

Modelo de Predição de Cancelamento de Reservas em Hotéis

SCC0230 Inteligência Artificial



Descrição da Problemática:

- Novos métodos de reserva de estadias em hotéis facilitam a experiência do usuário, porém podem trazer novos problemas para os hotéis.
- A comodidade pode fazer com que o cliente trate com menor importância a reserva.
- Reservas via site de domínio próprio ou aplicativos e/ou portais de reserva de terceiros distanciam o cliente da empresa.
- Maior comprometimento nas reservas mais “analógicas”, como via telefone ou até presencialmente.
- Cabe aos hotéis buscarem soluções para amenizar as chances do cliente não honrar com sua reserva, ou ao menos se preparar para isso.
- Busca-se mitigar o ônus dos hotéis utilizando modelos preditivos inteligentes.



Os dados:

- Booking ID
- no_of_adults
- no_of_childrens
- no_of_weekend_nights
- no_of_week_nights
- type_of_meal_plan
- required_car_parking_space
- room_type_reserved
- lead_time
- arrival_year
- arrival_month
- arrival_date
- market_segment_type
- repeated_guest
- no_of_previous_cancellations
- no_of_previous_bookings_not_canceled
- avg_price_per_room
- no_of_special_requests
- **booking_status**

Pré-processamento dos dados:

O dataset já se encontrava limpo, sem dados duplicados ou nulos, dispensando a necessidade de pré-processamento dos dados. Houve apenas uma preparação - distinta para cada modelo - antes de cada modelagem.

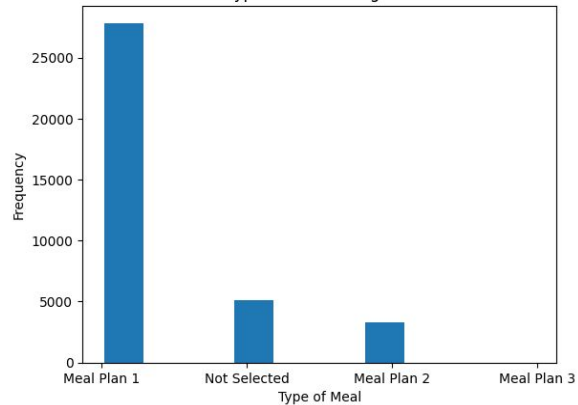


Análise dos dados:

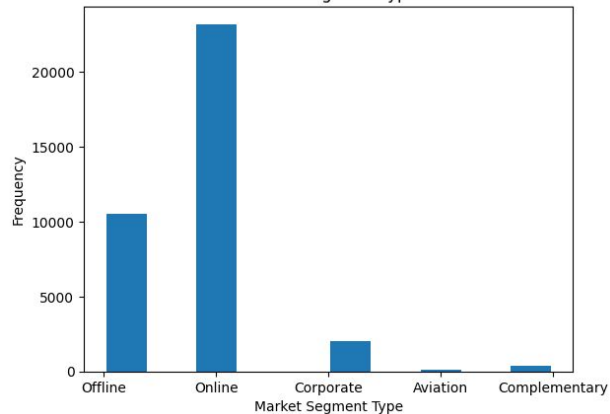
Histograma das variáveis categóricas

A investigação dessas variáveis visa fornecer insights sobre preferências, comportamentos e características distintas dos clientes. Os resultados dessa análise estão apresentados no slide seguinte.

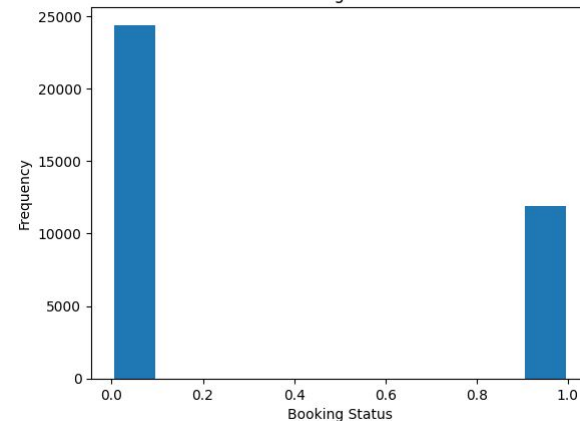
Type of Meal histogram



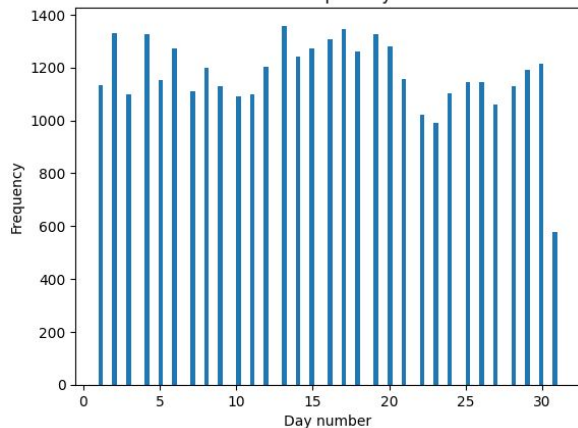
Market Segment Type



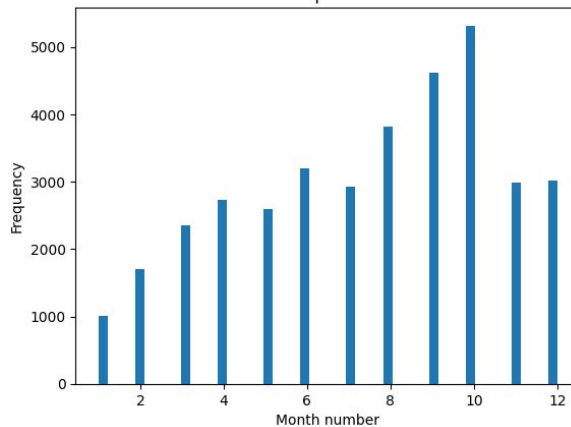
Booking Status



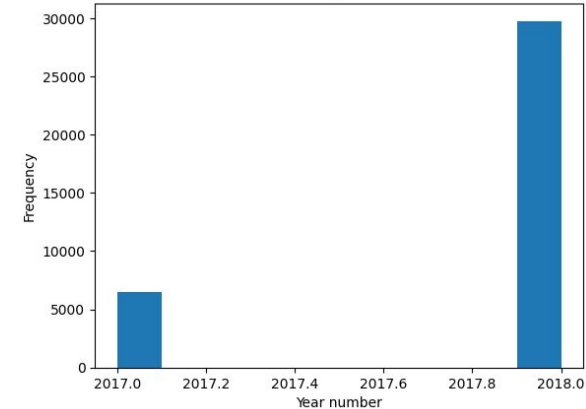
Arrival per Day



Arrival per Month



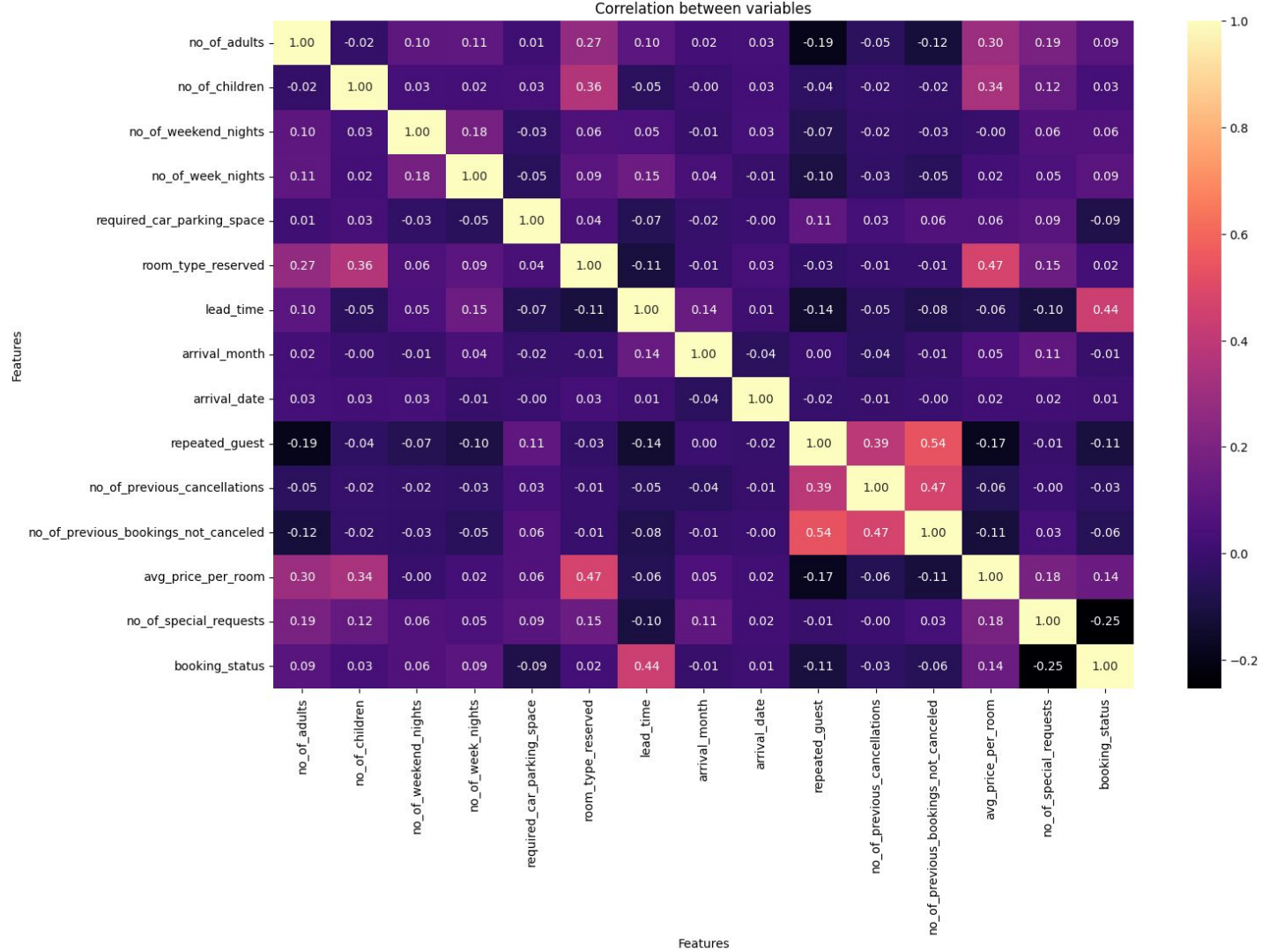
Arrival per Year



Análise dos dados:

Correlação entre as variáveis

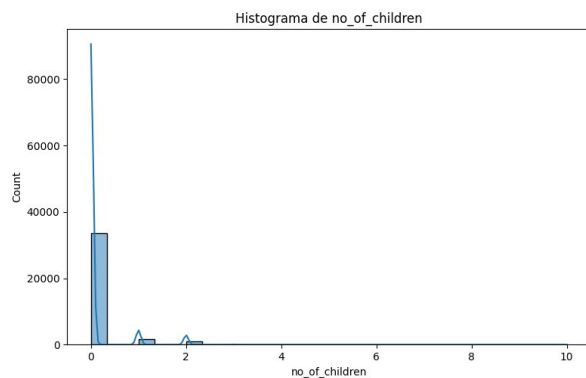
Conduzimos uma análise de correlação com o objetivo de explorar as relações entre as variáveis em nosso conjunto de dados de reservas de hotel, buscando entender as interconexões e padrões subjacentes. Os resultados estão detalhados no slide subsequente.



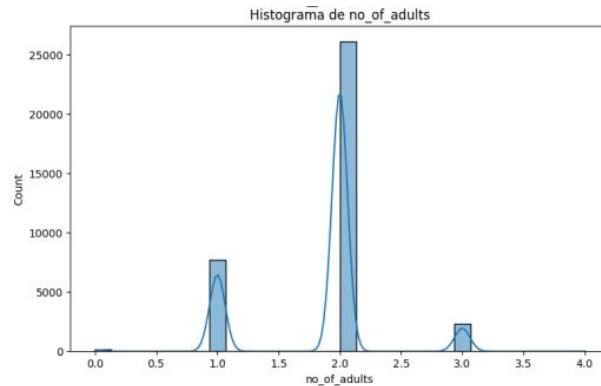
Análise dos dados:

Distribuição das variáveis contínuas

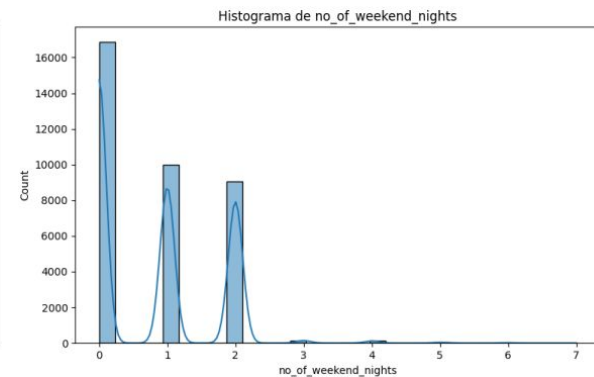
Fizemos uma análise exploratória com o intuito de obter uma visão abrangente das características de algumas das variáveis contínuas em nosso conjunto de dados de reservas de hotel, como apresentado no slide a seguir.



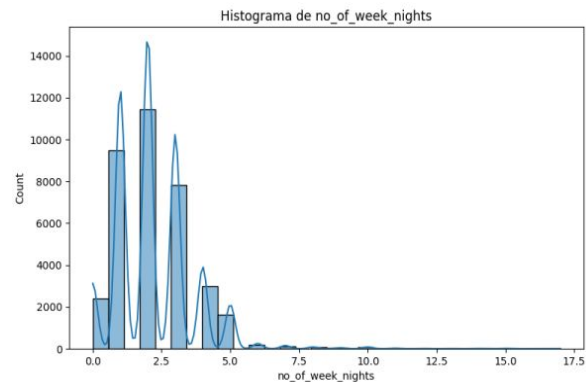
- Examina a distribuição da presença de crianças nas reservas.
- Ajuda a observar a frequência de reservas com ou sem crianças.



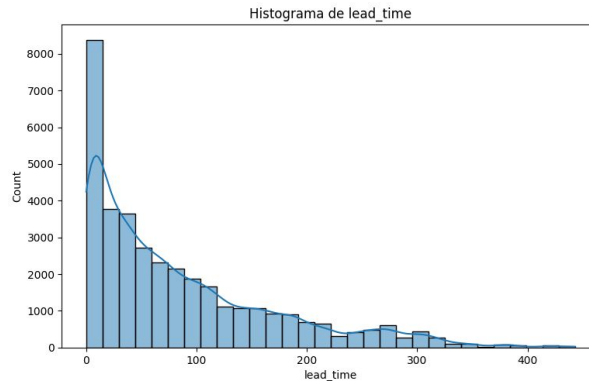
- O histograma revela a distribuição do número de adultos nas reservas.
- Permite identificar se a maioria das reservas é para casais, grupos ou viajantes individuais.



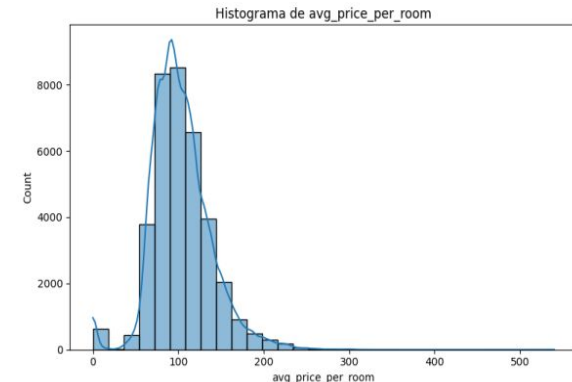
- Mostra como as reservas estão distribuídas em relação às noites de fim de semana.
- Pode indicar preferências por estadias durante a semana ou nos fins de semana.



- Analisa a distribuição das reservas em relação às noites de semana.
- Pode revelar se a maioria das reservas é para viagens de negócios ou lazer durante a semana.



- Exibe a distribuição do tempo entre a data da reserva e a data de chegada.
- Ajuda a entender padrões de reservas de última hora ou planejadas com antecedência.



- Mostra como os preços dos quartos estão distribuídos.
- Pode fornecer insights sobre a faixa de preço preferida pelos clientes.

Avaliação de modelos

Utilizamos algumas métricas para analisar qual o melhor modelo dentre os testados, que estão listadas no slide a seguir.

Métrica	O que Mede	Como é Calculada	Uso
Acurácia	Proporção de acertos do modelo.	$(\text{Verdadeiros Positivos} + \text{Verdadeiros Negativos}) / \text{Total de Previsões}$.	Avaliação geral de desempenho.
Pontuação F1	Equilíbrio entre precisão e recall.	$2 * (\text{Precisão} * \text{Recall}) / (\text{Precisão} + \text{Recall})$.	Útil em conjuntos de dados desbalanceados.
Precisão	Exatidão das previsões positivas.	$\text{Verdadeiros Positivos} / (\text{Verdadeiros Positivos} + \text{Falsos Positivos})$.	Importante quando os falsos positivos são custosos.
Recall	Capacidade de identificar casos positivos.	$\text{Verdadeiros Positivos} / (\text{Verdadeiros Positivos} + \text{Falsos Negativos})$.	Crucial quando os falsos negativos são custosos.

Resultados

Testamos 3 modelos diferentes, de 3 paradigmas diferentes, e avaliamos o desempenho de cada um, com o objetivo de averiguarmos qual modelo se adequa melhor aos dados apresentados.

1. **Random Forest**
2. **KNN**
3. **Naive Bayes**

Resultados:

1) Random Forest

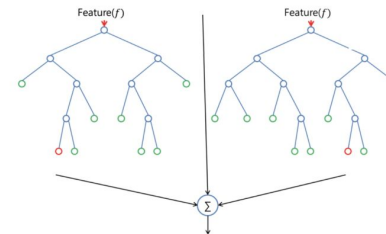
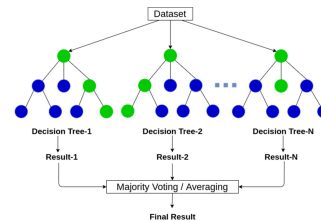
Paradigma: Aprendizado Simbólico + Ensemble

Ao combinar diversas árvores, cada uma treinada em subconjuntos diferentes dos dados, o modelo oferece estabilidade, reduzindo overfitting e proporcionando uma visão abrangente das relações nas variáveis.

Desempenho do Modelo:

	Model	Accuracy Score	F1 score	Precision	Recall
0	Random forest	0.904179	0.849263	0.893757	0.808989

Random Forest



Resultados:

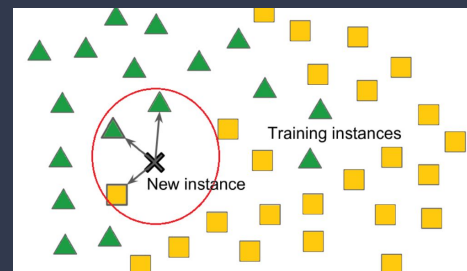
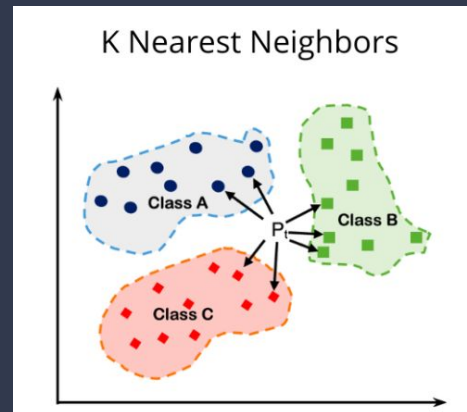
2) KNN

Paradigma: Aprendizado baseado em instâncias

O KNN, é uma técnica de aprendizado de máquina que se baseia na proximidade entre os pontos de dados. Ao classificar ou prever com base nos vizinhos mais próximos, o KNN destaca-se por sua simplicidade conceitual e eficácia em identificar padrões em conjuntos de dados.

Desempenho do Modelo:

	Model	Accuracy	Score	F1 score	Precision	Recall
0	KNN	0.876172	0.808329	0.835863	0.782551	



Resultados:

3) Naive Bayes: Categorical Naive Bayes

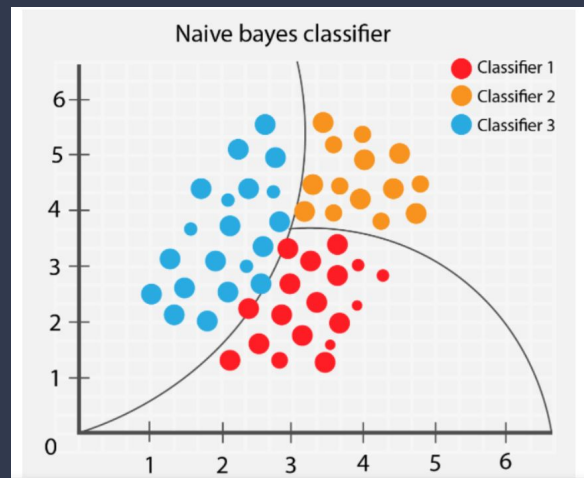
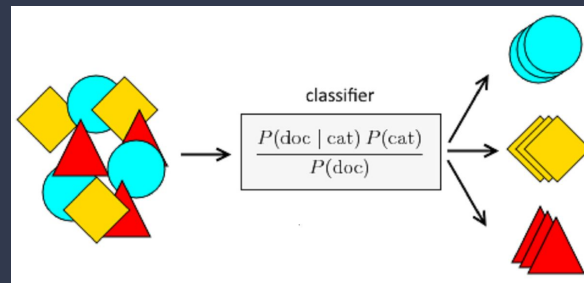
Paradigma: Aprendizado estatístico

Baseado no teorema de Bayes, o algoritmo calcula a probabilidade condicional de uma classe dado um conjunto de características. É especialmente útil quando as variáveis são independentes (hipótese ingênua).

O Naive Bayes categórico foi escolhido a partir da análise exploratória, onde foi percebida a grande quantidade de variáveis categóricas presentes nos dados.

Desempenho do Modelo:

	Model	Accuracy Score	F1 score	Precision	Recall
0	Categorical_nb	0.801522	0.671293	0.729365	0.621786

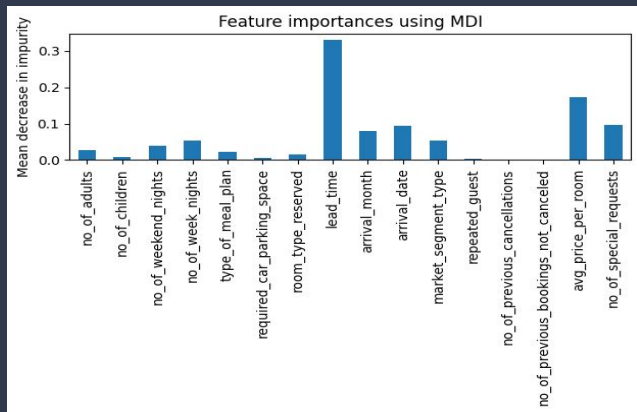


Modelo escolhido:

Random Forest

Interpretabilidade das Features:

As barras no gráfico representam a importância de cada feature na previsão do modelo, calculada usando o critério Mean Decrease in Impurity (MDI).



Esta decisão foi motivada por sua boa performance, expressa por uma Accuracy Score de 90.4%, indicando sua eficácia em prever com precisão se um hóspede irá cancelar ou não.

O Random Forest também se destacou com um F1 Score de 84.9%, uma métrica essencial quando lidamos com desequilíbrio nos dados, como é comum em previsões de cancelamentos. Esta pontuação leva em consideração tanto a precisão quanto o recall, equilibrando a preocupação com falsos positivos e falsos negativos.

Além de sua performance sólida, o Random Forest é conhecido por lidar bem com diferentes tipos de dados, oferecendo robustez contra o sobreajuste e se mostrando eficaz em tarefas de classificação. Dessa forma, acreditamos que o Random Forest é a escolha mais adequada para atender às nossas necessidades específicas de previsão de cancelamentos de reservas.

Aplicações práticas dos resultados:

- Munida das análises aqui mostradas as empresas hoteleiras podem tomar diversas ações:
 - Dada a predição da reserva em tempo real, um prazo, valor e multas adaptativas podem ser implementadas;
 - Avaliando as estatísticas de cancelamento das reservas, novos métodos de contato com o cliente podem ser pensados, retomando a proximidade de antigamente.
 - Pode-se explorar também diferentes fontes de reserva, uma vez que uma empresa pode ter mais de um canal de atendimento.

Fonte dos dados

<https://www.kaggle.com/datasets/ahsan81/hotel-reservations-classification-dataset/data>

Integrantes:

Ana Vitória Gouvea de Oliveira Freitas

Eduardo Vinícius Barbosa Rossi

Matheus Luis Oliveira da Silva

Pedro Augusto Ribeiro Gomes

Sofhia de Souza Gonçalves

Thiago Henrique dos Santos Cardoso