

# Small data: weak supervision, transfer & incorporation of priors

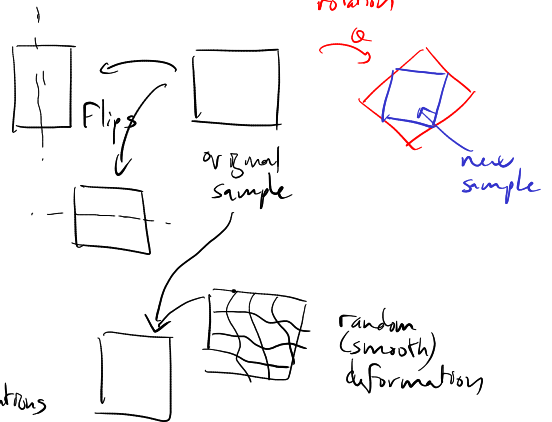
## I Dealing with small data

### Data augmentation

- For image classification tasks:

- deformation
- add rotations } spatial domain
- flips
- contrast
- color balance } intensity changes
- noise ...

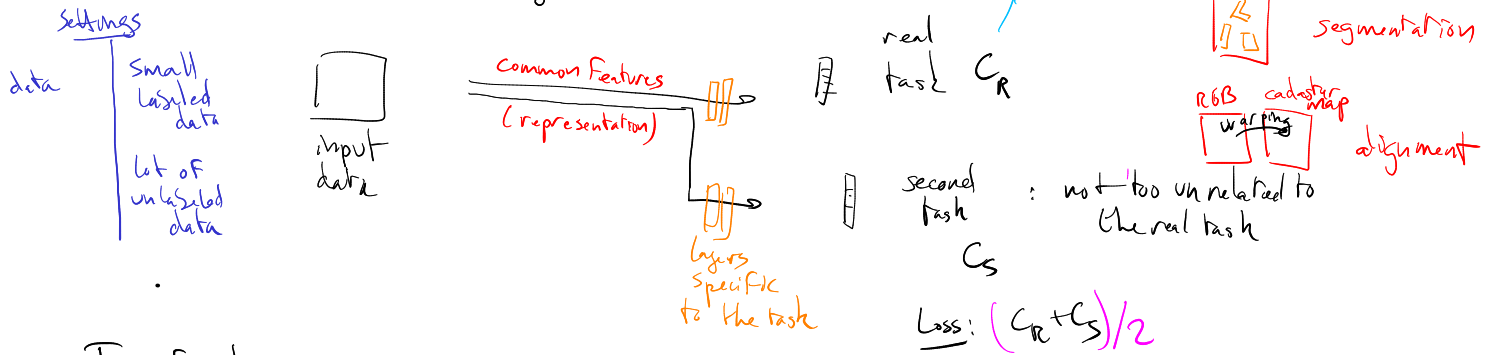
Model the "noise"  
non-meaningful transformations



- use a simulator  
↳ produce quantities of input data

### Multi-tasking

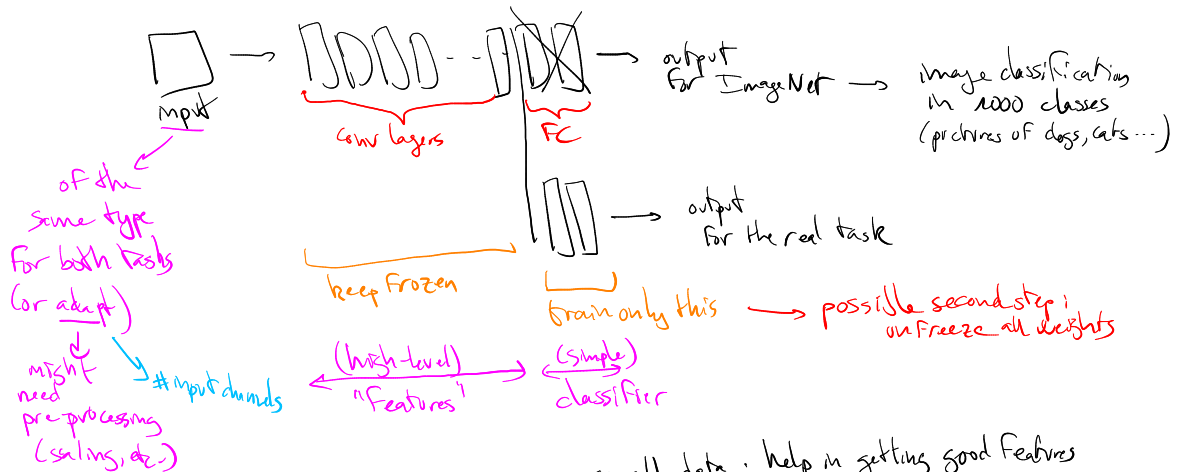
- consider (at the same training time) another task



### Transfer learning

- sequential training: - first  $C_S$   
- second  $C_R$

← pre-train a pre-trained network ex: for computer vision task, pre-trained VGG on ImageNet / ResNet



- analysis from [Rethinking ImageNet pre-training]

small data: help in getting good features  
big data: ≠ with training from scratch: big boost in training time (not necessary accuracy gain)

## II Forms of weak supervision

→ few labeled examples

### Semi-supervision

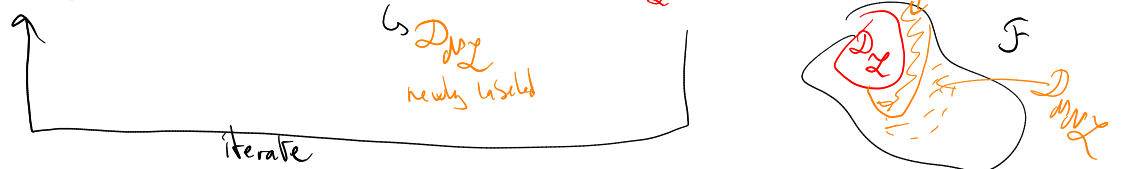


ex: when labeling is costly (requires time, expert...)

several approaches:

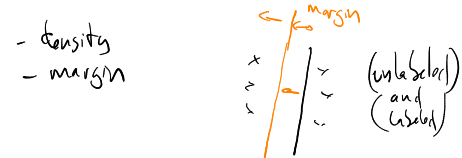
1) unsupervised training (on full set)  $\mathcal{F}$   $\rightarrow$  good representation  $\rightarrow$  supervised training  $\mathcal{D}_L$   
 ↳ Features  
 ↳ clustering  
 ↳ typically: auto-encoders

2) supervised training  $\rightarrow$  label some of the unlabeled samples  $\rightarrow$  bigger training set  $\mathcal{D}_L \cup \mathcal{D}_U$   
 $\mathcal{D}_U$   
 ↳  $\mathcal{D}_{new}$  newly labeled



issues: if mistakes, learn from wrongly labeled data

3) supervised training  $\rightarrow$  apply to full dataset  $\mathcal{F}$   $\rightarrow$  check some properties  
 ex: bias 50% / 50%  
 ↳ on  $\mathcal{F}$ : 60% / 60%  
 ↳ correct bias on  $\mathcal{F}$



Weak supervision

↳ more general: ex: labels could be noisy

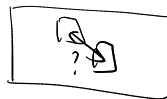
Self-supervision

$\rightarrow$  unsupervised (pre) training

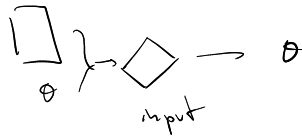
↳ supervised task (on  $\mathcal{F}$ )  $\rightarrow$  with a fake task with labels for all samples in  $\mathcal{F}$

ex: image classification

$\rightarrow$  image puzzle: extract patches from some image & ask for geometrical relation



$\rightarrow$  add a rotation to the image  $\rightarrow$  task: retrieve the angle  $\theta$  (random)



$\rightarrow$  data augmentation

↳ define "classes";  
 1 class = { all augmented input data coming from same sample }



ex: video classification

↳ auxiliary task:

$\rightarrow$  predict next frame: fully supervised

$\rightarrow$  give 3 frames, ask whether temporal order is correct

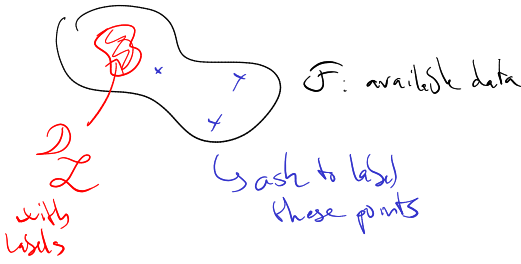


as many classes as original samples

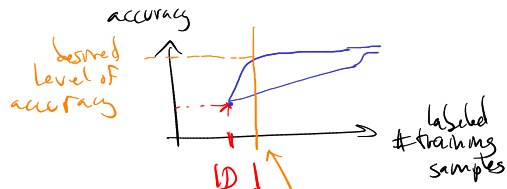
# Active learning

same setting as semi-supervision, + ask some samples to be labeled

(ex: costly labeling)



goal: increase accuracy as fast as possible (as a function of number of samples seen)



- large dataset  $\mathcal{F} = \{x_i\}$
- labels for few:  $(y_1, \dots, y_p)$  with  $p \ll |\mathcal{F}|$
- which  $x_i$  (with  $i \in [p, n]$ ) to pick? to be asked for labels

→ apply current model  $F$  to all samples → predictions  $\hat{y}_i = \begin{pmatrix} y_i^c \end{pmatrix}_{c \in \mathcal{C}}$  if classification task

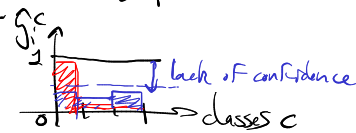
## Local methods

→ quantify the impact of the choice on the chosen sample only

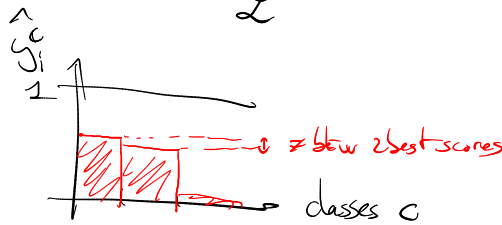
### \* uncertainty sampling:

pick  $x_i$  for which the model is the most uncertain → lowest prediction confidence

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



### \* margin sampling:



$$\arg \min_{i \in \mathcal{F} \setminus \mathcal{D}_L} \hat{y}_i^{c_1} - \hat{y}_i^{c_2}$$

where  $c_1, c_2$  are  $\hat{y}_i^c$

$$c_1 = \arg \max_c \hat{y}_i^c$$

$$c_2 = \dots$$

### \* entropy sampling:

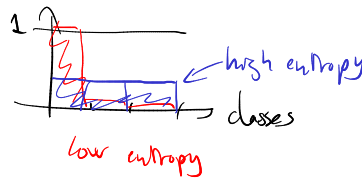
$$H(\hat{y}_i) = - \sum_{c \in \mathcal{C}} \hat{y}_i^c \log \hat{y}_i^c$$

entropy

proba. distrib. over classes

high: if probas. are well distributed over classes

low: if Dirac peak



### \* query by committee

→ predictor = ensemble of  $K$  models  $m_k \rightarrow \hat{y}_{i,k}$

→ do models agree?

→ pick  $x_i$  for which they most disagree

## Global methods

→ quantify impact of the choice over all dataset samples

### \* Expected model change

→ do one gradient-descent step with chosen sample

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

if one knew the class for sample  $i$

Dirac peak on  $\mathcal{E}_t$

$(x_i, c)$

$\mathbb{E} \nabla_{\theta} \hat{y}_i^c$

use predictions as class proba estimate

↳ expected parameter change  $\therefore \sum_{c \in \mathcal{C}} \hat{y}_i^c \left\| \nabla_{\theta} \text{Loss}(\hat{y}_i^c, \delta_c) \right\|$   
prio weight

↳ pick  $i$  leading to highest variation

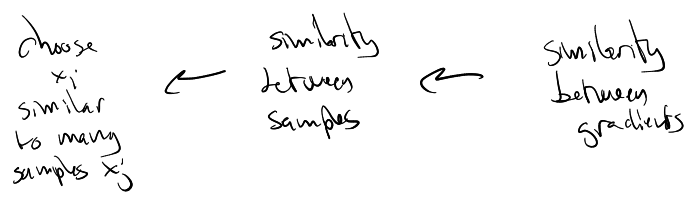
\* Expected error/reduction reduction

$\underset{i}{\text{arg min}} \sum_{c \in \mathcal{C}} \hat{y}_i^c \sum_j \text{prediction error for } x_j$   
 if trained with  $(x_i, \delta_c)$  also  
expectations over possible labels for  $x_i$

First-order dwp  
 $\hookrightarrow \hat{y}_j = F_{\theta_{\text{est}}}(x_j) \approx F_{\theta}(x_j) + (\delta\theta) \cdot \nabla_{\theta} F_{\theta}(x_j) + O(\delta\theta^2)$   
 $\theta_{\text{est}} - \theta$

$$-\eta \mathbb{E} \left[ \nabla_{\theta} \text{Loss}(x_i, \delta_c) \right]$$
  

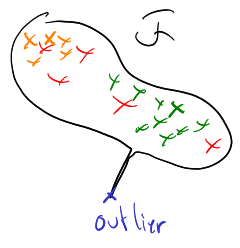
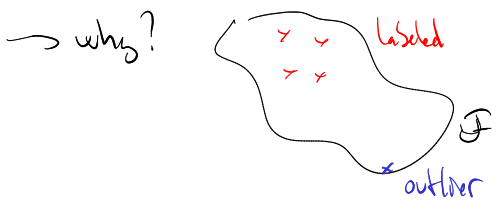
$$\nabla_{\theta} F_{\theta}(x_i) \cdot \nabla_{\theta} F_{\theta}(x_j)$$
  
 Sample chosen for active learning  
 any sample from dataset



→ density-weighted methods

↳ search for samples representative of many other ones

$\underset{i}{\text{arg max}} (\text{information brought by } (x_i, y_i)) \times \sum_j \text{similarity}(x_i, x_j)$   
x vs x      x vs x      x vs x



III Incorporation of priors

- small data
- help the training of the network by adding priors from physical knowledge

A) Invariance

Enforcement of invariance by design

- symmetry of the problem: group of transformations  $G$

$\forall g \in G, F(x) = F(gx)$

- no need to learn it
- easier training  $\approx$  data augmentation

- translation equi-variance: conv nets

$$F(\text{translated}(x)) = \text{translated}(F(x))$$

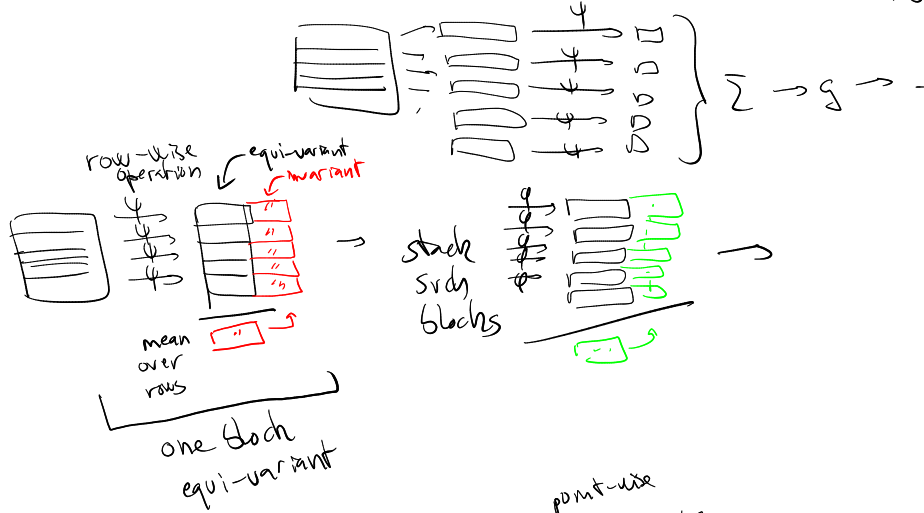
- permutation invariance:



[Deep Sets]

Theorem: any permutation-invariant function can be re-written as:

$$F(x) = g\left(\sum_r \psi(x^r)\right) \quad \forall F \dots \exists \psi, g \dots$$



Th: universal perm-equiv approximator

ex: point clouds → pointNet++



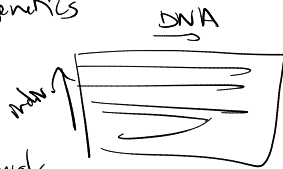
point-wise laser measurements on an object

→ 3D object classification

↳ point order -invariant

ex: population genetics

input: DNA from d individuals living now



mult-order -invariant

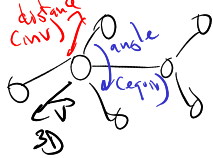
Task:

demography size inference

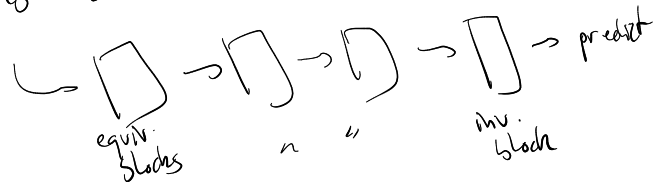
↳ how many people were living 5000 years ago?

- invariance to rotations  
equivariance

input: molecule



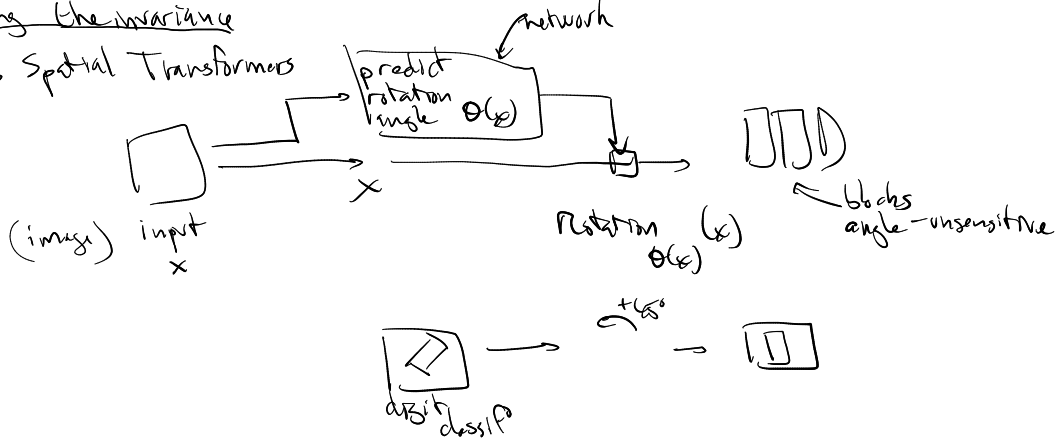
predict toxicity



$$\vec{v} \rightarrow \begin{cases} \|\vec{v}\| \text{ norm} \\ \text{direction } \frac{\vec{v}}{\|\vec{v}\|} \end{cases} \xrightarrow{\psi} \dots \rightarrow \psi(\|\vec{v}\|) \frac{\vec{v}}{\|\vec{v}\|}$$

Learning the invariance

→ Spatial Transformers



→ Capsule networks  
not no equiv

B) by task design → metrics  
→ data