



# **Advanced deep learning optimizers and convergence**

# Acknowledgement

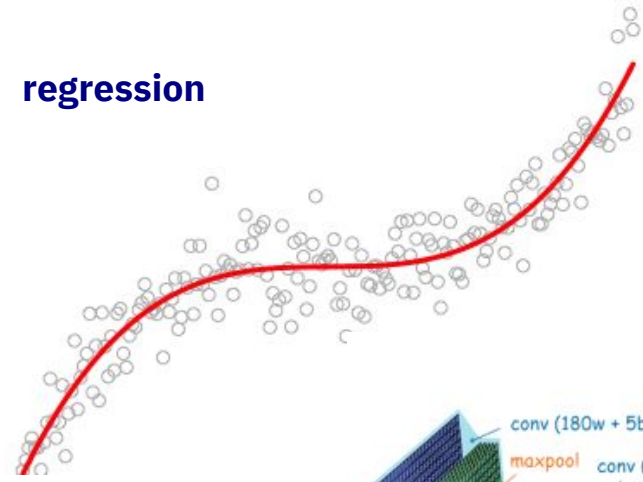
Special thanks to the work of Dr. Haitham Bou Ammar

Part of this lecture's content is based on the presentation of Machine Learning and AI Academy

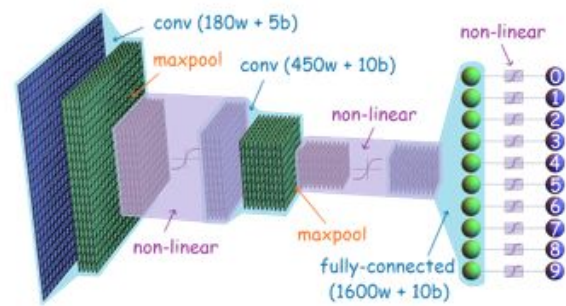
A group of experts in ML and AI with PhDs from top-tier schools and universities

# Why Optimization?

regression



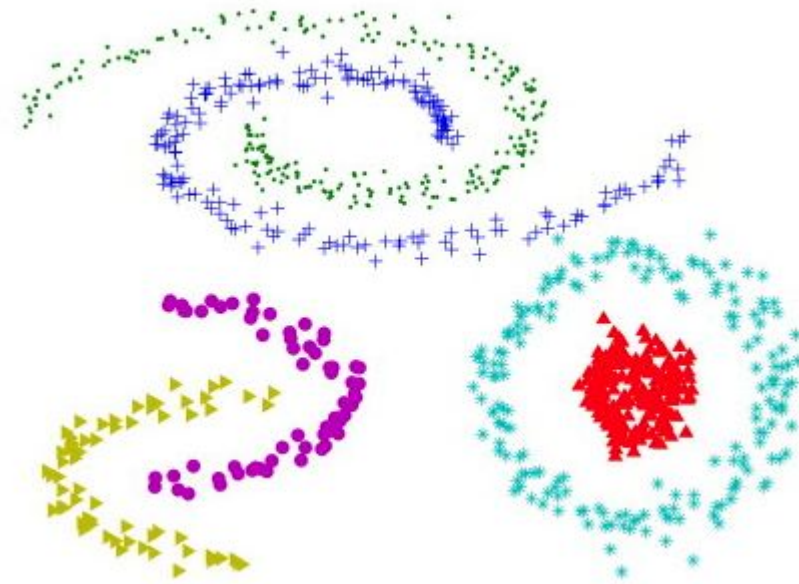
classification



## Supervised Learning

$$\min_{\theta} \frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\theta} \left( \mathbf{x}^{(i)}, y^{(i)} \right)$$

clustering/density estimation



## Unsupervised Learning

$$\min_{\theta} \frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\theta} \left( \mathbf{x}^{(i)} \right)$$

computer games



robotics

## Reinforcement Learning

$$\min_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} (\mathcal{R}_{\text{total}}(\tau))$$

... all these involve a minimization of some function ...

$$\min_{\theta \in \mathbb{R}^d} f(\theta)$$

# Optimization in Deep Learning

## Training Neural Network Models

Will your method hurt the convergence (optimization process)?

## Distributed Machine Learning

Will your method hurt the convergence (optimization process)?

## Gradients Compression

Will your method hurt the convergence (optimization process)?

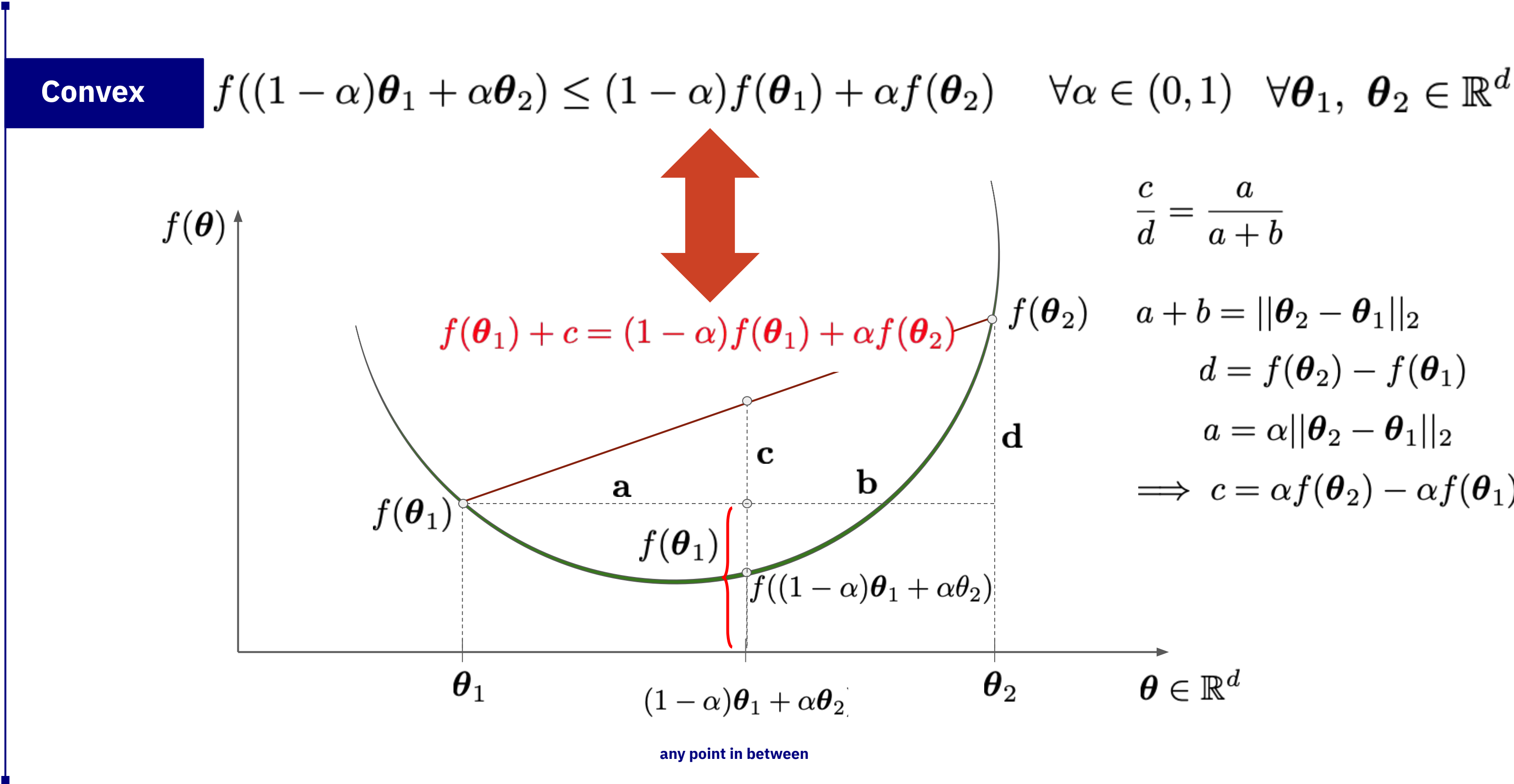
## Collaboration with application people

We are not allowed to change the model or dataset.

We can only improve the optimization part

# Function types, and what one can hope for ...

... optimizing for unknown parameters depends on the type of function under study ...



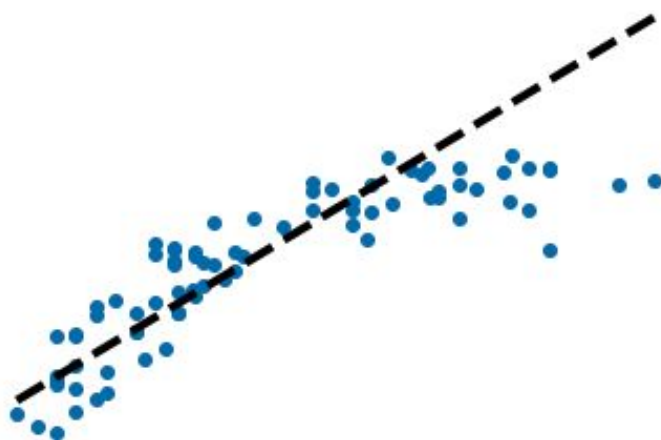


# Function types, and what one can hope for ...

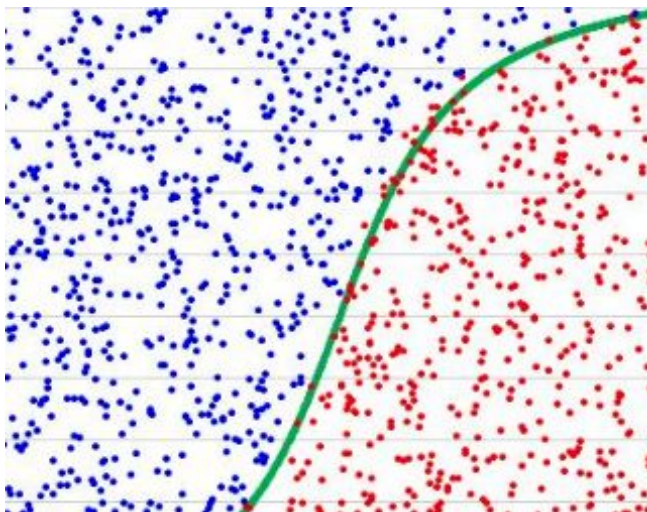
... optimizing for unknown parameters depends on the type of function under study ...

**Convex**

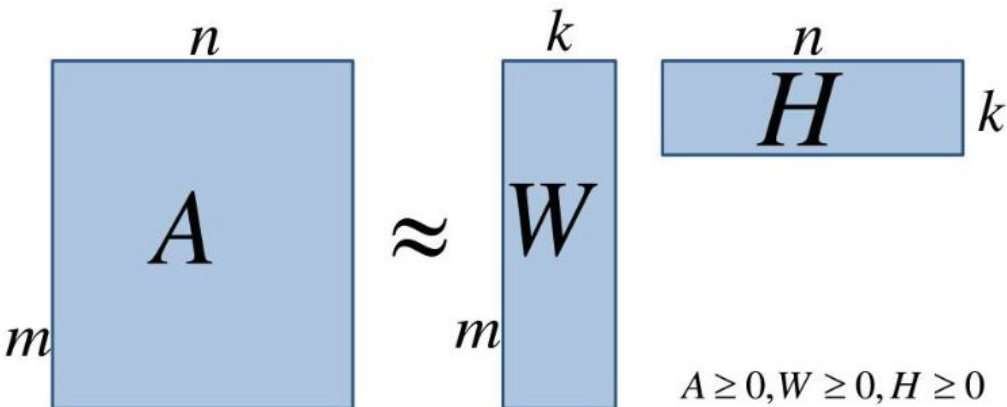
$$f((1-\alpha)\boldsymbol{\theta}_1 + \alpha\boldsymbol{\theta}_2) \leq (1-\alpha)f(\boldsymbol{\theta}_1) + \alpha f(\boldsymbol{\theta}_2) \quad \forall \alpha \in (0,1) \quad \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^d$$



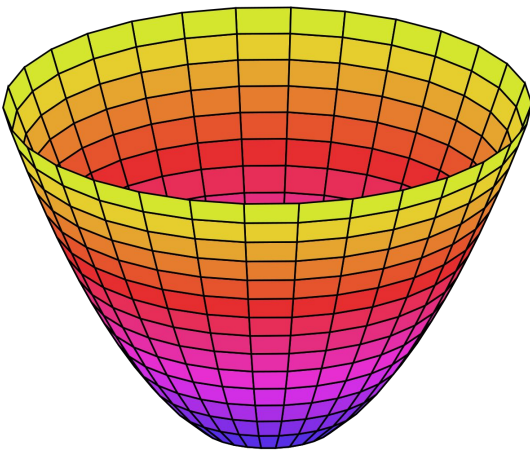
Linear Regression



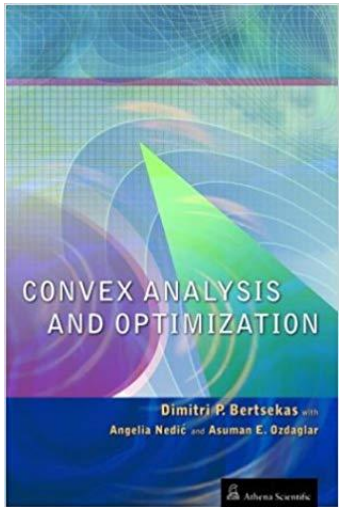
Classification with Hinge Loss

A diagram illustrating Non-Negative Matrix Factorisation. It shows a large blue square matrix A of size m by n, which is approximately equal to the product of a tall blue rectangle W of size m by k and a wide blue rectangle H of size k by n. Below the matrices, the conditions A ≥ 0, W ≥ 0, H ≥ 0 are specified.
$$A \approx WH \quad A \geq 0, W \geq 0, H \geq 0$$

Non-Negative Matrix Factorisation



Unique global minimum



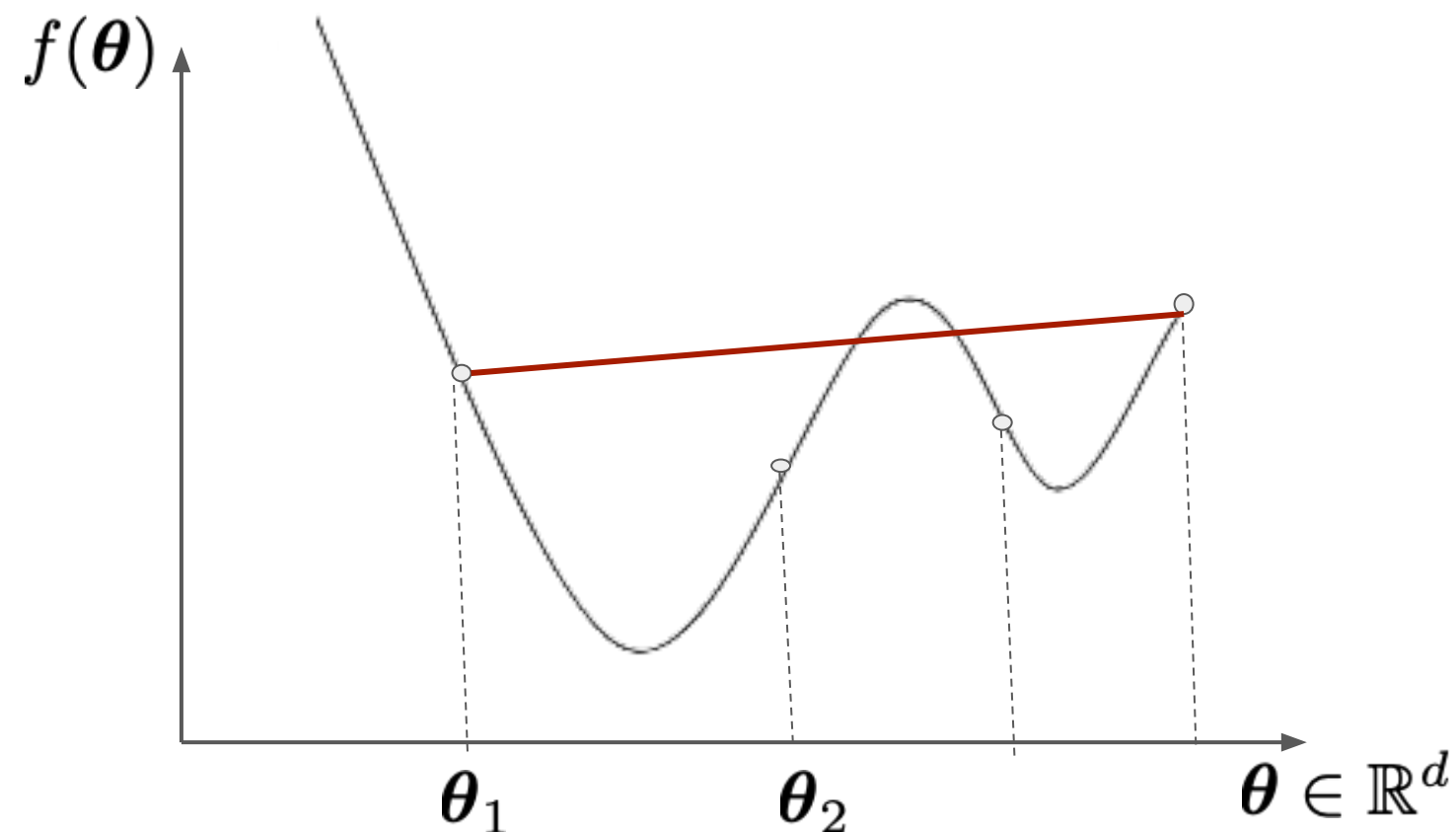
# Function types, and what one can hope for ...

... optimising for unknown parameters depends on the type of function under study ...

**Non-Convex**

... we want to negate the convex definition (and avoid concave definition) ...

$$\begin{aligned} \exists \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \text{ and } \alpha \in (0, 1) \text{ such that } f((1 - \alpha)\boldsymbol{\theta}_1 + \alpha\boldsymbol{\theta}_2) &> (1 - \alpha)f(\boldsymbol{\theta}_1) + \alpha f(\boldsymbol{\theta}_2) \\ \exists \tilde{\boldsymbol{\theta}}_1, \tilde{\boldsymbol{\theta}}_2, \text{ and } \tilde{\alpha} \in (0, 1) \text{ such that } f((1 - \tilde{\alpha})\tilde{\boldsymbol{\theta}}_1 + \tilde{\alpha}\tilde{\boldsymbol{\theta}}_2) &< (1 - \tilde{\alpha})f(\tilde{\boldsymbol{\theta}}_1) + \tilde{\alpha}f(\tilde{\boldsymbol{\theta}}_2) \end{aligned} \quad \&$$



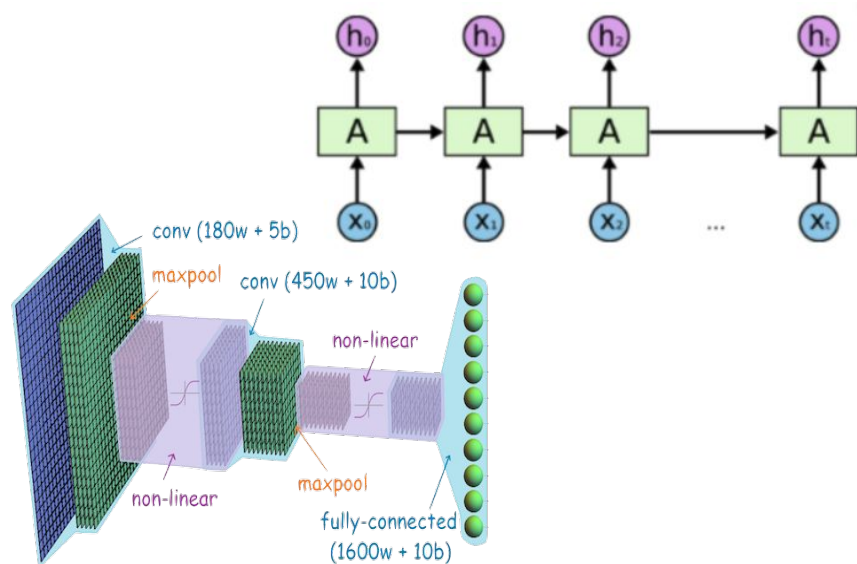
# Function types, and what one can hope for ...

... optimising for unknown parameters depends on the type of function under study ...

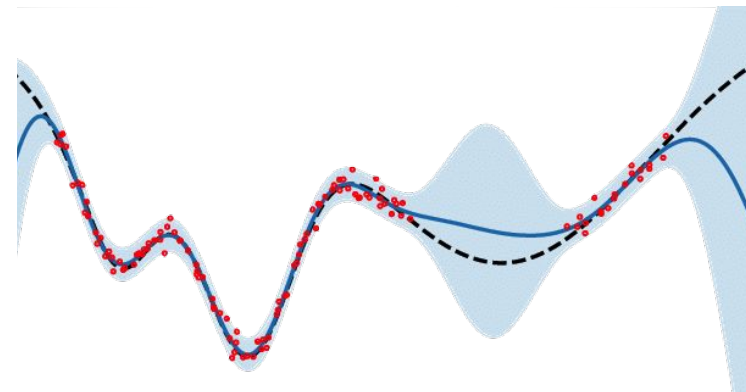
**Non-Convex**

... we want to negate the convex definition (and avoid concave definition) ...

$$\begin{aligned} \exists \theta_1, \theta_2, \text{ and } \alpha \in (0, 1) \text{ such that } f((1 - \alpha)\theta_1 + \alpha\theta_2) &> (1 - \alpha)f(\theta_1) + \alpha f(\theta_2) \\ \exists \tilde{\theta}_1, \tilde{\theta}_2, \text{ and } \tilde{\alpha} \in (0, 1) \text{ such that } f((1 - \tilde{\alpha})\tilde{\theta}_1 + \tilde{\alpha}\tilde{\theta}_2) &< (1 - \tilde{\alpha})f(\tilde{\theta}_1) + \tilde{\alpha}f(\tilde{\theta}_2) \end{aligned} \quad \&$$



**Deep Learning**



**Gaussian Processes & Bayesian Models**



**Reinforcement Learning**

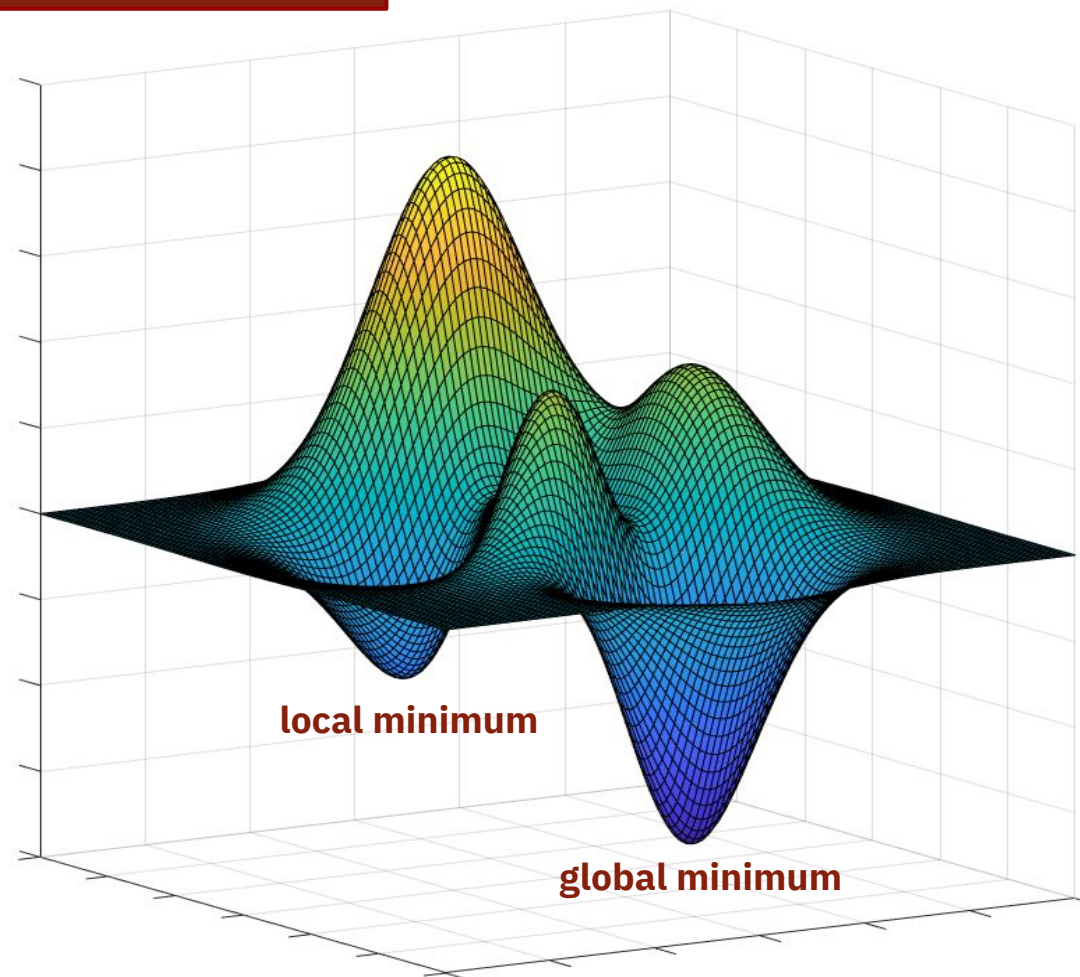


# Function types, and what one can hope for ...

... optimising for unknown parameters depends on the type of function under study ...

**Non-Convex**

... global and local minima (checking) are NP-Hard, we look for other types of points ...



$$\nabla_{\theta} f(\theta_{\text{stationary}}) = \mathbf{0}$$

... so instead, we are searching for stationary points ...

**1.  $\epsilon$ -First-Order-Stationary Point (FOSP):**  $\|\nabla_{\theta} f(\theta_{\text{FOSP}})\|_2 \leq \epsilon$

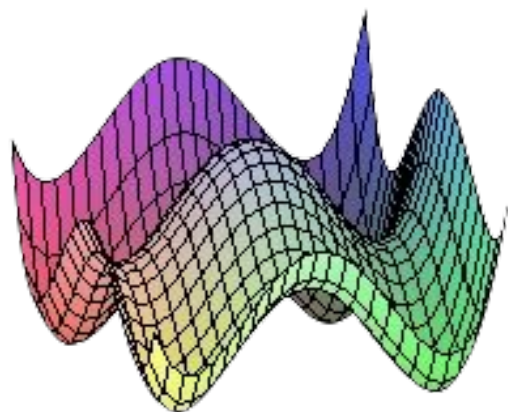
[e.g., all global and local minima, saddle points, plateau points]

**2.  $\epsilon$ - Second-Order-Stationary Point (SOSP):**

$$\|\nabla_{\theta} f(\theta_{\text{SOSP}})\|_2 \leq \epsilon \quad \text{and} \quad \lambda_{\min}(\nabla_{\theta, \theta}^2 f(\theta_{\text{SOSP}})) \geq -\sqrt{\epsilon}$$

[e.g., all global and local minima, plateau points]

# Algorithms vary in type of information used ...



## First-Order Methods

GD

SGD

ADAM

NAGD

AdaGrad

RMSProp

adaptive

Momentum

## Second-Order Methods

Newton Method

Regularised  
Newton Method

Stochastic  
Quasi-Newton



## Zero-Order Methods

Bayesian  
Optimisation

StoS00

Stroqu00L

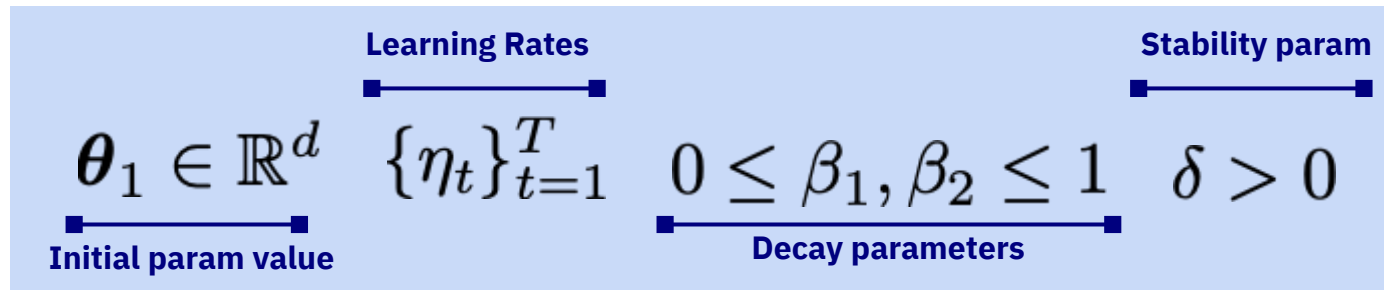
Non-Convex Optimisation

The background features a large, dark blue triangle on the right side, pointing towards the top right. To its left is a lighter blue triangle pointing towards the bottom left. The area between these two triangles is white. A diagonal gradient bar, transitioning from light blue at the top to light red at the bottom, runs from the bottom left towards the center.

# **A review on first-order optimizers**

# An example of loss function (ignore bias correction)

## Algorithm's Inputs:



## Update Procedure:

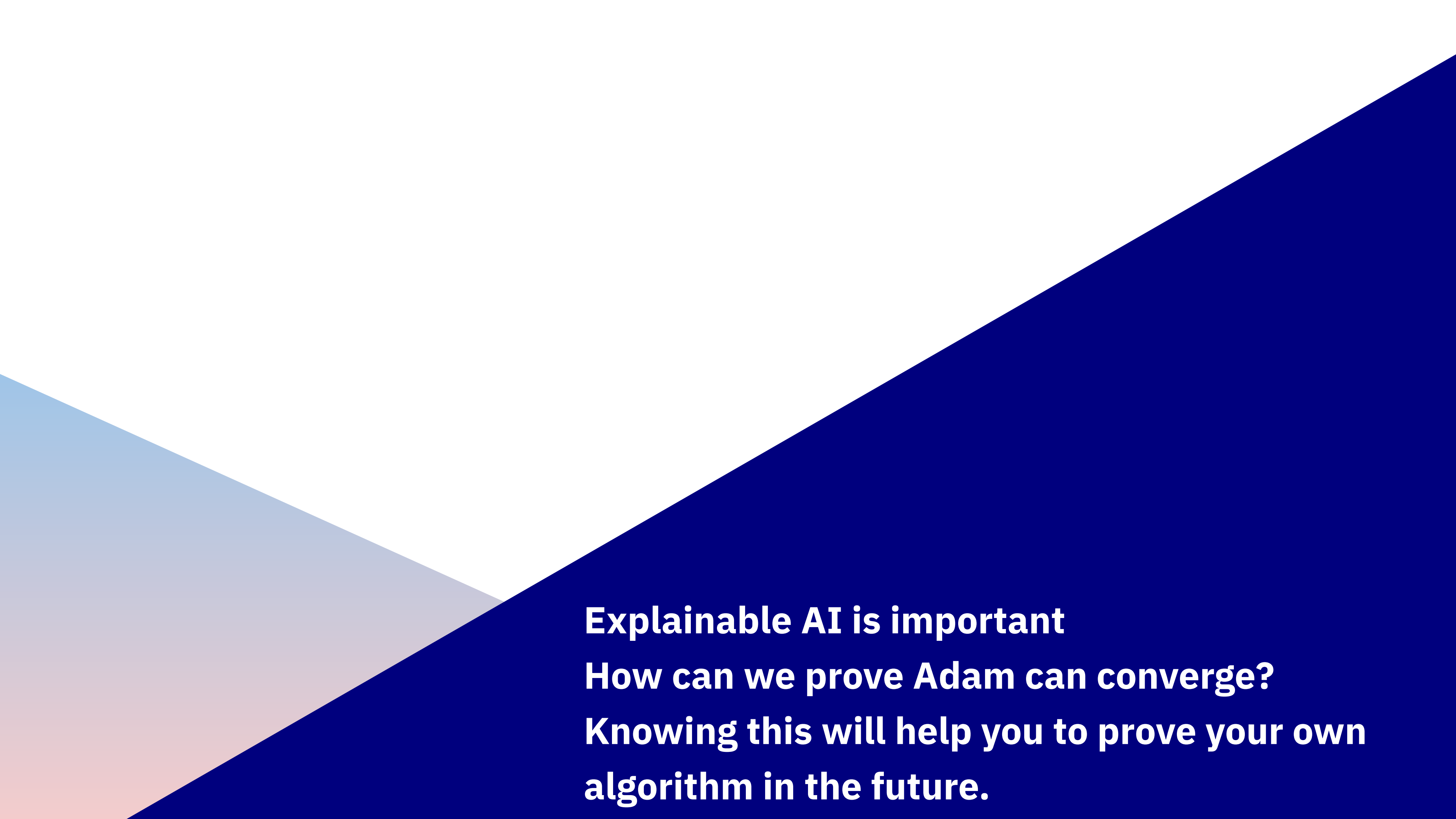
Set  $\mathbf{m}_0 = \mathbf{0}$ , and  $\mathbf{v}_0 = \mathbf{0}$   
**for**  $t = 1$  **to**  $T$  **do**  
  Draw a sample  $\xi_t$  from  $\mathbb{P}$   
  Compute  $\mathbf{g}_t = \nabla \mathcal{L}(\theta_t, \xi_t)$   
  Update  $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$   
  Update  $\mathbf{v}_t = \mathbf{v}_{t-1} - (1 - \beta_2)(\mathbf{v}_{t-1} - \mathbf{g}_t^2)$   
  Update  $\theta_{t+1} = \theta_t - \eta_t \frac{\mathbf{g}_t}{(\sqrt{\mathbf{v}_t} + \delta)}$   
**end for**

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(\mathbf{x}_i))^2$$

$$\text{Sample } \xi_t = i_t \in \{1, \dots, n\}$$

$$\Rightarrow \mathcal{L}(\theta, i_t) = (y_{i_t} - f_{\theta}(\mathbf{x}_{i_t}))^2$$

$$\begin{aligned} \nabla_{\theta} \mathcal{L}(\theta, i_t) &= \nabla_{\theta} (y_{i_t} - f_{\theta}(\mathbf{x}_{i_t}))^2 \\ &= -2(y_{i_t} - f_{\theta}(\mathbf{x}_{i_t})) \nabla f_{\theta}(\mathbf{x}_{i_t}) \end{aligned}$$



**Explainable AI is important  
How can we prove Adam can converge?  
Knowing this will help you to prove your own  
algorithm in the future.**



# From ML to ERM (Empirical risk minimization)

consider the following form of the objective function:  $\mathbb{E}_{\xi \sim \mathbb{P}} [\mathcal{L}(\boldsymbol{\theta}; \xi)]$

... for e.g., in regression

$$\xi \sim \text{Uniform}[1, n], \text{ then } \mathbb{E}_{\xi \sim \text{Uniform}}[(y_\xi - f_{\boldsymbol{\theta}}(\mathbf{x}_\xi))^2] = \frac{1}{n} \sum_{i=1}^n (y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i))^2$$

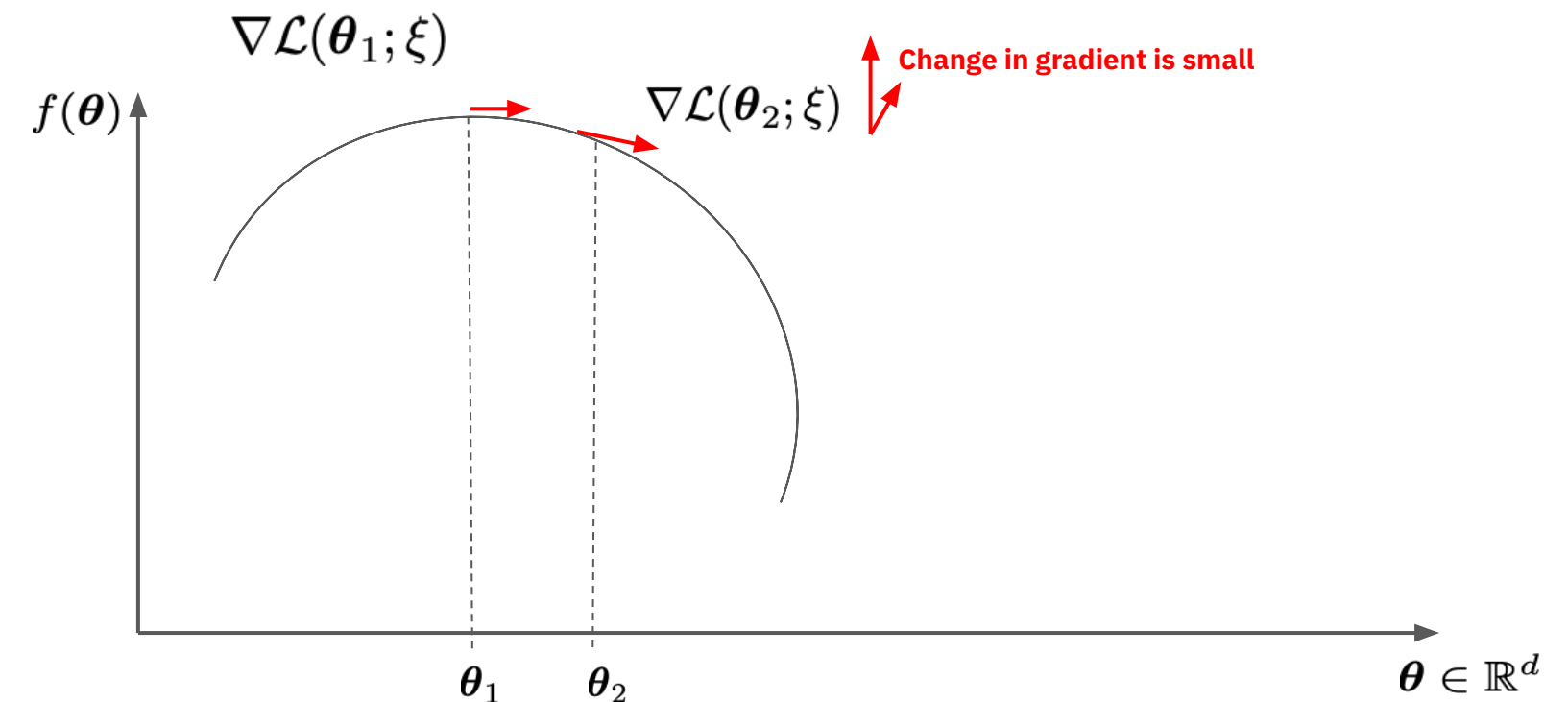
... now, our goal is to minimize the following

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{\xi \sim \mathbb{P}} [\mathcal{L}(\boldsymbol{\theta}; \xi)]$$

... using ADAM from the previous slide

Assumption I -- Loss Function is L-Smooth:

$$\|\nabla \mathcal{L}(\boldsymbol{\theta}_2; \xi) - \nabla \mathcal{L}(\boldsymbol{\theta}_1; \xi)\|_2 \leq L \|\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1\|_2 \quad \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \xi$$



# Proof Roadmap ...

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_t) + \dots +$$

$$\mathbb{E}[\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) \dots + \mathbb{E}\left[\frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} \middle| \boldsymbol{\theta}_t\right]$$

... we need to bound these ...

$$\dots + \mathbb{E}\left[\frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t\right]$$

**Objective Func. L-Smoothness**

... relation between 2 consecutive iterations ...

**True Components to Bound**

... consider randomness plug-in update rule, and achieve terms to bound ...

**Bounding the first term**

...

...

**Bound in terms of gradient norm**

**Choose hyper-params**



done ✓

**Bounding the second term**

...

...

**Bound in terms of batch-size**



# Convergence Proof ...

... as in any other optimization proof, we need to understand the change in function value between two consecutive iterations of the algorithm:

$$f(\boldsymbol{\theta}_{t+1}) \leq f(\boldsymbol{\theta}_t) - \Delta \implies \text{convergence to some point if the function is lower-bounded}$$

Some non-negative value

## Monotone Convergence Theorem (MCT)

... now, if we can say that the objective function is L-smooth, then we can have a relation between function values on two successive iterations:

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_t) + \nabla^\top \mathcal{L}(\boldsymbol{\theta}_t) (\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t) + \frac{L}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|_2^2$$

Relation between successive iterations

**But how to show that our objective function is L-Smooth**

Please see <https://xingyuzhou.org/blog/notes/Lipschitz-gradient>



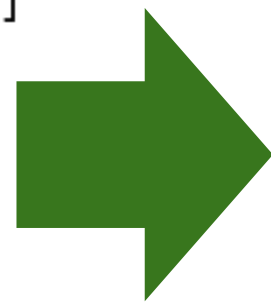
# Convergence Proof ...

... let us study the norm of the difference between the gradients of the objective function at any two given input points:

$$\begin{aligned} \|\nabla \mathcal{L}(\boldsymbol{\theta}_1) - \nabla \mathcal{L}(\boldsymbol{\theta}_2)\|_2 &= \|\nabla \mathbb{E}_{\xi}[\mathcal{L}(\boldsymbol{\theta}_1; \xi)] - \nabla \mathbb{E}_{\xi}[\mathcal{L}(\boldsymbol{\theta}_2; \xi)]\|_2 \\ &= \|\mathbb{E}_{\xi}[\nabla \mathcal{L}(\boldsymbol{\theta}_1; \xi)] - \mathbb{E}_{\xi}[\nabla \mathcal{L}(\boldsymbol{\theta}_2; \xi)]\|_2 \\ &= \|\mathbb{E}_{\xi}[\nabla \mathcal{L}(\boldsymbol{\theta}_1; \xi) - \nabla \mathcal{L}(\boldsymbol{\theta}_2; \xi)]\|_2 \\ &\leq \mathbb{E}_{\xi}[\|\nabla \mathcal{L}(\boldsymbol{\theta}_1; \xi) - \nabla \mathcal{L}(\boldsymbol{\theta}_2; \xi)\|_2] \\ &\leq \mathbb{E}_{\xi}[L\|\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1\|_2] \\ &= L\|\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1\|_2 \end{aligned}$$

Assumption I -- Loss Function is L-Smooth:

$$\|\nabla \mathcal{L}(\boldsymbol{\theta}_2; \xi) - \nabla \mathcal{L}(\boldsymbol{\theta}_1; \xi)\|_2 \leq L\|\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1\|_2 \quad \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \xi$$



**Objective function is L-Smooth**

# Convergence Proof ...

**... since we just proved that our objective is L-Smooth, now we can write that the objective value between two successive iterations abides by:**

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_t) + \nabla^\top \mathcal{L}(\boldsymbol{\theta}_t) (\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t) + \frac{L}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|_2^2$$

**... now, remember our update rules from the pseudo-code in the previous slides, we can write:**


$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t \implies \text{with } \beta_1 = 0, \text{ then } \mathbf{m}_t = \mathbf{g}_t \text{ then } \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \frac{\mathbf{g}_t}{(\sqrt{\mathbf{v}_t} + \delta)}$$

$$\text{... component-wise update } \theta_{i,t+1} = \theta_{i,t} - \eta_t \frac{\mathbf{g}_{i,t}}{(\sqrt{\mathbf{v}_{i,t}} + \delta)} \quad i \in \{1, \dots, d\}$$

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \left( [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2}$$



# Convergence Proof ...

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \left( [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2}$$


... now, taking the conditional expectation with respect to the sample at iteration  $t$  given a fixed random variable  $\boldsymbol{\theta}_t$ :

$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \underbrace{\mathcal{L}(\boldsymbol{\theta}_t)}_{\text{Fully known}} - \eta_t \sum_{i=1}^d \left( \underbrace{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i}_{\text{Fully known}} \times \mathbb{E} \left[ \underbrace{\frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta}}_{\text{Dependent RVs}} \middle| \boldsymbol{\theta}_t \right] \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

# Proof Roadmap ...

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_t) + \dots +$$

$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) \dots + \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} \middle| \boldsymbol{\theta}_t \right]$$

... we need to bound these ...

$$\dots + \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

**Objective Func. L-Smoothness**

... relation between 2 successive iterations ...

**True Components to Bound**

... consider stochasticity plug-in update rule, and realise terms to bound ...

**Bounding the first term**

...

...

**Bound in terms of  
gradient norm**



# Convergence Proof ...

$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \left( [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} \middle| \boldsymbol{\theta}_t \right] \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$



**How to deal with  
such a ratio**


$$= \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \left( [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} + \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$



**... adding and subtracting will allow us to deal with this ...**

# Convergence Proof ...

$$= \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \left( [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \mathbb{E} \left[ \underbrace{\frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta}}_a - \underbrace{\frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta}}_b + \underbrace{\frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta}}_c \middle| \boldsymbol{\theta}_t \right] \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

$\mathbb{E}[a - b + c] = \mathbb{E}[a - b] + \mathbb{E}[c]$ 


$$\mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] + \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right]$$

$$\frac{\mathbb{E}[\mathbf{g}_{i,t} | \boldsymbol{\theta}_t]}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} = \frac{[\nabla \mathcal{L}(\boldsymbol{\theta})]_i}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta}$$

$$= \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \left( [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \left[ \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} + \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] \right] \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

# Convergence Proof ...

$$= \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \left( [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \left[ \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} + \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] \right] \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

Diagram illustrating the first step of the convergence proof. A hand points to the term  $\frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta}$ , which is highlighted in a grey box. A grey arrow with a cross indicates a subtraction operation. A green arrow with a cross indicates an addition operation. The term  $[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right]$  is highlighted in a green box.

$$= \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \left( \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} + [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

Diagram illustrating the second step of the convergence proof. A large green arrow points downwards, indicating the next step in the derivation.


$$= \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} - \eta_t \sum_{i=1}^d [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$



# Convergence Proof ...


$$= \left[ \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \right] - \eta_t \sum_{i=1}^d [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \times \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

$$\left[ -\eta_t \sum_{i=1}^d a_i b_i \right] \leq \left| \eta_t \sum_{i=1}^d a_i b_i \right| \leq \eta_t \sum_{i=1}^d |a_i| |b_i|$$



$$\leq \left[ \eta_t \sum_{i=1}^d |[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i| \left| \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] \right| \right]$$

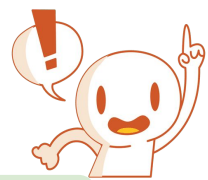
$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \left[ \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \right] + \left[ \eta_t \sum_{i=1}^d \left( |[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i| \left| \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] \right| \right) \right]$$



... our focus for now..


$$+ \left[ \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right] \right]$$

# Convergence Proof ...



$$|\mathbb{E}[x]| \leq \mathbb{E}[|x|]$$

$$\left| \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] \right| \leq \mathbb{E} \left[ \underbrace{\left| \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \right|}_{T_1} \middle| \boldsymbol{\theta}_t \right]$$

 ... our focus for now..



$$|\sqrt{a} - \sqrt{b}| = \frac{|a - b|}{\sqrt{a} + \sqrt{b}}$$

$$T_1 = \left| \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \right| = |\mathbf{g}_{i,t}| \underbrace{\left| \frac{1}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{1}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \right|}_{\times \text{ ... common denominator ..}} = \frac{|\mathbf{g}_{i,t}|}{(\sqrt{\mathbf{v}_{i,t}} + \delta)(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)} \overbrace{\left| \sqrt{\mathbf{v}_{i,t}} - \sqrt{\beta_2 \mathbf{v}_{i,t-1}} \right|}$$

$$= \frac{|\mathbf{g}_{i,t}|}{(\sqrt{\mathbf{v}_{i,t}} + \delta)(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)} \frac{|\mathbf{v}_{i,t} - \beta_2 \mathbf{v}_{i,t-1}|}{\sqrt{\mathbf{v}_{i,t}} + \sqrt{\beta_2 \mathbf{v}_{i,t-1}}}$$

# Convergence Proof ...

... update rule ...

$$\mathbf{v}_{i,t} = \beta_2 \mathbf{v}_{i,t-1} + (1 - \beta_2) \mathbf{g}_{i,t}^2$$

Remember  
Me?



... plug eq. in ...

$$\frac{|\mathbf{g}_{i,t}|}{(\sqrt{\mathbf{v}_{i,t}} + \delta)(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)} \frac{|\mathbf{v}_{i,t} - \beta_2 \mathbf{v}_{i,t-1}|}{\sqrt{\mathbf{v}_{i,t}} + \sqrt{\beta_2 \mathbf{v}_{i,t-1}}} = \frac{|\mathbf{g}_{i,t}|}{(\sqrt{\mathbf{v}_{i,t}} + \delta)(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)} \frac{(1 - \beta_2) \mathbf{g}_{i,t}^2}{\boxed{\sqrt{\mathbf{v}_{i,t}}} + \sqrt{\beta_2 \mathbf{v}_{i,t-1}}}$$



... plug eq. in ...

$$= \frac{|\mathbf{g}_{i,t}|}{(\sqrt{\mathbf{v}_{i,t}} + \delta)(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)} \frac{(1 - \beta_2) \mathbf{g}_{i,t}^2}{\boxed{\sqrt{\beta_2 \mathbf{v}_{i,t-1} + (1 - \beta_2) \mathbf{g}_{i,t}^2}} + \sqrt{\beta_2 \mathbf{v}_{i,t-1}}}$$

$$\mathbf{v}_{i,t} = \beta_2 \mathbf{v}_{i,t-1} + (1 - \beta_2) \mathbf{g}_{i,t}^2$$

Now what ...



# Convergence Proof ...



$$\frac{1}{a+b} \leq \frac{1}{a} \quad \text{for } a > 0 \text{ and } b \geq 0$$

$$= \frac{|g_{i,t}|}{(\sqrt{v_{i,t}} + \delta)(\sqrt{\beta_2 v_{i,t-1}} + \delta)} \frac{(1 - \beta_2)g_{i,t}^2}{\sqrt{\beta_2 v_{i,t-1}} + \underbrace{(1 - \beta_2)g_{i,t}^2}_{\text{non-negative}} + \sqrt{\beta_2 v_{i,t-1}}}$$

$$\leq \frac{|g_{i,t}|}{(\sqrt{v_{i,t}} + \delta)(\sqrt{\beta_2 v_{i,t-1}} + \delta)} \frac{(1 - \beta_2)g_{i,t}^2}{\underbrace{\sqrt{\beta_2 v_{i,t-1}}}_a + \underbrace{(1 - \beta_2)g_{i,t}^2}_b}$$



$$\sqrt{a+b} \geq \sqrt{b} \quad \text{if } a \geq 0 \Rightarrow \frac{1}{\sqrt{a+b}} \leq \frac{1}{\sqrt{b}}$$

$$\leq \frac{|g_{i,t}|}{(\sqrt{v_{i,t}} + \delta)(\sqrt{\beta_2 v_{i,t-1}} + \delta)} \frac{(1 - \beta_2)g_{i,t}^2}{\sqrt{(1 - \beta_2)g_{i,t}^2}}$$



... remember our focus ...



$$\mathbb{E} [\mathcal{L}(\theta_{t+1}) | \theta_t] \leq \dots + \eta_t \sum_{i=1}^d \left( \left| [\nabla \mathcal{L}(\theta_t)]_i \right| \mathbb{E} \left[ \left| \frac{g_{i,t}}{\sqrt{v_{i,t}} + \delta} - \frac{g_{i,t}}{\sqrt{\beta_2 v_{i,t-1}} + \delta} \right| \middle| \theta_t \right] \right) \dots$$

# Convergence Proof ...

$$\begin{aligned}
 T_1 &\leq \frac{|\mathbf{g}_{i,t}|}{(\underbrace{\sqrt{\mathbf{v}_{i,t}} + \delta}_b)(\underbrace{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta}_a)} \frac{(1 - \beta_2) \mathbf{g}_{i,t}^2}{\sqrt{(1 - \beta_2) \mathbf{g}_{i,t}^2}} \\
 &= \frac{1}{(\underbrace{\sqrt{\mathbf{v}_{i,t}} + \delta}_b)(\underbrace{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta}_a)} \sqrt{1 - \beta_2} \mathbf{g}_{i,t}
 \end{aligned}$$

...same trick...



$$\frac{1}{a+b} \leq \frac{1}{a} \text{ for } a > 0 \text{ and } b \geq 0$$



$$T_1 \leq \frac{\sqrt{1 - \beta_2} \mathbf{g}_{i,t}^2}{\delta(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)}$$



... now, we'll plug-back in  
the main bound ...

Remember  
Me?



... remember our focus ...



$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \dots + \eta_t \sum_{i=1}^d \left( \left| [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \right| \mathbb{E} \left[ \left| \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \right| \middle| \boldsymbol{\theta}_t \right] \right) \dots$$



# Plugging-Back in the main bound ...

$$\begin{aligned}\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] &\leq \dots + \eta_t \sum_{i=1}^d \left( |[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i| \left| \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] \right| \right) \dots \\ &\leq \dots + \eta_t \sum_{i=1}^d \left( |[\nabla \mathcal{L}(\boldsymbol{\theta})]_i| \underbrace{\mathbb{E} \left[ \left| \frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} - \frac{\mathbf{g}_{i,t}}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \right| \middle| \boldsymbol{\theta}_t \right]}_{T_1} \right) + \dots\end{aligned}$$

$$T_1 \leq \frac{\sqrt{1 - \beta_2} \mathbf{g}_{i,t}^2}{\delta(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)}$$



$$= \dots + \eta_t \sum_{i=1}^d (|[\nabla \mathcal{L}(\boldsymbol{\theta})]_i| \mathbb{E} [T_1 | \boldsymbol{\theta}_t]) + \dots = \dots + \eta_t \sum_{i=1}^d \left( |[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i| \frac{\sqrt{1 - \beta_2} \mathbb{E} [\mathbf{g}_{i,t}^2 | \boldsymbol{\theta}_t]}{\delta(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)} \right) + \dots$$

... hence, the overall bound ...

$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} + \eta_t \sum_{i=1}^d \left( |[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i| \frac{\sqrt{1 - \beta_2} \mathbb{E} [\mathbf{g}_{i,t}^2 | \boldsymbol{\theta}_t]}{\delta(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)} \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

# Bounding the gradient ...



$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} + \eta_t \sum_{i=1}^d \left( \boxed{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i} \frac{\sqrt{1 - \beta_2} \mathbb{E}[\mathbf{g}_{i,t}^2 | \boldsymbol{\theta}_t]}{\delta(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)} \right) + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

.... we can thus say ...

**Assumption II -- Loss functions has bounded gradient:**

$$\|\nabla \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{\xi})\| \leq G, \quad \forall \boldsymbol{\theta} \in \mathbb{R}^d, \quad \forall \boldsymbol{\xi}$$



$$\|\nabla \mathcal{L}(\boldsymbol{\theta})\| = \|\mathbb{E}_{\boldsymbol{\xi}} [\nabla \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{\xi})]\| \leq \mathbb{E}_{\boldsymbol{\xi}} [\|\nabla \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{\xi})\|] \leq G$$

$$\Rightarrow \boxed{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i \leq G}$$

e.g. infinity norm



$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} + \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

... now this ...

# Proof Roadmap ...

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_t) + \dots +$$

$$\mathbb{E}[\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) \dots + \mathbb{E}\left[\frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} \middle| \boldsymbol{\theta}_t\right]$$

... we need to bound these ...

$$\dots + \mathbb{E}\left[\frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t\right]$$

**Objective Func. L-Smoothness**

... relation between 2 successive iterations ...

**True Components to Bound**

... consider stochasticity plug-in update rule, and realise terms to bound ...

**Bounding the first term**

...

...

**Bound in terms of gradient norm**

**Choose hyper-params**

**Bounding the second term**



# Bounding the 3rd term ...

... update rule ...

$$\mathbf{v}_{i,t} = \beta_2 \mathbf{v}_{i,t-1} + (1 - \beta_2) \mathbf{g}_{i,t}^2$$

Remember Me?



... plug eq. in ...

$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \dots + \left[ \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right] \right] \Rightarrow \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{\left( \underbrace{\sqrt{\beta_2 \mathbf{v}_{i,t-1} + (1 - \beta_2) \mathbf{g}_{i,t}^2}}_{\text{non-negative}} + \delta \right)^2} \middle| \boldsymbol{\theta}_t \right]$$

$$\frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} + \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] + \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right]$$

Let's continue with the bound ...

$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \dots \quad \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] + \left[ \frac{L\eta_t^2}{2} \sum_{i=1}^d \mathbb{E} \left[ \frac{g_{i,t}^2}{(\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta)^2} \middle| \boldsymbol{\theta}_t \right] \right]$$

... same denominator ...

$$\leq \frac{L\eta_t^2}{2\delta} \sum_{i=1}^d \mathbb{E} \left[ \frac{g_{i,t}^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right]$$

➡

$$\left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \sum_{i=1}^d \mathbb{E} \left[ \frac{\mathbf{g}_{i,t}^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} \middle| \boldsymbol{\theta}_t \right] \leq \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \sum_{i=1}^d \frac{1}{\delta} \mathbb{E} [\mathbf{g}_{i,t}^2 | \boldsymbol{\theta}_t]$$

$$= \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \sum_{i=1}^d \mathbb{E} [\mathbf{g}_{i,t}^2 | \boldsymbol{\theta}_t]$$

Let's continue with the bound ...

$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \sum_{i=1}^d \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta} + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L \eta_t^2}{2\delta} \right) \mathbb{E} [\|\mathbf{g}_t\|^2 | \boldsymbol{\theta}_t]$$

$$\mathbf{v}_{i,t} \leq G^2 \quad \forall i, t \quad \Rightarrow \quad \sqrt{\beta_2 \mathbf{v}_{i,t-1}} + \delta \leq \sqrt{\beta_2} G + \delta \quad \Rightarrow \quad -\eta_t \sum_{i=1}^d \frac{[\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2}{\sqrt{\beta_2 \mathbf{v}_{i,t}} + \delta} \leq -\frac{\eta_t}{\sqrt{\beta_2} G + \delta} \sum_{i=1}^d [\nabla \mathcal{L}(\boldsymbol{\theta}_t)]_i^2 = -\frac{\eta_t}{\sqrt{\beta_2} G + \delta} \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2$$



$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \frac{\eta_t}{\sqrt{\beta_2} G + \delta} \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L \eta_t^2}{2\delta} \right) \mathbb{E} [\|\mathbf{g}_t\|^2 | \boldsymbol{\theta}_t]$$



—  $\Delta$  —

Some non-negative term



Let's continue with the bound ...



$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \frac{\eta_t}{\sqrt{\beta_2}G + \delta} \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \mathbb{E} [\|\mathbf{g}_t\|^2 | \boldsymbol{\theta}_t]$$

**Assumption III -- Variance of Loss is Bounded:**

$$\mathbb{E}_{\xi} [\|\nabla \mathcal{L}(\boldsymbol{\theta}; \xi) - \nabla \mathcal{L}(\boldsymbol{\theta})\|_2^2] \leq \sigma^2, \quad \forall \boldsymbol{\theta} \in \mathbb{R}^d, \quad \forall \xi$$



... if we use a mini-batch, we can write...

$$\mathbf{g}_t(\cdot) = \frac{1}{b_t} \sum_{\xi \in \mathcal{B}_t} \nabla \mathcal{L}(\cdot; \xi)$$



... then, we can prove ..

$$\mathbb{E} [\|\mathbf{g}_t\|_2^2 | \boldsymbol{\theta}_t] \leq \frac{1}{b_t} \left( \sigma^2 + \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right)$$

# Proof Roadmap ...

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_t) + \dots +$$

$$\mathbb{E}[\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) \dots + \mathbb{E}\left[\frac{\mathbf{g}_{i,t}}{\sqrt{\mathbf{v}_{i,t}} + \delta} \middle| \boldsymbol{\theta}_t\right]$$

... we need to bound these ...

$$\dots + \mathbb{E}\left[\frac{\mathbf{g}_{i,t}^2}{(\sqrt{\mathbf{v}_{i,t}} + \delta)^2} \middle| \boldsymbol{\theta}_t\right]$$

**Objective Func. L-Smoothness**

... relation between 2 successive iterations ...

**True Components to Bound**

... consider stochasticity plug-in update rule, and realise terms to bound ...

**Bounding the first term**

**Choose hyper-params**

**Bounding the second term**

...

...

...

...

**Bound in terms of gradient norm**

**Bound in terms of batch-size**



Therefore, we can write ...

$$\begin{aligned}\mathbb{E}[\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] &\leq \mathcal{L}(\boldsymbol{\theta}_t) - \frac{\eta_t}{\sqrt{\beta_2}G + \delta} \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} (\sigma^2 + \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2) \\ &= \mathcal{L}(\boldsymbol{\theta}_t) - \underbrace{\|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \left( \frac{\eta_t}{\sqrt{\beta_2}G + \delta} - \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \right)}_{\text{... has to be a constant ...}} + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \sigma^2\end{aligned}$$

... now, we need to handle each of these constants ...

$$\mathbb{E}[\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}) - \underbrace{\Delta}_{\text{... we want this to go to zero ...}} + \underbrace{\epsilon_t}_{\text{... has to be a constant ...}} + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \sigma^2$$

... let's start choosing free hyper-parameters (e.g., batch-sizes, learning rates ...) to get what we want ...



# Let's choose free hyper-parameters ...

## We'll make 3 choices:

1. Batch size:  $b_t$
2. Learning rate:  $\eta_t$
3. Free parameter:  $\beta_2$

$$\begin{aligned}\mathbb{E}[\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] &\leq \mathcal{L}(\boldsymbol{\theta}_t) - \frac{\eta_t}{\sqrt{\beta_2}G + \delta} \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \left( \sigma^2 + \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right) \\ &= \mathcal{L}(\boldsymbol{\theta}_t) - \underbrace{\|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \left( \frac{\eta_t}{\sqrt{\beta_2}G + \delta} - \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \right)}_{\text{... let's start with ... } \mathcal{A}} + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \sigma^2\end{aligned}$$

Choose  $b_t \geq 1$ , then we can say that:

$$\begin{aligned}\frac{1}{\delta b_t} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) &\leq \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \implies -\frac{1}{\delta b_t} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \geq -\frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \\ &\implies \mathcal{A} \geq \underbrace{\eta_t \left[ \frac{1}{\sqrt{\beta_2}G + \delta} - \frac{1}{\delta} \left( \frac{G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t}{2\delta} \right) \right]}_{\text{... let's call this ... } \mathcal{B}}\end{aligned}$$

# Let's choose free hyper-parameters ...

## We'll make 3 choices:

1. Batch size:  $b_t$
2. Learning rate:  $\eta_t$
3. Free parameter:  $\beta_2$

Choose  $b_t \geq 1$ , then we can say that:

$$\begin{aligned} \frac{1}{\delta b_t} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L \eta_t^2}{2\delta} \right) &\leq \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L \eta_t^2}{2\delta} \right) \implies -\frac{1}{\delta b_t} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L \eta_t^2}{2\delta} \right) \geq -\frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L \eta_t^2}{2\delta} \right) \\ &\implies \mathcal{A} \geq \eta_t \underbrace{\left[ \frac{1}{\sqrt{\beta_2} G + \delta} - \frac{1}{\delta} \left( \frac{G \sqrt{1 - \beta_2}}{\delta} + \frac{L \eta_t}{2\delta} \right) \right]}_{\text{... let's call this ... } \mathcal{B}} \end{aligned}$$

Choose  $\eta_t = \eta$ , such that  $\frac{L\eta}{2\delta} \leq \frac{G\sqrt{1 - \beta_2}}{\delta}$ , i.e.,  $\eta \leq \frac{2G\sqrt{1 - \beta_2}}{L}$ :

$$\frac{G\sqrt{1 - \beta_2}}{\delta} + \frac{L\eta}{2\delta} \leq \frac{2G\sqrt{1 - \beta_2}}{\delta} \implies -\left( \frac{G\sqrt{1 - \beta_2}}{\delta} + \frac{L\eta}{2\delta} \right) \geq -\frac{2G\sqrt{1 - \beta_2}}{\delta} \implies \mathcal{B} \geq \frac{1}{\sqrt{\beta_2} G + \delta} - \frac{2G\sqrt{1 - \beta_2}}{\delta^2}$$

Let's choose free hyper-parameters  $\frac{1}{G + \delta} \leq \frac{1}{\sqrt{\beta_2}G + \delta}$

**We'll make 3 choices:**

1. Batch size:  $b_t$
2. Learning rate:  $\eta_t$
3. Free parameter:  $\beta_2$

Further, choose  $\beta_2$  such that  $\frac{2G\sqrt{1-\beta_2}}{\delta^2} \leq \frac{1}{2} \left( \frac{1}{\sqrt{\beta_2}G + \delta} \right)$ , then:

---

Let us choose  $\beta_2$  such that:  $\frac{2G\sqrt{1-\beta_2}}{\delta^2} = \frac{1}{2} \frac{1}{(G + \delta)}$ , then:  $\beta_2 = 1 - \frac{\delta^4}{16G^2(G + \delta)}$

... should be close to one!

$$\Rightarrow \mathcal{B} \geq \frac{1}{2(\sqrt{\beta_2}G + \delta)} \Rightarrow \mathcal{A} \geq \eta \mathcal{B} \geq \frac{\eta}{2(\sqrt{\beta_2}G + \delta)} \Rightarrow -\mathcal{A} \leq -\frac{\eta}{2(\sqrt{\beta_2}G + \delta)}$$

$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \left( \frac{\eta_t}{\sqrt{\beta_2}G + \delta} - \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1-\beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \right) + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1-\beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \sigma^2$$

... let's call this ...  $\mathcal{C}$



# Let's choose free hyper-parameters ...

## We'll make 3 choices:

1. Batch size:  $b_t$
2. Learning rate:  $\eta_t$
3. Free parameter:  $\beta_2$


$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \left( \frac{\eta_t}{\sqrt{\beta_2}G + \delta} - \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \right) + \frac{1}{\delta} \left( \frac{\eta_t G \sqrt{1 - \beta_2}}{\delta} + \frac{L\eta_t^2}{2\delta} \right) \frac{1}{b_t} \sigma^2$$

■————■  
... let's call this ...  $\mathcal{C}$

Note, we chose  $\eta_t = \eta$  such that  $\frac{L\eta}{2\delta} \leq \frac{G\sqrt{1 - \beta_2}}{\delta}$  :

---


... then, we can say that  $\mathcal{C} \leq 2\eta \frac{G\sqrt{1 - \beta_2}}{\delta}$



$$\mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_{t+1}) | \boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \frac{\eta}{2(\sqrt{\beta_2}G + \delta)} + \frac{2\eta\sigma^2}{\delta^2 b_t} G \sqrt{1 - \beta_2}$$


Let's finalise the bound ...

$$\mathbb{E}[\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] \leq \mathcal{L}(\boldsymbol{\theta}_t) - \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \frac{\eta}{2(\sqrt{\beta_2}G + \delta)} + \frac{2\eta\sigma^2}{\delta^2 b_t} G \sqrt{1 - \beta_2}$$

$$\implies \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \frac{\eta}{2(\sqrt{\beta_2}G + \delta)} \leq \mathcal{L}(\boldsymbol{\theta}_t) - \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}_{t+1})|\boldsymbol{\theta}_t] + \frac{2\eta\sigma^2}{\delta^2 b_t} G \sqrt{1 - \beta_2}$$


$$\frac{1}{2(\sqrt{\beta_2}G + \delta)} \mathbb{E}_{\text{total}} \left[ \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right] \leq \frac{\mathbb{E}_{\text{total}}[\mathcal{L}(\boldsymbol{\theta}_t)] - \mathbb{E}_{\text{total}}[\mathcal{L}(\boldsymbol{\theta}_{t+1})]}{\eta} + \frac{2\sigma^2}{\delta_2 b_t} G \sqrt{1 - \beta_2}$$


$$\frac{1}{2(\sqrt{\beta_2}G + \delta)} \sum_{t=1}^T \mathbb{E}_{\text{total}} \left[ \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right] \leq \frac{\mathbb{E}_{\text{total}}[\mathcal{L}(\boldsymbol{\theta}_1)] - \mathbb{E}_{\text{total}}[\mathcal{L}(\boldsymbol{\theta}_{T+1})]}{\eta} + \frac{2\sigma^2}{\delta_2} G \sqrt{1 - \beta_2} \sum_{t=1}^T \frac{1}{b_t}$$


$$\frac{c_1}{T} \sum_{t=1}^T \mathbb{E}_{\text{total}} \left[ \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right] \leq \frac{\mathcal{L}(\boldsymbol{\theta}_1) - \mathcal{L}(\boldsymbol{\theta}_{\text{global-min}})}{T\eta} + c_2 \frac{1}{T} \sum_{t=1}^T \frac{1}{b_t}$$

# Let's finalise the bound ...

$$\frac{c_1}{T} \sum_{t=1}^T \mathbb{E}_{\text{total}} \left[ \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right] \leq \underbrace{\frac{\mathcal{L}(\boldsymbol{\theta}_1) - \mathcal{L}(\boldsymbol{\theta}_{\text{global-min}})}{T\eta} + c_2 \frac{1}{T} \sum_{t=1}^T \frac{1}{b_t}}_{\text{we want the RHS to be } \leq \epsilon c_1} \implies \frac{\mathcal{L}(\boldsymbol{\theta}_1) - \mathcal{L}(\boldsymbol{\theta}_{\text{global-min}})}{T\eta} + c_2 \frac{1}{T} \sum_{t=1}^T \frac{1}{b_t} \leq c_1 \epsilon$$

... with a constant batch-size ...

$$\left. \begin{aligned} b_t = b &\implies b = \lceil \frac{2c_2}{c_1\epsilon} \rceil \implies c_2 \frac{1}{T} \sum_{t=1}^T \frac{1}{b_t} = \frac{c_2}{b} \leq \frac{c_1\epsilon}{2} \\ T = \frac{2(\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}_{\text{global-min}}))}{\eta c_1 \epsilon} &\implies \frac{\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}_{\text{global-min}})}{T\eta} \leq \frac{c_1\epsilon}{2} \end{aligned} \right\} \implies \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\text{total}} \left[ \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right] \leq \epsilon$$

... but as T grows ...

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\text{total}} \left[ \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right] = \frac{c_2}{c_1 b} \neq 0 \quad \dots \text{we don't converge to a stationary point ...}$$

... how to fix that...



# Let's finalise the bound ...

$$\frac{c_1}{T} \sum_{t=1}^T \mathbb{E}_{\text{total}} \left[ \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right] \leq \underbrace{\frac{\mathcal{L}(\boldsymbol{\theta}_1) - \mathcal{L}(\boldsymbol{\theta}_{\text{global-min}})}{T\eta} + c_2 \frac{1}{T} \sum_{t=1}^T \frac{1}{b_t}}_{\text{we want the RHS to be } \leq \epsilon c_1} \implies \frac{\mathcal{L}(\boldsymbol{\theta}_1) - \mathcal{L}(\boldsymbol{\theta}_{\text{global-min}})}{T\eta} + c_2 \frac{1}{T} \sum_{t=1}^T \frac{1}{b_t} \leq c_1 \epsilon$$

... with an increasing batch-size ...      ... chose T such that...  $\frac{\ln T + \gamma}{T} \leq \frac{\epsilon}{2}$

$$\left. \begin{aligned} b_t = \lceil \frac{c_2}{c_1} \rceil t &\implies c_2 \frac{1}{T} \sum_{t=1}^T \frac{1}{b_t} \leq \frac{c_1}{T} \sum_{t=1}^T \frac{1}{t} = \frac{c_1}{T} (\ln T + \gamma) \leq \frac{c_1 \epsilon}{2} \\ T = \frac{2(\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}_{\text{global-min}}))}{\eta c_1 \epsilon} &\implies \frac{\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}_{\text{global-min}})}{T\eta} \leq \frac{c_1 \epsilon}{2} \end{aligned} \right\} \implies \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\text{total}} \left[ \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right] \leq \epsilon$$

... and as T grows ...

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\text{total}} \left[ \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|_2^2 \right] = 0$$

... we converge to a stationary point ...

# Summary

## Key idea

Partition the loss function, bound each term, select the hyper-parameters

## Assumptions

- subsampled loss function is  $L$ -smooth
- the gradient is upper bounded by a number  $G$
- the variance is bounded by  $\sigma^2$

## Steps

1. start from the property of  $L$ -smoothness and replace the equation
2. focus a ratio of random variables
3. add subtracted terms then we bound absolute values
4. plug back in and introduce assumptions
5. bound the third term

# References

1. Kingma, Diederik P., and Jimmy Ba. "Adam: A Method for Stochastic Optimization." In ICLR (Poster). 2015.
2. Reddi, Sashank J., Satyen Kale, and Sanjiv Kumar. "On the convergence of adam and beyond." arXiv preprint arXiv:1904.09237 (2019).
3. Reddi, S., Manzil Zaheer, Devendra Sachan, Satyen Kale, and Sanjiv Kumar. "Adaptive methods for nonconvex optimization." In Proceeding of 32nd Conference on Neural Information Processing Systems (NIPS 2018). 2018.
4. Sun, Ruoyu. "Optimization for deep learning: theory and algorithms." arXiv preprint arXiv:1912.08957 (2019).
5. <https://www.deeplearningbook.org/contents/optimization.html>
6. Smith, Samuel L., Pieter-Jan Kindermans, Chris Ying, and Quoc V. Le. "Don't decay the learning rate, increase the batch size." arXiv preprint arXiv:1711.00489 (2017).



# The Goal of This Lecture

- This example (Adam Convergence) is not perfect
- Hopefully you find a way to prove the convergence of your algorithms
  - In most of situations, you can't make a perfect proof
  - You need some assumptions