

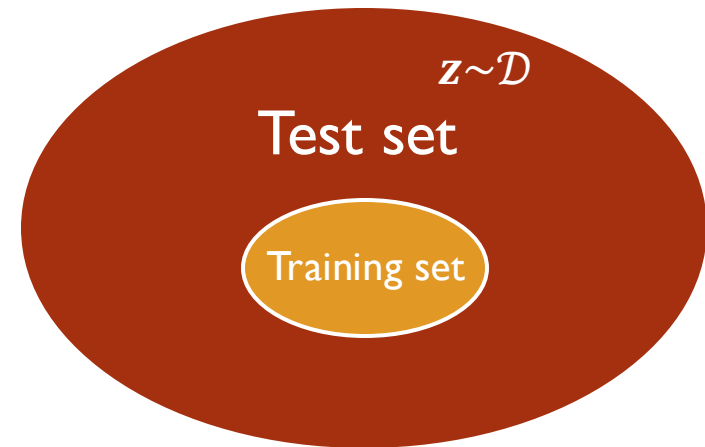
Stability of Gradient-Based Bilevel Algorithms for Hyperparameter Optimization

Rongzhen Wang

2024.3.28

Generalization and stability

- We care about:
 - Expected risk: $R = \mathbb{E}_{z \sim \mathcal{D}}[\ell(A(S); z)]$
- This can be estimated by:
 - Empirical risk: $R_S = \frac{1}{n} \sum_{i=1}^n \ell(A(S); z_i)$
- $R - R_S \rightarrow$ generalization error



Stability and generalization

- Stability:
 - The change in the performance of the algorithm's output if a single data point is replaced
 - $S = \{\mathbf{z}_1, \dots, \mathbf{z}_i, \dots, \mathbf{z}_n\}, \tilde{S}^i = \{\mathbf{z}_1, \dots, \tilde{\mathbf{z}}_i, \dots, \mathbf{z}_n\} \rightarrow \ell(A(\tilde{S}^i)) - \ell(A(S))$

- In expectation, generalization equals stability.

- $$\begin{aligned}\mathbb{E}_S[R - R_S] &= \mathbb{E}_S \left[\mathbb{E}_{\mathbf{z} \sim D}[\ell(A(S); \mathbf{z})] - \frac{1}{n} \sum_{i=1}^n \ell(A(S); \mathbf{z}_i) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{S, \tilde{S}^i} [\ell(A(\tilde{S}^i); \mathbf{z}_i) - \ell(A(S); \mathbf{z}_i)]\end{aligned}$$

Recap: Lipschitz conditions

- L -Lipschitz continuous:

$$\forall \mathbf{u}, \mathbf{v} \in \Omega, \|f(\mathbf{u}) - f(\mathbf{v})\| \leq L\|\mathbf{u} - \mathbf{v}\| \text{ or } \forall \mathbf{x} \in \Omega, \|\nabla f(\mathbf{x})\| \leq L.$$

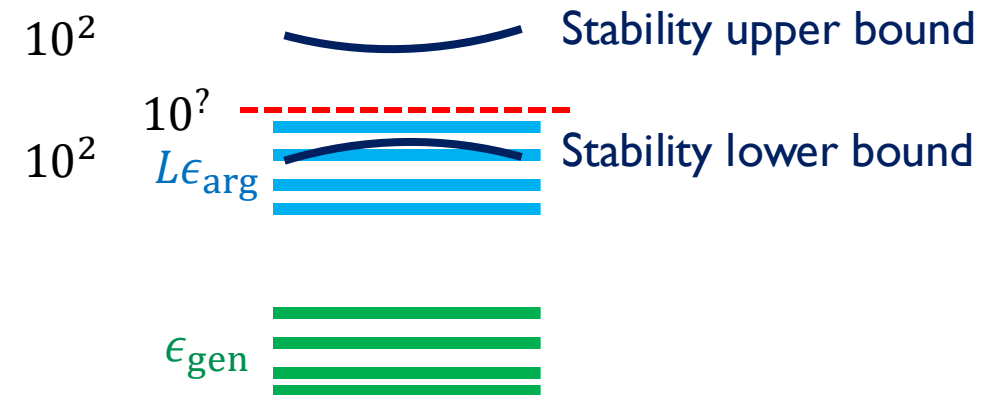
- γ -smooth (i.e., gradient Lipschitz continuous):

$$\forall \mathbf{u}, \mathbf{v} \in \Omega, \|\nabla f(\mathbf{u}) - \nabla f(\mathbf{v})\| \leq \gamma.$$

- $\mathbb{E}_S[R - R_S] \leq \sup_{S, \tilde{S}, \mathbf{z}} |\ell(A(\tilde{S}); \mathbf{z}) - \ell(A(S); \mathbf{z})| \rightarrow \text{uniform stability } \epsilon_{\text{arg}}$
 $\leq \sup_{S, \tilde{S}} L \|A(S) - A(\tilde{S})\| \rightarrow \text{uniform argument stability } \epsilon_{\text{stab}}$

Stability based generalization bound for Lipschitz loss functions

1. Generalization upper bound via stability: $\epsilon_{\text{gen}}(A) \leq \epsilon_{\text{stab}}(A) \leq L\epsilon_{\text{arg}}(A)$
2. Stability upper bound
 - About key factors: sample size, learning rate, number of steps
3. Stability lower bound
 - Construct an example



Outline

- Background on hyperparameter optimization (HO)
- Generalization and stability of HO algorithms
- Stability bounds of gradient-based HO algorithms
- Conclusion and discussion

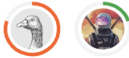




Hyperparameter optimization (HO)

- Hyperparameters:
 - `num_layers`, `lr` ...
 - Influential to training efficiency and results
 - Fixed during the training phase
- Three phases in machine learning
 - training
 - validation
 - test

Public **Private**

The private leaderboard is calculated with approximately 33% of the test data. This competition has completed. This leaderboard reflects the final standings.

■ Prize Winners

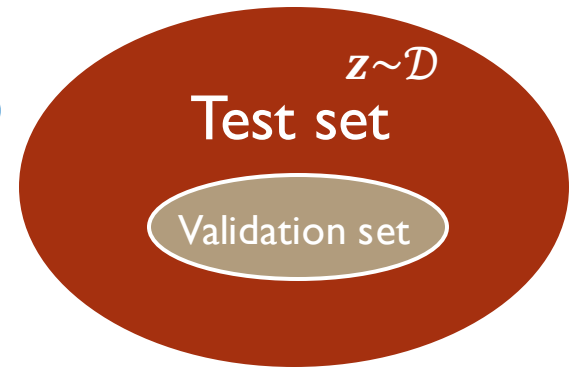
#	△	Team	Members
1	—	Clevert	
2	▲ 1022	ryo	
3	▲ 568	ForcewithMe	
4	▲ 1	Igor Krashenyi	
5	▲ 5	Ivan Panshin	

Generalization and stability of HO

- Recap: Single-level optimization

- Generalization via uniform argument stability

$$\epsilon_{\text{gen}}(A) \triangleq \mathbb{E}_S[R - R_S] \leq \sup_{S, \tilde{S}} L \mathbb{E}_A[\|A(S) - A(\tilde{S})\|] \triangleq L\epsilon_{\text{arg}}(A)$$



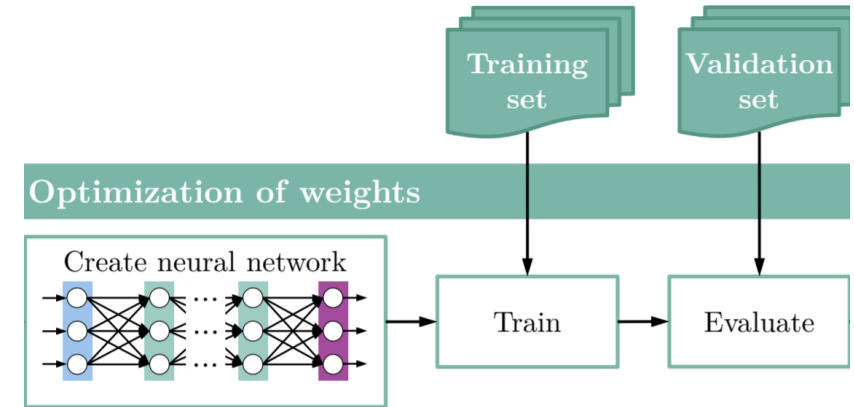
- Hyperparameter optimization:

- Generalization on validation via uniform argument stability

$$\epsilon_{\text{gen}}(A) \triangleq \mathbb{E}_{S^{\text{val}}, S^{\text{tr}}}[R - R_{S^{\text{val}}}] \leq \sup_{S^{\text{val}}, \tilde{S}^{\text{val}}, S^{\text{tr}}} L \mathbb{E}_A[\|A(S^{\text{val}}, S^{\text{tr}}) - A(\tilde{S}^{\text{val}}, S^{\text{tr}})\|] \triangleq L\epsilon_{\text{arg}}(A)$$

Formulation of HO

- Model parameters θ , hyperparameters λ
- Training loss ℓ^{tr} , validation loss ℓ^{val}
- Training set $S^{\text{tr}} = \{\mathbf{z}_i^{\text{tr}}\}_{i=1}^n$, validation set $S^{\text{val}} = \{\mathbf{z}_i^{\text{val}}\}_{i=1}^m$
- Example $\mathbf{z} = (\mathbf{x}, \mathbf{y}) \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$



Outer level: $\lambda^* = \operatorname{argmin}_{\lambda \in \Lambda} R_{S^{\text{val}}}(\lambda) \triangleq \frac{1}{m} \sum_{i=1}^m \ell^{\text{val}}(\lambda, \theta^*(\lambda); \mathbf{z}_i^{\text{val}}) \mathcal{L}^{\text{val}}(\lambda; \mathbf{z}_i^{\text{val}})$

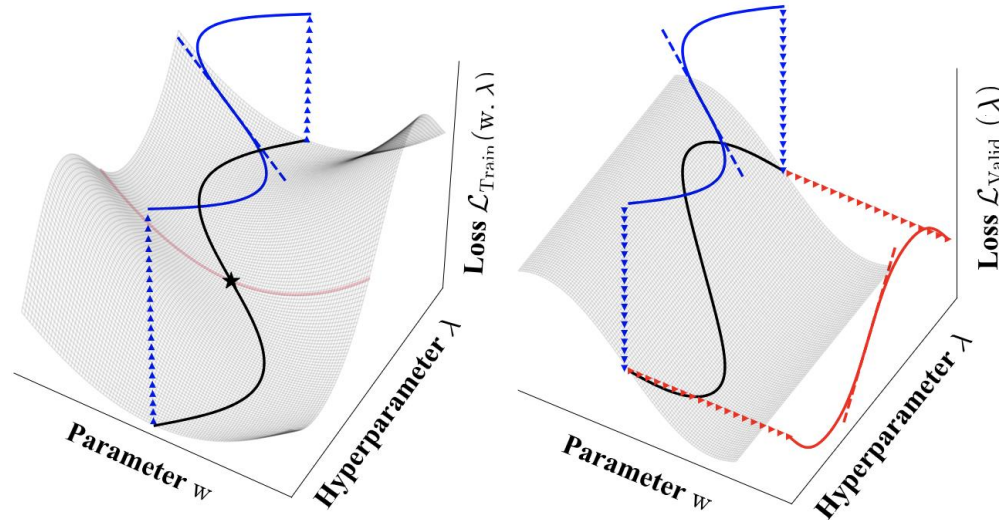
Inner level: where $\theta^*(\lambda) = \operatorname{argmin}_{\theta \in \Theta} R_{S^{\text{tr}}}(\lambda, \theta) \triangleq \frac{1}{n} \sum_{i=1}^n \ell^{\text{tr}}(\lambda, \theta; \mathbf{z}_i^{\text{tr}})$

$$\lambda^* = \operatorname{argmin}_{\lambda \in \Lambda} \frac{1}{m} \sum_{i=1}^m \mathcal{L}^{\text{val}}(\lambda; \mathbf{z}_i^{\text{val}})$$

Approaches for solving HO

Gradient-based methods

- How to get the hypergradient $\nabla_{\lambda} \mathcal{L}^{\text{val}}$
 - Modeling \mathcal{L}^{val} to be differentiable w.r.t. λ
 - e.g., architecture search, data reweighting
 - Calculate $\nabla_{\lambda} \mathcal{L}^{\text{val}}$
 - difficult for $\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda; \mathbf{z}) \neq \nabla_{\lambda} \ell^{\text{val}}(\lambda, \hat{\theta}; \mathbf{z})$



Algorithm 1 Gradient-based bilevel HO

Input: Initialization λ_0 and θ_0 ; training set S^{tr} and validation set S^{val} ; learning rate scheme α and η .

```

for  $t = 1$  to  $T$  do
    for  $k = 1$  to  $K$  do
        uniformly sampling  $j_k$  from  $[n]$ 
         $\theta_k \leftarrow \theta_{k-1} - \eta_k \nabla_{\theta} \ell^{\text{tr}}(\lambda_{t-1}, \theta_{k-1}; z_{j_k}^{\text{tr}})$ 
    end for
    uniformly sampling  $i_t$  from  $[m]$ 
     $\lambda_t \leftarrow \lambda_{t-1} - \alpha_t \nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; z_{i_t}^{\text{val}})$ 
end for
return  $\lambda_T$  and  $\theta_K$ 
  
```

Outer optimization

Figure credit: Lorraine J, Vicol P, Duvenaud D. Optimizing millions of hyperparameters by implicit differentiation[C]//International conference on artificial intelligence and statistics. PMLR, 2020: 1540-1552.

Approaches for solving HO

Gradient-based methods

- $\nabla_x f(\mathbf{u}(x), \mathbf{v}(x)) = \nabla_{\mathbf{u}} f(\mathbf{u}, \mathbf{v}) \nabla_x \mathbf{u}(x) + \nabla_{\mathbf{v}} f(\mathbf{u}, \mathbf{v}) \nabla_x \mathbf{v}(x)$
- $\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda; \mathbf{z}) = \nabla_{\lambda} \ell^{\text{val}}(\lambda, \hat{\boldsymbol{\theta}}(\lambda); \mathbf{z}) = \underbrace{\nabla_{\lambda} \ell^{\text{val}}(\lambda, \hat{\boldsymbol{\theta}}; \mathbf{z})}_{\text{Explicit, easy to be auto-calculated}} + \underbrace{\nabla_{\boldsymbol{\theta}} \ell^{\text{val}}(\lambda, \hat{\boldsymbol{\theta}}; \mathbf{z}) \nabla_{\lambda} \hat{\boldsymbol{\theta}}(\lambda)}_{\text{Implicit, involving the inner optimization}}$
- An example of optimizing the regularization coefficient in ridge regression
 - $\mathcal{L}^{\text{val}}(\lambda; \mathbf{z}) = \frac{1}{2} \left(y - \mathbf{x}^T \hat{\boldsymbol{\theta}}(\lambda) \right)^2 + \frac{1}{2} \lambda \|\hat{\boldsymbol{\theta}}(\lambda)\|_2^2$
 - $\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda; \mathbf{z}) = \frac{1}{2} \|\hat{\boldsymbol{\theta}}(\lambda)\|_2^2 + \left[\left(y - \mathbf{x}^T \hat{\boldsymbol{\theta}}(\lambda) \right) (-\mathbf{x}) + \lambda \hat{\boldsymbol{\theta}}(\lambda) \right] \nabla_{\lambda} \hat{\boldsymbol{\theta}}(\lambda)$
- A direct approach: unrolling differentiation
 - $\hat{\boldsymbol{\theta}}_k(\lambda) = \hat{\boldsymbol{\theta}}_{k-1}(\lambda) - \eta \nabla_{\boldsymbol{\theta}} \ell^{\text{tr}}(\lambda, \hat{\boldsymbol{\theta}}_{k-1}(\lambda)),$
 - $\nabla_{\lambda} \hat{\boldsymbol{\theta}}_k(\lambda) = \nabla_{\lambda} \hat{\boldsymbol{\theta}}_{k-1}(\lambda) - \eta \nabla^2_{\boldsymbol{\theta} \lambda} \ell^{\text{tr}}(\lambda, \hat{\boldsymbol{\theta}}_{k-1}) - \eta \nabla^2_{\boldsymbol{\theta} \boldsymbol{\theta}} \ell^{\text{tr}}(\lambda, \hat{\boldsymbol{\theta}}_{k-1}) \nabla_{\lambda} \hat{\boldsymbol{\theta}}_{k-1}(\lambda)$

Tight stability bounds

In the case of the constructed example

- When $\alpha_t = c/t$,

- **Thm 4** (Deformed argument stability bounds, [2]).

$\mathcal{L}^{\text{val}}(\lambda)$ is L -Lipschitz continuous and γ -smooth, the outer step size is $\alpha_t \leq \frac{c}{t}$, $\ell^{\text{tr}}(\lambda, \theta)$ is γ^{tr} -smooth w.r.t. θ and the inner step size $\eta_k = \eta$, then for all S^{tr} , S^{val} and \tilde{S}^{val} , we have

$$\Omega\left(\frac{T^{\ln(1+(1-\frac{1}{m})c\gamma)}}{m}\right) \leq \epsilon_{\text{arg}} \leq \mathcal{O}\left(\frac{T^{(1-\frac{1}{m})c\gamma}}{m}\right),$$

where $\gamma = \Theta((1 + \eta\gamma^{\text{tr}})^{2K})$.

10^2

$10^?$

10^2



Upper bound of ϵ_{arg}

ϵ_{arg}

Lower bound of ϵ_{arg}

Conclusion

- The fundamental relation between generalization and stability:
 - In expectation, generalization equals stability
 - $\epsilon_{\text{gen}}(A) \leq \epsilon_{\text{stab}}(A) \leq L\epsilon_{\text{arg}}(A)$
- Hyperparameter optimization (HO) problems and gradient-based algorithms
 - The core is to calculate the hypergradient
- Matching argument stability upper and lower bounds
 - $\Omega\left(\frac{T^{\ln\left(1+\left(1-\frac{1}{m}\right)c\gamma\right)}}{m}\right) \leq \epsilon_{\text{arg}} \leq \mathcal{O}\left(\frac{T^{\left(1-\frac{1}{m}\right)c\gamma}}{m}\right)$, where $\gamma = \Theta((1 + \eta\gamma^{tr})^{2K})$
 - The current upper bound cannot be improved

Stability upper bound

Recursion formula

- $\epsilon_{\arg}(A) = \sup_{S^{\text{val}}, \tilde{S}^{\text{val}}, S^{\text{tr}}} L \mathbb{E}_A[\|A(S^{\text{val}}, S^{\text{tr}}) - A(\tilde{S}^{\text{val}}, S^{\text{tr}})\|] = \sup_{S^{\text{val}}, \tilde{S}^{\text{val}}, S^{\text{tr}}} \mathbb{E}_A[\|\lambda_T - \tilde{\lambda}_T\|]$

- $\mathbb{E}_A[\|\lambda_T - \tilde{\lambda}_T\|] = \mathbb{E}_{\mathcal{A}}[\|\lambda_{t-1} - \alpha_t \nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - (\tilde{\lambda}_{t-1} - \alpha_t \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}}))\|]$
 $\leq \mathbb{E}_A[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \alpha_t \mathbb{E}_A[\|\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}})\|]$

If $\mathbf{z}_{i_t}^{\text{val}} = \tilde{\mathbf{z}}_{i_t}^{\text{val}}$, leverage smoothness
 $\leq \gamma \|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|$

If $\mathbf{z}_{i_t}^{\text{val}} \neq \tilde{\mathbf{z}}_{i_t}^{\text{val}}$, leverage continuity
 $\leq 2L$

$$\leq \mathbb{E}_A[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \alpha_t \left[\left(1 - \frac{1}{m}\right) \mathbb{E}_A[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \frac{1}{m} 2L \right]$$

$$= \left[1 + \alpha_t \left(1 - \frac{1}{m}\right) \gamma \right] \mathbb{E}_A[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \alpha_t \frac{1}{m} 2L$$

Stability upper bound

- Unroll the recursion:

- **Thm 1** (Argument stability upper bound, [1]).

Suppose $\mathcal{L}^{\text{val}}(\boldsymbol{\lambda})$ is L -Lipschitz continuous and γ -smooth, the outer step size is $\alpha_t \leq \frac{c}{t}$, the gradient of $\ell^{\text{tr}}(\boldsymbol{\lambda}, \boldsymbol{\theta})$ is γ^{tr} -smooth w.r.t. $\boldsymbol{\theta}$ and the inner step size $\eta_k = \eta$, then for all $S^{\text{tr}}, S^{\text{val}}$ and \tilde{S}^{val} , we have

$$\epsilon_{\text{arg}} \leq \sum_{t=1}^T \prod_{s=t+1}^{T+1} \left(1 + \left(1 - \frac{1}{m} \right) \alpha_s \gamma \right) \frac{2\alpha_t L}{m},$$

where $L = \mathcal{O}((1 + \eta\gamma^{\text{tr}})^K)$, $\gamma = \mathcal{O}((1 + \eta\gamma^{\text{tr}})^{2K})$.

- Considering T and m : When $\alpha_t = \alpha$, $\epsilon_{\text{arg}} = \mathcal{O} \left(\left(1 + \left(1 - \frac{1}{m} \right) \alpha \gamma \right)^T / m \right)$

Stability upper bound

- When $\alpha_t = c/t$,

- **Thm 2** (Deformed argument stability upper bound, [2]).

Suppose $\mathcal{L}^{\text{val}}(\boldsymbol{\lambda})$ is L -Lipschitz continuous and γ -smooth, the outer step size is $\alpha_t \leq \frac{c}{t}$, the gradient of $\ell^{\text{tr}}(\boldsymbol{\lambda}, \boldsymbol{\theta})$ is γ^{tr} -smooth w.r.t. $\boldsymbol{\theta}$ and the inner step size $\eta_k = \eta$, then for all $S^{\text{tr}}, S^{\text{val}}$ and \tilde{S}^{val} , we have

$$\epsilon_{\text{arg}} \leq \left(1 + \left(1 - \frac{1}{m}\right) c\gamma\right) T^{\left(1 - \frac{1}{m}\right) c\gamma} \frac{2L}{(m-1)\gamma} = \mathcal{O}\left(\frac{T^{\left(1 - \frac{1}{m}\right) c\gamma}}{m}\right),$$

where $\gamma = \mathcal{O}\left((1 + \eta\gamma^{\text{tr}})^{2K}\right)$.

Deformed stability upper bound

$$\begin{aligned}
\epsilon_{\text{arg}} &\leq \sum_{t=1}^T \prod_{s=t+1}^T (1 + \alpha_s(1 - 1/m)\gamma) \frac{2\alpha_t L}{m} \\
&\leq \sum_{t=1}^{T-1} \prod_{s=t+1}^T \exp\left[\left(1 - \frac{1}{m}\right) \frac{\gamma c}{s}\right] \frac{2cL}{tm} + \frac{2cL}{Tm} \\
&= \sum_{t=1}^{T-1} \exp\left[\left(1 - \frac{1}{m}\right) \gamma c \sum_{s=t+1}^T \frac{1}{s}\right] \frac{2cL}{tm} + \frac{2cL}{Tm} \\
&\leq \sum_{t=1}^{T-1} \exp\left[\left(1 - \frac{1}{m}\right) \gamma c \ln \frac{T}{t}\right] \frac{2cL}{tm} + \frac{2cL}{Tm} & (\forall t_2 > t_1 > 1, \sum_{t=t_1}^{t_2} \frac{1}{t} \leq \ln \frac{t_2}{t_1-1}) \\
&= \sum_{t=1}^T \left(\frac{T}{t}\right)^{(1-1/m)\gamma c} \frac{2cL}{tm} \\
&= \frac{2cL}{m} T^{(1-1/m)\gamma c} \sum_{t=1}^T t^{-(1-1/m)\gamma c-1} \\
&\leq \frac{2cL}{m} T^{(1-1/m)\gamma c} \left(1 + \int_1^T t^{-(1-1/m)\gamma c-1} dt\right) & (\forall a > 0, \sum_{t=1}^T t^{-a-1} \leq 1 + \int_1^T t^{-a-1} dt) \\
&= \frac{2cL}{m(1-1/m)c\gamma} \left[\left(1 + (1-1/m)\gamma c\right) T^{(1-1/m)\gamma c} - 1 \right] \\
&= \frac{2L}{(m-1)\gamma} \left[\left(1 + (1-1/m)\gamma c\right) T^{(1-1/m)\gamma c} - 1 \right] \\
&= \mathcal{O}\left(\frac{T^{(1-\frac{1}{m})c\gamma}}{m}\right),
\end{aligned}$$

Integral and exponential inequalities

Stability upper bound

- $\epsilon_{\arg}(A) = \sup_{S^{\text{val}}, \tilde{S}^{\text{val}}, S^{\text{tr}}} L \mathbb{E}_A[\|A(S^{\text{val}}, S^{\text{tr}}) - A(\tilde{S}^{\text{val}}, S^{\text{tr}})\|] = \sup_{S^{\text{val}}, \tilde{S}^{\text{val}}, S^{\text{tr}}} \mathbb{E}_A[\|\lambda_T - \tilde{\lambda}_T\|]$

- $\mathbb{E}_A[\|\lambda_T - \tilde{\lambda}_T\|] = \mathbb{E}_{\mathcal{A}}[\|\lambda_{t-1} - \alpha_t \nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - (\tilde{\lambda}_{t-1} - \alpha_t \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}}))\|]$

① $\leq \mathbb{E}_A[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \alpha_t \mathbb{E}_A[\|\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}})\|]$

If $\mathbf{z}_{i_t}^{\text{val}} = \tilde{\mathbf{z}}_{i_t}^{\text{val}}$, leverage smoothness

② $\leq \gamma \|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|$

If $\mathbf{z}_{i_t}^{\text{val}} \neq \tilde{\mathbf{z}}_{i_t}^{\text{val}}$, leverage continuity

③ $\leq 2L$

$$\leq \mathbb{E}_A[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \alpha_t \left[\left(1 - \frac{1}{m}\right) \mathbb{E}_A[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \frac{1}{m} 2L \right]$$

$$= \left[1 + \alpha_t \left(1 - \frac{1}{m}\right) \gamma \right] \mathbb{E}_A[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \alpha_t \frac{1}{m} 2L$$

Equality conditions

$$\begin{aligned} & \mathbb{E}_A \left[\left\| \lambda_{t-1} - \tilde{\lambda}_{t-1} - \alpha_t \left(\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}}) \right) \right\| \right] \\ & \leq \mathbb{E}_A \left[\left\| \lambda_{t-1} - \tilde{\lambda}_{t-1} \right\| \right] + \alpha_t \mathbb{E}_A \left[\left\| \nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}}) \right\| \right] \end{aligned}$$

- ① Triangle inequality

- $\lambda_{t-1} - \tilde{\lambda}_{t-1}$ and $-\left(\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}})\right)$ are parallel and in the same direction
 - $\lambda_0 = \tilde{\lambda}_0, \lambda_1 - \tilde{\lambda}_1 = -\alpha_1 \left(\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_0; \mathbf{z}_1^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_0; \tilde{\mathbf{z}}_1^{\text{val}}) \right) \doteq \mathbf{v}$
 - $\mathbf{v}_1 \doteq \mathbf{v}_2$, iff $\exists a \geq 0$, s. t. $\mathbf{v}_1 = a\mathbf{v}_2$
- \leftarrow For all $1 \leq t \leq T$, $-\left(\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}})\right) \doteq \mathbf{v}$, so that we also have $\lambda_t - \tilde{\lambda}_t \doteq \mathbf{v}$

Equality conditions

$$\text{If } \mathbf{z}_{i_t}^{\text{val}} = \tilde{\mathbf{z}}_{i_t}^{\text{val}}, \|\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \mathbf{z}_{i_t}^{\text{val}})\| \leq \gamma \|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|$$

- **② Smoothness:**

- Based on ①, $\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) = -\gamma(\lambda_{t-1} - \tilde{\lambda}_{t-1})$ for $\lambda_t - \tilde{\lambda}_t \triangleq \mathbf{v}$
 - The gradient of \mathcal{L}^{val} is linear w.r.t. λ
- $\leftarrow \mathcal{L}^{\text{val}}$ is quadratic w.r.t. λ with term: $\frac{1}{2} \lambda^T H \lambda$, where $H\mathbf{v} = -\gamma\mathbf{v}$

Equality conditions

$$\text{If } \mathbf{z}_{i_t}^{\text{val}} \neq \tilde{\mathbf{z}}_{i_t}^{\text{val}}, \|\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}})\| \leq 2L$$

• ③ Continuity:

- $\|\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}})\| = \|\nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}})\| = L = \sup_{\lambda, \mathbf{z}} \|\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda; \mathbf{z})\|$
 - Contradictory to ②: require \mathcal{L}^{val} to be linear
 - $\|\nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}})\| \geq 2L', \text{ where } L' \leq L$
- Based on ① and ②,

$$\begin{aligned} & \nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}}) \\ &= \nabla_{\lambda} \mathcal{L}^{\text{val}}(\lambda_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) + \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}}) \\ &= -\gamma(\lambda_{t-1} - \tilde{\lambda}_{t-1}) + \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \mathbf{z}_{i_t}^{\text{val}}) - \nabla_{\lambda} \mathcal{L}^{\text{val}}(\tilde{\lambda}_{t-1}; \tilde{\mathbf{z}}_{i_t}^{\text{val}}) \end{aligned}$$
 - \mathcal{L}^{val} has a linear cross term of λ and \mathbf{z}
- $\leftarrow \mathcal{L}^{\text{val}}$ has a linear term: $-y\mathbf{x}^T \lambda$ and $y_m \mathbf{x}_m - \tilde{y}_m \tilde{\mathbf{x}}_m \doteq \mathbf{v}$

Constructed example

- The composite validation loss: $\mathcal{L}^{\text{val}}(\boldsymbol{\lambda}; \mathbf{z}) = \frac{1}{2} \boldsymbol{\lambda}^T \mathbf{H} \boldsymbol{\lambda} - y \mathbf{x}^T \boldsymbol{\lambda}$
- \mathbf{H} is a real symmetric matrix
 - With the maximum absolute eigenvalue γ
 - With eigenvalue $-\gamma$ and a corresponding unit eigenvector \mathbf{v}
- S^{val} and \tilde{S}^{val} only differ in the last example, where $\mathbf{z}_m^{\text{val}} = (\mathbf{x}, y) = (\mathbf{v}, 1)$ and $\tilde{\mathbf{z}}_m^{\text{val}} = (\tilde{\mathbf{x}}, \tilde{y}) = (-\mathbf{v}, 1)$

Constructed example

- An HO problem with: $\ell^{\text{val}}(\boldsymbol{\lambda}, \boldsymbol{\theta}; \mathbf{z}) = \ell^{\text{tr}}(\boldsymbol{\lambda}, \boldsymbol{\theta}; \mathbf{z}) = \frac{1}{2} \boldsymbol{\theta}^T \mathbf{A} \boldsymbol{\theta} + \boldsymbol{\lambda}^T \boldsymbol{\theta} - y \mathbf{x}^T \boldsymbol{\theta}$
- \mathbf{A} is a real symmetric matrix
 - With the maximum absolute eigenvalue γ_A
 - With eigenvalue $-\gamma_A$ and a corresponding unit eigenvector \mathbf{v}
- S^{val} and \tilde{S}^{val} only differ in the last example, where $\mathbf{z}_m^{\text{val}} = (\mathbf{x}, y) = (\mathbf{v}, 1)$ and $\tilde{\mathbf{z}}_m^{\text{val}} = (\tilde{\mathbf{x}}, \tilde{y}) = (-\mathbf{v}, 1)$

Stability lower bound

$$\begin{aligned}
 \bullet \quad \mathbb{E}_{\mathcal{A}}[\|\lambda_t - \tilde{\lambda}_t\|] &= \mathbb{E}_{\mathcal{A}}[\|\lambda_{t-1} - \alpha_t(\mathbf{H}\lambda_{t-1} - y\mathbf{x}_{i_t}) - (\tilde{\lambda}_{t-1} - \alpha_t(\mathbf{H}\lambda_{t-1} - y\tilde{\mathbf{x}}_{i_t}))\|] \\
 &= \mathbb{E}_{\mathcal{A}}[\|\lambda_{t-1} - \tilde{\lambda}_{t-1} - \alpha_t[\mathbf{H}(\lambda_{t-1} - \tilde{\lambda}_{t-1}) + (y\mathbf{x}_{i_t} - y\tilde{\mathbf{x}}_{i_t})]\|] \\
 \textcircled{1} \quad &= \mathbb{E}_{\mathcal{A}}[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \underbrace{\alpha_t \mathbb{E}_{\mathcal{A}}[\|-\mathbf{H}(\lambda_{t-1} - \tilde{\lambda}_{t-1})\| + \|y\mathbf{x}_{i_t} - y\tilde{\mathbf{x}}_{i_t}\|]}_{\substack{\text{If } \mathbf{z}_{i_t}^{\text{val}} = \tilde{\mathbf{z}}_{i_t}^{\text{val}}, & \text{If } \mathbf{z}_{i_t}^{\text{val}} \neq \tilde{\mathbf{z}}_{i_t}^{\text{val}} \\ \textcircled{2} \quad = \gamma\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\| & \textcircled{3} \quad \geq 2L'}} \\
 &\geq \mathbb{E}_{\mathcal{A}}[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \alpha_t \left[\left(1 - \frac{1}{m}\right) \gamma \mathbb{E}_{\mathcal{A}}[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \frac{1}{m} 2L' \right] \\
 &= \left[1 + \alpha_t \left(1 - \frac{1}{m}\right) \gamma \right] \mathbb{E}_{\mathcal{A}}[\|\lambda_{t-1} - \tilde{\lambda}_{t-1}\|] + \alpha_t \frac{1}{m} 2L'
 \end{aligned}$$

Tight stability bounds

In the case of the constructed example

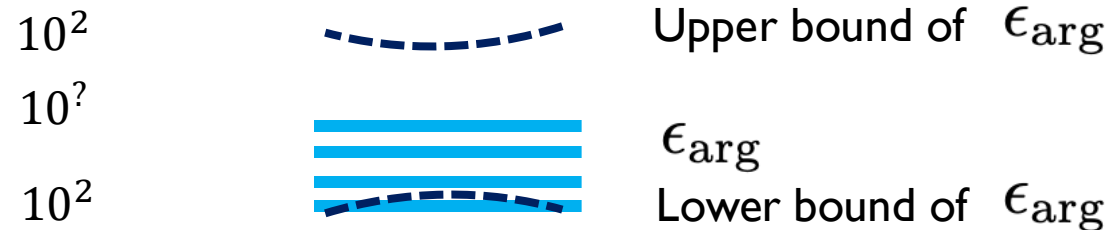
- Unroll the recursion:

- **Thm 3** (Argument stability bounds, [2]).

$\mathcal{L}^{\text{val}}(\boldsymbol{\lambda})$ is L -Lipschitz continuous and γ -smooth, the outer step size is α_t , $\ell^{\text{tr}}(\boldsymbol{\lambda}, \boldsymbol{\theta})$ is γ^{tr} -smooth w.r.t. $\boldsymbol{\theta}$ and the inner step size $\eta_k = \eta$, then for all $S^{\text{tr}}, S^{\text{val}}$ and \tilde{S}^{val} , we have

$$\sum_{t=1}^T \prod_{s=t+1}^{T+1} \left(1 + \left(1 - \frac{1}{m} \right) \alpha_s \gamma \right) \frac{2\alpha_t L'}{m} \leq \epsilon_{\text{arg}} \leq \sum_{t=1}^T \prod_{s=t+1}^{T+1} \left(1 + \left(1 - \frac{1}{m} \right) \alpha_s \gamma \right) \frac{2\alpha_t L}{m},$$

where $L' \simeq L = \Theta((1 + \eta\gamma^{\text{tr}})^K)$, $\gamma = \Theta((1 + \eta\gamma^{\text{tr}})^{2K})$.



Tight stability bounds

In the case of the constructed example

- When $\alpha_t = c/t$,

• **Thm 4** (Deformed argument stability bounds, [2]).

$\mathcal{L}^{\text{val}}(\boldsymbol{\lambda})$ is L -Lipschitz continuous and γ -smooth, the outer step size is $\alpha_t \leq \frac{c}{t}$, $\ell^{\text{tr}}(\boldsymbol{\lambda}, \boldsymbol{\theta})$ is γ^{tr} -smooth w.r.t. $\boldsymbol{\theta}$ and the inner step size $\eta_k = \eta$, then for all S^{tr} , S^{val} and \tilde{S}^{val} , we have

$$\Omega\left(\frac{T^{\ln(1+(1-\frac{1}{m})c\gamma)}}{m}\right) \leq \epsilon_{\text{arg}} \leq \mathcal{O}\left(\frac{T^{(1-\frac{1}{m})c\gamma}}{m}\right),$$

where $\gamma = \Theta((1 + \eta\gamma^{\text{tr}})^{2K})$.

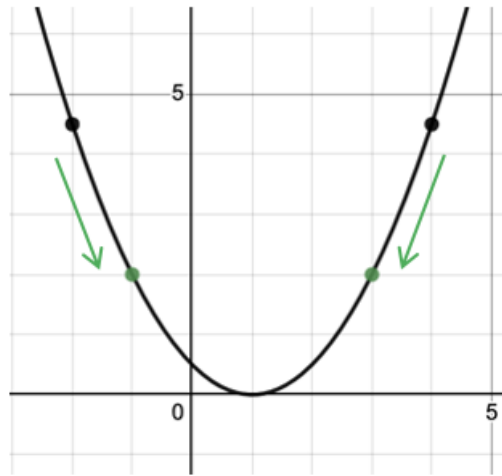
Deformed stability upper bound

$$\begin{aligned}
\epsilon_{\arg} &\geq \sum_{t=1}^T \prod_{s=t+1}^T (1 + \alpha_s(1 - 1/m)\gamma') \frac{2\alpha_t L'}{m} \\
&\geq \sum_{t=1}^T \prod_{s=t+1}^T \exp\left[r\left(1 - \frac{1}{m}\right)\alpha_s\gamma'\right] \frac{2\alpha_t L'}{m} \\
&= \sum_{t=1}^{T-1} \prod_{s=t+1}^T \exp\left[r\left(1 - \frac{1}{m}\right)\frac{c\gamma'}{s}\right] \frac{2cL'}{tm} + \frac{2cL'}{Tm} \quad (\alpha_t = c/t) \\
&= \sum_{t=1}^{T-1} \exp\left[r\left(1 - \frac{1}{m}\right)c\gamma' \sum_{s=t+1}^T \frac{1}{s}\right] \frac{2cL'}{tm} + \frac{2cL'}{Tm} \\
&\geq \sum_{t=1}^{T-1} \exp\left[r\left(1 - \frac{1}{m}\right)c\gamma' \ln \frac{T+1}{t+1}\right] \frac{2cL'}{tm} + \frac{2cL'}{Tm} \quad (\forall t_2 > t_1 > 0, \sum_{t=t_1}^{t_2} \frac{1}{t} \geq \ln \frac{t_2+1}{t_1}) \\
&= \sum_{t=1}^T \left(\frac{T+1}{t+1}\right)^{r(1-1/m)c\gamma'} \frac{2cL'}{tm} \\
&\geq \frac{2cL'}{m} (T+1)^{r(1-1/m)c\gamma'} \sum_{t=2}^{T+1} t^{-r(1-1/m)c\gamma'-1} \\
&\geq \frac{2cL'}{m} (T+1)^{r(1-1/m)c\gamma'} \int_2^{T+2} t^{-r(1-1/m)c\gamma'-1} dt \quad (\forall a > 0, \sum_{t=2}^{T+1} t^{-a-1} \geq \int_2^{T+2} t^{-a-1} dt) \\
&= \frac{2cL'}{m} (T+1)^{r(1-1/m)c\gamma'} \left[\frac{2^{-r(1-1/m)c\gamma'} - (T+2)^{-r(1-1/m)c\gamma'}}{r(1-1/m)c\gamma'} \right] \\
&= \frac{2L'}{r(m-1)\gamma'} (T+1)^{r(1-1/m)c\gamma'} \left[2^{-r(1-1/m)c\gamma'} - (T+2)^{-r(1-1/m)c\gamma'} \right] \\
&\geq \frac{2cL'}{m \ln(1 + (1-1/m)c\gamma')} \left[\left(\frac{T+1}{2}\right)^{\ln(1+(1-1/m)c\gamma')} - 1 \right], \quad (r = \frac{\ln(1+(1-1/m)c\gamma')}{(1-1/m)c\gamma'}) \\
&= \Theta\left(\frac{T^{\ln(1+(1-1/m)c\gamma')}}{m}\right).
\end{aligned}$$

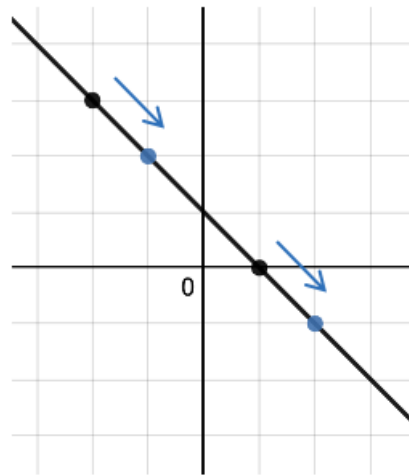
Conclusion

- The fundamental relation between generalization and stability:
 - In expectation, generalization equals stability
 - $\epsilon_{\text{gen}}(A) \leq \epsilon_{\text{stab}}(A) \leq L\epsilon_{\text{arg}}(A)$
- Hyperparameter optimization (HO) problems and gradient-based algorithms
 - The core is to calculate the hypergradient
- Matching argument stability upper and lower bounds
 - $\Omega\left(\frac{T^{\ln\left(1+\left(1-\frac{1}{m}\right)c\gamma\right)}}{m}\right) \leq \epsilon_{\text{arg}} \leq \mathcal{O}\left(\frac{T^{\left(1-\frac{1}{m}\right)c\gamma}}{m}\right)$, where $\gamma = \Theta((1 + \eta\gamma^{tr})^{2K})$
 - The current upper bound cannot be improved

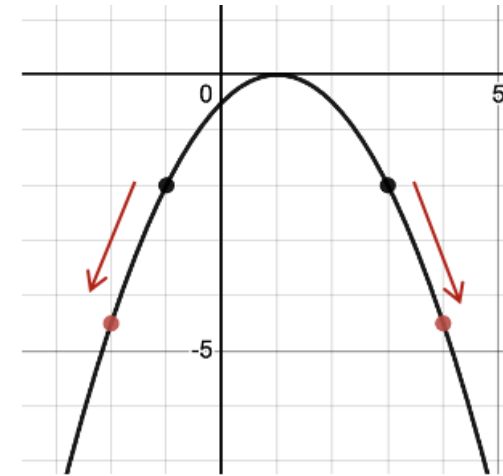
Recap: Uniform stability of SGD



strongly convex



convex



non-convex