

Importance Sampling for Nonlinear Models

Fred Roosta

School of Mathematics and Physics
University of Queensland

Example: Linear Regression

In simple linear regression:

- $f_i(\boldsymbol{\theta}) = \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle - y_i$
- $\ell(t) = t^2$ (squared loss)

$$\mathcal{L}(\boldsymbol{\theta}) = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2$$

where

- $\mathbf{X} \in \mathbb{R}^{n \times p}$ has rows \mathbf{x}_i^T
- $\mathbf{y} \in \mathbb{R}^n$ has entries y_i

RandNLA: Successes

- Major advances in core matrix problems:
 - matrix multiplication
 - least-squares / least-absolute deviation
 - low-rank approximation
 - and many more ...
- Broad toolkit developed:
 - oblivious methods (sketching, projections),
 - non-oblivious methods (leverage, row-norm sampling).
- Provides strong **approximation guarantees**: for any small ε ,

$$\mathcal{L}(\boldsymbol{\theta}_S^*) \leq \mathcal{L}(\boldsymbol{\theta}^*) + \mathcal{O}(\varepsilon)$$

$\boldsymbol{\theta}^*$: optimal parameter from **full** dataset

$\boldsymbol{\theta}_S^*$: optimal parameter from **subsampled** dataset

Revisit: Linear Regression

Consider linear least squares:

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2 = \sum_{i=1}^n \overbrace{\frac{1}{2} (\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle - y_i)^2}^{\ell_i(\boldsymbol{\theta})}$$

- **Hessians:**

$$\nabla^2 \ell_i(\boldsymbol{\theta}) = \mathbf{x}_i \mathbf{x}_i^T, \quad \nabla^2 \mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \mathbf{X}^T \mathbf{X}$$

- Suppose $\mathbf{X} \in \mathbb{R}^{n \times p}$ has full column rank.

- i^{th} **leverage score:**

$$\tau_i = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i = \text{Trace}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^T)$$

Going Beyond Linear Regression

- Assume each $\ell_i(\boldsymbol{\theta})$ is convex
- Nonlinear “Leverage” score (?):

$$\tau_i(\boldsymbol{\theta}) \triangleq \text{Trace}\left([\nabla^2 \mathcal{L}(\boldsymbol{\theta})]^\dagger \nabla^2 \ell_i(\boldsymbol{\theta})\right)$$

- Intuition: Consider the Riemannian manifold \mathcal{M} equipped with metric

$$g_{\boldsymbol{\theta}}(\mathbf{v}, \mathbf{v}) = \mathbf{v}^T \nabla^2 \mathcal{L}(\boldsymbol{\theta}) \mathbf{v}, \quad \mathbf{v} \in T_{\boldsymbol{\theta}}(\mathcal{M}).$$

So, $\tau_i(\boldsymbol{\theta})$ measures how much removing individual components ℓ_i changes curvature (locally)

- Properties:

$$0 \leq \tau_i(\boldsymbol{\theta}) \leq \text{Rank}(\nabla^2 \ell_i(\boldsymbol{\theta})), \quad \text{and} \quad \sum_{i=1}^n \tau_i(\boldsymbol{\theta}) = \text{Rank}(\nabla^2 \mathcal{L}(\boldsymbol{\theta})).$$

- This did **not** lead us to guarantees of the kind:

$$\mathcal{L}(\boldsymbol{\theta}_S^*) \leq \mathcal{L}(\boldsymbol{\theta}^*) + \mathcal{O}(\varepsilon).$$

Deriving the Nonlinear Adjoint

- Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be continuously differentiable.
- Mean value theorem gives

$$f(\boldsymbol{\theta}) = f(\mathbf{0}) + \int_0^1 \left\langle \frac{\partial}{\partial \boldsymbol{\theta}} f(t\boldsymbol{\theta}), \boldsymbol{\theta} \right\rangle dt.$$

- Equivalently,

$$f(\boldsymbol{\theta}) = f(\mathbf{0}) + \left\langle \boldsymbol{\theta}, \int_0^1 \frac{\partial}{\partial \boldsymbol{\theta}} f(t\boldsymbol{\theta}) dt \right\rangle.$$

- This looks like a dual pairing.

Definition: Nonlinear Adjoint Operator

Definition (Adjoint Operator)

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$. The adjoint operator is

$$\mathbf{f}^*(\boldsymbol{\theta}) \triangleq \int_0^1 \frac{\partial}{\partial \boldsymbol{\theta}} f(t\boldsymbol{\theta}) dt,$$

whenever $t \mapsto f(t\boldsymbol{\theta})$ is absolutely continuous.

- Absolute continuity allows a.e. differentiability.
- Extends beyond smooth settings.

Inner-Product Representation

- With the adjoint operator,

$$f(\boldsymbol{\theta}) = \langle \hat{\boldsymbol{\theta}}, \hat{\mathbf{f}}^*(\boldsymbol{\theta}) \rangle,$$

where

$$\hat{\mathbf{f}}^*(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{f}^*(\boldsymbol{\theta}) \\ f(\mathbf{0}) \end{bmatrix}, \quad \hat{\boldsymbol{\theta}} = \begin{bmatrix} \boldsymbol{\theta} \\ 1 \end{bmatrix}.$$

- If $f(\mathbf{0}) = 0$, this becomes $f(\boldsymbol{\theta}) = \langle \boldsymbol{\theta}, \mathbf{f}^*(\boldsymbol{\theta}) \rangle$.
- A nonlinear analogue of Riesz representation².

²Continuous linear functionals can be identified with inner products against a unique representer in the Hilbert space.

Computing the Adjoint Operator

- In general, $\mathbf{f}^*(\boldsymbol{\theta})$ involves an integral.
- Numerical approximation is possible if needed.
- Many ML models yield closed forms.

Proposition

Let $f = g \circ h$, where

- $g : \mathbb{R} \rightarrow \mathbb{R}$,
- $h : \mathbb{R}^p \rightarrow \mathbb{R}$ is positively homogeneous of degree α .

Then

$$\mathbf{f}^*(\boldsymbol{\theta}) = \begin{cases} \left(\frac{g(h(\boldsymbol{\theta})) - g(0)}{\alpha h(\boldsymbol{\theta})} \right) \frac{\partial}{\partial \boldsymbol{\theta}} h(\boldsymbol{\theta}), & h(\boldsymbol{\theta}) \neq 0, \\ \left(\frac{g'(0)}{\alpha} \right) \frac{\partial}{\partial \boldsymbol{\theta}} h(\boldsymbol{\theta}), & h(\boldsymbol{\theta}) = 0. \end{cases}$$

Example: Generalized Linear Predictors

- Generalized linear predictor models (single-index models):

$$f(\boldsymbol{\theta}) = \phi(\langle \boldsymbol{\theta}, \mathbf{x} \rangle), \quad \mathbf{x} \in \mathbb{R}^p.$$

- Here $f = g \circ h$ with $g(t) = \phi(t)$ and $h(\boldsymbol{\theta}) = \langle \boldsymbol{\theta}, \mathbf{x} \rangle$.
- Homogeneity degree: $\alpha = 1$.

- Adjoint:

$$\mathbf{f}^*(\boldsymbol{\theta}) = \left(\frac{\phi(\langle \boldsymbol{\theta}, \mathbf{x} \rangle) - \phi(0)}{\langle \boldsymbol{\theta}, \mathbf{x} \rangle} \right) \mathbf{x}.$$

- Lifted adjoint:

$$\widehat{\mathbf{f}}^*(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{f}^*(\boldsymbol{\theta}) \\ \phi(0) \end{bmatrix}.$$

- Scales \mathbf{x} by the local nonlinearity of ϕ .

Example: Logistic Regression

- Logistic function:

$$\phi(t) = \frac{1}{1 + e^{-t}}.$$

- Adjoint:

$$\mathbf{f}^*(\boldsymbol{\theta}) = - \left(\frac{\tanh(\langle \boldsymbol{\theta}, \mathbf{x} \rangle / 2)}{2 \langle \boldsymbol{\theta}, \mathbf{x} \rangle} \right) \mathbf{x}.$$

- Smooth, bounded scaling vs. linear regression's constant adjoint.

Example: ReLU Neural Networks (Single Neuron)

- Let

$$r(\boldsymbol{\theta}) = \phi(\psi(\boldsymbol{\theta})), \quad \psi(\boldsymbol{\theta}) = a \cdot \max\{\langle \mathbf{b}, \mathbf{x} \rangle, 0\},$$

with $\boldsymbol{\theta} = [a, \mathbf{b}]$.

- ψ has homogeneity degree $\alpha = 2$.
- Adjoint:

$$\mathbf{r}^*(\boldsymbol{\theta}) = \frac{\phi(a \max\{\langle \mathbf{b}, \mathbf{x} \rangle, 0\}) - \phi(0)}{2a \max\{\langle \mathbf{b}, \mathbf{x} \rangle, 0\}} \begin{bmatrix} \max\{\langle \mathbf{b}, \mathbf{x} \rangle, 0\} \\ a \mathbf{x} \mathbb{1}(\langle \mathbf{b}, \mathbf{x} \rangle > 0) \end{bmatrix}.$$

- For $\phi(z) = z$:

$$\mathbf{r}^*(\boldsymbol{\theta}) = \frac{1}{2} \begin{bmatrix} \max\{\langle \mathbf{b}, \mathbf{x} \rangle, 0\} \\ a \mathbf{x} \mathbb{1}(\langle \mathbf{b}, \mathbf{x} \rangle > 0) \end{bmatrix}.$$

ReLU Networks and Additivity

- Two-layer ReLU Network (m hidden neurons and a single output):

$$f(\boldsymbol{\theta}) = \sum_{j=1}^m \phi(a_j \max\{\langle \mathbf{b}_j, \mathbf{x} \rangle, 0\}).$$

- Adjoint is additive:

$$\mathbf{f}^*(\boldsymbol{\theta}) = \sum_{j=1}^m \mathbf{e}_j \otimes \mathbf{r}^*(\theta_j).$$

- Lifted adjoint:

$$\widehat{\mathbf{f}}^*(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{f}^*(\boldsymbol{\theta}) \\ m \phi(0) \end{bmatrix}.$$

From Linear to Nonlinear Least Squares

Consider nonlinear least-squares:

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n (f_i(\boldsymbol{\theta}))^2,$$

Since $f_i(\boldsymbol{\theta}) = \langle \hat{\boldsymbol{\theta}}, \hat{\mathbf{f}}_i^*(\boldsymbol{\theta}) \rangle$, we can write

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n \langle \hat{\boldsymbol{\theta}}, \hat{\mathbf{f}}_i^*(\boldsymbol{\theta}) \rangle^2 = \|\hat{\mathbf{F}}^*(\boldsymbol{\theta}) \hat{\boldsymbol{\theta}}\|^2,$$

where

$$\hat{\mathbf{F}}^*(\boldsymbol{\theta}) \triangleq \begin{bmatrix} \hat{\mathbf{f}}_1^*(\boldsymbol{\theta}) & \hat{\mathbf{f}}_2^*(\boldsymbol{\theta}) & \cdots & \hat{\mathbf{f}}_n^*(\boldsymbol{\theta}) \end{bmatrix}^T.$$

is the **nonlinear dual matrix**.

This mirrors the linear form $\|\mathbf{X}\boldsymbol{\theta}\|^2$, with $\hat{\mathbf{F}}^*(\boldsymbol{\theta})$ acting as the data matrix.

Key point: Tools from RandNLA can now be used.

Nonlinear Importance Sampling

We can define importance scores from the rows of $\widehat{\mathbf{F}}^*(\boldsymbol{\theta})$.

- **Nonlinear Leverage Scores:**

$$\tau_i(\boldsymbol{\theta}) \triangleq \frac{\langle \mathbf{e}_i, \widehat{\mathbf{F}}^*(\boldsymbol{\theta}) [\widehat{\mathbf{F}}^*(\boldsymbol{\theta})]^\dagger \mathbf{e}_i \rangle}{\text{Rank}(\widehat{\mathbf{F}}^*(\boldsymbol{\theta}))}.$$

- **Nonlinear Row-norm Scores:**

$$\tau_i(\boldsymbol{\theta}) \triangleq \frac{\|\widehat{\mathbf{f}}_i^*(\boldsymbol{\theta})\|_2^2}{\|\widehat{\mathbf{F}}^*(\boldsymbol{\theta})\|_F^2}.$$

Both score types define valid distributions:

$$\tau_i(\boldsymbol{\theta}) \geq 0, \quad \sum_{i=1}^n \tau_i(\boldsymbol{\theta}) = 1.$$

Remark: For linear models $f(\boldsymbol{\theta}) = \langle \boldsymbol{\theta}, \mathbf{x} \rangle$, these scores reduce to the classical ones.

Sampling and Approximate Loss

Let \mathcal{S} be a multiset of s indices sampled with replacement using $\{\tau_i(\boldsymbol{\theta})\}$.

Define the sampled loss

$$\mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) \triangleq \sum_{i \in \mathcal{S}} \frac{(f_i(\boldsymbol{\theta}))^2}{s \tau_i(\boldsymbol{\theta})}.$$

This estimator is unbiased:

$$\mathbb{E} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta}).$$

Let $\widehat{\mathbf{F}}_{\mathcal{S}}^*(\boldsymbol{\theta})$ be the sampled/rescaled rows of $\widehat{\mathbf{F}}^*(\boldsymbol{\theta})$.

Then

$$\mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) = \|\widehat{\mathbf{F}}_{\mathcal{S}}^*(\boldsymbol{\theta}) \widehat{\boldsymbol{\theta}}\|^2.$$

Approximation Guarantee

For fixed θ , if

$$s = \tilde{O}\left(p/\varepsilon^2\right),$$

then with probability at least $1 - \delta$, for all $\mathbf{v} \in \mathbb{R}^{p+1}$,

$$(1 - \varepsilon) \|\hat{\mathbf{F}}^*(\theta)\mathbf{v}\|^2 \leq \|\hat{\mathbf{F}}_S^*(\theta)\mathbf{v}\|^2 \leq (1 + \varepsilon) \|\hat{\mathbf{F}}^*(\theta)\mathbf{v}\|^2$$

As a direct consequence, with probability at least $1 - \delta$,

$$(1 - \varepsilon) \mathcal{L}(\theta) \leq \mathcal{L}_S(\theta) \leq (1 + \varepsilon) \mathcal{L}(\theta)$$

What Is Missing?

$$\theta^* = \arg \min_{\theta} \mathcal{L}(\theta)$$

$$\theta_S^* = \arg \min_{\theta} \mathcal{L}_S(\theta)$$

$$\mathcal{L}_S(\theta_S^*) \leq \mathcal{L}_S(\theta^*) \leq (1 + \varepsilon)\mathcal{L}(\theta^*)$$

But this is *not yet* the desired guarantee. Our goal is

$$\mathcal{L}(\theta_S^*) \leq \mathcal{L}(\theta^*) + \mathcal{O}(\varepsilon).$$

Two fundamental obstacles remain:

- 1 Importance scores depend on θ^* , which is unknown.
- 2 The bound involves $\mathcal{L}_S(\theta_S^*)$, not $\mathcal{L}(\theta_S^*)$.

Challenge 1: Unknown Sampling Scores

The inequality

$$\mathcal{L}_S(\boldsymbol{\theta}_S^*) \leq (1 + \varepsilon)\mathcal{L}(\boldsymbol{\theta}^*),$$

requires evaluating

$$\tau_i(\boldsymbol{\theta}^*),$$

which is infeasible since $\boldsymbol{\theta}^*$ is unknown.

Strategy: approximate nonlinear scores by *parameter-independent* quantities.

Near-optimal Sampling in RandNLA

From RandNLA: if

$$\beta \tau_i \leq \hat{\tau}_i,$$

for some $\beta \in (0, 1)$, then subspace embedding holds with sample size

$$\tilde{O}\left(p/(\beta\varepsilon^2)\right).$$

If we can approximate *nonlinear* scores similarly, guarantees carry over.

For specific model classes, the structure of the nonlinear adjoint operator allows parameter-independent upper bounds on nonlinear scores.

Example: Generalized Linear Predictors

Generalized linear predictor models $f(\boldsymbol{\theta}) = \phi(\langle \boldsymbol{\theta}, \mathbf{x} \rangle)$ with ϕ satisfying

$$\ell \leq \frac{\phi^2(t)}{t^2} \leq u \quad \text{for } t \in \mathcal{T}.$$

Then, for full-rank \mathbf{X} and $\mathbf{F}^*(\boldsymbol{\theta})$,

$$\ell \mathbf{X}^T \mathbf{X} \preceq [\mathbf{F}^*(\boldsymbol{\theta})]^T \mathbf{F}^*(\boldsymbol{\theta}) \preceq u \mathbf{X}^T \mathbf{X}.$$

As a consequence,

$$\tau_i(\boldsymbol{\theta}) \leq \frac{u}{\ell} \tau_i,$$

where τ_i is the *linear* leverage score from $\hat{\mathbf{X}}$.

Example: ReLU Neural Networks

Now consider a one-hidden-layer ReLU network with row-norm sampling.

Assume the activation satisfies

$$c_1 \leq \frac{(\phi(t) - \phi(0))^2}{t^2} \leq c_2 \quad \text{for } t \in \mathcal{T}.$$

Under mild assumptions on θ^* , define a compact set \mathcal{C} containing it.

For all $\theta \in \mathcal{C}$,

$$c_1 \ell \|\mathbf{x}_i\|^2 \leq \|\mathbf{f}_i^*(\theta)\|^2 \leq c_2 u \|\mathbf{x}_i\|^2,$$

where ℓ and u are used in the definition of \mathcal{C} .

Hence,

$$\tau_i(\theta) \leq \frac{\max\{c_2 u, 1\}}{\min\{c_1 \ell, 1\}} \tau_i,$$

where τ_i is a row-norm score from $\hat{\mathbf{X}}$.

Challenge 2: Relating Sampled and Full Loss

To complete the argument, we need a lower bound of the form

$$\mathcal{L}(\theta_S^*) \lesssim \mathcal{L}_S(\theta_S^*)$$

This requires uniform control over θ .

Consider a ball \mathcal{B}_R^* containing θ^* . Construct an ε -net $\mathcal{N}_\varepsilon \subset \mathcal{B}_R^*$ such that:

$$\forall \theta \in \mathcal{B}_R^*, \quad \exists \theta' \in \mathcal{N}_\varepsilon \text{ with } \|\theta - \theta'\| \leq \varepsilon R.$$

The net size satisfies $|\mathcal{N}_\varepsilon| = \Omega(1/\varepsilon^p)$. By choosing

$$\delta' = \delta/|\mathcal{N}_\varepsilon|,$$

a union bound ensures that

$$(1 - \varepsilon)\mathcal{L}(\theta) \leq \mathcal{L}_S(\theta) \leq (1 + \varepsilon)\mathcal{L}(\theta)$$

holds *simultaneously* for all $\theta \in \mathcal{N}_\varepsilon$ with probability $1 - \delta$.

Challenge 2: Relating Sampled and Full Loss

Assume \mathcal{L}_S is continuous on \mathcal{B}_R^* . By compactness, \mathcal{L}_S is Lipschitz:

$$|\mathcal{L}_S(\boldsymbol{\theta}) - \mathcal{L}_S(\boldsymbol{\theta}')| \leq L(f, \mathbf{X}, R) \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|.$$

This allows us to move from the net to arbitrary $\boldsymbol{\theta}$.

Using Lipschitz continuity and the net approximation,

$$\mathcal{L}_S(\boldsymbol{\theta}_S^*) \geq (1 - \varepsilon)\mathcal{L}(\boldsymbol{\theta}_S^*) - \varepsilon(2 - \varepsilon)R L(f, \mathbf{X}, R).$$

Putting it all together, we obtain

$$\mathcal{L}(\boldsymbol{\theta}_S^*) \leq \mathcal{L}(\boldsymbol{\theta}^*) + \frac{\varepsilon}{1 - \varepsilon} \left(\mathcal{L}(\boldsymbol{\theta}^*) + (2 - \varepsilon)R L(f, \mathbf{X}, R) \right).$$

Main Result

Theorem (Main Result)

Suppose $\theta^* \in \mathcal{C}$ and nonlinear scores satisfy

$$\beta \tau_i(\theta) \leq \tau_i \quad \forall \theta \in \mathcal{C}.$$

If

$$s = \tilde{\mathcal{O}}\left(p^2/(\beta\varepsilon^2)\right),$$

then, with probability at least $1 - \delta$,

$$\mathcal{L}(\theta_S^*) \leq \mathcal{L}(\theta^*) + \mathcal{O}(\varepsilon).$$

Discussion

- Recent work by Gajjar, Musco, and Hegde (2023) and Gajjar, Tai, Xingyu, Hegde, Musco, and Li (2024) studies leverage-score-based sampling for single-index models $f(\boldsymbol{\theta}) = \phi(\langle \boldsymbol{\theta}, \mathbf{x} \rangle)$.
- If ϕ is L -Lipschitz, Gajjar, Musco, and Hegde (2023) show that

$$\tilde{\mathcal{O}}(p^2/\varepsilon^4)$$

samples suffice to guarantee

$$\mathcal{L}(\boldsymbol{\theta}_S^*) \leq C \mathcal{L}(\boldsymbol{\theta}^*) + \mathcal{O}(\varepsilon),$$

for some constant $C > 1$.

- Gajjar, Tai, Xingyu, Hegde, Musco, and Li (2024) improve the sample complexity to

$$\tilde{\mathcal{O}}(p/\varepsilon^2),$$

matching linear rates up to logarithmic factors.

- However:** The constant C can be large (exceeding 10^3).

Discussion

	This work	Gajjar et al. (2024)
Applicability	Beyond single-index models	Single-index models
ε dependence	Same	Same
p dependence	$\mathcal{O}(p^2)$	$\mathcal{O}(p)$
Subproblem	Constrained	Unconstrained
Guarantee constant C	$C = 1$	Potentially large ($C \gg 1$)

SVHN - High Leverage (1 vs 0)



SVHN - Low Leverage (1 vs 0)



SVHN - High Leverage (1 vs 7)



SVHN - Low Leverage (1 vs 7)



NOTMNIST - High Leverage (A vs B)



Figure: Ambiguous or atypical character.

NOTMNIST - Low Leverage (A vs B)



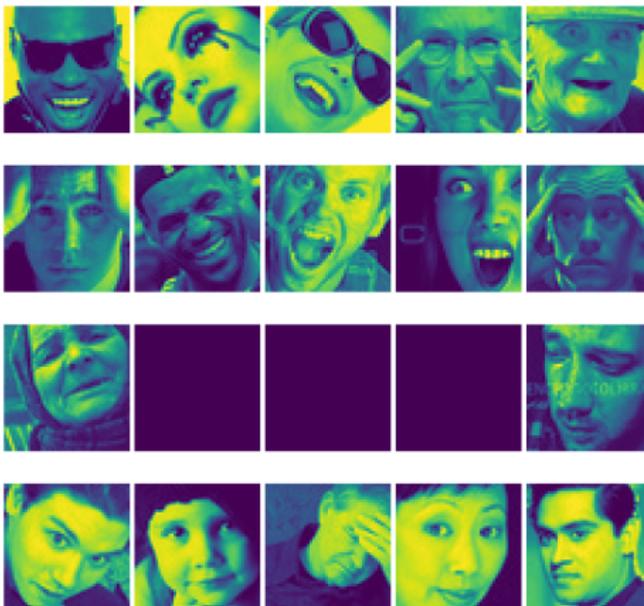
Figure: Canonical, easy-to-classify sample.

NOTMNIST - Low Leverage (B vs D)

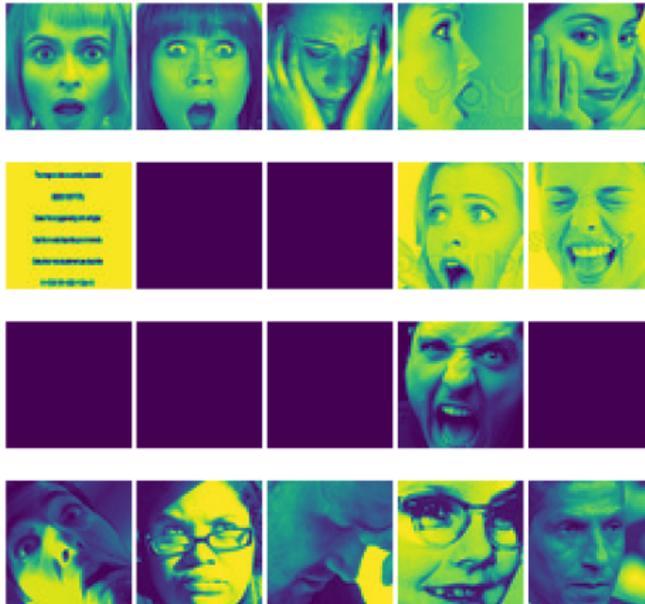


Figure: Low-information, prototypical sample.

FER - High Nonlinear Leverage



FER - Low Nonlinear Leverage





-  Avron, Haim, Michael Kapralov, Cameron Musco, Christopher Musco, Ameya Velingker, and Amir Zandieh (2017). “Random Fourier features for kernel ridge regression: Approximation bounds and statistical guarantees”. In: *International conference on machine learning*. PMLR, pp. 253–262.
-  — (2019). “A universal sampling method for reconstructing signals with simple fourier transforms”. In: *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pp. 1051–1063.
-  Erdélyi, Tamás, Cameron Musco, and Christopher Musco (2020). “Fourier sparse leverage scores and approximate kernel learning”. In: *Advances in Neural Information Processing Systems 33*, pp. 109–122.
-  Gajjar, Aarshvi and Cameron Musco (2021). “Subspace embeddings under nonlinear transformations”. In: *Algorithmic Learning Theory*. PMLR, pp. 656–672.
-  Gajjar, Aarshvi, Christopher Musco, and Chinmay Hegde (2023). “Active learning for single neuron models with lipschitz non-linearities”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR, pp. 4101–4113.

