

# A Comprehensive Survey of Continual Learning: Theory, Method and Application

## 인공지능과 유기체의 근본적 차이

- 인공지능은 일관적인 유기체를 처럼 학습 프레임워크나 인퍼런스 프레임워크가 통합된 구조가 아님.

## Continual learning의 시사하는 문제점 그리고 딜레마

- Continual learning의 핵심적인 문제는 "catastrophic forgetting".
- learning plasticity and memory stability 딜레마.
  - ↓
  - 비교적 중요성

## Continual learning의 방법론들

- continual learning은 usable resource를 최대한 활용하는 선에서 노력함. 따라서 continual learning이 추구하는 이상에 앞서는 방법론은 오직 새로운 training sample만 가지고 학습하는 방법.
- continual learning은 딥러닝 학습 프레임워크에서 어느 부분에서 솔루션을 찾아내며 여러 방법론을 크게 분류할 수 있음.
  - \* replay: 기존 task의 샘플을 다시 가져다와서 재학습에 활용.
  - \* architecture: 학습하는 네트워크 구조를 적절히 변형하여 필요한 부분에 학습.
  - \* representation: 다양한 상황에 대응해서 robust한 representation을 학습.
  - \* optimization: 기존 representation에 조금씩 더 나은 방향으로 최적화 진행.
  - \* regularization: 학습 시 정규화를 통해 학습 안정성 ↑.

## Continual learning의 다양한 시나리오들

### Instance-Incremental Learning (IIL)

- 모든 학습 sample은 같은 task에 속하며, 배치 단위로 들어옴.
- 이는 둘 일관적인 딥러닝 학습 알고리즘에서의 setting과 동일한 상황에서 input dataset만 continual하게 들어오는 상황을 가정하며, 단일 task를 가정하기 때문에 task specific identification 요구 X.

### Domain-Incremental Learning (DIL)

- IIL과 다르게 같은 task가 아닌 여러 tasks가 존재하는 상황을 전제로 함.
- 모든 task들은 동일한 data label space를 가지고 input distribution만 다름.
- 같은 task가 아님에도 불구하고 task specific identification 요구 X.

### Task-Incremental Learning (TIL)

- 각 task는 서로 disjoint한 data label space를 가짐.
- task specific identity에 대한 정보가 training 및 testing에서 제공됨.

### Class-Incremental Learning (CIL)

- 각 task는 서로 disjoint한 data label space를 가짐.
- 인퍼런스 단계에서 task specific identity 제한.

## Task-Free Continual Learning (TFCL)

- 각 task는 서로 disjoint한 data label space를 가짐.
- task specific identity 학습 단계에서 제한, 인스턴스 단계에서 optionally available.

## Online Continual Learning (OCL)

- 각 task는 서로 disjoint한 data label space를 가짐.
- 학습 시 각 task의 모든 샘플은 one-pass로 최적화에 사용됨 (한번 최적화에 사용된 데이터 배치는 다시 input으로 사용되지 않음). 따라서 줄이 task를 구분하지 않더라도 모든 배치는 single-stream framework를 사용하기 때문에 task에 무관한 학습이 됨.
- task specific identity 인스턴스 단계에서 optionally available.

## Blurred Boundary Continual Learning (BBCL)

- task의 경계가 애매(blurry)하고, data label space가 일부 겹침.
- 인스턴스 단계에서 task specific identity 제한.

## Continual Pre-training (CPT)

- 사전 학습되는 데이터가 continual (sequence)로 도착하는 상황.
- 즉각 목적은 downstream task의 성능을 올리는 것.
- pre-train 단계를 continual한 상황으로 가정하고, 때때로 testing 단계에서는 단일 downstream task를 기준으로 삼게 됨. 따라서 downstream task의 task label은 줄이 바뀌지 않음.

## Continual Learning에서의 Evaluation metric

- continual learning이 가지는 목적은 '아직 representation을 익히지 않은 새로운 task의 성능을 잘 올릴 수 있는가?'로 정리할 수 있음.

### ① 기존(old) task들의 성능 평가: Overall performance

- 기존 task들의 성능 평가 방법, average accuracy (AA)나 average incremental average (AIA)로 나타냄.
- \* AA: 여러 task들을 학습한 후, 각 task에 대한 평균 정확도를 계산.
- \* AIA: 새로운 task를 학습할 때마다, 이전에 학습된 모든 task들에 대한 평균 정확도를 계산.
- AA와 AIA를 계산할 때, 모델이 각 task에 대해 예측을 수행할 때 사용하는 output space를 고려해야 함. (IL의 경우, 모델을 현재까지 학습된 모든 클래스에 대해 예측 가능. TIL의 경우, 모델은 각 task에 대해 별도의 output space를 가질 수 있음).

### ② 기존(old) task에 대한 memory stability: Memory stability

- 기존 task를 익히거나 잘 기억하고 있는지 측정하는 지표로는 보통 forgetting measure (FM)이나 backward transfer (BWT)를 사용.
- \* FM: 각 task에 대한 forgetting 정도를 과거 학습한 성능 중 현재치와 현재 성능과의 차이로 계산.
- \* BWT: k번째 task가 이전 k-1번째 task까지의 성능에 오는 영향을 평균하여 계산.
- overall performance는 성능 자체를, memory stability는 성능 변동을 기록.

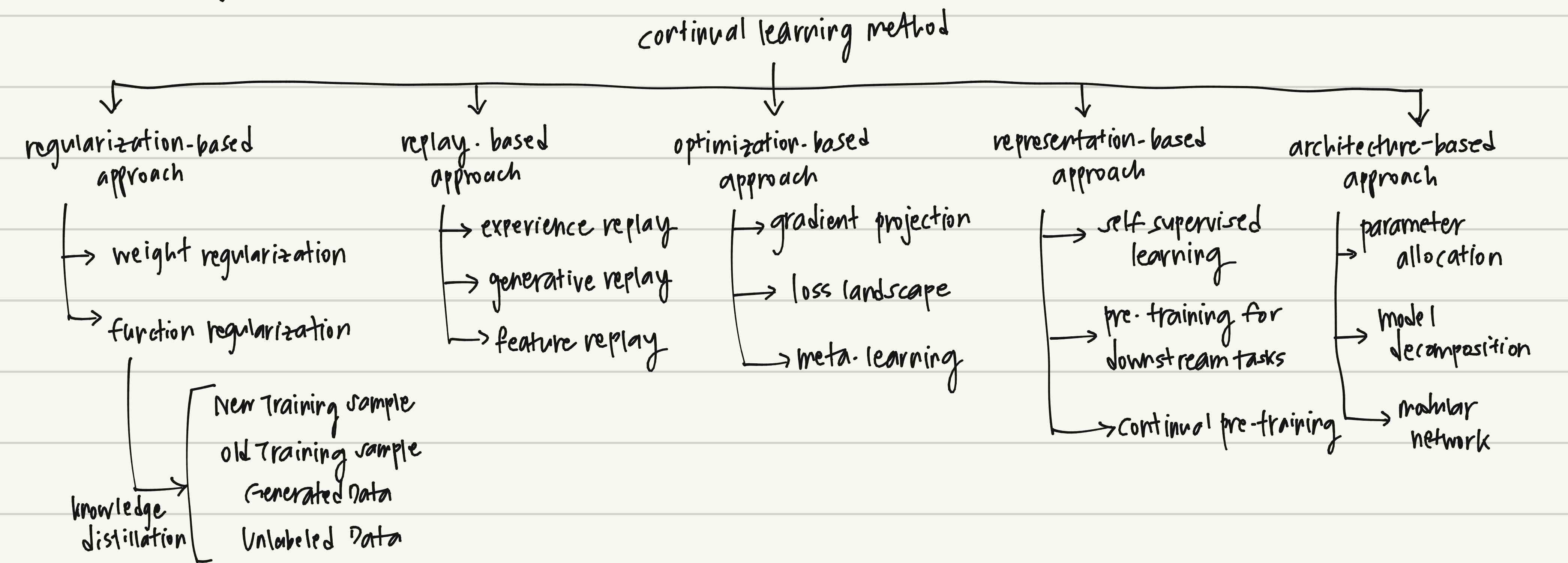


③ 새로운(New) task의 성능 평가: learning plasticity

- 새로운 task를 얼마나 빠르게 학습하는지에 대한 지표가 되는 learning plasticity는 일반적으로 instance measure (IM) 그리고 forward transfer (FT)로 측정.
- \* IM: 네트워크가 새로운 task를 잘 배우지 못하는 경우를 의미. joint training performance와 continual learning performance이 차이.
- joint training performance: 기존이 되는 모델은 k번째 task까지의 데이터를 jointly (함께) 학습시켰을 때의 성능.
- \* FT: 새로운 k번째 태스크를 학습할 때 이전에 학습된 태스크들이 새로운 태스크에 미치는 긍정적인 영향을 측정하는 방법. k번째 태스크를 학습한 후의 성능과 기존 성능 사이의 차이를 제한하여 측정.

Continual Learning Methods

- 각각 continual learning의 성능을 개선 방법은 k번째 task 각각에 대한 performance와, performance 변화 그리고 실제 continual learning 방식의 효과성에 대한 입증으로 구성됨.



Regularization based approach

- 일반적으로 생각할 수 있는 방법은 바로 old task와 new task 간의 균형을 위해 explicit한 regularization terms를 추가하는 것.
- 정규화가 consistency하게 불균라 하는 위치에 따라 weight regularization과 function regularization으로 구분됨.
- \* 정규화를 위해서 old task에 대해 학습된 frozen model을 reference로 가지고 있어야 한다는 점도 같음.

Weight Regularization

- 네트워크 파라미터의 변화를 선택적으로 규제하는 방법.
- 일반적인 구현은 loss function에 각 parameter의 변화에 대한 quadratic penalty를 추가한다면, 그에 따라 각 parameter의 중요도 (contribution / importance)에 따라 penalty를 주는 방법이 있음.
- 고차원의 뉴럴 네트워크를 다룰 때 ground truth가 되는 각 latent modality의 distribution을 직접 안 수렴하는 점이 real-world에 대한 근사적인 의미로 작용함. 하지만 만약 log-likelihood는 전체 dimension에 대해 가늠하지 않고 각 클래스 군사를 하게 되면, 적어도 해당 클래스의 local point에서의 분포 형태를 가늠한 분포를 구할 수 있음. 이것이 되면 각각 log-likelihood를 가늠한 군사 없으면 분포의 curvature 정보를 쉽게 구할 수, 이를 통해 중요한 parameter와 중요하지 않은 parameter를 구분할 수 있는 식도가 됨. (EWC | overcoming catastrophic forgetting in neural networks 논문)



- Online EWC (Progress & Compress: A scalable framework for continual learning)에서는 다수의 태스크를 학습하는 과정에서 각 태스크에 대한 공백도를 추적하고, 이를 통해 모델이 지속적으로 업데이트 될 때 중요한 가중치를 효과적으로 보존함. (task agnostic 하게 연속적으로 레이어 학습)
- \* EWC vs Online EWC: 기존 EWC는 소로 정서적으로 구분된 태스크들 사이의 연속 학습이 복잡함을 겪음.
- Synaptic Intelligence (Continual Learning Through Synaptic Intelligence): 각 파라미터 관위로 정의되는 공백도는 loss function이 레이어 parameter의 기울기를 구한 값과 catastrophic forgetting을 방지하기 위한 parameter 변화의 trajectory로 정해짐.
- 이러한 방법들은 모두 weight parameter에 대한 quadratic penalty term을 간접적으로 사용한다는 점에서 공통점을 가짐.
- penalty term을 사용하지 않고 factorized rotation을 기반으로 parameter space를 FIM(Fisher Information Matrix)에 diagonalize 하는 방법이나 (Rotate your Networks: Better Weight Consolidation and Less Catastrophic Forgetting 논문), Batch Normalization이 포함된 구조에 유사한 Kronecker-Factored Approximate Curvature (Continual Learning with Extended Kronecker-factored Approximate Curvature 논문)이 제안되기도 했음. 하지만 이런 방법들은 모두 "이전에 학습된 old task parameter"를 기준으로 한다는 점에서 parameter 변화를 막는다는 공통적인 constraint를 가지고, 이는 곧 새로운 task에 adaptation이 적용되는 과정에서 보충적인 효과를 가져옴.
- 이러한 문제점을 해결하고자 expansion 및 renormalization 방법이 제안되기도 했고, 이는 new task solution을 독립적으로 obtain한 뒤에, 이를 old model에 재배치하는 형태로 구현됨.
- IMM (Incremental Moment Matching / Overcoming Catastrophic Forgetting by Incremental Moment Matching 논문)이 초기 approach에 있는 논문인데, old task와 new task 간의 moment matching을 통해 점진적으로 새로운 task의 representation을 추가해가는 전략을 취함.
- \* moment matching: 학습된 볼록간의 유사성을 측정하고 조정하는데 사용됨.
- 이후 논문으로는 ResCL (IMM에 추가 combination coefficient를 학습 가능하게 바꾼 것 / Residual Continual Learning 논문), Online EWC가 있음.
- 추가 network에 학습된 task를 기존 network에 distillation 하는 과정인데, 이때 weight consolidation을 formulation 하는 과정에서 online EWC를 사용함.
- \* combination coefficient: 모델이 여러 태스크나 레이어들을 학습하면서 각각에 대한 강도를 개선하는데, 이러한 다양한 요소를 결합하거나 통합하는데 사용되는 계수.
- \* weight consolidation을 formulation 하는 것: 신경망의 파라미터를 조정하고 고정하여, 새로운 태스크를 학습하는 동안 이전 태스크들의 지식을 유지하는 것.
- AFEC (Active Forgetting of Negative Transfer in Continual Learning 논문)의 정세는 forgetting rate(시행자가 새로운 task를 받아들일 때 이전 정보와 conflict되는 정보를 active 하게 신경망)을 제한함. 이는 새로운 knowledge가 transfer 되는 과정에서 잠재적으로 기존 representation 과 공존할 때 발생하는 negative transfer 문제를 다뤄준다고 함.
- \* negative transfer: 새로운 task를 학습하는 과정에서 이전에 학습된 지식이 오히려 방해가 되어 학습 성능이 저하되는 현상.
- Towards Better Plasticity-Stability Trade-off in Incremental Learning - A Simple Linear Connector 논문: plasticity 및 stability 간의 trade-off (old task와 new task loss 사이의 trade-off) 균형을 위해 low-error path 간의 linear connector를 사용한 방법.
- \* low-error path: old task와 new task 각각에 대해 낮은 error를 보이는 신경망의 파라미터 집합.
- \* linear connector: 두 태스크 간의 가중치 공간에서 직접 경로를 형성하여, 이 경로를 따라 이동하면서 새로운 태스크를 학습하고 이전 태스크의 성능을 유지할 수 있음.



- Overcoming Catastrophic Forgetting by Neuron-level Plasticity Control 논문: parameter의 변화 자체를 극정하기 보라 learning rate를 줄이는 방법 제안.
- penalty를 통해 간접적으로 weight update를 막는 방법 대신 important neuron의 학습을 막음으로써 (freeze) hard regularization을 채택한 방법들을 찾아냄.

[function regularization]

- prediction function의 output 혹은 그 공간의 output을 기준으로 정제하는 간접적인 방법. 모델의 출력값이 일관성(consistency)을 유지해야 함.
- knowledge distillation 방법으로 훈련시키는데, teacher 모델은 기존의 지식을 바탕으로 일관된 출력을 제공하고, student 모델을 이를 참고해 새로운 데이터를 학습함.
- 대표적인 논문은 learning without forgetting.

Replay based approach

- 새로운 task를 학습할 때 old data distribution을 recover하고 approximate 하는 방법들도 생각해 볼 수 있는데, 각각 replay 하는 극제에 따라 차이가 있음.

[experience replay]

- 작은 memory buffer를 가지고 여기에 약간의 이전 task에 대한 training sample을 저장하는 형태로 구현됨.
  - \* reservoir sampling: 가장 간단하게 sample selection 과정마다 고정된 갯수만큼 old sample을 유지하고 나머지를 랜덤하게 대체하는 방식.
  - \* ring buffer sampling: reservoir sampling을 발전시켜서 old training sample을 새로 저장하는 과정마다 class마다 같은 샘플 수가 유지되도록.
  - \* mean-of-feature sampling: 각 class의 특징자 벡터의 평균을 일종의 prototype으로 간주하여 이미 가장 closest(유사한) 샘플을 저장하는 방식을 사용.
  - \* gradient based sample selection: gradient 기준으로 다양성을 최대화하는 방법.
  - \* Coresets via Bilevel Optimization for Continual Learning and Streaming 논문:
 

<u>entity</u> 간의 관계 조건인 ↓ 개미터의 샘플	<u>cardinality constraint</u> 를 back ↓ 샘플들 간의 관계 또는 상호작용을 나타내는 제약 조건
---	--

 performance와 연관짓는 방법
    - 샘플이 선정이 모델의 전반적인 성능, 특히 새로운 태스크를 학습할 때 이전 태스크의 정보를 잊거나 각 보존하는지에 대한 영향을 고려함.
  - \* batch 단위로 유사성을 보는 방법.
  - \* 최적화가 가능한 방법으로 latent decision boundary를 조정하는 등 고정되지 않은 인리듬으로서 정의된 approach가 발전하기 시작.
  - \* 동시에 샘플 재발 효과와 관련된 방법도 동시에 발전하기 시작했는데, vector quantize 방법을 포함하여 sample point를 압축하여 저장하기도 했는데, augmentation을 사용하여 한정된 샘플에 다양성을 보충하기도 함.
    - vector quantization: 연속적인 입력 레아터를 유한한 수의 벡터로 표현.
  - \* 의도적으로 forgetting boundary (잊지 쉬운 샘플들)을 adversarial하게 합성해 replay 효과를 높이는 방법으로 존재.

=> 굳이 샘플을 저장하지 않고 feature를 저장하면 안되는가?



[feature replay]

- experience replay에 비해 동량 측면에서 아쉽지 않으나. feature extractor가 업데이트되면서 같은 샘플에 대해서 representation이 달라지는 문제가 발생. 이를 곧 feature level에서의 replay가 catastrophic forgetting 위험성을 유발한다는 뜻.
- 몇가지 해결 방법:
  - \* old model과 new model 간의 feature distillation.
  - \* experience replay를 기준으로 feature representation (포장이나 포장재)이라 같은 용어를 쓰는 방법.
  - \* old sample: 1 representation을 저장해두고 이를 업데이트하면서 distribution shift를 예측하는 방법. (REP 논문)
  - \* 초기 layer를 고정시킨 채 중간 feature를 추출하여 이를 후반 layer를 업데이트 하는 다량의 replay sample도 사용하는 방법.

[generative replay]

- old task의 학습 데이터를 생성 모델을 사용하여 replay에 사용하는 전략.
  - \* ex) 새로운 task에 대해 생성된 데이터와 old task에 대해 생성된 데이터 간의 consistency를 사용하는 방법.

optimization based approach

- regularization 방법이나 replay는 결국 기호에 학습된 데이터를 implicit하게 혹은 explicit하게 활용하는 수 있는 방법.
- 이러한 "기호"라는 키워드에서 벗어나 explicit하게 학습한 구조를 바꾸거나 고정하는 방법이 소개되었으며, 이를 곧 optimization based approach라는 큰 틀로 묶을 수 있음.
- 이 접근법을 모델의 파라미터 업데이트 과정을 고정하여 새로운 태스크의 학습이 이전 태스크의 기억에 부정적인 영향을 최소화하도록 함.
- 가장 흔한 아이디어는 gradient projection. (새로운 태스크를 학습할 때 발생하는 그라디언트 업데이트가 이전 태스크의 학습 결과를 해치지 않도록 하는 것. 파라미터 업데이트 방향을 고정해 기존 태스크의 기억이 유지되도록 함)
  - \* replay-based approach의 변형: 기존 replay-based approach는 experience replay의 업데이트 방향을 기준으로 파라미터 업데이트를 고정.
  - \* OWM (Orthogonal Weight Modification): 파라미터 업데이트를 이전 input space에 orthogonal한 방향으로 수행. (기존 태스크의 입력 공간 유지)
  - \* OGD (Orthogonal Gradient Descent): 이전 태스크에서의 파라미터 업데이트를 보존하고, 새로운 태스크의 그라디언트 학습한 방향을 이전 태스크의 그라디언트와 수직인 방향으로 고정. (그라디언트 공간 유지 전략)
  - \* Bayesian weight regularization과 gradient projection의 결합: Bayesian 정제하는 과정이 가중치 불확실성을 고려하면서 데이터에 대해 일반화를 잘 할 수 있도록 돕고. 그라디언트 사영은 새로운 태스크의 학습이 이전 태스크의 성능에 부정적인 영향을 미치지 않도록 함.
- 여러 task에 대해 robust한 optimization이라고 하면 meta-learning을 떠올릴 수 있음.
  - \* meta learning: 학습 알고리즘 자체를 최적화. "학습하는 방법은 학습하는 것"
  - \* OML (Online Meta Learning): 메타러닝 학습 프레임워크를 제안함으로써 연속적으로 제시되는 샘플 input에 대해 학습 과정에서 얻은 지식을 online update에 도입하는 방법을 제안.
  - \* ANML (A Neuromodulated Meta-learning Algorithm): 점진적으로 공개하는 task에 대해, context-dependent 정보를 합체하여 특정 뉴런을 활성화하는 방식으로 학습 진행.



- generalization 관점 : loss landscape를 안정적으로 만드는 것. task 및 domain 관점에서 loss landscape가 안정적인 형태( curvature가 높음 . flat 한 local minima)를 가진수록 adaptation에 도움이 된다는 관점에서 중요함.

\* 이러한 loss landscape를 고려하는 방법들은 대부분 optimizer에 관한 연구로 구성되어 SGD나 Adam optimizer에 특화된 것임.

representation based approach

- generalization을 위한 meta learning, self-supervised learning, large scale pre-training과 같은 방법들 사용하여 continual learning에서 각 task에 특화하지 않고 넓은 범위의 한층 확장 가능 robust한 representation을 개발하는 전략.

- SSL을 포함한 방법들 :

- \* SSL approach를 사용한 논문들의 공통점은 "SSL로 학습된 representation"이 catastrophic forgetting에 강인하다는 점.
- \* LUMP: old task와 new task 간의 representation이 연속적이거나 가깝기에 interpolation을 진행하는 방식 사용.
- \* MindRed: 저장된 old training sample을 사용하여 replay experience를 기존 representation으로 복제 decorelate 시키는 방법 사용.
- \* self-supervised loss를 distillation 방법론과 결합 : 현재 representation을 이전 state로 mapping하거나 모델 간의 mapping을 하는 전략.

- pre-training을 포함한 방법들 :

- \* 기존 continual learning 과정에서 얻은 representation이 아키텍처는 사전 학습된 representation을 활용하여 여러 downstream continual setting에서 높은 성능을 나타내 하는 것이 목적.
- \* 레이블의 레이어셋에 대해 학습된 레전, 보라 큰 파라미터를 가지는 larger model에 레이어셋에 대해 학습된 레전. 아키텍처는 contrastive loss 등등 SSL 전략들과 함께 사용될 때 좋은 성과를 보임.

• contrastive loss : 모델이 레이어셋 간의 유사성을 학습하도록 하는 모션함수.

\* 이러한 방법론에서 공유한 것 사전 학습된 representation을 continual한 레전에서 어떻게 잘 관리하느냐.

① 사전 학습된 backbone 파라미터가 고정된 경우

- Side-Tuning : parameter-efficient tuning 전략. 기존 parameter와 parallel하게 학습될 수 있는 방법.
- Twf (Transfer without Forgetting) : backbone의 knowledge를 레이어셋으로 distillation 하는 전략.

② 사전 학습된 backbone 파라미터가 고정되지 않은 경우

- 생성 모델의 representation의 각 layer로부터 task-specific 파라미터를 학습하는 구조. (고정된 층도 있음, 안된 층도 있음)
- Adapter를 기반으로 사전 학습된 transformer를 tuning하는 방법.  
기존의 레전 사전 학습된 모델에 구조의 추가적인 레전이나 레이어를 삽입하여 특정 레전크에 대해 모델을 fine-tuning하는 방법.

architecture based approach

- 위의 방법들은 대체로 공유된 parameter (같은 model은 하나의 파라미터 단위로 생각했음)를 여러 task에서 잘 활용하는 방법에 대한 문제점.
- 하지만 매번 다른 파라미터를 여러 task에서 employ 한 경우 근본적으로 interference 문제를 해결하기 어렵다는 단점이 있음.
- 기존 방법들과 달리 task-specific parameter를 공유하는 parameter와 독립적으로 학습함으로써 간섭 문제를 해결하거나 한 방법이 바로 architecture based approach.

## [parameter allocation]

- 네트워크 수간선이 적어 각 task에 기어하는 parameter 부분집합 (subspace)를 분리. 이때 architecture는 고정되어있으므로 있고 dynamic하게 변경할 수 있음.
- 고정된 네트워크 구조를 활용한 경우 활성화하고자 하는 뉴런이나 파라미터를 각 task마다 분리해줄 수 있는 binary mask를 explicit하게 활용하는 형태.
  - \* binary mask: 네트워크의 각 파라미터에 대해 '0' 또는 '1'의 값을 가지는 마스크.

┐  
시용 X.

┐  
특정 레이어에 사용 0.
- 정확한 뉴런이나 파라미터 가치를 identify하는 방법들이 사용되기도 함. (iterative pruning, activation value, uncertainty estimation 등)
- 네트워크의 용량 한계로 인해 발생하는 성능 저하 문제를 해결하기 위해 네트워크 아키텍처를 dynamic하게 확장하는 방법들이 제안되었는데, 이는 특히 새로운 레이어를 효과적으로 추가하기 위해 기존 네트워크 파라미터의 용량이 부족한 데에 적합한 전략.

## [model decomposition]

- model을 task-sharing과 task-specific 성분으로 구분하는 approach.
- task-specific component는 확장 가능한 network를 가정하는 게 일반적임.
- task-specific components를 구성하는 구조:
  - \* ACL: parallel branches
  - \* GINCL: adaptive layer
  - \* 공간 feature map에 대한 mask를 생성하는 generator를 고안하는 방법.

## [modular network]

- parallel한 sub-network나 sub-module을 기존 네트워크 구조 안에 제안한 형태.
- progressive Network: 각 task마다 별도의 sub-network를 추가하는 방식 채택. adaptor connection으로부터 서로 다른 sub-network끼리 knowledge transfer를 수행하는 방법 제안.
- 여러 parallel branch를 통해 candidate path를 선택하고 최적의 경로를 선택하는 방법도 있음.
- 여러 sub-network 간의 앙상블을 사용하기도 함.