

8주차 (특강)

서비스향 ai 모델 개발하기

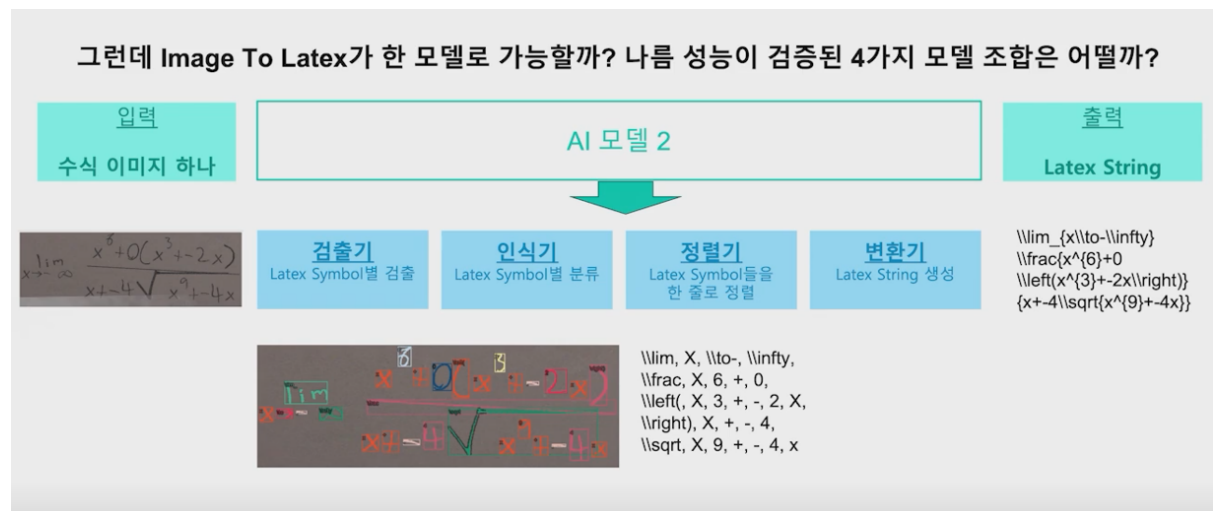
보통 수업/학교/연구에서는 정해진 데이터셋/평가 방식에 성능을 만드는 것을 목표로 가집니다. 다만, 서비스 개발 시에는 학습 데이터셋이 존재하지 않는다. 다만 서비스 요구 사항만이 존재한다. 따라서 첫번째로 학습 데이터셋을 정해야하고, 학습 데이터셋의 종류, 수량, 정답을 정하게 된다.

어떤 서비스를 개발해야할지 정확히 정의하고 시작하게 된다.

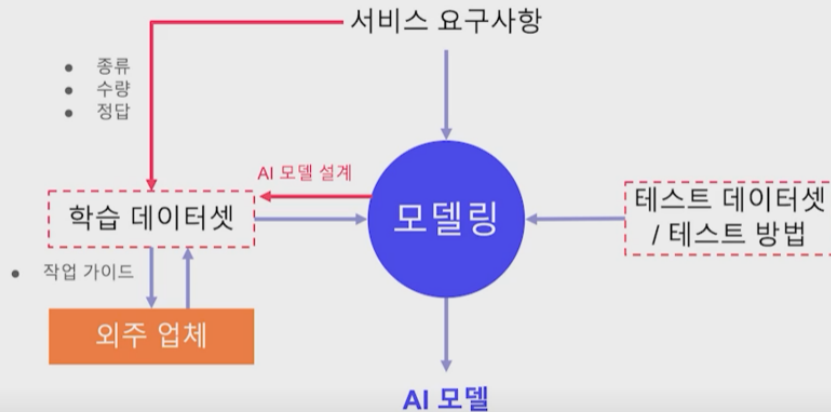
질의응답을 통해 데이터셋의 종류, 수량, 정답 관련 요구사항을 구체화한다.

그 표현과, 데이터셋의 종류를 구체화 하게 되고, 변화할 데이터셋 종류를 전부 대응해야하는지 확인하게 된다.

결국, 종류에 대해서도 정의를 해야한다. 어디까지 종류로 정의해서 각각 몇 장을 수집할 것인지 정해야 한다.



테스트 데이터셋은 학습 데이터셋에서 일부 사용한다고 하고, (사실은 이것도 할 얘기가 많지만..) 서비스 요구사항으로부터 테스트 방법을 도출해야 한다.



부드러운 워크플로우를 간단히 알 수 있었습니다.

실 서비스 전에 개발 환경에서 정량 평가와 실 서비스 적용시의 정량평가는 이질감이 굉장히 클 수 있다.

결국 서비스에서의 품질이 중요하기 때문에 오프라인 테스트 결과가 온라인 테스트 결과와 유사하게 오프라인 테스트를 설계해야 한다.

	Offline 서비스 적용 전 성능 평가	Online 서비스 적용 시 성능 평가
정량 평가	완벽하지 않기 때문에 AI 모델 후보 선택 목적으로 활용	해당 AI 모델을 서비스 시나리오에서 자동 정량 평가
정성 평가	각 후보 AI 모델에 대한 면밀 분석 후 서비스 출시 버전 선택	VOC (Voice Of Customer) AI 모델 개선 포인트 파악

또, 처리시간 또한 상당히 중요하다.

이미지가 이럽 후 수식 영역 정보가 출력될 때 까지의 시간

처리 시간 / 목표 정확도

● 처리 시간

처리 시간은 하나의 입력이 처리되어 출력이 나올 때까지의 시간

예) 주식 영역 검출의 경우,

OFFLINE TEST : 이미지 입력 후 주식 영역 정보가 출력될 때까지의 시간

ONLINE TEST : 이미지 촬영 후 이미지에서 주식 영역 정보가 화면 상에 표현되기까지의 시간

● 목표 정확도

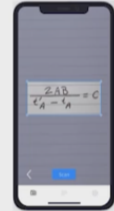
해당 기술 모듈의 정량적인 정확도

예) 신용카드 인식의 경우,

OFFLINE TEST : 입력된 이미지 내 카드 번호/유효기간에 대한 EDIT DISTANCE

ONLINE TEST : 사용자가 AI 모델의 결과값을 수정할 확률

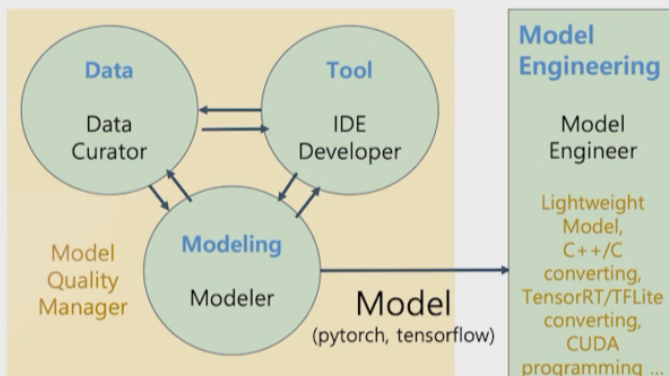
- 목표 qps
- Serving 방식
- 장비 사양



QPS: 1초당 몇번의 요청이 처리 가능하냐?

향상 방법 : 장비를 늘린다. 처리 시간을 줄인다. 모델 크기를 줄인다

우선 Serving HW향으로 모델을 최적화하는 인력이 필요하다.

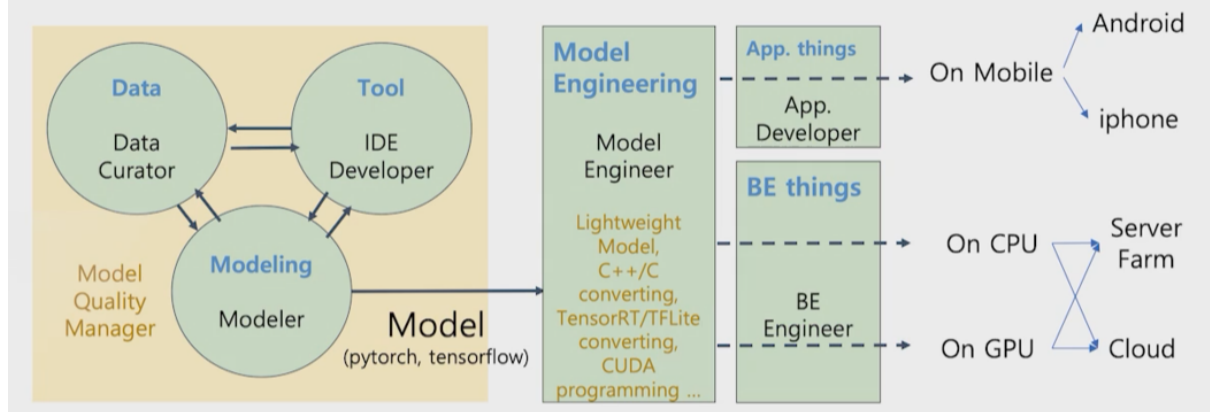


예)

- Mobile에서 구동하기 위해서
 - Pytorch → Tensorflow → TFLite로 변환
- Gpu server에서 구동하기 위해서
 - Pytorch → Tensorflow → TensorRT로 변환
- 각 toolkit에서 (TFLite, TensorRT) 지원되는 operation으로만 모델 구조 변경하여 재학습하거나 Custom Layer 구현
- 메모리를 줄이기 위한 Lightweight (경량화) 작업
 - Distillation, Quantization
- gpu 고속 처리를 위해 CUDA programming
- 아예 모든 연산을 C++/C로 변환

모델링 이후에도 참 많은 프로세스가 필요하군요?

마지막으로 모델을 실제 서빙하기 위한 추가 작업들이 end device에 맞춰 더 있다.



모델 엔지니어링/ 툴/ 서빙은 개발력이 많이 필요한 일이고, 앞으로 점차 니즈가 많아지기는 할테니 너무 AI 모델링 쪽으로 한번에 넘어가지 마라.

그 주변 기술도 많이 익히는 것이 중요하다. 대체 언제 익히란 말인가? 시간은 늘 부족하고 역량은 한정되어 있다.

ai시대의 커리어 빌딩

AI의 커리어

학교를 가야하나요? 회사를 가야하나요?

- 먼저 두 조직의 목표의 차이를 이해하는게 중요합니다.
 - 학교는 논문을 써서 연구 성과를 만드는 것이 목표
 - 회사는 서비스/상품을 만들어서 돈을 많이 버는 것이 목표
- “논문을 쓰고 싶어요”: 학교 > 회사
 - 회사는 논문을 쓸 시간적 여력이 부족할 수 있음
 - 학교는 논문을 쓰는 방식을 지도 받을 수 있음
- “상품/서비스를 만들고 싶어요”: 학교 < 회사
 - 회사는 상대적으로 데이터, 계산자원이 풍부함
 - 다 정제되어 있는 작은 규모의 토이 데이터가 아니라 어마어마한 양의 트래픽을 통해 발생하는 대규모의 리얼 데이터를 만질 수 있음 (따라서 굵은 일도 기꺼이 마다할 수도 있어야 함)

AI로 무엇을 더 잘하려는 회사가 있고,

AI로 새로운 비즈니스를 창출하려는 회사가 다르다는 것을 이해해야 합니다.

AI/ML 모델링은 팀 전체 업무의 일부

다양한 업무가 있는 만큼 팀 내에는 다양한 역할이 있음

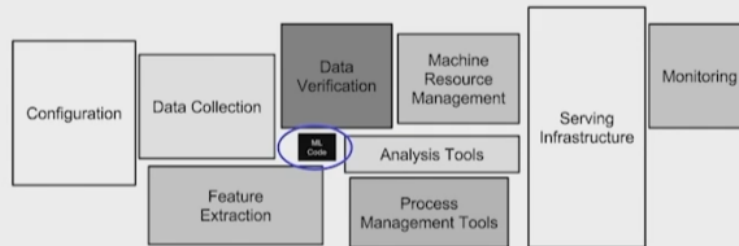


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

각 포지션에 대한 공통적인 표현은 아직 정립이 되지 않아서 모집 공고를 확실히 살펴봐야 한다.

일반적으로 내가 어떤 부분에 강점을 가진지 알고, 자신의 엷지를 어떻게 더 살릴 수 있을지 생각해봐야 한다. 내 자신을 이해하는것이 중요하지 않나?

내가 무엇에 관심이 있는지 어떻게 알아야 하나? ai 관련 인턴십을 추천한다고 합니다. AI 관련 인턴십을 조금 더 알아봐야 할 듯.

논문 재현 위주로 하는 워크샵

- 페이퍼스윛코드-rc2020

어떻게 커리어를 시작해야 하나?

일단은 인턴이 맞을 것 같습니다

자연어 처리를 위한 언어 모델의 학습과 평가

