# DISSECTING NON-VACUOUS GENERALIZATION BOUNDS BASED ON THE MEAN-FIELD APPROXIMATION

Pitas Konstantinos,
LTS2, EPFL,
Switzerland.
konstantinos.pitas@epfl.ch

# 1)OVERVIEW
## GENERALIZATION BOUNDS AND PAC-BAYES

$$\mathcal{L}(\rho) \leq \hat{\mathcal{L}}(\rho) + \text{complexity}(\rho)$$
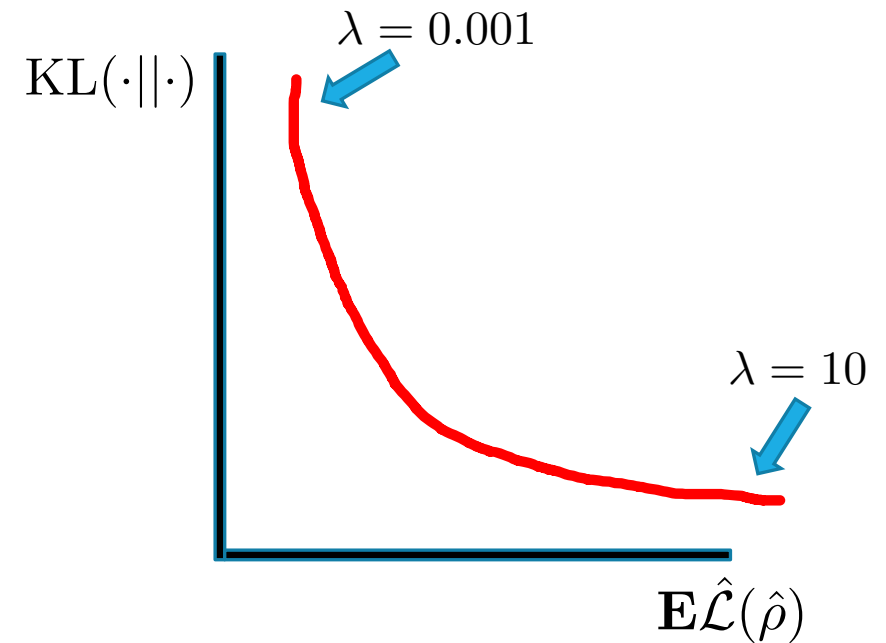
Risk    Empirical Risk

$$\mathbf{E}\mathcal{L}(\hat{\rho}) \leq \mathbf{E}\hat{\mathcal{L}}(\hat{\rho}) + \beta \mathrm{KL}(\hat{\rho}\|\pi)$$

# PAC-BAYES BOUNDS

$$\mathbf{E}\mathcal{L}(\hat{\rho}) \leq \mathbf{E}\hat{\mathcal{L}}(\hat{\rho}) + \beta \mathrm{KL}(\hat{\rho}||\pi)$$
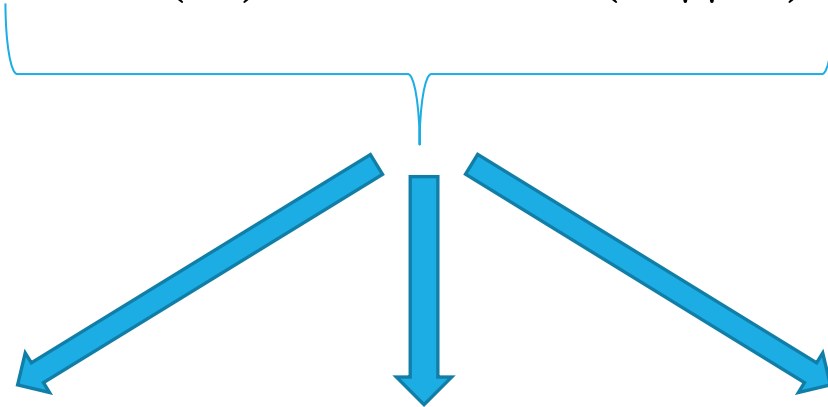
$$\hat{\rho}(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\hat{\rho}}, \lambda \mathbf{I})$$

$$\pi(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\mathrm{init}}, \lambda \mathbf{I})$$

# MODELING CHOICES

$$\mathbf{E}\mathcal{L}(\hat{\rho}) \leq \mathbf{E}\hat{\mathcal{L}}(\hat{\rho}) + \beta\mathrm{KL}(\hat{\rho}||\pi)$$

$$\hat{\rho}(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\hat{\rho}}, \lambda\mathbf{I})$$

$$\pi(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\mathrm{init}}, \lambda\mathbf{I})$$

Baseline

$$\hat{\rho}(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\hat{\rho}}, \boldsymbol{\sigma}_{\hat{\rho}})$$

$$\pi(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\mathrm{init}}, \lambda\mathbf{I})$$

VI

Invalid
Sanity Check

$$\hat{\rho}(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\hat{\rho}}, \boldsymbol{\sigma}_{\hat{\rho}})$$

$$\pi(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\mathrm{init}}, \lambda\boldsymbol{\sigma}_{\pi})$$

# RISK-COMPLEXITY PLOTS:
# A MORE INTUITIVE WAY OF COMPARING BOUNDS

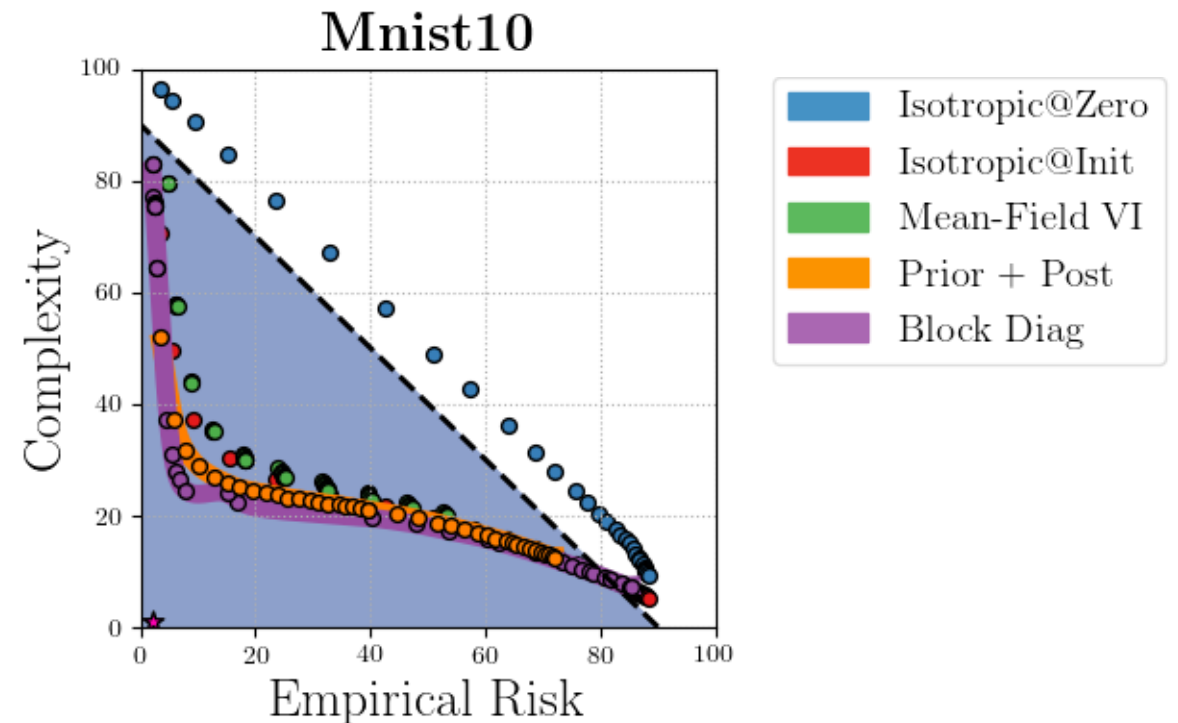$$40\% < 60\%\textcolor{red}{?}$$

Mnist10



$$\mathbf{E}\mathcal{L}(\hat{\rho}) \leq \mathbf{E}\hat{\mathcal{L}}(\hat{\rho}) + \beta\mathrm{KL}(\hat{\rho}\|\pi)$$

# MODELING CHOICES

Choosing the <u>prior</u> mean to be the random DNN initialization is very important!

# OVERVIEW OF RESULTS

- Baseline is non-vacuous in simple cases.

- Mean-Field VI of the **posterior covariance** yields marginal improvement (Problems with optimization?)

- Closed form optimization:
  - Optimizing the posterior diagonal covariance.
  - Optimizing the prior diagonal covariance.
  - Generalizing the posterior covariance to be **block diagonal**.



Mnist10

# 2)DETAILS
# MEAN-FIELD VI

$$\mathbf{E}_{\boldsymbol{\theta} \sim \hat{\rho}(\boldsymbol{\theta})} \hat{\mathcal{L}}_{X,Y}^{\ell_{01}}(f_{\boldsymbol{\theta}}) + \frac{1}{\be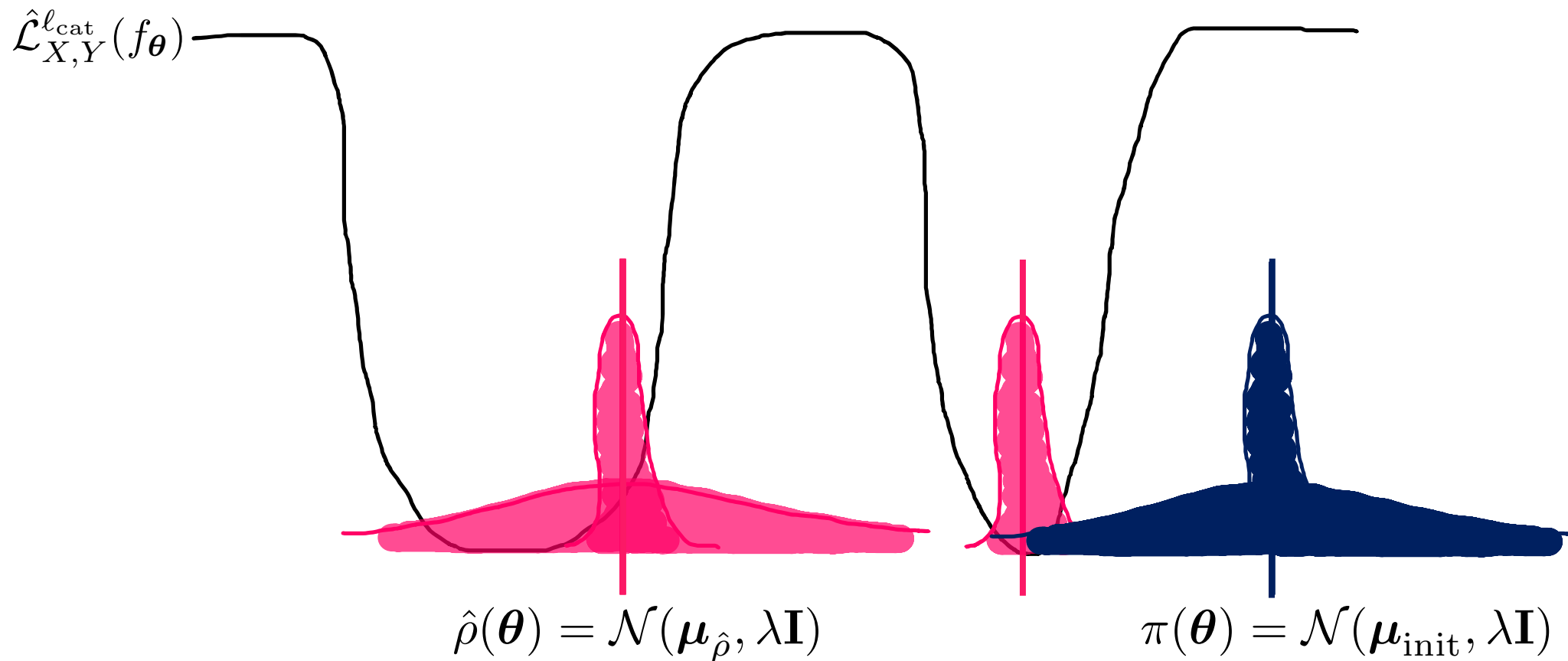ta n}(\mathrm{KL}(\hat{\rho}(\boldsymbol{\theta}) || \mathcal{N}(\boldsymbol{\mu}_{\mathrm{init}}, \lambda \mathbf{I})) + \ln \frac{1}{\delta})$$

$$\mathbf{E}_{\boldsymbol{\theta} \sim \hat{\rho}(\boldsymbol{\theta})} \hat{\mathcal{L}}_{X,Y}^{\ell_{\mathrm{cat}}}(f_{\boldsymbol{\theta}})$$

$$\sum_{i=0}^{T} \hat{\mathcal{L}}_{X,Y}^{\ell_{\mathrm{cat}}}(f_{\boldsymbol{\theta}_i}) \qquad \boldsymbol{\theta} = \boldsymbol{\mu}_{\hat{\rho}} + \sqrt{\boldsymbol{\sigma}_{\hat{\rho}}} \odot \mathcal{N}(\mathbf{0}, \mathbf{I}) \qquad \boxed{\boldsymbol{\mu}_{\hat{\rho}} \text{ remains fixed}}$$
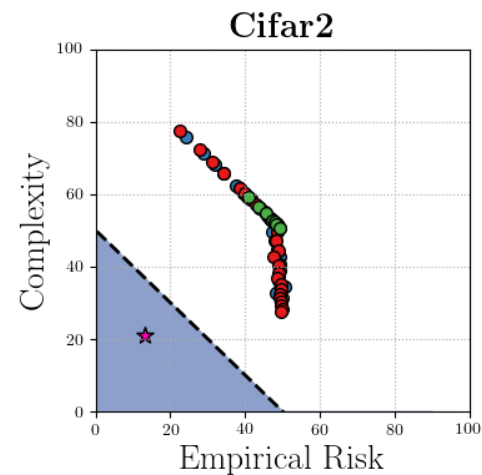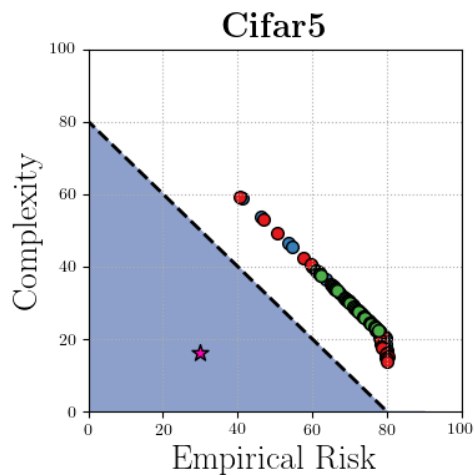
# MEAN-FIELD VI WITH FIXED POSTERIOR MEAN



$\hat{\mathcal{L}}_{X,Y}^{\ell_{\mathrm{cat}}}(f_{\boldsymbol{\theta}})$

$\hat{\rho}(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\hat{\rho}}, \lambda\mathbf{I})$

$\pi(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\mathrm{init}}, \lambda\mathbf{I})$

# MEAN-FIELD VI GRID SEARCH

$$\min_{\sigma_{\hat{\rho}}} \overbrace{\sum_{i=0}^{T} \hat{\mathcal{L}}_{X,Y}^{\ell_{\text{cat}}}(f_{\boldsymbol{\theta}_i})}^{} + \frac{1}{\beta n} \overbrace{(\text{KL}(\hat{\rho}(\boldsymbol{\theta}) || \mathcal{N}(\boldsymbol{\mu}_{\text{init}}, \lambda \mathbf{I})) + \ln \frac{1}{\delta})}^{}$$
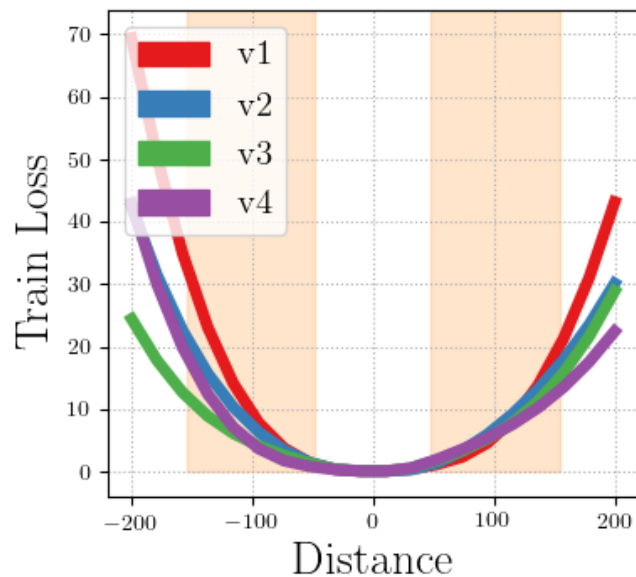
# MEAN-FIELD VI RESULTS

# QUADRATIC APPROXIMATION

$$C_\beta(X, Y; \hat{\rho}, \pi) = \mathbf{E}_{\boldsymbol{\theta} \sim \hat{\rho}(\boldsymbol{\theta})} \hat{\mathcal{L}}_{X,Y}^{\ell_{\mathrm{cat}}}(f_{\boldsymbol{\theta}}) + \beta \mathrm{KL}(\hat{\rho}(\boldsymbol{\theta}) \| \pi(\boldsymbol{\theta}))$$

$$\approx \mathbf{E}_{\boldsymbol{\eta} \sim \hat{\rho}'(\boldsymbol{\theta})}[\frac{1}{2} \boldsymbol{\eta}^T \nabla^2 \hat{\mathcal{L}}_{X,Y}^{\ell_{\mathrm{cat}}}(f_{\boldsymbol{\theta}}) \boldsymbol{\eta}] + \beta \mathrm{KL}(\hat{\rho}(\boldsymbol{\theta}) \| \pi(\boldsymbol{\theta}))$$

# QUADRATIC APPROXIMATION

**Lemma 4.1.** *The convex optimization problem* $\min_{\Sigma_{\hat{\rho}}} \mathbf{E}_{\eta \sim \hat{\rho}'(\theta)}[\frac{1}{2}\eta^T \mathbf{H}\eta] + \beta \mathrm{KL}(\hat{\rho}(\theta)\|\pi(\theta))$ *where* $\hat{\rho}(\theta) = \mathcal{N}(\mu_{\hat{\rho}}, \Sigma_{\hat{\rho}})$ *and* $\pi(\theta) = \mathcal{N}(\mu_{\pi}, \lambda\Sigma_{\pi})$ *is minimized at*

$$\Sigma_{\hat{\rho}}^* = \beta(\mathbf{H} + \frac{\beta}{\lambda}\Sigma_{\pi}^{-1})^{-1}, \qquad (6)$$

*where* $\mathbf{H} \equiv \nabla^2 \hat{\mathcal{L}}_{X,Y}^{\ell_{cat}}(f_{\theta})$ *captures the curvature at the minimum, while* $\Sigma_{\pi}$ *is the prior covariance.*

**Lemma 4.2.** *The optimal prior and posterior covariances for* $\min_{\sigma_{\hat{\rho}}, \sigma_{\pi}} \mathbf{E}_{\eta \sim \hat{\rho}'(\theta)}[\frac{1}{2}\eta^T \mathbf{H}\eta] + \beta \mathrm{KL}(\hat{\rho}(\theta)\|\pi(\theta))$ *with* $\hat{\rho}(\theta) = \mathcal{N}(\mu_{\hat{\rho}}, \sigma_{\hat{\rho}})$ *and* $\pi(\theta) = \mathcal{N}(\mu_{\pi}, \lambda\sigma_{\pi})$ *have elements*
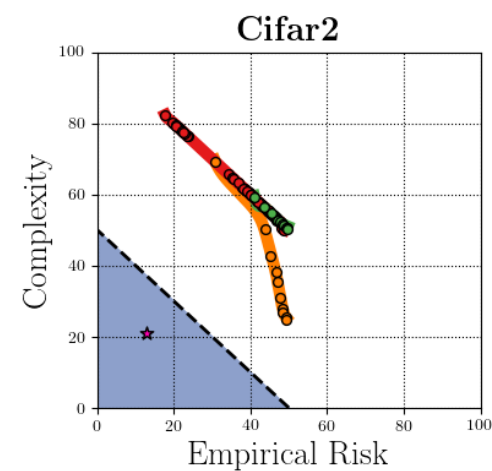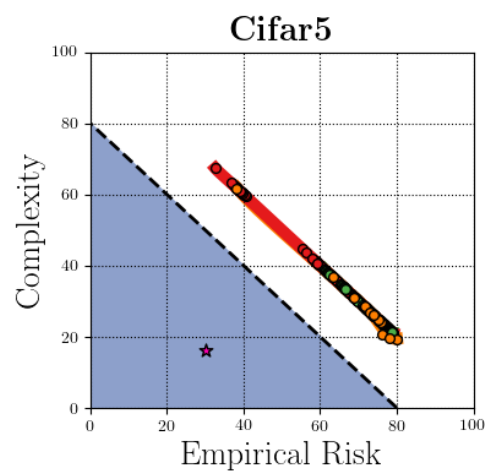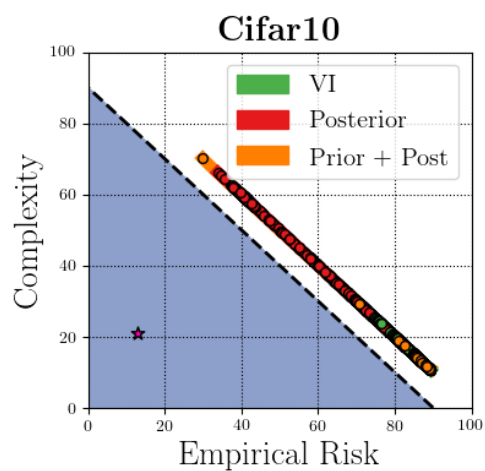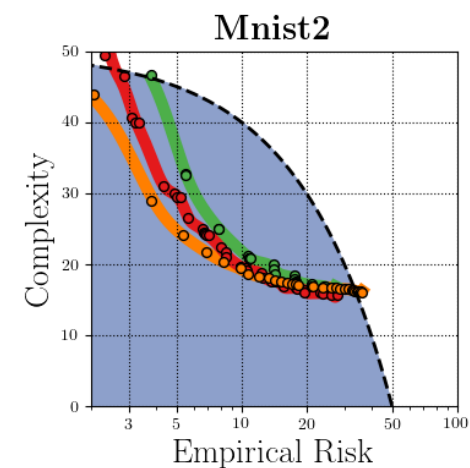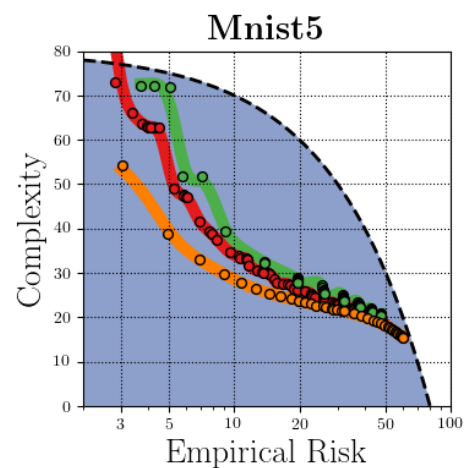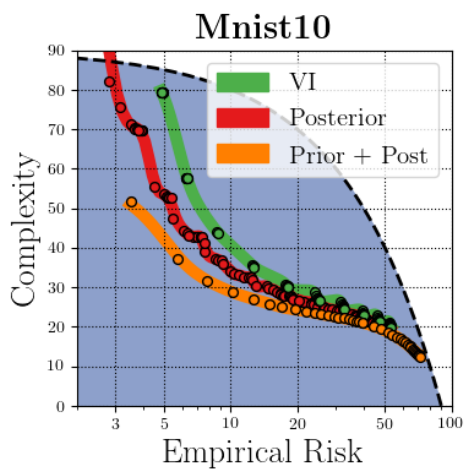
$$(\sigma_{\hat{\rho}i}^*)^{-1} = \frac{1}{2\beta}[h_i + \sqrt{h_i^2 + \frac{4\beta h_i}{(\mu_{i\hat{\rho}} - \mu_{i\pi})^2}}], \qquad (7)$$

$$(\sigma_{\pi i}^*)^{-1} = \frac{\lambda}{2\beta}[\sqrt{h_i^2 + \frac{4\beta h_i}{(\mu_{i\hat{\rho}} - \mu_{i\pi})^2}} - h_i], \qquad (8)$$
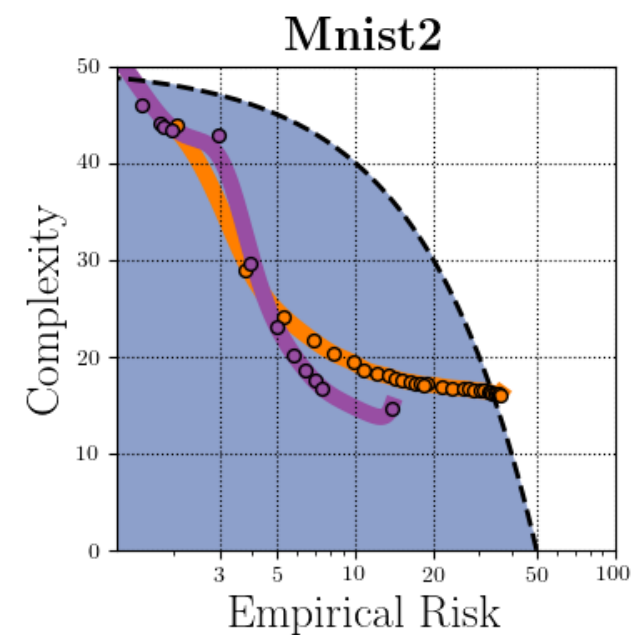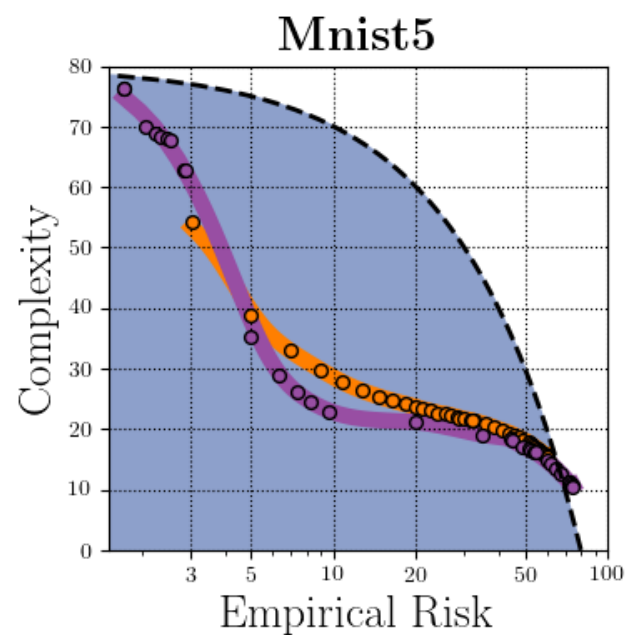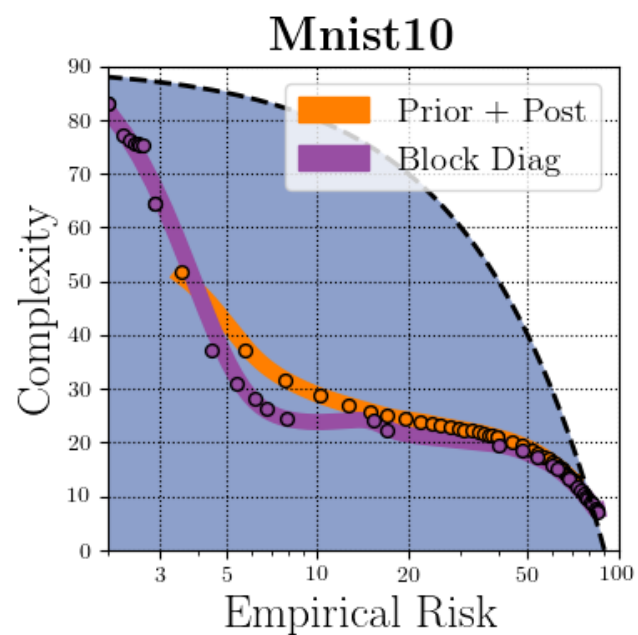
*where* $\mathbf{H} \equiv \nabla^2 \hat{\mathcal{L}}_{X,Y}^{\ell_{cat}}(f_{\theta})$ *captures the curvature at the minimum.*

Posterior

Prior + Posterior

# QUADRATIC APPROXIMATION RESULTS

# BLOCK DIAGONAL RESULTS

# THANK YOU FOR YOUR ATTENTION!

Arxiv preprint: https://arxiv.org/pdf/1909.03009.pdf

contact: konstantinos.pitas@epfl.ch

LTS2, EPFL, Switzerland