

Computer vision showcases

- Medicine
 - Instagram, snapchat filters, face id, images upscaling
- OCR (optical characters recognition)
 - car plates detections / recognition, monitoring safety violations, person identification
- CCTV
 - car plates detections / recognition, monitoring safety violations, person identification
- Satellites
 - agricultural segmentation (crops), area calculation
- AI in defects detection
 - welding seams inspection
- AI in agriculture
- AI in autonomous driving
 - traffic tracking, LIDAR data processing, predicts future states of cars, pedestrians, lane and road recognition, driver perception

Tips & tricks on libraries

torch - research-friendly library for CV

torchvision - contains datasets

datasets. MNIST - hand-written digits

Dataloader - special PyTorch class to load dataset

Neural network definition

in PyTorch

1. Extend nn.Module
2. Override __init__ to set up layers
3. Override forward(x)

loss function is called criterion

Train the model

1. Put model in train mode by calling net.train()
2. Run inference by calling evoke operator ()
3. Calculate loss by calling criterion
4. Set gradients to zero by calling zero_grad()
5. Call backward from loss to calculate gradients
6. Do update by calling optimizer.step()

net.eval() puts network in the inference mode

Computational graph

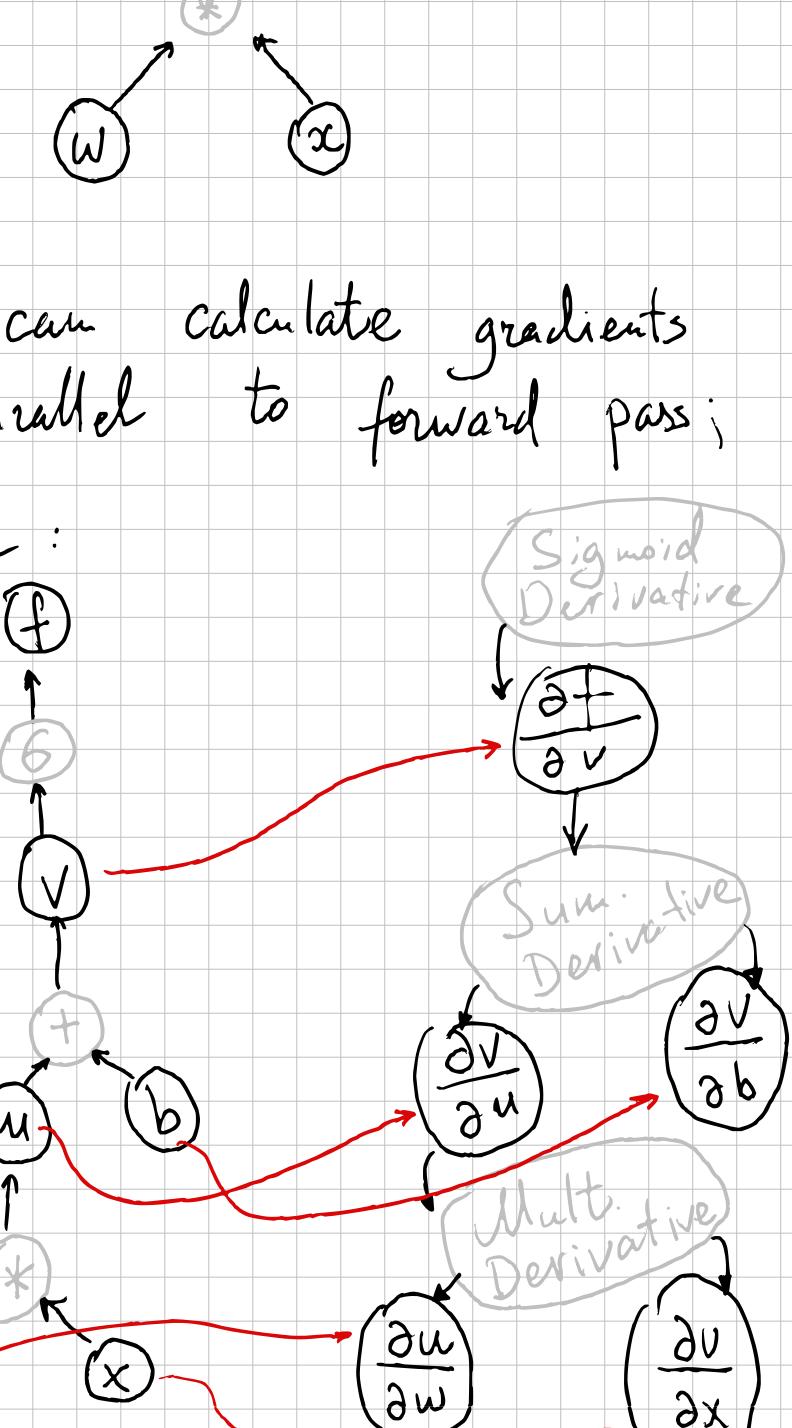
Idea: we can represent sequence of calculations as a graph;

let's take a sigmoid function as an example:

$$g(x) = \frac{1}{1 + \exp(-(w \cdot x + b))}$$

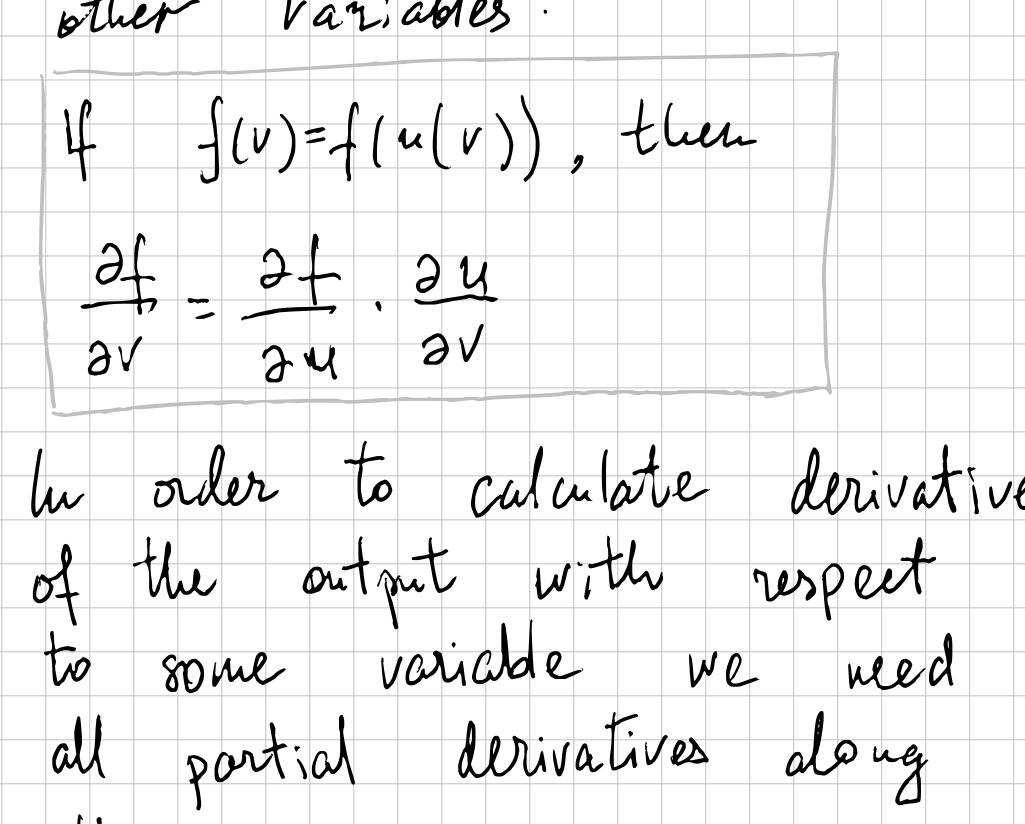
where w, b are parameters

Example



We can calculate gradients in parallel to forward pass;

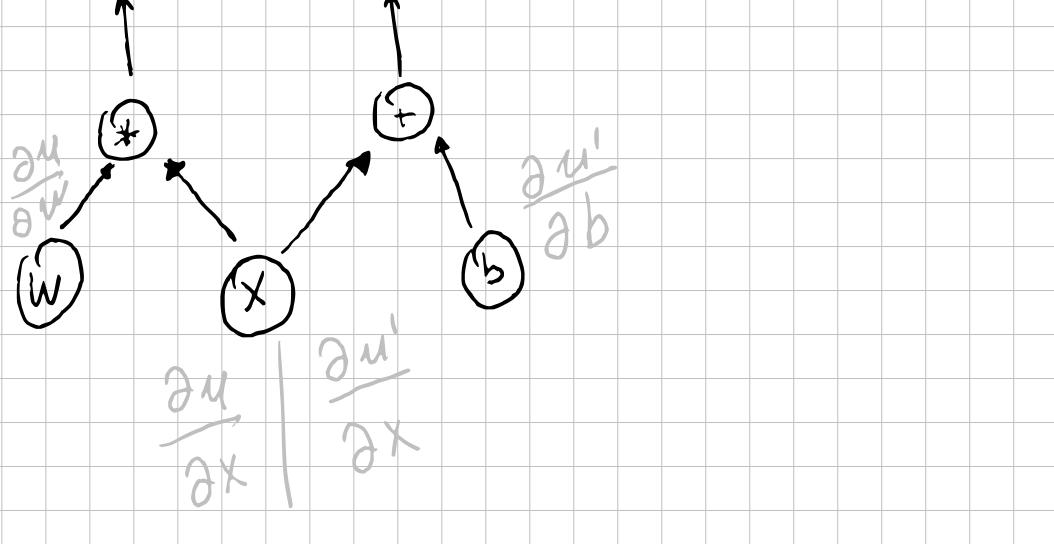
Example:



Chain rule helps us to calculate gradients of the output with respect to other variables:

If $f(v) = f(u(v))$, then

$$\frac{\partial f}{\partial v} = \frac{\partial f}{\partial u} \cdot \frac{\partial u}{\partial v}$$



Computer vision problems

Problem statement

Input image \underline{x}

We want to make a prediction \hat{y} based on an input image \underline{x}

$$\hat{y} = f(x)$$

We want to minimize empirical loss function

CV problems

- Classification : one & multiclass classification
- Detection
- Segmentation : semantic, panoptic, instance

CNNs

• ingredients :

1. An input image f of size $H \times W \times C$
2. Choose the receptive field F (trad. $S \times S \times C$)
3. For each pixel at position (i, j) calculate:

$$I'[i, j] = I[i:i+s, j:j+s, :] \odot F + b$$

4. Do above steps as many times as you want.

Properties of CNNs

- We use parameter sharing;
- Captures "local" relations;
- Translation invariants;

Dilated convolutions



Attention mechanism

- Squeeze-and-Excitation Attention

Main CV losses

Metrics are:

- evaluation from human perspective;
- can be used in optimization (ensembling)

Examples: accuracy, f1, auc-roc, mae, rmse, iou, dice

Losses are:

- evaluation from model perspective
- can be used to evaluate by human
 - can include regularizations
 - in neural networks should be differentiable

Examples: bce, mre, mae, focal, dice

Binary cross-entropy

\hat{y}, y

x -picture

θ -parameters

$$p(y=1 | x; \theta) = \hat{y}_\theta(x)$$

$$p(y=0 | x; \theta) = 1 - \hat{y}_\theta(x)$$

$$P(y/x; \theta) = \hat{y}_\theta(x)^y (1 - \hat{y}_\theta(x))^{1-y}$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$\rightarrow \max$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$\rightarrow \min$

Mean squared error

$$y = \hat{y}_\theta(x) + e, e \sim N(0, \sigma^2)$$

$$p(e) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{e^2}{2\sigma^2}\right) \rightarrow$$

$$\rightarrow p(y) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y-\hat{y})^2}{2\sigma^2}\right)$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n - \hat{y}_n)^2 \rightarrow \max$$

$$\text{Loss} = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2 \rightarrow \min$$

Other losses

Classification losses:

- Cross-entropy for multiclass classification

- Focal loss

- Dice loss

Regression losses:

- Mean absolute error

- Huber loss

Cross-entropy for multiclass

classification

$$CE = - \sum_{c=1}^C y_c \log(\hat{y}_c); C - \text{number of classes}$$

$$p(e) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{e^2}{2\sigma^2}\right) \rightarrow$$

$$\rightarrow p(y) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y-\hat{y})^2}{2\sigma^2}\right)$$

$$\log L = -C \sum_{n=1}^N (y_n - \hat{y}_n)^2 \rightarrow \max$$

$$\text{Loss} = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2 \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

$$L = \prod_{n=1}^N p(y_n) \rightarrow \max$$

$$\log L = -C \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow$$

$$\rightarrow \max$$

$$\text{Loss} = - \sum_{n=1}^N (y_n \log(\hat{y}) + (1-y_n) \log(1-\hat{y})) \rightarrow \min$$

CV metrics

Classification metrics:

- Accuracy
- Precision and recall
- F1-score, Dice score
- IOU, Jaccard index
- AU-ROC, RP - curve

Regression metrics:

- Mean Absolute Error
- Mean Squared Error
- Root Mean Squared Error
- R^2

Confusion matrix:

| | |
|----|----|
| TP | FN |
| FP | TN |

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FD}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{F1-score} = \frac{2 \times \text{precision} \times \text{recall}}{(\text{precision} + \text{recall})}$$

$$\text{IOU, Jaccard index} = \frac{TP}{TP + FP + FN}$$

ROC-AUC, RP - curve

Input: probabilities or scores

Sources of data and labels

Q: How to estimate amount of data?

A: Empirically

Depends on a task:

- Classification may require much more samples than, for example, segmentation

CIFAR10

60000

Image size: $32 \times 32 \times 3$

Number of classes: 10

Accuracy: 99%

CIFAR100

60000

Image size: $32 \times 32 \times 3$

Number of classes: 100

Accuracy: 96.8%

ImageNet - 1k

$1,3 \cdot 10^6$

Image size: $469 \times 387 \times 3$

Number of classes: 1000

Accuracy: 88.3%

LUNA16 (lung cancer dataset)

888 CT scans

Image size: $512 \times 512 \times 28$

Detection objects: 1186 nodules

Accuracy: 98.47% ROC-AUC

LITS2017

130 CT scans

Image size: $512 \times 512 \times 28$

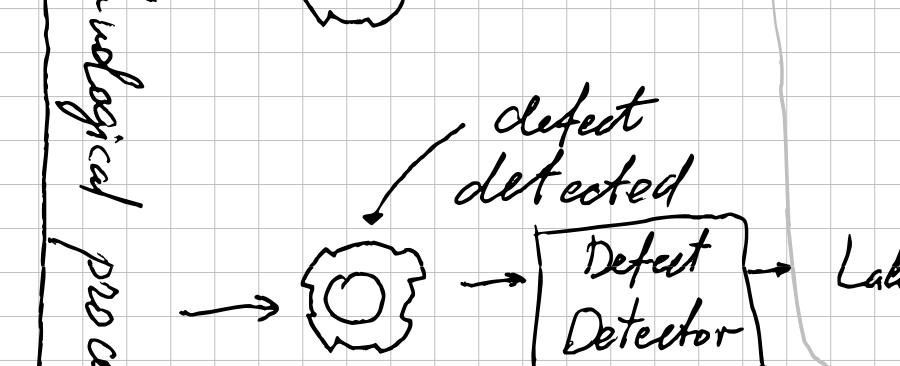
Task: Semantic segmentation

Accuracy: 89.85% IoU

Sources of data & labels

| Source | Pros. | Cons |
|-------------------------------------|---|---|
| Open datasets (LUNA, MS COCO) | <ul style="list-style-type: none"> + Free + Sometimes a lot of data + Sometimes very expensive and high quality labels | <ul style="list-style-type: none"> - Requires domain adaptation - Sometimes restrictive licensing |
| Crowdsource | <ul style="list-style-type: none"> + Cheap + Well-developed instruments | <ul style="list-style-type: none"> - Low quality - Requires a lot of supervision |
| Labeled by experts | <ul style="list-style-type: none"> + Full control of labeling process + Fits non-standard types of tasks | <ul style="list-style-type: none"> - Expensive - Requires a lot of supervision - Requires annotation tools |

Labels from another source



For example, person goes to the office with their magnetic pass: we can match data from surveillance system and the pass (badge).

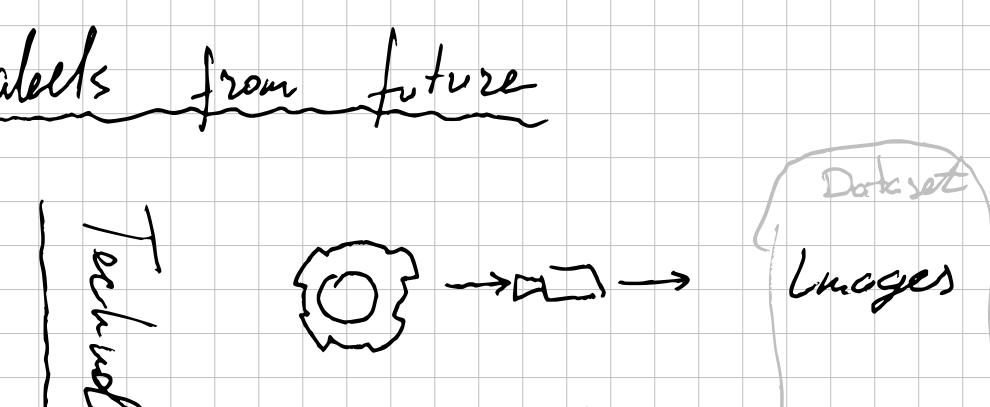
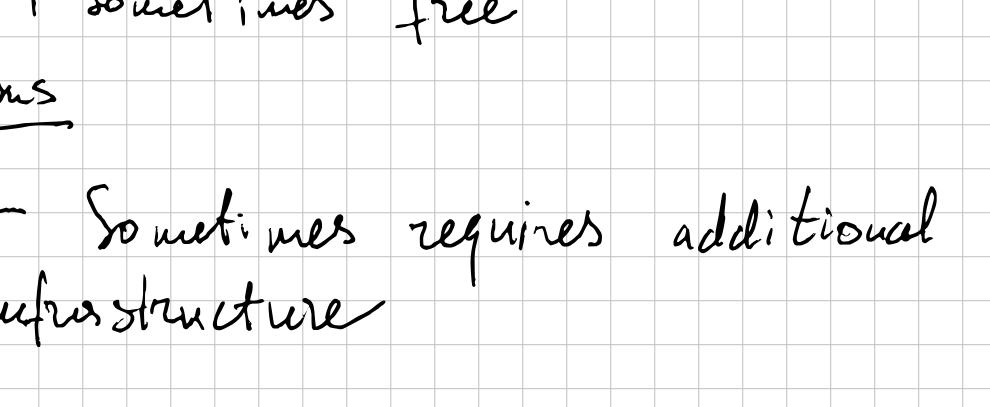
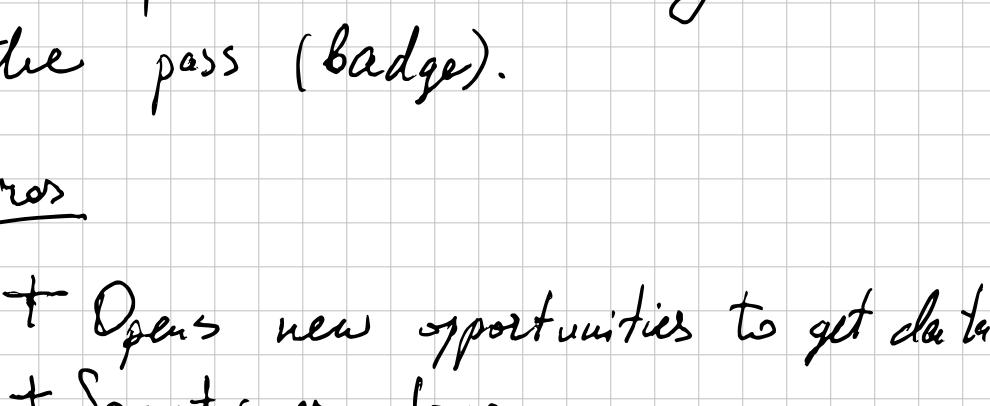
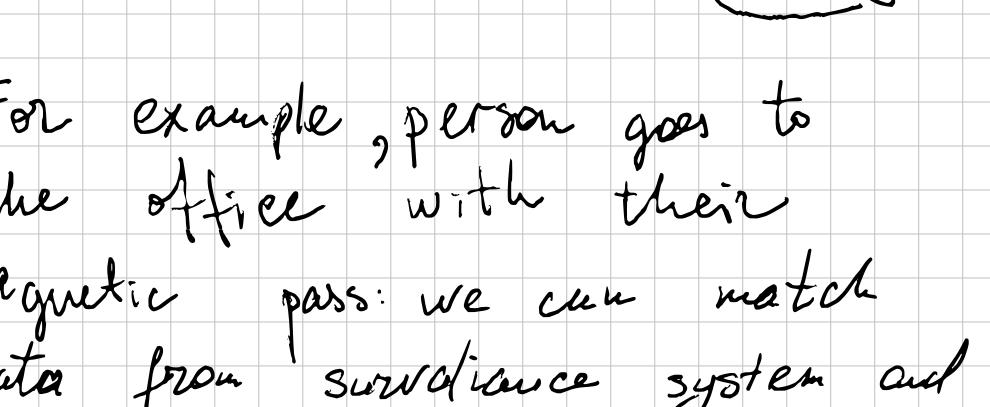
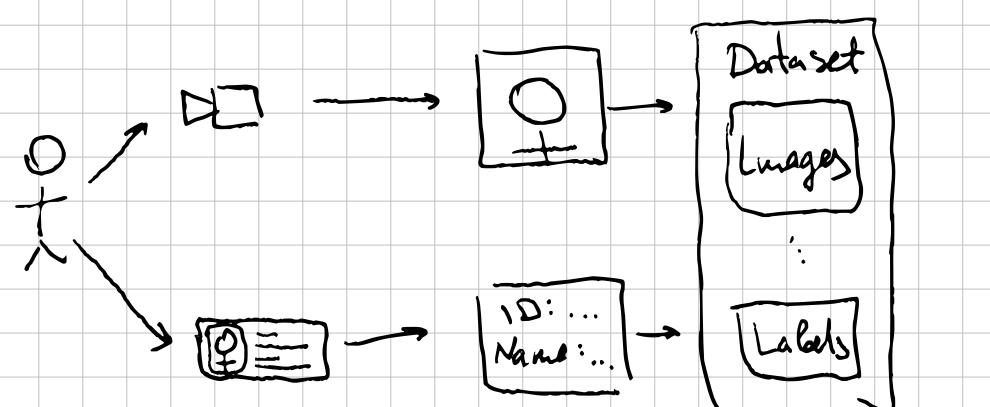
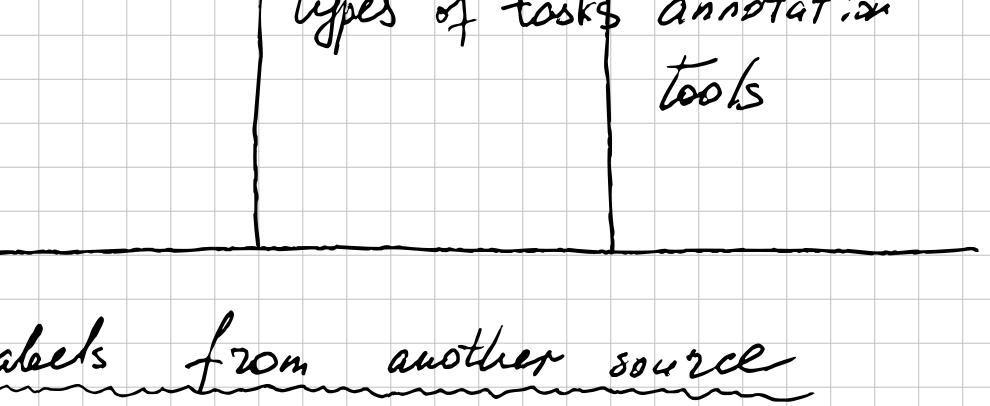
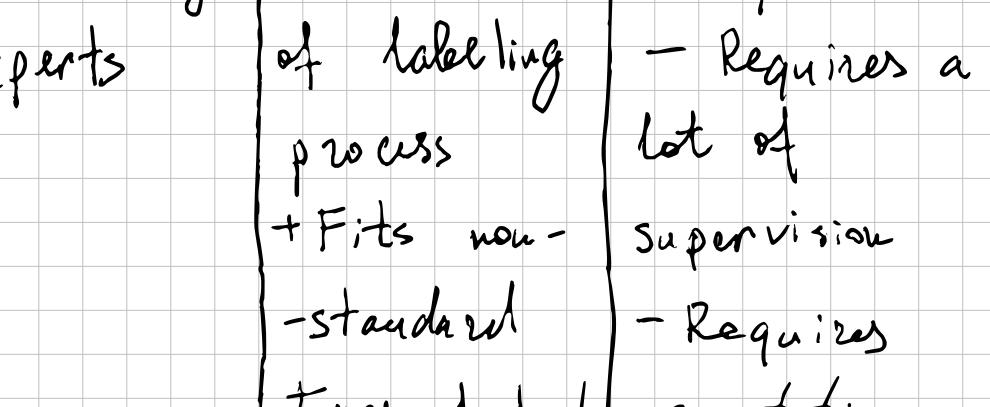
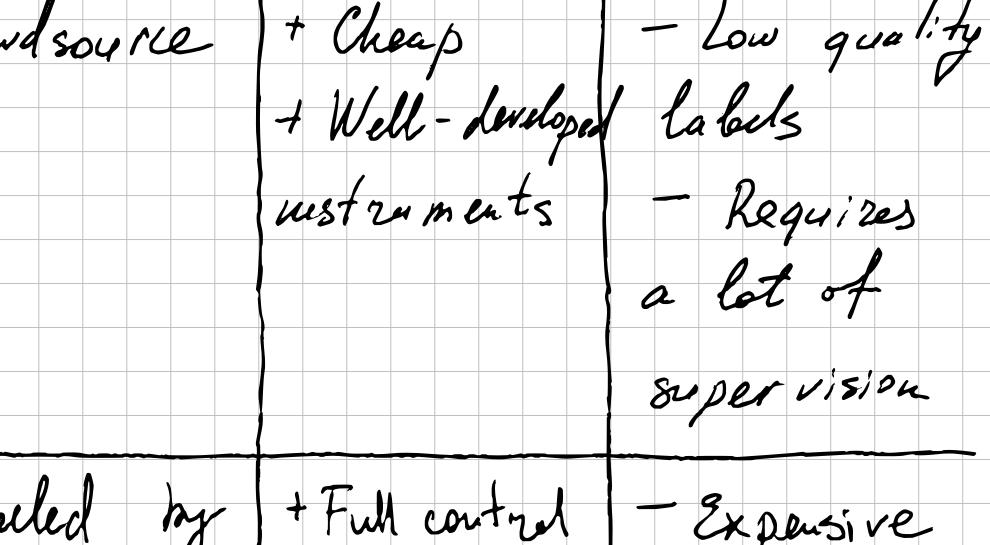
Pros

- + Opens new opportunities to get data
- + Sometimes free

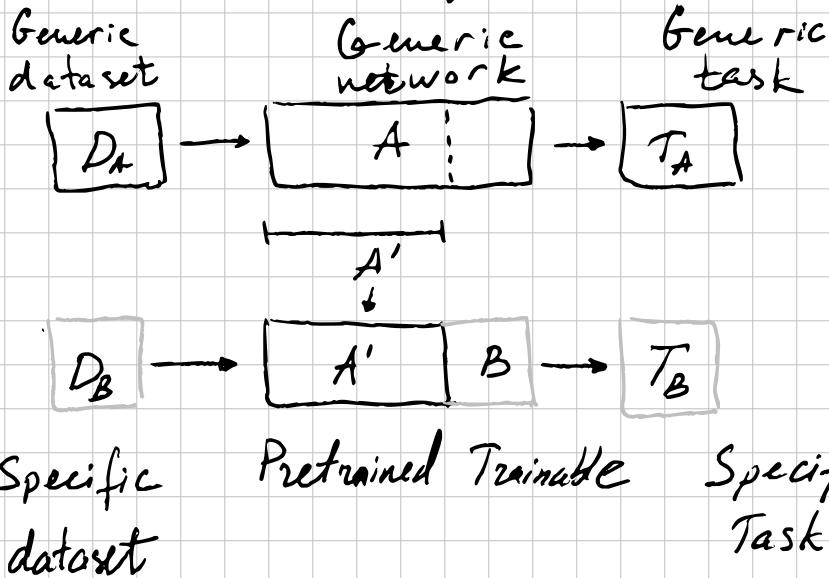
Cons

- Sometimes requires additional infrastructure

Labels from future



Transfer learning

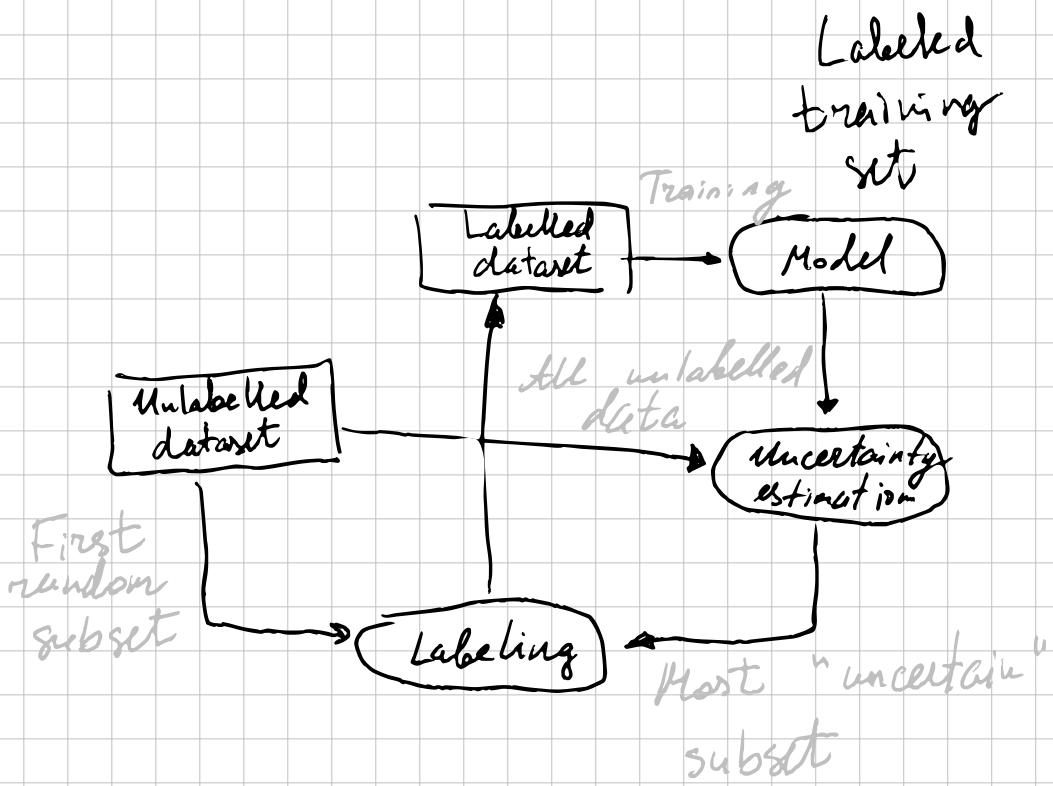
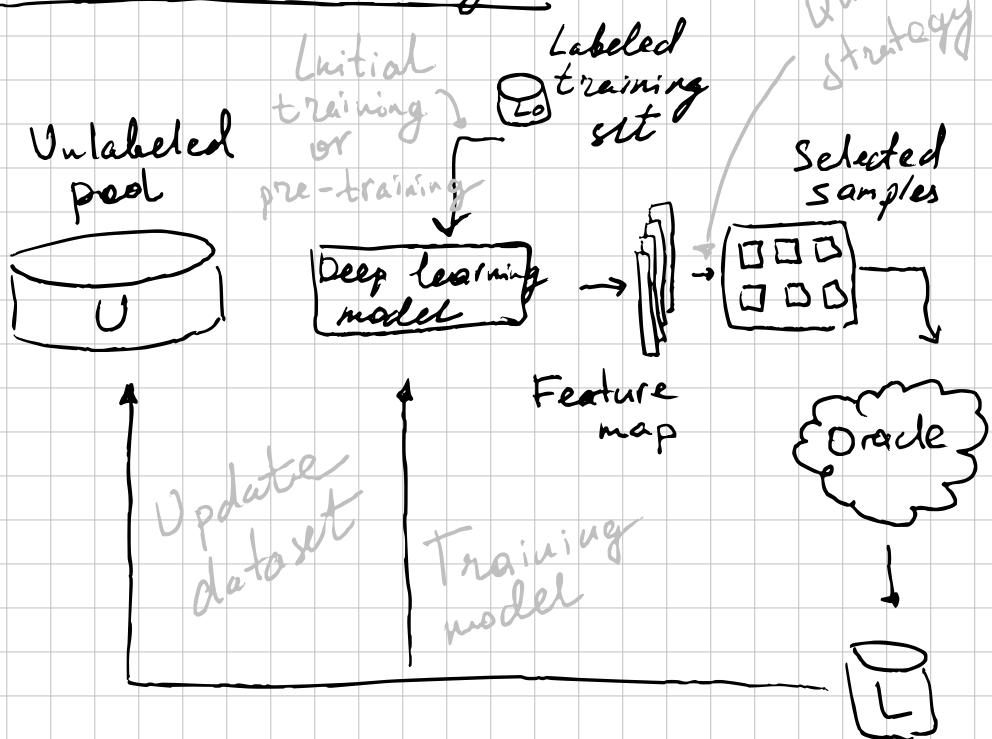


Usually, we can train model A efficiently on big generic dataset.

We assume A posses some knowledge about the system.

We train only our parameters of the network (B) .

Active learning



Pick examples that produces the most uncertain results, label them and return to model.

Estimation of dataset quality

1. Quality of a dataset

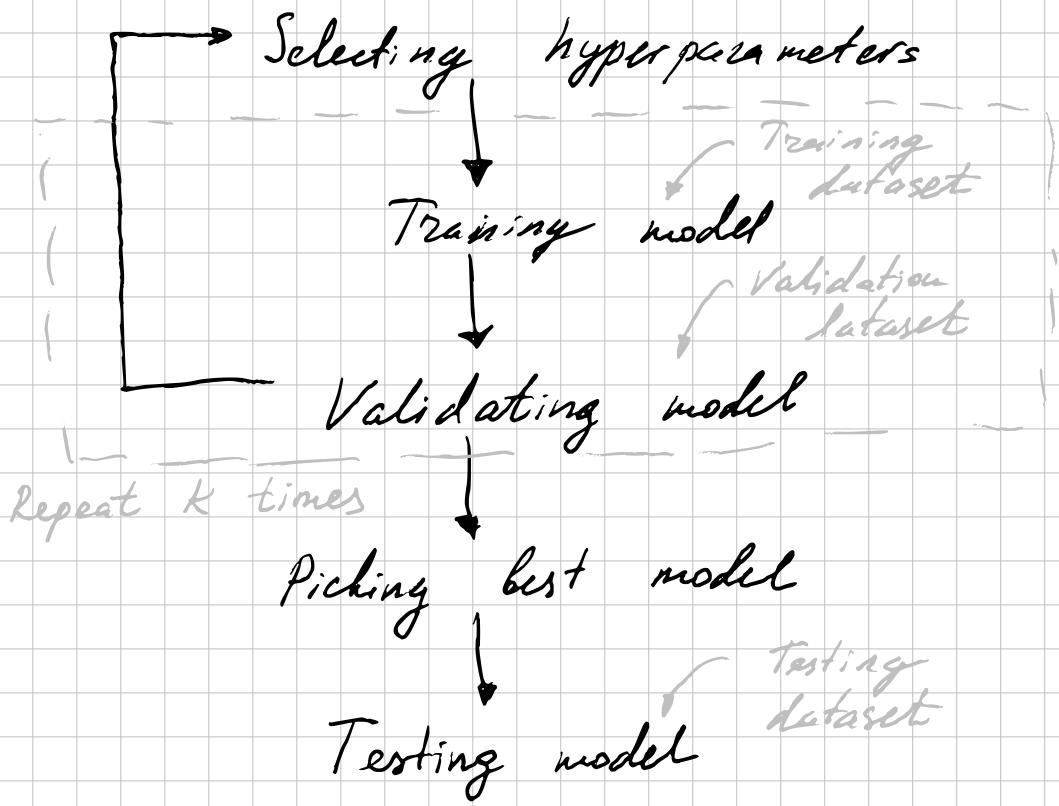
- Quantity: the more data we have - better
- Generalisation: how good data that we have represents general case

2. Quality of data

- Instrumental noise
- Exploratory data analysis : usually some general statistical error

3. Quality of labels

Model validation



- Test and validation sets should not have leaks (data from test set is in validation set)
- Test and validation sets should be close to production data;
no augmentations.
- If multiple separate test datasets - calculate separate metrics.
- Separate data into cohorts can be helpful

Learning rate optimizations and scheduling

We have a model $f(\cdot)$

let $\hat{y}_i = f(x_i)$ denote prediction

We penalise our model

using loss function $l(y_i, \hat{y}_i)$,
where y_i is ground truth

Empirical loss function is
defined as $L = \sum_i l(y_i, \hat{y}_i)$

We want to minimise
empirical loss function.

Stochastic gradient descent

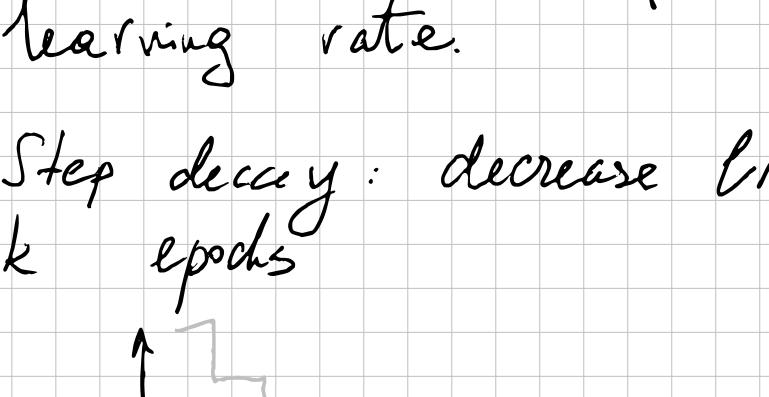
- Initialize parameters randomly
- Choose some learning rate λ
- Choose randomly k elements from the dataset;
- Update

$$W_{k+1} = W_k - \lambda \frac{1}{k} \sum_{i=1}^k \nabla_l(y_i, \hat{y}_i)$$

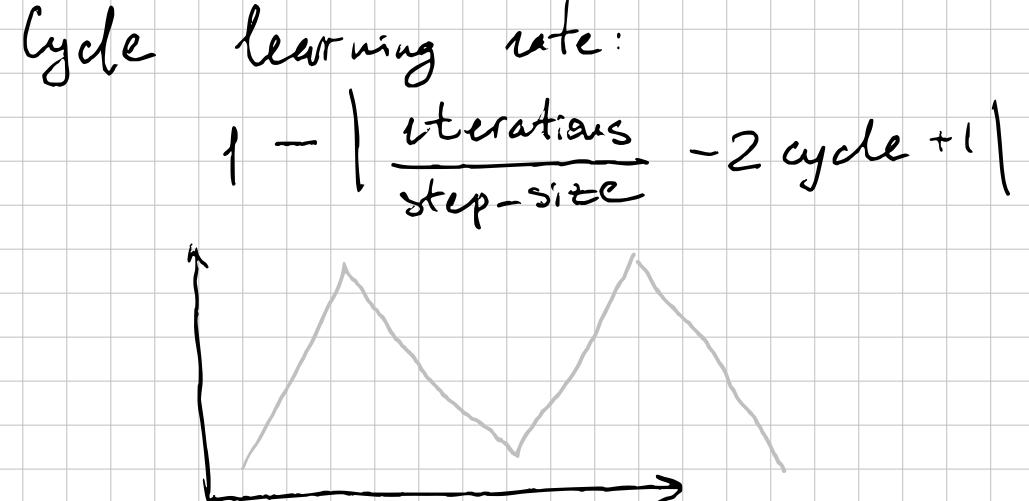
Update step:

$$\alpha := \alpha - \lambda \frac{d l(y_i, \hat{y}_i)}{da}$$

There are "good" functions
with easy to find local mins



Learning rate



Our goal is to find good learning rate.

Step decay: decrease lr every k epochs



Cycle learning rate:

$$1 - \left\lfloor \frac{\text{iterations}}{\text{step-size}} - 2 \right\rfloor \text{cycle} + 1$$

Benefits of cycle LR:

- Helps to avoid sharp minima;

- Helps to avoid saddle points;

Schedulers help with that.

Data augmentation

Motivation:

- NNs are overparametrized;
- NNs can produce very "sharp" predictions, i.e. in a vicinity of point x in parameter space it could predict different classes.

To solve the problem:

we can "enrich" our dataset by generating x' similar to x .

Pros:

- We will train our model on a dataset with more data samples;
- Hopefully the model will generalize better to unseen examples;

How to produce them:

1. Train generative model

$$p(x|y) \xleftarrow{\text{hard way}}$$

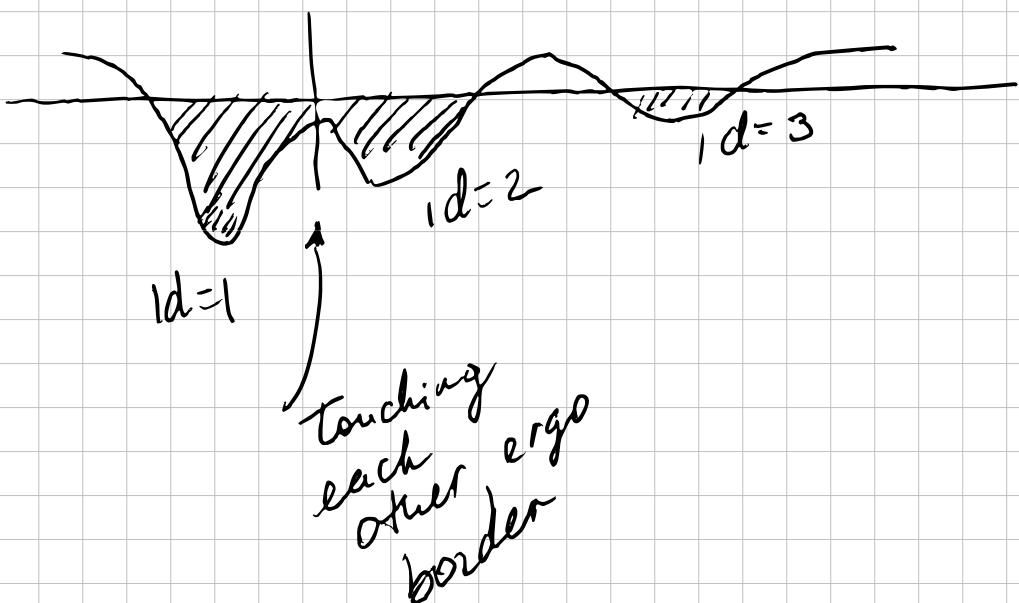
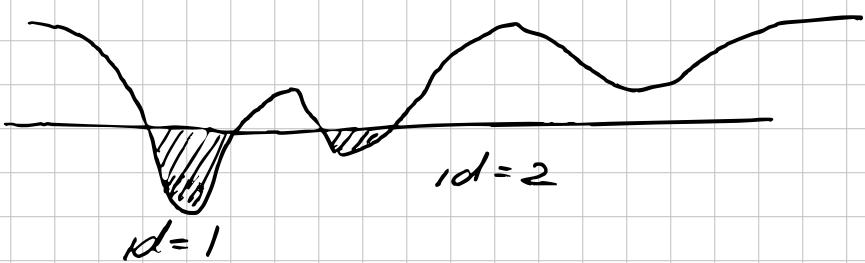
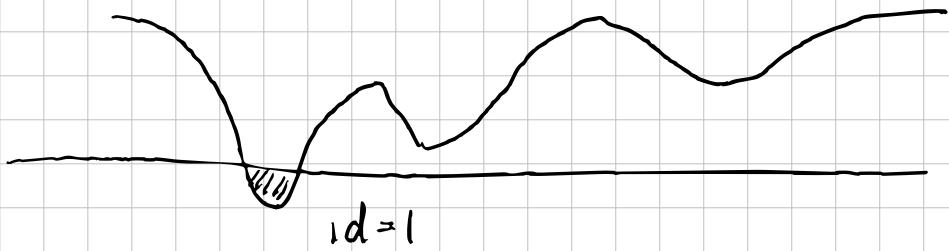
2. Data augmentations

Augmentations:

- Rotation
- Crop
- Coarse dropout
- Channel dropout
- Adding noise

Post processing

- Non-maximum suppression
- Watershed algorithm
 - 1. represent 1-channel image as height map
 - 2. Piling water onto it



- Erosion / Dilatation

The curse of dimensionality

Imagine that we have a function $f(x)$ and we want to calculate its values in uniformly distributed points

x_i

1d - 100 points

2d - 100^2 points

10d - 100^{10} points

etc

Hyperparameters optimisation

What if we want to optimise: learning rate, batch size, loss function, β in momentum.

First approach: Grid Search

We try every combination of hyperparameters.

• Bayesian optimization

① $x \in R$, $f(x) = \text{const.}$

$$\hat{y} = f(x) \in R$$

$$l_i = l(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2 = (y_i - c)^2$$

$$Z = \sum_i l_i = \sum_i (y_i - c)^2$$

$$Z_c' = \sum_i -2(y_i - c) = 0$$

$$\sum_i (y_i - c) = 0 \Leftrightarrow \sum y_i - \sum c = 0$$

$$\sum y_i - nc = 0 \Leftrightarrow nc = \sum y_i$$

$$c = \frac{\sum y_i}{n}$$

② $\hat{y}_i = f(x) = ax_i + b$

$$l_i = l(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2 =$$

$$= (y_i - ax_i - b)^2 = (ax_i + b - y_i)^2$$

$$Z = \sum_i l_i = \sum_i (ax_i + b - y_i)^2$$

$$Z_a' = \sum_i 2(ax_i + b - y_i) \cdot x_i$$

$$Z_b' = \sum_i 2(ax_i + b - y_i)$$

$$\sum ax_i + b - y_i = 0$$

$$nb = \sum y_i - ax_i$$

$$b = \frac{(\sum y_i - ax_i)}{n}$$

$$\bar{y} = \frac{\sum y_i}{n}$$

$$\bar{x} = \frac{\sum x_i}{n} \quad b = \bar{y} - a\bar{x}$$

$$\sum (ax_i + b - y_i)x_i = 0$$

$$a \sum x_i^2 = \sum x_i y_i - b \sum x_i$$

$$a = \frac{\sum x_i y_i - b \sum x_i}{\sum x_i^2}$$

Ensembling

- In order to produce ensemble of models that performs better than individual models, model should be:
 - Accurate (better than random guessing)
 - Diverse (f_1 & f_2 are diverse if $\ell(f_1(x), y) \neq \ell(f_2(x), y)$ on new data points)

Example :

10 model

$P = 2/3$

- we use majority voting
- error is

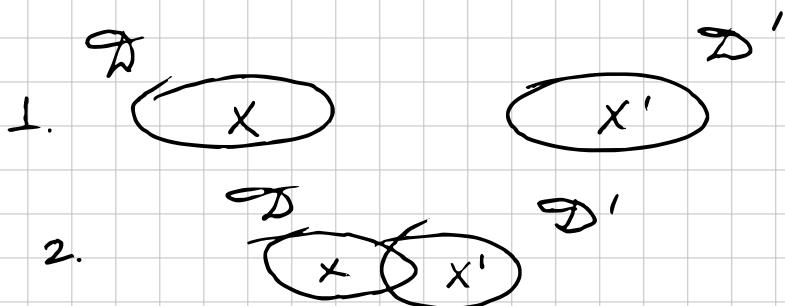
$$\sum_{i=0}^{10} \binom{10}{i} \frac{1}{3}^i \frac{2}{3}^{10-i} \approx 0.043$$

How to build ensembles:

- Bagging
- Cross-validation
- Stacking

Transfer learning

- Often we do not have enough data to train a neural network
- Pre-trained model was trained on a domain D
 - Domains could either intersect or not:



If new targets are not far away:

- Freeze a backbone, fine-tune classification head

New targets are far:

- Freeze some deep layers, and fine-tune top

Unsupervised learning

- Often we don't have an access to massive labeled data

Contrastive learning

- Split on positive and negative examples (that are far away)