Z_ROUR

[논문 정리] Outrageously Large

Neural Networks: The

Sparely-Gated

Mixture-of-Experts Layer ijnuvh · 2025년 4월 19일

TIL

▼ 목록 보기 11/14 () **Abstract** 정보를 흡수하는 신경망 능력 - 매개 변수의 수에 의해 제한.

통계 수정 삭제

네트워크의 일부가 예제 별로 활성화되는 조건부 계산 - 계산의 비례적인 증가 없이 모델 용량을 극적으로 증가시키는 방법으로 이론적으로 제안됨. But, 실제로는 상당한 알고리즘 및 성능 문제 존재.

결론:

위 문제를 해결하고 • 조건부 계산의 약속을 실현하여, 현대 GPU 클러스터에서 계산 효율성의 사소한 손실만으로 모델 용량에서 1000배 이상의 개선 달성.

• 최대 수천 개의 feed-forward sub-network로 구성된 a Sparsely-Gated Mixture-of-Experts layer(MoE) 도입.

• Trainable한 gating network는 각 예제에 사용될 MoE 조합을 결정.

• 여기서 모델 용량은 훈련 코퍼스에서 사용할 수 있는 방대한 양의 지식을 흡수하는 데 중요.

MoE - 언어 모델링 및 기계 번역 작업에 적용.

• 최대 1,370억 개의 매개 변수를 가진 MoE가 누적된 LSTM 레이어 사이에 Conv로 적용되는 모델 아키텍처 제시.

• 대규모 언어 모델링 및 기계 번역 벤치마크 성능 - 더 낮은 계산비용, 더 나은 결과 달성.

1. Introduction and Related Work 1-1. 조건부 계산

MoE layer

원래 조건부 계산할 때, parameter와 accuracy가 비례함.

컴퓨팅파워와 분산컴퓨팅은 이 수요를 감당할 수 없음.

 $G(x)_2$ G(x)_{n-1} Expert 1 Expert n-1 Expert n Expert 2 Gating Network

Figure 1: 반복 언어 모델 내에 포함된 MoE 계층. The sparse gating func. - 계산을 수행할 두 가지의 전문가를 선택. 그들의 출력은 Gating network output에 의해 변조된다. 위 아이디어는 이론적으로는 가능하지만, 실제로 현재까지 모델 용량, 훈련 시간 또는 모델 품질에서 대규모 개선을 입증한 작업은 아직 없음. 여러가지 난제가 있었음. 1-2. Our Approach: The Sparsely-Gated MoE Layer MoE: number of experts 로 구성, 각각은 간단한 feed-forward 신경망 & input의 처리를 위한 MoE 조합을 선택하는 trainable Gating network로 구성됨. 각 단계의 Network들은 모두 back-propagation에 의해 공동으로 train됨. "중심 Task: 언어 모델링 & 기계 번역" • 쌓여진 LSTM layer 사이에 MoE Conv 를 적용.

연구의 주제. SVMs (Collobert et al., 2002), Gaussian Processs (Tresp, 2001; Theis & Bethge, 2015; Deisenroth & Ng, 2015), Dirichlet Processs (Shahbaba & Neal, 2009) 같은 Deep Network들의 다양한 expert architecture 가 제안됨.

Jacobs et al., 1991; Jordan & Jacobs, 1994으로부터, 20년 전 도입된 동안 expert 접근 방식은 많은

Gating Network : Pre-trained 앙상블 NMT 모델에서 훈련. • NMT : *End-to-End 방식의 신경망 기반 기계번역 model*(Neural Machine Translation) 위의 논문들 - Top-Level MoE에 대한 것들. MoE = 전체 모델. Idea 1. Eigen et al.(2013)

• 복잡한 문제들은 각각 다른 expert를 필요로 하는 하위의 많은 문제들을 포함할 수도 있어서 이

• Multiple MoE 사용 & Deep model의 일부로 자체 gating network를 사용함.

• 결론에서 희소성을 도입할 것을 암시하며, MoE = 계산을 위한 수단으로 바꿈.

이 논문에서 Genral Neural Network 구성요소로서, MoEs의 이러한 사용을 기반으로 함. Idea 1의 논문은 두 set의 gating decision을 허용하는 두개의 stacked MoEs를 사용함. 반면, 우리의 MoE의 Conv적용은 MoE가 다른 gating decision을 텍스트의 각 위치에서 허용함. 또한, Sparse gating을 실현하고, model capacity를 게 늘리는 실용적인 방법으로 시연.

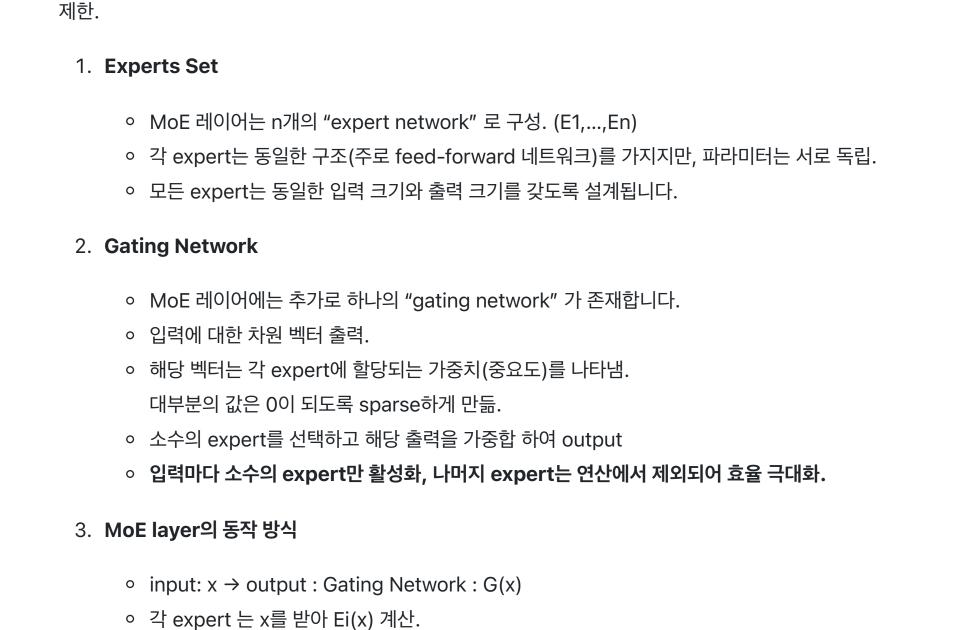
 $G(x)_2$ $G(x)_{n-1}$ Expert 1 Expert 3 Expert n-1 Expert n Expert 2 layer Gating Network

MoE layer: n개의 expert networks, (E1, ..., En), Gating network G (output이 흩어져있는 n차원

원칙적으로, expert들이 같은 크기의 input을 받아들이고, 같은 크기의 output을 생성하도록 요구함.

초기 조사에서 우리는 별도의 매개변수를 가졌지만, 동일한 구조를 가진 feed-forward network의 모델로

Expert들은 그 자신이 신경망이고, 그들은 각각 자신만의 parameter들을 가지고 있음.



방식.

3. Train 및 활용

와 같은 크기의 훨씬 적은 batch를 할당.

Expert 수의 증가에 따라 Navie MoE 구현이 비효율적이게 됨.

batch 축소에 대한 solution → original batch size 🄼

2. Taking Advantage of Convolutionality:

3. Increasing Batch Size for a Recurrent MoE:

각 디바이스의 파라미터는 파라미터 서버를 통해 동기화.

Mixing Data Parallelism and Model Parallelism

여러 디바이스에 모델 복사본을 두고 각기 다른 데이터 배치를 처리.

(k: 각 입력당 활성화되는 expert 수, n: 전체 expert 수)

• expert별 배치 크기를 d배로 늘릴 수 있음.

• 계층적 MoE의 경우,

conv. layer 활용하기

비동기적으로 처리.

데이터 병렬성:

○ MoE 레이어의 최종 출력:

2. Hierarchical MoE 확장 구조 • Expert가 너무 많을 경우, 2개의 계층으로 확장가능

• 1차 gating network가 그룹을 선택하고, 각 그룹에서 2차 gating network가 세부 expert를 선택하는

3. Addressing Performance Challenges 기존 - 매개변수 load 및 update의 overhead 때문에 큰 Batch size 가 필요함. Gating Network가 각 예제에 대해 n명의 전문가 중, K명을 선택하면 , b개의 예시 batch에 대해 kb/n<<B

MoE layer - 전체 Network와 함께 end-to-end 방식으로 사용

• gating network와 expert 모두 backpropagation으로 update.

• MoE 구조 ⇒ Transformer, LSTM 등 다양한 신경망 구조에 삽입해 사용가능.

최근의 MoE에서 batch size 증강 기존 분산 학습 방식: 분산 학습에서는 여러 디바이스(서버, GPU 등)에 모델 복사본을 두고, 각 디바이스가 서로 다른 데이터 배치를

○ 1차 gating network - 데이터 병렬 o 2차 MoE - model 병렬 • MoE: 일반 layer & Gating Network ~ 데이터 병렬로 처리 & expert는 각 디바이스에 1개씩, 모델 병렬로 처리.

디바이스 수를 늘리면 expert 수도 비례해 늘릴 수 있어, 모델 파라미터를 크게 확장할 수 있습니다. 전체 배치 크기는 늘어나지만, expert별 배치 크기, 디바이스별 메모리/대역폭 요구량, 학습 속도는 일정하게

• L_importance = importance coefficient of variation^2 × scaling fator 로드 밸런싱(Load Balancing) • 위의 중요도 균형만으로는 실제 예시 수 불균형 가능 • 해결책: expert별 예시 수를 균등하게 하는 추가 손실(L_load) 도입 5. BenchMark • MoE layer 대규모 언어 모델링 & 기계 번역 task 적용. • 최대 1370억(137 billion) parameter의 MoE layer -> LSTM stack 사이에 Conv 방식으로 삽입 후

• 해결책: 각 expert의 중요도(게이트 값 합계) 변동성을 최소화하는 추가 손실(L_importance) 도입

ijnuyh ad astra per aspera

Mask (opt.) Scale lain differences between Sparse and Soft MoE layers. While the router in Sparse to assign individual input tokens to each of the available slots, in Soft MoE layers (righ of a (different) weighted average of all the input tokens. Learning to make discrete a MatMu several optimization and implementation issues that Soft MoE sidesteps. [Paper] Attention is All You Need ... **MoE(Mixture of Experts)** YOLOv9 논문 리뷰 본 글은 Google Brain에서 2017 NIPS에 발표한 Outrageously Large Neural Networks: The YOLOv9: Learning What You Want to Learn Attention is All You Need 논문에 대한 리뷰입니 Sparsely-Gated Mixture-of-Experts Layer Using Programmable Gradient Information 다. RNN 모델의 장기 의존성 문제(long-term... Noam Shazeer, Azalia Mirhoseini*,... (2024, arxiv) 읽어보기

관심 있을 만한 포스트

댓글 작성

[논문 리뷰] Distilling the Knowledge ...

본 Paper Review는 고려대학교 스마트생산시스템

연구실 2024년 동계 논문 세미나 활동입니다.

• 0

2021년 3월 25일 · 1개의 댓글 2024년 2월 19일 · 0개의 댓글 by **뱅모** by 김재희 **♥** 3 Billion 단위의 언어모델을 학습시키기 위한 방법: Megatron-LM Billion 단위의 언어모델을 학습시키기 위한... NMT(Neural machine translation/... 최근 언어모델에서는 큰 트랜스포머 모델을 학습하는 인간의 신경을 모방한 머신러닝 기법이다. 감각기관 게 중요하다. 그러나 매우 큰 모델을 학습하는데는 많 을 통해서 정보가 들어오면, 신경세포인 뉴런이 그것 은 제약사항들이 따른다.Megatron-LM 논문에서... 을 뇌에 전달한다. 그러면 뇌는 정보를 종합해서 판... 2021년 11월 10일 · 0개의 댓글 2021년 9월 24일 · 0개의 댓글 by Sung.K by nawnoes

MoE layer 적용 시, conv 구조의 활용 장점 (시점별로 같은 연산을 적용하는 구조) 확장성: 유지됩니다. • In 언어 모델, MoE 레이어를 이전 레이어의 각 시점(time step) 출력에 동일하게 적용. • 이전 레이어가 끝난 뒤, 모든 시점의 출력을 **하나의 큰 Batch**로 묶어 한꺼번에 MoE 레이어에 입력으로 사용 시, input batch → time step 수만큼 커짐. ㅇ 연산 효율과 하드웨어 활용도 🔼

Increasing Batch Size for a Recurrent MoE

순환신경망(RNN, LSTM)에서의 적용.

"주요 성능 이슈 1: 네트워크 대역폭"

expert 활용 불균형

대역폭 보다 커야함. (GPU-수천:1)

한계 :

학습. • 기존 최고 성능(state-of-the-art) 모델 대비 현저히 더 우수한 결과 달성.

[TIL] vscode JAVA 기본 세팅 [TIL] npm error - log

2022년 12월 28일 · 0개의 댓글 2023년 8월 7일 · 1개의 댓글 2024년 3월 12일 · 0개의 댓글 by **Jomii** by kaeul **3** by choonsikmom 딥러닝 모델 경량화 딥러닝 모델 경량화는 크게 경량화 알고리즘과 알고 SPS Lab. 리즘 경량화로 나뉜다. 쉽게 말해 전자는 처음부터 가 벼운 알고리즘 쓰기이고 후자는 학습된 모델을 가볍 [Paper Review]

[논문 스터디 Week 4-5] Attention is ...

4 1

Powered by

2023년 4월 20일 · 0개의 댓글 by choonsikmom

Stellate

by choonsikmom ₩ 8 밑바닥부터 시작하는 딥러닝-8장 chapter8 딥러닝 8.1 더 깊은 신경망 8.1.1 손글씨 숫자를 인식하는 심층 CNN 손글씨심층CNN.png 3*3의 작은 필터를 사용한 합성곱 계층(Conv) 활성 화 함수(ReLU) 풀링 계층 추가해 중간 데이터의 공 간 크기를 줄여나감 완전연결 계층 뒤에 드롭 아웃 계 층 사용 완전 연결 신경망이란, 인접하는 계층의 모든 뉴런과 결합되어 있는 ... LLMs 전성시대! (Meta) LLaMA, (Sta... 2020년 1월 28일 · 1개의 댓글 LLaMA와 Alpaca 읽어보기

• 0

• 각 위치에서 potential한 다른 expert 조합 선택. • 이 potential 한 expeort는 구문과 의미론에 따라 고도로 전문화되는 경향을 가짐. 두 가지 task 벤치마크 모두, 연산은 더 적고, 성능은 더 나음. 1-3. Related work on Mixtrues of Experts

접근방식이 더 강력함.

Figure1. MoE module의 개요.

벡터) 로 구성.

• MoE - text의 각 위치에 대해 1번 호출.

다른 작업으로는 계층 구조(Yao et al., 2009), 무한한 수의 전문가(Rasmussen & Ghahramani, 2002), 순*차적 전문가 추가(Aljundi et al., 2016)*와 같은 다양한 expert 구성에 초점을 맞춤. Garmash & Monz (2016): 기계 번역을 위한 전문가들의 혼합 형식의 앙상블 모델 제안.

2. The Structure of the Mixture-of-Experts Layer MoE layer

오직 몇 개(예: 2~4개)만 활성화. 1. Gating Function & Sparsity 구현 • Gating Function : Softmax 함수와 noise를 결합한 Noisy Top-k Gating 방식 위주. → input 마다 상위 K개의 expert만 선택하여, sparsity를 강제함 sparse: 신경망의 가중치 행렬에서 대부분의 값이 0이거나 거의 0에 가까운 값일 때. • 해당 방식은 연산 효율 및 Expert간의 load balancing에도 기여. • Gating network가 특정 expert만 반복적으로 선택하는 현상을 방지하기 위해, expert 간 중요도와 부하 균형을 위한 추가 loss term을 도입.

image.png

• G(x)가 0인 expert는 아예 계산하지 않아 연산량을 크게 줄일 수 있음. \rightarrow 실제로는 수천 개의 expert 중

Batch size: forward & backward pass사이의 activation을 저장하는데에 있어 메모리 제한 존재. → Batch size를 늘리기 위한 본 논문에서의 제안 1. Mixing Data Parallelism and Model Parallelism: 데이터와 모델의 병렬성

모델 병렬성: 모델의 expert를 여러 디바이스에 나누어 배치. • 모델이 d개의 디바이스에 분산되고, 각 디바이스가 b 크기의 배치를 처리하면, 각 expert는 약 kdb/n개의 예시를 한 번에 받게 됨.

• 여러 데이터 병렬 배치에서 MoE 레이어로 들어오는 예시들을 합쳐서 각 expert가 한 번에 더 큰 배치를 처리. ○ 연산 효율 🔼 , 확장성 🔼 • 디바이스 수를 늘리면, expert 수도 늘릴 수 있음 ○ 모델 크기를 하드웨어에 맞게 확장가능. **Taking Advantage of Convolutionality:**

RNN/LSTM 가중치행렬 - MoE로 대체 시, 각 시점의 MoE입력이 이전 시점의 MoE 출력에 크게 의존 ○ Conv적 활용이 불가능함. (all time step output = 1 batch) ○ In 순환구조, MoE batch size 증강의 어려움 & 메모리 효율성 ☑ 3.2. Network Bandwidth)

• 이유 : 계산 효율 유지를 위해 expert의 연산량 대비 in/output 크기 비율이 계산 성능 대비 네트워크

• MoE 구조에서 네트워크 통신 대부분은 expert의 입력/출력 데이터 전송에 사용됨.

• 연산 효율 확보를 위해 expert의 연산량 대비 입출력 데이터 크기 비율이 커야 함.

hidden layer 크기 / 개수 → 연산 효율 / 네트워크 병목

○ L_importance : 모든 expert의 중요도가 비슷해지도록 유도.

4.Balancing Expert Utilization

● gating networ : 소수 expert에만 집중 → 불균형 발생

• 조건부 연산(Conditional Computation) 실제 구현

○ 모델 용량 1000배 이상 증가, 연산 비용 소폭 증가

○ 대규모 언어 모델링, 번역 벤치마크에서 SOTA 성능 달성

○ 데이터 크기·문제 복잡도 증가 시 모델 용량이 성능 향상 핵심

ㅇ 하드웨어 확장 시 트릴리언(1조) 파라미터 모델 학습 가능성

LSTM, Transformer 등 다양한 아키텍처 적용 가능

○ 수천 개 feed-forward expert 활용, 대규모 지식 효과적 흡수

ㅇ 입력별로 네트워크 일부만 활성화

ㅇ 데이터 병렬+모델 병렬 혼합 전략

• 자연어 처리 외 다양한 분야 활용 가능

• 모델 용량과 성능의 비례 관계 강화

• 대규모 분산 학습 확장성

• 실용적·범용 신경망 컴포넌트

0개의 댓글

MatMul

SoftMax

게 하기라고 할 수 있다. 합성곱 신경망은 처음으로 여러 개의 합성곱 층과 활성화 함수를 연속적으로 이

어 붙인 알렉스넷...

2022년 2월 8일 · 0개의 댓글

💇 🕽 by 마이클의 AI 연구소

glish to German: That is good.

cola sentence: The rse is jumping well."

ntence1: The rhino grazed grass. sentence2: A rhino grazing in a field."

rize: state authorities e damage after an onslaught weather in mississippi..."

> with a Unified Text-to-Text Transformer(2019) 읽어보기

2023년 4월 27일 · 0개의 댓글

WS by DSC W/S

댓글을 작성하세요

• 연산량(계산 비용) : 기존 모델과 비슷하거나 더 적으면서도, 모델 용량(파라미터 수) 1000배 이상 증가 • 언어 모델링 및 번역 task: 주요 지표에서 기존 대비 큰 폭의 개선을 보임. 6. Conclusion

'Das ist g not accept Τ5 "3.8" six people hospit a storm in attal T5(Text-to-Text Transfer Transfor... **Exploring the Limits of Transfer Learning**