

# Autoencoder - Derivations & Proofs

Paul F. Roysdon, Ph.D.

## Contents

<b>1 Mathematical Derivations &amp; Proofs</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 Data and Notation . . . . .	1
1.3 Model Formulation and Objective . . . . .	2
1.4 Backpropagation Derivation (One-Hidden-Layer AE) . . . . .	2
1.5 Linear Autoencoder and PCA (Equivalence Proof) . . . . .	3
1.6 Regularized Variants . . . . .	3
1.7 Algorithm (Autoencoder Training) . . . . .	4
1.8 Summary of Variables and Their Dimensions . . . . .	4
1.9 Summary . . . . .	5

## 1 Mathematical Derivations & Proofs

### 1.1 Introduction

An autoencoder (AE) is a neural architecture that learns an *encoder* mapping inputs to a low-dimensional representation (latent code) and a *decoder* that reconstructs the inputs from that code. Training minimizes a reconstruction loss, optionally augmented with regularizers that shape the representation (sparsity, contraction, denoising). We derive: (i) the AE objective from empirical risk minimization, (ii) complete backpropagation gradients for a one-hidden-layer AE (generalizes to deep AEs by recursion), and (iii) the connection between the *linear* AE and PCA, giving a proof (via projection and Eckart–Young) that the optimal linear AE spans the principal subspace. We also state variants (denoising, sparse, contractive) with their objectives and gradients.

### 1.2 Data and Notation

Given training data

$$\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n, \quad \mathbf{x}_i \in \mathbb{R}^d,$$

we learn encoder/decoder functions  $f_{\phi} : \mathbb{R}^d \rightarrow \mathbb{R}^r$  and  $g_{\theta} : \mathbb{R}^r \rightarrow \mathbb{R}^d$  with  $r \ll d$  (bottleneck). For the canonical one-hidden-layer AE.

The Encoder is defined

$$\mathbf{h} = f_{\phi}(\mathbf{x}) = \sigma_e(\mathbf{W}_e \mathbf{x} + \mathbf{b}_e) \in \mathbb{R}^r,$$

the Decoder

$$\hat{\mathbf{x}} = g_{\theta}(\mathbf{h}) = \sigma_d(\mathbf{W}_d \mathbf{h} + \mathbf{b}_d) \in \mathbb{R}^d,$$

with parameters:

$$\phi = \{\mathbf{W}_e \in \mathbb{R}^{r \times d}, \mathbf{b}_e \in \mathbb{R}^r\}, \quad \theta = \{\mathbf{W}_d \in \mathbb{R}^{d \times r}, \mathbf{b}_d \in \mathbb{R}^d\},$$

where the weights are  $\mathbf{W}_e, \mathbf{W}_d$  and biases  $\mathbf{b}_e, \mathbf{b}_d$ .

For a mini-batch  $\mathcal{B}$  we stack columns in matrices  $\mathbf{X} \in \mathbb{R}^{d \times |\mathcal{B}|}$ ,  $\mathbf{H} \in \mathbb{R}^{r \times |\mathcal{B}|}$ ,  $\hat{\mathbf{X}} \in \mathbb{R}^{d \times |\mathcal{B}|}$ .

### 1.3 Model Formulation and Objective

Let  $\ell(\mathbf{x}, \hat{\mathbf{x}})$  measure reconstruction error. The empirical risk (with optional weight decay) is

$$\min_{\phi, \theta} \hat{\mathcal{L}}(\phi, \theta) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, g_\theta(f_\phi(\mathbf{x}_i))) + \frac{\lambda}{2} (\|\mathbf{W}_e\|_F^2 + \|\mathbf{W}_d\|_F^2). \quad (1)$$

Common choices:

$$\ell_{\text{MSE}}(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{2} \|\hat{\mathbf{x}} - \mathbf{x}\|_2^2, \quad \ell_{\text{CE}}(\mathbf{x}, \hat{\mathbf{x}}) = - \sum_{j=1}^d \left( x_j \log \hat{x}_j + (1 - x_j) \log (1 - \hat{x}_j) \right)$$

(MSE for real-valued inputs with linear decoder; cross-entropy for binary inputs with sigmoid decoder).

### 1.4 Backpropagation Derivation (One-Hidden-Layer AE)

For a single example, define pre-activations

$$\mathbf{a}^{(e)} = \mathbf{W}_e \mathbf{x} + \mathbf{b}_e, \quad \mathbf{h} = \sigma_e(\mathbf{a}^{(e)}), \quad \mathbf{a}^{(d)} = \mathbf{W}_d \mathbf{h} + \mathbf{b}_d, \quad \hat{\mathbf{x}} = \sigma_d(\mathbf{a}^{(d)}).$$

Gradients follow by chain rule. Let  $\odot$  denote Hadamard product and  $\sigma'_e, \sigma'_d$  the elementwise derivatives.

**Decoder layer.** For general  $\ell$ ,

$$\frac{\partial \ell}{\partial \mathbf{a}^{(d)}} = (\nabla_{\hat{\mathbf{x}}} \ell) \odot \sigma'_d(\mathbf{a}^{(d)}).$$

For MSE:  $\nabla_{\hat{\mathbf{x}}} \ell = \hat{\mathbf{x}} - \mathbf{x}$ , thus

$$\boldsymbol{\delta}^{(d)} \triangleq \frac{\partial \ell}{\partial \mathbf{a}^{(d)}} = (\hat{\mathbf{x}} - \mathbf{x}) \odot \sigma'_d(\mathbf{a}^{(d)}). \quad (2)$$

Gradients:

$$\nabla_{\mathbf{W}_d} \ell = \boldsymbol{\delta}^{(d)} \mathbf{h}^\top + \lambda \mathbf{W}_d, \quad \nabla_{\mathbf{b}_d} \ell = \boldsymbol{\delta}^{(d)}. \quad (3)$$

**Encoder layer.** Backpropagate through decoder:

$$\frac{\partial \ell}{\partial \mathbf{h}} = \mathbf{W}_d^\top \boldsymbol{\delta}^{(d)}, \quad \boldsymbol{\delta}^{(e)} \triangleq \frac{\partial \ell}{\partial \mathbf{a}^{(e)}} = (\mathbf{W}_d^\top \boldsymbol{\delta}^{(d)}) \odot \sigma'_e(\mathbf{a}^{(e)}).$$

Gradients:

$$\nabla_{\mathbf{W}_e} \ell = \boldsymbol{\delta}^{(e)} \mathbf{x}^\top + \lambda \mathbf{W}_e, \quad \nabla_{\mathbf{b}_e} \ell = \boldsymbol{\delta}^{(e)}. \quad (4)$$

**Proof (chain rule).** *Proof.* For MSE,  $\ell = \frac{1}{2} \|\hat{\mathbf{x}} - \mathbf{x}\|^2$  so  $d\ell = (\hat{\mathbf{x}} - \mathbf{x})^\top d\hat{\mathbf{x}}$ . With  $\hat{\mathbf{x}} = \sigma_d(\mathbf{a}^{(d)})$  and  $\mathbf{a}^{(d)} = \mathbf{W}_d \mathbf{h} + \mathbf{b}_d$ ,

$$d\ell = \underbrace{((\hat{\mathbf{x}} - \mathbf{x}) \odot \sigma'_d(\mathbf{a}^{(d)}))^\top}_{\boldsymbol{\delta}^{(d)}} (d\mathbf{W}_d \mathbf{h} + \mathbf{W}_d d\mathbf{h} + d\mathbf{b}_d).$$

Collect coefficients of  $d\mathbf{W}_d$  and  $d\mathbf{b}_d$  to obtain Eqn. (3). Then propagate  $d\mathbf{h}$  back through  $\mathbf{h} = \sigma_e(\mathbf{a}^{(e)})$  with  $\mathbf{a}^{(e)} = \mathbf{W}_e \mathbf{x} + \mathbf{b}_e$  to obtain Eqn. (4).  $\blacksquare$

**Mini-batch form.** Stacking columns for batch  $\mathcal{B}$ ,

$$\begin{aligned} \boldsymbol{\Delta}^{(d)} &= (\hat{\mathbf{X}} - \mathbf{X}) \odot \sigma'_d(\mathbf{A}^{(d)}), \quad \nabla_{\mathbf{W}_d} \hat{\mathcal{L}}_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \boldsymbol{\Delta}^{(d)} \mathbf{H}^\top + \lambda \mathbf{W}_d, \\ \boldsymbol{\Delta}^{(e)} &= (\mathbf{W}_d^\top \boldsymbol{\Delta}^{(d)}) \odot \sigma'_e(\mathbf{A}^{(e)}), \quad \nabla_{\mathbf{W}_e} \hat{\mathcal{L}}_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \boldsymbol{\Delta}^{(e)} \mathbf{X}^\top + \lambda \mathbf{W}_e. \end{aligned}$$

**Tied weights (optional).** If we enforce  $\mathbf{W}_d = \mathbf{W}_e^\top$  (parameter sharing), then the total derivative w.r.t.  $\mathbf{W}_e$  includes both paths:

$$\nabla_{\mathbf{W}_e} \ell \Big|_{\text{tied}} = \boldsymbol{\delta}^{(e)} \mathbf{x}^\top + \left( \nabla_{\mathbf{W}_d} \ell \right)^\top + \lambda \mathbf{W}_e.$$

## 1.5 Linear Autoencoder and PCA (Equivalence Proof)

Assume centered data ( $\frac{1}{n} \sum_i \mathbf{x}_i = \mathbf{0}$ ), linear activations ( $\sigma_e, \sigma_d$  identity), no biases, and tied weights with orthonormal rows:

$$\mathbf{h} = \mathbf{W}\mathbf{x}, \quad \hat{\mathbf{x}} = \mathbf{W}^\top \mathbf{h} = \mathbf{W}^\top \mathbf{W}\mathbf{x}, \quad \mathbf{W} \in \mathbb{R}^{r \times d}, \quad \mathbf{W}\mathbf{W}^\top = \mathbf{I}_r.$$

The reconstruction error is

$$\frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \frac{1}{2} \text{tr}(\mathbf{S} - 2\mathbf{SP} + \mathbf{P}\mathbf{SP}), \quad \mathbf{S} = \frac{1}{n} \sum_i \mathbf{x}_i \mathbf{x}_i^\top, \quad \mathbf{P} = \mathbf{W}^\top \mathbf{W}. \quad (5)$$

Since  $\mathbf{P}$  is an orthogonal projector onto the row-space of  $\mathbf{W}$ , it satisfies  $\mathbf{P} = \mathbf{P}^\top = \mathbf{P}^2$ . Then  $\mathbf{P}\mathbf{SP} = \mathbf{PS}$ , and Eqn. (5) reduces to

$$\frac{1}{2} \text{tr}(\mathbf{S}) - \text{tr}(\mathbf{SP}).$$

The first term is constant; minimizing reconstruction error is equivalent to maximizing  $\text{tr}(\mathbf{SP})$  over rank- $r$  projectors  $\mathbf{P}$ .

**Eckart–Young/Principal subspace argument.** *Proof.* Let  $\mathbf{S} = \mathbf{U}\Lambda\mathbf{U}^\top$  with eigenvalues  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ . Any rank- $r$  projector can be written  $\mathbf{P} = \mathbf{U} \begin{bmatrix} \mathbf{I}_r & \\ & \mathbf{0} \end{bmatrix} \mathbf{U}^\top$  after choosing an  $r$ -dimensional subspace.

Then

$$\text{tr}(\mathbf{SP}) = \sum_{j=1}^r \lambda_{i_j} \leq \sum_{j=1}^r \lambda_j,$$

with equality when the projector spans the top- $r$  eigenvectors. Therefore the linear AE with orthonormal-tied weights learns the PCA principal subspace, and the reconstruction  $\hat{\mathbf{x}} = \mathbf{Px}$  is the orthogonal projection onto that subspace. ■

**Untied weights, no orthonormal constraint.** Even without explicit orthonormality, at the optimum (with linear decoder/encoder and centered data), the column space of  $\mathbf{W}_d$  equals the row space of  $\mathbf{W}_e$  and spans the top- $r$  eigenvectors; the factorization  $\mathbf{W}_d \mathbf{W}_e$  is the best rank- $r$  approximation to the identity in the metric induced by  $\mathbf{S}$ , again by Eckart–Young for low-rank approximation of the data matrix.

## 1.6 Regularized Variants

**Denoising Autoencoder (DAE).** Corrupt inputs by  $\tilde{\mathbf{x}} \sim q(\tilde{\mathbf{x}} | \mathbf{x})$  (e.g., Gaussian noise or masking) and train to reconstruct the clean  $\mathbf{x}$ :

$$\min_{\phi, \theta} \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tilde{\mathbf{x}}_i \sim q(\cdot | \mathbf{x}_i)} \ell(\mathbf{x}_i, g_\theta(f_\phi(\tilde{\mathbf{x}}_i))) + \frac{\lambda}{2} (\|\mathbf{W}_e\|_F^2 + \|\mathbf{W}_d\|_F^2). \quad (6)$$

Gradients are as in Section 1.4 (replace  $\mathbf{x}$  by  $\tilde{\mathbf{x}}$  in the forward pass, keep the target  $\mathbf{x}$ ). For small additive Gaussian noise of variance  $\sigma^2$  and MSE loss,  $g_\theta(f_\phi(\mathbf{x})) - \mathbf{x} \approx \sigma^2 \nabla_{\mathbf{x}} \log p(\mathbf{x})$  (score-matching connection).

**Sparse Autoencoder.** Encourage low average activation in the code by a KL penalty

$$\Omega_{\text{sparse}}(\mathbf{H}) = \beta \sum_{j=1}^r \text{KL}(\rho \parallel \hat{\rho}_j), \quad \hat{\rho}_j = \frac{1}{n} \sum_{i=1}^n h_{ij}, \quad \text{KL}(\rho \parallel \hat{\rho}) = \rho \log \frac{\rho}{\hat{\rho}} + (1-\rho) \log \frac{1-\rho}{1-\hat{\rho}}. \quad (7)$$

The gradient contribution to the encoder pre-activation for unit  $j$  adds

$$\frac{\partial \Omega_{\text{sparse}}}{\partial a_j^{(e)}} \approx \beta \left( -\frac{\rho}{\hat{\rho}_j} + \frac{1-\rho}{1-\hat{\rho}_j} \right) \sigma'_e(a_j^{(e)}),$$

broadcast to all samples; in practice one uses mini-batch estimates of  $\hat{\rho}_j$ .

**Contractive Autoencoder (CAE).** Penalize the encoder Jacobian to enforce local robustness:

$$\Omega_{\text{contr}}(\phi) = \frac{\gamma}{n} \sum_{i=1}^n \left\| \nabla_{\mathbf{x}} f_{\phi}(\mathbf{x}_i) \right\|_F^2. \quad (8)$$

For elementwise  $\sigma_e$  with rows  $\mathbf{w}_{e,j}^\top$  of  $\mathbf{W}_e$ ,

$$\left\| \nabla_{\mathbf{x}} f_{\phi}(\mathbf{x}) \right\|_F^2 = \sum_{j=1}^r (\sigma'_e(a_j^{(e)}))^2 \|\mathbf{w}_{e,j}\|_2^2,$$

so CAE adds a data-dependent weight decay scaled by  $(\sigma'_e(a_j^{(e)}))^2$ .

## 1.7 Algorithm (Autoencoder Training)

1. **Input:** data  $\{\mathbf{x}_i\}$ , code dimension  $r$ , activations  $(\sigma_e, \sigma_d)$ , loss  $\ell$ , regularization hyperparameters  $(\lambda, \beta, \gamma)$ .
2. **Initialize:** parameters  $\phi, \theta$  (e.g., Xavier/He).
3. **Iterate (SGD/Adam on mini-batches):**
  - (a) (Optional corruption) draw  $\tilde{\mathbf{x}}$  for DAE.
  - (b) Forward:  $\mathbf{h} = \sigma_e(\mathbf{W}_e \mathbf{x} + \mathbf{b}_e)$ ;  $\hat{\mathbf{x}} = \sigma_d(\mathbf{W}_d \mathbf{h} + \mathbf{b}_d)$ .
  - (c) Compute loss  $\hat{\mathcal{L}}_{\mathcal{B}}$  (ERM + regularizers).
  - (d) Backprop: compute  $\Delta^{(d)}, \Delta^{(e)}$ ; form gradients Eqns. (3)–(4) plus regularizer terms.
  - (e) Update parameters.
4. **Output:** trained encoder  $f_{\phi}$  (features) and decoder  $g_{\theta}$  (reconstruction).

## 1.8 Summary of Variables and Their Dimensions

- $\mathbf{x}_i \in \mathbb{R}^d$ : input vector (column).
- $r$ : code dimension (bottleneck),  $1 \leq r \leq d$ .
- Encoder parameters:  $\mathbf{W}_e \in \mathbb{R}^{r \times d}$ ,  $\mathbf{b}_e \in \mathbb{R}^r$ ; pre-activation  $\mathbf{a}^{(e)} \in \mathbb{R}^r$ ; code  $\mathbf{h} \in \mathbb{R}^r$ .
- Decoder parameters:  $\mathbf{W}_d \in \mathbb{R}^{d \times r}$ ,  $\mathbf{b}_d \in \mathbb{R}^d$ ; pre-activation  $\mathbf{a}^{(d)} \in \mathbb{R}^d$ ; reconstruction  $\hat{\mathbf{x}} \in \mathbb{R}^d$ .
- Loss  $\ell(\mathbf{x}, \hat{\mathbf{x}}) \in \mathbb{R}_{\geq 0}$ ; regularization  $\lambda, \beta, \gamma \in \mathbb{R}_{\geq 0}$ .
- (Linear AE/PCA)  $\mathbf{S} \in \mathbb{R}^{d \times d}$  covariance; projector  $\mathbf{P} = \mathbf{W}^\top \mathbf{W}$  with  $\mathbf{W} \in \mathbb{R}^{r \times d}$ ,  $\mathbf{W} \mathbf{W}^\top = \mathbf{I}_r$ .

## 1.9 Summary

From first principles, an autoencoder minimizes empirical reconstruction error Eqn. (1) under an encoder–decoder factorization. Exact gradients are obtained by the chain rule Eqns. (2)–(4), yielding standard backprop updates (with tied-weights handled by parameter sharing). In the linear, centered, tied, orthonormal case, the optimal AE computes the orthogonal projection onto the PCA principal subspace, as shown by maximizing  $\text{tr}(\mathbf{SP})$  via Eckart–Young. Regularized variants (denoising, sparse, contractive) add principled penalties or noise models that shape the learned representation while preserving the same backpropagation structure.