de forma permanente em memória.

As características dos fluxos contínuos de dados implicam as seguintes restrições aos algoritmos que os processam (??):

1. É impossível armazenar todos os dados do fluxo. Somente uma pequena parcela pode ser processada e armazenada, enquanto o restante é descartado;
2. A velocidade de chegada dos eventos exige que os elementos sejam processados em tempo real;
3. A distribuição dos dados pode mudar com o tempo. Assim, os dados do passado podem se tornar irrelevantes ou mesmo prejudiciais para a descrição dos conceitos atuais.

A área de Aprendizado de Máquina (AM) estuda algoritmos que melhoram o seu desempenho conforme ganham experiência (??). Esses algoritmos se dividem em duas categorias principais: não supervisionados (agrupamento ou *clustering*) e supervisionados (classificação ou regressão). Algoritmos de ambas as categorias foram adaptados para que pudessem ser aplicados em cenários com fluxos contínuos de dados. As principais características de cada categoria e as especializações propostas serão discutidas a seguir.

As técnicas não supervisionadas realizam o agrupamento automático de dados segundo o seu grau de semelhança. Essas técnicas têm como objetivo a formação de grupos com alta similaridade intragrupo e baixa similaridade intergrupo (??). Os seguintes algoritmos são exemplos de técnicas não supervisionadas para cenários em lote: K-Means (??), DBSCAN (??), PAM (??) e OPTICS (??).

De acordo com ??), a principal dificuldade ao aplicar técnicas não supervisionadas em cenários com fluxos contínuos é a manutenção da qualidade e consistência dos grupos formados conforme novos dados são observados. Portanto, é necessário que os algoritmos atuem de forma incremental, evoluindo os grupos formados ao longo do tempo (??). Sendo assim, foram desenvolvidos métodos não supervisionados especializados para fluxos contínuos de dados. Os seguintes trabalhos são exemplos dessas especializações: CluStream (??), StreamKM++ (??), DenStream (??), D-Stream (??) e ClusTree (??).

Os algoritmos supervisionados realizam previsões para novos exemplos utilizando um modelo criado a partir de uma base de treinamento (??). Se a previsão é categórica, entende-se como um problema de classificação. Se a previsão resulta em um valor numérico, trata-se de uma tarefa de regressão. Exemplos de algoritmos supervisionados para cenários em lote, incluem: árvores de decisão (??), métodos baseados em regras, redes neurais e máquinas de vetores suporte (SVM) (??).

Segundo ??), as técnicas supervisionadas tradicionais não podem ser aplicadas a contextos com fluxos contínuos de dados, pois estes métodos não contemplam as severas restrições de uso de memória e de tempo de execução desses cenários. Dessa forma, novos algoritmos supervisionados foram propostos para esses contextos (??????????).

As especializações mencionadas buscam atender às restrições de uso de memória e de tempo de execução dos contextos com fluxos contínuos de dados. Contudo, não consideram que na maioria desses cenários as informações são geradas de forma não estacionária, e que a distribuição dos dados ou o contexto do processo gerador podem sofrer variações, alterando os resultados esperados. Na literatura, essas variações são denominadas mudanças de conceito e a sua ocorrência pode impactar a acurácia da técnica aplicada (??).

Neste projeto de mestrado, considera-se que os dados são obtidos a partir de fluxos contínuos de dados com ocorrência de mudanças de conceito. Na próxima seção, o fenômeno mudança de conceito será discutido em detalhes.

## 2.3 MUDANÇA DE CONCEITO

Técnicas de Aprendizado de Máquina aplicadas a cenários com fluxos contínuos de dados devem ser capazes de lidar com alterações na distribuição dos dados ou no contexto do processo gerador. Essas alterações são denominadas Mudanças de Conceito (*Concept Drift* ou *Concept Shift*, em inglês) e podem alterar os resultados esperados (conceitos-

alvo) dos algoritmos, prejudicando sua acurácia (??).

Na literatura, é comum utilizar a Teoria Bayesiana de Decisão (??) para descrever a tarefa de classificação. Esta descrição será utilizada como base para formalização do fenômeno de mudança de conceito.

Sendo  $X \in \mathbb{R}^p$  uma instância em um espaço  $p$ -dimensional de atributos e  $X \in c_i$  onde  $c_1, c_2, \dots, c_k$  é o conjunto de classes, o classificador ótimo para classificar  $x \rightarrow c_i$  é determinado a partir das probabilidades a priori das classes  $P(c_i)$  e pela função de densidade de probabilidade condicionada às classes  $p(X|c_i), i = 1, \dots, k$ . É possível definir um conceito como um conjunto de probabilidades a priori e condicionais das classes, como mostra a Equação ??:

$$S = \{(P(c_1), P(X|c_1)), (P(c_2), P(X|c_2)), \dots, (P(c_k), P(X|c_k))\} \quad (2.1)$$

Ainda segundo a Teoria Bayesiana, a classificação de uma instância  $X$  baseada na máxima probabilidade a posteriori pode ser obtida através da Equação ??:

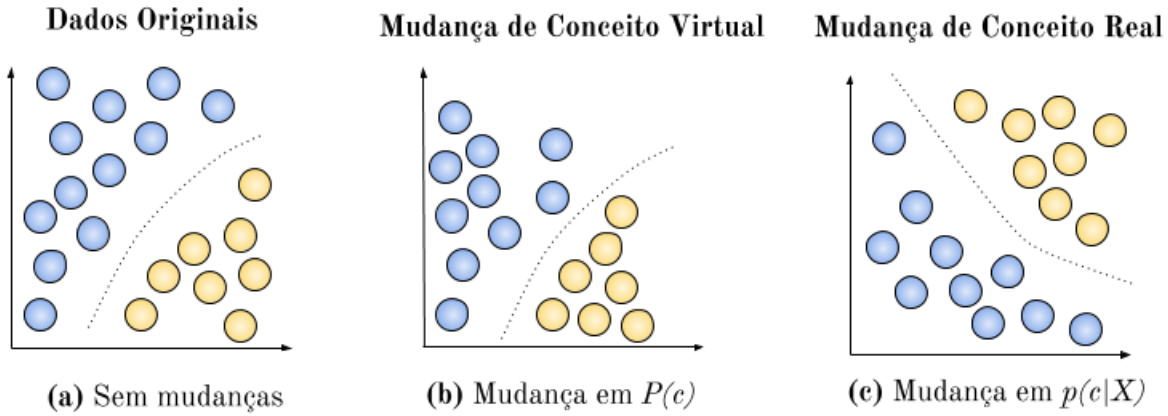
$$p(c_i|X) = \frac{p(c_i) * p(X|c_i)}{p(X)} \quad (2.2)$$

Assim, é possível afirmar que há mudança de conceito entre os instantes  $t_0$  e  $t_1$  se:

$$\exists X : p_{t_0}(X, c) \neq p_{t_1}(X, c) \quad (2.3)$$

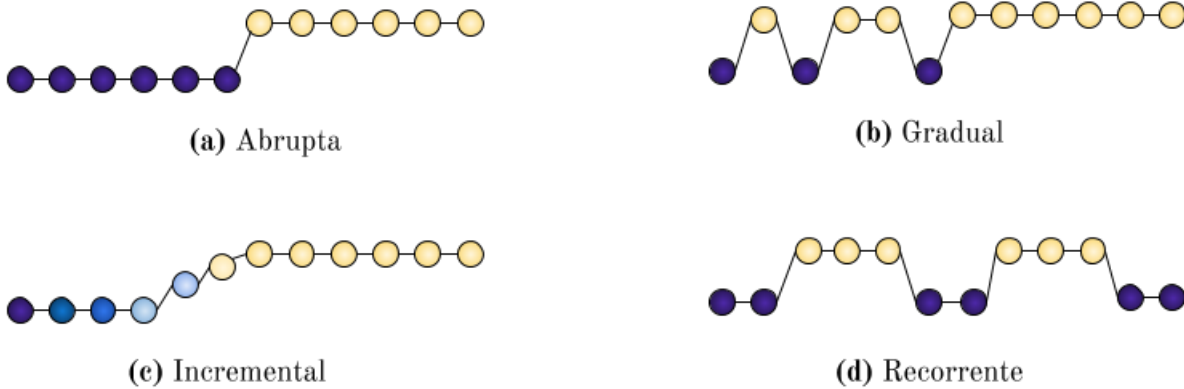
onde  $p_{t_0}$  e  $p_{t_1}$  denotam as distribuições de probabilidades conjuntas nos instantes  $t_0$  e  $t_1$ , respectivamente, para  $X$  e  $c$  (??). Em outras palavras, um conjunto de dados possui resultados esperados legítimos em  $t_0$ , mas este mesmo conjunto passa a ter resultados esperados diferentes, também legítimos, em  $t_1$  (??).

De acordo com ??), as mudanças de conceito podem ser categorizadas como virtuais ou reais. As mudanças virtuais são causadas por alterações na probabilidade a priori das classes,  $P(c)$ , e não alteram os conceitos-alvo. Enquanto que as mudanças de conceito reais surgem a partir de alterações na probabilidade a posteriori,  $p(c|X)$ , e modificam os resultados esperados. Os dois tipos de mudança de conceito estão representados na Figura ??.



**Figura 2.1** Mudança de Conceito Virtual vs. Mudança de Conceito Real

Conforme ??), as mudanças de conceito podem ocorrer de forma abrupta, gradual, incremental ou recorrente. A Figura ?? ilustra estes padrões, utilizando círculos na cor azul para representar o conceito  $A$  e círculos na cor bege para o conceito  $B$ :



**Figura 2.2** Padrões de ocorrência de Mudanças de Conceito

Na mudança abrupta, o conceito  $A$  é repentinamente substituído pelo conceito  $B$  (Figura ?? (a)).

Na mudança gradual, ocorre uma transição mais suave entre os conceitos  $A$  e  $B$ . Inicialmente, eventos pertencentes a ambos os conceitos coexistem. Com o passar do tempo, os eventos pertencentes ao conceito  $A$  diminuem de frequência, até pararem de ocorrer. Por fim, os eventos pertencentes a  $B$  se tornam predominantes (Figura ?? (b)).

A mudança incremental descreve a evolução de um único conceito ao longo do tempo. Essa evolução pode ser discretizada como uma sequência de conceitos consecutivos. Nesta sequência, cada conceito intermediário difere pouco dos seus conceitos antecessor e sucessor. Portanto, as mudanças são notáveis apenas à longo prazo (Figura ?? (c)).

A mudança recorrente acontece quando um conceito anteriormente ativo reaparece após um determinado período de tempo. Contudo, não se trata de uma sazonalidade periódica, pois não é evidente o momento no qual o conceito voltará a ser ativo (Figura ?? (d)).

Este trabalho de mestrado propõe um método baseado em Redes de Função de Base Radial para detecção de mudanças de conceito em fluxos contínuos de dados, independente do padrão de ocorrência. Na próxima subseção, será apresentada a terminologia do fenômeno mudança de conceito.

### 2.3.1 Terminologia

O fenômeno mudança de conceito tem sido estudado em diferentes comunidades de pesquisa, incluindo Mineração de Dados, Aprendizado de Máquina, Estatística e Recuperação de Informação (??). Contudo, o tema apresenta diferentes nomenclaturas em cada comunidade. Na Tabela ?? são listados os termos correspondentes a mudança de conceito em cada área de pesquisa.

**Tabela 2.1** Terminologia - Mudança de Conceito (??)

Área	Termos
Mineração de Dados	Mudança de Conceito
Aprendizado de Máquina	Mudança de Conceito, Mudança de Covariável
Computação Evolucionária	Ambiente Evolutivo, Ambiente em Mudança
IA e Robótica	Ambiente Dinâmico
Estatística, Séries Temporais	Não Estacionário
Recuperação de Informação	Evolução Temporal

Outra fonte comum de equívocos são os termos detecção de *outliers*, detecção de novidade, detecção de *change points* e detecção de mudança de conceito. Estes termos são muitas vezes utilizados de forma indistinta, mas, para o contexto deste trabalho, é importante distingui-los.

As técnicas para detecção de *outliers* têm como objetivo identificar padrões de dados em desacordo com o comportamento esperado. Estes padrões são geralmente classificados como anomalias ou ruídos (??).

Os métodos para detecção de novidade identificam padrões ainda não observados, mas que se enquadram no comportamento esperado. Estes métodos se diferenciam das técnicas para detecção de *outliers* pois os novos padrões são incorporados ao modelo (??).

As estratégias para detecção de *change points* identificam variações abruptas de valor, que podem representar transições entre estados, em séries temporais unidimensionais estacionárias (??).

Por fim, cabe definir o que é detecção de mudança de conceito. Este termo abrange técnicas que monitoram a distribuição dos dados ou indicadores produzidos pelas técnicas de aprendizado aplicadas, com o objetivo de identificar variações que possam impactar a acurácia do modelo em uso (??).

Na próxima subseção, os principais algoritmos para detecção de mudança de conceito serão discutidos.

### 2.3.2 Algoritmos para Detecção de Mudança de Conceito

Os algoritmos para detecção de mudança de conceito caracterizam e quantificam as mudanças de conceito através da delimitação dos instantes ou intervalos de tempo em que as mudanças ocorrem (??).

Esses algoritmos se dividem em duas categorias, conforme a necessidade de rotulação dos dados (??):

**Algoritmos Explícitos/Supervisionados** Dependem da rotulação dos dados por um especialista. Estes rótulos são utilizados no cálculo de medidas de performance como taxa de erro e acurácia, que são monitoradas ao longo do tempo. Mudanças de conceito são sinalizadas quando essas medidas atingem um limite previamente definido.

**Algoritmos Implícitos/Não Supervisionados** Independem da rotulação por especialistas. As mudanças de conceito são detectadas a partir do monitoramento de características dos próprios dados ou de indicadores produzidos pelas técnicas de aprendizado aplicadas. Apesar de serem mais propensos a alarmes falsos, são uma alternativa para cenários onde a obtenção de rótulos é dispendiosa, demorada ou inviável.

Segundo ??), os algoritmos *explícitos* / *supervisionados* podem ser segmentados em três subcategorias:

**Métodos Baseados em Análise Sequencial** Avaliam continuamente os indicadores de performance (por exemplo: taxa de erro) do algoritmo de aprendizado aplicado. A mudança de conceito é detectada quando esses indicadores atingem um limite pré-definido. Os algoritmos *Cumulative Sum (CUSUM)*, *PageHinkley (PH)* (??) e *Geometric Moving Average (GMA)* (??) são representantes desta subcategoria.

**Abordagens baseadas em Estatística** Identificam mudanças de conceito através da análise de parâmetros estatísticos como média e desvio padrão associados aos resultados das predições. Os métodos *Drift Detection Method (DDM)* (??), *Early Drift Detection Method (EDDM)* (??), *Exponentially Weighted Moving Average (EWMA)* (??) e *Reactive Drift Detection Method (RDDM)* (??) são exemplos desta subcategoria.

**Métodos baseados em Janelas** Utilizam uma janela de tamanho fixo para sumarizar informações passadas e uma janela deslizante para sumarizar os dados recentes. Uma diferença significativa entre as distribuições dessas janelas implica na



ocorrência de mudança de conceito. Esta diferença é verificada a partir de testes estatísticos ou desigualdades matemáticas, considerando como hipótese nula a igualdade das distribuições. Os algoritmos *Adaptive Windowing (ADWIN)* (??), *SeqDrift* (??), *HDDMA* e *HDDMW* (??) pertencem a esta subcategoria.

De forma similar, os algoritmos *implícitos / não supervisionados* também foram divididos em três subcategorias (??):

**Detecção de Novidade / Métodos de Agrupamento** Utilizam técnicas derivadas dos métodos de agrupamento e de detecção de *outliers* para identificar padrões ainda não observados. A partir dessa identificação, são realizados cálculos de distância e/ou densidade para confirmar a ocorrência de mudança de conceito (??). Os métodos *OLINDDA* (??), *MINAS* (??), *Woo* (??), *DETECTNOD* (??), *ECSMiner* (??) e *GC3* (??) fazem parte desta subcategoria.

**Monitoramento de distribuição multivariada** Monitoram diretamente a distribuição dos dados para cada atributo. A distribuição de um conjunto de treinamento é resumizada e utilizada como referência. Esta referência é, então, comparada à distribuição dos dados do conjunto atual. Diferenças significativas entre esses conjuntos indicam a ocorrência de mudança de conceito. Os algoritmos *CoC* (??), *HDDDM* (??), *PCA-detect* (??) são representantes desta subcategoria.

**Monitoramento dependente de modelo** Requerem a aplicação de um algoritmo de classificação probabilístico, pois as mudanças de conceito são detectadas a partir do monitoramento da probabilidade a posteriori calculada (??). Esta restrição permite reduzir a ocorrência de falsos positivos e torna o processo computacionalmente eficiente, pois apenas um único fluxo univariado de valores é observado. Os métodos *A-distance* (??), *CDBD* (??) e *Margin* (??) integram esta subcategoria.

Por fim, a Tabela ?? resume as categorias, as subcategorias e as técnicas apresentadas. O método de detecção proposto neste trabalho se enquadra na categoria de algoritmos *implícitos / não supervisionados*, mais especificamente na subcategoria *detecção de novidades / métodos de agrupamento*. Na próxima seção, as ferramentas utilizadas para implementação e validação deste método serão apresentadas.

### 2.3.3 Ferramentas

Nesta seção, os frameworks *Massive Online Analysis* (MOA) e *Tornado* serão apresentados. Estes frameworks serão utilizados para implementar e validar a técnica de detecção proposta neste trabalho. Além disso, essas ferramentas possuem implementações para os principais algoritmos de detecção de mudança de conceito presentes na literatura, possibilitando a comparação deste trabalho com o estado da arte.

**Tabela 2.2** Resumo - Algoritmos de detecção (??)

Categoria	Subcategoria	Algoritmos
Algoritmos Explícitos/Supervisionados	Métodos Baseados em Análise Sequencial	Cumulative Sum (CUSUM)
		PageHinkley (PH) (??)
		Geometric Moving Average (GMA) (??)
	Abordagens baseadas em Estatística	Drift Detection Method (DDM) (??)
		Early Drift Detection Method (EDDM) (??)
		Exponentially Weighted Moving Average (EWMA) (??)
		Reactive Drift Detection Method (RDDM) (??)
	Métodos baseados em Janelas	Adaptive Windowing (ADWIN) (??)
		SeqDrift (??)
		HDDMA/HDDMW (??)
Algoritmos Implícitos/Não Supervisionados	Detecção de Novidade / Métodos de Agrupamento	OLINDDA (??)
		MINAS (??)
		Woo (??)
		DETECTNOD (??)
		ECSMiner (??)
		GC3 (??)
		CoC (??)
	Monitoramento de distribuição multivariada	HDDDM (??)
		PCA-detect (??)
		A-distance (??)
	Monitoramento dependente de modelo	CDBD (??)
		Margin (??)

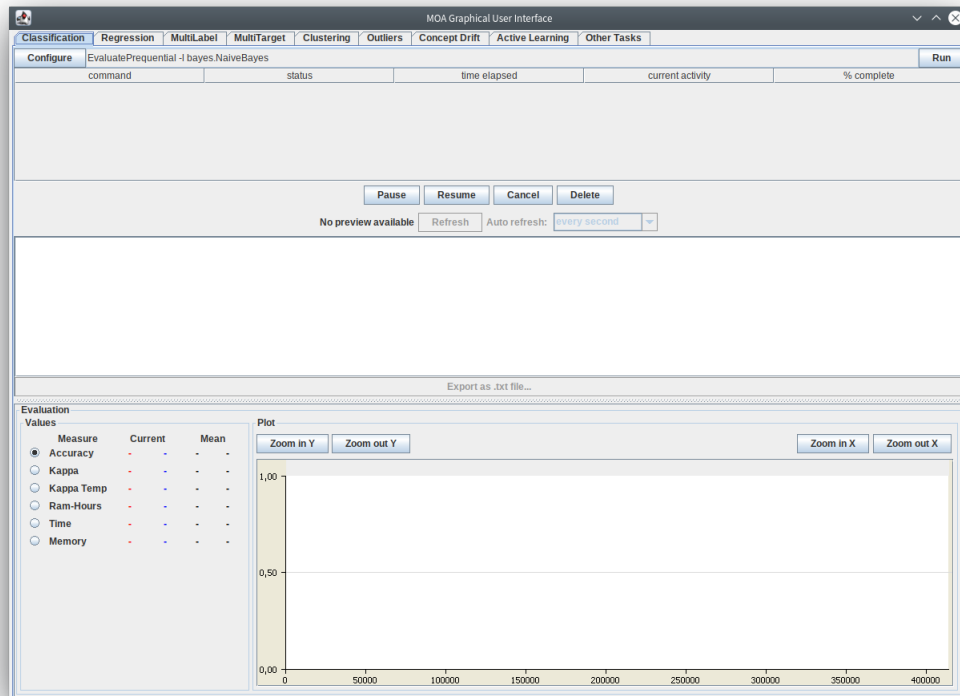
### 2.3.4 MOA

Atualmente, o *MOA – Massive Online Analysis*<sup>1</sup> é o principal framework para mineração de dados em fluxos contínuos (??). O framework é de código-aberto<sup>2</sup> e apresenta uma comunidade bastante ativa e crescente. A aplicação é composta por uma ampla coleção de algoritmos da área de Aprendizado de Máquina, contemplando técnicas de classificação, regressão, agrupamento, busca por padrões, detecção de *outliers*, detecção de mudanças de conceito e sistemas de recomendação. Além das implementações, também estão disponíveis rotinas para avaliação dessas técnicas. A aplicação é desenvolvida em Java, o que permite a sua execução nos principais sistemas operacionais e a integração com o projeto WEKA (??).

O MOA divide as suas funcionalidades em tarefas (*tasks*). Estas tarefas podem ser executadas a partir da interface gráfica (GUI) ou por linha de comando. A interface gráfica permite executar múltiplas tarefas de forma concorrente, controlar suas execuções e visualizar os resultados parciais. A tela principal da aplicação é demonstrada na Figura ??.

<sup>1</sup><https://moa.cms.waikato.ac.nz/>

<sup>2</sup><https://github.com/Waikato/moa>



**Figura 2.3** MOA - Tela Inicial

A aplicação é capaz de ler arquivos em formato ARFF<sup>3</sup>, além de permitir a produção de fluxos de dados dinamicamente, através de geradores. Alguns dos geradores de fluxo disponíveis no MOA são: *Random Trees* (??), *SEA* (??), *STAGGER* (??), *Rotating Hyperplane* (??), *Random RBF*, *LED* (??), *Waveform* (??), e *Function* (??).

Outra funcionalidade importante do framework é a possibilidade de adicionar mudanças de conceito a fluxos estacionários existentes. Esse processo é realizado através de uma função sigmóide, que modela o evento de mudança de conceito como uma combinação balanceada de duas distribuições homogêneas, que caracterizam os conceitos-alvo antes e depois da mudança. Além destes conceitos, o usuário também pode definir o momento da mudança e a sua duração (??).

A arquitetura do framework é modular, o que permite a implementação de novos detectores de forma trivial. Para tanto, basta estender a classe abstrata *AbstractChangeDetector* e implementar o algoritmo desejado. A janela de configuração deste detector, similar a Figura ??, será criada dinamicamente, a partir dos atributos definidos na classe.

<sup>3</sup>Attribute-Relation File Format



**Figura 2.4** MOA - Configuração detector

Além de possuir implementações para diversos algoritmos da área de Aprendizado de Máquina e de áreas correlatas, o framework também dispõe de classes para avaliar a acurácia e a performance dessas técnicas. Neste trabalho de mestrado, a classe *BasicConceptDriftPerformanceEvaluator* será utilizada para avaliar a acurácia e a performance do método proposto. Na próxima subseção, o framework *Tornado* será apresentado.

### 2.3.5 Tornado

O *Tornado* é um framework desenvolvido para analisar algoritmos de detecção de mudanças de conceito (??). O projeto é construído na linguagem Python e o seu código está disponível<sup>4</sup>. O framework se diferencia do *MOA*, pois apresenta um cenário de avaliação específico: analisar a execução, em paralelo, de pares formados por um classificador e um detector de mudanças, para identificar o par ótimo ao longo do tempo, em relação ao fluxo de dados.

---

<sup>4</sup><https://github.com/alipsggh/tornado>

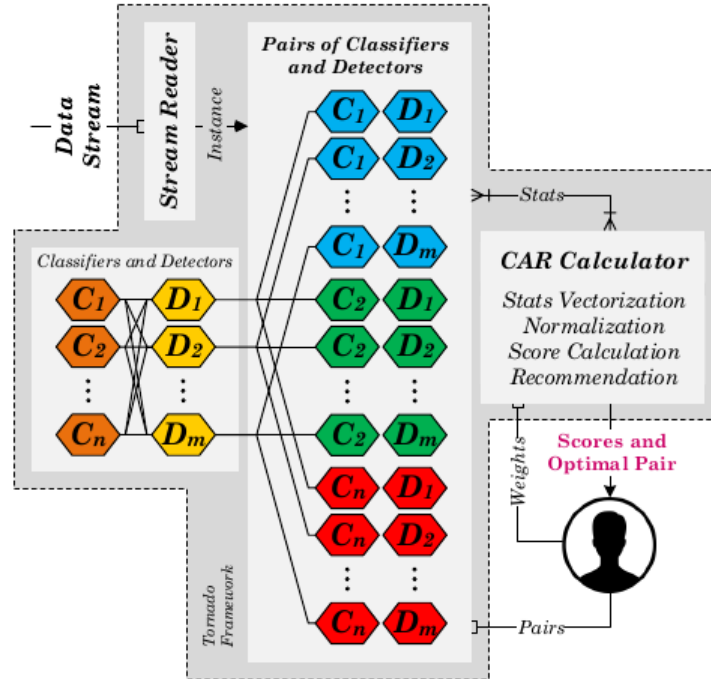
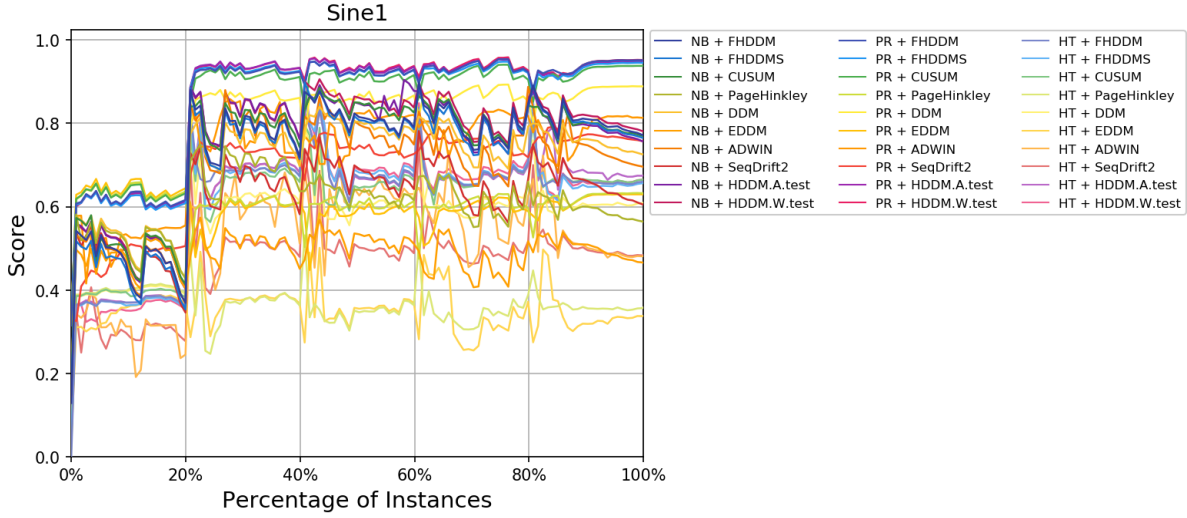


Figura 2.5 Framework Tornado (??).

Conforme apresentado na Figura ??, os principais componentes do framework são: *Stream Reader*, *Classifiers*, *Detectors*, *Classifier-Detector Pairs* e *CAR Calculator*. A entrada de dados é composta por um fluxo (*Stream*), uma lista de pares (classificador, detector) e um vetor com pesos. O componente *Stream Reader* recebe as instâncias e as encaminha para construção do modelo de forma incremental. Por seguir a abordagem *prequential*, cada instância é primeiramente utilizada para testes e depois como treinamento. Simultaneamente, os classificadores enviam suas estatísticas aos detectores, para que a mudança de conceito possa ser sinalizada. Por fim, o componente *CAR Calculator* calcula uma pontuação para cada par, baseada na taxa de erro, atraso para detecção (*delay*), falsos positivos, falsos negativos, quantidade de memória utilizada e tempo de execução (??).

O resultado produzido identifica o par ótimo para cada instante da execução. Esta abordagem de avaliação se mostra relevante, pois o par ideal pode mudar ao longo do tempo, devido ao aprendizado incremental ou às mudanças de conceito. A Figura ?? apresenta um exemplo de resultado da ferramenta.

Neste trabalho, o algoritmo proposto será implementado e testado nas duas ferramentas apresentadas. Os detalhes de implementação e os resultados desses testes serão discutidos no Capítulo ??. A seguir, as Redes de Função de Base Radial são detalhadas.



**Figura 2.6** Tornado - Exemplo de resultado (??)

## 2.4 REDES DE FUNÇÃO DE BASE RADIAL

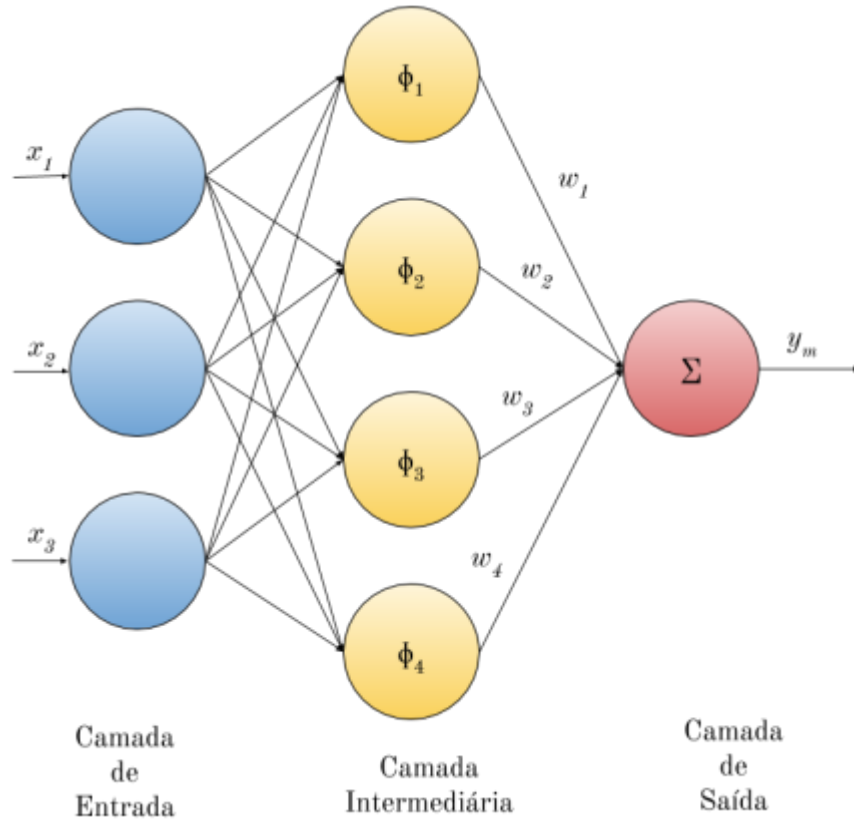
Uma Rede de Função de Base Radial (redes RBF, em português, ou *radial basis function network*, *RBF network*, em inglês) pode ser definida como um modelo de múltiplas camadas alimentadas adiante (*feedforward*), capaz de analisar padrões complexos e resolver problemas não-linearmente separáveis, utilizando uma abordagem de aproximação de funções. Estas redes têm como principal diferencial a sua forma de ativação, realizada através do cálculo da distância entre o dado e um centro definido (??).

A arquitetura de uma rede de função de base radial, em sua forma mais básica, envolve três camadas. A camada de entrada representa os atributos do problema. As unidades dessa camada não realizam processamento e simplesmente distribuem as entradas para os neurônios na camada intermediária. A camada intermediária, única camada oculta da rede, é constituída por funções de ativação de base radial, que atuam como os neurônios da rede. Por fim, a camada de saída pondera, através de pesos, os resultados da camada intermediária, agregando-os linearmente para compor a resposta final da rede (??). A Figura ?? demonstra essa arquitetura.

Na literatura, as funções Gaussianas (Função ??) são as funções de ativação mais usuais em redes RBF, sendo definidas por seus centros e larguras:

$$\varphi(v_i) = e^{-(\sigma r)^2} \quad (2.4)$$

onde,  $r$  é a distância euclidiana ( $\|\mathbf{v}_i - \mathbf{c}_i\|$ ),  $v$  é o valor de entrada,  $c_i$  representa o centro e  $\sigma$  é o parâmetro limitador do raio.

**Figura 2.7** Arquitetura RBF

A modelagem das redes RBF define a entrada como um vetor de números reais  $x \in R^n$  e o resultado da rede como uma função escalar desse vetor  $\varphi : R^n \rightarrow R$ . Esta função é obtida a partir da combinação linear dos resultados das funções de ativação (Equação ??).

$$\varphi(x) = \sum_{i=1}^N a_i \rho(\|x - c_i\|) \quad (2.5)$$

onde  $N$  é o número de neurônios na camada intermediária,  $c_i$  é o centro do neurônio  $i$ ,  $\|\dots\|$  é a norma euclidiana e  $a_i$  é o peso (bias) atribuído ao neurônio para realização da combinação linear.

O algoritmo proposto neste trabalho utiliza as camadas inicial e intermediária das Redes de Função de Base Radial para compôr um novo método de detecção de mudanças de conceito. Isto é possível, pois a camada intermediária cria, de forma implícita, agrupamentos durante o processo de ativação. Dessa forma, a criação de novos centros e a mudança do centro ativo dos neurônios podem sinalizar a ocorrência de mudanças de conceito. A técnica desenvolvida será discutida em detalhes no Capítulo ??. Na próxima seção, os trabalhos relacionados encontrados na literatura são apresentados.

## 2.5 TRABALHOS RELACIONADOS

Além das referências básicas apresentadas neste capítulo, foi realizada uma pesquisa na literatura em busca de trabalhos que propõem métodos para identificação de mudanças de conceito, de forma online e independente de rótulos, em fluxos contínuos de dados. Também foram estudadas técnicas que pudessem subsidiar o desenvolvimento de novos algoritmos que atendam a esses requisitos.

Inicialmente, foram analisados os algoritmos *Implícitos / Não Supervisionados* pertencentes a subcategoria *Deteção de Novidade / Métodos de Agrupamento*. Estes algoritmos são descritos a seguir.

O algoritmo *OnLine Novelty and Drift Detection Algorithm* (OLINDDA) utiliza a técnica de agrupamento *K-means* para monitorar e se adaptar ao surgimento de novos padrões (??). Os exemplos recepcionados são incluídos em uma fila e são periodicamente agrupados. Os grupos resultantes desse processamento definem se os exemplos devem ser incorporados a um grupo existente ou formar um novo grupo.

O algoritmo MINAS, proposto por ??), atua de forma similar ao algoritmo OLLINDA, contudo utiliza *micro clusters* obtidos a partir da aplicação de uma técnica de agrupamento específica para fluxos contínuos de dados (CluStream) e estende a abordagem para permitir a aplicação em problemas multiclasse.

O algoritmo DETECTNOD (??) utiliza um modelo de clusterização para delimitar os exemplos observados. Uma representação compacta desses grupos é produzida através da Transformada Discreta de Cosseno. Esta representação é utilizada, em conjunto com o algoritmo de vizinho mais próximo, para aproximar novos exemplos. Estes exemplos são categorizados como desvios ou novidades, conforme o grau de similaridade calculado.

Os algoritmos Woo (??) e ECSMiner (??) também se baseiam no conceito de *micro clusters*. A técnica Woo atribui um classificador para cada cluster formado. Estes classificadores são utilizados para verificar a similaridade de novos exemplos em relação ao cluster. Exemplos que não se enquadram em nenhum grupo são marcados como suspeitos e passam a ter sua densidade monitorada. Um número crescente de vizinhos no raio dos exemplos monitorados indica o surgimento de um novo conceito, induzindo o retreino do modelo e ajuste do centróide do agrupamento. O ECSMiner (??) utiliza o conceito de Outliers Filtrados, que se refere a amostras que estão fora do limite de todos os clusters existentes, em uma abordagem similar ao algoritmo Woo.

O algoritmo GC3 (??) estende a ideia de *micro clusters* ao aplicar um algoritmo de agrupamento baseado em grade e densidade. Dessa forma, novos padrões são identificados a partir do surgimento de novas áreas de alta densidade.

Os algoritmos baseados em técnicas de detecção de novidade utilizam métodos de agrupamento para identificar padrões emergentes. Portanto, também apresentam dificuldade para lidar com alta dimensionalidade e/ou dados binários, por dependerem do cálculo de distância, e são computacionalmente custosos.

A segunda etapa da pesquisa analisou técnicas para detecção de mudanças de conceito em séries temporais que atuam de forma online. Estas técnicas também são comumente denominadas como técnicas para detecção de *change points*. A maioria dos trabalhos identificados se baseia no monitoramento dos valores residuais da aplicação de um modelo.



Nesta abordagem, um modelo estatístico ou de Aprendizado de Máquina é aplicado a um conjunto de observações extraídas da série temporal, representando um conceito. Assume-se que se não houver mudança de conceito, os valores residuais dessa aplicação devem ser estacionários. Os trabalhos identificados são apresentados a seguir.

??) propõem um método de detecção online capaz de identificar outliers e mudanças de conceito em séries temporais. A abordagem aplica um modelo autoregressivo aos dados e atualiza os parâmetros de forma incremental, para que o peso de exemplos passados diminua ao longo do tempo. Atribui-se uma pontuação para cada observação, calculada com base na função de perda. Esta pontuação indica o grau de desvio do dado em relação ao modelo aplicado. Pontuações elevadas indicam alta possibilidade de que o dado seja um outlier. A detecção de mudança de conceito é feita através do monitoramento da média das funções de perda ao longo da série temporal. Todavia, esta abordagem só é aplicável a séries temporais com comportamento autoregressivo, o que limita o seu escopo de uso.

??) propõem um método diferente, também online, para detecção de mudanças em séries temporais autoregressivas, baseado no monitoramento dos parâmetros do modelo. Uma adaptação do método CUSUM (??) é aplicada para identificar alterações nos parâmetros. A hipótese nula do algoritmo assume que estes parâmetros devem permanecer estacionários ao longo do tempo. Esta hipótese é testada para cada nova observação recepcionada e a mudança de conceito é sinalizada quando a diferença ultrapassa um limiar estabelecido. Esta abordagem também apresenta limitação de escopo, pois é aplicável apenas a séries temporais autoregressivas.

A utilização de valores residuais para detecção de mudanças de conceito apresenta desvantagens. Por serem baseados na acurácia do modelo aplicado, estes métodos sofrem influência de problemas ocorridos durante a parametrização ou fase de treinamento. Em teoria, o monitoramento direto de características dos dados pode permitir uma identificação de mudanças de conceito mais robusta.

Poucos estudos utilizam as características das séries temporais para identificar mudanças de conceito. Neste contexto, o trabalho proposto por ??) se destaca por apresentar um método de detecção de mudanças de conceito em séries temporais que atua de forma online, baseando-se no monitoramento da característica de auto-similaridade da série. Entretanto, o método proposto tem como principal limitação o fato de ser especializado para séries que apresentam auto-similaridade e periodicidade. Além disso, a abordagem utiliza uma considerável quantidade de memória, pois armazena uma sequência de dados que representa o comportamento estável ou regular da série temporal.

Para cenários nos quais os fluxos de dados não apresentam as propriedades esperadas pelas técnicas discutidas anteriormente, ??) propõem um novo método para detecção de mudanças de conceito que atua de forma online. Este trabalho se diferencia por utilizar um algoritmo de agrupamento hierárquico estável, conforme a propriedade proposta por ??). A abordagem cria janelas consecutivas de dados. A técnica de agrupamento é aplicada em cada janela e os dendogramas resultantes comparados, a fim de identificar mudanças de conceito na distribuição dos dados.

As técnicas para detecção de mudanças de conceito são, geralmente, aplicadas para minimizar o intervalo entre a ocorrência da mudança e a atualização do modelo aplicado.

Diante disso, um método de detecção ideal deve:

- ser transparente para o usuário, detectando e informando explicitamente sobre a ocorrência de mudanças;
- ser computacionalmente eficiente;
- atuar de forma online, buscando minimizar o atraso de atualização do modelo; e
- basear-se em características dos dados que reflitam os conceitos presentes.

Por fim, a pesquisa realizada evidenciou a inexistência de técnicas que atendam a esses requisitos e que sejam computacionalmente eficientes, independentes de rótulos e aplicáveis a fluxos de dados de qualquer natureza. Este projeto de trabalho de mestrado tem como objetivo tentar preencher esta lacuna na literatura, através da proposição de um novo método, baseado em Redes de Função de Base Radial, para identificação de mudanças de conceito em fluxos contínuos de dados.

## **2.6 CONSIDERAÇÕES FINAIS**

Neste capítulo foram apresentados os principais conceitos utilizados na pesquisa. Foram discutidos conceitos de Fluxos Contínuos de Dados, técnicas de Aprendizado de Máquina, Mudança de Conceito, Técnicas de Detecção e Redes de Função de Base Radial. Por fim, foram apresentados os trabalhos relacionados encontrados na literatura. No próximo capítulo, o plano de pesquisa será detalhado.

## **PLANO DE PESQUISA**

### **3.1 CONSIDERAÇÕES INICIAIS**

Este capítulo descreve como este projeto de mestrado será desenvolvido. A abordagem proposta utiliza Redes de Função de Base Radial para detectar mudanças de conceito em fluxos contínuos de dados. Espera-se que a utilização de redes RBF permita detectar mudanças em tempo de execução, de forma computacionalmente eficiente e independente de rótulos. A seguir, são apresentados detalhes sobre cada etapa do desenvolvimento do projeto.

### **3.2 DESCRIÇÃO DO PROBLEMA**

Nos últimos anos, a quantidade de dados produzidos por sistemas computacionais tem crescido de forma exponencial (??). Parte significativa dessas informações é produzida na forma de fluxos contínuos de dados, que são sequências potencialmente infinitas e de alta frequência (??).

Devido a esse crescimento, pesquisadores passaram a utilizar técnicas de Aprendizado de Máquina para extrair informações úteis de grandes volumes de dados. Essas técnicas precisaram ser adaptadas para contextos com fluxos contínuos, pois estes cenários apresentam severas restrições de tempo de execução e de uso dos recursos computacionais.

Contudo, as adaptações propostas não tratam alterações na distribuição dos dados ou no contexto do processo gerador do fluxo. Estas alterações são denominadas mudanças de conceito e podem afetar negativamente a acurácia do algoritmo (??).

Para mitigar este problema, métodos de detecção de mudanças de conceito foram desenvolvidos. Estes métodos identificam o momento da mudança com maior precisão, permitindo que o modelo de decisão seja atualizado de forma eficiente.

Entretanto, as técnicas de detecção encontradas na literatura apresentam limitações ao serem aplicadas em cenários com fluxos contínuos de dados. Os métodos de detecção supervisionados/explicitos necessitam que o rótulo correto de cada exemplo processado seja informado, o que pode torná-los inviáveis, por causa do custo e/ou do tempo para

rotulação. Enquanto que as técnicas não supervisionadas/implícitas têm dificuldade para atender às restrições de tempo de execução e de uso dos recursos computacionais desses cenários. Diante disso, este trabalho propõe uma nova abordagem para detecção de mudanças de conceito baseada em Redes de Função de Base Radial.

As redes RBF podem ser definidas como redes neurais multicamadas, alimentadas adiante, capazes de analisar padrões complexos e resolver problemas que não são linearmente separáveis (??). A arquitetura básica dessas redes é composta por três camadas. A camada de entrada recebe os dados e repassa para a camada seguinte. A camada intermediária realiza a ativação dos neurônios através de funções de base radial. E, por fim, a camada de saída produz o resultado final da rede através de uma combinação linear dos resultados da camada intermediária (??).

Este trabalho de mestrado utiliza apenas as camadas inicial e intermediária dessa arquitetura. Isto ocorre, pois a camada intermediária forma, implicitamente, grupos (*clusters*) durante o processo de ativação. Este agrupamento possui um centro ativo que muda conforme o valor processado. O algoritmo desenvolvido nesta pesquisa monitora o agrupamento formado para identificar quando esse centro é alterado. Estas alterações são utilizadas como indícios de possíveis mudanças de conceito.

O método proposto, denominado *RBFDriftDetector*, utiliza a função Gaussiana para ativação dos neurônios e requer a definição de dois parâmetros:  $\sigma$ , responsável por limitar o raio da radial, e  $\lambda$ , que define um limiar para ativação de um centro. A técnica é descrita na forma de pseudocódigo em Algoritmo ??, tendo sido também implementada na linguagem Java<sup>1</sup>, para realização dos experimentos na plataforma MOA.

**Entrada:** *valor*,  $\sigma$ ,  $\lambda$

**Saída:** booleano indicando a ocorrência ou não de mudança de conceito

**início**

```

centros  $\leftarrow$  (); centroAtual  $\leftarrow$  null; centroAtivo  $\leftarrow$  null
mudanca  $\leftarrow$  falso
para todo centro faça
    | ativacao  $\leftarrow$  gaussiana(valor, centro,  $\sigma$ )
    | se ativacao  $\geq \lambda$  então
    | | centroAtivo  $\leftarrow$  centro
    | |  $\lambda \leftarrow$  ativacao
se centroAtivo == null então
    | centros  $\leftarrow$  valor
    | centroAtivo  $\leftarrow$  valor
se centroAtual  $\neq$  centroAtivo então
    | centroAtual  $\leftarrow$  centroAtivo
    | mudanca  $\leftarrow$  verdadeiro
retorna mudanca

```

#### Algoritmo 1: RBFDRIFTDETECTOR

Para exemplificar a execução do algoritmo proposto neste projeto, considere o conjunto  $S = \{0.11, 0.12, 0.13, 0.34, 0.45, 0.47, 0.33, 0.25, 0.14, 0.10\}$  como fonte de dados.

---

<sup>1</sup><https://git.io/fjGuv>

Para este exemplo, os parâmetros foram definidos de forma empírica. O parâmetro  $\sigma$  foi definido com o valor 0.2 e o parâmetro  $\lambda$  foi fixado em 0.6. A seguir, será descrito o comportamento do algoritmo para cada dado recebido a partir de  $S$ .

No instante  $T1$ , o valor 0.11 é recepcionado. Como não existem centros estabelecidos, o valor é definido como centro e ativado. Em  $T2$  e  $T3$  são recebidos, respectivamente, os valores 0.12 e 0.13 que são imediatamente vinculados ao centro atualmente ativo (0.11), pois seus valores de ativação foram maiores que 0.6 ( $\lambda$ ).

No instante  $T4$ , o valor 0.34 não atinge o valor mínimo de ativação ( $\lambda$ ) para ser vinculado ao centro ativo (0.11). Dessa forma, o valor é definido como um novo centro e ativado. Devido a alteração do centro ativo, o algoritmo sinaliza a ocorrência de mudança de conceito.

Os valores dos instantes  $T5$ ,  $T6$ ,  $T7$  e  $T8$  são vinculados ao último centro ativo (0.34). Contudo, em  $T9$ , o valor 0.14 é recebido e apresenta maior valor de ativação para o centro 0.11, que é reativado. Assim, uma nova mudança de conceito é sinalizada. Finalmente, o valor do instante  $T10$ , 0.10, é vinculado ao centro ativo (0.11).

A Figura ?? apresenta esse comportamento graficamente. Os círculos na cor amarela representam valores vinculados ao primeiro centro estabelecido (0.11 no momento  $T1$ ), enquanto os de cor lilás foram ativados pelo segundo centro definido (0.34 em  $T4$ ). As linhas verticais tracejadas, de cor vermelha, indicam a ocorrência de mudança de conceito.

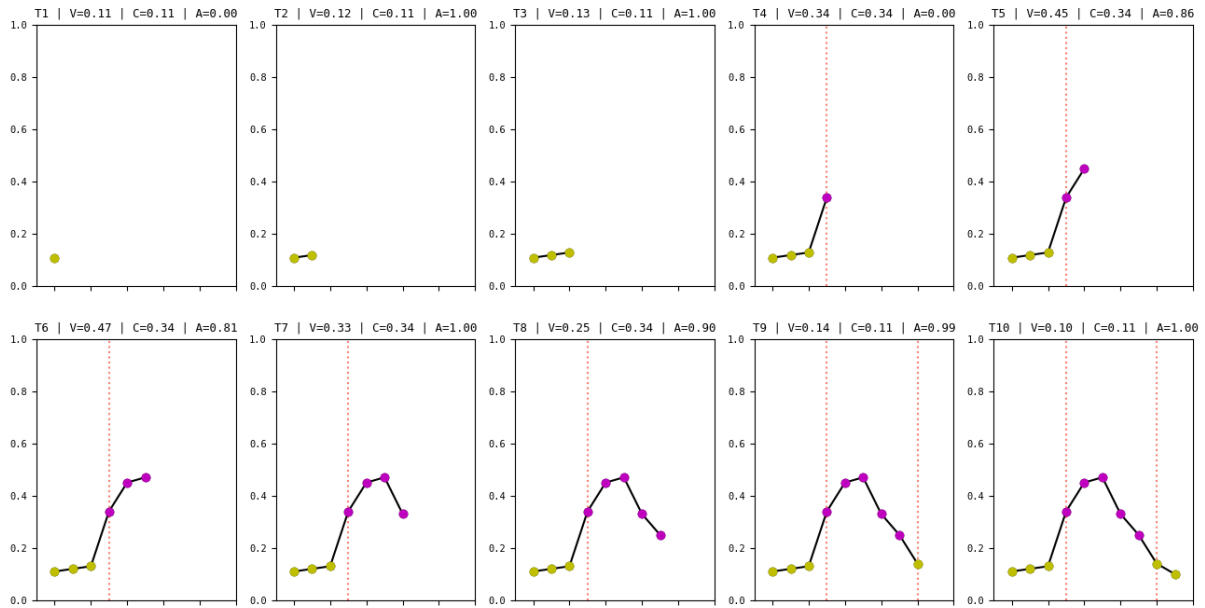


Figura 3.1 Exemplo de funcionamento do algoritmo

### 3.3 ATIVIDADES DE PESQUISA

A Tabela ?? apresenta o cronograma das atividades planejadas para a realização da pesquisa. Atividades concluídas são representadas pelo símbolo  $X$  e as futuras por  $\bullet$ .

**Tabela 3.1** Cronograma de atividades

Atividades	Meses																							
	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24
1-Disciplinas	X	X	X	X	X	X	X	X	X	X														
2-Revisão da Literatura							X	X	X	X	X	X												
3-Experimentos									X	X	X	X	•	•	•	•								
4-Análise dos Resultados											X	X				•	•	•						
5-Escrita da qualificação										X	X	X												
6-Estágio docente																•	•	•	•	•	•			
7-Pesquisa Orientada							X	X	X	X	X	X	•	•	•	•	•	•	•	•	•	•	•	•
8-Apresentação da qualificação													•											
9-Escrita de artigos																		•	•					
10-Escrita da dissertação										X	X	X					•	•	•	•	•	•	•	•
11- Defesa da dissertação																								•

Para a conclusão da atividade 1, o Programa de Pós-Graduação em Ciência da Computação (PGCOMP) da Universidade Federal da Bahia (UFBA) exige que um mestrando obtenha um total de 18 créditos em disciplinas. Nos semestres 2018.1 e 2018.2, foram obtidos os créditos exigidos ao cursar as disciplinas: MATD74 – Algoritmos e Grafos, MATE32 – Tópicos em Inteligência Computacional II, MATE64 – Seminários Científicos, MATE65 – Fundamentos de Pesquisa em Ciência da Computação I, MATE10 – Tópicos em Inteligência Computacional I e MATE84 – Tópicos em Fundamentos da Computação IV.

A segunda atividade planejada neste cronograma foi realizada a partir da disciplina MATE65 – Fundamentos de Pesquisa em Ciência da Computação I e nos meses subsequentes. Além disso, durante a execução desta tarefa foi realizada a prova de proficiência em inglês.

As atividades 3 e 4 do cronograma consistem na realização dos experimentos e análise dos resultados. Essas atividades foram divididas em duas partes. A primeira contém apenas experimentos preliminares que foram realizados para validar esta proposta de trabalho. Nesta parte, conjuntos de dados sintéticos foram analisados conforme apresentado no Capítulo ???. A segunda parte dos experimentos e suas análises serão realizadas após a qualificação.

As atividades 5 e 8 estão relacionadas com o componente curricular MATD75 – Exame de qualificação. A atividade 5 refere-se à escrita deste texto e as atividades 6 e 7 representam os componentes curriculares MATA32 – Estágio Docente e MATA31 – Pesquisa Orientada, respectivamente. A atividade 8 está relacionada à apresentação desta qualificação de mestrado.

A escrita de artigos, listada no item 9 do cronograma, será realizada com base nos resultados gerados com os experimentos (Atividades 3 e 4) e nas contribuições obtidas com a apresentação da qualificação. Por fim, como requisito para a defesa de dissertação, fica pendente a atividade MATE93 - Defesa de Proposta de Mestrado, a qual se refere aos itens 10 e 11 da Tabela ??.

### **3.4 CONSIDERAÇÕES FINAIS**

Neste capítulo, apresentou-se de maneira detalhada o projeto de pesquisa, o plano de atividades e o cronograma planejado para a conclusão do mestrado. No próximo capítulo, serão discutidos os experimentos iniciais, realizados com o objetivo de analisar a viabilidade da proposta de mestrado.





## EXPERIMENTOS INICIAIS

### 4.1 CONSIDERAÇÕES INICIAIS

Este capítulo apresenta um conjunto de experimentos preliminares realizados com dados sintéticos, cujo objetivo foi demonstrar a viabilidade da proposta de mestrado. Os resultados obtidos se mostraram promissores, indicando que o tema de pesquisa deve continuar a ser investigado. A próxima seção apresenta como os dados sintéticos foram produzidos para condução dos experimentos.

### 4.2 CONFIGURAÇÃO DOS EXPERIMENTOS

Os fluxos de dados sintéticos foram produzidos através de classes geradoras do framework MOA e salvos em arquivos no formato ARFF. Cada fluxo produzido representa até dois conceitos e é composto por 2.500 observações, com valores entre 0 e 1. O restante desta seção apresenta as classes geradoras utilizadas, as modificações realizadas e os parâmetros aplicados.

A classe geradora *AbruptChangeGenerator* produz fluxos de dados sintéticos com mudanças de conceito abruptas. Os fluxos gerados simulam as mudanças através da intercalação de sequências com o valor 0.2 e sequências com o valor 0.8, conforme o tamanho definido para os conceitos. Com o objetivo de tornar os dados produzidos mais próximos da realidade, essa classe foi alterada<sup>1</sup> para permitir a adição de um ruído randômico, com amplitude configurável, para cada exemplo produzido. A versão com suporte a ruído foi utilizada nos experimentos para produção do fluxo com mudanças abruptas. O gerador foi parametrizado para produzir exemplos não binários, com conceitos formados com 400 instâncias e ruído limitado ao intervalo  $[-0.1, 0.2]$  seguindo uma distribuição uniforme.

A classe *GradualChangeGenerator* produz fluxos sintéticos com mudanças de conceito graduais. Durante sua configuração, percebeu-se uma limitação, pois a classe só gera uma mudança gradual, independente do tamanho dos conceitos. Para superar esta limitação, a

---

<sup>1</sup><https://git.io/fjGEj>

classe foi modificada<sup>2</sup> para gerar uma quantidade de mudanças coerente com o tamanho definido. A classe modificada foi utilizada para gerar o fluxo sintético com mudanças graduais utilizado nos experimentos, sendo parametrizada para produzir exemplos não binários e conceitos compostos por 400 instâncias.

Por fim, a classe *NoChangeGenerator*, que produz fluxos sem ocorrência de mudanças de conceito, também foi modificada<sup>3</sup> para permitir a incidência de ruído nos resultados gerados. Esta versão foi utilizada para produzir um fluxo sintético sem mudanças, com ruído entre  $[-0.1, 0.1]$ .

### 4.3 CRITÉRIOS DE AVALIAÇÃO

A classe de avaliação *BasicConceptDriftPerformanceEvaluator*, pertencente ao framework MOA, foi utilizada para avaliar o método proposto neste trabalho. Esta classe permite mensurar a acurácia e o desempenho de algoritmos para detecção de mudanças de conceito. Para utilizá-la, é necessário construir uma tarefa do tipo *EvaluateConceptDrift* através da aba *Concept Drift*, presente na tela inicial da aplicação. Para condução dos experimentos deste trabalho, a tarefa *EvaluateConceptDrift* foi configurada conforme descrito na Tabela ??.

**Tabela 4.1** Configuração da classe *BasicConceptDriftPerformanceEvaluator*

Parâmetro	Valor	Observação
learner	ChangeDetectorLearner	O algoritmo de detecção de mudanças de conceito a ser testado é definido no atributo <i>driftDetectionMethod</i> da classe <i>ChangeDetectorLearner</i> .
stream	ARFFFileStream	Caminho para um dos arquivos <i>ARFF</i> descrito na seção anterior. O atributo <i>classIndex</i> deve ser definido como 0, pois não existem rótulos nestes conjuntos de dados.
instanceLimit	-1	Desabilita o limite de instâncias a serem processadas.
timeLimit	-1	Desabilita o limite de tempo de execução.
sampleFrequency	1	Uma linha de resultado do avaliador deve ser gerada para cada instância processada.

<sup>2</sup><https://git.io/fjGue>

<sup>3</sup><https://git.io/fjGu3>

A classe de avaliação utilizada produz diversos indicadores referentes a acurácia e performance do algoritmo executado. A Tabela ?? apresenta os indicadores analisados neste trabalho.

**Tabela 4.2** Indicadores analisados

Indicador	Observação
Tempo de Processamento	Tempo médio (seg.) de processamento por instância.
Mudanças Existentes	Quantidade de mudanças existentes.
Mudanças Detectadas	Quantidade de mudanças detectadas corretamente.
Falso-positivos	Quantidade de mudanças detectadas erroneamente.
Atraso de Detecção	Quantidade média de instâncias até a detecção.

#### 4.4 Pettitt

Durante a revisão da literatura, investigou-se a aplicabilidade das técnicas para detecção de *changing points* em problemas de detecção de mudança de conceito. Para tanto, foi realizada a implementação, no framework MOA<sup>4</sup>, de um novo algoritmo baseado no método estatístico proposto por ??).

O método de Pettitt é um teste não paramétrico, no qual se verifica se duas amostras  $Y_1, \dots, Y_t$  e  $Y_{t+1}, \dots, Y_T$  são da mesma população. A estatística  $U_{t,T}$  faz uma contagem do número de vezes que um membro da primeira amostra é maior que um membro da segunda. A estatística  $U_{t,T}$  é obtida através da Equação ??:

$$U_{t,T} = U_{t-1,T} + \sum_{j=1}^T \text{sgn}(Y_t - Y_j) \quad (4.1)$$

Onde  $t = 2, \dots, T$ ;  $\text{sgn}(x) = 1$  para  $x > 0$ ;  $\text{sgn}(x) = 0$  para  $x = 0$  e  $\text{sgn}(x) = -1$  para  $x < 0$ . A estatística  $U_{t,T}$  é então calculada para os valores de  $1 \leq t \leq T$  e a estatística  $k(t)$  do teste de Pettitt é o máximo valor absoluto de  $U_{t,T}$ .

Os testes foram executados com os conjuntos de dados sintéticos produzidos na Seção ?. Os resultados obtidos são apresentados a seguir, na Tabela ?.

**Tabela 4.3** Resultados - Método de Pettitt

Conjunto de Dados	Tempo de processamento	Mudanças Existentes	Mudanças Detectadas	Falso-positivos	Atraso de Detecção
Sem mudanças	3.19	0	0	4	–
Mudanças Abruptas	2.09	6	3	2	132
Mudanças Incrementais	1.48	6	4	2	133

<sup>4</sup><https://git.io/fj8xu>

A análise dos resultados demonstra que o algoritmo se mostrou propenso à produção de falsos-positivo e, em comparação aos outros algoritmos abordados neste trabalho, computacionalmente ineficiente.

Por não apresentar resultados competitivos em relação ao estado da arte, esta linha de pesquisa foi suspensa, permitindo que os esforços se voltassem para a análise da aplicação de Redes de Função de Base Radial para detecção de mudanças de conceito.

#### 4.5 RBFDriftDetector

Esta seção apresenta um conjunto de experimentos realizados com o objetivo de validar o algoritmo desenvolvido nesta pesquisa, denominado RBFDriftDetector. Embora os experimentos ainda não sejam suficientes para comprovar a hipótese proposta, a análise empírica aqui apresentada evidencia a importância de analisar a aplicação de Redes de Função de Base Radial para detecção de mudanças de conceito em fluxos contínuos de dados.

Para execução dos experimentos, foram utilizados três conjuntos de dados sintéticos, conforme descrito na seção ???. Além disso, os seguintes algoritmos foram utilizados para comparação: CUSUM, PageHinkley e ADWIN. A Tabela ?? lista os parâmetros utilizados para cada algoritmo.

**Tabela 4.4** Parâmetros utilizados para cada algoritmo

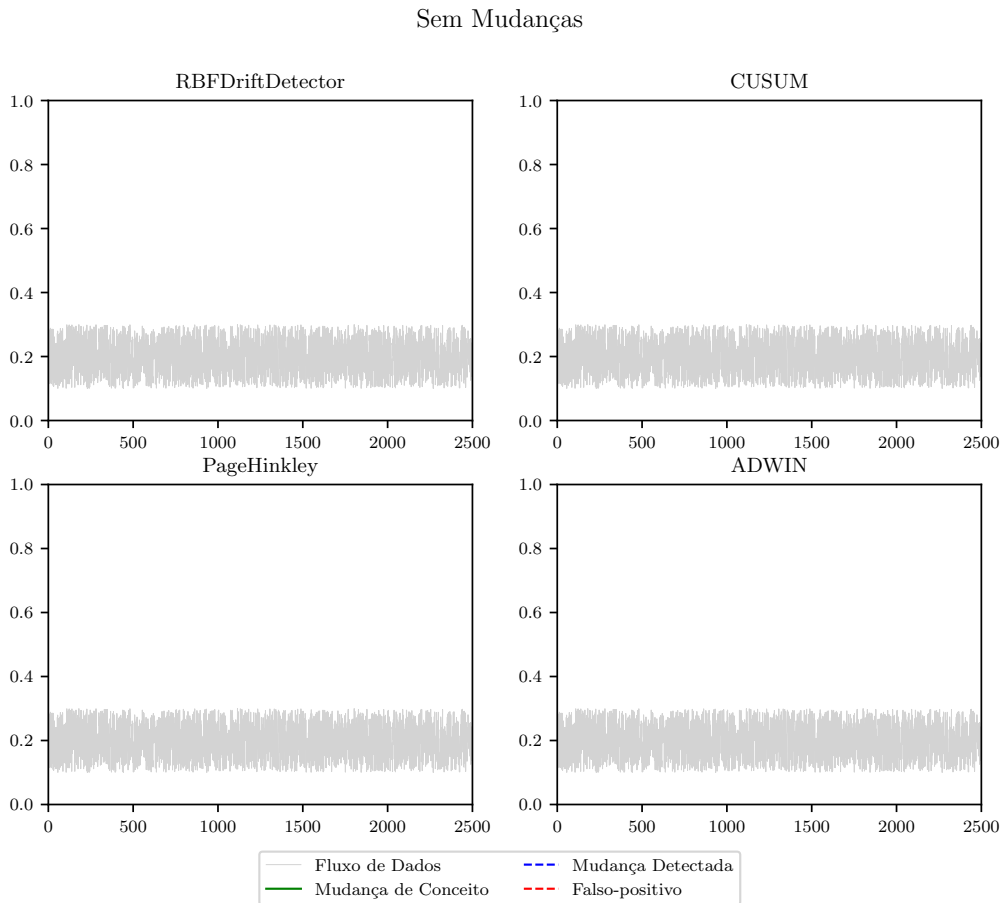
Algoritmo	Parâmetros
RBFDriftDetector	$\sigma = 2; \lambda = 0.5$
CUSUM	$MinNumInstances = 30; \delta = 0.005; \lambda = 50$
PageHinkley	$MinNumInstances = 30; \delta = 0.005; \lambda = 50; \alpha = 1$
ADWIN	$\delta = 0.002$

O primeiro experimento utilizou o fluxo de dados sintético sem mudanças de conceito, a fim de avaliar a tendência de produção de falsos-positivo. O resultado desta análise pode ser visto na Tabela ??.

**Tabela 4.5** Experimento 1 - Fluxo sem mudanças de conceito

Algoritmo	Tempo de processamento	Mudanças Existentes	Mudanças Detectadas	Falso-positivos	Atraso de Detecção
RBFDriftDetector	0.22	0	0	0	—
CUSUM	0.31	0	0	0	—
PageHinkley	0.24	0	0	0	—
ADWIN	0.21	0	0	0	—

Como pode ser observado, todos algoritmos testados demonstraram tolerância a ruídos e não indicaram nenhum falso positivo. Quanto a performance, o algoritmo proposto obteve a segunda melhor média em tempo de processamento. Superado pelo algoritmo ADWIN, por uma pequena margem. O comportamento dos algoritmos e do conjunto de dados utilizado é apresentado graficamente na Figura ??.



**Figura 4.1** Representação Gráfica - Sem mudanças de conceito

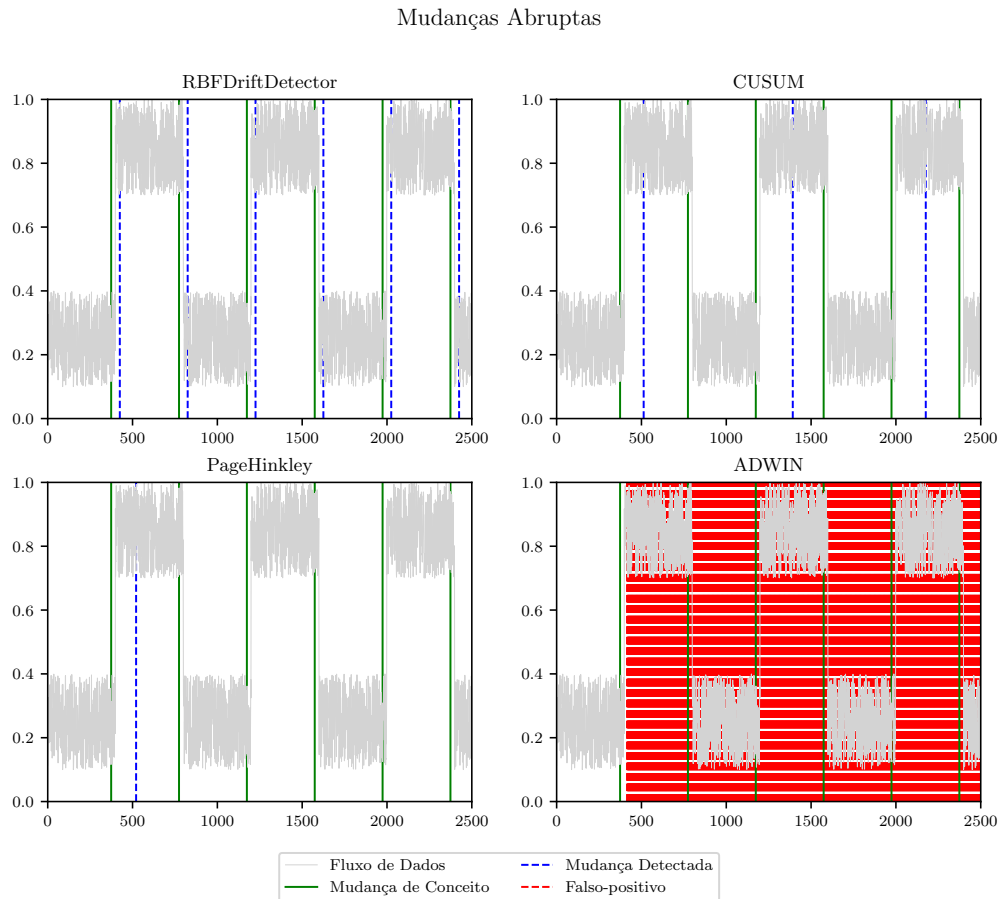
Em seguida, utilizou-se a base de dados com mudanças abruptas como base de testes. Os resultados obtidos são apresentados na Tabela ??.

**Tabela 4.6** Experimento 2 - Fluxo com mudanças de conceito abruptas

Algoritmo	Tempo de processamento	Mudanças Existentes	Mudanças Detectadas	Falso-positivos	Atraso de Detecção
RBFDriftDetector	0.23	6	6	0	1
CUSUM	0.29	6	3	0	68
PageHinkley	0.22	6	1	0	17
ADWIN	0.21	6	2052	2046	9

Ao analisar este resultado, verifica-se que o algoritmo proposto neste trabalho identificou todas as mudanças de conceito existentes no fluxo sem indicar nenhum falso-positivo, sendo também o algoritmo com menor atraso de detecção. No quesito performance, a técnica desenvolvida obteve o terceiro melhor resultado. Paralelamente, os outros algoritmos analisados apresentaram baixa acurácia. O algoritmo CUSUM detectou apenas metade das mudanças e apresentou a maior taxa de atraso. Enquanto o método PageHinkley detectou apenas uma mudança. Por fim, o algoritmo ADWIN apresentou maior eficiência computacional e apesar de detectar as 6 mudanças existentes, mostrou-se hipersensível, pois foram detectados 2046 falsos-positivos.

Os resultados do segundo experimento são demonstrados graficamente na Figura ?? . Nesta imagem, as linhas sólidas na cor verde representam as mudanças de conceito existentes no conjunto de dados analisado. As linhas tracejadas azuis representam mudanças de conceito identificadas corretamente, enquanto que as linhas tracejadas de cor vermelha representam falsos-positivos.



**Figura 4.2** Representação Gráfica - Mudanças Abruptas

O último experimento realizado utilizou o fluxo sintético com mudanças de conceito graduais. Seu resultado é demonstrado na Tabela ??.

**Tabela 4.7** Experimento 3 - Fluxo com mudanças de conceito graduais

Algoritmo	Tempo de processamento	Mudanças Existentes	Mudanças Detectadas	Falso-positivos	Atraso de Detecção
RBFDriftDetector	0.24	6	5	1	171
CUSUM	0.65	6	3	0	32
PageHinkley	0.26	6	1	0	4
ADWIN	0.27	6	2244	2238	1

Neste experimento, o algoritmo RBFDriftDetector identificou 5 das 6 mudanças existentes e sinalizou um falso-positivo. Ainda assim, obteve a melhor acurácia, apesar de apresentar a maior taxa de atraso. Os outros algoritmos analisados apresentaram comportamentos similares aos do experimento anterior. O algoritmo CUSUM identificou metade das mudanças. O método PageHinkley identificou apenas uma alteração de conceito. E, finalmente, o algoritmo ADWIN voltou a se mostrar muito sensível, ao identificar 2238 falsos-positivos.

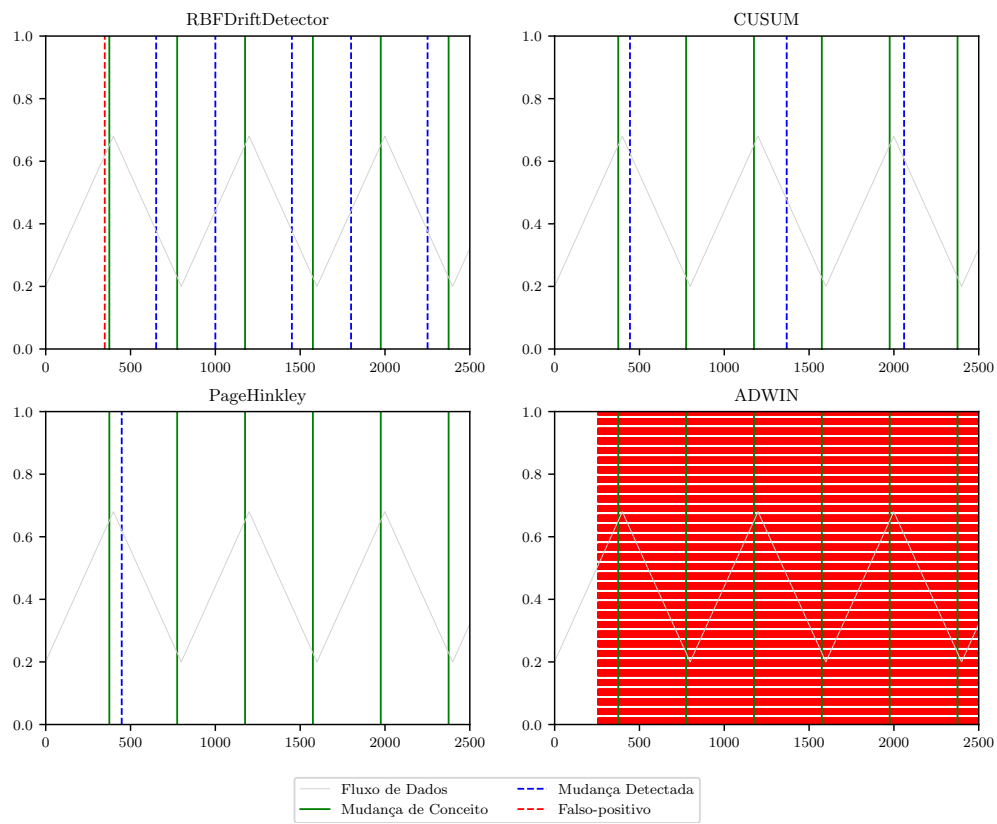
Por fim, o comportamento do conjunto de dados e dos algoritmos neste último experimento é apresentado na Figura ?. Na imagem, as linhas sólidas na cor verde representam as mudanças de conceito existentes no conjunto de dados analisado. As linhas tracejadas azuis representam mudanças de conceito identificadas corretamente, enquanto que as linhas tracejadas de cor vermelha representam falsos-positivos.

## 4.6 CONSIDERAÇÕES FINAIS

Nesta seção foram apresentados os resultados dos experimentos iniciais deste trabalho de pesquisa, os quais visavam confirmar a aplicabilidade das Redes de Função de Base Radial na detecção de mudanças de conceito em fluxos contínuos de dados.

As seguintes etapas estão previstas para serem executadas: i) realizar experimentos com fluxos recorrentes e incrementais; ii) integrar uma cadeia de Markov ao método proposto, permitindo a análise do comportamento das mudanças e o aprimoramento da acurácia das detecções; iii) implementar os algoritmos OLINDDA, MINAS e DETECNOD para serem comparados; iv) implementar e validar o algoritmo no framework Tornado; e v) alterar o algoritmo para permitir que o centro se desloque dentro do grupo formado. Espera-se ainda realizar uma aplicação prática da abordagem proposta, utilizando-a em um conjunto de dados oriundo de um sistema do mundo real.

### Mudanças Graduais



**Figura 4.3** Representação Gráfica - Mudanças Graduais