

Regularized Empirical Risk Minimization

Problem Setting:

Unknown distribution : rv $Z \in \mathcal{Z}$ with distribution ρ

Parameter $\theta \in \mathcal{H}$, \mathcal{H} a Hilbert space

Problem: Minimize an expected loss:

$$\min_{\theta \in \mathcal{H}} L(\theta) := \mathbb{E}[\ell_Z(\theta)], \quad \ell_Z(\theta) \text{ loss function} \quad (1)$$

Well-specified assumption $\theta^* \in \operatorname{argmin}_{\theta \in \mathcal{H}} L(\theta)$ exists

Statistical performance: $L(\theta) - L(\theta^*)$

Data: access to ρ through n i.i.d observations $(z_i)_{1 \leq i \leq n} \in \mathcal{Z}^n$ from Z

Regularized ERM:

$$\hat{\theta}_\lambda = \operatorname{argmin}_{\theta \in \mathcal{H}} \hat{L}(\theta) + \frac{\lambda}{2} \|\theta\|_{\mathcal{H}}^2, \quad \hat{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell_{z_i}(\theta)$$

Basic result : slow rates

$$L(\hat{\theta}_\lambda) - L(\theta^*) \leq \frac{\|\nabla L\|_\infty^2 + \|\theta^*\|^2}{\sqrt{n}}, \quad \lambda = \frac{1}{\sqrt{n}}$$

Bias-variance trade-off for least-squares

Loss: $\ell_{x,y}(\theta) = \|y - \theta \cdot x\|^2$, $L(\theta) = \mathbb{E}[\|Y - \theta \cdot X\|^2]$

Covariance operator: $\Sigma = \mathbb{E}[X \otimes X]$

$$\forall \theta \in \mathcal{H}, L(\theta) - L(\theta^*) = \|\Sigma^{1/2}(\theta - \theta^*)\|^2 = \|\theta - \theta^*\|_\Sigma^2.$$

Two main terms :

Effective dimension: $\operatorname{df}_\lambda = \operatorname{Tr}(\Sigma_\lambda^{-1} \Sigma)$, $\Sigma_\lambda = \Sigma + \lambda I$

Bias term: $b_\lambda = \lambda \|\Sigma_\lambda^{-1} \theta^*\|$

Bias-Variance trade-off :

$$L(\hat{\theta}_\lambda) - L(\theta^*) \leq b_\lambda^2 + \frac{\operatorname{df}_\lambda}{n}.$$

Parametrization and optimal rates

Effective dimension \leftrightarrow **spectrum of covariance matrix**

$(\lambda_i)_i$ eigenvalues of Σ in decreasing order.

assumption: $\operatorname{df}_\lambda \leq Q^2 \lambda^{-1/\alpha} \leftrightarrow \lambda_i = O(i^{-\alpha})$

Bias term \leftrightarrow **difficulty of the learning problem**

assumption : $b_\lambda \leq L \lambda^{1/2+r} \leftrightarrow \|\Sigma^{-r} \theta^*\| < \infty$

Optimal fast rates for $\lambda = (Q/L)^2 n^{-\alpha/(1+2r)\alpha+1}$

$$L(\hat{\theta}_\lambda) - L(\theta^*) \leq Q^{2\gamma} L^{2(1-\gamma)} n^{-\gamma}, \quad \gamma = (1+2r)\alpha/(1+(1+2r)\alpha)$$

Acknowledgments and References

We acknowledge support from the ERCIM Alain Bensoussan Fellowship and the European Research Council (SEQUOIA 724063)

• A. Caponnetto and E. De Vito. Optimal rates for the regularized least-squares algorithm. Found. Comput. Math., 7(3):331–368, July 2007.

• Dmitrii Ostrovskii and Francis Bach. Finite-sample analysis of M-estimators using self-concordance. Technical Report 1810.06838, arXiv, 2018.

Generalized bias-variance trade-off

• **Hessian at optimum**: $\mathbf{H}(\theta^*) = \nabla^2 L(\theta^*) = \mathbb{E}[\nabla^2 \ell_Z(\theta^*)]$

Main terms :

• **Effective dimension**: $\operatorname{df}_\lambda = \mathbb{E}[\|\mathbf{H}_\lambda^{-1/2}(\theta^*) \nabla \ell_Z(\theta^*)\|^2]$

• **Bias term**: $b_\lambda = \lambda \|\mathbf{H}_\lambda^{-1}(\theta^*) \theta^*\| = \|\nabla L_\lambda(\theta^*)\|_{\mathbf{H}_\lambda^{-1}(\theta^*)}$

• **Dikin radius**: $r_\lambda = \sqrt{\lambda}/R$

Bias-Variance trade-off :

$$L(\hat{\theta}_\lambda) - L(\theta^*) \leq b_\lambda^2 + \frac{\operatorname{df}_\lambda}{n}, \quad b_\lambda, \sqrt{\operatorname{df}_\lambda/n} \leq r_\lambda$$

Generalized optimal rates

Effective dimension \leftrightarrow **renormalized Fischer information**

assumption: $\operatorname{df}_\lambda \leq Q^2 \lambda^{-1/\alpha}$.

For GLMs in the well specified case, $\operatorname{df}_\lambda = \operatorname{Tr}(\mathbf{H}_\lambda^{-1}(\theta^*) \mathbf{H}(\theta^*))$.

Bias term \leftrightarrow **difficulty of the learning problem**

assumption : $b_\lambda \leq L \lambda^{1/2+r} \leftrightarrow \|\mathbf{H}(\theta^*)^{-r} \theta^*\| < \infty$

Optimal fast rates for $\lambda = (Q/L)^2 n^{-\alpha/((1+2r)\alpha+1)}$

$$L(\hat{\theta}_\lambda) - L(\theta^*) \leq Q^{2\gamma} L^{2(1-\gamma)} n^{-\gamma}$$

Sketch of proof

Define $\theta_\lambda^* = \operatorname{argmin}_{\theta \in \mathcal{H}} L(\theta) + \frac{\lambda}{2} \|\theta\|^2$.

Idea: Decompose the statistical performance $\theta^* \rightarrow \theta_\lambda^* \rightarrow \hat{\theta}_\lambda$.

Bias term: $\theta^* \leftrightarrow \theta_\lambda^*$ Using localisation on L_λ :

$$b_\lambda \leq r_\lambda \implies \begin{cases} R \|\theta_\lambda^* - \theta^*\| \leq \log 2 \\ \|\theta_\lambda^* - \theta^*\|_{\mathbf{H}_\lambda} \leq 2b_\lambda \end{cases}$$

Equivalence of norms: $\mathbf{H}_\lambda \sim \hat{\mathbf{H}}_\lambda$ Concentration inequality.

Variance term: $\hat{\theta}_\lambda \leftrightarrow \theta_\lambda^*$ Localization + Concentration inequality :

1. Localization on \hat{L}_λ + Equivalence of norms

$$\|\nabla \hat{L}_\lambda(\theta_\lambda^*)\|_{\mathbf{H}_\lambda^{-1}} \leq r_\lambda \implies \begin{cases} R \|\theta_\lambda^* - \hat{\theta}_\lambda\| \leq \log 2 \\ \|\theta_\lambda^* - \hat{\theta}_\lambda\|_{\mathbf{H}_\lambda} \leq 2\|\nabla \hat{L}_\lambda(\theta_\lambda^*)\|_{\mathbf{H}_\lambda^{-1}} \end{cases}$$

2. Concentration inequality

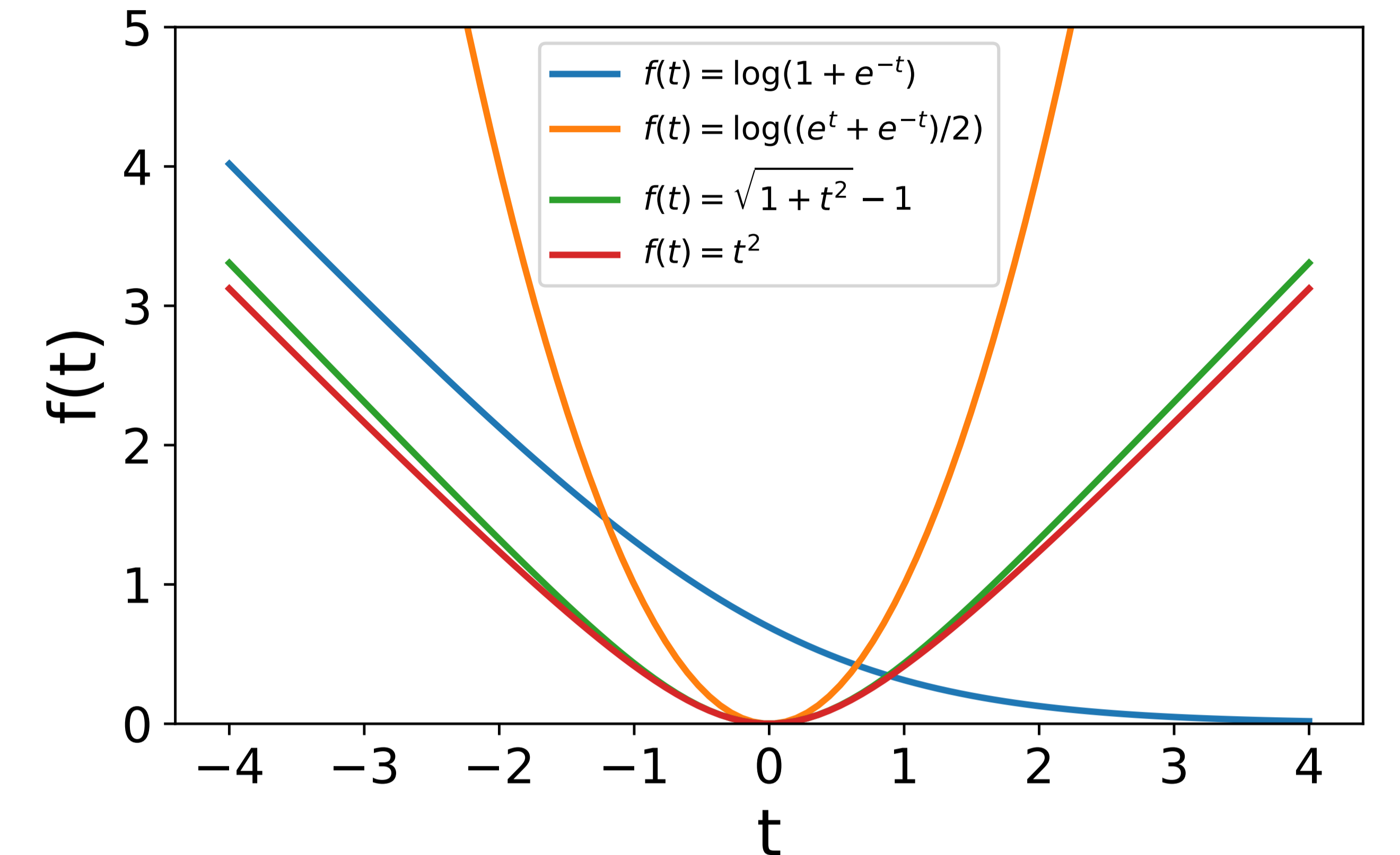
$$\|\nabla \hat{L}_\lambda(\theta_\lambda^*)\|_{\mathbf{H}_\lambda^{-1}} \leq \sqrt{\operatorname{df}_\lambda/n}$$

Putting things together

$$\sqrt{\operatorname{df}_\lambda/n}, b_\lambda \leq r_\lambda \implies \begin{cases} R \|\theta^* - \hat{\theta}_\lambda\| \leq 2 \log 2 \\ \|\theta^* - \hat{\theta}_\lambda\|_{\mathbf{H}_\lambda} \leq \sqrt{\operatorname{df}_\lambda/n} + b_\lambda \end{cases}$$

Quadratic approximation: $L(\hat{\theta}_\lambda) - L(\theta^*) \leq 4(\sqrt{\operatorname{df}_\lambda/n} + b_\lambda)^2$

Generalized Self-Concordance



Definition (GSC function)

$$\nabla^{(3)} F(\theta)[h, k, k] \leq R \|h\| \nabla^2 F(\theta)[k, k]$$

Assumption: the $(\ell_z)_{z \in \operatorname{supp}(\rho)}$ are all GSC functions for R .

Consequence: $L, L_\lambda = L + \frac{\lambda}{2} \|\cdot\|^2, \hat{L}_\lambda$ are GSC for R .

Examples from Supervised Learning

Distribution: $Z = (X, Y)$, where we want to predict Y from X .

Predictor: $\theta \cdot \Phi(x)$ ($\theta \cdot \Phi(x, y)$ for multiclass)

Assumption: Y is bounded, $\Phi(X)$ is bounded

Regression: $\ell_z(\theta) = \psi(y - \theta \cdot \Phi(x))$

• **Square loss**: $\psi(t) = \frac{1}{2} t^2$

• **Huber loss 1**: $\psi(t) = \sqrt{1 + t^2} - 1$

• **Huber loss 2**: $\psi(t) = \log \frac{e^t + e^{-t}}{2}$

Classification:

• **Logistic loss**: $\ell_z(\theta) = \log(1 + e^{-y \cdot \theta \cdot \Phi(x)})$

• **GLMs**: $\ell_z(\theta) = -\theta \cdot \Phi(x, y) + \log \int_Y \exp(\theta \cdot \Phi(x, y')) d\mu(y')$

Properties of GSC functions

Assume F is GSC, with minimizer θ^* . Let $\theta \in \mathcal{H}$, $t = R \|\theta - \theta^*\|$.

Quadratic approximation

$$F(\theta) - F(\theta^*) \leq e^t \|\theta - \theta^*\|_{\nabla^2 F(\theta)}^2$$

Localization using gradients Define $F_\lambda = F + \frac{\lambda}{2} \|\cdot\|^2$.

$$\|\nabla F_\lambda(\theta)\|_{\nabla^2 F_\lambda(\theta)^{-1}} \leq r_\lambda \implies \begin{cases} t \leq \log 2 \\ \|\theta - \theta^*\|_{\nabla^2 F_\lambda(\theta)} \leq 2\|\nabla F_\lambda(\theta)\|_{\nabla^2 F_\lambda(\theta)^{-1}} \end{cases}$$