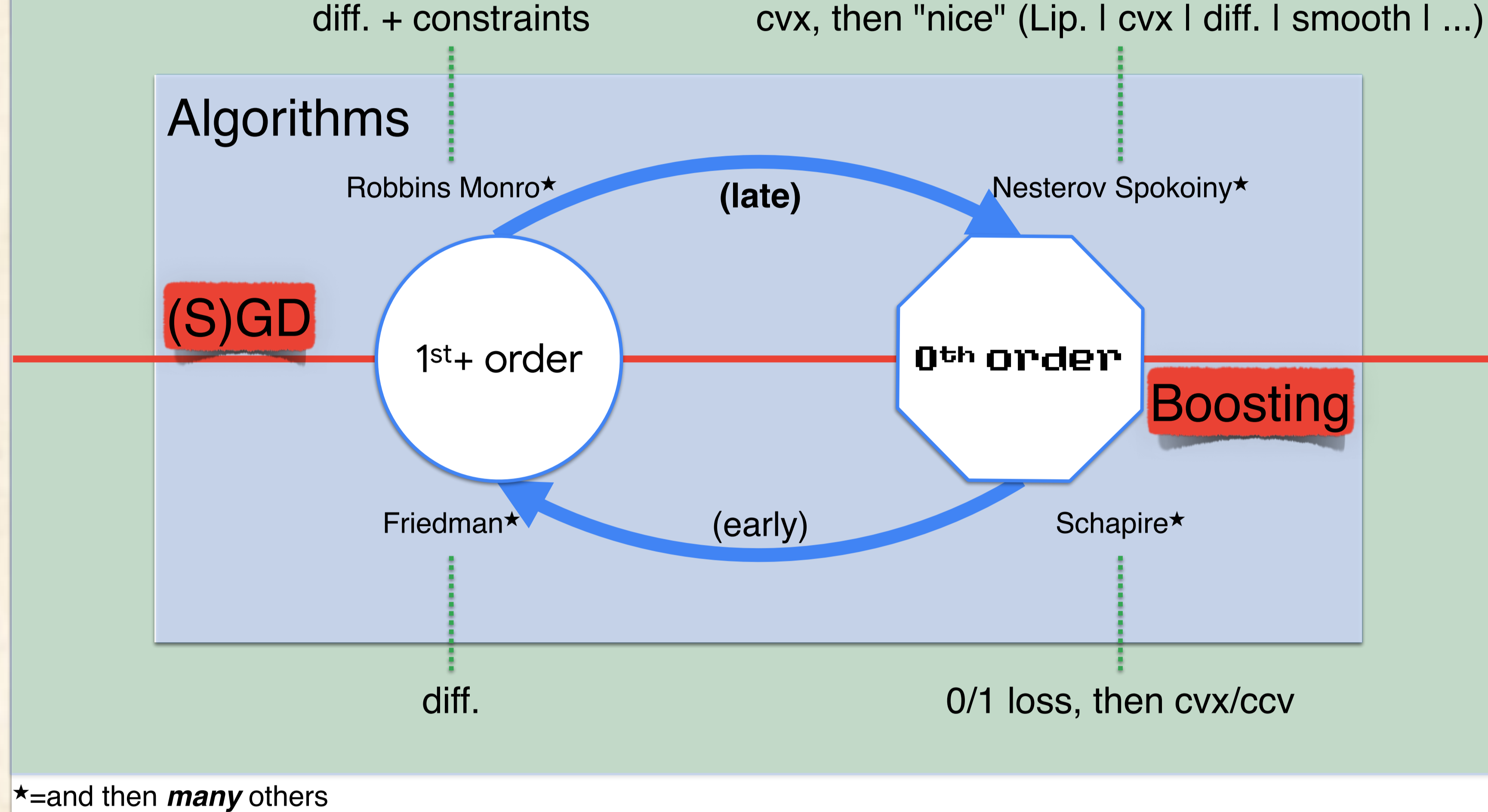


Summary

- Recent evolution of **(S)GD** → 0th order opt.: gradient-free, only loss queries (not any loss: they need to be somehow nice)

Losses



- Boosting** is gradient-free by design (Kearns/Valiant) w/ "nasty" 0/1 loss and then evolved → 1st order opt. w/ differentiable losses

- We question the power of the original 0th order framework*: *what losses can it directly optimize under the weak learning assumption?*
- Answer:** any loss whose set of discontinuities has 0 Lebesgue measure - computer-wise, this means **any loss** + our technique is constructive: *we give an algorithm*

*=analysis of boosting-compliant convergence on training, since generalization entails restrictions on losses w/ SOTA toolbox (no different from (S)GD → 0th order's mainstream analysis)

- (S)GD → 0th order "natively" operates on (m)any architectures
- Boosting implies finding the architecture (how "blocks" from weak learner are assembled), so (still) restricted from this standpoint

Algorithm[☆]

[☆]simplified, see paper for full presentation

Algorithm 1 SECBOOST(S, T)

Input sample $\mathcal{S} = \{(x_i, y_i), i = 1, 2, \dots, m\}$, number of iterations T , initial (h_0, v_0) (constant classification and offset).
 Step 1 : let $H_0 \leftarrow 1 \cdot h_0$ and $w_1 = -\delta_{v_0} F(h_0) \cdot 1$;
 Step 2 : **for** $t \in [T]$
 Step 2.1 : let $h_t \leftarrow \text{WEAK_LEARNER}(\mathcal{S}_t, |w_t|)$;
 Step 2.2 : compute leveraging coefficient α_t , params $\varepsilon_t > 0, \bar{w}_{2,t} > 0$;
 Step 2.3 : let $H_t \leftarrow H_{t-1} + \alpha_t \cdot h_t$;
 Step 2.4 : **for** $i \in [m]$, let $v_{ti} \leftarrow \text{OFFSET_ORACLE}(t, i, \varepsilon_t \cdot \alpha_t^2 M_t^2 \bar{w}_{2,t})$;
 Step 2.5 : **for** $i \in [m]$, let $w_{(t+1)i} \leftarrow -\delta_{v_{ti}} F(y_i H_t(x_i))$;
 Step 2.6 : **if** $w_{t+1} = 0$ **then** break;
Return H_T .

$$M_t \doteq \max_i |h_t(x_i)|$$

Step 2.2

- Two possibilities to get $\alpha_t, \varepsilon_t, \bar{w}_{2,t}$, where $\bar{w}_{2,t}$ is any > 0 real s.t. $\mathbb{E}_{i \sim [m]} \left[\delta_{\{\alpha_t y_i h_t(x_i), v_{(t-1)i}\}} F(y_i H_{t-1}(x_i)) \cdot \left(\frac{h_t(x_i)}{M_t} \right)^2 \right] \leq \bar{w}_{2,t}$

tricky bit!

If offsets were → 0, this would be a second-order derivative

➡ **Possibility 1:** F is "nice" ⇒ easy bound: we just have to pick $\varepsilon_t > 0, \pi_t \in (0, 1)$ and α_t as:
 (example: F β -smooth ⇒ $\bar{w}_{2,t} = 2\beta$)

➡ **Possibility 2:** no niceness ⇒ Cf paper for efficient algorithm providing all params $(\alpha_t, \varepsilon_t, \bar{w}_{2,t} \& \pi_t)$

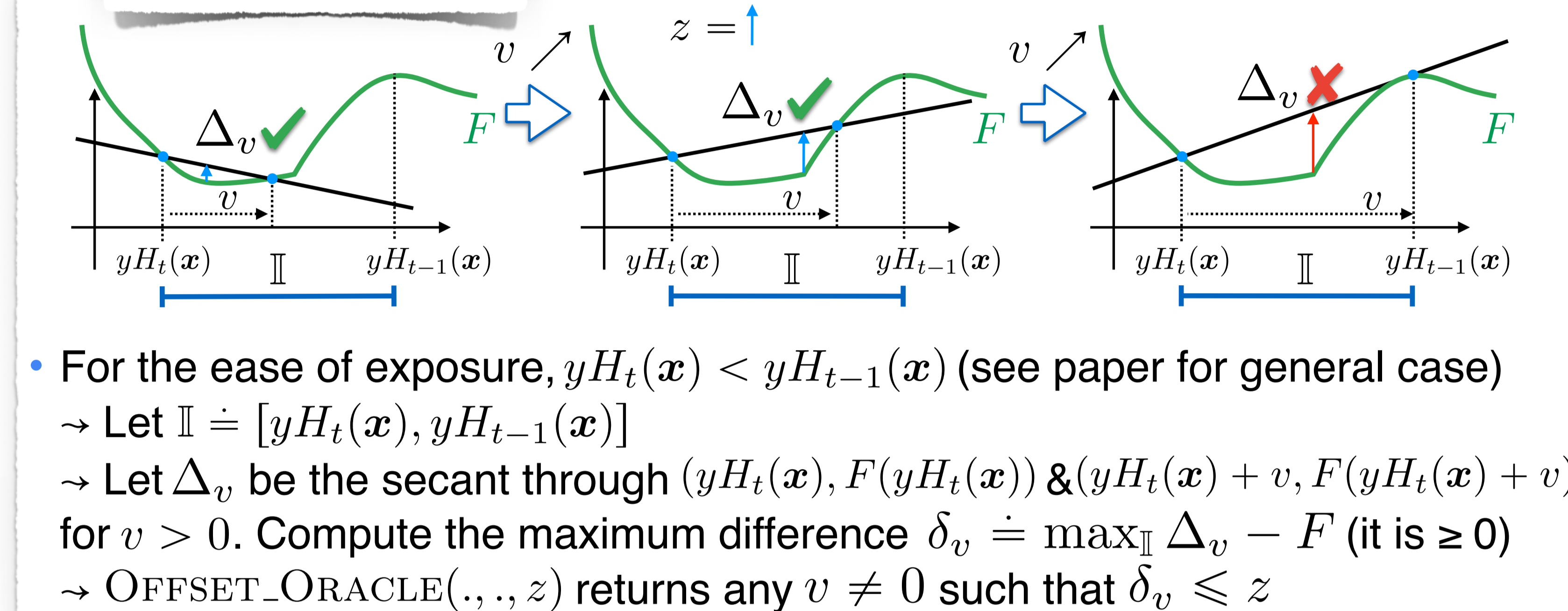
$$\alpha_t \in \frac{\eta_t}{2(1 + \varepsilon_t) M_t^2 \bar{w}_{2,t}} \cdot [1 - \pi_t, 1 + \pi_t]$$

$$\eta_t \doteq (1/m) \cdot \sum_i w_{ti} y_i h_t(x_i) \quad \text{important for boosting rate}$$

Notable generalizations with respect to "boosting-à-la-Valiant"

- Examples for WEAK_LEARNER can be *label flipped*: $\mathcal{S}_t \doteq \{(x_i, y_i \cdot \text{sign}(w_{ti}))\}$
- Need an "oracle" giving *offsets* (implementation generic or loss dependent)

The offset oracle



- For the ease of exposure, $y_{H_t}(x) < y_{H_{t-1}}(x)$ (see paper for general case)
 → Let $\mathbb{I} \doteq [y_{H_t}(x), y_{H_{t-1}}(x)]$
 → Let Δ_v be the secant through $(y_{H_t}(x), F(y_{H_t}(x)))$ & $(y_{H_t}(x) + v, F(y_{H_t}(x) + v))$, for $v > 0$. Compute the maximum difference $\delta_v \doteq \max_{\mathbb{I}} \Delta_v - F$ (it is ≥ 0)
 → $\text{OFFSET_ORACLE}(\cdot, \cdot, z)$ returns any $v \neq 0$ such that $\delta_v \leq z$

Toolbox

- Generalization of quantum calculus' (\neq quantum computation) v -derivative:

$$\delta_{\mathcal{V}} F(z) \doteq \begin{cases} F(z) & \text{if } \mathcal{V} = \emptyset \\ \delta_v F(z) & \text{if } \mathcal{V} = \{v\} \\ \delta_{\{v\}}(\delta_{\mathcal{V} \setminus \{v\}} F)(z) & \text{otherwise } \mathcal{V} = \{v, w, \dots\} \text{ (eventually multiset)} \end{cases}$$

- Singleton $\mathcal{V} = \{v\} \Rightarrow$ classical secant's slope

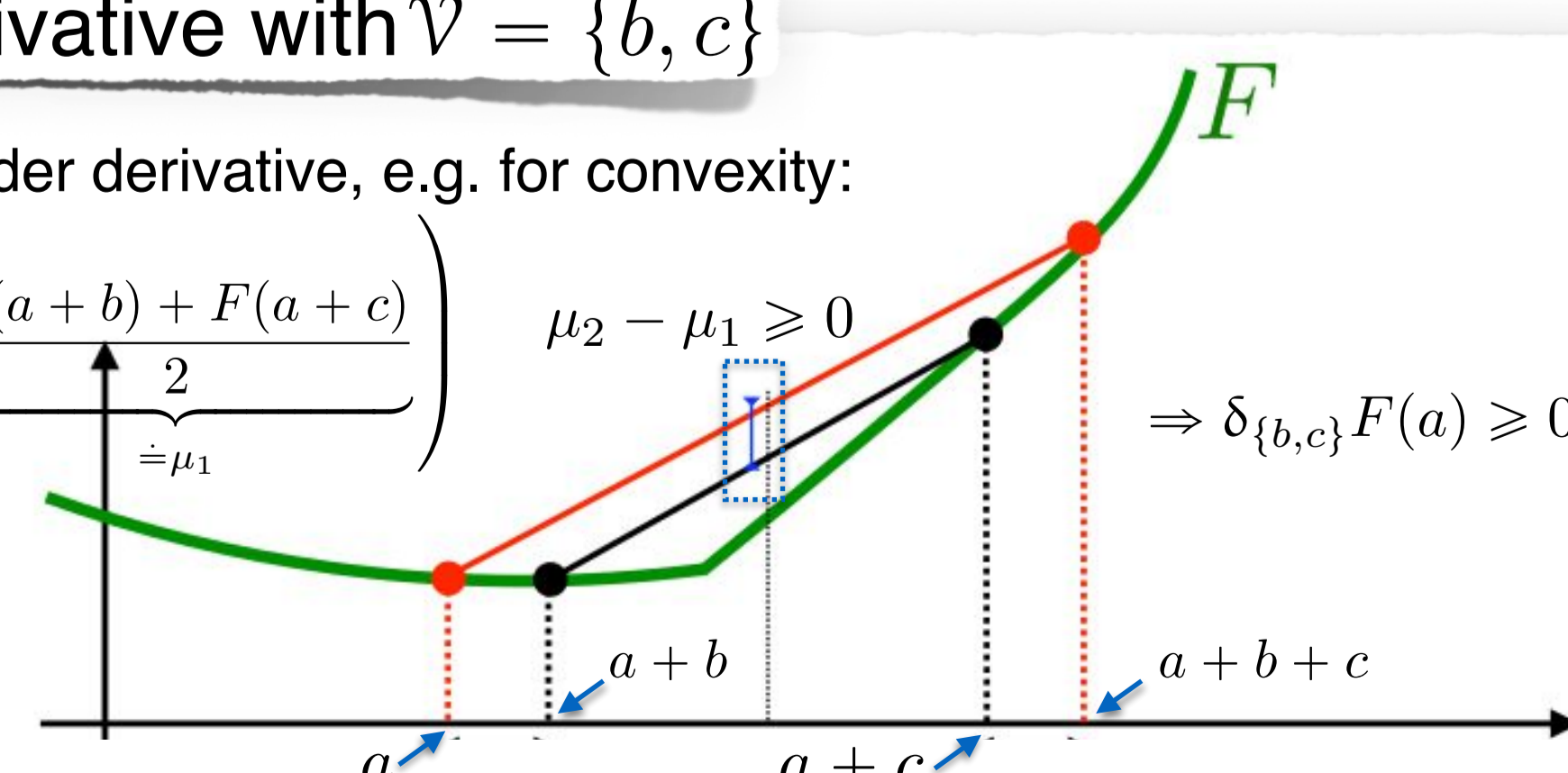
$$\delta_v F(z) \doteq \frac{F(z+v) - F(z)}{v} \quad \text{called } h\text{-derivative, with } \mathcal{V} = \{v, v, \dots\} \text{ in "Quantum calculus", Kac \& Cheung, 2002}$$

Example: second-order v -derivative with $\mathcal{V} = \{b, c\}$

generalizes some properties of second-order derivative, e.g. for convexity:

$$\delta_{\{b,c\}} F(a) = \frac{2}{b} \cdot \frac{1}{c} \cdot \left(\frac{F(a+b+c) + F(a)}{2} - \frac{F(a+b) + F(a+c)}{2} \right) \quad \mu_2 - \mu_1 \geq 0 \Rightarrow \delta_{\{b,c\}} F(a) \geq 0$$

(wlog $c > b > 0$, see Lemma 5.2 in paper for more & general case(s))



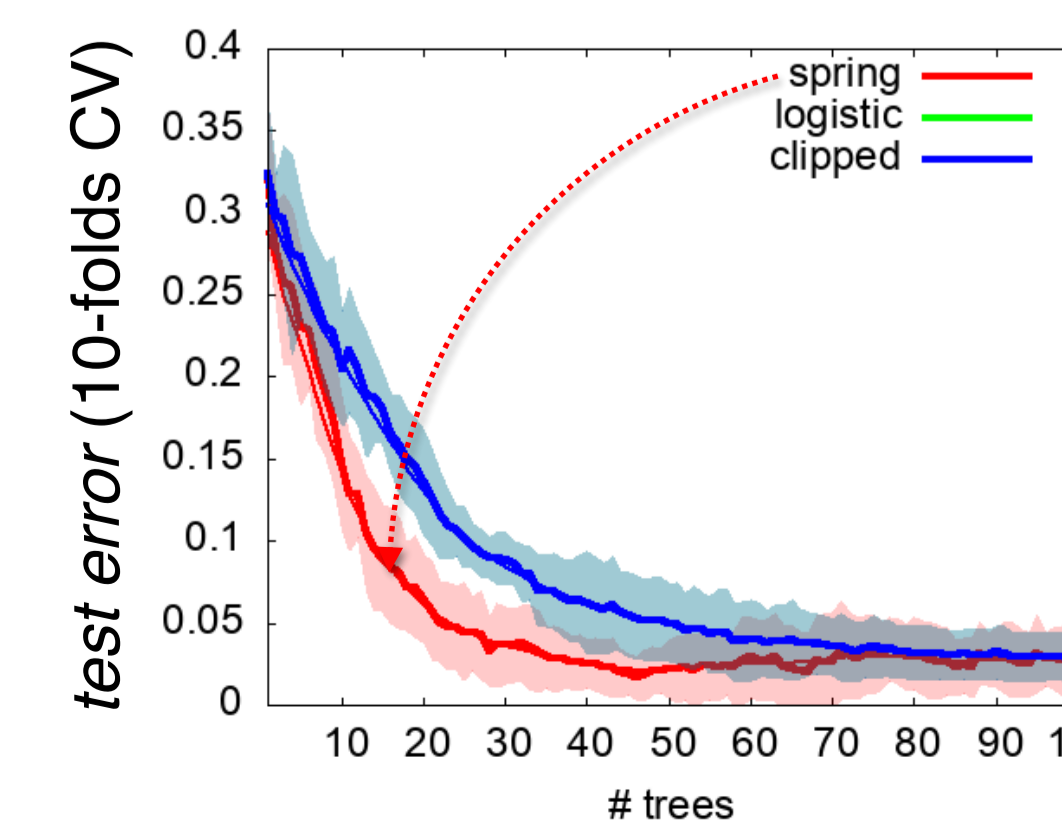
Boosting!

- Weak Learning Assumption: $\left| \mathbb{E}_{w_t} \left[\tilde{y}_{ti} \cdot \frac{h_t(x_i)}{M_t} \right] \right| \geq \gamma > 0$
 ➡ $\tilde{y}_{ti} \doteq y_i \cdot \text{sign}(w_{ti})$ label eventually *flipped*
 ➡ $\tilde{w}_t = |w_t|$ normalized to unit
- Weak Convergence Regime: $\frac{\mathbb{E}_{i \sim [m]} [w_{ti}]^2}{\bar{w}_{2,t}} \geq \rho > 0$
 ➡ numerator ← 1st order v -derivative, expected *signed* weights
 ➡ denominator ← 2nd-order v -derivative, loss "jiggling"

Theorem. Let the expected empirical loss of classifier H be $F(\mathcal{S}, H) \doteq \mathbb{E}_{i \sim [m]} [F(y_i H(x_i))]$ and its initial value $F_0 \doteq F(\mathcal{S}, h_0)$. Suppose WLA+WCR hold. Then for any $z \in \mathbb{R}$ such that $F(z) \leq F_0$, if SecBoost is run for a number of iterations

$$T \geq \frac{4(F_0 - F(z))}{\gamma^2 \rho} \cdot \frac{1 + \max_t \varepsilon_t}{1 - \max_t \pi_t^2}$$

then $F(\mathcal{S}, H_T) \leq F(z)$.



- Example implementation (details: Cf paper)
- Weak Classifiers = size-20 DTs
- Losses: **logistic** & two variations: **clipped logistic** and a non-[cvx, Lip, diff] loss ("**spring loss**")

