

## Capítulo

# 1

## **Processamento e análise de *Big Data* para aplicação de algoritmos de *Machine Learning* com interface *Myriad* através da utilização da plataforma *HPCC Systems***

Mauro Donato Marques (LexisNexis Risk Solutions)

### *Abstract*

*Throughout the mini-course, participants will have the opportunity to learn the essential concepts of processing and analyzing massive volumes of data (Big Data) and the process of developing a query service using the open-source platform High Performance Computing Cluster (HPCC Systems) and also, the application of Machine Learning algorithms with Myriad interface, as well as having the possibility to apply the knowledge acquired in a training environment made available in the classroom.*

### *Resumo*

*Ao longo do minicurso, os participantes terão a oportunidade de conhecer os conceitos essenciais de processamento e análise de volumes massivos de dados (Big Data) e o processo de desenvolvimento de um serviço de consulta através da utilização da plataforma open-source composta por um Cluster Computacional de Alto Desempenho (HPCC Systems) e, também, a aplicação de algoritmos de Aprendizado de Máquina com interface Myriad, bem como terão a possibilidade de aplicar os conhecimentos adquiridos em um ambiente de treinamento disponibilizado em sala de aula.*

**Informações Técnicas:** Curso de nível básico. O curso requer apenas um computador com acesso à Internet e uma conta no [GitHub](#).

GitHub repositório: [https://github.com/mauromarx/SSCAD\\_2024](https://github.com/mauromarx/SSCAD_2024)

## 1.1. A Plataforma HPCC Systems

HPCC Systems (High Performance Computing Cluster) é uma plataforma para solução de desafios de Big Data com as seguintes características:

- Supercomputação: processamento paralelo e dados distribuídos;
- Open source: código aberto e gratuita;
- Completa: gestão compreensiva e simplificada do fluxo de dados.

HPCC Systems oferece o melhor de ambos os mundos no que tange ao processamento e análise de volumes massivos de dados (Big Data), ou seja, combina o rápido desempenho de um “Data Warehouse” para a entrega de informações com a capacidade de tratar os dados como se estivessem em um “Data Lake”. HPCC Systems usa arquitetura de dados distribuída e uma metodologia de processamento paralelo para trabalhar com grandes conjuntos de dados.

### 1.1.1. Um breve histórico

A plataforma de Big Data que se tornaria HPCC Systems foi desenvolvida em 2001 por uma equipe interna de engenharia da LexisNexis Risk Solutions. Os primórdios da base tecnológica para a plataforma HPCC Systems foi desenvolvida pela primeira vez na Seisint em 1999, a fim de gerenciar um número significativo de conjuntos de dados. Em 2000, a equipe da Seisint precisava coletar e analisar grandes quantidades de informações brutas de consumo de dados financeiros de uma ampla variedade de fontes para um cliente responsável por determinação da “score” de crédito ao consumidor nos Estados Unidos. Isto levou à criação da linguagem de programação, conhecida como ECL (Enterprise Control Language), que HPCC Systems usa atualmente. Em 2004, a LexisNexis Risk Solutions (LNRS) adquiriu Seisint junto com sua tecnologia, incluindo a linguagem ECL usada para programação em clusters Thor e Roxie. Durante esse período, a plataforma HPCC Systems era empregada, principalmente, como uma ferramenta interna da LexisNexis Risk Solutions, e ainda não havia atingido o seu pleno potencial. Em 2008, a LexisNexis Risk Solutions adquiriu a ChoicePoint, uma seguradora e provedora de análise e, nos próximos três anos, o portfólio de produtos da ChoicePoint foi integrado à plataforma HPCC Systems. Combinando a plataforma HPCC Systems com os produtos de seguros ChoicePoint criou-se negócios poderosos e vantajosos - uma solução de processamento de Big Data extremamente eficiente, capaz de gerenciar grandes quantidades de dados, gerando produtos de tratamento de dados muito mais eficientes no já volumoso mercado de seguros.

Em 2011, a LexisNexis decidiu lançar a plataforma HPCC Systems sob uma licença de código aberto. Desde o seu lançamento, a HPCC Systems desenvolveu uma rica comunidade de usuários e uma rede global de clientes. Além disso, a HPCC Systems continua a inovar a plataforma, adicionando suporte para as arquiteturas baseadas em nuvem; em desenvolvimento de novas ferramentas de administração, governança e orquestração de dados; e adicionando novos recursos de dashboard.

Os usuários atuais da plataforma HPCC Systems que podem ser mencionados publicamente incluem a Quod, um Bureau de Crédito no Brasil, o qual a LNRS ajudou a criar, e muitas universidades proeminentes, como: Clemson University, Florida Atlantic University, Kennesaw State University, RV College of Engineering, Universidade de São

Paulo, Universidade Federal de Santa Catarina, Universidade Federal do Pará e várias outras universidades em todo o mundo.



**Figura 1.1: A evolução da plataforma HPCC Systems.**

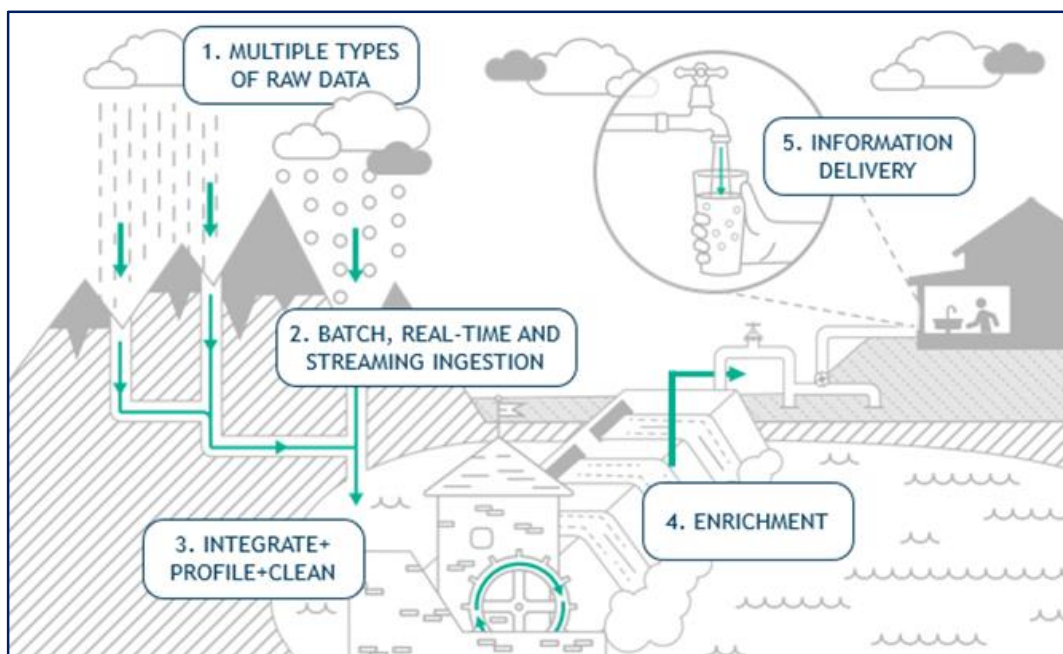
### 1.1.2. O que é a plataforma HPCC Systems?

HPCC Systems é uma plataforma de “Data Lake” de código aberto (open source) projetada para obter dados de diversas fontes de dados em formatos estruturados e não estruturados. Os dados geralmente são armazenados em arquivos simples, como arquivos básicos de disco, ou em armazenamentos de objetos como Amazon S3 e armazenamento BLOB do Azure. Em outras palavras, não há um esquema de formatação pré-definida, na qual os dados precisam ser ajustados antes do armazenamento.

A implementação de uma típica plataforma HPCC Systems começa com apenas algumas fontes de dados, alguns processos de análises iniciais e ferramentas de relatório, mas o tamanho, a complexidade e a capacidade do “Data Lake” pode crescer rapidamente. Depois que os dados são ingeridos no “Data Lake” e, refinados através de limpeza e padronização dos dados, começa o processo de enriquecimento desses dados. Esse enriquecimento de dados é um processo iterativo e evolutivo que extrai tanto conhecimento quanto possível das fontes de dados. Uma vez que esse conhecimento é extraído, ele fica disponível para outros usuários do “Data Lake”, que precisam dos mesmos por meio de um processo conhecido como entrega de dados. Durante a entrega de dados, a plataforma HPCC Systems garante que os dados sejam transferidos aos usuários do “Data Lake” de maneira responsiva e segura. Uma analogia pode ajudar a ilustrar o que acontece com os dados à medida que eles entram em um sistema de “Data Lake” do HPCC Systems. Nesta analogia, a água (dados brutos) é coletada em um reservatório (Data Lake) onde em seguida, é processado para torná-lo adequado ao consumo do público.

Para continuar entregando água à população, a usina de beneficiamento não pode parar a coleta ou processamento de água; o processo deve fornecer água de forma confiável 24 horas por dia, sete dias por semana para acompanhar a demanda do consumidor. Um “Data Lake” deve oferecer o mesmo nível de disponibilidade. Não importa quantos dados sejam adicionados ao sistema, o processo de catalogação e análise desses dados deve operar continuamente em níveis de serviço de “cinco noves” (o “Data Lake” deve estar ativo e operacional pelo menos 99,999 por cento do tempo).

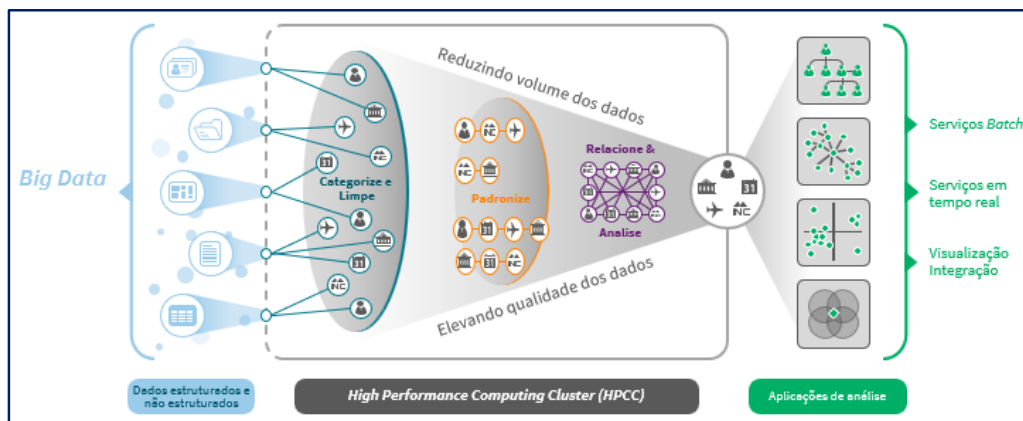
A figura abaixo registra o ciclo de vida dos dados em um sistema real de “Data Lake” do HPCC Systems atualmente em uso por um cliente da plataforma. Movendo da esquerda para a direita, as fontes de dados entregam os dados ao “Data Lake” do HPCC Systems para ingestão, refinamento, enriquecimento, indexação e análise. O HPCC Systems pode gerar relatórios ou dashboards sobre os dados em qualquer etapa do processo, dependendo de qual informação o consumidor necessita. Todos esses processos ocorrem dentro do ambiente de “Data Lake” para produzir resultados consistentes.



**Figura 1.2:** Um “Data Lake” deve ser capaz de ingerir, formatar e enriquecer dados com disponibilidade 24 horas por dia, 7 dias por semana.

### 1.1.3. O Pipeline de enriquecimento de dados no HPCC Systems

O pipeline de dados no HPCC Systems segue os dados desde a origem até sua ingestão no cluster do HPCC Systems, onde é formatado, enriquecido e depois disponibilizado para aplicações hospedadas no Cluster.



**Figura 1.3:** Fluxo de dados no HPCC Systems – “Funil” de dados.

### 1.1.4. Os componentes da plataforma HPCC Systems

O HPCC Systems é composto pelos seguintes componentes:

- A linguagem de programação ECL (Enterprise Control Language) é uma linguagem de programação declarativa e orientada a dados desenvolvida para uso em “Data Lakes” na plataforma HPCC Systems;
- Thor é um cluster de processamento de dados em massa, o qual limpa, padroniza e indexa dados de entrada para uso pelo “Data Lake”. Depois que os dados forem refinados pelo Thor, podem, então, serem usados pelo cluster Roxie;
- Roxie é um cluster de API/consulta (query) em tempo real para consultar dados após refinamento por Thor. As consultas Roxie são executadas em menos de um segundo e fornecem resultados de forma concorrente.

Os clusters Thor e Roxie apresentam objetivos específicos e, assim, fazendo uma analogia simplista com o mundo oceânico seria algo como mostrado na figura abaixo:



Figura 1.4: Objetivos dos clusters Thor e Roxie.

### 1.1.5. A plataforma HPCC Systems

Um diagrama da plataforma HPCC Systems apresentando um cluster Thor (para processamento de dados em massa) e um cluster Roxie (para lidar com consultas de dados) e, ainda, o chamado Power Trio:

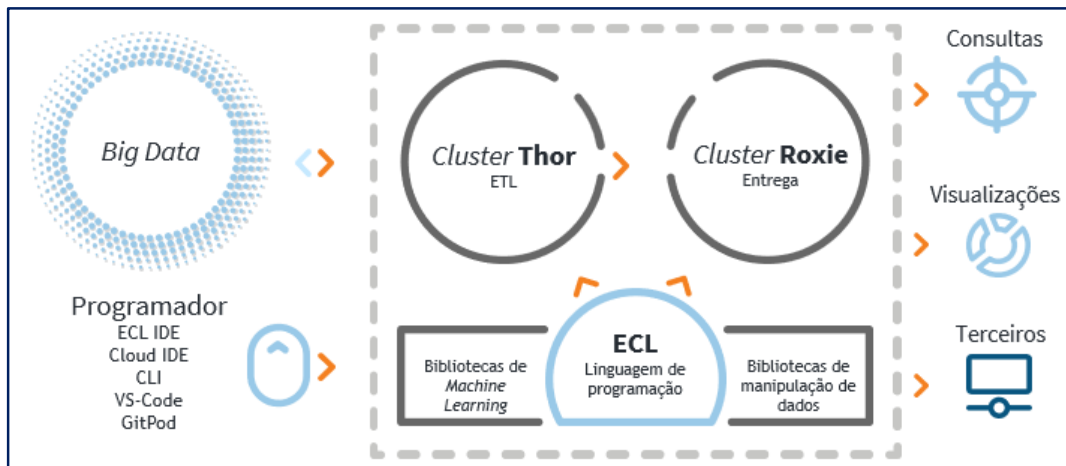


Figura 1.5: Clusters Thor e Roxie.

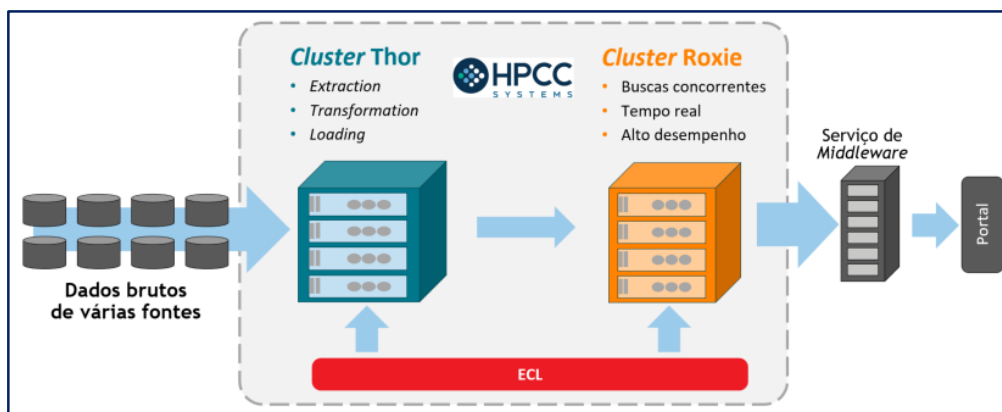


Figura 1.6: Power Trio: A) Thor, B) Roxie e C) ECL.

### A) Thor: Arquitetura

Thor é um cluster projetado especificamente para executar processos de manipulação massiva de dados (ETL). Thor é um cluster de preparação de dados de back-office e não se destina a consultas de nível de produção do usuário final.

Os clusters Thor são usados para fazer todo o trabalho "pesado" de preparação de dados para processar dados brutos em formatos padrão. Uma vez concluído o processo, os usuários finais podem consultar esses dados padronizados para coletar informações reais.

No entanto, os usuários finais geralmente desejam ver os seus resultados "imediatamente" e, ainda, geralmente mais de um usuário final deseja obter seus resultados ao mesmo tempo. O cluster Thor só funciona em uma consulta por vez, o que o torna inviável para o usuário final, por isso foi criado o cluster Roxie.

### B) Roxie: Arquitetura

Roxie é um cluster projetado especificamente para atender as consultas padrão, fornecendo uma taxa de transferência de mais de mil respostas por segundo (a taxa de resposta real para qualquer consulta, logicamente, depende de sua complexidade). Roxie é um cluster de nível de produção projetado para aplicações de missão crítica.

Os clusters Roxie podem lidar com milhares de usuários finais simultâneos e fornecer a todos eles a percepção dos resultados "imediatamente". Ele faz isso permitindo apenas que os usuários finais executem consultas padrão pré-compiladas que foram desenvolvidas especificamente para uso do usuário final no cluster Roxie.

Normalmente, essas consultas usam índices e, portanto, fornecem desempenho extremamente rápido. No entanto, o cluster Roxie é inviável para uso como ferramenta de desenvolvimento, pois todas as suas consultas devem ser pré-compiladas e os dados utilizados devem ter sido implantados anteriormente.

### C) A linguagem ECL

ECL é uma linguagem de programação declarativa, que apresenta inúmeras vantagens sobre o modelo de programação mais convencional. A programação declarativa permite ao programador expressar a lógica de uma computação sem descrever seu controle de

fluxo. Em termos leigos, a linguagem ECL permite que os desenvolvedores digam ao sistema o que eles precisam, mas deixa para o sistema determinar a melhor maneira para fazer isso.



Além de simplificar o design e a implementação de algoritmos complexos, também melhora a qualidade do código, minimizando ou eliminando efeitos colaterais no código do “Data Lake”, o que facilita o teste de código e simplifica a manutenção do código. O código ECL é mais fácil de entender e verificar, mesmo por pessoas que não estão familiarizadas com o design original, o que ajuda encurtar a curva de aprendizado para novos programadores.

ECL é a única linguagem necessária para expressar algoritmos de dados em toda a plataforma HPC Systems. No Thor, a linguagem ECL expressa “workflows” de dados que consistem em carregamento dados, transformação, vinculação, indexação etc. No Roxie, a linguagem ECL define consultas de dados. Isso significa que os analistas de dados e programadores da plataforma HPC Systems só precisam aprender uma linguagem para definir o ciclo de vida completo dos dados.

A linguagem ECL é implicitamente paralela, portanto, o mesmo código ECL desenvolvido para ser executado em um cluster de um nó pode ser executado com a mesma facilidade em um cluster com milhares de nós. O programador não precisa se preocupar em implementar a paralelização, e a linguagem ECL possui uma função otimizadora que garante o melhor desempenho para uma arquitetura específica.

A linguagem ECL foi projetada desde o início para ser uma linguagem de programação orientada a dados. Ao contrário de outras linguagens, funções primitivas de alto nível para dados, como JOIN, TRANSFORM, PROJECT, SORT, DISTRIBUTE, MAP, NORMALIZE etc.; são funções de primeira classe, então operações básicas de dados podem ser implementadas em uma única linha de código. Isto torna a linguagem ECL uma linguagem de programação ideal para análise de dados, pois pode ser usada para expressar algoritmos de dados diretamente, eliminando a necessidade de escrever as especificações do software. Em essência, com o uso da linguagem ECL, os “Data Lakes” no HPC Systems precisam de menos programadores para entregar mais projetos em menor quantidade de tempo.

- Conceitos básicos de ECL
  - Paradigma declarativo (não-procedural);
  - ECL “não” é sensível a caixa alta/baixa (uppercase/lowercase)
  - Espaço em branco é ignorado para uma melhor leitura;
  - Comentários em linha (`//`) e em bloco (`/* e */`);
  - ECL utiliza sintaxe Objeto.Propriedade:
    - `Dataset.Campo // um campo em um dataset`
    - `NomedoDiretorio.Definicao // uma definição em outro módulo`
- Definições vs. Ações

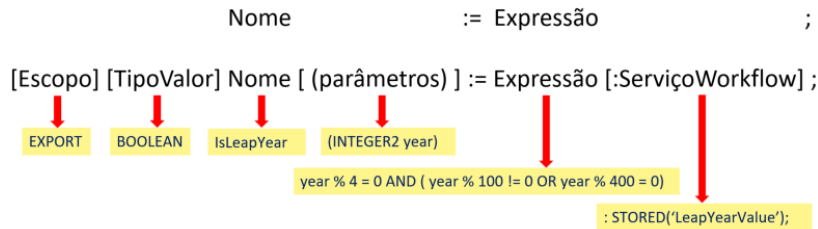
O código ECL é constituído de “Definições” e “Ações”.

  - Definições estabelecem o que as coisas são (arquivos de definição ECL):
    - `MyDef := 'Hello World'; // não inicia uma WorkUnit`



- Ações em ECL resultam em compilação e execução (arquivos BWR):
  - OUTPUT(MyDef); // inicia uma WorkUnit
  - OUTPUT('Hello World, again...'); // inicia uma WorkUnit

### Sintaxe completa de uma Definição ECL



### 1.1.6. Outras ferramentas da plataforma HPCC Systems

A plataforma HPCC Systems fornece para os desenvolvedores um ambiente de desenvolvimento integrado (IDE) denominado como “ECL IDE”, a fim de facilitar o desenvolvimento de código ECL. O “ECL IDE” é uma aplicação para o sistema operacional Windows. Há também uma extensão de linguagem ECL disponível para VS Code que alguns desenvolvedores preferem usar.

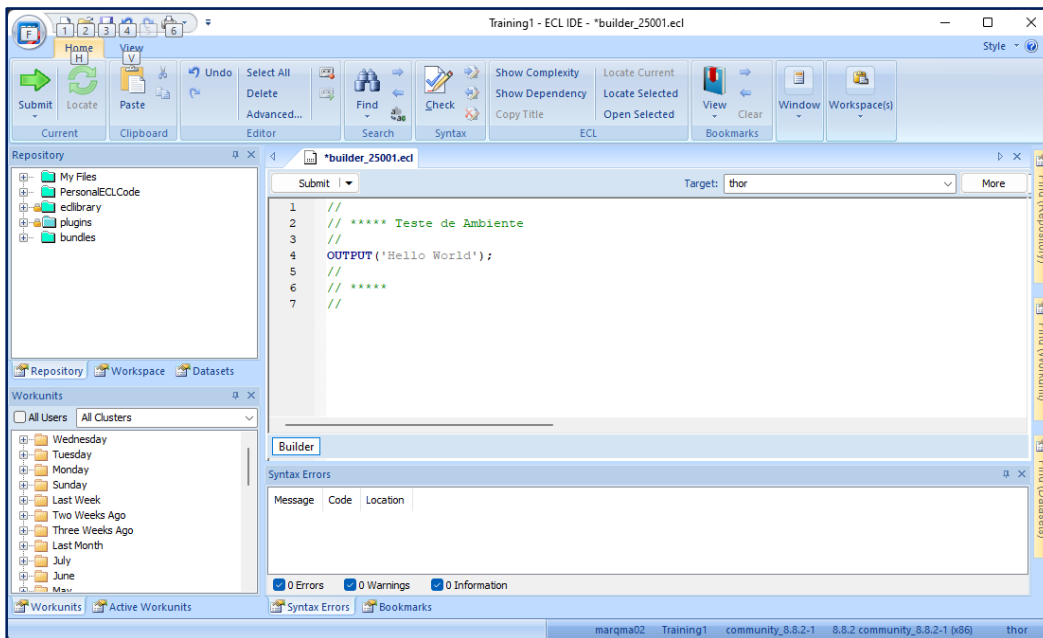


Figura 1.7: ECL Integrated Development Environment (IDE).



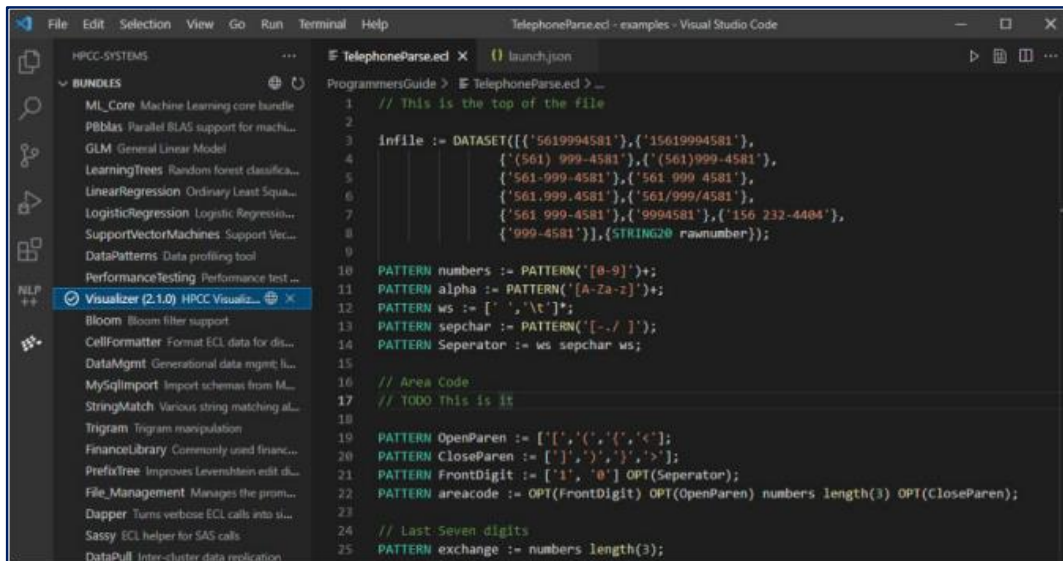


Figura 1.8: VS Code IDE.

## 1.2. Algoritmos de Aprendizagem de Máquina para a plataforma HPCC Systems

A característica principal do Aprendizagem de Máquina (Machine Learning) é a capacidade de inferir sobre relacionamentos e, visa prever uma resposta razoável quando apresentado com dados nunca vistos.

O "aprendizado" do Machine Learning (ML) tem várias categorias:

- Supervisionado - O tipo mais comum de ML. Esse método envolve o treinamento do sistema em que os recordsets, juntamente com o padrão de saída de destino, são fornecidos ao sistema para executar uma tarefa;
- Não Supervisionado - Este método não envolve a saída de destino, o que significa que nenhum treinamento é fornecido ao sistema. O sistema precisa aprender por meio da determinação e adaptação de acordo com as características estruturais nos padrões de entrada.
- Deep Learning - Move-se para a área dos métodos ML de Redes Neurais (Neural Networks - NNs). O Deep Learning implica várias camadas maiores que '2' e, também, implica em técnicas utilizadas com dados complexos, como análise de vídeo ou áudio.

### 1.2.1. Aprendizagem Supervisionado

Essa será a categoria abordada nesse estudo com foco no método de Árvores de Decisão (Learning Trees - Random Forests).

A premissa básica do ML Supervisionado é que, dado um conjunto de "amostras de dados" (registros) e um conjunto de "valores-alvo" em campo e formato de registro, que ele aprenda como prever "valores-alvo para novas amostras".

- Amostras de dados: conhecidas como variáveis "Independentes", porque são as informações fornecidas e não dependem de nenhum outro dado. Variáveis Independentes também são conhecidas como 'Características' (features) dos dados.

- Valores-alvo: conhecidos como variáveis “Dependentes”, porque eles são de alguma forma dependentes das amostras de dados.

As variáveis ‘Independentes’ e ‘Dependentes’ juntas são conhecidas como “conjunto de treinamento”.

Dois tipos básicos de Modelos de Análise em Aprendizado Supervisionado:

- Quantitativo - conhecido como “Regressão” - implica um valor numérico;
- Qualitativo - conhecido como “Classificação” - implica uma categoria ou, às vezes, um resultado binário.

### 1.3. Bundles de Machine Learning da plataforma HPCC Systems

Os bundles de produção (excluindo os bundles de suporte ML\_Core e PBblas) fornecem uma interface principal muito semelhante para o Machine Learning. No entanto, cada algoritmo tem suas próprias peculiaridades (suposições e restrições) que devem ser levadas em consideração. Portanto, é importante ler a documentação que acompanha cada bundle para usá-lo efetivamente.

- Link definitivo para todos os bundles de produção, sua documentação e tutoriais, quando aplicável:

<https://hpccsystems.com/download/free-modules/machine-learning-library>

a) Bundles Principais:

- ML\_Core - Machine Learning Core

Fornecer as principais definições de dados para ML. É um pré-requisito para todos os outros pacotes configuráveis de produção.

Mais informações: [https://github.com/hpcc-systems/ML\\_Core](https://github.com/hpcc-systems/ML_Core)

- PBblas - Parallel Block Basic Linear Algebra Subsystem

Fornecer operações de matriz escalonáveis e distribuídas usadas por vários dos outros pacotes configuráveis. Também pode ser usado diretamente sempre que as operações da matriz estiverem em ordem. Essa é uma dependência para vários dos outros pacotes configuráveis.

Mais informações: <https://github.com/hpcc-systems/PBblas>

b) Bundles de Aprendizado Supervisionado:

- LearningTrees (baseado no algoritmo clássico “ML RandomForest”)

Classificação e Regressão baseadas em Árvores de Decisão. Um dos melhores métodos de ML "prontos para uso", pois faz poucas suposições sobre a natureza dos dados e é resistente ao “overfitting”<sup>(\*)</sup>. Capaz de lidar com muitas variáveis independentes. Cria uma “floresta” de diversas Árvores de Decisão e calcula a média das respostas das diferentes Árvores.

Mais informações: <https://github.com/hpcc-systems/LearningTrees>

<sup>(\*)</sup> Superajuste (overfitting): Muitos algoritmos tendem a superestimar o conjunto de treinamento. Isso significa que ele está reduzindo o erro de previsão ajustando-se ao ruído que ocorre nesse conjunto de dados. Um modelo de excesso de ajuste tratará esse ruído como sinal e usará o que aprendeu para prever os próximos dados apresentados. Infelizmente, esse próximo conjunto de dados estará sujeito a um conjunto de ruído completamente diferente, o que não fornecerá bons resultados.

## 1.4. Interface Myriad

Todos os bundles de ML de produção suportam a interface Myriad, que é uma maneira de executar muitas ações semelhantes em diferentes conjuntos de dados, com uma única invocação da interface. Por exemplo, pode-se querer criar um modelo separado para cada grande área metropolitana (cidade/estado) e, então, usar esse conjunto de modelos para prever dados para cada área usando o seu próprio modelo exclusivo. A interface Myriad permite o processamento dessas atividades em paralelo.

Este material pressupõe que se tenha um conhecimento básico dos conceitos e terminologia de Machine Learning, que se esteja familiarizado com os conceitos básicos de uso dos bundles de ML nativos na plataforma HPCC Systems e que se tenha alguma experiência com programação em ECL.

### 1.4.1. Conceito

Conceitualmente, a interface Myriad...

- É um conjunto de padrões incorporados aos bundles de ML de produção nativos na plataforma HPCC Systems;
- É desenhada para executar várias atividades de Aprendizado de Máquina em uma única invocação da interface: Isso significa que se pode invocar um único processo de treinamento, o qual poderá retornar vários modelos, possibilitando fazer uma única previsão em dados desses vários modelos. Logo, pode-se avaliar a precisão de todos os modelos com uma única chamada para a função “Accuracy(...)”;
- É uma forma de maximizar a utilização dos recursos do cluster da plataforma HPCC Systems: Ao executar as atividades ao mesmo tempo, permite aproveitar totalmente o paralelismo do cluster, realizando todas as atividades em uma única invocação, contrastando com a execução serial, onde muitos nós podem ficar ociosos para qualquer atividade;
- Permite que todas as atividades sejam balanceadas em todos os nós do cluster;
- Fornece sincronização mais eficiente, necessária para muitos algoritmos de Machine Learning.

### 1.4.2. Usando a Interface Myriad

Como um exemplo de uso da interface Myriad, conforme já mencionado, será considerado que se deseje criar um modelo separado para cada grande área metropolitana (cidade/estado), a fim de que se possa fazer previsões mais granulares baseadas em dados de cada área.

Com a interface Myriad pode-se invocar um único treinamento, aqui processado através da função “GetModel(...)”, de maneira que retorne “n” modelos diferentes compactados em um só. Então, com uma única invocação de “Predict(...)”, pode-se fazer previsões sobre dados de todas as “n” áreas metropolitanas, cada uma em relação ao modelo aprendido apropriado para essa área. Pode-se, então, avaliar a precisão dos “n” modelos com uma única chamada para “Accuracy(...)”. Isso é conveniente do ponto de vista da programação, mas, mais importante ainda é poder fazer o melhor uso dos recursos do cluster da plataforma HPCC Systems.

Executar as atividades ao mesmo tempo permite aproveitar mais completamente o paralelismo do cluster da plataforma HPCC Systems realizando todas as atividades em uma invocação. Compare isso com a invocação serial de cada atividade. Se cada atividade não puder usar o conjunto completo de nós da plataforma, porque é de tamanho moderado, ou o algoritmo específico não pode utilizar todos os nós de uma vez, muitos dos nós ficarão ociosos para cada atividade. Executar as atividades uma de cada vez, portanto, permite que muitos nós sejam subutilizados durante toda a operação. Executar as atividades em paralelo com a interface Myriad, possibilita que as atividades sejam balanceadas em todos os nós do cluster, maximizando assim o desempenho (ou seja, minimizando o tempo de execução).

Além disso, muitos dos algoritmos de ML exigem várias etapas sequenciais para concluir uma atividade. Cada etapa requer sincronização entre os nós no cluster. Quanto mais sincronizações forem necessárias, menos otimizado será o processamento paralelo dentro do cluster. Executar várias atividades em paralelo não aumenta o número de sincronizações em relação a uma única atividade, embora cada sincronização se torne maior (ou seja, mais dados). Embora as sincronizações maiores afetem o desempenho, o impacto é pequeno em comparação ao aumento do número de etapas de sincronização. Isso ocorre devido à latência (atrasos) na comunicação de rede. Ao esperar por respostas de um nó para outro, nenhum trabalho pode ser feito. Requisições maiores não criam atrasos tão improdutivos porque não há espera, exceto pelo tempo real da transferência de dados na rede.

A estratégia da interface Myriad está incorporada nas principais estruturas de registro usadas pelo ML. Lembrando que, o ML utiliza vários tipos principais de registro em suas interfaces:

- NumericField é usado para fornecer dados de valor real em forma de matriz;
- DiscreteField é usado para transmitir dados discretos (de valor inteiro); e
- Layout\_Model é usado para codificar o Modelo que armazena tudo o que é aprendido sobre os dados.

Cada um desses layouts tem um campo conhecido como “work-item-id” (‘wi’). É esse parâmetro “id” que fornece isolamento das várias atividades independentes na interface Myriad.

Quando se fornece os dados independentes, pode-se usar um ‘wi’ diferente para os dados de cada área metropolitana. O mesmo pode ser feito para os dados de treinamento dependentes. Quando for chamado o “GetModel(...)”, o conjunto de dados retornado conterá todos os modelos separados, cada um identificado pelo seu respectivo ‘wi’. Por exemplo, o modelo que foi treinado por dados independentes e dependentes com ‘wi=1’ será rotulado com ‘wi=1’. O mesmo ocorre para cada ‘wi’. Agora, usando o(s) modelo(s) para prever um novo valor dependente, por meio de “Predict(...)” (para Regressão) ou “Classify(...)” (para Classificação), as previsões serão baseadas no modelo com o mesmo ‘wi’ que os itens de dados independentes. O mesmo é verdadeiro quando se avalia os modelos usando “Accuracy(...)”.

Na verdade, cada interface dentro dos bundles de ML de produção pode lidar com várias atividades independentes por meio do uso de “work-item-ids”.

A figura mostrada abaixo, ilustra esse fluxo de dados.

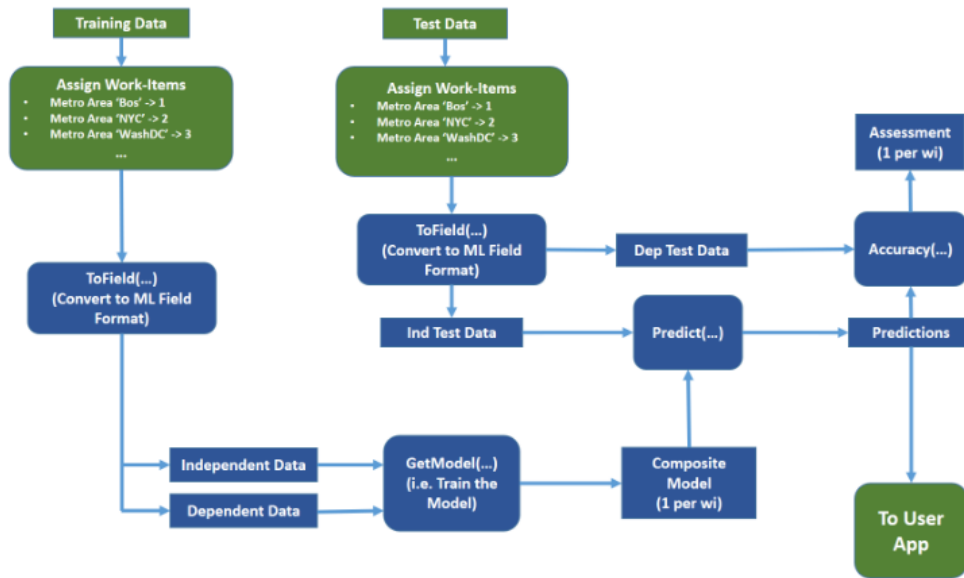


Figura 1.9: Como a interface Myriad funciona.

Mais adiante será apresentado um estudo de caso como exemplo da aplicação de algoritmos de Machine Learning fazendo uso da interface Myriad.

### 1.5. ML Supervisionado: Árvores de Decisão (Learning Trees - Random Forests)

As Árvores de Decisão têm sido utilizadas, pelo menos, desde a década de 1930 como uma forma de estruturar o conhecimento usando um conjunto de regras em cascata. Elas são conceitualmente simples e razoavelmente fáceis de entender e interpretar.

O LearningTrees bundle fornece uma implementação eficiente e escalável dos métodos Learning Trees. Atualmente, fornece algoritmos de "Decision Trees", "Random Forest", "Gradient Boosted Trees" e "Boosted Forest".

Diante dos diversos algoritmos disponíveis, qual algoritmo se deve escolher?

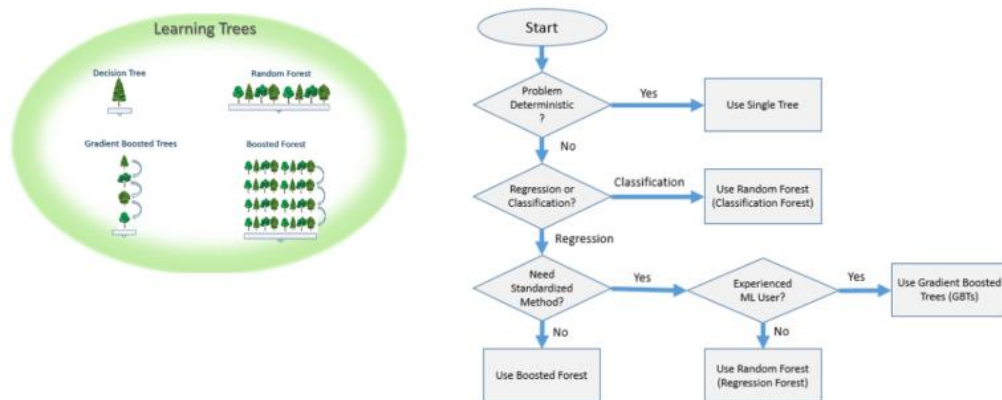
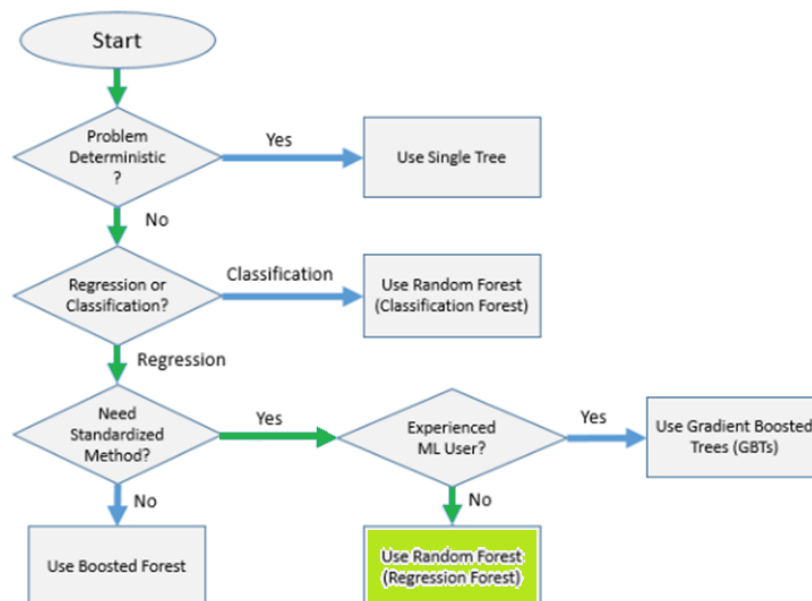
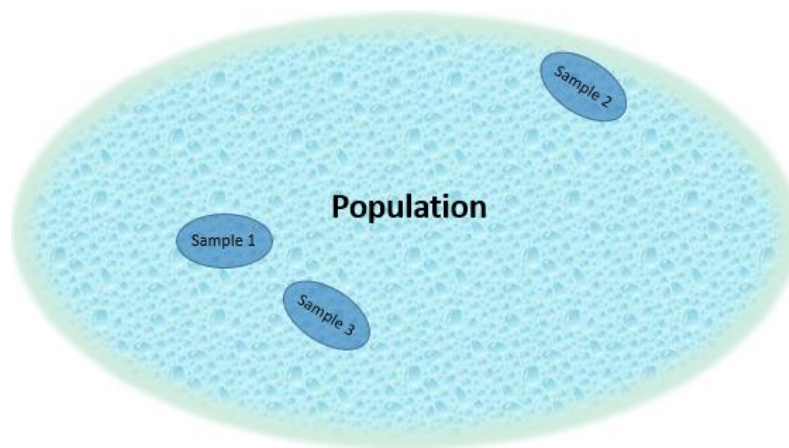


Figura 1.10: Fluxograma a ser usado para auxiliar nessa escolha.



**Figura 1.11: No caso em estudo “preço de imóveis” – Escolha: Random Forest - Regression Forest.**

Isso funciona muito bem desde que o problema seja ‘determinístico’ - ou seja, as mesmas variáveis (features) sempre produzem os mesmos resultados e, há dados de treinamento suficientes, desencadeando no uso de uma Árvore Simples. Se o problema for ‘estocástico’ - incorpora aleatoriedade nos dados ou resultados - uma Árvore de Decisão provavelmente não “generalizará” bem. Quanto ao conceito de “generalização” é importante entender a natureza da ‘População’ versus ‘Amostra’ (sample). Os dados de treinamento para um modelo de ML são quase sempre uma pequena amostra da população-alvo:



**Figura 1.12: Generalização.**

Em uma Random Forest, uma série de Árvores de Decisão são geradas, com alguma aleatoriedade adicionada ao processo de geração para garantir que cada árvore use um processo de decisão diferente na separação dos pontos. Então, para cada novo ponto de dados, todas as árvores são consultadas para sua previsão. Essas previsões individuais são agregadas para formar a previsão final.

Se o objetivo da Random Forest é escolher uma das várias classes (ou seja, uma ClassificationForest), então a agregação é feita por votação: a “classe prevista pelo maior número de árvores” torna-se a previsão final. Se o objetivo é prever um valor numérico (como o preço de imóveis), então temos uma RegressionForest, que forma uma previsão final pela “média das previsões das várias árvores”.

Um dos aspectos interessantes do Random Forest é que há muito poucos parâmetros para ajustar, e a escolha desses parâmetros geralmente tem muito pouca influência na precisão dos resultados. Bons resultados geralmente podem ser alcançados usando os parâmetros "default", tornando-o um dos algoritmos de ML mais fáceis de usar.

Outra grande característica é que o Random Forest pode lidar com muitas variáveis. Não é incomum usá-lo com milhares de variáveis.

Por fim, Random Forest são facilmente “paralelizáveis”, sendo, portanto, um algoritmo ideal para uso em clusters do HPC Systems.

### 1.5.1. O Fluxo de Aprendizado de Máquina Supervisionado

No estudo de caso foi selecionado o bundle “LearningTrees ML” pelos seguintes motivos:

1. Fácil de usar:
  - Faz muito pouca suposição sobre a distribuição dos dados ou seus relacionamentos;
  - Ele pode lidar com muitos registros e muitos campos;
  - Ele pode lidar com relações não lineares e descontínuas;
  - Quase sempre funciona bem usando os parâmetros padrão e sem nenhum ajuste.
2. Escalabilidade
  - Escala bem em clusters HPC Systems de quase qualquer tamanho.
3. Sua precisão de previsão é competitiva com os melhores algoritmos de última geração.



Figura 1.13: Fluxo de ML Supervisionado.

#### 1.5.1.1. Definição do problema

“Dado um conjunto de atributos de uma propriedade (localização, metragem quadrada, ano de construção/aquisição, número de cômodos, etc.), como prever o seu valor real de venda?”



propertyid	house_numbr	house_m	predir	street	street	postdir	apt	city	state	zip	total_value	assessed_value	year_acquired	land_square_foot	living_square_fee	bedrooms	full_bath
028195	144			MCKERNAN	DR			WALNUT CREEK	CA	94597	62614	52614	2006	20418	2485	3	2
1144455	281			CENTER	ST			BALTIMORE	MD	21136	105500	10550	2007	4807	1368	0	0
1494347	483			NEWTON	RD			FLAGSTAFF	AZ	86001	2220	2220	0	5654	1811	3	1
1910847	882			HITCHCOCK	CT			WOODLAND	MA	98674	156800	156800	0	6094	0	2	1
4267962	5807		E	ROY ROGERS	RD			TROY	HI	48069	227253	227253	2007	3464	0	3	0
4888682	7687			PERRLESTONE	DR		000009	KERIVILLE	CA	92239	721179	721179	2010	19597	6132	6	6
48725	4			LONG	AVE			SUNRISE	FL	33325	271000	271000	2008	6880	2392	4	2
83520	6			TRILLIUM	LN			WAYLAND	MA	02153	79889	79889	2007	7657	1657	4	1
94604	7			PARIENTER	AVE			PLYMOUTH	MI	55441	23800	23800	2005	19994	1754	3	2
220326	17			TINGER	RD			LOS ANGELES	CA	90003	89000	89000	2008	7840	954	3	1
994609	212			FREYER	DR	NE		PHILOMONT	VA	20131	59800	59800	2009	11199	1241	3	0
1836173	724			EASTER	ST			ALLENTOUN	PA	18102	191600	191600	0	9100	2534	4	2
2910797	1903			SADDLE BROOK	DR			CLIO	CA	96106	61610	11610	2007	0	0	0	0
3083959	2158			RIVERSIDE	DR			UPPER MORELA	PA	19006	90300	0	0	1235	3	2	0
3952189	4040			GRAND VIEW	BLVD		000054	RIO LINDA	CA	95673	0	0	0	2700720	0	0	0
4186238	4726			LAS PALMAS	CT			WAELEDER	TX	78959	18816	18816	2009	2159	1320	0	0
4597143	6213			WILSON	RD			ZOLFO SPRINGS	FL	33880	72600	0	0	8496	0	3	1
4624905	6321			STONEWALL	LN			PATERSON	NJ	07514	139880	139880	2008	10454	1391	4	2
92326	7			KNOLLCREST	DR			NARANJA	FL	33032	76214	76214	2008	4800	930	2	0
1792852	784			ERIN	DR			TRABUCO	CA	92678	28010	28010	2007	5200	0	3	1
1843977	728		S	ARLINGTON HE.	RD			BLOOMING GRL.	TX	76626	130400	130400	2007	36154	1629	3	1
4114872	4811			HORTLE GUY	DR		000015	SAN FERRANDO	CA	91324	18887	18887	2007	92654	0	0	0

Figura 1.14: Preço do imóvel – Valor real de venda (total\_value).

### 1.5.1.2. Extração dos dados

Será utilizado o dataset “Property.csv” no formato “CSV” contendo 1.662.959 registros. Usando os campos selecionados, tentaremos prever o preço da propriedade com base em outros dados (por exemplo: localização, metragem quadrada, ano de aquisição/construção, número de cômodos, código postal, etc.).

Arquivo lógico: “~CLASS::XYZ::ML::Property”.

- Códigos ECL utilizados:  
modProperty.ecl  
BWR\_BrowseData.ecl

### 1.5.1.3. Preparação dos dados

A preparação dos dados é a etapa mais importante ao usar qualquer algoritmo ML. Como em qualquer processo de programação, a qualidade dos dados recebidos terá um grande efeito na qualidade dos resultados.

Aqui estão algumas regras importantes a serem consideradas:

- Os dados devem conter todos os valores numéricos;
- O primeiro campo no registro deve ser um identificador exclusivo (geralmente um ID de registro sem sinal - UNSIGNED);
- Para facilitar a implementação, mova seu campo “dependente” para o final da estrutura RECORD;
- Randomize seus dados para criar um recordset de treinamento e um teste mais preciso. Isso pode ser feito adicionando um campo com um número aleatório (RANDOM);
- Faça previamente a limpeza dos dados (cleasing);
- Durante a preparação dos dados, atribua “work-item-ids” com base no campo ‘state’ através da função MAP, a fim de criar um modelo de ML separado para cada estado norte-americano existente no dataset “Property”;
- Usando o campo aleatório gerado, classifique e segregue os dados nos dados iniciais de Treinamento e de Teste. Os dados de Treinamento serão usados para treinar seu modelo de ML e os dados de Teste serão usados para avaliar (ou analisar) a eficácia do modelo. É fundamental que se reserve alguns dos dados para teste, pois é uma péssima ideia testar o modelo com os mesmos dados nos quais se treinou.

- Códigos ECL utilizados:  
isCleanFilter.ecl  
CleanProperty.ecl  
modPrepData.ecl  
BWR\_ViewData1.ecl

##	propertyid	zip	assessed_value	year_acquired	land_square_footage	living_square_feet	bedrooms	full_baths	half_baths	year_built	total_value
1	79784	33424	76440	2015	4299	1255	3	2	0	2010	76440
2	3924129	20601	95900	2013	11224	1468	3	2	1	2007	95900
3	413843	8803	76000	2015	57000	1858	3	2	0	1970	76000
4	608224	98370	39340	2012	7405	1066	3	1	1	1967	39340
5	942963	72032	278400	2008	9600	2459	3	2	0	1963	278400
6	2237271	79935	143600	2011	8430	1008	2	1	1	1961	143600
7	4443742	84065	166934	2013	9317	1700	4	2	0	1991	166934
8	3834707	66227	348350	2012	15300	2663	4	2	1	2002	348350
9	3592739	19606	54000	2015	15060	2292	4	2	1	1980	90000
10	2916349	34639	119050	2015	6947	1709	3	2	0	2009	140950

**Figura 1.15: Limpeza, padronização e consolidação de registros.**

#### 1.5.1.4. Segregação dos dados

Ao segregar os dados é importante coletar amostras aleatoriamente do seu dataset, ao invés de usar os primeiros ‘N’ registros no conjunto, porque pode existir alguma ordem oculta no processo que originalmente gerou os dados. A segregação dos dados de Treinamento e de Teste é feita facilmente usando ECL.

Como regra geral, cerca de 20 a 30% dos seus dados devem ser reservados para Teste. Entretanto, para fins didáticos e, principalmente, objetivando uma melhor performance durante a realização desse estudo, esse percentual será reduzido, considerando uma proporção de 5.000 registros para a amostra de Treinamento e 2.000 registros para a amostra de Teste:

```
// Considerando os primeiros 5.000 registros como amostra de Treinamento
MyriadTrainData := PROJECT(MyriadPrepDataSort[1..5000],
    $.modPrepData.ML_Prop) :PERSIST('~CLASS::XYZ::ML::MyriadTrain');
// Considerando os 2.000 registros seguintes como amostra de Teste
MyriadTestData := PROJECT(MyriadPrepDataSort[5001..7000],
    $.modPrepData.ML_Prop) :PERSIST('~CLASS::XYZ::ML::MyriadTest');
```

O próximo passo é converter os dados para o formato usado pelos pacotes configuráveis ML. Para operar genericamente com qualquer dado, o ML requer que os dados estejam em um “layout de matriz orientado a célula” conhecido como “NumericField”, considerando o parâmetro “wifield” associado ao novo campo ‘wi\_id’. O bundle configurável ML\_Core facilita essa tarefa:

```
IMPORT ML_Core;
// Conversão Matricial dos campos numéricos
ML_Core.ToField(MyriadTrainData, MyriadTrainDataNF, wifield := wi_id);
// "idfield" parâmetro não fornecido, assume o 1º campo 'propertyid'
ML_Core.ToField(MyriadTestData, MyriadTestDataNF, wifield := wi_id);
// "idfield" parâmetro não fornecido, assume o 1º campo 'propertyid'
```

A etapa final antes de aplicar o modelo ML é separar os dados “independentes” dos dados “dependentes”, definindo os dados do ML RECORD para colocar o campo de dados dependente no “final” do layout de registro. Um filtro simples nos dados de treinamento e de teste conclui esta etapa.

Sempre deve ser excluído o ID do registro exclusivo e, conte apenas seus campos independentes e dependentes. No caso, os dados de treinamento do dataset “Property”, tem nove (9) campos independentes e o último campo (10º) é o campo dependente (o valor que será previsto).

O uso do PROJECT definindo o campo numérico = 1 não é estritamente necessário. Isso indica que é o primeiro campo dos dados dependentes. Como existe apenas um campo dependente, ele foi numerado de acordo:

```
EXPORT MyriadIndTrainDataNF := MyriadTrainDataNF(number < 10);
EXPORT MyriadDepTrainDataNF := PROJECT(MyriadTrainDataNF(number =
10),
TRANSFORM(RECORDOF(LEFT), SELF.number := 1, SELF := LEFT));
EXPORT MyriadIndTestDataNF := MyriadTestDataNF(number < 10);
EXPORT MyriadDepTestDataNF := PROJECT(MyriadTestDataNF(number = 10),
TRANSFORM(RECORDOF(LEFT), SELF.number := 1, SELF := LEFT));
```

- Códigos ECL utilizados:
  - modSegConvData.ecl
  - BWR\_ViewData2.ecl

##	wi	id	number	value
1	22	2292044	1	6751.0
2	22	2292044	2	75837.0
3	22	2292044	3	2013.0
4	22	2292044	4	3484.0
5	22	2292044	5	763.0
6	22	2292044	6	1.0
7	22	2292044	7	1.0
8	22	2292044	8	1.0
9	22	2292044	9	2002.0
10	16	3675939	1	12546.0

##	wi	id	number	value
1	22	2292044	1	75837.0
2	16	3675939	1	131500.0
3	3	4346206	1	78954.0
4	3	4168683	1	65000.0
5	18	2956615	1	168821.0
6	41	2622288	1	65568.0
7	18	1196844	1	36100.0
8	43	114945	1	19796.0
9	6	961534	1	156893.0
10	41	882691	1	71786.0

**Figura 1.16: Nove (9) campos independentes e o último campo (10º) é o campo dependente, com o work-item-id (wi) associado aos estados norte-americanos.**

### 1.5.1.5. Treinamento e avaliação do modelo

O primeiro passo é selecionar o modelo "learner". Para Regressão (quantitativa), será usado o módulo RegressionForest:

```
IMPORT $;
IMPORT LearningTrees AS LT;
/* Selecionando o algoritmo...
  Sintaxe:
  RegressionForest(UNSIGNED numTrees=100, UNSIGNED featuresPerNode=0,
  UNSIGNED maxDepth=100, SET OF UNSIGNED nominalFields=[]);
*/
```

```
MyriadLearnerR := LT.RegressionForest(); // parâmetros default
```

Em seguida, o “learner” será usado para treinar e recuperar o modelo:

```
MyriadModelR :=  
  MyriadLearnerR.GetModel($.modSegConvData.MyriadIndTrainDataNF,  
    $.modSegConvData.MyriadDepTrainDataNF);
```

Observe que o modelo obtido é uma estrutura de dados opaca (interpretável apenas pelo bundle) que encapsula todos os resultados do treinamento. O primeiro parâmetro fornecido para a função “GetModel” é o dataset de treinamento “independente” e o segundo parâmetro é o dataset de treinamento “dependente”.

Em seguida, o modelo será usado para fazer previsões, prevendo o campo dependente com base nos campos de teste independentes:

```
MyriadPredictedDeps := MyriadLearnerR.Predict(MyriadModelR,  
  $.modSegConvData.MyriadIndTestDataNF);
```

Portanto, aplicando o modelo de treinamento e prevendo os resultados dependentes com base nos dados independentes, permitirá realmente produzir e visualizar esses resultados para análise.

Após o treinamento, será testado o modelo de previsão RegressionForest comparando-o com o dataset de teste dependente [MyriadDepTestDataNF], que foi extraído anteriormente dos registros de teste segregados. ML\_Core tem uma ótima maneira de fazer isso no módulo “Analysis”. Esse módulo fornece uma avaliação genérica do poder preditivo do modelo. Também há uma grande variedade de outras métricas de avaliação fornecidas pelos pacotes individuais, específicas para os algoritmos desse pacote. No Modelo de Regressão do LearningTrees, basta uma linha de ECL para realizar essa análise:

```
MyriadAssessmentR :=  
  ML_Core.Analysis.Regression.Accuracy(MyriadPredictedDeps,  
    $.modSegConvData.MyriadDepTestDataNF);
```

A função “Accuracy” usa o resultado da função “Predict” como primeiro parâmetro, comparando-o com os dados de teste dependentes que foi criado anteriormente. Aqui estão os resultados dos dados de treinamento para o dataset “Property” (vide Nota):

- Métricas chaves:
  - $r^2$  (R-Quadrado - O "coeficiente de determinação". Aproximadamente, a proporção da variação nos dados originais que foram capturados pela regressão. O R-Quadrado pode variar ligeiramente de negativo a 1. Um R-Quadrado 1 significa que as previsões correspondem perfeitamente às reais. R-Quadrado zero ou abaixo significa que o modelo não tem valor preditivo. R-Quadrado 0,5 significa que metade da variação dos dados foi explicada pelo modelo.
  - MSE (Mean Squared Error - Erro quadrático médio) - O desvio médio quadrático entre os valores previstos e reais.
  - RMSE (Root Mean Squared Error - Raiz do Erro quadrático médio) - A raiz quadrada do MSE. Essa é uma expectativa do erro médio em uma determinada amostra.

##	wi	regressor	r2	mse	rmse
1	48	1	1.0	7899654400.0	88880.0
2	40	1	0.9311376324090863	86100656.52500002	9279.043944555928
3	6	1	0.8407434005563046	1188713430.574167	34477.72368608703
4	53	1	0.8065354715574288	2619422585.198289	51180.29489166987
5	37	1	0.8051321121617373	3681390045.728787	60674.45958332705
6	41	1	0.8033767773856457	2630432563.605217	51287.74282033883
7	32	1	0.8003100933540183	1700812118.14	41240.90345930846
8	22	1	0.798261028979762	953955898.739406	30886.17649919468
9	34	1	0.7728380339814155	522134100.9433332	22850.25384855348
10	17	1	0.7696693188373955	2411024752.625271	49102.18684157836

Figura 1.17: Métricas chaves – Ordenação pelo Coeficiente de Determinação (r2).

Nota: Para o conjunto de treinamento aplicado em “Property” (dataset bruto), a precisão foi um pouco abaixo de 30% na primeira tentativa, onde para melhorar a precisão foi necessário revisar os campos independentes. Após uma análise mais aprofundada, pode-se verificar que o campo ‘zip’ não é realmente quantitativo, mas qualitativo (variável categórica). Então, a seguinte modificação no modelo de aprendizado foi realizada:

```
MyriadLearnerR := LT.RegressionForest(,,[1]);
```

A etapa final foi revisar as amostras e remover os registros com valores de campo “nulos” das amostras de treinamento e de teste, tendo sido usado para tal procedimento o código [CleanProperty].

➤ Códigos ECL utilizados:

```
BWR_TrainReg.ecl
```

### 1.5.1.6. Implantação do modelo

Essa etapa corresponde ao desenvolvimento de um serviço de consulta através da codificação de uma “função” (estrutura Function) da linguagem ECL, seguida da compilação do código no cluster Roxie:

The screenshot shows the Roxie web interface for a dynamic form. The form is titled "FN\_MYRIADGETPRICEREGREQUEST" and contains several input fields for variables. The variables and their values are: \_state\_ (CA), assess\_val (118720), bedrooms (3), full\_baths (2), half\_baths (1), land\_sq\_ft (14774), living\_sq\_ft (1437), year\_acq (2011), year\_built (1968), and zip (95451). At the bottom of the form, there are checkboxes for "Capture Log Info." and "No Timeout", and buttons for "Call Query", "Output Tables", "FORM POST", "Submit", and "Clear All".

roxie

fn\_myriadgetpricereg\_web Dynamic Form

**FN\_MYRIADGETPRICEREG\_WEBREQUEST**

\_state\_01: CA

\_state\_02: NY

\_state\_03: FL

assess\_val: 118720

bedrooms: 3

full\_baths: 2

half\_baths: 1

land\_sq\_ft: 14774

living\_sq\_ft: 1437

year\_acq: 2011


year\_built: 1968

zip: 95451

Capture Log Info. Trace Level:   No Timeout

Call Query Output Tables FORM POST Submit Clear All

Figura 1.18: Carregamento de dados e disponibilização de serviços de consulta através de ESP (Enterprise Service Platform).



## SSCAD 2024

### \$ Preço do Imóvel por Estado 🏠

Estados (EUA):

Valor avaliado:

Nº de quartos:

Nº de banheiros:

Nº de lavabos:

Área construída (ft<sup>2</sup>):

Área interna (ft<sup>2</sup>):

Ano de aquisição [AAAA]:

Ano de construção [AAAA]:

CEP:


Consulta

Reset

Mauro Marques

**Estados dos EUA disponíveis na base de dados:**

AK	Alasca	MT	Montana
AL	Alabama	NC	Carolina do Norte
AR	Arkansas	ND	Dakota do Norte
AZ	Arizona	NE	Nebraska
CA	Califórnia	NH	Nova Hampshire
CO	Colorado	NJ	Nova Jérsei
CT	Connecticut	NM	Novo México
DE	Delaware	NV	Nevada
FL	Flórida	NY	Nova Iorque
GA	Geórgia	OH	Ohio
HI	Havaí	OK	Oklahoma
IA	Iowa	OR	Oregon
ID	Idaho	PA	Pensilvânia
IL	Illinois	RI	Rhode Island
IN	Indiana	SC	Carolina do Sul
KS	Kansas	SD	Dakota do Sul
KY	Kentucky	TN	Tennessee
LA	Louisiana	TX	Texas
MA	Massachusetts	UT	Utah
MD	Maryland	VT	Vermont
ME	Maine	VA	Virgínia
MI	Michigan	WA	Washington
MN	Minnesota	WI	Wisconsin
MO	Missouri	WV	Virgínia Ocidental
MS	Mississippi	WY	Wyoming



**Treinamento na plataforma**  
HPCC Systems

Copyright © 2024  
LexisNexis Risk Solutions

**\$ Preço do Imóvel por Estado** 💰

Estados (EUA): 01 02 03

Valor avaliado:

Nº de quartos:

Nº de banheiros:

Nº de lavabos:

Área construída (ft<sup>2</sup>):

Área interna (ft<sup>2</sup>):

Ano de aquisição [AAAA]:

Ano de construção [AAAA]:

CEP:

**Consulta**

**Reset**

*Mauro Marques*

**Figura 1.19: Carregamento de dados e disponibilização de serviços de consulta através de interface Web.**

➤ Códigos ECL utilizados:

fn\_MyriadGetPriceReg.ecl

fn\_MyriadGetPriceReg\_Web.ecl (código usado na interface Web)

BWR\_TestQueriesReg.ecl



### 1.5.2. Conclusão

Algoritmos baseados em Árvores de Decisão são provavelmente os algoritmos mais poderosos e fáceis de usar no arsenal de ML. Se os dados forem numéricos, houver uma quantidade razoável de dados e o objetivo for “Previsão” ao invés de “Explicação”, será difícil encontrar um melhor algoritmo do que Learning Trees.

Através desse estudo foi possível saber tudo o que se precisa saber para criar, treinar, prever, avaliar e testar modelos de Machine Learning fazendo uso da interface Myriad com base nos dados disponíveis e usá-los para prever valores quantitativos (Regressão). Caso se deseje usar um bundle de ML diferente, se concluirá que todos os bundles operam de maneira muito semelhante, com algumas variações menores.

No exemplo, foi treinado um número de modelos associados a quantidade de estados norte-americanos existentes no dataset “Property”, correspondente à cardinalidade de 44 atributos distintos para o campo ‘state’, sendo que todos os modelos eram homogêneos — com base no mesmo layout de registro e relacionados aos mesmos dados dependentes (preço de imóveis). Poderia ter sido usado datasets completamente diferentes que não tinham nada a ver um com o outro, ou poderia ter sido treinado modelos em dados dependentes completamente diferentes se assim fosse desejado

A simplicidade conceitual de ML é um tanto enganadora. Cada algoritmo tem suas próprias peculiaridades (ou premissas e restrições) que precisam ser levadas em consideração para maximizar a precisão preditiva. Além disso, quantificar a eficácia das previsões requer habilidades em ciência de dados e estatística que a maioria não possui. Por esses motivos, é muito importante que não seja utilizada a exploração de ML para produzir produtos ou reivindicar habilidades sem antes consultar especialistas no campo.

### 1.5.3. Apêndice

#### ➤ Myriad RoadMap

- Procedimento: Sequência de criação, execução (submit) e compilação (compile) dos códigos:  
modProperty.ecl  
BWR\_BrowseData.ecl (submit)  
isCleanFilter.ecl  
CleanProperty.ecl  
modPrepData.ecl  
BWR\_ViewData1.ecl (submit)  
modSegConvData.ecl  
BWR\_ViewData2.ecl (submit)  
BWR\_TrainReg.ecl (submit)  
fn\_MyriadGetPriceReg.ecl (compile em hThor & Roxie)  
fn\_MyriadGetPriceReg\_Web.ecl (compile em hThor & Roxie)  
BWR\_TestQueriesReg.ecl (submit)

#### ➤ Códigos ECL

```
[modProperty.ecl]
//
EXPORT modProperty := MODULE
EXPORT Layout := RECORD
  UNSIGNED8 personid;
```

```

INTEGER8 propertyid;
STRING10 house_number;
STRING10 house_number_suffix;
STRING2 predir;
STRING30 street;
STRING5 streettype;
STRING2 postdir;
STRING6 apt;
STRING40 city;
STRING2 state;
STRING5 zip;
UNSIGNED4 total_value;
UNSIGNED4 assessed_value;
UNSIGNED2 year_acquired;
UNSIGNED4 land_square_footage;
UNSIGNED4 living_square_feet;
UNSIGNED2 bedrooms;
UNSIGNED2 full_baths;
UNSIGNED2 half_baths;
UNSIGNED2 year_built;
END;
EXPORT File := DATASET('~CLASS::XYZ::ML::Property',Layout,CSV);
END;
//

```

#### [BWR\_BrowseData.ecf]

```

IMPORT $;
//
// Visualização da extração dos dados
OUTPUT($.modProperty.File);
COUNT($.modProperty.File); // 1.662.959 registros
//

```

#### [isCleanFilter.ecf]

```

IMPORT $;
//
Property := $.modProperty.File;
//
// Limpando os dados...
EXPORT isCleanFilter := Property.zip <> " AND
    Property.assessed_value <> 0 AND
    Property.year_acquired <> 0 AND
    Property.land_square_footage <> 0 AND
    Property.living_square_feet <> 0 AND
    Property.bedrooms <> 0 AND
    Property.full_baths <> 0 AND
    Property.year_Built <> 0;
//

```

#### [CleanProperty.ecf]

```

IMPORT $;
//
Property := $.modProperty.File;
//
EXPORT CleanProperty := Property($.isCleanFilter);
//
// CleanProperty := Property($.isCleanFilter);

```

```

// OUTPUT(CleanProperty);
// COUNT(CleanProperty);    // 575.814 registros
//

[modPrepData.ecl]
IMPORT $;
//
Property := $.modProperty.File;
//
EXPORT modPrepData := MODULE
EXPORT ML_Prop := RECORD
UNSIGNED8 propertyid;    // corresponde ao campo "idfield" usado na conversão matricial pelo
ML_Core.ToField
UNSIGNED3 zip;          // variável Categórica
UNSIGNED4 assessed_value;
UNSIGNED2 year_acquired;
UNSIGNED4 land_square_footage;
UNSIGNED4 living_square_feet;
UNSIGNED2 bedrooms;
UNSIGNED2 full_baths;
UNSIGNED2 half_baths;
UNSIGNED2 year_built;
UNSIGNED4 total_value;    // variável Dependente (a ser determinada)
UNSIGNED4 wi_id;          // Work-Item ID usado na Interface Myriad
END;
//
EXPORT ML_PropExt := RECORD(ML_Prop)
UNSIGNED4 rnd;            // número randômico
END;
//
EXPORT MyriadPrepData := PROJECT($.CleanProperty, TRANSFORM(ML_PropExt,
SELF.rnd := RANDOM(),
SELF.zip := (UNSIGNED3)LEFT.zip,
SELF.wi_id := MAP(LEFT.state = 'AE' => 1, LEFT.state = 'WY' => 2,
LEFT.state = 'TN' => 3, LEFT.state = 'MI' => 4,
LEFT.state = 'RI' => 5, LEFT.state = 'SC' => 6,
LEFT.state = 'VA' => 7, LEFT.state = 'NH' => 8,
LEFT.state = 'AP' => 9, LEFT.state = 'NM' => 10,
LEFT.state = 'IL' => 11, LEFT.state = 'KS' => 12,
LEFT.state = 'ME' => 13, LEFT.state = 'AL' => 14,
LEFT.state = 'WA' => 15, LEFT.state = 'NY' => 16,
LEFT.state = 'MO' => 17, LEFT.state = 'PA' => 18,
LEFT.state = 'HI' => 19, LEFT.state = 'GU' => 20,
LEFT.state = 'VT' => 21, LEFT.state = 'CT' => 22,
LEFT.state = 'NC' => 23, LEFT.state = 'OR' => 24,
LEFT.state = 'IA' => 25, LEFT.state = 'DE' => 26,
LEFT.state = 'VI' => 27, LEFT.state = 'ID' => 28,
LEFT.state = 'MT' => 29, LEFT.state = 'AR' => 30,
LEFT.state = 'MS' => 31, LEFT.state = 'UT' => 32,
LEFT.state = 'NE' => 33, LEFT.state = 'IN' => 34,
LEFT.state = 'GA' => 35, LEFT.state = 'WV' => 36,
LEFT.state = 'NJ' => 37, LEFT.state = 'LA' => 38,
LEFT.state = 'WI' => 39, LEFT.state = 'AK' => 40,
LEFT.state = 'CA' => 41, LEFT.state = 'NV' => 42,
LEFT.state = 'FL' => 43, LEFT.state = 'MA' => 44,
LEFT.state = 'CO' => 45, LEFT.state = 'AZ' => 46,
LEFT.state = 'SD' => 47, LEFT.state = 'DC' => 48,
LEFT.state = 'KY' => 49, LEFT.state = 'MP' => 50,

```

```

LEFT.state = 'ND' => 51, LEFT.state = 'AS' => 52,
LEFT.state = 'TX' => 53, LEFT.state = 'PR' => 54,
LEFT.state = 'MD' => 55, LEFT.state = 'OH' => 56,
LEFT.state = 'MN' => 57, LEFT.state = 'OK' => 58,
LEFT.state = 'AA' => 59, 0),
        SELF := LEFT))
:PERSIST('~CLASS::XYZ::ML::MyriadPrepData');
//
END;
//

[BWR_ViewData1.ecf]
IMPORT $;
//
// Preparação dos dados
OUTPUT($.modPrepData.MyriadPrepData);
COUNT($.modPrepData.MyriadPrepData); // 575.814 registros
//

[modSegConvData.ecf]
IMPORT $;
IMPORT ML_Core;
//
MyriadPrepData := $.modPrepData.MyriadPrepData;
//
// Torne os dados aleatórios, ordenando os registros pelo número randômico
MyriadPrepDataSort := SORT(MyriadPrepData(wi_id <> 0), rnd);
//
//
// Segregação dos dados - Considerando os primeiros 5.000 registros como amostra de Treinamento
MyriadTrainData := PROJECT(MyriadPrepDataSort[1..5000], $.modPrepData.ML_Prop)
:PERSIST('~CLASS::XYZ::ML::MyriadTrain'); // layout sem o campo rnd
//
// Segregação dos dados - Considerando os 2.000 registros seguintes como amostra de Teste
MyriadTestData := PROJECT(MyriadPrepDataSort[5001..7000], $.modPrepData.ML_Prop)
:PERSIST('~CLASS::XYZ::ML::MyriadTest'); // layout sem o campo rnd
//
//
// Conversão Matricial dos campos numéricos
ML_Core.ToField(MyriadTrainData, MyriadTrainDataNF, wifield := wi_id); // "idfield" não fornecido,
assume 1º campo = propertyid
ML_Core.ToField(MyriadTestData, MyriadTestDataNF, wifield := wi_id); // "idfield" não fornecido,
assume 1º campo = propertyid
//
//
EXPORT modSegConvData := MODULE
EXPORT MyriadIndTrainDataNF := MyriadTrainDataNF(number < 10); // excluindo o campo
propertyid
//
EXPORT MyriadDepTrainDataNF := PROJECT(MyriadTrainDataNF(number = 10),
TRANSFORM(RECORDOF(LEFT),
        SELF.number := 1,
        SELF := LEFT));
//
EXPORT MyriadIndTestDataNF := MyriadTestDataNF(number < 10); // excluindo o campo
propertyid
//

```

```
EXPORT MyriadDepTestDataNF := PROJECT(MyriadTestDataNF(number = 10),
TRANSFORM(RECORDOF(LEFT),
```

```
SELF.number := 1,
SELF := LEFT));
```

```
END;
```

```
//
```

```
[BWR_ViewData2.ecl]
```

```
IMPORT $;
```

```
//
```

```
// Segregação dos dados e Conversão matricial dos campos numéricos
```

```
OUTPUT($.modSegConvData.MyriadIndTrainDataNF, NAMED('MyriadIndTrainData'));
```

```
OUTPUT($.modSegConvData.MyriadDepTrainDataNF, NAMED('MyriadDepTrainData'));
```

```
OUTPUT($.modSegConvData.MyriadIndTestDataNF, NAMED('MyriadIndTestData'));
```

```
OUTPUT($.modSegConvData.MyriadDepTestDataNF, NAMED('MyriadDepTestData'));
```

```
//
```

```
[BWR_TrainReg.ecl]
```

```
IMPORT $;
```

```
IMPORT ML_Core;
```

```
IMPORT LearningTrees AS LT;
```

```
//
```

```
// Selecionando o algoritmo
```

```
// Sintaxe: RegressionForest(UNSIGNED numTrees=100, UNSIGNED featuresPerNode=0,
```

```
// UNSIGNED maxDepth=100, SET OF UNSIGNED nominalFields=[])
```

```
// MyriadLearnerR := LT.ReggressionForest(); // parâmetros default
```

```
// MyriadLearnerR := LT.ReggressionForest(,,,[1]); // zip não é realmente quantitativo, mas qualitativo
```

```
MyriadLearnerR := LT.ReggressionForest(10,,10,[1]);
```

```
//
```

```
//
```

```
// Obtendo o modelo composto (composto por 45 modelos = números de estados existentes no dataset)
```

```
MyriadModelR :=
```

```
MyriadLearnerR.GetModel($.modSegConvData.MyriadIndTrainDataNF,$.modSegConvData.MyriadDep
TrainDataNF);
```

```
OUTPUT(MyriadModelR, '~CLASS::XYZ::ML::MyriadModelR',
```

```
NAMED('Myriad_Modelo_Treinado'), OVERWRITE);
```

```
OUTPUT(COUNT(DEDUP(SORT((MyriadModelR), wi), wi)), NAMED('Numero_de_Modelos'));
```

```
//
```

```
//
```

```
// Testando o modelo
```

```
MyriadPredictedDeps :=
```

```
MyriadLearnerR.Predict(MyriadModelR,$.modSegConvData.MyriadIndTestDataNF);
```

```
OUTPUT(MyriadPredictedDeps, NAMED('Myriad_Valores_Previstos'));
```

```
//
```

```
//
```

```
// Avaliando o modelo
```

```
MyriadAssessmentR :=
```

```
ML_Core.Analysis.Reggression.Accuracy(MyriadPredictedDeps,$.modSegConvData.MyriadDepTestDa
taNF);
```

```
OUTPUT(MyriadAssessmentR, NAMED('Myriad_Avaliacao_dos_Modelos'));
```

```
OUTPUT(SORT(MyriadAssessmentR, -r2), NAMED('Ordenacao_pela_acuracia_dos_Modelos'));
```

```
//
```

```
[fn_MyriadGetPriceReg.ecl]
```

```
IMPORT $,STD,Visualizer;
```

```
IMPORT ML_Core;
```

```
IMPORT LearningTrees AS LT;
```

```

//
// Função de predição de preços de imóveis
EXPORT fn_MyriadGetPriceReg(zip,
    assess_val,
    year_acq,
    land_sq_ft,
    living_sq_ft,
    bedrooms,
    full_baths,
    half_baths,
    year_built,
    // state_code = 0) := FUNCTION
    STRING2 _state_ = 'XX') := FUNCTION
//
WUID := WORKUNIT; // obter o Id da WorkUnit
//
// Transformação dos parâmetros de entrada no formato de ML data frame
MyriadInSet := [zip, assess_val, year_acq, land_sq_ft, living_sq_ft, bedrooms, full_baths, half_baths,
year_built];
MyriadInDS := DATASET(MyriadInSet, {REAL8 MyriadInValue});
ML_Core.Types.NumericField PrepDataReg(RECORDOF(MyriadInDS) Le, INTEGER cnt) :=
TRANSFORM
// SELF.wi := state_code,
    SELF.wi := CASE(STD.Str.ToUpperCase(_state_), 'AE' => 1, 'WY' => 2, 'TN' => 3, 'MI' => 4, 'RI' =>
5, 'SC' => 6,
        'VA' => 7, 'NH' => 8, 'AP' => 9, 'NM' => 10, 'IL' => 11, 'KS' => 12,
        'ME' => 13, 'AL' => 14, 'WA' => 15, 'NY' => 16, 'MO' => 17, 'PA' => 18,
        'HI' => 19, 'GU' => 20, 'VT' => 21, 'CT' => 22, 'NC' => 23, 'OR' => 24,
        'IA' => 25, 'DE' => 26, 'VI' => 27, 'ID' => 28, 'MT' => 29, 'AR' => 30,
        'MS' => 31, 'UT' => 32, 'NE' => 33, 'IN' => 34, 'GA' => 35, 'WV' => 36,
        'NJ' => 37, 'LA' => 38, 'WI' => 39, 'AK' => 40, 'CA' => 41, 'NV' => 42,
        'FL' => 43, 'MA' => 44, 'CO' => 45, 'AZ' => 46, 'SD' => 47, 'DC' => 48,
        'KY' => 49, 'MP' => 50, 'ND' => 51, 'AS' => 52, 'TX' => 53, 'PR' => 54,
        'MD' => 55, 'OH' => 56, 'MN' => 57, 'OK' => 58, 'AA' => 59, 0),
SELF.id := 1,
SELF.number := cnt,
SELF.value := Le.MyriadInValue;
END;
MyriadNewIndDataReg := PROJECT(MyriadInDS, PrepDataReg(LEFT, COUNTER));
//
// Predição e retorno do valor do imóvel consultado
MyriadModelR :=
DATASET('~CLASS::XYZ::ML::MyriadModelR', ML_Core.Types.Layout_Model2, FLAT, PRELOAD);
MyriadLearnerR := LT.RegressionForest(10,,10,[1]);
MyriadPredictDeps := MyriadLearnerR.Predict(MyriadModelR, MyriadNewIndDataReg);
//
//
Action01 := OUTPUT(MyriadPredictDeps, {preco:=ROUND(value)}, NAMED('Preco_Imovel'));
Action02 := OUTPUT(MyriadPredictDeps, {_state_, preco:=ROUND(value)},
NAMED('Choropleth_USStates'));
Action03 := Visualizer.Choropleth.USStates('Myriad,, 'Choropleth_USStates');
Action04 := OUTPUT(WUID, NAMED('WUID'));
//
//
// RETURN OUTPUT(MyriadPredictDeps, {preco:=ROUND(value)});
RETURN SEQUENTIAL(Action01, PARALLEL(Action02, Action03), Action04);
END;
//

```

```

[fn_MyriadGetPriceReg_Web.ecl]
IMPORT $,STD,Visualizer;
IMPORT ML_Core;
IMPORT LearningTrees AS LT;
//
// Função de predição de preços de imóveis
EXPORT fn_MyriadGetPriceReg_Web(zip,
    assess_val,
    year_acq,
    land_sq_ft,
    living_sq_ft,
    bedrooms,
    full_baths,
    half_baths,
    year_built,
    // state_code = 0) := FUNCTION
    STRING2 _state_01 = 'XX',
    STRING2 _state_02 = 'YY',
    STRING2 _state_03 = 'ZZ') := FUNCTION
//
WUID := WORKUNIT; // obter o Id da WorkUnit
//
// Transformação dos parâmetros de entrada no formato de ML data frame
MyriadInSet := [zip, assess_val, year_acq, land_sq_ft, living_sq_ft, bedrooms, full_baths, half_baths,
year_built];
MyriadInDS := DATASET(MyriadInSet, {REAL8 MyriadInValue});
//
ML_Core.Types.NumericField PrepDataReg01(RECORDOF(MyriadInDS) Le, INTEGER cnt01) :=
TRANSFORM
// SELF.wi := state_code,
    SELF.wi := CASE(STD.Str.ToUpperCase(_state_01), 'AE' => 1, 'WY' => 2, 'TN' => 3, 'MI' => 4, 'RI' =>
5, 'SC' => 6,
        'VA' => 7, 'NH' => 8, 'AP' => 9, 'NM' => 10, 'IL' => 11, 'KS' => 12,
        'ME' => 13, 'AL' => 14, 'WA' => 15, 'NY' => 16, 'MO' => 17, 'PA' => 18,
        'HI' => 19, 'GU' => 20, 'VT' => 21, 'CT' => 22, 'NC' => 23, 'OR' => 24,
        'IA' => 25, 'DE' => 26, 'VI' => 27, 'ID' => 28, 'MT' => 29, 'AR' => 30,
        'MS' => 31, 'UT' => 32, 'NE' => 33, 'IN' => 34, 'GA' => 35, 'WV' => 36,
        'NJ' => 37, 'LA' => 38, 'WI' => 39, 'AK' => 40, 'CA' => 41, 'NV' => 42,
        'FL' => 43, 'MA' => 44, 'CO' => 45, 'AZ' => 46, 'SD' => 47, 'DC' => 48,
        'KY' => 49, 'MP' => 50, 'ND' => 51, 'AS' => 52, 'TX' => 53, 'PR' => 54,
        'MD' => 55, 'OH' => 56, 'MN' => 57, 'OK' => 58, 'AA' => 59, 0),
SELF.id := 1,
SELF.number := cnt01,
SELF.value := Le.MyriadInValue;
END;
MyriadNewIndDataReg01 := PROJECT(MyriadInDS, PrepDataReg01(LEFT, COUNTER));
//
ML_Core.Types.NumericField PrepDataReg02(RECORDOF(MyriadInDS) Le, INTEGER cnt02) :=
TRANSFORM
// SELF.wi := state_code,
    SELF.wi := CASE(STD.Str.ToUpperCase(_state_02), 'AE' => 1, 'WY' => 2, 'TN' => 3, 'MI' => 4, 'RI' =>
5, 'SC' => 6,
        'VA' => 7, 'NH' => 8, 'AP' => 9, 'NM' => 10, 'IL' => 11, 'KS' => 12,
        'ME' => 13, 'AL' => 14, 'WA' => 15, 'NY' => 16, 'MO' => 17, 'PA' => 18,
        'HI' => 19, 'GU' => 20, 'VT' => 21, 'CT' => 22, 'NC' => 23, 'OR' => 24,
        'IA' => 25, 'DE' => 26, 'VI' => 27, 'ID' => 28, 'MT' => 29, 'AR' => 30,
        'MS' => 31, 'UT' => 32, 'NE' => 33, 'IN' => 34, 'GA' => 35, 'WV' => 36,

```



```

'NJ' => 37,'LA' => 38,'WI' => 39,'AK' => 40,'CA' => 41,'NV' => 42,
'FL' => 43,'MA' => 44,'CO' => 45,'AZ' => 46,'SD' => 47,'DC' => 48,
'KY' => 49,'MP' => 50,'ND' => 51,'AS' => 52,'TX' => 53,'PR' => 54,
'MD' => 55,'OH' => 56,'MN' => 57,'OK' => 58,'AA' => 59,0),

SELF.id := 1,
SELF.number := cnt02,
SELF.value := Le.MyriadInValue;
END;
MyriadNewIndDataReg02 := PROJECT(MyriadInDS, PrepDataReg02(LEFT, COUNTER));
//
ML_Core.Types.NumericField PrepDataReg03(RECORDOF(MyriadInDS) Le, INTEGER cnt03) :=
TRANSFORM
// SELF.wi := state_code,
SELF.wi := CASE(STD.Str.ToUpperCase(_state_03), 'AE' => 1,'WY' => 2,'TN' => 3,'MI' => 4,'RI' =>
5,'SC' => 6,

'VA' => 7,'NH' => 8,'AP' => 9,'NM' => 10,'IL' => 11,'KS' => 12,
'ME' => 13,'AL' => 14,'WA' => 15,'NY' => 16,'MO' => 17,'PA' => 18,
'HI' => 19,'GU' => 20,'VT' => 21,'CT' => 22,'NC' => 23,'OR' => 24,
'IA' => 25,'DE' => 26,'VI' => 27,'ID' => 28,'MT' => 29,'AR' => 30,
'MS' => 31,'UT' => 32,'NE' => 33,'IN' => 34,'GA' => 35,'WV' => 36,
'NJ' => 37,'LA' => 38,'WI' => 39,'AK' => 40,'CA' => 41,'NV' => 42,
'FL' => 43,'MA' => 44,'CO' => 45,'AZ' => 46,'SD' => 47,'DC' => 48,
'KY' => 49,'MP' => 50,'ND' => 51,'AS' => 52,'TX' => 53,'PR' => 54,
'MD' => 55,'OH' => 56,'MN' => 57,'OK' => 58,'AA' => 59,0),

SELF.id := 1,
SELF.number := cnt03,
SELF.value := Le.MyriadInValue;
END;
MyriadNewIndDataReg03 := PROJECT(MyriadInDS, PrepDataReg03(LEFT, COUNTER));
//
//
// Predição e retorno do valor do imóvel consultado - State 01
MyriadModel01 :=
DATASET('~CLASS::XYZ::ML::MyriadModelR',ML_Core.Types.Layout_Model2,FLAT,PRELOAD);
MyriadLearner01 := LT.RegressionForest(10,,10,[1]);
MyriadPredictDeps01 := MyriadLearner01.Predict(MyriadModel01, MyriadNewIndDataReg01);
//
// Predição e retorno do valor do imóvel consultado - State 02
MyriadModel02 :=
DATASET('~CLASS::XYZ::ML::MyriadModelR',ML_Core.Types.Layout_Model2,FLAT,PRELOAD);
MyriadLearner02 := LT.RegressionForest(10,,10,[1]);
MyriadPredictDeps02 := MyriadLearner02.Predict(MyriadModel02, MyriadNewIndDataReg02);
//
// Predição e retorno do valor do imóvel consultado - State 03
MyriadModel03 :=
DATASET('~CLASS::XYZ::ML::MyriadModelR',ML_Core.Types.Layout_Model2,FLAT,PRELOAD);
MyriadLearner03 := LT.RegressionForest(10,,10,[1]);
MyriadPredictDeps03 := MyriadLearner03.Predict(MyriadModel03, MyriadNewIndDataReg03);
//
//
Action01a := OUTPUT(MyriadPredictDeps01,{preco:=ROUND(value)}, NAMED('Preco_Imovel_01'));
Action01b := OUTPUT(MyriadPredictDeps02,{preco:=ROUND(value)}, NAMED('Preco_Imovel_02'));
Action01c := OUTPUT(MyriadPredictDeps03,{preco:=ROUND(value)}, NAMED('Preco_Imovel_03'));
//
Action02a := OUTPUT(MyriadPredictDeps01,{states:=_state_01, preco:=ROUND(value)},
NAMED('Choropleth_USStates'),EXTEND);
Action02b := OUTPUT(MyriadPredictDeps02,{states:=_state_02, preco:=ROUND(value)},
NAMED('Choropleth_USStates'),EXTEND);

```

```

Action02c := OUTPUT(MyriadPredictDeps03,{states:=_state_03, preco:=ROUND(value)},
NAMED('Choropleth_USStates'),EXTEND);
//
Action03 := Visualizer.Choropleth.USStates('Myriad', 'Choropleth_USStates');
Action04 := OUTPUT(WUID, Named('WUID'));
//
//
// RETURN OUTPUT(MyriadPredictDeps, {preco:=ROUND(value)});
RETURN SEQUENTIAL(Action01a,Action01b,Action01c,
Action02a,Action02b,Action02c,
Action03,
Action04);
END;
//

```

[BWR\_TestQueriesReg.ecl]

```

IMPORT $;
//
// Teste da Função
// Zip | Assess_val | Year_acq | Land_sq_ft | Living_sq_ft | Bedrooms | Full_baths | Half_baths | Year_built
| States 01/02/03
//
// $.fn_MyriadGetPriceReg(95451,118720,2011,14774,1437,3,2,1,1968,'CA');
$.fn_MyriadGetPriceReg_Web(95451,118720,2011,14774,1437,3,2,1,1968,'CA','NY','FL');
//
/*
_state_01: CA
_state_02: NY
_state_03: FL
assess_val: 118720
bedrooms: 3
full_baths: 2
half_baths: 1
land_sq_ft: 14774
living_sq_ft:1437
year_acq: 2011
year_built: 1968
zip: 95451
*/
//

```

## 1.6. Referências

Introduction to HPCC Systems Open Source Big Data Platform. Disponível em:

[https://cdn.hpccsystems.com/whitepapers/wp\\_introduction\\_HPCC.pdf](https://cdn.hpccsystems.com/whitepapers/wp_introduction_HPCC.pdf)

Acesso em: 25 set. 2024.

Machine Learning Demystified. Disponível em:

<https://hpccsystems.com/resources/machine-learning-demystified/>

Acesso em: 25 set. 2024.

HPCC Systems Machine Learning Library. Disponível em:

<https://hpccsystems.com/download/free-modules/hpcc-systems-machine-learning-library/>

Acesso em: 25 set. 2024.

Using HPCC Systems Machine Learning. Disponível em:

<https://hpccsystems.com/resources/using-hpcc-systems-machine-learning/>

Acesso em: 25 set. 2024.

Introducing the new, improved HPCC Systems Machine Learning Library | HPCC Systems. Disponível em:

<https://hpccsystems.com/resources/introducing-the-new-improved-hpcc-systems-machine-learning-library/>

Acesso em: 25 set. 2024.

ECL-ML Machine Learning Module. Disponível em:

<https://hpccsystems.com/resources/ecl-ml-machine-learning-module/>

Acesso em: 25 set. 2024.

ML\_Core Documentation. Disponível em:

[https://cdn.hpccsystems.com/pdf/ml/ML\\_Core.pdf](https://cdn.hpccsystems.com/pdf/ml/ML_Core.pdf)

Acesso em: 25 set. 2024.

Source code: HPCC Systems ML\_Core repository on GitHub. Disponível em:

[https://github.com/hpcc-systems/ML\\_Core](https://github.com/hpcc-systems/ML_Core)

Acesso em: 25 set. 2024.

Introduction to using PBblas on HPCC Systems. Disponível em:

<https://hpccsystems.com/resources/introduction-to-using-pbblas-on-hpcc-systems/>

Acesso em: 25 set. 2024.

Documentation: PBblas Documentation. Disponível em:

<https://cdn.hpccsystems.com/pdf/ml/PBblas.pdf>

Acesso em: 25 set. 2024.

Source code: HPCC Systems PBblas repository on GitHub. Disponível em:

<https://github.com/hpcc-systems/PBblas>

Acesso em: 25 set. 2024.

Learning Trees — A guide to Decision Tree based Machine Learning. Disponível em:

<https://hpccsystems.com/resources/learning-trees-a-guide-to-decision-tree-based-machine-learning/>

Acesso em: 25 set. 2024.

Documentation: LearningTrees Documentation. Disponível em:

<https://cdn.hpccsystems.com/pdf/ml/LearningTrees.pdf>

Acesso em: 25 set. 2024.

Source code: LearningTrees repository on GitHub. Disponível em:

<https://github.com/hpcc-systems/LearningTrees>

Acesso em: 25 set. 2024.

Understanding the Myriad Interface feature of HPCC Systems Machine Learning | HPCC Systems. Disponível em:

<https://hpccsystems.com/resources/understanding-the-myriad-interface-feature-of-hpcc-systems-machine-learning/>

Acesso em: 25 set. 2024.