# GT Deep Learning: 03

## Vincent Perrollaz

### Institut Denis Poisson

## 17 Février 2026

## Summary

Hyperparameters  network architecture, layer sizes, activation functions

Human Decision, Methodology to come.

Parameters  previous session $\implies$ belong to

$$\arg\min_{p\in\mathbb{R}^D} J(p). \tag{1}$$

Remark  first $J \in \mathcal{C}^1$, later specific form

Existence  AI literature $\cap$ Perron's Paradox $= \emptyset$

## Computation

Equation Classical calculus $\implies$ compute

$$\{p \in \mathbb{R}^D \, : \, \mathrm{d}J(p) = 0\}. \qquad (2)$$

Successive Approximation Construct

$$\text{recursively } (p_n)_{n \geq 0} \in (\mathbb{R}^D)^{\mathbb{N}}$$
$$n \to +\infty \implies J(p_n) \to \inf J. \quad (3)$$

$$\boxed{\text{Main Idea}}$$

$$\begin{aligned}
\text{Start} \quad & \text{No information} \implies p_0 \sim \mathcal{N}(0,1) \\
\text{Analysis} \quad & p_{n+1} = p_n + (p_{n+1} - p_n) \\
.\text{Expansion} \quad & J(p_{n+1}) \approx J(p_n) + (p_{n+1} - p_n)J'(p_n) \\
\text{Sign condition} \quad & \mathrm{sign}(p_{n+1} - p_n) = -\mathrm{sign}(J'(p_n)) \\
\text{Size condition} \quad & p_{n+1} - p_n = -\eta J'(p_n) \\
\text{Parameter} \quad & \eta \text{ learning rate, small}
\end{aligned}$$

# ODE Link

Rewriting $\dfrac{p_{n+1} - p_n}{\eta} = -\nabla J'(p_n)$

Limit $\eta \to 0 \implies \dot{p}(t) = -J'(p(t))$

Lyapunov $\dfrac{\mathrm{d}}{\mathrm{d}t} J(p(t)) = -\left(J'(p(t))\right)^2$

Necessity $p(t) \to p^* \implies J'(p^*) = 0$

Sufficiency $J$ convex $\implies p^* \in \arg\min J.$

$$\boxed{\text{Gradient Descent}}$$

$$\begin{cases} p_{n+1} = p_n - \eta \nabla J(p_n), & \forall n \geq 0 \\ p_0 = \mathcal{N}(0,1)(\omega). \end{cases} \tag{4}$$

$\boxed{\text{Łojasiewicz inequality}}$

1. $J$ $\mu$-PL when $\exists p^*$

   $$\forall p \in \mathbb{R}^D, \quad \|\nabla J(p)\|^2 \geq 2\mu(J(p) - J(p^*))$$

2. $J$ uniformly convex $\implies$ $J$ $\mu$-PL
3. Locally true for real analytic functions!
4. ODE $\implies$

   $$\frac{\mathrm{d}}{\mathrm{d}t} J(p(t)) \leq -2\mu(J(p(t)) - J(p^*)) \qquad (5)$$

5. Gronwall $\implies$

   $$J(p(t)) - J(p^*) \leq (J(p_0) - J(p^*))e^{-2\mu t} \qquad (6)$$

$$\boxed{\text{Momentum}}$$

$$\begin{cases} m_{n+1} = \beta m_n + (1 - \beta)\nabla J(p_n) \\ p_{n+1} = p_n - \eta m_{n+1} \\ p_0 = \mathcal{N}(0, 1)(\omega) \\ m_0 = 0. \end{cases} \tag{7}$$

$$\boxed{\text{RMSProp}}$$

$$\begin{cases} s_{n+1} = \beta s_n + (1-\beta)\|\nabla J(p_n)\|^2 \\ p_{n+1} = p_n - \eta \frac{\nabla J(p_n)}{\epsilon + \sqrt{s_{n+1}}} \\ p_0 = \mathcal{N}(0,1)(\omega) \\ s_0 = 0. \end{cases} \qquad (8)$$

$$\boxed{\text{Adam}}$$

$$\begin{cases} m_{n+1} = \beta_1 m_n + (1 - \beta_1)\nabla J(p_n) \\ s_{n+1} = \beta_2 s_n + (1 - \beta_2)\|\nabla J(p_n)\|^2 \\ \hat{m}_{n+1} = \frac{m_{n+1}}{1 - \beta_1^{n+1}} \\ \hat{s}_{n+1} = \frac{s_{n+1}}{1 - \beta_2^{n+1}} \\ p_{n+1} = p_n - \eta \frac{\hat{m}_{n+1}}{\epsilon + \sqrt{\hat{s}_{n+1}}} \\ p_0 = \mathcal{N}(0, 1)(\omega) \\ m_0 = 0 \\ s_0 = 0. \end{cases} \tag{9}$$

## Remarks

- Solvers $\implies$ additional hyperparameters
- e.g. Adam: $\beta_1$, $\beta_2$, $\eta$, $\epsilon$.
- ⚠ Theoretical Garantees ⚠
- When do we stop?

$$\boxed{\text{Loss Function}}$$

$$\forall p \in R^D, \quad J(p) := \frac{1}{N} \sum_{0 \le k < N} J_k(p). \qquad (10)$$

1. Specific Form to Machine Learning
2. $N$ : sample size, very large (internet size)
3. $D$ : parameter size $10^9$ to $10^{12}$.
4. Classical Gradient Descent expensive steps!

$$\boxed{\text{Stochastic Gradient Descent}}$$

1. $N = B \cdot S$
2. $E$ : number of epochs
3. $1 \leq e \leq E$ : $\sigma_e$ random in $\mathbb{S}(N)$
4. $0 \leq s < S$ : loss function

$$J_{s,e}(p) := \frac{1}{B} \sum_{sB \leq k < (s+1)B} J_{\sigma_e(k)}(p). \qquad (11)$$

5. One step of some gradient algorithm for $J_{s,e}$.