

Excess Risk in Hyperparameter Selection

$$\hat{\lambda} = \arg \min_{\lambda \in \Lambda} E_{n_v}^v(\hat{g}_{n_t, \lambda}) \quad \hat{g}_{n_t, \lambda} = \arg \min_{f \in \mathcal{F}_\lambda} E_{n_t}(f)$$

Outer DFO

Inner ERM

- Hyperparameter selection: Bilevel Optimization Problem
- Inner level:
 - Empirical risk minimization (ERM) problem
 - Objective depends on *training split* of the data set
- Outer level:
 - Derivative-free optimization (DFO) problem
 - Objective depends on *held-out validation split* of the data set

$$\mathcal{E} = E(\hat{g}_{n_t, \hat{\lambda}}) - E(f^*)$$

Excess Risk

- Excess risk of model learned via ERM for hyperparameter selected via DFO
- Both levels of optimization use *empirical estimates* of true risk

$$E(f) = \int \ell(y, f(x)) dP(x, y), \quad E_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

True Risk

Empirical Risk

Additional Sources of Excess Risk

- Inner ERM approximated, implying following hyperparameter selection

$$\hat{\lambda} = \arg \min_{\lambda \in \Lambda} E_{n_v}^v(\tilde{g}_{n_t, \lambda}), \quad \tilde{g}_{n_t, \lambda} \in \left\{ g \in \mathcal{F}_\lambda : E_{n_t}(g) \leq \min_{f \in \mathcal{F}_\lambda} E_{n_t}(f) + \rho_{in} \right\}$$

Approximate Inner ERM

- After hyperparameter selection, *final model trained on data that combines training and validation splits*, leading to model discrepancy

$$\hat{f}_{n, \hat{\lambda}} = \arg \min_{f \in \mathcal{F}_{\hat{\lambda}}} E_n(f)$$

Final ERM with selected HP on full data

Contributions

- Provide novel excess risk bounds for above scenarios
- Propose novel data-driven practical heuristics for improved performance

Excess Risk with Exact ERM

No Model Discrepancy

Theorem 1 Let $L = 2|\Lambda| + 2$. Then, with probability at least $1 - \delta$ for any $\delta > 0$, the excess risk $\mathcal{E} = E(\hat{g}_{n_t, \hat{\lambda}}) - E(f^*)$ is bounded from above as:

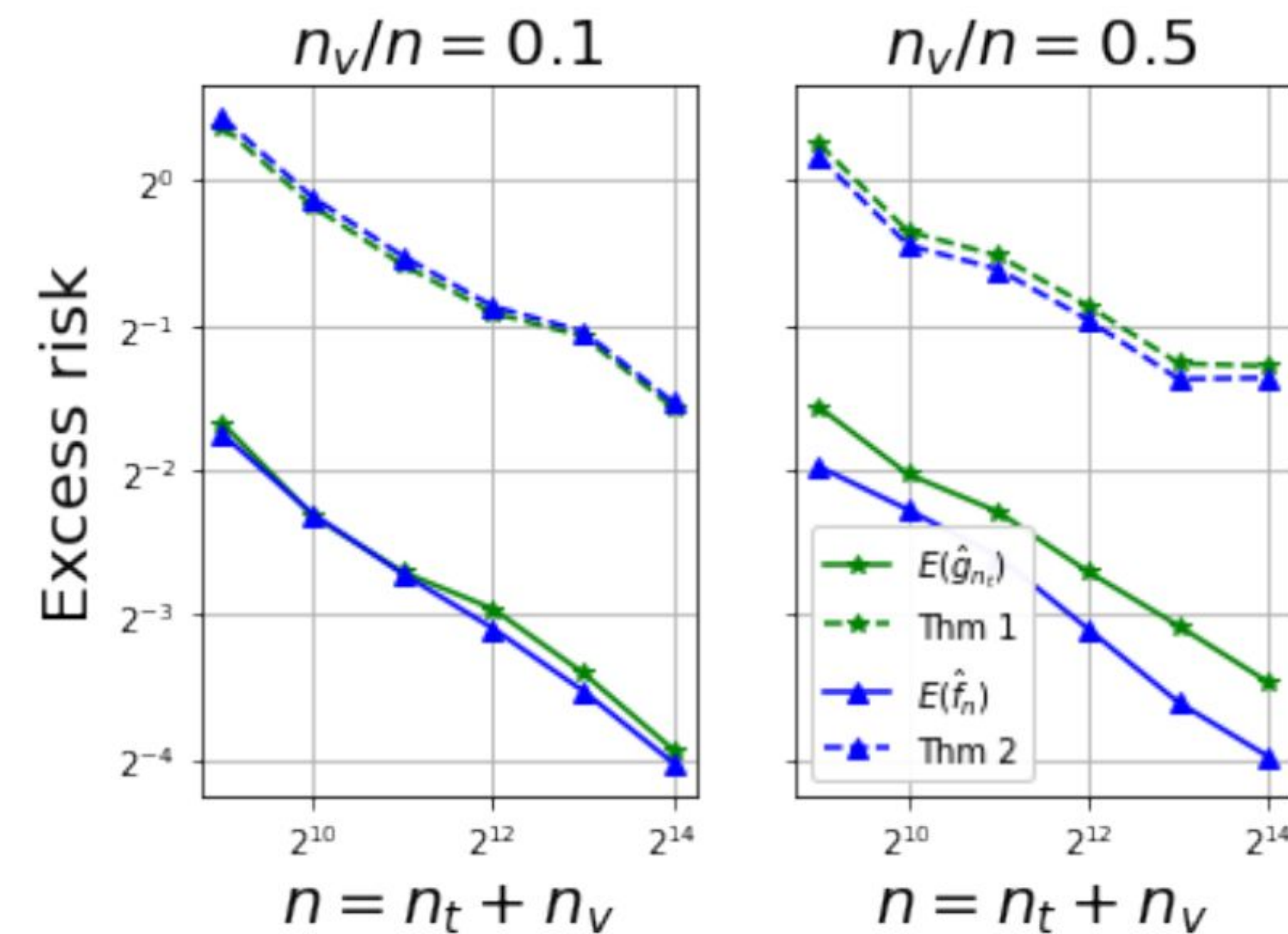
$$\mathcal{E} \leq \min_{\lambda \in \Lambda} \{ 2\Delta(\mathcal{F}_\lambda, n_t, \delta/L_1) + \mathcal{E}_{app}(\lambda) \} + B\sqrt{2\log(L_1/\delta)/n_v}$$

Model Discrepancy

Theorem 2 Let $L_2 = 2|\Lambda| + 3$. Let $\mathcal{I}_{n, n_t, \hat{\lambda}} = E_n(\hat{g}_{n_t, \hat{\lambda}}) - E_n(\hat{f}_{n, \hat{\lambda}})$ denote the “empirical risk improvement” obtained by refitting the model on the full training set. Then, with probability at least $1 - \delta$ for any $\delta > 0$, the excess risk $\mathcal{E} = E(\hat{f}_{n, \hat{\lambda}}) - E(f^*)$ is bounded from above by:

$$\mathcal{E} \leq \min_{\lambda \in \Lambda} \{ 2\Delta(\mathcal{F}_\lambda, n_t, \delta/L_2) + \mathcal{E}_{app}(\lambda) \} - \mathcal{I}_{n, n_t, \hat{\lambda}} + B\sqrt{2\log(L_2/\delta)}(1/\sqrt{n} + 1/\sqrt{n_v})$$

Comparing Bounds to True Excess Risk

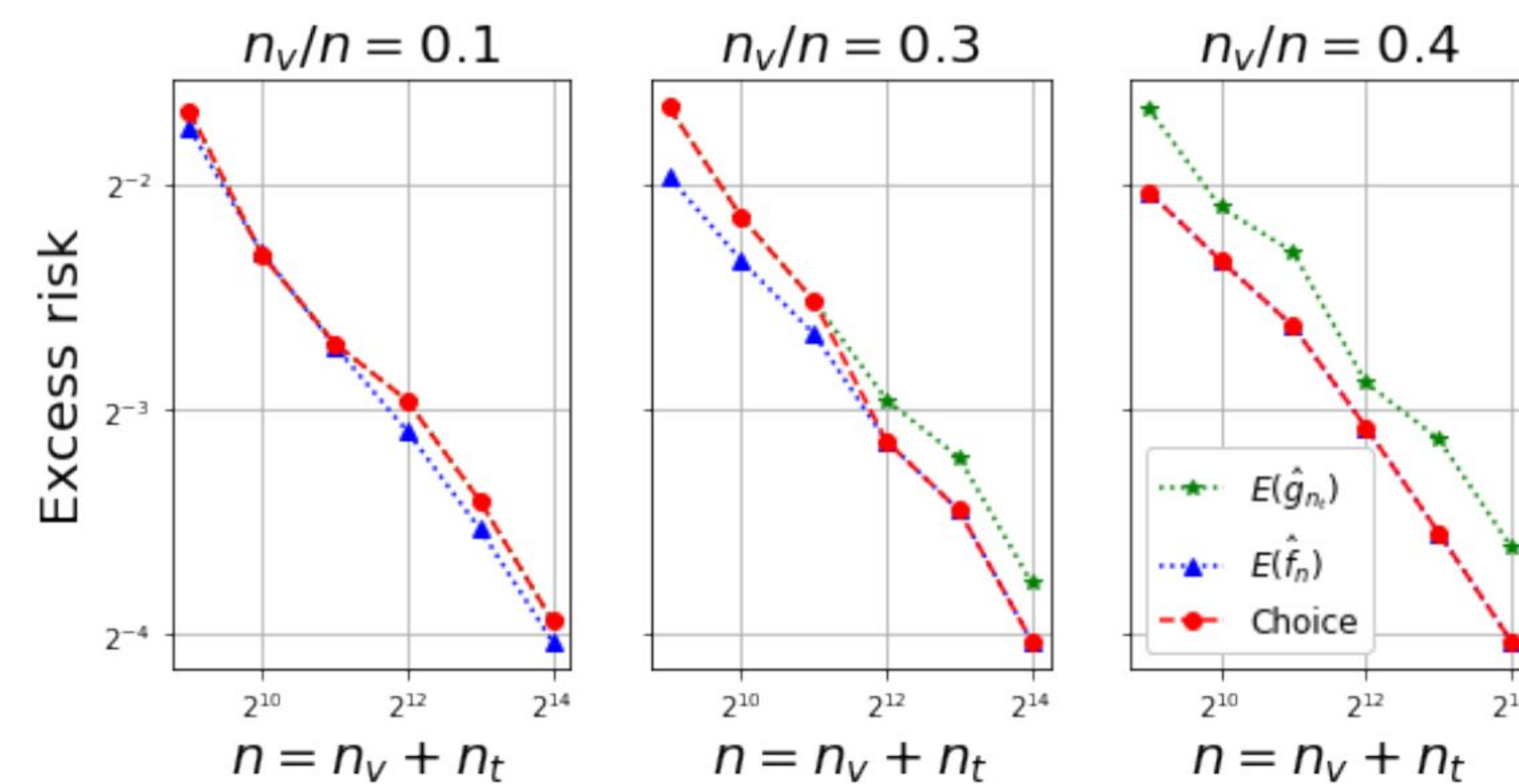


Risk bounds reflect relative performance

Data-driven Choice for Final Model

Heuristic 1 Let us define the following data-dependent scalars α, β based on the quantities in Theorems 1 & 2, and we select $\hat{f}_{n, \hat{\lambda}}$ as the final model if $\alpha \geq \beta$, or select $\hat{g}_{n, \hat{\lambda}}$ otherwise:

$$\alpha = B\sqrt{2\log(L_1/\delta)/n_v}, \quad \beta = -\mathcal{I}_{n, n_t, \hat{\lambda}} + B\sqrt{2\log(L_2/\delta)}(1/\sqrt{n} + 1/\sqrt{n_v})$$



Data-driven heuristic able to match best in most cases

Excess Risk with Approximate ERM

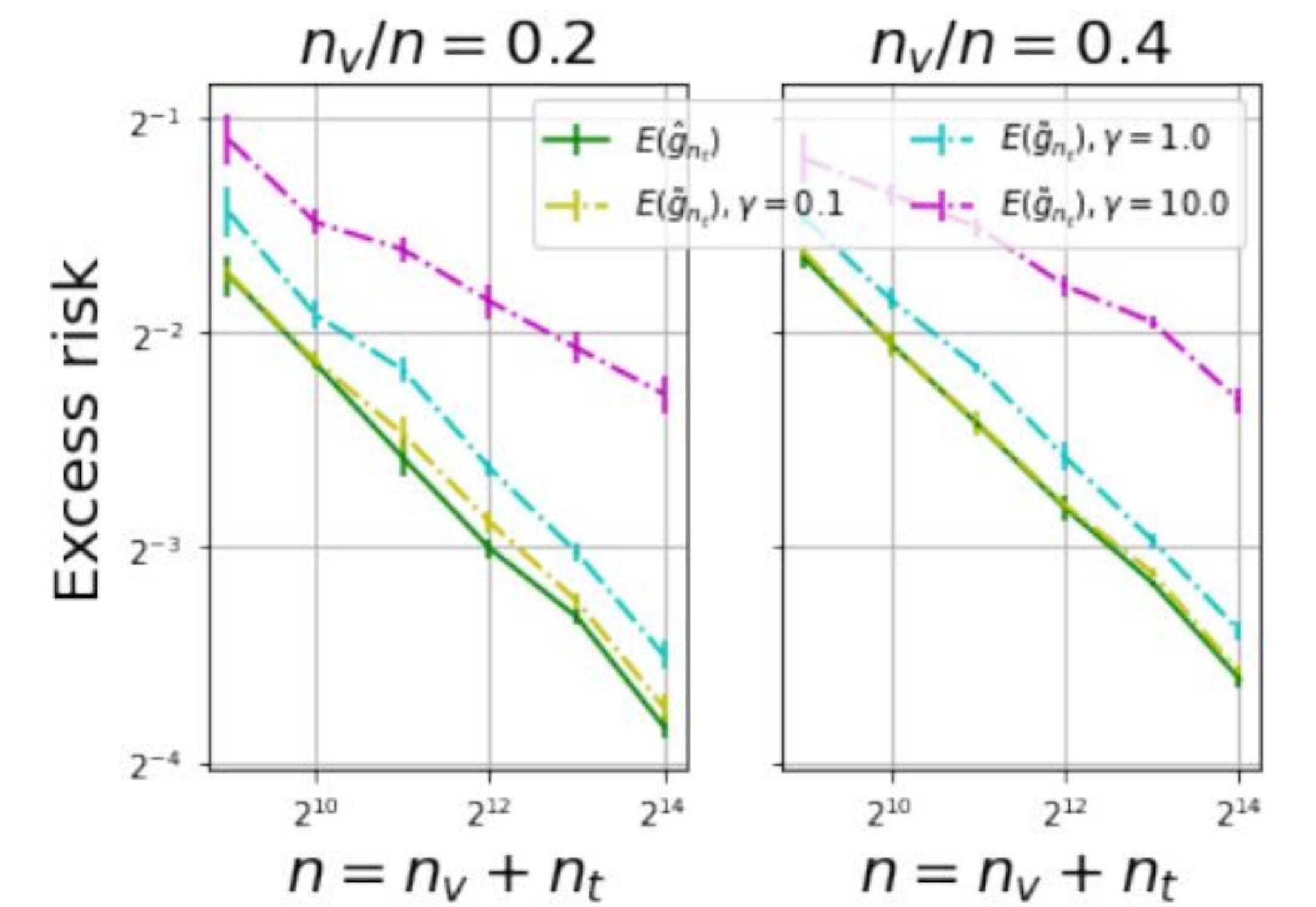
No Model Discrepancy

Theorem 3 The excess risk $\mathcal{E} = E(\tilde{g}_{n_t, \hat{\lambda}}) - E(f^*)$ can be bounded from above with probability at least $1 - \delta$ for any $\delta > 0$ and $L_3 = (2 + 3|\Lambda|)$:

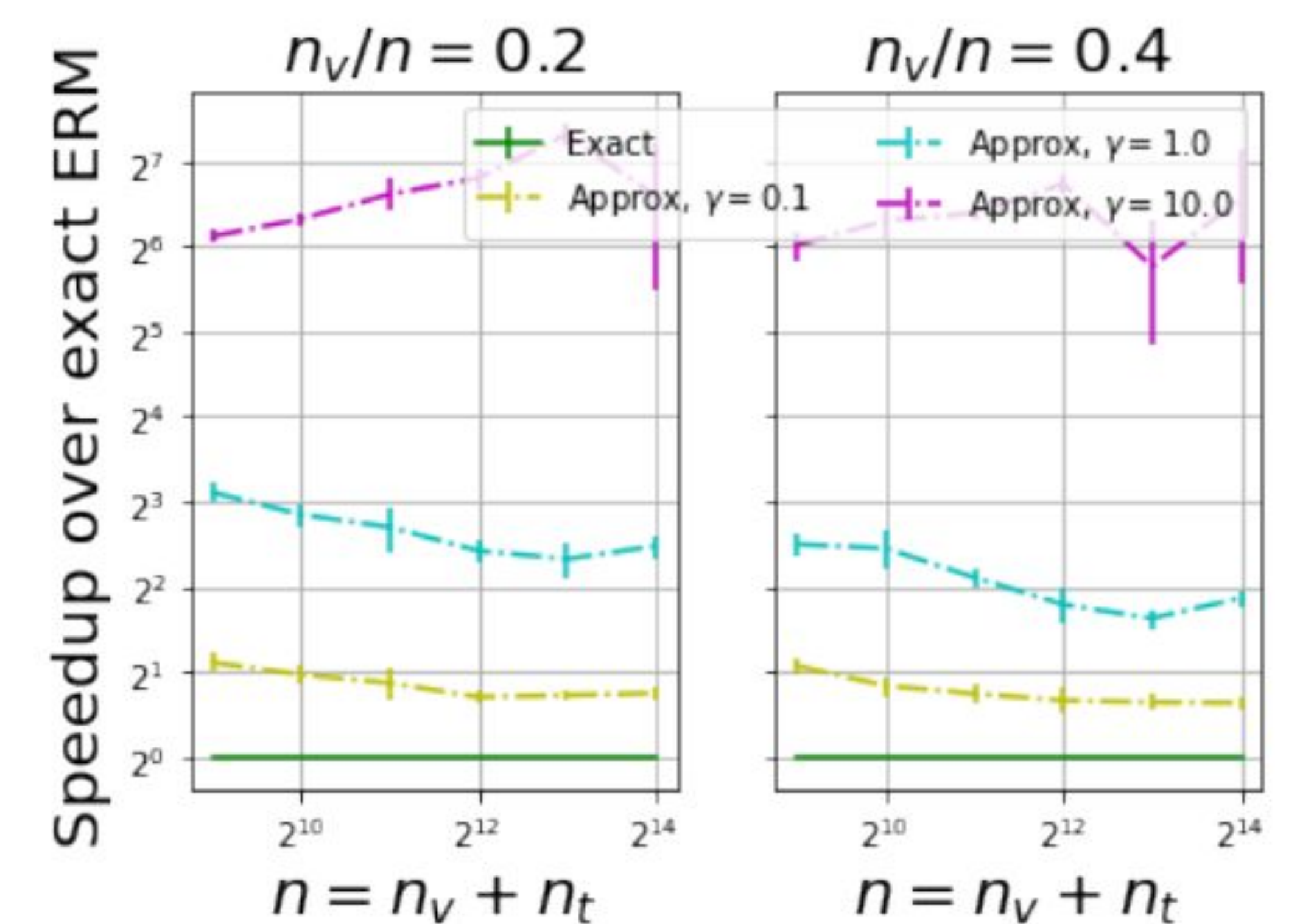
$$\mathcal{E} \leq \min_{\lambda \in \Lambda} \{ 2\Delta(\mathcal{F}_\lambda, n_t, \delta/L_3) + \mathcal{E}_{app}(\lambda) \} + B\sqrt{2\log(L_3/\delta)/n_v} + \rho_{in}$$

Data-driven Choice for ERM Approximation

Heuristic 3 Based on the terms defined in Theorem 3, select a scaling parameter $\gamma > 0$ and set ρ_{in} as $\rho_{in} = \gamma B\sqrt{2\log(L_3/\delta)/n_v}$ such that $\rho_{in} \sim o(B\sqrt{2\log(L_3/\delta)/n_v})$. A value of $\gamma = 0.1$ suffices in our experience.



Data-driven choice of approximation in inner ERM does not increase excess risk significantly over exact inner ERM



Data-driven choice of approximation in inner ERM provides 2X speedup over exact ERM with no additional excess risk, and can provide 4-6X speedup with slight increase in excess risk

Notations

$$f^* : \mathcal{X} \rightarrow \mathcal{Y} \text{ such that } f^*(x) = \arg \min_{\hat{y} \in \mathcal{Y}} \mathbb{E}[\ell(y, \hat{y}) | x]$$

Bayes optimal model

$$\bar{f}_\lambda = \arg \min_{f \in \mathcal{F}_\lambda} E(f)$$

True risk minimizer

$$\mathcal{E}_{app}(\lambda) = E(\bar{f}_\lambda) - E(f^*)$$

Approximation risk

$$\sup_{f \in \mathcal{F}} |E_n(f) - E(f)| \leq \Delta(\mathcal{F}, n, \delta)$$

Estimation risk bound