

Small or noisy data: forms of weak supervision

Academic setting:

- Large dataset
- Focus on architectures / training options

Real world / data (amount)

I Small data

- Best case: simulator = train on simulated data, then fine-tune on real data

- data augmentation:

- vision:
 - rotate, translate, crop, flip
 - spatial deformations



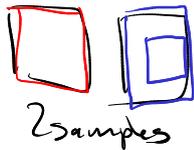
- pixel intensities:
 - Contrast
 - color balance
 - noise: suited to the task



Random Transform



- CutMix



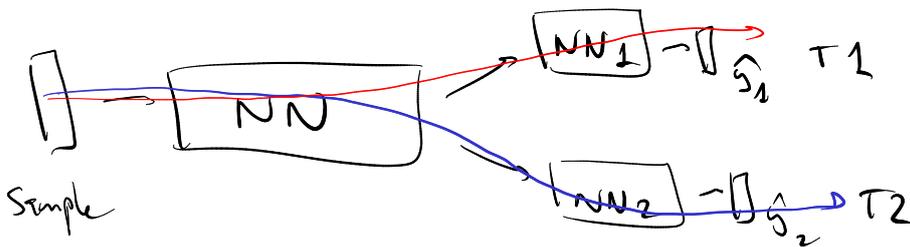
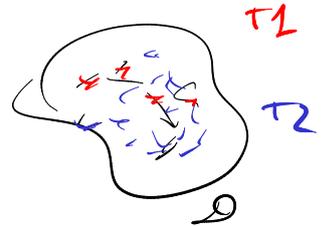
random location & size

- multi-tasking

Task 1: small data

Task 2: lot of data

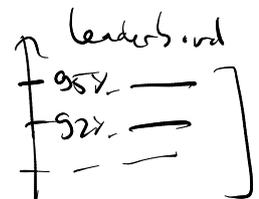
} same data → few labels
→ many labels



⚠ beware = similarity of tasks

Transfer learning

- learning from scratch: building architecture ... ≈ 6 months
- build on previous models:
 - huggingface
 - benchmarks / challenges



1) copy the architecture



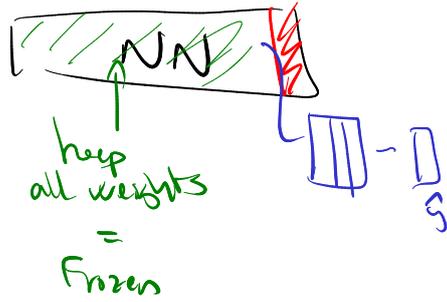
retrain from scratch

or a few hours/days

+ copy the full processing/optim² chain:

- pre-processing (data augm²)
- drop-out
- early stopping
- learning rate scheduler

or 2) copy the weights: "pre-trained" model

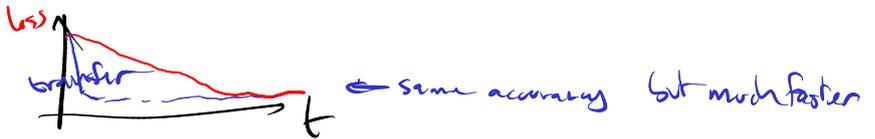


- a few seconds/minutes

or 3) after training the new layers, unfreeze the pretrained model & fine-tune
→ hours

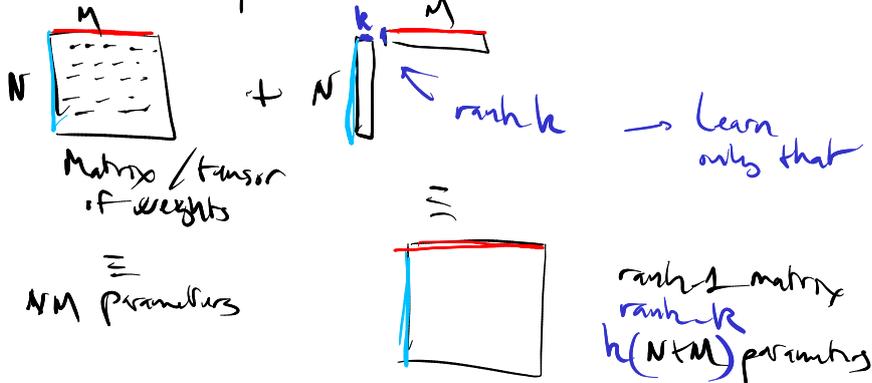
Effect of transfer learning

- if small data: \hookrightarrow help in getting relevant features
- if not small data: \hookrightarrow \hookrightarrow some accuracy but much faster

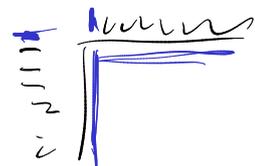
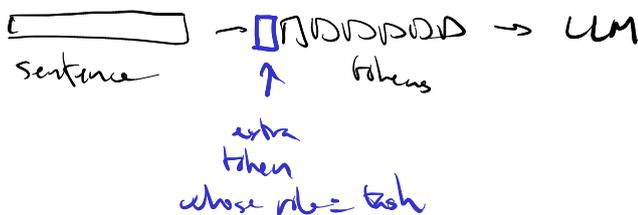


Other ways to fine-tune: (more specific to ULM)

- LoRA: Low-Rank Adaptation



- add a F₀ token



- prompting
- RAG

Remark about transfer / Foundation models

- tend to transfer from large, generic models
- end up with huge model
- ⇒ reduce models

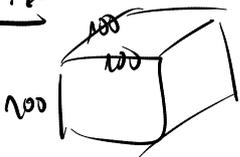
→ pruning: weights
neurons
layers

↳ group-Lasso
 L_1

how: - magnitude of weights
- " of weight gradient
- variation of weight over iterations

then: Fine-tune
usually: iterative \rightarrow Prune \rightarrow Fine-tune

→ tensorize



100 100 100 parameters → low-rank tensor?

→ CP decomposition



rank-1 tensor

→ Tucker:

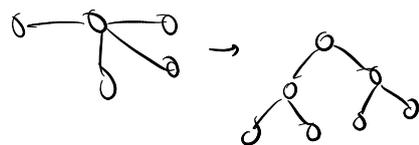


core

→ Tensor train sequence of multiplications of tensors



→ combine these approaches + lots of variations
↳ eg: Hierarchical Tucker



→ quantize:

$w = 2.551783 \dots$

Float: 32 vs 16 vs ...

≈ 2.6

extreme case: Sparse networks

$-1, 0, +1$

list of possible weights: $\{-2.5, -0.2, 0.15, \dots\}$

eg. 5 values

↳ during trainings

cluster the weights and push towards the centers

→ "knowledge distillation": "teacher-student"

Task: 



 initiate the outputs of NN

\neq :
↳ Feedback
as many Feedbacks as # classes

- Lottery ticket ~~AD~~:



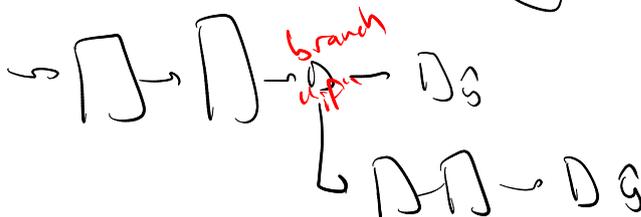
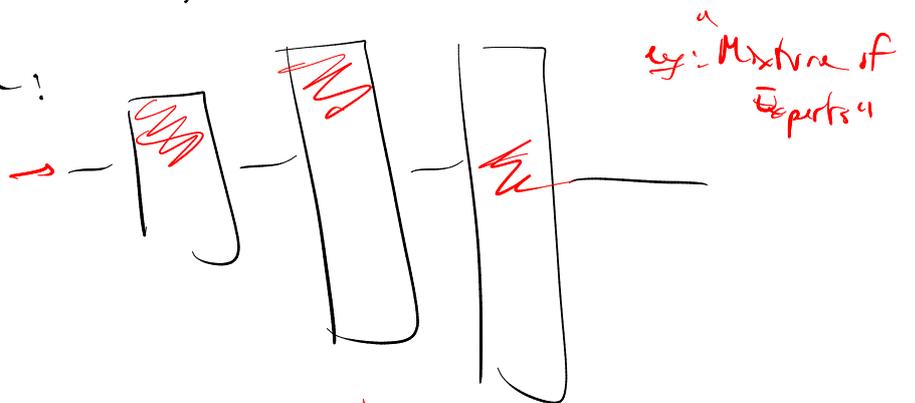
a subnetwork not even trained might do the job

eye of small models:

- MobileNet
- SqueezeNet
- EfficientNet

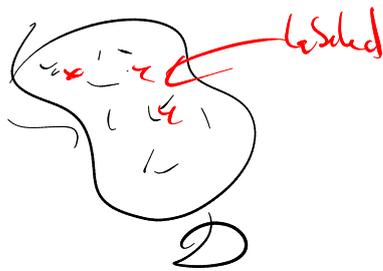
"tiny ML"

↳ dynamic architecture:



- Neural Architecture Growth for Frugal AI

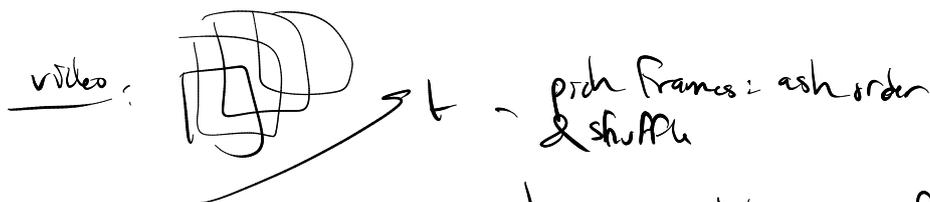
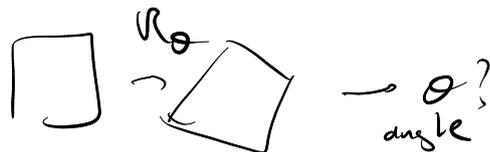
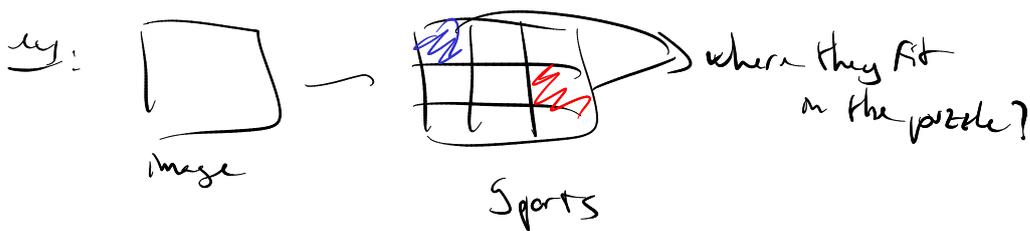
II Few labeled ex: weak supervision



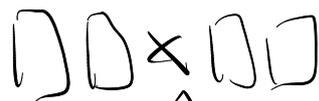
- 1) unsupervised pre-training (on all data)
then fine-tune on real task with labeled samples
- 2) train supervisedly
then label a few more samples (the most reliable ones)
 & re-train \uparrow mistake
- 3) train supervisedly
but check global properties on the whole dataset
 ex. bias, density --
 ↓ classes classif. pb
 20%, 80%,
 A B
 → calibration

"Self-supervision"

- create "dummy" task \uparrow supervised with labels automatically generated



- harder: predict a missing frame

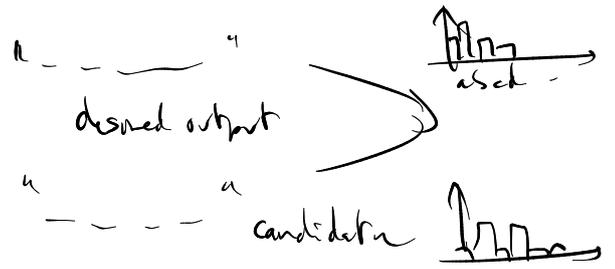


- hardest: predict next frame

→ useful in NLP

eg: translation

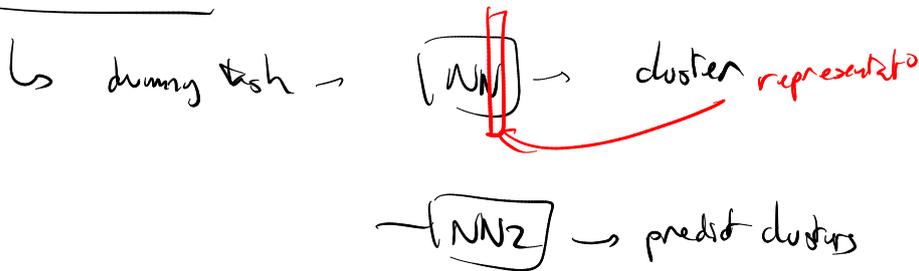
- a) not many ex of translation
- b) good metric? → BLEU / BOWE



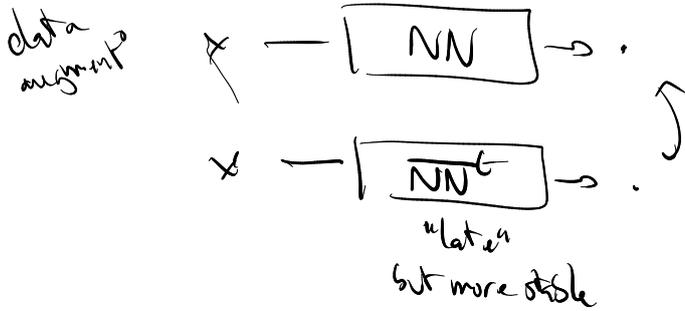
Self-supervised:

Time is $\frac{\emptyset}{\dots}$
 missing word

ClusterFit / Dino:



↳ Dino:



Sample x

$x^c = \text{augmented version}(x)$

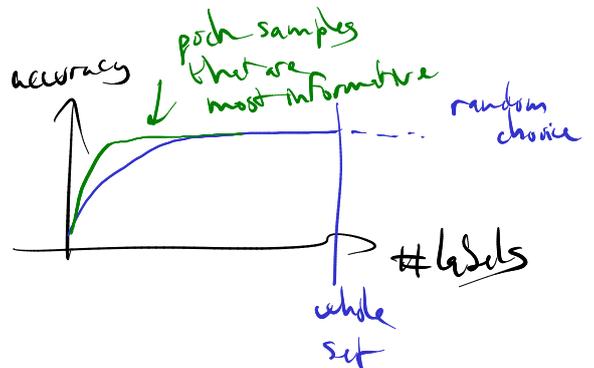
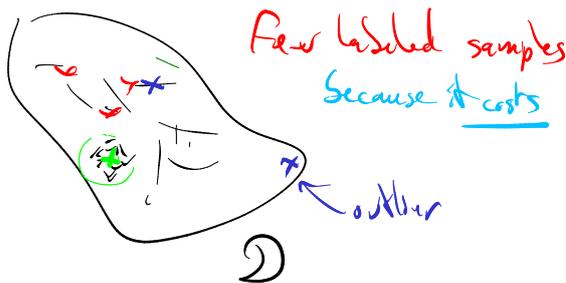
Loss:

$$\|g(x) - \hat{g}(x^c)\|$$

Learn Features that are invariant to transform group

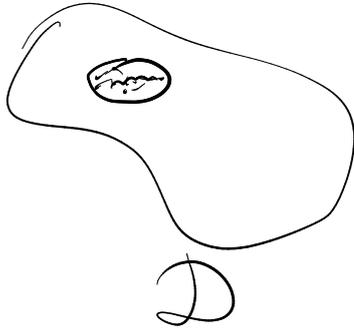
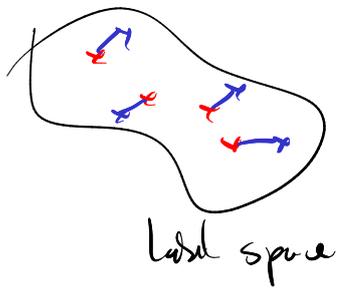
Ask for variance on outputs

Active learning:



Noisy data:

potentially big data, but labels are noisy



Same Sample regression

real label

if noise centered:

$$\min_{\mathbf{g}} \sum_i \| \hat{\mathbf{g}} - \mathbf{y}_i \|^2$$

argmin \mathbf{g}

noisy labels

$$\hat{\mathbf{g}}^* = \frac{1}{N} \sum_i \mathbf{y}_i$$

$$\| \hat{\mathbf{g}}^* - \mathbf{g}_{\text{real}} \|^2 \sim \frac{1}{N} \sigma^2$$

noise variance

classified

ImageNet: change labels

10% → random labels

20%

50%

99% →

x100 frames set size ⇒ in order to keep # rightly-labeled $\frac{1}{N}$: