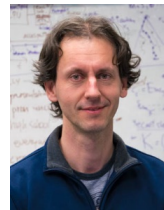


# Characterization of Convex Objective Functions and Optimal Expected Convergence Rates of SGD

---

**Phuong Ha Nguyen<sup>1</sup>**

Marten van Dijk<sup>1</sup>, Lam M. Nguyen<sup>2</sup> and Dzung T. Phan<sup>2</sup>



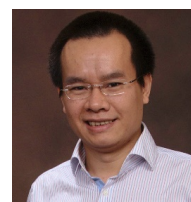
Marten



Lam



P. Ha



Dzung



1. Secure Computation Laboratory, ECE, University of Connecticut
2. IBM Research, Thomas J. Watson Research Center



International Conference on Machine Learning (ICML)  
Long Beach, California, 2019

# Problem Setting

---

- Solve

$$\min_{w \in \mathbb{R}^d} \{F(w) = E_{\xi}[f(w; \xi)]\}$$



- **Solve**

$$\min_{w \in \mathbb{R}^d} \{F(w) = E_{\xi}[f(w; \xi)]\}$$

- **Assumptions**

- **Convex:**

$$f(w; \xi) - f(w'; \xi) \geq \langle \nabla f(w'; \xi), (w - w') \rangle$$

- **Smooth:**

$$\|\nabla f(w; \xi) - \nabla f(w'; \xi)\| \leq L\|w - w'\|$$



- **Solve**

$$\min_{w \in R^d} \{F(w) = E_{\xi}[f(w; \xi)]\}$$

- **Assumptions**

- **Convex:**

$$f(w; \xi) - f(w'; \xi) \geq \langle \nabla f(w'; \xi), (w - w') \rangle$$

- **Smooth:**

$$\|\nabla f(w; \xi) - \nabla f(w'; \xi)\| \leq L\|w - w'\|$$

- **Find a  $w_t$  close to**

$$W^* = \{w_* \in R^d : \forall_{w \in R^d}, F(w) \geq F(w_*)\}$$



# Problem Setting

- **Solve**

$$\min_{w \in R^d} \{F(w) = E_{\xi}[f(w; \xi)]\}$$

- **Assumptions**

- **Convex:**

$$f(w; \xi) - f(w'; \xi) \geq \langle \nabla f(w'; \xi), (w - w') \rangle$$

- **Smooth:**

$$\|\nabla f(w; \xi) - \nabla f(w'; \xi)\| \leq L\|w - w'\|$$

- **Find a  $w_t$  close to**

$$W^* = \{w_* \in R^d : \forall w \in R^d, F(w) \geq F(w_*)\}$$

## Stochastic Gradient Descent (SGD):

**Initialize:**  $w_0$

**Iterate:**

**for**  $t = 0, 1, 2, \dots$ , **do**

  Choose  $\eta_t > 0$

  Generate random  $\xi_t$

  Compute  $\nabla f(w_t; \xi_t)$

  Update  $w_{t+1} = w_t - \eta_t \nabla f(w_t; \xi_t)$

**end for**



# Problem Setting

- **Solve**

$$\min_{w \in R^d} \{F(w) = E_{\xi}[f(w; \xi)]\}$$

- **Assumptions**

- **Convex:**

$$f(w; \xi) - f(w'; \xi) \geq \langle \nabla f(w'; \xi), (w - w') \rangle$$

- **Smooth:**

$$\|\nabla f(w; \xi) - \nabla f(w'; \xi)\| \leq L\|w - w'\|$$

- **Find a  $w_t$  close to**

$$W^* = \{w_* \in R^d : \forall_{w \in R^d}, F(w) \geq F(w_*)\}$$

- **Problem: Characterize Expected Convergence Rates**

$$E \left[ \inf_{w_* \in W^*} \|w_t - w_*\|^2 \right] \text{ and } E[F(w_t) - F(w_*)]$$

**Stochastic Gradient Descent (SGD):**

**Initialize:**  $w_0$

**Iterate:**

**for**  $t = 0, 1, 2, \dots$ , **do**

    Choose  $\eta_t > 0$

    Generate random  $\xi_t$

    Compute  $\nabla f(w_t; \xi_t)$

    Update  $w_{t+1} = w_t - \eta_t \nabla f(w_t; \xi_t)$

**end for**

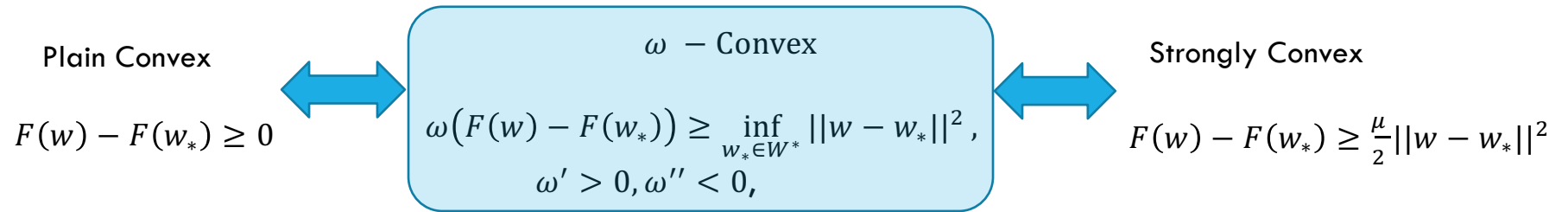


# UCONN Beyond convex and strongly convex functions

---

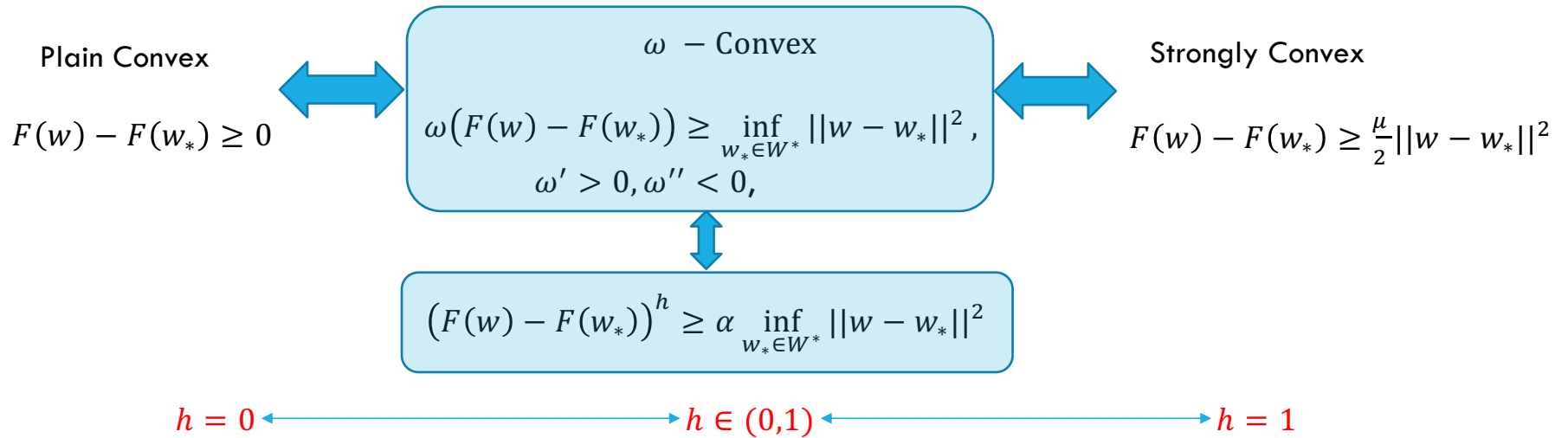


# $\omega$ -Convexity

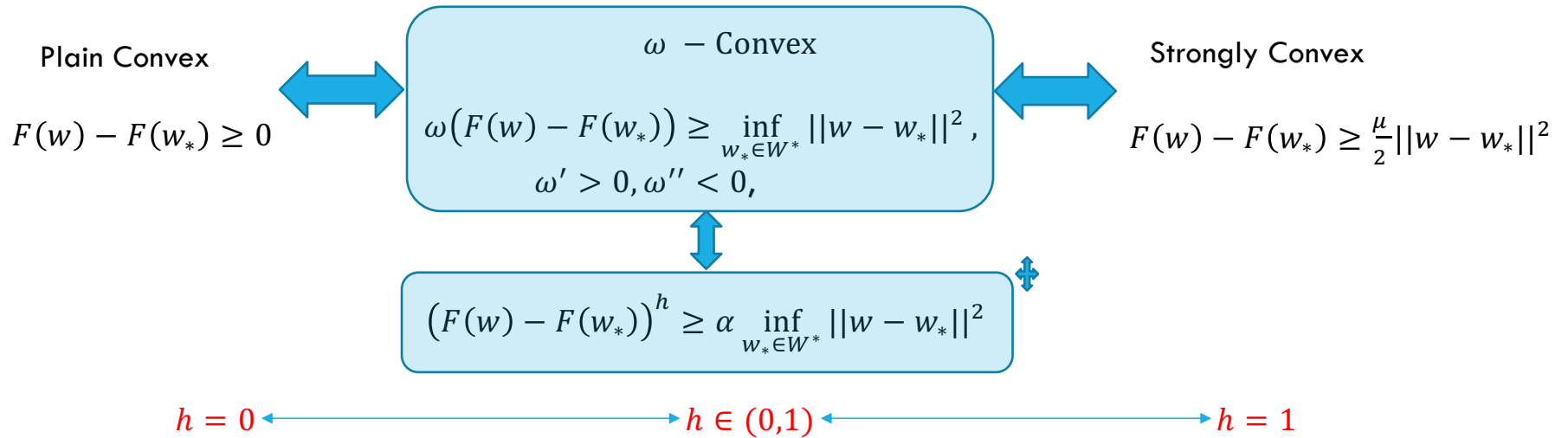




# $\omega$ -Convexity with curvature $h \in [0,1]$



# $\omega$ -Convexity with curvature $h \in [0,1]$



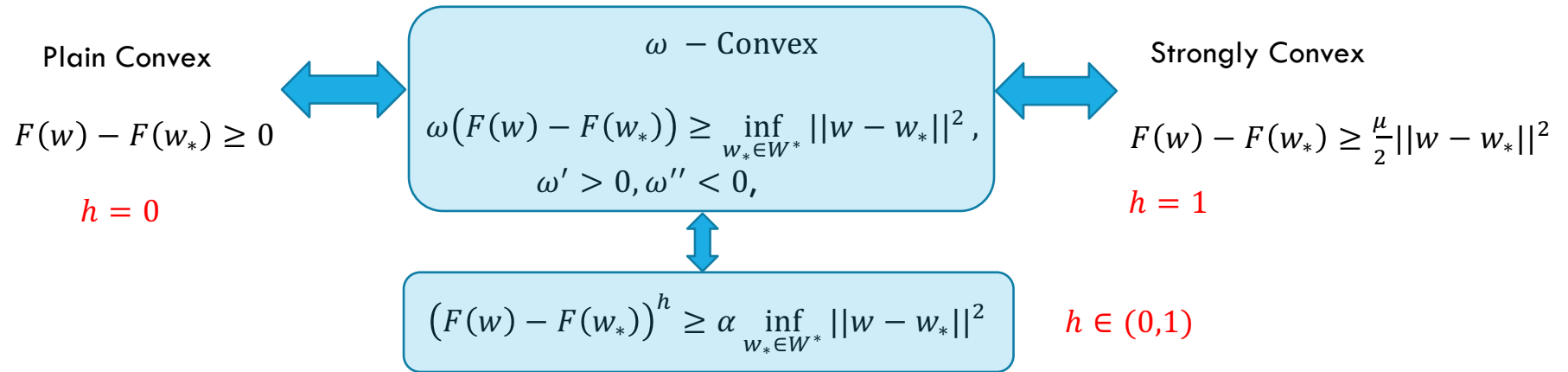
⚡ HEB (Holderian Error Bound):  $(F(w) - F(w_*))^h \geq \alpha \inf_{w_* \in W^*} \|w - w_*\|^2,$  where  $h \in (0,2]$ .

HEB and  $\omega$ -convexity are not subclasses of one another but they do intersection for  $h \in (0,1]$ .

[Bolte, J., Nguyen, T. P., Peypouquet, J., and Suter, B. W. From error bounds to the complexity of first order descent methods for convex functions. *Mathematical Programming*, 165(2):471–507, Oct 2017]



# Close to optimal stepsize



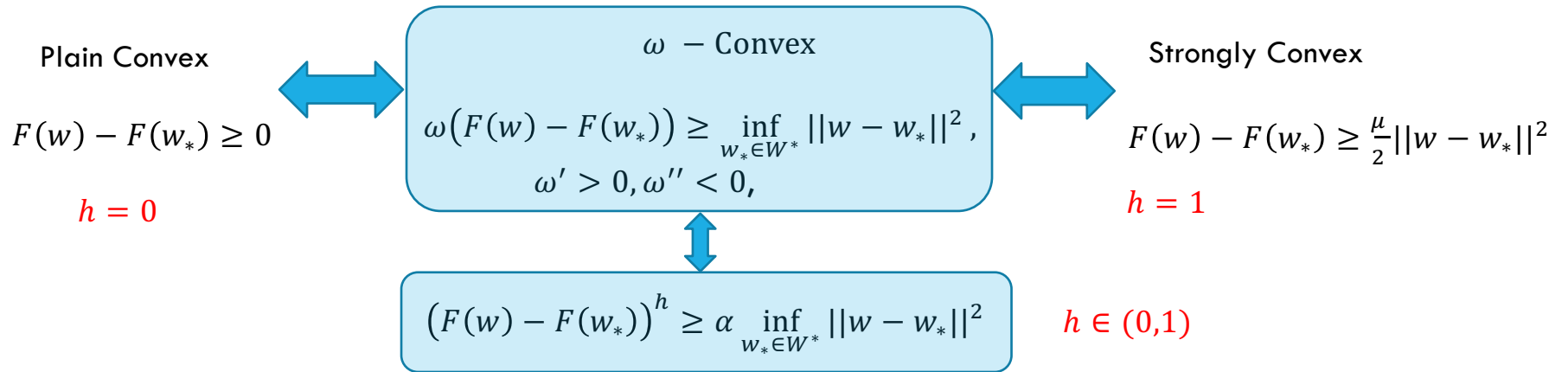
SGD

$$\eta_t = \frac{c}{(t+\Delta)^{1/(2-h)}}$$

Close to optimal stepsize



# Convergence Rate of SGD



SGD

$$\eta_t = \frac{c}{(t+\Delta)^{1/(2-h)}}$$

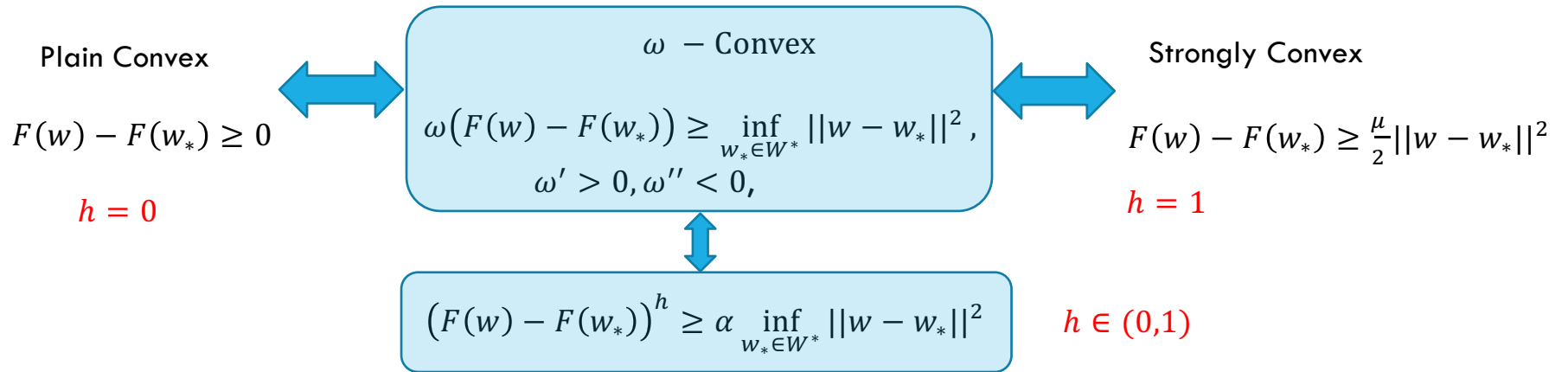
Close to optimal stepsize

$$E \left[ \inf_{w_* \in W^*} \|w_t - w_*\|^2 \right] = O(t^{-h/(2-h)})$$

$$\frac{1}{t} \sum_{i=t+1}^{2t} E[F(w_i) - F(w_*)] = O(t^{-1/(2-h)})$$



# Convergence Rate of SGD



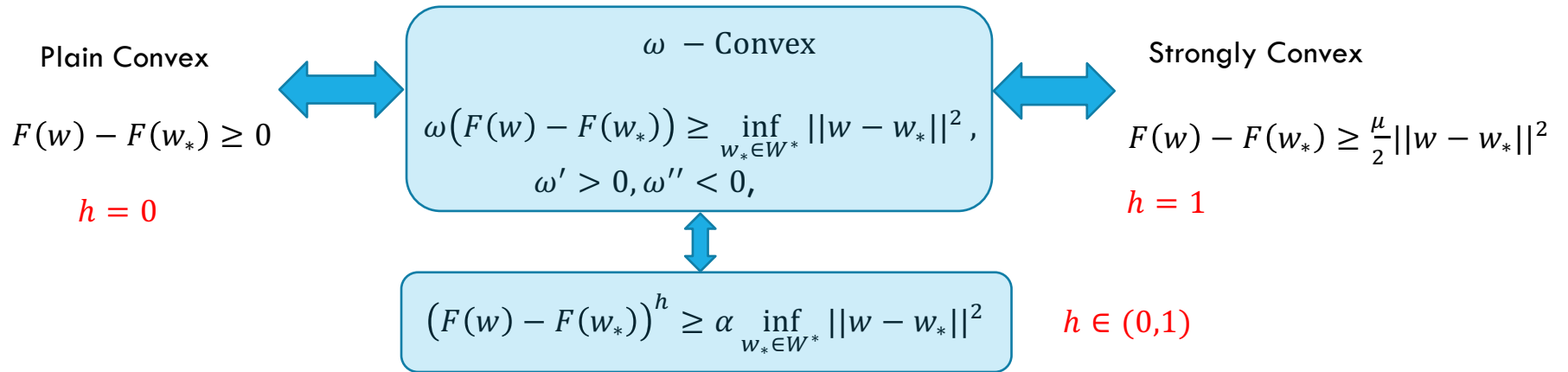
$$E \left[ \inf_{w_* \in W^*} \|w_t - w_*\|^2 \right] = O(t^{-h/(2-h)}) \quad [\text{Useless}, 0] \longleftrightarrow [\text{Useful}, 1]$$

$0 \leftarrow h \rightarrow 1$

$$\frac{1}{t} \sum_{i=t+1}^{2t} E[F(w_i) - F(w_*)] = O(t^{-1/(2-h)}) \quad [\text{Useful}, 0] \longleftrightarrow [\text{Useful}, 1]$$



# Convergence Rate of SGD



$$E \left[ \inf_{w_* \in W_*} \|w_t - w_*\|^2 \right] = O(t^{-h/(2-h)})$$

$$\frac{1}{t} \sum_{i=t+1}^{2t} E[F(w_i) - F(w_*)] = O(t^{-1/(2-h)})$$

$h = 1/2$

$$F(w) = H(w) + \lambda G(w), H(w) - \text{convex}$$

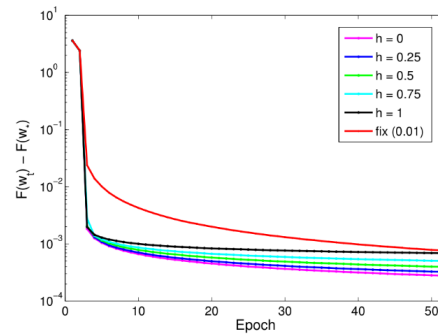
$$G(w) = \sum_{i=1}^d [e^{w_i} + e^{-w_i} - 2 - w_i^2]$$



# Experiment

Curvature 0 (convex)

$$f_i(w) = \log(1 + \exp(-y_i x_i^T w))$$

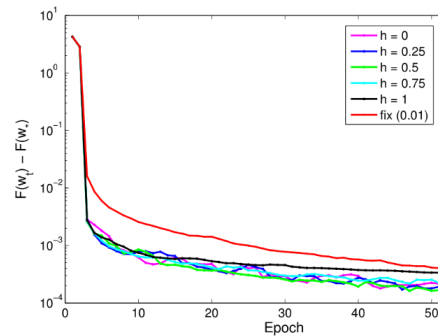


(a)

Curvature 1/2

$$f_i^a(w) = f_i(w) + \lambda G(w)$$

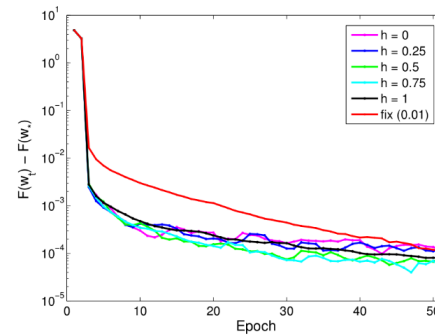
$$G(w) = \sum_{i=1}^d [e^{w_i} + e^{-w_i} - 2 - w_i^2]$$



(c)

Curvature unknown

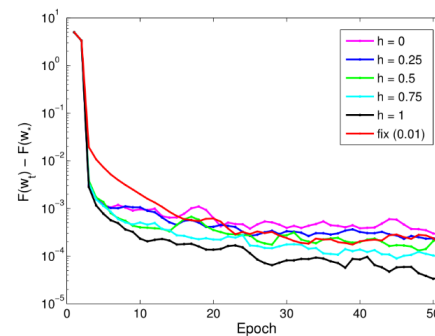
$$f_i^a(w) = f_i(w) + \lambda ||w||$$



(b)

Curvature 1 (strongly convex)

$$f_i^c(w) = f_i(w) + \frac{\lambda}{2} ||w||^2$$



(d)



# Conclusion

---

- $\omega$ -convexity notion: plain convex, strongly convex and something in between
- SGD with  $\omega$ -convex objective functions

**Thank you for your attention! 😊**

<https://arxiv.org/abs/1810.04100>

Poster Number: #193 – Pacific Ballroom. – 06:30—09:00PM – 06/11

