# The Mathematics of Neural Networks

Gabriel Peyré

**www.numerical-tours.com**

# Overview

- **Empirical Risk Minimization**

- Perceptrons

- Optimization

- Convolutional Networks

- Residual Networks

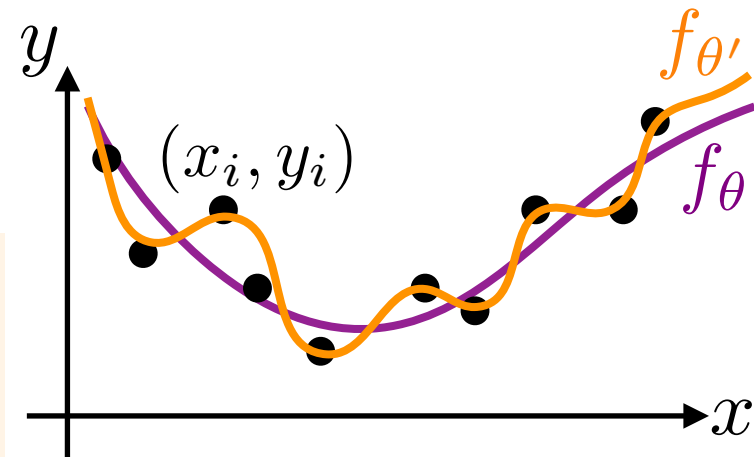- Transformers

# Empirical Risk Minimization

Dataset: $(x_i, y_i)_{i=1}^n$.     $x_i \in \mathbb{R}^d$     $y_i \in \mathbb{R}$

Prediction:    $y \approx f_\theta(x)$
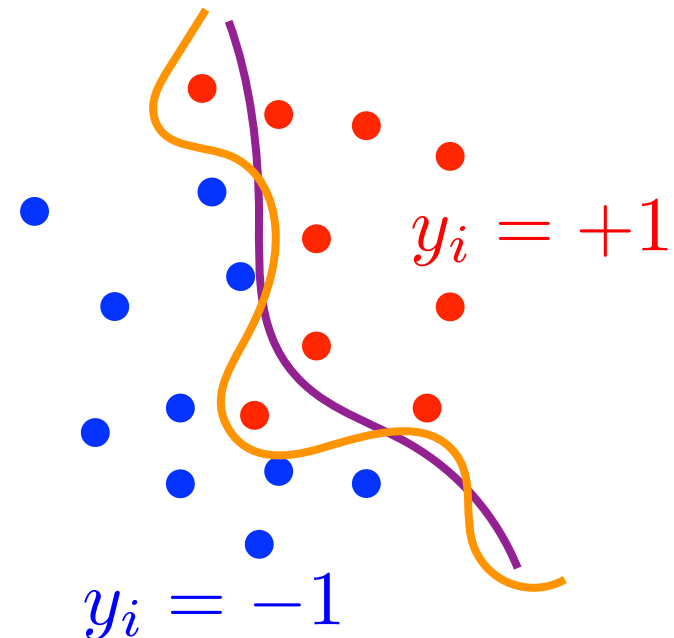
Empirical risk minimization:

$$\min_\theta \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i)$$

*Least square:* $\ell(y, y') = (y - y')^2$

# Empirical Risk Minimization

Dataset: $(x_i, y_i)_{i=1}^n$. $\quad x_i \in \mathbb{R}^d \quad\quad y_i \in \mathbb{R}$

Prediction: $\quad y \approx f_\theta(x)$

Empirical risk minimization:

$$\min_\theta \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i)$$

*Least square:* $\ell(y, y') = (y - y')^2$

Classification: $\quad y_i \in \{-1, 1\}$

$y \approx \text{sign}(f_\theta(x))$

*Logistic:* $\ell(y, y') = \log(1 + e^{-yy'})$

# Empirical Risk Minimization

Dataset: $(x_i, y_i)_{i=1}^n$.     $x_i \in \mathbb{R}^d$     $y_i \in \mathbb{R}$

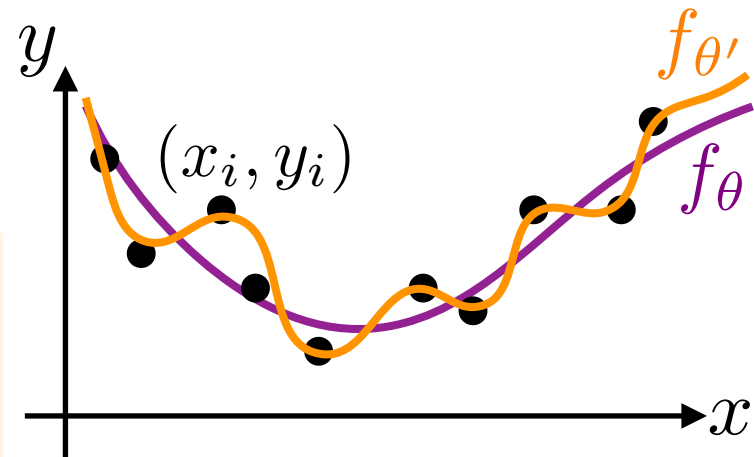Prediction:    $y \approx f_\theta(x)$

Empirical risk minimization:

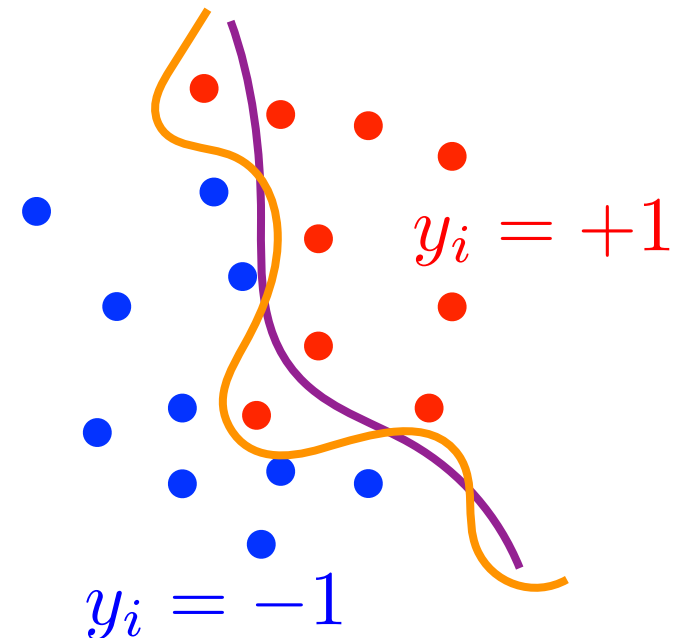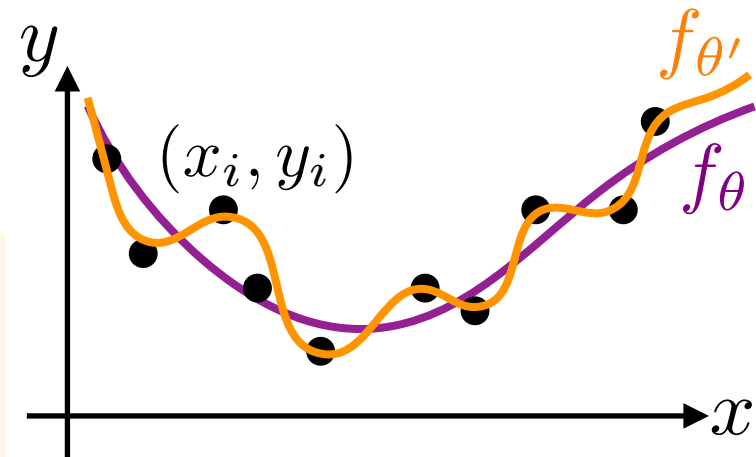$$\min_\theta \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i)$$

*Least square:* $\ell(y, y') = (y - y')^2$

Classification:    $y_i \in \{-1, 1\}$
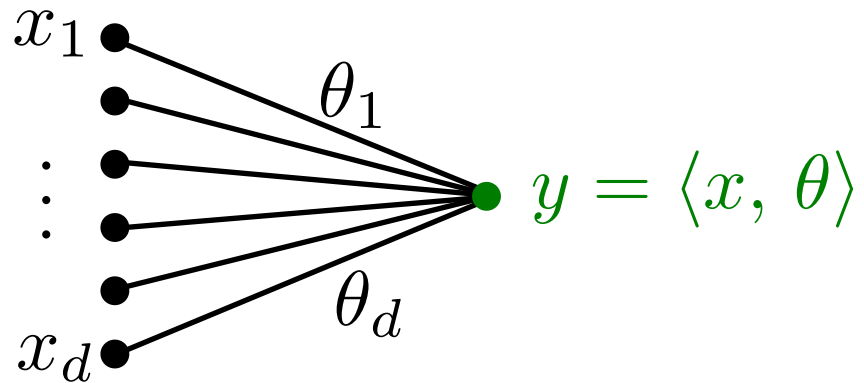
$y \approx \text{sign}(f_\theta(x))$

*Logistic:* $\ell(y, y') = \log(1 + e^{-yy'})$
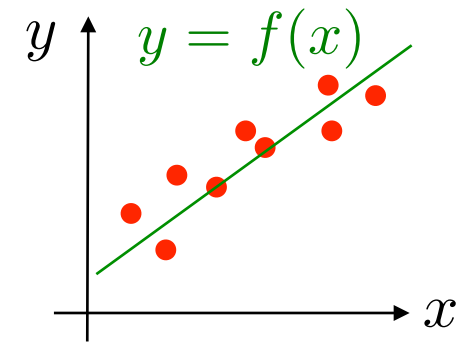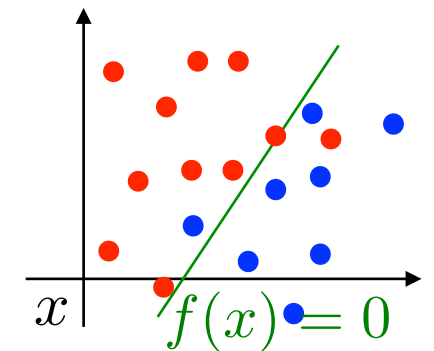
Overfitting, regularization, ...

# Linear model (1 layer)

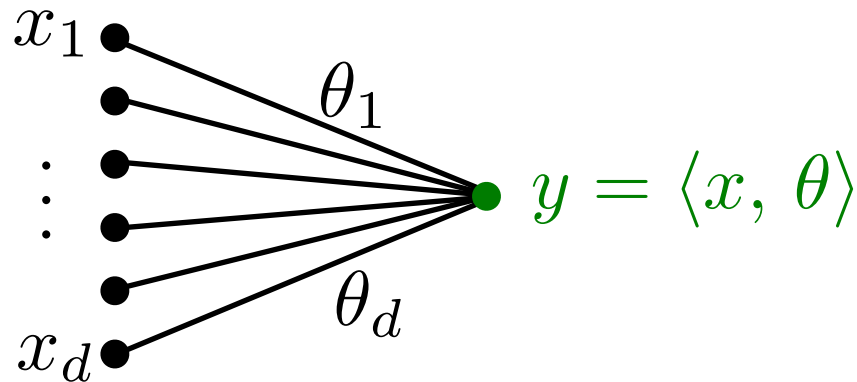$$f_\theta(x) = \langle x, \theta \rangle = \sum_k x_k \theta_k$$



$x_1$

$\theta_1$

$\vdots$

$y = \langle x, \theta \rangle$

$\theta_d$

$x_d$

Regression

$y$

$y = f(x)$

$x$

Classification:

$x$

$f(x) = 0$

# Linear model (1 layer)

$$f_\theta(x) = \langle x, \theta \rangle = \sum_k x_k \theta_k$$



$$y = \langle x, \theta \rangle$$

Regression



$y = f(x)$

Convex optimization:

$$\min_\theta \triangleq \frac{1}{n} \sum_{i=1}^{n} \ell(\langle x_i, \theta_i \rangle, y_i)$$

Classification:



$f(x) = 0$

# Linear model (1 layer)

$$f_\theta(x) = \langle x, \theta \rangle = \sum_k x_k \theta_k$$

Regression

$y = f(x)$

$y = \langle x, \theta \rangle$

Convex optimization:

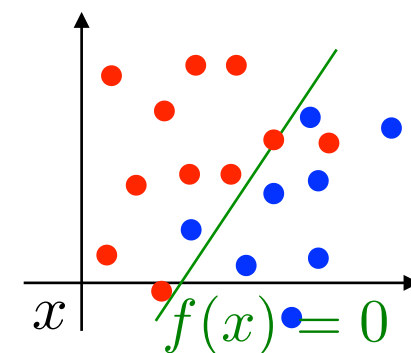$$\min_\theta \triangleq \frac{1}{n} \sum_{i=1}^{n} \ell(\langle x_i, \theta_i \rangle, y_i)$$

Classification:

$f(x) = 0$

*Kernel methods:* replace $x$ by $\varphi(x) \in \mathbb{R}^D$

$(D \gg d$, even $D = \infty!)$

*Deep learning methods:* learn $\varphi(x)$!

# Overview

- Empirical Risk Minimization

- **Perceptrons**

- Optimization

- Convolutional Networks

- Residual Networks

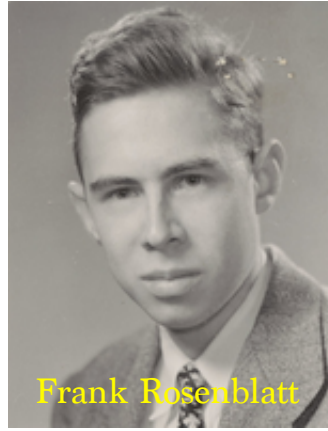- Transformers

# Multi-layer Perceptron



$$x_0 \leftarrow x$$

$$x_{k+1} \triangleq \sigma(W_k x_k + b_k)$$

$$\theta = \{(W_k, b_k)\}_{k=0}^{D-1}$$

$$W_k \in \mathbb{R}^{d_{k+1} \times d_k}$$

$$b_k \in \mathbb{R}^{d_{k+1}}$$

Frank Rosenblatt

# Multi-layer Perceptron



$$x_0 \leftarrow x$$

$$x_{k+1} \triangleq \sigma(W_k x_k + b_k)$$

Frank Rosenblatt

$$\theta = \{(W_k, b_k)\}_{k=0}^{D-1}$$

$$W_k \in \mathbb{R}^{d_{k+1} \times d_k}$$

$$b_k \in \mathbb{R}^{d_{k+1}}$$

$$f_\theta(x_0) = x_D$$

*Non-linearity:* $\sigma$ must be non-polynomial to increase expressivity.



Sigmoid

ReLu

# Multi-layer Perceptron



$$x_0 \leftarrow x$$
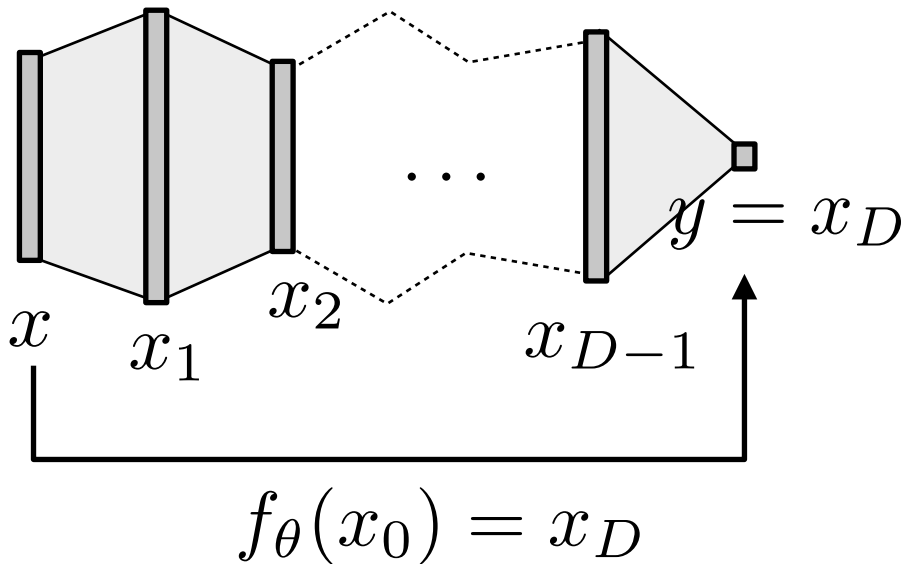
$$x_{k+1} \triangleq \sigma(W_k x_k + b_k)$$

$$f_\theta(x_0) = x_D$$

$$\theta = \{(W_k, b_k)\}_{k=0}^{D-1}$$

$$W_k \in \mathbb{R}^{d_{k+1} \times d_k}$$

$$b_k \in \mathbb{R}^{d_{k+1}}$$

Frank Rosenblatt

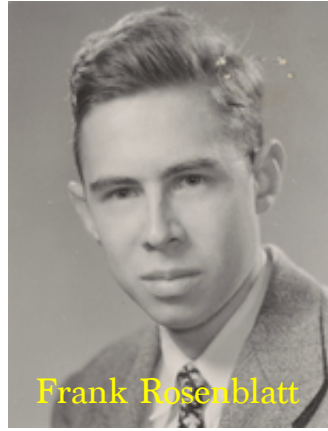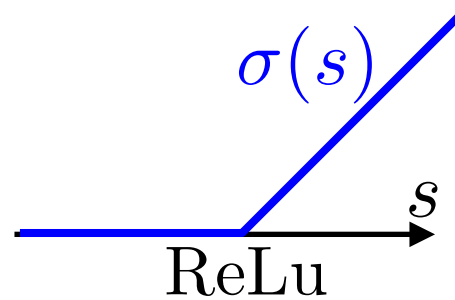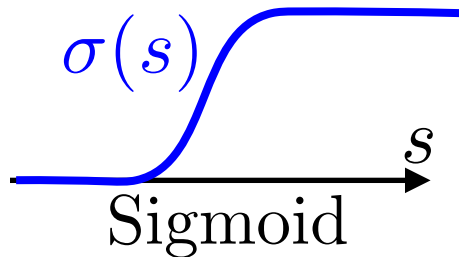*Non-linearity:* $\sigma$ must be non-polynomial to increase expressivity.



Sigmoid

ReLu

*Weight matrix:* needs extra constraints (e.g. convolution & sub-sampling)

# Two Layers Perceptron



$$f_\theta(x) \triangleq \sum_{s=1}^{p} a_s \sigma(\langle x, w_s \rangle + b_s)$$

$\rightarrow$ sum of "ridge" functions $\sigma(\langle x, w \rangle + b)$

Input $y = f(x)$     $p = 6$ neurons     $p = 30$ neurons     $p = 100$ neurons

# Two Layers Perceptron



$$f_\theta(x) \triangleq \sum_{s=1}^{p} a_s \sigma(\langle x, w_s \rangle + b_s)$$

$\rightarrow$ sum of "ridge" functions $\sigma(\langle x, w \rangle + b)$

Input $y = f(x)$ | $p = 6$ neurons | $p = 30$ neurons | $p = 100$ neurons

# Universality of Perceptrons

*Theorem:* If $f$ is continuous on a compact $\Omega$, for all $\varepsilon > 0$
for $p$ large enough, there exists $\theta$ such that
$$\forall\, x \in \Omega,\ |f_\theta(x) - f(x)| \leqslant \varepsilon$$

$\rightarrow$ non quantitative ... no free lunch.

George Cybenko

# Universality of Perceptrons

*Theorem:* If $f$ is continuous on a compact $\Omega$, for all $\varepsilon > 0$
for $p$ large enough, there exists $\theta$ such that
$$\forall\, x \in \Omega,\ |f_\theta(x) - f(x)| \leqslant \varepsilon$$

$\rightarrow$ non quantitative ... no free lunch.

Barron's functions:   $\|f\|_B \triangleq \int_{\mathbb{R}^d} \|\omega\| |\hat{f}(\omega)| \mathrm{d}\omega < +\infty$

*Theorem:* for $p$ large, there exists $\theta$ such that
$$\|f - f_\theta\|_{L^2(\Omega)} \leqslant \frac{2\mathrm{diam}(\Omega)\|f\|_B}{\sqrt{p}}$$

George Cybenko

Andrew Barron

# Universality of Perceptrons

*Theorem:* If $f$ is continuous on a compact $\Omega$, for all $\varepsilon > 0$ for $p$ large enough, there exists $\theta$ such that
$$\forall\, x \in \Omega,\ |f_\theta(x) - f(x)| \leqslant \varepsilon$$

$\rightarrow$ non quantitative ... no free lunch.

Barron's functions: $\quad \|f\|_B \triangleq \int_{\mathbb{R}^d} \|\omega\| |\hat{f}(\omega)| \mathrm{d}\omega < +\infty$

*Theorem:* for $p$ large, there exists $\theta$ such that
$$\|f - f_\theta\|_{L^2(\Omega)} \leqslant \frac{2\operatorname{diam}(\Omega)\|f\|_B}{\sqrt{p}}$$

$\rightarrow$ non-constructive.



$\rightarrow$ for $p$ "large enough" gradient descent works [Chizat-Bach 2018]

George Cybenko

Andrew Barron

# Overview

- Empirical Risk Minimization

- Perceptrons

- **Optimization**

- Convolutional Networks

- Residual Networks

- Transformers

# Gradient Descent

$$\min_\theta \mathcal{E}(\theta) \triangleq \frac{1}{n} \sum_{i=1}^{n} \ell(f_\theta(x_i), y_i)$$

Gradient descent:
$$\theta_{\ell+1} = \theta_\ell - \tau_\ell \nabla \mathcal{E}(\theta_\ell)$$



Small $\tau_\ell$      Large $\tau_\ell$      Optimal $\tau_\ell = \tau_\ell^\star$

# Gradient Descent

$$\min_{\theta} \mathcal{E}(\theta) \triangleq \frac{1}{n} \sum_{i=1}^{n} \ell(f_\theta(x_i), y_i)$$

Gradient descent:
$$\theta_{\ell+1} = \theta_\ell - \tau_\ell \nabla \mathcal{E}(\theta_\ell)$$

Small $\tau_\ell$

Large $\tau_\ell$

Optimal $\tau_\ell = \tau_\ell^\star$

Stochastic gradient descent:

$$\theta_{\ell+1} = \theta_\ell - \frac{\tau}{\ell} \nabla \mathcal{E}_\ell(\theta_\ell)$$

$$i \leftarrow \texttt{rand}$$

$$\mathcal{E}_\ell(\theta) \triangleq \ell(f_\theta(x_i), y_i)$$

Herbert Robbins

Sutton Monro

# The Complexity of Gradient Computation

**Setup:** $\mathcal{E} : \mathbb{R}^d \to \mathbb{R}$ computable in $K$ operations.

```python
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \dots)$

and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^d \to \mathbb{R}^d$?

# The Complexity of Gradient Computation

**Setup:** $\mathcal{E} : \mathbb{R}^d \to \mathbb{R}$ computable in $K$ operations.

```python
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \dots)$

and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^d \to \mathbb{R}^d$?

Finite differences:
$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon}(\mathcal{E}(\theta + \varepsilon\delta_1) - \mathcal{E}(\theta), \dots \mathcal{E}(\theta + \varepsilon\delta_d) - \mathcal{E}(\theta))$$

$K(d+1)$ operations, intractable for large $d$.

# The Complexity of Gradient Computation

**Setup:** $\mathcal{E} : \mathbb{R}^d \to \mathbb{R}$ computable in $K$ operations.

```python
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \ldots)$
and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^d \to \mathbb{R}^d$?

Finite differences:
$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon}(\mathcal{E}(\theta + \varepsilon\delta_1) - \mathcal{E}(\theta), \ldots \mathcal{E}(\theta + \varepsilon\delta_d) - \mathcal{E}(\theta))$$
$K(d+1)$ operations, intractable for large $d$.

*Theorem:* there is an algorithm to compute $\nabla \mathcal{E}$ in $O(K)$ operations.
[Seppo Linnainmaa, 1970]

This algorithm is reverse mode automatic differentiation

```python
def BackwardNN(A,b,X):
    gx = lossG(X[R],Y) # initialize the gradient
    for r in arange(R-1,-1,-1):
        M = rhoG( A[r].dot(X[r]) + tile(b[r],[1,n]) ) * gx
        gx = A[r].transpose().dot(M)
        gA[r] = M.dot(X[r].transpose())
        gb[r] = MakeCol(M.sum(axis=1))
    return [gA,gb]
```



Seppo Linnainmaa

# Computational Graph

Computer program $\Leftrightarrow$ directed acyclic graph $\Leftrightarrow$ linear ordering of nodes $(\theta_r)_r$

# Example

$$\ell(\theta_1, \theta_2) \stackrel{\text{def.}}{=} \theta_2 e^{\theta_1} \sqrt{\theta_1 + \theta_2 e^{\theta_1}}$$



$\theta_1$

$\theta_2$

input

$g_3$  $\theta_3 \stackrel{\text{def.}}{=} e^{\theta_1}$

$g_4$  $\theta_4 \stackrel{\text{def.}}{=} \theta_2 \theta_3$

$g_5$  $\theta_5 \stackrel{\text{def.}}{=} \theta_1 + \theta_4$

$g_6$  $\theta_6 \stackrel{\text{def.}}{=} \sqrt{\theta_5}$

$g_7$  $\theta_7 \stackrel{\text{def.}}{=} \theta_4 \theta_6$

output  $\ell$

# Example

$$\ell(\theta_1, \theta_2) \stackrel{\text{def.}}{=} \theta_2 e^{\theta_1} \sqrt{\theta_1 + \theta_2 e^{\theta_1}}$$



input

output $\ell$

Chain rules:

$$``\frac{\partial \theta_j}{\partial \theta_1} = \sum_{i \in \text{Parent}(j)} \frac{\partial \theta_j}{\partial \theta_i} \frac{\partial \theta_i}{\partial \theta_1}"$$

$$\partial_i g_j(\theta)$$

"Classical" evaluation: **forward**.
Complexity $\sim$ #inputs.

# Example

$$\ell(\theta_1, \theta_2) \overset{\text{def.}}{=} \theta_2 e^{\theta_1} \sqrt{\theta_1 + \theta_2 e^{\theta_1}}$$



Chain rules:

$$\text{``}\frac{\partial \theta_j}{\partial \theta_1} = \sum_{i \in \text{Parent}(j)} \frac{\partial \theta_j}{\partial \theta_i} \frac{\partial \theta_i}{\partial \theta_1}\text{''}$$

$$\partial_i g_j(\theta)$$

"Classical" evaluation: **forward**.
Complexity $\sim$ #inputs.

$$\text{``}\frac{\partial \theta_N}{\partial \theta_j} = \sum_{k \in \text{Child}(j)} \frac{\partial \theta_N}{\partial \theta_k} \frac{\partial \theta_k}{\partial \theta_j}\text{''}$$

$$\nabla_j \ell(\theta) \qquad \nabla_k \ell(\theta) \qquad \partial_j g_k(\theta)$$

**Backward** evaluation.
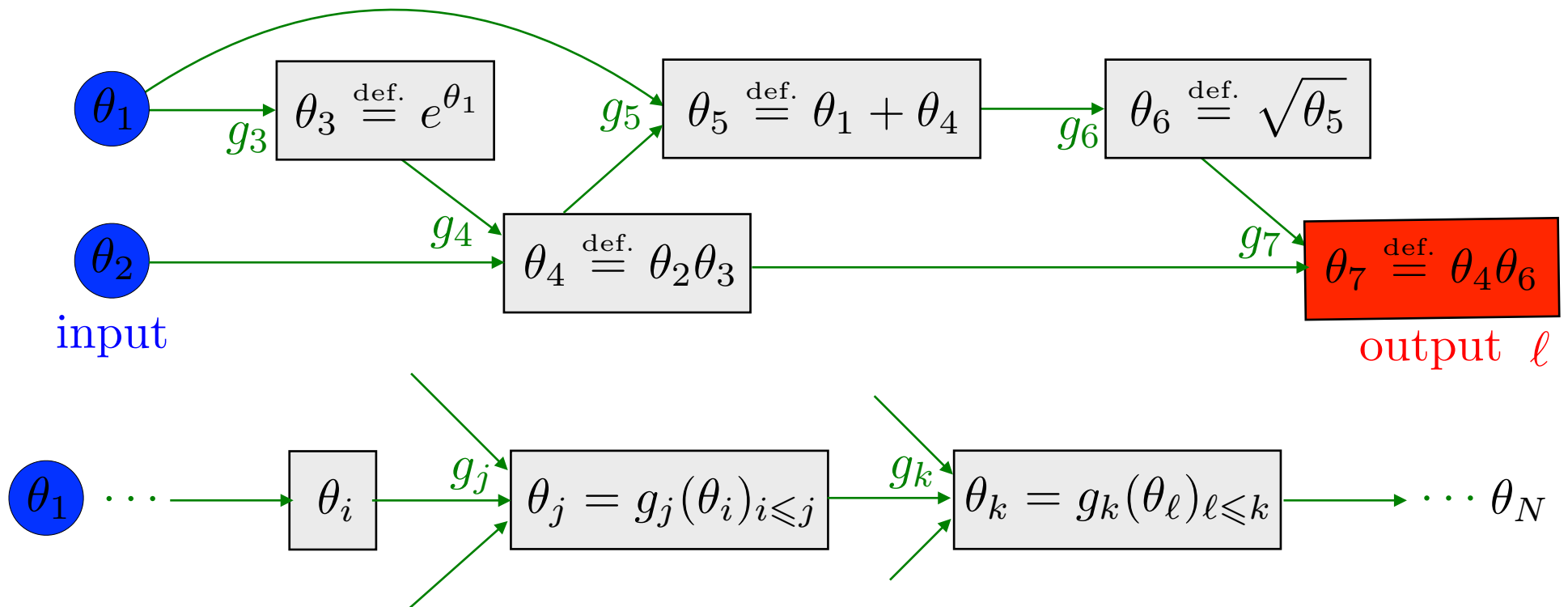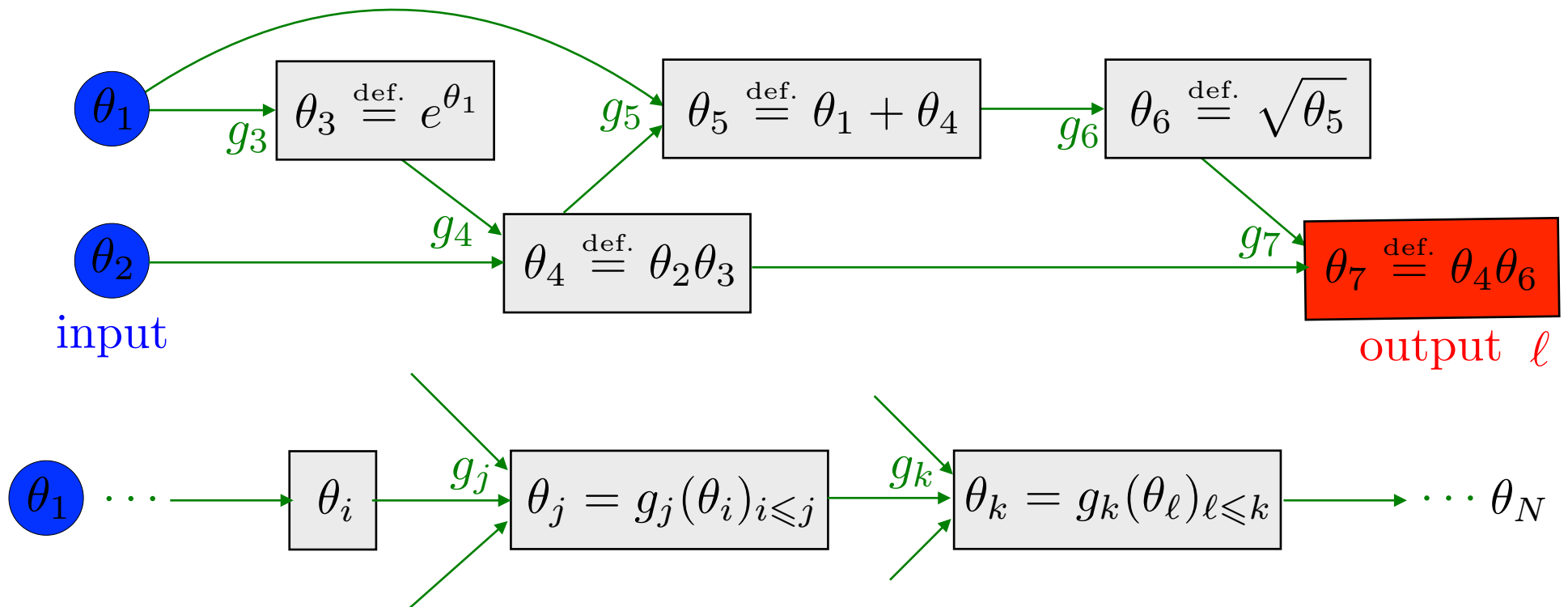Complexity $\sim$ #outputs (1 for grad).

# Backward Automatic Differentiation

$$\ell(\theta_1, \theta_2) \stackrel{\text{def.}}{=} \theta_2 e^{\theta_1} \sqrt{\theta_1 + \theta_2 e^{\theta_1}}$$

computing $\ell$



computing $\nabla\ell$

forward

```
function ℓ(θ₁, ..., θ_M)
    for r = M + 1, ..., R
    |   θ_r = g_r(θ_Parents(r))
    return θ_R
```

backward

```
function ∇ℓ(θ₁, ..., θ_M)
    ∇_R ℓ = 1
    for r = R − 1, ..., 1
    |   ∇_r ℓ = Σ_{s∈Child(r)} ∂_r g_s(θ) ∇_s ℓ
    return (∇₁ℓ, ..., ∇_M ℓ)
```

# Overview

- Empirical Risk Minimization

- Perceptrons

- Optimization

- **Convolutional Networks**

- Residual Networks

- Transformers

# Convolutional CNN



$$y = x_D \qquad x_{k+1} \triangleq \sigma(W_k x_k + b_k)$$

$\rightarrow$ Leverage translation invariance of images.

$\rightarrow$ Sub-sampling: breaks invariance but increase receptive fields.

# Convolutional CNN



$$y = x_D \qquad x_{k+1} \triangleq \sigma(W_k x_k + b_k)$$

$\to$ Leverage translation invariance of images.

$\to$ Sub-sampling: breaks invariance but increase receptive fields.



$227 \times 227 \times 3$

$55 \times 55 \times 96$

$27 \times 27 \times 256$

$13 \times 13 \times 384$  $13 \times 13 \times 384$  $13 \times 13 \times 256$

AlexNet, 2011

4096  4096  1000

# Example of Activations

# Overview

- Empirical Risk Minimization

- Perceptrons

- Optimization

- Convolutional Networks

- **Residual Networks**

- Transformers

ResNet-34

image

# ResNet-type Architectures [He et al' 16]

ResNet-34

image
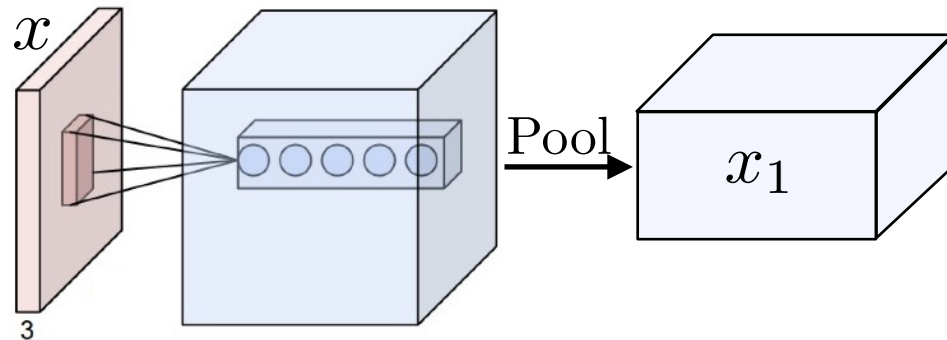
changes dimension

skip-connexion

$$x_t = x_{t-1} + v_{\theta_t}(x_{t-1})$$

# ResNet-type Architectures [He et al' 16]



ResNet-34

image

changes dimension

skip-connexion

$$x_t = x_{t-1} + v_{\theta_t}(x_{t-1})$$

$\rightarrow$ Makes the "infinite depth" limit non-degenerate.

$\rightarrow$ Enable $v_\theta = 0$ initialization, i.e. identity map.

$x_0$   $x_1$   $x_T$

$x_0$   $x_1$   $x_{2T}$

$x_0$   $x_1$   $x_{4T}$

ResNet [He et al, 2016]

$$\Phi_\theta(x_0) \triangleq x_T \quad \text{where}$$

$$x_{t+1} = x_t + \frac{1}{T} v_{\theta_t}(x_t)$$

# Infinite Depth and Neural-ODEs

# Infinite Depth and Neural-ODEs



ResNet [He et al, 2016]

$$\Phi_\theta(x_0) \triangleq x_T \quad \text{where}$$

$$x_{t+1} = x_t + \frac{1}{T} v_{\theta_t}(x_t)$$

$$T \to +\infty$$

Neural ODE [Chen et al, 2018]

$$\Phi_\theta(x(0)) \triangleq x(1) \quad \text{where}$$

$$\frac{\mathrm{d}x(t)}{\mathrm{d}t} = v_{\theta(t)}(x(t))$$

$x_0$ $x_1$ $x_T$

$x(0)$ $x(1)$

Trajectories cannot cross: $\Phi_\theta$ defines a diffeomorphism.

$T \to +\infty$ is a singular limit ($\theta$ can "explodes" during training)

# On the importance of scale and initialization

$$x_{t+1} = x_t + \frac{1}{T} v_{\theta_t}(x_t)$$

Zero/smooth initialization of $(\theta_t)_t$

$$\downarrow \quad T \to +\infty$$

Deterministic ODE

$$\frac{\mathrm{d}x(t)}{\mathrm{d}t} = v_{\theta(t)}(x(t))$$

$x(0)$
$x(1)$

$$x_{t+1} = x_t + \frac{1}{\sqrt{T}} v_{\theta_t}(x_t)$$

Random initialization of $(\theta_t)_t$

$$\downarrow \quad T \to +\infty$$

Stochastic ODE

$$\mathrm{d}x(t) = v_{\theta(t)}(x(t))\mathrm{d}t + \mathrm{d}W(t)$$

$x(0)$
$x(1)$
$x(1)$

[R. Cont, A. Rossier, R. Xu, 2022]

[P. Marion, Fermanian, Biau, Vert, 2022]

# Training Dynamic



Training: $\min_{\theta} f(\theta) \triangleq \frac{1}{N} \sum_{i=1}^{N} \| B \Phi_\theta (A x^i) - y^i \|^2$

Gradient descent: $\theta^{(k+1)} = \theta^{(k)} - \tau \nabla f(\theta^{(k)})$

$\rightarrow$ **No explicit regularization!**

# Training Dynamic



Training: $\quad \min_\theta f(\theta) \triangleq \frac{1}{N} \sum_{i=1}^N \| B\Phi_\theta(Ax^i) - y^i \|^2$

Gradient descent: $\quad \theta^{(k+1)} = \theta^{(k)} - \tau \nabla f(\theta^{(k)})$

$\rightarrow$ **No explicit regularization!**

*Question:* convergence of $\theta^k$ toward global minimum?

Neural tangent kernel [Jacot et al'18]:   local linear expansion.

Polyak-Łojasiewicz inequality [Liu, Zhu, Belkin 2021]:

$\quad \rightarrow$ conditionning might explodes as $T \rightarrow +\infty$.

$\quad \rightarrow$ find a suitable limit model and show "implicit" regularization effect.

Simplified analysis: [Barboni et al. 2022]

# Overview

- Empirical Risk Minimization

- Perceptrons

- Optimization

- Convolutional Networks

- Residual Networks

- **Transformers**

# Transformers

# Transformers

Token extraction
Vector encoding

Positional encoding

$x_1$
$x_2$

Points cloud
$\{x_i\}_i$

Demain, dès l'aube, à l'heure où blanchit la campagne, Je partirai. Vois-tu, je sais que tu m'attends. J'irai par la forêt, j'irai par la montagne. Je ne puis demeurer loin de toi plus longtemps. Je marcherai les yeux fixés sur mes pensées,…

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs

# Transformers

Token extraction
Vector encoding

Positional encoding

$x_1$
$x_2$
...

Points cloud
$\{x_i\}_i$

Demain, dès l'aube, à l'heure où blanchit la campagne, Je partirai. Vois-tu, je sais que tu m'attends. J'irai par la forêt, j'irai par la montagne. Je ne puis demeurer loin de toi plus longtemps. Je marcherai les yeux fixés sur mes pensées,…

Replace convolution by attention:

Residual

$x_j$
$x_i$

$$\tilde{x}_i = x_i + \frac{\sum_j A_{i,j}\, x_j}{\sum_j A_{i,j}}$$

Attention matrix: $A_{i,j} \triangleq e^{\langle K x_i, Q x_j \rangle}$

trained

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head Attention

N×

Positional Encoding

Input Embedding

Positional Encoding

Output Embedding

Inputs

Outputs

# Conclusion

Strong connexion with mathematical concepts:

– Going wider $\sim$ function approximation.

– Going deeper $\sim$ differential equations.

– Attention $\sim$ interacting particles.
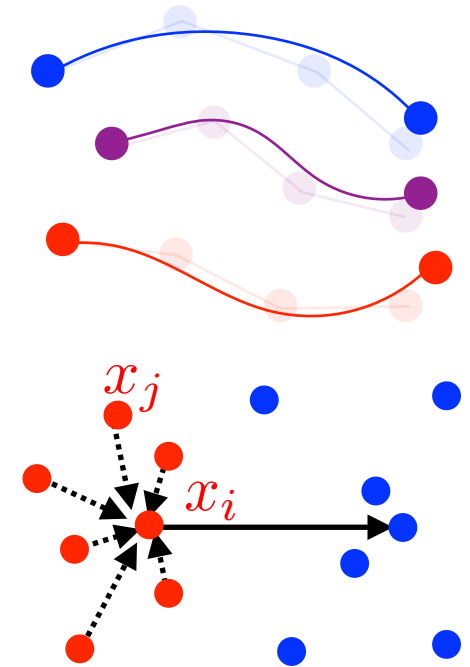
# Conclusion

Strong connexion with mathematical concepts:

– Going wider $\sim$ function approximation.

– Going deeper $\sim$ differential equations.

– Attention $\sim$ interacting particles.

Very limited theoretical understanding:

– Why gradient descent works?

– Implicit bias of architectures.
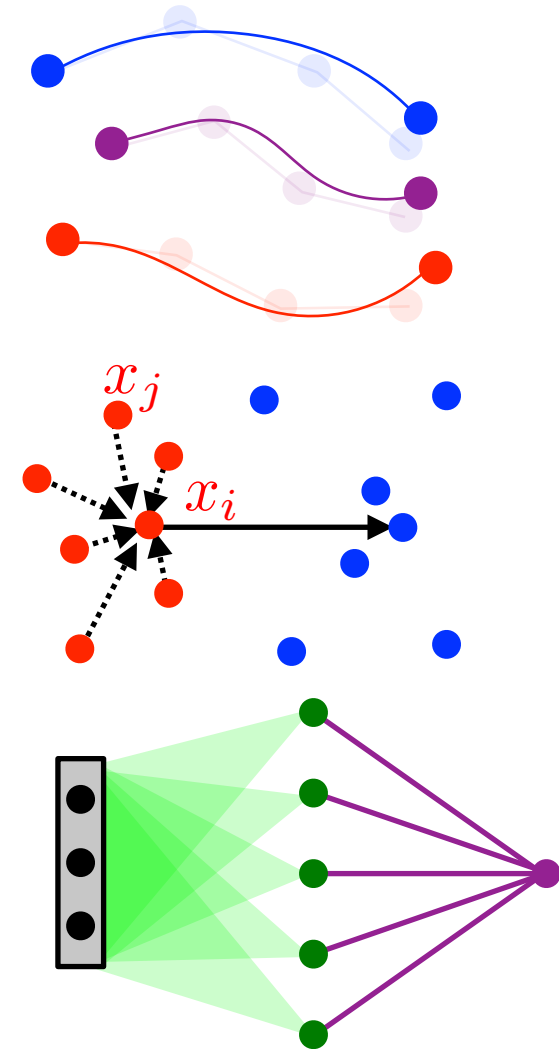
– Implicit bias of optimizers.

# Conclusion

Strong connexion with mathematical concepts:

– Going wider $\sim$ function approximation.

– Going deeper $\sim$ differential equations.

– Attention $\sim$ interacting particles.

Very limited theoretical understanding:

– Why gradient descent works?

– Implicit bias of architectures.

– Implicit bias of optimizers.

Examples of open problems:

– Why some optimizers (e.g. Adam) works for transformers?

– Mean field analysis of transformers (optimal transport?)