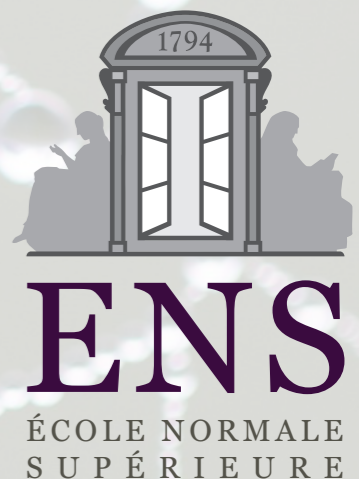


Les mathématiques de l'intelligence artificielle

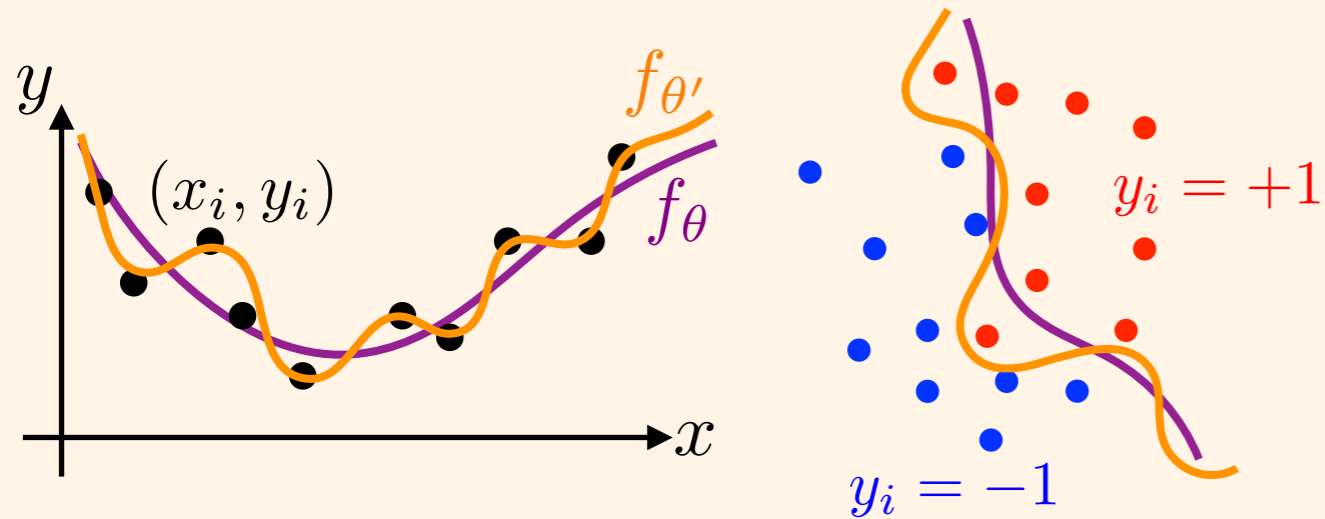
Gabriel Peyré

www.numerical-tours.com

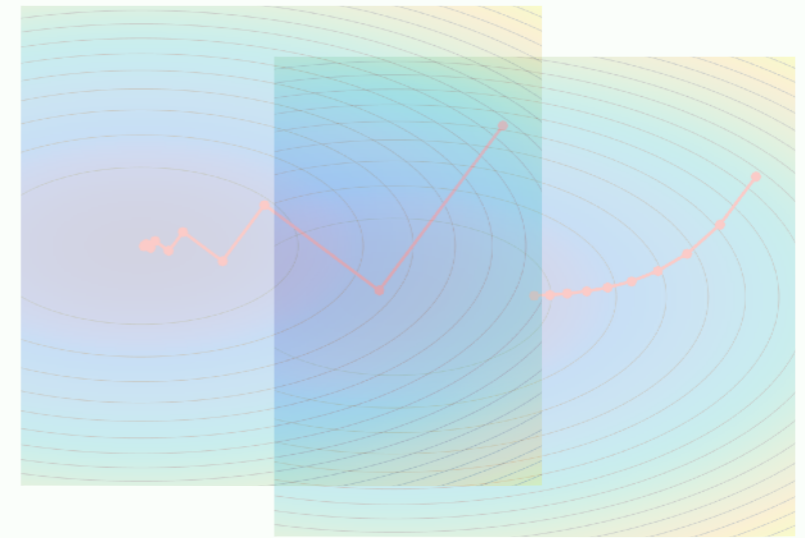


Overview

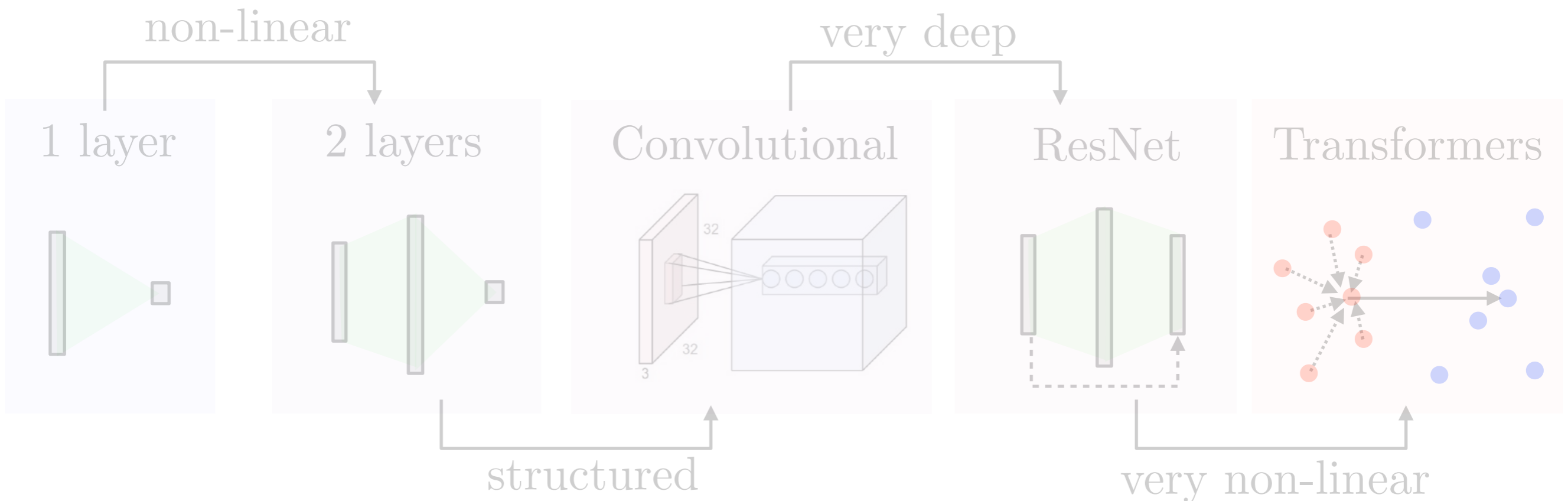
Machine Learning



Optimization

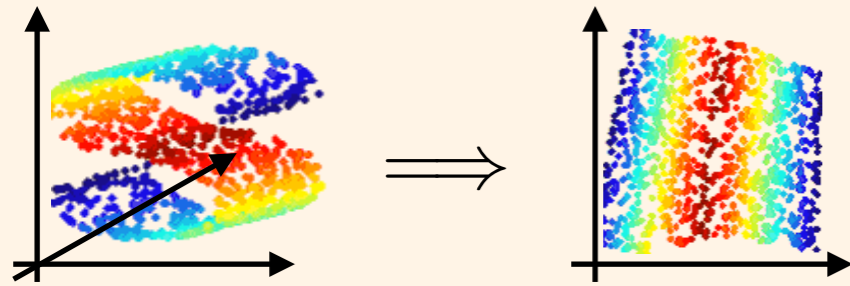


Networks Architectures



Supervised vs Unsupervised Learning

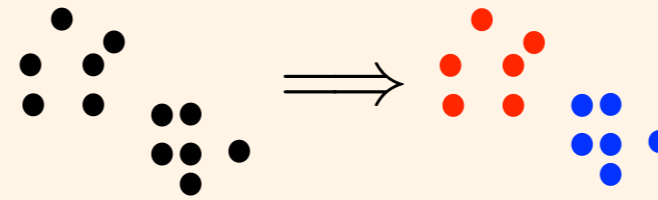
Dimensionality reduction



PCA, t-SNE, UMAP

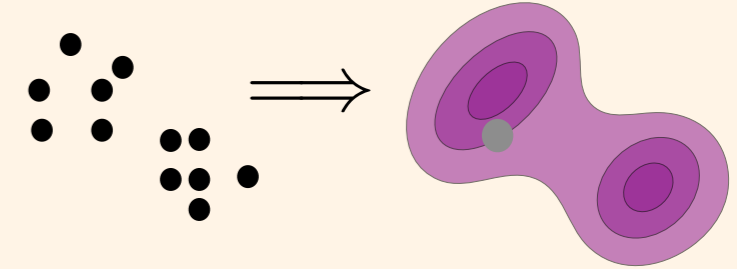
Un-supervised learning:

Clustering



k-means, DBScan

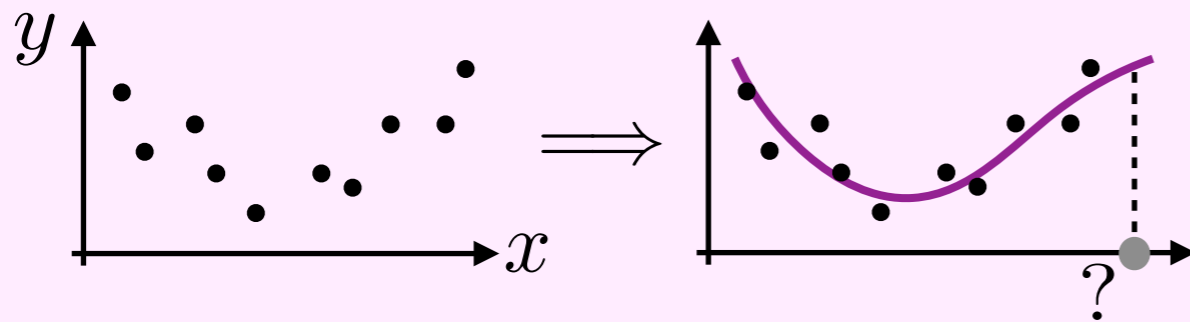
Generative modeling



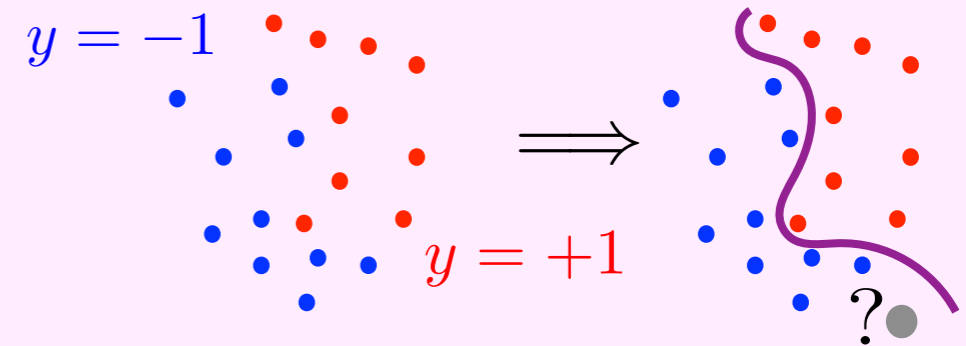
VAE, GAN, diffusion

Supervised learning:

Regression

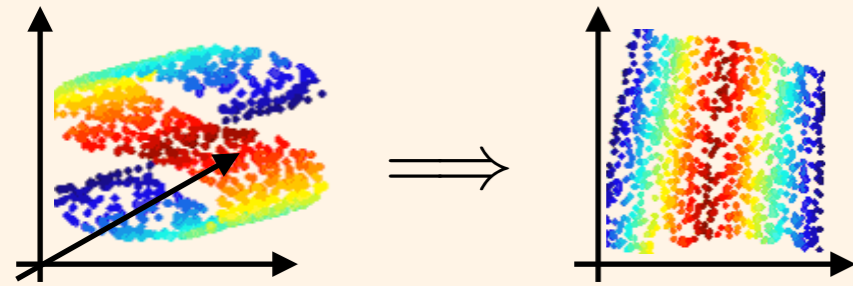


Classification



Supervised vs Unsupervised Learning

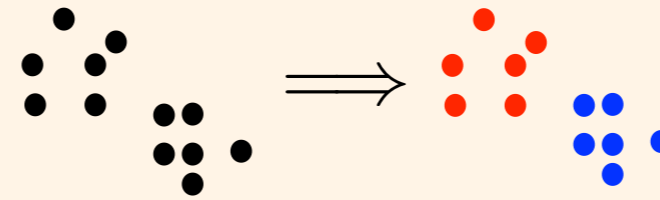
Dimensionality reduction



PCA, t-SNE, UMAP

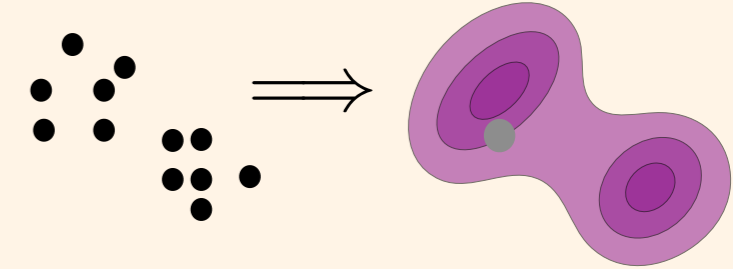
Un-supervised learning:

Clustering



k-means, DBScan

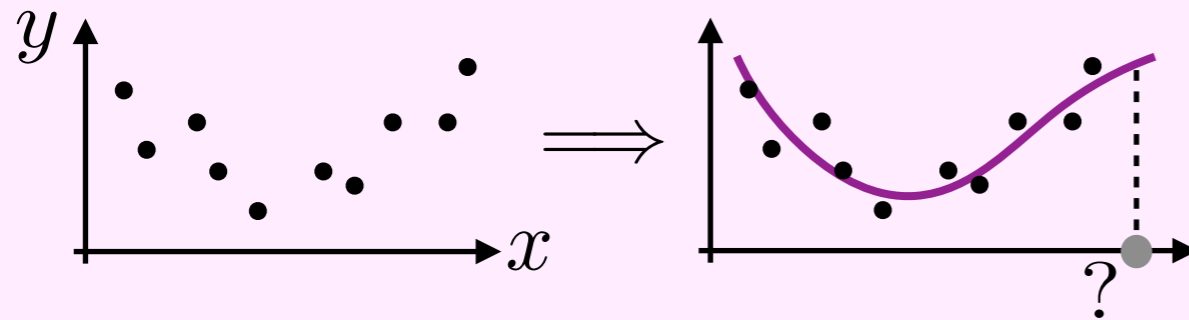
Generative modeling



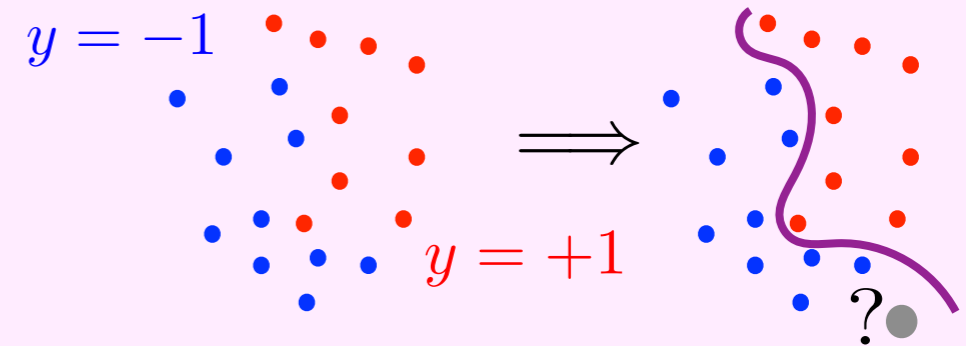
VAE, GAN, diffusion

Supervised learning:

Regression



Classification



Self-supervised learning:



L'École normale supérieure⁶, appelée aussi «de la rue d'Ulm», «Normale Sup'», est l'une des institutions universitaires et de recherche les plus prestigieuses et les plus sélectives de France, spécialisée en lettres et en sciences.

L'École normale supérieure⁶, appelée aussi «de la rue d'Ulm», «Normale Sup'», est l'une des institutions universitaires et de recherche les plus prestigieuses et les plus sélectives de France, spécialisée en lettres et en sciences.

L'École normale supérieure⁶, appelée aussi «de la rue d'Ulm», «Normale Sup'», est l'une des institutions universitaires et de recherche les plus prestigieuses et les plus sélectives de France, spécialisée en lettres et en sciences.

Add noise

Denoise

Masking

Next token prediction

Empirical Risk Minimization

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

Empirical risk minimization:

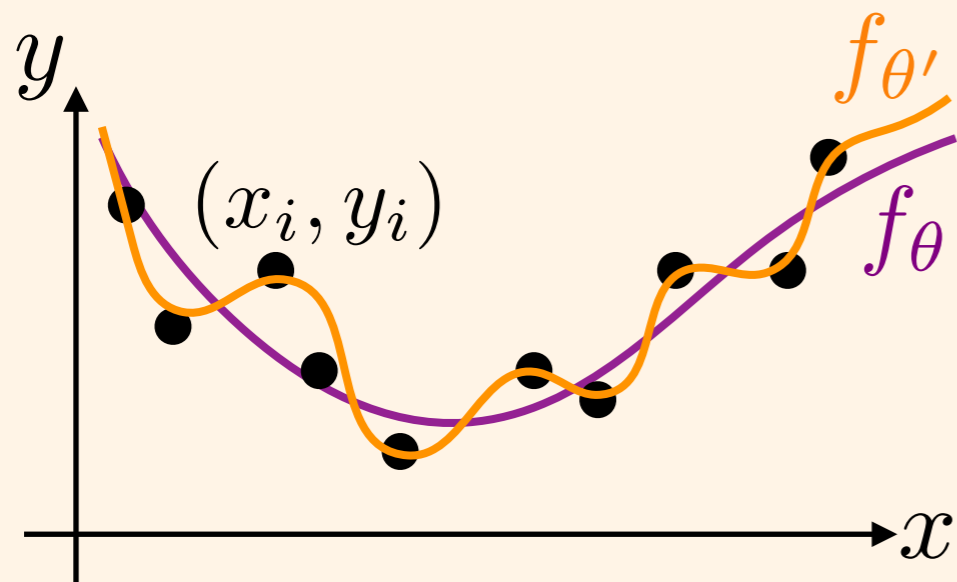
$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i)$$

Regression: $y_i \in \mathbb{R}$

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

Least square:

$$\ell(y, y') = (y - y')^2$$



Empirical Risk Minimization

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

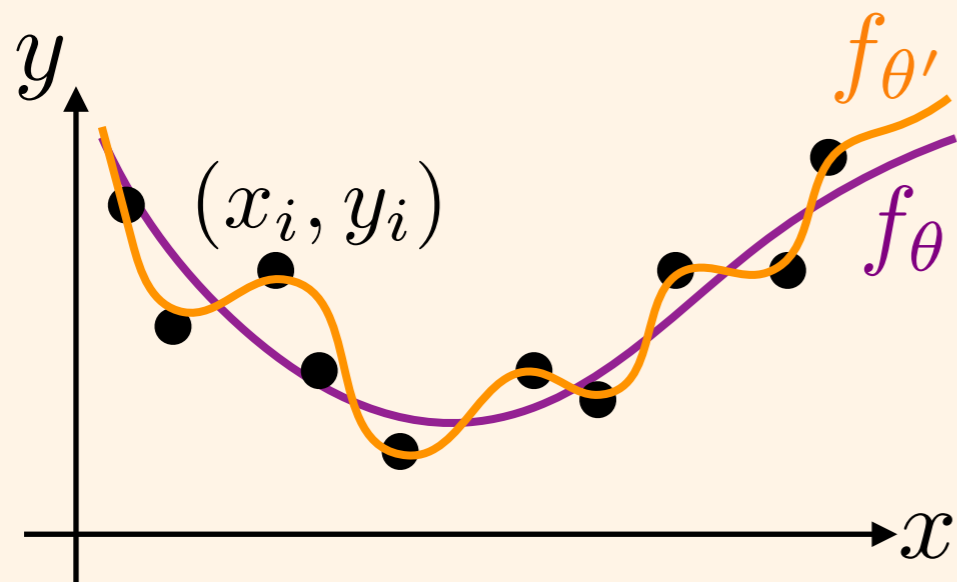
Empirical risk minimization:

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

Regression: $y_i \in \mathbb{R}$
 $y \approx f_{\theta}(x)$

Least square:

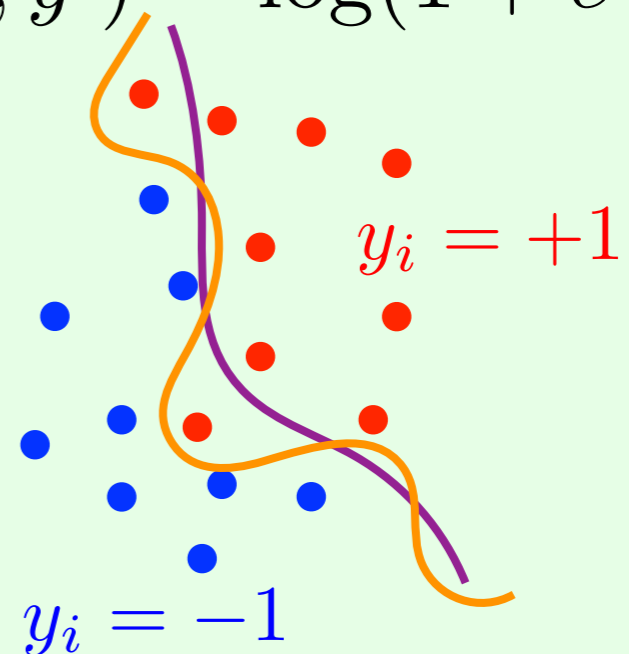
$$\ell(y, y') = (y - y')^2$$



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

Logistic:

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



Empirical Risk Minimization

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

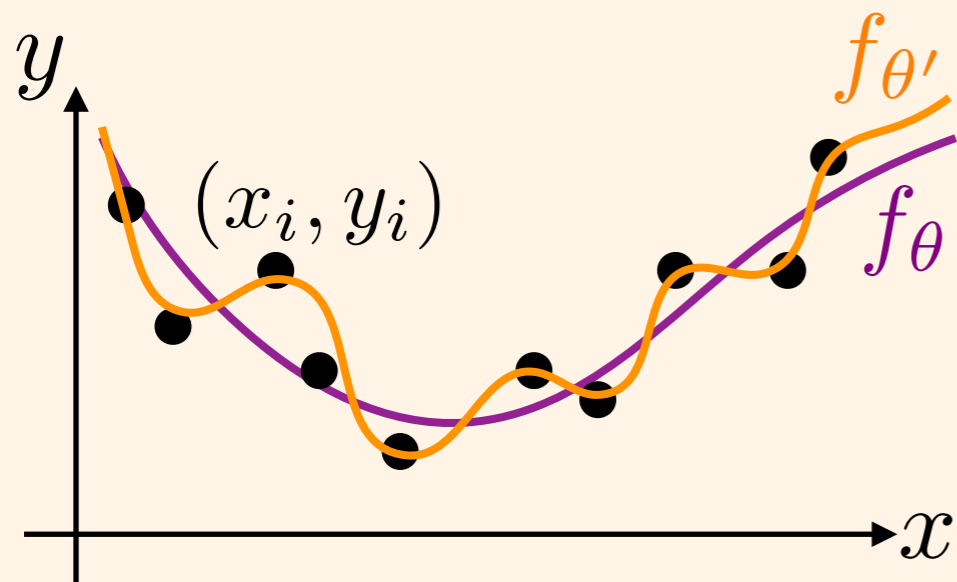
Empirical risk minimization:

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

Regression: $y_i \in \mathbb{R}$
 $y \approx f_{\theta}(x)$

Least square:

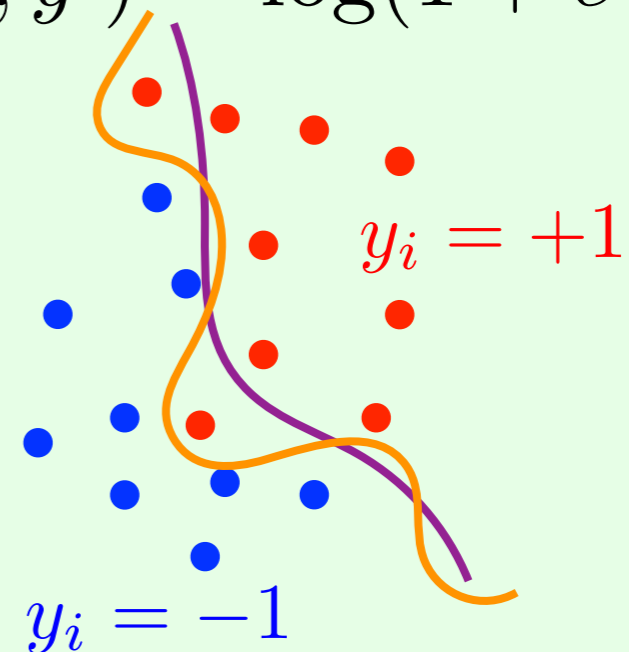
$$\ell(y, y') = (y - y')^2$$



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

Logistic:

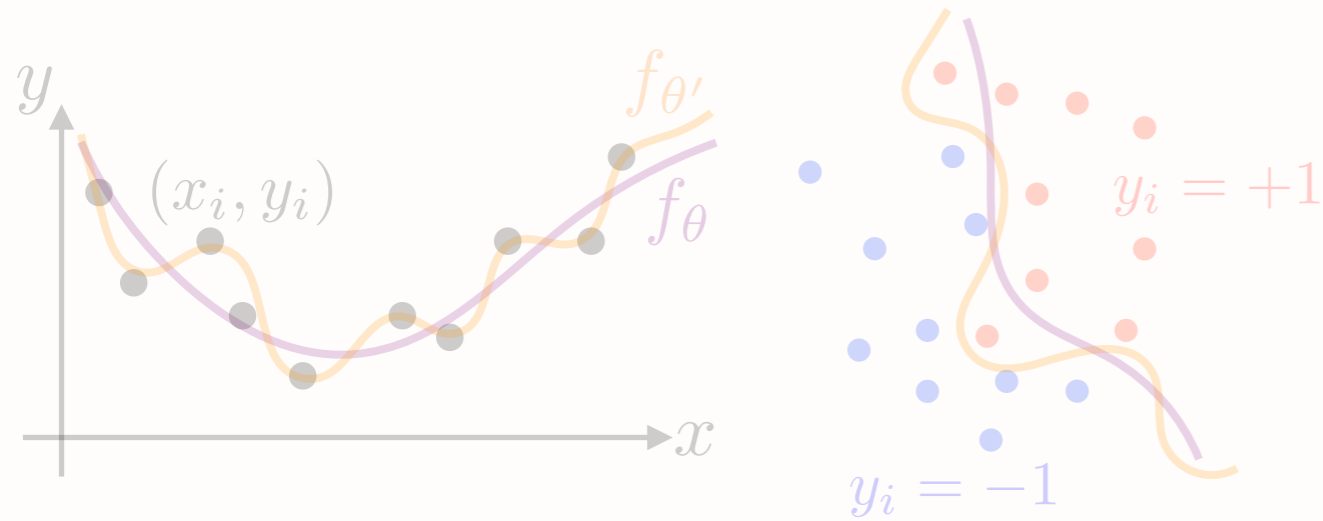
$$\ell(y, y') = \log(1 + e^{-yy'})$$



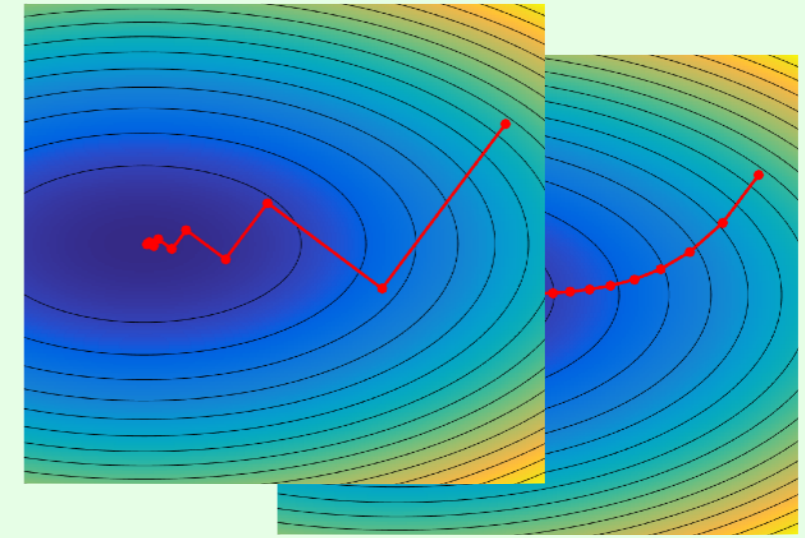
Overfitting, regularization, ...

Overview

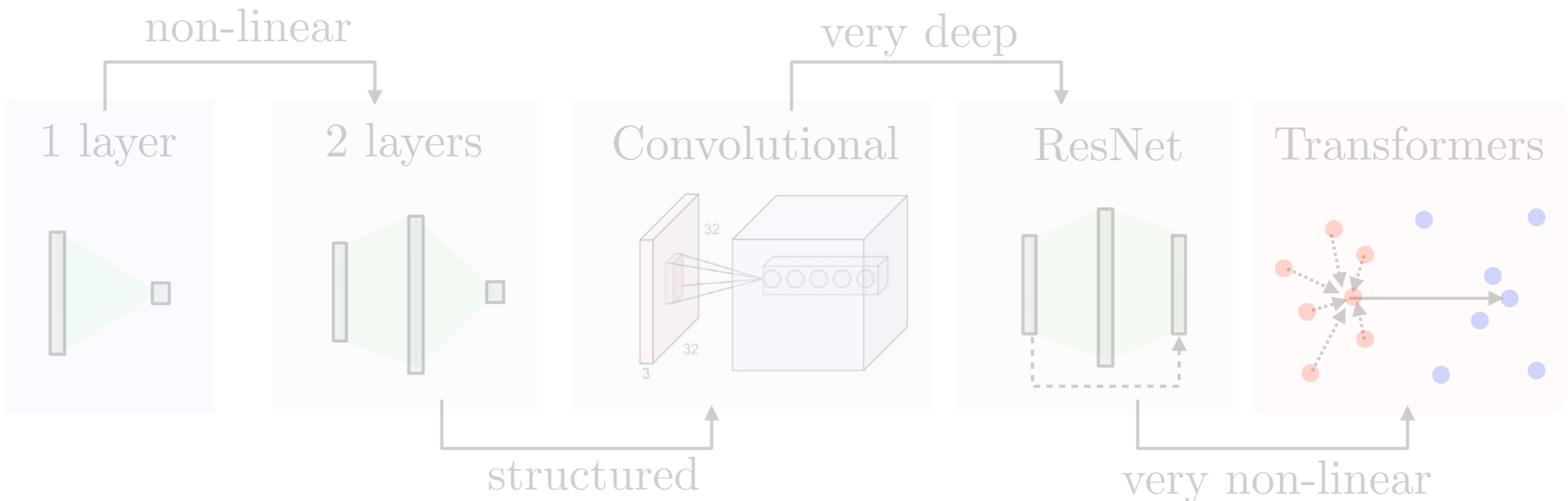
Machine Learning



Optimization



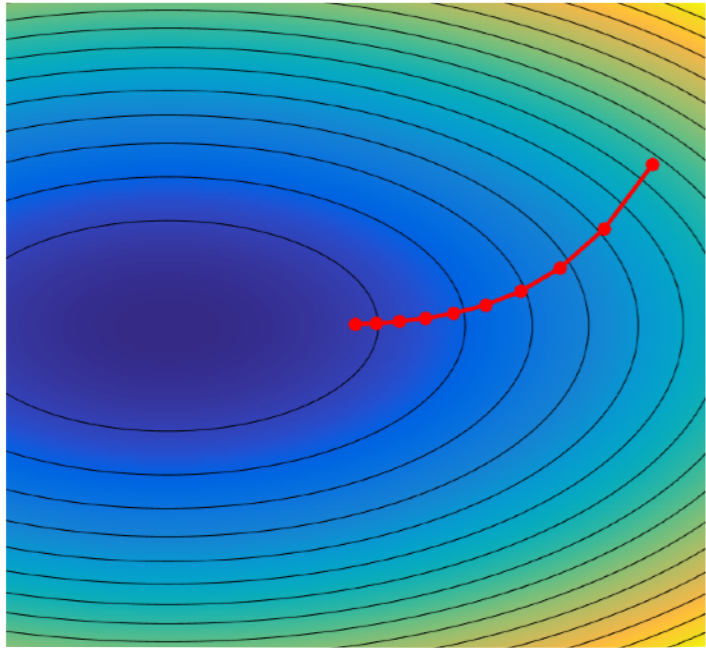
Networks Architectures



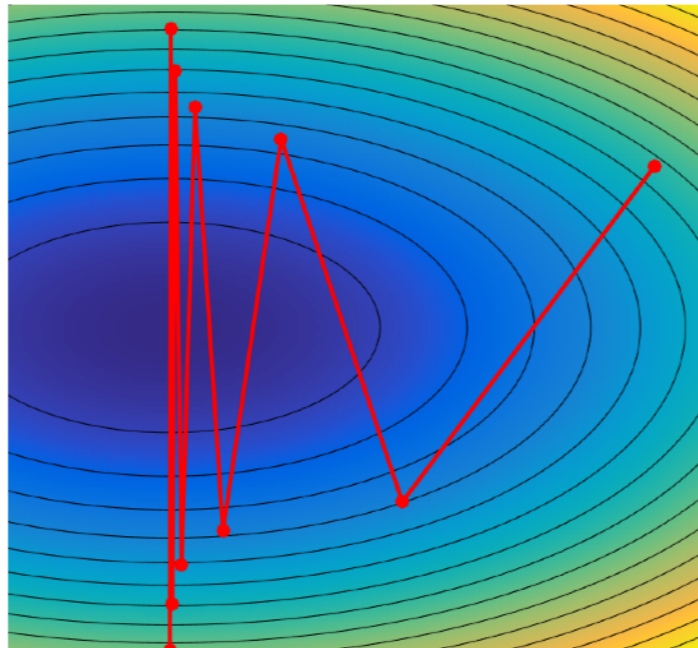
Gradient Descent

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

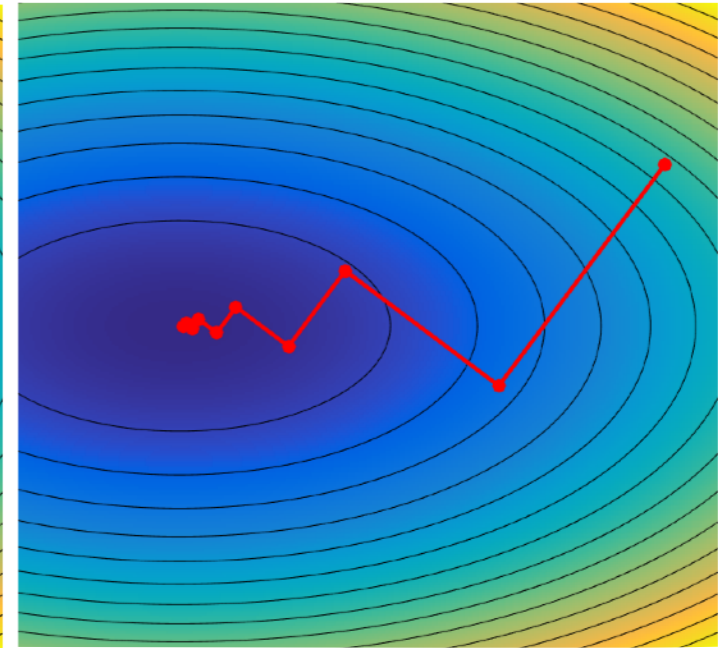
Gradient descent:
$$\theta_{\ell+1} = \theta_{\ell} - \tau_{\ell} \nabla \mathcal{E}(\theta_{\ell})$$



Small τ_{ℓ}



Large τ_{ℓ}

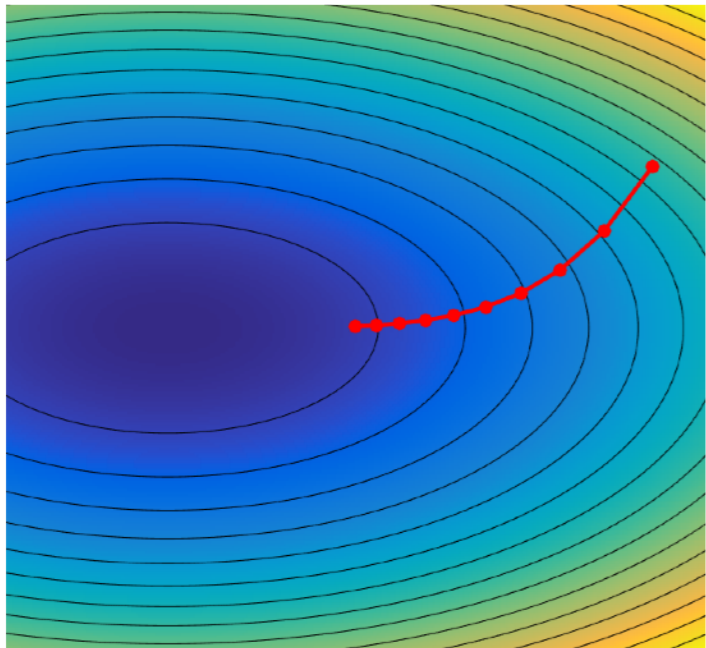


Optimal $\tau_{\ell} = \tau_{\ell}^*$

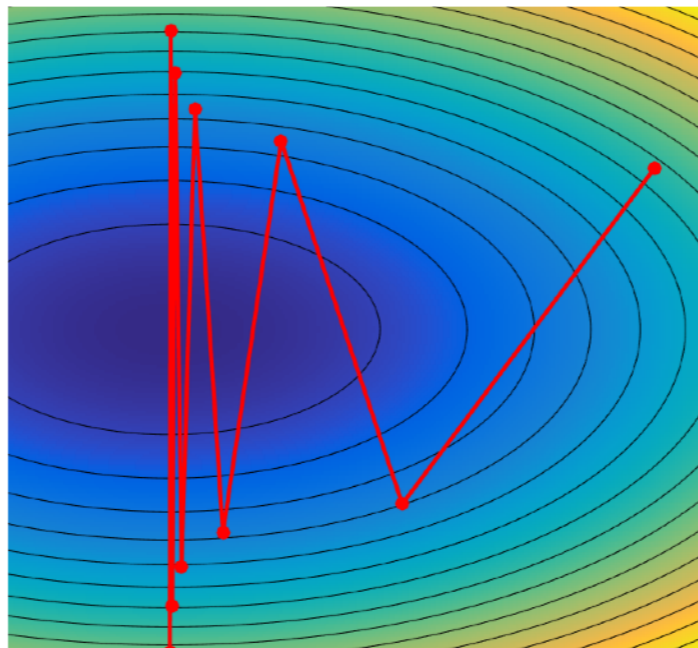
Gradient Descent

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

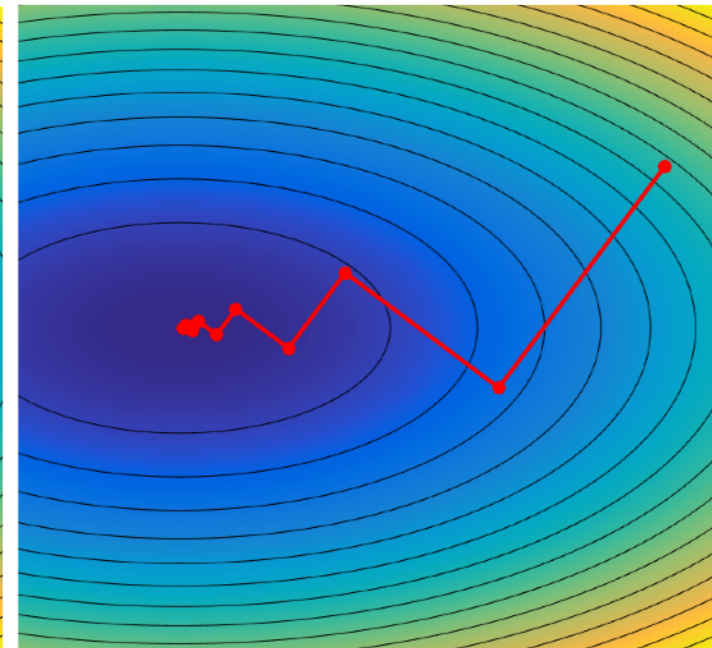
$$\text{Gradient descent:}$$
$$\theta_{\ell+1} = \theta_{\ell} - \tau_{\ell} \nabla \mathcal{E}(\theta_{\ell})$$



Small τ_{ℓ}



Large τ_{ℓ}



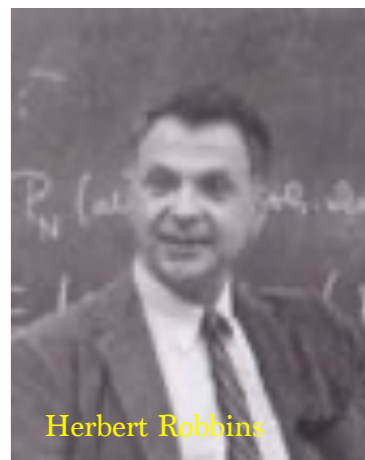
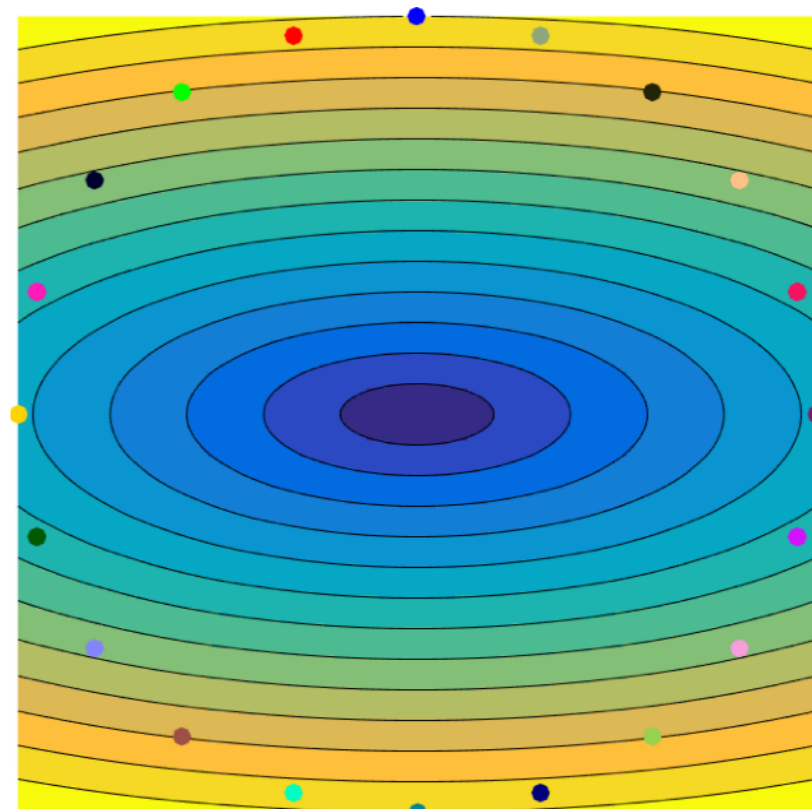
Optimal $\tau_{\ell} = \tau_{\ell}^*$

Stochastic gradient descent:

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

$i \leftarrow \text{rand}$

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



Herbert Robbins

Sutton Monro

The Complexity of Gradient Computation

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

```
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} ...) and their derivatives cost $O(1)$.

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

The Complexity of Gradient Computation

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

```
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} ...) and their derivatives cost $O(1)$.

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

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

The Complexity of Gradient Computation

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

```
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} ...) and their derivatives cost $O(1)$.

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

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

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

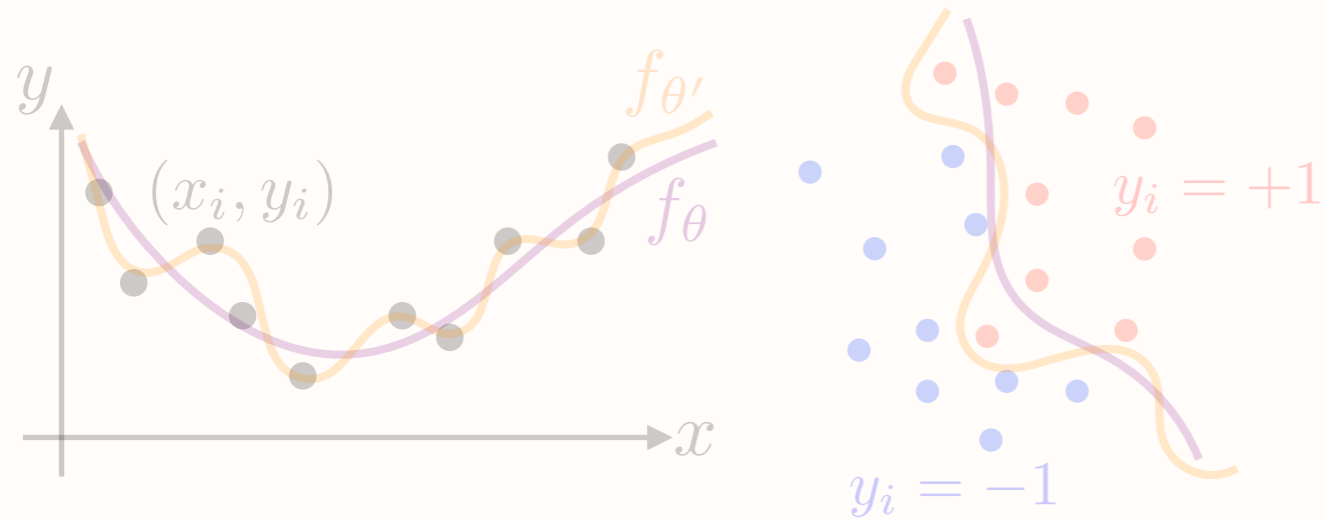
This algorithm is reverse mode automatic differentiation

```
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]
```

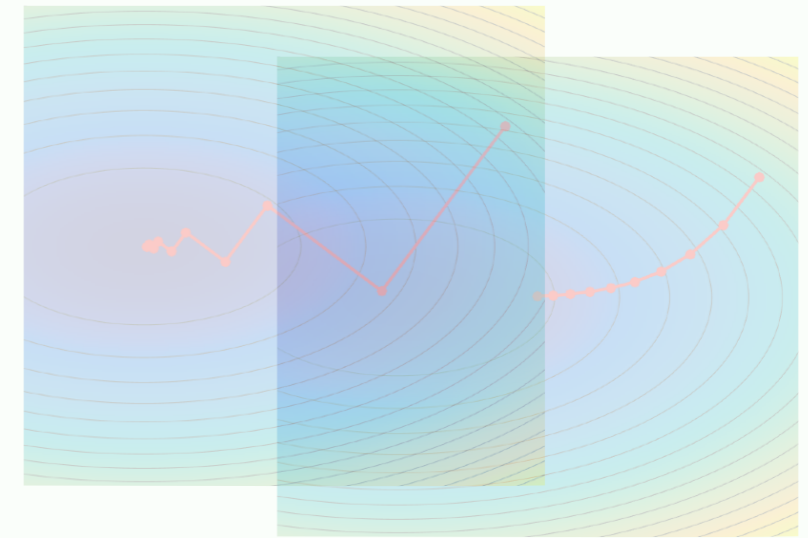


Overview

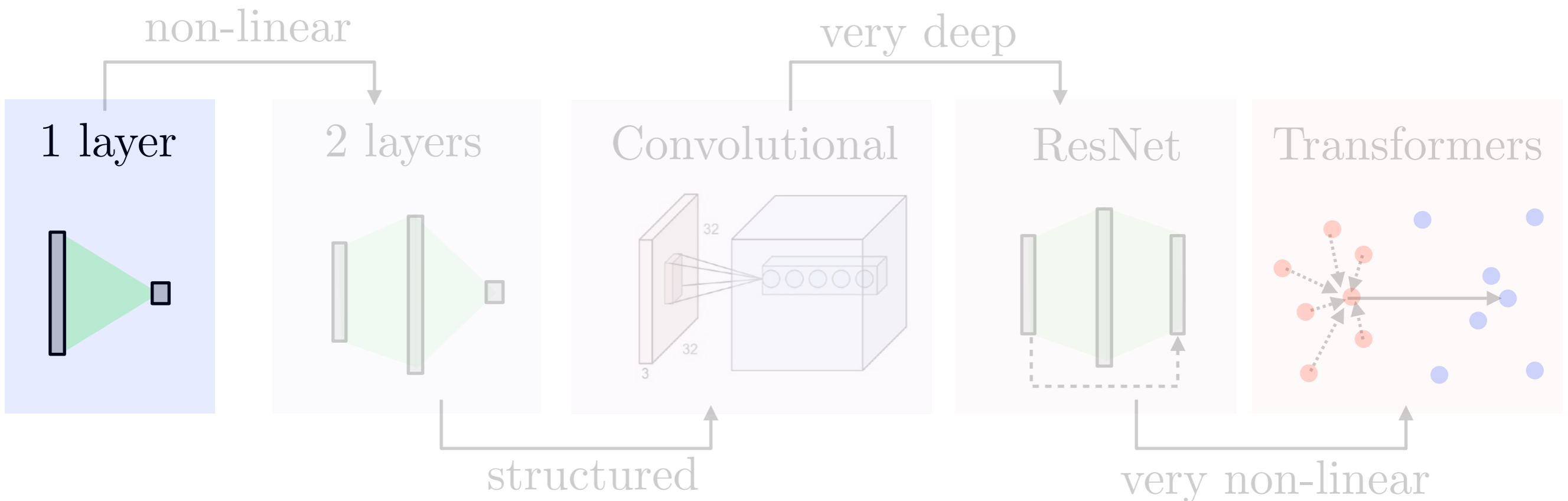
Machine Learning



Optimization

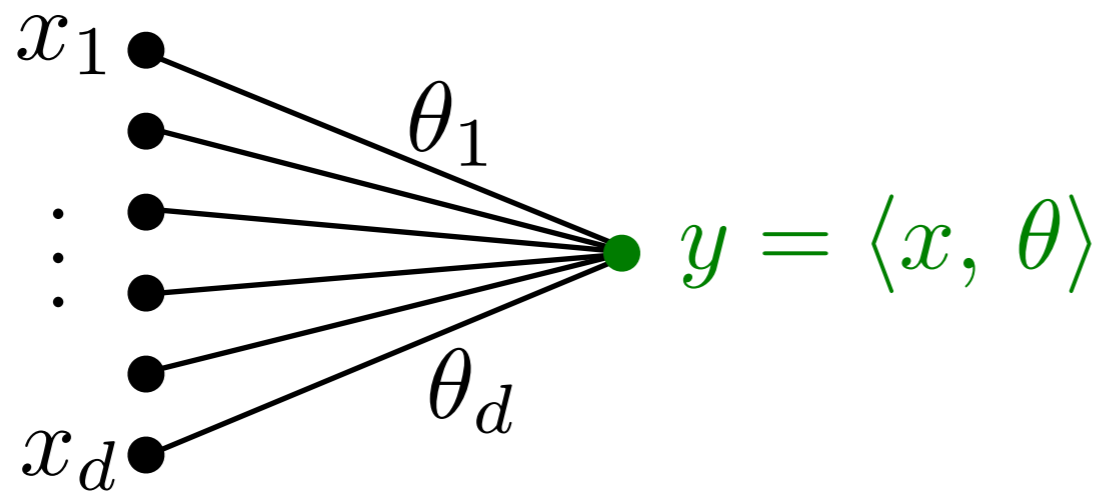


Networks Architectures

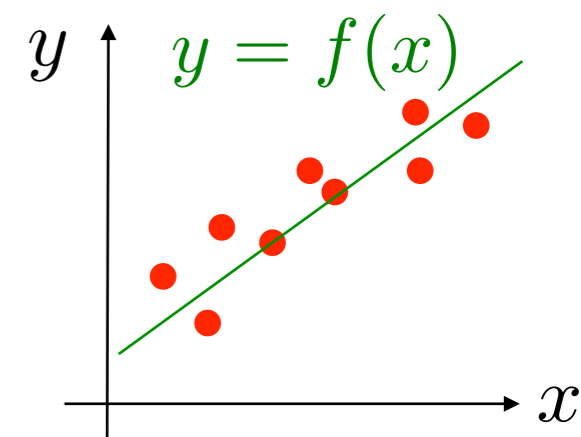


Linear model (1 layer)

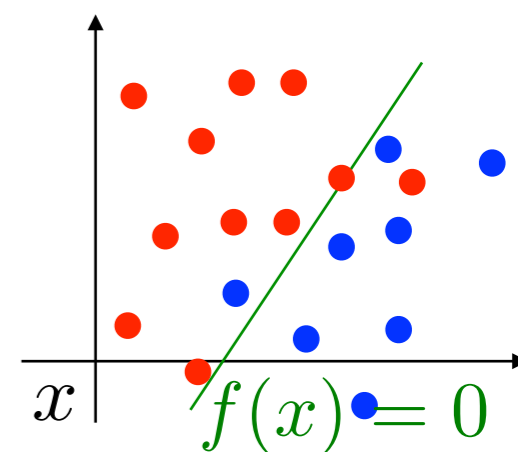
$$f_{\theta}(x) = \langle x, \theta \rangle = \sum_k x_k \theta_k$$



Regression

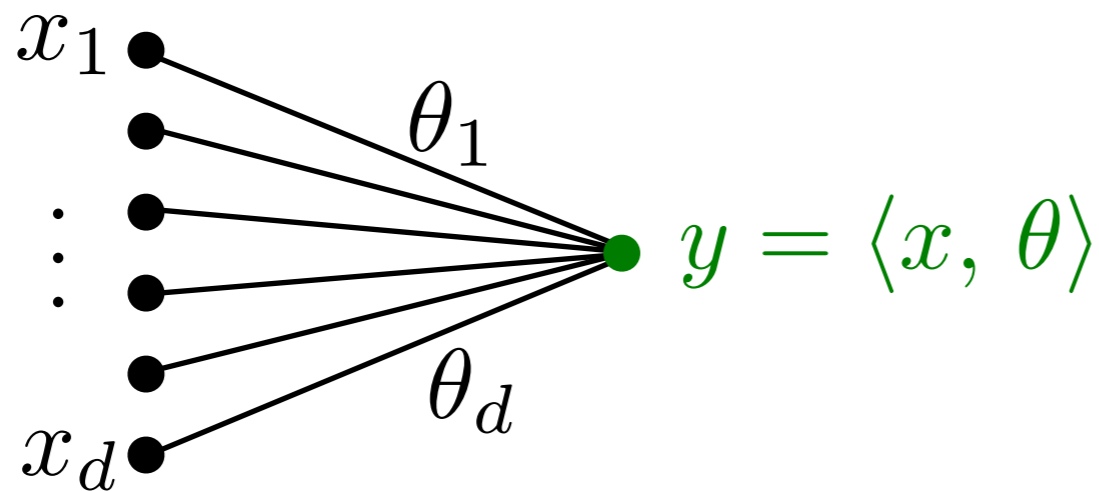


Classification:



Linear model (1 layer)

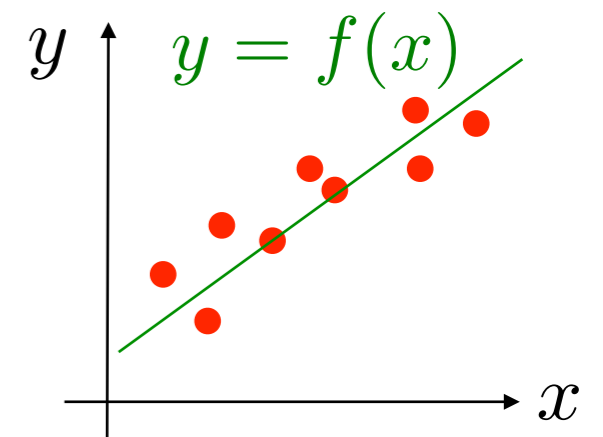
$$f_{\theta}(x) = \langle x, \theta \rangle = \sum_k x_k \theta_k$$



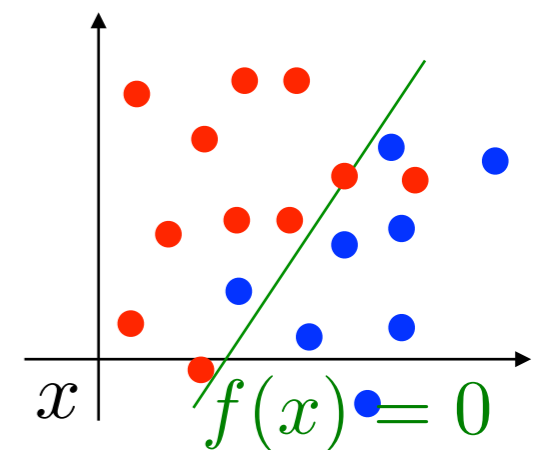
Convex optimization:

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

Regression

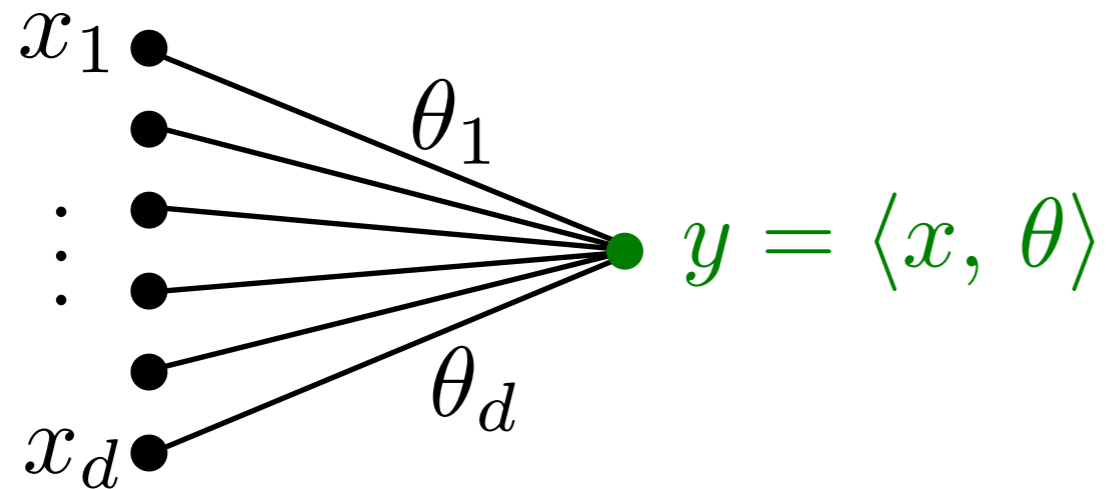


Classification:

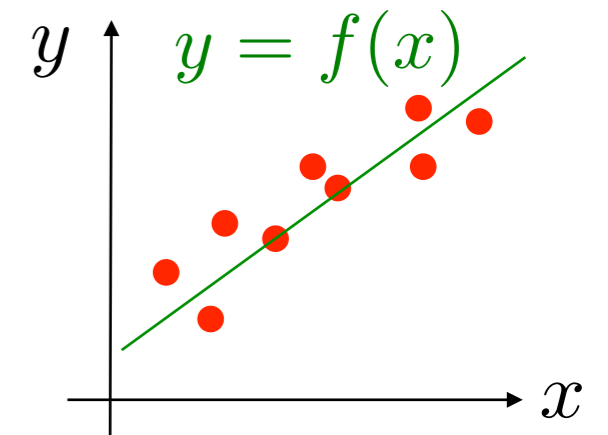


Linear model (1 layer)

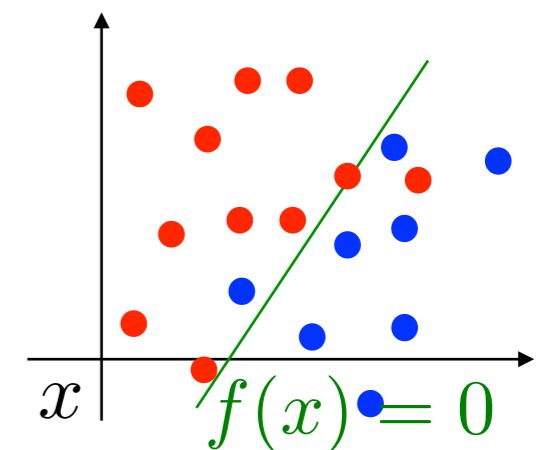
$$f_{\theta}(x) = \langle x, \theta \rangle = \sum_k x_k \theta_k$$



Regression



Classification:



Convex optimization:

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

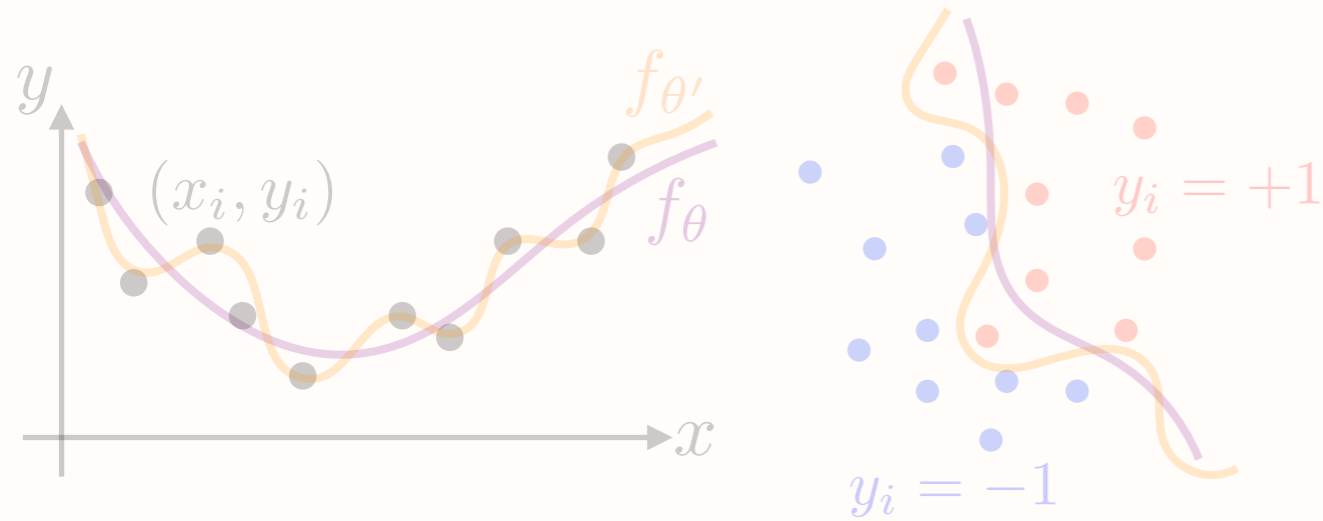
Kernel methods: replace x by $\varphi(x) \in \mathbb{R}^D$

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

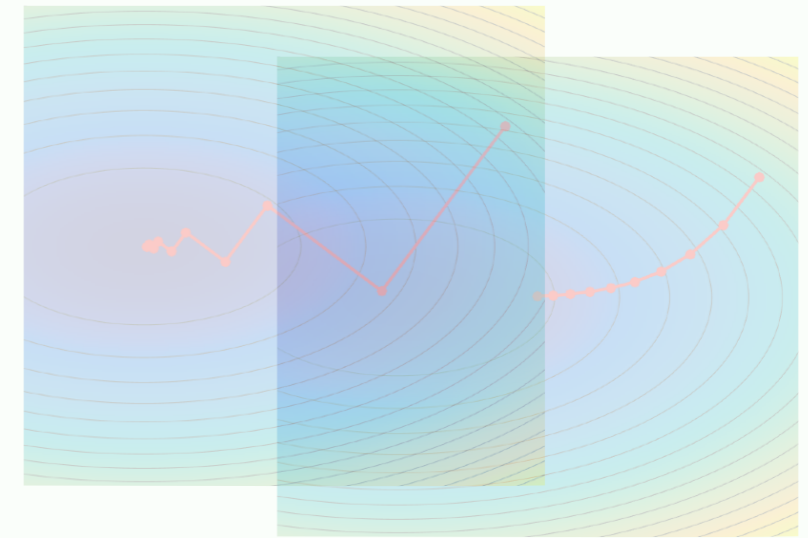
Deep learning methods: learn $\varphi(x)$!

Overview

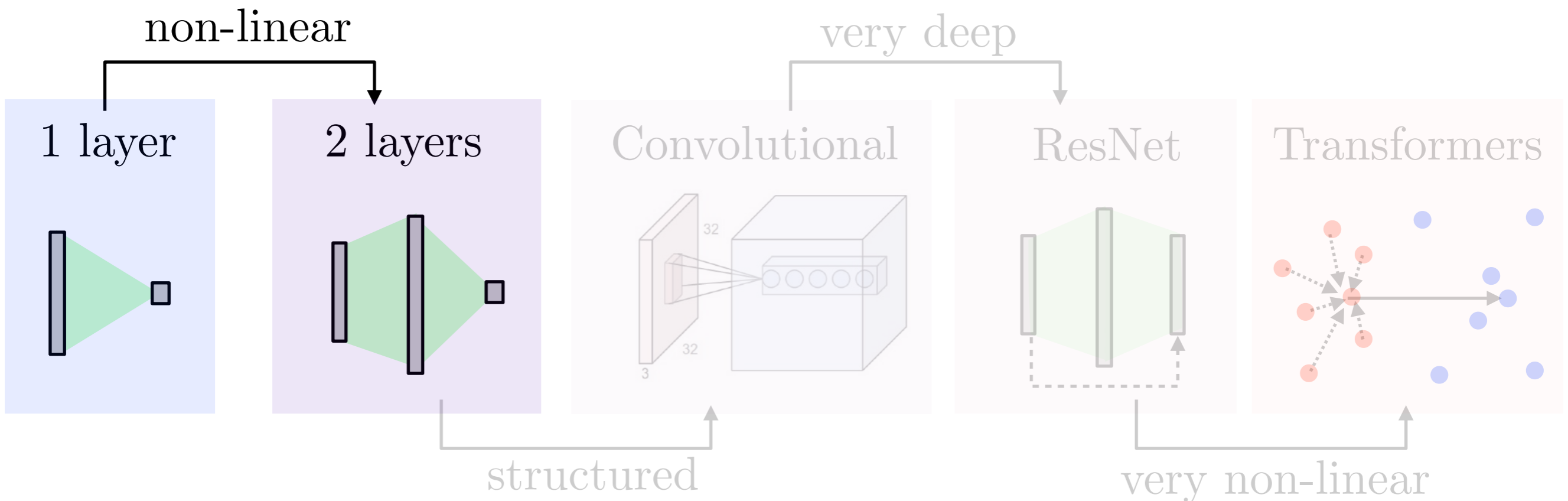
Machine Learning



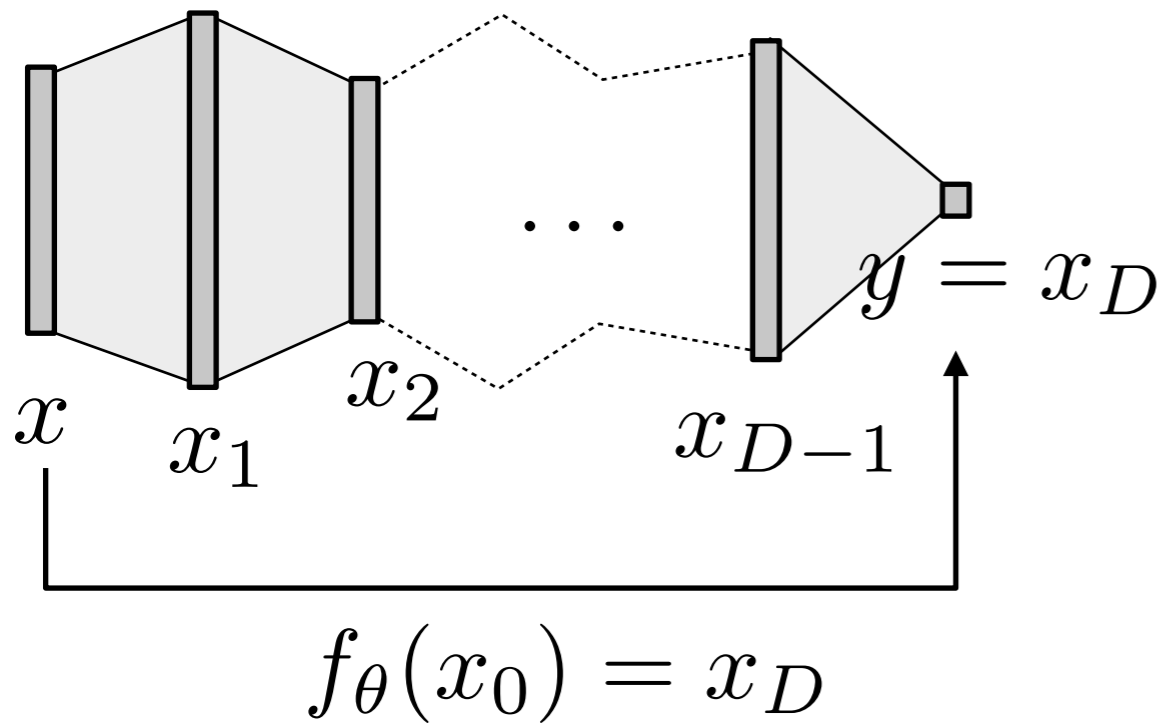
Optimization



Networks Architectures



Multi-layer Perceptron



$$x_0 \leftarrow x$$

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

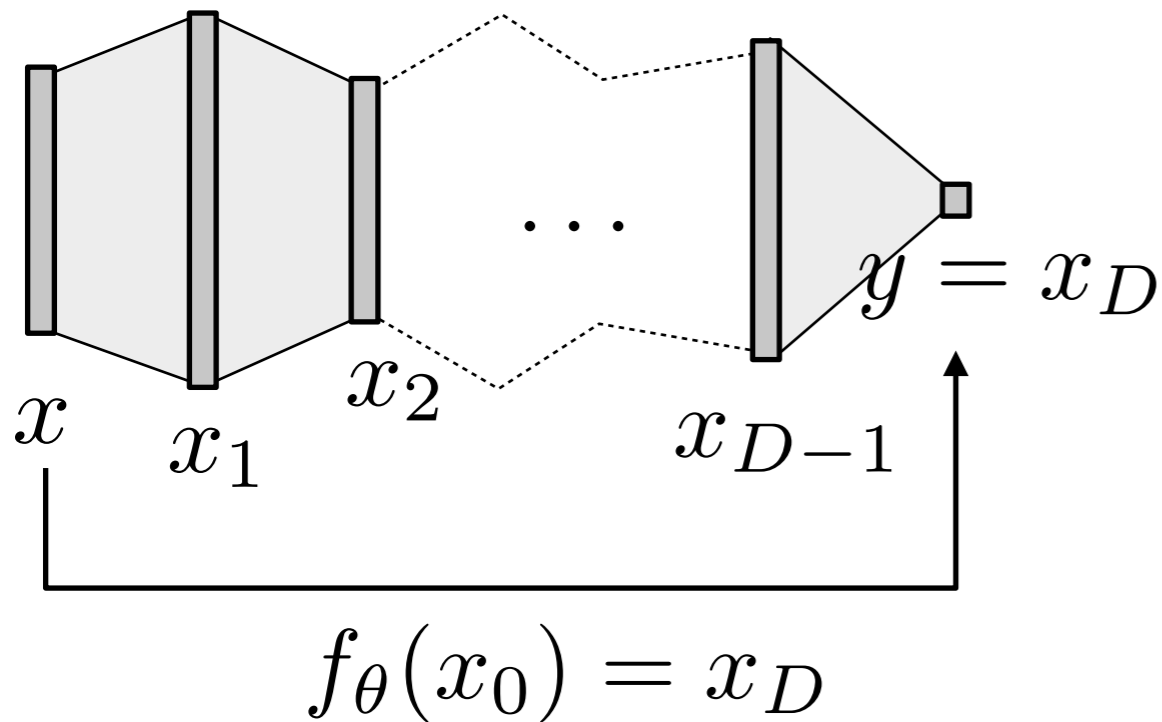
$$\theta = \{(W_k, b_k)\}_{k=0}^{D-1}$$

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

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



Multi-layer Perceptron



$$x_0 \leftarrow x$$

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

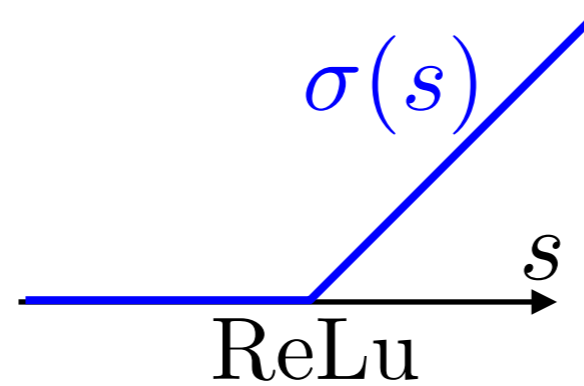
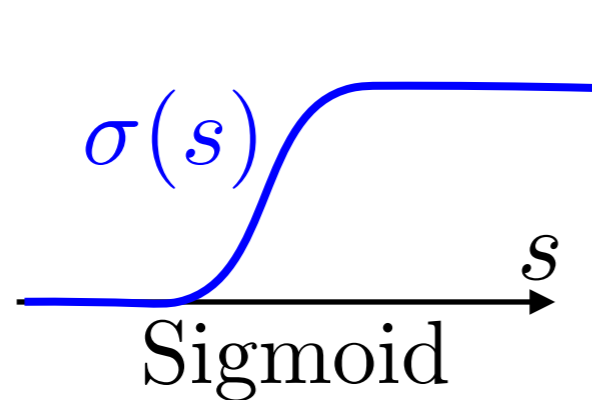
$$\theta = \{(W_k, b_k)\}_{k=0}^{D-1}$$

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

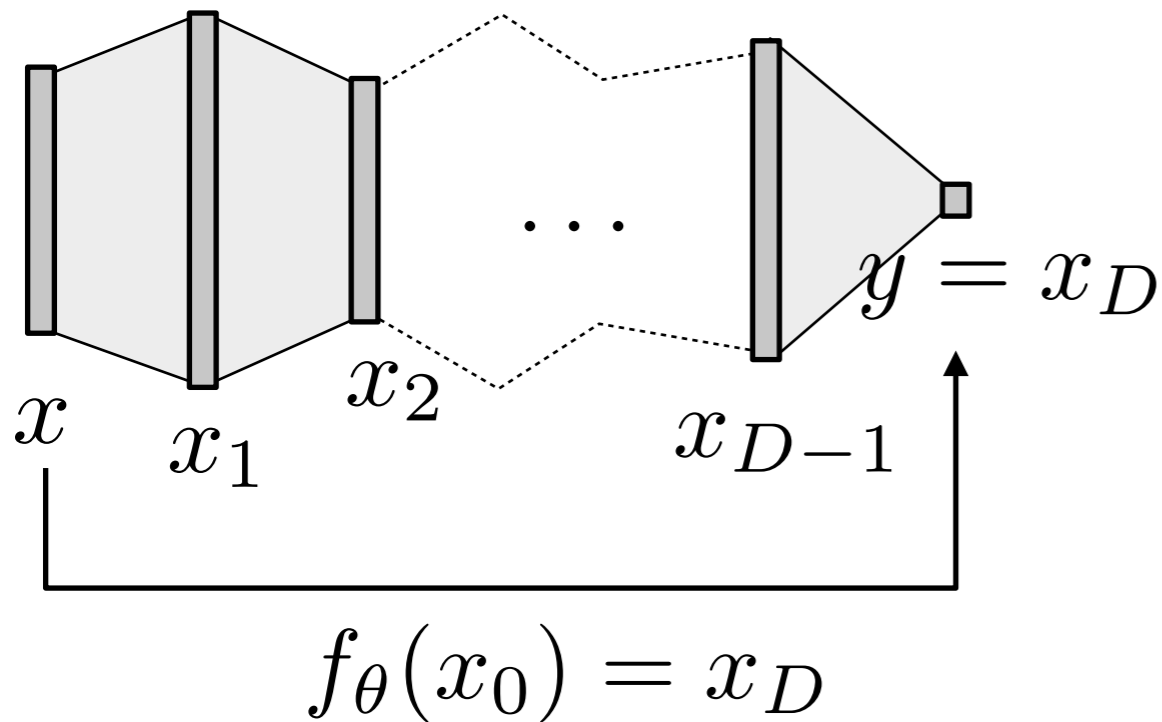
$$b_k \in \mathbb{R}^{d_{k+1}}$$



Non-linearity: σ must be non-polynomial to increase expressivity.



Multi-layer Perceptron

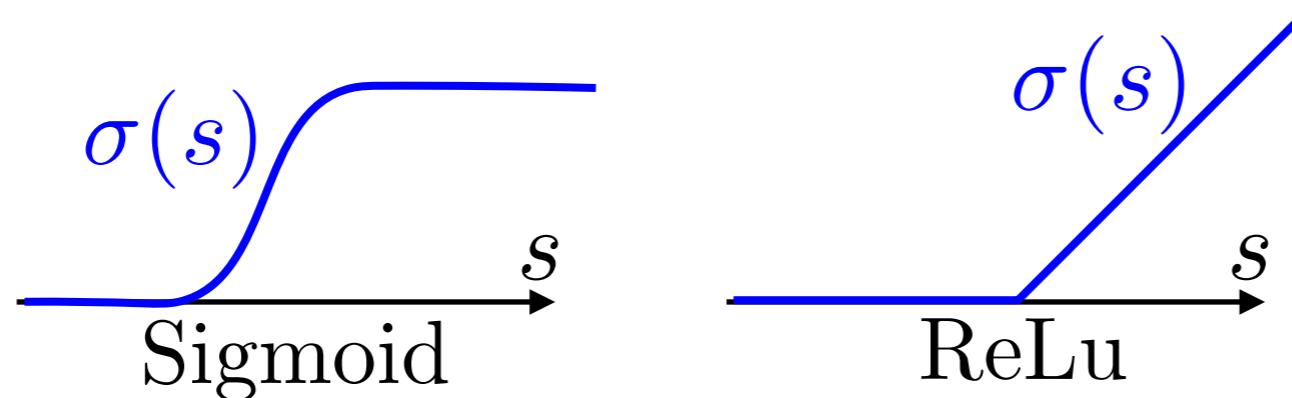


$$x_0 \leftarrow x$$
$$x_{k+1} \triangleq \sigma(W_k x_k + b_k)$$

$$\theta = \{(W_k, b_k)\}_{k=0}^{D-1}$$
$$W_k \in \mathbb{R}^{d_{k+1} \times d_k}$$
$$b_k \in \mathbb{R}^{d_{k+1}}$$

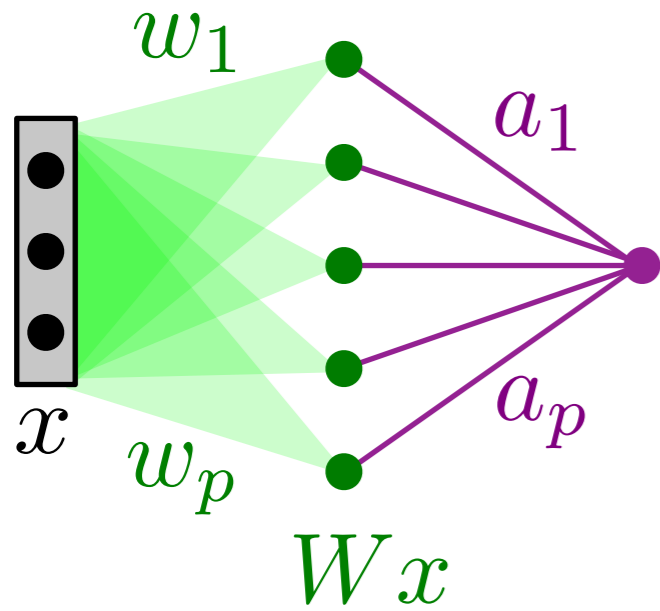


Non-linearity: σ must be non-polynomial to increase expressivity.



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

Two Layers Perceptron



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

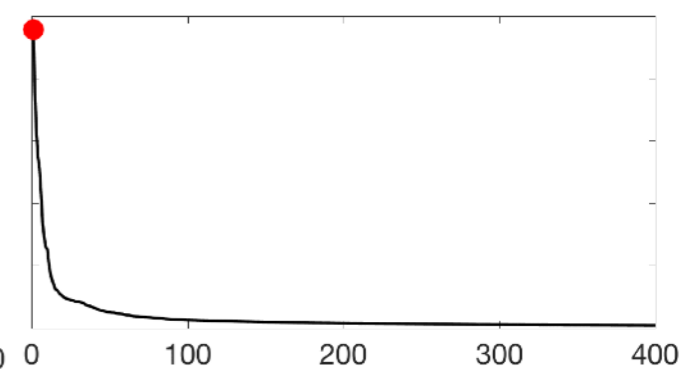
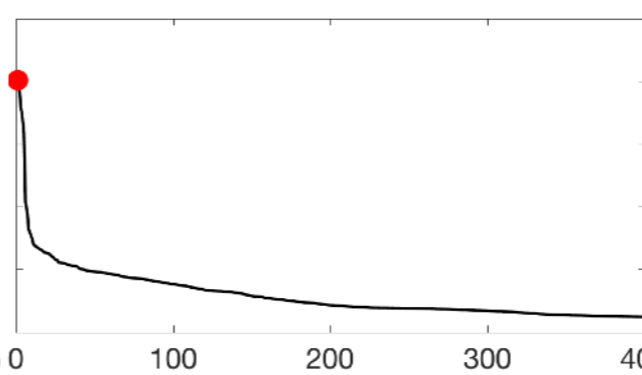
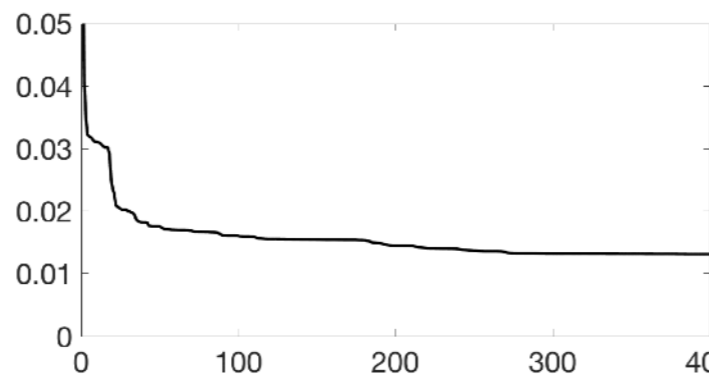
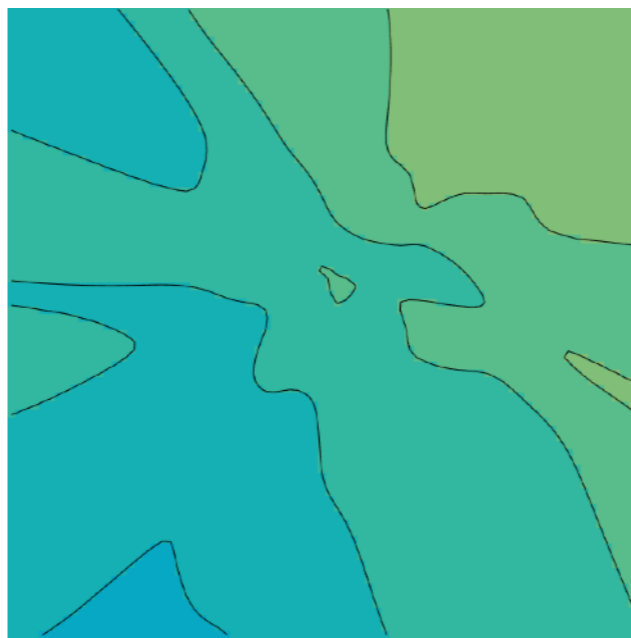
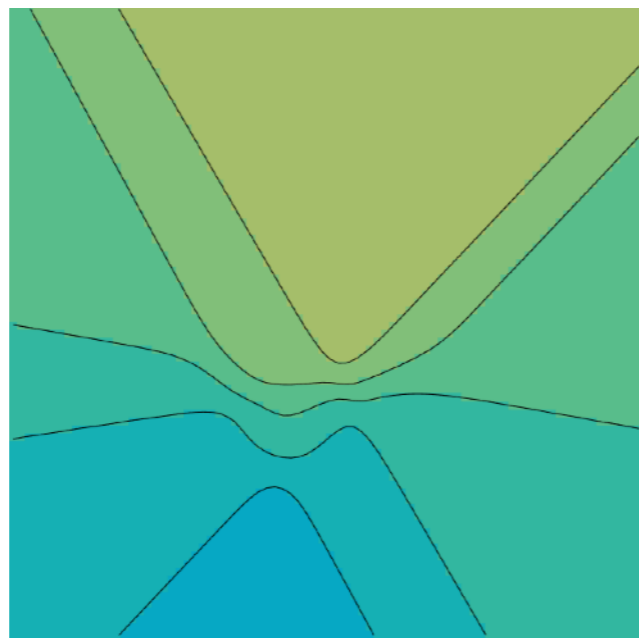
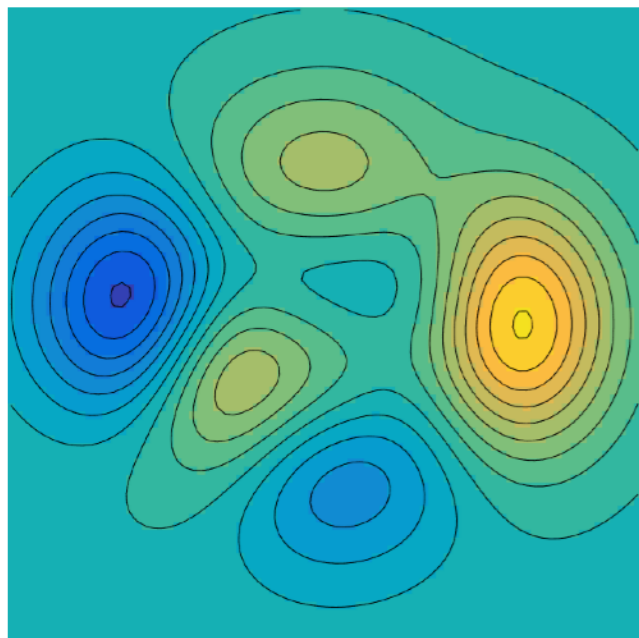
→ sum of “ridge” functions $\sigma(\langle x, w \rangle + b)$

Input $y = f(x)$

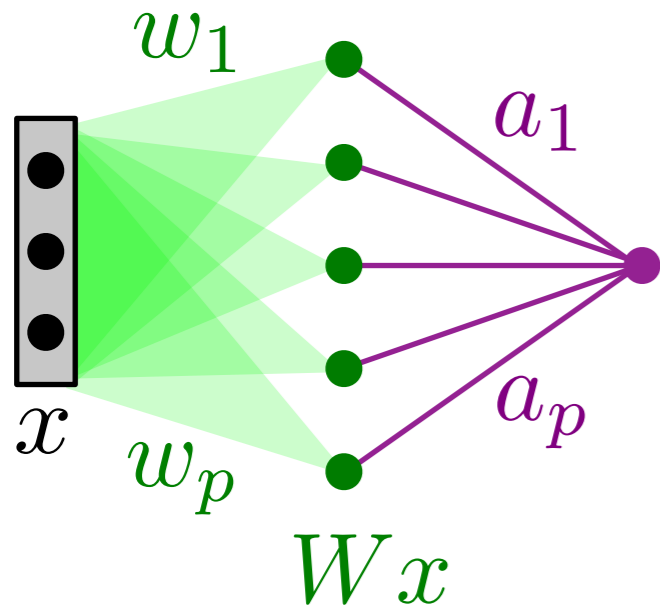
$p = 6$ neurons

$p = 30$ neurons

$p = 100$ neurons



Two Layers Perceptron



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

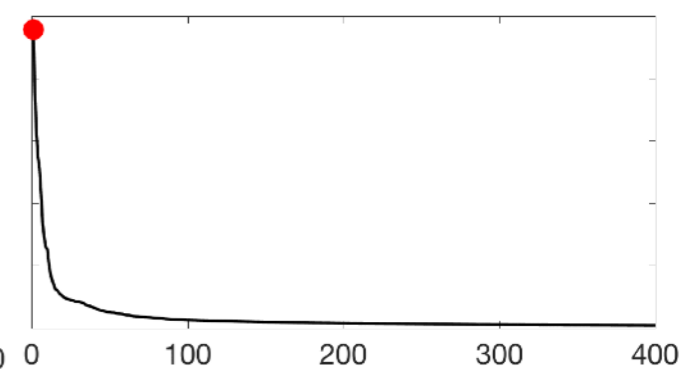
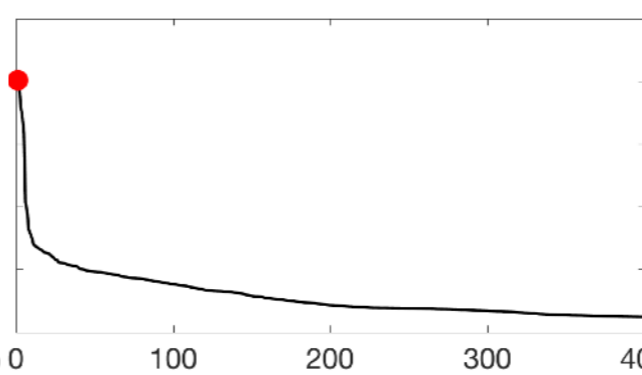
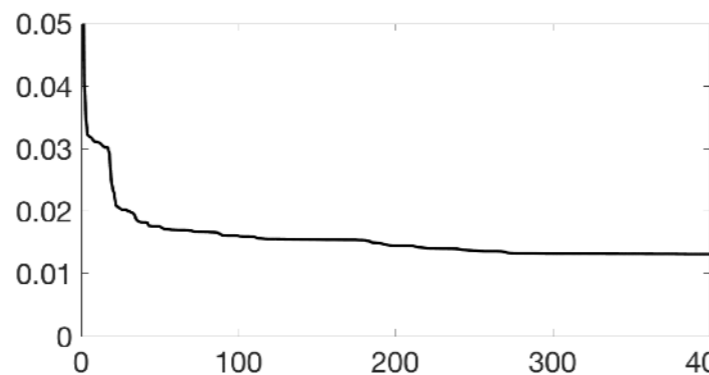
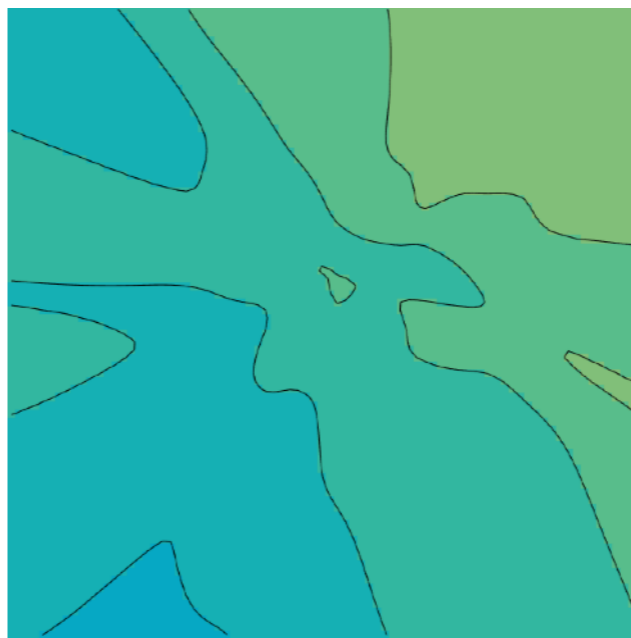
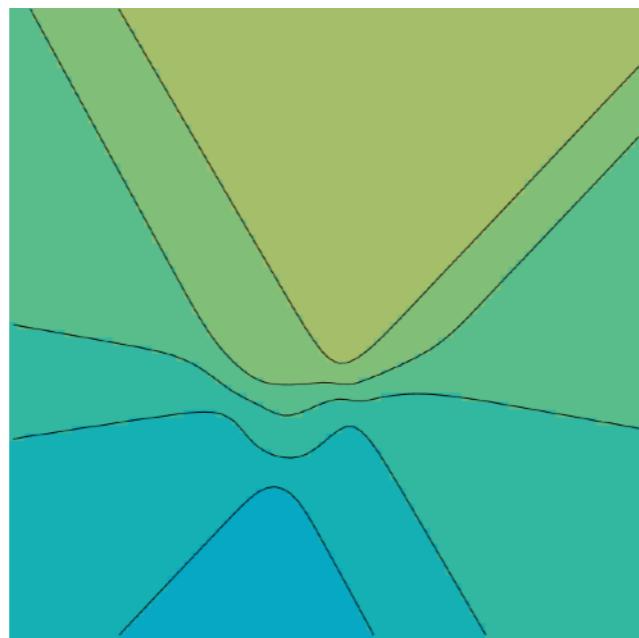
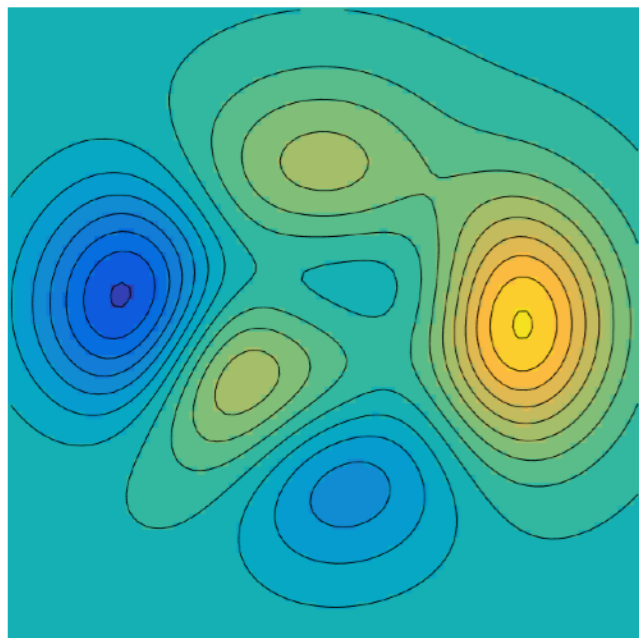
→ sum of “ridge” functions $\sigma(\langle x, w \rangle + b)$

Input $y = f(x)$

$p = 6$ neurons

$p = 30$ neurons

$p = 100$ neurons



Universality of Perceptrons

Theorem: If f is continuous on a compact Ω , for all $\varepsilon > 0$ for p large enough, there exists θ such that

$$\forall x \in \Omega, |f_{\theta}(x) - f(x)| \leq \varepsilon$$

→ non quantitative ... no free lunch.



Universality of Perceptrons

Theorem: If f is continuous on a compact Ω , for all $\varepsilon > 0$ for p large enough, there exists θ such that

$$\forall x \in \Omega, |f_\theta(x) - f(x)| \leq \varepsilon$$

→ non quantitative ... no free lunch.

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

Theorem: for p large, there exists θ such that

$$\|f - f_\theta\|_{L^2(\Omega)} \leq \frac{2\text{diam}(\Omega)\|f\|_B}{\sqrt{p}}$$



George Cybenko



Andrew Barron

Universality of Perceptrons

Theorem: If f is continuous on a compact Ω , for all $\varepsilon > 0$ for p large enough, there exists θ such that

$$\forall x \in \Omega, |f_\theta(x) - f(x)| \leq \varepsilon$$

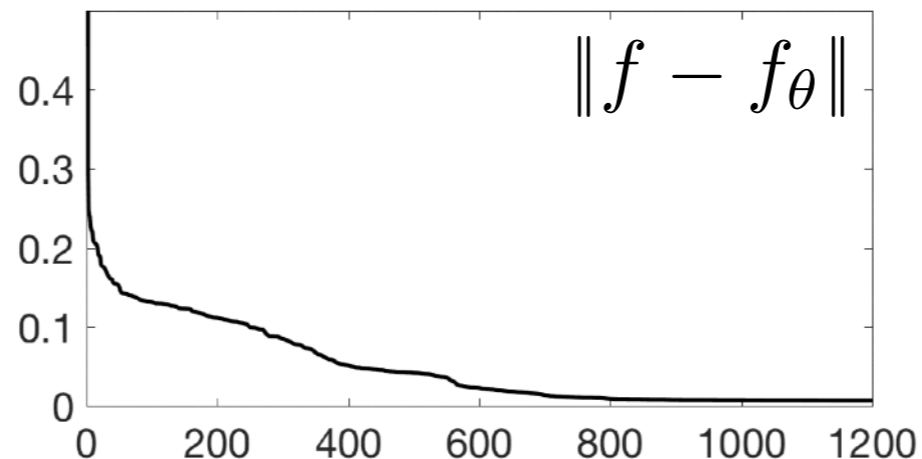
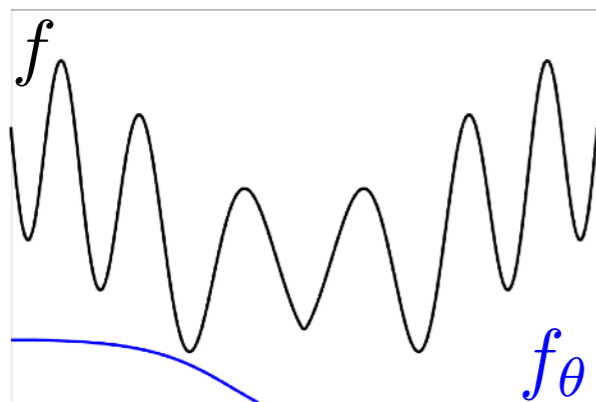
→ non quantitative ... no free lunch.

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

Theorem: for p large, there exists θ such that

$$\|f - f_\theta\|_{L^2(\Omega)} \leq \frac{2\text{diam}(\Omega)\|f\|_B}{\sqrt{p}}$$

→ non-constructive.



→ for p “large enough” gradient descent works
[Chizat-Bach 2018]



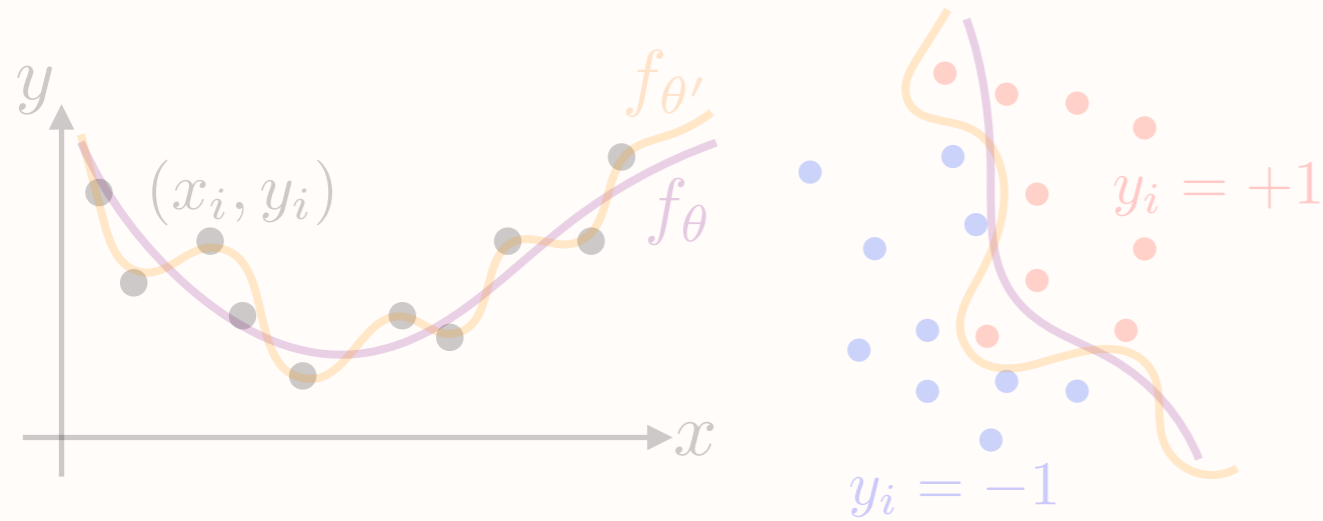
George Cybenko



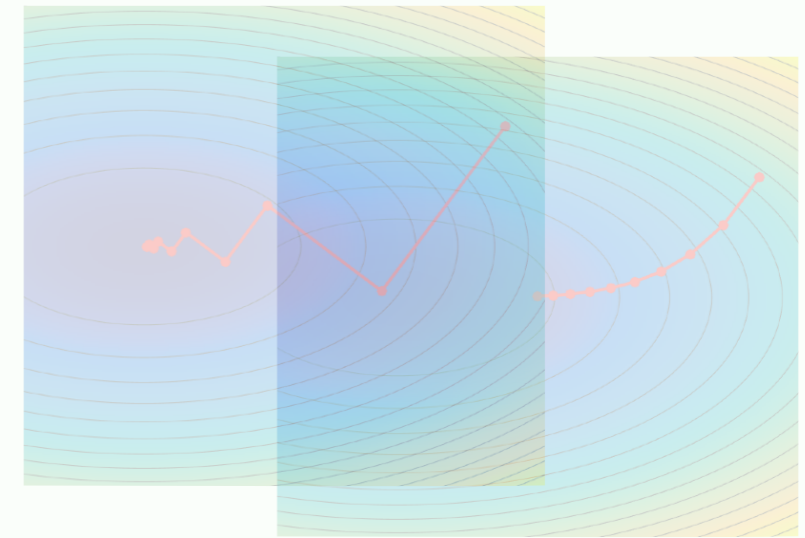
Andrew Barron

Overview

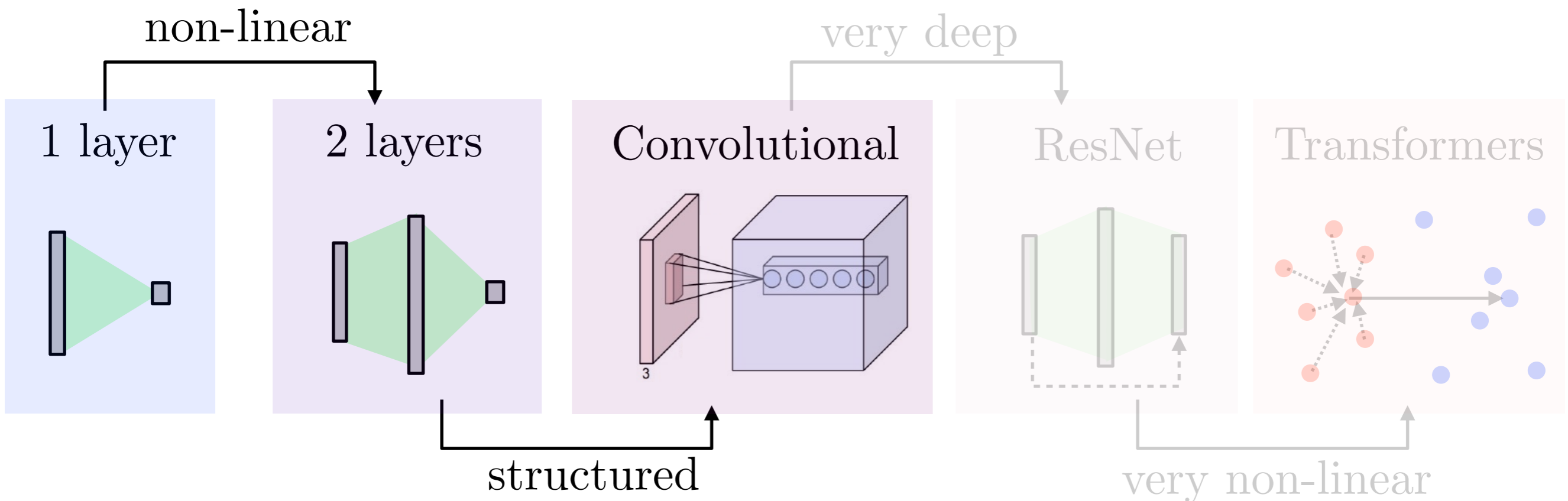
Machine Learning



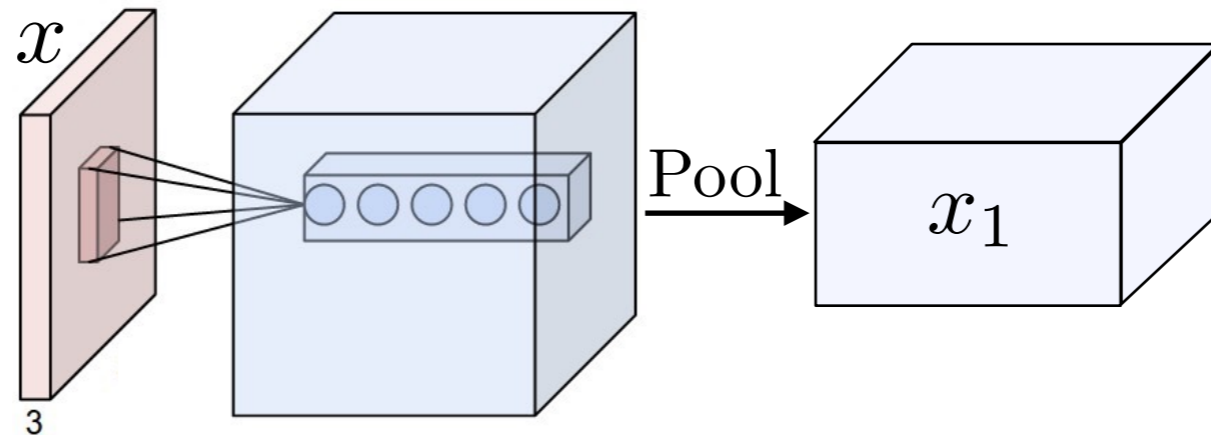
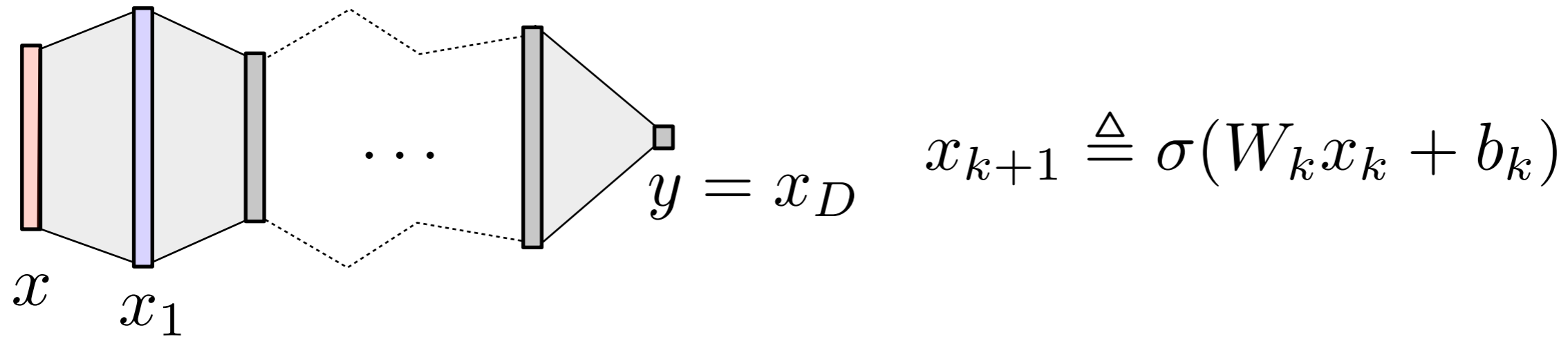
Optimization



Networks Architectures



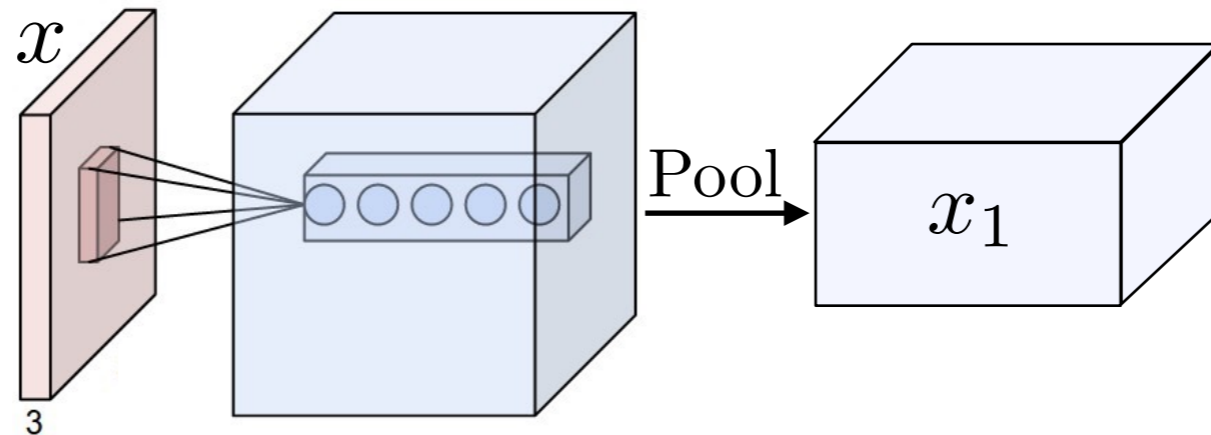
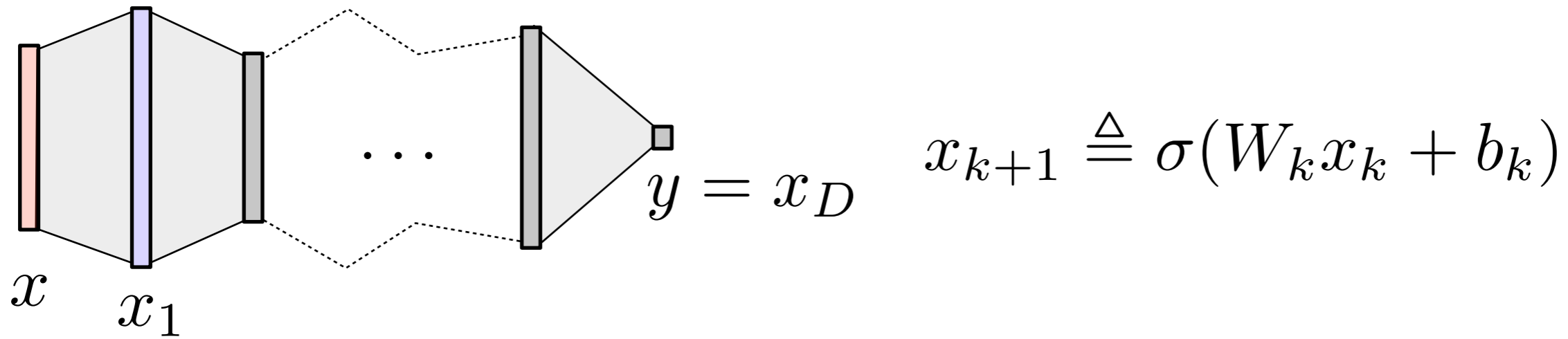
Convolutional CNN



→ Leverage translation invariance of images.

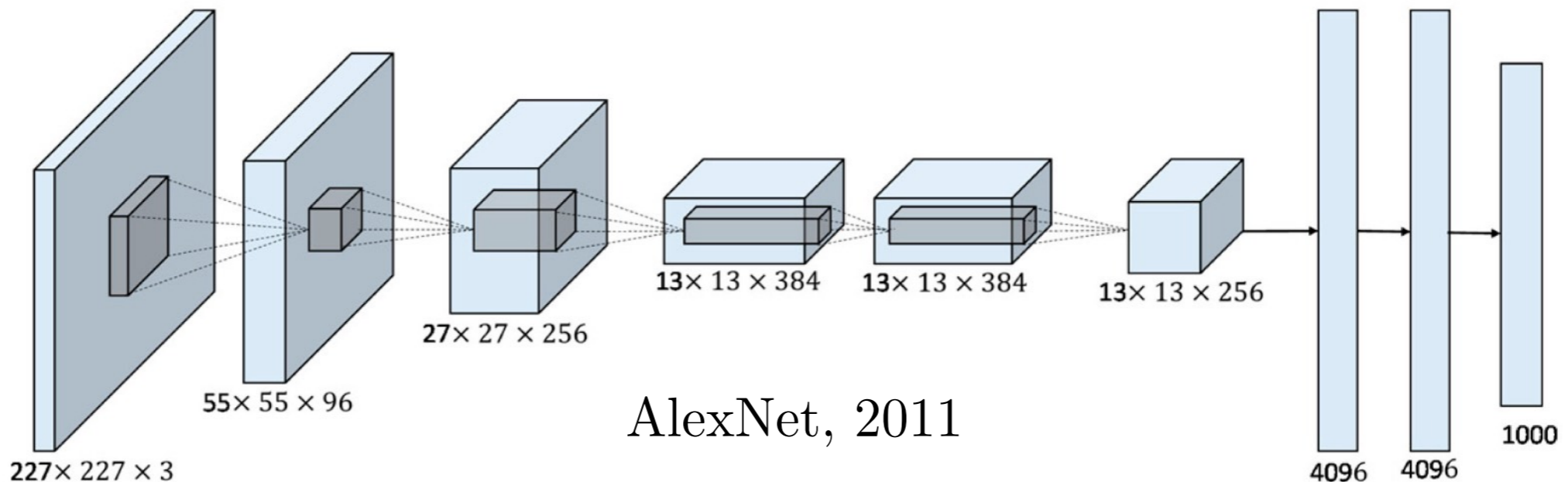
→ Sub-sampling: breaks invariance but increase receptive fields.

Convolutional CNN

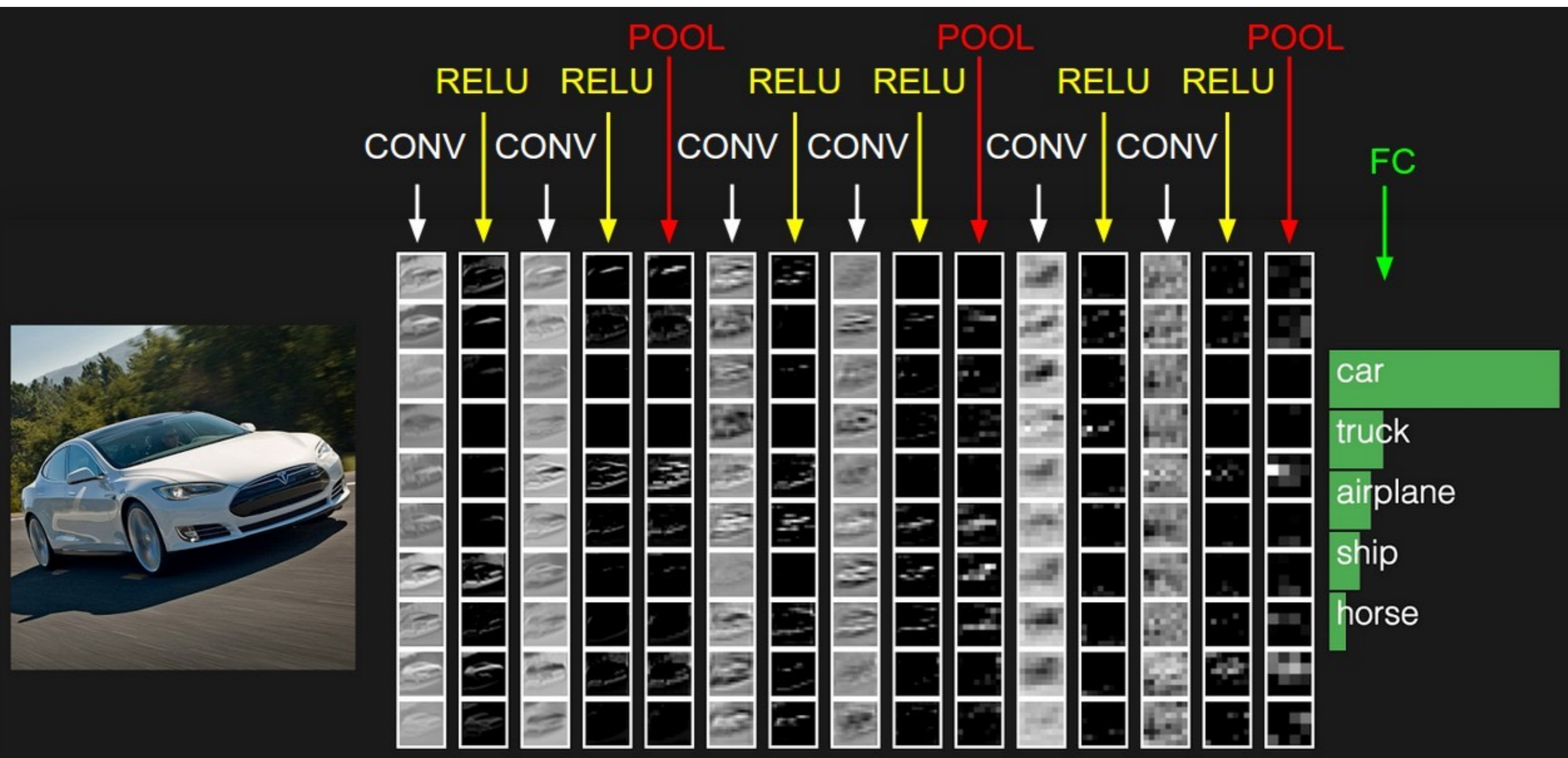


→ Leverage translation invariance of images.

→ Sub-sampling: breaks invariance but increase receptive fields.

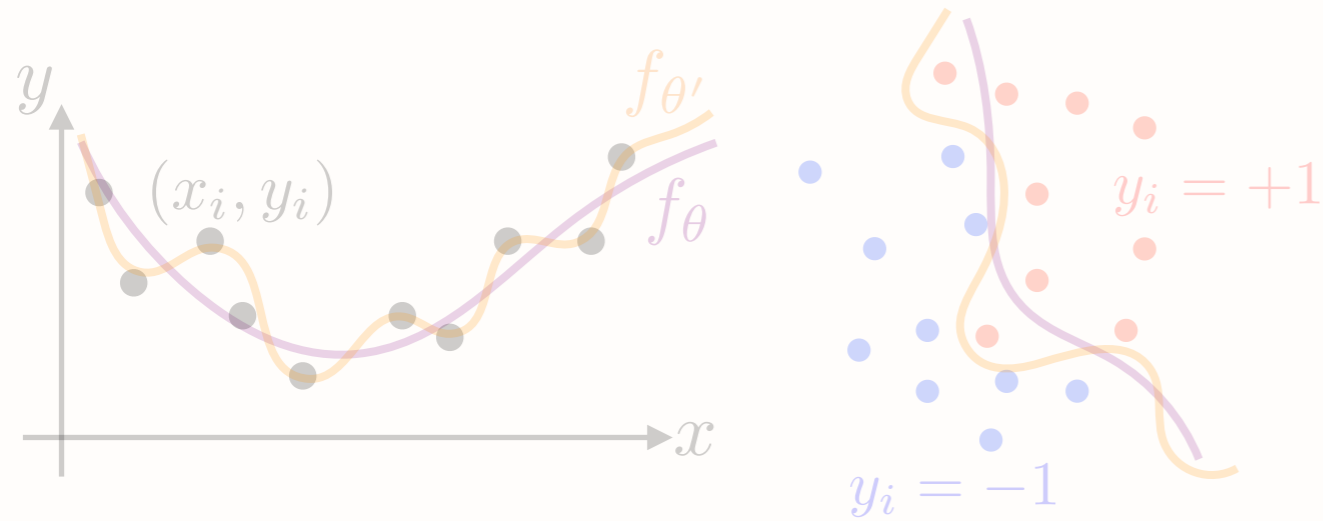


Example of Activations

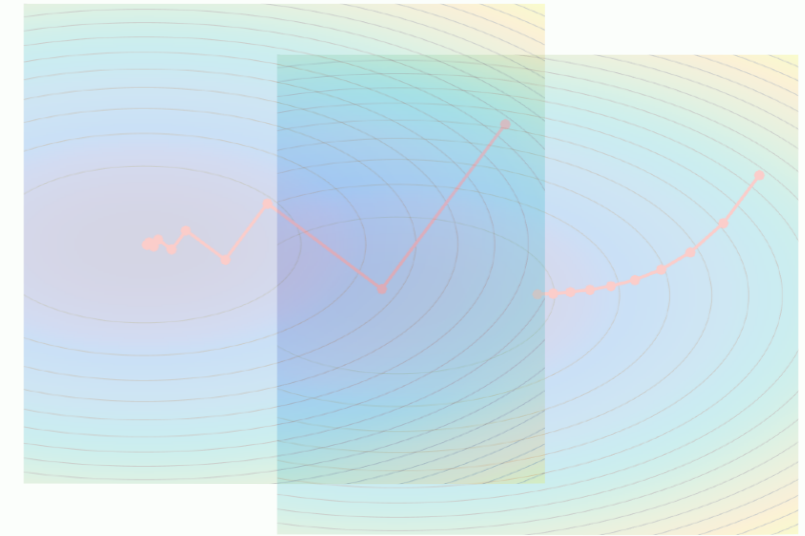


Overview

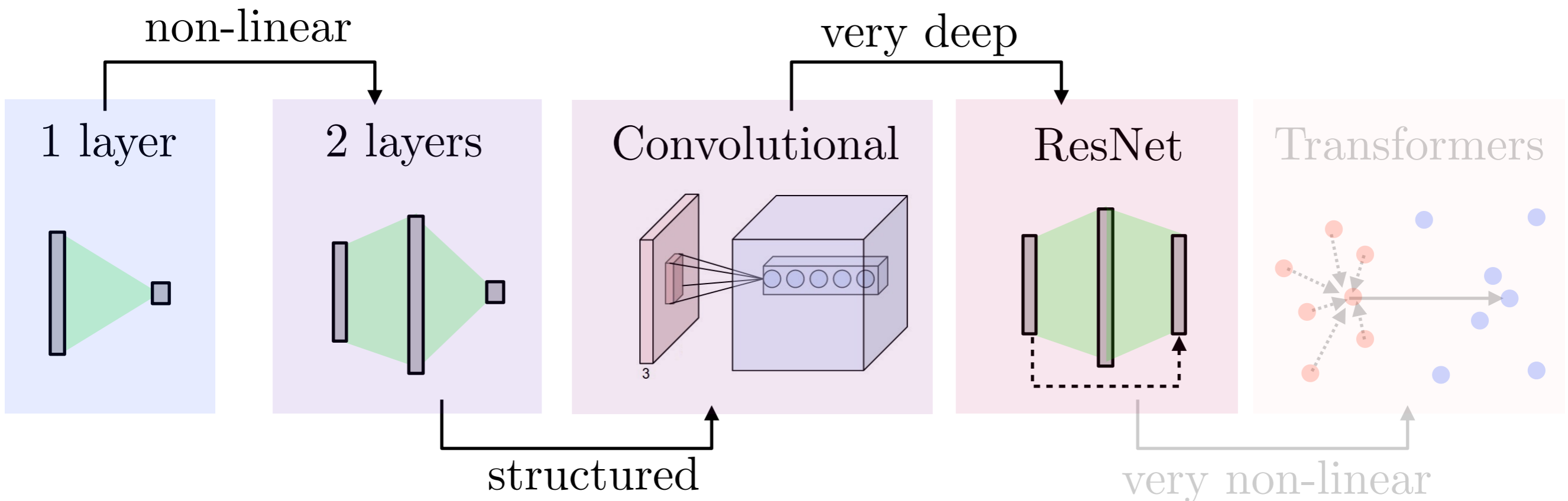
Machine Learning



Optimization

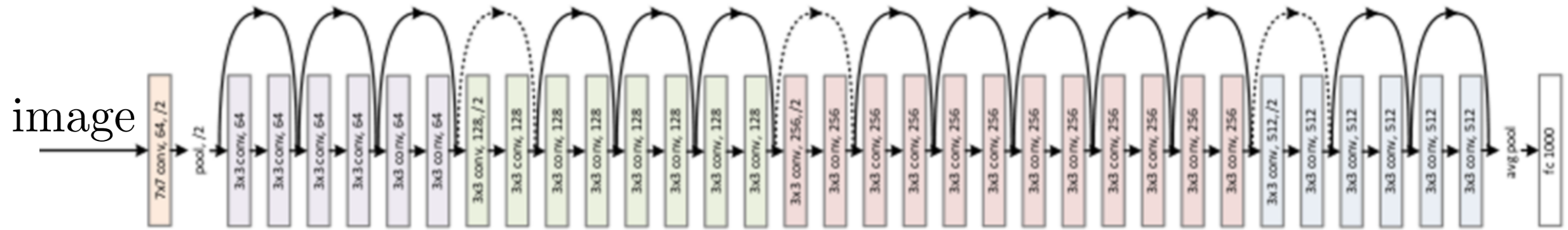


Networks Architectures



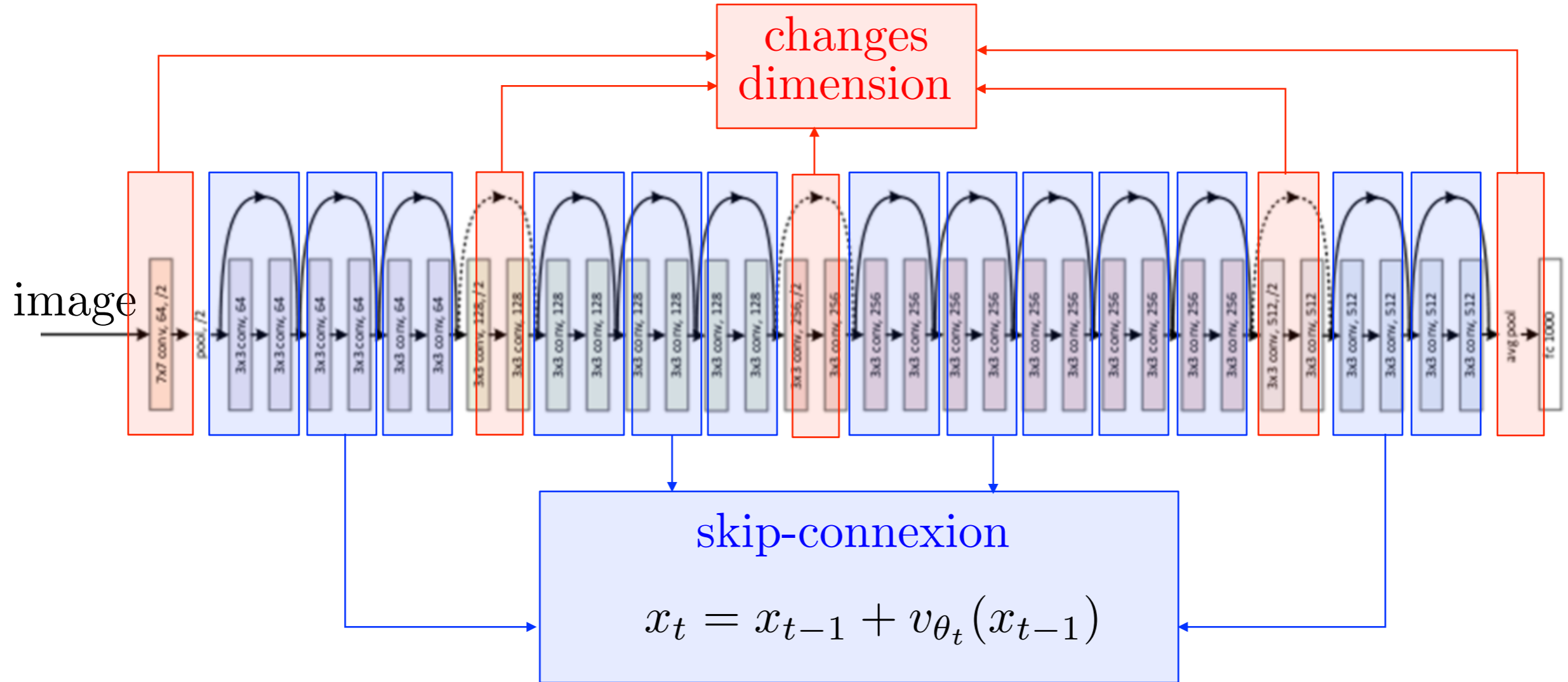
ResNet-type Architectures [He et al' 16]

ResNet-34

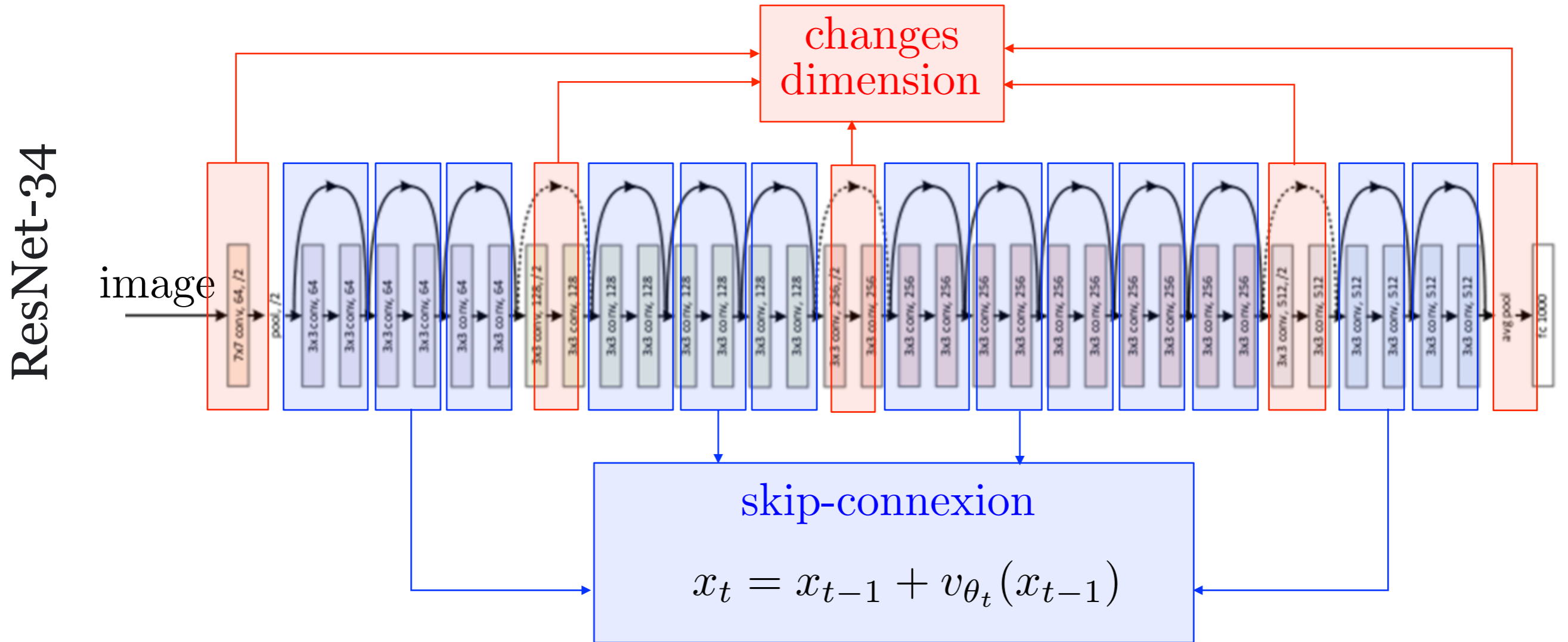


ResNet-type Architectures [He et al' 16]

ResNet-34

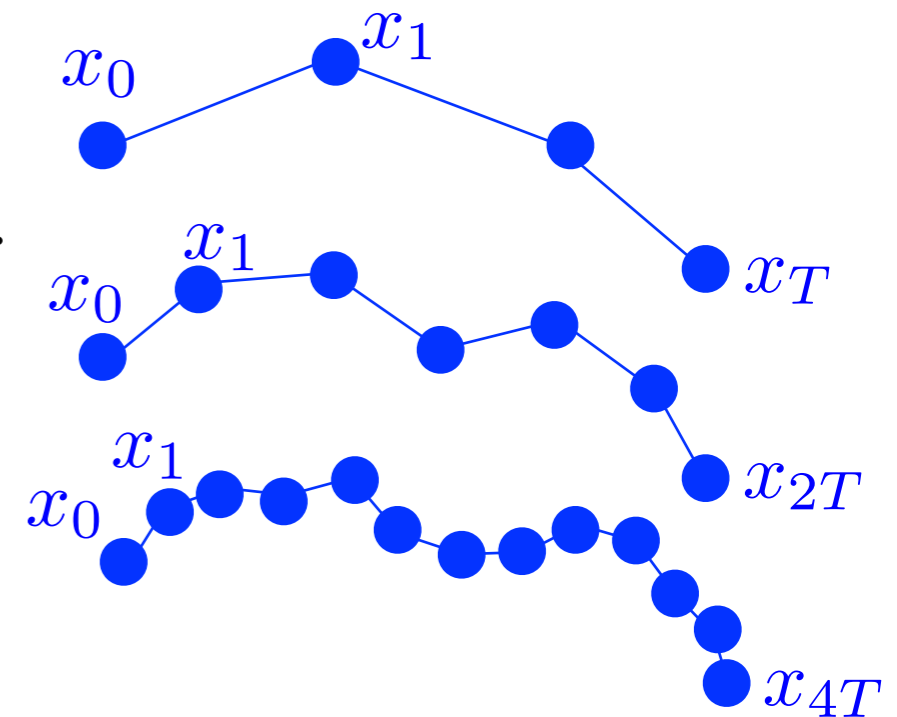


ResNet-type Architectures [He et al' 16]



→ Makes the “infinite depth” limit non-degenerate.

→ Enable $v_{\theta} = 0$ initialization, i.e. identity map.

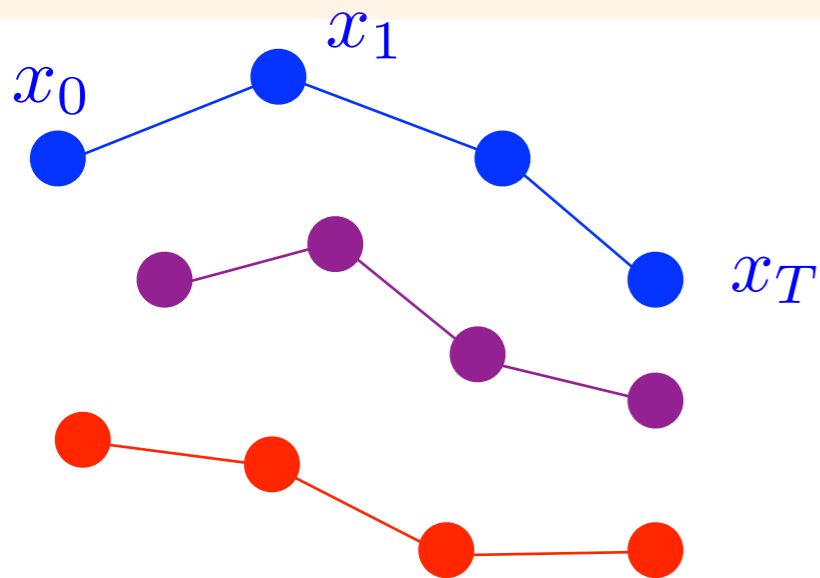


Infinite Depth and Neural-ODEs

ResNet [He et al, 2016]

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

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



Infinite Depth and Neural-ODEs

ResNet [He et al, 2016]

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

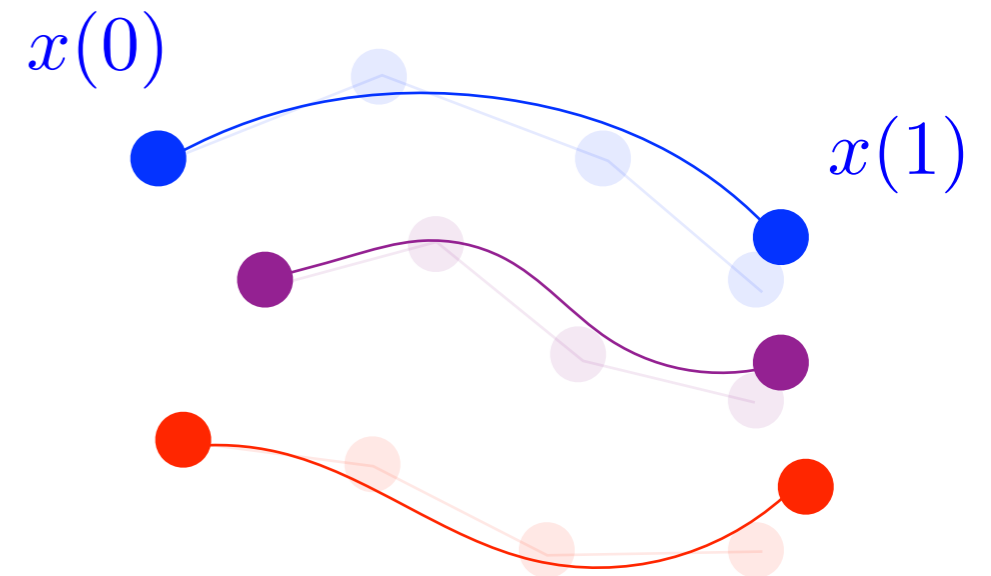
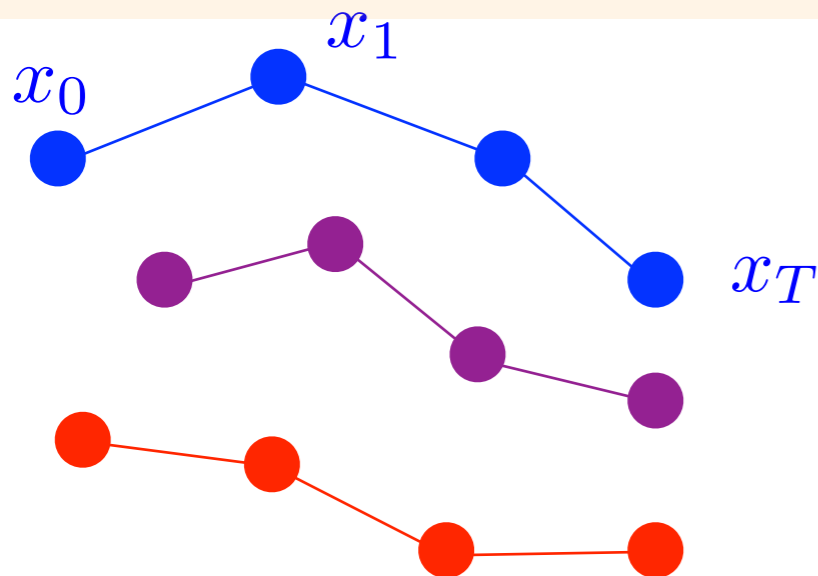
$$x_{t+1} = x_t + \frac{1}{T} v_{\theta_t}(x_t)$$

$T \rightarrow +\infty$

Neural ODE [Chen et al, 2018]

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

$$\frac{dx(t)}{dt} = v_{\theta(t)}(x(t))$$



Infinite Depth and Neural-ODEs

ResNet [He et al, 2016]

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

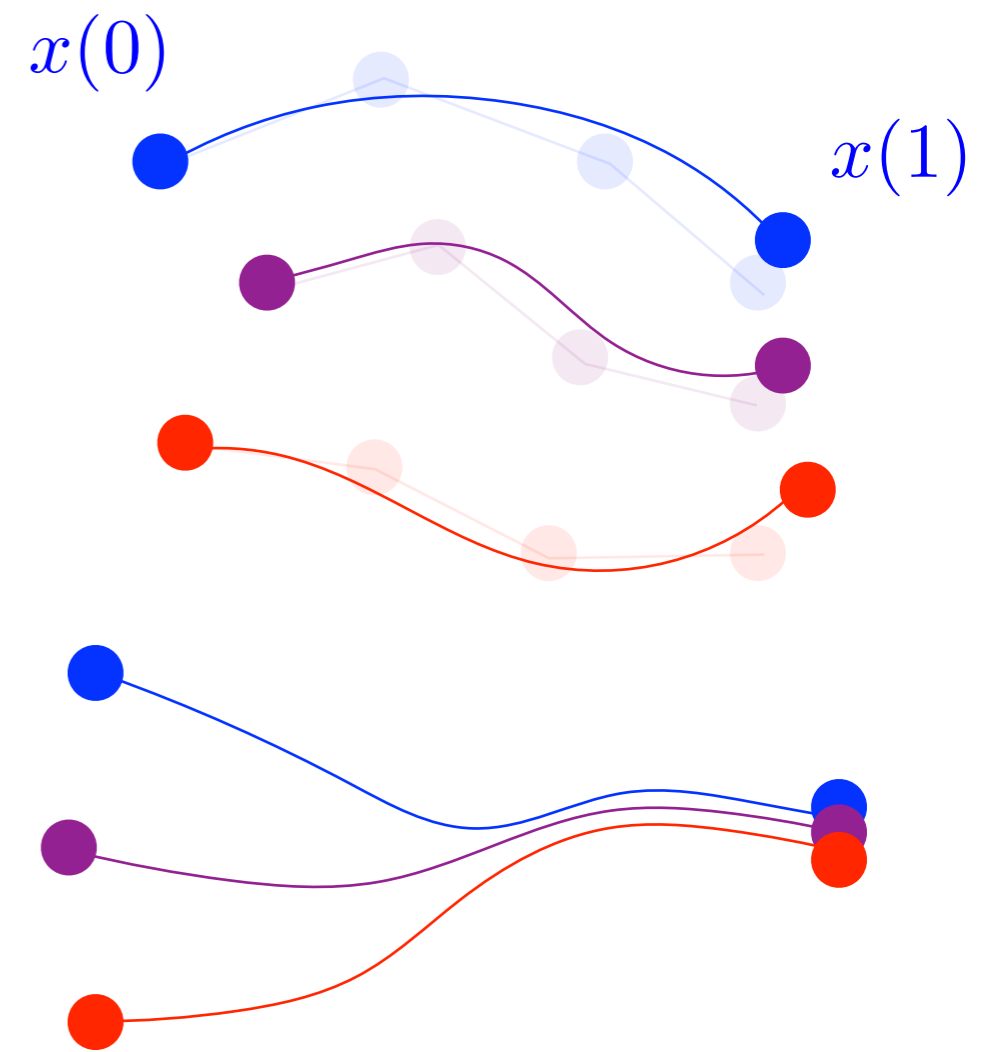
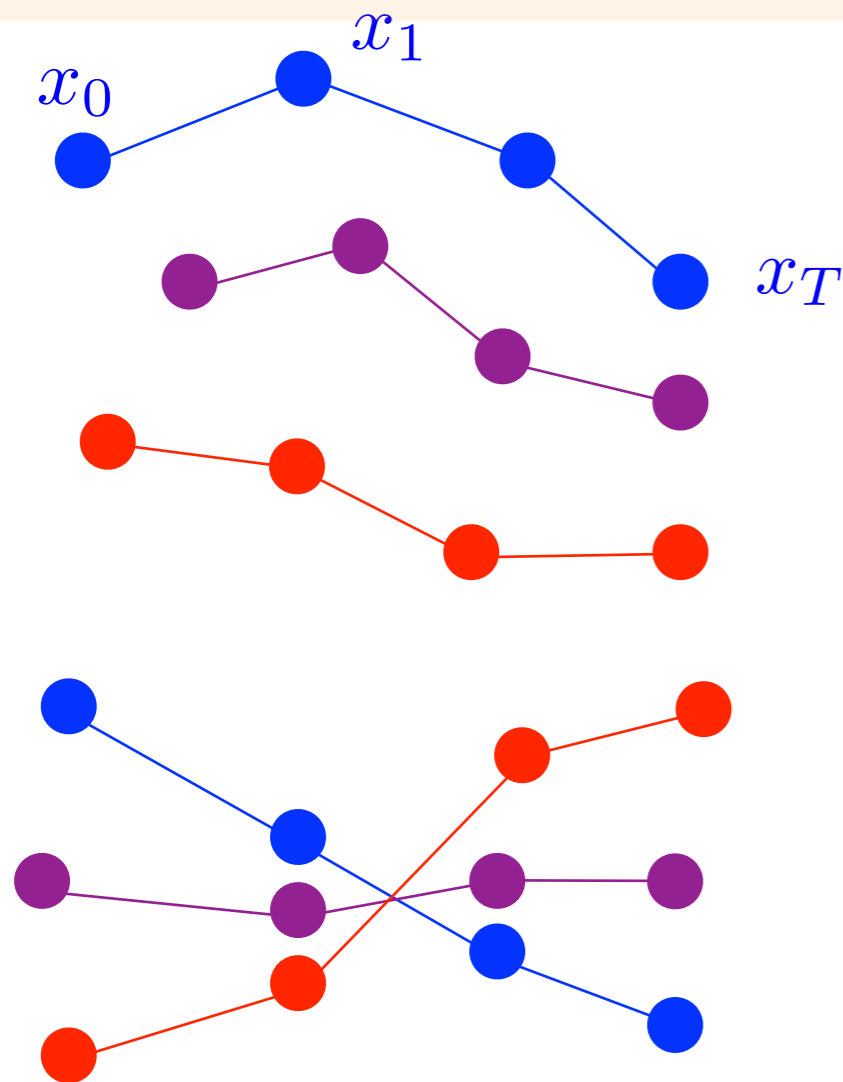
$$x_{t+1} = x_t + \frac{1}{T} v_{\theta_t}(x_t)$$

$$T \rightarrow +\infty$$

Neural ODE [Chen et al, 2018]

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

$$\frac{dx(t)}{dt} = v_{\theta(t)}(x(t))$$

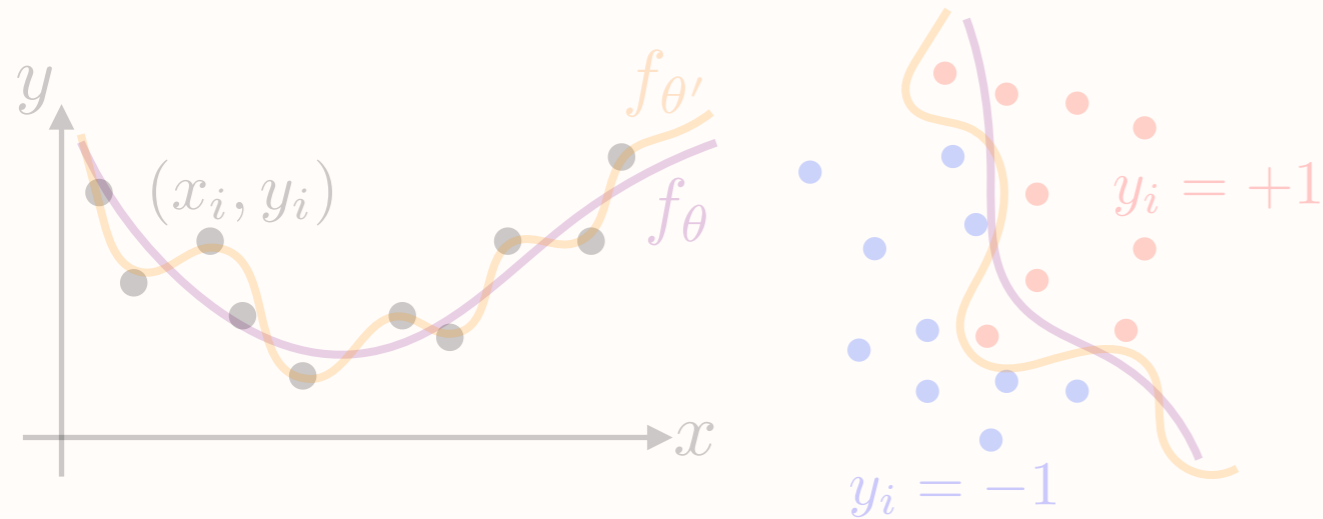


Trajectories cannot cross: Φ_{θ} defines a diffeomorphism.

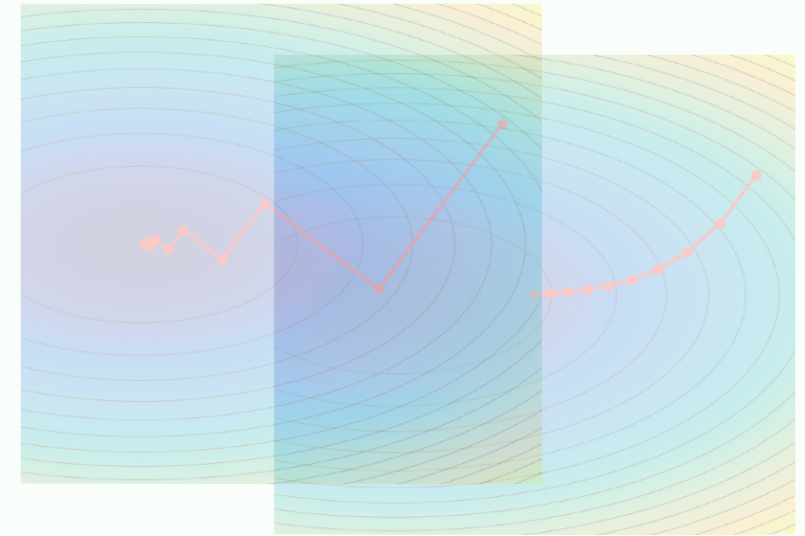
$T \rightarrow +\infty$ is a singular limit (θ can “explode” during training)

Overview

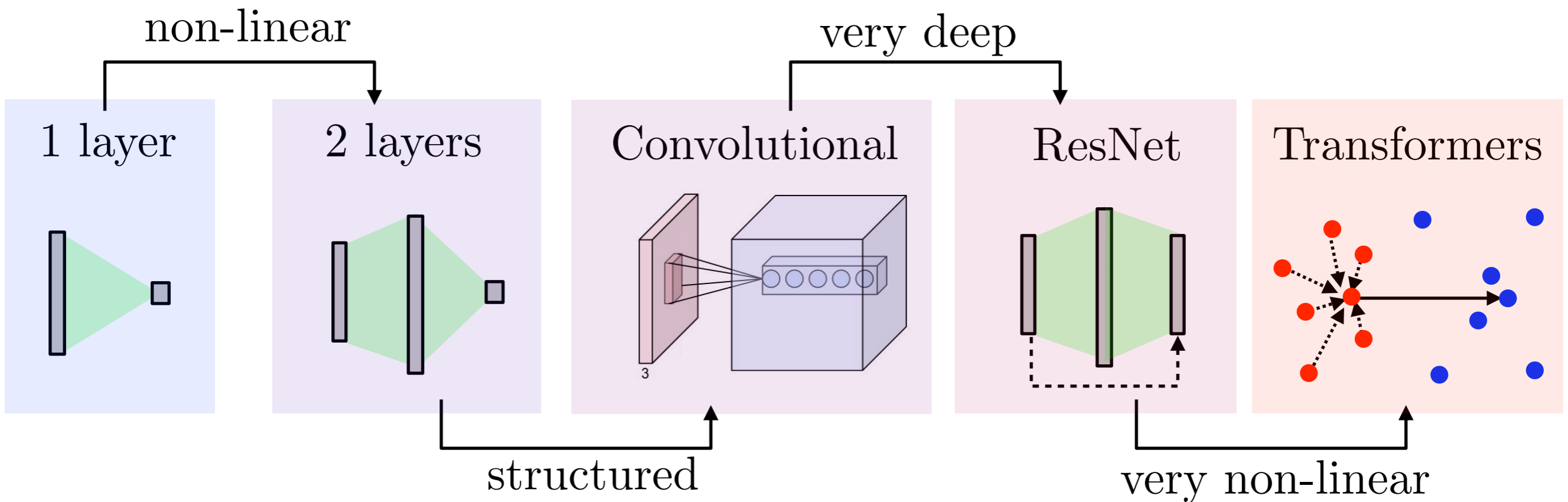
Machine Learning



Optimization

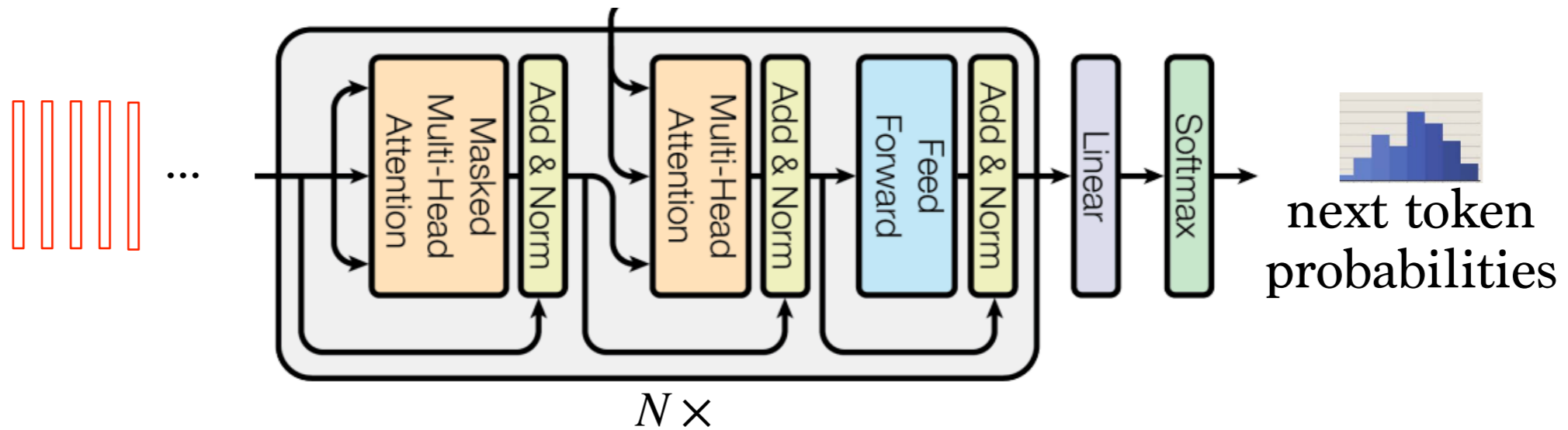
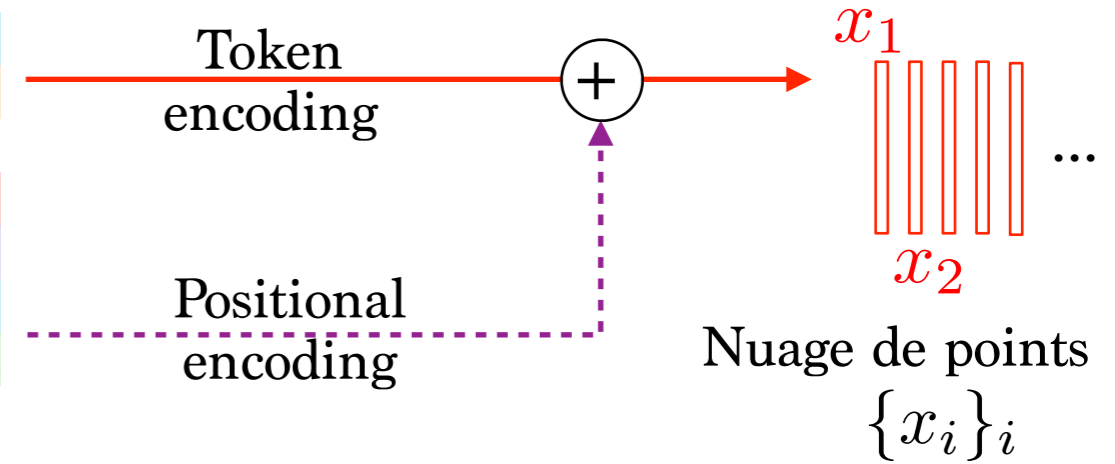
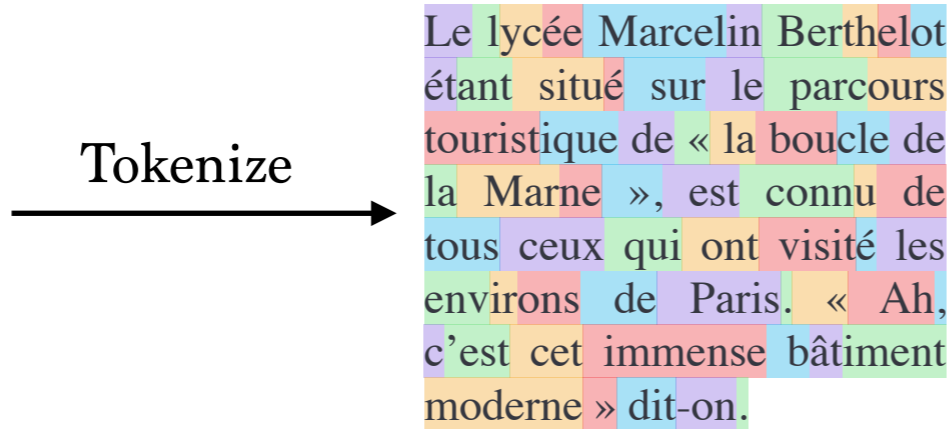


Networks Architectures

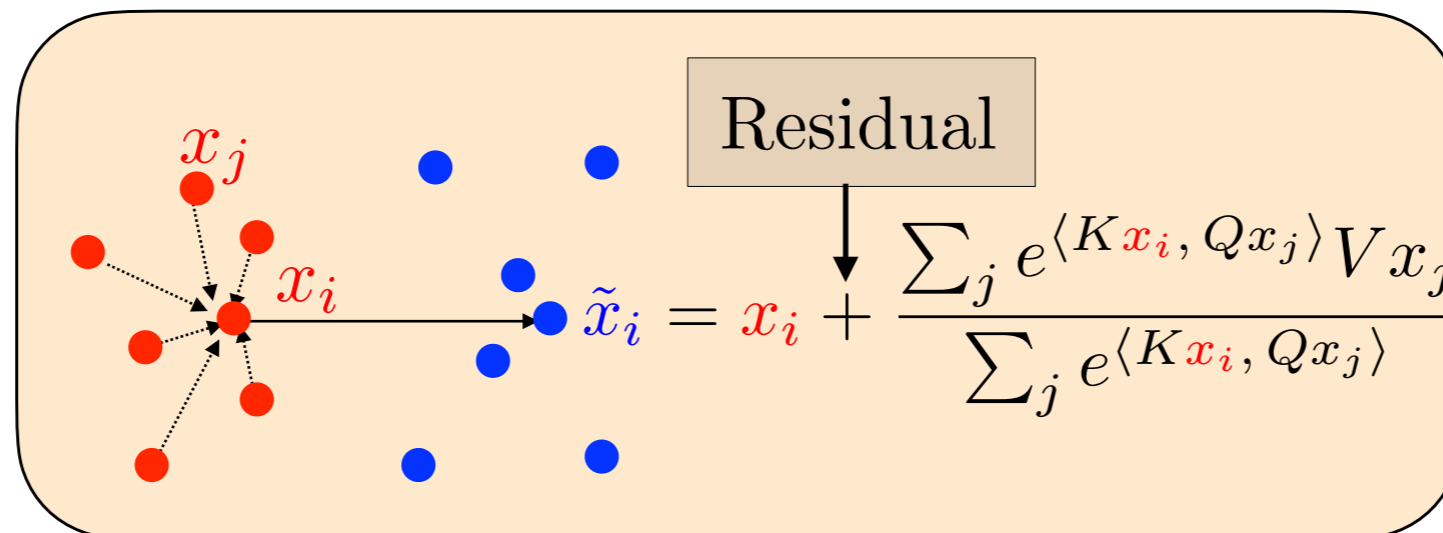


Les transformers et le mécanisme d'attention

Le lycée Marcelin Berthelot étant situé sur le parcours touristique de « la boucle de la Marne », est connu de tous ceux qui ont visité les environs de Paris. « Ah, c'est cet immense bâtiment moderne » dit-on.



Convolutions remplacées par des attentions



Large Language Models (LLM)

Pas seulement générer des données!

But: Répondre à des questions
Faire « tourner » des « algorithmes ».

Question:

Un transformer entraîné uniquement à faire du « next token prediction » est-il capable de raisonner?

Enjeu mathématique:

Modéliser et comprendre le phénomène de « in context learning ».



You

The Frobenius norm over matrices is $\|S, T\|_F \coloneqq \sqrt{\sum_{i,j} S_{i,j}^2 + T_{i,j}^2}$. Show that the Frobenius inner product can be written as $\langle S, T \rangle_F = \text{trace}(ST^\top)$.



ChatGPT

To demonstrate that the Frobenius inner product can be expressed as $\langle S, T \rangle_F = \text{trace}(ST^\top)$, let's start by defining the Frobenius inner product and the trace of a matrix.

- Frobenius Inner Product:** The Frobenius inner product of two matrices S and T is defined as $\langle S, T \rangle_F := \sum_{i,j} S_{i,j} T_{i,j}$, where $S_{i,j}$ and $T_{i,j}$ are the elements of matrices S and T at the i^{th} row and j^{th} column respectively.
- Trace of a Matrix:** The trace of a square matrix A is defined as the sum of its diagonal elements, i.e., $\text{trace}(A) = \sum_i A_{i,i}$.

Now, let's prove that $\langle S, T \rangle_F = \text{trace}(ST^\top)$:

- The product ST^\top results in a matrix where each element $(ST^\top)_{i,j}$ is the dot product of the i^{th} row of S and the j^{th} row of T . That is, $(ST^\top)_{i,j} = \sum_k S_{i,k} T_{j,k}$.

14/20

Large Language Models (LLM)

Pas seulement générer des données!

But: Répondre à des questions
Faire « tourner » des « algorithmes ».

Question:

Un transformer entraîné uniquement à faire du « next token prediction » est-il capable de raisonner?

Enjeu mathématique:

Modéliser et comprendre le phénomène de « in context learning ».



You

The Frobenius norm over matrices is $\|S, T\|_F \coloneqq \sqrt{\sum_{i,j} S_{i,j}^2 + T_{i,j}^2}$. Show that the Frobenius inner product can be written as $\langle S, T \rangle_F = \text{trace}(ST^\top)$.



ChatGPT

To demonstrate that the Frobenius inner product can be expressed as $\langle S, T \rangle_F = \text{trace}(ST^\top)$, let's start by defining the Frobenius inner product and the trace of a matrix.

- Frobenius Inner Product:** The Frobenius inner product of two matrices S and T is defined as $\langle S, T \rangle_F := \sum_{i,j} S_{i,j} T_{i,j}$, where $S_{i,j}$ and $T_{i,j}$ are the elements of matrices S and T at the i^{th} row and j^{th} column respectively.
- Trace of a Matrix:** The trace of a square matrix A is defined as the sum of its diagonal elements, i.e., $\text{trace}(A) = \sum_i A_{i,i}$.

Now, let's prove that $\langle S, T \rangle_F = \text{trace}(ST^\top)$:

- The product ST^\top results in a matrix where each element $(ST^\top)_{i,j}$ is the dot product of the i^{th} row of S and the j^{th} row of T . That is, $(ST^\top)_{i,j} = \sum_k S_{i,k} T_{j,k}$.

14/20

Theorem [Sander, Peyré 2024]:

If one trains a transformer on geometric sequences $\{(x_1, x_2 = Ax_1, \dots, x_t = Ax_{t-1}) : A, x_1\}$ then for an input $(x_1, x_2 = Bx_1, \dots, x_t = Bx_{t-1})$, the transformer:

1) estimates B

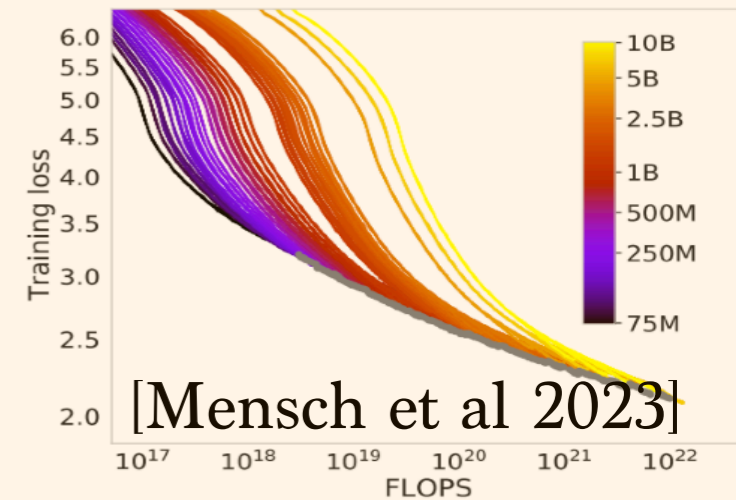
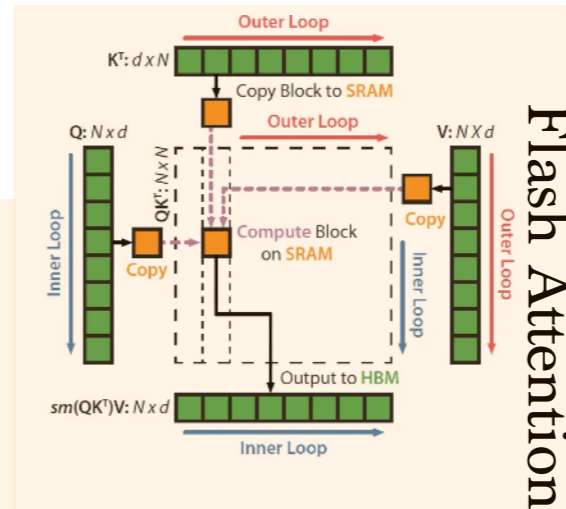
2) predicts $x_{t+1} = Bx_t$

Le futur ?

Le machine learning (IA) change à toute vitesse. Impossible de faire des prédictions à 5 ans ...

Numérique

High performance computing
Au delà de l'attention
(meilleures scaling laws) ?

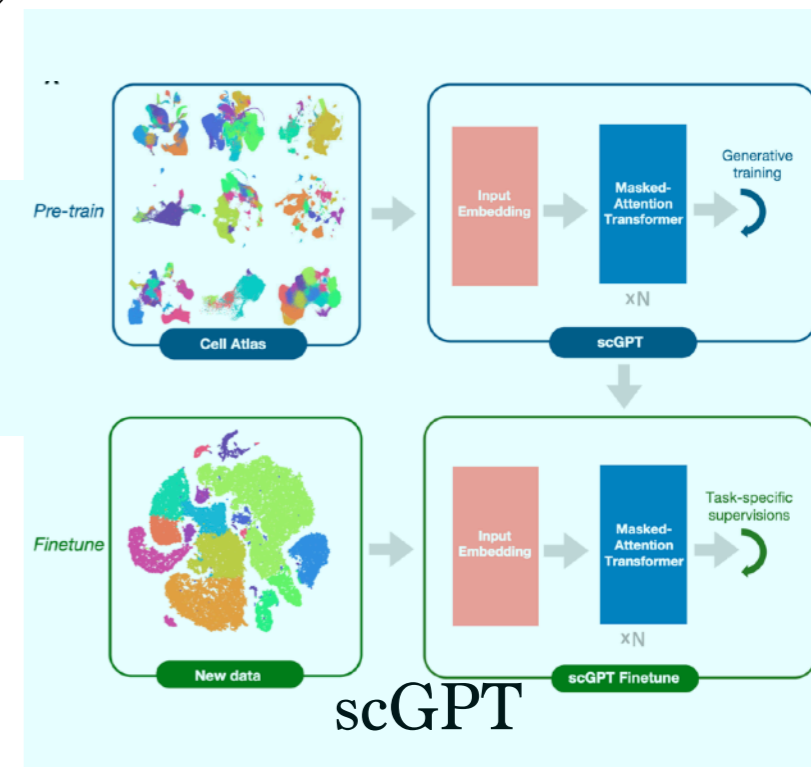


Théorie

Connecter modèles de diffusion (continue) et LLM (discret)
Comprendre le processus de « in context learning ».

Applications

AI for science : au delà du next token prediction
Vers le multimodal (video, 3D, etc).



Industrie

Open weights v.s. open source (données d'entraînement?)
Interface avec la vie quotidienne?