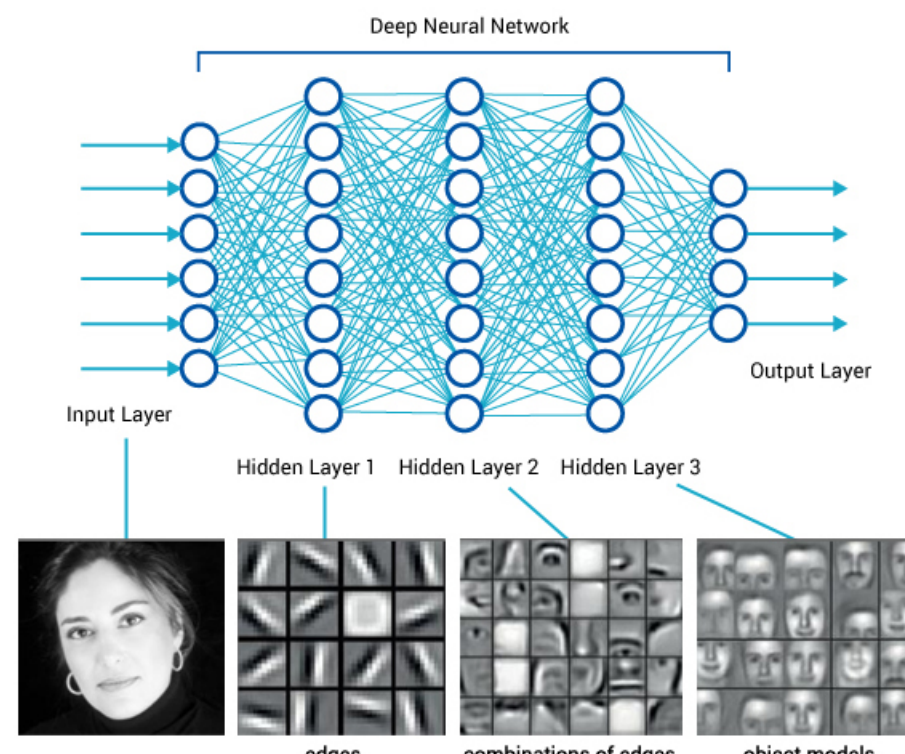
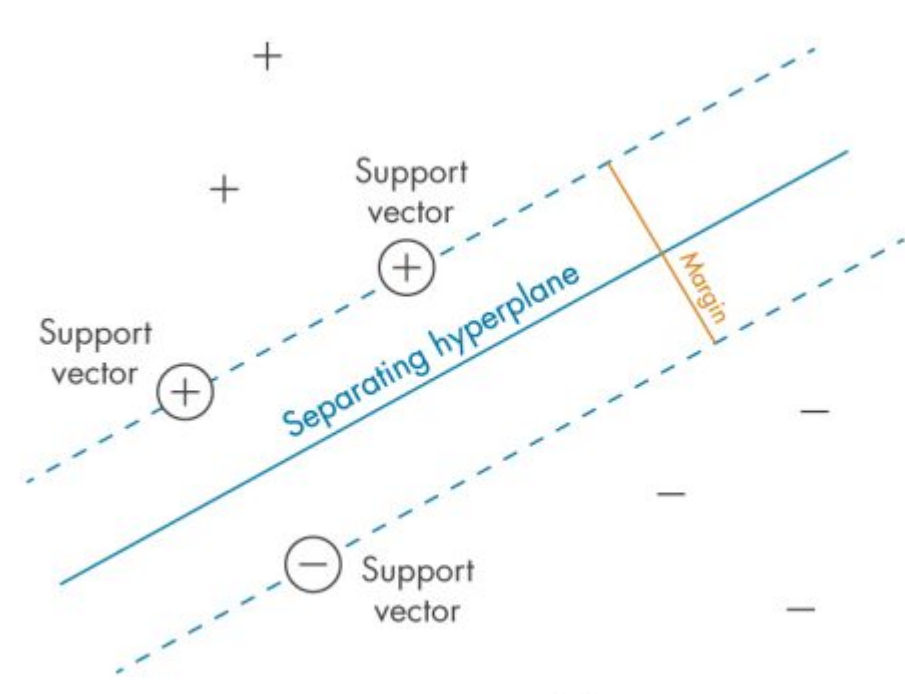


Motivation

Tractable models for deep learning



What we would like to understand



What we actually have theory for

- no quantitative theory for deep learning yet
- exact understanding of simple models (e.g., linear SVM)
- simple models sometimes capture deep learning phenomenology

Need for tractable, realistic surrogate models for deep learning

Ingredients for a surrogate model

- learning architecture** (ridge regression, support vector machine, ...)
- data/feature model** (i.i.d. Gaussian, non-diagonal covariance ...)
- training algorithm (not here, direct focus on estimators)

Examples

- Instances of ridge regression with i.i.d. coordinates captures the so-called **double descent** [BHMM19] phenomenon
- Gaussian mixtures** are appropriate models for GAN data [SLTC20]
- Convex Generalized Linear Models** (GLM) with correlated Gaussian designs [LGC⁺21] capture a wide range of single task regression problems, with structured data/feature maps (kernels, GAN data, ...)

Objective

Can we have a realistic benchmark for multiclass classification problems ?

Contributions

- study classification of a high-dimensional **K-Gaussian mixture with a convex Generalized Linear Model (GLM)**
- generic means and covariances for the clusters**
- exact asymptotic distribution of the estimator**
- study of both random design and real data problems**

The generative model : a K-Gaussian mixture

Consider the Gaussian mixture density with K cluster $\{C_k\}_{1 \leq k \leq K}$:

$$P(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^K y_k \rho_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (1)$$

- means $\boldsymbol{\mu}_k \in \mathbb{R}^d$, covariances $\boldsymbol{\Sigma}_k \in \mathbb{R}^{d \times d}$ positive definite
- cluster membership $\rho_k \in [0, 1]$ with $\sum_k \rho_k = 1$
- labels y_k are **one-hot-encoded**, i.e. cluster C_k is denoted by the k -th basis vector in \mathbb{R}^K
- sample n pairs $(\mathbf{x}^\nu, \mathbf{y}^\nu)$ from Eq.(1)
- design matrix denoted $\mathbf{X} \in \mathbb{R}^{n \times d}$

Learn **K separating hyperplanes** in \mathbb{R}^d : $\mathbf{W}^* \in \mathbb{R}^{K \times d}$

The learning method : a convex GLM

Estimator obtained by minimising the empirical risk:

$$\mathcal{R}(\mathbf{W}, \mathbf{b}) \equiv \sum_{\nu=1}^n \ell\left(\mathbf{y}^\nu, \frac{\mathbf{W} \mathbf{x}^\nu}{\sqrt{d}} + \mathbf{b}\right) + \lambda r(\mathbf{W}), \quad (2)$$

$$(\mathbf{W}^*, \mathbf{b}^*) \equiv \underset{\mathbf{W} \in \mathbb{R}^{K \times d}, \mathbf{b} \in \mathbb{R}^K}{\operatorname{argmin}} \mathcal{R}(\mathbf{W}, \mathbf{b}), \quad (3)$$

- $\mathbf{W} \in \mathbb{R}^{K \times d}$, $\mathbf{b} \in \mathbb{R}^K$ are the weights and bias to be learned
- ℓ is a convex loss function
- r is a convex regularisation function with strength $\lambda \in \mathbb{R}$

Examples : least-squares, logistic loss, ℓ_2 or ℓ_1 penalty, ...

Goal : asymptotic properties of \mathbf{W}^*

High-dimensional limit : $n, d \rightarrow \infty$ with fixed $\alpha = n/d$
We characterise the asymptotic distribution of the estimator $(\mathbf{W}^*, \mathbf{b}^*)$.

In particular, we are interested in:

- the average training loss

$$\epsilon_\ell = \frac{1}{n} \sum_{\nu=1}^n \ell\left(\mathbf{y}^\nu, \frac{\mathbf{W}^* \mathbf{x}^\nu}{\sqrt{d}} + \mathbf{b}^*\right), \quad (4)$$

- the average training error ϵ_t and generalisation error ϵ_g :

$$\epsilon_t = \frac{1}{n} \sum_{\nu=1}^n \mathbb{I}\left[\mathbf{y}^\nu \neq \hat{\mathbf{y}}\left(\frac{\mathbf{W}^* \mathbf{x}^\nu}{\sqrt{d}} + \mathbf{b}^*\right)\right], \quad (5)$$

$$\epsilon_g = \mathbb{E}_{(\mathbf{x}^{\text{new}}, \mathbf{y}^{\text{new}})} \left[\mathbb{I}\left[\mathbf{y}^{\text{new}} \neq \hat{\mathbf{y}}\left(\frac{\mathbf{W}^* \mathbf{x}^{\text{new}}}{\sqrt{d}} + \mathbf{b}^*\right)\right] \right],$$

where $(\mathbf{x}^{\text{new}}, \mathbf{y}^{\text{new}})$ is a new sample from Eq. (1), and $\hat{\mathbf{y}}_k(\mathbf{x}) = \mathbb{I}(\max_{k'} x_{k'} = x_k)$.

Useful notation

Suppose that the matrix $\mathbf{G} = (G_{ki})_{ki} \in \mathbb{R}^{K \times d}$ is given, alongside the four-index tensor $\mathbf{A} = (A_{kik'j'})_{kik'j'} \in \mathbb{R}^{K \times d} \otimes \mathbb{R}^{K \times d}$. We will use the notation $\mathbf{G} \odot \mathbf{A} = \sum_{ki} G_{ki} A_{kik'j'} \in \mathbb{R}^{K \times d}$. Similarly, given a four-index tensor \mathbf{A} , we will define $\sqrt{\mathbf{A}}$ as the tensor such that $\mathbf{A} = \sqrt{\mathbf{A}} \odot \sqrt{\mathbf{A}}$.

Main result : exact asymptotics [LSG⁺21]

- Let $\boldsymbol{\xi}_{k \in [K]} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ be collection of K -dimensional standard normal vectors independent of other quantities.
- let $\{\boldsymbol{\Xi}_k\}$ a set of K matrices, $\boldsymbol{\Xi}_k \in \mathbb{R}^{K \times d}$, with i.i.d. standard normal entries, independent of other quantities.
- let $\mathbf{Z}^* = \frac{1}{\sqrt{d}} \mathbf{W}^* \mathbf{X} \in \mathbb{R}^{K \times n}$

Under mild feasibility and regularity assumptions, for any pseudo-Lispchitz functions $\phi_1 : \mathbb{R}^{K \times d} \rightarrow \mathbb{R}$, $\phi_2 : \mathbb{R}^{K \times n} \rightarrow \mathbb{R}$:

$$\phi_1(\mathbf{W}^*) \xrightarrow[n, d \rightarrow +\infty]{P} \mathbb{E}_{\boldsymbol{\Xi}}[\phi_1(\mathbf{G})], \quad \phi_2(\mathbf{Z}^*) \xrightarrow[n, d \rightarrow +\infty]{P} \mathbb{E}_{\boldsymbol{\xi}}[\phi_2(\mathbf{H})],$$

where we have introduced the proximal for the loss:

$$\mathbf{h}_k = \mathbf{V}_k^{1/2} \operatorname{Prox}_{\ell(\mathbf{e}_k, \mathbf{V}_k^{1/2} \cdot)}(\mathbf{V}_k^{-1/2} \boldsymbol{\omega}_k) \in \mathbb{R}^K$$

$$\boldsymbol{\omega}_k \equiv \mathbf{m}_k + \mathbf{b} + \mathbf{Q}_k^{1/2} \boldsymbol{\xi}_k,$$

and $\mathbf{H} \in \mathbb{R}^{K \times n}$ is obtained by concatenating each \mathbf{h}_k , $\rho_k n$ times.

We have also introduced the matrix proximal $\mathbf{G} \in \mathbb{R}^{K \times d}$:

$$\mathbf{G} = \mathbf{A}^{\frac{1}{2}} \odot \operatorname{Prox}_{r(\mathbf{A}^{\frac{1}{2}} \odot \cdot)}(\mathbf{A}^{\frac{1}{2}} \odot \mathbf{B}), \quad \mathbf{A}^{-1} \equiv \sum_k \hat{\mathbf{V}}_k \otimes \boldsymbol{\Sigma}_k,$$

$$\mathbf{B} \equiv \sum_k \left(\boldsymbol{\mu}_k \hat{\mathbf{m}}_k^\top + \boldsymbol{\Xi}_k \odot \sqrt{\hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k} \right).$$

The collection of parameters $(\mathbf{Q}_k, \mathbf{m}_k, \mathbf{V}_k, \hat{\mathbf{Q}}_k, \hat{\mathbf{m}}_k, \hat{\mathbf{V}}_k)_{k \in [K]}$ is given by the fixed point of the following self-consistent equations:

$$\begin{cases} \mathbf{Q}_k = \frac{1}{d} \mathbb{E}_{\boldsymbol{\Xi}}[\mathbf{G} \boldsymbol{\Sigma}_k \mathbf{G}^\top] \\ \mathbf{m}_k = \frac{1}{\sqrt{d}} \mathbb{E}_{\boldsymbol{\Xi}}[\mathbf{G} \boldsymbol{\mu}_k] \\ \mathbf{V}_k = \frac{1}{d} \mathbb{E}_{\boldsymbol{\Xi}} \left[\left(\mathbf{G} \odot (\hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k) \right)^{-\frac{1}{2}} \odot (\mathbf{I}_K \otimes \boldsymbol{\Sigma}_k) \right] \boldsymbol{\Xi}_k^\top \\ \hat{\mathbf{Q}}_k = \alpha \rho_k \mathbb{E}_{\boldsymbol{\xi}}[\mathbf{f}_k \mathbf{f}_k^\top] \\ \hat{\mathbf{V}}_k = -\alpha \rho_k \mathbf{Q}_k^{-\frac{1}{2}} \mathbb{E}_{\boldsymbol{\xi}}[\mathbf{f}_k \boldsymbol{\xi}^\top] \\ \hat{\mathbf{m}}_k = \alpha \rho_k \mathbb{E}_{\boldsymbol{\xi}}[\mathbf{f}_k] \end{cases}$$

where $\mathbf{f}_k \equiv \mathbf{V}_k^{-1}(\mathbf{h}_k - \boldsymbol{\omega}_k)$, and the vector \mathbf{b}^* is such that $\sum_k \rho_k \mathbb{E}_{\boldsymbol{\xi}}[\mathbf{V}_k \mathbf{f}_k] = \mathbf{0}$ holds.

Important remarks

- very generic statement
- proximal operators are easy to compute, summarize the effect of loss and penalty
- greatly simplifies with assumptions on covariances, separability of functions, ...
- in most cases reduces to low/one dimensional statement

Corollary : training and generalisation error

The training loss, the training error and the generalisation error are given by

$$\epsilon_\ell = \sum_{k=1}^K \rho_k \mathbb{E}_{\boldsymbol{\xi}}[\ell(\mathbf{e}_k, \mathbf{h}_k)], \quad (6)$$

$$\epsilon_t = 1 - \sum_{k=1}^K \rho_k \mathbb{E}_{\boldsymbol{\xi}}[\hat{y}_k(\mathbf{h}_k)], \quad (7)$$

$$\epsilon_g = 1 - \sum_{k=1}^K \rho_k \mathbb{E}_{\boldsymbol{\xi}}[\hat{y}_k(\boldsymbol{\omega}_k)]. \quad (8)$$

Results on a synthetic dataset

- multiclass logistic regression with ridge penalty
- effect of sample complexity, number of clusters and regularisation strength
- recover and extend previous result on separability transition [MKL⁺20]

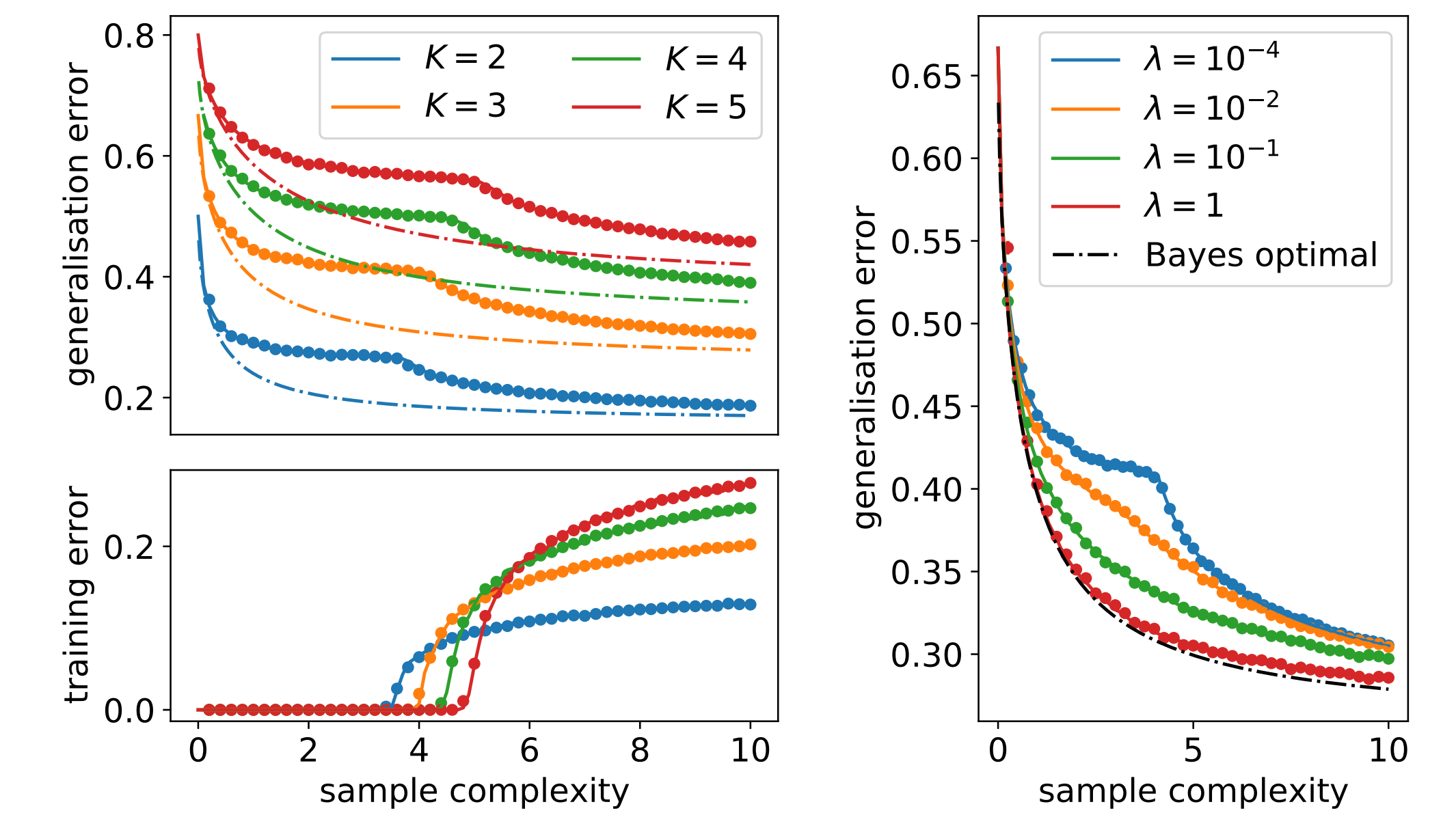


Figure: Gaussian means and $\boldsymbol{\Sigma}_k \equiv \Sigma = 1/2 \mathbf{I}_d$. (Left) Generalisation error ϵ_g (top) and training error ϵ_t (bottom) as function of α at $\lambda = 10^{-4}$. Theoretical predictions (full lines) are compared with the results of numerical experiments (dots). Dash-dotted lines of the corresponding color represent, for comparison, the Bayes-optimal error. (Right) Dependence of the generalisation error on the regularization λ for $K = 3$ and $\Delta = 1/2$, $\rho_k = 1/K$

Results on a real dataset

- binary classification with the logistic loss on MNIST/Fashion-MNIST
- comparison between the estimator obtained with real data and a synthetic (Gaussian) approximation with matching covariances
- the real learning curve is captured by the synthetic one**

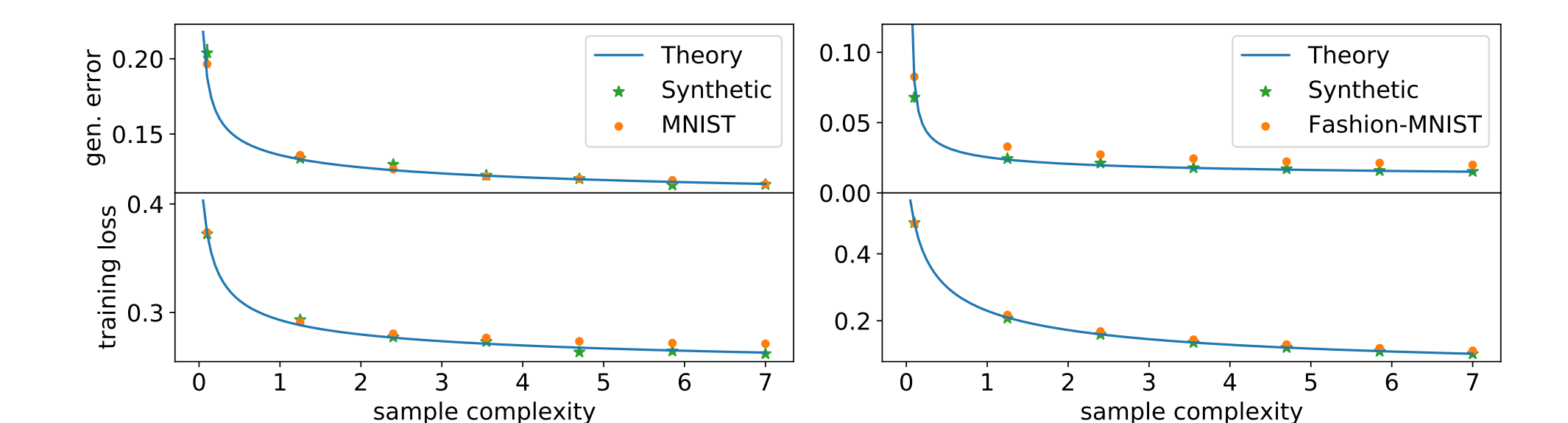


Figure: Generalisation error and training loss on MNIST with $\lambda = 0.05$ (left) and on Fashion-MNIST with $\lambda = 1$ (right)

References

- [BHMM19] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [LGC⁺21] Bruno Loureiro, Cédric Gerbelot, Hugo Cui, Sebastian Goldt, Florent Krzakala, Marc Mézard, and Lenka Zdeborová. Capturing the learning curves of generic features maps for realistic data sets with a teacher-student model. *arXiv preprint arXiv:2102.08127*, 2021.
- [LSG⁺21] Bruno Loureiro, Gabriele Sicuro, Cédric Gerbelot, Alessandro Pacco, Florent Krzakala, and Lenka Zdeborová. Learning gaussian mixtures with generalised linear models: Precise asymptotics in high-dimensions. *arXiv preprint arXiv:2106.03791*, 2021.
- [MKL⁺20] Francesca Mignacco, Florent Krzakala, Yue Lu, Pierfrancesco Urbani, and Lenka Zdeborová. The role of regularization in classification of high-dimensional noisy Gaussian mixture. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 6874–6883. PMLR, 13–18 Jul 2020.
- [SLTC20] Mohamed El Amine Seddik, Cosme Louart, Mohamed Tamaazousti, and Romain Couillet. Random matrix theory proves that deep learning representations of GAN-data behave as Gaussian mixtures. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 8573–8582. PMLR, 13–18 Jul 2020.

Contact: cedric.gerbelot@ens.fr