

Covariate-Guided Clusterwise Linear Regression for Generalization to Unseen Data

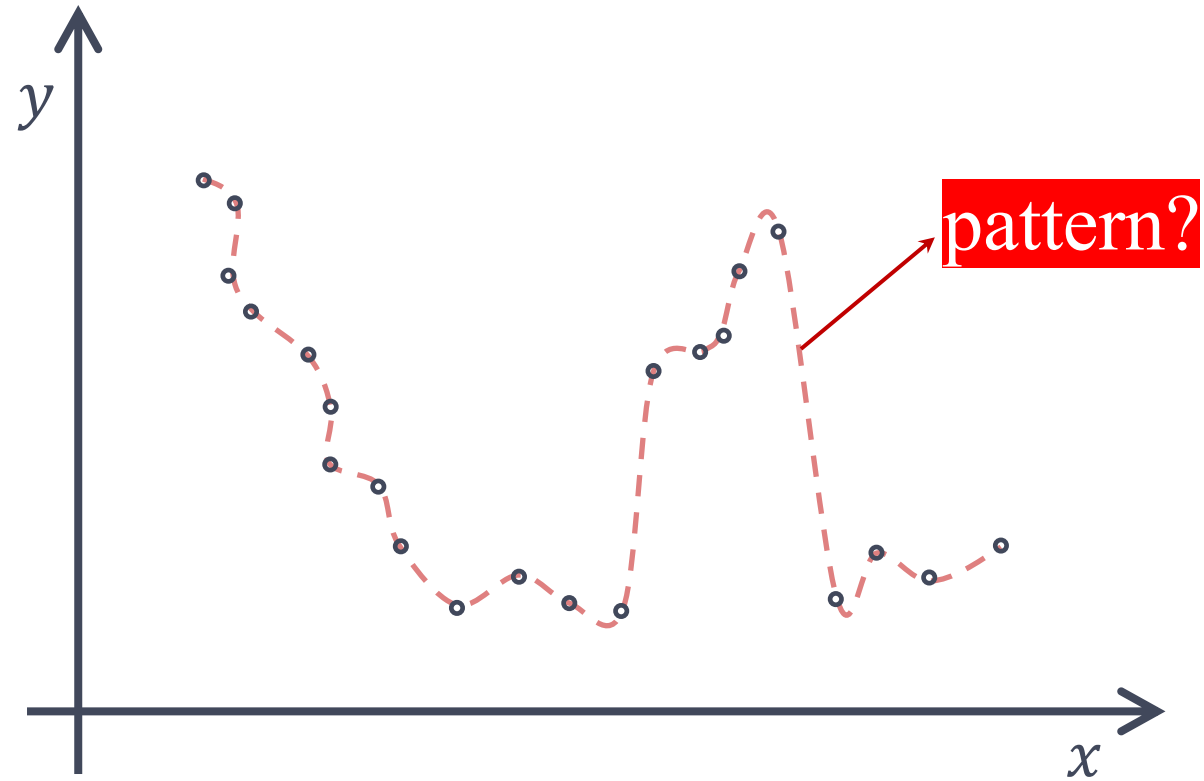
Dohyun Bu

Hyunho Kim

Jong-Seok Lee

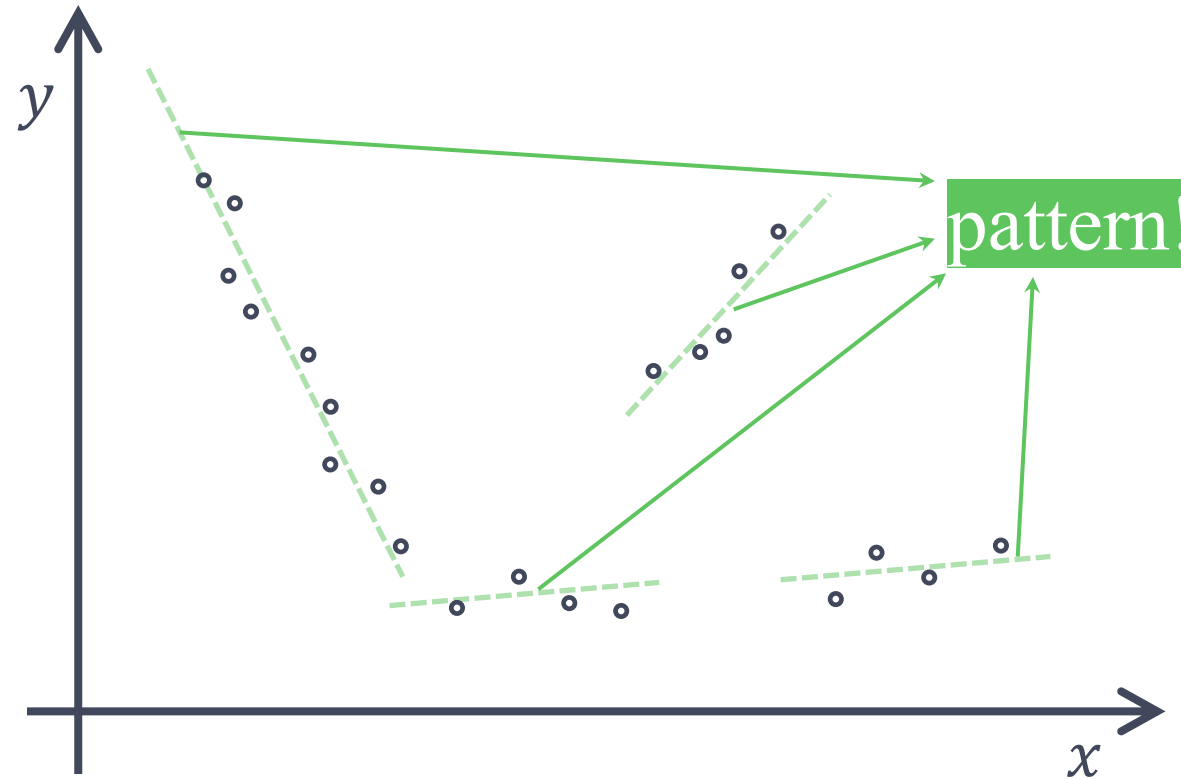


Why Local Linear Models?



Many real-world data exhibit **globally nonlinear**, but **locally simple** patterns.

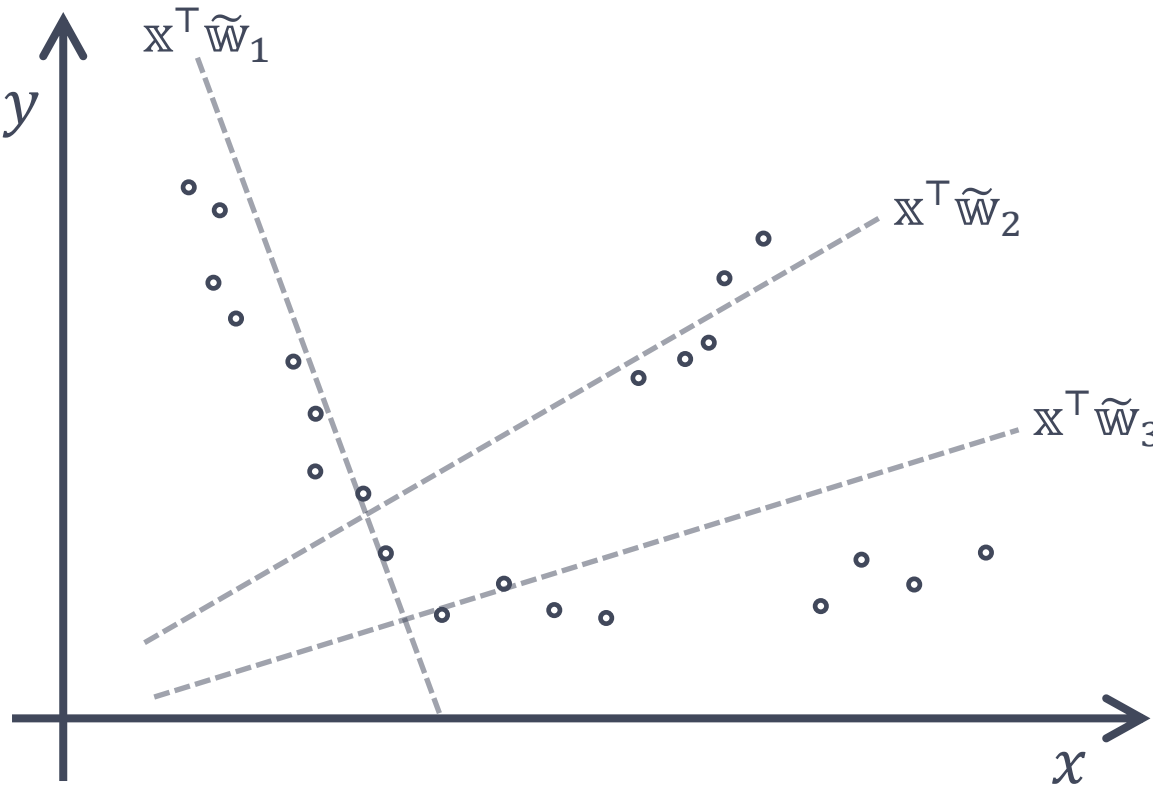
Why Local Linear Models?



Many real-world data exhibit **globally nonlinear**, but **locally simple** patterns.

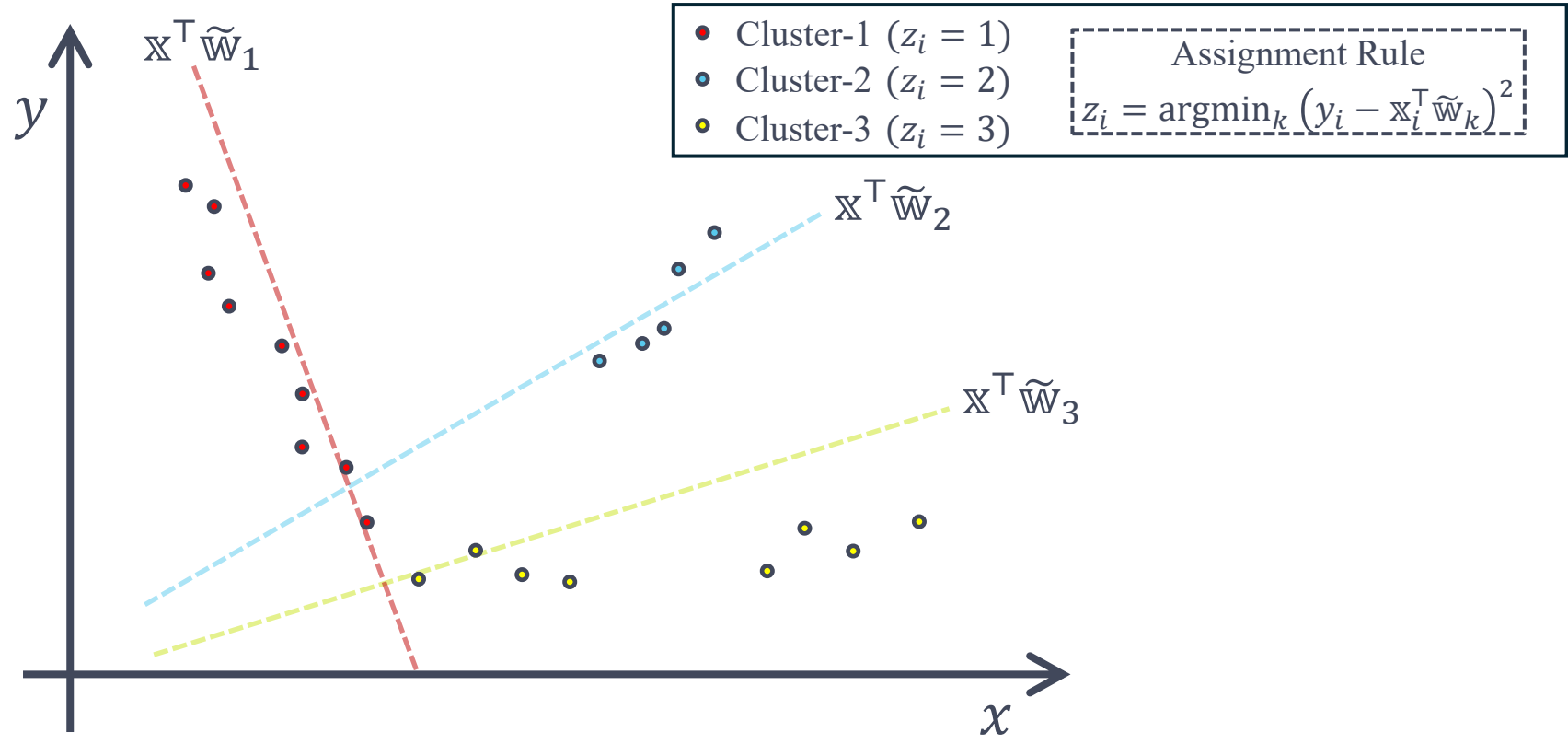
→ This motivates **Clusterwise Linear Regression (CLR)**

CLR Learns Local Experts,



{ Clusterwise Linear Regression (CLR) }

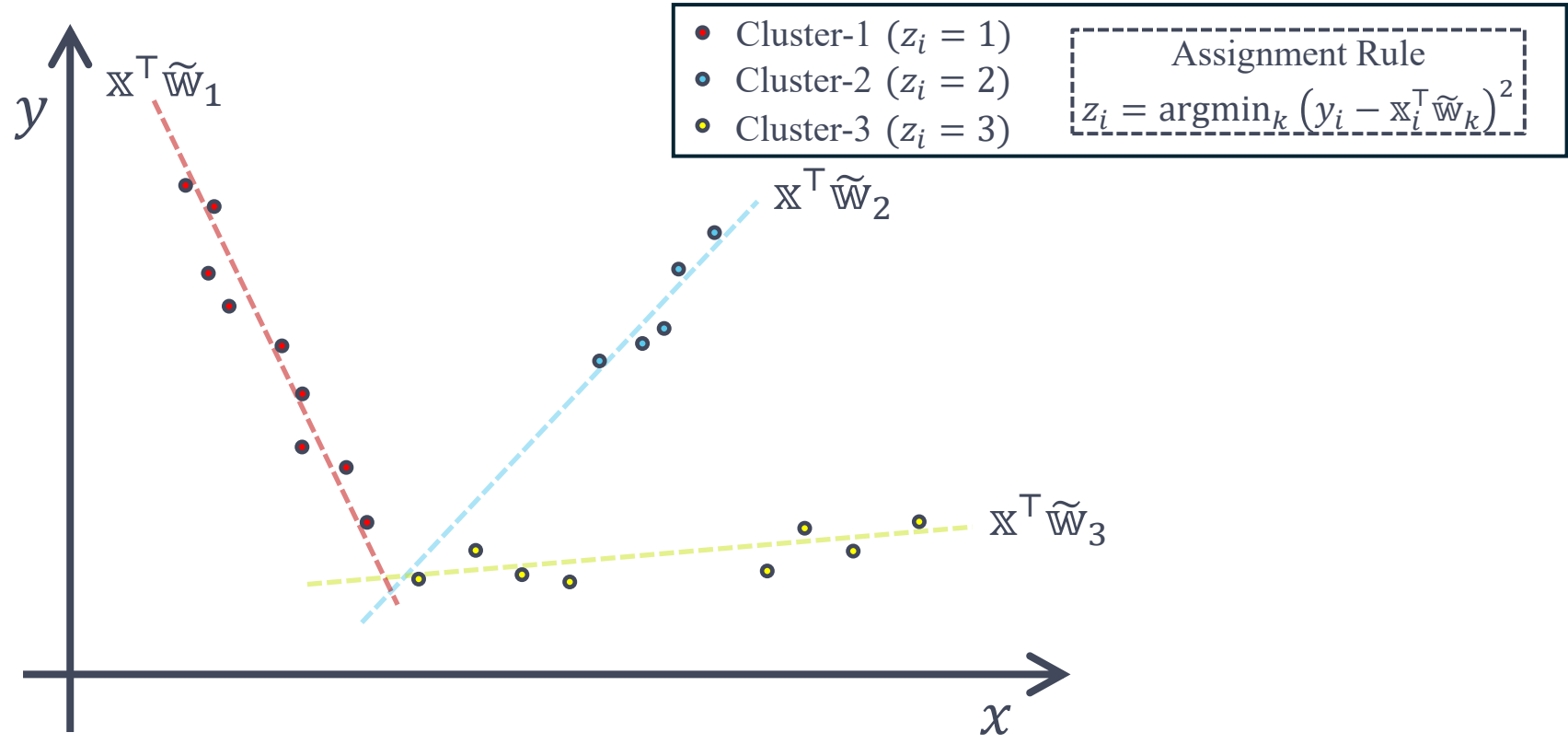
CLR Learns Local Experts,



{ Clusterwise Linear Regression (CLR) }

Step 1. Response-aware assignment

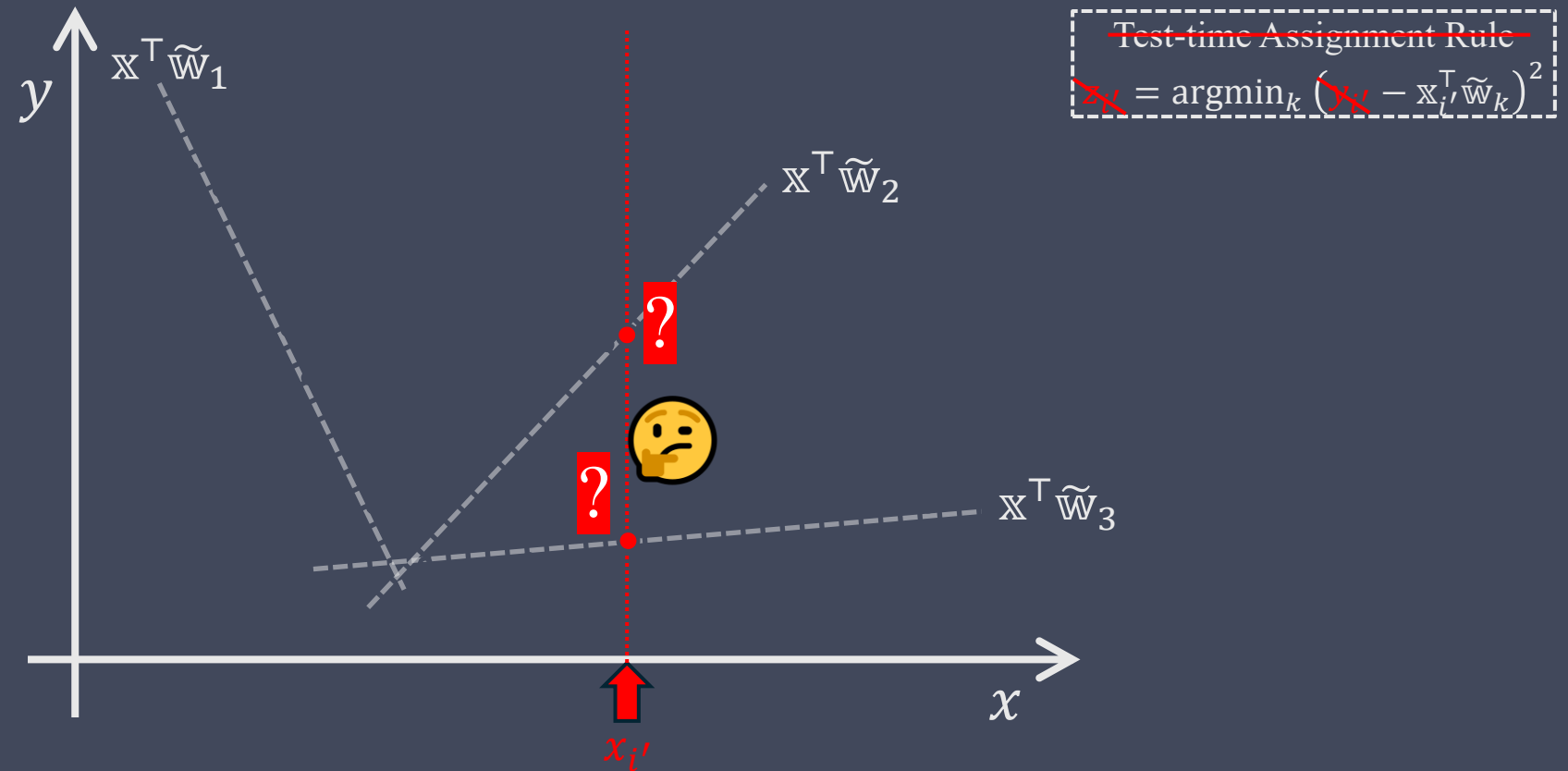
CLR Learns Local Experts,



{ Clusterwise Linear Regression (CLR) }

Step 2. Update each expert using assigned points **Stable convergence!** 😊

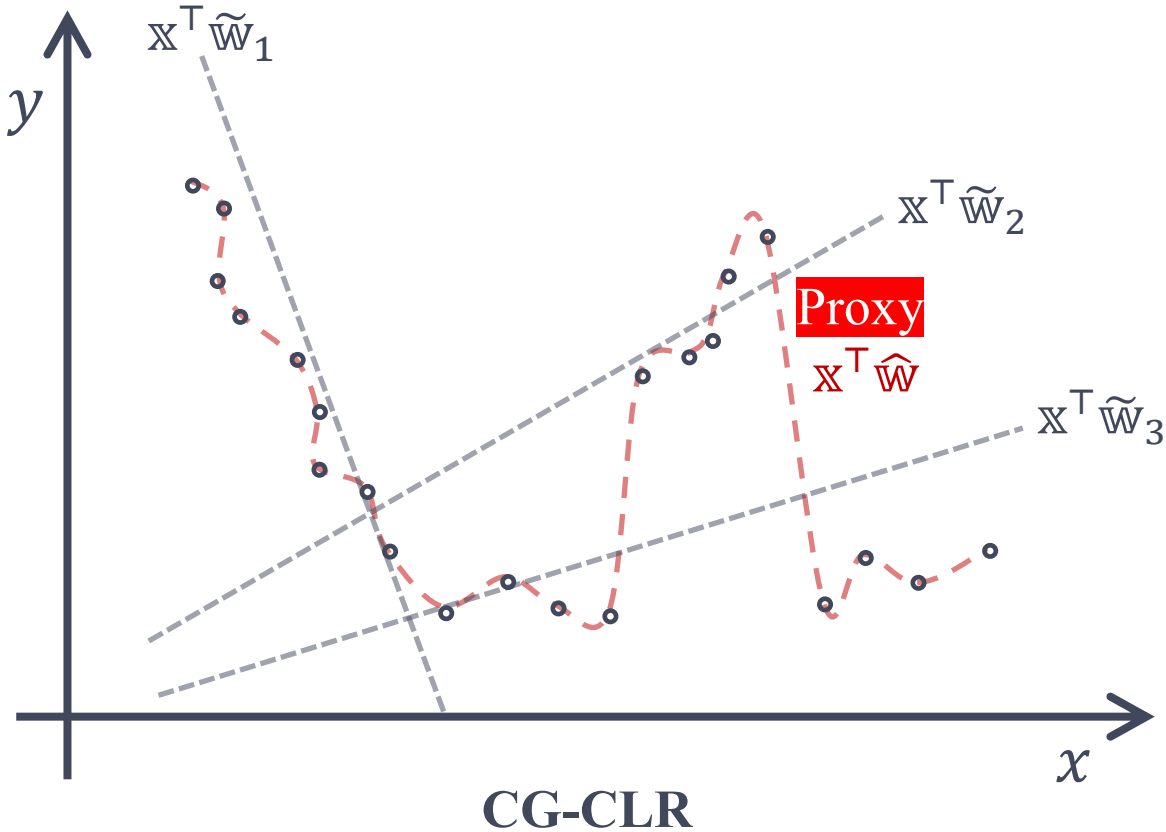
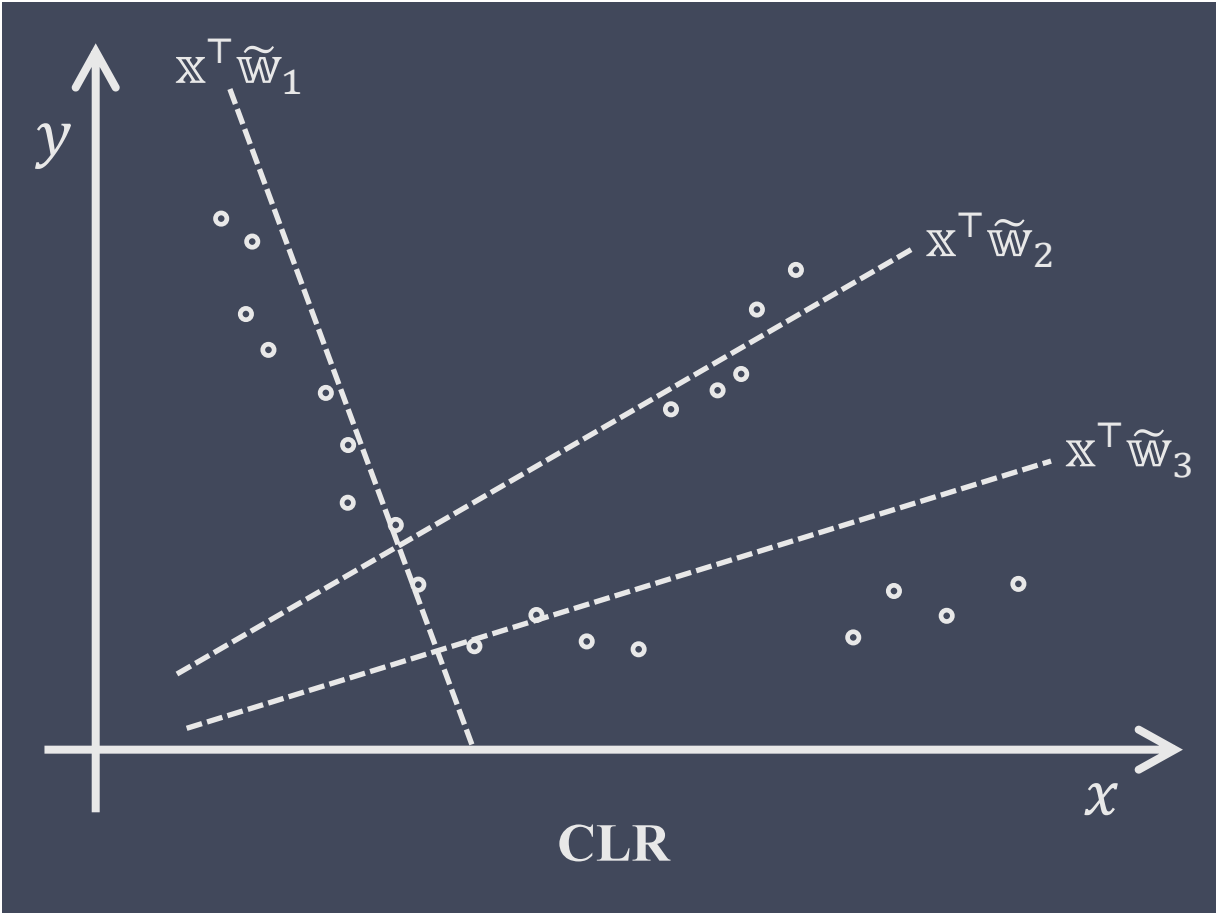
CLR Learns Local Experts, But Not Test-time Routing



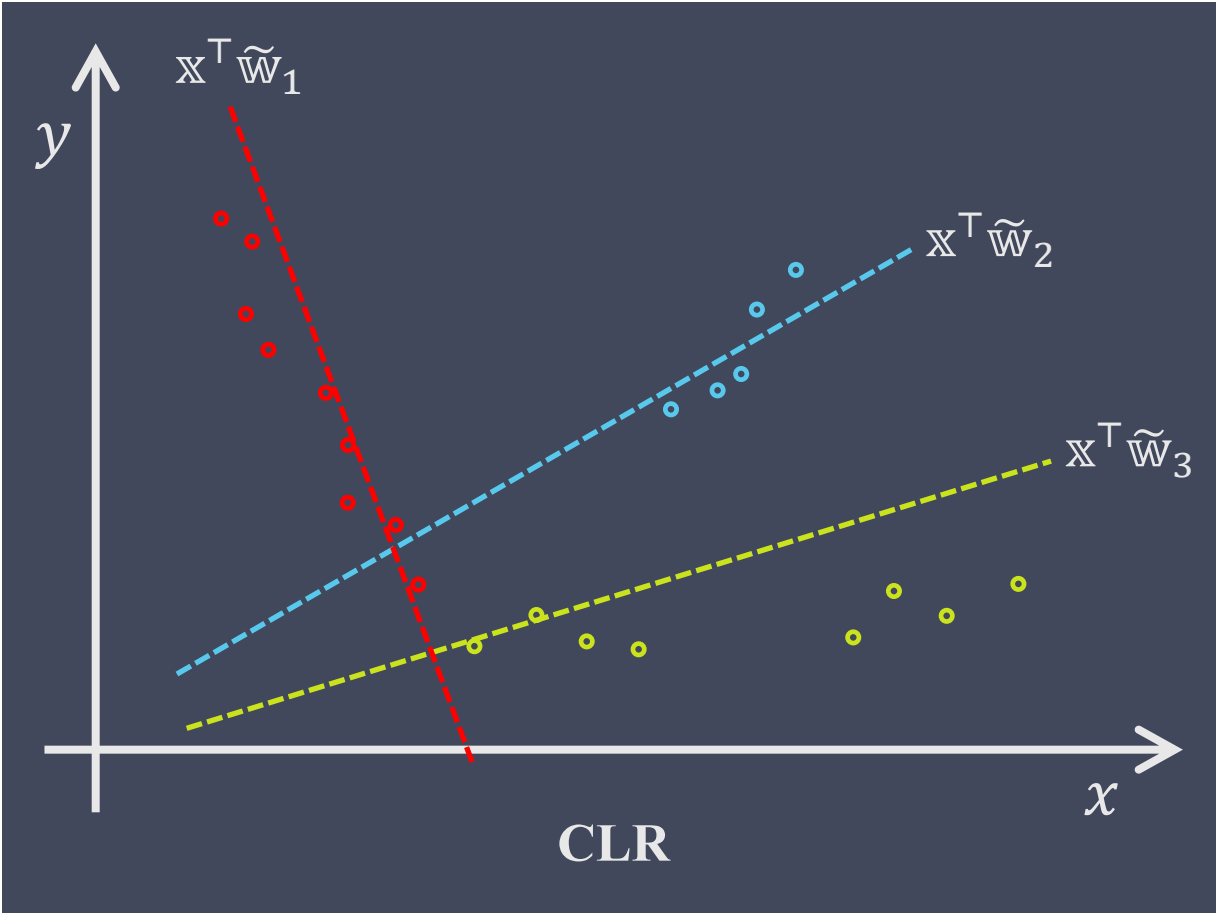
For a new covariate $x_{i'}$ and unknown $y_{i'}$, we cannot assign the cluster- $z_{i'}$.

“Response-aware training, **infeasible** test-time prediction”

CG-CLR Adds Test-time Routing via a Proxy!



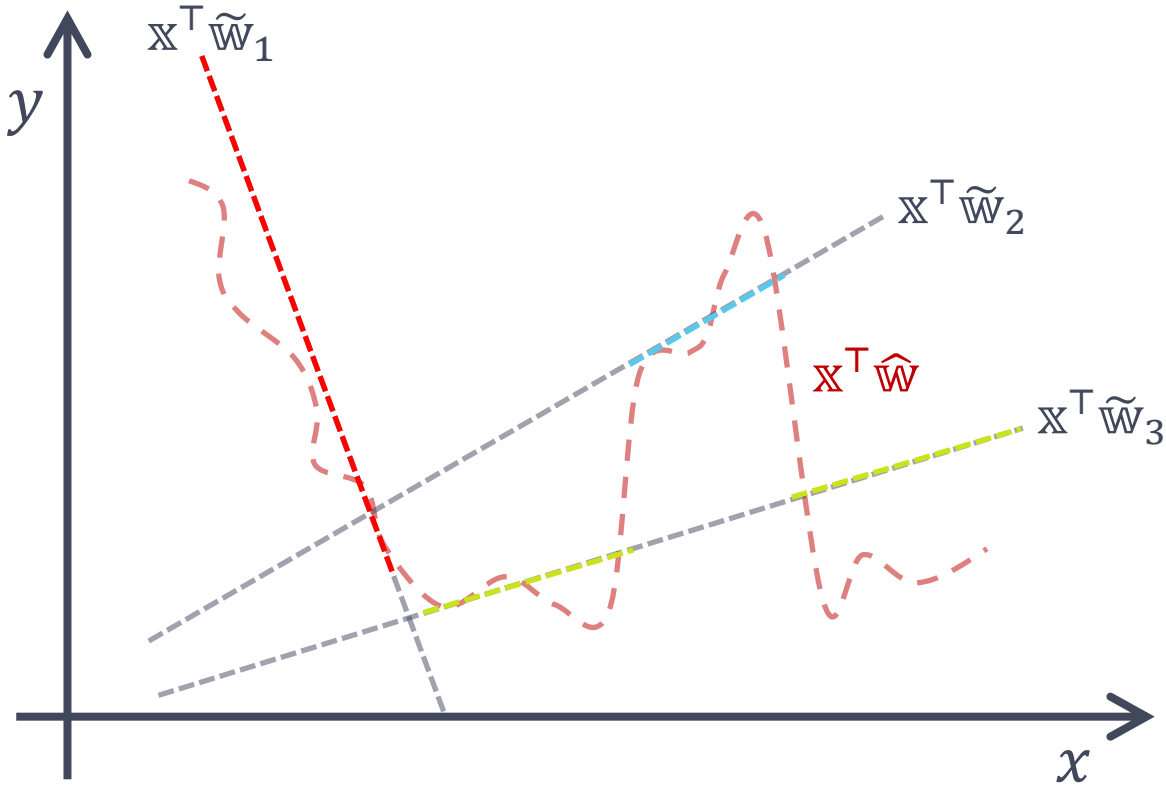
CG-CLR Adds Test-time Routing via a Proxy!



CLR

Response-aware Assignment

$$z_i = \operatorname{argmin}_k (y_i - \mathbf{x}_i^T \tilde{\mathbf{w}}_k)^2$$

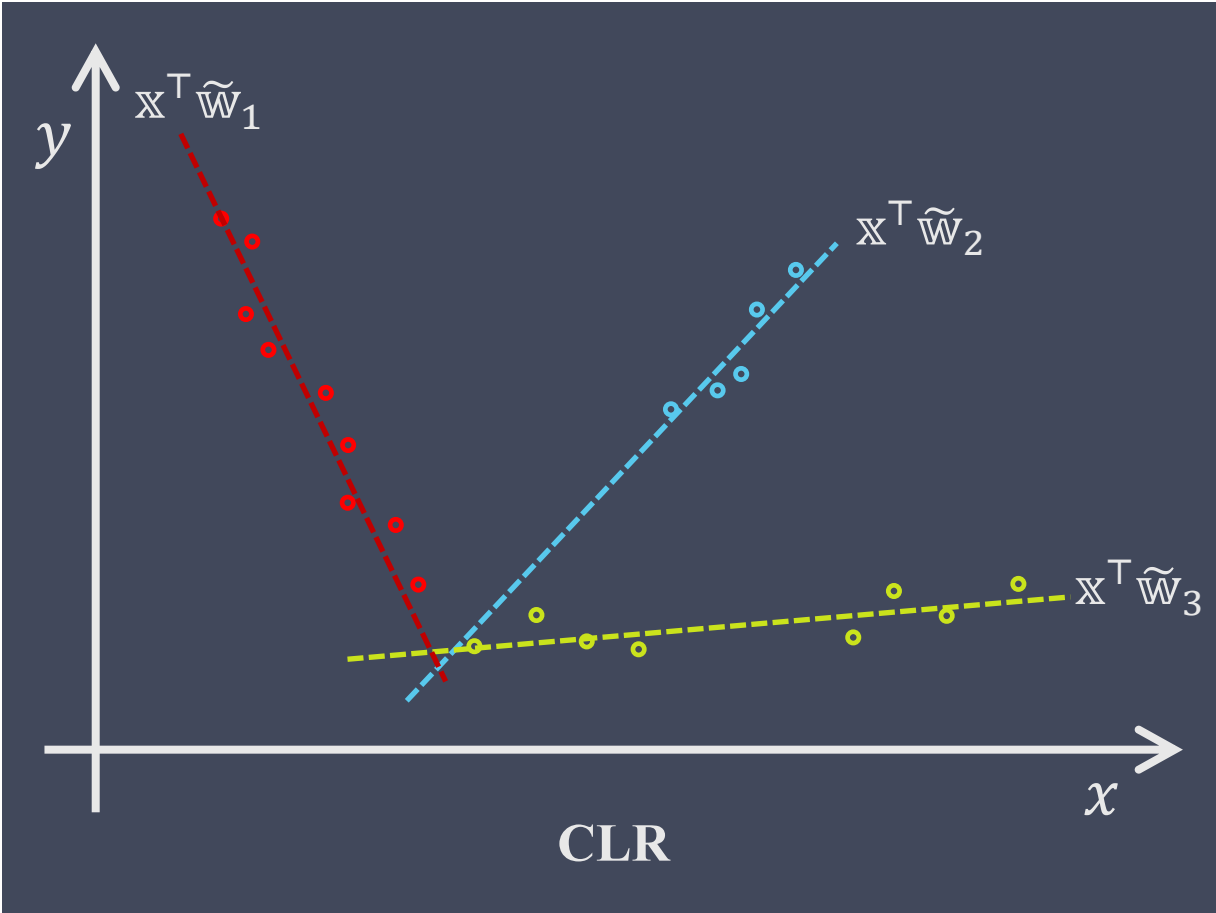


CG-CLR

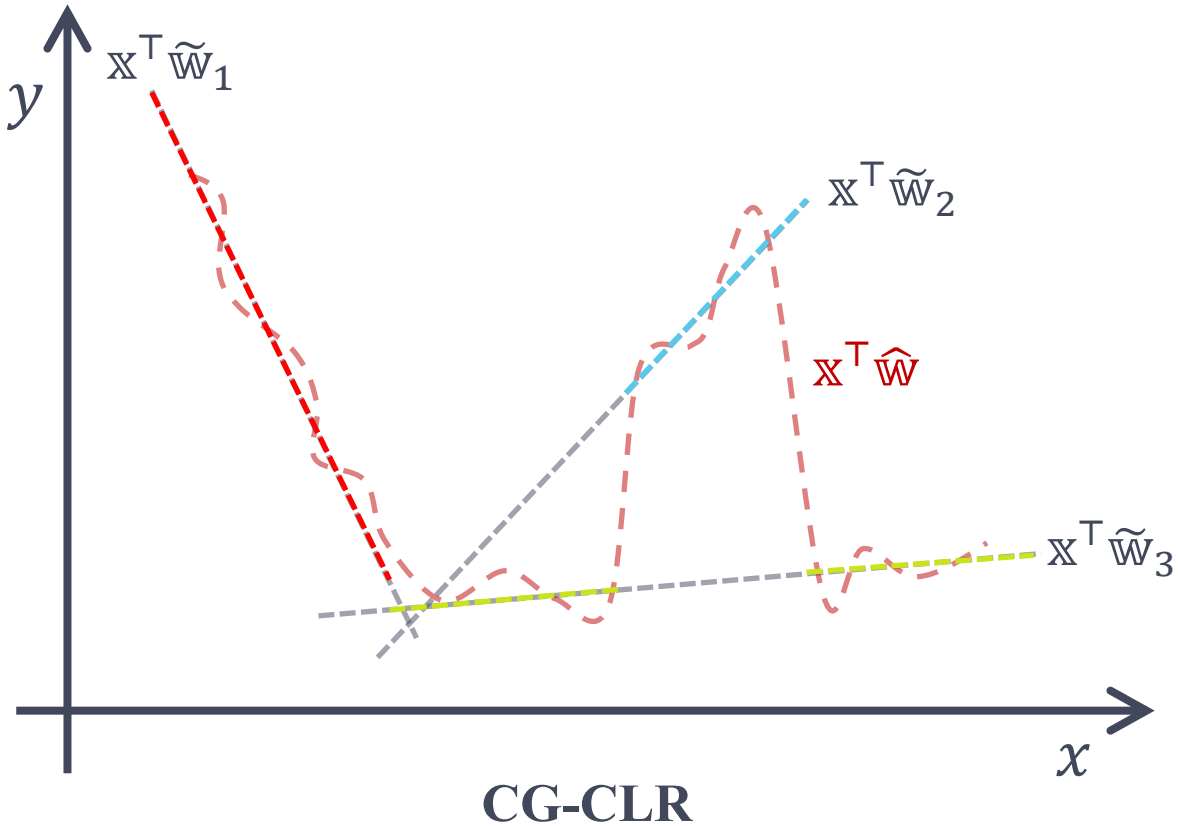
Proxy-based Assignment

$$z_i = \operatorname{argmin}_k (\mathbf{x}_i^T \hat{\mathbf{w}}_i - \mathbf{x}_i^T \tilde{\mathbf{w}}_k)^2$$

CG-CLR Adds Test-time Routing via a Proxy!



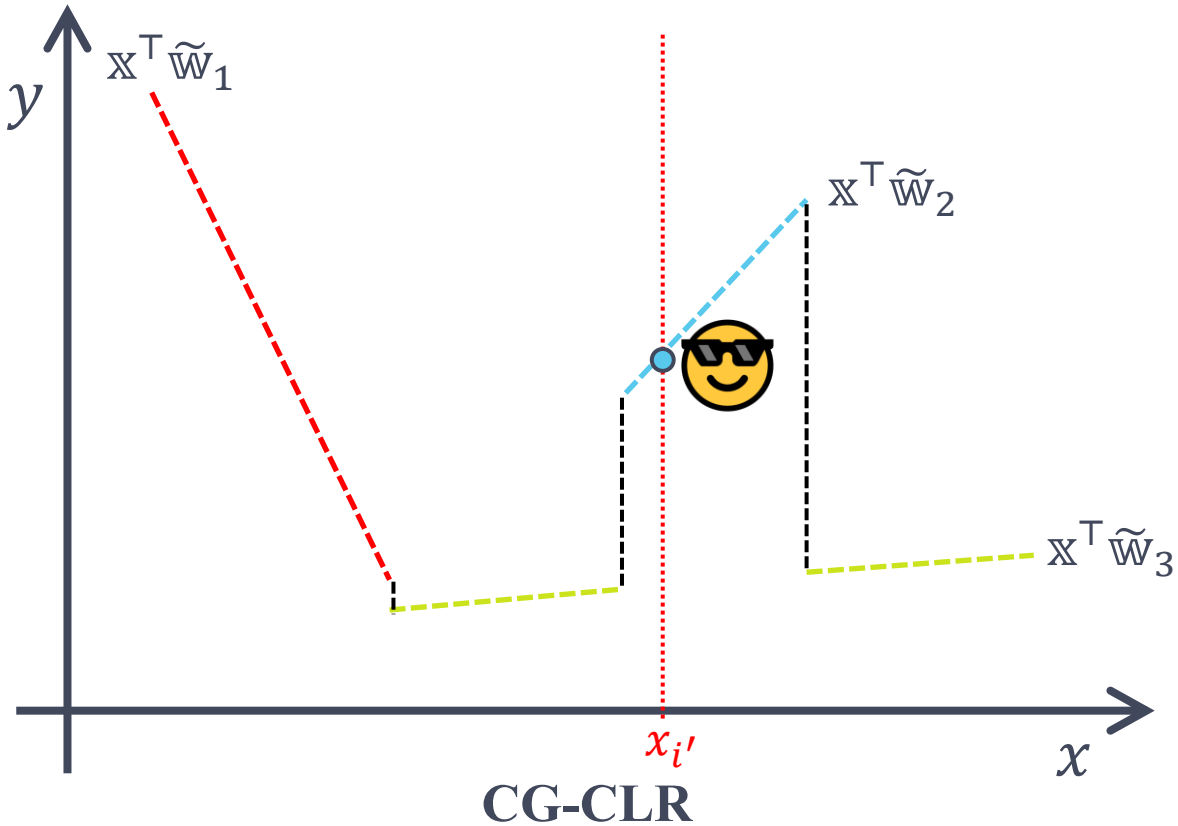
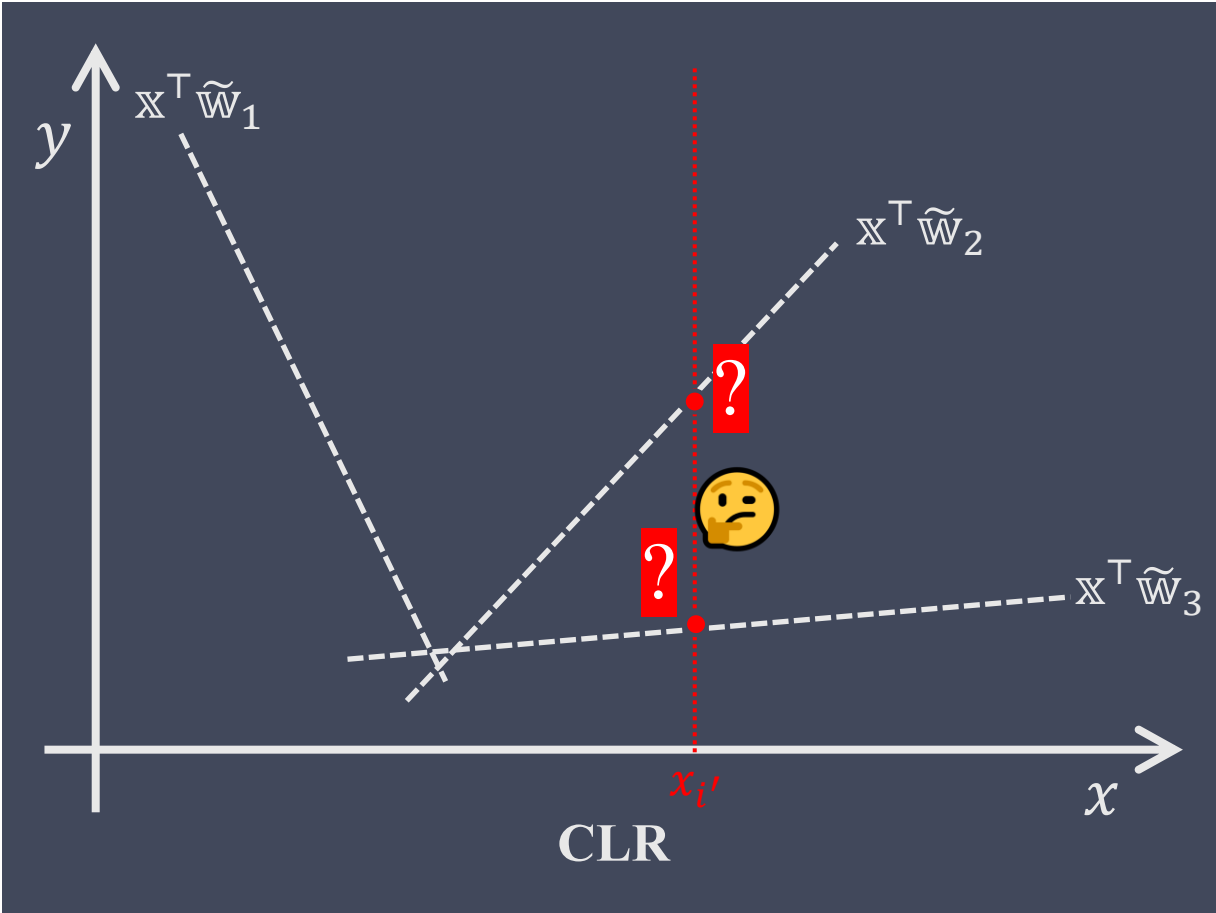
Update each expert using assigned points



Update each expert using assigned **proxy outputs**

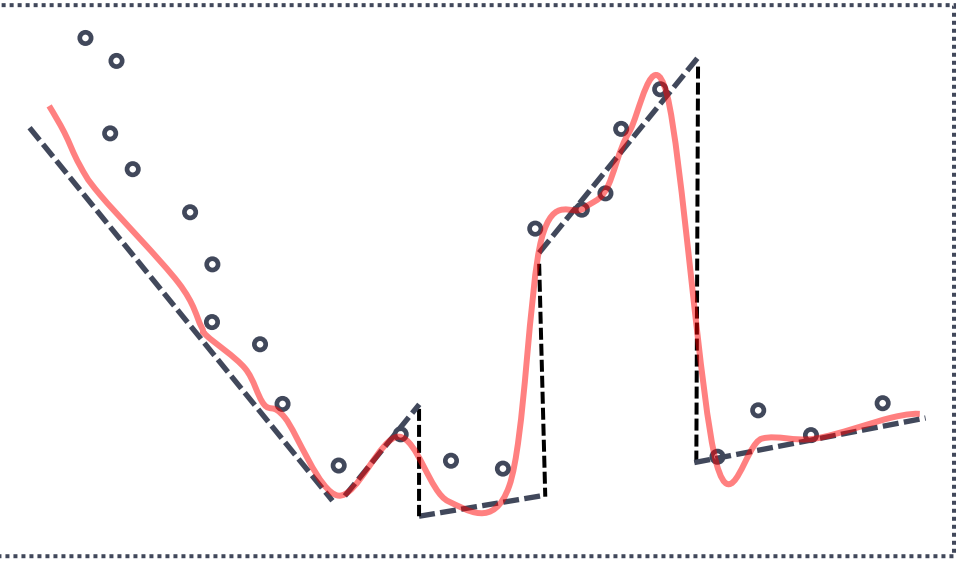
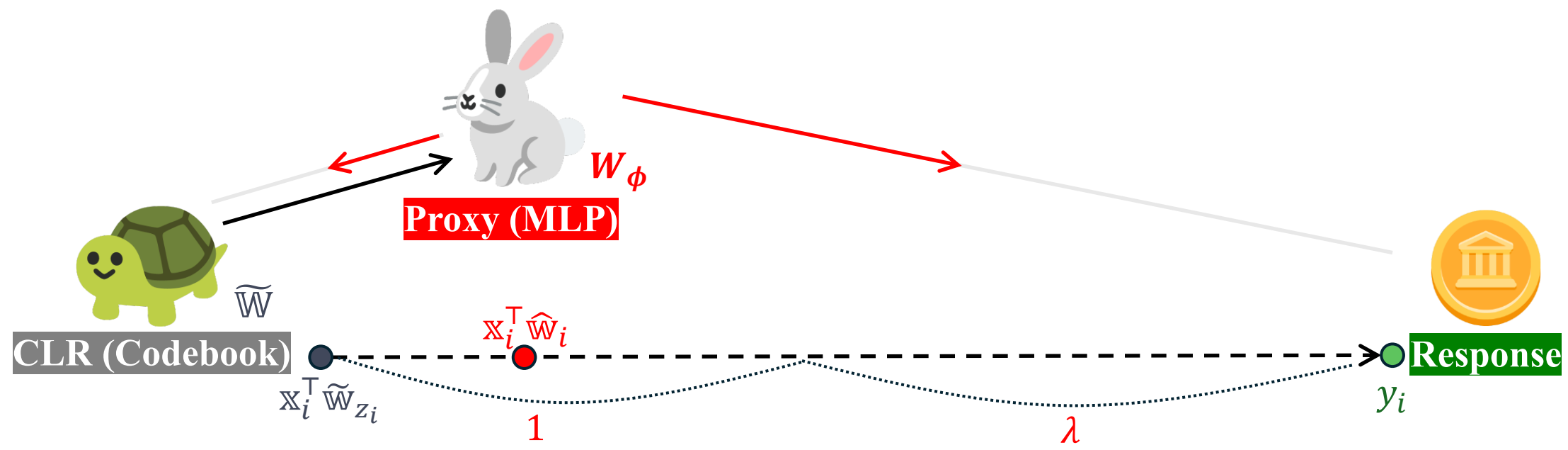
Stable convergence! 😊

CG-CLR Adds Test-time Routing via a Proxy!



Covariate-guided test-time prediction!
without response $y_{i'}$

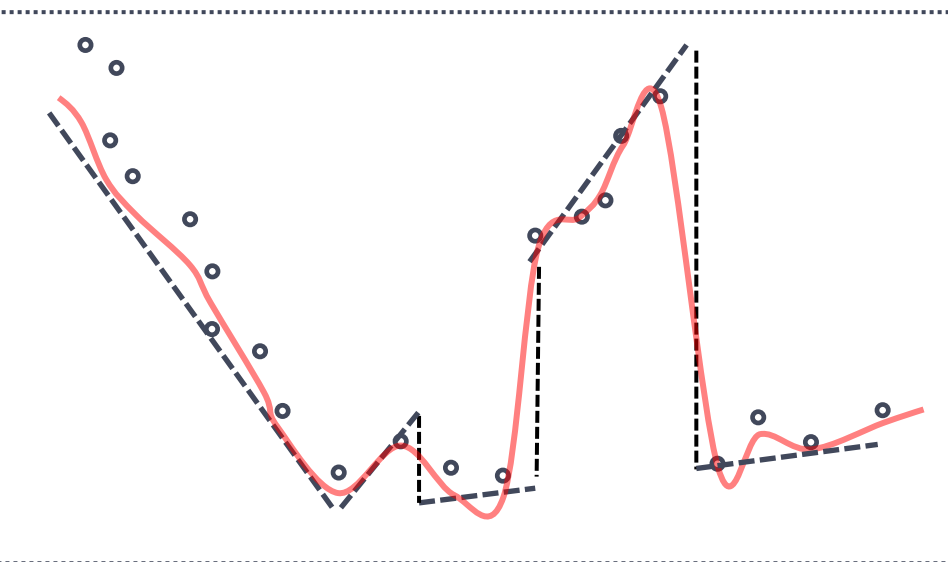
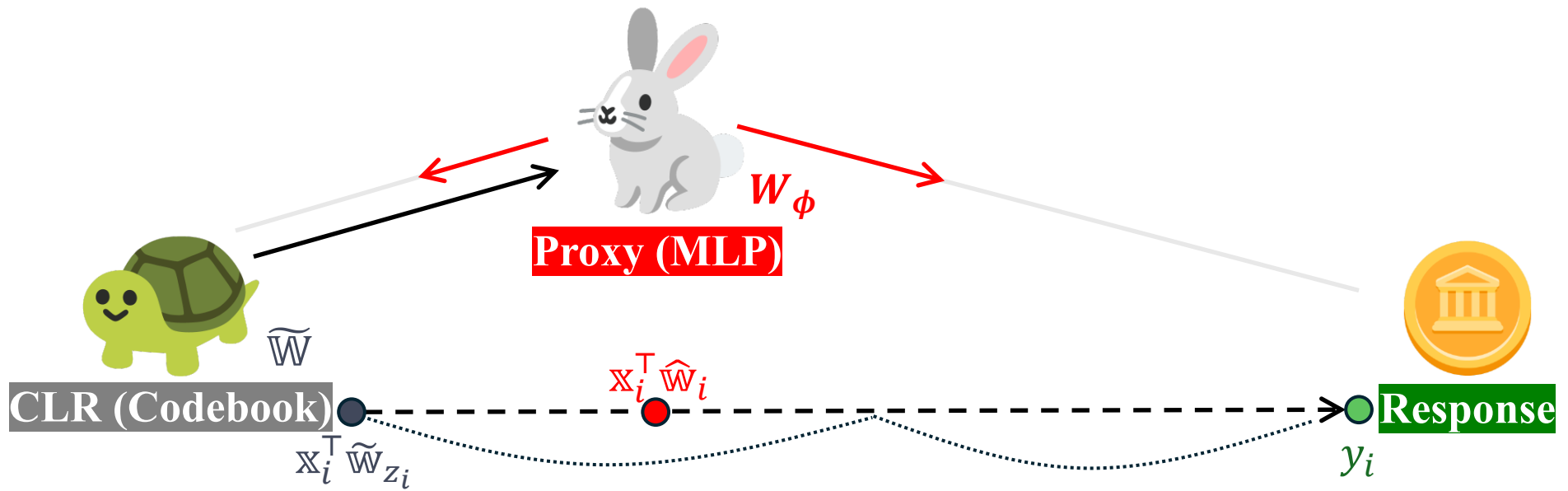
Why does the Optimization stay Stable?



$$\phi^{(t+1)} \leftarrow \phi^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i=1}^N \left(\frac{y_i + \lambda \mathbf{x}_i^\top \tilde{\mathbf{w}}_{z_i}^{(t)}}{1+\lambda} - \mathbf{x}_i^\top \hat{\mathbf{w}}_i^{(t)} \right) (\nabla_{\phi} \hat{\mathbf{w}}_i^{(t)})^\top \mathbf{x}_i,$$

$$\tilde{\mathbf{w}}_j^{(t+1)} \leftarrow \tilde{\mathbf{w}}_j^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i \in S_j^{(t)}} \mathbf{x}_i \mathbf{x}_i^\top (\hat{\mathbf{w}}_i^{(t)} - \tilde{\mathbf{w}}_j^{(t)}), \quad \forall j \in [K].$$

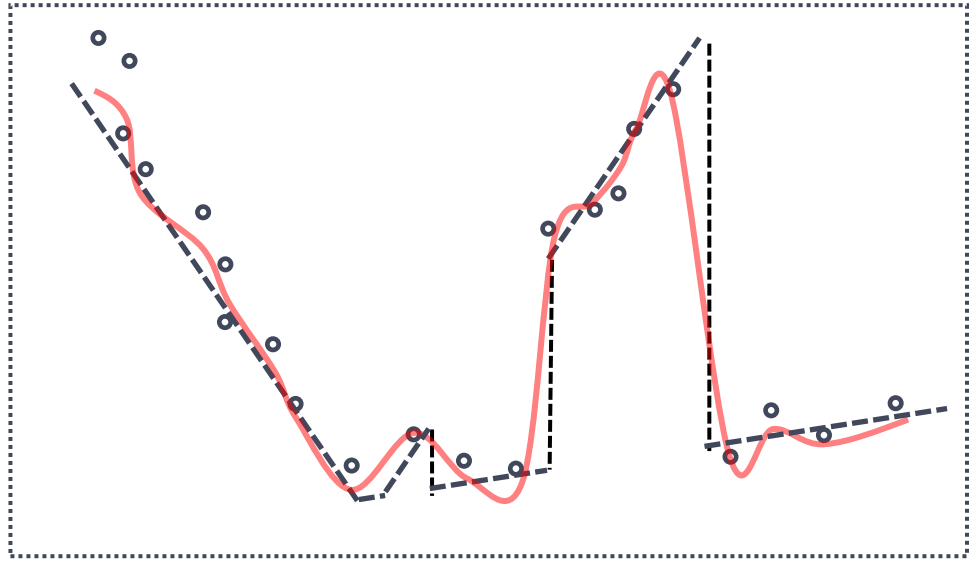
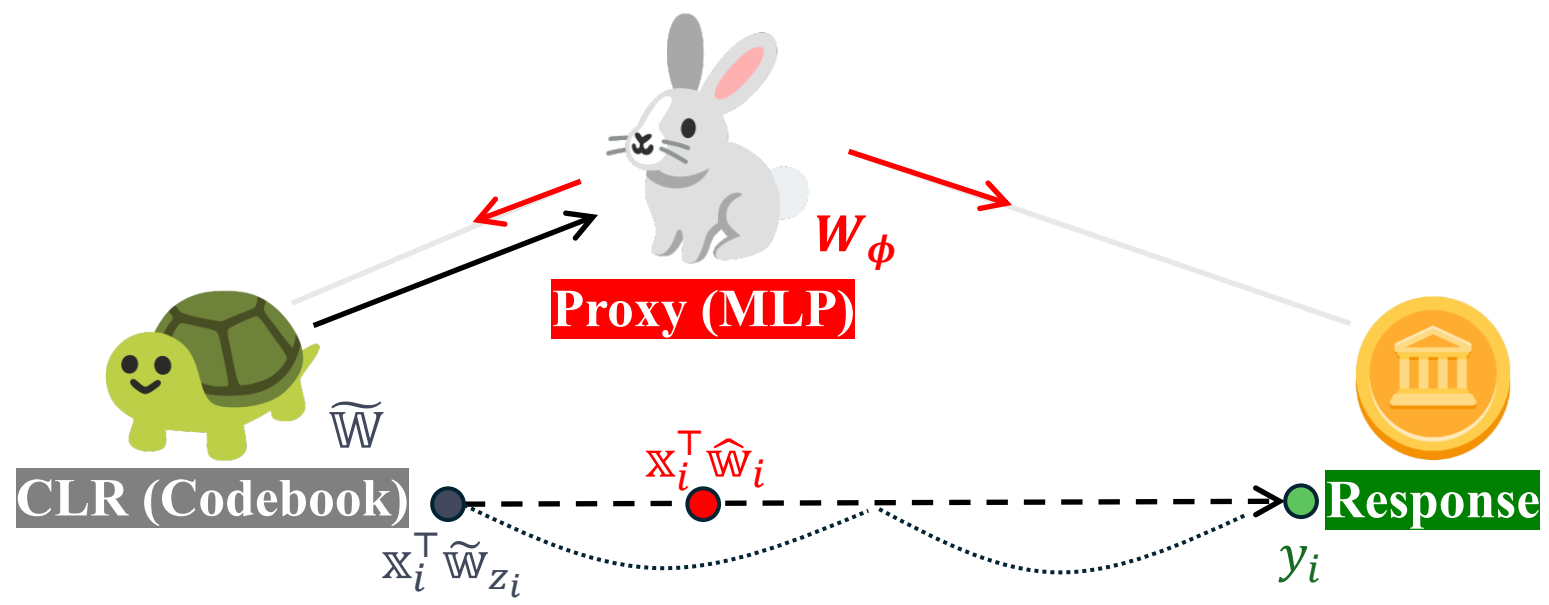
Why does the Optimization stay Stable?



$$\phi^{(t+1)} \leftarrow \phi^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i=1}^N \left(\frac{y_i + \lambda \mathbf{x}_i^\top \tilde{W}_{z_i}^{(t)}}{1+\lambda} - \mathbf{x}_i^\top \hat{W}_i^{(t)} \right) (\nabla_{\phi} \hat{W}_i^{(t)})^\top \mathbf{x}_i,$$

$$\tilde{W}_j^{(t+1)} \leftarrow \tilde{W}_j^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i \in S_j^{(t)}} \mathbf{x}_i \mathbf{x}_i^\top (\hat{W}_i^{(t)} - \tilde{W}_j^{(t)}), \quad \forall j \in [K].$$

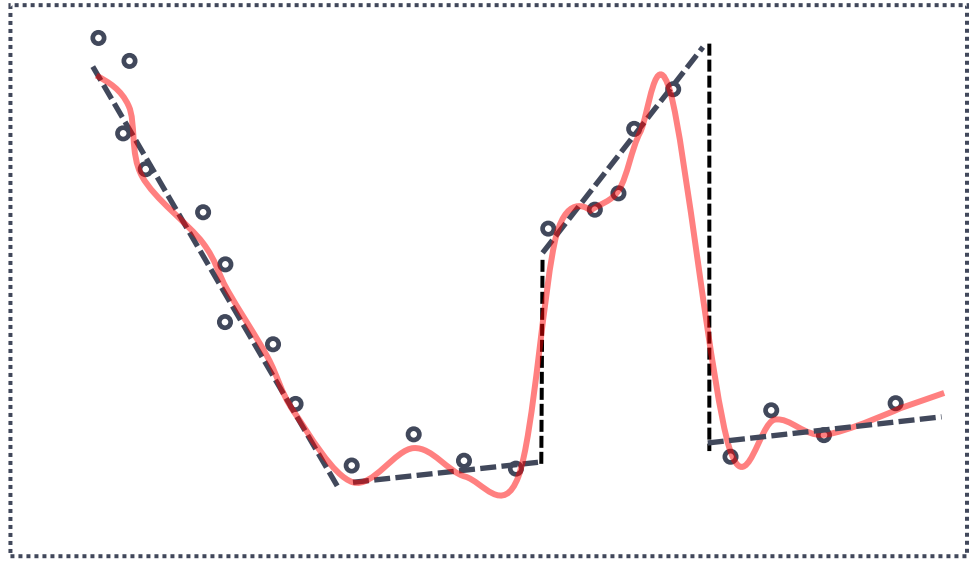
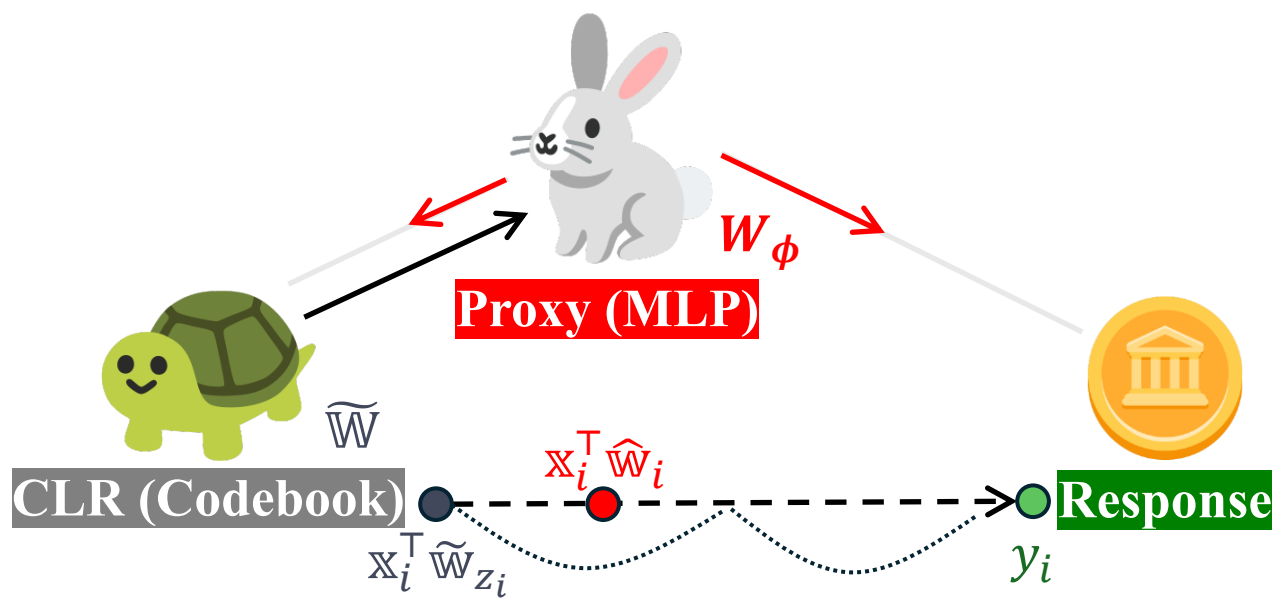
Why does the Optimization stay Stable?



$$\phi^{(t+1)} \leftarrow \phi^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i=1}^N \left(\frac{y_i + \lambda \mathbf{x}_i^T \tilde{W}_{z_i}^{(t)}}{1+\lambda} - \mathbf{x}_i^T \hat{W}_i^{(t)} \right) (\nabla_{\phi} \hat{W}_i^{(t)})^T \mathbf{x}_i,$$

$$\tilde{W}_j^{(t+1)} \leftarrow \tilde{W}_j^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i \in S_j^{(t)}} \mathbf{x}_i \mathbf{x}_i^T (\hat{W}_i^{(t)} - \tilde{W}_j^{(t)}), \quad \forall j \in [K].$$

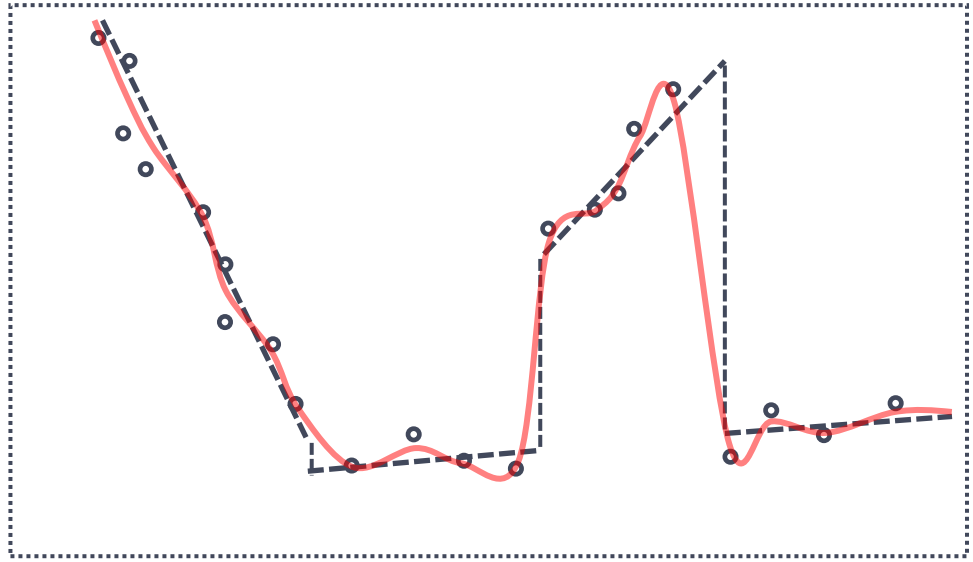
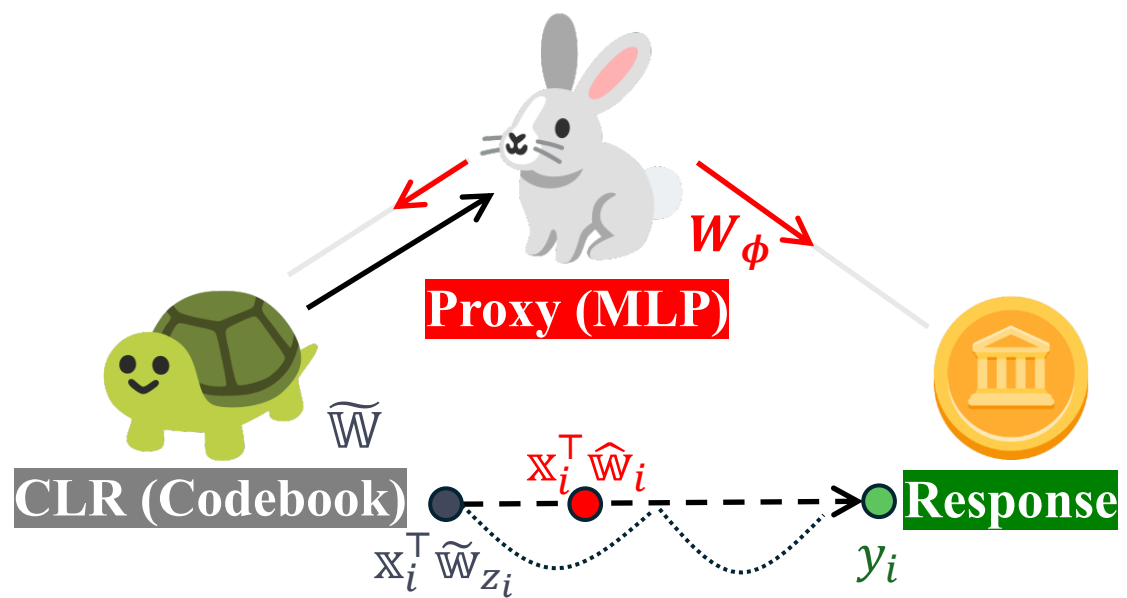
Why does the Optimization stay Stable?



$$\phi^{(t+1)} \leftarrow \phi^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i=1}^N \left(\frac{y_i + \lambda \mathbf{x}_i^T \tilde{W}_{z_i}^{(t)}}{1+\lambda} - \mathbf{x}_i^T \hat{W}_i^{(t)} \right) (\nabla_{\phi} \hat{W}_i^{(t)})^T \mathbf{x}_i,$$

$$\tilde{W}_j^{(t+1)} \leftarrow \tilde{W}_j^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i \in S_j^{(t)}} \mathbf{x}_i \mathbf{x}_i^T (\hat{W}_i^{(t)} - \tilde{W}_j^{(t)}), \quad \forall j \in [K].$$

Why does the Optimization stay Stable?



$$\phi^{(t+1)} \leftarrow \phi^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i=1}^N \left(\frac{y_i + \lambda \mathbf{x}_i^\top \tilde{W}_{z_i}^{(t)}}{1+\lambda} - \mathbf{x}_i^\top \hat{W}_i^{(t)} \right) (\nabla_{\phi} \hat{W}_i^{(t)})^\top \mathbf{x}_i,$$

$$\tilde{W}_j^{(t+1)} \leftarrow \tilde{W}_j^{(t)} + \frac{2\eta(1+\lambda)}{N} \sum_{i \in S_j^{(t)}} \mathbf{x}_i \mathbf{x}_i^\top (\hat{W}_i^{(t)} - \tilde{W}_j^{(t)}), \quad \forall j \in [K].$$

Why does the Optimization stay Stable?

Proposition 3.1 (One-epoch descent). *Under Assumptions 1 and 2, suppose the cluster assignment $z^{(t)}$ remains fixed during epoch t . Define*

$$L_V := \max\left\{(1 + \lambda)L, (4 + 2\lambda)L_x^2 \mathbb{X}_{\max}^2\right\}.$$

For any step-size $0 < \eta \leq 1/L_V$, the simultaneous gradient updates guarantee a strict decrease of the Lyapunov function:

$$V_\lambda(\phi^{(t+1)}, \tilde{W}^{(t+1)}) \leq V_\lambda(\phi^{(t)}, \tilde{W}^{(t)}) - \frac{\eta}{2} \left(\|\phi^{(t+1)} - \phi^{(t)}\|^2 + \|\tilde{W}^{(t+1)} - \tilde{W}^{(t)}\|^2 \right).$$

This explicitly verifies that V_λ serves as a valid Lyapunov function, rigorously ensuring the convergence of the CG-CLR algorithm.

Theorem 3.2 (Linear convergence). *Let Assumptions 1–4 hold. Define*

$$\mu_V := \min\left\{2m_x^2 \mathbb{X}_{\min}^2, (1 + \lambda)\mu\right\}, \quad L_V := \max\left\{(1 + \lambda)L, (4 + 2\lambda)L_x^2 \mathbb{X}_{\max}^2\right\}, \quad q := \frac{L_V - \mu_V}{L_V + \mu_V}.$$

Gradient descent with the optimal step-size $\eta_ = 2/(\mu_V + L_V)$ ensures that the stacked parameter vector $\theta = (\phi, \tilde{W})$ converges linearly at the rate q , while the Lyapunov gap contracts at the rate q^2 .*

CLR (Codebook)

Proxy (MLP)



**“The Proxy follows the interior point of Response and CLR,
and CLR follows the Proxy”**

Check out our paper & Join our poster!

 Scan Me !

