

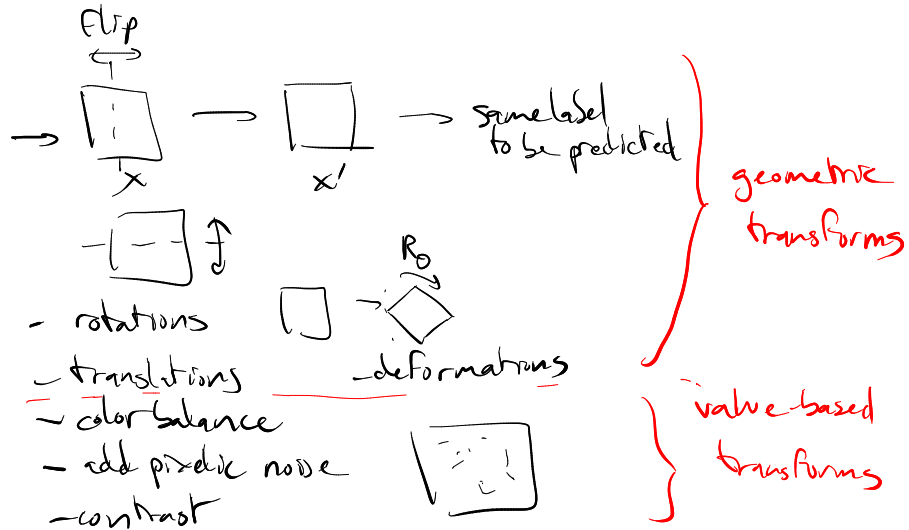
# Forms of weak supervision

## I Small data

### Data augmentation

ex: image classification task  
add transformations

Label:  $\swarrow$  invariant  $\searrow$  equivariant



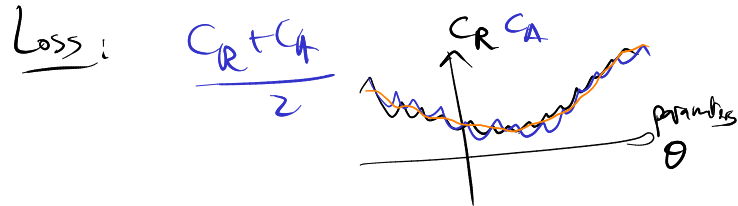
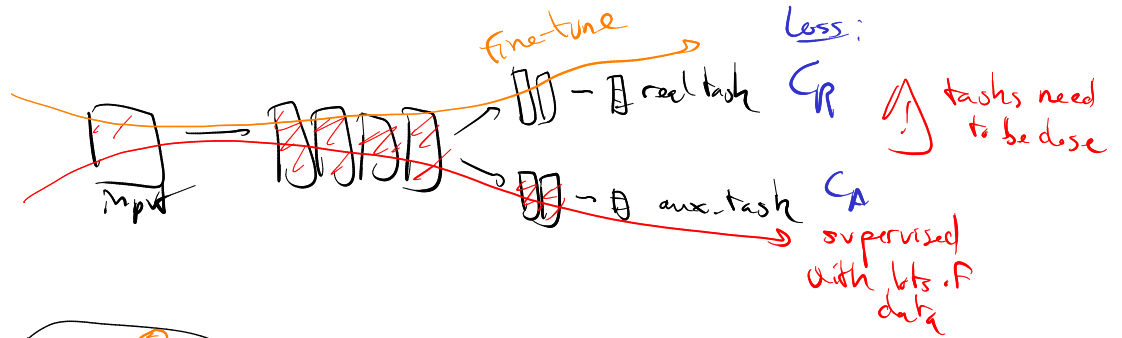
- use a simulator
- generate lots of data
- how realistic?

## Multi-tasking

- one real task + one auxiliary task

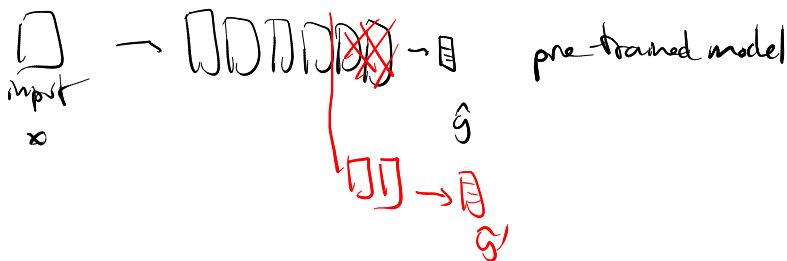
real task: few labeled data

auxiliary task: lots of labeled data

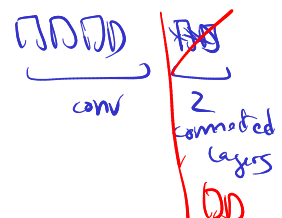


## Transfer learning

- sequential training: first on auxiliary task  $\leftarrow$  pick a pre-trained model  
then on real task



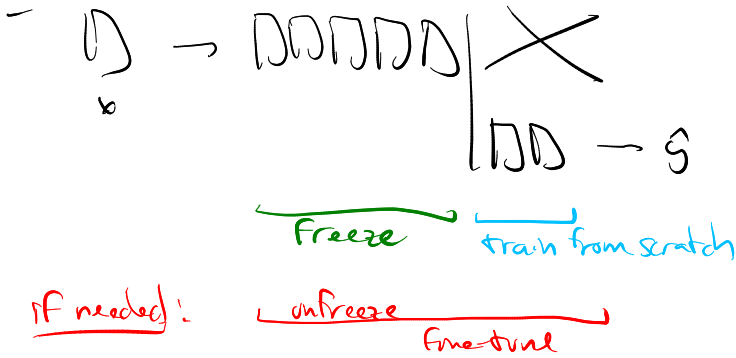
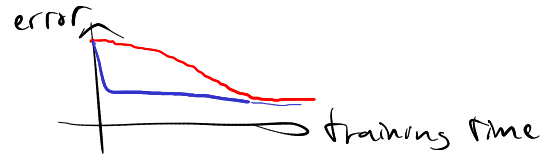
conv: ResNet / VGG  
ImageNet



- analysis from [Rethinking ImageNet pre-training]

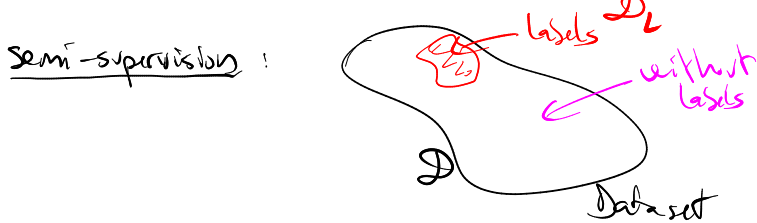
small data:  
helps in getting good features

big data:  
≠ from training from scratch,  
a big boost in training time



II Forms of weak supervision

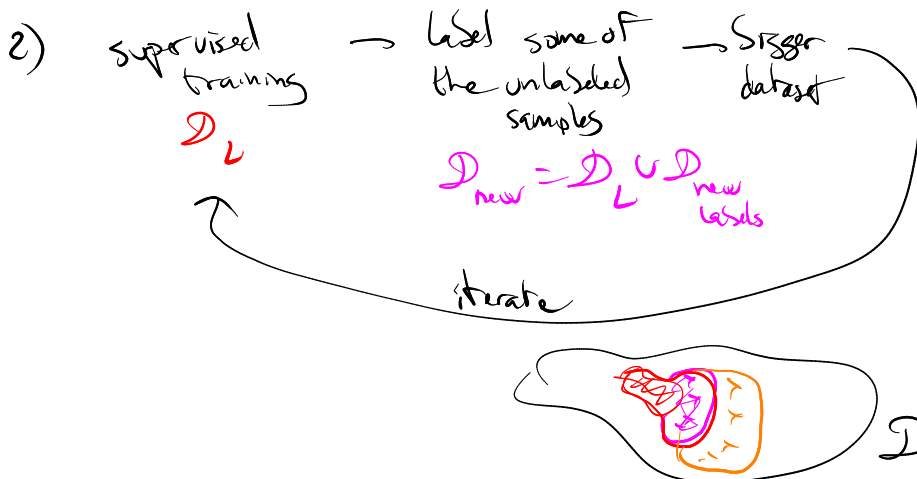
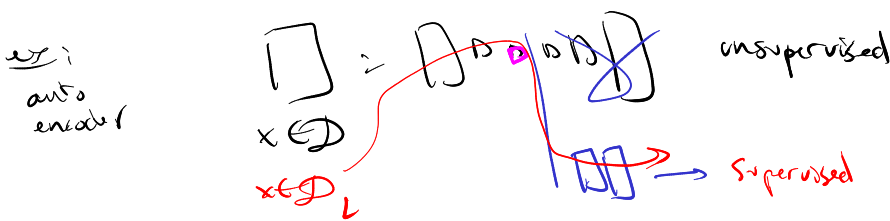
Amount of data, but few are labeled



ex: when labeling is costly  
(requires time, expertise, ...)

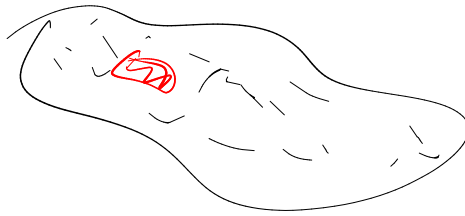
Several approaches:

1) unsupervised training on full dataset  $D$  → good representation → supervised task on  $D_L$



Issue: what if mistakes?

3) supervised training  $\rightarrow$  apply to full dataset  $\rightarrow$  check some properties & adjust parameters



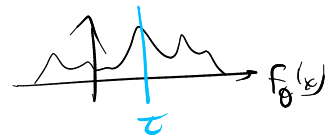
ex: global bias

target: A / B  
50% / 50%,  
80% / 70%

adjust the decision threshold

$$f_{\theta}(x) < 0$$

$$f_{\theta}(x) < \tau$$



properties:  
- bias  
- density  
- margin (SVM)

### Weak supervision

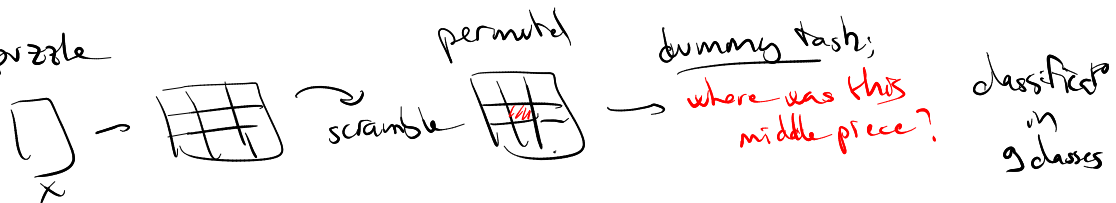
$\hookrightarrow$  more general:  
ex: labels could be noisy

### Self-supervision

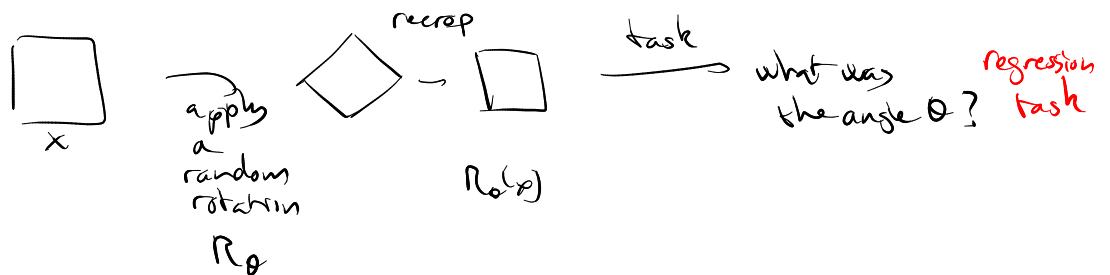
$\hookrightarrow$  used for pre-training  
 $\hookrightarrow$  no supervision (no label given by hand)  
 $\hookrightarrow$  unsupervised task formulated as a supervised one with automatic labels

ex: image classification

- image puzzle



- image rotations



ex: video classification

- predict next frame  
- give 3 frames:  $\square \square \square$  ask temporal order

# Building on teacher-student techniques

- "ClusterFit".



train a new network:  $\boxed{x} \rightarrow \text{DNDNDND} \rightarrow \hat{y}''$  task: predict this label, i.e. which cluster  $x$  belongs to student

- DINO: student  $\boxed{x} \rightarrow \text{DNDNDND} \rightarrow$  randomly initialized

teacher  $\boxed{x} \rightarrow \text{D} - \dots - \text{D}$  average of the past students (recent history) moving average

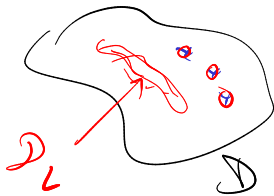
data augmentation  $\rightarrow$  learning to be invariant to the chosen group of transforms

+ something against collapse

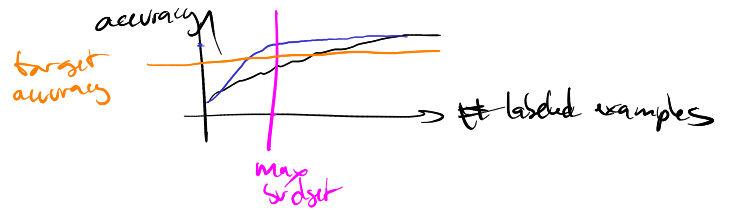
$\hookrightarrow$  just train a linear classifier on top  $\Rightarrow$  almost as good as best supervised technique ever on ImageNet

## Active learning

same setting as semi-supervision + ask some samples to be labeled  $\hookrightarrow$  costly labeling



goal: increase global accuracy as fast as possible (in terms of # of labeled samples)



- large dataset:  $\mathcal{D} = \{x_i\}$

- labels for now:  $\mathcal{D}_L = \{y_1, \dots, y_p\}$  with  $p \ll |\mathcal{D}|$

$\leftarrow$  which  $x_i \in \mathcal{D}$  to pick? ( $i \in ]p, |\mathcal{D}|$ ) to ask to be labeled

\* apply the current model  $f_\theta$  to all samples  $\rightarrow$  predictions  $\hat{y}_i = (y_i^c)_{c \in \mathcal{C}}$  if classification task   
  $\mathcal{C}$  classes

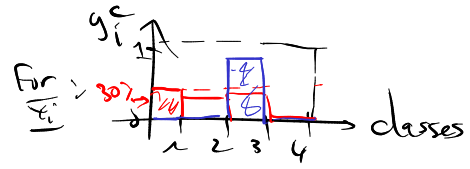
## Local methods

$\rightarrow$  quantify the impact of the choice of  $x_i$  on the predictions for that point only

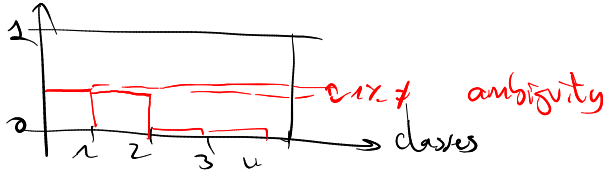


uncertainty: pick  $x_i$  for which  $f_\theta$  is the most uncertain

$$\arg \min_{x \in \mathcal{D} \setminus \mathcal{D}_L} \sup_{c \in \mathcal{C}} \hat{y}_i^c$$



margin:



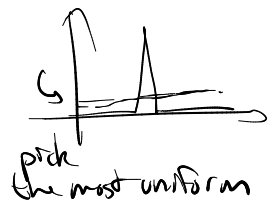
$$\arg \min_i \sup_{c_1 \in \mathcal{C}} \hat{y}_i^{c_1} - \hat{y}_i^{c_2}$$

where  $c_1 = \arg \max_c \hat{y}_i^c$   
 $c_2 = \text{second arg max}$

entropy:  $H(\hat{y}_i) = -\sum_{c \in \mathcal{C}} \hat{y}_i^c \log \hat{y}_i^c$

$\arg \max_i H(\hat{y}_i)$

$\hookrightarrow$  quantifies how spread the distribution over classes is



query by committee:

if predictor = ensemble of  $K$  models  $\hookrightarrow \hat{y}_{i,k}$

$\hookrightarrow$  do models agree? pick  $i$  where models disagree most

Global methods

$\hookrightarrow$  quantify impact of the sample choice over all dataset examples

Expected model change

$F_\theta \xrightarrow[\text{new labeled sample}]{\text{retrain with}} f_{\theta+\theta} \xrightarrow[\text{all unlabeled samples}]{\text{apply to}} \Rightarrow \text{impact?}$

just one training step  $(\nabla \downarrow)$

$$\theta_t \longrightarrow \theta_{t+1} = \theta_t - \eta \nabla_{\theta} \text{Loss}(\hat{y}_i, \delta_{c^*})$$

quantify the information gain as  $\|\theta_{t+1} - \theta_t\|$

$$\mathbb{E} \|\nabla_{\theta} \text{Loss}(\hat{y}_i, \delta_{c^*})\|$$

True label not known  $\Rightarrow$  average over possibilities  $c^*$

$$\mathbb{E}_{c^* \in \mathcal{C}} \left[ \|\nabla_{\theta} \text{Loss}(\hat{y}_i, \delta_{c^*})\| \right]$$

$\hookrightarrow p(c^*) = \hat{y}_i^{c^*}$

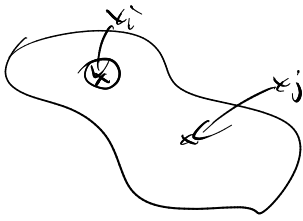
$$\sum_{c \in \mathcal{C}} \hat{y}_i^c \|\nabla_{\theta} \text{Loss}(\hat{y}_i, \delta_c)\|$$

# Expected error reduction

$$\arg \min_i \sum_{c \in C} c_j \hat{f}_j \sum_j \text{prediction variation for sample } x_j \text{ if trained also with } (x_i, c) \text{ all samples}$$

density-based method

or similarity  $\downarrow$  between  $x_i$  &  $x_j$



$$\theta_t \rightarrow \theta_{t+1} = \delta \theta = \eta \nabla_{\theta} \text{Loss}(x_i, c)$$

$$F_{\theta_{t+1}}(x_j) = F_{\theta_t}(x_j) + \underbrace{\delta \theta \cdot \nabla_{\theta} F_{\theta_t}(x_j)}_{\text{prediction variation}} + O(\delta \theta^2)$$

$$-\eta \nabla_{\theta} \text{Loss}(x_i, c) \cdot \nabla_{\theta} F_{\theta_t}(x_j)$$

Cons: computational power

pros: avoid choosing outliers

Sub focus on yet-unlabeled clusters of similar points

$$\nabla_{\theta} F_{\theta}(x_i) \frac{\partial L}{\partial f_{\theta}(x_i)} \nabla_{\theta} F_{\theta_t}(x_j)$$

(chain rule)

$\approx$  influence functions

$$h(x_i, x_j) = \nabla_{\theta} F_{\theta}(x_i) - \nabla_{\theta} F_{\theta}(x_j)$$

similarity

