



## Universidade do Minho

Licenciatura em Engenharia Informática

Aprendizagem e Decisão Inteligentes

3º Ano, 2º Semestre

Ano letivo 2024/2025

Guião prático nº 5

Fevereiro, 2025

### Tema

Exploração avançada e preparação de dados com KNIME

### Objetivos de aprendizagem

Com a realização desta ficha prática pretende-se que os estudantes:

- Explore a plataforma de análise de dados KNIME;
- Experimentem tarefas de exploração avançada de dados;
- Realizem tarefas de preparação de dados;

### Enunciado

Os dados respeitam a voos e respetivos atrasos nas chegadas (*target*) de diversas companhias aéreas.

Realizar as tarefas seguintes:

1. Carregar o *dataset* «flights» com os dados dos voos;
2. Aplicar nodos para exploração de dados que respondam à resolução das seguintes situações:
  - a. Em que colunas os valores são sempre iguais?
  - b. Que colunas têm variância elevada?
  - c. Que colunas representam o mesmo conhecimento?
  - d. Existe alguma relação relevante entre as *features*?
3. Aplicar nodos de exploração e tratamento de dados para:
  - a. Remover colunas que possam ser irrelevantes ou que não traduzam conhecimento novo;
  - b. Verificar a existência de entradas duplicadas;
  - c. Tratar valores em falta (*missing values*);
  - d. Remover *outliers*;
  - e. Realizar *encoding* das *features* (transformar dados em colunas);
  - f. Fazer a normalização de valores;
4. Criar modelos de previsão:
  - a. Treinar um modelo de regressão linear para prever o atraso na chegada;
  - b. Analisar o resultado da previsão.



## Descrição do *dataset* FLIGHTS

ATRIBUTO	DESCRIÇÃO
<b>id</b>	Identificador do voo
<b>year</b>	Ano em que o voo se realizou
<b>month</b>	Mês em que o voo se realizou (1 = janeiro; 12 = dezembro).
<b>day</b>	Dia em que o voo se realizou
<b>dep_time</b>	Hora de partida real (formato “hhmm”)
<b>sched_dep_time</b>	Hora de partida agendada (formato “hhmm”)
<b>dep_delay</b>	Atraso na partida do voo (em minutos) - diferença entre horário real e agendado - valores positivos são atraso e negativos são antecipação
<b>arr_time</b>	Hora de chegada real (formato “hhmm”)
<b>sched_arr_time</b>	Hora de chegada agendada (formato “hhmm”)
<b>arr_delay</b>	Atraso na chegada do voo (em minutos) - diferença entre horário real e agendado - valores positivos são atraso e negativos são antecipação
<b>carrier</b>	Código da companhia aérea (duas letras)
<b>flight</b>	Número do voo
<b>tailnum</b>	Identificador do avião
<b>origin</b>	Código do aeroporto de partida (três letras)
<b>dest</b>	Código do aeroporto de destino (três letras)
<b>air_time</b>	Duração do voo (em minutos)
<b>distance</b>	Distância entre os aeroportos de origem e destino (em milhas)
<b>hour</b>	Hora da partida agendada
<b>minute</b>	Minutos da partida agendada
<b>time_hour</b>	Hora agendada de partida do voo (“aaaa-mm-dd hh:00:00”)
<b>name</b>	Nome completo da companhia aérea responsável pelo voo

Mais detalhes sobre estes dados podem ser encontrados neste link: [kaggle.com/datasets/flights](https://www.kaggle.com/datasets/flights)