# Unsupervised dataset distillation

## Project

- General task: train a (shallow) feature extractor via PSRs (on a learned dataset) that work well with a task-specific head

$$\min_{\mathcal{D}} \sum_j \min_{w_{\text{head}}^j} \mathcal{L}_{\text{tr}}^j(w_{\text{head}}^j \circ w_{\text{base}}^*) \qquad \text{s.t. } w_{\text{base}}^* = \arg\min_w \mathcal{L}_{\text{PSR}}(w; \mathcal{D})$$

- Why shallow feature extractor?
  Inner problem should be "easy" and quick to solve
  Head part of NN can be arbitrary
- E.g. use shallow LVM such as GRBM

$$\log p_\theta(x) \doteq -\frac{1}{2\sigma^2}\|x - a\|^2 + \sum_k S(w_k^\top x + b_k)$$

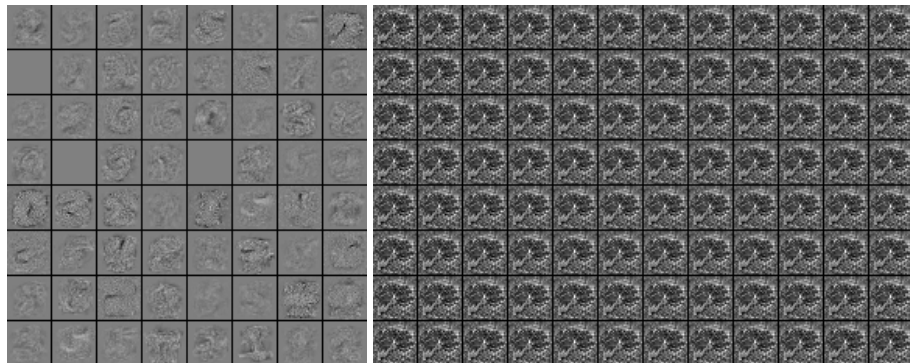$$\nabla_x \log p_\theta(x) = \frac{1}{\sigma^2}\|a - x\|^2 + \sum_k w_k \sigma(w_k^\top x + b_k)$$

Feature extractor $S'(w_k^\top x + b_k) = \sigma(w_k^\top x + b_k)$

- Dataset $\mathcal{D} = \{x_1, \ldots, x_N\}$ of unlabeled data

Quadratic penalty method

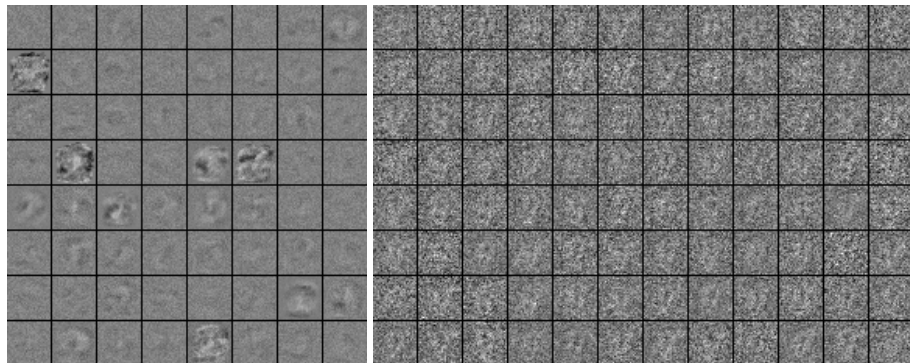| Epoch | Test accuracy | $\|X - AE(x)\|$ |
|---|---|---|
| 1 | 46.1% | 5.4 |
| 100 | 92.1% | 0.067 |
| 200 | 94.1% | 0.01 |



$w_{\text{base}}$ and distilled dataset $\mathcal{D}$

# Unsupervised dataset distillation: Reference implementation

Quadratic penalty method (now with "diversity term" $-\lambda_{\text{divers}}\|\mathcal{D} - \text{mean}(\mathcal{D})\|^2$)

| Epoch | Test accuracy | $\|X - AE(x)\|$ |
|---|---|---|
| 1 | 24.0% | 5.6 |
| 100 | 92.3% | 0.22 |
| 200 | 93.2% | 0.20 |



$w_{\text{base}}$ and distilled dataset $\mathcal{D}$

Optimal value reformulation (diversity term has little impact)

| Epoch | Test accuracy | $\|X - AE(x)\|$ |
|-------|---------------|-----------------|
| 1     | 32.6%         | 4.8             |
| 100   | 81.4%         | 2.6             |
| 200   | 94.6%         | 3.3             |



$w_{\text{base}}$ and distilled dataset $\mathcal{D}$

# Unsupervised dataset distillation: What you should do

- Choose meta-learning approach
  - Quadratic penalty method (probably the easiest to make it work)
  - Optimal value reformulation
  - Implicit differentiation & approximate Hessian matrix inversion
  - Or any other method that approximately solves bilevel programs
- Baseline code provided
  - Available on Canvas
  - Implements penalty methods using AE as unsupervised learner
  - Implemented in Julia + Flux
- Extend what is implemented in the baseline, e.g.
  - Different NN / DEM architectures, different target dataset and/or different $|\mathcal{D}|$
  - Multiple training datasets for outer loss (MTL setting)
  - Use different PSRs / unsupervised learning methods (instead of simple AE)
  - Investigate other ways to prevent mode collapse and/or to make $\mathcal{D}$ look more natural
    - E.g. pre-train a prior $\log p_\psi(x)$ to represent a desired image statistics and add to loss
    - Note: there is an ambiguity between $w_{\text{base}}$ and $\mathcal{D}$: $w_{\text{base}} \mathcal{D} = w_{\text{base}} Q Q^{-1} \mathcal{D}$
  - Apply the approach layer-wise / meta-learn multiple layers at once
  - You only need to work on some of these (and your own) suggestions, not all!
    - Your time spent should be worth 2 ECTS ($\approx$ 40 hours, i.e. 56 hours minus lecture time)

# Unsupervised dataset distillation: Hints

- If you opt for the quadratic penalty method, you need to calculate $\nabla_w \mathcal{L}_{\text{inner}}(w, \theta)$
  - Because most auto-diff packages cannot do $\nabla_w f(\nabla_x g(x, w))$ well (but JAX might be able to do it)
- In the case of AE and GBRM:

$$\nabla_{W,b,a} \left( \frac{1}{2} \left\| \frac{1}{\sigma^2}(a - x) + W\sigma(W^\top x + b) \right\|^2 \right)$$

- Define $\mathbf{r} := \frac{1}{\sigma^2}(a - x) + W\sigma(W^\top x + b)$
- Main task: $\nabla_{w_k}$, where $w_k$ is $k$-th column of $W$

$$\nabla_{w_k} \left( \frac{1}{2} \left\| \frac{1}{\sigma^2}(a - x) + \sum w_k \sigma(w_k^\top x + b_k) \right\|^2 \right)$$

- Hint: $\nabla$ yields column vector, hence chain rule is transposed

$$\nabla_{w_k} = \left( \sigma(w_k^\top x + b_k)\mathsf{I} + \sigma'(w_k^\top x + b_k)w_k x^\top \right) \mathbf{r}$$

- Efficient implementation in Julia:

```julia
as = m.W' * Xs .+ m.b
hiddens = my_σ.(as); dhiddens = ∇my_σ.(as)
XXs = m.W * hiddens .+ m.a / σ2; residual = (XXs - Xs / σ2) / N
W_r = m.W' * residual

gs_W = residual * hiddens' + Xs * (W_r .* dhiddens)'
gs_b = sum(W_r .* dhiddens, dims=2)
gs_a = sum(residual, dims=2) / σ2

return sum(abs2.(gs_W)) + sum(abs2.(gs_b)) + sum(abs2.(gs_a))
```

# Unsupervised dataset distillation: Hints

- Consult the "Matrix Cookbook" for vector and matrix-based calculus
- Penalty: verify numerically that $\|\nabla_w \mathcal{L}_{\text{inner}}(w, \theta)\|^2$ matches with and w/o auto-diff
  - Math: Jacobian $\sigma'(w_k^\top x + b_k)$ is a diagonal matrix (tensor for a batch)
  - Implementation: `dhiddens` $= \nabla\text{my}\_\sigma.\text{(W'} \ast \text{x + b)}$ is a vector (matrix for a batch)
  - Efficient summation over entire batch
- OVR: make sure $\min_{w'} \mathcal{L}_{\text{inner}}(w', \theta)$ works well!

## Don't panic

- Solving bilevel programs is difficult & brittle!
  - Especially for non-convex inner losses
  - Results may vary significantly with different hyper-parameters
- No need to win a benchmark!

# Unsupervised dataset distillation: What to submit

- Short report (pdf!) describing what you have done beyond the baseline code
  - And some visualization of the distilled dataset and the training progress
  - Also document unexpected problems/surprises you encountered
- Your Python / Julia / Matlab / . . . code
- Email to zach@chalmers.se by November 30, 2024
  - With subject containing "Project LFR-2024"
- Contact me with a proposal for a project if you cannot do this one
  - E.g. no access to sufficiently capable PC or no coding experience

  You can work in pairs (and submit one pdf and implementation)!

# Examples of last year's project reports

# Examples of last year's project reports