

The Breakdown of Gaussian Universality in Classification of High-Dimensional Linear Factor Mixtures

Xiaoyi MAI^{*}, Zhenyu LIAO[†]

^{*}University of Toulouse - Jean Jaurès, Toulouse Mathematics Institute, France

[†]EIC Huazhong University of Science & Technology, China

ICLR, 2025

Introduction: Empirical Risk Minimization

Supervised ML: building a classifier from a training set of $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with feature vectors $\mathbf{x}_i \in \mathbb{R}^p$ and class labels $y_i \in \{-1, 1\}$.

Introduction: Empirical Risk Minimization

Supervised ML: building a classifier from a training set of $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with feature vectors $\mathbf{x}_i \in \mathbb{R}^p$ and class labels $y_i \in \{-1, 1\}$.

Empirical Risk Minimisation (ERM) : many supervised algorithms (e.g., SVM, LR, ANN with pretrained hidden layers) can be summarized by the following ERM

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(y_i, \beta^\top \mathbf{x}_i)}_{\text{Empirical loss}} + \underbrace{\lambda \|\beta\|^2}_{\text{regularization term}} .$$

Introduction: Empirical Risk Minimization

Supervised ML: building a classifier from a training set of $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with feature vectors $\mathbf{x}_i \in \mathbb{R}^p$ and class labels $y_i \in \{-1, 1\}$.

Empirical Risk Minimisation (ERM) : many supervised algorithms (e.g., SVM, LR, ANN with pretrained hidden layers) can be summarized by the following ERM

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(y_i, \beta^\top \mathbf{x}_i)}_{\text{Empirical loss}} + \underbrace{\lambda \|\beta\|^2}_{\text{regularization term}} .$$

- Popular choices of ℓ : square loss for least squares method, logistic loss for LR, hinge loss for SVM, etc.

Introduction: Performance Analysis in the Big Data Regime

Modern ML: comparably numerous features and samples, i.e., $p \sim n \gg 1$.



Figure: Image Classification

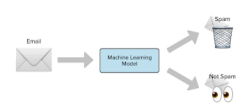


Figure: Spam Detection

Introduction: Performance Analysis in the Big Data Regime

Modern ML: comparably numerous features and samples, i.e., $p \sim n \gg 1$.



Figure: Image Classification

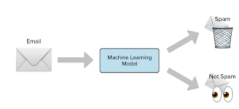


Figure: Spam Detection

Analysis of Modern ML:

- Complications of $n \sim p$:
 - Performance **sensitive** to the sample size n/p , and the hyperparameters ℓ, λ .
 - **Random** classifier $\hat{\beta}$ depending on $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ in a **non-linear** and **implicit** manner.

Introduction: Performance Analysis in the Big Data Regime

Modern ML: comparably numerous features and samples, i.e., $p \sim n \gg 1$.



Figure: Image Classification

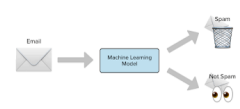


Figure: Spam Detection

Analysis of Modern ML:

- Complications of $n \sim p$:
 - Performance **sensitive** to the sample size n/p , and the hyperparameters ℓ, λ .
 - **Random** classifier $\hat{\beta}$ depending on $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ in a **non-linear** and **implicit** manner.
- Technical conveniences of $n, p \gg 1$:
 - **Convergence of performance curve** as a function of sample ratio n/p .
 - **Gaussian universality** (GU) induced by **Central Limit Theorem** (e.g., $\hat{\beta} = \text{avg}(y_i \mathbf{x}_i)$).

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the first two moments.

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the **first two moments**.

- ▶ Advantage: use of tools requiring data Gaussianity , e.g., CGMT, AMP, replica.

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the **first two moments**.

- ▶ Advantage: use of tools requiring data Gaussianity , e.g., CGMT, AMP, replica.
- ▶ Limitation: no insight into the learning of **higher order data statistics**.

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the **first two moments**.

- ▶ Advantage: use of tools requiring data Gaussianity , e.g., CGMT, AMP, replica.
- ▶ Limitation: no insight into the learning of **higher order data statistics**.

Previous results on GU of ERM: proofs of GU under the assumptions of **pointwise normality** $\hat{\beta}^T \mathbf{x}$ on the ERM solution $\hat{\beta}$ [MS22; Dan+24], **not verifiable** from the data distribution.

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the **first two moments**.

- ▶ Advantage: use of tools requiring data Gaussianity , e.g., CGMT, AMP, replica.
- ▶ Limitation: no insight into the learning of **higher order data statistics**.

Previous results on GU of ERM: proofs of GU under the assumptions of **pointwise normality** $\hat{\beta}^T \mathbf{x}$ on the ERM solution $\hat{\beta}$ [MS22; Dan+24], **not verifiable** from the data distribution.

Our contributions [ML25]:

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the **first two moments**.

- ▶ Advantage: use of tools requiring data Gaussianity , e.g., CGMT, AMP, replica.
- ▶ Limitation: no insight into the learning of **higher order data statistics**.

Previous results on GU of ERM: proofs of GU under the assumptions of **pointwise normality** $\hat{\beta}^T \mathbf{x}$ on the ERM solution $\hat{\beta}$ [MS22; Dan+24], **not verifiable** from the data distribution.

Our contributions [ML25]:

- ▶ Sharp performance analysis of ERM under **linear factor mixture models** (LFMM).
 - ▶ The development of **leave-one-out** approach to characterize performance beyond GU.

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the **first two moments**.

- ▶ Advantage: use of tools requiring data Gaussianity , e.g., CGMT, AMP, replica.
- ▶ Limitation: no insight into the learning of **higher order data statistics**.

Previous results on GU of ERM: proofs of GU under the assumptions of **pointwise normality** $\hat{\beta}^T \mathbf{x}$ on the ERM solution $\hat{\beta}$ [MS22; Dan+24], **not verifiable** from the data distribution.

Our contributions [ML25]:

- ▶ Sharp performance analysis of ERM under **linear factor mixture models** (LFMM).
 - ▶ The development of **leave-one-out** approach to characterize performance beyond GU.
- ▶ **Condition** of GU depending on LFMM. .
 - ▶ Explicit condition allowing insight into the impact of **data structure**.

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the **first two moments**.

- ▶ Advantage: use of tools requiring data Gaussianity , e.g., CGMT, AMP, replica.
- ▶ Limitation: no insight into the learning of **higher order data statistics**.

Previous results on GU of ERM: proofs of GU under the assumptions of **pointwise normality** $\hat{\beta}^T \mathbf{x}$ on the ERM solution $\hat{\beta}$ [MS22; Dan+24], **not verifiable** from the data distribution.

Our contributions [ML25]:

- ▶ Sharp performance analysis of ERM under **linear factor mixture models** (LFMM).
 - ▶ The development of **leave-one-out** approach to characterize performance beyond GU.
- ▶ **Condition** of GU depending on LFMM. .
 - ▶ Explicit condition allowing insight into the impact of **data structure**.
- ▶ Implication of GU breakdown on the **optimal choice of loss**.

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the **first two moments**.

- ▶ Advantage: use of tools requiring data Gaussianity , e.g., CGMT, AMP, replica.
- ▶ Limitation: no insight into the learning of **higher order data statistics**.

Previous results on GU of ERM: proofs of GU under the assumptions of **pointwise normality** $\hat{\beta}^T \mathbf{x}$ on the ERM solution $\hat{\beta}$ [MS22; Dan+24], **not verifiable** from the data distribution.

Our contributions [ML25]:

- ▶ Sharp performance analysis of ERM under **linear factor mixture models** (LFMM).
 - ▶ The development of **leave-one-out** approach to characterize performance beyond GU.
- ▶ **Condition** of GU depending on LFMM. .
 - ▶ Explicit condition allowing insight into the impact of **data structure**.
- ▶ Implication of GU breakdown on the **optimal choice of loss**.
 - ▶ Consequence: square loss no longer optimal as in GMM [TPT20; ML20].

Introduction: Gaussian Universality of Empirical Risk Minimization

Implication of GU : universal performance depending on the **first two moments**.

- ▶ Advantage: use of tools requiring data Gaussianity , e.g., CGMT, AMP, replica.
- ▶ Limitation: no insight into the learning of **higher order data statistics**.

Previous results on GU of ERM: proofs of GU under the assumptions of **pointwise normality** $\hat{\beta}^T \mathbf{x}$ on the ERM solution $\hat{\beta}$ [MS22; Dan+24], **not verifiable** from the data distribution.

Our contributions [ML25]:

- ▶ Sharp performance analysis of ERM under **linear factor mixture models** (LFMM).
 - ▶ The development of **leave-one-out** approach to characterize performance beyond GU.
- ▶ **Condition** of GU depending on LFMM. .
 - ▶ Explicit condition allowing insight into the impact of **data structure**.
- ▶ Implication of GU breakdown on the **optimal choice of loss**.
 - ▶ Consequence: square loss no longer optimal as in GMM [TPT20; ML20].
 - ▶ Sharp results on the advantage of non-square losses in learning high-order data statistics.

Linear Factor Mixture Model

Definition (Linear Factor Mixture Model (LFMM))

A data-label pair $(\mathbf{x}, y) \sim \mathcal{D}_{(\mathbf{x}, y)}$ with class label $y \in \{\pm 1\}$ is said to follow a linear factor mixture model if $\mathbf{x} \in \mathbb{R}^p$ is the linear combination of p **factors** z_1, \dots, z_p

$$\mathbf{x} = \sum_{k=1}^p z_k \mathbf{v}_k = \sum_{k=1}^p (y s_k + e_k) \mathbf{v}_k, \quad (1)$$

for linearly independent deterministic $\mathbf{v}_1, \dots, \mathbf{v}_p \in \mathbb{R}^p$ and standardized noises $e_1, \dots, e_p \in \mathbb{R}$ independent of y with bounded fourth moments.

- q **informative factors** z_1, \dots, z_q with **deterministic signals** $s_k > 0$, $k \in \{1, \dots, q\}$;

Linear Factor Mixture Model

Definition (Linear Factor Mixture Model (LFMM))

A data-label pair $(\mathbf{x}, y) \sim \mathcal{D}_{(\mathbf{x}, y)}$ with class label $y \in \{\pm 1\}$ is said to follow a linear factor mixture model if $\mathbf{x} \in \mathbb{R}^p$ is the linear combination of p **factors** z_1, \dots, z_p

$$\mathbf{x} = \sum_{k=1}^p z_k \mathbf{v}_k = \sum_{k=1}^p (y s_k + e_k) \mathbf{v}_k, \quad (1)$$

for linearly independent deterministic $\mathbf{v}_1, \dots, \mathbf{v}_p \in \mathbb{R}^p$ and standardized noises $e_1, \dots, e_p \in \mathbb{R}$ independent of y with bounded fourth moments.

- ▶ q **informative factors** z_1, \dots, z_q with **deterministic signals** $s_k > 0$, $k \in \{1, \dots, q\}$;
- ▶ $p - q$ **noise factors** z_{q+1}, \dots, z_p with $s_k = 0$, $k \in \{q + 1, \dots, p\}$.

Linear Factor Mixture Model

Definition (Linear Factor Mixture Model (LFMM))

A data-label pair $(\mathbf{x}, y) \sim \mathcal{D}_{(\mathbf{x}, y)}$ with class label $y \in \{\pm 1\}$ is said to follow a linear factor mixture model if $\mathbf{x} \in \mathbb{R}^p$ is the linear combination of p **factors** z_1, \dots, z_p

$$\mathbf{x} = \sum_{k=1}^p z_k \mathbf{v}_k = \sum_{k=1}^p (y s_k + e_k) \mathbf{v}_k, \quad (1)$$

for linearly independent deterministic $\mathbf{v}_1, \dots, \mathbf{v}_p \in \mathbb{R}^p$ and standardized noises $e_1, \dots, e_p \in \mathbb{R}$ independent of y with bounded fourth moments.

- ▶ q **informative factors** z_1, \dots, z_q with **deterministic signals** $s_k > 0$, $k \in \{1, \dots, q\}$;
- ▶ $p - q$ **noise factors** z_{q+1}, \dots, z_p with $s_k = 0$, $k \in \{q+1, \dots, p\}$.
- ▶ $\text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_q\}$ orthogonal to $\text{Span}\{\mathbf{v}_{q+1}, \dots, \mathbf{v}_p\}$

Linear Factor Mixture Model

Definition (Linear Factor Mixture Model (LFMM))

A data-label pair $(\mathbf{x}, y) \sim \mathcal{D}_{(\mathbf{x}, y)}$ with class label $y \in \{\pm 1\}$ is said to follow a linear factor mixture model if $\mathbf{x} \in \mathbb{R}^p$ is the linear combination of p **factors** z_1, \dots, z_p

$$\mathbf{x} = \sum_{k=1}^p z_k \mathbf{v}_k = \sum_{k=1}^p (y s_k + e_k) \mathbf{v}_k, \quad (1)$$

for linearly independent deterministic $\mathbf{v}_1, \dots, \mathbf{v}_p \in \mathbb{R}^p$ and standardized noises $e_1, \dots, e_p \in \mathbb{R}$ independent of y with bounded fourth moments.

- ▶ q **informative factors** z_1, \dots, z_q with **deterministic signals** $s_k > 0$, $k \in \{1, \dots, q\}$;
- ▶ $p - q$ **noise factors** z_{q+1}, \dots, z_p with $s_k = 0$, $k \in \{q+1, \dots, p\}$.
- ▶ $\text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_q\}$ orthogonal to $\text{Span}\{\mathbf{v}_{q+1}, \dots, \mathbf{v}_p\}$
- ▶ class-conditional means and covariances of \mathbf{x} :

$$\mathbb{E}[\mathbf{x}|y] = y\boldsymbol{\mu}, \quad \text{Cov}[\mathbf{x}|y] = \boldsymbol{\Sigma} \quad (2)$$

with $\boldsymbol{\mu} = \sum_{k=1}^p s_k \mathbf{v}_k$ and $\boldsymbol{\Sigma} = \sum_{k=1}^p \mathbf{v}_k \mathbf{v}_k^\top$.

Theorem (Asymptotic distribution of predicted scores under LFMM)

For ERM classifier $\hat{\beta}$ obtained on $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ of size n drawn i.i.d. from an LFMM, we have that, for any bounded Lipschitz function $f: \mathbb{R} \rightarrow \mathbb{R}$,

$$\text{Testing score: } \boxed{\mathbb{E}[f(\hat{\beta}^\top \boldsymbol{\nu})] - \mathbb{E}[f(\tilde{\beta}^\top \boldsymbol{\nu})] \rightarrow 0}, \quad \forall \text{ deterministic } \boldsymbol{\nu} \in \mathbb{R}^p$$

$$\text{Training score: } \boxed{\mathbb{E}[f(\hat{\beta}^\top \mathbf{x}_i)] - \mathbb{E}[f(\text{prox}_{\kappa, \ell(\cdot, y_i)}(\tilde{\beta}^\top \mathbf{x}_i))] \rightarrow 0}, \quad \forall i \in \{1, \dots, n\},$$

where

$$\tilde{\beta} = (\lambda \mathbf{I}_p + \theta \boldsymbol{\Sigma})^{-1} \left(\eta \boldsymbol{\mu} + \sum_{k=1}^q \omega_k \mathbf{v}_k + \gamma \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{u} \right), \quad (3)$$

for Gaussian vector $\mathbf{u} \sim \mathcal{N}(\mathbf{0}_p, \mathbf{I}_p/n)$ independent of $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and constants $\theta, \eta, \gamma, \omega_1, \dots, \omega_q$ determined by a (known) self-consistent system of equations

$$[\theta, \eta, \gamma, \omega_1, \dots, \omega_q] = G_{\ell, \lambda, n/p, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{v}_1, \dots, \mathbf{v}_q, \mathcal{D}_{(z_1, \dots, z_q)}}([\theta, \eta, \gamma, \omega_1, \dots, \omega_q]).$$

Definition of Gaussian Universality

Definition (Equivalent Gaussian mixture model (Equivalent GMM))

For an LFMM $\mathcal{D}_{(\mathbf{x}, y)}$, we define its equivalent Gaussian mixture model $\mathcal{D}_{(\mathbf{g}, y)}$ as the GMM with the same class-conditional means $\boldsymbol{\mu}$ and covariances $\boldsymbol{\Sigma}$ as the LFMM in (2):

$$\mathbf{g} \sim \mathcal{N}(y\boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (4)$$

We denote $\hat{\beta}^{\mathbf{g}}$ the ERM solution obtained on n i.i.d. $(\mathbf{g}_1, y_1), \dots, (\mathbf{g}_n, y_n) \sim \mathcal{D}_{(\mathbf{g}, y)}$.

Definition of Gaussian Universality

Definition (Equivalent Gaussian mixture model (Equivalent GMM))

For an LFMM $\mathcal{D}_{(\mathbf{x}, y)}$, we define its equivalent Gaussian mixture model $\mathcal{D}_{(\mathbf{g}, y)}$ as the GMM with the same class-conditional means $\boldsymbol{\mu}$ and covariances $\boldsymbol{\Sigma}$ as the LFMM in (2):

$$\mathbf{g} \sim \mathcal{N}(y\boldsymbol{\mu}, \boldsymbol{\Sigma}). \quad (4)$$

We denote $\hat{\beta}^{\mathbf{g}}$ the ERM solution obtained on n i.i.d. $(\mathbf{g}_1, y_1), \dots, (\mathbf{g}_n, y_n) \sim \mathcal{D}_{(\mathbf{g}, y)}$.

Definition (Gaussian universality under LFMM)

For an ERM solution $\hat{\beta}$ on LFMM $\mathcal{D}_{(\mathbf{x}, y)}$ and $\hat{\beta}^{\mathbf{g}}$ on the equivalent GMM, we say **Gaussian universality** holds if

$$\Pr(y_i \mathbf{x}_i^T \hat{\beta} > 0) \simeq \Pr(y_i \mathbf{g}_i^T \hat{\beta}^{\mathbf{g}} > 0)$$

$$\Pr(y' \mathbf{x}'^T \hat{\beta} > 0) \simeq \Pr(y' \mathbf{g}'^T \hat{\beta}^{\mathbf{g}} > 0)$$

for $(\mathbf{x}', y') \sim \mathcal{D}_{(\mathbf{x}, y)}$ independent of $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, and $(\mathbf{g}', y') \sim \mathcal{D}_{(\mathbf{g}, y)}$ independent of $\{(\mathbf{g}_i, y_i)\}_{i=1}^n$

Conditional Gaussian universality under LFMM

Conditional Gaussian universality under LFMM

Corollary (Conditional Gaussian universality under LFMM)

*Under LFMM, the Gaussian universality of ERM holds if and only if the **informative factors** z_1, \dots, z_q of LFMM are class-conditional **Gaussian**.*

Conditional Gaussian universality under LFMM

Corollary (Conditional Gaussian universality under LFMM)

Under LFMM, the Gaussian universality of ERM holds if and only if the *informative factors* z_1, \dots, z_q of LFMM are class-conditional *Gaussian*.

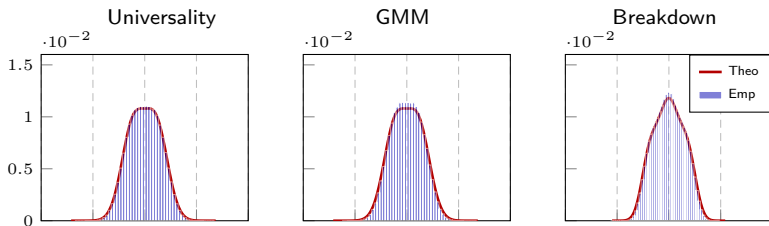


Figure: Theoretical and empirical distribution of predicted scores $\hat{\beta}^T \mathbf{x}'$ for some fresh test data $(\mathbf{x}', y') \sim \mathcal{D}_{(\mathbf{x}, y)}$ independent of $\hat{\beta}$. The theoretical probability densities (red), and the empirical histograms (blue) are the values of $\hat{\beta}^T \mathbf{x}'$ over 10^6 independent copies of \mathbf{x}' , for three different LFMMs with $n = 600$, $p = 200$, $\rho = 0.5$, $s = \lceil \sqrt{2} \rceil; \mathbf{0}_{p-1}$ (so that $q = 1$), and Haar distributed \mathbf{V} . **Left:** normal e_1 and uniformly distributed e_2, \dots, e_p ; **Middle:** normal e_1, \dots, e_p ; **Right:** uniformly distributed e_1 , and normal e_2, \dots, e_p .

Performance under Gaussian Universality Breakdown

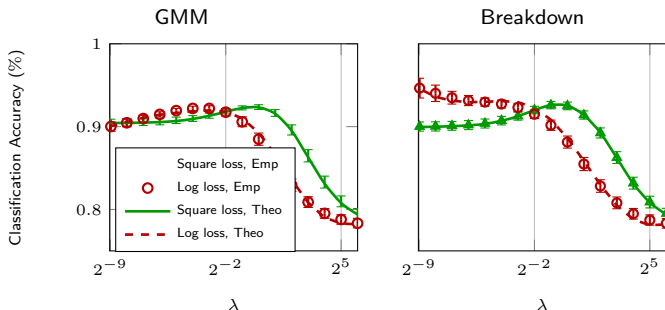


Figure: Empirical classification accuracy of $\hat{\mathbf{w}}_{\ell, \lambda}$ averaged over 100 trials with a width of ± 1 standard deviation, versus theoretical curve given by the square loss and the logistic loss on $n = 800$ training samples. **Left:** GMM under with $p = 200$, $\rho = 0.5$, $s = [1, 5; 0.5; \mathbf{0}_{p-2}]$ (so that $q = 2$), and $\mathbf{V} = \text{diag}(2, \mathbf{1}_{p-1})\mathbf{H}$ with Haar distributed \mathbf{H} . **Right:** LFMM identical to the GMM in the left, but with Rademacher e_1 .

Performance under Gaussian Universality Breakdown

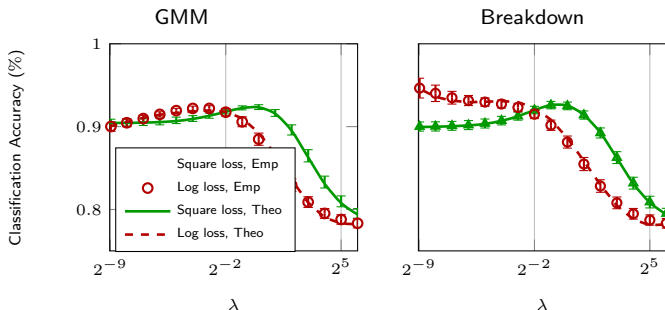
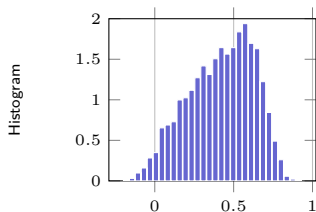


Figure: Empirical classification accuracy of $\hat{\mathbf{w}}_{\ell, \lambda}$ averaged over 100 trials with a width of ± 1 standard deviation, versus theoretical curve given by the square loss and the logistic loss on $n = 800$ training samples. **Left:** GMM under with $p = 200$, $\rho = 0.5$, $s = [1, 5; 0.5; \mathbf{0}_{p-2}]$ (so that $q = 2$), and $\mathbf{V} = \text{diag}(2, \mathbf{1}_{p-1})\mathbf{H}$ with Haar distributed \mathbf{H} . **Right:** LFMM identical to the GMM in the left, but with Rademacher e_1 .

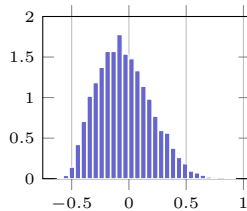
Remark: under LFMM, the square loss is **no longer optima** as in GMM [TPT20; ML20].

Experiments on Real Image Data

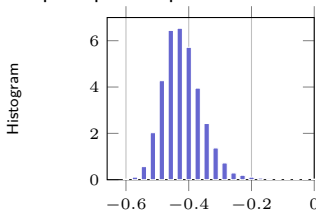
Case 1: Classes 3&7 of Fashion-MNIST data, for which *approximately Gaussian* informative factors (estimated by the principal components in PCA) can be observed.



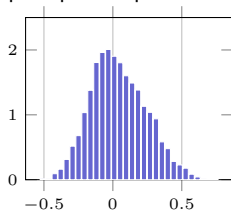
First principal component of Class 3



Second principal component of Class 3



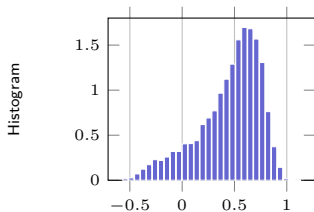
First principal component of Class 7



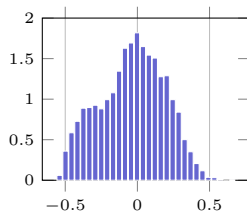
Second principal component of Class 7

Experiments on Real Image Data

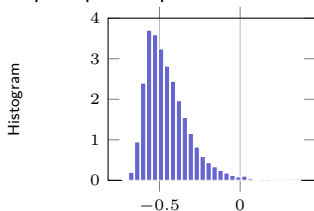
Case 2: Classes 4&5 of Fashion-MNIST data, as an example of *non-Gaussian* informative factors (estimated by the principal components in PCA).



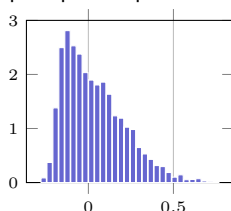
First principal component of Class 4



Second principal component of Class 4



First principal component of Class 5



Second principal component of Class 5

Experiments on Real Image Data

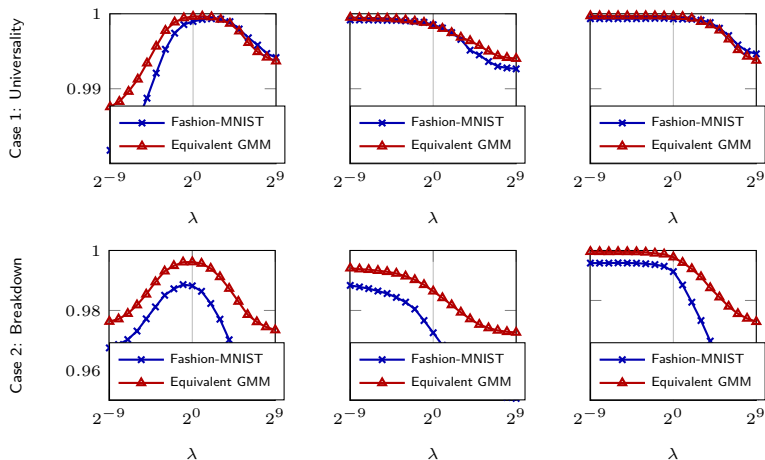


Figure: Classification accuracies as a function of the regularization penalty γ , for Fashion-MNIST data and Equivalent GMM of sample size $n = 512$, with square (**left**), logistic (**middle**), and square hinge (**right**) losses.

Experiments on Real Image Data

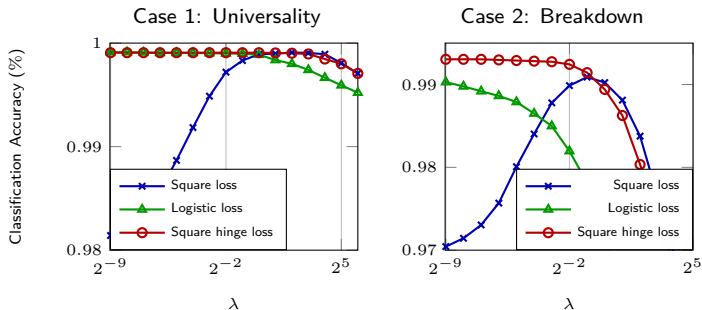


Figure: Classification accuracies as a function of the regularization penalty, for square, logistic, and square hinge loss, on Fashion-MNIST data of sample size $n = 512$. **Left:** Class 3 versus 7, as an example of (close-to) Gaussian information factors. **Right:** Class 4 versus 5, as an example of non-Gaussian information factors.

- ▶ Previous proofs of GU in ERM [MS22; Dan+24] required assumptions on the unknown statistical behaviour of $\hat{\beta}$, **not verifiable** from the data distribution.

- ▶ Previous proofs of GU in ERM [MS22; Dan+24] required assumptions on the unknown statistical behaviour of $\hat{\beta}$, **not verifiable** from the data distribution.
- ▶ Our results demonstrate a breakdown of GU in ERM **conditioned** on the **data structure**: under LFMM, GU **only holds** in the case of class-conditional **Gaussian informative factors**.

- ▶ Previous proofs of GU in ERM [MS22; Dan+24] required assumptions on the unknown statistical behaviour of $\hat{\beta}$, **not verifiable** from the data distribution.
- ▶ Our results demonstrate a breakdown of GU in ERM **conditioned** on the **data structure**: under LFMM, GU **only holds** in the case of class-conditional **Gaussian informative factors**.
- ▶ Our sharp performance analysis provides insight into the advantages of non-square losses in learning **higher order data statistics** under the GU breakdown.

- ▶ Previous proofs of GU in ERM [MS22; Dan+24] required assumptions on the unknown statistical behaviour of $\hat{\beta}$, **not verifiable** from the data distribution.
- ▶ Our results demonstrate a breakdown of GU in ERM **conditioned** on the **data structure**: under LFMM, GU **only holds** in the case of class-conditional **Gaussian informative factors**.
- ▶ Our sharp performance analysis provides insight into the advantages of non-square losses in learning **higher order data statistics** under the GU breakdown.
- ▶ Our results help predict GU on real-world data from the distributions of their **principal components**.

- ▶ Previous proofs of GU in ERM [MS22; Dan+24] required assumptions on the unknown statistical behaviour of $\hat{\beta}$, **not verifiable** from the data distribution.
- ▶ Our results demonstrate a breakdown of GU in ERM **conditioned** on the **data structure**: under LFMM, GU **only holds** in the case of class-conditional **Gaussian informative factors**.
- ▶ Our sharp performance analysis provides insight into the advantages of non-square losses in learning **higher order data statistics** under the GU breakdown.
- ▶ Our results help predict GU on real-world data from the distributions of their **principal components**.

Our Poster Session: Thu 24 Apr 10 a.m. SGT – 12:30 p.m. SGT.

- [MS22] Andrea Montanari and Basil N Saeed. “Universality of empirical risk minimization”. In: *Conference on Learning Theory*. PMLR. 2022, pp. 4310–4312.
- [Dan+24] Yatin Dandi, Ludovic Stephan, Florent Krzakala, Bruno Loureiro, and Lenka Zdeborová. “Universality laws for gaussian mixtures in generalized linear models”. In: *Advances in Neural Information Processing Systems* 36 (2024).
- [ML25] Xiaoyi Mai and Zhenyu Liao. “The Breakdown of Gaussian Universality in Classification of High-dimensional Mixtures”. In: *(accepted) 2025 International Conference on Learning Representations (ICLR)*. 2025.
- [TPT20] Hossein Taheri, Ramtin Pedarsani, and Christos Thrampoulidis. “Optimality of Least-squares for Classification in Gaussian-Mixture Models”. In: *2020 IEEE International Symposium on Information Theory (ISIT)* 00 (2020), pp. 2515–2520.
- [ML20] Xiaoyi Mai and Zhenyu Liao. *High Dimensional Classification via Regularized and Unregularized Empirical Risk Minimization: Precise Error and Optimal Loss*. Nov. 2020. arXiv: 1905.13742 [stat].
- [Bea+13] Derek Bean, Peter J Bickel, Noureddine El Karoui, and Bin Yu. “Optimal M-estimation in high-dimensional regression”. In: *Proceedings of the National Academy of Sciences* 110.36 (2013), pp. 14563–14568.