

Modelo de Previsão de Vendas Multiproduto

Autor

Ricardo Castro (rc@rcpt)

Orientadores

João Alves (Metyis Porto)

Nelson Pires (ISLA)

Firmino Silva (ISLA)

INTRODUÇÃO

O projeto apresentado, foi realizado em colaboração com a Metyis e destina-se à "The Greenery" uma cooperativa de agricultores holandeses dedicada ao comércio e venda de frutas e vegetais, com origem na Holanda e não só, e destinados à distribuição nacional e internacional, essencialmente para supermercados e grandes superfícies. Mantendo o foco na rentabilidade do negócio e resultados financeiros crescentes, a "The Greenery", pretende ajustar a oferta à procura reduzindo o desperdício e evitando rotura de stocks.

A operar num mercado cada vez mais exigente e concorrencial, e sentindo a necessidade de otimizar os seus processos pela inovação e digitalização, a "The Greenery", forneceu via Metyis, os dados da sua atividade, de modo a que aplicando técnicas e modelos de data science e business analytics, lhes fosse desenvolvida uma ferramenta de utilização simples, capaz de fazer previsões de vendas dos vários produtos com que trabalha.

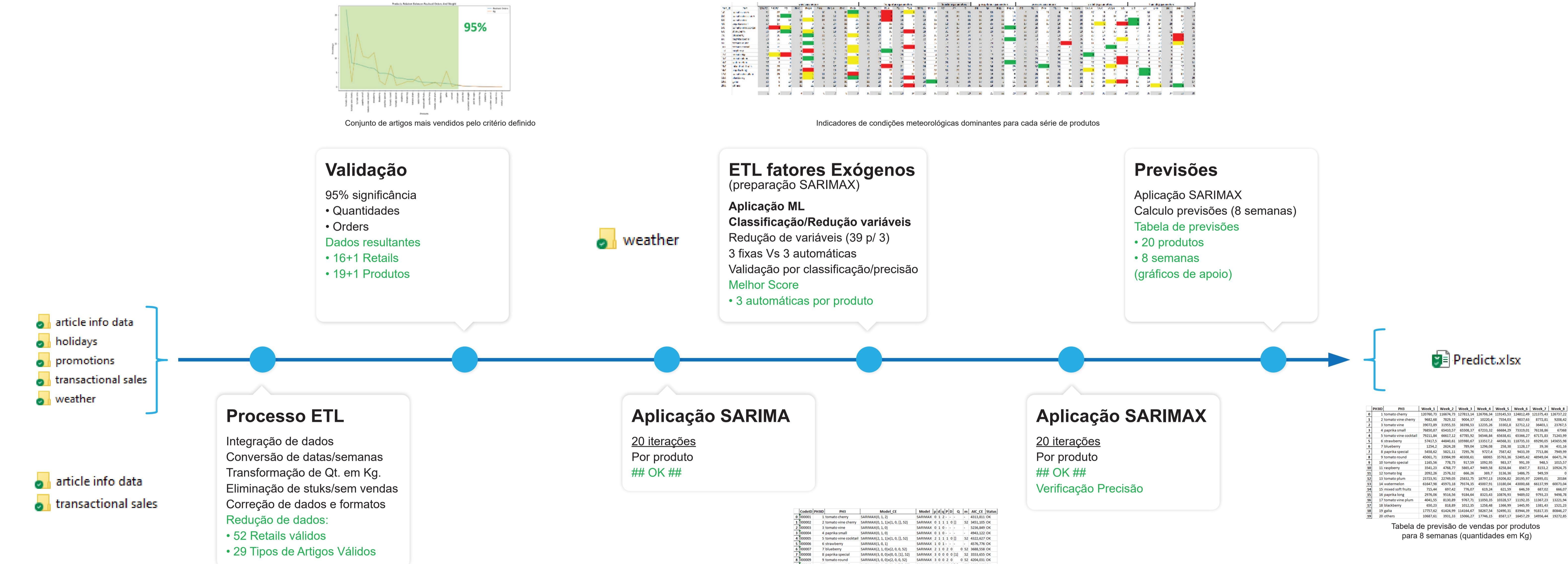
Palavras-chave:

CRISP-DM, SARIMAX, MACHINE LEARNING; PYTHON; THE GREENERY.



IMPLEMENTAÇÃO DO PROJECTO

A implementação do projeto seguiu várias etapas conforme workflow que se apresenta, passando no essencial pela verificação da qualidade dos dados, ingestão no processo, e tratamento cuidado e criterioso, não só para aproveitar ao máximo a informação recebida, mas também garantir que a mesma apresenta formatos, coerentes, legíveis e integráveis com todo o processo de modelação. A modelação teve por base a utilização de modelos autorregressivos integrando variáveis móveis (SARIMAX), aplicado a séries temporais, integrando a influência das condições climáticas como variáveis exógenas ao processo de previsão. Os modelos criados seguiram técnicas matemáticas, estatísticas e de Machine Learning aprendidas.



Etapa 1 - Processo ETL - Os dados recolhidos foram lidos, tratados e integrados, com uniformização de grandezas, eliminação de redundâncias ou duplicados, e tratamento de valores nulos ou elegíveis. Foram recebidos dados desde janeiro 2018 até início de maio de 2022, num total de 315.474 registos de vendas, 59.244 produtos, traduzindo-se depois do ETL em 29 produtos únicos e válidos, e 52 clientes válidos e com operações registadas.

Etapa 2 - Validação - Foi definido fazer uma análise detalhada para os produtos que percentualmente tivessem maior representatividade e cuja soma percentual representasse pelo menos 95% das vendas globais, em quantidades vendidas (Kg), mas também em número de ordens de compra recebidas. Resultaram 19 produtos significativos a tratar individualmente, sendo os restantes tratados como um conjunto.

CONCLUSÃO

A solução desenvolvida para a "The Greenery" em contexto real, consiste numa aplicação em Python, capaz de prever vendas de vários produtos num horizonte temporal de 8 semanas recorrendo aos dados históricos de vendas ocorridas. O resultado desta aplicação é uma tabela de fácil interpretação e utilização, conforme foi solicitado.

O trabalho apresentado, e partindo de uma necessidade concreta do mercado, é o resultado de um processo iterativo de aprendizagem e aplicação de conhecimento, traduzido numa solução possível de fácil aplicação na "The Greenery".

BIBLIOGRAFIA & SITES CONSULTADOS

 1) <https://www.datascience-pm.com/crisp-dm> - Informação relativa à metodologia CRISP-DM apresentação e implementação;

 2) <https://www.thegreenery.com/en> - Informação geral e operacional da empresa em estudo a "The Greenery";

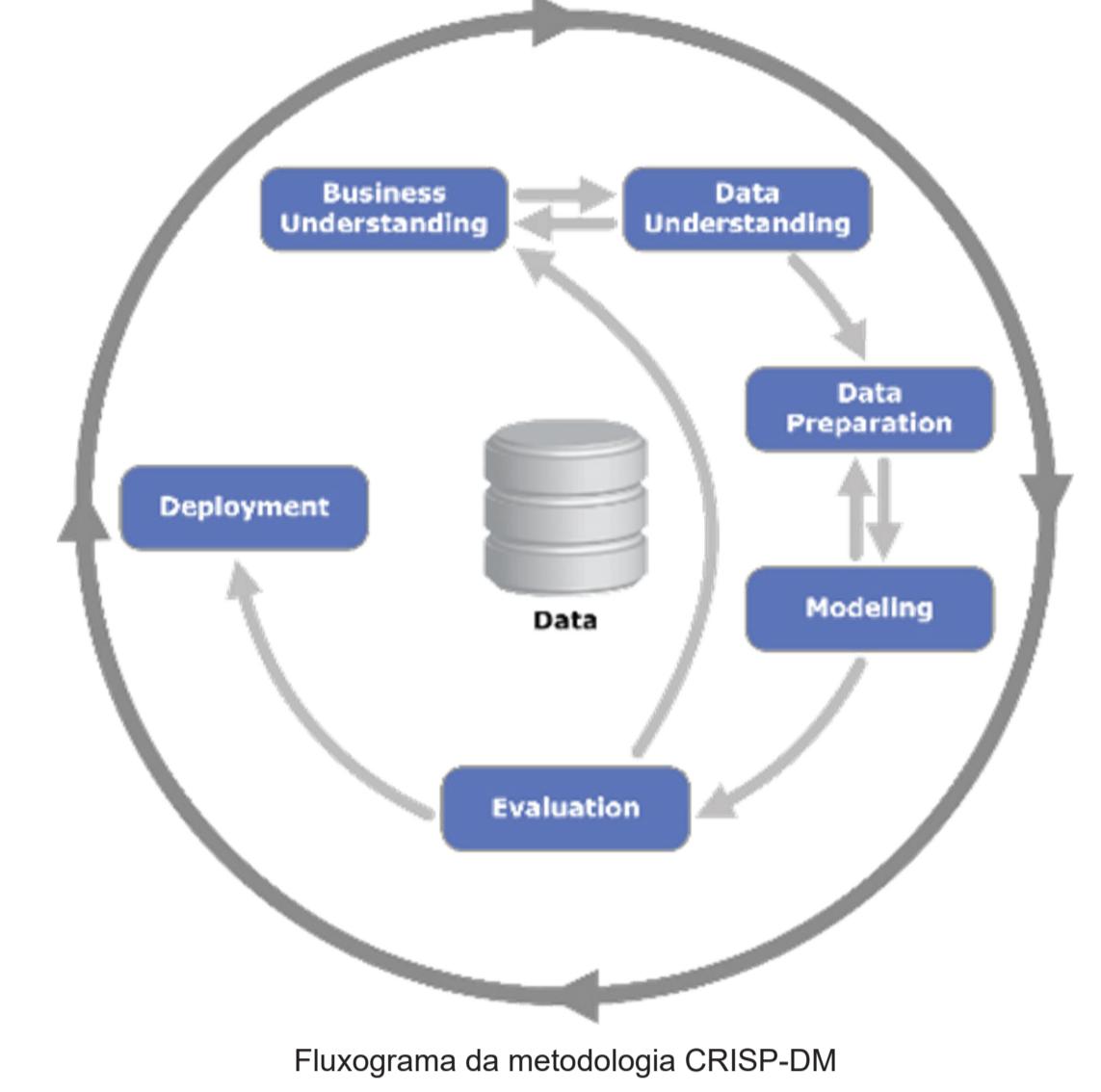
 3) <https://stackoverflow.com> - Exemplos e rotinas em python utilizadas no desenvolvimento do processo e programa realizado;

 4) https://scikit-learn.org/stable/user_guide.html - Exemplos e algoritmos de Machine Learning utilizados no desenvolvimento do processo e programa realizado;

 5) <https://www.statsmodels.org/dev/examples> - Exemplos e aplicações de modelos de previsão aplicados a séries temporais;

OBJETIVO

Desenvolver uma ferramenta de previsão baseada em linguagem Phyton, aplicando a metodologia CRISP-DM e utilizando técnicas de data Science, com modelos matemáticos, estatísticos e de Machine Learning, como algoritmos de classificação, redução de variáveis e modelos autorregressivos integrando variáveis móveis (SARIMAX), com o objetivo último de produzir previsões de venda para os produtos considerados mais relevantes, numa janela temporal de 8 semanas, considerando a influência de fatores exógenos, no caso as condições meteorológicas.



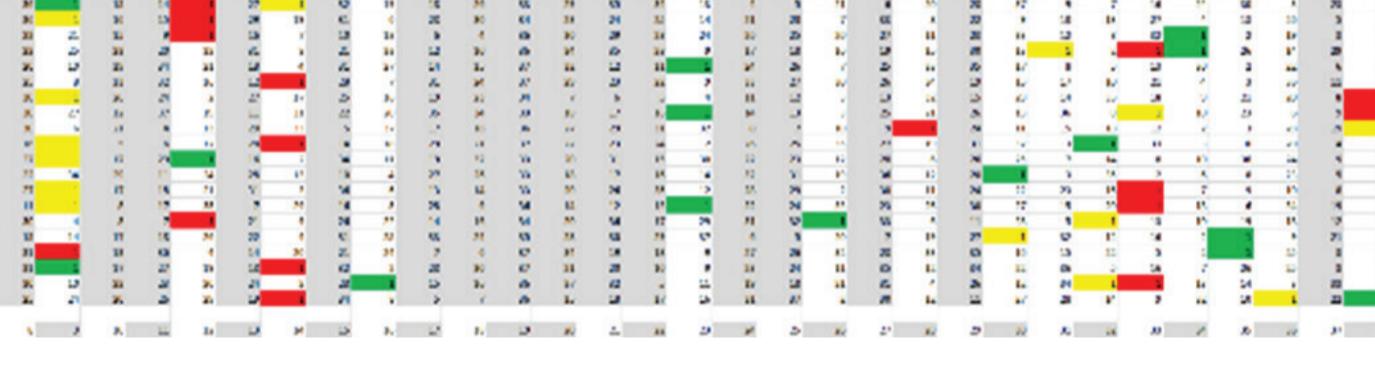
A metodologia seguida neste projeto foi a metodologia CRISP-DM - CRoss Industry Standard Process for Data Mining, nas 6 fases conhecidas:

- Business understanding;
- Data understanding;
- Data Preparation;
- Modeling;
- Evaluation;
- Deployment.

Etapa 3 - Aplicação SARIMA - Foi utilizada uma função automática em Python, para determinação dos parâmetros adequados a cada modelo SARIMA, aplicado a cada série temporal dos produtos definidos. Os modelos foram treinados e testados fornecendo resultados consistentes que validam o processo desenvolvido, pese embora os erros verificados rondem os 26%.

Etapa 4 - ETL Fatores Exógenos - Como fatores exógenos selecionamos as condições meteorológicas compostas por 39 parâmetros diferentes. Numa primeira abordagem, intuitivamente, foram selecionados de forma fixa a Temperatura diária, a Duração da luz solar e a Humidade do ar. Simultaneamente e aplicado um algoritmo de redução de variáveis, RFE – Recursive Feature Elimination, tentámos pelo método RFE automático encontrar os 3 parâmetros principais que mais influenciam a venda de cada produto. Verificou-se que os 3 parâmetros fixos previamente selecionados não eram os dominantes, assim como se verificou, que os 3 principais parâmetros obtidos pelo processo RFE variavam de produto para produto.

Indicadores de condições meteorológicas dominantes para cada série de produtos



Conjunto de artigos mais vendidos pelo critério definido

95%

95%

Indicadores de condições meteorológicas dominantes para cada série de produtos

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%

95%