

Setting

Teacher function:

$$y^\nu = f(\mathbf{x}^\nu, \mathbf{W}^*) + \sqrt{\Delta}\zeta,$$

with $\mathbf{W}^* \in \mathbb{R}^{d \times k}$ and

$$f(\mathbf{x}^\nu, \mathbf{W}^*) = \frac{1}{k} \sum_{r=1}^k \sigma \left(\frac{\mathbf{w}_r^{*\top} \mathbf{x}^\nu}{\sqrt{d}} \right).$$

Learned by a *student* two-layer neural network with weights $\mathbf{W} \in \mathbb{R}^{d \times p}$

$$\hat{f}(\mathbf{x}, \mathbf{W}) = \frac{1}{p} \sum_{j=1}^p \sigma \left(\frac{\mathbf{w}_j^\top \mathbf{x}}{\sqrt{d}} \right)$$

through Stochastic Gradient Descent (SGD):

$$\mathbf{w}_j^{\nu+1} = \mathbf{w}_j^\nu - \gamma \nabla_{\mathbf{w}_j} (y^\nu - \hat{f}(\mathbf{x}^\nu, \mathbf{W}))^2$$

SGD aims to directly minimize the *population risk* \mathcal{R} :

$$\mathcal{R}(\mathbf{W}, \mathbf{W}^*) \equiv \mathbb{E}_x \left[\left(f(\mathbf{x}, \mathbf{W}^*) - \hat{f}(\mathbf{x}, \mathbf{W}) \right)^2 \right]$$

Local fields and overlaps

Everything in f and \hat{f} happens through the *local fields*

$$\lambda_r^* = \frac{\mathbf{w}_r^{*\top} \mathbf{x}^\nu}{\sqrt{d}}, \quad \lambda_j = \frac{\mathbf{w}_j^\top \mathbf{x}^\nu}{\sqrt{d}}.$$

If $x \sim \mathcal{N}(0, 1)$ everything is characterized through the *order parameters*

$$\mathbf{Q}^\nu \equiv \mathbb{E} [\lambda^\nu \lambda^{\nu\top}] = \frac{1}{d} \mathbf{W}^{\nu\top} \mathbf{W}^\nu,$$

$$\mathbf{M}^\nu \equiv \mathbb{E} [\lambda^\nu \lambda^{*\nu\top}] = \frac{1}{d} \mathbf{W}^{\nu\top} \mathbf{W}^*,$$

$$\mathbf{P} \equiv \mathbb{E} [\lambda^{*\nu} \lambda^{*\nu\top}] = \frac{1}{d} \mathbf{W}^{*\top} \mathbf{W}^*,$$

collected into the overlap matrix

$$\Omega^\nu = \begin{pmatrix} \mathbf{Q}^\nu & \mathbf{M}^\nu \\ \mathbf{M}^{\nu\top} & \mathbf{P} \end{pmatrix}$$

Update equations for the overlaps

The updates for the overlap matrix read

$$q_{j\ell}^{\nu+1} - q_{j\ell}^\nu = \frac{\gamma_{\text{eff}}}{d} \underbrace{\left(\sigma'(\lambda_j^\nu) \lambda_\ell^\nu + \sigma'(\lambda_\ell^\nu) \lambda_j^\nu \right) \mathcal{E}^\nu}_{\text{learning}} + \frac{\gamma_{\text{eff}}^2}{d} \underbrace{\sigma'(\lambda_j^\nu) \sigma'(\lambda_\ell^\nu) (\mathcal{E}^\nu)^2}_{\text{variance}}, \quad (1)$$

$$m_{j_r}^{\nu+1} - m_{j_r}^\nu = \frac{\gamma_{\text{eff}}}{d} \underbrace{\sigma'(\lambda_j^\nu) \mathcal{E}^\nu \lambda_r^{*\nu}}_{\text{learning}},$$

where $\gamma_{\text{eff}} = \gamma/p$ and

$$\mathcal{E}^\nu = y^\nu - \hat{f}(\mathbf{x}^\nu, \mathbf{W}^\nu)$$

Learning term \iff Gradient flow approximation

Scaling of γ_{eff} \iff Relative weight of learning and variance terms

Theorem: Rigorous ODE approximation

Define

$$\delta t = \frac{\gamma_{\text{eff}} \vee \gamma_{\text{eff}}^2}{d},$$

and let $\psi : \mathbb{R}^{(p+k) \times (p+k)} \rightarrow \mathbb{R}^{(p+k) \times (p+k)}$ be the expectation of the RHS of (1):

$$\psi(\Omega)_{ij} = \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(0, \Omega)} \left[\frac{\Omega_{ij}^{\nu+1} - \Omega_{ij}^\nu}{\delta t} \right]$$

Then Ω converges to the solution $\bar{\Omega}$ of

$$\frac{d\bar{\Omega}}{dt} = \psi(\bar{\Omega}),$$

with rate

$$\|\Omega^\nu - \bar{\Omega}(\nu \delta t)\|_\infty \leq \frac{C(t) \ln(p) \sqrt{\gamma_{\text{eff}} \vee \gamma_{\text{eff}}^2}}{\sqrt{d}}$$

Extension of Saad & Solla for $p \gg 1$ with nonasymptotic bound

References

- [1] Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová. Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [2] David Saad and Sara A. Solla. On-line learning in soft committee machines. *Phys. Rev. E*, 52:4225–4243, Oct 1995.

Phase transitions in γ_{eff}

$\gamma_{\text{eff}} \ll 1$, perfect learning:

$$\psi(\Omega) = \psi_{\text{GF}}(\Omega) + o(1)$$

Equivalent to gradient flow approximations

$\gamma_{\text{eff}} \gg 1$, variance dominates:

$$\psi(\Omega) = \psi_{\text{var}}(\Omega) + o(1)$$

No learning terms: $\mathbf{M}^\nu \approx \mathbf{M}_0$

$\gamma_{\text{eff}} \propto 1$: Saad & Solla line

$$\psi(\Omega) = \psi_{\text{GF}}(\Omega) + \psi_{\text{var}}(\Omega)$$

Some learning, then a plateau with asymptotic risk $\propto \gamma \Delta$

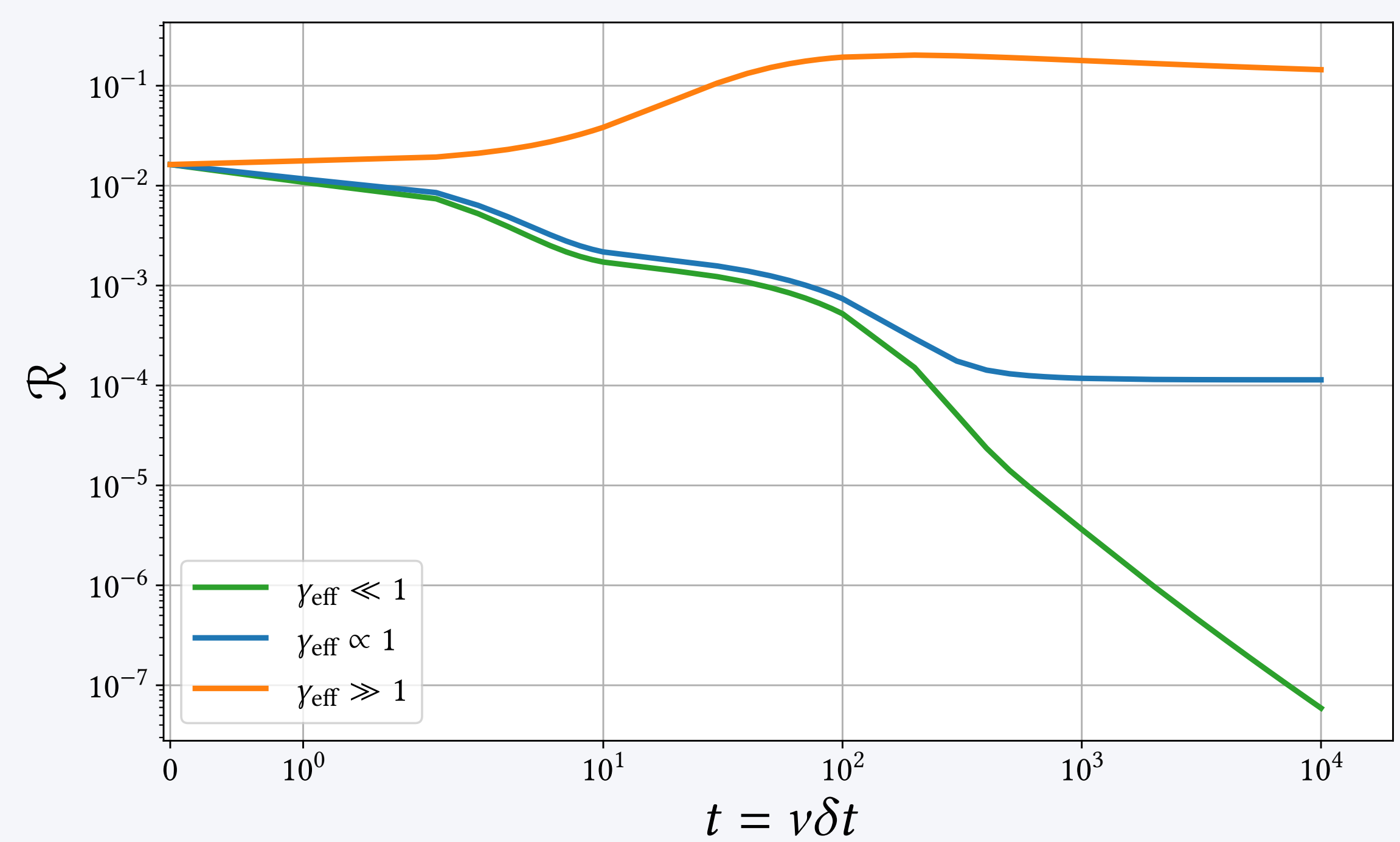


Figure 1. Illustration of the different regimes of ODEs

Interplay between γ and p

Phase transitions happen in term of $\gamma_{\text{eff}} = \gamma/p$, so we can make a two-dimensional phase diagram

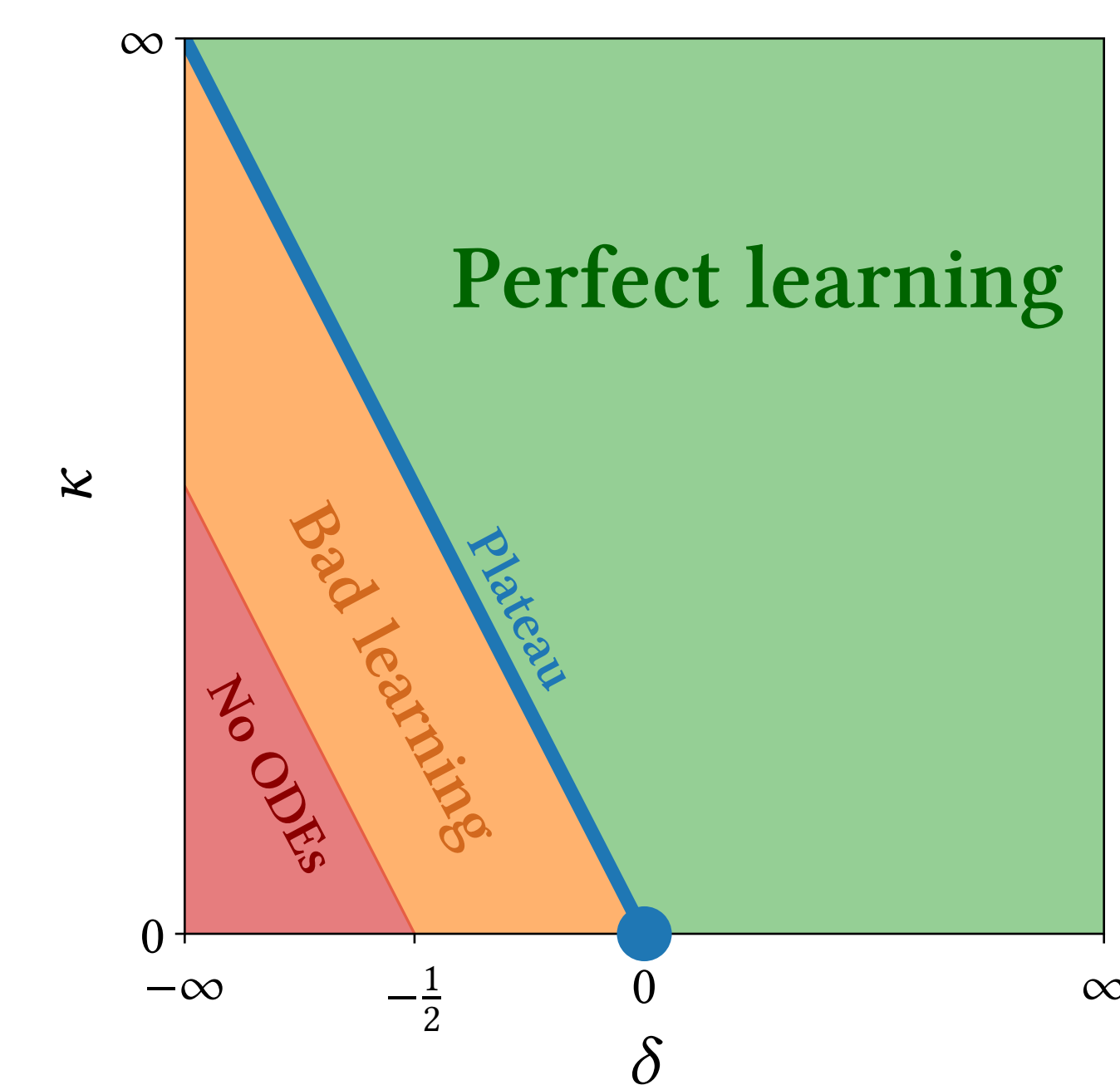


Figure 2. Phase diagram with $\gamma \sim d^{-\delta}$, $p \sim d^{\kappa}$

Overparametrization \iff Tuning the learning rate

Sample complexity

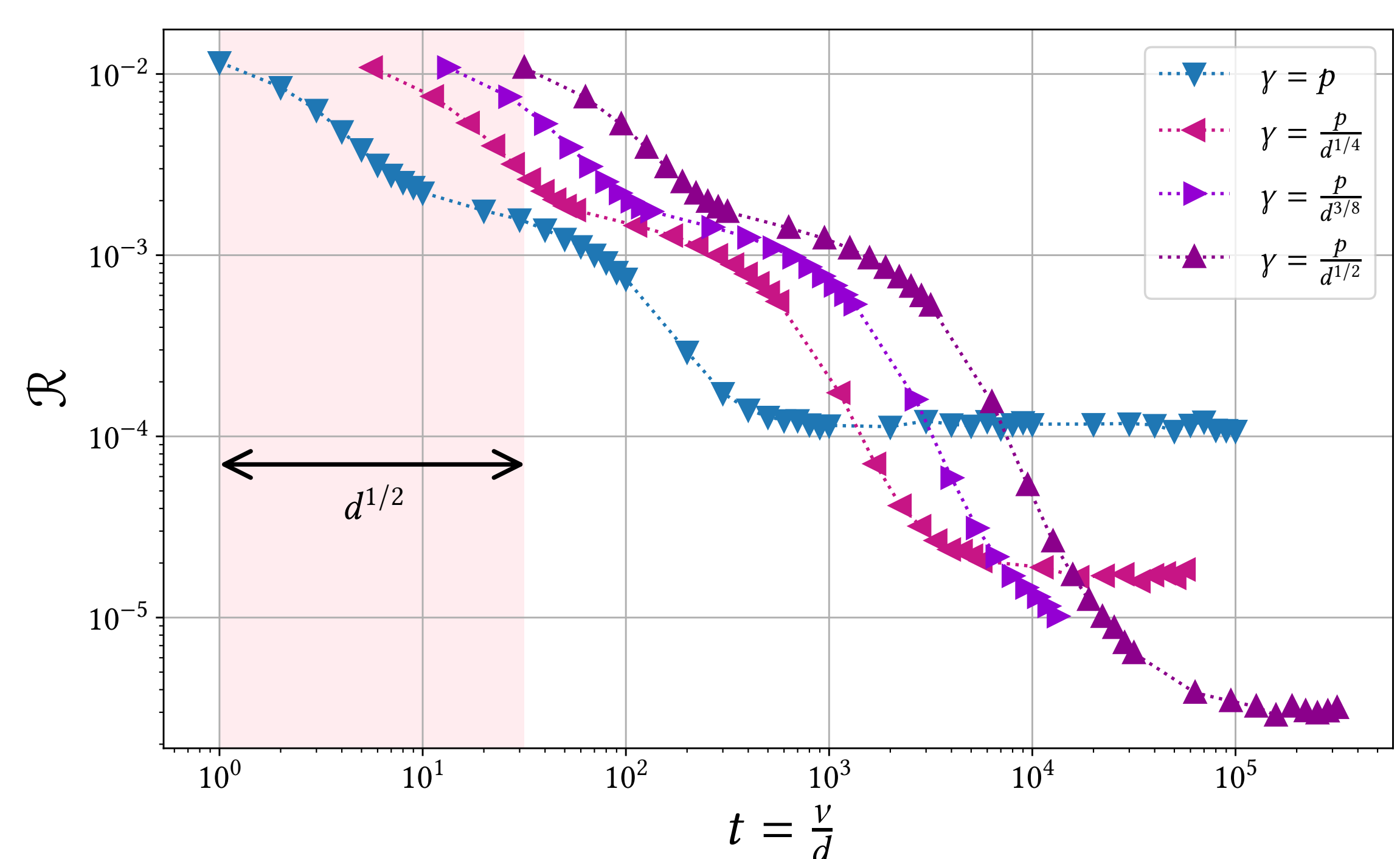


Figure 3. Effect of γ on convergence times

Tradeoff between achieved minima and sample complexity