# Continual Learning and Catastrophic Forgetting

by Paul Hand
Northeastern University
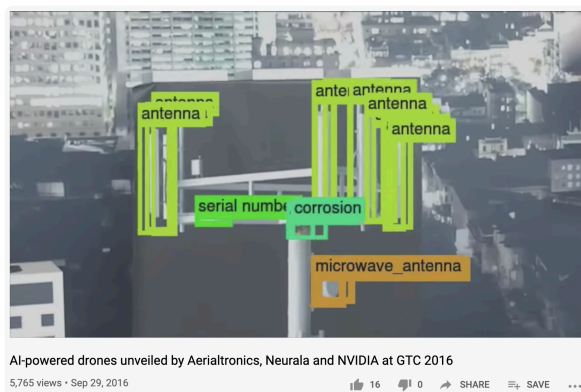
Outline:
    Context + initial approaches
    Evaluating algorithms
    Algorithms for CL

Example context for continual learning





AI-powered drones unveiled by Aerialtronics, Neurala and NVIDIA at GTC 2016
5,765 views · Sep 29, 2016     👍 16   👎 0   ➤ SHARE   ≡+ SAVE   ...

Other example: autonomous vehicles

**What are examples of situations where continual learning is desirable?**

Can you simply train on new data?

Task A: $D_A = \{(x_i, y_i)\}$        Task B: $D_B = \{(x_i, y_i)\}$

First,

$$\min_{\theta} \sum_{x_i, y_i \in D_A} L(\hat{y}_\theta(x_i), y_i) \quad \text{initialize randomly}$$

Then,

$$\min_{\theta} \sum_{x_i, y_i \in D_B} L(\hat{y}_\theta(x_i), y_i) \quad \begin{array}{l}\text{initialize w/ soln to}\\ \text{above task}\end{array}$$

Failure mode: catastrophic forgetting / interference

Typically, good performance at B
worse performance at A

**Visualization in parameter space of catastrophic forgetting**

**Demonstration of catastrophic forgetting using linear regression in 2-d**

$$D_A = \{ (\begin{pmatrix} 1 \\ 0 \end{pmatrix}), 1) \}$$

$$D_B = \{ (\begin{pmatrix} 1 \\ 1 \end{pmatrix}), 2) \}$$

$$y = (1 \ 1) \cdot x$$

$$\underset{\text{true response}}{|}$$

Model: $y = x^t \theta + \varepsilon$

Linear regression

$$\underset{\theta}{\min} \ \| y - x^t \theta \|^2 \qquad \text{inialized at } \theta = \theta_0$$

Soln $\theta = P_{y=x^t\theta} \theta_0$

How can you mitigate forgetting?

Train from scratch w/ new data and old data

Drawbacks

- Many industry nets take days - weeks to train
- Waste of power and compute
- Original training data may be unavailable

**When might original training data be unavailable?**

- Need access to GPUs/Cloud / steady internet

**When would cloud/gpu access be an issue?**

Humans can learn incrementally, so it is possible to do

Replay old training data w/ new data

- requires storing old data
- Storage costs grow linearly w/ tasks

Dilemma: plasticity - stability

Reviews: Parisi et al. 2019

Chen and Liu 2018

# Evaluating Continual Learning Algorithms

(Kemker et al. 2017)

## Data permutation tasks

img — class 0-9

Task 1: Training data ~ MNIST $(X_i, Y_i)$
Given $X_i$, predict $Y_i$

Task 2: Fix an image permutation $P_2$

Training data ~ MNIST $(P_2 X_i, Y_i)$
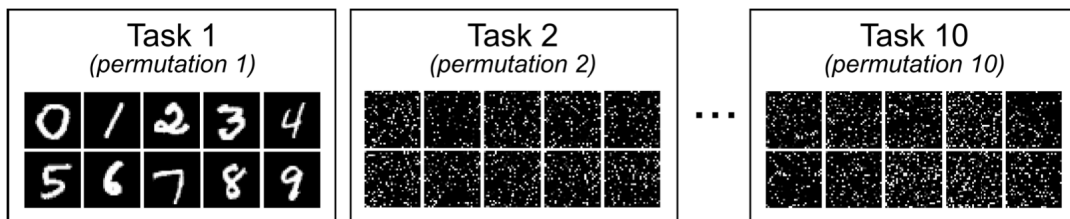Given $P_2 X_i$, predict $Y_i$



Figure 2: Schematic of permuted MNIST task protocol.
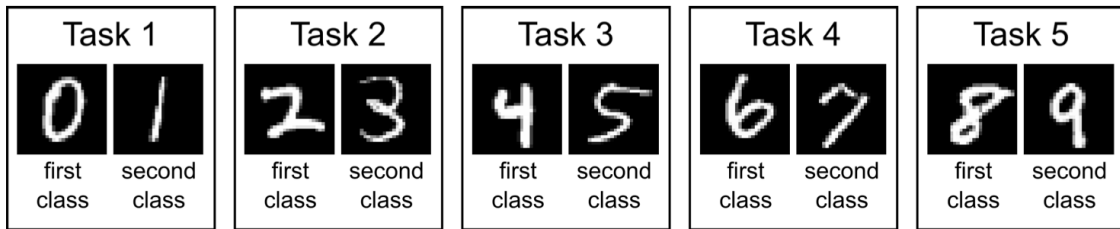
( Van de Ven and Tolias 2019)

## Comments:
Each task is equally difficult

**When/why are these tasks equally difficult?**

# Incremental Class Learning
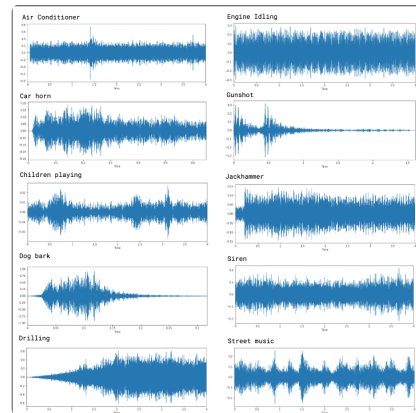
## Learn a base task set, then learn additional classes

## Shared features w/ new classes

**Are these tasks equally difficult?**

# Multimodal learning

Learn an image classification task
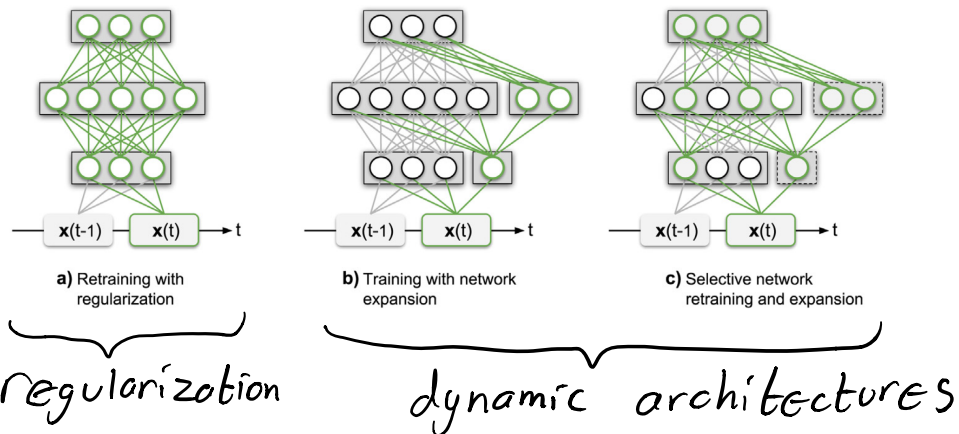then learn audio classification





Different features must be learned

\

# Approaches for Continual learning

- Train whole net w/ regularization
- Dynamic architectures (add neurons)
- Complementary Learning Systems (memory + replay)

(and more)

**a)** Retraining with regularization

**b)** Training with network expansion

**c)** Selective network retraining and expansion

regularization          dynamic architectures

# Regularization approaches

Update network weights but penalize changes in order to minimize forgetting

# Learning without Forgetting (LwF)

Consider predictor with shared parameters across tasks and some task specific parameters

At new task, update:

Shared params, new params, AND old params

So that

Output of old task on new data doesnt change too much.

LEARNINGWITHOUTFORGETTING:

Start with:
$\theta_s$: shared parameters
$\theta_o$: task specific parameters for each old task
$X_n, Y_n$: training data and ground truth on the new task

Initialize:
$Y_o \leftarrow \text{CNN}(X_n, \theta_s, \theta_o)$   // compute output of old tasks for new data
$\theta_n \leftarrow \text{RANDINIT}(|\theta_n|)$    // randomly initialize new parameters

Train:
Define $\hat{Y}_o \equiv \text{CNN}(X_n, \hat{\theta}_s, \hat{\theta}_o)$    // old task output
Define $\hat{Y}_n \equiv \text{CNN}(X_n, \hat{\theta}_s, \hat{\theta}_n)$    // new task output
$$\theta_s^*, \theta_o^*, \theta_n^* \leftarrow \underset{\hat{\theta}_s, \hat{\theta}_o, \hat{\theta}_n}{\arg\min} \left( \lambda_o \mathcal{L}_{old}(Y_o, \hat{Y}_o) + \mathcal{L}_{new}(Y_n, \hat{Y}_n) + \mathcal{R}(\hat{\theta}_s, \hat{\theta}_o, \hat{\theta}_n) \right)$$

parameter for plasticity-stability

modified cross-entropy loss

cross-entropy loss

weight decay

Note: does not require seeing old data!

**Why optimize the task specific parameters for the previous task instead of leaving them fixed?**
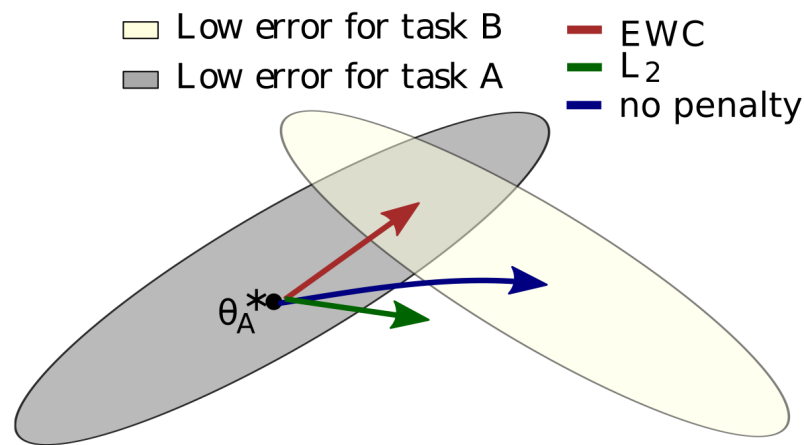
**Learning without forgetting is a "regularization" based method for continual learning.  Where is the regularization?**

# Elastic Weight Consolidation (EWC)

When training on task B, identify
weights that were important to A
and penalize updates to those weights

Try to stay in low error region for A



☐ Low error for task B    ▬ EWC
■ Low error for task A    ▬ L$_2$
                                   ▬ no penalty

$\theta_A^*$

Minimize

$$L(\theta) = L_B(\theta) + \sum_i \frac{\lambda}{2} F_i \left( \theta_i - \theta_{A,i}^* \right)^2$$

parameter

Solution to task A

where

$$F_i = \mathbb{E}_X \left( \partial_{\theta_i} \log p_\theta(x) \right)^2$$

diagonal entry of Fisher information
matrix of predictor at $\theta_A^*$

Motivated by a Bayesian learning perspective

**What is Bayesian learning?**

$(x_i, y_i) \sim$ dist of data iid

$D = \{(x_i, y_i)\}_{i=1\cdots n}$

model: $x \mapsto$ distribution on $y$    params $\theta$

$\max_{\theta} P(D|\theta)$    vs.    $\max_{\theta} P(\theta|D)$

MLE                                  MAP

frequentist                    max a posteriori estimaten

prior $P(\theta)$
collect $D$
update prior $P(\theta|D)$

Bayesian

Bayes Theorem:

$$\log P(\theta | D) = \log P(D|\theta) + \log P(\theta) - \log P(D)$$

If $P(\theta) \equiv$ const $\Rightarrow$ equivalent MLE $\Leftrightarrow$ MAP

    uninformative
    prior

Does there exist a prior $P(\theta)$ uninformative over $\mathbb{R}^d$?

Improper prior

**Example of Bayesian perspective on learning:**

MLE training of a NN — **MAP** Estimation
                           w/ noninformative improper prior
Weight Decay

$$\min_{\theta} \quad -\log P(D|\theta) + \|\theta\|_2^2$$

Bayesian perspective : prior $\log P(\theta) = -\|\theta\|_2^2$
$$P(\theta) \sim N(0, I)$$

# Derivation of Fisher Information

Dist $\quad P(X|\theta) \quad P_\theta(x)$

How sensitive is it to changes in $\theta$?

$$\nabla_\theta \log P(X|\theta) = 0 \quad \text{at a soln to training}$$
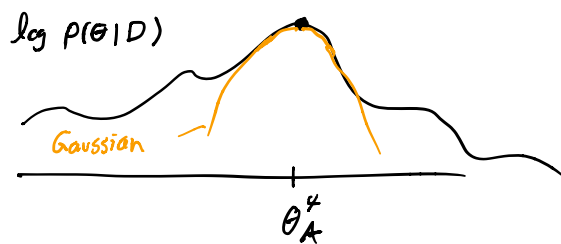
Instead, look at $\quad D_\theta^2 \log P(X|\theta)$

$$\mathbb{E}_{X \sim P_\theta} D_\theta^2 \log P(X|\theta) = \underbrace{\mathbb{E}_X \nabla \log P(X|\theta) \nabla \log P(X|\theta)^t}$$
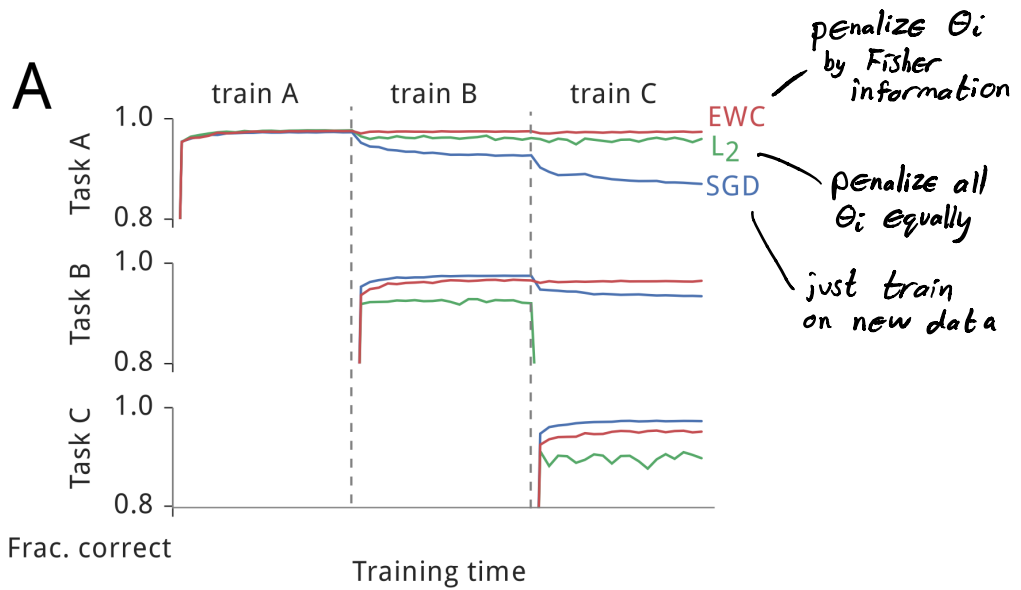
Fisher Information of $P_\theta$

large diagonal entries have large information content

Fisher info is Covariance$^{-1}$ of $\nabla \log P$

Fisher Information Matrix provides locally Gaussian approximation to a posterior distribution

$\log P(\theta)$

_____

_____

prior

$\log P(\theta|D)$



Gaussian

$\theta_A^*$

A



penalize $\theta_i$
by Fisher
information

penalize all
$\theta_i$ equally

just train
on new data

EWC
L2
SGD

Task A · Task B · Task C

1.0 · 0.8 · 1.0 · 0.8 · 1.0 · 0.8

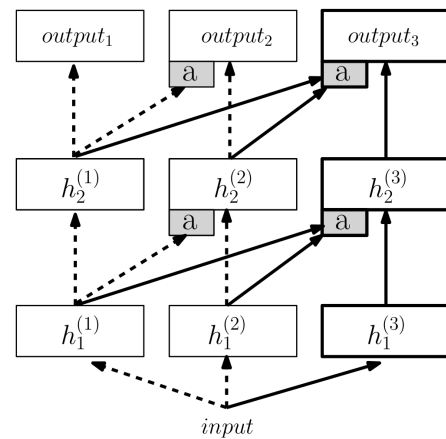train A · train B · train C

Frac. correct

Training time

# Progressive Neural Networks

(Rusu et al. 2016)

For each new task,
- add neurons
- add output layer
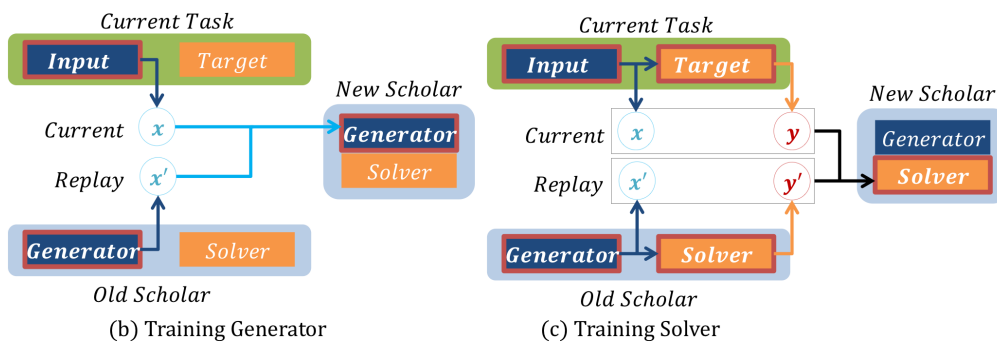- add lateral connections
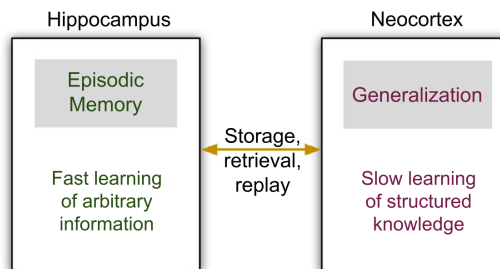- dont modify old weights

# Generative Replay

Train a generative model to output synthetic data that follows same distribution as training data.

Replay synthetic data along w/ new data



(b) Training Generator

(c) Training Solver

Takes inspiration from human learning

**b)** Complementary Learning Systems (CLS) theory

**Does generative replay avoid the data storage concerns that motivated continual learning methods?**

**Does generative replay avoid the data privacy concerns that motivated continual learning methods?**