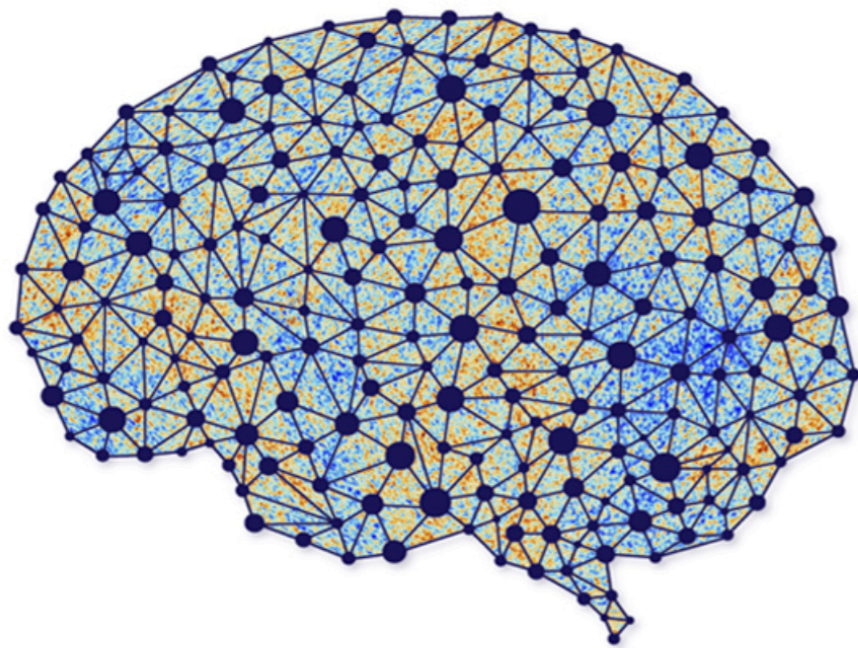


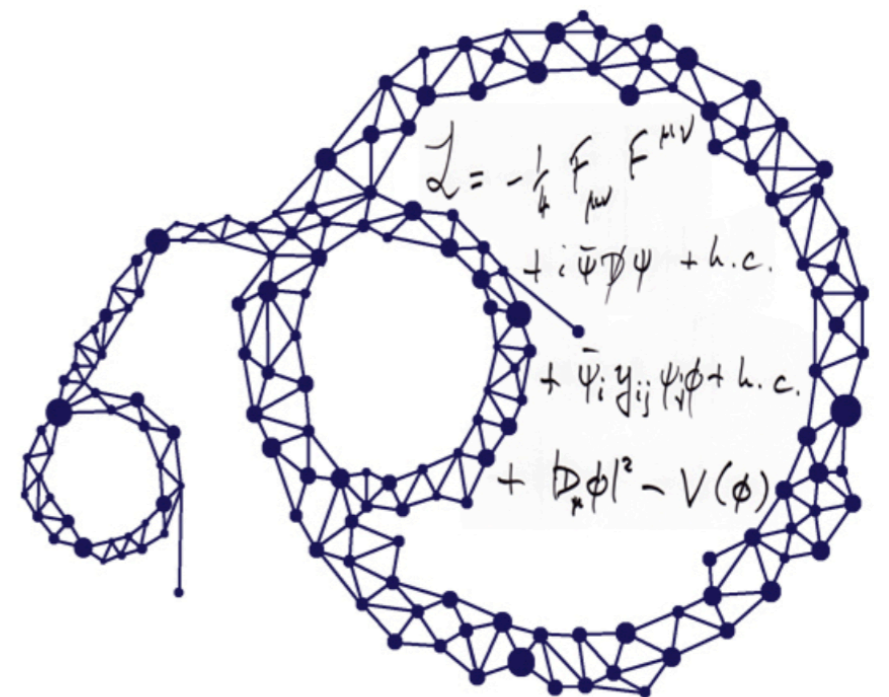
PHY 835: Machine Learning in Physics

Lecture 11: Unsupervised Learning (PCA)

February 27, 2024



AI
∩
Universe



Outline for today

- Unsupervised learning
- Challenges of High-dimensional data
- Principal component analysis (PCA)
- Multi-dimensional scaling (MDS)

References: 1803.08823, Deep Learning Book

Unsupervised Learning

- Discovering structure in unlabelled data.
- Two ways: 1) some appropriate numerical measure (e.g. distance in some representation space). 2) with visualizations.
- Need to dimensionally reduce data as it is impractical for datasets involving large number of measured features (e.g. images)
- We call the dimensionally reduced space **latent space**.
- By dimensional reduction we often loose information. This is not necessarily bad. By loosing only irrelevant information, we can find good representations.

Challenges of High-dimensional Data

- **High-dimensional data lives near the edge of sample space.**
- Consider data distributed uniformly at random in a D-dimensional hypercube $C = [-e/2, e/2]^D$. Probability of a data point inside a D-dimensional hypersphere S of radius $e/2$ centered at the origin:

$$p(\|\mathbf{x}\|_2 < e/2) \sim (1/2)^D \rightarrow 0 \text{ exponentially as } D \rightarrow \infty$$

- Most of the data will concentrate outside the hypersphere, in the corners of the hypercube.
- Recall this property underlies properties of statistical systems such as the Maxwell distribution.

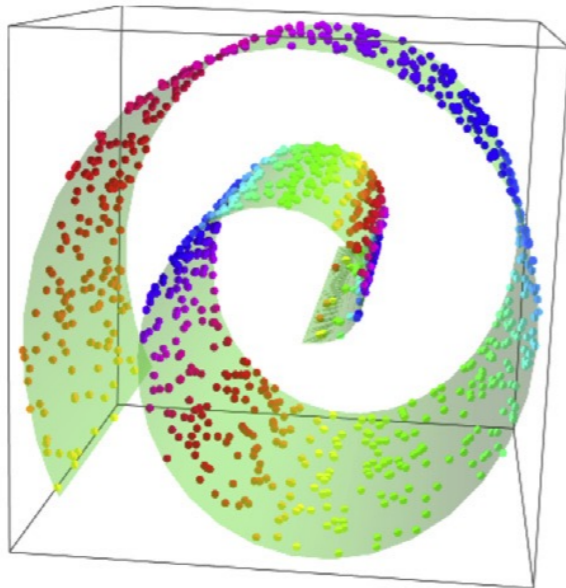
Challenges of High-dimensional Data

- **Real-world data is usually not random or uniformly distributed** (data lives in a lower-dim. space compared with original space).
- “Blessing of non-uniformity”: Data will typically be locally smooth (local variation will not incur a change in the target variable).
- Examples: thermodynamics variables (temperature, pressure, etc) are not sensitive to variations of the dynamical variables (position and momentum of individual particles); small number of order parameters in statistical systems with large number of dofs.
- Objective: preserve relative pairwise distances between data points when going to latent space.

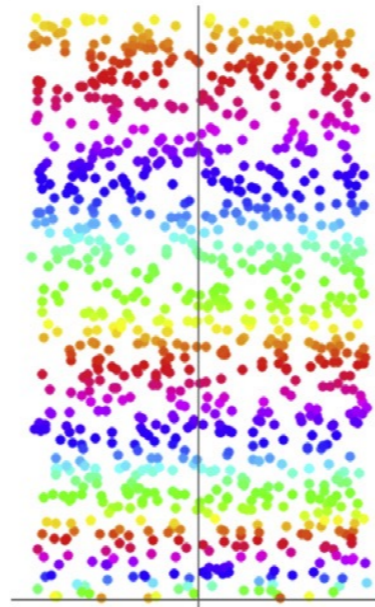
Challenges of High-dimensional Data

- **Intrinsic dimensionality and the crowding problem:**

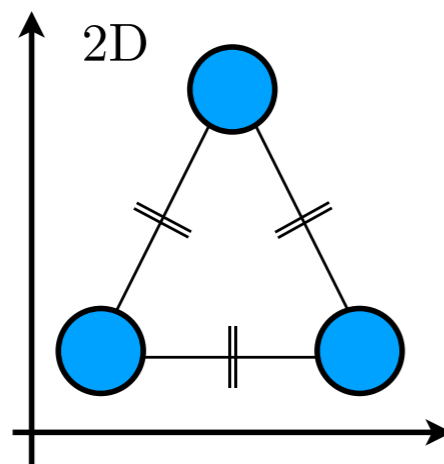
a)



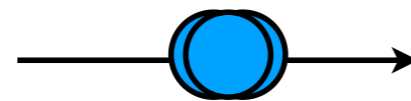
b)



Intrinsic dim = min. # parameters to parametrize the data.



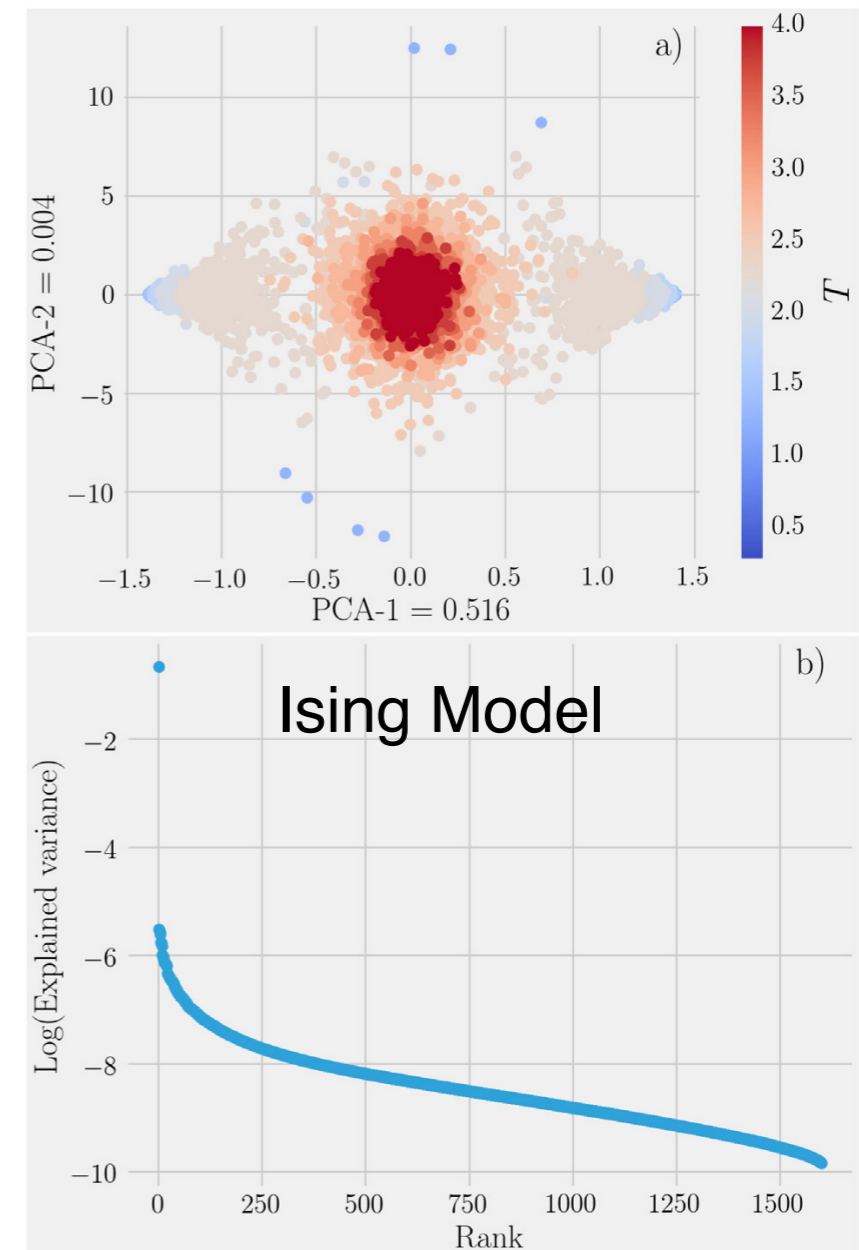
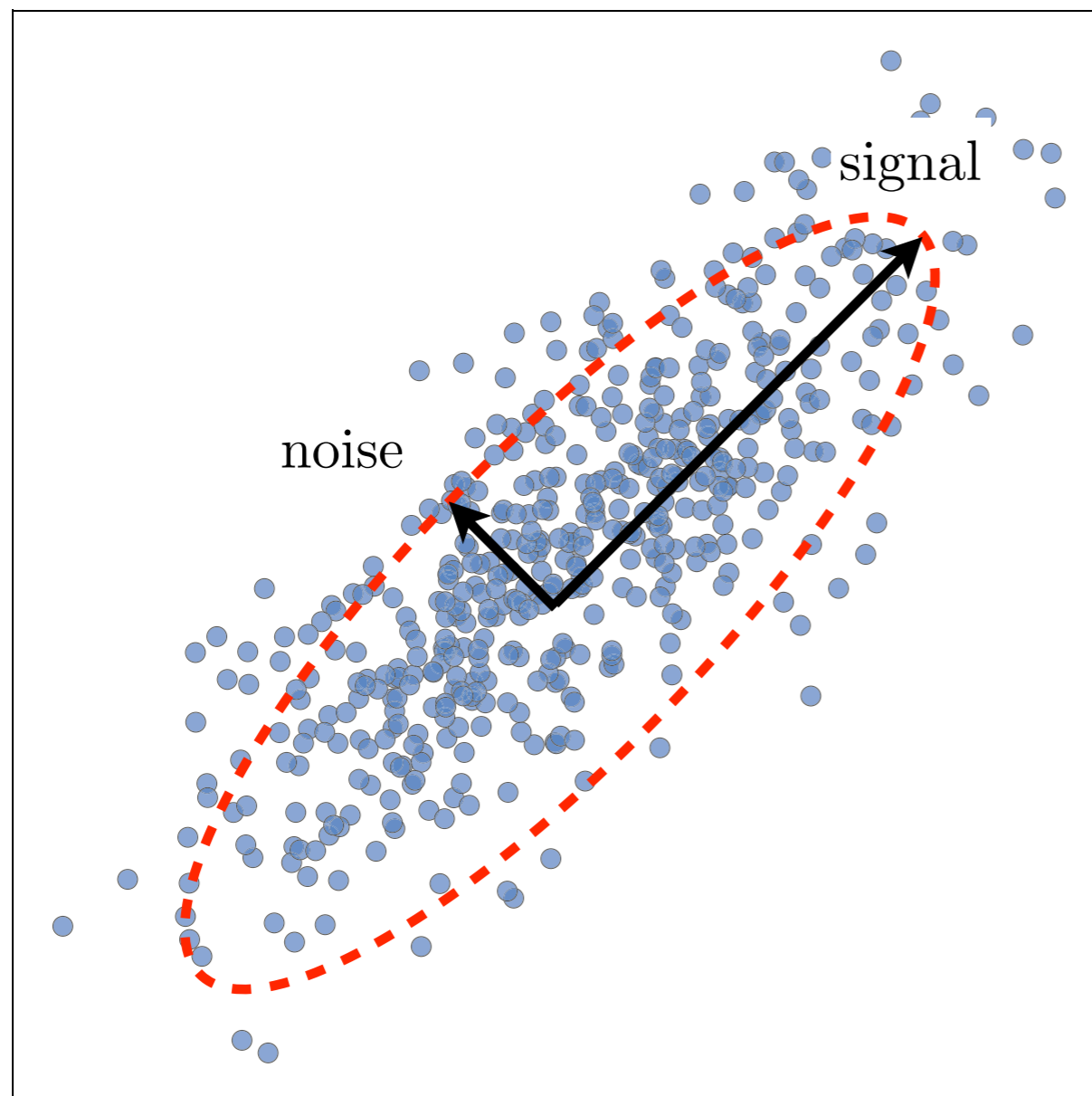
1D



Attempts to represent data in a space with $\text{dim} < \text{intrinsic dimensionality}$ lead to a “crowding” problem.

Principal Component Analysis (PCA)

- Perform an orthogonal transformation of the data to find the high variance directions \Leftrightarrow minimizing the error in projection.



PCA — Minimizing Decoding Error

- Suppose we have a collection of N points $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ in \mathbb{R}^n .
- Compress them into code vectors $\{\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(N)}\}$ in \mathbb{R}^l with $l < n$
- Encoding function: $f(\mathbf{x}) = \mathbf{c}$; Decoding function: $\mathbf{x} = g(\mathbf{c})$
- Encoding + decoding: $\tilde{\mathbf{x}} = g(f(\mathbf{x}))$
- A measure of goodness for your compression is how accurate is this encoding+decoding:

$$\|\mathbf{x} - \tilde{\mathbf{x}}\| \ll 1 \quad (\text{good compression})$$

PCA — Minimizing Decoding Error

- Let $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$ where $\mathbf{D} \in \mathbb{R}^{n \times l}$ is a matrix defining the decoding.
- Columns of \mathbf{D} are orthogonal to each other and have unit norm.
- Minimizing the loss: $\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2$

or equivalently (and more conveniently): $\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2^2$

- The function to be minimized: $(\mathbf{x} - g(\mathbf{c}))^\top (\mathbf{x} - g(\mathbf{c}))$
 $= \mathbf{x}^\top \mathbf{x} - \mathbf{x}^\top g(\mathbf{c}) - g(\mathbf{c})^\top \mathbf{x} + g(\mathbf{c})^\top g(\mathbf{c}) = \mathbf{x}^\top \mathbf{x} - 2\mathbf{x}^\top g(\mathbf{c}) + g(\mathbf{c})^\top g(\mathbf{c})$
- Omit the first term which does not depend on \mathbf{c} :

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} -2\mathbf{x}^\top g(\mathbf{c}) + g(\mathbf{c})^\top g(\mathbf{c})$$

PCA — Minimizing Decoding Error

- Using the definition of the decoding function:

$$\begin{aligned} \mathbf{c}^* &= \arg \min_{\mathbf{c}} -2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{D}^\top \mathbf{D}\mathbf{c} \\ &= \arg \min_{\mathbf{c}} -2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{I}_l \mathbf{c} \end{aligned}$$

because the columns of \mathbf{D} are orthogonal and have unit norm.

- The optimization problem has the solution:

$$\begin{aligned} \nabla_{\mathbf{c}}(-2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{c}) &= \mathbf{0} \\ -2\mathbf{D}^\top \mathbf{x} + 2\mathbf{c} &= \mathbf{0} \\ \mathbf{c} &= \mathbf{D}^\top \mathbf{x}. \end{aligned}$$

- The encoding function: $f(\mathbf{x}) = \mathbf{D}^\top \mathbf{x}$
- PCA reconstruction operation: $r(\mathbf{x}) = g(f(\mathbf{x})) = \mathbf{D}\mathbf{D}^\top \mathbf{x}$

PCA — Minimizing Decoding Error

- Since we use the same matrix \mathbf{D} to decode all the points, we minimize the Frobenius norm of the matrix of errors computed over all dimensions and all points:

$$\mathbf{D}^* = \arg \min_{\mathbf{D}} \sqrt{\sum_{i,j} \left(x_j^{(i)} - r(\mathbf{x}^{(i)})_j \right)^2} \text{ subject to } \mathbf{D}^\top \mathbf{D} = \mathbf{I}_l.$$

- Consider $l = 1$ (exercise: generalization to other l), then $\mathbf{D} = \mathbf{d}$

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{d}\mathbf{d}^\top \mathbf{x}^{(i)}\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1.$$

- Some cosmetic changes (noting $\mathbf{d}^\top \mathbf{x}^{(i)}$ is a scalar, and so its transpose is equal to itself) give:

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{x}^{(i)\top} \mathbf{d}\mathbf{d}\|_2^2 \text{ subject to } \|\mathbf{d}\|_2 = 1.$$

PCA — Minimizing Decoding Error

- Introduce compact notation by defining the matrix \mathbf{X} :

$$\mathbf{X} \in \mathbb{R}^{m \times n} \quad \mathbf{X}_{i,:} = \mathbf{x}^{(i)\top}$$

- The decoding error is minimized when:

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top\|_F^2 \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1.$$

- The Frobenius norm part:

$$\begin{aligned} \arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top\|_F^2 &= \arg \min_{\mathbf{d}} \text{Tr} \left(\left(\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top \right)^\top \left(\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top \right) \right) \\ &= \arg \min_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top - \mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} + \mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \\ &= \arg \min_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X}) - \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) - \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}) + \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \end{aligned}$$

does not depend on \mathbf{d}

PCA — Minimizing Decoding Error

- Cycle the order of the matrices inside a trace, the Frobenius norm:

$$= \arg \min_{\mathbf{d}} -2 \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top)$$

$$= \arg \min_{\mathbf{d}} -2 \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \mathbf{d} \mathbf{d}^\top)$$

- The constraint $\mathbf{d}^\top \mathbf{d} = 1$ gives:

$$= \arg \min_{\mathbf{d}} -2 \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1$$

- Thus minimizing decoding error is the same as maximizing variance:

$$= \arg \min_{\mathbf{d}} - \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1$$

$$= \arg \max_{\mathbf{d}} \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1$$

$$= \arg \max_{\mathbf{d}} \operatorname{Tr}(\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d}) \text{ subject to } \mathbf{d}^\top \mathbf{d} = 1.$$

PCA — Maximizing Variance

- The covariance matrix of data matrix \mathbf{X} is defined as:

$$\Sigma(\mathbf{X}) = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$$

- $\Sigma(\mathbf{X})_{jj}$ corresponds to the variance of the j -th feature while $\Sigma(\mathbf{X})_{ij}$ measures the covariance (correlation) between feature i & feature j .
- Find a new basis that emphasizes highly variable directions while reducing redundancy between basis vectors. Perform SVD:

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T,$$
$$\begin{aligned} \Sigma(\mathbf{X}) &= \frac{1}{N-1} \mathbf{V}\mathbf{S}\mathbf{U}^T \mathbf{U}\mathbf{S}\mathbf{V}^T \\ &= \mathbf{V} \left(\frac{\mathbf{S}^2}{N-1} \right) \mathbf{V}^T \\ &\equiv \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T. \end{aligned}$$

- The eigenvalues λ_i of $\mathbf{\Lambda}$ are given by $\lambda_i = s_i^2 / (N-1)$.

PCA — Maximizing Variance

- To reduce the dimensionality of data from n to l , construct the $n \times l$ projection matrix \mathbf{V}_l by selecting the singular components with the l largest singular values. The projection is then

$$\mathbf{Y} = \mathbf{X}\mathbf{V}_l$$

- The singular vector with the largest singular value (largest variance) is the first principal component; the singular vector with the second largest variance is the second principal component, etc.
- Common in data visualization is to project on the first few principal components (as long as a large part of the variance is explained in those components, e.g., Ising Model).
- Low explained variance may imply that the intrinsic dimensionality of the data is high, or it cannot be captured by a linear representation.

Multidimensional Scaling (MDS)

- Non-linear dimensional reduction technique which preserves the pairwise distance (or dissimilarity) d_{ij} between data points.
- **Metric MDS:** the latent coordinates are obtained by minimizing:

$$\tilde{\mathbf{Y}} = \arg \min_{\mathbf{Y}} \sum_{i < j} w_{ij} |d_{ij}(\mathbf{X}) - d_{ij}(\mathbf{Y})|,$$

- w_{ij} specifies the level of confidence (precision) in the value of $d_{ij}(\mathbf{X})$. If Euclidean metric is used, MDS=PCA; known as **classical scaling**.

<https://stats.stackexchange.com/questions/14002/whats-the-difference-between-principal-component-analysis-and-multidimensional>

- MDS (metric or non-metric) is a generalization of PCA.
- **Non-metric MDS:** d_{ij} can be any distance matrix that preserves the ordination, i.e., if $d_{12}(\mathbf{X}) < d_{13}(\mathbf{X})$ then $d_{12}(\mathbf{Y}) < d_{13}(\mathbf{Y})$.

Multidimensional Scaling (MDS)

- Both MDS and PCA can be implemented using standard Python packages such as Scikit.
- **MDS** algorithms have a scaling of $\mathcal{O}(N^3)$ where $N = \#$ data points.
- Sample-based methods can reduce this scaling to $\mathcal{O}(N \log N)$.
- **PCA** has a scaling of $\mathcal{O}(Np^2 + p^3)$ for a complete decomposition.
 - computation of covariance matrix
 - eigenvalue decomposition where $p = \#$ features
- Can be improved to give a $\mathcal{O}(Np^2 + p)$ scaling for PCA if only a few principal components are desired.