

특허문서의 IPC 분류기 생성을 위한 데이터 전처리

박수현¹, 김진²

¹상명대학교 빅데이터융합전공 학부생

²상명대학교 빅데이터융합전공 교수

202110794@sangmyung.kr, jinkim@smu.ac.kr

Data Pre-processing for Create IPC Classifiers for Patent Documents

Su-Hyun Park¹, Jin Kim²

¹Big Data Convergence Major, Sangmyung University

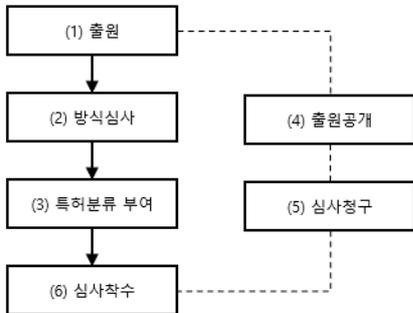
²Big Data Convergence Major, Sangmyung University

요약

특허심사절차는 짧지 않은 과정으로 이루어져 있는데, 현재 모든 절차가 사람이 직접 관여하여 진행되고 있다. 특허심사절차의 효율적 시간 분배를 위해, 특허문서 분류 과정의 자동화 처리 필요성을 느끼게 되었다. 따라서, 본 논문에서는 해당 분류기 생성을 위한 데이터의 전처리 과정을 다루었다.

1. 서론

특허청의 특허심사절차 일반[1]에 따르면, 특허심사 과정 중 특허문서의 분류 과정은 아래의 절차에 따라 진행된다.



(그림 1) 특허문서의 분류 과정.

특허문서의 분류 과정을 지나면 특허결정, 의견서/보정서, 재심사 청구 등의 과정을 거쳐 심판/소송/판결 단계까지 진행된다. 이처럼 특허심사절차의 모든 과정에 사람이 직접 관여하여, 오랜 시간이 소요된다.

이러한 배경에서 효율적인 절차 진행에 대해 고민하게 되었고, 특허문서의 분류 과정을 자동화 하는 분류기의 필요성을 느끼게 되었다.

다만, 특허문서를 분류하는 코드인 International Patent Classification (IPC)의 구조는 (그림 2)[2]의

구성처럼 방대한 분류로 이루어져 있기 때문에, 서브그룹까지 판단하는 분류기를 생성하기에는 현실적인 어려움이 있을 것으로 예상된다.

섹션(Section)	A	B	C	D	E	F	G	H
클래스(Class)	F01	F02	...	F45	F46	...	F42	F99
서브클래스(Subclass)	F01B	F01C	F01L	F01P
메인그룹(Maingroup)	F01C1/00	F01C3/00	F01C19/00	F01C20/00
서브그룹(Subgroup)	F01C1/02	F01C1/04	...	F01C1/44	F01C1/46	...	F01C1/352	F01C1/356

(그림 2) IPC의 구조.

따라서 IPC의 구조 중 서브클래스까지 판단하는 분류기 생성을 목표로, 본 논문에서는 해당 분류기 생성을 위한 데이터 전처리에 대해 논할 예정이다.

2. 사용 데이터와 데이터 전처리

특허문서 분류의 자동화를 위한 분류기에 사용할 데이터는, 기존의 특허문서 데이터가 포함된 xml 형식의 데이터이다.

‘특허문서 필드의 기능적 특성을 활용한 IPC 다중 레이블 분류’[3]를 보면, 청구항 필드는 타 필드에 비해 데이터가 많아 분류기 성능을 떨어트리는 요소로 작용하는 것을 알 수 있다. 다만, 청구항 중 독립청구항만 포함했을 경우의 분류기 성능은 확인되지 않았다.

따라서 독립청구항 여부에 따른 분류기 성능 비

