

@삼성리서치

포트폴리오

2023.07.21
김진웅

순서

주요 개발&연구 경험

1. 머신러닝 플랫폼 스케줄러 및 모니터링 시스템 개발
2. 자동 하이퍼파라미터 최적화 시스템 개발
3. GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

입사 후 하고 싶은 일

1. 최적화
2. 자동화
3. 서빙 시스템 구축

#1 머신러닝 플랫폼 스케줄러 및 모니터링 시스템 개발

배경 및 문제점

- 클로바는 자체 머신러닝 플랫폼 NSML 을 운영하고 있었음
- 점점 유저들이 늘어남에 따라 기존 스케줄러의 문제점이 들어났음
- 단순 테이블 & 인메모리 스케줄링을 하고 있었기 때문에, 1) GPU 모니터링의 부재, 2) 큐의 부재, 3) 같은 GPU에 여러 작업 할당 등 여러 문제들이 있었음
 - 특히 동일한 GPU에 여러 학습이 수행되는 경우, 단순히 학습이 지연되는게 아니라, 메모리 침범으로 인해 이상한 결과를 도출하기도 하였음

#1 머신러닝 플랫폼 스케줄러 및 모니터링 시스템 개발

해결

- 스케줄러 및 모니터링 시스템 구축
 - 스케줄링 알고리즘 직접 디자인
 - 1) bin-packing, buddy cell allocation 응용하여 fragmentation 최소화
 - 2) 도커 이미지, 빌드 타임, 데이터셋 크기 등 고려하여 스코어링
 - 3) 거부된 작업들에 서브밋 가이드 제공
 - Hot, Warn Standby 적용하여 HA 구성
- 모니터링 시스템
 - nvidia-smi query 를 이용하여 agent 구현하여 GPU 서버들에 배포
- DB 도입을 통해 큐 구현 및 모니터링 대쉬보드 셋업

#1 머신러닝 플랫폼 스케줄러 및 모니터링 시스템 개발

결과

- 스케줄링 알고리즘 개선을 통한 작업 할당 시간 감소, 동일 GPU에 작업 이중 할당 제거, 큐 도입으로 인한 예약 학습 가능
- 모니터링 시스템 도입으로 인한 장애 감지 및 분석, 서비스 제외 자동화 도입
- 개인적으로는,
 - 수백명의 유저가 24시간, 365일 이용해야 하는 시스템 만드는 것이 매우 힘들다는 것을 알게 됨
 - 단순 업데이트나 버그 픽스가 치명적인 결과를 초래 할 수 있어서 테스트 강화 및 복구 프로세스 셋업 등 수행함
 - 실제로 약 2년동안 다운타임 10분 이내 었음

#2. 자동 하이퍼파라미터 최적화 시스템 개발

문제

- 조직내 평균 GPU 사용률이 매우 낮아서 자원들(GPU, CPU, MEMORY 등)이 낭비되고 있었음.
- 하지만, 최대 GPU 사용률은 높아서 지속적인 GPU 구매 요청이 들어왔음.
- 이에 따른 원인 분석 및 해결책 제시가 시급했음

#2. 자동 하이퍼파라미터 최적화 시스템 개발

분석&해결방법 (1/2)

- 모니터링 로그를 분석해보니, 이런 **순간적인 GPU 사용량 증가**의 많은 부분이 하이퍼파라미터 튜닝이었음
 - 많은 수의 GPU를 이용하여 탐색하고, 결과를 보고 다시 탐색을 하는데, 주말이나 연휴가 있으면 그 사이에 유휴 GPU가 많이 발생하는 상황이었음
- 이를 해결하기 위해, 하이퍼 파라미터를 자동화하는 시스템을 구축함
- 사용 편리성을 위해, 기존 시스템에 아래 예시처럼 json으로 된 파일 하나만 더 던져주면 자동으로 하이퍼파라미터를 튜닝 할 수 있게 함

Configuration

- Dictionary-based file

```
config = {  
  'n_params': {  
    'lr': {'parameters': [0.001, 0.005], 'distribution': 'log_uniform', 'type': 'float', 'p_range': [0.001, 0.005]},  
    'activation': {'parameters': ['relu', 'sigmoid', 'tanh'], 'distribution': 'categorical', 'type': 'str', 'p_range': []}  
  },  
  'entry': 'main.py',  
  'dataset': 'mnist',  
  'population': 5,  
  'step': 500,  
  'order': 'descending',  
  'measure': 'test/accuracy',  
  'tune': {'opt': {'exploit': 'truncation', 'explore': 'perturb', 'fork': False}},  
  'resources': {  
    'gpus': 1, 'cpus': 2, 'memory': '24G', 'shm-size': '1G'  
  },  
  'termination': {'max_session_number': 50}  
}
```

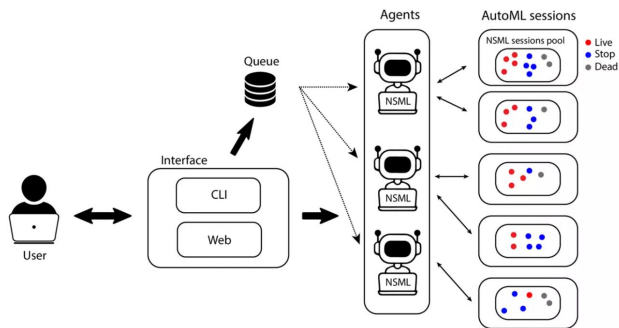
What to tune?

How to tune?

#2. 자동 하이퍼파라미터 최적화 시스템 개발

분석&해결방법 (2/2)

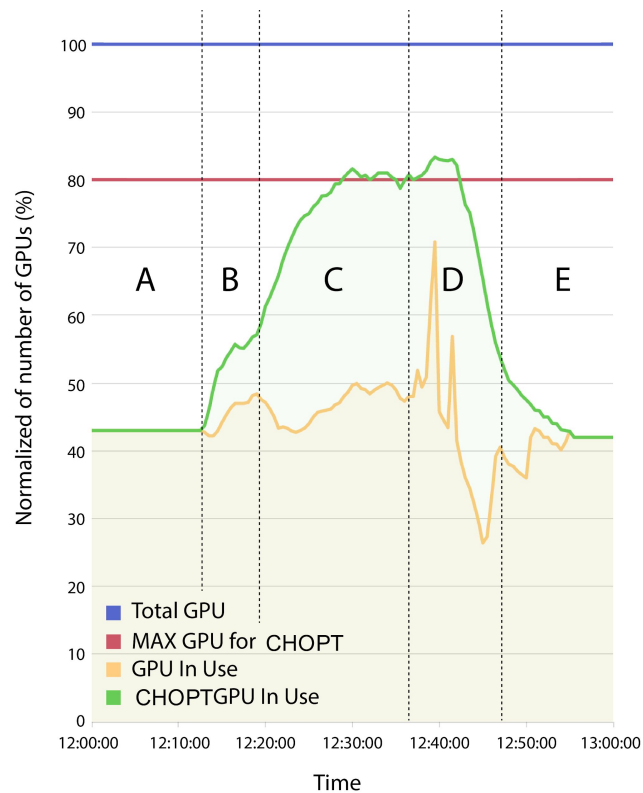
- DeepMind의 PBT 알고리즘을 참고하여 구현하였으나, 랜덤 알고리즘만으로도 엄청 좋은 결과를 얻어냈음.
- 이유를 분석해보니, 랜덤이나 PBT 같은 알고리즘이 **scalability**가 좋아서 다수의 GPU를 효과적으로 사용 할 수 있었기 때문이었음. 베이지안 같은 알고리즘은 비교적 소수의 GPU 만 사용 가능 할 때 활용하는 편이 더 좋은 결과를 보여주었음



#2. 자동 하이퍼파라미터 최적화 시스템 개발

결과

- 유휴 GPU를 활용하여 GPU 구매량을 수배 낮출 수 있었음



#2. 자동 하이퍼파라미터 최적화 시스템 개발

결과

- 유휴 GPU를 활용하여 GPU 구매량을 수배 낮출 수 있었음
- 또한, 보통 사람이 몇 주 걸릴 작업을 몇 시간 안에 끝내주게 되어 인적, 물적 자원을 수십배 절약 할 수 있게 되었음
 - Image Classification with CIFAR-100^[1]
 - Reasoning-QA with SQuAD 1.1^[2]

Task	Model	Best Score in Paper	CHOPT
Image Classification	ResNet ^[3]	76.27	77.75
	WRN ^[4]	81.51	81.66
	ResNet with RE ^{[3][5]}	77.9	79.45
	WRN with RE ^{[4][5]}	82.77	83.1
QuestionAnswering	BiDAF ^[6]	77.3	77.93

[1]: Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.

[2]: Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. Squad: 100,000+ questions for machine comprehension of text. arXiv preprint arXiv:1606.05250, 2016.

[3]: He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770 - 778, 2016.

[4]: Zagoruyko, S. and Komodakis, N. Wide residual networks. arXiv preprint arXiv:1605.07146, 2016.

[5]: Zhong, Z., Zheng, L., Kang, G., Li, S., and Yang, Y. Random erasing data augmentation. arXiv preprint arXiv:1708.04896, 2017.

[6]: Seo, M., Kembhavi, A., Farhadi, A., and Hajishirzi, H. Bidirectional attention flow for machine comprehension. arXiv preprint arXiv:1611.01603, 2016.

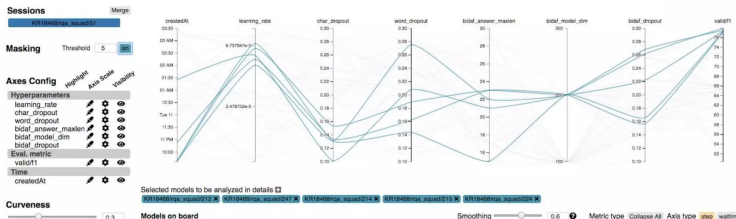
#2. 자동 하이퍼파라미터 최적화 시스템 개발

결과

- 유휴 GPU를 활용하여 GPU 구매량을 수배 낮출 수 있었음
- 또한, 보통 사람이 몇 주 걸릴 작업을 몇 시간 안에 끝내주게 되어 인적, 물적 자원을 수십배 절약 할 수 있게 되었음
- 네이버 개발자 컨퍼런스에서 발표하였으며, MLSys에서 데모 발표하였음

Select best hyperparameters

DEVIEW
2018



#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

배경 및 문제점

- 20년 하반기 OpenAI가 발표한 GPT3가 센세이셔널한 성능을 보여줌
- 기존 검색시스템 서비스들의 품질 향상을 위해 적용해보기로 함
- 하지만, OpenAI의 GPT3는 한국어 데이터가 너무 적었음
- 따라서, 한국어 및 네이버 만의 도메인 특성(블로그, 카페등)이 적용된 우리만의 언어 모델을 학습하기 시작했음
- 그런데, 학습뿐만 아니라 데이터 전송, 모델 배포, 테스트 등에 아주 많은 시간이 걸렸음

#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

해결방법

- 자동 머신러닝 파이프라인을 개발 및 최적화 하였음
- 기존 서비스에 쉽게 접목 할 수 있도록, 다양한 테스트 후 가이드 하였음
- 또한, 성능 향상을 위한 파인 튜닝 및 배포 자동화 하였음

2.1 전처리/배포 자동화



실제 서비스에 적용하기 전까지 적게는 수십, 많게는 수백 번 반복
대규모 모델의 경우 전처리와 배포 자동화가 절실 하였음



#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

해결방법

- 자동 머신러닝 파이프라인을 개발 및 최적화 하였음

2.1 전처리/배포 자동화



실제 서비스에 적용하기 전까지 적게는 수십, 많게는 수백 번 반복
대규모 모델의 경우 전처리와 배포 **자동화가 절실** 하였음



#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

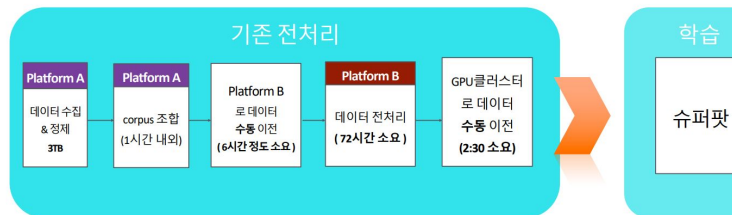
해결방법

- 자동 머신러닝 파이프라인을 개발 및 최적화 하였음

2.2 기존 전처리 방식의 문제

기존 방식에는 플랫폼간 이동과 수동 작업이 빈번했음
또한 전처리 최적화가 덜 되어 있어 많은 시간을 필요로 했음

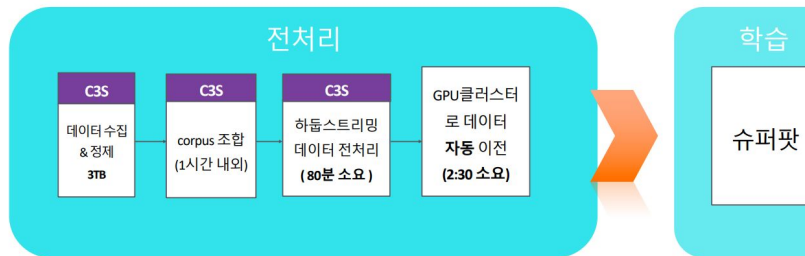
N DEVIEW
2021



2.3 기존 전처리 방식의 개선

단일 플랫폼(C3S)로 통일하여 자동화 파이프라인 구축
C3S의 하둡 스트리밍 병렬화를 이용해 전처리 시간 58배 단축

N DEVIEW
2021



#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

해결방법

- 자동 머신러닝 파이프라인을 개발 및 최적화 하였음

2.4 학습 종료 후 배포

N DEVVIEW
2021

학습 중에도 3시간마다 1TB(82B 모델 기준) 모델 저장
수시로 가져와서 성능 비교 필요하기 때문에 자동화
AiSuite*는 Kubeflow기반의 사내 머신 러닝 플랫폼



Devview 2021 'AiSuite': Kubeflow를 통해 더 나은 AI 모델 서빙과 MLOps 실현하기 - 원완규, 박정욱 참고

#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

해결방법

- 자동 머신러닝 파이프라인을 개발 및 최적화 하였음
- 기존 서비스에 쉽게 접목 할 수 있도록, 다양한 테스트 후 가이드 하였음

3.3.2 궁금증



모델 사이즈가 두배가 되면,

Q. 실행시간이 두배가 될까? → 3.4.7

Q. GPU 메모리 사용량도 두배가 될까? → 3.4.8

Q. 품질도 두배가 될까? → 4.3.4

Q. GPU 타입별 성능은 어떻게 다를까? → 4.2.3

...

Q. 입력 쿼리가 길어지면 생성 단어 속도도 느려질까? → 싱글배치

Q. 입력 쿼리가 여러개면 어떻게 될까? → 멀티배치

...

#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

해결방법

- 자동 머신러닝 파이프라인을 개발 및 최적화 하였음
- 기존 서비스에 쉽게 접목 할 수 있도록, 다양한 테스트 후 가이드 하였음

9 lessons we learned (1)

전처리/배포 파이프라인

1. MLOps는 더 이상 선택이 아닌 필수이며, 나아가서 **최적화 필요**

GPT3 Performance Study

2. Latency에 영향을 많이 주는 요소는 생성 단어 수
 - 입력 토큰 수와 모델 사이즈에 latency는 sub-linear하게 증가
 - 특정 토큰 생성시 조기 종료하는 기법은 latency 감소에 도움 될 수 있음
 - 워밍업은 초기 메모리 할당 오버헤드를 감추는데 도움이 됨
3. GPU memory에 영향을 많이 주는 요소는 입력 토큰 수, 배치 사이즈
 - 한 배치내 입력 토큰 간의 길이 차이가 심하면 가장 긴 토큰으로 성능 수렴

9 lessons we learned (2)

GPT3 성능 최적화

5. 인퍼런스 전용 트랜스포머(FasterTransformer) 사용시 2~4배 성능 향상
6. V100 → A100 변경 시 1.5~2배 성능 향상
 - 메모리 대역폭 및 NVLink 성능 향상이 많은 영향 줌
7. 파인 튜닝은 때때로 좋은 선택!

정합성 관련

- +8. 배치 사이즈, GPU 종류, 트랜스포머에 따라 생성되는 단어가 다를 수 있다
 - 내부 matrixMul 연산 및 호출 알고리즘이 달라지기 때문. 품질은 비슷
- +9. 모든 인퍼런스는 FP16으로 진행하고 있고, so far so good!

#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

해결방법

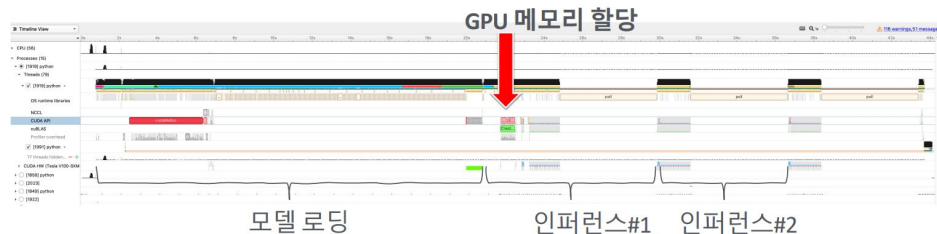
- 자동 머신러닝 파이프라인을 개발 및 최적화 하였음
- 기존 서비스에 쉽게 접목 할 수 있도록, 다양한 테스트 후 가이드 하였음

3.4.9 인퍼런스는 워밍업이 필요하다

N DEVVIEW
2021

처음 인퍼런스 수행 시 **메모리 할당에 많은 시간 소모**

메모리 할당 후 해제를 안하므로, 추후 메모리 할당 오버헤드를 피할 수 있다



#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

해결방법

- 자동 머신러닝 파이프라인을 개발 및 최적화 하였음
- 기존 서비스에 쉽게 접목 할 수 있도록, 다양한 테스트 후 가이드 하였음
- 또한, 성능 향상을 위한 파인 튜닝 및 배포 자동화 하였음

5.1. we are serving now

DEVVIEW
2021



서빙중



준비중



배치성



실시간성



대화형



요약형



제안형

#3 GPT3를 네이버 검색 서비스에 효율적으로 적용하기 위한 머신 러닝 파이프 개발

결과

- 서비스 기획자들이 쉽게 자신들의 기존 서비스에 **GPT3**를 적용 할 수 있게 되었고, 이러한 신속성으로 대한민국에서 최초로 실서비스에 **GPT3**를 적용 할 수있게 되었음. (21년 5월)
 - 기존 룰베이스로 처리 안 되던 까다로운 검색어들을 **GPT3**로 처리하여 품질을 향상 시킬 수 있게 되었음
- 새로운 **GPT3** 관련 서비스 약 10개 정도 출시함 (21년도)

입사 후 희망 업무

최적화

도메인은 계속 바뀌었지만 제가 박사때부터 한 것은 .. 최적화

- 최적화를 잘 하기 위해선 문제파악 부터
- 학습 속도 저하의 원인이 알고리즘 때문일수도 있지만 인프라 문제일수도
- 모든 부분을 분석 할 수 있어야 함
- 인프라 모니터링& 분석 부터, GPU 프로파일링까지 가능

예) 학습 속도 최적화, 인퍼런스 성능 최적화, 자원 효율화 극대 등



자동화, 서빙 시스템 구축

자동화

- 머신러닝 업무는 반복성이 매우 높으나 대부분은 자동화를 하지 않음
- 생산성을 늘리기 위해선 업무 자동화가 필수

예) 모델 서빙 시스템 구축 및 자동화, 파이프라인 자동화, 그 외 각종 잡무 자동화

서빙 시스템 구축

- 데이터가 계속 변함에 따라 인공지능 모델은 효율적이고 지속적인 업데이트 필요
- 안정적인 서빙을 위해선 **다양한 시스템적 기능이 필요**
 - 로드 밸런싱, HA, 멀티 모델, 플랫폼 지원 등
- 고성능 서빙 시스템 구축을 위해선 인공지능 뿐만 아니라 시스템 전반적인 지식 & 경험 필요

저의 주요 역량은,

- 폭 넓은 시야
 - 다양한 분야 및 R&R
- 빠른 학습 속도
 - 질문,질문,질문...
- 높은 생산성
 - 효율의 극대화, 수레바퀴를 두 번 만들지 말자, 쓸데없는 고집 및 오버 엔지니어링 지양
 - 무엇이 중요한가를 스스로, 동료 및 상사들과 끊임없이 논의하는 성향