

INDICE

1 INTRODUCTION

MSE

BCE + Suggested Reading
Generative vs Discriminative

2 SHALLOW AND DEEP NETWORKS

- Shallow Network
- Universal Approximation Theorem
- Activation Functions
- Deep Networks
- VAT: Depth Version

3 TRAINING & LOSS FUNCTION

- Loss Function, Maximum likelihood
- Log likelihood, Negative Log likelihood
- Recipe for loss
 - Root Square Loss
 - Heteroscedastic vs Homoscedastic
- Multiclass Classification, Multiclass Cross Entropy
- Cross Entropy

4 MODEL TRAINING, FITTING AND BACKPROP

- Gradient Descent
- Gabor Model and SGD
 - SGD Scheduler LA, Momentum, Nesterov Momentum, Adam
- Backpropagation
- Computation Graph
- Parameter Initialization
- Vanishing & Exploding gradient

5 PERFORMANCE AND REGULARIZATION

- Sources of Errors
- Noise, Bias, Praise
- Variance, Bias
- Double Descent
- Overfitting
- Regularization
 - Early Stopping
 - As a prior
 - Dropout
 - L2
 - Data Augmentation

6 CNN - RES-NET

- CNN Geometric view
- Convolution Theorem
- Discrete Conv
- Conv Layer, Receptive Field
- Resnet: the depth
- Exploding gradient & Batch Norm
- U-Net
- Highway Networks

7 SELF SUP. LEARNING

- Siamese NN
- Neural Collaborative Filtering
- Self Supervision, Proxy Task
- Contrastive Learning, Triplet Loss
- SIMCLR
- BYOL + Evolution
- BARLOW TWINS + Evolution
- Deep Metric

8 TRANSFORMER

- RNN
- Seq 2 Seq
- Autoregressive Models
- ELMo, Word2Vec
- Attention, Self-Attention, MHA
- Transformer
- BERT
- Visual Transformer
- Transfer Learning, Neural IDBS, T5

9 GENERATIVE AI

- Discriminative vs Generative
- GAN, Lotka, Convergence, Mode collapse
- Wasserstein GANs
- Autoencoders, Denoising Autoencoders
- VA, ELBO, Reparametrization Trick

10 GRAPH NEURAL NETWORK

- GNN Basic
- Tasks
- Graph CNN
- Batching
- Aggregation
- Spectral Graph
- Graph Deletion
- Spectral Graph convolution
- Chebnet

11 GEOMETRIC DEEP LEARNING

- Group & Actions
- Geom. deep learning
- Graph NN, Equivariance
- Deep Sets
- Group Equivariant

12 MODEL COMPRESSION

- Pruning NN
- Lottery Ticket Hypothesis
- Weight Quantization

14 RAG

15 NOISY LABEL

13 META LEARNING

- Meta Learning Problem, Few Shot
- Supervised vs meta learning
- Metric-Based Meta Learning
- Model-Based Meta Learning
- Optimization-Based Meta Learning
- Learning to Learn
- MAML

CAPITOLO 1

INTRODUCTION

MEAN SQUARED ERROR Used for regression

$$L(\theta) = \sum_{i=1}^n [f(x, \theta) - y_i]^2 = \sum_{i=1}^n [(\theta_0 + \theta_1 x_i) - y_i]^2$$

BINARY CROSS ENTROPY : For classification Problems

$$L(\theta) = \sum_{i=0}^n - (1 - y_i) \log \left[1 - \frac{1}{1 + e^{-h}} \right] - y_i \log \left[\frac{1}{1 + e^{-h}} \right]$$

Optimization of BCE, computation of gradient

$$h = w^T X \quad z = \text{sig}(h) = \frac{1}{1 + e^{-h}}$$

$$J(w) = - (y \log(z) + (1 - y) \log(1 - z))$$

The Chain Rule

$$\frac{\partial J(w)}{\partial w} = \frac{\partial J(w)}{\partial z} \frac{\partial z}{\partial h} \frac{\partial h}{\partial w}$$

$$\text{Gradient Descent} = w - \alpha \frac{\partial J(w)}{\partial w}$$

DISCRIMINATIVE JS GENERATIVE

$$y = f(x, \theta)$$

$$x = g(y, \theta)$$

In this setting the inference to get y would require the inversion of the equation
 $y = g^{-1}(x, \theta)$

CAPITOLO 2

SHALLOW AND DEEP NEURAL NETWORKS

- Shallow Network
- Universal Approximation Theorem
- Activation Functions
- Deep Networks
- VAT : Depth Version

BASICS OF NEURAL NETWORKS

$f[x, \phi]$

input parameters output

x input scalar and y is a scalar

ϕ is composed of 10 parameters $(\phi_0, \phi_1, \phi_2, \phi_3, \theta_{10}, \theta_{11}, \theta_{20}, \theta_{21}, \theta_{30}, \theta_{31})$

Then we can express

$$y = f(x, \phi)$$

$$= \phi_0 + \phi_1 \square + \phi_2 \square + \phi_3 \square$$

Now is a multivariate (3) linear function

$$= \phi_0 + \phi_1 a(\square) + \phi_2 a(\square) + \phi_3 a(\square)$$

LINEAR LAYER

$$= \phi_0 + \phi_1 a(\theta_{10} + \theta_{11}x) + \phi_2 a(\theta_{20} + \theta_{21}x) + \phi_3 a(\theta_{30} + \theta_{31}x)$$

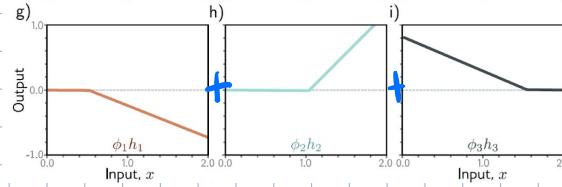
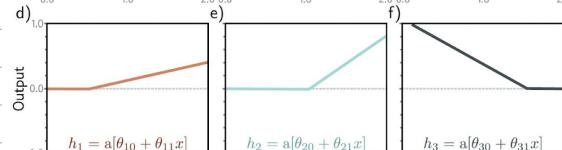
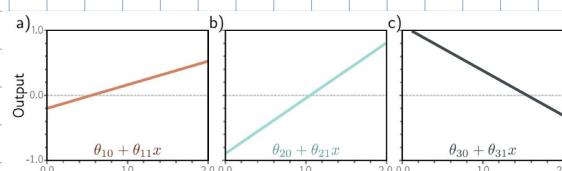
Now let's see this function by the use of hidden units

$$= \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

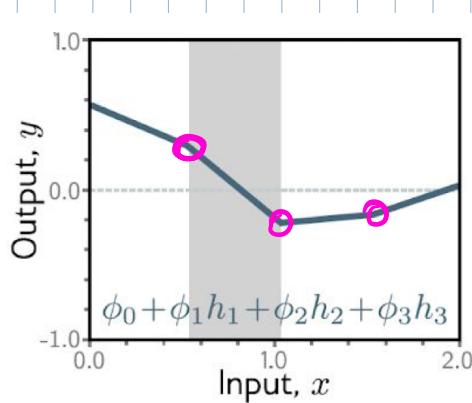
$$h_i = [\theta_{i0} + \theta_{i1}x]$$

Let's take $a = \text{ReLU}(x)$

$$= \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$



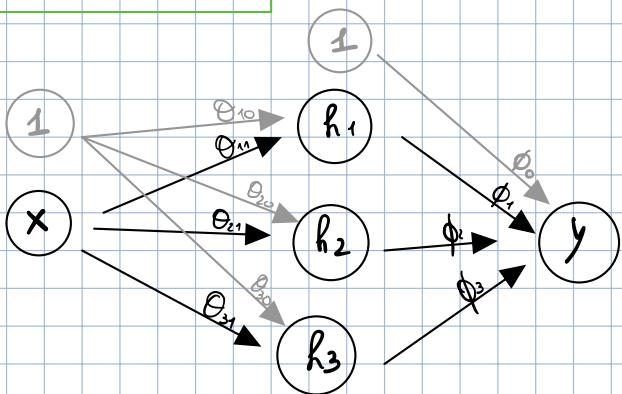
Sum of functions \rightarrow the expression



3 hidden units
= 3+1 splits

Now let's see the function with the neuron scheme

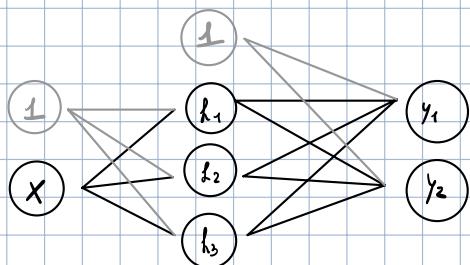
$$y = f(x, \phi) = \\ = \phi_0 + \phi_1 a(\theta_{10} + \theta_{11} x) + \phi_2 a(\theta_{20} + \theta_{21} x) + \phi_3 a(\theta_{30} + \theta_{31} x)$$



Considering a multivariate output function

$$y_1 = \phi_{10} + \phi_{11} a(\theta_{110} + \theta_{111} x) + \phi_{12} a(\theta_{120} + \theta_{121} x) + \phi_{13} a(\theta_{130} + \theta_{131} x)$$

$$y_2 = \phi_{20} + \phi_{21} a(\theta_{210} + \theta_{211} x) + \phi_{22} a(\theta_{220} + \theta_{221} x) + \phi_{23} a(\theta_{230} + \theta_{231} x)$$

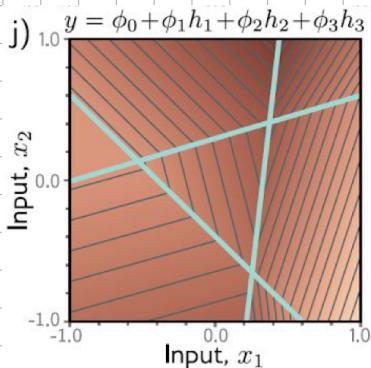


Multivariate Input function

$$\text{Now } x = [x_1, x_2]$$

$$\text{then } y = \phi_0 + \phi_1 a(\theta_{10} + \theta_{11} x_1 + \theta_{12} x_2) + \phi_2 a(\theta_{20} + \theta_{21} x_1 + \theta_{22} x_2) + \phi_3 a(\theta_{30} + \theta_{31} x_1 + \theta_{32} x_2)$$

linear function of the $\left\{ h_i \right\}$
multivariate input



3 hidden units (3 lines)

2 input (2 dimension - plot)

Multivariate input and output

$$x \in \mathbb{R}^{d_i} \quad y \in \mathbb{R}^{d_o} \quad h \in \mathbb{R}^d \quad h_d = a \left(\theta_{d0} + \sum_i \theta_{di} x_i \right) \quad y = \phi_0 + \sum_i \phi_i h_i$$

UNIVERSAL APPROXIMATION THEOREM

D hidden units h_1, \dots, h_D $h_i = a \left(\theta_{i0} + \theta_{ii} x \right)$

The non linear function will be $y = \theta_0 + \sum_i \theta_i h_i$

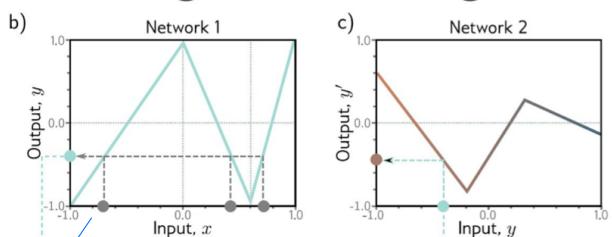
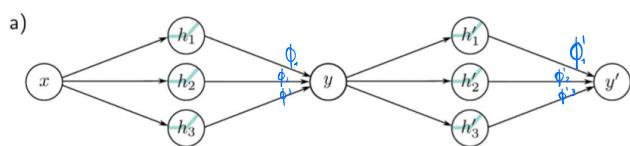
With D hidden units we have $D+1$ splits

Make $D \rightarrow \infty$ and you can approximate (any) complex functions

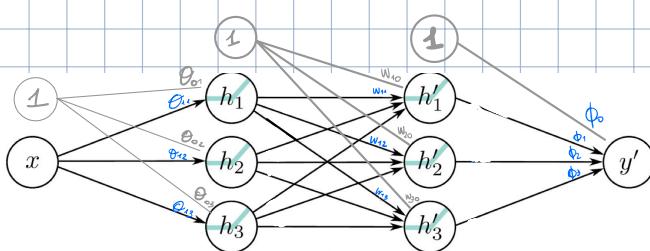
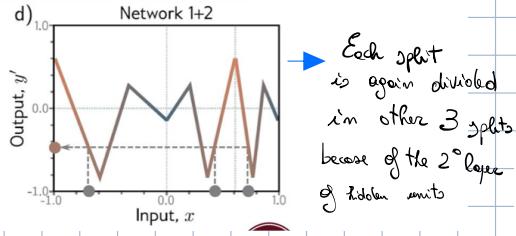
THEOREM: For any continuous function \exists a shallow NN that approximates that function

DEEP NEURAL NETWORKS

Case 1



h then 3 splits



CASE 2 : MLP

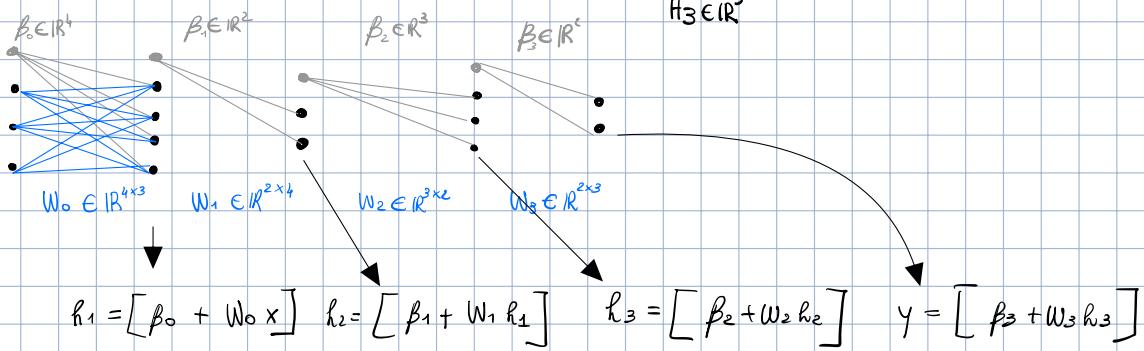
$$y' = \phi_0 + \phi_1 h'_1 + \phi_2 h'_2 + \phi_3 h'_3$$

$$\begin{aligned} &= \phi_1 a(w_{10} + w_{11} h_1 + w_{12} h_2 + w_{13} h_3) + \phi_2 a(w_{20} + w_{21} h_2 + w_{22} h_2 + w_{23} h_3) + \phi_3 (w_{30} + w_{31} h_1 + w_{32} h_2 + w_{33} h_3) \\ &= \phi_1 a(w_{10} + w_{11} a(\theta_{10} + \theta_{11} x) + w_{12} a(\theta_{20} + \theta_{21} x) + w_{13} a(\theta_{30} + \theta_{31} x)) \\ &\quad + \phi_2 a(w_{20} + w_{21} a(\theta_{20} + \theta_{21} x) + \dots) \end{aligned}$$

MATRIX NOTATION

$$x \in \mathbb{R}^3 \quad y \in \mathbb{R}^2 \quad H = \{H_1, H_2, H_3\}$$

$$H_1 \in \mathbb{R}^4 \\ H_2 \in \mathbb{R}^2 \\ H_3 \in \mathbb{R}^3$$



EXERCISE MLP PARAMETERS COUNT

$$x \in \mathbb{R}^{6220800} \quad D_1 = 512 \quad D_2 = 1024 \quad D_3 = 512 \quad y \in \mathbb{R}^1$$

$$\begin{aligned} \# \text{ of parameters} &= (\beta_0 + W_0) + (\beta_1 + W_1) + (\beta_2 + W_2) + (\beta_3 + W_3) \\ &= (512 + (512 \cdot 6220800)) + (1024 + (1024 \cdot 512)) + (512 + (512 \cdot 1024)) + (1 + 512) = \dots \end{aligned}$$

CAPITOLO 3

TRAINING & Loss Function

Loss

Is a function that returns a measure of the mismatch between prediction of the model and ground truth

Training objective is to find the best parameters of $f(x, \theta)$, minimizing the mismatch

$$\hat{\theta} = \underset{\phi}{\operatorname{argmax}} L(\phi)$$

MAXIMUM LIKELIHOOD

We can see the ML model as a method to emit a conditional probability over the possible output, given the x .

$$P(y|x), \text{ then we have to consider the prob. to get an output given the input and the parameters of the model}$$

$$P_c \{y | x; \theta\} \rightarrow \text{likelihood}$$

What the model should do is to find the parameters θ^* such that the likelihood is maximum: the combined probability over all the dataset is maximum

Maximum Likelihood Criterion

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \left(\prod_{i=1}^D P_c(y_i | f(x_i, \theta)) \right)$$

the model will generate an output probability for each possible output value given each input x in D

Maximum Log Likelihood (Transform P in Σ)

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \left[\sum_{i=1}^D \log \left(P_c(y_i | f(x_i, \theta)) \right) \right]$$

Negative Log Likelihood (Convert the problem into minimization problem)

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \left[- \sum_{i=1}^D \log \left(P_c(y_i | f(x_i, \theta)) \right) \right] = \underset{\theta}{\operatorname{argmin}} L(\theta)$$

Inference (Return y with highest prob. according to the model)

$$\hat{y} = \underset{y}{\operatorname{argmax}} P_c(y | x_i, \hat{\theta})$$

- Loss Function, Maximum Likelihood
- Log Likelihood, Negative Log Likelihood
- Recipe for loss
 - Root Square Loss
 - Heteroscedastic vs Homoscedastic
- Multiclass Classification, Multiclass Cross Entropy
- Cross Entropy

BUILD A LOSS in 4 steps

- ① Choose a prob. distribution $P_2(y|\theta)$ θ parameters of the distribution
- ② Set the ML model $f(x, \phi)$ ϕ parameters of the model predict $\hat{\theta} = f(x, \phi)$
and $P_2(y|\theta) = P(y|f(x, \phi))$
- ③ Train the model finding ϕ such that minimize the loss (the NLL Loss)

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[- \sum_{i=1}^D \log \left(P(y_i|f(x_i, \phi)) \right) \right]$$
- ④ Inference return $P_2(y|f(x, \phi))$ or $\underset{y}{\operatorname{argmax}} P_2(y|f(x, \phi))$

DERIVATION OF THE MSE Loss

- ① Consider Gaussian distribution with y, σ^2 $P_2(y|y, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y-y)^2}{2\sigma^2} \right) =$
 - ② Set 2 ML models compute both y, σ^2
then $y = f_1[x, \phi]$ $\sigma^2 = f_2[x, \phi]^2$
 - ③ $\hat{\phi} = \underset{\phi, \sigma^2}{\operatorname{argmin}} \left[- \sum_{i=1}^D \log \left[\frac{1}{\sqrt{2\pi f_2(x|\phi)}} \exp \left(-\frac{(y-f_1(x|\phi))^2}{2f_2(x|\phi)^2} \right) \right] \right]$
- or using one model to get both (set as differentiable also σ^2) $\hat{\phi} = \underset{\phi, \sigma^2}{\operatorname{argmin}} \left[- \sum_{i=1}^D \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y-f(x|\phi))^2}{2\sigma^2} \right) \right] \right]$

BINARY CROSS ENTROPY

$$P_2(y|\lambda) = \begin{cases} 1-\lambda & y=0 \\ \lambda & y=1 \end{cases} \rightarrow P_2(y|\lambda) = (1-\lambda)^{1-y} \cdot \lambda^y$$

- ① Choose prob distribution Bernoulli
 - ② Set $\lambda = f(x, \phi)$ ③ likelihood $\rightarrow P_2(y|x) = (1 - \operatorname{sig}(f(x, \phi)))^{1-y} \cdot \operatorname{sig}(f(x, \phi))^y$
- Negative log likelihood $L(\phi) = \sum_{i=1}^D - \left((1-y_i) \log(\operatorname{sig}(f(x_i, \phi))) + y_i \log(1 - \operatorname{sig}(f(x_i, \phi))) \right)$

MULTICLASS CLASSIFICATION

- ① $y \in \{1, \dots, k\}$ there is a prob. for each class, then k prob.
 $P_2(y=k) = \lambda_k$
 λ_i has to be in $[0, 1]$ and $\sum_{i=1}^k \lambda_i = 1$
- ② We set $\lambda_i \quad i=1 \dots k = f(x, \phi)$
Set the constraint using $\operatorname{softmax}_k(z) = \frac{\exp(z_i)}{\sum_{k=1}^k \exp(z_k)}$
- ③ NLL Loss
Likelihood $P_2(y=k|x) = \operatorname{softmax}(f(x, \phi))$
NLL $L(\phi) = - \sum_{i=1}^D \log \left[\operatorname{softmax}_{y_i}(f(x_i, \phi)) \right] =$
 $- \sum_{i=1}^D \left(f_{y_i}(x_i, \phi) - \log \left(\sum_{k=1}^k \exp [f_k(x_i, \phi)] \right) \right)$

CAPITOLO 4

MODEL - TRAINING

GRADIENT DESCENT

Find parameters that leads to lower loss values

- Gradient Descent

- Gabor Model and SGD

- SGD Scheduler LR, Momentum, Nesterov Momentum, Adam

- Backpropagation

- Computation Graph

- Parameter Initialization

- Vanishing & Exploding gradient

① Derivative of loss wrt ϕ

$$L(\phi) = \sum_{i=1}^D \ell_i$$

$$\frac{\partial L}{\partial \phi} = \begin{bmatrix} \frac{\partial L}{\partial \phi_0} \\ \vdots \\ \frac{\partial L}{\partial \phi_N} \end{bmatrix}$$

② Update the parameters

$$\phi = \phi - \alpha \frac{\partial L}{\partial \phi}$$

learning rate

linear regression

$$y = f(x, \phi) = \phi_0 + \phi_1 x$$

$$\text{Loss}(\phi) = \sum_{i=1}^D \ell_i = \sum_{i=1}^D (f(x, \phi) - y_i)^2 = \sum_{i=1}^D (\phi_0 + \phi_1 x_i - y_i)^2$$

Derivative of the loss

$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^D \ell_i = \sum_{i=1}^D \frac{\partial \ell_i}{\partial \phi}$$

this means

$$\frac{\partial \ell_i}{\partial \phi} = \begin{bmatrix} \frac{\partial \ell_i}{\partial \phi_0} \\ \frac{\partial \ell_i}{\partial \phi_1} \end{bmatrix} = \begin{bmatrix} 2(\phi_0 + \phi_1 x_i - y_i) \\ 2x_i(\phi_0 + \phi_1 x_i - y_i) \end{bmatrix}$$

STOCHASTIC GRADIENT DESCENT

Compute the gradient descent on a minibatch

Pros

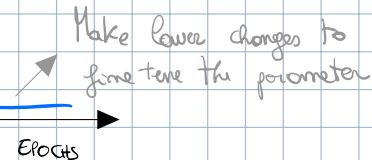
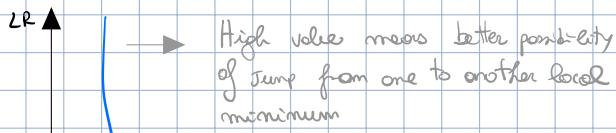
- Less expensive
- Escape local minima
- Allows generalization over unseen data

Update Rule

$$\phi_{t+1} = \phi_t - \alpha \sum_{i=1}^{B_c} \frac{\partial \ell_i[\phi_t]}{\partial \phi}$$

Batch instead of the entire dataset

SGD and Scheduler of LR



$$\text{Remember } LR = \alpha \quad \phi_{t+1} = \phi_t - \alpha \cdot \frac{\partial L}{\partial \phi_t}$$

SGD with Momentum Effect: Smooth out the trajectory to the minimum and makes the convergence fast

Add a Momentum Term m_t instead of ∂L

The new update rule

$$\phi_{t+1} = \phi_t - \alpha m_{t+1}$$

the overall direction depends from the previous momentum

$$m_{t+1} = \beta m_t + (1-\beta) \sum_{i=1}^B \frac{\partial l_i[\phi_t]}{\partial \phi}$$

$\in [0,1]$ weight that gives more power to the previous or current update Gradient Computed on the current batch

SGD NESTEROV ACCELERATED MOMENTUM

$$\phi_{t+1} = \phi_t - \alpha m_{t+1}$$

$$m_{t+1} = \beta m_t + (1-\beta) \sum_{i=1}^B \frac{\partial l_i[\phi_t - \alpha m_t]}{\partial \phi}$$

Effect provide a better path to reach the minimum

SGD - ADAM

Problem until now:

- big gradient lead to big parameters update
- choose a LR that makes the model stable is difficult

Solution 1:

Normalize the momentum by the scale of the gradient to avoid huge updates, then:

$$m_{t+1} = \frac{\partial L(\phi_t)}{\partial \phi} \quad v_{t+1} = \left(\frac{\partial L(\phi_t)}{\partial \phi} \right)^2$$

$$\phi_{t+1} = \phi_t - \alpha \frac{m_{t+1}}{\sqrt{v_{t+1}} + \epsilon}$$

Adam: Have an even smoother trajectory convergence

$$m_{t+1} = \beta m_t + (1-\beta) \frac{\partial L(\phi_t)}{\partial \phi}$$

$$v_{t+1} = \gamma v_t + (1-\gamma) \left(\frac{\partial L(\phi_t)}{\partial \phi} \right)^2$$

At each timestep m_{t+1} is $\tilde{m}_{t+1} = \frac{m_{t+1}}{1-\beta^{t+1}}$

$$v_{t+1} \text{ is } \tilde{v}_{t+1} = \frac{v_{t+1}}{1-\gamma^{t+1}}$$

$$\text{Update Rule} \quad \phi_{t+1} = \phi_t - \alpha \frac{\tilde{m}_{t+1}}{\sqrt{\tilde{v}_{t+1}} + \epsilon}$$

FITTING THE MODEL WITH BACK PROPAGATION

Let's consider

x : Multivariate input $(x_1 \dots x_m) \in \mathbb{R}^m$

ϕ : Parameters

h_1, h_2, h_3 Hidden layers

$$h_1 = \alpha (\beta_0 + \phi_0 x)$$

$$h_2 = \alpha (\beta_1 + \phi_1 h_1)$$

$$h_3 = \alpha (\beta_2 + \phi_2 h_2)$$

$$f(x, \phi) = \beta + \phi_3 h_3$$

loss for each sample l_i

loss for the entire dataset

$$L(\phi) = \sum_{i=1}^n l_i$$

Parameter update using SGD
 $\phi_{t+1} = \phi_t - \alpha \sum_{i=1}^B \frac{\partial l_i(\phi_t)}{\partial \phi}$

This algorithm is called Back propagation

It includes a forward pass

Computing the derivatives $\frac{\partial l_i}{\partial \phi_k}$ and $\frac{\partial l_i}{\partial \phi_k}$

A backward pass that uses the derivatives to update the weights

An Example 8 parameter for a function $f(x, \phi)$

$$\phi = \{\beta_0, w_0, \beta_1, w_1, \beta_2, w_2, \beta_3, w_3\}$$

3 non linear functions sin, exp, cos

$$f[x, \phi] = \beta_3 + w_3 \cdot \cos[\beta_2 + w_2 \cdot \exp[\beta_1 + w_1 \cdot \sin[\beta_0 + w_0 \cdot x]]]$$

$$\ell_i = (f[x_i, \phi] - y_i)^2$$

compute

$$\frac{\partial \ell_i}{\partial \beta_0, \partial \beta_1, \partial \beta_2, \partial \beta_3}$$

and

$$\frac{\partial \ell_i}{\partial w_0, \partial w_1, \partial w_2, \partial w_3}$$

$$\frac{\partial \ell_i}{\partial w_0} = -2 \left(\beta_3 + w_3 \cdot \cos[\beta_2 + w_2 \cdot \exp[\beta_1 + w_1 \cdot \sin[\beta_0 + w_0 \cdot x_i]]] - y_i \right) \cdot w_1 w_2 w_3 \cdot x_i \cdot \cos[\beta_0 + w_0 \cdot x_i] \cdot \exp[\beta_1 + w_1 \cdot \sin[\beta_0 + w_0 \cdot x_i]] \cdot \sin[\beta_2 + w_2 \cdot \exp[\beta_1 + w_1 \cdot \sin[\beta_0 + w_0 \cdot x_i]]].$$

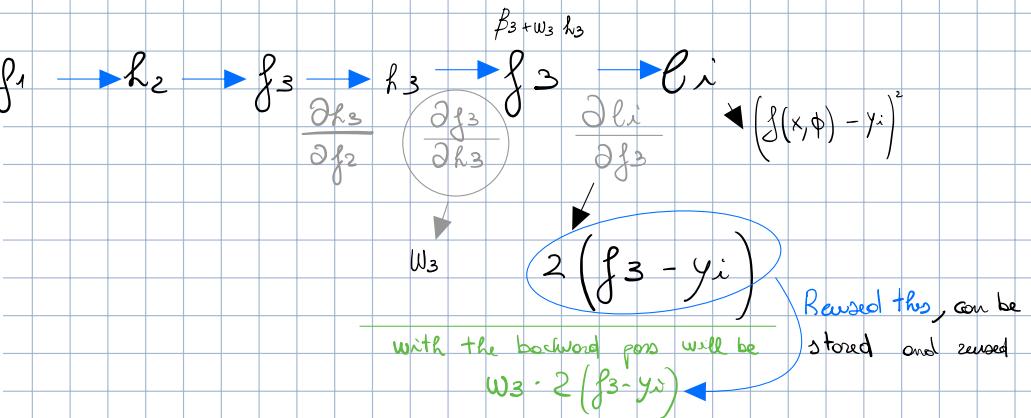
Computing derivatives you face a lot redundant computation

COMPUTATION GRAPH

Avoid the redundant computation in the computation of derivative using a graph

$$x_i \rightarrow f_0 \rightarrow h_1 \rightarrow f_1 \rightarrow h_2 \rightarrow f_2 \rightarrow h_3 \rightarrow f_3 \rightarrow h_3 \rightarrow \ell_i$$

$$\begin{aligned}\frac{\partial \ell_i}{\partial f_2} &= \frac{\partial h_3}{\partial f_2} \left(\frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right) \\ \frac{\partial \ell_i}{\partial h_2} &= \frac{\partial f_2}{\partial h_2} \left(\frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right) \\ \frac{\partial \ell_i}{\partial f_1} &= \frac{\partial h_2}{\partial f_1} \left(\frac{\partial f_2}{\partial h_2} \frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right) \\ \frac{\partial \ell_i}{\partial h_1} &= \frac{\partial f_1}{\partial h_1} \left(\frac{\partial h_2}{\partial f_1} \frac{\partial f_2}{\partial h_2} \frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right) \\ \frac{\partial \ell_i}{\partial f_0} &= \frac{\partial h_1}{\partial f_0} \left(\frac{\partial f_1}{\partial h_1} \frac{\partial h_2}{\partial f_1} \frac{\partial f_2}{\partial h_2} \frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3} \right).\end{aligned}$$



PARAMETER INITIALIZATION

Depend from where you start from a place where the variance of the pesos is small to do not oscillate too much or to low

$$\text{Consider } f_k = \beta_k + \Omega_k h_k$$

$$= \beta_k + \Omega_k \alpha(f_{k-1})$$

With Ω very small variance there will be no update.

Because RELU

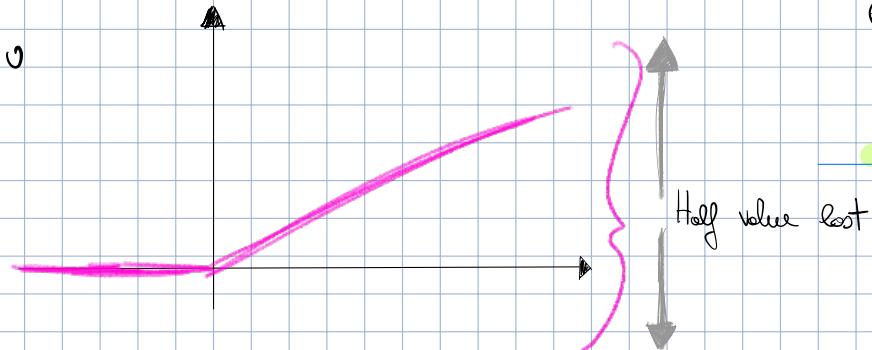
IMPOSE SOME VARIANCE: a good balance
 Ω initialized $\in N(y, \sigma^2)$

Going through layers becomes much smaller

SMALL WEIGHTS

leads to

VANISHING GRADIENT



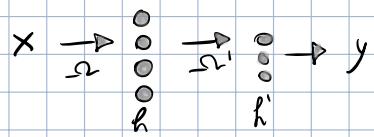
High variance σ^2 = Big Weights, huge jump in update thus leads to **Exploding Gradients**

VARIANCE CONSTRAINT : Initialization

$$\sigma^2 \Omega = \frac{2}{D_h}$$

If we have 2 adjacent hidden layer

Consider on initializ.



The total initialization

$$\sigma^2 \Omega' = \frac{2}{D_{h'}}$$

$$\sigma^2 \Omega = 4(D_h + D_{h'})$$

CAPITOLO 5

PERFORMANCE & REGULARIZATION

- Overfitting

- Regularization

- Explicit

- As a prior

- L2

- Early Stopping

- Dropout

- Data Augmentation

SOURCE OF ERRORS

NOISE

Noise on the data

BIAS

Chosen model not flexible/powerful enough for learning the true function

$$\mathbb{E}_{\mathcal{D}} [\mathbb{E}_y [L[x]]] = \underbrace{\mathbb{E}_{\mathcal{D}} [(f[x, \phi[\mathcal{D}]] - f_\mu[x])^2]}_{\text{variance}} + \underbrace{(f_\mu[x] - \mu[x])^2}_{\text{bias}} + \sigma^2.$$

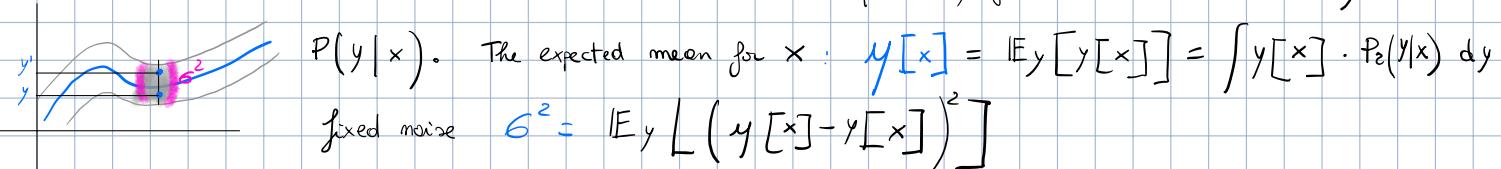
Variance in the learning process leads to different local minimum. Not enough variant dataset will not be a good set for the model to learn the ground truth distribution

VARIANCE

NOISE, BIAS, VARIANCE, MATHEMATICAL PROOF

Consider 1D Regression problem. Data x have noise with variance σ^2 .

For each x there might be multiple y value. We then consider a probability for x to have a certain y



REDUCE ERRORS

NOISE:

Nothing to do.

Noise cannot be removed, is intrinsic of the data generation process

REDUCE THE VARIANCE

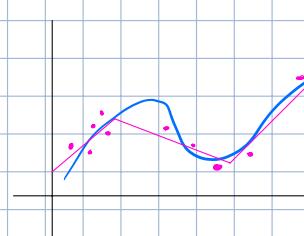
- Add data to train from makes the model more confident on the prediction that is giving. less data, more uncertainty.
- Noise is symmetric, add data makes the model even more robust w.r.t data noise

THE BIAS-VARIANCE TRADE OFF

Chosen a certain type of model to fit on training data, we are imposing a bias.

We are essentially limiting the flexibility of the model. To reduce the bias we can use a more powerful model.

With a lot of sample the model will try to get the lowest error for each of them.



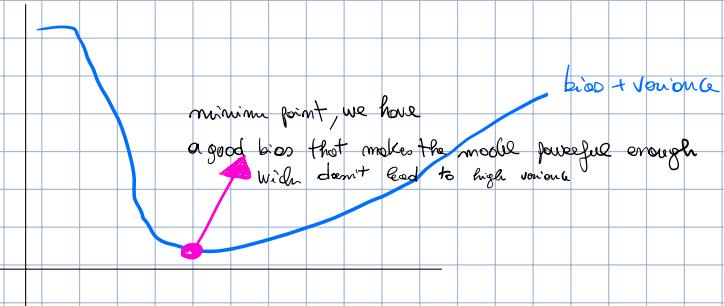
Some hidden units in each plot

Increasing the picking
We have many more
functions to pick from

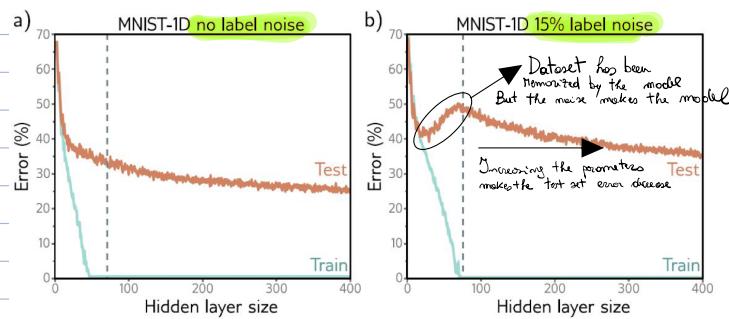
= We are increasing
the Variance



The capacity issue pick: How much data to use to fit the model?



DOUBLE DESCENT



The double descent

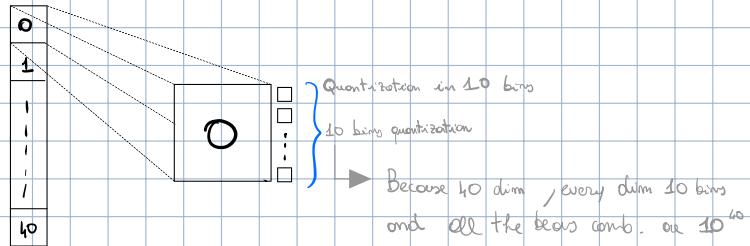
- ① Training, the model learns
 - ② Adding parameters, the model memorize the data, but make mistake because of the noise
 - ③ Mistake grows
 - ④ Still add parameter
Mistake decrease
- mistake on the test set which does not contain the noise

COURSE OF DIMENSIONALITY

In low dimensionality is easy to know if two objects are close or far. Increasing the dimension everything starts to get complicated.

Add hidden unit let the model to be powerful.

BUT still, we have to consider high dimensional data (40 dim)

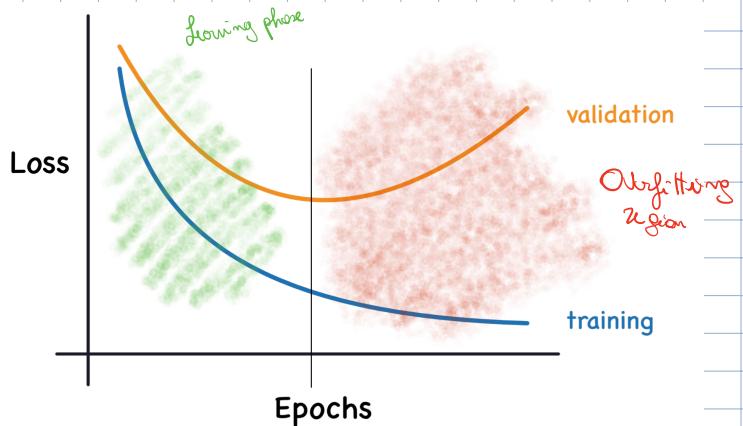
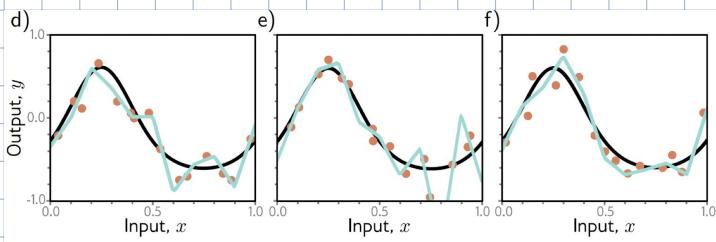


With 10^k example (10^5)
we are observing a small portion of this high dimensional space w.r.t 10^{40}

There will be a sample every 10^{35} bins

OVERFITTING

Increasing the number of hidden units (Universal theorem approximator) lead the model to perfectly represent the data. But this doesn't mean that the model is actually learning the true distribution underlying the data.



REGULARIZATION

Explicit Regularization

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left(\sum_{i=1}^D l_i + \lambda \cdot g(\phi) \right)$$

Importance to give to this

given a weight value return a value to influence then the loss. Big value when ϕ has value that is not perform

Regularization as a Prior

Remember the Maximum likelihood

$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \prod_{i=1}^D P(y_i | x_i)$$

$$= \underset{\phi}{\operatorname{argmax}} \prod_{i=1}^D P_2(y_i | f[x, \phi])$$

With regularization term

$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \left[\prod_{i=1}^D P_2(y_i | f(x, \phi)) P_e(\phi) \right]$$

L2 Regularization

Explicit (L2 Regularization)

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \sum_{i=1}^D l_i(x_i, y_i) + \lambda \sum_j \phi_j^2$$

$g(\phi)$ here is penalizing weights with high values

As a prior L2 Regularization

Consider prior dist₂. over the weights $\in N(\mu=0, \sigma^2)$ Then:

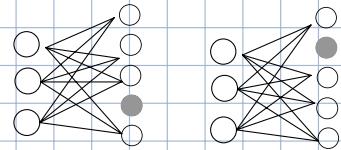
$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \prod_{i=1}^D P_2(y_i | f(x, \phi)) \cdot \prod_{j=1}^J \text{Norm } \phi_j [0, \sigma^2]$$

EARLY STOPPING

Stop the training procedure before that overfits the data

DROPOUT

Some hidden layer's unit weight are set to 0, both in coming and outgoing. We may see this as an ensemble of sub-networks



Note: dropout to works need a lot of labeled training example

CAPITOLO 6

CNN - RESNET

- Principles Idea
- CNN Geometric view
- Convolution Theorem
- Discrete Conv
- Conv Layer, Receptive Field
- Resnet : the depth
- Exploding gradient & Batch Norm
- U-Net
- Highway Networks

PRINCIPLES & IDEA

MLP are universal approximator BUT

Complexity Difficult optimization limited Generalization

We would like to leverage structural similarities

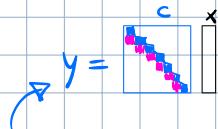
SELF SIMILARITY TRANSLATION INVARIANCE DEFORMATION INVARIANCE HIERARCHY COMPOSITIONALITY

CNN: GEOMETRIC VIEW

let's define a Continuous Convolution

$$(f * g)(x) = \int_{-\infty}^{+\infty} f(t) g(x-t) dt$$

Discrete time Convolution $(x * w) = \sum_{k=0}^{m-1} w_{i-k} \text{ mod } m X_k$
 (2 vector x, w)



Or as a linear operator $y = C(w)x$ where $C(w)$ is a Circulant Matrix
 (stack the vector w m times with shift)

let's say something regarding Circulant Matrix

With W vector $[0, 1, 0, \dots, 0]$ we have a Shift-Matrix (a type of Circulant Matrix)

$$y = \begin{matrix} y \\ 2 \\ 3 \\ 4 \end{matrix} = \begin{matrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{matrix} \begin{matrix} x \\ 1 \\ 2 \\ 3 \\ 4 \end{matrix} \quad \left. \begin{matrix} \text{Shift right} \\ \text{by } \pm \end{matrix} \right\} \quad \begin{matrix} X \rightarrow 1, 2, 3, 4 \\ Y \rightarrow 2, 3, 4, 1 \end{matrix}$$

PROPERTY

Matrix is circulant
 if (matrix) with shift

C is circulant if
 $C(w) \cdot S(w) = S(w) \cdot C(w)$

FACT: If Matrices commute, they have the same Eigenvector: they are joint diagonalizable
 The eigen values of $C(w)$ is the Fourier Transform of w

Formally: Convolution in the time domain (sliding dot product between kernel and signal)
 is the same of multiplication in the frequency domain

THE CONVOLUTION THEOREM

Compute the convolution between $x * w$ in 2 ways:

Spatial Domain

Take a Circulant Matrix $C(w)$
 and apply it to x

Spectral Domain

Compute Fourier Transform of x (signal)
 Compute FT of w (kernel)
 Multiply FT x with FT w
 Compute the inverse FT of the multiplication (go back in time domain)

DISCRETE 2D CONVOLUTION

$$(f * g) [m, n] = \sum_k \sum_e g[k, e] \cdot g[m - k, n - e]$$

As in 1D case we would have
 $f(t) - g(x-t)$

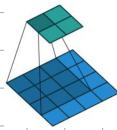
on image: $m \times m$ matrix

this is like a sliding window over the image



• Padding

- Zero padding (valid)



- Full Padding (same)

The size of the output = Size of input

- Stride

- Dilation (skip within a convolution certain offset of the current position)

- DIMENSION AFTER CONVOLUTION

$$\begin{array}{ll} W_{in} & F \text{ Size of Kernel} \\ H_{in} & P \text{ padding} \\ & S \text{ stride} \end{array}$$

$$W_{out} = \frac{W_{in} - F + 2P}{S} + 1$$

$$H_{out} = \frac{H_{in} - F + 2P}{S} + 1$$

Convolutional Layer of Receptive Field

Convolutional layer

$$\left. \begin{array}{l} \text{Input} \\ \text{Convolve Input} \\ \text{Add bias} \\ \text{Activation function} \end{array} \right\} = h_i = a \left[\beta + w_1 x_{i-1} + w_2 x_i + w_3 x_{i+1} \right]$$

$$= a \left[\beta + \sum_{j=1}^3 w_j x_{i+j-2} \right]$$

$F = \text{kernel size} = 3$

Fully connected VS convolution

$$\begin{array}{l} 6 \text{ size of input} \\ 6 \text{ hidden unit} \\ (6 \times 6) + 6 = 42 \text{ param} \\ \text{Bias} \end{array}$$

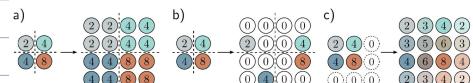
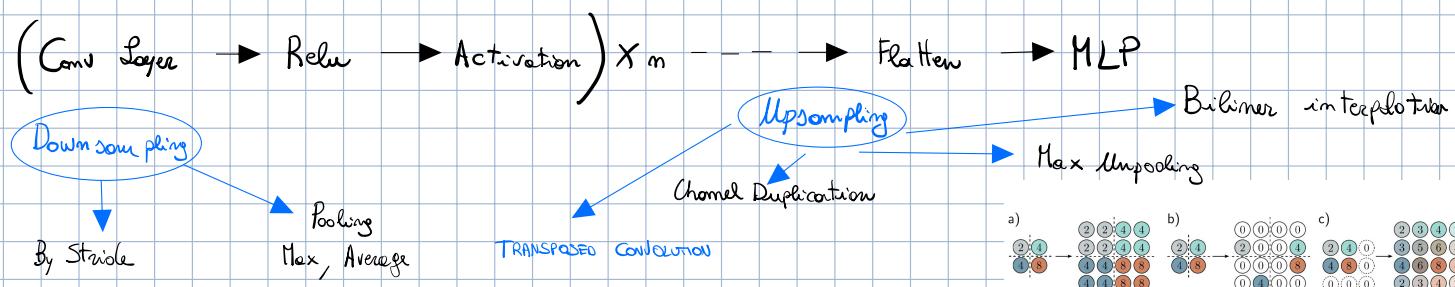
$$\begin{array}{l} F \text{ kernel size} = 3 \\ (3 \times 1 + 6) \times 1 = 9 \\ \beta \end{array}$$

Number of parameters

$$20 (F \times F \times C_{in} \times C_{out}) + C_{out}$$

$$1D (F \times (C_{in} + Bias)) \times \text{Number of filters}$$

Convolutional Network Structure



PARAMETERS OPTIMIZATION AND GRADIENT FLOW

RESNET - BATCH NORM

The problem: More layer, better approximation (Universal Approximation Theorem)

But more layer = more difficult search of parameters

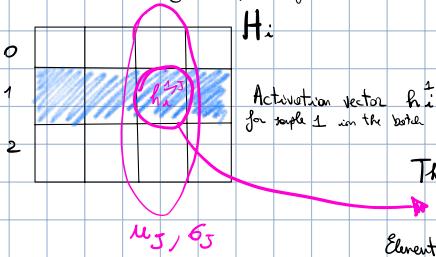
Network do not converge
Exploding / Vanishing Gradients

BATCH NORMALIZATION

Adaptive parameterization

$$x \xrightarrow{w_1} f_1 \xrightarrow{w_2} f_2 \dots \xrightarrow{w_k} f_k \quad h_i = \underbrace{\sum_{j=1}^k w_j}_{\text{num layers}} \cdot \underbrace{h_i}_{\text{unit of weights}}$$

Now, let's suppose $H_i = \langle h_i^1 \dots h_i^k \rangle$ the activation of layer i for the minibatch of samples from 1 to k



$$H_i^1 = \frac{h_i^1 - \mu_j}{\sigma_j}$$

The normalization
 $\hat{h}_i^{1,j} = (h_i^{1,j} - \mu_j) / \sigma_j$
Element 1 of the both j element of the hidden vector
for the hidden layer i

While training

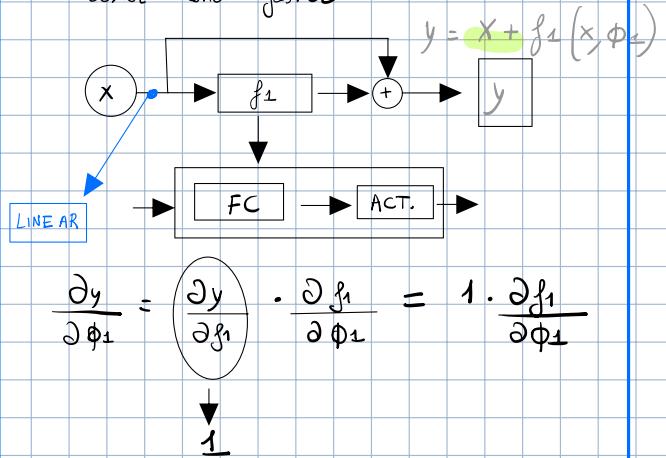
$$y = \frac{1}{k} \sum_i H_i$$

Mean vector y computed for each layer i

$$\begin{bmatrix} \mu_1 & \mu_2 & \dots & \mu_k \\ \vdots & \vdots & \ddots & \vdots \\ \mu_k & \end{bmatrix}$$

RESIDUAL BLOCK

The introduction of residual connection make the problem of parameters optim. easier and faster

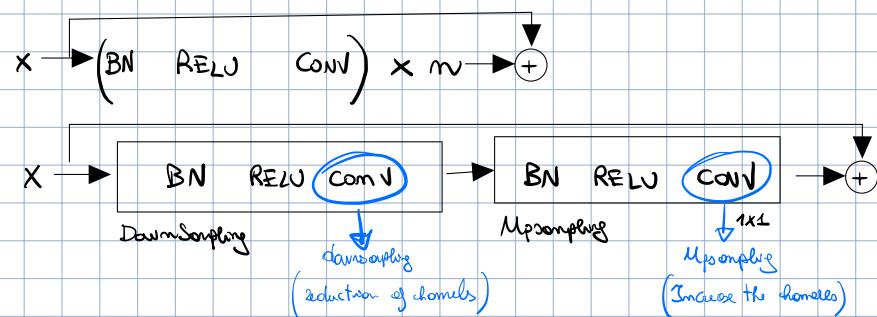


Both y and σ are trainable parameters

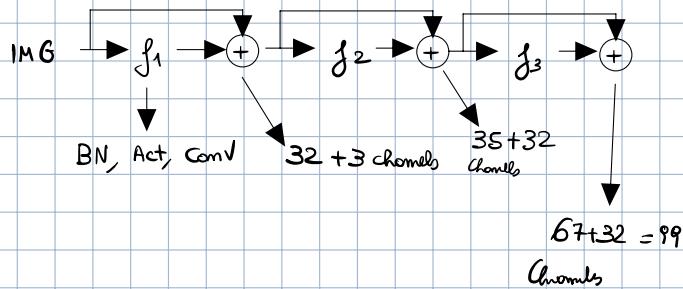
$$\sigma = \sqrt{s + \frac{1}{k} \sum_i (H_i - y)^2}$$

RESNET + BATCH NORM ARCHITECTURE

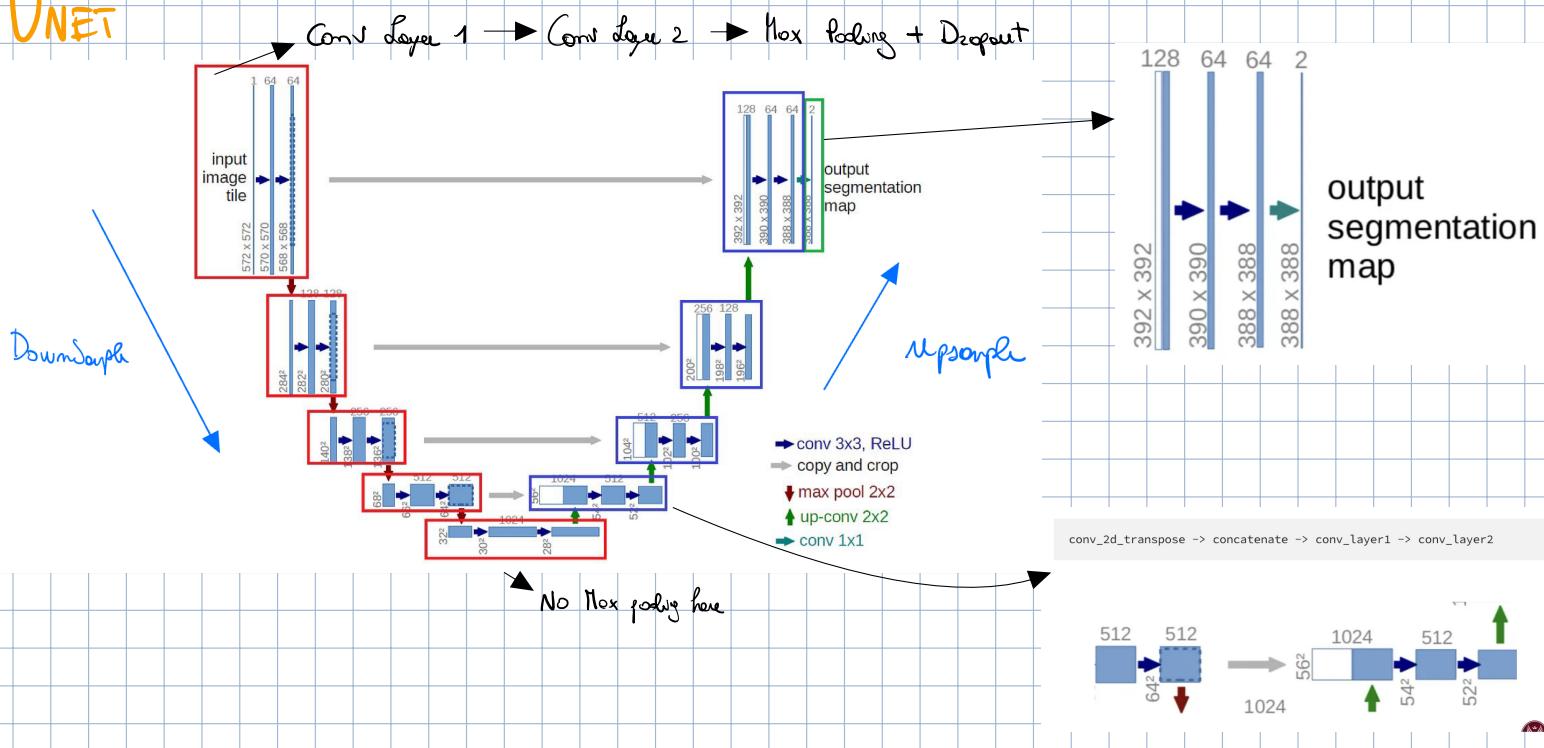
RESNET



DENSE NET



UNet



HIGHWAY NETWORKS

$$\text{given } y = f(x, w_f)$$

on highway nn also have other 2 transformation of the input x : $T(x, w_T)$

$$y = f(x, w_f) \cdot T(x, w_T) + x \cdot C(x, w_C)$$

Usually set as $1 - T(x, w_T)$

$$C(x, w_C)$$

CAPITOLO 7

SELF SUPERVISED LEARNING

Siamese NN

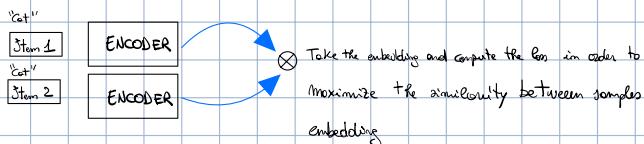
- Neural Collaborative Filtering
- Self Supervision, Proxy Task
- Contrastive Learning, Triplet Loss

- SIM CLR
- BYOL + Evolution
- BARLOW TWINS + Evolution
- Deep Metric

SIAMESE NEURAL NETWORK

Two or more identical subnetworks.

IDEA: Some items should have some representation

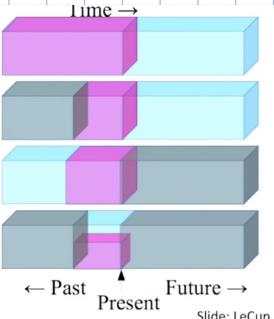


Use data as label for themselves to learning their represent.

SELF SUPERVISION PROXY TASK

We don't care about performance of specific application, we want to learn input representation such that we can then use them in downstream task. We learn by the use of PROXY TASK.

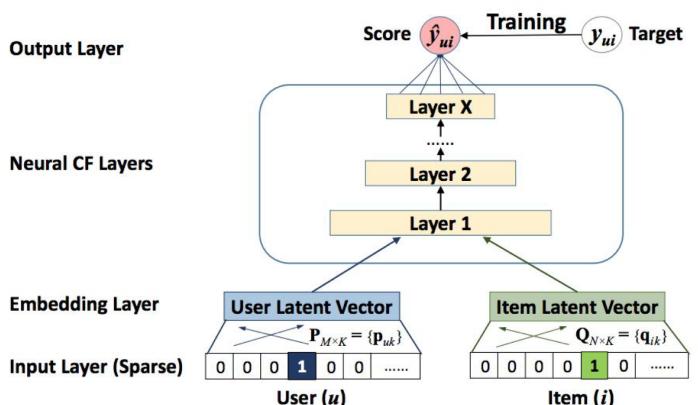
- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the **occluded** from the **visible**
- ▶ Pretend there is a part of the input you don't know and predict that.



NEURAL COLLABORATIVE FILTER

User = {User 1 -- User n} Item = {item 1 -- item m} Preferences = {User i, item j}

Goal: get a ranked list of items (**unseen**) for a certain user
 $\hat{y}_{ui} = f(u, i | \theta)$ predict the score of interaction for (u, i)



CONTRASTIVE LEARNING

IDEA: Similar samples close, while dissimilar ones far apart

Sample from some class \rightarrow Similar embeddings

$x_1, x_2 \in \text{class } X \rightarrow$ minimize embedding distance

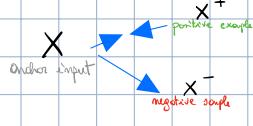
$x_1 \in X, x_2 \in Y \rightarrow$ maximize embedding distance

With a lower bound distance between different classes

CONTRASTIVE LOSS (x_1, x_2)

$$L(X_i, X_s, \theta) = \mathbb{I}[y_i = y_s] \|f_\theta(x_i) - f_\theta(x_s)\|_2^2 + \mathbb{I}[y_i \neq y_s] \max(0, \epsilon \|f_\theta(x_i) - f_\theta(x_s)\|_2)$$

TRIPLET LOSS



$$L_{\text{Triplet}}(x, x^+, x^-) = \sum_{x \in D} \max(0, \|f(x) - f(x^+)\|_2^2 - \|f(x) - f(x^-)\|_2^2 + \epsilon)$$

To make this loss decrease the positive distance has to be at least the negative distance + ϵ

$$L(x, x^+, x^-) = 2 - \lambda + \epsilon = 2 - \max(0, -2) = 0 \text{ Good}$$

IF $x^+ \sim x$
 $x^- \sim x$ $\max(0, 1-2) = 2 \text{ Bad}$

INGREDIENTS

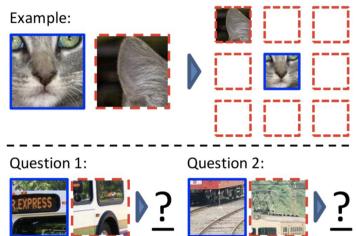
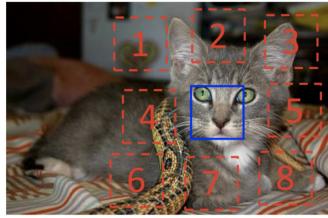
- Create positive samples with data augmentation
- Large Batch size to get a good negative samples coverage
- Create a proxy strategy to get hard negative (by labels or by an algorithm)

IMAGE BASED SELF - SUPERVISED LEARNING

WORKFLOW

- Training with patent task (a proxy task)
- Use a layer to train a logistic regression over the ImageNet classes evaluating the representation given

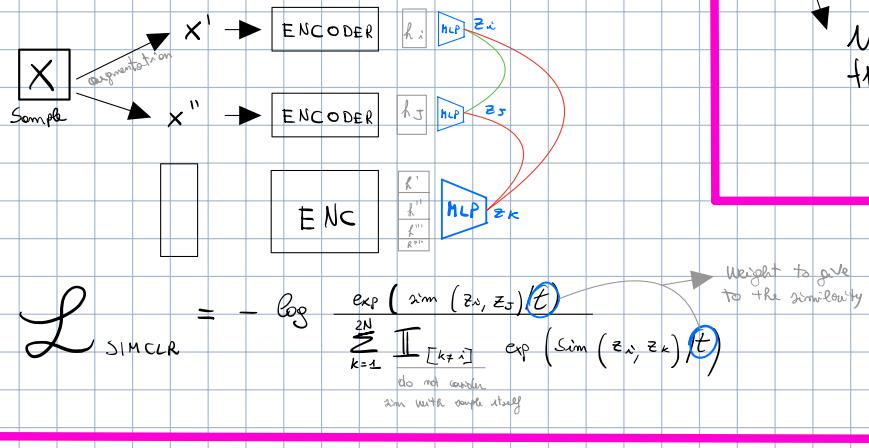
Have some representation for patches corrupted but starting from an original one (Exemplar-CNN)
Given: the degree of rotation of an image
Predict Relation between patches



$$X = (\text{cat face}, \text{tail}); Y = 3$$

SIM CLR

(x, x') , (x, x'') positive samples all the rest are negative



Use the 9 patches to guess which permutation of the patches has been applied

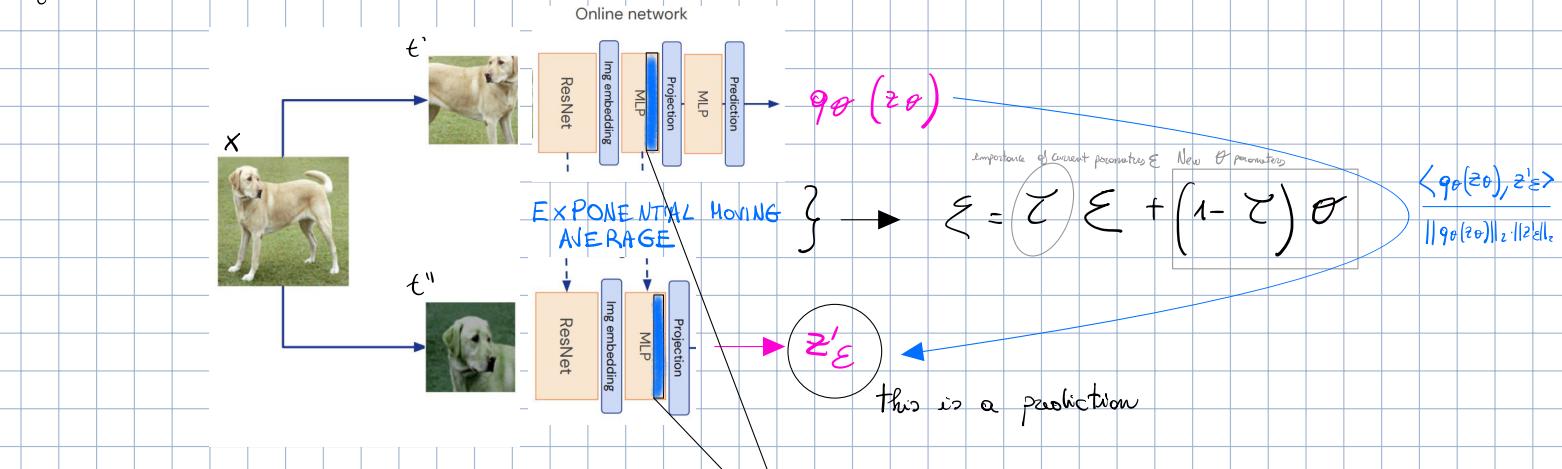
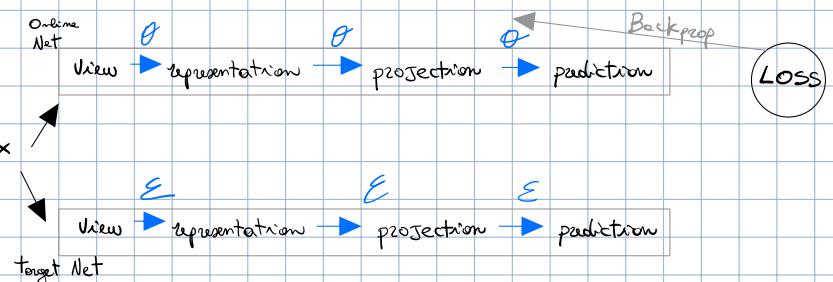
BYOL

Can we get rid of negatives?

Getting 2 identical network with 3 stages

Train online net and update the θ parameter with GD
But update the target with MOVING AVERAGE

This avoid the 2 represent. embedding from online and target network to not be equal hence the function to be constant



LOSS: Minimize the distance between representations given by source network with representation given by target network

BATCH NORM
Create contrastive samples

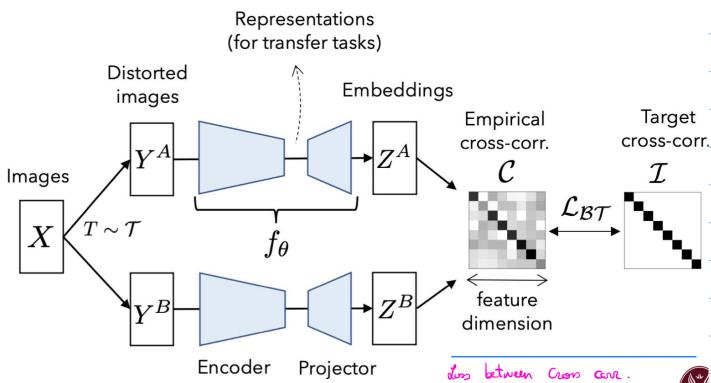
$$\mathcal{L}_{\theta, \varepsilon} = \| q_{\theta}(z_{\theta}) - z^{\dagger} \varepsilon \|_2^2 = 2 - 2 \cdot \frac{\langle q_{\theta}(z_{\theta}) - z^{\dagger} \varepsilon \rangle}{\| q_{\theta}(z_{\theta}) \|_2 \cdot \| z^{\dagger} \varepsilon \|_2}$$

} Magnitude of 2 vectors }
} a kind of cos similarity }

BARLOW TWINS

Measure cross-correlation matrix between the embeddings of two identical networks. GOAL \rightarrow Make this matrix close to the identity

EFFECT: the embedding of distorted version of the same sample have to be similar



Gross Correlation matrix between the output of two identical networks along the batch dimension

$$\mathcal{L}_{BT} = \sum_i (1 - C_{ii})^2 + \lambda \sum_i \sum_{j \neq i} C_{ij}^2$$

Invariance term Redundancy Reduction term

Try to make the off-diagonal elements equal to 0 output units contain non-redundant information for the sample

Cross Correlation Matrix computed as:

$$C_{ij} = \frac{\sum_b z_{b,i}^A \cdot z_{b,j}^B}{\sqrt{\sum_b (z_{b,i}^A)^2} \sqrt{\sum_b (z_{b,j}^B)^2}}$$

SUPERVISED DEEP METRIC LEARNING

Again, put sample of the same class near in the latent space using a feature extractor $f_{\theta}: \mathbb{X}^D \rightarrow \mathbb{R}^N$.

TECHNIQUE \rightarrow Contrastive learning

Expansion issue: it is not guaranteed that if sample $s \in Y_i$ the distance between sample is small

Sampling issue: How to define negative samples?

Center loss: Take the representation, use a classifier using softmax.

Add a term to loss to ensure that samples are pushed toward the centroid of the class

$$\mathcal{L}_{center} = \mathcal{L}_{softmax} + \frac{\lambda}{2} \sum_i^N \| z_i - c_i \|_2^2$$

} distance from the centroid of class sample i }

Sphere Face: Modify the softmax to allow a different choice of the "center" in the center loss

Consider θ_{y_i} as the angle between z_i and centroid of class y

Starting from the

$$\text{Softmax } L_{\text{mod softmax}} = -\frac{1}{N} \sum_i^N \log \frac{\exp(W_{y_i}^T z_i + b_{y_i})}{\sum_j^N \exp(W_j^T z_i + b_j)} = -\frac{1}{N} \sum_1^N \log \frac{\exp(\|z_i\| \cos(\theta_{y_i}))}{\sum_j^N \exp(\|z_i\| \cos(\theta_{y_j}))}$$

We assign to sample i the class y if the projection z_i to the center of class y is the largest or if the cosine angle is the smallest between all the centroids

CAPITOLO 8

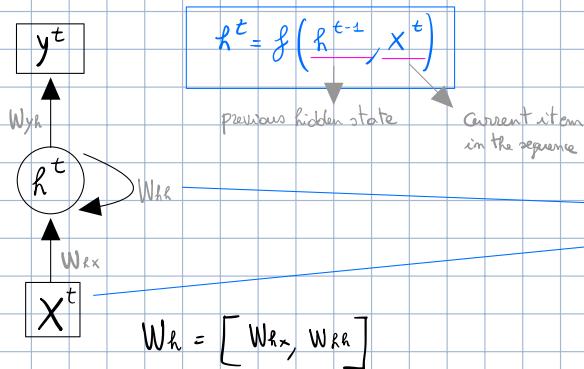
TRANSFORMERS

RNN

$x = \text{Input}$, composed of sequence data

Only one architecture to each item of the sequence. The weight are shared for all the tokens.

RNN can in theory deal with an arbitrary long sequences



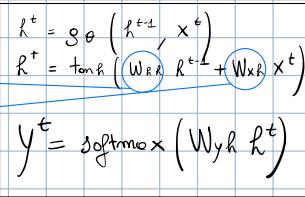
An Example of RNN

Starting from $h^0 = f(h^{t-1}, x^0)$

$$h^t = g_\theta(h^{t-1}, x^t)$$

$$h^t = \tanh((W_{RA} h^{t-1} + W_{RB} x^t))$$

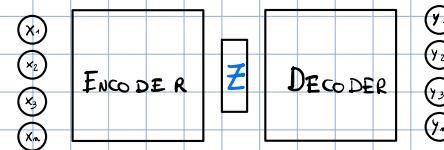
$$y^t = \text{softmax}(W_{yh} h^t)$$



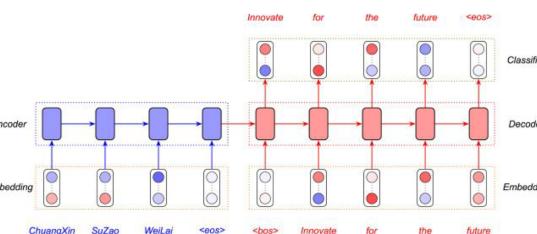
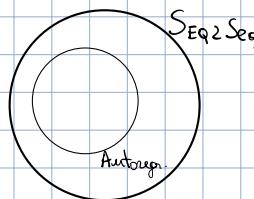
SEQUENCE TO SEQUENCE

Many to one \rightarrow Classification

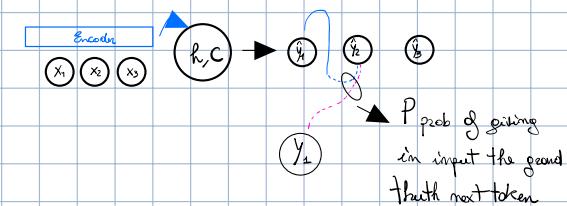
Many to many One to Many



AUTOREGRESSIVE MODELS



TEACHER FORCING



WORD 2 VEC



CBOW

Product target word from the context

$P(\text{Is} | \text{The, Cat, On, The})$

$$P(w_t | w_{t-h} \dots w_{t-1}, w_{t+1} \dots w_{t+h}) = \text{softmax}(w_t U V^T) \quad U, V \in \mathbb{R}^{m \times d}$$

ISSUE

Static word embedding: the embedding vector does not change for ambiguous words

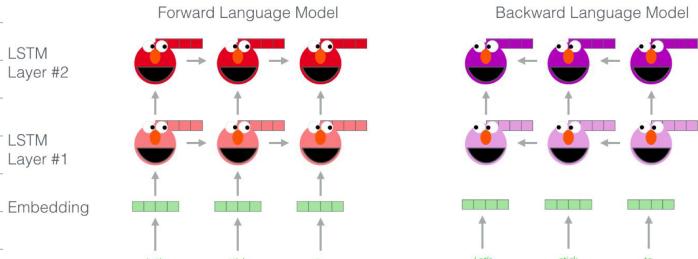
SKIPGRAM

The, Cat, Is, On, The, Table
Predict the context from the target

$P(\text{The} | \text{Is})$
 $P(\text{Cat} | \text{Is})$
 $P(\text{On} | \text{Is})$

ELMO

Embedding of "stick" in "Let's stick to" - Step #1



Embedding of "stick" in "Let's stick to" - Step #2

1- Concatenate hidden layers

2- Multiply each vector by a weight based on the task

$\times s_2$

$\times s_1$

$\times s_0$

3- Sum the (now weighted) vectors

ELMo embedding of "stick" for this task in this context

Forward Language Model

stick

Backward Language Model

stick



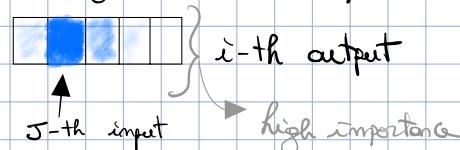
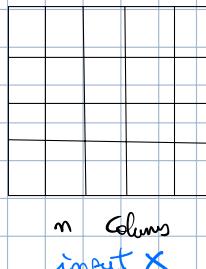
THE ATTENTION MECHANISM

THE ATTENTION: Association for each output a vector long as the input sequence
let's take a vector for one output

$$\text{Formally Input } x = [x_1 \dots x_n] \quad y = [y_1 \dots y_m]$$

The attention is a matrix

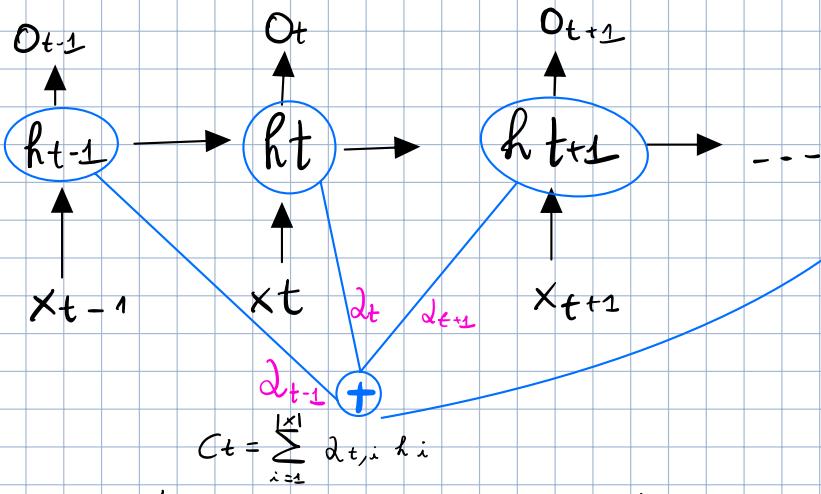
m rows
output y



$|x| = m$

$|y| = m$

RNN ENCODER



DECODER

the decoder works on

s_t that is a function of
 c_t

c_t context = sum of hidden states of
the encoder weighted
by the attention score α

α represent the alignment

"Align the output to the input"
of the output token

$$\text{SCORE}(s_t, h_i) = \sqrt{\alpha} \tanh \left(W_a [s_t, h_i] \right)$$

to be learnt

α is a kind of multiclass classification

$$\text{softmax}(y_t, x_i)$$

the scores α are learnt by another NN

SCALED DOT PRODUCT ATTENTION

$$\text{SCORE}(s_t, h_i) = \frac{s_t^T h_i}{\sqrt{m}}$$

normalized by seq. length

THE SELF ATTENTION

Input $X = [x_1 \dots x_m]$ x_i is a vector (input embedding)

The self-attention: given each word which are the other word in the same sentence that contribute the most to explain the semantic of that word

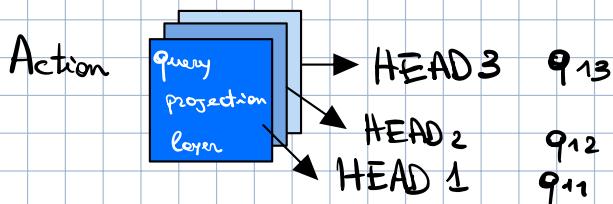
$$\text{Attention}(Q, k, v) = \text{softmax} \left(\frac{Q k^T}{\sqrt{m}} \right) \quad \begin{matrix} \checkmark & \text{scaled dot product} \\ // \end{matrix}$$

Example of self attention "Action gets result"

	query	key	value	score	score / \sqrt{m}	Softmax	Softmax * V	Sum
Action	q_1	k_1	v_1	$q_1 \cdot k_1$	$q_1 \cdot k_1 / 8$	x_{11}	$x_{11} \cdot v_1$	z_1
gets		k_2	v_2	$q_1 \cdot k_2$	$q_1 \cdot k_2 / 8$	x_{12}	$x_{12} \cdot v_2$	
Result		k_3	v_3	$q_1 \cdot k_3$	$q_1 \cdot k_3 / 8$	x_{13}	$x_{13} \cdot v_3$	

SCALED DOT Product

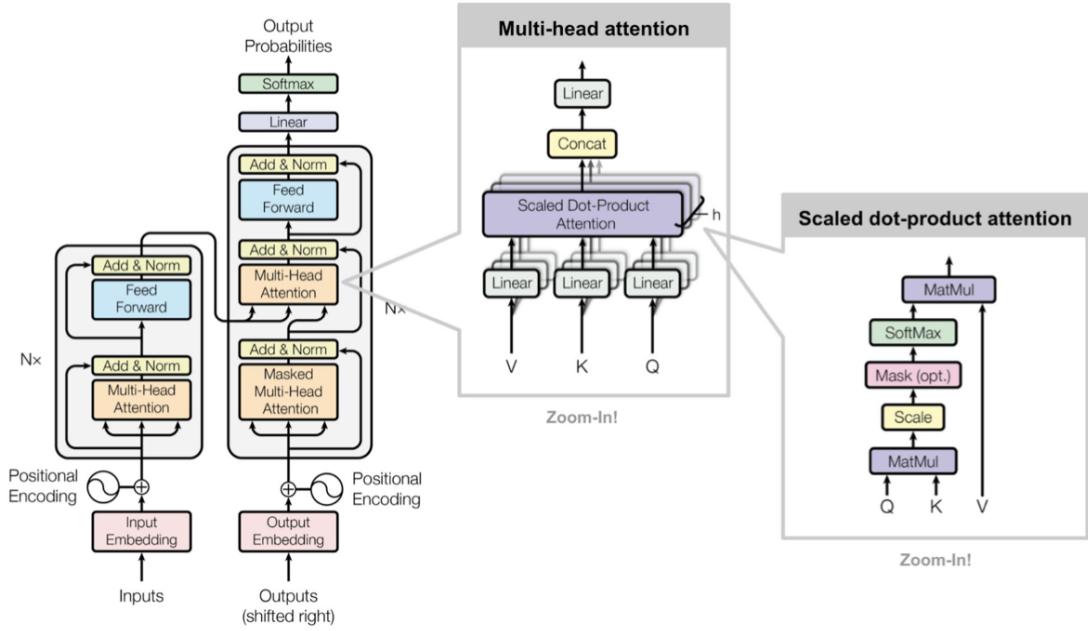
MULTI HEAD ATTENTION



$$\text{Multihead}(Q, k, v) = \left[\text{head}_1 \dots \text{head}_h \right]$$

head $i = \text{Attention} [Q W_i^Q, k W_i^k, v W_i^v]$

TRANSFORMER



BERT Bi-directional Encoder Representation from Transformer

Multi-layer bidirectional transformer

Base 110 M / Large 340 M

Proxy Tasks

Masked Language Modeling

Some input is masked and the model has to predict the right token

80% [MASK] token
10% Random token
10% Some token

Next Sentence Prediction

PRE - TRAINING

[CLS] + Input Embedding + Positional Emb + Segment embedding + [SEP] + Input Embedding + Positional Emb + Segment embedding
(Wordpiece Embedding) (Sentence 1) (Wordpiece Embedding) (Sentence 1)

FINE - TUNING (None Entity recognition, Grounded common sense inference, Textual Entailment, Question Answering) → Sentence 1 + Continuation → FFN → Softmax → Label

If head, If token, 1 Hidden state → Hidden state is a tensor 12/24, seq_len, 768/1024 → Pooling → FFN → Label

TRANSFER LEARNING : Neural IDBS

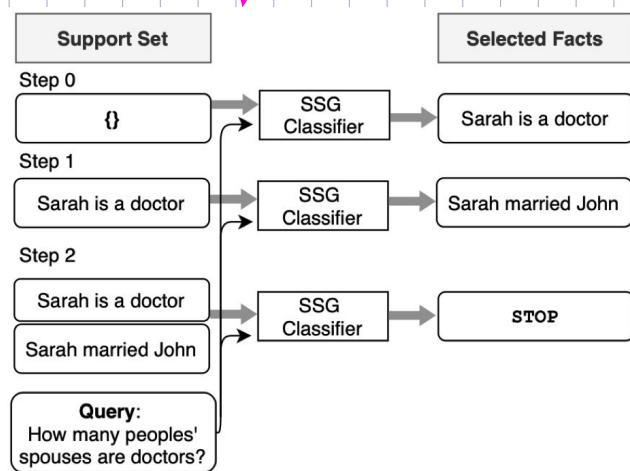
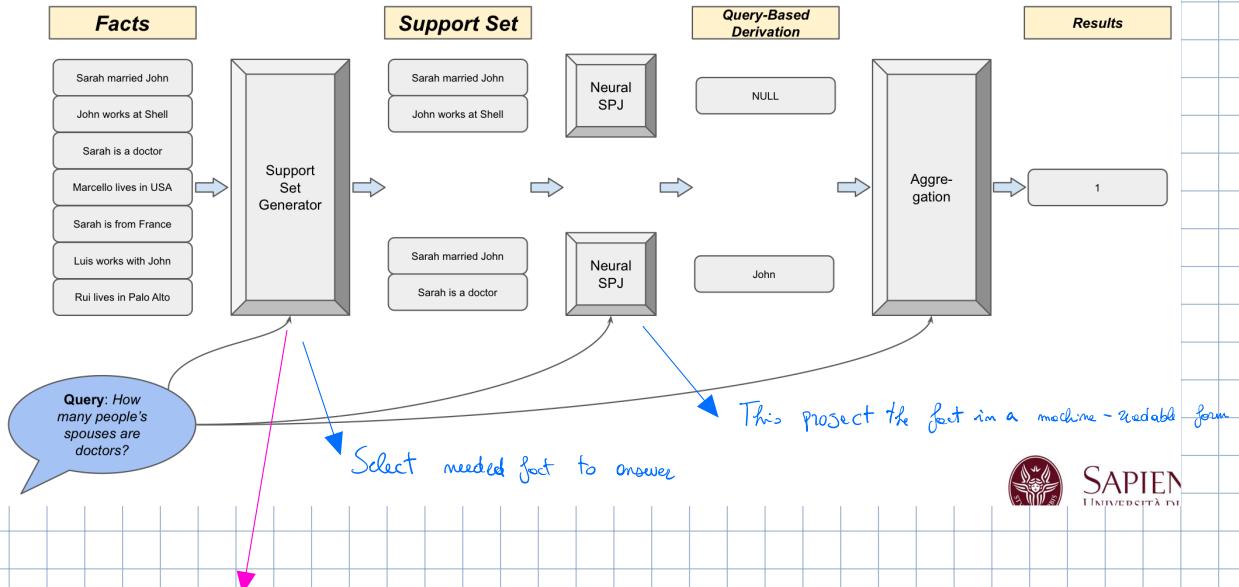
Make use of large pretrained model to adapt its ability on performing downstream task.

Convert or basically replace the SCHEMA of the database with the only unstructured content.

Neural Query Processing

Given a query and a small number of facts from the database, can the T5 accurately answer queries that are posed in natural language, whose answer may require projection (i.e., extracting part of a sentence), join, and aggregation?

This can be used by a Neural Processor which can interpret the content and reply to answer based on that content.



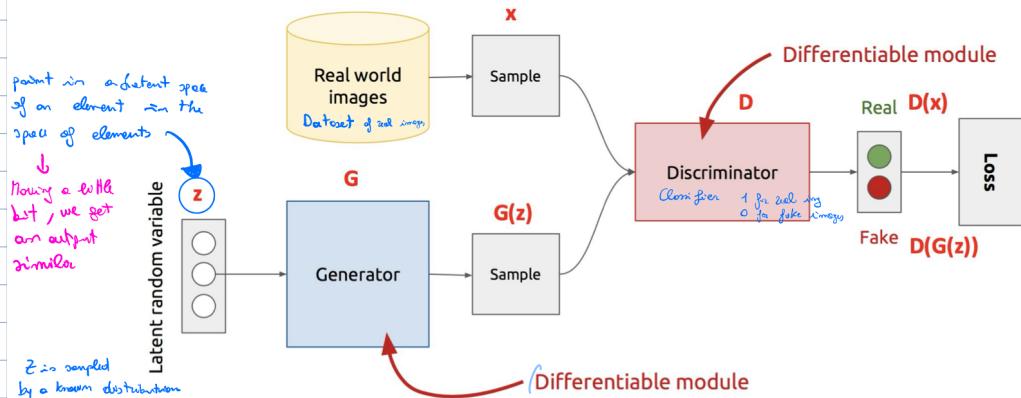
CAPITOLO 9

GENERATIVE AI

GAN

Considering the prob. of $X \sim P_{\theta}\{X\}$. We assume the generation of $x \in X$ conditioning on a variable z , $P_{\theta}\{x|z\}$ and z comes from a normal distib.

- Discriminative vs Generative
- GAN, Lotter, Convergence, Mode collapse
- Wasserstein GANs
- Autoencoders, Denoising Autoencoders
- VA, ELBO, Reparameterization Trick



GAN TRAINING : SETUP

for epoch

Train discriminator, while the generator is fixed

for k step:

- obtain $P_{\theta}(x)$ sampling m samples from Dataset
- obtain m samples generated

$$\text{MAX}_D \frac{1}{m} \sum_{i=0}^m \log D(x_i) + \frac{1}{m} \sum_{i=0}^m \left[\log (1 - D(G(z_i))) \right]$$

Train generator

- obtain m generated samples

$$\text{MIN}_G \frac{1}{m} \sum_{i=1}^m \log [1 - D(G(z_i))]$$

this is what the generator want to minimize

if $D(G(z)) = 1$:

the discriminator says that the image is fake, then

$\log \rightarrow \infty$

while if $D(G(z))$ miss classifies the sample, saying that is real then $\log \rightarrow 0$

TRAIN GAN

- 1 Freeze Generator, Train discriminator
- 2 Freeze Discriminator, Train generator

ALL MATHEMATICAL PROOF FOR GAN

- The value of generator depends on the discriminator

This is a ZERO-SUM GAME
(if one win, the other lose)

When the value of the discriminator is low (is committing error), the value of the generator is high

When the value of the discriminator is high means that the generator is no able to fool it, hence the discriminator detect all the fake images

$$\begin{matrix} \text{MIN} & \text{MAX} \\ G & D \end{matrix} \quad V(D, G)$$



Nash Equilibrium: the moment of convergence of a GAN
When is the moment in which the generator is working at best?

Basically when is totally uncertain about a sample

$$D(x) = \frac{1}{2} \quad \text{This means that } P_{\text{data}}(x) = P_{\text{gen}}(x) \quad \text{the generator has learnt the distribution of the data}$$

1 Real, 1 Fake
MAX UNCERTAINTY

there is a possible set of strategies (in this case the set of weight for the discriminator and for the generator) such that changing one strategy will make worse one of the strategy in the set.

Obtain the best Discriminator

$$\underset{D}{\text{MAX}} \quad V(G, D) = \mathbb{E}_{x \sim P_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim q(z)} [\log (1 - D(G(z)))]$$

Maximization part

Express the prob. of x to come from P_{data} or P_D

A pre defined distib such as $N(0, \sigma)$

Minimization part

Basically this means that we want to maximize the probability for the discriminator to predict $x \sim P_{\text{data}}$ and minimize the probability that discriminator predicts $G(z) \sim P_{\text{data}}$

Obtain the best generator

We want for the best generator that is able to fool the discriminator: is able to maximize the probability that $D(G(z)) = \text{True}$

$$\underset{G}{\text{MAX}} \quad V(D, G) = \mathbb{E}_{z \sim q(z)} [\log D(G(z))]$$

Equivalent to

$$\underset{G}{\text{MIN}} \quad V(D, G) = \mathbb{E}_{z \sim q(z)} [\log (1 - D(G(z)))]$$

Then we can reach the Nash Equilibrium merging the best discriminator and the best generator in min max settings

$$\underset{G}{\text{min}} \quad \underset{D}{\text{max}} \quad V(D, G) = \mathbb{E}_{x \sim P_{\text{data}}} \log [D(x)] + \mathbb{E}_{z \sim q(z)} \log [1 - D(G(z))]$$

Law of UNCONSCIOUS STATISTICIAN (lotus)

$$z \sim X \text{ sampled from } X$$

$$E(g(X)) = \int_{-\infty}^{\infty} g(z) f_X(z) dx$$

↓
a random variable

↓
probability density function of X

1 Let's rewrite the second part of the loss

$$\mathbb{E}_{z \sim q(z)} \log [1 - D(G(z))] = \mathbb{E}_{z \sim p_G(z)} \log [1 - D(G(z))]$$

↑
Corridor samples coming out from the generator directly

2 Apply lotus for the two terms of the GAN loss

$$\int_X P_{\text{data}}(x) \log D(x) dx + \int_Z P_G(z) \log [1 - D(G(z))] dz$$

over the random variable x over random variable z

Prob. density function PDF

unifying the two integrals

$$\int_X P_{\text{data}}(x) \log D(x) + P_G(x) \log [1 - D(G(x))] dx$$

3 Use the Liebniz Rule to obtain the optimal D given G

$$\frac{d}{dx} \left(\int_a^b f(x,t) dt \right) = \int_a^b \frac{\partial}{\partial x} f(x,t) dt$$

Take the derivative inside and because x is multivariate we must take the partial derivative ∂

A point where the previous integral is minimized

hence where the derivative is 0

Let's write again the integral using a function $f(y)$

$$f(y) = a \log y + b \log (1-y) \quad \text{and} \quad \text{the derivative of } f(y)$$

↓
P_{data} D(x) P_G

$$f'(y) = \frac{a}{y} - \frac{b}{(1-y)}$$

4

to minimize we set the derivative to 0

$$f'(y) = 0 \Rightarrow \frac{a}{y} - \frac{b}{(1-y)} = y = \frac{a}{a+b} \Rightarrow D(x) = \frac{P_{\text{data}}}{P_{\text{data}} + P_G}$$

Knowing that the Nash Equilibrium $D(x) = \frac{1}{2}$ lead to the best generator, then if $P_{\text{data}} = P_G$ we have $D(x) = \frac{1}{2}$

Theorem: the global minimum of the criterion $C(G) = \max_D V(G, D)$ is reached when $P_G = P_{\text{data}}$

Proof: • Assumption : we know $D^* = \frac{1}{2}$
and $P_G = P_{\text{data}}$

Just plug $\frac{1}{2}$

$$\begin{aligned} V(G, D^*) &= \int_X P_{\text{data}}(x) \log \frac{1}{2} + P_G(x) \log \left(1 - \frac{1}{2}\right) dx \\ &= -\log 2 \int_X P_{\text{data}}(x) dx - \log 2 \int_X P_G(x) dx = 2 - \log 2 = -\log 4 \end{aligned}$$

• Assumption $C(G) = \max_D V(G, D)$

$$D_G^* = P_{\text{data}} / (P_{\text{data}} + P_G)$$

Let's plug D_G^* into $C(G) = \max_D V(G, D)$

$$C(G) = \int_X P_{\text{data}}(x) \log \left(\frac{P_{\text{data}}}{P_{\text{data}} + P_G} \right) + P_G(x) \log \left(\frac{P_G}{P_{\text{data}} + P_G} \right) dx$$

Now let's write again $C(G)$ using a trick

$$C(G) = \int_X (\log 2 - \log 2) P_{\text{data}}(x) + P_{\text{data}}(x) \log \left(\frac{P_{\text{data}}}{P_{\text{data}} + P_G} \right) + (\log 2 - \log 2) P_G(x) + P_G(x) \log \left(\frac{P_G}{P_{\text{data}} + P_G} \right) dx$$

$$\begin{aligned} C(G) &= \int_X \cancel{\log 2 P_{\text{data}}(x)} - \cancel{\log 2 P_{\text{data}}} + P_{\text{data}}(x) \log \left(\frac{P_{\text{data}}}{P_{\text{data}} + P_G} \right) + \cancel{\log 2 P_G(x)} - \cancel{\log 2 P_G} + P_G(x) \log \left(\frac{P_G}{P_{\text{data}} + P_G} \right) dx \\ &= \int_X -\log 2 (P_{\text{data}}(x) + P_G(x)) dx + \boxed{\int_X P_{\text{data}}(x) \left(\log 2 + \log \left(\frac{P_{\text{data}}}{P_{\text{data}} + P_G} \right) \right) dx} + \boxed{\int_X P_G(x) \left(\log 2 + \log \left(\frac{P_G}{P_{\text{data}} + P_G} \right) \right) dx} \end{aligned}$$

$$= -\log 2 \int_X P_{\text{data}}(x) + P_G(x) dx + \boxed{\int_X P_{\text{data}}(x) \left(\log 2 + \log \left(\frac{P_{\text{data}}}{P_{\text{data}} + P_G} \right) \right) dx} + \boxed{\int_X P_G(x) \left(\log 2 + \log \left(\frac{P_G}{P_{\text{data}} + P_G} \right) \right) dx}$$

by definition of probability densities
the integration = 1
 $-\log_2(1+1) = -2 \log 2$

$$\log \left(\frac{2 P_{\text{data}}}{P_{\text{data}} + P_G} \right)$$

the 2 integrals are KL divergences

$$C(G) = -\log 4 + \text{KL} \left(P_{\text{data}} \mid \frac{P_{\text{data}} + P_G}{2} \right) + \text{KL} \left(P_G \mid \frac{P_{\text{data}} + P_G}{2} \right)$$

or written as

$$C(G) = -\log 4 + 2 \cdot \text{JSD} \left(P_{\text{data}} \mid P_G \right)$$

Because these are non-negative
we see that $-\log 4$ is the
global minimum of this function

or written as

$$C(G) = -\log 4 + 2 \cdot \text{JSD} \left(P_{\text{data}} \mid P_G \right)$$

CANDIDATE FOR
GLOBAL MINIMUM

$$\frac{1 - P_{\text{data}}}{P_{\text{data}} + P_G} = \frac{P_{\text{data}} + P_G - P_{\text{data}}}{P_{\text{data}} + P_G}$$

CONVERGENCE OF GAN

In the min max formulation with $D_G^* = \frac{1}{2}$ the PG converges to Polar

$$\min_G \max_D V(G, D) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D_G^*(x)] + \mathbb{E}_{x \sim p_G} [\log 1 - D_G^*(x)]$$

And we know that $p_G = p_{\text{data}}$ when the mesh equilibrium is reached and when $D(x) = \frac{1}{2}$, the discriminator is totally untrained.

PROBLEMS OF GANS

Saturation of Loss G

In beginning of training, G is poor. It rejects sample with high confidence because samples are really different from all ones. Instead of minimize $\log [1 - D(G(z))]$ that saturates we can maximize $\log [D(G(z))]$.

Mode Collapse: The generator learns the subset (of very similar) samples that fools the discriminator. If the discriminator get stuck in a local minimum and do not find that the subset of generated is fake, then the G will never learn to generate other samples.

When we have $\nabla_{\theta_G} \frac{1}{m} \sum_{i=1}^m \log [1 - D(G(z_i))]$ we will obtain $x^* = \operatorname{argmax}_x D(x)$ and x^* will be almost independent from z .

Considering the gradient $\frac{\partial \text{Loss}}{\partial z} \approx 0$ and no update will be performed.

D is too powerful that G 's gradient vanish and doesn't learn

WASSERSTEIN GAN

$$\begin{array}{cccccc} P_1 = 3 & P_2 = 2 & P_3 = 1 & P_4 = \frac{1}{4} \\ Q_1 = 1 & Q_2 = 2 & Q_3 = \frac{1}{4} & Q_4 = 3 \end{array} \left. \begin{array}{l} \text{How do } \mathcal{J} \\ \text{match these 2 distributions?} \end{array} \right.$$

$$\begin{array}{ll} \text{Move } (2, P_1, P_2) & P_1 = Q_1 \\ \text{Move } (2, P_2, P_3) & P_2 = Q_2 \\ \text{Move } (1, Q_3, Q_4) & P_3 = Q_3, P_4 = Q_4 \end{array} \quad \begin{array}{ll} P_1 = 1 & P_2 = \frac{1}{4} \\ P_2 = Q_2 & P_1 = 1 \\ P_3 = 1 & P_2 = 2 \\ Q_1 = 1 & Q_2 = 2 \end{array} \quad \begin{array}{ll} P_3 = 1 & P_4 = \frac{1}{4} \\ P_3 = 3 & P_4 = \frac{1}{4} \\ Q_3 = 3 & Q_4 = \frac{1}{4} \end{array}$$

$$W = \sum |\delta_i| = 2+2+1 = 5 \quad \text{This is the discrete case}$$

For the continuous case

$$W(P_G, P_{\text{data}}) = \inf_{\gamma \sim \Pi(P_G, P_{\text{data}})} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|]$$

Take the smallest amount of sand to move
all the possible sand movement PLAN between the 2 dist.

We can use the wasserstein distance as a loss function

This distance is more smoothly, not like KL or JSD

If we want to use W distance in the loss, we have to consider that we cannot enumerate \mathcal{G} (all the possible set of flows of data to move)

We can use $W(p_2, p_g) = \frac{1}{k} \sup_{\|g\|_L \leq k} |\mathbb{E}_{x \sim p_2}[g(x)] - \mathbb{E}_{x \sim p_g}[g(x)]|$

f should be k -Lipschitz continuous

Instead of take infimum over joint distributions, now we take sup of a function

Now we parametrize f with a set of parameters that we optimize. f will be a Neural Net

We will have an objective that directly maximize the similarity between p_{data} and p_G

$$\text{Loss}(p_2, p_g) = W(p_2, p_g) = \max_{w \in \mathcal{W}} |\mathbb{E}_{x \sim p_2}[f_w(x)] - \mathbb{E}_{z \sim p_g}[f_w(g_\theta(z))]|$$

TRAIN A GAN with WASSERSTEIN DISTANCE

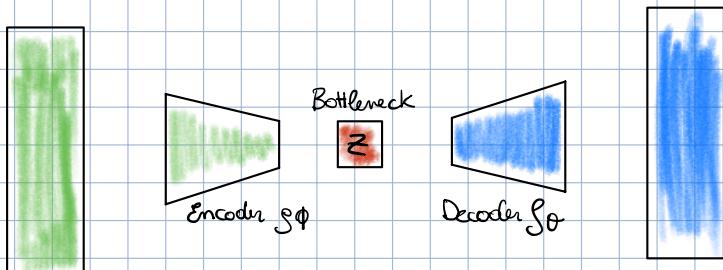
While θ not converged

for iteration

- update the param to enforce $p_{\text{data}} = p_g$ using W distance
- Sample $\{x_i\}_{i=1}^m$ Real data
 - Sample $\{z_i\}_{i=1}^m$ Noise Samples
 - $\nabla_w = \nabla_w \left[\frac{1}{m} \sum_{i=1}^m f_w(x_i) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z_i)) \right]$ this is to enforce $p_{\text{data}} = p_g$ by the using of a function f_w
 - $w = w + \alpha \text{ RMS Prop}(w, \nabla_w)$ using the parameter of generator
 - $w = \text{clip}(w, -c, c)$ this is to enforce k -Lipschitz

- Update G params
- Sample $\{z_i\}_{i=1}^m$ from $q(z)$ a Noise Samples
 - $g_\theta = -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z_i))$ gradient w.r.t θ params of generator
 - $\theta = \theta - \alpha \text{ RMS Prop}(\theta, g_\theta)$

AUTOENCODERS



Note: this has a problem, might learn an identity function

Denoising autoencoder

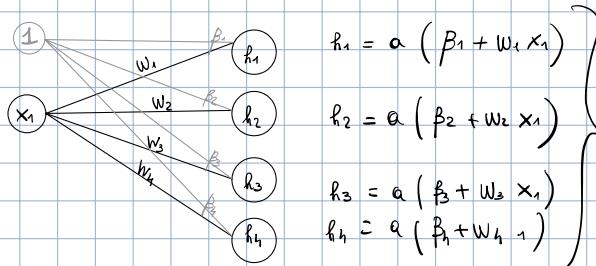
Make the input noisy $\tilde{x} = M_D^\text{noise}(x)$

$$x' = f_\phi(g_\phi(\tilde{x}))$$

$L(\phi, \theta) = \sum_{i=1}^m (x_i - x'_i)^2 = \sum_{i=1}^m (x_i - f_\phi(g_\phi(x_i)))^2$

SPARSE AUTOENCODER

Avoid overfitting, improve robustness.
Enforce activation of few neurons.



Enforce small activations
with a term in the loss

$$L(\phi, \theta) = L_{\text{MSE}}(\phi, \theta) + \beta D_{KL}(\phi || \hat{\phi})$$

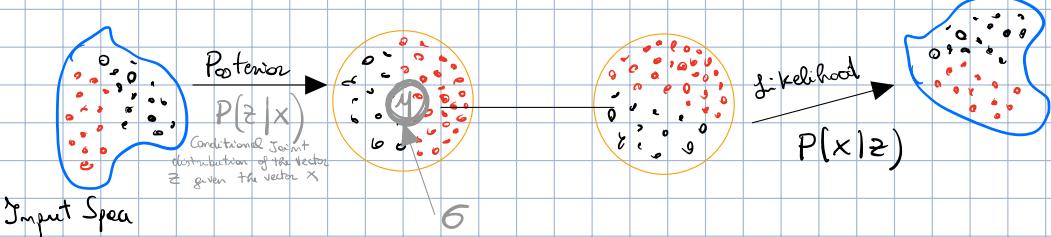
target activation
for a certain layer

VARIATIONAL AUTOENCODERS

Idea Behind: the input space follows a distribution. We want to learn a represent. for the whole input space

Instead of learn a fixed vector for a given sample x , we learn a distribution or a latent space. (How sample are distributed in a lower dimension latent space)

- $P_{prior}(z)$ Prob. of latent space
- Likelihood $p(x|z)$ Prob. of getting a sample x from a latent variable z)
- Posterior $p(z|x)$ Prob. of obtaining the latent variable z from x sample



Again, we want to learn the distribution of the input data. We do it by observe the input and use it to estimate the parameters of another latent variable \mathbf{z} .



f is a NN

How do we learn the weight of f

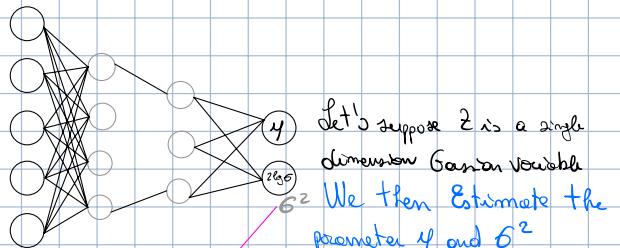
We need a loss which should measure the DISTANCE between our prob. distribution of the input, w_2 to \mathbf{z} .

We may use KL Divergence or (Expected value of log likelihood)

$$D_{KL}(q_\phi || p_\theta) = \mathbb{E}_{q_\phi} \log \left[\frac{q_\phi(z|x)}{p_\theta(z|x)} \right]$$

↑ prediction
↑ ground truth prob. distribution

$$= \int_{\mathbb{R}} q_\phi(z|x) \log \left[\frac{q_\phi(z|x)}{p_\theta(z|x)} \right]$$



we change this $\log \sigma^2$ then to $2 \log \sigma$ to have the output of the net to be a positive Real Number

$$(z \sim N(y, \exp(0.5 \cdot 2 \log \sigma)))$$

↑ THIS IS GIVING THE PREDICTION
 $q_\phi(z|x)$

this is the GT
but we do not have this

For this reason we use

ELBO Evidence Lower Bound Is the loss of the variational Autoencoder

Derivation of the loss starting from KL Divergence

$$D_{KL}(q_\phi || p_\theta) = \mathbb{E}_{q_\phi} \left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)} \right]$$

Just rewrite the fraction as subtraction using log property

$$= \mathbb{E}_{q_\phi} [\log q_\phi(z|x)] - \mathbb{E}_{q_\phi} [\log p_\theta(z|x)]$$

By Bayes Rule
 $p_\theta(z|x) = p_\theta(x|z) p_\theta(z)$
 $p_\theta(x)$

$$= \mathbb{E}_{q_\phi} [\log q_\phi(z|x)] - \mathbb{E}_{q_\phi} \left[\log \frac{p_\theta(z|x)}{p_\theta(x)} \right]$$

$$= \mathbb{E}_{q_\phi} [\log q_\phi(z|x)] - \mathbb{E}_{q_\phi} [\log p_\theta(z,x)] + \mathbb{E}_{q_\phi} [\log p_\theta(x)]$$

$$= \mathbb{E}_{q_\phi} [\log q_\phi(z|x)] - \mathbb{E}_{q_\phi} [\log p_\theta(z,x)] + \int_{\mathbb{R}} q_\phi(z|x) \cdot \log p_\theta(x) dz$$

$$= \mathbb{E}_{q_\phi} [\log q_\phi(z|x)] - \mathbb{E}_{q_\phi} [\log p_\theta(z,x)] + \int_{\mathbb{R}} (\log p_\theta(x)) q_\phi(z|x) dz = 1$$

Moved out because doesn't depend on z

$$= \mathbb{E}_{q_\phi} [\log q_\phi(z|x)] - \mathbb{E}_{q_\phi} [\log p_\theta(z,x)] + \boxed{\log p_\theta(x)} \text{ Marginal Log likelihood Evidence}$$

We cannot do nothing then we rewrite every thing like follow

Marginal Log likelihood

$$\log p_\theta(x) = \boxed{-\mathbb{E}_{q_\phi} [\log q_\phi(z|x)] + \mathbb{E}_{q_\phi} [\log p_\theta(z,x)]} + D_{KL}(q_\phi || p_\theta)$$

if This is MAXIMIZED

This has to be MINIMIZED

to hold the equality

we know that this ≥ 0

Let's reason a bit on this

Then we can say that at minimum we can have

$$\log p_\theta(x) \geq \boxed{-\mathbb{E}_{q_\phi} [\log q_\phi(z|x)] + \mathbb{E}_{q_\phi} [\log p_\theta(z,x)]} \text{ Joint prob of } x$$

We call this ELBO

$$\begin{aligned} \text{ELBO} &= -\mathbb{E}_{q_\phi} [\log q_\phi(z|x)] + \mathbb{E}_{q_\phi} [\log p_\theta(x|z)] + \mathbb{E}_{q_\phi} [\log p_\theta(z)] \\ &= \mathbb{E}_{q_\phi} [\log p_\theta(x|z)] - \mathbb{E}_{q_\phi} [\log \frac{q_\phi(z|x)}{p_\theta(z)}] \end{aligned}$$

Expected reconstruction error

D_{KL} Divergence Between posterior and prior

Using Elbo we want to maximize it

$$\text{ELBO} = \mathbb{E}_{q_\phi} [\log p_\theta(x|z)] - \mathbb{E}_{q_\phi} [\log \frac{q_\phi(z|x)}{p_\theta(z)}]$$

Here we have problems

We have this

We have this

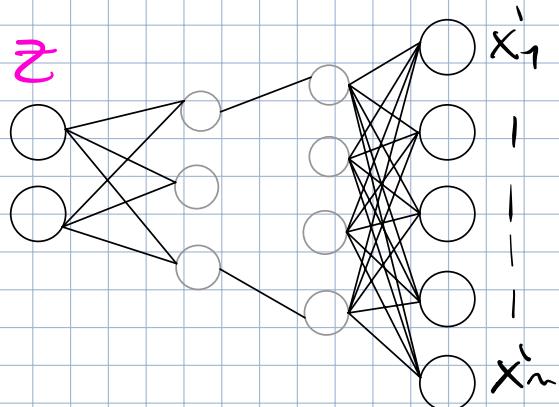
Again we can apply the same reason of before, we

try to find the parameter for a function that allows to

obtain x' giving as input z .

$z = \text{encoder}(x)$

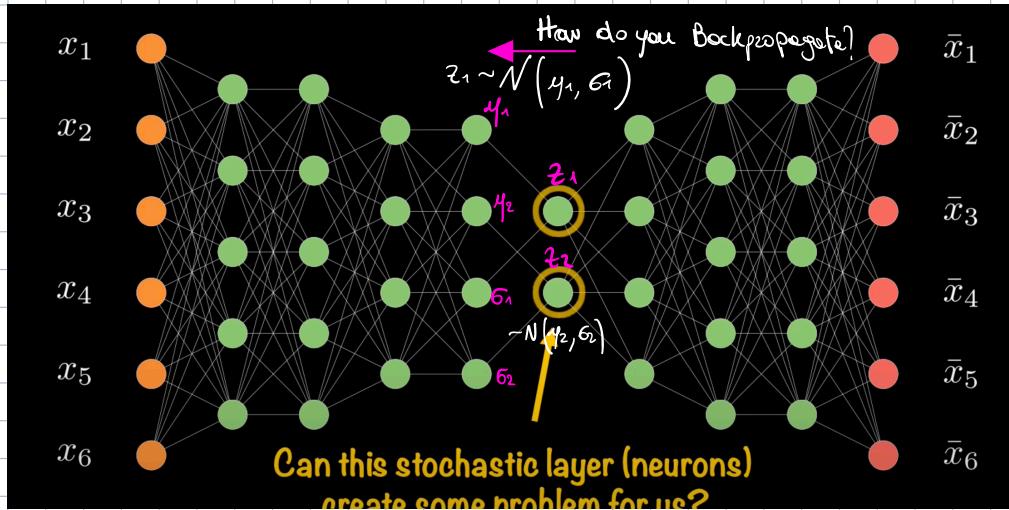
$x' = \text{decoder}(z)$



with x' we can now substitute this term just with

$$\text{MSE}(x, x')$$

REPARAMETERIZATION TRICK



Remember the Elbo Loss

$$\mathcal{L}_{\phi, \theta}(x) = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{q_{\phi}(z|x)}{p_{\theta}(z)} \right]$$

↑
parameters of the encoder

↑
this is the set that contains μ, σ for z

Let's rewrite some stuff

$$\begin{aligned} &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) - \log \frac{q_{\phi}(z|x)}{p_{\theta}(z)} \right] \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) - \log q_{\phi}(z|x) + \log p_{\theta}(z) \right] \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) - \log q_{\phi}(z|x) \right] \end{aligned}$$

And now we can take gradient WRT ϕ and θ , for θ ok, for ϕ we cannot because it's not possible to apply the Liebniz Rule

Change of variable to sample from a distrib.

$x \sim p_{\theta}(X)$ and we take $x = g(\epsilon, \theta)$ where $\epsilon \sim p(\epsilon)$

Example $p_{\theta}(x) = N(x; \mu, \sigma)$ where $\theta = \{\mu, \sigma\}$

Sample on $X = \mu + \sigma \cdot \epsilon$ where $\epsilon \sim p(\epsilon)$ base distribution

How do we sample from the Base Distribution?
How the Base Distrib. is characterized?

$$p(\epsilon) = N(\epsilon, 0, 1)$$

Now Just a Reminder of Lutus

$$\mathbb{E}_{p_\theta(x)} [f(x)] = \int f(x) p_\theta(x) dx$$

then because $x \sim p_\theta(x)$ and $x = g(\epsilon, \theta)$ where $\epsilon \sim p(\epsilon)$ then

$$\nabla_\theta \mathbb{E}_{p_\theta(x)} [f(x)] = \nabla_\theta \int f(g(\epsilon, \theta)) p(\epsilon) d\epsilon \quad \} \text{ Lutus}$$

this is the def
of expectation
of continuous RV

$$= \left\{ \nabla_\theta f(g(\epsilon, \theta)) p(\epsilon) d\epsilon \right\} \text{ Liebniz}$$

$$= \mathbb{E}_{p(\epsilon)} [\nabla_\theta f(g(\epsilon, \theta))]$$

$$\text{by} \quad \approx \frac{1}{N} \sum_i^N \nabla_\theta f(g(\epsilon^i, \theta)) \quad \epsilon^i \sim p(\epsilon)$$

stochastic
Batch

Now we can do the same on ELBO

$$\text{ELBO} = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)]$$

set $\epsilon \sim p(\epsilon)$ sample from a base dist.
 $z = g(\phi, x, \epsilon)$ deterministic func of ϕ
 (differentiable) by encoder

Here is
applied change
of variable

Using Lutus

$$\mathcal{L} = \mathbb{E}_{\epsilon \sim p(\epsilon)} [\log p_\theta(x, z) - \log q_\phi(z|x)]$$

Thanks to Liebniz Rule we move in the gradient

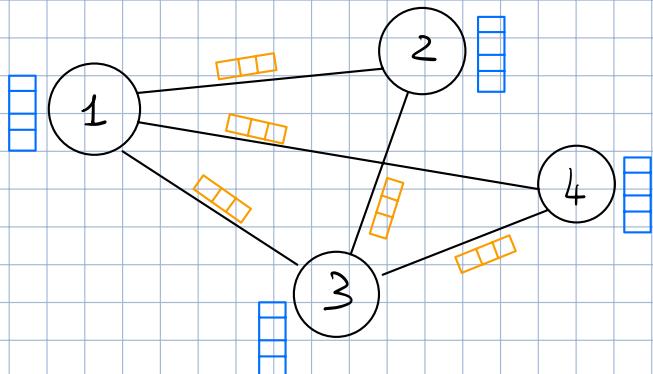
$$\nabla_\phi \mathcal{L}_{\theta, \phi}(x) = \nabla_\phi \mathbb{E}_{p(\epsilon)} [\log p_\theta(x, z) - \log q_\phi(z|x)]$$

CAPITOLO 10

GRAPH NEURAL NETWORKS

- GNN Basic
- Tasks
- Graph CNN
- Bathing
- Aggregation
- Spectral Graph
- Graph convolution
- Chebnet

GNN Basic



To describe a Graph G we need 3 matrices

- A Adjacency Matrix

$|v| \times |v|$ $i,j = 1$ if
node i linked to node j

	1	2	3	4
1	0	1	1	1
2	1	0	1	0
3	1	1	0	1
4	1	0	1	0

- X Node Feature Matrix

$d^0 \times |v|$ num of nodes
 $\times d^0$ num of features

	1	2	3	4
(1,2)	blue	blue	blue	blue
(1,3)	blue	blue	blue	blue
(1,4)	blue	blue	blue	blue
(2,3)	blue	blue	blue	blue
(2,4)	blue	blue	blue	blue
(3,4)	blue	blue	blue	blue

- E Edge feature Matrix

$d^e \times |E|$ dim of features
 $\times |E|$ num of edges

	1	2	3	4
(1,2)	orange	orange	orange	orange
(1,3)	orange	orange	orange	orange
(1,4)	orange	orange	orange	orange
(2,3)	orange	orange	orange	orange
(2,4)	orange	orange	orange	orange
(3,4)	orange	orange	orange	orange

PERMUTATION OF NODES

→ P Matrix that change the ordering of node

$$X' = PX$$

$$A' = P^T A P$$

TASKS

Graph Classification

$$\Pr(y=1 | X, A) = \text{SIG} \left[\beta_k + W_k H_k \frac{1}{N} \right]$$

Probability of graph to be classified
as 1 given node feature Matrix
And Adj matrix

$$\text{hidden} \quad \text{last hidden layer}$$

$$\beta_k + W_k H_k \frac{1}{N}$$

Column Vector 1

1xD

this gives a
single value

GNN Layer

$$h_0 = \text{layer } (\mathbf{x})$$

Each column represent the hidden
representation of a node

	1	2	3	4
h_0	blue	blue	blue	blue
1	blue	blue	blue	blue
2	blue	blue	blue	blue
3	blue	blue	blue	blue
4	blue	blue	blue	blue

Node Classification

$$\Pr(y^m=1 | X, A) = \text{SIG} \left(\beta_k + W_k h^m \right)$$

Just consider the hidden state of the node
to be classified (m)

Edge Classification

$$\Pr(y^{mm}=1 | X, A) = \text{SIG} \left[\begin{pmatrix} h^m \\ h^m \end{pmatrix}^T \begin{pmatrix} h^m & h^m \end{pmatrix} \right]$$

Make use of the
emb. of nodes

GRAPH CNN

Graph Convolutional Layer

$$H_1 = F(X, A, \Phi_0)$$

$$H_i = F(H_{i-1}, A, \Phi_{i-1})$$

depends on previous hidden state and previous set of Φ

$$\nabla H_k = F(H_{k-1}, A, \Phi_{k-1})$$

If mode m , in layer i we sum the aggregator $AGG[m, i] = \sum_{m \in N(e)} h_i^m$

We sum up all the hidden states of this layer from the neighbors mode e of mode m

How do we get h_i^m (Hidden state of mode m for layer $i+1$)

$$h_i^m = \alpha [B_i^T + W_i h_i + W_i AGG[m, i]]$$

Proj. of current hidden State information Aggregated information projected

If we want to do it for all the nodes

$$H_{i+1} = \alpha [B_i^T + W_i H_i + W_i H_k A]$$

$$= \alpha [B_i^T + W_i H_i (A + I)]$$

\hookrightarrow Identity

Then, A full network of GCL

$$H_1 = \alpha [\beta_0 + W_0 X (A + I)]$$

$$H_i = \alpha [\beta_{i-1} + W_i H_{i-1} (A + I)]$$

$$H_k = \alpha [\beta_{k-1} + W_k H_{k-1} (A + I)]$$

GRAPH CLASSIF.

$$f(X, A, \Phi) = \alpha [\beta_k + W_k H_k 1/N]$$

NODE CLASS.

$$f[X, A, \Phi] = \text{sig} [\beta_k 1^T + W_k H_k]$$

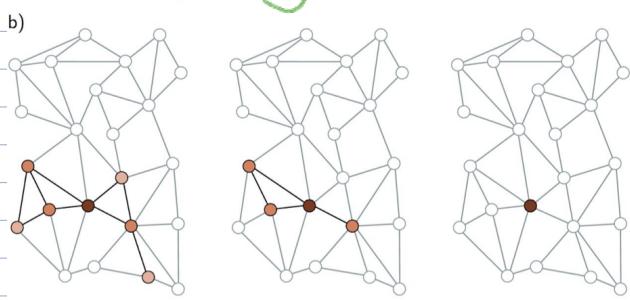
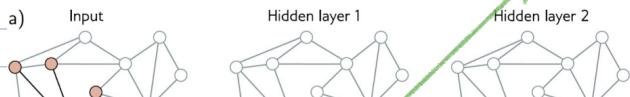
AGGREGATION

BATCHING

NODE SUBSET

Choose a random subset of nodes at training step and then we select the neighbors from previous layers

One layer means we are taking only the 1-hop Nodes



$$h_1^x = \alpha [\beta_0 + W_0 X + W_1 \text{Agg}(X, 1)]$$

$$\text{Or a random subset of } N(x) \quad \sum_{m \in N(x)} h_m$$

$$h_2 = \alpha [\beta_1 + W_0 h_1 + W_1 \text{Agg}(X, z)]$$

Another Mechanism Is the GRAPH PARTITIONING

AGGREGATION MECHANISM

Summation

$$H_{k+1} = \alpha \left[\beta_k \mathbf{1}^\top + W_k H_k (A + I) \right]$$

Diagonal Enhancement

$$H_{k+1} = \alpha \left[\beta_k \mathbf{1}^\top + W_k H_k \left(A + (1 + \varepsilon_k) I \right) \right]$$

$$\text{AGG}(x, i) = \sum_{m \in N(x)} (1 + \varepsilon_i)^{h_m^i}$$

Each node is multiplied by a factor of $(1 + \varepsilon_k)$ before contributing to the sum, ε_k is learnt and differ for each layer

Generalized Diagonal Enhancement

$$H_{k+1} = \alpha \left[\beta_k \mathbf{1}^\top + \underbrace{\Omega_k}_{\text{AGG}} H_k A + W_k H_k \right]$$

Different weight for the sub. of the current Node

$$\Omega_k = \begin{bmatrix} H_k A \\ H_k \end{bmatrix}$$

Residual Connection

$$H_{k+1} = \left[\alpha \left[\beta_k + W_k H_k A \right], H_k \right]$$

① Aggregate, linear transf., activation
② Conc. to the projection of current node

Mean Aggregation

$$\text{AGG}(x, i) = \frac{1}{|N(x)|} \sum_{m \in N(x)} h_m$$

Smooth the contr. of aggregation, less inform. about the structure of the graph

$$H_{k+1} = \alpha \left[\beta_k \mathbf{1}^\top + W_k H_k (AD^{-1} + I) \right]$$

KIPF Normaliz.

$$\text{AGG}(x, i) = \sum_{m \in N(x)} \frac{h_m}{\sqrt{|N(x)| \cdot |N(m)|}}$$

Smooth the contribution of the neighbor nodes that have a lot of neighbors.
(They have less unique information)

$$H_{k+1} = \alpha \left[\beta_k \mathbf{1}^\top + W_k H_k (D^{-1/2} A D^{-1/2} + I) \right]$$

MAX Pooling

$$\text{AGG}(n, i) = \max_{m \in N(x)} h_m$$

Take the max embedding of the neighbors.

Attention in Aggregation

$$H'_k = \beta_k \mathbf{1}^\top + W_k H_k$$

Then consider similarity score $s_{xm} = \alpha \left[\Phi_k^\top \begin{bmatrix} h'_x \\ h'_m \end{bmatrix} \right]$

S is a $N \times N$ Matrix of similarity Scores

$$H_{k+1} = \alpha \left[H'_k \cdot \text{Softmax}(S, A + I) \right]$$

SPECTRAL GRAPH CONV. FRAMEWORK

A GRAPH CONVOLUTIONAL LAYER is composed of

• PATCH FUNCTIONS

Shape of the convolutional filter:
selection of the nodes that have to interact with each other in each layer. Matrix $|N| \times |N|$, one \downarrow layer is $\downarrow i=1 \dots k$
 $|N| \times |N| \rightarrow (f_1(W, H^e) \dots f_k(W, H^e))$

Matrix of node features
at layer l

k can be:

- Spectral: defining the number of neighbors
- Spectral: with the rank of matrix

• CONVOLUTION FILTERS WEIGHTS

In layer we have K matrices $d_e \times d_{e+1}$ ($\Theta_1^e, \dots, \Theta_K^e$) and this are convolved with every patch H^e

$$m_k^{l+1} = f_k(W, H^e) H^e \Theta_k^e \text{ for } 1 \leq k \leq K$$

• MERGING FUNCTIONS

Combine the output of multiple convolution into one single representation

$$H^{l+1} = h(m_1^{l+1}, \dots, m_K^{l+1})$$

Method	Model	$g_k(\cdot)$	$h(m_1, \dots, m_k)$
Spectrum-based: $\tilde{L} = U \Lambda U^\top, f_k(W, H) = g_k(U)$	SCNN	$g_k(U) = u_k u_k^\top$	$\sigma(\sum_k m_k)$
Spectrum-free: $f_k(W, H) = g_k(W, D)$	ChebNet	$g_k(W, D) = T_k(\frac{2(I - D^{-1/2}WD^{-1/2})}{\lambda_{max}(I - D^{-1/2}WD^{-1/2})} - I)$	$\sigma(\sum_k m_k)$
	GCN	$g_1(W, D) = (D + I)^{-1/2}(W + I)(D + I)^{-1/2}$	$\sigma(m_1)$
Spatial: $f_k(W, H) = g_k(W, D)$	SAGE-mean	$g_1(W, D) = I, g_2(W, D) \sim \mathcal{U}_{\text{norm}}(D^{-1}W, q)$	$\sigma(m_1 + m_2)$
	GGNN	$g_1(W, D) = I, g_2(W, D) = W$	GRU(m_1, m_2)
Attention: $f_k(W, H) = \alpha(W \cdot g_k(H))$	MoNet	$g_k(U^s) = \exp(-\frac{1}{2}(U^s - \mu_k)^\top \Sigma_k^{-1}(U^s - \mu_k))$	$\sigma(\sum_k m_k)$
	GAT	$g_k(H) = \text{LeakyReLU}(HB^\top b_0 \oplus b_1^\top BH^\top)$	$\sigma(m_1 \dots m_k)$

GRAPH LAPLACIAN

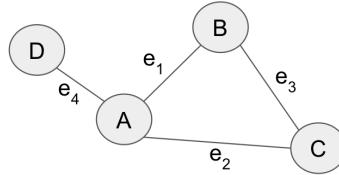
Given a graph $G := (V, E)$

- $V := \{v_1, v_2, \dots, v_n\}$ is the set of nodes/vertices
- $E := \{e_1, e_2, \dots, e_m\}$ is the set of edges
- adjacency matrix $A \in \{0, 1\}^{n \times n}$

$$A_{ij} := \begin{cases} 1, & \text{if there is an edge between } v_i \text{ and } v_j \\ 0, & \text{otherwise} \end{cases}$$

and degree matrix $D \in \mathbb{Z}^{n \times n}$

$$D_{ij} := \begin{cases} \text{degree}(v_i), & \text{if } i = j \\ 0, & \text{otherwise} \end{cases}$$

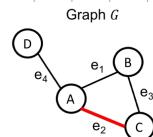


Gradient It is the easy difference of f for 2 nodes $e_k = (v_i, v_j) \in E$ edges

$$g(e_k) = f(v_i) - f(v_j) \quad i < j \quad \text{and with ordered edges we can write } g = [g(e_1) \dots g(e_m)]$$

How to do this in a Matrix Form? Incidence matrix K .

$$K_{ij} := \begin{cases} 1, & \text{if } i == k \text{ and } i > h \\ 1, & \text{if } i == h \text{ and } i > k \\ -1, & \text{if } i == k \text{ and } i < h \\ -1, & \text{if } i == h \text{ and } i < k \\ 0, & \text{otherwise} \end{cases}$$



$$\begin{array}{c} \text{Incidence Matrix } K \\ \hline \begin{matrix} & A & B & C & D \\ e_1 & 1 & -1 & 0 & 0 \\ e_2 & 1 & 0 & -1 & 0 \\ e_3 & 0 & 1 & -1 & 0 \\ e_4 & 1 & 0 & 0 & -1 \end{matrix} \end{array}$$

Functions map each point into real numbers $f: V \rightarrow \mathbb{R}$

Define the notion of $f: [f(v_1) \dots f(v_m)]$

- point
- gradient
- functions
- divergence

Point Each vertex $v \in V$ is a point

Then we can compute

$$\begin{array}{ccccc} K & & f & & g \\ \left[\begin{array}{cccc} 1 & -1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & -1 & 0 \\ 1 & 0 & 0 & -1 \end{array} \right] & \left[\begin{array}{c} f(A) \\ f(B) \\ f(C) \\ f(D) \end{array} \right] & = & \left[\begin{array}{c} f(A) - f(B) \\ f(A) - f(C) \\ f(B) - f(C) \\ f(A) - f(D) \end{array} \right] \end{array}$$

Divergence given mode A Consider divergence $(A) = g(e_1) + g(e_2) + g(e_4)$

$$\begin{array}{ccccc} K^T & & g & & d \\ \left[\begin{array}{cccc} 1 & 1 & 0 & 1 \\ -1 & 0 & 1 & 0 \\ 0 & -1 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{array} \right] & \left[\begin{array}{c} g(e_1) \\ g(e_2) \\ g(e_3) \\ g(e_4) \end{array} \right] & = & \left[\begin{array}{c} g(e_1) + g(e_2) + g(e_4) \\ -g(e_1) + g(e_3) \\ -g(e_2) - g(e_3) \\ -g(e_4) \end{array} \right] \end{array}$$

$$\begin{array}{c} \text{The Laplacian is thus } L := K^T K \\ K^T K = D - A \end{array}$$

CAPITOLO 11

GEOMETRIC DEEP LEARNING

- Group & Actions
- Comm. deep learning
- Graph NN, Equivariance
- Deep Sets

- Group Equivariant

Branch that aims to develop geometric invariant learning

Euclidean : function defined over a given number of coordinates (x, y, z)

In NON Euclidean : the function is not defined over the same number of coordinate & input of the function
In Graph we have node A linked with 2 nodes (2 coordinates) and node B linked with 3 coordinates.

We want invariance also in this type of Spaces

Goal of Geometric Deep learning : build representation of object defined by a particular structure / geometry.

We define a Geometry as the set of possible transformation that applied over objects, preserve certain property

Euclidean Geometry \rightarrow the set of transformations that preserve area, angle, distance

Given a certain type of data

- ① Identify the invariant property (translation, rotation, symmetry...)
- ② Create a model which is robust to a given invariant (maintain a given invariant)
(choose or optimize)
- Leverage the INDUCTIVE BIAS

Given the dense operation in a Deep learning model : input is a vector, layer is a Matrix of weight

We can describe transformation (this matrices) as element of a Group

Representation $f: \rho: G \rightarrow \text{General Linear Group}(\mathbb{V})$

I need a tool that takes an element of this group and transform it into a Matrix

Invariance

Given a Group of transformation G

I have represent of the element in G $\rho(G)$

a function f is Invariant to G if

The outcome of f

$f(\rho(G)x) = f(x)$ is the same both if
 x is not transformed

or is transformed by

an element of the group $f(\rho(G)x)$

Equivariance

f is Equivariant w.r.t

$$f(f(\rho(G)x)) = f(\rho(G)f(x))$$

Transform input and get result =
Transform the representation

Question: Can we emulate the receptive field mechanism of the convolution?

Hence, can we look at features of the input at a different SCALES (different levels of granularity) ?

The Local stability is the property of preserving a certain type of feature at different scales

Example $f\left(\begin{array}{|c|} \hline \text{Tree} \\ \hline 100 \times 100 \text{ px} \\ \hline \end{array}\right) = \text{True}$ $f\left(\begin{array}{|c|} \hline \text{Tree} \\ \hline 200 \times 200 \text{ px} \\ \hline \end{array}\right) = \text{True}$ The TREE Property is locally stable in f

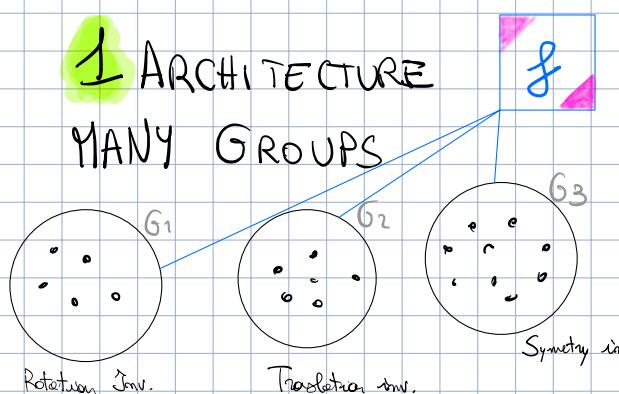
GEOMETRIC DEEP LEARNING

Input $x \rightarrow$ Equivariant layer $f(f(G)x) = f(G)f(x)$ \rightarrow Local Pooling Is an aggregation function \rightarrow Equivariant layer \rightarrow Single global Representation

this is why you take an Elmt of the group of transformation
ASSUME A CERTAIN GEOMETRY OF THE INPUT.

This is because then you need the group of certain transformation adapt for certain geometry

What's the advantage?



$$f(f(G_i)x) = f(G_i)f(x)$$

A function that choose a certain element of a given group

Ideal characteristic of GNN 2 isomorphic graphs \rightarrow Same representation.
2 non-isomorphic \rightarrow Different represent.

To make a graph at ∞ power of WL test
the aggregation funct should be

$$h_v^{(k)} = \text{MLP}^{(k)} \left(\left(1 - \varepsilon^{(k)} \right) h_v^{(k-1)} + \sum_{i \in N(v)} h_i^{(k-1)} \right)$$

Self loop weight
Hetero layer func.

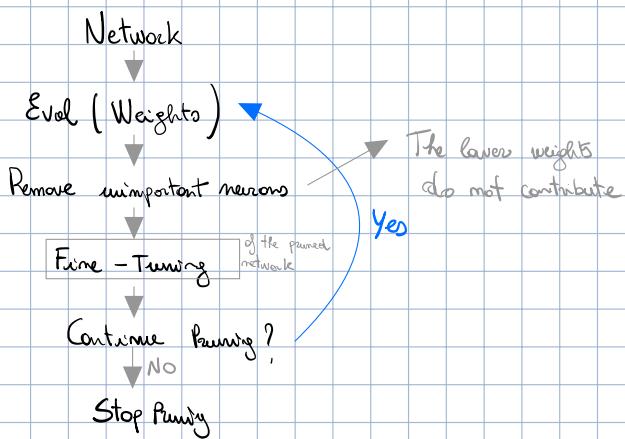
This makes it injective

CAPITOLO 12

- Pruning NN
- Lottery Ticket Hypothesis
- Weight Quantization

MODEL COMPRESSION

PRUNING



LOTTERY TICKET HYPOTHESIS

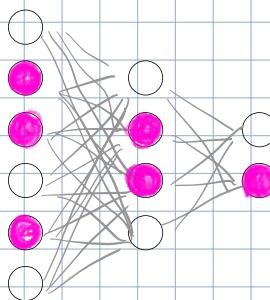
Net (original) $\rightarrow X\%$ of accuracy
 Net (pruned) $\rightarrow Y\%$ of accuracy with $Y \leq X$

The retraining idea

Smaller pruned network perform worse than the original network because of course they are less powerful.

The lottery ticket hypothesis

A random initialized network contains a subnetwork (using the original weight variables) such that if you just train this small network you can get the same performance of the big original network.



Given $f(x, \theta)$ init weight $\theta = \theta_0$

Train f and it gets loss l at step J accuracy a

Now take $f(x, m \circ \theta)$ with init $\theta = m \circ \theta_0$

then you get l_0 at J_0 with accuracy a_0

Lottery ticket: $\exists m / J_0 < J, a_0 \geq a, \|m\| \ll \theta$

The algorithm works as follows:

- Randomly initialize a neural network $f(x; \theta_0)$ (where $\theta_0 \sim D_\theta$).
- Train the network for J iterations, arriving at parameters θ_J .
- Prune $p\%$ of the parameters in θ_J , creating a mask m .
- Reset the remaining parameters to their values in θ_0 , creating the winning ticket $f(x; m \circ \theta_0)$.

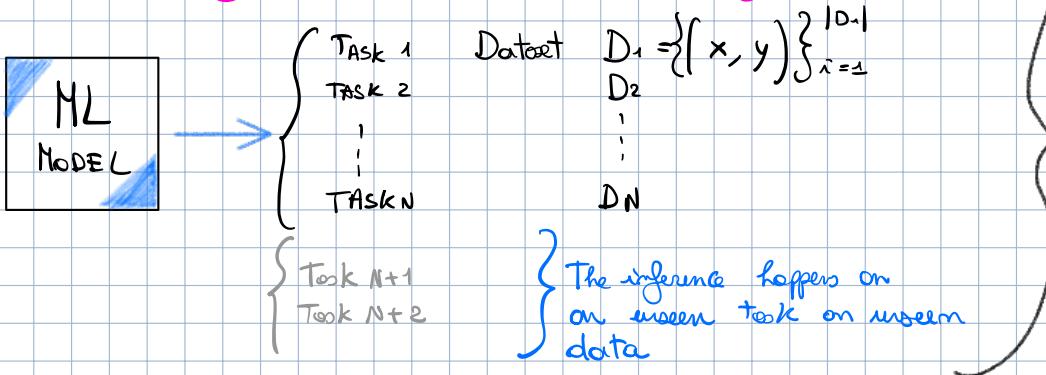
CAPITOLO 13

META LEARNING

- Meta Learning Problem, Few Shot
- Supervised VS meta learning
- Metric - Based Meta Learning
- Model - Based Meta Learning
- Optimization - Based Meta Learning
- Learning to Learn
- MAML

Design a ML model with human learning abilities:
Learning new concept and tasks with few training examples.

Meta Learning in Supervised Learning



Learning will be

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{D \sim p(D)} [\ell_\theta(D)]$$

Find the best θ by optimize over a distribution of task

FEW-SHOT LEARNING k-SHOT N-CLASS

Meta learning supervised learning where each Dataset $D_i = \langle S, B \rangle$

S is the subset of D_i of k sample used for training

N classes tells the number of Dataset available

In terms of likelihood

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{(x, y) \in D} [\mathbb{P}_\theta(y|x)]$$

with minibatches $\underset{\theta}{\operatorname{argmax}} \mathbb{E}_{B \in D} \left[\sum_{(x, y) \in B} \mathbb{P}_\theta(y|x) \right]$

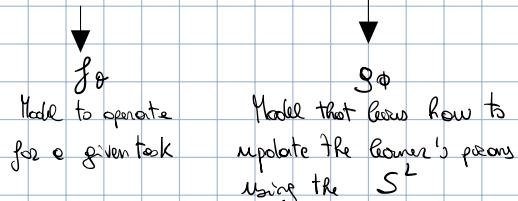
SUPERVISED VS META-LEARNING

$$\theta = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{L \in \mathcal{L}} \left[\mathbb{E}_{S^L \in D, B^L \in D} \left[\sum_{(x, y) \in B^L} \mathbb{P}_\theta(y|x, S^L) \right] \right]$$

Given a sample from batch use the support set S^L to compute the likelihood

a subset of the total labeled pairs
The whole dataset

LEARNER & META-LEARNER



Given the problem, we can now look at the approaches

How $\mathbb{P}_\theta(y|x)$ is modeled

Kernel Function that stores memory

METRIC-BASED $\sum_{(x, y) \in S^L} k_\theta(x, x_i) \cdot y_i (*)$

MODEL-BASED $f_\theta(x, S^L)$

GRADIENT-DESCENT $P_{f_\theta(\theta, S^L)}(y|x)$

$$\theta^* = g(\theta, S^L)$$

The final optimization would be

$$\mathbb{E}_{\mathcal{L} \in \mathcal{L}} \left[\mathbb{E}_{S^L \in D, B^L \in D} \left[\sum_{(x, y) \in B^L} \mathbb{P}_{g(\theta, S^L)}(y|x) \right] \right]$$

METRIC - BASED META-LEARNING

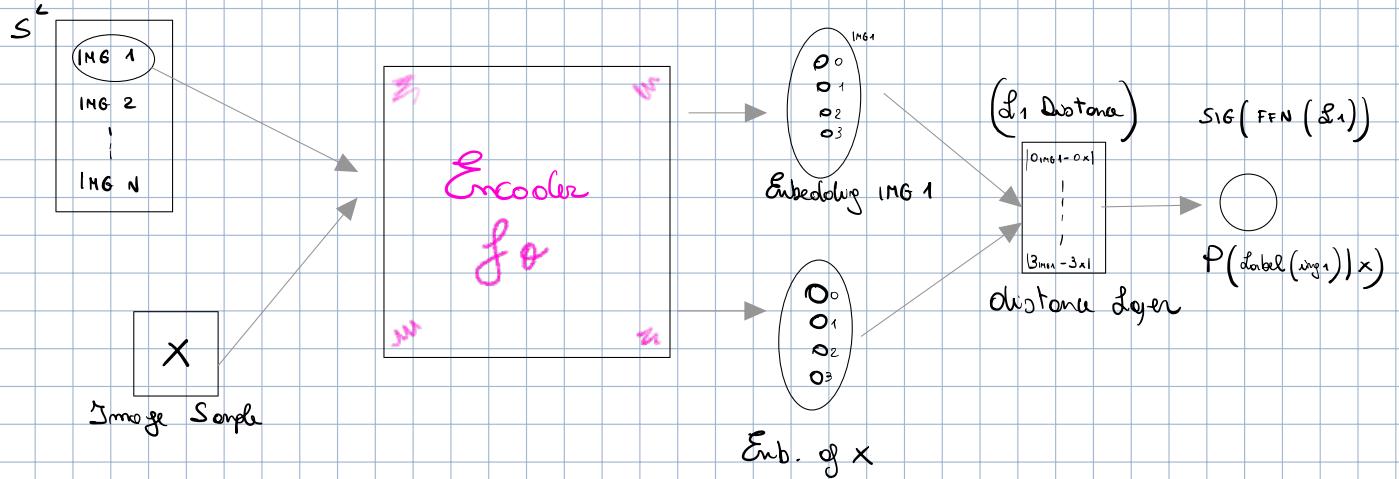
$$P_\theta(y|x) = \sum_{(x_i, y_i) \in S^L} K_\theta(x, x_i) \cdot y_i$$

The probability of y given a sample x is the weighted sum of a kernel function the measure the similarity between x , x_i (element of S^L) as weight

Then to implement this we need

- ① Learn the embedding of samples
- ② Learn how to compute the similarity with the kernel's parameters

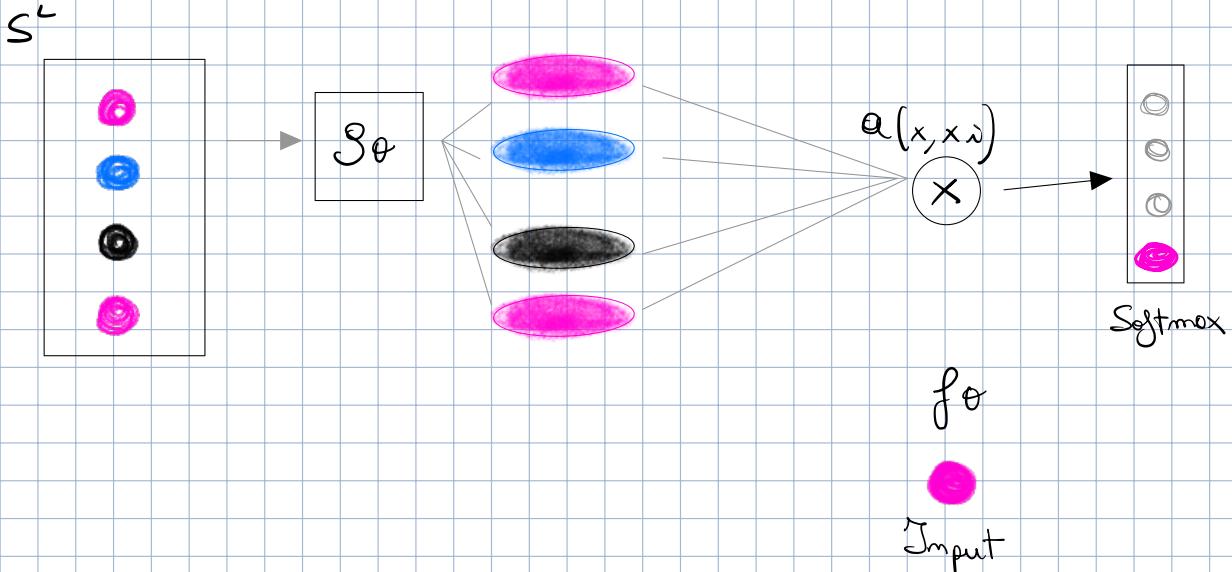
1 CONVOLUTIONAL SIAMESE NEURAL NETWORKS



$$\text{Prediction } \text{Class}_{S^L}(x) = \text{class}\left(\underset{x_i \in S^L}{\operatorname{argmax}} P(x, x_i)\right)$$

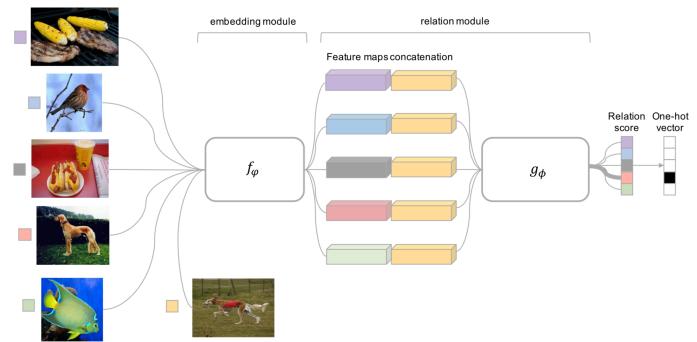
2 MATCHING NETWORKS

$$\text{Prediction } \text{class}_{S^L}(x) = P(y|x, S^L) = \sum_{i=1}^K a(x, x_i) y_i$$



3

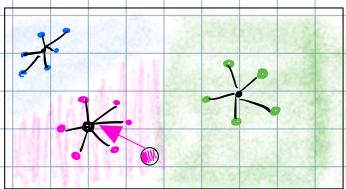
RELATION NETWORK



4

PROTO TYPICAL NETWORKS

Compute embeddings & class sample in S^c
take the mean

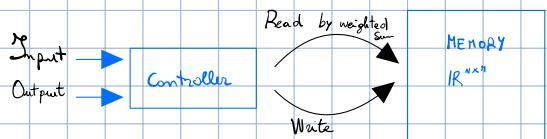


$$\text{Centroid for class } c \quad V_c = \frac{1}{|S^c \in S^L|} \sum_{(x_i, y_i) \in S^c} f_\theta(x_i)$$

MODEL BASED META-LEARNING

START FROM NEURAL TURING MACHINES

2 Components : CONTROLLER (NN)



MEMORY ($N \times M$) Matrix

Read by weighted sum

$$r_t = \sum_{i=1}^N w_t(i) M_t(i)$$

attention Score

$$\text{The sum of the attention over the memory} = 1 \quad \sum_{i=1}^N w_t(i) = 1 \quad \forall i \quad 0 \leq w_t(i) \leq 1$$

Write like LSTM

① Erase some content from M according to an erase vector

$$\tilde{M}_t(i) = M_{t-1}(i) [1 - w_t(i) e_t]$$

② Add Some content by an add vector

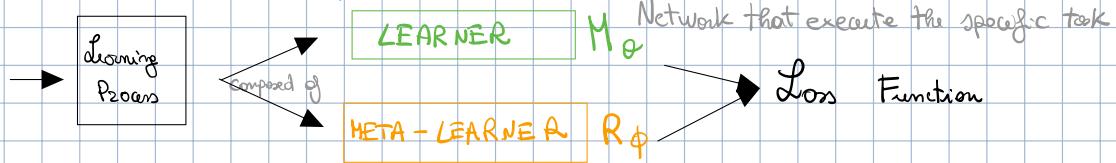
$$M_t(i) = \tilde{M}_t(i) + w_t(i) a_t$$

OPTIMIZATION - BASED META LEARNING

Gradient - Descent needs
 → A lot of samples to learn
 → A lot of steps to converge

- How can a GD approach work with few examples?

The Meta-Learner Approach



LSTM - MODEL
 Has the role of update
 θ parameters (learner's params)
 using a small support set
 such that the learner lead
 to high performance

the meta learner introduce a function $g_t(\nabla f(\theta_t), \phi)$ this function depends on the gradient of the learner and the parameter of the meta learner itself

Update the learner by GD using the g function to define how to actually update the weights

$$\theta_{t+1} = \theta_t + g_t(\nabla f_t(\theta_t), \phi)$$

This means that the ^(best) ~~parameters~~ of the learner θ^*
 Depends on the parameters of the learner itself but
 also on the parameters of the meta learner ϕ

$$\theta^*(f_\phi, \phi)$$

How we define the goodness of the meta-learner?

$$\text{Considering its loss } L(\phi) = \mathbb{E}_f [f(\theta^*(f, \phi))]$$

This means that the good meta learner is the one that give a low loss computed by the outcome that the learner function gives when makes use of the best parameters θ^* that depends on f and ϕ

g_t is an RNN parametrized by ϕ
 in RNN we have

- h_t : HIDDEN STATES
- trajectory length : history T (horizon) of evolution of gradient of f

$$L(\phi) = \mathbb{E}_f \left[\sum_{t=1}^T w_t f(\theta_t) \right]$$

arbitrary weight for time t

$$\text{update } \theta_{t+1} = \theta_t + g_t$$

$$\text{and } \begin{bmatrix} g_t \\ h_{t+1} \end{bmatrix} = \text{meta learner Model } (\nabla f_t, h_t, \phi)$$

↑
 the current
 gradient of f
 $\nabla \theta f(\theta_t)$

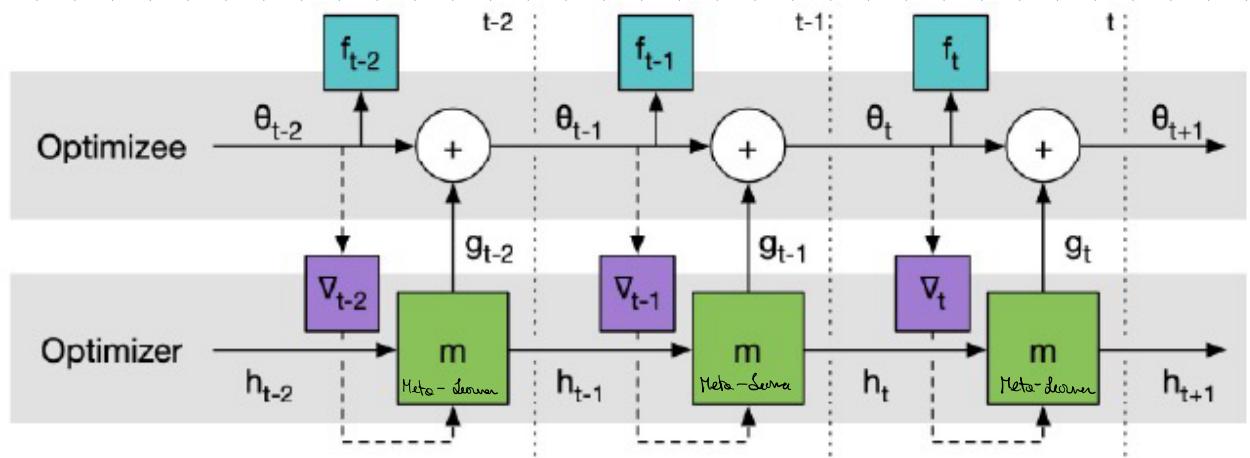


Figure 2: Computational graph used for computing the gradient of the optimizer.

MAML : Model Agnostic Meta-Learner

Really simple idea : the parameters of the learner θ are updated using a batch of training for a given task \mathcal{Z}_i . With f and a task \mathcal{Z}_i

$$\text{we can update } \theta_i' = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{Z}_i}(f_{\theta})$$

the new set of parameters of the learner for the given task i

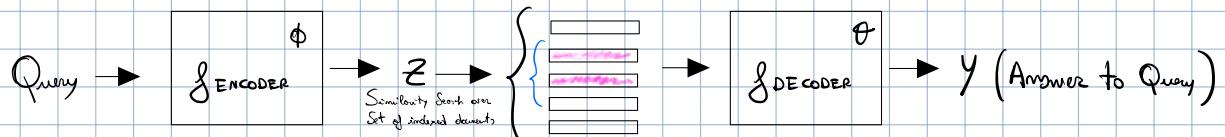
When all the θ_i' are computed for all the tasks, then we can update

$$\theta = \theta - \beta \nabla_{\theta} \sum_{\mathcal{Z}_i \sim P(\mathcal{Z})} \mathcal{L}_{\mathcal{Z}_i}(f_{\theta_i'})$$

CAPITOLO 14

• Make an LLM Access external information

Retrieval-Augmented Generation



$$\text{PRAG } (y|x) = \sum_{z \in \text{Top } k \text{ sim}(\cdot, x)} p_{\phi}(z|x) \cdot p_{\theta}(y|z)$$

$$p_{\phi}(z|x) = \exp(d(z)^T \cdot f_{\phi}(x))$$

$$\mathcal{L}_{\text{RAG}}(x, z^+, z_i^- \dots z_j^-, \phi, \theta) = \frac{\text{sim}(x, z^+)}{\text{sim}(x, z^+) + \sum_{j=1}^m \text{sim}(x, z_j^-)}$$

- DPR → Backprop f_{ENCODER} (only for query), use a Frozen Generator
- RAG → Backprop f_{ENCODER} query, freeze doc encoder,
- REALM → Doc Encoder and f_{ENCODER} both trained

CAPITOLO 15

Deal with Noisy Label

SOURCE OF NOISE

- Human Error
- Ambiguity of the task
- Multiple correct labels

The training procedure has as objective minimize the empirical risk ($E_{\mathcal{D}}$) over the dataset \mathcal{D}

$$\text{Minimize } R_{\mathcal{D}}(\hat{f}) = E_{\mathcal{D}} \mathbb{I}(\hat{y} \neq \underset{x \in \{1..C\}}{\text{argmax}} f(x), y)$$

the best classifier will be

$$\hat{f}^* \in \arg \min_{\hat{f}} R_{\mathcal{D}}(\hat{f}) = \frac{1}{m} \sum_{i=1}^m \ell(f(x_i), y_i)$$

ERM for noisy dataset

$$\hat{R}_{\mathcal{D}, \epsilon}(f) = \frac{1}{m} \sum_{i=1}^m \ell(f(x_i), \tilde{y}_i)$$

Quantify the noise with the Noise Transition Matrix

Observation = Noise transition Matrix $\times T(x)$

$$\begin{bmatrix} P_1(\tilde{y}=1|x=x) \\ \vdots \\ P_2(\tilde{y}=C|x=x) \end{bmatrix} = \begin{bmatrix} P_1(\tilde{y}=1|y=1, x=x) & \dots & P_1(\tilde{y}=1|y=C, x=x) \\ \vdots & \ddots & \vdots \\ P_2(\tilde{y}=C|y=1, x=x) & \dots & P_2(\tilde{y}=C|y=C, x=x) \end{bmatrix} \cdot \begin{bmatrix} P(y=1|x=x) \\ \vdots \\ P(y=C|x=x) \end{bmatrix}$$

$$T(x) = \begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.2 & 0.7 & 0.1 \\ 0.1 & 0.2 & 0.7 \end{bmatrix}$$

APPROACHES TO DEAL WITH NOISE

• Estimate the Transition matrix

• Robot Loss

• Sample Weighting / Sample Selection

• Ensemble

• Meta-Learning

More than one network

Classifier trained on noisy dataset

Lead to the same even of the classifier trained on clean data

• GENERALIZED CROSS ENTROPY

$$L_g(f(x), e_j) = -(1 - f_j(x))^2$$

One hot of class j
J-th component of the output prob. distribution

• EARLY-LEARNING REGULARIZ. LOSS

$$ELR = ECE + \lambda \sum_{i=1}^m \log \left(1 - \langle p_i^j, q_i^j \rangle \right)$$

steps in training
prediction target sample i

• Regularization

Mix-Up Data Augment

$$\begin{aligned} x &= \lambda x_1 + (1-\lambda) x_2 && \text{overlap} \\ y &= \lambda y_1 + (1-\lambda) y_2 && \text{images and labels} \end{aligned}$$

$$\lambda \approx 0.5$$

Sample Selection

• Co-Teaching (2 Models w_f w_g)

• Get loss of samples for model w_f and model w_g

• Select Sample with loss $< R(T)$ D_f and \tilde{D}_g

• Update $w_f = w_f - \alpha \nabla \ell(f, D_f)$

$$w_g = w_g - \alpha \nabla \ell(g, \tilde{D}_g)$$

• Update the threshold $R(T)$

• Ensemble

More than one network

• Early Learning + Memorization

Learning easy pattern (clean samples) ↗

The model starts to memorize the noise