

# Information Retrieval and Text Mining

Nuno Escudeiro ([nfe@isep.ipp.pt](mailto:nfe@isep.ipp.pt))

Ricardo Almeida ([ral@isep.ipp.pt](mailto:ral@isep.ipp.pt))

# Planning

# Planning

S	T (2 HORAS/SEMANA)	TP (1 HORA/SEMANA)	PL (2 HORAS/SEMANA)
1 24 - fev	<ul style="list-style-type: none"> <li>- Introdução à Recuperação de Informação (IR), Filtragem de Informação (IF), Extração de Informação (IE) e Mineração de Texto</li> <li>- Processo geral da recuperação de informação.</li> <li>- Indexação e consulta/pesquisa de informação não estruturada</li> </ul>	Componentes de um sistema de recuperação de informação	Lançamento Projeto e Seminário. Criação dos grupos. Seleção de temas. Preparação do artigo sobre os temas selecionados
2 3 - mar	<ul style="list-style-type: none"> <li>- Avaliação em recuperação e extração de informação. Precision, recall, F1-score</li> <li>- ROC curve, AUC, Precision/Recall curve</li> </ul>	Avaliação em recuperação e extração de informação	Preparação do artigo
3 10 - mar	<ul style="list-style-type: none"> <li>- Boolean Retrieval Model: Basics of Boolean algebra and operations in information retrieval</li> <li>- Conjunctive and Disjunctive Queries: Understanding AND, OR, and NOT operators</li> </ul>	Apresentação de informação, visualização	R R packages for IR R packages for TM
4 17 - mar	<ul style="list-style-type: none"> <li>- Técnicas de indexação: Inverted index creation and compression</li> <li>- Posting lists: Representação e manipulação de posting lists</li> </ul>	Construção de índices	Lucene
5 24 - mar	<ul style="list-style-type: none"> <li>- Retrieval Models: Probabilistic and language modeling approaches</li> <li>- Ranking Algorithms: Okapi BM25 and Language Model for Information Retrieval (LMIR)</li> <li>- Relevance Feedback: Techniques for improving search results based on user feedback.</li> <li>- Query Expansion: Expanding queries to enhance retrieval effectiveness</li> </ul>	Modelação de textos e conjuntos de textos (corpora)	Métricas de avaliação: Precision, recall, F1-score, and their significance Evaluation Metrics: ROC curve, AUC, Precision/Recall curve Entrega dos artigos - IR Seminar

# Planning

S	T (2 HORAS/SEMANA)	TP (1 HORA/SEMANA)	PL (2 HORAS/SEMANA)
1 24 - fev	<ul style="list-style-type: none"> <li>- Introdução à Recuperação de Informação (IR), Filtragem de Informação (IF), Extração de Informação (IE) e Mineração de Texto</li> <li>- Processo geral da recuperação de informação.</li> <li>- Indexação e consulta/pesquisa de informação não estruturada</li> </ul>	Componentes de um sistema de recuperação de informação	<b>Lançamento Projeto e Seminário. Criação dos grupos. Seleção de temas. Preparação do artigo sobre os temas selecionados</b>
2 3 - mar	<ul style="list-style-type: none"> <li>- Avaliação em recuperação e extração de informação. Precision, recall, F1-score</li> <li>- ROC curve, AUC, Precision/Recall curve</li> </ul>	Avaliação em recuperação e extração de informação	Preparação do artigo
3 10 - mar	<ul style="list-style-type: none"> <li>- Boolean Retrieval Model: Basics of Boolean algebra and operations in information retrieval</li> <li>- Conjunctive and Disjunctive Queries: Understanding AND, OR, and NOT operators</li> </ul>	Apresentação de informação, visualização	R R packages for IR R packages for TM
4 17 - mar	<ul style="list-style-type: none"> <li>- Técnicas de indexação: Inverted index creation and compression</li> <li>- Posting lists: Representação e manipulação de posting lists</li> </ul>	Construção de índices	Lucene
5 24 - mar	<ul style="list-style-type: none"> <li>- Retrieval Models: Probabilistic and language modeling approaches</li> <li>- Ranking Algorithms: Okapi BM25 and Language Model for Information Retrieval (LMIR)</li> <li>- Relevance Feedback: Techniques for improving search results based on user feedback.</li> <li>- Query Expansion: Expanding queries to enhance retrieval effectiveness</li> </ul>	Modelação de textos e conjuntos de textos (corpora)	Métricas de avaliação: Precision, recall, F1-score, and their significance Evaluation Metrics: ROC curve, AUC, Precision/Recall curve  <b>Entrega dos artigos - IR Seminar</b>

# Planning

S	T (2 HORAS/SEMANA)	TP (1 HORA/SEMANA)	PL (2 HORAS/SEMANA)
6 31 - mar	Apresentação dos temas escolhidos pelos grupos para o sistema de IR/TM a implementar.	Avaliação dos artigos - IR Seminar	Avaliação dos artigos - IR Seminar Iniciar fase 2 (projeto IR/TM).
7 7 - abr	<ul style="list-style-type: none"> <li>- Processo de recuperação de informação na web, motores de pesquisa</li> <li>- Algoritmo PageRank; Algoritmos: Introdução à análise de hiperlinks e grafos web.</li> <li>- Recuperação de informação centrada no utilizador</li> <li>- Sistemas de recomendação Collaborative filtering</li> </ul>	Modelação de textos e conjuntos de textos (corpora). Avaliação de sistemas de recuperação de informação. Avaliação de Sistemas de Recomendação: Métricas para avaliar a qualidade das recomendações	Projeto IR/TM
8 14- abr	<ul style="list-style-type: none"> <li>- Processo de recuperação de informação na web, motores de pesquisa</li> <li>- Algoritmo PageRank; Algoritmos: Introdução à análise de hiperlinks e grafos web.</li> <li>- Recuperação de informação centrada no utilizador</li> <li>- Sistemas de recomendação Collaborative filtering</li> </ul> <p>Páscoa (MINTRI tem aulas à 6ªf)</p>	Modelação de textos e conjuntos de textos (corpora). Avaliação de sistemas de recuperação de informação. Avaliação de Sistemas de Recomendação: Métricas para avaliar a qualidade das recomendações	Páscoa (MINTRI tem aulas à 6ªf)
21-abr	Páscoa	Páscoa	Páscoa
9 28 - abr	<ul style="list-style-type: none"> <li>- Modelação de dados semiestruturados e não estruturados. Modelação de textos e conjuntos de textos (corpora). Dicionários, Estruturas de dados, Tolerância de termos</li> <li>- Term Frequency-Inverse Document Frequency (TF-IDF)</li> <li>- Vector Space Model</li> </ul>	Medidas de similaridade entre documentos e consultas	Caso de estudo de mineração de texto

# Planning

S	T (2 HORAS/SEMANA)	TP (1 HORA/SEMANA)	PL (2 HORAS/SEMANA)
6 31 - mar	<b>Apresentação dos temas escolhidos pelos grupos para o sistema de IR/TM a implementar.</b>	<b>Avaliação dos artigos - IR Seminar</b>	<b>Avaliação dos artigos - IR Seminar</b> <b>Iniciar fase 2 (projeto IR/TM).</b>
7 7 - abr	<ul style="list-style-type: none"> <li>- Processo de recuperação de informação na web, motores de pesquisa</li> <li>- Algoritmo PageRank; Algoritmos: Introdução à análise de hiperlinks e grafos web.</li> <li>- Recuperação de informação centrada no utilizador</li> <li>- Sistemas de recomendação Collaborative filtering</li> </ul>	Modelação de textos e conjuntos de textos (corpora). Avaliação de sistemas de recuperação de informação. Avaliação de Sistemas de Recomendação: Métricas para avaliar a qualidade das recomendações	Projeto IR/TM
8 14- abr	<ul style="list-style-type: none"> <li>- Processo de recuperação de informação na web, motores de pesquisa</li> <li>- Algoritmo PageRank; Algoritmos: Introdução à análise de hiperlinks e grafos web.</li> <li>- Recuperação de informação centrada no utilizador</li> <li>- Sistemas de recomendação Collaborative filtering</li> </ul> <p>Páscoa (MINTRI tem aulas à 6ªf)</p>	Modelação de textos e conjuntos de textos (corpora). Avaliação de sistemas de recuperação de informação. Avaliação de Sistemas de Recomendação: Métricas para avaliar a qualidade das recomendações	Páscoa (MINTRI tem aulas à 6ªf)
21-abr	Páscoa	Páscoa	Páscoa
9 28 - abr	<ul style="list-style-type: none"> <li>- Modelação de dados semiestruturados e não estruturados. Modelação de textos e conjuntos de textos (corpora). Dicionários, Estruturas de dados, Tolerância de termos</li> <li>- Term Frequency-Inverse Document Frequency (TF-IDF)</li> <li>- Vector Space Model</li> </ul>	Medidas de similaridade entre documentos e consultas	Caso de estudo de mineração de texto

# Planning

S	T (2 HORAS/SEMANA)	TP (1 HORA/SEMANA)	PL (2 HORAS/SEMANA)
5 - mai	Queima das Fitas	Queima das Fitas	Queima das Fitas
10 12 -mai	- Pré-processamento de texto (stop words, stemming, seleção de atributos, thesauri, etc) - Pré-processamento de texto (stop words, stemming, seleção de atributos, thesauri, etc)	Ferramentas, linguagens e bibliotecas (R, Wordnet, , etc)	Projeto IR/TM
11 19 -mai	- Operações sobre texto (POS tagging, NER, ...) - Operações sobre texto (POS tagging, NER, etc)	Ferramentas, linguagens e bibliotecas (R, Wordnet, etc)	Projeto IR/TM
12 26 -mai	- Aprendizagem automática na recuperação de informação - Aplicações (classificação, clustering, sentiment analysis, topic drift, etc) - Classificação de texto	Tendências na área da recuperação de informação e mineração de texto Considerações éticas	Projeto IR/TM
13 2 - jun	Projeto IR/TM	Projeto IR/TM	Projeto IR/TM.
14 9 - jun	<b>Entrega do projeto.</b> <b>Avaliação</b>	<b>Avaliação</b>	<b>Avaliação</b>
15 16 -Jun	<b>Avaliação</b>	<b>Avaliação</b>	<b>Avaliação</b>
16 23 -Jun	<b>Avaliação</b>	<b>Avaliação</b>	<b>Avaliação</b>

# Evaluation

Brief description, purpose



# Evaluation

Components	Type	Weight	Min	Repeatable 1 <sup>st</sup> Exam	Repeatable 2 <sup>nd</sup> Exam
M1 – Seminar	Paper	25	9,5 (weighted average)	NR	NR
M2 – Project	Project	45		NR	NR
EG – Exam	Written	30	8.0	Yes	Yes

- Final Classification (CF)
  - $CF = M1 * 0,25 + M2 * 0,45 + EG * 0,3$
- Partial Improvement Repeatable Component (MP) - EXAM
  - $CF = MNR * 0,7 + MP * 0,3$