

UNIVERSIDADE FEDERAL DE SANTA CATARINA
CENTRO TECNOLÓGICO
DEPARTAMENTO DE ENGENHARIA ELÉTRICA E ELETRÔNICA
CURSO DE GRADUAÇÃO EM ENGENHARIA ELÉTRICA

Gustavo Henrique Angeoleti Lopes

**Predição do Preço de Liquidação das Diferenças do Submercado Sul a partir
de Variáveis Hidrológicas**

Florianópolis
2022

Gustavo Henrique Angeoleti Lopes

**Predição do Preço de Liquidação das Diferenças do Submercado Sul a partir
de Variáveis Hidrológicas**

Trabalho de Conclusão de Curso de Graduação em Engenharia Elétrica do Centro Tecnológico da Universidade Federal de Santa Catarina para a obtenção do título de Bacharel em Engenharia Elétrica.
Orientador: Prof. Eduardo Luiz Ortiz Batista, Dr.

Florianópolis

2022

Ficha de identificação da obra elaborada pelo autor,
através do Programa de Geração Automática da Biblioteca Universitária da UFSC.

Lopes, Gustavo Henrique Angeoleti
Predição do Preço de Liquidação das Diferenças do
Submercado Sul a partir de Variáveis Hidrológicas / Gustavo
Henrique Angeoleti Lopes ; orientador, Eduardo Luiz Ortiz
Batista, 2022.
75 p.

Trabalho de Conclusão de Curso (graduação) -
Universidade Federal de Santa Catarina, Centro Tecnológico,
Graduação em Engenharia Elétrica, Florianópolis, 2022.

Inclui referências.

1. Engenharia Elétrica. 2. PLD Horário. 3. Aprendizado
de Máquina. 4. Séries Temporais. I. Batista, Eduardo Luiz
Ortiz. II. Universidade Federal de Santa Catarina.
Graduação em Engenharia Elétrica. III. Título.

Gustavo Henrique Angeoleti Lopes

Predição do Preço de Liquidação das Diferenças do Submercado Sul a partir de Variáveis Hidrológicas

Este Trabalho Conclusão de Curso foi julgado adequado para obtenção do Título de Bacharel em Engenharia Elétrica e aprovado em sua forma final pelo Curso de Engenharia Elétrica.

Florianópolis, 14 de abril de 2022.



Documento assinado digitalmente
Miguel Moreto
Data: 20/04/2022 16:18:59-0300
CPF: 948.850.100-63
Verifique as assinaturas em <https://v.ufsc.br>

Prof. Miguel Moreto, Dr.
Coordenador do Curso

Banca Examinadora:



Documento assinado digitalmente
Eduardo Luiz Ortiz Batista
Data: 20/04/2022 09:31:26-0300
CPF: 036.521.889-85
Verifique as assinaturas em <https://v.ufsc.br>

Prof. Eduardo Luiz Ortiz Batista, Dr.
Orientador
Universidade Federal de Santa Catarina



Documento assinado digitalmente
Erlon Cristian Finardi
Data: 20/04/2022 10:23:08-0300
CPF: 020.364.749-18
Verifique as assinaturas em <https://v.ufsc.br>

Prof. Erlon Cristian Finardi, Dr.
Avaliador
Universidade Federal de Santa Catarina

Enga. Isla Almeida Oliveira
Avaliadora
Fundação CERTI

AGRADECIMENTOS

À minha família, pelo apoio e companheirismo ao longo dos últimos anos.

Ao meu orientador Eduardo, pelo suporte, sugestões, paciência e amizade.

À Isla e ao Erlon, por terem aceitado contribuir com este trabalho.

A todos os professores que tive a sorte de cruzar ao longo dos últimos 20 anos.

A todos que produzem conhecimento aberto.

Ao Joinville Esporte Clube, por me tornar habituado às derrotas, tão comuns ao longo da graduação em Engenharia Elétrica.

“A new scientific truth does not triumph by convincing its opponents and making them see the light, but rather because its opponents eventually die, and a new generation grows up that is familiar with it.”
(Max Planck)

RESUMO

A partir dos anos 1990, o mercado de energia elétrica brasileiro passou por um intenso período de abertura econômica, tornando o ambiente mais competitivo e, em certo aspecto, eficiente. Nesse contexto, uma mudança notável que ocorreu ainda mais recentemente é a implementação do Preço de Liquidação das Diferenças (PLD) horário. Tradicionalmente, o PLD define, a partir de modelos de otimização, o preço da energia no mercado de curto prazo. Com a mudança, passou-se a ser divulgado o PLD com discretização horária, permitindo maior liquidez e dinamismo ao mercado. Assim, diversas novas oportunidades começam a poder ser exploradas, em geral visando à redução de exposição a riscos que tal parâmetro, agora mais volátil, pode causar. Nesse contexto, este trabalho propõe utilizar modelos de aprendizado de máquina para prever, com uma semana de antecedência, o valor do PLD. Para fazê-lo, são implementados diversos algoritmos, que utilizam como dados de entrada séries temporais históricas de variáveis hídricas, com a expectativa de obter um valor acurado do preço futuro da energia no curto prazo. Adicionalmente, o trabalho desenvolve um breve estudo de caso, que visa avaliar um eventual ganho financeiro que um agente do mercado poderia ter, caso em posse de um modelo de predição. Com os resultados obtidos, é possível observar que, não apenas alguns modelos conseguiram atingir um valor satisfatório de predição, como também o estudo de caso indicou que há uma significativa redução na exposição a riscos a partir de modelos com predição.

Palavras-chave: PLD horário, aprendizado de máquina, séries temporais

ABSTRACT

Since the 1990s, the Brazilian electricity market has been experiencing an intense period of economic opening, making its environment more competitive and, in a certain way, more efficient. In this context, a notable change that has taken place in recent years is the implementation of the hourly energy spot price. Thus, that price, which defines, based on optimization models, the price of energy in the short-term market, started to be defined hourly, allowing greater liquidity and dynamism to the market. Thus, several new market opportunities can be explored, generally aiming to reduce exposure to the risks that the new and more volatile spot price can cause. In this context, this work proposes using machine learning models to predict the value of the spot price one week in advance. To this end, several models are implemented, using time series of water variables as input data, hoping to obtain an accurate value of the future energy price. Additionally, the work develops a brief case study to evaluate an eventual financial gain that a market agent could have, if it had a forecast model. With the acquired results, it is possible to observe that some models were able to reach a satisfactory prediction value, but the the case study indicated that there is a significantly lower exposure to risks of models with predictions similar to the one obtained.

Keywords: energy spot price, machine learning, time series

LISTA DE FIGURAS

Figura 1 – Estrutura organizacional de alguns dos principais agentes do setor elétrico	27
Figura 2 – Fluxograma comparando os dois ambientes de comercialização de energia elétrica no Brasil	30
Figura 3 – Balanço energético a ser contabilizado e liquidado no Mercado de Curto Prazo (MCP)	32
Figura 4 – Contratos de Comercialização de Energia no Ambiente de Contratação Livre (CCEAL) com modulação	32
Figura 5 – CCEAL <i>flat</i>	33
Figura 6 – Fluxograma ilustrando o problema do despacho em sistemas hidrotérmicos	34
Figura 7 – Funções de Custo Imediato e Futuro da água	35
Figura 8 – Ávore de decisão simplificada	44
Figura 9 – Funcionamento simplificado de um algoritmo baseado em <i>boosting</i> . . .	45
Figura 10 – Funcionamento simplificado de um algoritmo baseado em <i>bagging</i> . . .	46
Figura 11 – Funcionamento simplificado de uma Redes Neurais Artificiais (RNA) .	47
Figura 12 – Fluxo de trabalho definido para desenvolvimento deste trabalho	51
Figura 13 – Valor do PLD no Submercado Sul em R\$/MWh	54
Figura 14 – Valor da Energia Armazenada (EAR) na Bacia do Iguaçu em MWmês .	55
Figura 15 – Valor da Energia Natural Afluente (ENA) na Bacia do Uruguai em MWmês	55
Figura 16 – Precipitação diária observada em cada ponto no dia 15 de agosto de 2020	56
Figura 17 – Precipitação diária observada total em todos os pontos em mm	56
Figura 18 – Estrutura simplificada dos modelos	57
Figura 19 – Dinâmica de tabularização dos dados para treinamento	59
Figura 20 – Comparação entre o PLD real e o predito pelo modelo de regressão linear	61
Figura 21 – Comparação entre o PLD real e o predito pelo modelo de floresta aleatória	62
Figura 22 – Comparação entre o PLD real e o predito pelo modelo de <i>Gradient Boosting</i>	62
Figura 23 – Comparação entre o PLD real e o predito pelo modelo de Redes Neurais	63
Figura 24 – Diagrama de blocos do cálculo de economia do estudo de caso	64
Figura 25 – Distribuição dos horários considerados ótimos para início da produção no estudo de caso	66

LISTA DE TABELAS

Tabela 1 – Resumo dos dados adquiridos	53
Tabela 2 – Resumo do desempenho de cada modelo	63

LISTA DE ABREVIATURAS E SIGLAS

ACL	Ambiente de Contratação Livre
ACR	Ambiente de Contratação Regulada
AED	Análise Exploratória de Dados
ANEEL	Agência Nacional de Energia Elétrica
ANP	Agência Nacional do Petróleo
CCEAL	Contratos de Comercialização de Energia no Ambiente de Contratação Livre
CCEE	Câmara de Comercialização de Energia Elétrica
CGCE	Câmara da Gestão da Crise de Energia Elétrica
CMO	Custo Marginal de Operação
CMSE	Comitê de Monitoramento do Setor Elétrico
EAR	Energia Armazenada
ENA	Energia Natural Afluente
EPE	Empresa de Pesquisa Energética
FCF	Função de Custo Futuro
FCI	Função de Custo Imediato
MAE	<i>Mean Absolut Error</i>
MAPE	<i>Mean Absolute Percentage Error</i>
MCP	Mercado de Curto Prazo
MME	Ministério de Minas e Energia
ONS	Operador Nacional do Sistema Elétrico
PCH	Pequenas Centrais Hidrelétricas
PD	Pesquisa e Desenvolvimento
PDO	Programação Diária da Operação
PLD	Preço de Liquidação das Diferenças
PMAE	Preço do Mercado Atacadista de Energia Elétrica
PO	Precipitação Diária Observada
RMSE	<i>Root Mean Squared Error</i>
RNA	Redes Neurais Artificiais
RSS	<i>Residual Sum of Squares</i>
SEB	Setor Elétrico Brasileiro
SIN	Sistema Interligado Nacional

SUMÁRIO

1	INTRODUÇÃO	21
1.1	OBJETIVOS	22
1.1.1	Objetivo Geral	22
1.1.2	Objetivos Específicos	22
2	FUNDAMENTAÇÃO TEÓRICA	25
2.1	MERCADO DE ENERGIA ELÉTRICA BRASILEIRO	25
2.1.1	Contexto Histórico	25
2.1.2	Modelo Atual	26
2.1.2.1	Agentes do mercado	27
2.1.2.2	Comercialização de Energia Elétrica	29
2.1.3	Mercado Livre de Energia	30
2.1.3.1	Contratos de energia no mercado livre	31
2.1.4	Planejamento da Operação	33
2.1.4.1	Problema do despacho em sistemas hidrotérmicos	33
2.1.4.2	Modelos de otimização e formação de preços	35
2.1.4.3	PLD	37
2.2	APRENDIZADO DE MÁQUINA	37
2.2.1	Notação	38
2.2.2	Fundamentos	38
2.2.3	Conceitualização	39
2.2.3.1	Conjunto de dados de treinamento e teste	39
2.2.3.2	Aprendizado supervisionado e não supervisionado	39
2.2.3.3	Classificação e regressão	40
2.2.3.4	Função perda	40
2.2.3.5	Validação cruzada	41
2.2.3.6	Hiperparâmetros	41
2.2.4	Modelos	41
2.2.4.1	Regressão Linear	42
2.2.4.2	Métodos Baseados em Árvore	43
2.2.4.3	Métodos de <i>Ensemble</i>	44
2.2.4.4	Redes Neurais Artificiais	47
2.2.5	Métricas de avaliação	48
3	DESENVOLVIMENTO	51
3.1	AQUISIÇÃO DOS DADOS	52
3.2	DETERMINAÇÃO DAS VARIÁVEIS	53
3.3	TREINAMENTO	57
4	RESULTADOS	61

4.1	AVALIAÇÃO DOS MODELOS	61
4.2	ESTUDO DE CASO	63
5	CONCLUSÃO	67
	REFERÊNCIAS	71

1 INTRODUÇÃO

Desde os anos 1990, o setor de energia elétrica no Brasil passou por diversas reestruturações e novas modelagens, que envolveram longos processos de desestatização e abertura de mercado. Esses processos contribuíram para uma maior eficiência do setor e conseqüentemente o crescimento econômico do país (SCHOUCHANA, 2010).

Outra consequência de tais mudanças foi um aumento da competitividade no mercado de energia, acelerado pela criação do então chamado Novo Modelo do Setor Elétrico, através da Lei nº 10.848 de 2004, em que foram estabelecidos conceitos relevantes como a modicidade tarifária, garantia de suprimento e estabilidade do marco regulatório (BRASIL, 2004). Todos esses movimentos contribuíram para a criação de um ambiente competitivo e favorável a novos investimentos (SCHOUCHANA, 2010).

Ainda no contexto das mudanças mencionadas, destaca-se também a criação de diversos agentes, responsáveis por regular o mercado, formular as novas políticas, planejar a expansão do setor, coordenar a operação, desenvolver novos estudos, entre outras atribuições. Com a inserção dessas entidades, o sistema até os dias atuais vem se desenvolvendo, tornando-se cada vez mais competitivo e com a presença de novos agentes do setor privado.

A partir disso, é relevante a discussão acerca do Preço de Liquidação das Diferenças (PLD), valor que determina o preço vigente de toda a energia elétrica que foi produzida, mas ainda não contratada pelos agentes do mercado (CCEE, 2021a). Esse parâmetro, definido pela Câmara de Comercialização de Energia Elétrica (CCEE) a partir de diversos modelos computacionais de otimização, tem a relevante função de definir o preço que será pago pela energia no Mercado de Curto Prazo (MCP) para liquidar a diferença entre os contratos estabelecidos e a carga consumida real, tornando o mercado mais dinâmico e balanceando a relação entre oferta e demanda (SANTOS, 2019).

Dessa forma, novamente no contexto das mudanças que vem permitindo ao Setor Elétrico Brasileiro (SEB) rápido avanço, o PLD passou, a partir de 1º de janeiro de 2021, a ser calculado com base nos resultados de execução do modelo DESSEM, apresentado com mais profundidade na Seção 2.2.4 (CCEE, 2021a). Aqui, cabe pontuar a relevância dessa mudança, uma vez que permitiu ao PLD ter uma discretização horária, a qual até então era semanal. Essa transição permitiu ainda maior dinamismo ao mercado, uma vez que, sendo este valor definido atualmente para cada hora do dia, há uma maior volatilidade e, portanto, maior espaço para operações mais rentáveis por parte dos agentes envolvidos.

Como breves exemplos das possíveis vantagens que os consumidores têm a partir dessa mudança, pode-se citar o caso de uma empresa que opta por replanejar sua operação para horários cujo preço de energia é mais baixo, tornando os custos atrelados à produção mais econômicos. Ademais, podem ser citadas novas oportunidades do mercado, considerando o atual cenário, em que há, cada vez mais, abertura e demanda de instrumentos financeiros que podem oferecer proteção para mitigar os riscos atrelados à volatilidade dos preços

(ENGIE, 2021).

Em contrapartida, é importante citar também a complexidade com a qual o PLD é definido, uma vez que a matriz energética brasileira é predominantemente hídrica. Sendo assim, os modelos que geram o PLD precisam considerar relações complexas, como o benefício entre o uso imediato da água e o uso futuro da água. Ainda, os modelos precisam considerar outras variáveis, como preço de combustível para usinas térmicas, aspectos técnicos dos equipamentos de transmissão e geração, demandas de energia, novos projetos, dentre outros (MEDEIROS, 2004).

Considerando todos os pontos mencionados, Medeiros (2004) destaca a relevância do desenvolvimento de modelos de previsão como um dos grandes desafios da indústria de energia elétrica, com o intuito de oferecer melhores ferramentas aos tomadores de decisões para definição das suas estratégias, tanto em questão de comercialização e consumo, como também na avaliação de novos investimentos no setor. Assim, o tópico abordado ao longo deste trabalho é de grande relevância, considerando, em especial, o interesse dos agentes do mercado de energia por modelos de predição do preço de energia no MCP e as recentes mudanças na dinâmica do mesmo.

1.1 OBJETIVOS

Diante dos pontos expostos na Seção anterior, é possível estabelecer os objetivos gerais e específicos deste trabalho, os quais serão expostos a seguir.

1.1.1 Objetivo Geral

O objetivo do presente trabalho é, a partir de dados históricos das variáveis hídricas que impactam o SEB, desenvolver modelos de predição do PLD, visando à maior acurácia possível. Para tal, objetiva-se utilizar modelos de aprendizado de máquina para obter uma predição do PLD do submercado Sul para um período posterior de sete dias, com discretização horária. Em relação às variáveis hídricas, a abordagem será detalhada na Seção 3; porém, de forma geral, são considerados os dados que impactam os modelos de otimização, especificamente aqueles que definem os volumes dos reservatórios e os relacionados à precipitação. A partir disso, espera-se utilizar os resultados obtidos para comparar o desempenho de cada modelo de aprendizado de máquina, especificamente em questão de proximidade entre os valores preditos e os valores reais.

1.1.2 Objetivos Específicos

Além do objetivo geral, o trabalho envolve alguns objetivos específicos. Entre eles, cabe citar as expectativas de

- a) estudar as dinâmicas e as recentes mudanças do mercado de energia brasileiro;

-
- b) estudar com profundidade os modelos de aprendizado de máquina, bem como avaliar quais são e como influenciam os hiperparâmetros;
 - c) desenvolver modelos que façam a predição do PLD com um prazo de sete dias de antecedência;
 - d) discutir a eficácia de cada modelos de aprendizado da máquina no contexto do problema estudado;
 - e) quantificar a influência das variáveis hidrológicas utilizadas como entrada na formação do PLD;
 - f) avaliar os potenciais ganhos que um agente do mercado de energia pode ter, caso em posse de um modelo de predição como estes a serem desenvolvidos.

2 FUNDAMENTAÇÃO TEÓRICA

2.1 MERCADO DE ENERGIA ELÉTRICA BRASILEIRO

Muito em razão do recente e acelerado desenvolvimento tecnológico do setor elétrico, bem como pelas mudanças regulatórias que o país vem passando nos últimos anos, o sistema elétrico do Brasil está em constante evolução (ANEEL, 2008). Nesse contexto, esta seção visa tratar brevemente do contexto histórico e da estrutura do setor, buscando detalhar as questões relativas ao mercado de energia elétrica, aspecto fundamental para o entendimento das potenciais aplicações deste trabalho, para enfim desenvolver os conceitos que impactam diretamente no PLD.

2.1.1 Contexto Histórico

Destacam-se, dentro do contexto dos avanços no setor elétrico observados ao longo dos anos 1990, duas relevantes mudanças ligadas diretamente ao modelo institucional: a primeira é relativa à privatização das companhias geradoras na segunda metade da década de 90, enquanto a segunda, mais recente, se deu com a introdução do Novo Modelo do Setor de Energia Elétrica, em 2004 (ANEEL, 2008).

A respeito desse processo de desestatização, tem-se como relevante o marco inicial da liberalização do setor elétrico, feita em 1993, através dos Decretos nº 915 e nº 1.009 (BRASIL, 1993b) (BRASIL, 1993a). A partir deles, o então presidente Itamar Franco decretou, entre outros pontos, a autorização de formação de consórcios entre concessionárias e autoprodutores para exploração de aproveitamentos, permitindo livre acesso à malha de transmissão. Adicionalmente, em 1995, a Lei nº 8.987 definiu o novo regime de concessões e permissões dos serviços públicos, enquanto a Lei nº 9.074 criou a figura do produtor independente de energia e do consumidor livre, estabelecendo as normas para facilitar a privatização das empresas do setor elétrico (FARIAS, 2006).

De forma geral, tais mudanças foram a base para uma relevante abertura do setor no país. Neste contexto, nota-se também uma intensificação dos mecanismos competitivos do setor de energia elétrica a partir das privatizações e desverticalizações, que se sucederam com estas mudanças regulatórias. Cabe também citar as então recém iniciadas interações entre agentes de diferentes regiões, permitida com o novo modelo, que colocou fim à reserva geográfica de mercado, estabelecida até então (SILVA, 2006).

Na direção contrária, o evento que ficou conhecido como apagão do setor elétrico, ocorrido em maio de 2001, após um intenso período de seca, e que resultou em uma drástica redução da geração hídrica, predominante na matriz energética nacional, fez com que o então Presidente Fernando Henrique Cardoso tomasse a iniciativa de fazer diversas mudanças. Entre elas, houve a criação da Comissão de Análise do Sistema Hidrotérmico de Energia Elétrica e da Câmara da Gestão da Crise de Energia Elétrica (CGCE), ambas

na mesma semana, com o intuito de mapear as possíveis causas do apagão.

Após alguns meses de trabalho, foi divulgado um relatório final, indicando pontos relevantes acerca das motivações da crise recém ocorrida. Segundo Landi (2011), os resultados indicaram que a crise do apagão se devia, em grande parte, às mudanças estruturais implementadas no processo de reestruturação do setor elétrico, não apenas em virtude das questões climáticas, motivo que acentuou a necessidade por reformas no setor.

A partir disso, tem-se a implementação do chamado Novo Modelo do Setor de Energia Elétrica, que teve suas bases definidas pela Lei nº 10.848 de 2004, e permitiu inovações para o produtor independente de energia, agente que detém a possibilidade de produzir a própria energia, e para o comercializador de energia. Também, a recém criada Agência Nacional de Energia Elétrica (ANEEL) fica definida como o órgão responsável pela autorização de atuação de tais agentes (BRASIL, 2004). Além disso, destacam-se, como transformações do anterior modelo, a criação da CCEE, com o intuito de substituir o Mercado Atacadista de Energia, e a Empresa de Pesquisa Energética (EPE), responsável por desenhar a expansão do setor elétrico.

Ainda como resultado deste conjunto de mudanças, Brito (2010) destaca a Lei nº 10.848/2004 como relevante instrumento prático para práticas que garantissem a expansão da oferta de energia elétrica. Destas, cabe citar as definições de que toda demanda dos agentes deve estar contratada, que todo contratado deve ser respaldado por capacidade firme de geração e que toda contratação das distribuidoras deve ser realizada por meio de leilões, medidas que visavam aumentar a segurança do fornecimento de energia elétrica no Brasil.

Adicionalmente, tem-se, conforme indicado pelo trabalho de Silva (2006), uma significativa mudança que resultou na definição de um novo modelo de comercialização de energia, a partir de diversos instrumentos normativos assinados à época. Com eles, tem-se, então, formas detalhadas de: a) negociação e contratação de energia, já segmentadas pelos dois ambientes de regulação - o Ambiente de Contratação Regulada (ACR) e o Ambiente de Contratação Livre (ACL); b) regras para outorga de concessões; e c) regras para os leilões de energia elétrica. Tais mudanças são fundamentais para a consolidação do mercado como hoje ele o é.

2.1.2 Modelo Atual

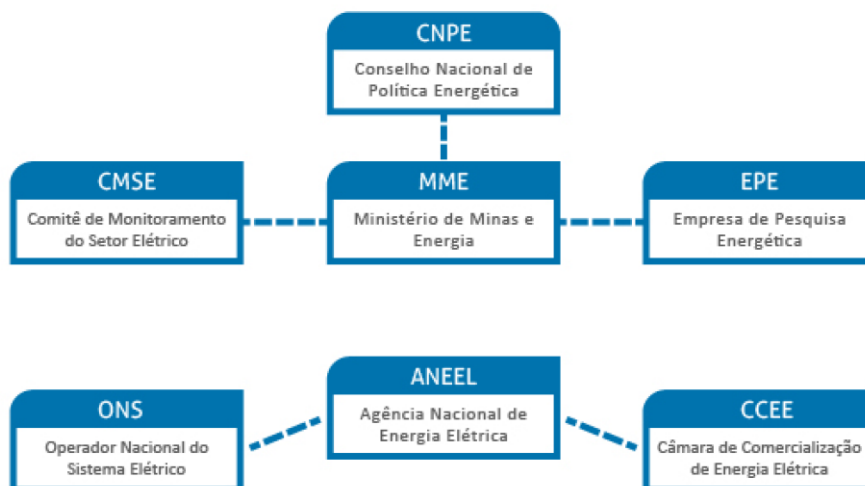
Todas as mudanças discutidas na seção 2.1.1 tornaram o mercado de energia elétrica brasileiro atual robusto, sob supervisão de uma rigorosa regulação, contando ainda com a presença de inúmeros agentes, complexos modelos de otimização e uma estrutura dinâmica de contratação de energia. Esta seção tem como objetivo discutir os principais fatores que, de alguma forma, influenciam este mercado, em especial no que tange sua precificação.

2.1.2.1 Agentes do mercado

Dentro da já mencionada complexidade do atual modelo vigente do mercado de energia, chama a atenção o elevado número de agentes atuantes, sejam eles reguladores ou os participantes propriamente ditos, que interagem entre si de diversas maneiras.

Com o intuito de compreender melhor a dinâmica atual, é relevante observar de forma detalhada o papel de cada um destes no modelo em vigor. Assim, tem-se, na Figura 1 um organograma que relaciona alguns destes agentes. Dentre eles, alguns até já citados na Seção 2.1.1, cabe destacar os apresentados a seguir.

Figura 1 – Estrutura organizacional de alguns dos principais agentes do setor elétrico



Fonte: (CCEE, 2021b)

- *MME*

O Ministério de Minas e Energia (MME) tem, entre outras funções, a incumbência de formular e assegurar a execução das políticas necessárias para a gestão sustentável dos recursos energéticos do país (FEDERAL, s.d.). Assim, no contexto do sistema elétrico nacional, o órgão é responsável por formular, planejar e implementar todas as ações do Governo Federal que tangem a política de energia do país.

- *ANEEL*

A autarquia, vinculada por regime especial ao MME, foi criada, conforme descrito na Seção 2.1.1, no contexto das mudanças relevantes do setor elétrico, com o intuito de fazer sua regulação no Brasil.

Com suas atividades iniciadas apenas no final de 1997, a ANEEL tem como principais atribuições a) regular a geração, transmissão, distribuição e comercialização de energia elétrica; b) fiscalizar as concessões, as permissões e os serviços de energia elétrica; c) implementar políticas e diretrizes definidas para exploração de energia elétrica definidas pelo governo federal; d) estabelecer as tarifas; e) dirimir as divergências; e f) promover as atividades de outorgas de concessão (ANEEL, s.d.).

- *ONS*

O Operador Nacional do Sistema Elétrico (ONS) foi mais um dos órgãos criado no mesmo contexto das reformas do setor elétrico. Sua atuação envolve a coordenação e controle da operação das instalações de geração e transmissão de energia elétrica do Sistema Interligado Nacional (SIN), bem como o planejamento dos sistemas isolados do país, sob fiscalização da ANEEL.

Como objetivos do ONS tem-se a) promover a otimização da operação do sistema eletroenergético; b) garantir acesso não discriminatório de todos os agentes à rede de transmissão; e c) contribuir para a expansão do SIN sob as melhores condições (ONS, s.d.). Estes pontos são fundamentais para a compreensão da formação dos preços da energia elétrica e, por consequência, o escopo principal deste trabalho descrito na Seção 3.

- *EPE*

A EPE visa ao desenvolvimento de estudos e pesquisas voltadas a suportar o planejamento do setor energético, sendo uma prestadora de serviços do MME. Caracteriza-se por ser uma empresa pública federal, e, portanto, é dependente do Orçamento Geral da União.

Sua criação, também pela Lei nº 10.847 de 2004, deu-se com o intuito de resgatar a responsabilidade constitucional do Estado em assegurar um desenvolvimento sustentável da infraestrutura energética do país. Atualmente, o órgão consolida-se como parte fundamental das atividades de definição de políticas e diretrizes, que, posteriormente, materializam-se em estudos e pesquisas que orientam o desenvolvimento do setor de energia no país (EPE, s.d.).

- *CCEE*

A CCEE é o órgão responsável por tornar possível a comercialização de energia no SIN, mediante supervisão da ANEEL, que a regula e fiscaliza. Conforme será descrito na Seção 3, a CCEE tem como uma das mais relevantes atribuições a realização dos leilões públicos de energia no ACR, também mediante delegação da ANEEL (CCEE, s.d.).

Ademais, cabe à CCEE o registro de todos os contratos de comercialização do ACR, do ACL e os de contratações de ajustes, bem como a contabilização e liquidação das transações de curto prazo.

- *CMSE*

O Comitê de Monitoramento do Setor Elétrico (CMSE), conforme descrito na Seção 2.1.1, foi criado com a função de acompanhar e avaliar a continuidade e segurança do fornecimento e suprimento de energia elétrica em todo o Brasil.

A partir do decreto 5.174/2004, o comitê é presidido pelo Ministro de Estado de Minas e Energia, sendo composto, ainda, por outros quatro representantes do MME, além dos titulares indicados pela ANEEL, Agência Nacional do Petróleo (ANP), CCEE, EPE e ONS.

- *CEPEL*

O CEPEL é um centro vinculado às empresas Eletrobras que realiza pesquisas na área de engenharia elétrica. O órgão realiza projetos de Pesquisa e Desenvolvimento (PD), além de serviços tecnológicos e laboratoriais especializados, além de suporte ao MME e outras entidades.

No conceito dos tópicos apresentados neste trabalho, o CEPEL tem um papel fundamental no desenvolvimento dos modelos de otimização que fazem o planejamento de operação do SIN, conforme descrito na Seção 2.1.4.1.

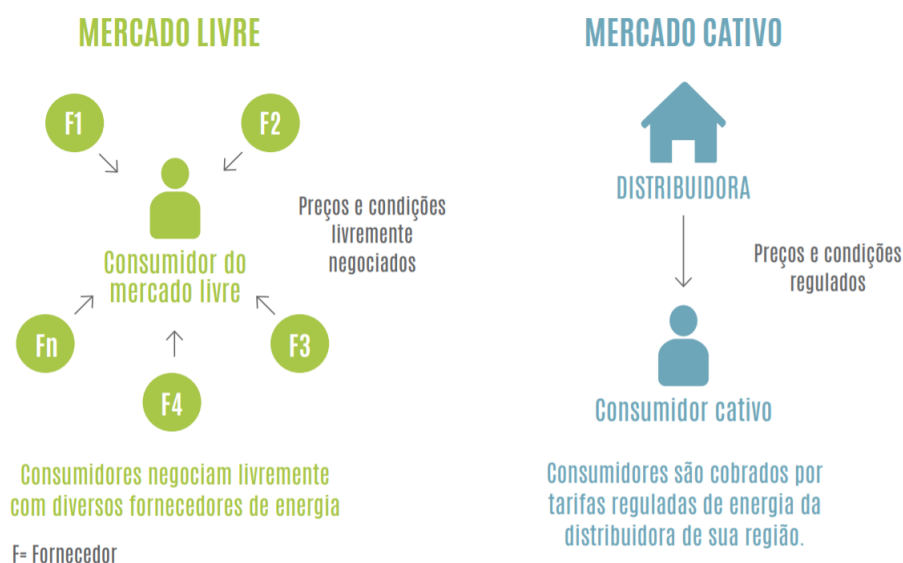
2.1.2.2 Comercialização de Energia Elétrica

A comercialização de energia elétrica no contexto brasileiro acontece de duas formas: livre ou a partir da definição ou limitação de preços e quantidades pelo Poder Público (ANEEL, 2018). Ambas as formas são operacionalizadas pela CCEE no âmbito do SIN, sendo a ANEEL o órgão que estabelece os regulamentos a serem seguidos.

Estas duas formas distintas de se comercializar energia elétrica no Brasil deram origem aos dois ambientes de contratação hoje em vigor no país, também consequência do novo modelo instituído em 2014: o ACL e o ACR, também conhecido como mercado cativo. Enquanto no primeiro são realizadas operações de compra e venda de energia a partir de contratos bilaterais livremente negociados, conforme as regras anteriormente citadas, no segundo estas operações entre agentes são negociadas a partir de leilões de energia. A Figura 2 apresenta um fluxograma simplificado, que indica a diferença básica entre ambos os ambientes.

Conforme analisado pelo trabalho de Brito (2010), tem-se que a forma de comercialização definida neste novo modelo priorizou, essencialmente, a contratação regulada de energia, com prazos previamente estabelecidos e os valores de compras definidos a partir de leilões estruturados, cuja energia é proveniente de empreendimentos de geração, podendo ser existentes ou novos. Assim, tem-se um cenário no qual a atuação dos agentes é bem estabelecida, sendo obrigatória a participação das distribuidoras no ACR, com a

Figura 2 – Fluxograma comparando os dois ambientes de comercialização de energia elétrica no Brasil



Fonte: (ABRACEEL, 2019)

possibilidade de atuação dos geradores em ambos os ambientes, enquanto os consumidores livres têm a permissão de atuar apenas no ACL.

Além disso, foi definida a permanência no MCP da contabilização e da liquidação das diferenças no ambiente livre. Assim, tem-se no MCP o segmento no qual a diferença entre energia elétrica gerada e de fato consumida é contabilizada e liquidada a PLD, que substitui o Preço do Mercado Atacadista de Energia Elétrica (PMAE) (BRITO, 2010).

Dessa forma, levando em conta que o PLD, conceito fundamental deste trabalho, é uma ferramenta que está inserida apenas no ACL, cabe um olhar detalhado sobre a dinâmica deste ambiente livre.

2.1.3 Mercado Livre de Energia

O ACL, também conhecido como o mercado livre de energia elétrica, é o ambiente no qual os consumidores têm a possibilidade de escolher de forma livre os seus fornecedores de energia, negociando, entre ambos, as condições de contratação.

Hoje, cerca de 60% da energia que as indústrias consomem no Brasil são adquiridas através deste ambiente, o que proporciona, desde 2013, uma economia média de 18% em comparação com o mercado cativo (ABRACEEL, 2019).

Em relação aos potenciais consumidores deste ambiente, tem-se que há dois possíveis tipos: os consumidores livres e os especiais. Os primeiros, cuja demanda deve ser de até 3.000 kW, podem contratar energia gerada por qualquer fonte, enquanto os consumidores especiais, limitados a contratação de energia proveniente das fontes chamadas especiais,

devem ter uma demanda entre 500 e 3.000 kW. Estas fontes especiais de energia elétrica incluem usinas eólicas, solares, a biomassa, Pequenas Centrais Hidrelétricas (PCH) ou hidráulica, desde que o empreendimento tenha potência de até 50.000 kW.

Outro aspecto importante a se discutir relativo ao mercado livre de energia é a compreensão de quem são os agentes dos quais se pode comprar energia. A respeito disso, a definição do ACL indica que a energia elétrica pode ser disponibilizada através de agentes comercializadores, importadores, autoprodutores (apenas excedentes da autogeração), geradores, além dos próprios consumidores livres e especiais, a partir de cessão de contratos, contanto que estes estejam cadastrados como agentes na CCEE (ABRACEEL, 2019).

2.1.3.1 Contratos de energia no mercado livre

Conforme discutido na Seção 2.1.2.1, cabe à CCEE o registro dos contratos de energia elétrica firmados no mercado livre (CCEE, s.d.). Assim, é através deste órgão que são concretizados, no SIN, todos os processos de comercialização de energia elétrica no ACL. Para este ambiente, os contratos são conhecidos como CCEAL.

Dentro do ACL, ainda, alguns pontos relevantes devem ser respeitados quando firmados os contratos. Nas negociações, por exemplo, é imprescindível que os agentes apresentem 100% de lastro para a venda de energia, materializado pela garantia física, seja por geração própria ou não. Nesse sentido, cabe à CCEE também a aplicação de penalidades, caso este lastro não seja comprovado, mediante regras da própria câmara (SANTOS, 2019).

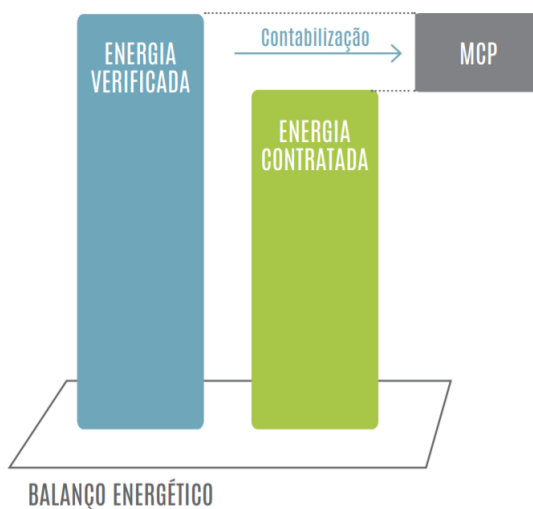
Ademais, um ponto fundamental para compreensão das análises deste trabalho é relativo ao acerto das diferenças. Considerando que não há uma ligação direta entre os contratos de compra e venda com geração e consumo, mesmo havendo tais lastros, pode acontecer de uma usina ter um valor de geração superior ou inferior ao que estava previsto no contrato. Analogamente, a parte do contrato que fez a aquisição de energia pode também consumir uma quantidade diferente daquela estabelecida. Para liquidação destas diferenças, existe o MCP (ABRACEEL, 2019).

Neste cenário, entra novamente em ação a CCEE, que faz a comparação entre os valores estabelecidos pelos contratos, e os valores posteriormente medidos como gerados e consumidos. Havendo esta diferença, os agentes devem utilizar o MCP para contabilizar e liquidar estes montantes. A Figura 3 apresenta de forma simplificada essa situação.

Além disso, é de fundamental importância para o escopo deste trabalho pontuar que este valor excedido dos contratos é sempre negociado no MCP ao valor vigente do PLD, preço determinado a partir dos modelos apresentados na Seção 2.2.4.

Ainda, o estabelecimento dos CCEAL também deve contemplar a definição dos montantes e vigência. Para isso, os agentes envolvidos devem definir o volume de energia contratado em MW médio para uma determinada vigência, a partir de um perfil de entrega acordado. Essa definição pode envolver um contrato com uma única vigência, em que os

Figura 3 – Balanço energético a ser contabilizado e liquidado no MCP

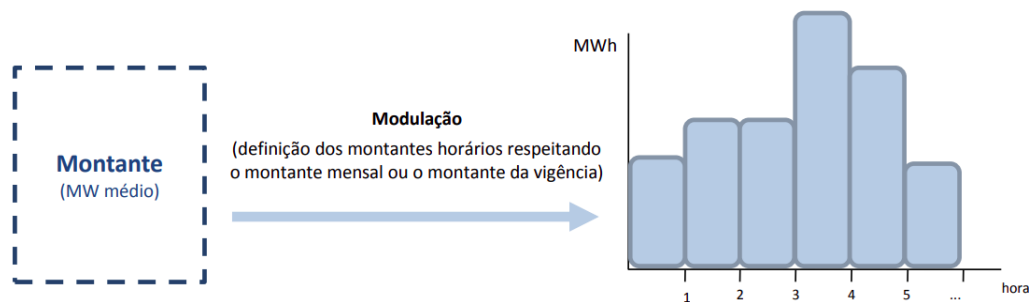


Fonte: (ABRACEEL, 2019)

montantes são fixos, ou variando os montantes para cada vigência (CCEE, 2017).

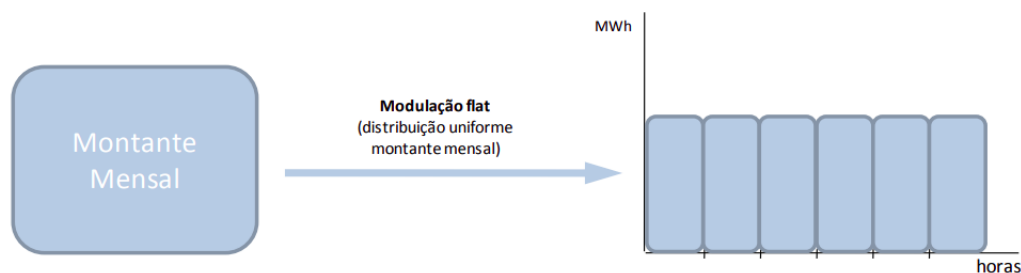
Adicionalmente, os agentes podem ou não optar por firmar um contrato com ou sem modulação. Se o fizerem, os agentes distribuem os montantes os valores de energia do contrato em base horária. Uma representação simplificada dessa dinâmica pode ser observada na Figura 4.

Figura 4 – CCEAL com modulação



Fonte: (CCEE, 2017)

Em contrapartida, caso os agentes não optem por estabelecer a modulação do contrato, tem-se o chamado contrato *flat*. Nesse caso, o CCEAL é modulado automaticamente de forma uniforme, dividindo-se o total de energia para cada hora da vigência do contrato, conforme apresentado na Figura 5.

Figura 5 – CCEAL *flat*

Fonte: (CCEE, 2017)

2.1.4 Planejamento da Operação

Conforme já descrito anteriormente, uma das incumbências do ONS é realizar o planejamento da operação eletroenergética. Este se dá através de estudos que avaliam as condições futuras de suprimento de energia, a partir de critérios de otimização, estudos de recomposição, reserva e segurança operativa e análises de continuidade do suprimento eletroenergético (ONS, s.d.).

Assim, este planejamento da operação, que ocorre para diversos horizontes de tempo, tem como entrada principal as condições hidrológicas e os possíveis cenários de cargas futuros, naturalmente todos expostos a altos níveis de incerteza. Dessa forma, espera-se, como resultado, a definição de estratégias adequadas para melhor utilização dos recursos disponíveis, subsidiando a programação da operação do SIN e a pré-operação para elaboração dos despachos, visando atender o mercado com segurança e custo minimizado, a partir dos recursos disponíveis.

Nesse contexto, tem-se que há dois possíveis modelos para esse despacho centralizado: o *Tight Pool*, adotado no Brasil, e o *Loose Pool*. Cada modelo é definido de acordo com o nível controle da entidade independente.

Em geral, esta dinâmica adotada no Brasil costuma ser implementada também em países no qual a matriz energética é formada majoritariamente por usinas hidrelétricas, e, por consequência, cuja previsibilidade de geração é menor, caso oposto ao de países como Colômbia, Inglaterra e Itália (SCHOUCHANA, 2010).

2.1.4.1 Problema do despacho em sistemas hidrotérmicos

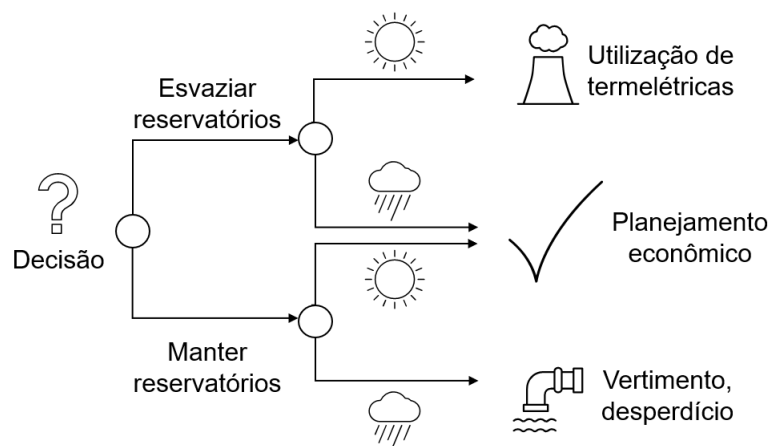
A partir dessa contextualização do planejamento da operação, é relevante discutir os modelos que, de fato, determinam essa operação. Contudo, dada a complexidade envolvida em um sistema predominantemente hidrotérmico, como é o caso do SIN, alguns conceitos iniciais devem ser discutidos previamente.

De forma geral, esta complexidade se dá, uma vez que, diferentemente de um sis-

tema puramente térmico, a geração hidráulica tem, a cada decisão tomada na fase de planejamento, uma consequência futura e um risco associado, uma vez que a afluência futura é uma variável de difícil previsibilidade, além de envolver um limite na capacidade dos reservatórios das usinas (MEDEIROS, 2004).

Para ilustrar tal complexidade, Medeiros (2004) ilustra algumas potenciais consequências das tomadas de decisão em situações comuns no planejamento da operação. No primeiro cenário, toma-se a decisão de esvaziar os reservatórios do sistema, tendo dois possíveis desdobramentos: (i) as aflúncias futuras são altas, e o resultado é uma operação econômica; e (ii) as aflúncias futuras são baixas, e, com os reservatórios vazios, há uma operação deficitária ou a necessidade de geração térmica, cujos custos são muito maiores. Já no segundo cenário, tem-se a decisão inversa, a de manter os reservatórios cheios. As consequências agora envolvem (i) novamente aflúncias futuras altas, necessitando verter água, desperdiçando um potencial de geração, não resultando no cenário mais otimizado; e (ii) as aflúncias vêm abaixo do esperado, mas, como os reservatórios estavam cheios, tem-se, aqui, uma operação econômica. A Figura 6 ilustra este exemplo apresentado.

Figura 6 – Fluxograma ilustrando o problema do despacho em sistemas hidrotérmicos

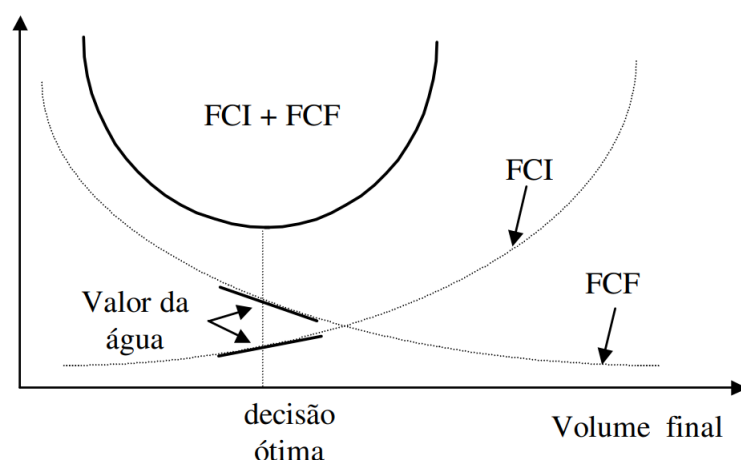


Fonte: Próprio autor a partir de Medeiros (2004)

A complexidade deste cenário aumenta à medida que são consideradas as inúmeras bacias envolvidas no SIN, além do efeito cascata, ou seja, usinas cuja disponibilidade de água dependem da tomada de decisão de outras usinas. Dessa forma, é evidente que há um *tradeoff* entre a utilização imediata da geração hidroelétrica com o benefício de armazenamento do reservatório ao longo do período de planejamento (SANTOS, 2019). A forma analítica de resolver este problema envolve a Função de Custo Futuro (FCF) e a Função de Custo Imediato (FCI) da água, ilustradas na Figura 7.

A interpretação da FCF indicada na Figura 7 consiste em compreender que, quanto maior é o volume armazenado final, menor é o custo futuro, pois haverá uma maior

Figura 7 – Funções de Custo Imediato e Futuro da água



Fonte: (SIMÕES; GOMES, 2017)

disponibilidade de água, independente das afluências. Já para o FCI, entende-se que um volume final muito alto acarretará em um custo imediato também elevado, uma vez que será necessária a geração térmica para suprir a demanda no momento presente (MEDEIROS, 2004). O uso otimizado da água será, para cada caso, aquele em que o somatório do FCF e do FCI é minimizado, isto é, o ponto de derivada zero na parábola indicada na Figura 7. Este é um dos pontos fundamentais a ser resolvido pelos modelos de otimização que fazem o planejamento da operação do SIN.

2.1.4.2 Modelos de otimização e formação de preços

Novamente trazendo para o contexto do mercado brasileiro, o Custo Marginal de Operação (CMO) é definido pelo ONS para cada um dos quatro submercados que compõem o SIN - Norte, Nordeste, Sul e Sudeste/Centro-Oeste - a partir dos modelos que consideram tais conceitos discutidos na Seção 2.1.4.1. Assim, tem-se o CMO como resultado de saída dos modelos de otimização desenvolvidos pelos agentes citados na Seção 2.1.2.1.

Neste contexto, é evidente a relevância de se discutir alguns destes modelos de otimização que resultam no CMO e, por consequência, no PLD.

- *NEWAVE*

O NEWAVE é primeiro modelo da cadeia, desenvolvido pelo CEPEL com o intuito de otimizar o planejamento de operação de sistemas hidrotérmicos no longo e médio prazo, uma vez que o seu horizonte de planejamento é de 5 anos. Seu objetivo básico é calcular a política de operação que estima os valores de água armazenada nos reservatórios, além de indicar as metas de geração de cada usina mensalmente, com o intuito de atender à

demanda e minimizar os custos esperados da operação, a partir de um critério de risco pré-determinado (CEPEL, 2018b).

É importante pontuar que, por ser um modelo que visa um horizonte de tempo mais longo, algumas simplificações se fazem necessárias, como o conceito de reservatório equivalente. Para períodos de estudo mais longo, esse conceito é mais amplamente utilizado, ao passo que, conforme os estudos focam em planejamento de mais curto prazo, há uma individualização das usinas.

Adicionalmente, o NEWAVE, discretizado mensalmente, utiliza 20 cenários hidrológicos para cada período, fazendo com que, no total dos 5 anos, tenha-se 10^{78} . Assim, mesmo que o algoritmo implementado no NEWAVE seja baseado em conceitos de programação dinâmica dual estocástica, não são percorridos todos os conjuntos da árvore, apenas um subconjunto de cenários, escolhidos da distribuição original da variável aleatória (CEPEL, 2018b).

Além dos cenários de afliências, outras variáveis são utilizadas como entrada do NEWAVE, como previsão de demanda por energia, previsão dos preços dos combustíveis utilizados nas termelétricas, custos de novas usinas, entre outros (SANTOS, 2019). Ainda, é relevante que o NEWAVE tem um papel fundamental na cadeia de modelos, uma vez que é o responsável por definir a FCF a ser utilizada nos demais.

- *DECOMP*

Posterior ao NEWAVE, o DECOMP é o segundo modelo da cadeia de programas. O seu objetivo é resolver o problema de planejamento de operação no curto prazo, desagregando, para cada reservatório, agora individualizado, as FCF resultantes da etapa anterior (MEDEIROS, 2004).

Além de estar acoplado ao NEWAVE e receber como entrada algumas de suas saídas, o DECOMP também tem como características gerais os cenários de afliências, os patamares de carga, os limites de interligação entre os sistemas, os contratos de exportação e importação de energia, além de restrições elétricas e de transmissão, características particulares de usinas hídricas e termelétricas, entre outros (CEPEL, 2018a).

- *DESSEM*

O DESSEM é o modelo responsável pela Programação Diária da Operação (PDO) do sistema. Para isso, trabalha com técnicas e ferramentas que podem modelar e resolver o problema de otimização diária da operação do SIN, considerando não apenas aspectos da rede elétrica, mas também das operações das usinas do sistema, buscando fazê-lo da forma mais acurada possível (CEPEL, 2018a).

Análoga à situação anterior, aqui também o DESSEM é acoplado ao DECOMP, através de um conversor de dados chamado DECODESS. Além disso, é importante destacar

a relevância do DESSEM no contexto deste trabalho, uma vez que o CMO é uma de suas saídas.

Mais recente dos modelos, o DESSEM passou por uma longa etapa de validações, e, recentemente em janeiro de 2020, foi implementado pelo ONS, para realizar o planejamento do despacho hidrotérmico no curtíssimo prazo, etapa fundamental para que para o início do PLD horário, implementado no início de 2021.

2.1.4.3 PLD

Finalmente, a partir de toda a discussão teórica apresentada anteriormente, é possível introduzir o conceito de PLD, também conhecido como preço *spot*. Este parâmetro, obtido diretamente do CMO, aplicadas algumas restrições de limite superior e inferior, é utilizado pela CCEE para liquidar e contabilizar as diferenças dos contratos no MCP, conforme apresentado na Seção 2.1.3.1.

Ainda, a partir do que já foi apresentado na Seção sobre o DESSEM, o CMO passou recentemente a ser calculado com discretização semi-horária. Dessa forma, com uma granularidade muito maior de dados, também o PLD passou a ter seu valor apresentado a cada hora, mudança relevante para trazer dinamismo ao mercado livre de energia. Além disso, essa mudança é uma das principais motivações deste trabalho, uma vez que, com uma disponibilidade maior de dados, é possível o desenvolvimento de modelos de predição com maior acurácia.

2.2 APRENDIZADO DE MÁQUINA

Antes um nicho limitado dentro dos departamentos de computação, engenharia e matemática, a base de usuários de aprendizado de máquina vem crescendo de forma excepcional nas últimas décadas, incluindo agora pesquisadores e estudantes, tanto no meio acadêmico, como também em diversos segmentos do mercado (WATT; BORHANI; KAT-SAGGELOS, 2016). De forma geral, é importante notar que, com um olhar simplificado, esses modelos visam à identificação de padrões a partir de um conjunto de dados. De acordo com Bishop (2006), ambos o reconhecimento de padrões e o aprendizado de máquina podem ser vistos como duas facetas de uma mesma área de pesquisa, mesmo que este tenha tido um desenvolvimento muito mais recente, enquanto aquele é um problema com diversos registros ao longo da história. Neste sentido, esta seção é dedicada a introduzir conceitos básicos para a compreensão dos fundamentos do aprendizado de máquina, para que, com essa base consolidada, seja possível iniciar a discussão dos modelos utilizados na Seção 3, visando ao entendimento da fundamentação matemática dos algoritmos, bem como o significado e o impacto dos parâmetros envolvidos.

2.2.1 Notação

Embora ocorrida essa consolidação dos conceitos de aprendizado de máquina na última década, ainda há uma certa divergência em questão de padronização da notação utilizada, muito em razão de ser um campo de pesquisa bastante abrangente. Nesse sentido, este trabalho seguirá o padrão descrito e utilizado por Bishop (2006).

Assim, segue como notação deste trabalho que os vetores são representados por algarismos romanos, como x , sendo também todos os vetores considerados vetores coluna. Além disso, utiliza-se o sobrescrito T para indicar o transposto. Por fim, algarismos romanos escritos em maiúsculo e negrito representam matrizes, como em \mathbf{M} , enquanto a notação (w_1, \dots, w_M) indica um vetor linha que contém M elementos.

Além disso, por conveniência, padroniza-se nesta seção e ao longo do restante do trabalho a variável n para indicar o número de variáveis preditoras de um modelo e m como o número de amostras de cada conjunto de dados. O contexto dessas nomenclaturas é apresentado ao longo da Seção 2.2.2.

2.2.2 Fundamentos

De forma geral, tem-se dois tipos principais de variáveis nos modelos de aprendizado de máquina. As variáveis de entrada, muitas vezes chamadas também variáveis independentes, *features* ou preditoras, costumam ser denotadas pelo símbolo \mathbf{X} , junto a um subscrito para distinção, já que os modelos costumam utilizar diversas variáveis de entrada. Além destas, há também as variáveis de saída, ocasionalmente chamadas de variáveis dependentes ou *targets*, usualmente denotadas por y (JAMES *et al.*, 2013). O conceito fundamental é que, para cada entrada ou observação do conjunto de dados, haverá um vetor de variáveis de entrada e, eventualmente, um valor da variável de saída correspondente, conforme discutido com mais profundidade na Seção 2.2.3.2.

Generalizando a relação entre elas, tem-se a observação de uma variável quantitativa y , feita a partir de n diferentes *features*, \mathbf{X} . Assume-se que a relação pode ser genericamente expressada como

$$y = f(\mathbf{X}) + \epsilon, \quad (1)$$

onde ϵ é o termo de erro aleatório desta equação, enquanto f representa a informação sistemática que \mathbf{X} provê acerca de y .

No contexto deste trabalho, em que se busca trabalhar com modelos de predição, esta mesma equação pode ser generalizada, assumindo que a média do termo ϵ tende a zero, o que resulta em

$$\hat{y} = \hat{f}(\mathbf{X}), \quad (2)$$

onde \hat{f} indica a estimativa de f , enquanto \hat{y} representa a predição resultante de y . Nesta configuração, e especialmente para algoritmos de aprendizado de máquina, \hat{f} é muitas vezes tratado como uma caixa preta, no sentido de que não se tem, tipicamente, o intuito de saber qual a sua exata função, contanto que sua aplicação gere predições acuradas de y (JAMES *et al.*, 2013).

Outra característica consequente de (1) e (2), discutida em James *et al.* (2013) e fundamental também para o escopo deste trabalho, é o interesse que se pode ter no entendimento da forma com a qual y é afetada a partir das mudanças em (X_1, \dots, X_M) . Conforme citado na Seção 1, mesmo que o objetivo geral do trabalho seja, de fato, a predição da variável de saída, também espera-se compreender, por exemplo, de que forma ela é impactada pelas variáveis de entrada. Dessa forma, tem-se a possibilidade de não apenas predizer os valores do *target*, como também discutir aspectos mais qualitativos do modelo.

2.2.3 Conceitualização

A partir desta fundamentação inicial, é possível discutir outros conceitos relevantes do campo de aprendizado de máquinas.

2.2.3.1 Conjunto de dados de treinamento e teste

Há duas fases distintas no processo de execução do modelo de aprendizado de máquina: a fase de treinamento e a de testes. A partir disso, também é relevante pontuar que é utilizado um conjunto de dados diferente para cada uma delas. Assim, há, durante a fase de pré-processamento, a necessidade de dividir o conjunto de dados disponíveis, formando, com isso, o conjunto de treino e de teste.

A fase de treino, também conhecida como etapa de aprendizado, é aquela na qual os parâmetros das funções que modelam \hat{f} são determinados. A forma como isso é aprendido pelo algoritmo é melhor detalhada na Seção 2.2.4, uma vez que o processo difere a depender do modelo. Posteriormente, há a etapa de teste, na qual se utiliza novos exemplos de dados de entrada, para categorizar novos exemplos, a partir das variáveis de entrada, almejando-se o processo conhecido como generalização (BISHOP, 2006).

2.2.3.2 Aprendizado supervisionado e não supervisionado

Ainda relativo ao conjunto de dados, as aplicações podem ser classificadas como de aprendizado supervisionado e não supervisionado. Naquelas em que o conjunto de treinamento contém os vetores de entrada junto aos seus respectivos vetores *target*, tem-se os modelos de aprendizado supervisionado.

Embora esta seja a abordagem mais tradicional, e que segue, de forma geral, os fundamentos discutidos na Seção 2.2.2, há também situações na qual o conjunto de

treinamento contém apenas as variáveis de entrada, sem valores de saída correspondente. Nesse caso, o objetivo dos modelos de aprendizado não supervisionados é fazer a chamada clusterização, ou seja, agrupar cada um dos dados em grupos a partir dos valores das suas variáveis de entrada. Também estas aplicações podem ter como objetivo determinar a distribuição dos dados dentro de um espaço de entrada, ou até mesmo projetar os dados em algum espaço de alta dimensionalidade (BISHOP, 2006).

De forma complementar, há também a técnica de aprendizado por reforço, em que o objetivo é encontrar decisões adequadas para cada situação a partir de uma situação em que se almeja maximizar uma dada recompensa (BISHOP, 2006). Como o escopo deste trabalho se limita a aplicações de aprendizado supervisionado, a Seção 2.2.4 focará nestas abordagens, não aprofundando as demais.

2.2.3.3 Classificação e regressão

Outra relevante segmentação entre os modelos de aprendizado de máquina é dado pelo problema que se almeja resolver. Embora ambos sejam de certa forma semelhantes, Watt, Borhani e Katsaggelos (2016) indica que a diferença chave é que os modelos de classificação têm como objetivo a predição de valores discretos ou das chamadas classes. Como exemplo clássico, tem-se que um problema no qual se objetiva distinguir imagens de algarismos decimais escritos à mão, cujas classes, naturalmente, podem assumir 10 valores discretos.

Em contrapartida, os problemas de regressão não possuem essa limitação, podendo, portanto, assumir valores contínuos dentro de um domínio. Trazendo este conceito para o escopo deste trabalho, um problema de predição de preços, no qual estes, definidos como *target*, podem assumir quaisquer valores, é um exemplo de regressão. Novamente, considerando que os modelos desenvolvidos ao longo deste trabalho e apresentados na Seção 3 são classificados como de regressão, a Seção 2.2.4 foca em desenvolver os modelos com esta perspectiva.

2.2.3.4 Função perda

Fundamental ainda para os modelos de regressão é o entendimento da chamada função perda. Voltando a (2), tem-se que o objetivo da função perda é definir uma métrica para atribuir valores que possam mensurar a assertividade da predição de \hat{y} , dado que o valor correto que se almeja é y . Dessa notação, tem-se que a função perda representada como $L(y, \hat{y})$.

Como exemplos clássicos de função perda, temos o erro absoluto, também conhecido como l_1 , definido como

$$L(y, \hat{y}) = |\hat{y} - y|, \quad (3)$$

e erro quadrático, l_2 , como

$$L(y, \hat{y}) = (\hat{y} - y)^2. \quad (4)$$

Este conceito passa a ser relevante a partir da ideia de que, para mensurar o quão bom está o modelo, utiliza-se os cálculos da função perda para todas as amostras do conjunto de treinamento, conforme citado na Seção 2.2.3.1. De forma bastante simplificada, o objetivo dos diversos algoritmos de treinamento, descritos com mais detalhes na Seção 2.2.4, é minimizar a função custo, cada um a partir de uma abordagem própria.

2.2.3.5 Validação cruzada

A validação cruzada é uma estratégia utilizada para avaliar o desempenho dos modelos e o quão bem eles conseguem generalizar os dados. Um dos exemplos mais comuns utilizados em diversos projetos é a Validação Cruzada *k-Fold*, abordagem que envolve dividir o conjunto de treino aleatoriamente em k subconjuntos de tamanho aproximadamente igual.

A partir disso, o modelo é treinado com um dos subconjuntos de fora, que, por sua vez, será utilizado posteriormente para teste, resultando em valor de erro. Repete-se esse processo k vezes, em cada uma delas deixando um dos subsets de fora (JAMES *et al.*, 2013). Ao fim, tem-se como resultado k valores de erro, que, tomando a média entre eles, resulta em um parâmetro mais representativo acerca da capacidade do modelo receber novos dados.

2.2.3.6 Hiperparâmetros

Ao longo do desenvolvimento teórico dos modelos apresentados na Seção 2.2.4, diversas vezes serão discutidos valores chamados de hiperparâmetros. Estes nada mais são do que constantes a serem definidas no momento de implementação do modelo, com o intuito de otimizar seu resultado, seja evitando *overfitting*, melhorando a acurácia ou a partir de qualquer outra estratégia aplicada à abordagem específica.

2.2.4 Modelos

Assim como já descrito na Seção 2.2.2, cada modelo de aprendizado de máquina possui uma abordagem particular, com formulação bem definida e, para os modelos mais tradicionais, uma literatura extensa e já bastante consolidada. Com isso, é conveniente dedicar uma seção para discutir os fundamentos dos algoritmos utilizados neste trabalho e apresentados na Seção 3. Ainda, conforme os modelos forem sendo discutidos, também são inseridos alguns conceitos e técnicas chave constantemente presentes nos modelos de aprendizado de máquina.

2.2.4.1 Regressão Linear

Dentre as abordagens discutidas nesta seção, a regressão linear é a mais simples, sendo um tópico já tradicional na literatura de predições quantitativas (JAMES *et al.*, 2013). Partindo do mais simples dos modelos, no qual uma única variável preditora x é utilizada para obter uma estimativa de \hat{y} da saída y , usando a mesma notação de (2), tem-se

$$\hat{y} = f(x) = w_0 + w_1 \cdot x, \quad (5)$$

onde w_0 e w_1 são coeficientes cujos valores se almeja descobrir, chamados de interseção e inclinação, respectivamente, uma vez que formam uma reta.

Assim, para encontrar o valor desses coeficientes, a regra geral indica a busca por valores que tornem a reta mais próxima possível dos valores de *target*. Embora existam diversos métodos para tal, um dos mais simples envolve o critério dos mínimos quadrados. Nele, utiliza-se os valores da *Residual Sum of Squares* (RSS), definidos como o somatório de todos os quadrados das diferenças entre \hat{y} e y para cada observação, para, com o auxílio de ferramentas de cálculo, encontrar w_0 e w_1 que minimizam esse valor. De maneira geral, a RSS é definida como

$$RSS = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_m - \hat{y}_m)^2. \quad (6)$$

Para modelos com mais de uma variável preditora, a estratégia que pode ser utilizada é análoga. Aqui, as variáveis e os coeficientes se relacionam da seguinte forma:

$$\hat{y} = f(\mathbf{X}) = w_0 + w_1 \cdot x_1 + \dots + w_n \cdot x_n. \quad (7)$$

Novamente, busca-se encontrar os valores de w_0, w_1, \dots, w_n que levem à minimização de uma função custo como a RSS. Naturalmente, trata-se agora de um problema mais complexo, mas que, a partir de equações de álgebra matricial, pode ser resolvido com certa facilidade.

- *Funções de base*

Para muitas situações, como é, de fato, o escopo deste trabalho, funções lineares, como a (5) e a (7), têm dificuldade em modelar os problemas envolvidos, os quais via de regra são não-lineares. A partir disso, uma estratégia comumente utilizada é aplicar uma transformação nas variáveis predictoras. Assim, ao invés de utilizar um modelo linear em X , tem-se um modelo na forma

$$\hat{y} = w_0 + w_1 \cdot b_1(x_1) + \dots + w_n \cdot b_n(x_n), \quad (8)$$

onde b_1, b_2, \dots, b_n são funções fixas e conhecidas, chamadas de funções de base. Para regressões polinomiais, por exemplo, tem-se que $b_j(x_i) = x_i^j$. Esse tipo de transformação se mostra bastante útil no desenvolvimento deste trabalho, visto que, conforme descrito na Seção sobre o DESSEM, os complexos modelos de otimização que fazem a obtenção do PLD não são lineares.

2.2.4.2 Métodos Baseados em Árvore

Os métodos baseados em árvore são abordagens de aprendizado de máquina, cujas aplicações utilizadas podem se dar tanto para problemas de regressão, como de classificação. De forma geral, estas abordagens envolvem estratificação ou segmentação do preditor para um número mais simples de regiões (JAMES *et al.*, 2013). Assim, há regras definidas para serem responsáveis por essas separações nestes modelos, tornando a tomada de decisão por subconjuntos, motivo pelo qual os modelos são conhecidos como árvores de decisão.

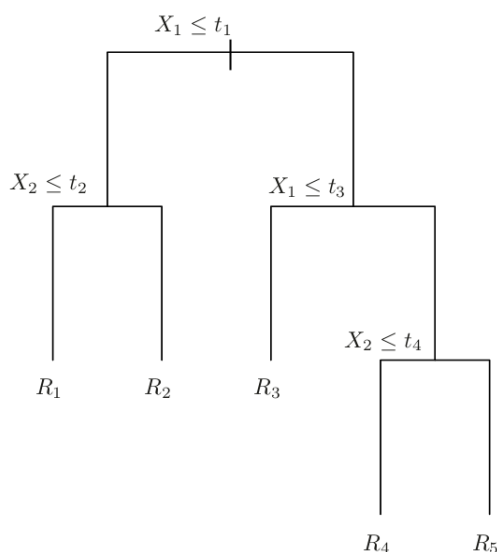
Embora estes modelos sejam simples e de visualização bastante intuitiva, especialmente para casos com menos variáveis, os métodos baseados em árvore não costumam atingir resultados competitivos em termos de acurácia de predição (JAMES *et al.*, 2013). Nesse sentido, estes modelos tornam-se também relevante no contexto em que diversos conceitos apresentados aqui são relevantes para compreender os métodos de *Ensemble*, apresentados na Seção 2.2.4.3, e que, por sua vez, apresentam resultados mais satisfatórios.

Por conveniência, esta seção se dedica a explorar os fundamentos de uma abordagem bastante simples dentre deste conjunto: a árvore de regressão. O seu funcionamento básico pode ser dividido conforme James *et al.* (2013) em duas etapas. Na primeira, divide-se o espaço preditor - ou seja, os possíveis valores das variáveis de entrada (X_1, \dots, X_n) - em J regiões distintas e sem interseção, (R_1, \dots, R_J). Posteriormente, para cada observação que entra dentro da região R_J , faz-se a mesma predição, o que é simplesmente a média dos valores das respostas para as variáveis de treinamento.

Apesar da descrição formal aparentar ser complexa, a intuição desse modelo é bastante simples. Para fazê-lo, observa-se a Figura 8. Nela, tem-se um exemplo simplificado de uma árvore de regressão, em que há duas variáveis de entrada, X_1 e X_2 . Para a etapa de testes, avalia-se se a variável de entrada X_1 é maior ou igual ao limiar estabelecido t_1 . Em caso afirmativo, passa-se para o nó seguinte à direita, repetindo o processo para o limiar t_3 , e, em caso negativo, à esquerda, tornando agora a comparar o valor de X_2 com t_2 . Quando exauridos os nós, define-se, enfim, a qual região, entre R_1, R_2, R_3, R_4 o vetor testado pertence; esta região, onde não há mais divisões, é chamada de folha. Para cada, tem-se um diferente valor preditor da variável de saída.

Clarificada a etapa de testes, a questão fundamental aqui é a definição, durante a etapa de treinamento, dos limiares t_1, t_2, t_3, t_4 . Novamente, utiliza-se o conceito de valores residuais, com o objetivo de encontrar, para cada divisão, os limiares que fornecem o menor RSS. Assim, mais uma vez há a aplicação de cálculo para, a cada nó, definir o

Figura 8 – Ávore de decisão simplificada



Fonte: (JAMES *et al.*, 2013)

limiar que minimiza o erro médio quadrático das predições.

Também para a árvore de regressão, há de se evitar que haja *overfitting*, uma vez que, novamente seguindo o exemplo de Figura 8, é possível ir fazendo novos nós que reduzam o RSS, até que se tenha uma observação para cada folha, fazendo o modelo zerar o erro para treinamento, porém com problemas para generalizar no teste. Nesse sentido, a estratégia mais usual utilizada é limitar essa segmentação, definindo o número mínimo de amostras dentro do espaço para criação de uma nova folha, hiperparâmetro definido na implementação do modelo.

Outra possível abordagem nesse sentido é o de profundidade máxima, que também visa evitar que o modelo vá se desdobrando de forma muito extensa, criando uma nova folha para cada observação. Novamente, esse limite é estabelecido como um hiperparâmetro.

2.2.4.3 Métodos de *Ensemble*

O conjunto de métodos conhecido como *Ensemble* são aqueles cuja abordagem combina diversos modelos, em geral com desempenho menos competitivo, visando à obtenção de um único e mais poderoso modelo (JAMES *et al.*, 2013). De forma geral, estes métodos podem ser organizados de duas formas: *bagging*, na qual os modelos se organizam de forma paralela, e *boosting*, de forma sequencial.

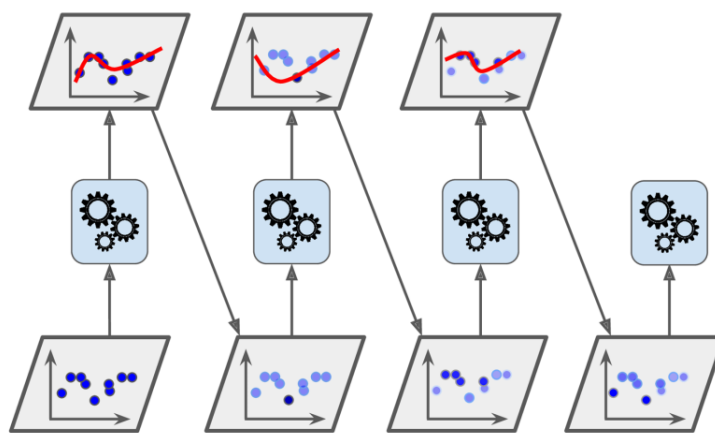
O *bagging* é um procedimento que, de forma geral, objetiva reduzir a variância do modelo, a partir dessa estrutura paralela de agregação de modelos. Isso se dá a partir da lógica que, em um conjunto de n amostras independentes, cada uma com variância

σ^2 , tem-se que a variância da média dessas observações é σ^2/n . Essa observação matemática relevante indica que a obtenção da média de um conjunto reduz a variância. Como consequência, uma forma natural de diminuir a variância e aumentar a acurácia de um modelo é simplesmente aumentar o número de conjuntos de treino, construir modelos independentes e tirar a média dos resultados preditos (JAMES *et al.*, 2013). Essa é a lógica dos modelos de *bagging*.

De forma semelhante, o *boosting* é uma estratégia que pode ser utilizada para diversos métodos de aprendizado de máquina, embora muito comumente o feito a partir de árvores de decisão, e é dessa perspectiva que convém explicá-lo. Ainda a partir de analogias à abordagem anterior, enquanto o *bagging* roda modelos para diversos conjuntos de dados e os combina no final, o *boosting* trabalha de forma semelhante, à exceção de que as árvores são desenvolvidas de forma sequencial, em que cada uma utiliza os dados da anterior (JAMES *et al.*, 2013).

A lógica de ambas as abordagens é ilustrada nas Figuras 9 e 10. A partir desta introdução conceitual, é válido discutir os modelos implementados neste trabalho que fazem uso de tais estratégias.

Figura 9 – Funcionamento simplificado de um algoritmo baseado em *boosting*

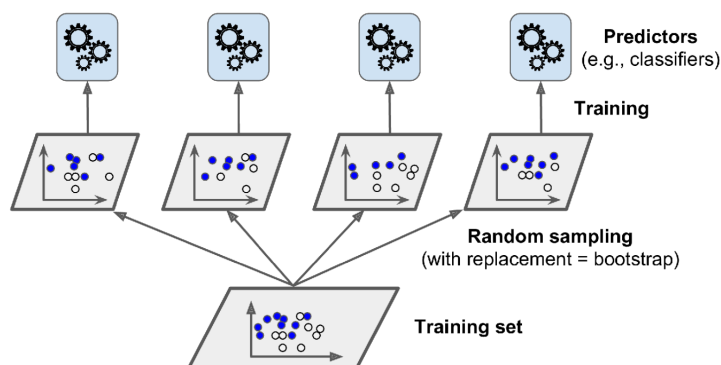


Fonte: (GÉRON, 2019)

- *Florestas Aleatórias*

No contexto da introdução apresentada, o modelo de florestas aleatórias é uma abordagem que utiliza a estratégia de *bagging*, agregando, de forma geral, árvores de decisão.

Para sua construção, são aplicadas diversas árvores de decisão à amostra de treinamento, com o grande diferencial o fato de que, na construção dessas árvores, a cada divisão feita em um nó, uma amostra aleatorizada de m preditoras é escolhida dentro do

Figura 10 – Funcionamento simplificado de um algoritmo baseado em *bagging*

Fonte: (GÉRON, 2019)

conjunto total de *features*. A cada nova divisão, um novo subconjunto com m preditoras é definido (JAMES *et al.*, 2013).

A justificativa para a escolha de apenas um subconjunto dentre todas as variáveis preditoras é bastante intuitivo. Em situações em que há uma variável forte, mesmo agregando diversos modelos paralelamente com o *bagging*, eles se desenvolvem de forma semelhante, uma vez que esta variável forte inevitavelmente estará posicionada nos nós superiores das árvores. A consequência lógica disso é o aumento da correlação entre os modelos, acabando fundamentalmente com o objetivo de se utilizar *bagging*.

Para as Florestas Aleatórias, os hiperparâmetros são o número de estimadores, que pode ser interpretado como o número de árvores na floresta do modelo, a profundidade máxima, que limita a expansão dos nós, e a mínima divisão de amostras, que define um valor mínimo de amostras necessárias para segmentar um nó.

- *AdaBoost e Gradient Boosting*

Os métodos de *boosting* têm uma intuição bastante simples, ao passo que costumam apresentar resultados bastante competitivos.

O primeiro deles, o *AdaBoost*, foi introduzido pela primeira vez a partir do trabalho de Freund e Schapire (1997). O princípio fundamental é aqui é que, a cada iteração, são aplicados pesos para cada amostra de treinamento. Inicialmente, pesos são atribuídos com um mesmo valor padrão, fazendo com que o algoritmo se inicie com um simples modelo fraco. A partir disso, a cada nova iteração, o peso atribuído às amostras é modificado, aumentando para aquelas que foram classificadas de forma incorreta. Esse processo se repete até que o último modelo, em que cada amostra tem um peso diferente, é compilado, obtendo-se o resultado final.

Similarmente, o algoritmo de *Gradient Boosting* é também bastante popular, e funciona de forma sequencial, adicionando parâmetros a cada nova iteração. A diferença

reside no fato de que, ao invés de alterar os pesos das observações, esta abordagem treina novamente o modelo a partir do erro residual da iteração anterior (GÉRON, 2019).

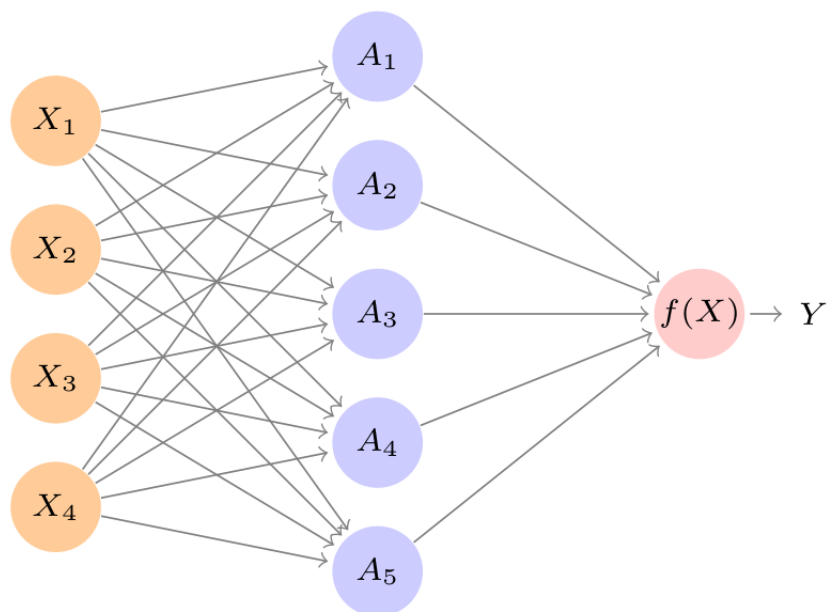
Para este modelo, os hiperparâmetros são o número de estimadores, referente ao número de estágios a serem iterados, a taxa de aprendizado, que determina a contribuição de cada modelo fraco e a profundidade máxima, que limita a expansão dos modelos fracos.

2.2.4.4 Redes Neurais Artificiais

Adicionalmente, cabe também apresentar o conceito de RNA: modelos de aprendizado de máquina cuja dinâmica é inspirada nos neurônios biológicos do nosso cérebro. Estas RNA são os elementos fundamentais de outro conceito relevante utilizado nos principais projetos de aprendizado de máquina atualmente: o aprendizado profundo (GÉRON, 2019).

Novamente, para explicar o funcionamento de uma RNA, é válido começar pelo mais simples dos modelos. Na Figura 11, é apresentada a estrutura de uma simples RNA *feed-forward*. Com uma dinâmica similar aos modelos anteriores, aqui tem-se um vetor de entrada \mathbf{X} , que contém 4 observações, e uma saída y . Dentro da terminologia de aprendizado profundo, tem-se que cada um dos círculos é um neurônio artificial e que o conjunto deles forma uma camada: em amarelo a camada de entrada, em roxo a camada oculta e em rosa, a de saída. As flechas indicam que há uma comunicação entre cada neurônio, ou seja, uma função determina a relação entre eles (JAMES *et al.*, 2013).

Figura 11 – Funcionamento simplificado de uma RNA



Fonte: (JAMES *et al.*, 2013)

Nesta camada oculta, entende-se que o número de neurônios seja um hiperparâmetro a ser definido. Sendo K o número de neurônios nesta camada, a ideia geral aqui é que as K ativações A_k sejam calculadas a partir das variáveis de entrada \mathbf{X} pela relação

$$A_k = h_k(\mathbf{X}) = g(w_{k0} + \sum_{j=1}^p w_{kj} \mathbf{X}_j), \quad (9)$$

onde $g(z)$ é conhecida como função de ativação, sendo definida previamente por uma função não-linear. É possível pensar que cada A_k é uma transformação de $h_k(X)$ das variáveis originais de entrada (JAMES *et al.*, 2013).

Posteriormente, as K ativações da camada oculta alimentam a camada de saída, resultado na função

$$f(\mathbf{X}) = \beta_0 + \sum_{k=1}^K \beta_k A_k. \quad (10)$$

Em (9) e 10, os parâmetros β_k e w_{kj} são estimados a partir dos dados de treinamento.

Como nos demais modelos, a etapa de treinamento aqui visa à definição de tais parâmetros, de forma a reduzir ao máximo a função de perda do conjunto, a partir de conjuntos de entrada e saída de dados.

Generalizando este modelo, é típico que as RNA tenham mais de uma camada oculta, além de um número maior de neurônios por cada camada. A teoria fundamental que rege essa dinâmica é que, matematicamente, uma camada que contenha um número suficientemente grande de neurônios pode aproximar grande parte das funções existentes (JAMES *et al.*, 2013).

Nas redes neurais, os hiperparâmetros definidos incluem o número de épocas, que é a quantidade de vezes que o conjunto de treinado passa pelo modelo, as camadas convolucionais e o tamanho da máscara, que definem o comportamento das operações de convolução da rede neural, e o *dropout*, que consiste na exclusão de algumas amostras do conjunto de treinamento para evitar *overfitting*.

2.2.5 Métricas de avaliação

Ao longo das seções 3 e 4, o desempenho do modelo é avaliado a partir de duas diferentes métricas. A primeira delas é o *Root Mean Squared Error* (RMSE), tipicamente utilizada para mensurar modelos de regressão. Sua utilização é em geral recomendada para situações em que se deseja penalizar com um peso maior os erros mais significativos (GÉRON, 2019). A partir da notação apresentada na Seção 2.2.2, o RMSE pode ser equacionado da seguinte forma:

$$RMSE = \sqrt{\frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2}. \quad (11)$$

Em outros contextos, contudo, é mais conveniente a utilização do *Mean Absolut Error* (MAE), definido como:

$$MAE = \frac{1}{m} \sum_{i=1}^m |y_i - \hat{y}_i|. \quad (12)$$

Isto se dá em situações contendo diversos *outliers* na predição, que fariam com que a métrica de erro anterior fosse muito grande se utilizada. Neste contexto, inclusive, é conveniente utilizar em certos contextos ambas as métricas, visando observar se o modelo está com erros uniformes, com pequenas variâncias ou com muitos pontos fora da curva.

Para o presente modelo, contudo, a métrica escolhida é o *Mean Absolute Percentage Error* (MAPE), valor cuja unidade é apresentada em porcentagem, permitindo uma interpretação mais tangível do resultado do modelo. Seu equacionamento é definido como:

$$MAPE = \frac{100\%}{m} \sum_{i=1}^m \left| \frac{y_i - \hat{y}_i}{y_i} \right|. \quad (13)$$

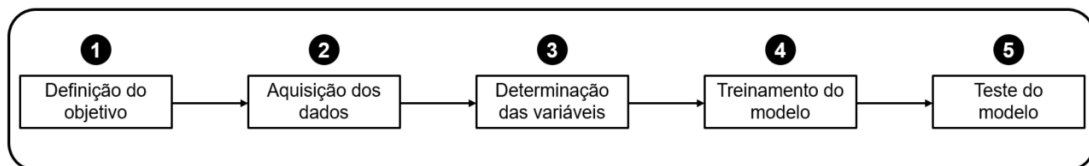
3 DESENVOLVIMENTO

A partir da fundamentação teórica apresentada, é possível iniciar a apresentação do desenvolvimento dos modelos em si. Para isso, também cabe elencar previamente as premissas definidas para tal. São elas:

- a) Os modelos devem prever a série temporal dos **7 dias seguintes de PLD**, ou seja, os 168 valores seguintes ao momento da predição;
- b) As entradas, também séries temporais, são definidas para os **7 dias anteriores do momento da predição**;
- c) As estratégias de otimização dos modelos devem buscar obter o **menor valor possível de MAPE**.

A metodologia seguida no presente trabalho é ilustrada na Figura 12, seguindo a sugestão proposta por Watt, Borhani e Katsaggelos (2016). Seguindo com este fluxo, tem-se o primeiro bloco como Definição do Objetivo, no qual se deve definir de forma clara quais são os objetivos do problema a ser resolvido. Dado que esta questão já foi abordada na Seção 1.1, não há necessidade de aprofundar novamente tais pontos.

Figura 12 – Fluxo de trabalho definido para desenvolvimento deste trabalho



Fonte: Próprio autor a partir de Watt, Borhani e Katsaggelos (2016)

Com isso, o único objetivo adicional determinado nesta etapa inicial foi relativo às questões técnicas de implementação do modelo. Neste sentido, as ferramentas utilizadas no desenvolvimento deste trabalho são as seguintes:

- a) *Jupyter Notebook* como plataforma para aplicação dos códigos;
- b) *Python* como a linguagem de programação;
- c) *Numpy* e *Pandas* como as bibliotecas principais para o tratamento e estruturação dos dados;
- d) *Scikit-Learn* e *Darts* como as bibliotecas de séries temporais e aprendizado de máquina para implementação dos modelos;
- e) *Matplotlib* e *Seaborn* como as bibliotecas de visualização de dados para a análise exploratórias das variáveis e de desempenho dos modelos.

Eventualmente, outras ferramentas foram utilizadas como suporte ao longo do desenvolvimento do trabalho. Contudo, estas acima citadas as foram fundamentais para

a obtenção dos resultados. Além disso, a última etapa do projeto, na qual o conjunto de dados de teste é aplicado ao modelo para verificar os resultados obtidos e avaliado o desempenho resultante, também é, por conveniência, apresentada compilando todos os resultados no Capítulo 4. Dessa maneira, esta seção foca em apresentar os três passos restantes: a aquisição de dados, a definição das *features* e a etapa de teste dos modelos, sendo esta última também segmentada para discutir a implementação de cada um dos algoritmos implementados.

3.1 AQUISIÇÃO DOS DADOS

Ainda seguindo a metodologia definida a partir de Watt, Borhani e Katsaggelos (2016), a indicação é que a aquisição de dados objetive coletar a maior quantidade possível de dados para a resolução do problema proposto, buscando fazê-la também de forma diversa. Esta seção trata da coleta de dados e também são abordados os pontos relativos ao pré-processamento desses dados, posterior à coleta.

Assim, o primeiro passo consiste na obtenção dos dados relativos ao PLD. Conforme discutido na seção 2.1.4.3, o PLD passou, a partir de janeiro de 2021, a ser estabelecido em base horária. Contudo, o modelo DESSEM já vinha sendo testado como operação sombra, ou seja, sem implementação prática, desde 2018, porém inaugurado pelo ONS apenas em janeiro de 2020 (CEPEL, 2018a). Assim, para aumentar a disponibilidade de dados, e, por consequência, melhorar o desempenho dos modelos, opta-se por trabalhar com os dados de PLD que vão do 1º de janeiro de 2020 até o dia 31 de junho de 2021. A data final do conjunto coincide com o início do desenvolvimento desse trabalho.

Para fazer a obtenção desses dados históricos, a CCEE disponibiliza o Painel de Preços, no qual diversas aquisições relacionadas ao PLD são possíveis, permitindo fazer inclusive essa filtragem por data. Algumas rotinas de pré-processamento foram necessárias para adequar a base de dados baixada para um formato que posteriormente fosse possível adicionar as demais variáveis.

Ainda como efeito de simplificação, também foram filtrados nesta etapa os dados apenas para o PLD do Submercado Sul, definido convenientemente por ser a região onde está localizada a UFSC. A partir dessa definição do espaço temporal no qual as análises seriam feitas, os demais dados de hidrologias foram obtidos através do portal Portal de Dados Abertos do ONS, disponibilizado pelo operador do sistema. O portal, lançado em 2021 para disponibilizar diversos dados históricos do setor elétrico no Brasil, faz parte da estratégia de digitalização e transparência do ONS. Nesse portal, foram adquiridos dois parâmetros segmentados para cada bacia hidro energética do SIN: o ENA e o EAR. O primeiro deles representa a energia produzível de todas as usinas que fazem parte da bacia, sendo calculada a partir do produto das vazões naturais aos reservatórios com as produtividades a 65% dos volumes úteis (ONS, 2021b). Já o EAR é uma medida que

indica a energia associada ao volume de água disponível nos reservatórios e que pode ser convertida em geração, seja na própria usina ou nas seguintes que estejam à jusante (ONS, 2021a).

Para o pré-processamento dessas duas variáveis, foi necessário também replicar cada observação por 24, uma vez que a discretização desses parâmetros é diária, enquanto o que virá a ser o *target* do modelo, o PLD, tem suas observações horárias. Ademais, rotinas em *Python* foram novamente implementadas para padronizar os dados para um formato possível de se trabalhar com as bibliotecas estabelecidas para este trabalho.

Ainda, também no portal do ONS, foi adquirida a série histórica de Precipitação Diária Observada (PO), relativa às chuvas mensuradas nas estações meteorológicas e postos pluviométricos. O operador fornece estes dados para cada dia em todas as estações de observação. O conjunto de dados obtido fornece, além do valor total diário de precipitação em mm, a localização geográfica da estação em questão, a partir dos dados de latitude e longitude de cada uma. Novamente, estando esses valores obtidos em discretização diária, faz-se necessário uma etapa anterior de pré-processamento dos dados.

A Tabela 1 sintetiza os dados adquiridos nesta etapa inicial de desenvolvimento.

Tabela 1 – Resumo dos dados adquiridos

	Variável	Unidade	Fonte
PLD	Preço de Liquidação das Diferenças	R\$/MWh	CCEE
ENA	Energia Natural Afluentes por Bacia	MWmês	ONS
EAR	Energia Armazenada por Bacia	MWmês	ONS
PO	Precipitação Observada	mm	ONS

Fonte: Próprio autor

3.2 DETERMINAÇÃO DAS VARIÁVEIS

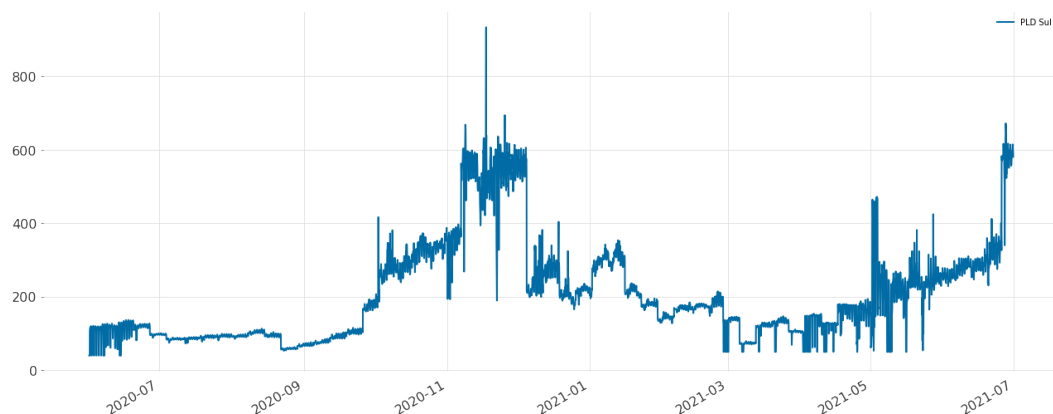
Segundo Watt, Borhani e Katsaggelos (2016), esta etapa consiste em analisar quais são as *features* que entram para o modelo. Neste sentido, a etapa seguinte à aquisição dos dados citados na Seção 3.1 consiste em desenvolver a Análise Exploratória de Dados (AED). Segundo Hartwig e Dearing (1979), as técnicas de AED são utilizadas de forma vital para visualizar e revelar informações acerca dos dados utilizados. A partir disso, foram utilizadas as bibliotecas de dados citadas no início deste capítulo para a obtenção de visualizações que pudessem ajudar na compreensão das bases de dados à disposição.

Seguindo a mesma ordem da aquisição de dados descrita na Seção 3.1, inicialmente foram feitas visualizações para compreender os dados do PLD no Submercado Sul. A observação dos dados de PLD permite algumas observações relevantes. A primeira delas

é a respeito da volatilidade dos dados, que tende a dificultar a precisão dos modelos de predição. Em contrapartida, para alguns períodos, há uma certa estabilização do valor ao longo de ao menos algumas semanas.

Como consequência dessas observações, os dados de PLD assumem dois papéis distintos neste trabalho: variável de saída, a partir do objetivo definido de prever os seus valores para os 7 dias seguintes, e como variável de entrada, utilizando-se dos 7 dias anteriores de dados observados. Em razão dessa volatilidade discutida, uma visualização limitada em 7 meses, de dezembro de 2020 a julho de 2021, é mais conveniente para ilustrar as conclusões, conforme apresentado na Figura 13.

Figura 13 – Valor do PLD no Submercado Sul em R\$/MWh



Fonte: Próprio autor a partir dos dados de ONS (s.d.)

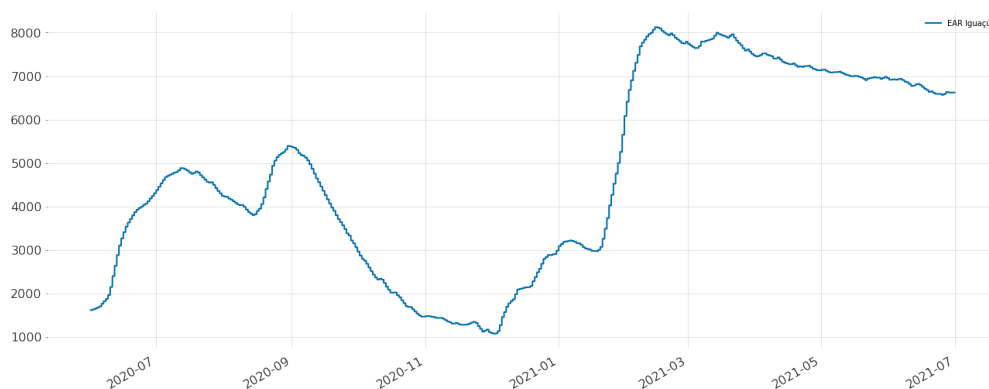
Posteriormente, trabalha-se na visualização dos dados relativos às bacias. Aqui, conforme observado nos exemplos obtidos para duas bacias, do Uruguai, para o EAR, e do Iguazú, para o ENA, apresentados nas Figuras 14 e 15. Um ponto importante é que há uma evidente sazonalização para algumas bacias do sistema brasileiro.

Além disso, uma observação minuciosa dos dados das diversas bacias cujos valores foram adquiridos indica que essa sazonalidade não é comum para todas elas e que há tendências divergentes de padrões para cada uma, muito em razão da extensão territorial do SIN. Assim, julga-se necessária a inclusão dos dados de todas as bacias como variáveis de entrada. Para fazê-lo de forma simplificada, calcula-se, previamente, o somatório dos valores de todas as bacias para obter um parâmetro diários único de ENA e EAR.

As bacias cujos dados de ENA e EAR foram incluídos nos modelos são as do Amazonas, Araguari, Capivari, Doce, Grande, Iguazu, Jacuí, Jequitinhonha, Paraguaçu, Paraguai, Paraíba do Sul, Paraná, Parnaíba, Paranapanema, São Francisco, Tietê, Tocantins e Uruguai.

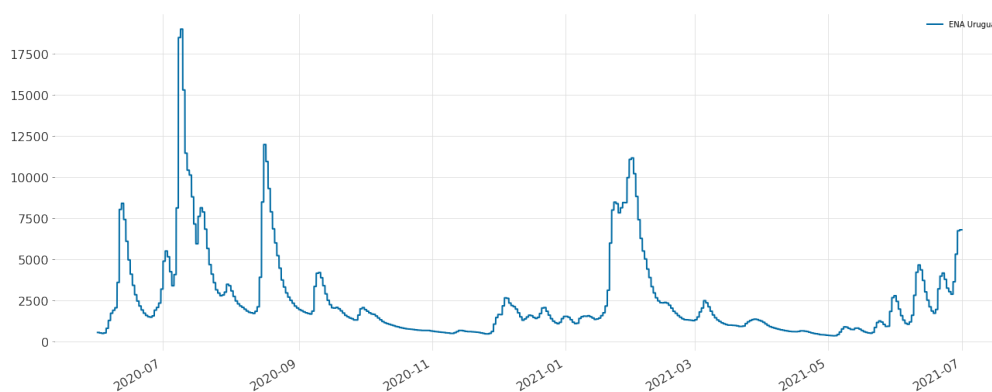
Por fim, os últimos dados para os quais foram feitas visualizações foram os de precipitação. Com os dados segmentados para cada um dos 1.440 pontos de observação, a

Figura 14 – Valor da EAR na Bacia do Iguazú em MWmês



Fonte: Próprio autor a partir dos dados de ONS (2021a)

Figura 15 – Valor da ENA na Bacia do Uruguai em MWmês



Fonte: Próprio autor a partir dos dados de ONS (2021b)

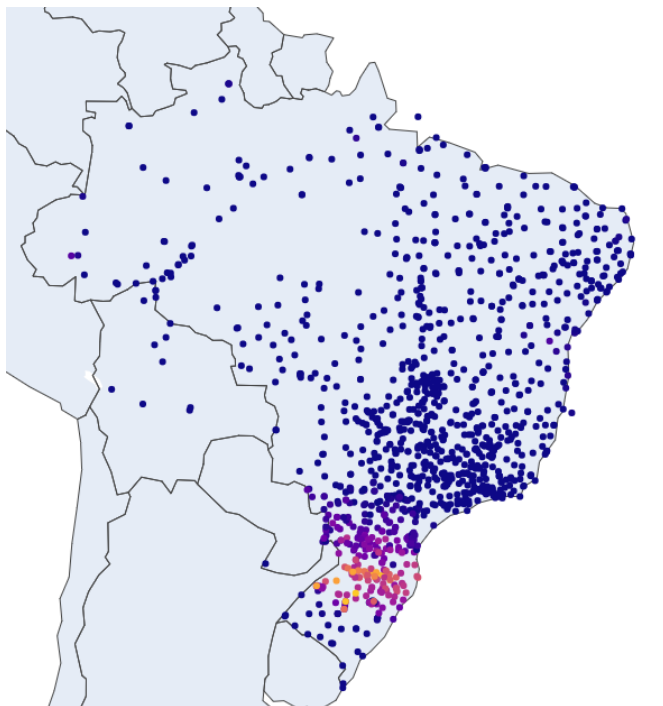
estratégia inicial consistia em selecionar alguns pontos para inclusão no modelo.

A partir da execução da AED desses dados, porém, nota-se que a precipitação observada em um dia pode diferir de forma excessiva para pontos muito próximos, conforme observado na figura 16, referente à PO no dia 15 de agosto de 2020, data definida de forma aleatória para ilustração.

Assim, para evitar a situação na qual se perderiam dados relevantes para o modelo, opta-se por não selecionar alguns pontos, mas incluir como variável de entrada a soma de todos eles para cada dia. O resultado, embora volátil, pode representar bem a precipitação total nas bacias do SIN. Também com uma evidente sazonalização, a série é apresentada na Figura 17.

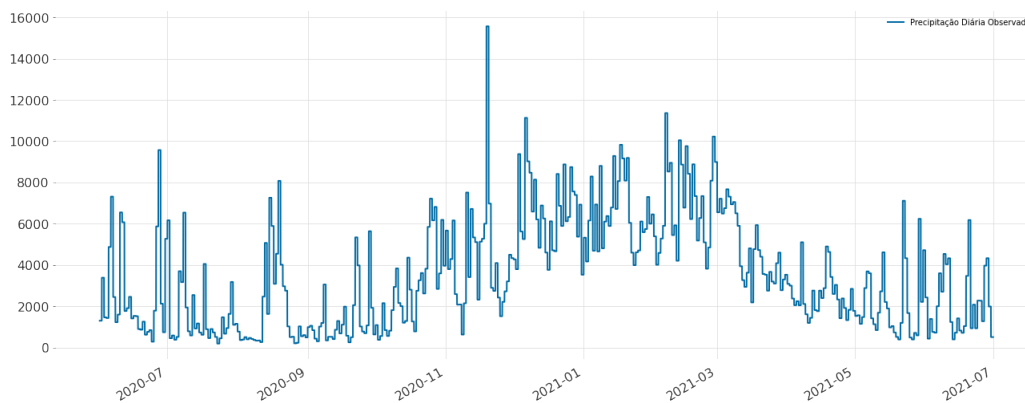
A partir dos pontos expostos, utiliza-se todas as variáveis adquiridas para o treinamento do modelo, melhor descrito na seção a seguir. Adicionalmente, como há uma diferença significativa do preço do PLD dentro de um mesmo dia - de forma geral, mais

Figura 16 – Precipitação diária observada em cada ponto no dia 15 de agosto de 2020



Fonte: Próprio autor

Figura 17 – Precipitação diária observada total em todos os pontos em mm

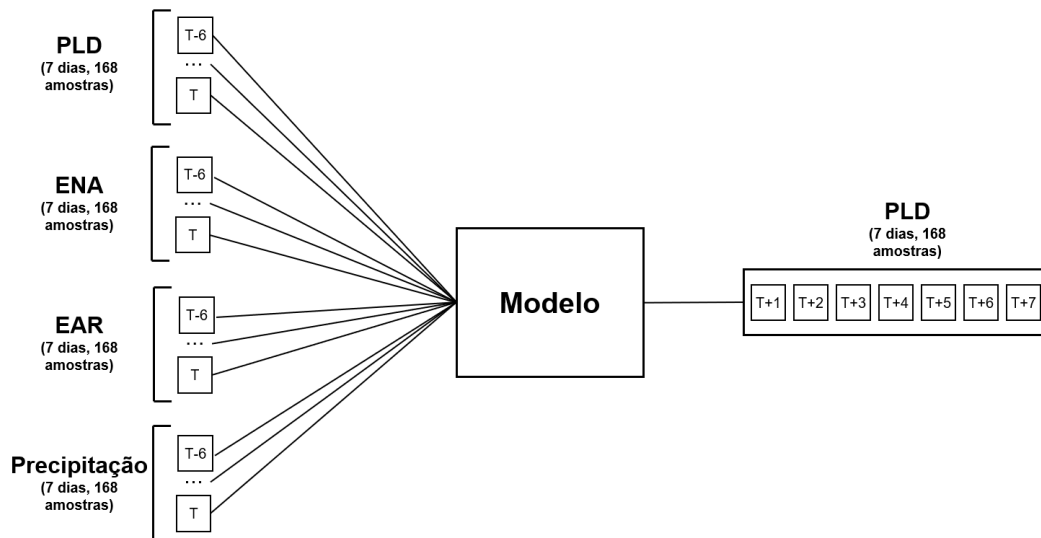


Fonte: Próprio autor

caro em horários de consumo mais alto - tem-se também a hora da observação com variável de entrada, variando de 0 a 23.

Como síntese, a Figura 18 apresenta a estrutura simplificada do modelo, com as variáveis de entrada à esquerda, todas alimentando o modelo, que prediz o PLD.

Figura 18 – Estrutura simplificada dos modelos



Fonte: Próprio autor

3.3 TREINAMENTO

Na etapa de treinamento, tem-se, enfim, a implementação do modelo, com o intuito de encontrar os parâmetros que fazem a predição dos dados de teste. Para fazê-lo, porém, é necessário previamente a definição dos hiperparâmetros adequados para cada algoritmo.

Para encontrar a combinação de hiperâmetros que apresentam melhor desempenho, uma opção é iterar manualmente diversos arranjos e comparar seus resultados. Essa possibilidade, contudo, toma obviamente um esforço mecânico bastante custoso e obviamente maçante, conforme discute Géron (2019). Ainda de acordo com o autor, uma alternativa bastante viável é utilizar a ferramenta *GridSearch* para fazê-lo automaticamente.

Nesta função, define-se uma série de valores para cada um dos hiperparâmetros do modelo, e, automaticamente, o algoritmo utiliza a técnica de validação cruzada para explorar todas as possíveis combinações (GÉRON, 2019). Como possível saída desta função, por exemplo, há a possibilidade de se obter o valor de MAPE para cada combinação, permitindo, assim, seguir com a implementação do modelo com os melhores hiperparâmetros. Todos os modelos de validação cruzada foram feitos com $k = 5$, conforme Seção 2.2.3.5.

Em relação à escolha dos modelos, opta-se por trabalhar com 4 deles: regressão linear simples, floresta aleatória, *Gradient Boosting* e redes neurais convolucionais. A escolha se dá pela possibilidade de comparar um modelo mais simples que pode servir de referência para os demais, dois modelos de *Ensemble*, sendo um de *Bagging* e outro de *Boosting*, e finalmente um modelo de redes neurais.

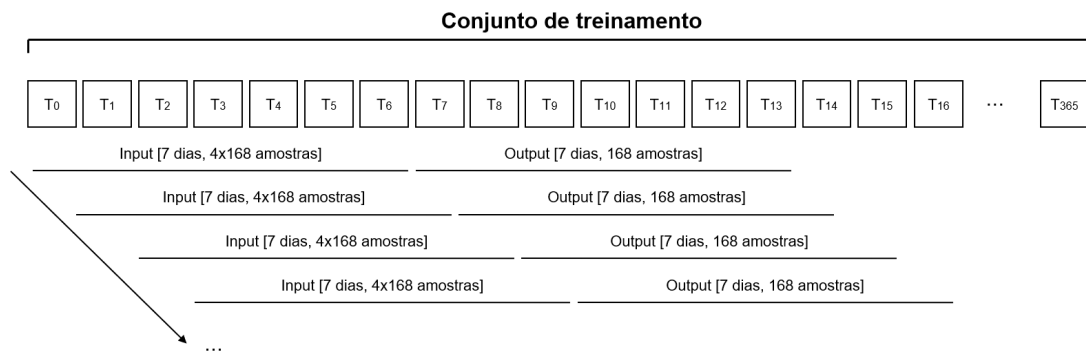
Sendo essa a abordagem escolhida para a etapa de treinamento dos modelos deste trabalho, cabe descrever (i) os modelos implementados e (ii) os hiperparâmetros iterados e

os selecionados para cada um deles. Os resultados são apresentados na lista abaixo, em que os hiperparâmetros do conjunto que melhor teve desempenho são destacados em negrito.

- **Regressão Simples** - sem hiperparâmetros;
- **Florestas Aleatórias**
 - Número de Estimadores = [40, **80**, 100].
 - Profundidade Máxima = [10, **25**, Indefinida].
 - Mínima divisão de amostras = [2, **3**, 5].
- **Gradient Boosting**
 - Número de Estimadores = [60, **70**, 100, 120].
 - Taxa de Aprendizado = [0.01, **0.1**, 0.3].
 - Profundidade Máxima = [2, **3**, 5].
 - Mínima Divisão de Amostras = [1, **2**, 3].
- **Redes Neurais**
 - Número de Épocas = [**200**].
 - Camadas Convolucionais = [**1**].
 - Tamanho da Máscara = [1, **3**, 5].
 - *Dropout* = [0, **0.1**, 0.2].

Ainda relativo à etapa de treinamento, a Figura 19 apresenta como são tabulados os dados de treinamento do modelo. Essa dinâmica é feita de forma automática pela biblioteca *Darts*, e visa à criação dos pares de entrada, que incluem todas as observações das 4 variáveis de entrada citadas na Seção anterior para 7 dias, e saída, no caso o valor do PLD para os 7 dias posteriores. Para ser possível realizar a validação do modelo em etapa posterior, o conjunto de treinamento é segmentado como os primeiros 12 meses do conjunto de dados, ou seja, contempla todo o período de 2020.

Figura 19 – Dinâmica de tabularização dos dados para treinamento



Fonte: Próprio autor

4 RESULTADOS

Para avaliar os resultados, divide-se esse capítulo em duas seções: primeiramente são apresentadas métricas obtidas para cada modelo quando aplicado o conjunto de teste, e, posteriormente, é brevemente discutida uma possível aplicação de um modelo de predição do PLD para um agente do SIN, mensurando também os seus benefícios.

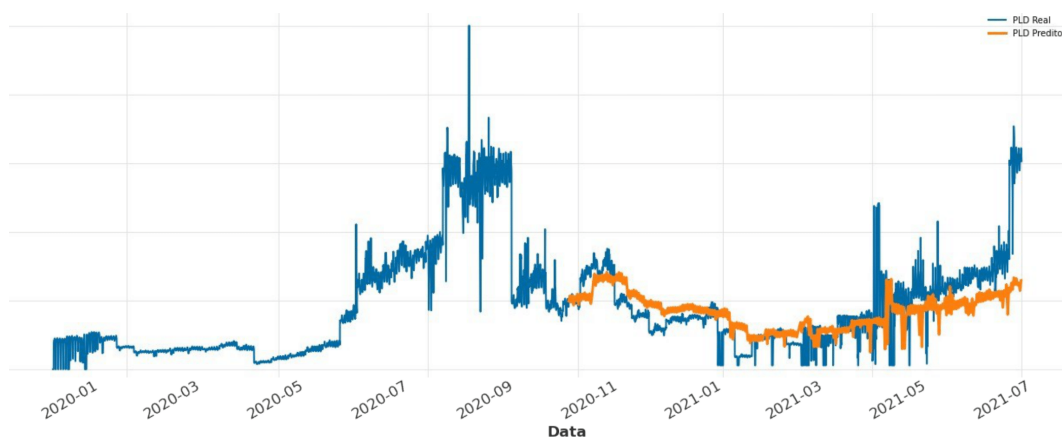
4.1 AVALIAÇÃO DOS MODELOS

Definidos os parâmetros, hiperparâmetros e os modelos, e treinados os algoritmos, foi enfim possível aplicar o conjunto de dados de teste no modelo para avaliar a eficácia de cada um. Para fazê-lo, foi novamente aplicada uma função da biblioteca *Darts* para aplicar uma espécie de *Backtesting*, isto é, calcula-se os valores que teriam sido preditos no passado, se houvesse sido aplicado o modelo à época. Essa dinâmica permite uma avaliação rápida e precisa do comportamento da predição frente aos dados reais.

Para melhorar a visualização e poder facilmente comparar os modelos, todos os gráficos apresentados nesta seção são apresentados com os 6 primeiros meses contendo apenas os dados reais, enquanto nos 6 meses restantes o valor real do PLD em azul e a predição em cor alaranjada.

A partir disso, para o modelo de regressão linear, mesmo que bastante simples, já foi possível obter resultados interessantes, os quais indicam uma convergência do modelo para as tendências de PLD, conforme figura 20. O MAPE obtido foi de 38,77%.

Figura 20 – Comparação entre o PLD real e o predito pelo modelo de regressão linear

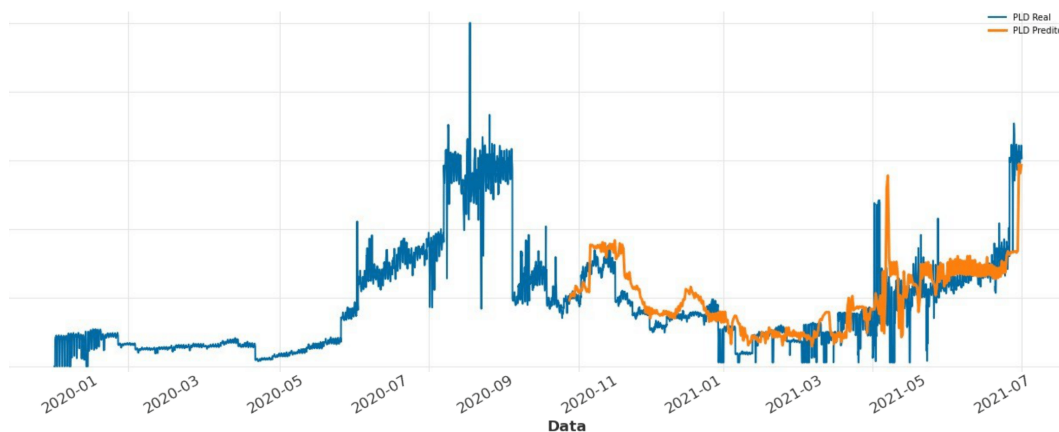


Fonte: Próprio autor

Posteriormente, foram obtidos os valores para os modelos de *Ensemble*, tanto o de árvore aleatória e o *Gradient Boosting*, apresentados nas Figuras 21 e 22, respectivamente.

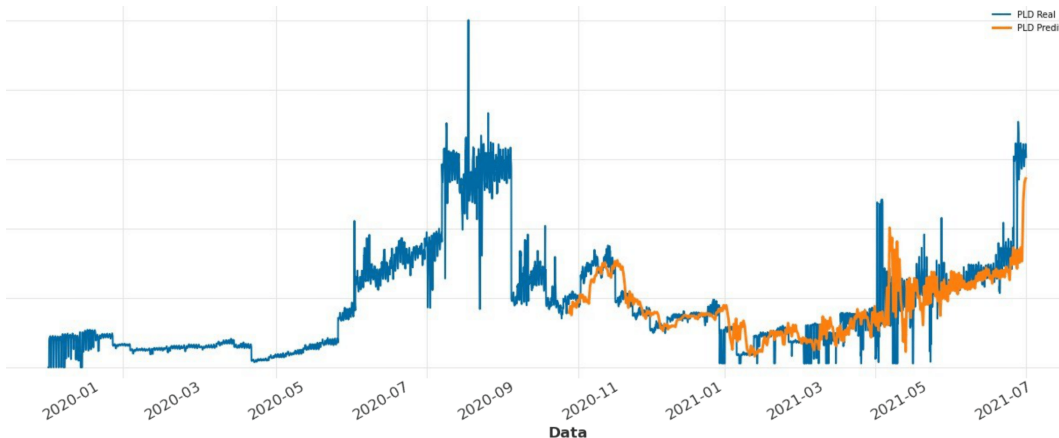
O erro médio absoluto do primeiro foi de 27,59%, e 28,39% no segundo. Conforme esperado, os modelos apresentaram um resultado superior ao anterior.

Figura 21 – Comparação entre o PLD real e o predito pelo modelo de floresta aleatória



Fonte: Próprio autor

Figura 22 – Comparação entre o PLD real e o predito pelo modelo de *Gradient Boosting*



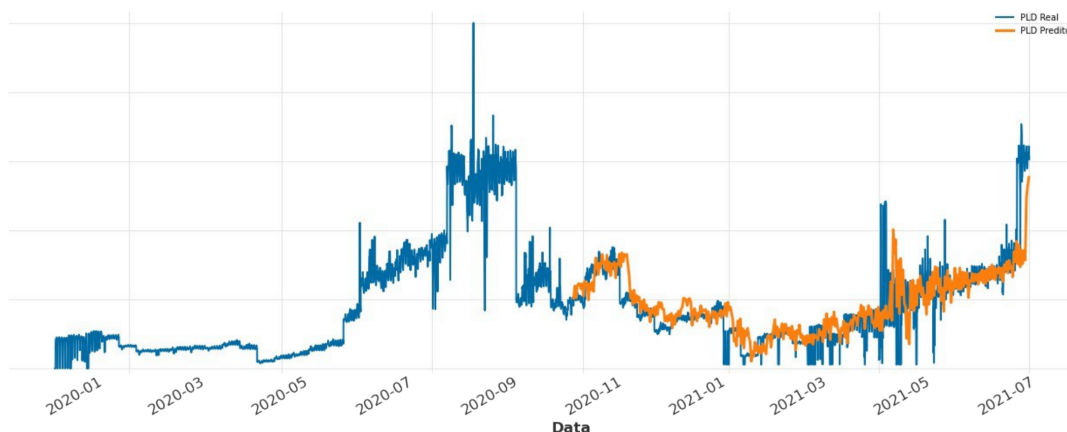
Fonte: Próprio autor

Finalmente, o modelo de RNA foi testado com a mesma dinâmica, apresentando resultados ligeiramente melhores que os modelos implementados até então, com um MAPE de 24,70%. A visualização para o resultado do modelo é apresentada na figura 23.

Para evidenciar a diferença entre o desempenho de cada modelo, a Tabela 2 sintetiza o MAPE de todos os modelos testados.

A partir dessa tabela, é interessante notar que os resultados estão alinhados com a teoria apresentada. Apesar disso, no início do trabalho havia a expectativa de uma diferença mais significativa entre eles. É plausível acreditar que, a partir de mais testes e trabalhando

Figura 23 – Comparação entre o PLD real e o predito pelo modelo de Redes Neurais



Fonte: Próprio autor

Tabela 2 – Resumo do desempenho de cada modelo

Modelo	MAPE
Regressão Simples	38,77%
Floresta Aleatória	27,59%
<i>Gradient Boosting</i>	28,39%
Redes Neurais	24,70%

Fonte: Próprio autor

com uma janela de tempo maior, essa disparidade se tornaria mais evidente. Outra possível alteração sugerida para trazer melhoria aos resultados é a criação de modelos que tenham variáveis preditoras limitadas à região do submercado Sul, ao invés da utilização de dados agregados do SIN. Como os modelos de otimização são desenvolvidos tendem a melhorar os modelos. Ademais, também se nota um certo atraso dos valores preditos para os valores reais. Uma possível interpretação disso é que as variáveis preditoras conseguem de forma satisfatória aprender as tendências do PLD, mas não são suficientes para antecipar subidas ou descidas bruscas no preço da energia.

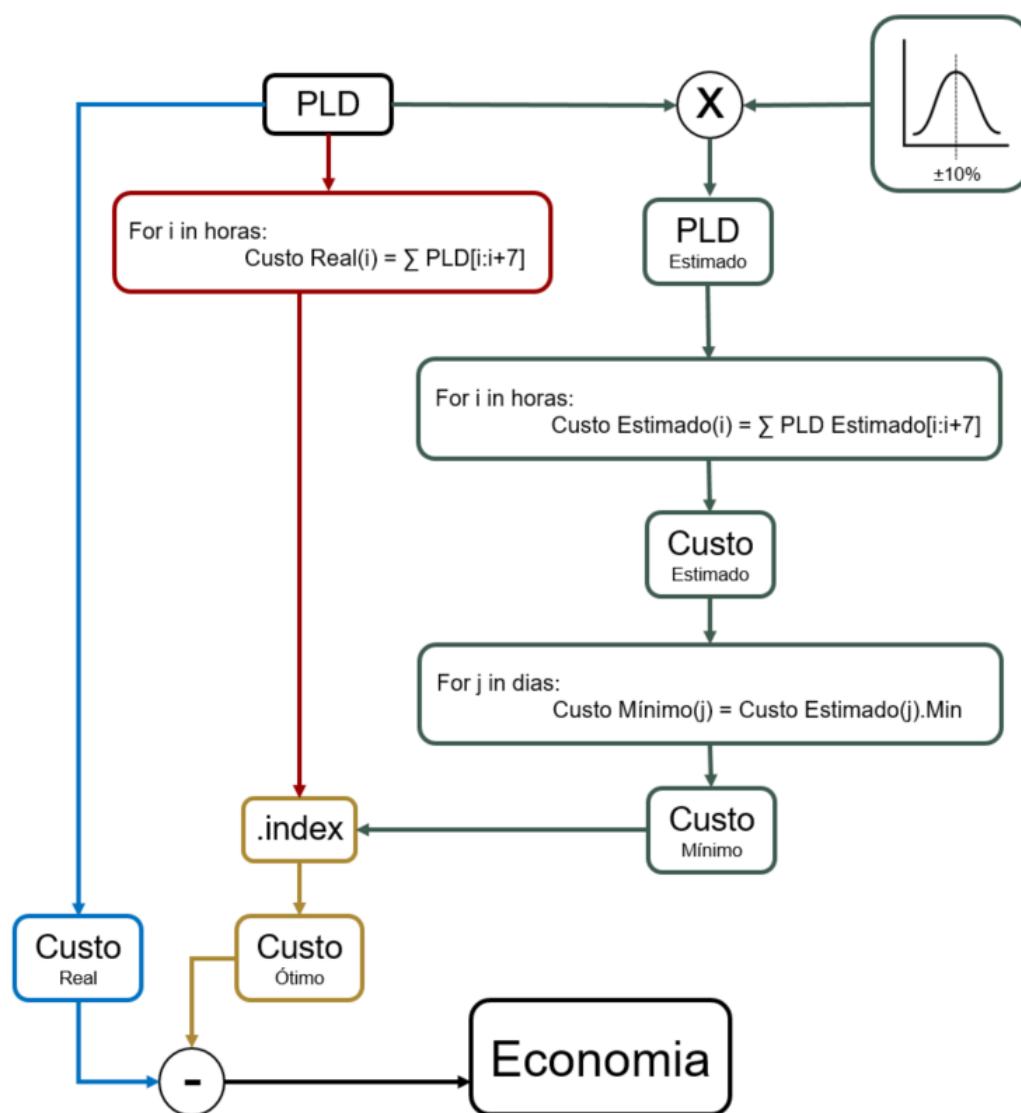
4.2 ESTUDO DE CASO

Conforme discutido exhaustivamente ao longo deste trabalho, a mudança na discretização do PLD é uma enorme inovação para o mercado brasileiro de energia elétrica, permitindo diversos novos produtos que tendem a aumentar a eficiência do setor. Nesta direção, esta Seção apresenta uma das inúmeras possíveis aplicações que podem ser exploradas, tanto a partir da mudança para o PLD horário, como também a partir da posse de

um modelo que faça previsões do seu valor, mesmo não o sendo com alta acurácia.

O desenvolvimento deste estudo de caso parte da seguinte pergunta: quanto poderia uma empresa, com flexibilidade para deslocar o horário de início das 8 horas diárias de produção, poderia ter de economia, caso em posse de um modelo de previsão do PLD. Para isso, foi simulado, usando apenas os dados históricos PLD, uma nova rotina de códigos, com as mesmas ferramentas descritas na Seção 3. Um diagrama de blocos com alguns pseudo-códigos apresentando a lógica utilizada é apresentado na figura 24.

Figura 24 – Diagrama de blocos do cálculo de economia do estudo de caso



Fonte: Próprio autor

O princípio aqui é considerar um consumidor livre de energia, cujo contrato de energia estabelecido com a comercializadora é do tipo *flat*, e que possui a flexibilidade para alterar o horário do seu consumo de carga.

Para isso, a rotina de códigos consiste, de forma geral, em multiplicar o vetor de PLDs por valores distribuídos em uma gaussiana com média de 10%, pseudo-aleatoriamente gerados para valores negativos e positivos. O intuito da aquisição desse novo vetor, apresentado na figura como PLD Estimado no laço em verde, é simular a aplicação de um modelo de predição, com erro médio absoluto percentual de 10%.

Outra premissa importante de ser colocada é que o contrato deste agente é apenas *Flat*, de acordo com o discutido na Seção 2.1.3.1. Fosse um contrato modulado, a aplicação não teria tanto resultado.

Em posse desse vetor, o algoritmo adquire quais os horários ótimos para início da produção, ou seja, qual o momento do dia em que a soma dos 8 PLDs estimados seguintes resultam no menor somatório. Esse valor é apresentado como Custo Mínimo. Como este parâmetro é resultado ainda da aleatorização, encontra-se qual o valor real de PLD com o início nos horários obtidos como ótimos. O racional deste ponto é saber qual o custo na prática que teria a empresa, caso iniciada a produção no horário recomendado pelo modelo.

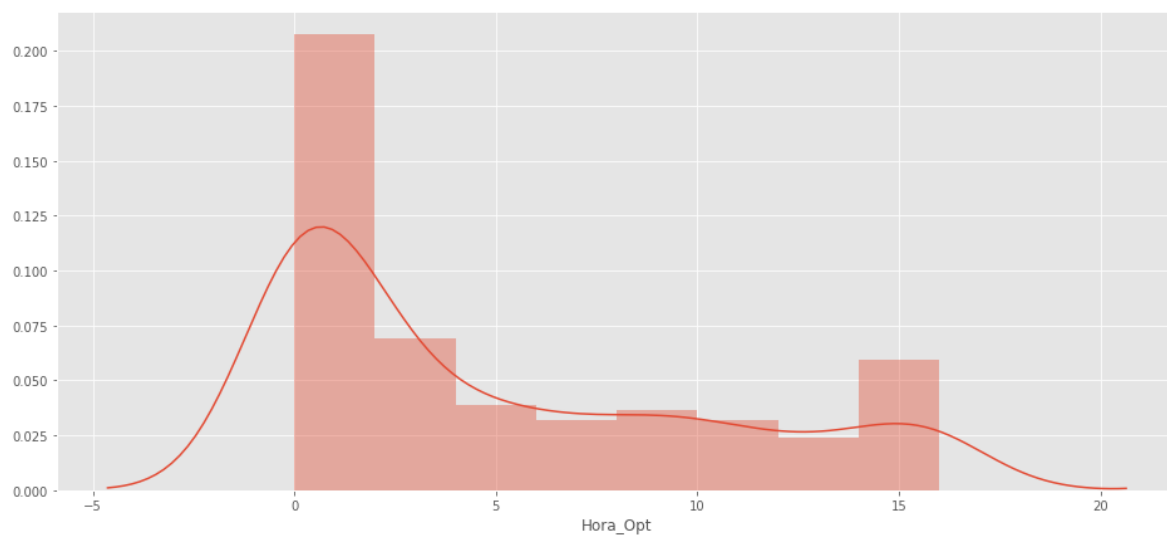
Por fim, um terceiro laço, em azul na figura 24, obtém o custo da empresa, caso rodando a produção em um horário mais usual; aqui, definido como entre 10 e 18 horas. Esse custo é apresentado como real, e comparado com o custo ótimo.

O resultado desse algoritmo é bastante pertinente. Aplicando essa rotina para o período de PLD entre janeiro de 2020 e julho de 2021, a economia média diária é de 7,3%.

A Figura 25 evidencia um ponto importante a partir dessa análise: de forma geral, o horário para início ótimo da produção é durante a madrugada, período no qual o consumo de energia é evidentemente menor. O motivo pelo qual todos os horários iniciais ótimos de início de produção estão entre as 00 e as 16 horas é que uma das premissas do estudo de caso é que a produção deveria iniciar e terminar no mesmo dia.

Importante observar que essa economia não impacta diretamente no custo total de energia consumida, somente aquela adquirida a PLD no MCP. Contudo, com essa redução no preço médio pago, há a possibilidade de o agente aumentar sua exposição neste próprio mercado de curto prazo, sem necessariamente aumentar seus riscos. Para indústrias do setor de consumo cíclico, por exemplo, em que a demanda é fortemente volátil e influenciada por fatores externos, poder reduzir a contratação de energia e aumentar sua exposição no mercado de energia, pode ser uma interessante vantagem competitiva.

Figura 25 – Distribuição dos horários considerados ótimos para início da produção no estudo de caso



Fonte: Próprio autor

5 CONCLUSÃO

Diversas mudanças de dinâmica e regulação pautaram o mercado de energia elétrica no Brasil nas últimas décadas. A transição ocorrida em 2020 quanto à discretização horária do PLD, contudo, pode ser citada com uma das mais relevantes e impactantes.

Ao longo deste trabalho, em especial na Seção 4, foi discutida com mais profundidade uma possível vantagem competitiva que pode ser desenvolvida a partir da implementação do PLD horário. Contudo, há uma infinidade de novas ferramentas e produtos que podem ser criadas conforme o mercado for ganhando mais maturidade.

Também no contexto de exposição, por exemplo, a quantidade de produtos financeiros que podem ser utilizados como proteção para as mais diversas operações serão ainda pautas de inúmeros estudos. Para torná-los mais eficientes, algoritmos de predição de tendência do preço spot serão fundamentais e, agentes que conseguirem fazê-lo com a melhor acurácia possível, conseguirão largar na frente.

Ainda, mesmo que se tenha discutido até aqui as vantagens competitivas a partir um viés financeiro, é importante também estender para a questão ambiental atrelada à eficiência energética. Principalmente em um contexto de transição para fontes renováveis, permitir que os agentes consumam energia elétrica em períodos em que a demanda tende a ser menor, pode aumentar de forma substancial a eficiência de todo setor, diminuindo, por exemplo, a necessidade de geração a partir de fontes térmicas.

De forma geral, os resultados obtidos neste trabalho são bastante satisfatórios, principalmente na circunstância de se discutir mudanças ainda bastante recentes, o que torna o tópico passível de inúmeros aprofundamentos.

Sem dúvida alguma os modelos de predição do PLD horário serão amplamente discutidos ao longo dos próximos anos, sendo a expectativa deste trabalho tornar-se uma abordagem inicial bastante simples para referência de trabalhos mais robustos. Nesse sentido, os futuros projetos podem tanto focar no desenvolvimento de novos algoritmos, como também na implementação de variáveis adicionais de entrada.

Em relação ao primeiro, e novamente citando algoritmos de aprendizado profundo, muito se tem trabalhado no contexto de séries temporais nos modelos de redes neurais recorrentes, por exemplo. Esse modelo teria como vantagem adicional a possibilidade de se trabalhar com um comportamento dinâmico, a partir da utilização da camada oculta anterior para o cálculo da saída. Nesse sentido, é como se houvesse uma memória, em que a cada período a rede pudesse utilizar os dados armazenados do passado.

Naturalmente, esse tipo de aplicação demandaria idealmente uma quantidade muito maior de dados, mas a expectativa é que, com o passar do tempo, a disponibilidade de informações coletadas sobre o PLD seja maior, permitindo, de fato, a construção de um modelo desse tipo.

Já quanto aos dados adicionais de entrada, é importante citar que a proposta deste

trabalho era de fato desenvolver uma abordagem simples, avaliando exclusivamente as variáveis hídricas. Contudo, os modelos de otimização citados na Seção 2.2.4 envolvem entradas adicionais. Dessa forma, envolver também dados de geração térmica, eólica e solar, demanda atual e prevista de carga, previsão do tempo, entre outros, tendem a melhorar a acurácia.

Ainda acerca da discussão das sugestões de trabalhos futuros, saindo desse viés de modelos de predição, a implementação do PLD horário abre espaço para outras discussões relevantes. Como exemplos, são relevantes as discussões sobre a razoabilidade desse novo parâmetro, ou seja, se de fato ele otimiza o planejamento da operação, ou até mesmo se ele reflete de forma satisfatória como o sistema está operando e quais outras variáveis poderiam auxiliar na modelagem e fazê-lo.

REFERÊNCIAS

- ABRACEEL. **Mercado Livre de Energia Elétrica - Um guia básico para consumidores potencialmente livres e especiais.** [S.l.], 2019. Acessado em: 12 jan. 22. Disponível em: https://www.abraceel.com.br/archives/files/Abraceel_Cartilha_MercadoLivre_V9.pdf.
- ANEEL. **Atlas de Energia Elétrica no Brasil.** [S.l.], 2008. Acessado em: 15 ago. 21. Disponível em: <http://www2.aneel.gov.br/arquivos/PDF/atlas3ed.pdf>.
- ANEEL. **Institucional: A Aneel.** [S.l.]. Acessado em: 10 jan. 22. Disponível em: <https://www.aneel.gov.br/a-aneel>.
- ANEEL. **Regulação do Mercado de Energia Elétrica - Comercialização.** [S.l.], 2018. Acessado em: 10 jan. 22. Disponível em: <https://www.aneel.gov.br/mercado-de-eletricidade>.
- BISHOP, Christopher. **Pattern Recognition and Machine Learning.** [S.l.]: Springer, 2006. Disponível em: <https://www.microsoft.com/en-us/research/publication/pattern-recognition-machine-learning>.
- BRASIL. DECRETO N° 1.009. **Diário Oficial [da] República Federativa do Brasil,** 1993a. DOI: http://www.planalto.gov.br/ccivil_03/decreto/1990-1994/d1009.htm.
- BRASIL. DECRETO N° 915. **Diário Oficial [da] República Federativa do Brasil,** 1993b. DOI: http://www.planalto.gov.br/ccivil_03/decreto/1990-1994/D0915.htm.
- BRASIL. Lei N° 10.848. **Diário Oficial [da] República Federativa do Brasil,** 2004. DOI: http://www.planalto.gov.br/ccivil_03/_ato2004-2006/2004/lei/110.848.htm.
- BRITO, Erico. Revisão tarifária e diferenças regionais: Um estudo de concessões de distribuição de energia elétrica no Brasil. **Dissertação de mestrado,** 2010. DOI: 10.11606/D.86.2010.tde-30092010-153101.
- CCEE. **Conceitos de Preço.** [S.l.], 2021a. Acessado em: 14 jan. 22. Disponível em: <https://www.ccee.org.br/web/guest/precos/conceitos-precos>.

CCEE. **O que fazemos: Com quem se relaciona.** [S.l.], 2021b. Acessado em: 10 jan. 22. Disponível em: <https://www.ccee.org.br>.

CCEE. **O que fazemos: Informações ao mercado.** [S.l.]. Acessado em: 10 jan. 22. Disponível em: <https://www.ccee.org.br>.

CCEE. **Regras de Comercialização: Contratos.** 2017.1.0. ed. [S.l.], 2017. Acessado em: 22 jan. 22. Disponível em: https://www.ccee.org.br/ccee/documentos/ccee_377240.

CEPEL. **DESSEM - Modelo de Despacho Hidrotérmico de Curto Prazo.** [S.l.], 2018a. Acessado em: 12 jan. 22. Disponível em: http://srvlumis02.cepel.br/pt_br/produtos/dessem-modelo-de-despacho-hidrotermico-de-curto-prazo.htm.

CEPEL. **Mercado Livre de Energia Elétrica - Um guia básico para consumidores potencialmente livres e especiais.** [S.l.], 2018b. Acessado em: 12 jan. 22. Disponível em: http://srvlumis02.cepel.br/pt_br/produtos/newave-modelo-de-planejamento-da-operacao-de-sistemas-hidrotermicos-interligados-de-longo-e-medio-prazo.htm.

ENGIE. **PLD Horário: saiba o que o primeiro ano ensinou ao mercado.** [S.l.], 2021. Disponível em: <https://www.alemdaenergia.engie.com.br/pld-horario-o-que-os-primeiros-seis-meses-mostram-ao-mercado/>.

EPE. **A EPE - Quem Somos.** [S.l.]. Acessado em: 10 jan. 22. Disponível em: <https://www.epe.gov.br/pt/a-epe/quem-somos>.

FARIAS, Regina. **Atuação Estatal e a Privatização do Setor Elétrico Brasileiro. Dissertação de mestrado,** 2006. DOI: <https://portal.tcu.gov.br/biblioteca-digital/atuacao-estatal-e-a-privatizacao-do-setor-eletrico-brasileiro.htm>.

FEDERAL, Governo. **Ministério de Minas e Energia - MME.** [S.l.]. Acessado em: 10 jan. 22. Disponível em: <https://dados.gov.br/organization/about/ministerio-de-minas-e-energia-mme>.

FREUND, Yoav; SCHAPIRE, Robert. A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting. **Journal of Computer and System Sciences**, v. 55, 1997. DOI: <https://doi.org/10.1006/jcss.1997.1504>.

GÉRON, Aurélien. **Hands-on Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems**. 2. ed. [S.l.]: O'Reilly Media, 2019.

HARTWIG, Frederick; DEARING, Brian E. **Exploratory data analysis**. [S.l.]: Sage, 1979.

JAMES, Gareth; WITTEN, Daniela; HASTIE, Trevor; TIBSHIRANI, Robert. **An Introduction to Statistical Learning with Applications in R**. 2. ed. [S.l.]: Springer, 2013. Disponível em: <https://www.statlearning.com/>.

LANDI, Mônica. Evolução do setor elétrico brasileiro no contexto econômico nacional: uma análise histórica e econométrica de longo prazo. **Dissertação de mestrado**, 2011. DOI: <https://www.teses.usp.br/teses/disponiveis/86/86131/tde-12032012-091848/en.php>.

MEDEIROS, Lucio. Previsão do Preço Spot do Mercado de Energia Elétrica. **Dissertação de doutorado**, 2004. DOI: <https://doi.org/10.17771/PUCRio.acad.4777>.

ONS. **EAR Diário por Bacia**. [S.l.], 2021a. Acessado em: 21 jan. 22. Disponível em: <https://dados.ons.org.br/dataset/ear-diario-por-bacia>.

ONS. **ENA Diário por Bacia**. [S.l.], 2021b. Acessado em: 21 jan. 22. Disponível em: <https://dados.ons.org.br/dataset/ena-diario-por-bacia>.

ONS. **Sobre ONS - O que é ONS**. [S.l.]. Acessado em: 10 jan. 22. Disponível em: <http://www.ons.org.br/paginas/sobre-o-ons/o-que-e-ons>.

SANTOS, Guilherme. Uma aplicação de Redes Neurais Recorrentes do tipo LSTM à previsão dos preços de curto prazo do mercado de energia elétrica brasileiro. **Dissertação de mestrado**, 2019. DOI: <https://hdl.handle.net/10438/28069>.

SCHOUCHANA, Felipe. Decisão ótima em Swaps de energia elétrica no Brasil utilizando a medida Omega. **Dissertação de mestrado**, 2010. DOI: <https://doi.org/10.17771/PUCRio.acad.16911>.

SILVA, Bruno. Energia elétrica e políticas públicas: a experiência do setor elétrico brasileiro no período de 1934 a 2005. **Dissertação de doutorado**, 2006. DOI: 10.11606/T.86.2006.tde-10112011-102906.

SIMÕES, Mario; GOMES, Leonardo. Decisão de sazonalização de contratos de fornecimento de energia elétrica no Brasil através da otimização da medida Ômega, 2017. DOI: <https://doi.org/10.1590/S1413-23112011000100007>.

WATT, Jeremy; BORHANI, Reaz; KATSAGGELOS, Aggelos. **Machine Learning Refined: Foundations, Algorithms and Applications**. 1. ed. New York: Cambridge University Press, 2016.