

**UNIVERSIDADE DE SÃO PAULO  
INSTITUTO DE CIÊNCIAS MATEMÁTICAS E DE COMPUTAÇÃO**

**Henrique Cursino Vieira**

**Classificação de patentes utilizando random forest**

**São Carlos**

**2020**



**Henrique Cursino Vieira**

## **Classificação de patentes utilizando random forest**

Trabalho de conclusão de curso apresentado ao Centro de Ciências Matemáticas Aplicadas à Indústria do Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo, como parte dos requisitos para conclusão do MBA em Ciências de Dados.

Área de concentração: Ciências de Dados

Orientador: Prof. Dr. Nikolai Valtchev Kolev

**Versão original**

**São Carlos**

**2020**



*“Nenhuma grande descoberta foi feita jamais sem um  
palpite ousado.”  
Isaac Newton*



## LISTA DE ABREVIATURAS E SIGLAS

IDE	Integrated Development Environment
SVM	Support Vector Machine
TF-IDF	Term Frequency - Inverse Document Frequency
FPO	Free Patents Online
DTM	Document-Term Matrix
LDA	Latent Dirichlet allocation





## SUMÁRIO

<b>1</b>	<b>INTRODUÇÃO</b>	<b>9</b>
<b>1.1</b>	<b>Apresentação</b>	<b>9</b>
<b>1.2</b>	<b>Justificativa</b>	<b>10</b>
<b>1.3</b>	<b>Problema</b>	<b>10</b>
<b>1.4</b>	<b>Objetivo geral</b>	<b>10</b>
<b>1.5</b>	<b>Objetivos específicos</b>	<b>10</b>
<b>1.6</b>	<b>Metodologia</b>	<b>11</b>
<b>2</b>	<b>DESENVOLVIMENTO</b>	<b>13</b>
<b>2.1</b>	<b>Revisão sistemática</b>	<b>13</b>
2.1.1	Descrição do objeto de estudo	13
2.1.2	Delineamento da pesquisa	13
<b>2.2</b>	<b>Materiais e métodos</b>	<b>13</b>
2.2.0.1	Extração dos dados	13
2.2.0.2	Construção do dicionário	13
2.2.1	Validação do dicionário	15
<b>2.3</b>	<b>Resultados parciais</b>	<b>15</b>
2.3.1	Extração de dados	15
2.3.2	Construção do dicionário	15
2.3.2.1	Levantamento de tópicos	15
2.3.2.2	Validação dos tópicos	16
2.3.2.3	Expansão do dicionário	16
2.3.2.4	Modelo	17
<b>2.4</b>	<b>Discussão parcial</b>	<b>18</b>
	<b>REFERÊNCIAS</b>	<b>19</b>



# 1 INTRODUÇÃO

## 1.1 Apresentação

No desenvolvimento geral de produtos, a pesquisa por documentos de patentes visa garantir o não infringimento de propriedades intelectuais que ainda não estão em domínio público (BREITZMAN; MOGEE, 2002). O sistema de patentes é um conjunto de medidas utilizados para visar o retorno do valor privado investido ao valor social de suas invenções, fornece aos inventores um período temporário de poder de mercado, recuperando os custos de seus investimentos na pesquisa (WILLIAMS, 2017). De acordo com o *World Intellectual Property Indicator 2017*, em 2016, o número de documentos de patente excedeu 3 milhões pela primeira vez, um aumento de 8.3% (LI, 2018). Em uma pesquisa de documentos de patente, documentos relacionados a tecnologia, economia e jurídico são tratadas, classificadas e analisadas para se obter um alto vantagem técnica e comercial (LI, 2018).

A classificação de documentos é o processo de classificação de um documento em uma categoria predefinida, desempenhando um papel importante no gerenciamento e busca de temas (ANNE et al., 2017). A automatização da classificação de documentos a partir de aprendizado de máquina, pode rotular documentos de um tema único e a rotulagem em vários temas é relativamente desafiador (ANNE et al., 2017).

De acordo com Shahid et al (2019), a classificação de documentos de patente em temas e a atribuição de valor de relevância para estes temas, permitem ao pesquisador filtrar as patentes que o interessa e reduzindo o escopo de análise. Nesse trabalho, realizou a construção de uma matriz de valores de term frequency - inverse document frequency (TF-IDF), notações e peso ponderado por BM25, que posteriormente foi testado em diferentes classificadores, classificando os documentos de patente em cada assunto. Vide Anne et al (2017), identificou uma matriz de métodos a serem aplicados com os modelos k-Nearest Neighbors (kNN), Support Vector Machine (SVM), Random Forest e J48. Os principais passos para essa pesquisa foram técnicas de seleção de características, com uso de ganho de informação e correlação para efetividade do classificadores.

Destes dois estudos, foi observado que a adição de mais características para os modelos de classificação utilizados, a acurácia foi melhorada (SHAHID et al., 2020). E que obstáculos, como o desbalanceamento dos dados foram atenuados pela adição de novas características (ANNE et al., 2017). Balancear a relação entre esses dois pontos é um desafio quanto a classificação de documentos de patente.

## 1.2 Justificativa

Como tratar, classificar e analisar documentos de patente havendo algumas centenas de documentos sobre um assunto específico? O método tradicional necessita de tempo e equipe para realizá-lo, apresentando um resultado com deficiências devido ao alto volume de documentos de patente a serem analisadas (LI, 2018). Hoje, já há portais web que oferecem ferramentas das quais algumas auxiliam ao pesquisador a reduzir essa pesquisa (ABBAS; ZHANG; KHAN, 2014), mas classificam os documentos em uma relevância geral. Esse resultado somente demonstra que dentro daquela amostra de documentos, uma visão macro sobre o tema que muitas vezes o pesquisador está em busca de um subtema, como quais mercados essa tecnologia está presente, quais os processos de produção desta tecnologia ou qual a formulação desse composto.

## 1.3 Problema

Com o rápido crescimento de documentos de patente, torna-se urgente a questão de automatização da classificação de documentos de patente de forma acurada e rápida (ZHU et al., 2020). Os documentos de patente contem um potencial conhecimento tecnológico na resolução de problemas no processo de fabricação, nos quais são de grande valor científico e tecnológico, no entanto, esse conhecimento está implícito em longos textos (LI, 2018; WANG et al., 2016). A classificação de documentos de patente em temas e subtemas utilizando de modelos de aprendizado de máquina se beneficiaria do uso da extração de características úteis vindas do próprio documento (ANNE et al., 2017). Observa-se que mais de 90% das informações de científicas e tecnológicas estão em documentos de patente, e sua análise resultaria em decisões de negócio de sucesso (LI, 2018).

## 1.4 Objetivo geral

Este projeto se propõe a classificar documentos de patente por tema específico e subtemas de interesse do pesquisador, reduzindo o escopo de documentos de patentes a serem estudados à somente os mais relevantes para o que se procura.

## 1.5 Objetivos específicos

A classificação de documentos de patente envolverá o uso técnicas de processamento de linguagem natural para o tratamento e preparação dos dados que serão usados no modelo de classificação por relevância que será desenvolvido. Este modelo usará inicialmente a medida estatística TF-IDF e avaliaremos outras medidas. Haverá a necessidade de criação de dicionários que auxiliem na classificação dos documentos de patente. E então será treinado um algoritmo para classificar os documentos de acordo com o tema.

## 1.6 Metodologia

Realizaremos a obtenção de um conjunto de documentos de patente aplicado a agricultura através da ferramenta Free Patents Online - FPO (<https://www.freepatentsonline.com/>). Não foi encontrado artigos ou materiais que fizessem essa aplicação para patentes relacionadas ao setor agrônomo, para gerenciamento de patentes, desenvolvimento de produtos e descoberta de mercados. Faremos o uso do modelo de classificação baseado em florestas aleatórias, a vantagem desse modelo, é a flexibilidade para o uso em regressão e classificação, além da sua facilidade de interpretação do resultado obtido. A construção de dicionários será a partir de técnicas de Processamento de Linguagem Natural, elencando as palavras mais relacionadas a área. A análise, construção de dicionários e modelagem do modelos de regressão e classificação será feita na linguagem de programação Python.



## 2 DESENVOLVIMENTO

### 2.1 Revisão sistemática

O estudo de Revisão Sistemática da Literatura seguiu as recomendações Preferred Reporting Items for Systematic Reviews and Meta-Analysis – PRISMA. Foram buscados os termos: “patent mining”, “patent”, “random forest”, “machine learning” - nas seguintes bases de dados: Periodicos CAPES, Microsoft research, Semantic Scholar e Google Scholar. O intervalo de publicação dos artigos selecionados estão entre 2012 a 2020 e restrito a somente artigos escritos em inglês.

#### 2.1.1 Descrição do objeto de estudo

Foi realizado a extração de dados de documentos de patentes no site Free Patents Online - FPO (<https://www.freepatentsonline.com/>). Este site contem os dados dos documentos de patentes de forma pública.

#### 2.1.2 Delineamento da pesquisa

Foi buscado o termo “agronomy” e filtrado para somente documentos de patentes registrados nos Estados Unidos. Foi totalizado 12906 patentes, dos quais selecionamos uma amostragem das 200 primeiras patentes. Construímos uma aplicação de webscraping na linguagem Python para realizar a extração dos dados de documentos de patentes. Os dados extraídos foram armazenados em um banco de dados.

### 2.2 Materiais e métodos

#### 2.2.0.1 Extração dos dados

A aplicação de webscraping dos dados de documentos de patentes foi escrita na linguagem de programação Python, com uso das bibliotecas *requests* e *BeautifulSoup*. Essa aplicação é modular o suficiente para que seja definido quantos documentos de patentes terão suas informações extraídas, como também quais informações serão extraídas, como demonstrado na figura 1. Os dados são organizados na forma de tabela e armazenado em um pequeno banco de dados feito em *SQLite*.

#### 2.2.0.2 Construção do dicionário

A construção do dicionário que será utilizado no projeto é composto pelas seguintes etapas, figura 2, geração de um corpora de documentos de patentes, pré processamento do corpora, obtenção da matriz de documento-termo (Document-Term Matrix – DTM) e aplicação do modelo Latent Dirichlet allocation (LDA). A partir dos tópicos apresentados

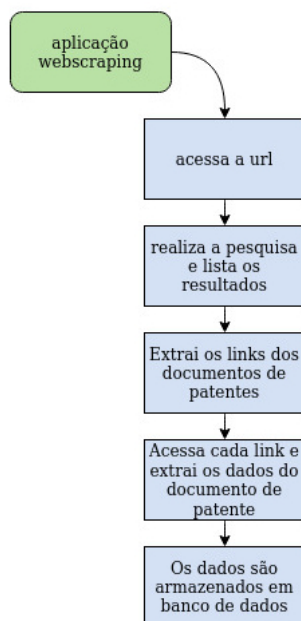


Figura 1: Fluxo de captura dos dados de documentos de patente

pelo resultado do LDA, são adicionados ao tópicos, palavras relacionadas, tais como sinônimos, hiperônimos e hipônimos através do banco de dados wordnet.

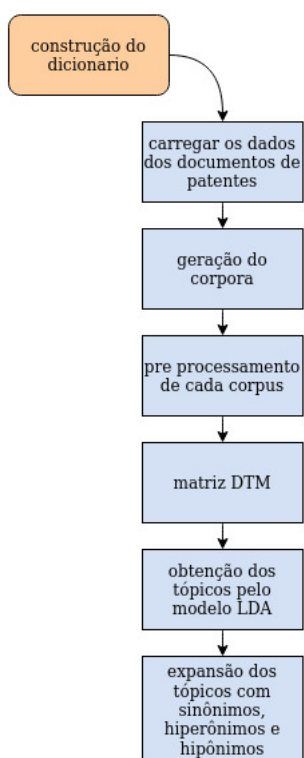


Figura 2: Fluxo de criação do dicionário



### 2.2.1 Validação do dicionário

A avaliação do dicionário obtido, consiste em observar se o valor  $k$  utilizado para geração de tópicos conseguiu separar adequadamente os assuntos contidos no corpora.

## 2.3 Resultados parciais

### 2.3.1 Extração de dados

Extraímos uma amostra no total de 200 documentos de patentes através do uso da técnica de webscraping. Destes documentos, os dados de **Título** e **Resumo** foram pré processados, removendo as quebras de linhas, espaços no início e fim da frase, uso de somente um espaço como separador e transformação em minúsculo. Estes dados foram concatenados e usados para a montagem do corpora de documentos de patentes, que poderá ser utilizado para outros projetos. Podemos visualizar na figura 3 como é distribuída a relação de palavras.

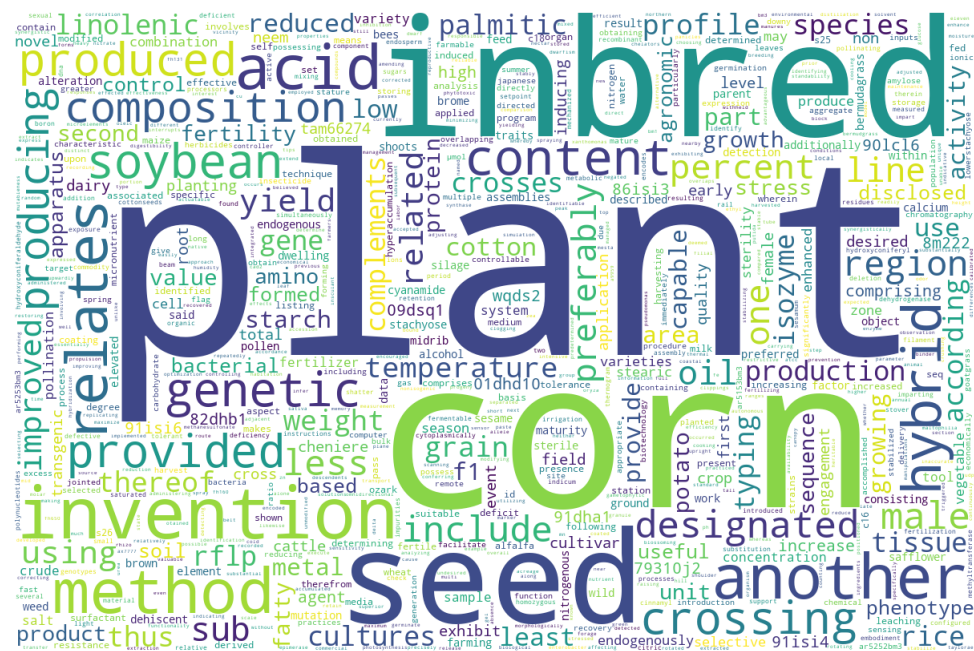


Figura 3: Construção nuvem de palavras dos termos mais representativos para este corpora.

### 2.3.2 Construção do dicionário

A construção do dicionário engloba o levantamento de tópicos, expansão dos termos, avaliação dos tópicos e expansão dos termos.

#### 2.3.2.1 Levantamento de tópicos

O corpora de documentos de patentes foi carregado, removido as stopwords, removido também caracteres numéricos e especiais, depois separados em palavras (tokens) e

desflexionados para a sua raiz (lemmas). Resultando em 200 conjuntos de palavras normalizadas representando cada documento de patente e que esta pronto para ser utilizado em modelos de Processamento de Linguagem Natural e em modelos de Aprendizado de Máquina. Aplicamos o modelo LDA, com os seguintes parâmetros - random\_state igual a 100, update\_every igual a 1, chunksize igual a 100, passes igual a 10 e alpha automático. Para definir a quantidade de tópicos  $k$ , usamos um laço de 40 interações e anotamos o valor da métrica Coherence.

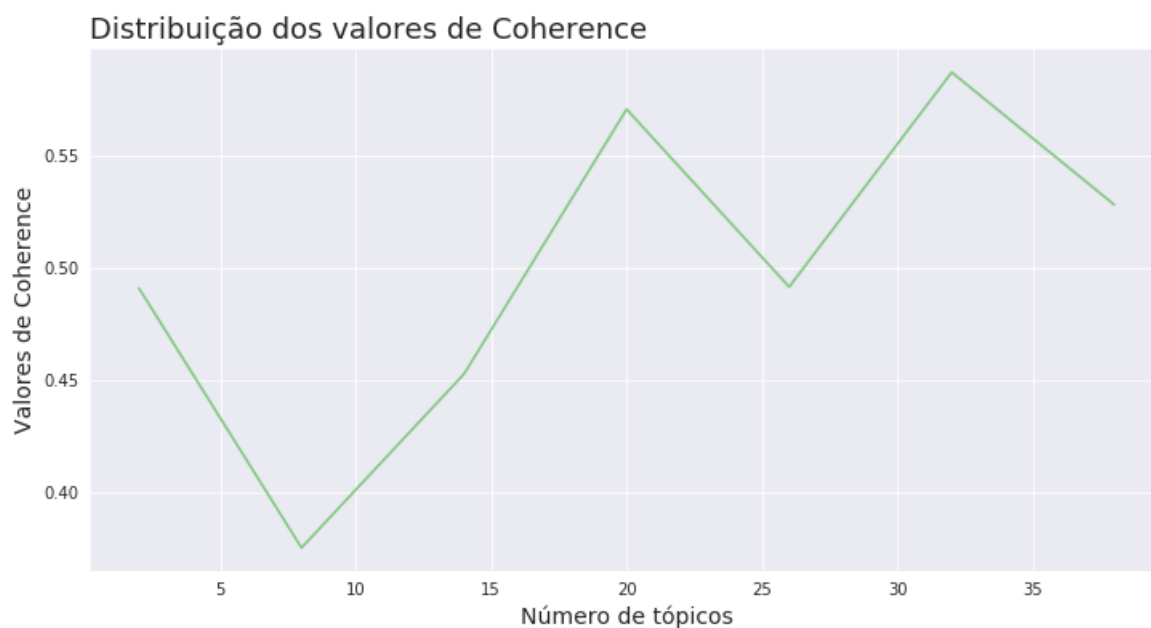


Figura 4: A distribuição dos valores de Coherence ao longo da variação do parametro  $k$ , permite que observemos qual a quantidade de tópicos devemos ter.

O gráfico aponta que um  $k$  igual a 20 resulta no segundo mais alto valor de Coherence, como visto na figura 4. Utilizaremos este valor para  $k$ , pois é um número menor de conjuntos de palavras que precisará serem interpretadas para se definir qual o título do tópico se referem.

#### 2.3.2.2 Validação dos tópicos

Examinamos o tópicos obtidos através da ferramenta pyLDAvis, figura 5. Os termos que compõe os tópicos gerados representam bem o corpora usado. Temos pouca sobreposição, com exceção do tópico 18, e os termos de cada tópico possuem uma alta relevância com o tema agronomia.

#### 2.3.2.3 Expansão do dicionário

Antes de expandir o dicionário, realizamos a remoção dois tópicos que estavam muito similares. Os tópicos geraram no total de 98 palavras únicas que foram submetidas ao wordnet e adicionado os sinônimos, hiperônimos e hipônimos destes termos, totalizando

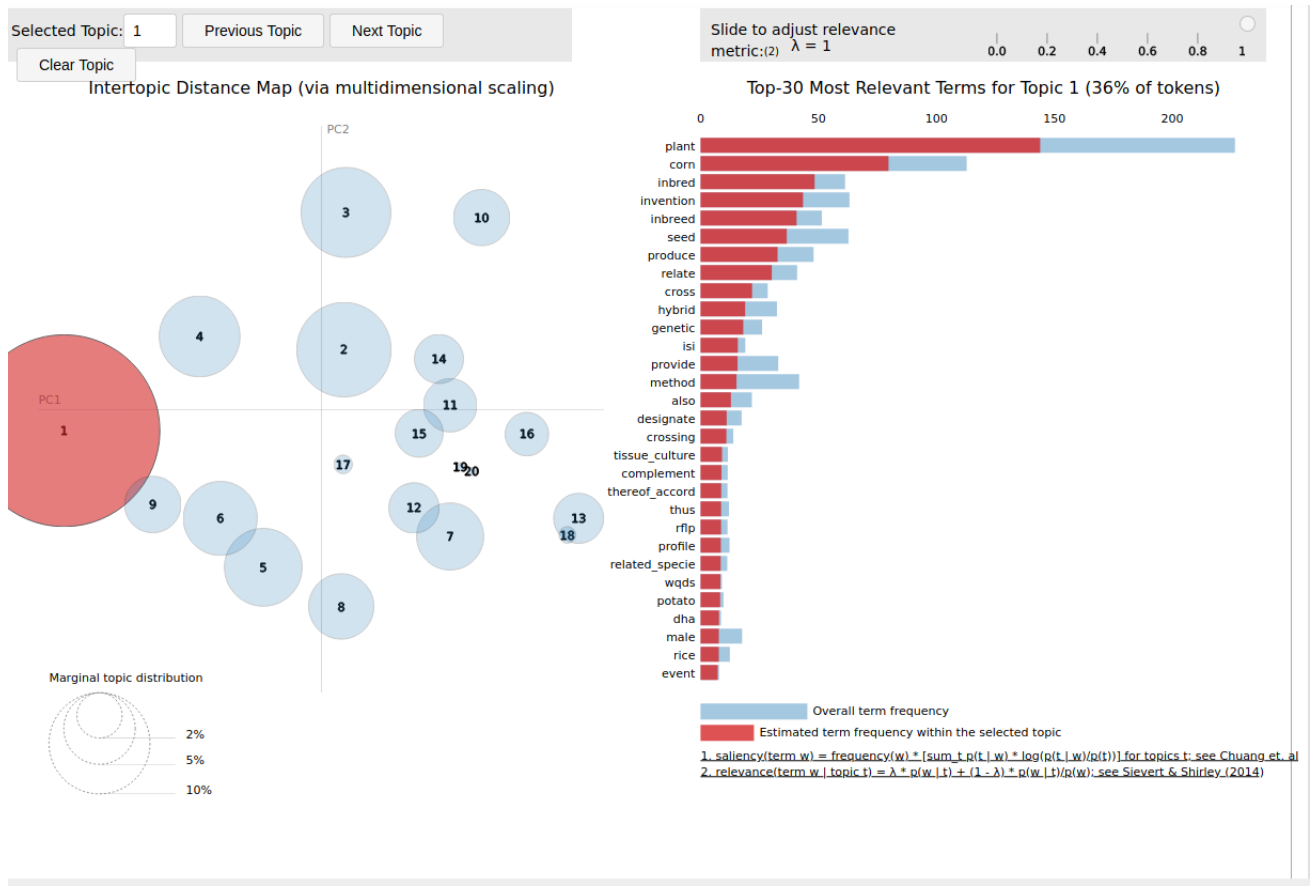


Figura 5: O grafico de bolhas, cada bolha representa um tópicos, o tamanho da bolha representa a prevalencia do tópicos e a sobreposição de bolhas aponta a similaridade entre os tópicos. O gráfico da direita, as barras representam a relevância do termo para o tópicos observado.

433 termos que representam cada tópicos. A estrutura do dicionário criado é composta por três colunas, a primeira é o tópicos, a segunda são os termos que estão atrelada ao tópicos e a terceira coluna são as palavras derivadas dos termos.

Obtivemos no final um dicionário com 954 linhas e três colunas, que foi utilizado para fazer uma classificação inicial dos documentos de patentes.

#### 2.3.2.4 Modelo

Para a construção do modelo, realizamos um pre processamento convertendo o conteúdo do documento de patente em uma matriz documento-termo, esta matriz tem a estrutura da seguinte forma:

- colunas: palavras de relevância
- linhas: documentos

- valores: correspondem ao valor de TF-IDF obtido, quando a palavra não consta na entrada, o valor será igual a zero.

O modelo testado foi o Random Forest, com critério de separação Gini, obtendo um score igual a 0,56.

## **2.4 Discussão parcial**

Conseguimos realizar a extração bem sucedida de uma amostra de documentos de patentes, do qual pré processamos e criamos um corpora que poderá ser usado não somente para este trabalho como para outros trabalhos com documentos de patentes. Fizemos um primeiro levantamento dos tópicos usando o modelo LDA e obtivemos 20 tópicos que serão avaliados e rotulados corretamente. Utilizando termos simples para rotular os tópicos neste primeiro momento percebemos que dois tópicos puderam ser removidos por serem semelhantes. Com os termos, expandimos com novos termos que são sinônimos, hiperônimos e hipônimos dos termos associados aos tópicos fazendo com que o nosso dicionário cubra uma área maior do assunto. Realizamos a construção de um modelo a partir do Random Forest, onde obtivemos o valor de score de 0,56, este é um valor ruim, o que era esperado era algo acima de 0,70.

Para os próximos passos para a finalização do trabalho, esperamos aumentar o tamanho da base de dados, testar um conjunto de parâmetros que auxiliem ao modelo resultar em um score mais alto e testar o modelo Naive Bayes e SVM.

## REFERÊNCIAS

- ABBAS, A.; ZHANG, L.; KHAN, S. U. A literature review on the state-of-the-art in patent analysis. **World Patent Information**, Elsevier Ltd, v. 37, p. 3–13, 2014. ISSN 01722190. Disponível em: <<http://dx.doi.org/10.1016/j.wpi.2013.12.006>>.
- ANNE, C. et al. Multiclass patent document classification. **Artificial Intelligence Research**, v. 7, n. 1, p. 1, 2017. ISSN 1927-6974.
- BREITZMAN, A. F.; MOGEE, M. E. The many applications of patent analysis. **Journal of Information Science**, v. 28, n. 3, p. 187–205, 2002. ISSN 01655515.
- LI, G. A Literature Review on Patent Texts Analysis Techniques. **International Journal of Knowledge and Language Processing**, v. 9, n. 3, p. 1–15, 2018.
- SHAHID, M. et al. Automatic patents classification using supervised machine learning. In: SPRINGER. **International Conference on Soft Computing and Data Mining**. [S.l.], 2020. p. 297–307.
- WANG, G. et al. Extraction of Principle Knowledge from Process Patents for Manufacturing Process Innovation. **Procedia CIRP**, The Author(s), v. 56, p. 193–198, 2016. ISSN 22128271. Disponível em: <<http://dx.doi.org/10.1016/j.procir.2016.10.053>>.
- WILLIAMS, H. L. How Do Patents Affect Research Investments? **Annual Review of Economics**, v. 9, n. 1, p. 441–469, 2017. ISSN 1941-1383.
- ZHU, H. et al. Patent automatic classification based on symmetric hierarchical convolution neural network. **Symmetry**, v. 12, n. 2, p. 1–12, 2020. ISSN 20738994.