Principal Bit Analysis: Autoencoding with Schur-Concave Loss Sourbh Bhadane, Aaron B. Wagner, Jayadev Acharya,

Autoencoders and their Limitations

▶ Given a centered dataset $x_1, x_2 \cdots x_n \in \mathbb{R}^d$ with an empirical covariance matrix K, an **autoencoder** consists of an encoder $f : \mathbb{R}^d \mapsto \mathbb{R}^k$ and a decoder $g: \mathbb{R}^k \mapsto \mathbb{R}^d$. In a linear autoencoder (LAE), f and g are linear maps. Conventional LAEs do not identify principal directions of the dataset. ► Autoencoders are sometimes described as "compressing" the data.

LAE with Schur-Concave constraint

Linear Encoder	$W^{ op}x_i$	Quantizer	$W^{\top}x_i + \varepsilon$	Linear Decoder
(W)				(7)

Figure: Compression Block Diagram

► Quantization assumptions:

 \triangleright Dither: If Q() maps a real number to its nearest integer, then $Q(z + \varepsilon) - \varepsilon \stackrel{d}{\sim} z + \varepsilon$ for $\varepsilon \sim \text{Unif}[-0.5, 0.5]$

> To constrain number of bits, clip $w_i^{\top}x$ to the interval

▷ Number of bits to represent $w_i^\top x$ is $\frac{1}{2} \log (4a^2 w_i^\top K w_j + 1)$.

• **Goal:** Minimize mean squared error (MSE) when $W^{\top}x$ is quantized, subject to a rate constraint on the number of bits required to represent quantized $W^{\top}x$. ► Observe $\{w_j^\top K w_j\}_{j=1}^d \mapsto \frac{1}{2} \log (4a^2 w_j^\top K w_j + 1)$ is **Schur-concave**.

Solve a more general nonconvex problem with any Schur-concave constraint ρ .

$$\inf_{W,T} \quad \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}_{\varepsilon} \left[\left\| x_{i} - T \left(W^{\top} x_{i} + \varepsilon \right) \right\|_{2}^{2} \right]$$

ubject to $R \geq \rho \left(\left\{ w_{j}^{\top} K w_{j} \right\}_{j=1}^{d} \right).$

Theorem (Optimal LAE with Schur-concave constraint)

For Schur-concave $\rho : \mathbb{R}^d_{>0} \to \mathbb{R}_{\geq 0}$ and R > 0, the set of matrices whose nonzero columns are eigenvectors of the covariance matrix K is optimal. If ρ is strictly Schur-concave and K contains distinct eigenvalues, this set contains all optimal solutions.

Remarks

 \blacktriangleright W = US, where U is eigenvector matrix and S is an unknown diagonal matrix. ▶ Let $\rho(\mathbf{x}) = \sum_{i=1}^{n} \rho_{si}(\mathbf{x}_i)$, where $\rho_{si} : \mathbb{R}_{\geq 0} \mapsto \mathbb{R}_{\geq 0}$.

Cornell University





 Kw_j+1

(1)

Conventional LAE

- ► Conventional LAE: $W, T \in \mathbb{R}^{d \times k}$ where k is a parameter. > Principal component analysis (PCA): $W = T = U_k$, where U_k is a matrix with top
- k eigenvectors as columns. PCA is a global optimal solution, but not unique. We recover conventional LAEs by penalizing dimension,

 $\rho_{sl}(x) = 1 [x > 0].$ The optimal encoder is given by

> ∞ ∞ S =(2) 0 11 A 44

where top min (|R|, d) entries are ∞ . Since latent variables are quantized in our formulation, PCA with parameter k is equivalent to S with top k diagonal entries equal to ∞ .

Principal Bit Analysis (PBA)

Solve original problem by choosing

 $\rho_{sl}(\mathbf{x}) = \frac{1}{2} \log \left(\frac{\gamma}{\sigma^2} \mathbf{x} + 1 \right)$

 $\triangleright \gamma = 1$: classical waterfilling solution

 $\succ \gamma \in (1, 2]$: convex optimization problem

 $\triangleright \gamma > 2$: nonconvex optimization problem

PBA takes as input $\lambda > 0$ and outputs the optimal rate-distortion Lagrangian solution. By sweeping λ , we obtain the rate-distortion curve.

Variable-Rate Compression

Ballé et al. [1] defined an autoencoder-based variable-length compressor objective called nonlinear transform coding (NTC),

 $\inf_{f,g} \mathbb{E}_{x,\varepsilon} \left[\|x - g(Q(f(x) + \varepsilon) - \varepsilon)\|_2^2 \right] + \lambda$

Assuming a Gaussian source, linear f and g, and under independent encoding of each dimension, the NTC Lagrangian is

 $\inf_{W,T} \mathbb{E}_{\boldsymbol{x},\varepsilon} \left[\left\| \boldsymbol{x} - T \left(\boldsymbol{W}^{\top} + \varepsilon \right) \right\|_{2}^{2} \right] + \lambda \cdot \right]$

Theorem: Under the above assumptions, any W that achieves the infimum has all its nonzero columns as eigenvectors of K.

$$H(Q(f(x) + \varepsilon) - \varepsilon | \varepsilon).$$
 (3)

$$\sum_{i=1}^{d} h\left(w_i^{\top} x + [\varepsilon]_i\right).$$
 (4)

Experiments

- We compare across three metrics. 1) SNR = $10 \cdot \log_{10} (P/MSE)$, 2) Structural similarity index (SSIM/MS-SSIM) 3) Performance on downstream tasks, specifically classification accuracy.



SSIM Performance



Classification Accuracy Performance





We compare a PBA-based fixed-rate compressor with PCA. For images, we compare against JPEG, JPEG2000. For audio, we compare against AAC.

