

Today: Meta-learning  
Forgetting  
Generative Models: VAE

Announce: Fill out survey  
Extra Credit for everyone (3%)  
IF 75% of class does the survey

Recall Approaches To Adapt A Model to a new task:

0) IF something promptable, simply prompt it (Potentially using a prompt optimizer)

1) Use the pretrained model as an embedder / Feature extractor  
Called Linear Probing:

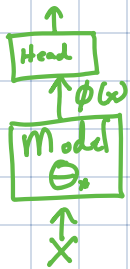


Train a new task-specific head for this task.  
e.g. linear classifier, regression

Advantage: Easy to do. No need for data beyond task.

Disadvantage: Linear Probing might not work as well.

2) Full fine-tune: Train everything — use pretrained model as initial condition for a part of the model.



Advantage: Typically higher performance on new task.

Disadvantage: Can be far bigger of a training job.  
Risk of overfitting.

Practical Tip: Don't just initialize a random head & then fine-tune.

Better: Initialize Head to 0.

Even Better: Use some data to train just the head first.

3) LoRAs or Soft-prompting combined with a new head...

Meta-learning: Making a model better at being fine-tuned for tasks  
(Think catigues 2&3 above)

What's a good baseline approach?

- 0) Do Nothing: Random Initialization
- 1) General Foundation Model
- 2) MAML: Model Agnostic Meta-Learning

What do we need?

- A) A collection of tasks from the family.  
i.e. Training Data for these different tasks  
+ Loss Function.
- B) Approach to finetuning.  
e.g. Use a LoRA and the SGD optimizer on training data.
- C) Approach to evaluation  
e.g. Eval performance on held-out set

Key Insight: In ML, default is train like you'll be tested.

Second Insight: Learning Process of SGD is like an RNN.

Let's be precise.—

A task  $i$  has Training Data  $D_i$ , Training Loss  $L_i$ , Test Loss  $\tilde{L}_i$

Want to train to do well on all tasks from this family, including held-out tasks.

Note: tasks may or may not share output cardinalities.

If they do, or have nested/overlapping outputs, we can consider learning a good initialization for the output heads too.

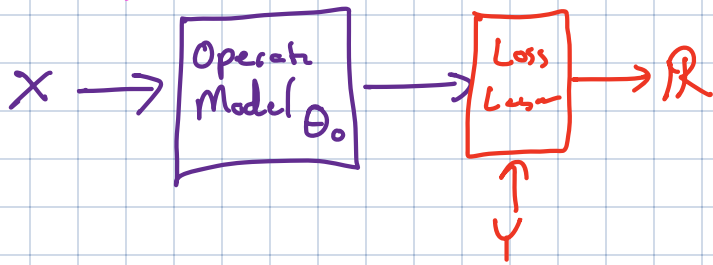
If they don't, we can use zero for the initialization of task-specific heads.

Assume we have a good initialization  $\Theta_0$ . How would we use it?

- 1) Start Model at  $\Theta_0$   $\leftarrow$  All <sup>learnable</sup> parameters in the model
- 2) Do gradient descent steps using  $D, L$  to get to  $\Theta_{\text{final}}$
- 3) Evaluate  $\Theta_{\text{final}}$  using held-out  $\tilde{D}$  and Test Loss  $\tilde{L}$  to get  $h_{\text{test}}$

(1,2,3)

Fit this into our standard form:



What is  $X$ ?  
 $Y$ ?

opert  
mult?

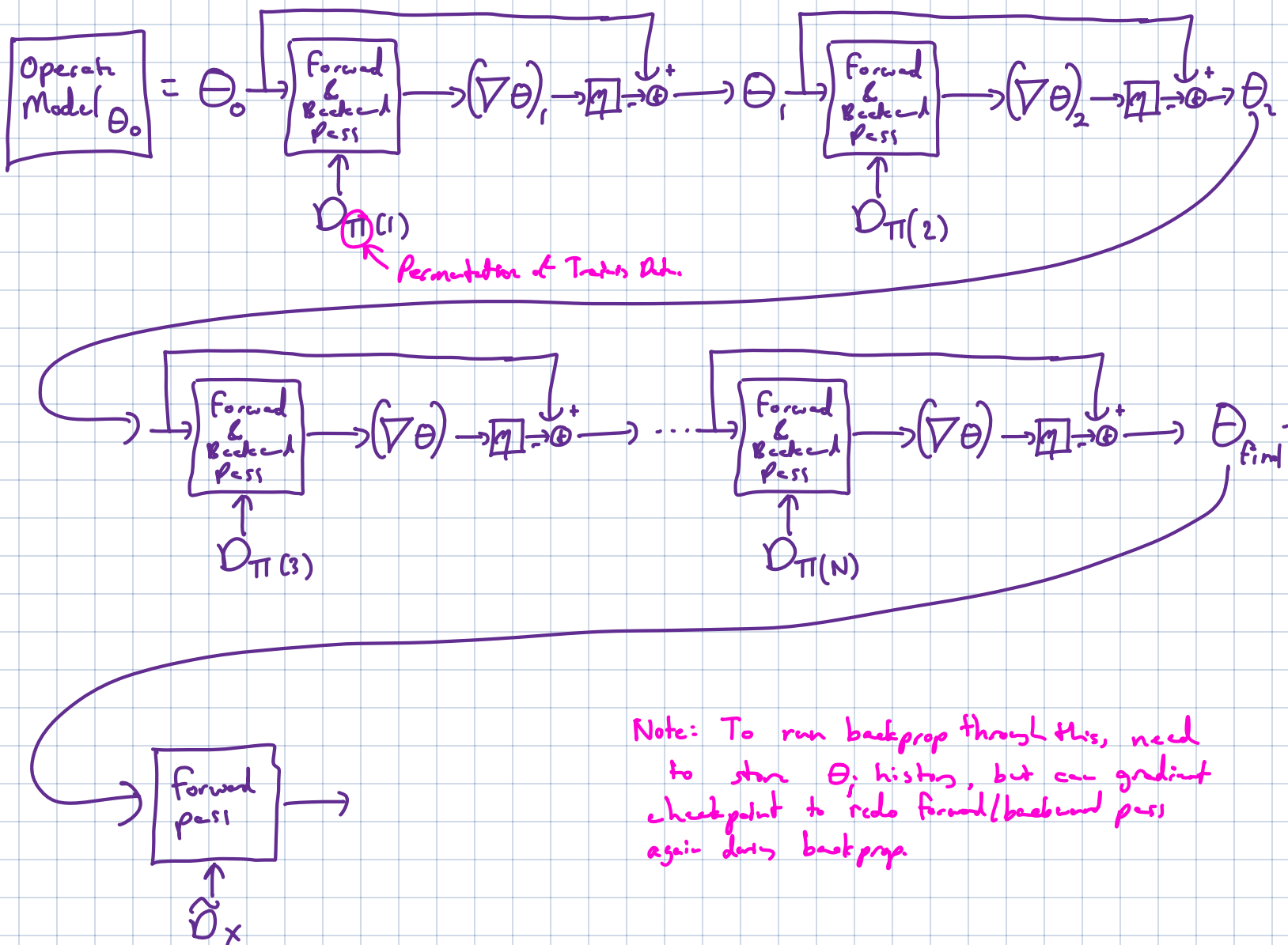
Loss Layer?

Just Pattern Match:

$$X \leftrightarrow (D, \tilde{D}_x)$$

$$Y \leftrightarrow \text{Label info in } \tilde{D}_y$$

$$\text{Loss Layer} \leftrightarrow \tilde{L}$$



MAML insight: That's just a deep model. Can backprop through it

Gives a gradient  $\nabla \Theta_0$  on initial condition.

Can take a step in that direction for a better initial condition

Repeat for a new task (randomly drawn).....

Hope: This gives us a much better initial condition.

Schematic Picture: Outer-loop on tasks, learning-rate  $\eta_{out}$

Inner-loop on batches/examples, learning-rate  $\eta_{inner}$

Inner Iterations  $\rightarrow$

Task 1:  $\Theta_{0,0} \rightarrow \Theta_{1,0} \xrightarrow{\text{Backprop}} \Theta_{1,1} \rightarrow \Theta_{1,2} \rightarrow \dots \rightarrow \Theta_{1,N_1} \rightarrow \boxed{L_1} \rightarrow \mathbb{R}$   
           $\downarrow$  *Applies outer step*  
2:  $\Theta_{0,1} \rightarrow \Theta_{2,0} \rightarrow \Theta_{2,1} \rightarrow \dots \rightarrow \Theta_{2,N_2} \rightarrow \boxed{L_2} \rightarrow \mathbb{R}$

3:  $\vdots$

M:  $\Theta_{0,M} \rightarrow \Theta_{M,0} \rightarrow \dots$

$\Theta_{M,N_M} \rightarrow \boxed{L_M} \rightarrow \mathbb{R}$

$\downarrow$   
 $\Theta_{0,final}$

outer-loop  $\rightarrow$

## Two Variants:

- 1) Reptile: Avoid backprop through backprop  
Approximate  $\nabla \Theta \approx \Theta_{\text{ANIL}} - \Theta$
- 2) ANIL / Meta Opt Net / RZDZ: Optimize for Linear Probs



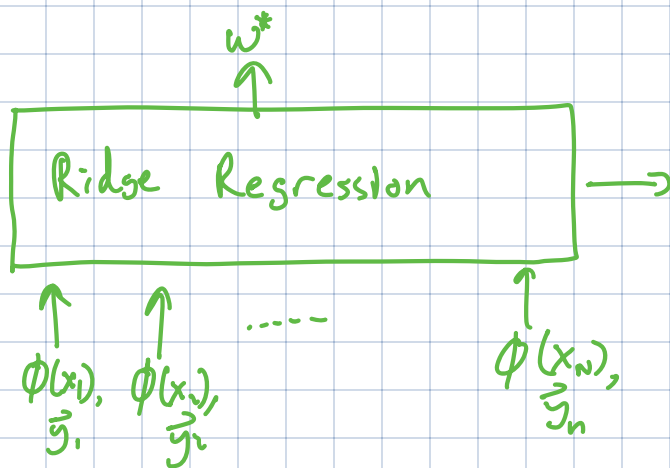
e.g. consider regression problems.

Have closed form formulas for Head.  
(Instead of Gradient updates...)

Take a gradient step through that to the parameters of feature extractor

## Explicitly:

(Added dummy office hours)



$$\text{Let } \Phi = \begin{bmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix}$$

$$Y = \begin{bmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_n^T \end{bmatrix}$$

$$W^* = \underbrace{\left[ \Phi^T \Phi + \lambda I \right]^{-1} \Phi^T Y}_{\text{Differentiable}}$$

$$W^* \phi(x_{\text{test}}) = \hat{y}_{\text{test}}^T$$

↑  
output output point

$\Rightarrow$  Sends gradients to  $\Phi$  and hence to  $\phi(x_1), \dots, \phi(x_n)$  during backprop.

Note: Convex Problem (even not in closed form) can send gradients via natural iteration solutions. e.g. A few Newton Steps

So loss  $\tilde{L}$  on  $\hat{y}_{\text{train}}$  relative to  $y_{\text{test}}$  sends gradients to both  $w^*$  and  $\phi(x_{\text{test}})$

We end up with parallel gradients flowing back thru weighted  $\phi(x_{\text{test}}), \phi(x_1), \dots, \phi(x_n)$

These can be used to get gradients on  $\Theta \rightarrow$  param definitions  $\phi$

$\Rightarrow$  Updates on  $\Theta$  for meta-learning

Note: Can use gradient checkpointing ideas to avoid having to store all intermediate activations. Do forward pass twice and backward once.

Catastrophic Forgetting: When fine-tuning, model forgets how to do what it knew how to do.

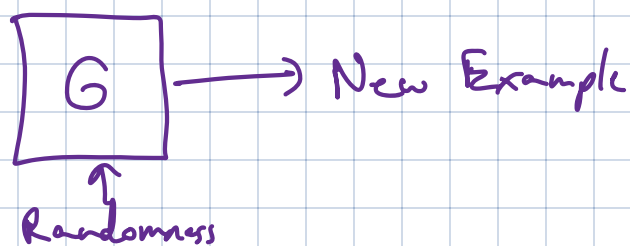
Key Practical Solution: During fine-tuning, mix in some pretraining-style data.  
Like 10%.

Added in Office Hours:

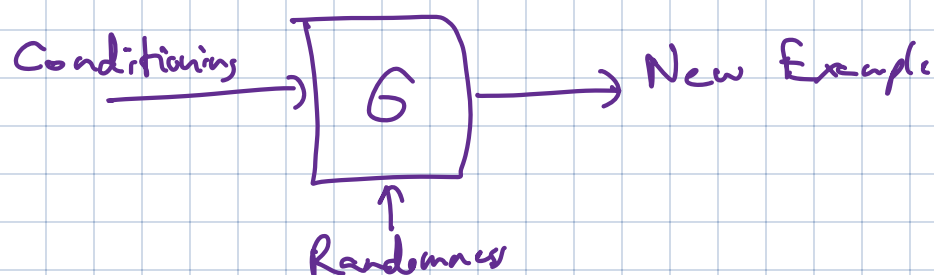
- A) If the pretraining-task(s) had distinct heads, keep them during fine-tuning — they will (along with their losses) send gradients into the shared part of the model to prevent/reduce forgetting. These heads should update — they are not frozen.
- B) Forgetting can be thought of as a type of overfitting.

# Generative Models

## Unconditional

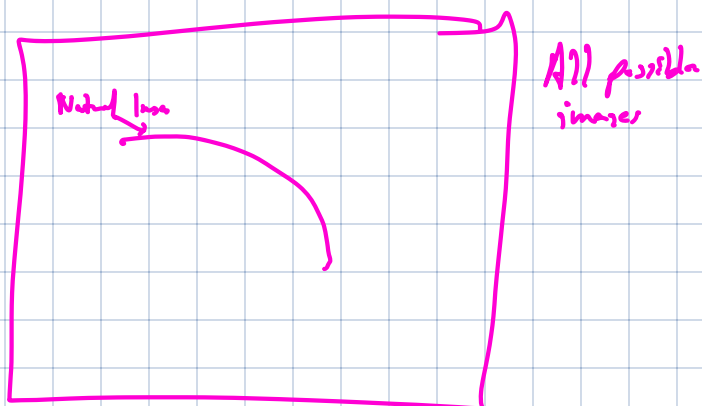


## Conditional



## Ideas that don't work

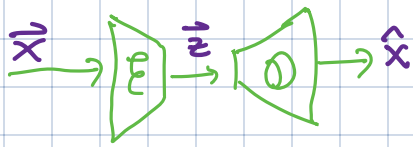
A) Use a classifier



## B) Use an autoencoder

Core Ingredients: Labels are  $\vec{x}_i$  itself.

Architecture has an encoder followed by a decoder  
Bottleneck in the middle.



Traditional Perspective: Decoder is scaffolding.

Try using  $\mathcal{D}$  to generate samples. — IF  $\vec{z}$  too small, get blurry junk  
IF  $\vec{z}$  big, "his"

## VAE Approach

- 3 key ingredients:
- 1) Make  $\vec{z}$  random during training
  - 2) Add a loss on distribution of  $\vec{z}$
  - 3) Make this work with SGD