
Understanding the loss landscape of DNNs

Student: Luca Zancato
Supervisor: Prof. Alessandro Chiuso

Control Days 2019
9 May 2019, Padova



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Outline of the presentation:

- 1** Machine Learning Framework

Outline of the presentation:

- 1 Machine Learning Framework
- 2 Optimization for Deep Neural Networks

Outline of the presentation:

- 1 Machine Learning Framework
- 2 Optimization for Deep Neural Networks
- 3 Simplified model design

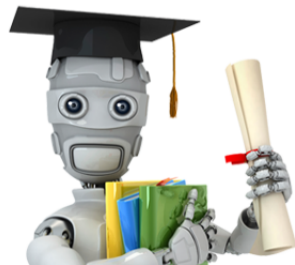
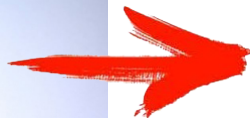
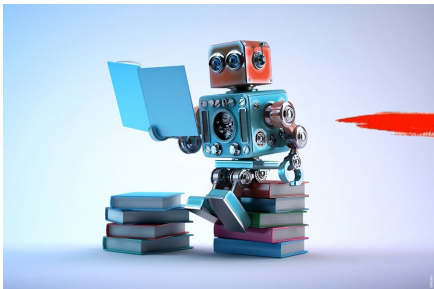
Outline of the presentation:

- 1 Machine Learning Framework
- 2 Optimization for Deep Neural Networks
- 3 Simplified model design
- 4 SGD Asymptotic and transient analysis

Experience



Expertise



Ingredients of Machine Learning:

- Training data: $(y_i, x_i) \quad 1, \dots, m$
- Model: $h_\theta \in \mathcal{H}$

Ingredients of Machine Learning:

- Training data: $(y_i, \mathbf{x}_i) \quad 1, \dots, m$
- Model: $h_\theta \in \mathcal{H}$

↳ Empirical Loss function: $L(\theta) := \sum_{i=1}^m l(y_i, h_\theta(\mathbf{x}_i))$



Ingredients of Machine Learning:

- Training data: $(y_i, \mathbf{x}_i) \quad 1, \dots, m$
- Model: $h_\theta \in \mathcal{H}$

Empirical Loss function: $L(\theta) := \sum_{i=1}^m l(y_i, h_\theta(\mathbf{x}_i))$

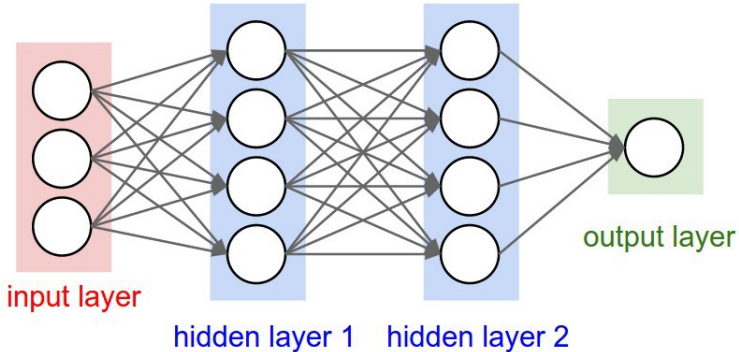


$$\theta_{opt} := \arg \min_{\theta \in \mathbb{R}^d} L(\theta)$$

Deep Neural Networks (DNNs)



1 What is a DNN? \implies Weighted directed acyclic graph



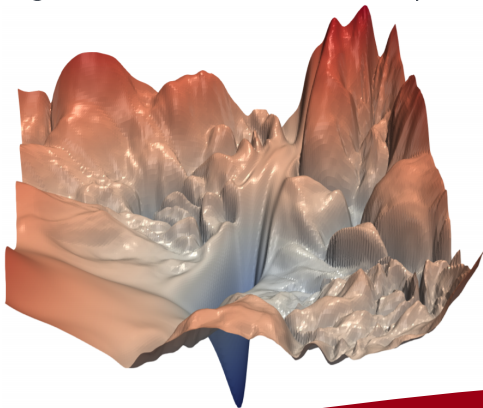
Composed by **many** layers performing non-linear transformations

- 1 What is a DNN? \implies Weighted directed acyclic graph
- 2 How to choose the parameters? \implies Gradient Descent

- 1 What is a DNN? \implies Weighted directed acyclic graph
- 2 How to choose the parameters? \implies Gradient Descent

Optimization problems:

- High dimensional **non convex** empirical loss function



- 1 What is a DNN? \implies Weighted directed acyclic graph
- 2 How to choose the parameters? \implies Gradient Descent

OPTIMIZATION PROBLEMS:

- High dimensional **non convex** empirical loss function
- **Computations**

- 1 What is a DNN? \implies Weighted directed acyclic graph
- 2 How to choose the parameters? \implies Gradient Descent

OPTIMIZATION PROBLEMS:

- High dimensional **non convex** empirical loss function
- **Computations**
- **Generalization** properties of θ_{opt} ?

Literature focuses on:

- 1 **Why** and **when** is deep learning **effective**?
- 2 Good **optimization** \implies Good **generalization**?

Literature focuses on:

- 1 **Why** and **when** is deep learning **effective**?
- 2 Good **optimization** \implies Good **generalization**?

By means of:

- **Stochastic optimization algorithms**
 - SGD
 - RMSProp
 - Adam

Literature focuses on:

- 1 **Why** and **when** is deep learning **effective**?
- 2 Good **optimization** \implies Good **generalization**?

By means of:

- **Stochastic optimization algorithms**
 - SGD
 - RMSProp
 - Adam
- **Geometry of the empirical loss landscape**



Gradient Descent

$$\theta_{k+1} = \theta_k - \eta \sum_{i=1}^m \underbrace{\nabla l(\theta_k, \mathbf{x}_i)}_{\text{gradient loss}}$$

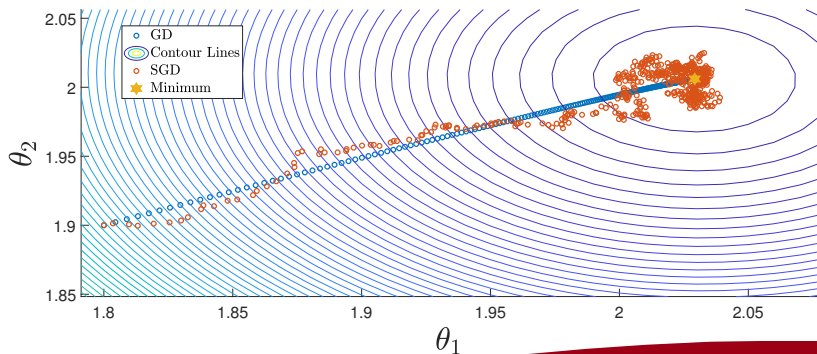
SGD

$$\theta_{k+1} = \theta_k - \eta \sum_{i=1}^B \nabla l(\theta_k, \mathbf{x}_i)$$

$$m \gg B$$

$$\text{SGD} \implies \theta_{k+1} = \theta_k - \underbrace{\eta \sum_{i=1}^m \nabla l(\theta_k, \mathbf{x}_i)}_{\text{GD dynamics}} + \underbrace{v_k}_{\text{noise term}}$$

GD vs SGD

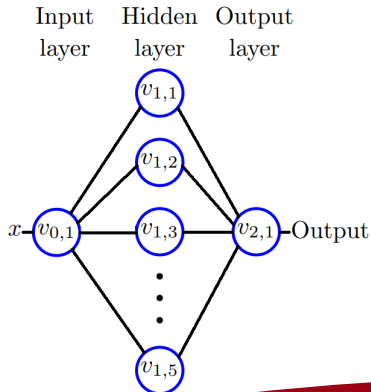


Simplified model and Local minima



Consider a regression task from \mathbb{R} to \mathbb{R} :

$$h_{\alpha, \mathbf{w}, \mathbf{b}}(\mathbf{x}) = \sum_{i=1}^d \alpha_i \tanh(\mathbf{w} \cdot (\mathbf{x} - \mathbf{b}_i))$$



Consider a regression task from \mathbb{R} to \mathbb{R} :

$$h_{\alpha, \mathbf{w}, \mathbf{b}}(\mathbf{x}) = \sum_{i=1}^d \alpha_i \tanh(\mathbf{w} \cdot (\mathbf{x} - \mathbf{b}_i))$$

$$L(\alpha, \mathbf{w}, \mathbf{b}, \mathbf{X}) = \sum_{i=1}^m (y_i - h_{\alpha, \mathbf{w}, \mathbf{b}}(\mathbf{x}_i))^2 = \|\mathbf{Y} - \Phi(\mathbf{w}, \mathbf{b}, \mathbf{X})\alpha\|^2$$

Consider a **regression** task from \mathbb{R} to \mathbb{R} :

$$h_{\alpha, \mathbf{w}, \mathbf{b}}(\mathbf{x}) = \sum_{i=1}^d \alpha_i \tanh(\mathbf{w} \cdot (\mathbf{x} - \mathbf{b}_i))$$

Fix the centers \mathbf{b}_i :

$$L(\alpha, \mathbf{w}, \mathbf{X}) = \|\mathbf{Y} - \Phi(\mathbf{w}, \mathbf{X})\alpha\|^2$$

And α ?

$$\arg \min_{\alpha, \mathbf{w}} L(\alpha, \mathbf{w}, \mathbf{X}) = \arg \min_{\mathbf{w}} \left\| \mathbf{Y} - \Phi(\mathbf{w}, \mathbf{X})\Phi^\dagger(\mathbf{w}, \mathbf{X})\mathbf{Y} \right\|^2$$

Design of a simplified model:

$$L(\underbrace{\theta}_{\mathbb{R}^d}) \implies L(\underbrace{w}_{\mathbb{R}^p}) \quad d \gg p$$

Simplified model and local minima



Design of a simplified model:

$$L(\underbrace{\theta}_{\mathbb{R}^d}) \implies L(\underbrace{w}_{\mathbb{R}^p}) \quad d \gg p$$

Is the geometry of the simplified model rich enough?

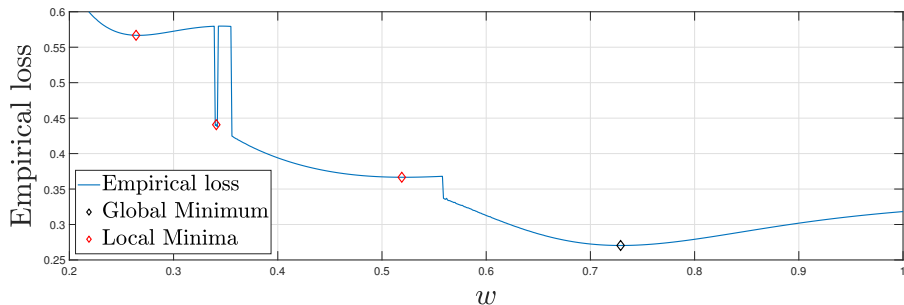
Simplified model and local minima



Design of a simplified model:

$$L(\underbrace{\theta}_{\mathbb{R}^d}) \implies L(\underbrace{w}_{\mathbb{R}^p}) \quad d \gg p$$

Is the geometry of the simplified model rich enough?



How to **approximate** the following stochastic difference equation?

$$\theta_{k+1} = \underbrace{\theta_k - \eta \nabla L(\theta_k)}_{\text{GD dynamics}} + \underbrace{v_k(\theta_k)}_{\text{colored noise}}$$



$$d\theta(t) = - \underbrace{\nabla L(\theta)}_{\text{drift term}} dt + \underbrace{\sqrt{\beta D(\theta)}}_{\text{diffusion matrix}} dW(t)$$

With $\beta = \frac{\eta}{B}$ **Temperature**

Let $\rho(\theta, t)$ be the probability distribution over the parameter space. It is well known that this is ruled by the **Fokker-Planck** equation:

$$\frac{\partial \rho}{\partial t} = \nabla \cdot \left(\nabla L(\theta) \rho + \beta \nabla \cdot (D(\theta) \rho) \right)$$

Let $\rho(\theta, t)$ be the probability distribution over the parameter space. It is well known that this is ruled by the **Fokker-Planck** equation:

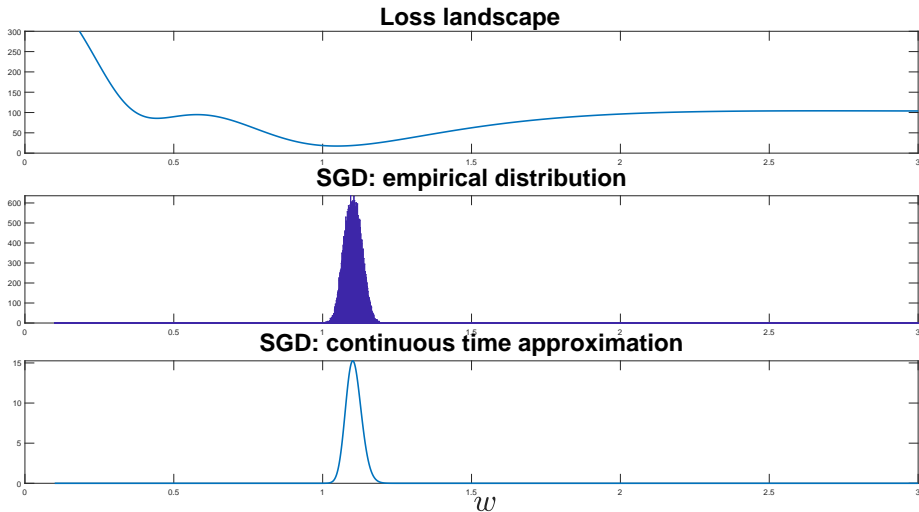
$$\frac{\partial \rho}{\partial t} = \nabla \cdot \left(\nabla L(\theta) \rho + \beta \nabla \cdot (D(\theta) \rho) \right)$$

Let the **Stationary distribution** $\rho^{SS}(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\rho^{SS}(\theta) \propto e^{-\beta^{-1} \Phi(\theta)} = e^{-\frac{E}{\eta} \Phi(\theta)}$$

where $\Phi(\theta)$ s.t. $\frac{\partial \rho^{SS}}{\partial t} = 0$

Asymptotic distribution approximation



Deterministic gradient vs Stochastic gradient

$$\nabla L(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla l_i(\theta) \implies \nabla G(\theta) = \frac{1}{B} \sum_{i \in \mathcal{B}} \nabla l_i(\theta)$$

with $i_1, \dots, i_B \in \mathcal{B}$ i.i.d. r.v. in the set $1, 2, \dots, m$

Deterministic gradient vs Stochastic gradient

$$\nabla L(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla l_i(\theta) \implies \nabla G(\theta) = \frac{1}{B} \sum_{i \in \mathcal{B}} \nabla l_i(\theta)$$

with $i_1, \dots, i_B \in \mathcal{B}$ i.i.d. r.v. in the set $1, 2, \dots, m$

$$\mathbb{E}_{i_1, i_2, \dots, i_B} [\nabla G(\theta)] = \nabla L(\theta)$$

Deterministic gradient vs Stochastic gradient

$$\nabla L(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla l_i(\theta) \implies \nabla G(\theta) = \frac{1}{B} \sum_{i \in \mathcal{B}} \nabla l_i(\theta)$$

with $i_1, \dots, i_B \in \mathcal{B}$ i.i.d. r.v. in the set $1, 2, \dots, m$

$$\mathbb{E}_{i_1, i_2, \dots, i_B} [\nabla G(\theta)] = \nabla L(\theta)$$

$$\text{Var}_{i_1, i_2, \dots, i_B} [\nabla G(\theta)] = \frac{1}{B} \left(\frac{\sum_{j=1}^m \nabla l_j(\theta) \nabla l_j(\theta)^T}{m} - \nabla G(\theta) \nabla G(\theta)^T \right)$$

Assumptions:

1 Quadratic loss:

$$l_i(\theta) = (y_i - h_\theta(\mathbf{x}_i))^2 \implies \nabla l_i(\theta) = -2 \underbrace{(y_i - h_\theta(\mathbf{x}_i))}_{:=e_\theta(\mathbf{x}_i)} \nabla h_\theta(\mathbf{x}_i)$$

Assumptions:

1 Quadratic loss:

$$l_i(\theta) = (y_i - h_\theta(\mathbf{x}_i))^2 \implies \nabla l_i(\theta) = -2 \underbrace{(y_i - h_\theta(\mathbf{x}_i))}_{:=e_\theta(\mathbf{x}_i)} \nabla h_\theta(\mathbf{x}_i)$$

2 Noisy data generative process:

$$y_i = h_{\theta^*}(\mathbf{x}_i) + \epsilon_i \quad \forall i$$

with ϵ_j i.i.d. with $\mathbb{E}[\epsilon_j] = 0$ and $\text{Var}[\epsilon_j] = \sigma^2$

Assumptions:

1 Quadratic loss:

$$l_i(\theta) = (y_i - h_\theta(\mathbf{x}_i))^2 \implies \nabla l_i(\theta) = -2 \underbrace{(y_i - h_\theta(\mathbf{x}_i))}_{:=e_\theta(\mathbf{x}_i)} \nabla h_\theta(\mathbf{x}_i)$$

2 Noisy data generative process:

$$y_i = h_{\theta^*}(\mathbf{x}_i) + \epsilon_i \quad \forall i$$

with ϵ_i i.i.d. with $\mathbb{E}[\epsilon_i] = 0$ and $\text{Var}[\epsilon_i] = \sigma^2$

3 θ_{loc} is a local minimum i.e. $\nabla L(\theta_{loc}) = 0 \not\Rightarrow \nabla G(\theta_{loc}) = 0$

Then:

$$e_{\theta}(x_i) = y_i - h_{\theta}(x_i) \underbrace{=}_{hyp1} h_{\theta^*}(x_i) + \epsilon_i - h_{\theta}(x_i) = \underbrace{h_{\theta^*}(x_i) - h_{\theta}(x_i)}_{:=\Delta_{\theta}(x_i)} + \epsilon_i$$

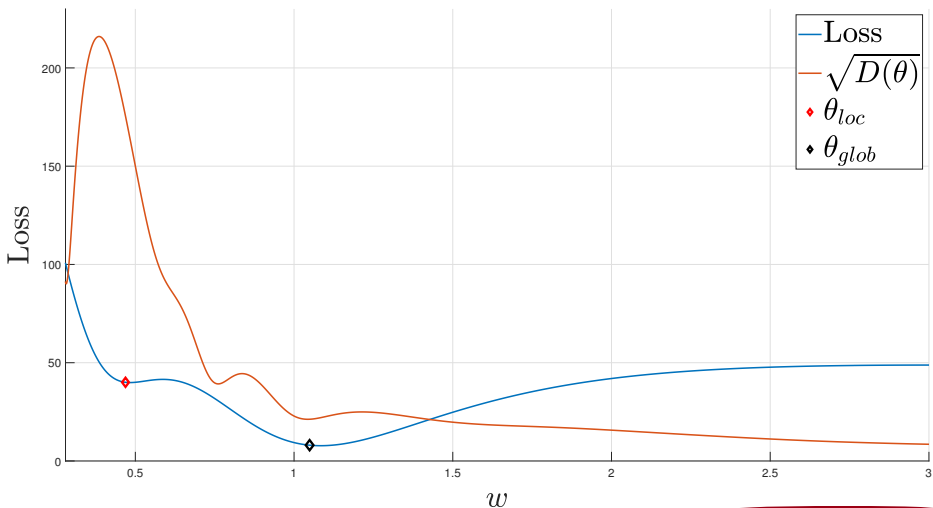
Then:

$$e_{\theta}(x_i) = y_i - h_{\theta}(x_i) \underbrace{=}_{hyp1} h_{\theta^*}(x_i) + \epsilon_i - h_{\theta}(x_i) = \underbrace{h_{\theta^*}(x_i) - h_{\theta}(x_i)}_{:=\Delta_{\theta}(x_i)} + \epsilon_i$$

$$\mathbb{E}_e[\text{Var}_w[G(\theta)]] \propto$$

$$\left[\underbrace{\sigma^2 \sum_{i=1}^m \nabla h_{\theta}(x_i) \nabla h_{\theta}(x_i)^T}_{\text{intrinsic noise term}} + \underbrace{\sum_{i=1}^m \Delta_{\theta}^2(x_i) \nabla h_{\theta}(x_i) \nabla h_{\theta}(x_i)^T}_{\text{exploration noise term}} \right]$$

Variance decomposition





Up to now we described **asymptotic** behaviours!

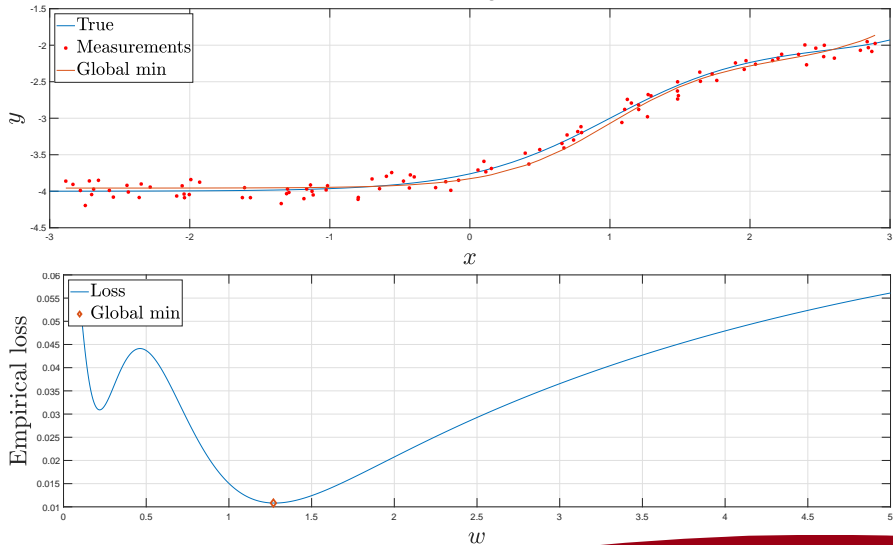
In practice we are also interested in reaching as quickly as possible the stationary distribution.

How to study the **transient dynamic**?

Genesis of local minima



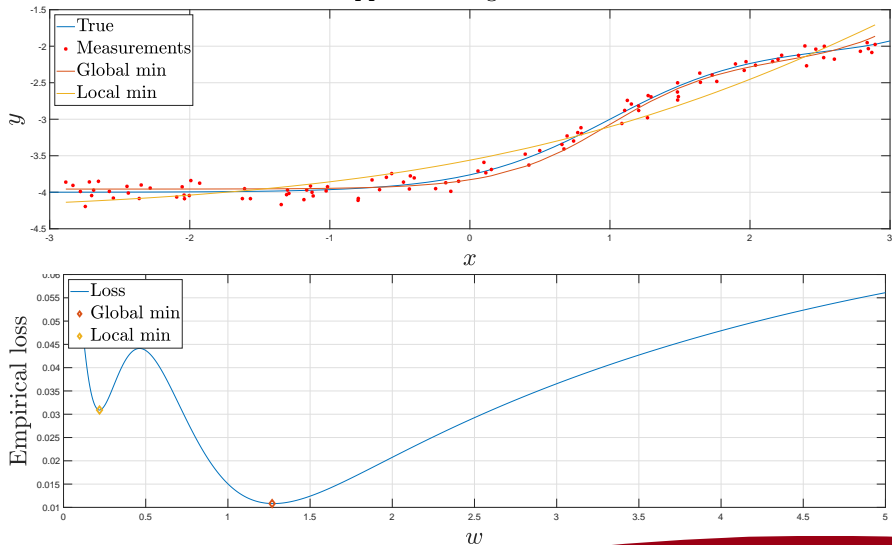
Approximating functions



Genesis of local minima



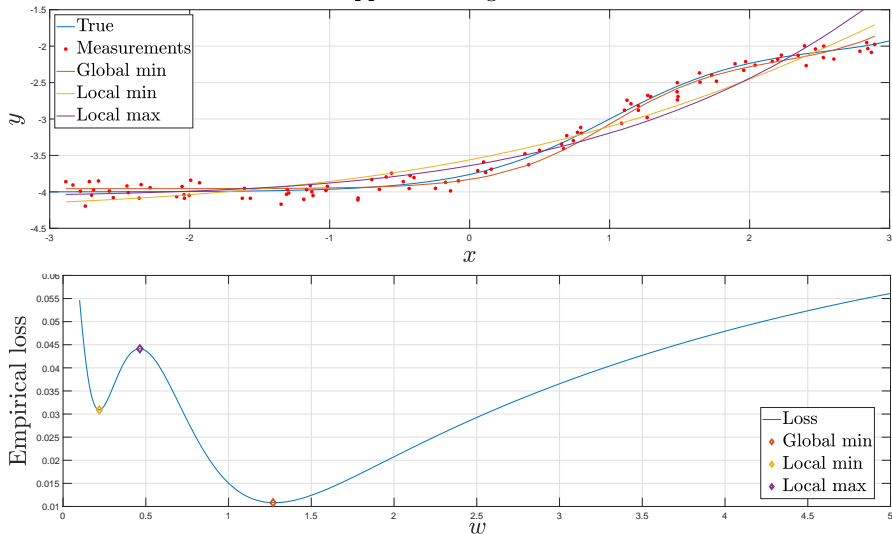
Approximating functions



Genesis of local minima



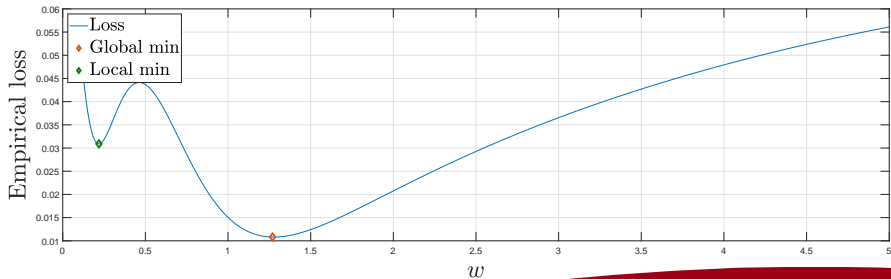
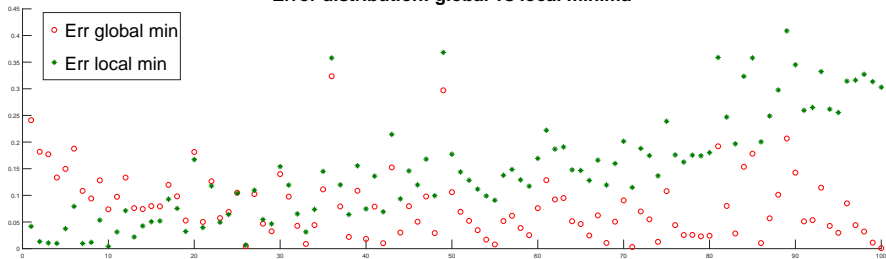
Approximating functions



A Local Analysis on the errors distribution



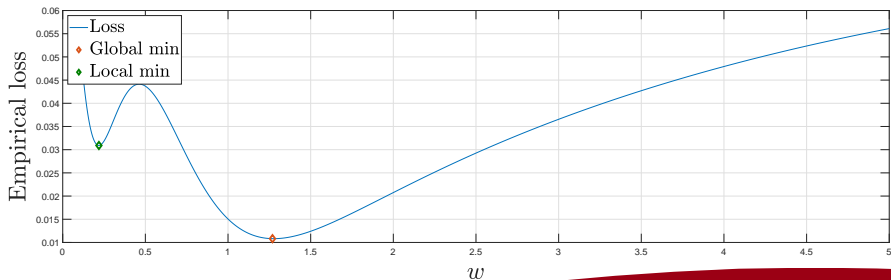
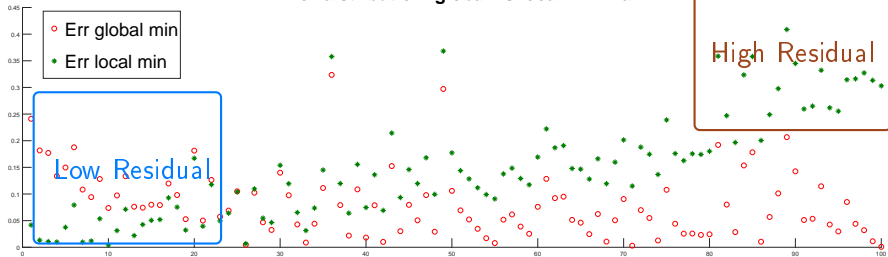
Error distribution: global vs local minima



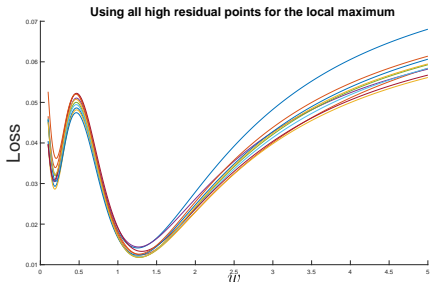
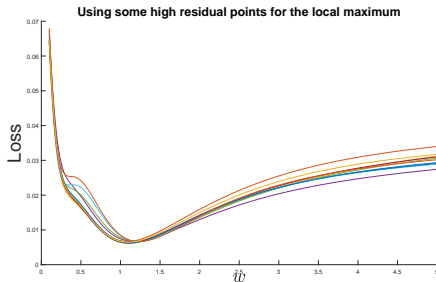
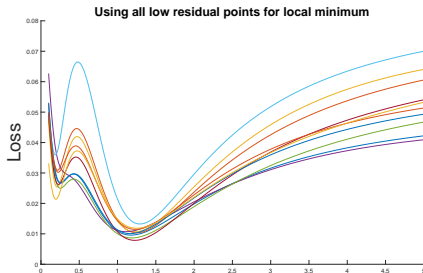
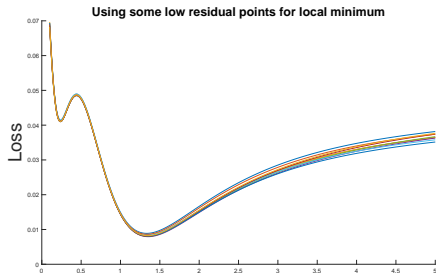
A Local Analysis on the errors distribution



Error distribution: global vs local minima



A Local Analysis on the errors distribution





Local maxima are slowing the optimization!

If caused by high residual \implies we can **easily identify** them



Local maxima are slowing the optimization!

If caused by high residual \implies we can **easily identify** them

Dynamic outliers: data points that are not well fitted by the model under the current parameter θ_k

Local maxima are slowing the optimization!

If caused by high residual \implies we can **easily identify** them

Dynamic outliers: data points that are not well fitted by the model under the current parameter θ_k

Importance sampling strategy:

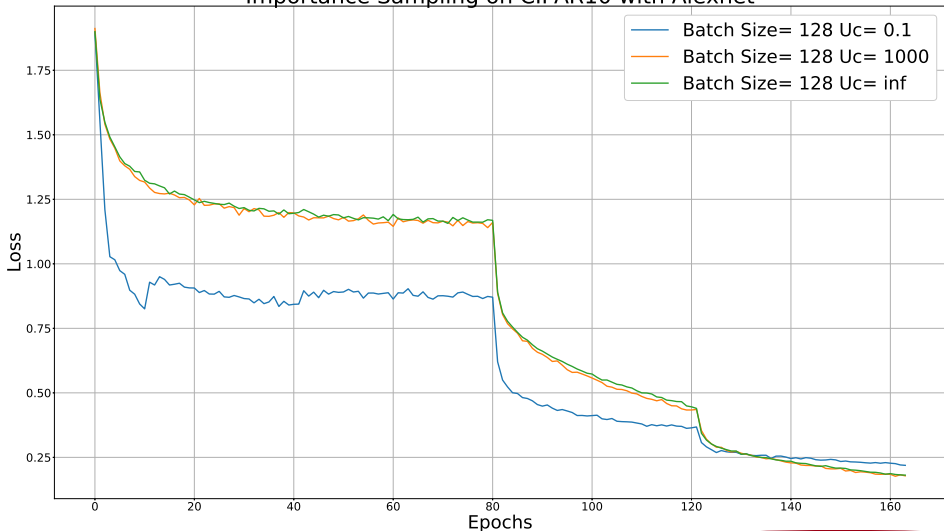
$$p(\mathbf{x}_i) := e^{-\frac{l_i(\theta_k)}{Uc}}$$

Note if $Uc \uparrow \implies$ reducing importance sampling

A new perspective: importance sampling



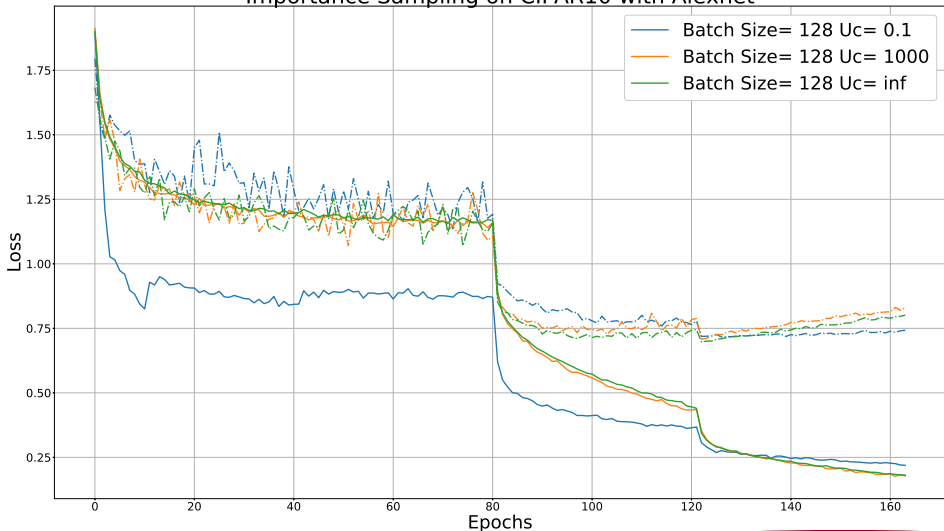
Importance Sampling on CIFAR10 with Alexnet



A new perspective: importance sampling



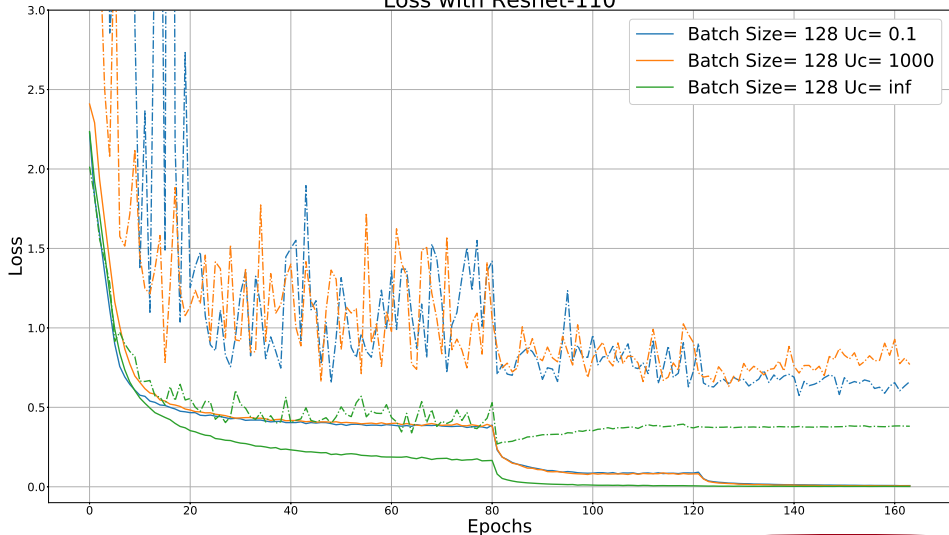
Importance Sampling on CIFAR10 with Alexnet



A new perspective: importance sampling



Loss with Resnet-110



Conclusions:

- 1 Stochasticity \implies escaping local minima

Conclusions:

- 1 Stochasticity \implies escaping local minima
- 2 SDE can be used to describe SGD

Conclusions:

- 1 Stochasticity \implies escaping local minima
- 2 SDE can be used to describe SGD
- 3 Instability of high potential local minima

Conclusions:

- 1 Stochasticity \implies escaping local minima
- 2 SDE can be used to describe SGD
- 3 Instability of high potential local minima
- 4 Genesis of local minima

Conclusions:

- 1 Stochasticity \implies escaping local minima
- 2 SDE can be used to describe SGD
- 3 Instability of high potential local minima
- 4 Genesis of local minima
- 5 Importance sampling \implies smoothing

Conclusions:

- 1 Stochasticity \implies escaping local minima
- 2 SDE can be used to describe SGD
- 3 Instability of high potential local minima
- 4 Genesis of local minima
- 5 Importance sampling \implies smoothing

Future work:

- 1 Local minima vs structure of DNNs?

Conclusions:

- 1 Stochasticity \implies escaping local minima
- 2 SDE can be used to describe SGD
- 3 Instability of high potential local minima
- 4 Genesis of local minima
- 5 Importance sampling \implies smoothing

Future work:

- 1 Local minima vs structure of DNNs?
- 2 Avoiding/Escaping local minima in higher dimension?

Conclusions:

- 1 Stochasticity \implies escaping local minima
- 2 SDE can be used to describe SGD
- 3 Instability of high potential local minima
- 4 Genesis of local minima
- 5 Importance sampling \implies smoothing

Future work:

- 1 Local minima vs structure of DNNs?
- 2 Avoiding/Escaping local minima in higher dimension?
- 3 Different importance sampling schemes

Thanks



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Thank you for your attention!