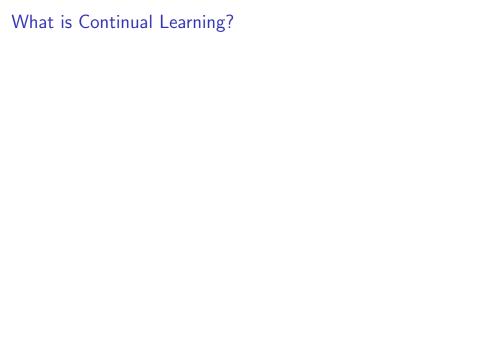
Continual Zero-Shot Learning

Ivan Skorokhodov

March 17, 2020



▶ Modern neural networks are prone to *catastrophic forgetting*: they forget previous tasks while are learning new ones.

- ▶ Modern neural networks are prone to *catastrophic forgetting*: they forget previous tasks while are learning new ones.
- ► Continual Learning tries to find ways to make model learn skills one by one in such a way that we do not forget previous skills

- ▶ Modern neural networks are prone to *catastrophic forgetting*: they forget previous tasks while are learning new ones.
- Continual Learning tries to find ways to make model learn skills one by one in such a way that we do not forget previous skills
 - Example 1: a robot that travels the world and learn new skills. We want it not to forget previous skills while he is acquiring new ones.

- ▶ Modern neural networks are prone to *catastrophic forgetting*: they forget previous tasks while are learning new ones.
- Continual Learning tries to find ways to make model learn skills one by one in such a way that we do not forget previous skills
 - Example 1: a robot that travels the world and learn new skills. We want it not to forget previous skills while he is acquiring new ones.
 - Example 2: a classification model is learning datasets one by one: we do not want its performance on previously learned datasets to decrease.

Modern CL techniques can be divided into three groups:

Modern CL techniques can be divided into three groups:

Regularization-based ([4], [1], etc): detect the weights which are important for previous tasks and do not change them much in the future.

Modern CL techniques can be divided into three groups:

- ▶ Regularization-based ([4], [1], etc): detect the weights which are important for previous tasks and do not change them much in the future.
- ▶ Rehearsal-based ([2], [6], etc): store a part of previous data to replay it in the future.

Modern CL techniques can be divided into three groups:

- Regularization-based ([4], [1], etc): detect the weights which are important for previous tasks and do not change them much in the future.
- ▶ Rehearsal-based ([2], [6], etc): store a part of previous data to replay it in the future.
- Component-based ([5], [3], etc): divide your network into components, and let future tasks not to break components which are important for previous tasks.

What data do we have:

▶ We will consider classification tasks from now on...

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in C$ we are given an *attribute vector* $a_c \in A$ which describes the class:

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in \mathcal{C}$ we are given an *attribute vector* $a_c \in \mathcal{A}$ which describes the class:
 - Imagine that we are classifying birds

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in \mathcal{C}$ we are given an *attribute vector* $a_c \in \mathcal{A}$ which describes the class:
 - Imagine that we are classifying birds
 - Then for each bird a_c includes bird's characteristics: color of a tail, body size, length of a beak, etc

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in C$ we are given an *attribute vector* $a_c \in A$ which describes the class:
 - Imagine that we are classifying birds
 - Then for each bird a_c includes bird's characteristics: color of a tail, body size, length of a beak, etc
- ▶ All classes are divided into seen and unseen:

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in \mathcal{C}$ we are given an *attribute vector* $a_c \in \mathcal{A}$ which describes the class:
 - Imagine that we are classifying birds
 - Then for each bird a_c includes bird's characteristics: color of a tail, body size, length of a beak, etc
- ▶ All classes are divided into *seen* and *unseen*:
 - Seen dataset: $D^s = \{X^s, Y^s, A^s\}$

- We will consider classification tasks from now on...
- ▶ For each class $c \in C$ we are given an *attribute vector* $a_c \in A$ which describes the class:
 - Imagine that we are classifying birds
 - Then for each bird a_c includes bird's characteristics: color of a tail, body size, length of a beak, etc
- ▶ All classes are divided into *seen* and *unseen*:
 - Seen dataset: $D^s = \{X^s, Y^s, A^s\}$
 - ▶ Unseen dataset: $D^u = \{X^u, Y^u, A^u\}$

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in C$ we are given an *attribute vector* $a_c \in A$ which describes the class:
 - Imagine that we are classifying birds
 - Then for each bird a_c includes bird's characteristics: color of a tail, body size, length of a beak, etc
- ▶ All classes are divided into *seen* and *unseen*:
 - Seen dataset: $D^s = \{X^s, Y^s, A^s\}$
 - ▶ Unseen dataset: $D^u = \{X^u, Y^u, A^u\}$
- ▶ During training we have an access only to seen dataset D^s .

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in \mathcal{C}$ we are given an *attribute vector* $a_c \in \mathcal{A}$ which describes the class:
 - Imagine that we are classifying birds
 - Then for each bird a_c includes bird's characteristics: color of a tail, body size, length of a beak, etc
- ▶ All classes are divided into *seen* and *unseen*:
 - Seen dataset: $D^s = \{X^s, Y^s, A^s\}$
 - ▶ Unseen dataset: $D^u = \{X^u, Y^u, A^u\}$
- \triangleright During training we have an access only to seen dataset D^s .
 - Our goal is to learn to match images with class descriptions

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in \mathcal{C}$ we are given an *attribute vector* $a_c \in \mathcal{A}$ which describes the class:
 - Imagine that we are classifying birds
 - Then for each bird a_c includes bird's characteristics: color of a tail, body size, length of a beak, etc
- ▶ All classes are divided into *seen* and *unseen*:
 - Seen dataset: $D^s = \{X^s, Y^s, A^s\}$
 - ▶ Unseen dataset: $D^u = \{X^u, Y^u, A^u\}$
- \triangleright During training we have an access only to seen dataset D^s .
 - Our goal is to learn to match images with class descriptions
 - I.e. model learns to detect "blue tails", "large heads", "short beaks", etc and not only concrete bird species

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in \mathcal{C}$ we are given an *attribute vector* $a_c \in \mathcal{A}$ which describes the class:
 - Imagine that we are classifying birds
 - Then for each bird a_c includes bird's characteristics: color of a tail, body size, length of a beak, etc
- ▶ All classes are divided into *seen* and *unseen*:
 - Seen dataset: $D^s = \{X^s, Y^s, A^s\}$
 - ▶ Unseen dataset: $D^u = \{X^u, Y^u, A^u\}$
- \triangleright During training we have an access only to seen dataset D^s .
 - Our goal is to learn to match images with class descriptions
 - I.e. model learns to detect "blue tails", "large heads", "short beaks", etc and not only concrete bird species
- \triangleright At test time we evaluate model performance on unseen dataset D^u

- ▶ We will consider classification tasks from now on...
- ▶ For each class $c \in C$ we are given an *attribute vector* $a_c \in A$ which describes the class:
 - Imagine that we are classifying birds
 - Then for each bird a_c includes bird's characteristics: color of a tail, body size, length of a beak, etc
- ▶ All classes are divided into *seen* and *unseen*:
 - Seen dataset: $D^s = \{X^s, Y^s, A^s\}$
 - ▶ Unseen dataset: $D^u = \{X^u, Y^u, A^u\}$
- \triangleright During training we have an access only to seen dataset D^s .
 - Our goal is to learn to match images with class descriptions
 - I.e. model learns to detect "blue tails", "large heads", "short beaks", etc and not only concrete bird species
- ightharpoonup At test time we evaluate model performance on unseen dataset D^u
- ▶ Using the knowledge about how inputs and attributes correspond to each other we can detect birds that we have not seen before just based on their class description a_c .

▶ Embedding-based: build two embedder models:

- ▶ Embedding-based: build two embedder models:
 - ▶ Model f(x) to embed images

- ▶ Embedding-based: build two embedder models:
 - ▶ Model f(x) to embed images
 - ▶ Model h(a) to embed attributes

- ▶ Embedding-based: build two embedder models:
 - ▶ Model f(x) to embed images
 - Model h(a) to embed attributes
 - ▶ Compute classification logits just by computing $d(f(x), h(a_c))$ for each class c (here d is some distance function.

- ► Embedding-based: build two embedder models:
 - ▶ Model f(x) to embed images
 - Model h(a) to embed attributes
 - ▶ Compute classification logits just by computing $d(f(x), h(a_c))$ for each class c (here d is some distance function.
 - I.e. we just measure distance between an image embedding and attribute embeddings

- ► Embedding-based: build two embedder models:
 - Model f(x) to embed images
 - Model h(a) to embed attributes
 - ▶ Compute classification logits just by computing $d(f(x), h(a_c))$ for each class c (here d is some distance function.
 - ▶ I.e. we just measure distance between an image embedding and attribute embeddings
 - ► Challenges: how to embed images properly, how to embed attributes properly, what distance function to use, etc

- ► Embedding-based: build two embedder models:
 - Model f(x) to embed images
 - Model h(a) to embed attributes
 - ▶ Compute classification logits just by computing $d(f(x), h(a_c))$ for each class c (here d is some distance function.
 - I.e. we just measure distance between an image embedding and attribute embeddings
 - Challenges: how to embed images properly, how to embed attributes properly, what distance function to use, etc
- Generative-based

- ► Embedding-based: build two embedder models:
 - \blacktriangleright Model f(x) to embed images
 - ▶ Model *h*(*a*) to embed attributes
 - ▶ Compute classification logits just by computing $d(f(x), h(a_c))$ for each class c (here d is some distance function.
 - ▶ I.e. we just measure distance between an image embedding and attribute embeddings
 - Challenges: how to embed images properly, how to embed attributes properly, what distance function to use, etc
- Generative-based
 - Train a conditional generative model that will learn to generate images based on class descriptions

- ► Embedding-based: build two embedder models:
 - \blacktriangleright Model f(x) to embed images
 - Model h(a) to embed attributes
 - ▶ Compute classification logits just by computing $d(f(x), h(a_c))$ for each class c (here d is some distance function.
 - I.e. we just measure distance between an image embedding and attribute embeddings
 - Challenges: how to embed images properly, how to embed attributes properly, what distance function to use, etc
- Generative-based
 - Train a conditional generative model that will learn to generate images based on class descriptions
 - ► I.e. GAN model that can generate a pigeon given description "grey feather, white head, short legs, etc"

- Embedding-based: build two embedder models:
 - ▶ Model f(x) to embed images
 - Model h(a) to embed attributes
 - ▶ Compute classification logits just by computing $d(f(x), h(a_c))$ for each class c (here d is some distance function.
 - I.e. we just measure distance between an image embedding and attribute embeddings
 - Challenges: how to embed images properly, how to embed attributes properly, what distance function to use, etc

Generative-based

- Train a conditional generative model that will learn to generate images based on class descriptions
- ► I.e. GAN model that can generate a pigeon given description "grey feather, white head, short legs, etc"
- At test time generate a lot of synthetic images, then train a classifier based on this synthetic dataset

- ► Embedding-based: build two embedder models:
 - ▶ Model f(x) to embed images
 - Model h(a) to embed attributes
 - ▶ Compute classification logits just by computing $d(f(x), h(a_c))$ for each class c (here d is some distance function.
 - I.e. we just measure distance between an image embedding and attribute embeddings
 - Challenges: how to embed images properly, how to embed attributes properly, what distance function to use, etc

Generative-based

- Train a conditional generative model that will learn to generate images based on class descriptions
- ► I.e. GAN model that can generate a pigeon given description "grey feather, white head, short legs, etc"
- At test time generate a lot of synthetic images, then train a classifier based on this synthetic dataset
- Challenges: how to train a good conditional generative model?

Modern ZSL techniques (for classification)

- ► Embedding-based: build two embedder models:
 - ▶ Model f(x) to embed images
 - Model h(a) to embed attributes
 - ▶ Compute classification logits just by computing $d(f(x), h(a_c))$ for each class c (here d is some distance function.
 - I.e. we just measure distance between an image embedding and attribute embeddings
 - Challenges: how to embed images properly, how to embed attributes properly, what distance function to use, etc

Generative-based

- Train a conditional generative model that will learn to generate images based on class descriptions
- I.e. GAN model that can generate a pigeon given description "grey feather, white head, short legs, etc"
- At test time generate a lot of synthetic images, then train a classifier based on this synthetic dataset
- Challenges: how to train a good conditional generative model?
- Currently performs better than embedding-based approaches

▶ Project: let's use class attributes to improve the performance on future tasks (and this also should improve past performance)

- ▶ Project: let's use class attributes to improve the performance on future tasks (and this also should improve past performance)
- ► Why?

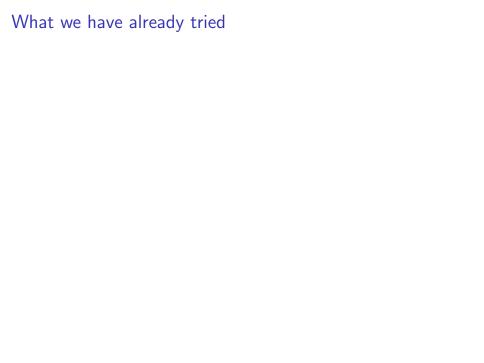
- ▶ Project: let's use class attributes to improve the performance on future tasks (and this also should improve past performance)
- ► Why?
 - A robot should be able to understand things it has never seen before but only heard about.

- Project: let's use class attributes to improve the performance on future tasks (and this also should improve past performance)
- ► Why?
 - A robot should be able to understand things it has never seen before but only heard about.
 - ▶ And it shouldn't forget previously seen objects while doing so...

- Project: let's use class attributes to improve the performance on future tasks (and this also should improve past performance)
- ► Why?
 - A robot should be able to understand things it has never seen before but only heard about.
 - ► And it shouldn't forget previously seen objects while doing so...
 - ▶ And the set of attribute descriptions can grow over time...

- Project: let's use class attributes to improve the performance on future tasks (and this also should improve past performance)
- ► Why?
 - A robot should be able to understand things it has never seen before but only heard about.
 - And it shouldn't forget previously seen objects while doing so...
 - And the set of attribute descriptions can grow over time...
 - ▶ In some sense, it is "the next" step of ZSL

- Project: let's use class attributes to improve the performance on future tasks (and this also should improve past performance)
- ► Why?
 - A robot should be able to understand things it has never seen before but only heard about.
 - ► And it shouldn't forget previously seen objects while doing so...
 - And the set of attribute descriptions can grow over time...
 - ▶ In some sense, it is "the next" step of ZSL
- I.e. build a zero-shot model that is trained in a continual learning fashion



► Attempt #1: use some simple (but working) generative-based ZSL approach

- ► Attempt #1: use some simple (but working) generative-based ZSL approach
 - Problem: generative models are hard to train in CL scenario (it is slow and they forget things fast)

- ► Attempt #1: use some simple (but working) generative-based ZSL approach
 - Problem: generative models are hard to train in CL scenario (it is slow and they forget things fast)
- Attempt #2: use some simple (but working) embeddings-based ZSL approach

- ► Attempt #1: use some simple (but working) generative-based ZSL approach
 - Problem: generative models are hard to train in CL scenario (it is slow and they forget things fast)
- Attempt #2: use some simple (but working) embeddings-based ZSL approach
 - Problem #1: It was already explored in previous literature (in A-GEM paper)

- ▶ Attempt #1: use some simple (but working) generative-based ZSL approach
 - Problem: generative models are hard to train in CL scenario (it is slow and they forget things fast)
- Attempt #2: use some simple (but working) embeddings-based ZSL approach
 - Problem #1: It was already explored in previous literature (in A-GEM paper)
 - ▶ Problem #2: Attributes do not have optimal "discriminative" power

References



Rahaf Aljundi et al. "Memory Aware Synapses: Learning what (not) to forget". In: CoRR abs/1711.09601 (2017).



Arslan Chaudhry et al. "Efficient Lifelong Learning with A-GEM". In: *International Conference on Learning Representations*. 2019.



Chrisantha Fernando et al. "PathNet: Evolution Channels Gradient Descent in Super Neural Networks". In: CoRR abs/1701.08734 (2017).



James Kirkpatrick et al. "Overcoming catastrophic forgetting in neural networks". In: *Proceedings of the National Academy of Sciences* 114.13 (2017), pp. 3521–3526.



Joan Serra et al. "Overcoming Catastrophic Forgetting with Hard Attention to the Task". In: *ICML*. Vol. 80. Proceedings of Machine Learning Research. PMLR, Oct. 2018, pp. 4548–4557.



Chenshen Wu et al. "Memory Replay GANs: Learning to Generate New Categories without Forgetting". In: *Advances in Neural Information Processing Systems* 31. Curran Associates, Inc., 2018, pp. 5962–5972.