

빠르게 정확도 높여가는 인공지능 번역

최근 자동 번역에 도입된 인공지능경망 기술이 영어 고민을 덜어주고 있다. 해외여행 때마다 모르는 언어로 쓰인 메뉴판을 보고 음식을 주문해야 하는 상황이나 길을 찾거나 각종 안내 표지판을 확인할 때 도움을 주고 있는 것. 실시간 번역이 가능해 외국인도 두렵지 않다.

김형자 과학칼럼니스트



인공신경망 번역기, 문장 통째로 파악해 오류 적어

최근 구글의 번역 서비스를 이용해 보면 놀랄 만큼 자연스러운 문장을 만들어낸다. 인공지능(AI) 덕분에 가능해진 변화다. 인공지능에 기반을 둔 '인공신경망 기계 번역'(NMT)의 등장이다. 인공신경망 기술은 기계학습 기법 중 하나다. 인간 뇌의 신경망을 본뜬 인공신경망을 통해 기계가 스스로 학습하는 것이다. 마치 아기가 시행착오를 겪으면서 무언가를 배우듯, 컴퓨터가 주어진 데이터를 반복적으로 분석해 의미를 찾아내고 스스로 배운다.

예전부터 구글 번역기는 기계학습으로 언어들을 번역해 왔다. 이전에 구글이 이용한 방법은 '통계 기반 번역'이다. 번역할 문장이 들어오면 단어와 구 형식으로 각각 나눠 과거에 사람이 번역해 놓은 문서에서 어떻게 번역됐는지 통계에 따라 번역 결과를 내보낸 것. 단어 하나하나를 번역하고 단어들을 조합해서 문장을 만들었다는 얘기가. 예를 들어 'do homework=숙제를 하다', 'do study=공부를 하다'를 바탕으로 'do X=X를 하다'라고 번역하는 방식이다. 단어와 문구에 치우치다 보니 말투가 이상하거나 단어를 잘못 번역하는 경우가 많다.

그래서 등장한 것이 인공신경망이다. 인공신경망 번역은 지금까지 번역된 글과 원문을 찾아서 단어 사이의 연결성을 보고 인공지능 스스로 번역의 규칙을 찾는 것이다. 단어를 파악하는 게 아니라 기계학습을 통해 문장을 통째로 파악하고 어순, 의미, 문맥 별 의미 차이 등을 반영해 번역 결과물을 내놓는다. 예를 들어 '저녁에 집에 돌아와 저녁을 먹는다'라는 말을 번역할 때 'evening'과 'dinner'를 구분할 수 있게 된다는 의미다. 이처럼 단어나 문구가 아닌 문장 전체를 분석하기 때문에 맥락에 대한 이해도가 높아져 번역 결과가 더욱 정확하고 자연스럽다. 그렇다면 기계학습이 어떻게 이뤄지기에 이런 결과가 나올 수 있는 걸까?

사실 인공지능에 기반을 둔 '인공신경망'의 언어 습득은 사람이 배우는 것과 비슷한 점이 있다. 어린아이가 부모 행동을 보고 똑같이 따라하는 것처럼 인공지능도 결국 자신이 배운 것을 번역한다. 지

금까지 배우지 못한 말은 번역하지 못하고, 편견이 담긴 말들을 가르치면 그런 말들을 똑같이 번역한다. 실제로 마이크로소프트가 내놓은 '테이'라는 인공지능은 인종차별적이고 혐오적인 발언들을 내뱉은 적이 있다. 사람들이 그런 나쁜 말들을 가르쳤기 때문이다.

하지만 인공지능이 언어를 배우는 과정은 우리의 학습 과정과는 다르다. 인공지능망은 문장의 모든 요소를 잘게 나누어 알고리즘(일련의 계산법)에 따라 분석한다. 이를테면 고양이 가 뭐지 모르면서도 고양이의 사진을 한없이 작은 부분들로 나누어 그 관계를 학습하며 고양이를 구분하는 법을 배우는 것과 비슷하다. 오직 수학과 확률만 따져 학습하는 것이다. 그 알고리즘은 얼마나 좋은 특징을 뽑아내느냐에 따라 성능이 크게 좌우된다. 인간이 사물을 구분하듯 데이터를 완벽히 분류하게 되면 번역은 더 정확해진다. 사람이 바둑 규칙을 일일이 입력하지 않았음에도 알파고가 이세돌 9단을 이긴 것도 이 같은 학습능력 덕분이다.

인공지능망 기법으로 학습하기 위해서는 많은 데이터가 필요하다. 과거에 사람들이 번역해 놓은 원문과 번역문이 많아야 그것을 보고 배울 수 있기 때문이다. 통계 기반 번역도 마찬가지. 많은 데이터베이스가 구축되면 번역의 정확도가 높아지지만, 사용 빈도가 낮은 언어에서는 문법 정확도가 떨어진다. 그렇다면 데이터가 충분히 학습할 수 있는 양에 미치지 못하는 영어-한국어 번역을 구글은 어떻게 극복했을까.

구조 비슷한 언어 묶어 학습해 세계 언어 번역

구글은 '다중 언어 트레이닝'이라는 새로운 기법을 도입해 이를 해결했다. 이것은 비슷한 구조를 가진 언어들을 한 번에 묶어 학습시키는 방식이다. 한국어의 경우 일본어, 터키어와 언어 구조

가 비슷하다. 그래서 구글은 어순이 비슷한 이 세 언어를 함께 학습했다. 다시 말해 한국어-영어의 데이터 부족을, 번역 데이터가 많은 일본어를 통해 인공지능이 스스로 학습하여 한국어 번역에 적용한 것이다.

이 말은 무엇을 의미할까. 영어와 한국어, 영어와 터키어 사이의 비교 데이터가 충분하다면 이를 기반으로 데이터가 부족한 한국어와 터키어도 번역할 수 있다는 얘기다. 이런 방식으로 구글은 현재 세계 인구의 3분의 1이 쓰는 11개 언어에 인공지능망 번역을 제공하고 있다.

구글은 지난 10년 동안 번역 서비스를 계속 진화시켜 왔다. 인공지능망 방식이 세상에 등장한 지는 불과 1~2년. 그럼에도 지난 10년 이상 발전해 온 '통계 기반 번역'보다 인공지능망을 통한 번역이 훨씬 좋은 결과가 나타났다. 인공지능망 기술을 적용한 후 번역 오류

가 55%에서 최대 85% 가까이 줄어든 것이다. 참으로 대단한 성과다.

구글은 인공지능망 번역과 함께 증강현실(AR) 기술도 번역에 활용하고 있다. 증강현실 번역은 카메라로 촬영한 문자를 바로 번역해 주는 기술이다. 영어로 표현된 간판 등 텍스트를 촬영해 손가락으로 문지르면 해당 부분을 번역해 준다. 바로 찍은 사진은 물론이고 이전에 촬영해 둔 사진도 불러와 손가락으로 번역할 문자를 문지르면 해당 문자를 바로 번역해 준다.

물론 인공지능망 번역이 아직은 사람 수준에 미치지 못하는 게 사실이다. 짧은 문장은 인식률이 높지만 문장이 길어질수록 오차가 크게 나타난다. 하지만 인공지능망은 방대한 데이터를 모아 스스로 학습해 가며 인간과 비슷한 수준의 번역 결과물을 곧 내놓게 될 것이다. 어쩌면 인공지능의 언어와 번역 능력이 완벽에 가까워질 때, 우리는 도리어 인간의 고유한 영역에 대한 확신을 잃게 될지도 모를 것이다. TTA

