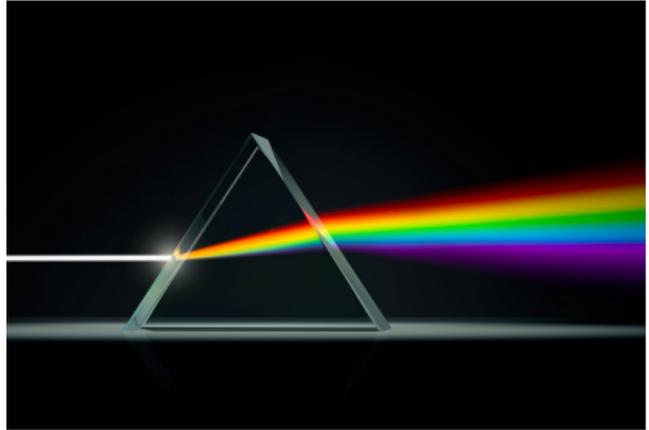
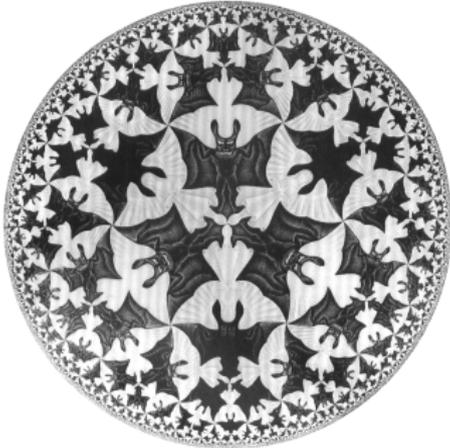


Manifold Learning and Spectral Methods

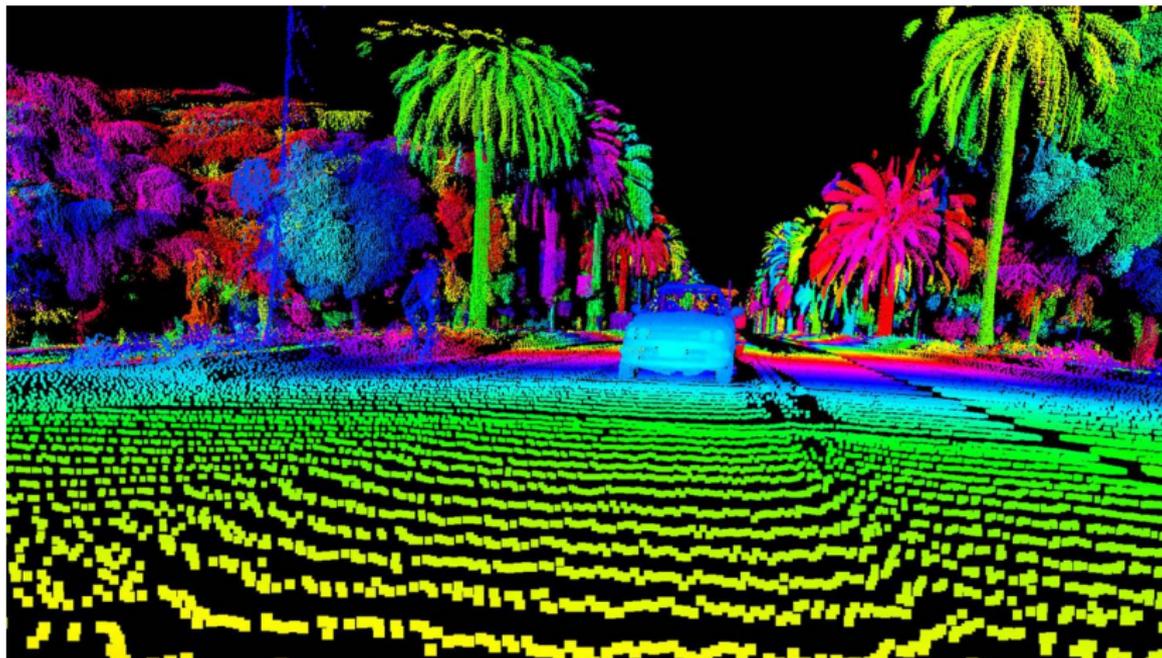


David Pfau



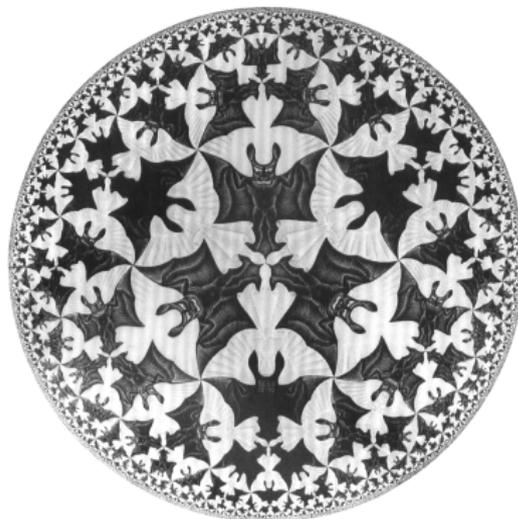
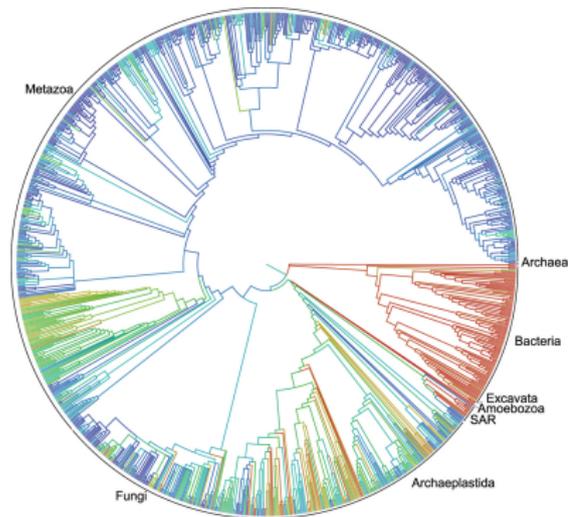
Buenos Aires MLSS, 28 June 2018

What this tutorial is about



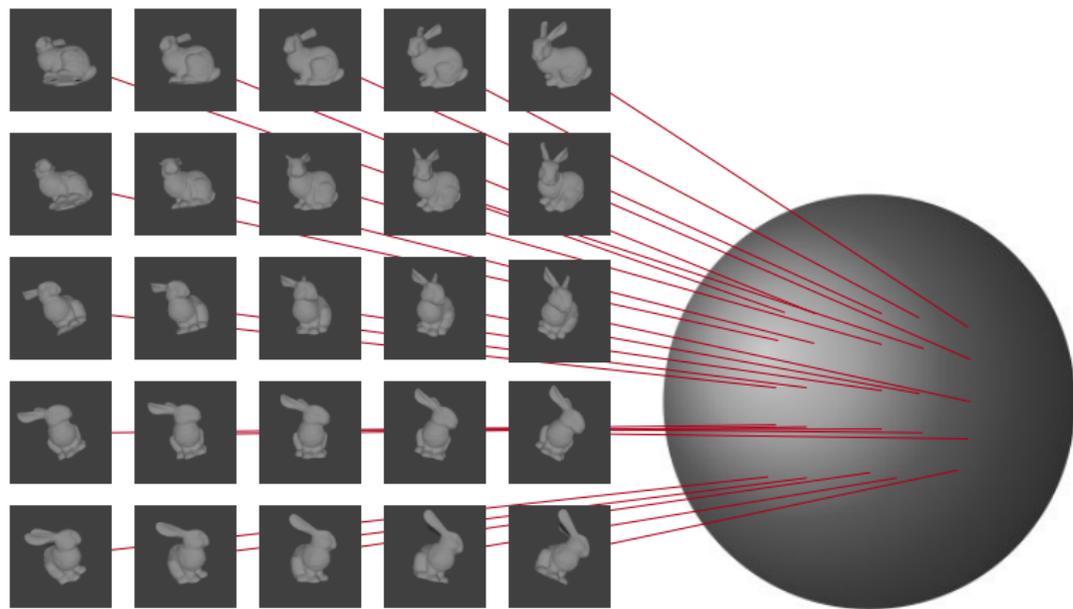
Designing models for graph- or manifold-structured [input](#)

What this tutorial is about



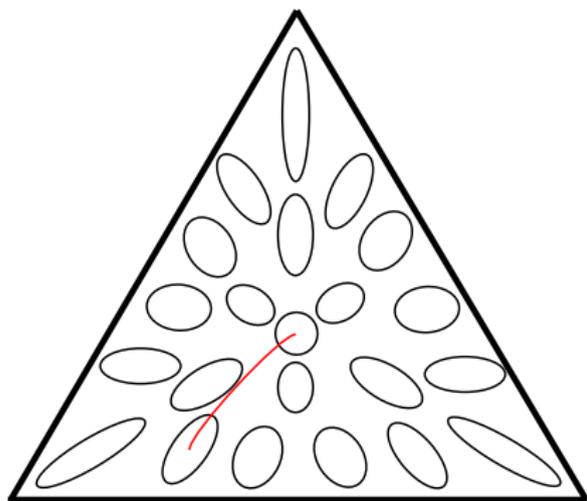
Designing models for data with **latent** graph- or manifold-structure

What this tutorial is about



Discovering latent manifold structure in data

What this tutorial is *not* about



$$\mathcal{F}(\theta) = \mathbb{E}_x [\nabla_{\theta} \log p(x|\theta) \nabla_{\theta} \log p(x|\theta)^T]$$

$$\theta_{t+1} \leftarrow \theta_t + \mathcal{F}^{-1} \nabla_{\theta} \log p(x_t|\theta_t)$$

Information geometry - the manifold structure of **parameter** space

Spectral Learning

Any learning algorithm where the usual tools (SGD, EM, MCMC...) are replaced with **eigendecomposition or SVD**

Spectral Learning

Any learning algorithm where the usual tools (SGD, EM, MCMC...) are replaced with **eigendecomposition** or **SVD**

Solves a **nonconvex** optimization exactly!

$$\mathbf{A} = \mathbf{\Phi} \mathbf{\Lambda} \mathbf{\Phi}^T$$

$$\min_{\substack{\mathbf{\Phi} \\ \mathbf{\Phi}^T \mathbf{\Phi} = \mathbf{I}}} \text{Tr}(\mathbf{\Phi}^T \mathbf{A} \mathbf{\Phi})$$

Spectral Learning

Any learning algorithm where the usual tools (SGD, EM, MCMC...) are replaced with **eigendecomposition or SVD**

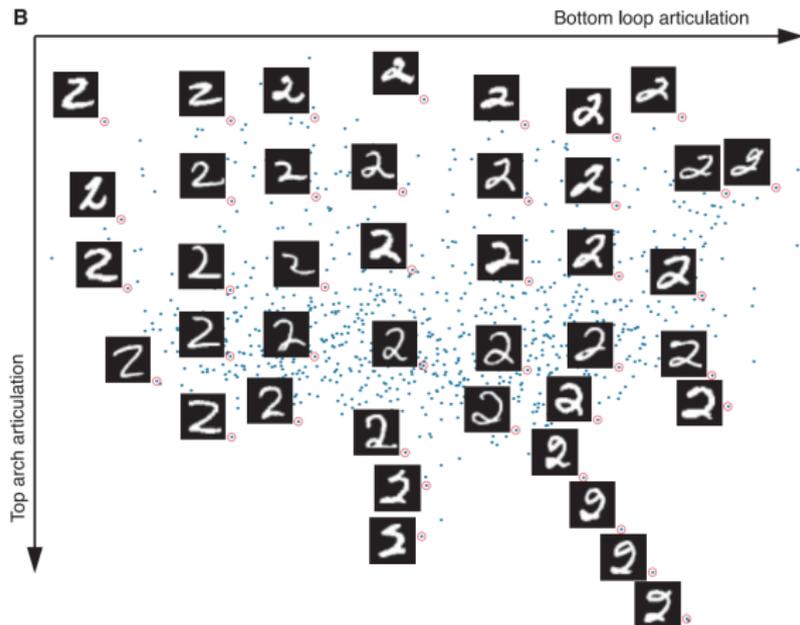
Solves a **nonconvex** optimization exactly!

$$\mathbf{A} = \Phi \Lambda \Phi^T$$

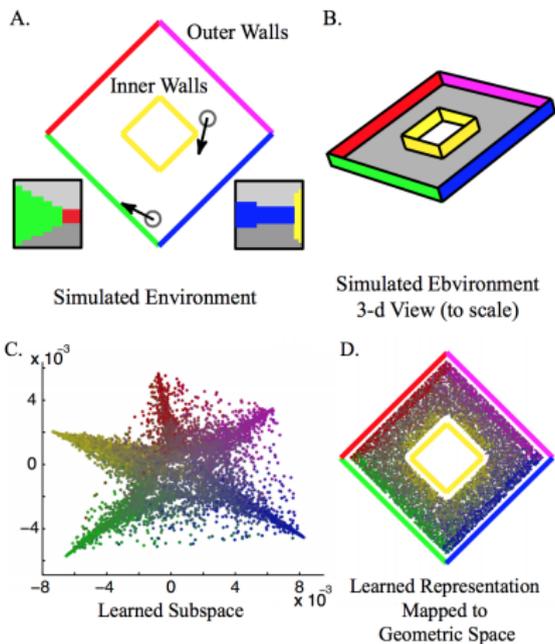
$$\min_{\substack{\Phi \\ \Phi^T \Phi = \mathbf{I}}} \text{Tr}(\Phi^T \mathbf{A} \Phi)$$

Functions on \mathbb{R}^n can be decomposed with Fourier basis. Spectral analysis **generalizes Fourier analysis to manifolds and graphs**

What is this good for?



What is this good for?



Outline

Theory

Graphs

Manifolds and Differential Geometry

Spectral Theory

Outline

Theory

Graphs

Manifolds and Differential Geometry

Spectral Theory

The Geometry of Data

Classic Manifold Learning

Embedding Hierarchies in Hyperbolic Space

Analyzing the Geometry of Deep Generative Models

Outline

Theory

Graphs

Manifolds and Differential Geometry

Spectral Theory

The Geometry of Data

Classic Manifold Learning

Embedding Hierarchies in Hyperbolic Space

Analyzing the Geometry of Deep Generative Models

Spectral Deep Learning

Convolutions on Graphs and Manifolds

Spectral Graph Convolutional Neural Networks

Inference in Spectral Learning with Deep Networks

Part I

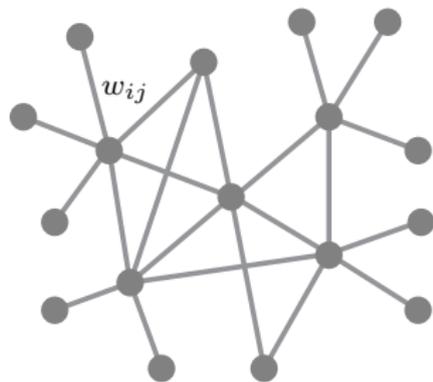
Theory

Part Ia

Graph Theory (in three slides)

Graphs

Weighted undirected graph \mathcal{G} with vertices $\mathcal{V} = \{1, \dots, n\}$, edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ and edge weights $w_{ij} \geq 0$ for $(i, j) \in \mathcal{E}$

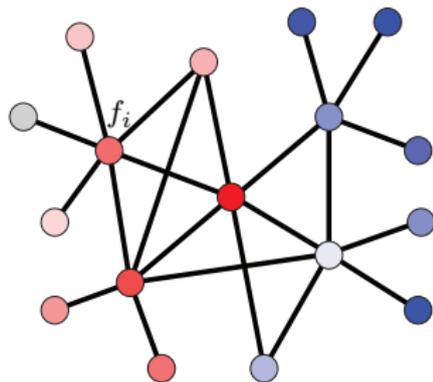


Graphs

Weighted undirected graph \mathcal{G} with vertices $\mathcal{V} = \{1, \dots, n\}$, edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ and edge weights $w_{ij} \geq 0$ for $(i, j) \in \mathcal{E}$

Functions over the vertices

$$L^2(\mathcal{V}) = \{f : \mathcal{V} \rightarrow \mathbb{R}\}$$

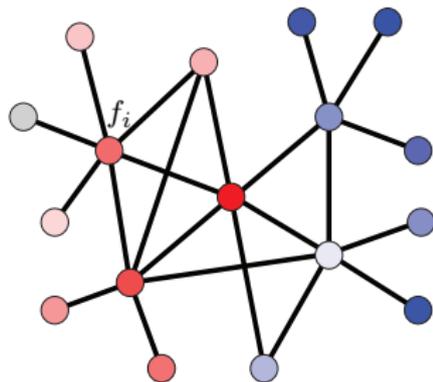


Graphs

Weighted undirected graph \mathcal{G} with vertices $\mathcal{V} = \{1, \dots, n\}$, edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ and edge weights $w_{ij} \geq 0$ for $(i, j) \in \mathcal{E}$

Functions over the vertices

$L^2(\mathcal{V}) = \{f : \mathcal{V} \rightarrow \mathbb{R}\}$ represented as vectors $\mathbf{f} = (f_1, \dots, f_n)$



Graphs

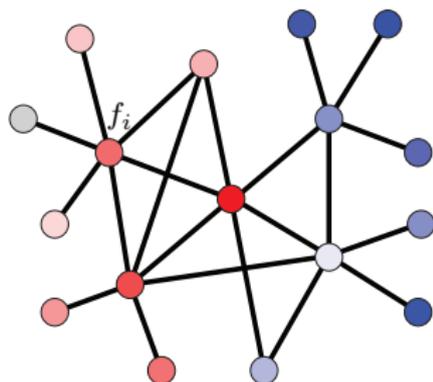
Weighted undirected graph \mathcal{G} with vertices $\mathcal{V} = \{1, \dots, n\}$, edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ and edge weights $w_{ij} \geq 0$ for $(i, j) \in \mathcal{E}$

Functions over the vertices

$L^2(\mathcal{V}) = \{f : \mathcal{V} \rightarrow \mathbb{R}\}$ represented as vectors $\mathbf{f} = (f_1, \dots, f_n)$

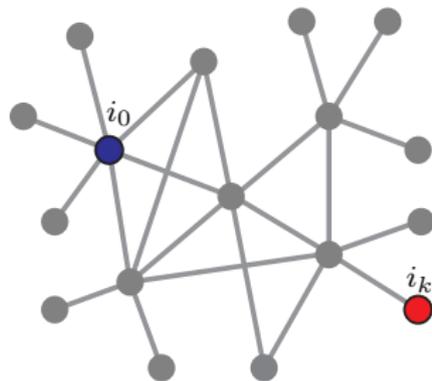
Inner product

$$\langle f, g \rangle_{L^2(\mathcal{V})} = \sum_{i \in \mathcal{V}} f_i g_i = \mathbf{f}^\top \mathbf{g}$$



Shortest paths on graphs

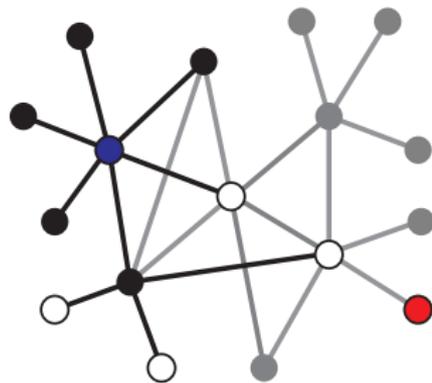
Dijkstra's Algorithm: Find sequence $i_0, \dots, i_k \in \mathcal{V}$ from **source** i_0 to **sink** i_k that minimizes $\sum_t w_{i_t i_{t+1}}$



Shortest paths on graphs

Dijkstra's Algorithm: Find sequence $i_0, \dots, i_k \in \mathcal{V}$ from **source** i_0 to **sink** i_k that minimizes $\sum_t w_{i_t i_{t+1}}$

Create queue of **active** nodes and **shortest-path tree**, starting from source



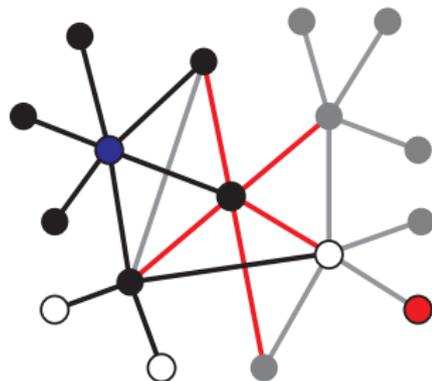
Shortest paths on graphs

Dijkstra's Algorithm: Find sequence $i_0, \dots, i_k \in \mathcal{V}$ from **source** i_0 to **sink** i_k that minimizes $\sum_t w_{i_t i_{t+1}}$

Create queue of **active** nodes and **shortest-path tree**, starting from source

Remove active node i closest to source, compute distance to neighbors

$$\text{dist}[j] = \text{dist}[i] + w_{ij}.$$



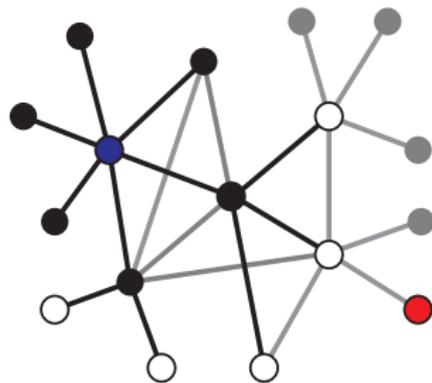
Shortest paths on graphs

Dijkstra's Algorithm: Find sequence $i_0, \dots, i_k \in \mathcal{V}$ from **source** i_0 to **sink** i_k that minimizes $\sum_t w_{i_t i_{t+1}}$

Create queue of **active** nodes and **shortest-path tree**, starting from source

Remove active node i closest to source, compute distance to neighbors
 $\text{dist}[j] = \text{dist}[i] + w_{ij}$.

If $\text{dist}[j]$ is less than before, set i as parent node of j in shortest path tree.



Shortest paths on graphs

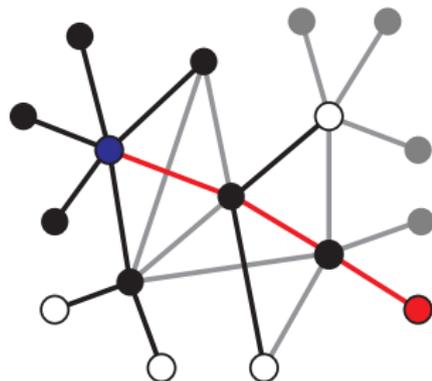
Dijkstra's Algorithm: Find sequence $i_0, \dots, i_k \in \mathcal{V}$ from **source** i_0 to **sink** i_k that minimizes $\sum_t w_{i_t i_{t+1}}$

Create queue of **active** nodes and **shortest-path tree**, starting from source

Remove active node i closest to source, compute distance to neighbors
 $\text{dist}[j] = \text{dist}[i] + w_{ij}$.

If $\text{dist}[j]$ is less than before, set i as parent node of j in shortest path tree.

Stop when sink is reached. Path to the root of the tree is the shortest path.



Shortest paths on graphs

Dijkstra's Algorithm: Find sequence $i_0, \dots, i_k \in \mathcal{V}$ from **source** i_0 to **sink** i_k that minimizes $\sum_t w_{i_t i_{t+1}}$

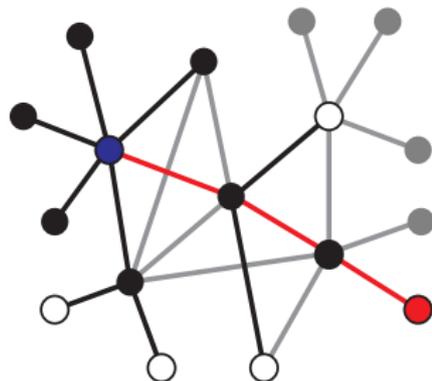
Create queue of **active** nodes and **shortest-path tree**, starting from source

Remove active node i closest to source, compute distance to neighbors
 $\text{dist}[j] = \text{dist}[i] + w_{ij}$.

If $\text{dist}[j]$ is less than before, set i as parent node of j in shortest path tree.

Stop when sink is reached. Path to the root of the tree is the shortest path.

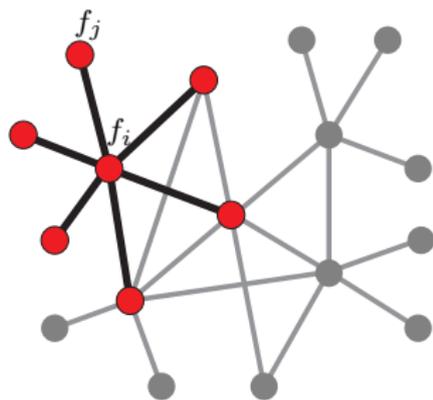
Complexity: $\mathcal{O}(|\mathcal{E}| + |\mathcal{V}|\log|\mathcal{V}|)$



Graph Laplacian

Unnormalized Laplacian $\Delta : L^2(\mathcal{V}) \rightarrow L^2(\mathcal{V})$

$$(\Delta f)_i = \sum_{j:(i,j) \in \mathcal{E}} w_{ij}(f_i - f_j)$$

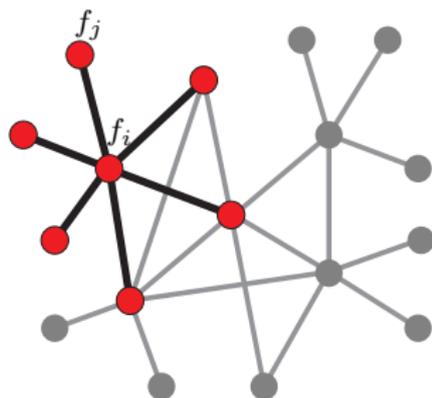


Graph Laplacian

Unnormalized Laplacian $\Delta : L^2(\mathcal{V}) \rightarrow L^2(\mathcal{V})$

$$\begin{aligned}(\Delta f)_i &= \sum_{j:(i,j) \in \mathcal{E}} w_{ij}(f_i - f_j) \\ &= f_i \sum_{j:(i,j) \in \mathcal{E}} w_{ij} - \sum_{j:(i,j) \in \mathcal{E}} w_{ij} f_j\end{aligned}$$

(up to scale) difference between f and its local average



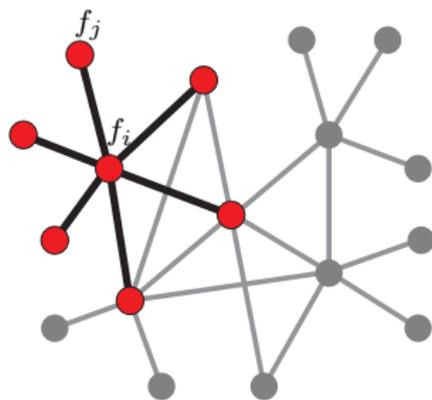
Graph Laplacian

Unnormalized Laplacian $\Delta : L^2(\mathcal{V}) \rightarrow L^2(\mathcal{V})$

$$(\Delta f)_i = \sum_{j:(i,j) \in \mathcal{E}} w_{ij}(f_i - f_j)$$

(up to scale) difference between f and its local average

Represented as a **positive semi-definite** $n \times n$ matrix $\Delta = \mathbf{D} - \mathbf{W}$ where $\mathbf{W} = (w_{ij})$ and $\mathbf{D} = \text{diag}(\sum_{j \neq i} w_{ij})$



Graph Laplacian

Unnormalized Laplacian $\Delta : L^2(\mathcal{V}) \rightarrow L^2(\mathcal{V})$

$$(\Delta f)_i = \sum_{j:(i,j) \in \mathcal{E}} w_{ij}(f_i - f_j)$$

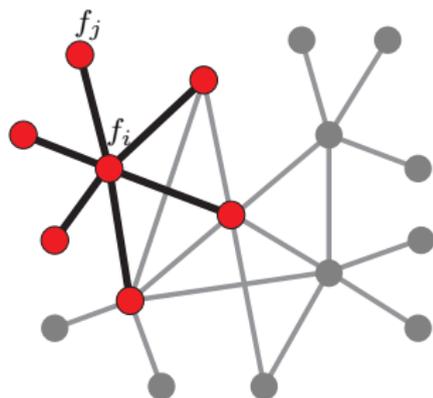
(up to scale) difference between f and its local average

Represented as a **positive semi-definite** $n \times n$ matrix $\Delta = \mathbf{D} - \mathbf{W}$ where $\mathbf{W} = (w_{ij})$ and $\mathbf{D} = \text{diag}(\sum_{j \neq i} w_{ij})$

Dirichlet energy of f

$$\|f\|_{\mathcal{G}}^2 = \frac{1}{2} \sum_{i,j=1}^n w_{ij}(f_i - f_j)^2 = \mathbf{f}^\top \Delta \mathbf{f}$$

measures the **smoothness** of f (how fast it changes locally)



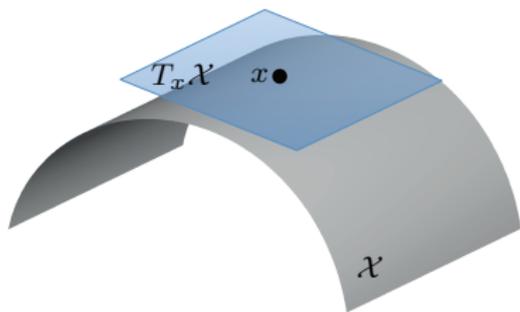
Part Ib

Manifolds and Differential Geometry

Riemannian manifolds

Manifold \mathcal{X} = locally flat space (no non-differentiable corners)

Tangent plane $T_x\mathcal{X}$ = local Euclidean representation of manifold \mathcal{X} around x



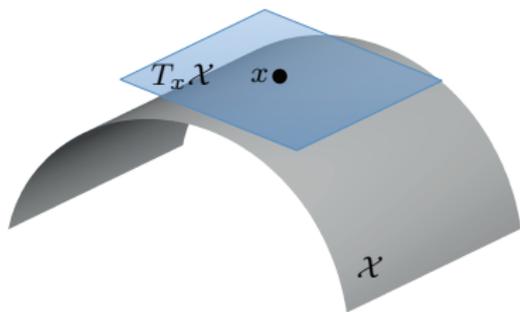
Riemannian manifolds

Manifold \mathcal{X} = locally flat space (no non-differentiable corners)

Tangent plane $T_x\mathcal{X}$ = local Euclidean representation of manifold \mathcal{X} around x

Riemannian metric describes the local intrinsic structure at x

$$\langle \cdot, \cdot \rangle_{T_x\mathcal{X}} : T_x\mathcal{X} \times T_x\mathcal{X} \rightarrow \mathbb{R}$$



Riemannian manifolds

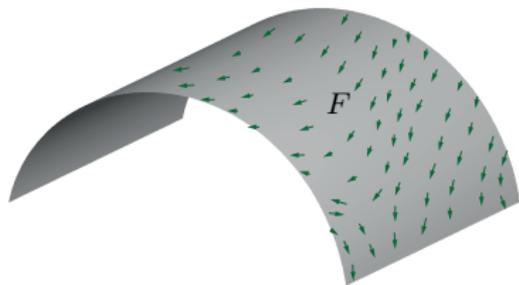
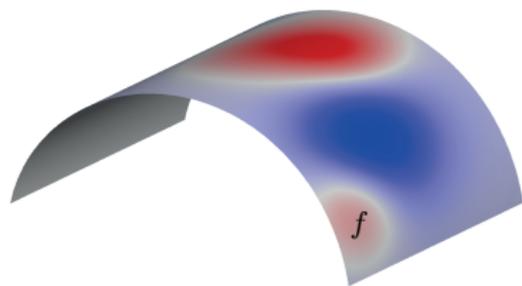
Manifold \mathcal{X} = locally flat space (no non-differentiable corners)

Tangent plane $T_x\mathcal{X}$ = local Euclidean representation of manifold \mathcal{X} around x

Riemannian metric describes the local intrinsic structure at x

$$\langle \cdot, \cdot \rangle_{T_x\mathcal{X}} : T_x\mathcal{X} \times T_x\mathcal{X} \rightarrow \mathbb{R}$$

Scalar fields $f : \mathcal{X} \rightarrow \mathbb{R}$ and **vector fields** $F : \mathcal{X} \rightarrow T\mathcal{X}$



Riemannian manifolds

Manifold \mathcal{X} = locally flat space (no non-differentiable corners)

Tangent plane $T_x\mathcal{X}$ = local Euclidean representation of manifold \mathcal{X} around x

Riemannian metric describes the local intrinsic structure at x

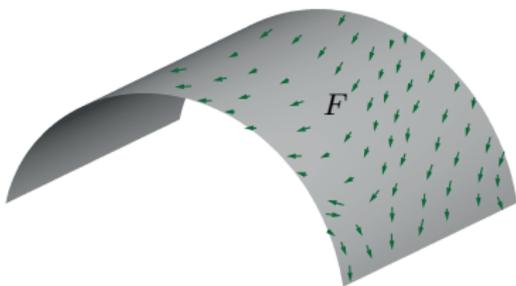
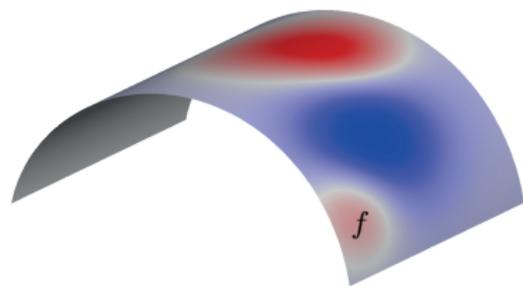
$$\langle \cdot, \cdot \rangle_{T_x\mathcal{X}} : T_x\mathcal{X} \times T_x\mathcal{X} \rightarrow \mathbb{R}$$

Scalar fields $f : \mathcal{X} \rightarrow \mathbb{R}$ and **vector fields** $F : \mathcal{X} \rightarrow T\mathcal{X}$

Inner products

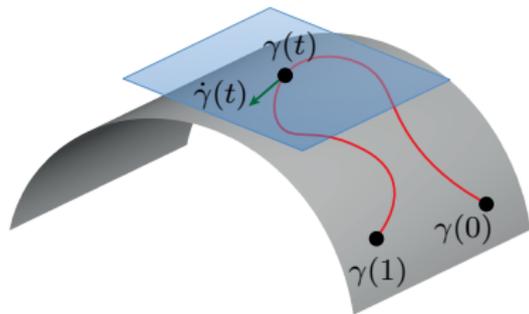
$$\langle f, g \rangle_{L^2(\mathcal{X})} = \int_{\mathcal{X}} f(x)g(x)dx$$

$$\langle F, G \rangle_{L^2(T\mathcal{X})} = \int_{\mathcal{X}} \langle F(x), G(x) \rangle_{T_x\mathcal{X}} dx$$



Shortest paths on manifolds

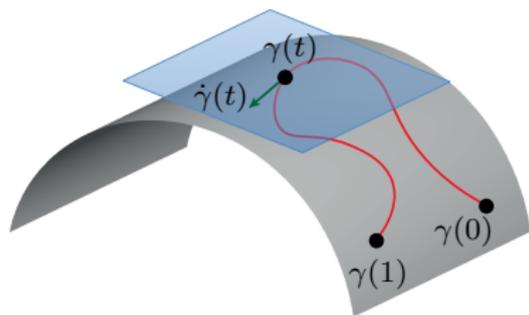
Curve $\gamma(\cdot) : [0, 1] \rightarrow \mathcal{X}$ = smooth path
along manifold



Shortest paths on manifolds

Curve $\gamma(\cdot) : [0, 1] \rightarrow \mathcal{X} =$ smooth path along manifold

Velocity $\dot{\gamma}(\cdot) : [0, 1] \rightarrow T_x \mathcal{X} =$ local direction of change along the curve



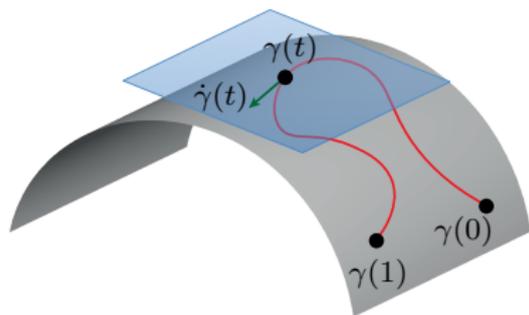
Shortest paths on manifolds

Curve $\gamma(\cdot) : [0, 1] \rightarrow \mathcal{X} =$ smooth path along manifold

Velocity $\dot{\gamma}(\cdot) : [0, 1] \rightarrow T_x \mathcal{X} =$ local direction of change along the curve

Curve energy integrates length at constant velocity

$$\mathcal{S}[\gamma] = \int_0^1 dt \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle_{T_{\gamma(t)} \mathcal{X}}$$



Shortest paths on manifolds

Curve $\gamma(\cdot) : [0, 1] \rightarrow \mathcal{X} =$ smooth path along manifold

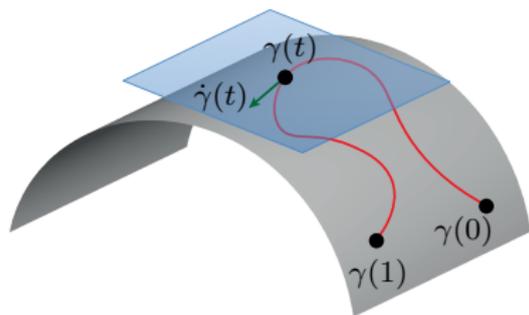
Velocity $\dot{\gamma}(\cdot) : [0, 1] \rightarrow T_x \mathcal{X} =$ local direction of change along the curve

Curve energy integrates length at constant velocity

$$\mathcal{S}[\gamma] = \int_0^1 dt \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle_{T_{\gamma(t)} \mathcal{X}}$$

Geodesic shortest curve between points

$$\gamma^* = \min_{\substack{\gamma \\ \gamma(0)=x_0 \\ \gamma(1)=x_1}} \mathcal{S}[\gamma]$$



Shortest paths on manifolds

Curve $\gamma(\cdot) : [0, 1] \rightarrow \mathcal{X} =$ smooth path along manifold

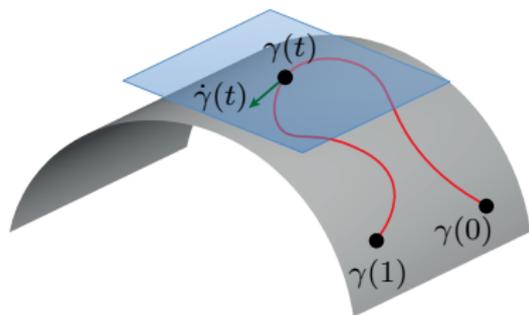
Velocity $\dot{\gamma}(\cdot) : [0, 1] \rightarrow T_x \mathcal{X} =$ local direction of change along the curve

Curve energy integrates length at constant velocity

$$\mathcal{S}[\gamma] = \int_0^1 dt \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle_{T_{\gamma(t)} \mathcal{X}}$$

Geodesic shortest curve between points

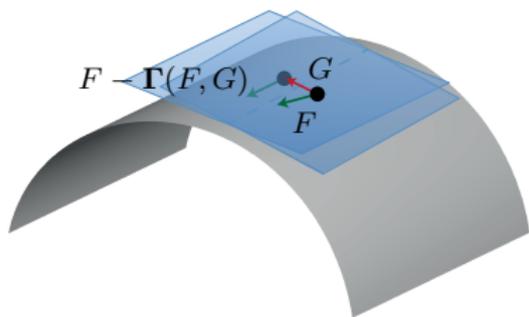
$$\gamma^* = \min_{\substack{\gamma \\ \gamma(0)=x_0 \\ \gamma(1)=x_1}} \mathcal{S}[\gamma]$$



Geodesic is a path such that the velocities are **locally parallel**

Parallel Transport

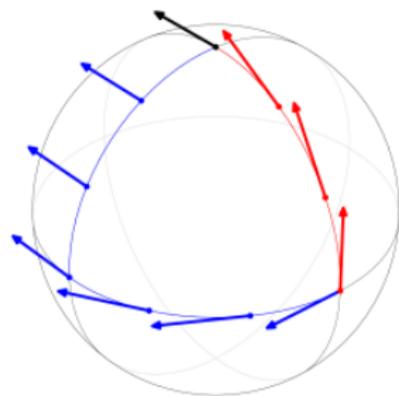
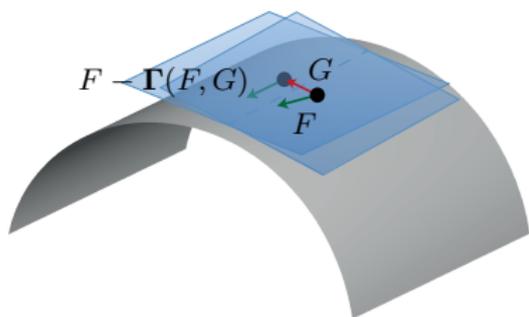
Connection $\Gamma_x(F, G)$: infinitesimal change to vector $F \in T_x\mathcal{X}$ that keeps it locally parallel when moved in the direction $G \in T_x\mathcal{X}$.



Parallel Transport

Connection $\Gamma_x(F, G)$: infinitesimal change to vector $F \in T_x\mathcal{X}$ that keeps it locally parallel when moved in the direction $G \in T_x\mathcal{X}$.

Parallel transport Sequence of vectors $F(t) \in T_{\gamma(t)}\mathcal{X}$ along curve γ that are all locally parallel.

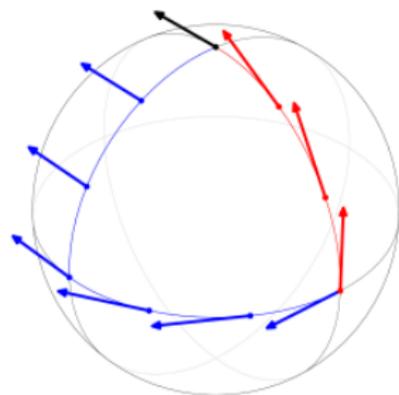
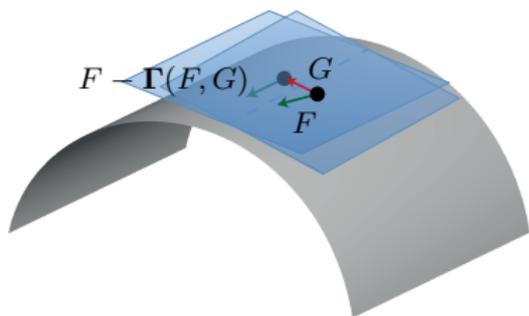


Parallel Transport

Connection $\Gamma_x(F, G)$: infinitesimal change to vector $F \in T_x\mathcal{X}$ that keeps it locally parallel when moved in the direction $G \in T_x\mathcal{X}$.

Parallel transport Sequence of vectors $F(t) \in T_{\gamma(t)}\mathcal{X}$ along curve γ that are all locally parallel.

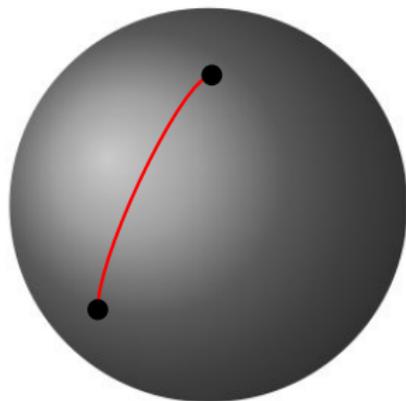
$$\dot{F}(t) + \Gamma_{\gamma(t)}(F(t), \dot{\gamma}(t)) = 0$$



Geodesic equation

Parallel transport:

$$\dot{F}(t) + \mathbf{\Gamma}_{\gamma(t)}(F(t), \dot{\gamma}(t)) = 0$$



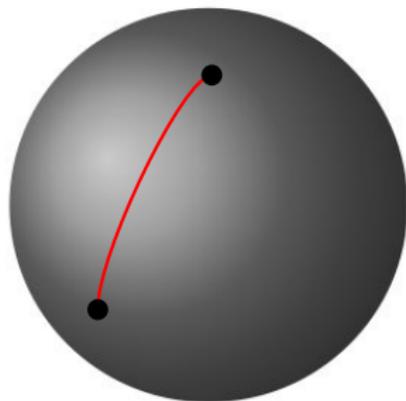
Geodesic equation

Parallel transport:

$$\dot{F}(t) + \mathbf{\Gamma}_{\gamma(t)}(F(t), \dot{\gamma}(t)) = 0$$

The geodesic equation:

$$\ddot{\gamma}(t) + \mathbf{\Gamma}_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t)) = 0$$



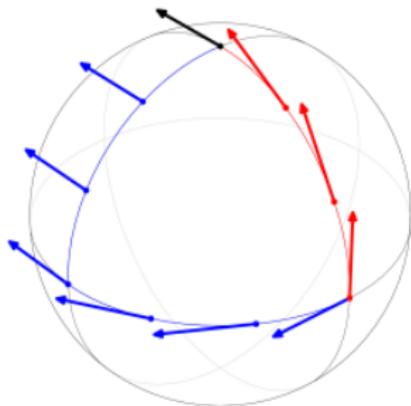
Geodesic equation

Parallel transport:

$$\dot{F}(t) + \mathbf{\Gamma}_{\gamma(t)}(F(t), \dot{\gamma}(t)) = 0$$

The geodesic equation:

$$\ddot{\gamma}(t) + \mathbf{\Gamma}_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t)) = 0$$



A geodesic is a curve such that if you parallel transport its velocity at a point along itself, the vector always remains tangent to the curve

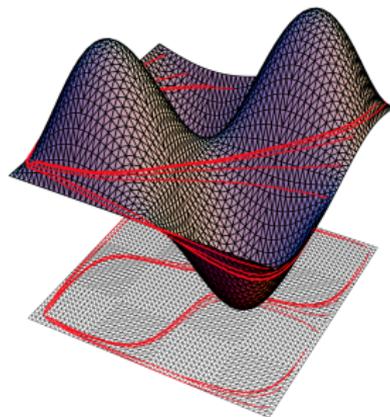
Geodesic equation

Parallel transport:

$$\dot{F}(t) + \mathbf{\Gamma}_{\gamma(t)}(F(t), \dot{\gamma}(t)) = 0$$

The geodesic equation:

$$\ddot{\gamma}(t) + \mathbf{\Gamma}_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t)) = 0$$



A geodesic is a curve such that if you **parallel transport its velocity** at a point along itself, the vector always **remains tangent** to the curve

Solve numerically on meshes with **fast marching method**

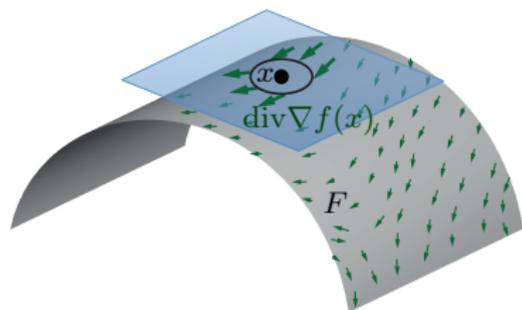
Manifold Laplacian

Laplacian $\Delta : L^2(\mathcal{X}) \rightarrow L^2(\mathcal{X})$

$$\Delta f(x) = -\operatorname{div} \nabla f(x)$$

where **gradient** $\nabla : L^2(\mathcal{X}) \rightarrow L^2(T\mathcal{X})$
and **divergence** $\operatorname{div} : L^2(T\mathcal{X}) \rightarrow L^2(\mathcal{X})$
are adjoint operators

$$\langle F, \nabla f \rangle_{L^2(T\mathcal{X})} = \langle -\operatorname{div} F, f \rangle_{L^2(\mathcal{X})}$$



Manifold Laplacian

Laplacian $\Delta : L^2(\mathcal{X}) \rightarrow L^2(\mathcal{X})$

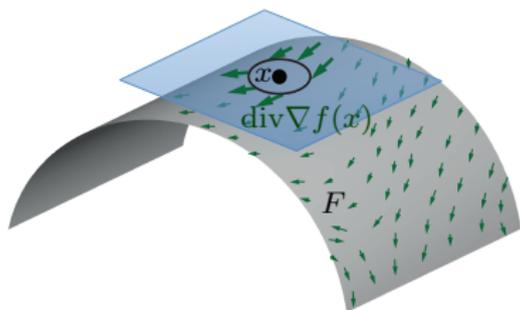
$$\Delta f(x) = -\operatorname{div} \nabla f(x)$$

where **gradient** $\nabla : L^2(\mathcal{X}) \rightarrow L^2(T\mathcal{X})$
and **divergence** $\operatorname{div} : L^2(T\mathcal{X}) \rightarrow L^2(\mathcal{X})$
are adjoint operators

$$\langle F, \nabla f \rangle_{L^2(T\mathcal{X})} = \langle -\operatorname{div} F, f \rangle_{L^2(\mathcal{X})}$$

Laplacian is self-adjoint

$$\langle \Delta f, f \rangle_{L^2(\mathcal{X})} = \langle f, \Delta f \rangle_{L^2(\mathcal{X})}$$



Manifold Laplacian

Laplacian $\Delta : L^2(\mathcal{X}) \rightarrow L^2(\mathcal{X})$

$$\Delta f(x) = -\operatorname{div} \nabla f(x)$$

where **gradient** $\nabla : L^2(\mathcal{X}) \rightarrow L^2(T\mathcal{X})$
and **divergence** $\operatorname{div} : L^2(T\mathcal{X}) \rightarrow L^2(\mathcal{X})$
are adjoint operators

$$\langle F, \nabla f \rangle_{L^2(T\mathcal{X})} = \langle -\operatorname{div} F, f \rangle_{L^2(\mathcal{X})}$$

Laplacian is self-adjoint

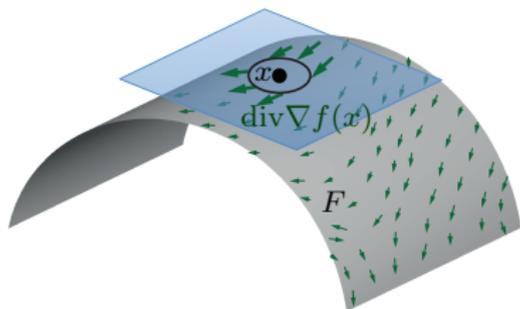
$$\langle \Delta f, f \rangle_{L^2(\mathcal{X})} = \langle f, \Delta f \rangle_{L^2(\mathcal{X})}$$

Continuous limit of graph Laplacian
under some conditions

Dirichlet energy of f

$$\langle \nabla f, \nabla f \rangle_{L^2(T\mathcal{X})} = \int_{\mathcal{X}} f(x) \Delta f(x) dx$$

measures the **smoothness** of f (how fast it changes locally)



Part Ic

Spectral Theory for Graphs and Manifolds

Orthogonal bases on graphs

Find the **smoothest orthogonal basis** $\{\phi_1, \dots, \phi_n\} \subseteq L^2(\mathcal{V})$

$$\min_{\phi_1} E_{\text{Dir}}(\psi_1) \quad \text{s.t.} \quad \|\phi_1\| = 1$$

$$\min_{\phi_k} E_{\text{Dir}}(\psi_k) \quad \text{s.t.} \quad \|\phi_k\| = 1, \quad k = 2, 3, \dots, n$$

$$\phi_k \perp \text{span}\{\phi_1, \dots, \phi_{k-1}\}$$

Orthogonal bases on graphs

Find the **smoothest orthogonal basis** $\{\phi_1, \dots, \phi_n\} \subseteq L^2(\mathcal{V})$

$$\min_{\Phi \in \mathbb{R}^{n \times n}} \text{trace}(\Phi^\top \Delta \Phi) \quad \text{s.t.} \quad \Phi^\top \Phi = \mathbf{I}$$

Orthogonal bases on graphs

Find the **smoothest orthogonal basis** $\{\phi_1, \dots, \phi_n\} \subseteq L^2(\mathcal{V})$

$$\min_{\Phi \in \mathbb{R}^{n \times n}} \text{trace}(\Phi^\top \Delta \Phi) \quad \text{s.t.} \quad \Phi^\top \Phi = \mathbf{I}$$

Solution: $\Phi =$ Laplacian eigenvectors

Laplacian eigenvectors and eigenvalues

Eigendecomposition of a graph Laplacian

$$\Delta = \Phi \Lambda \Phi^\top$$

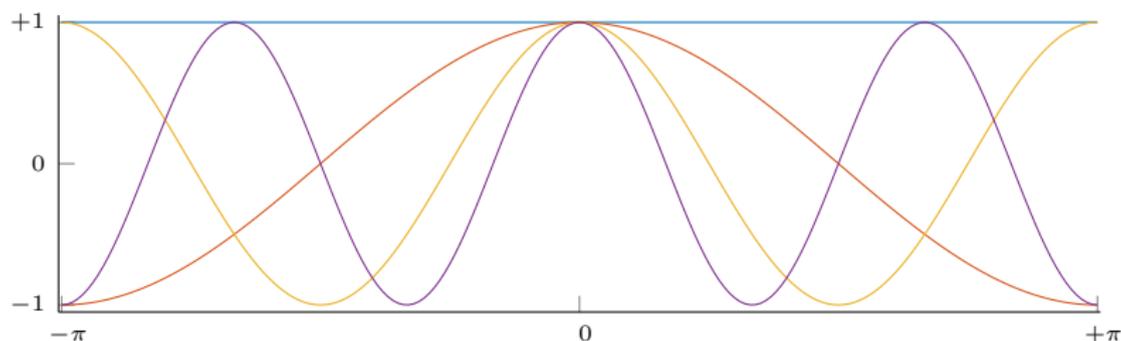
where $\Phi = (\phi_1, \dots, \phi_n)$ are **orthogonal eigenvectors** ($\Phi^\top \Phi = \mathbf{I}$) and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ the corresponding **non-negative eigenvalues**

Laplacian eigenvectors and eigenvalues

Eigendecomposition of a graph Laplacian

$$\Delta = \Phi \Lambda \Phi^\top$$

where $\Phi = (\phi_1, \dots, \phi_n)$ are **orthogonal eigenvectors** ($\Phi^\top \Phi = \mathbf{I}$) and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ the corresponding **non-negative eigenvalues**



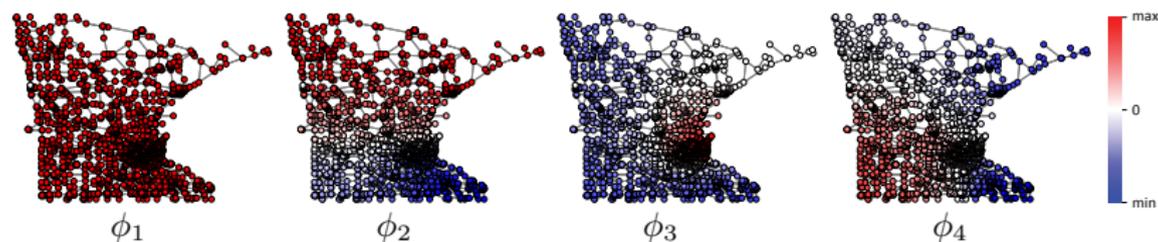
First eigenfunctions of 1D Euclidean Laplacian

Laplacian eigenvectors and eigenvalues

Eigendecomposition of a graph Laplacian

$$\Delta = \Phi \Lambda \Phi^\top$$

where $\Phi = (\phi_1, \dots, \phi_n)$ are **orthogonal eigenvectors** ($\Phi^\top \Phi = \mathbf{I}$) and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ the corresponding **non-negative eigenvalues**



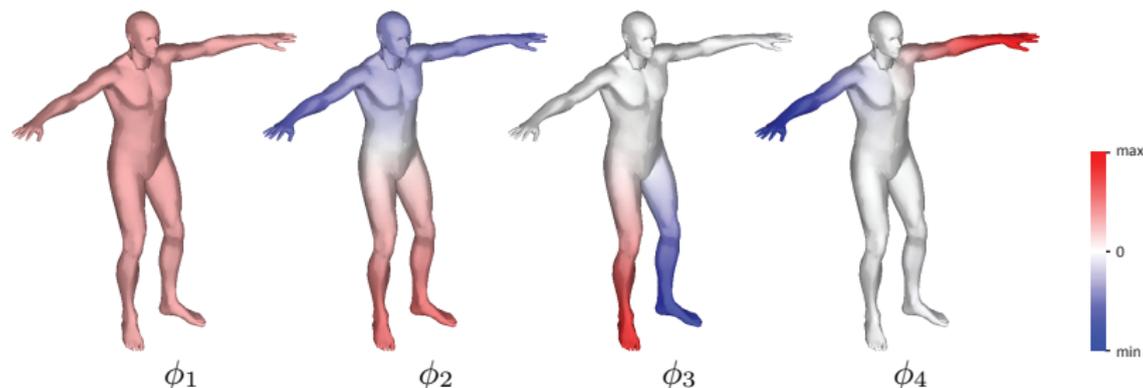
First eigenfunctions of a graph Laplacian

Laplacian eigenvectors and eigenvalues

Eigendecomposition of a graph Laplacian

$$\Delta = \Phi \Lambda \Phi^\top$$

where $\Phi = (\phi_1, \dots, \phi_n)$ are **orthogonal eigenvectors** ($\Phi^\top \Phi = \mathbf{I}$) and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ the corresponding **non-negative eigenvalues**



First eigenfunctions of a manifold Laplacian

Fourier analysis on Euclidean spaces

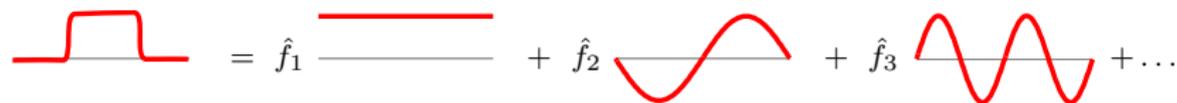
A function $f : [-\pi, \pi] \rightarrow \mathbb{R}$ can be written as a **Fourier series**

$$f(x) = \sum_{k \geq 0} \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x') e^{-ikx'} dx' e^{ikx}$$

Fourier analysis on Euclidean spaces

A function $f : [-\pi, \pi] \rightarrow \mathbb{R}$ can be written as a **Fourier series**

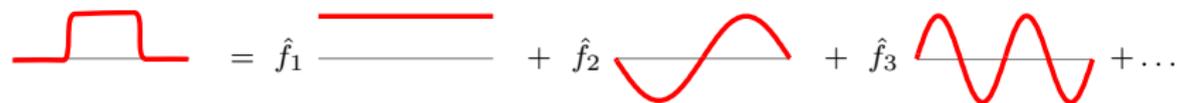
$$f(x) = \sum_{k \geq 0} \langle f, e^{ikx} \rangle_{L^2([-\pi, \pi])} e^{ikx}$$



Fourier analysis on Euclidean spaces

A function $f : [-\pi, \pi] \rightarrow \mathbb{R}$ can be written as a **Fourier series**

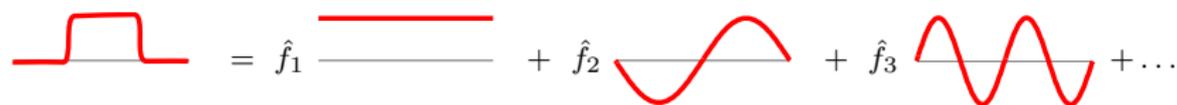
$$f(x) = \sum_{k \geq 0} \underbrace{\langle f, e^{ikx} \rangle_{L^2([-\pi, \pi])}}_{\hat{f}_k \text{ Fourier coefficient}} e^{ikx}$$



Fourier analysis on Euclidean spaces

A function $f : [-\pi, \pi] \rightarrow \mathbb{R}$ can be written as a **Fourier series**

$$f(x) = \sum_{k \geq 0} \underbrace{\langle f, e^{ikx} \rangle_{L^2([-\pi, \pi])}}_{\hat{f}_k \text{ Fourier coefficient}} e^{ikx}$$



Fourier basis = **Laplacian eigenfunctions**: $-\frac{d^2}{dx^2} e^{ikx} = k^2 e^{ikx}$

Fourier analysis on graphs and manifolds

A function $f : \mathcal{V} \rightarrow \mathbb{R}$ can be written as **Fourier series**

$$f = \sum_{k=1}^n \underbrace{\langle f, \phi_k \rangle_{L^2(\mathcal{V})}}_{\hat{f}_k} \phi_k$$

Fourier basis = **Laplacian eigenfunctions**: $\Delta \phi_k = \lambda_k \phi_k$

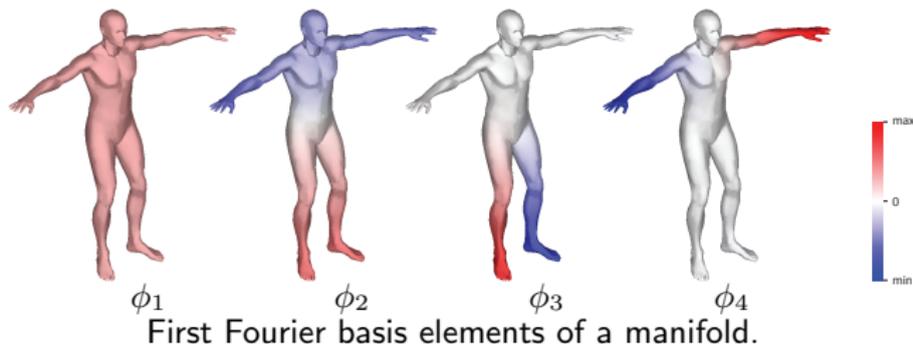
Fourier analysis on graphs and manifolds

A function $f : \mathcal{V} \rightarrow \mathbb{R}$ can be written as **Fourier series**

$$f = \sum_{k=1}^n \underbrace{\langle f, \phi_k \rangle_{L^2(\mathcal{V})}}_{\hat{f}_k} \phi_k$$

Fourier basis = **Laplacian eigenfunctions**: $\Delta \phi_k = \lambda_k \phi_k$

$\lambda_k =$ **frequency**



Summary

Manifolds and graphs are natural extensions of vector spaces

Summary

Manifolds and **graphs** are natural extensions of vector spaces

Lines can be generalized to **shortest paths** and **geodesics**

Summary

Manifolds and **graphs** are natural extensions of vector spaces

Lines can be generalized to **shortest paths** and **geodesics**

The **Laplacian** operator defines **smoothness** of a function

Summary

Manifolds and **graphs** are natural extensions of vector spaces

Lines can be generalized to **shortest paths** and **geodesics**

The **Laplacian** operator defines **smoothness** of a function

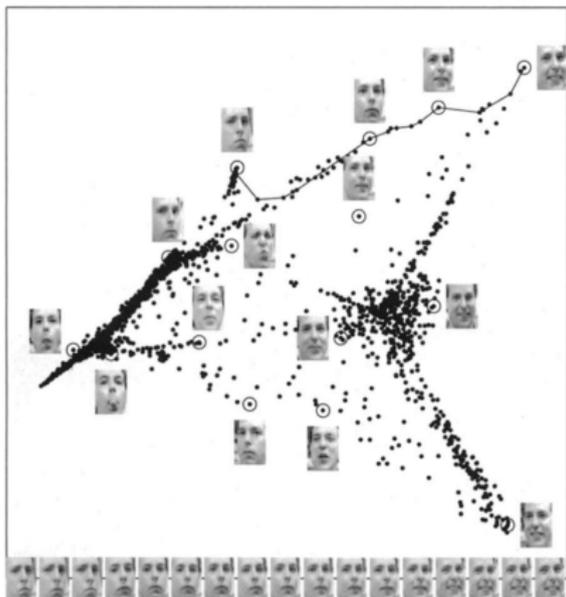
Spectral decompositions generalize **Fourier** analysis

Part II

The Geometry of Data

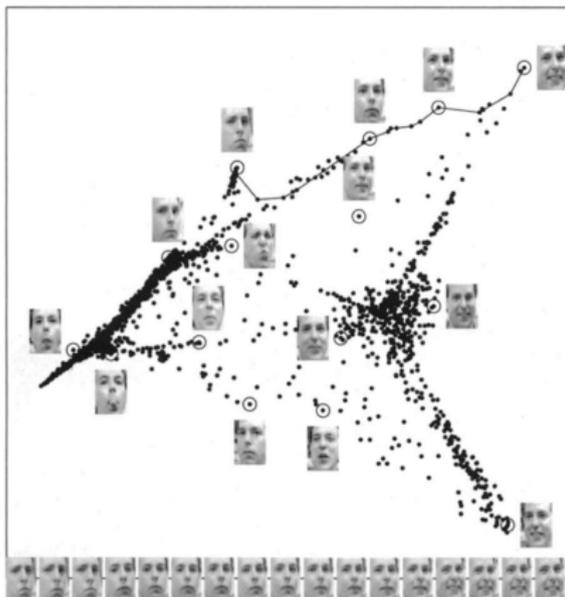
Classic Manifold Learning: A Time Before Deep Learning

Manifold Learning



Learn **nonlinear** embedding of data into low dimensional space that **preserves distances locally**

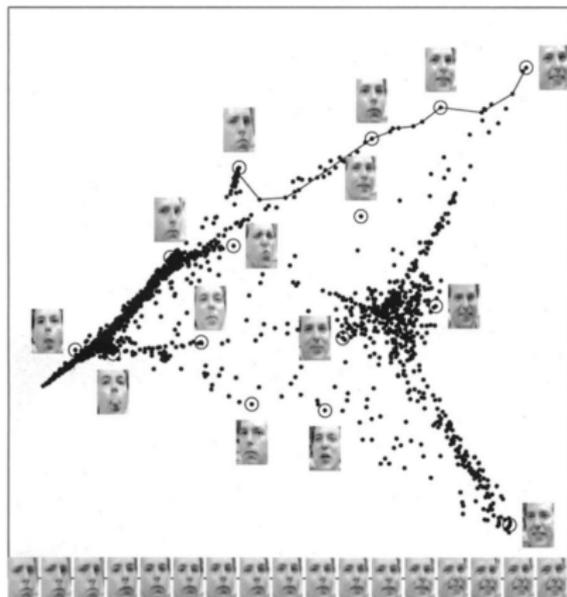
Manifold Learning



Learn **nonlinear** embedding of data into low dimensional space that **preserves distances locally**

The data itself does **not** have to be manifold structured

Manifold Learning



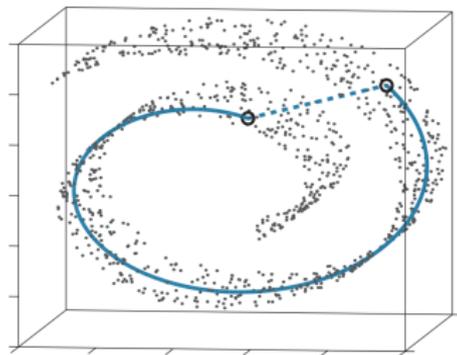
Learn **nonlinear** embedding of data into low dimensional space that **preserves distances locally**

The data itself does **not** have to be manifold structured

The **dataset** has some latent manifold structure

IsoMap

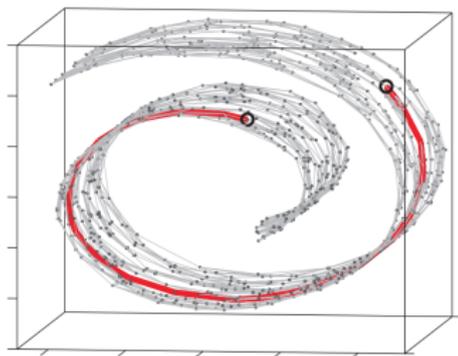
Compute **nearest neighbors** graph of dataset



IsoMap

Compute **nearest neighbors** graph of dataset

Construct matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ of **geodesic** distances between pairs of data by **Dijkstra's algorithm**

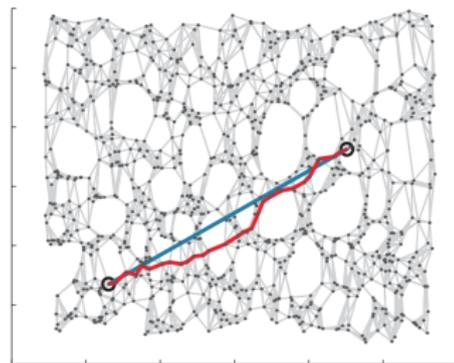


IsoMap

Compute **nearest neighbors** graph of dataset

Construct matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ of **geodesic** distances between pairs of data by **Dijkstra's algorithm**

Embed data in low dimensional space with **multidimensional scaling** (MDS):



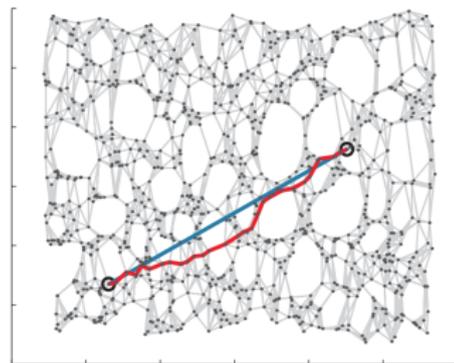
IsoMap

Compute **nearest neighbors** graph of dataset

Construct matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ of **geodesic** distances between pairs of data by **Dijkstra's algorithm**

Embed data in low dimensional space with **multidimensional scaling** (MDS):

$$\mathbf{B} = -\frac{1}{2} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) \mathbf{D} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right)$$



IsoMap

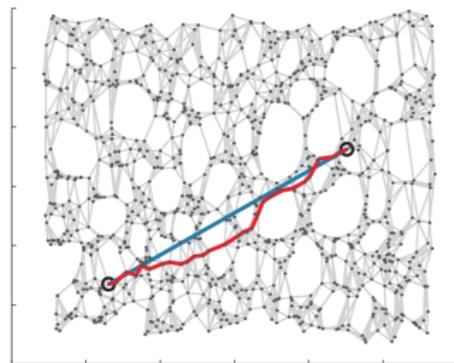
Compute **nearest neighbors** graph of dataset

Construct matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ of **geodesic** distances between pairs of data by **Dijkstra's algorithm**

Embed data in low dimensional space with **multidimensional scaling** (MDS):

$$\mathbf{B} = -\frac{1}{2} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) \mathbf{D} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right)$$

$\lambda_1, \dots, \lambda_k, \phi_1, \dots, \phi_k$ **top** eigenvalues and eigenvectors of \mathbf{B}



IsoMap

Compute **nearest neighbors** graph of dataset

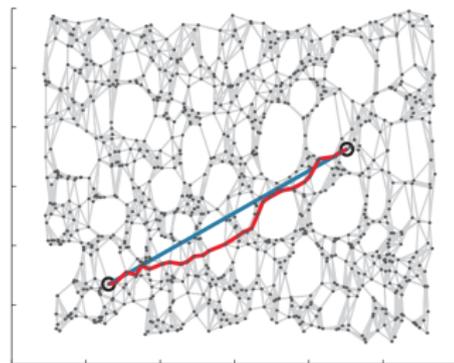
Construct matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ of **geodesic** distances between pairs of data by **Dijkstra's algorithm**

Embed data in low dimensional space with **multidimensional scaling** (MDS):

$$\mathbf{B} = -\frac{1}{2} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) \mathbf{D} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right)$$

$\lambda_1, \dots, \lambda_k, \phi_1, \dots, \phi_k$ **top** eigenvalues and eigenvectors of \mathbf{B}

Embedding: $(\sqrt{\lambda_1}\phi_1, \dots, \sqrt{\lambda_k}\phi_k)$



IsoMap

Compute **nearest neighbors** graph of dataset

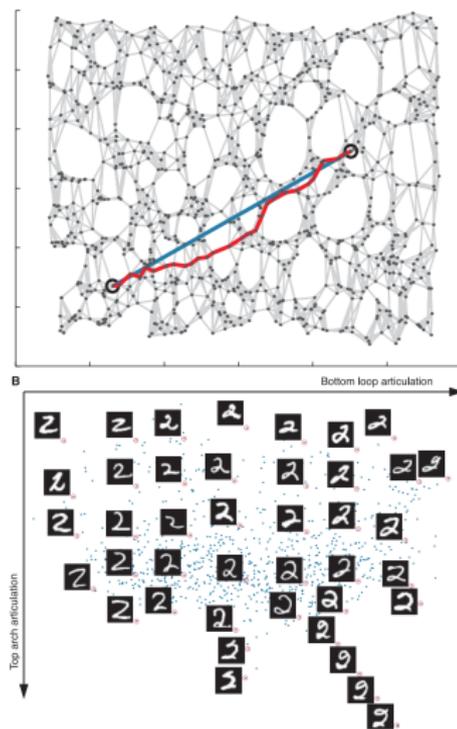
Construct matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ of **geodesic** distances between pairs of data by **Dijkstra's algorithm**

Embed data in low dimensional space with **multidimensional scaling** (MDS):

$$\mathbf{B} = -\frac{1}{2} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) \mathbf{D} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right)$$

$\lambda_1, \dots, \lambda_k, \phi_1, \dots, \phi_k$ **top** eigenvalues and eigenvectors of \mathbf{B}

Embedding: $(\sqrt{\lambda_1} \phi_1, \dots, \sqrt{\lambda_k} \phi_k)$



IsoMap

Compute **nearest neighbors** graph of dataset

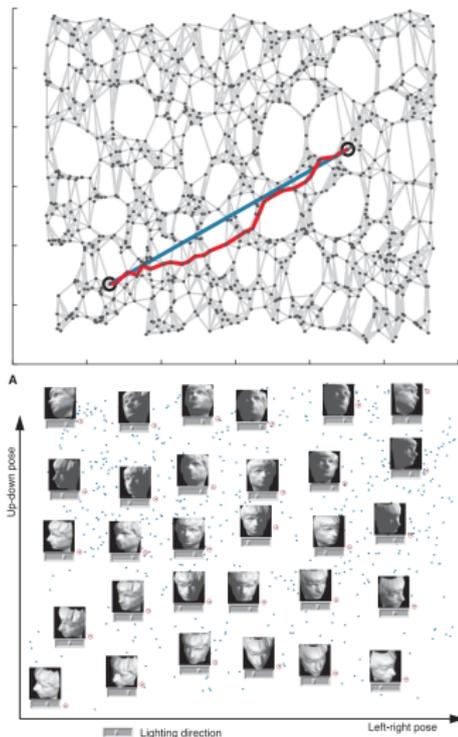
Construct matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ of **geodesic** distances between pairs of data by **Dijkstra's algorithm**

Embed data in low dimensional space with **multidimensional scaling** (MDS):

$$\mathbf{B} = -\frac{1}{2} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) \mathbf{D} \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right)$$

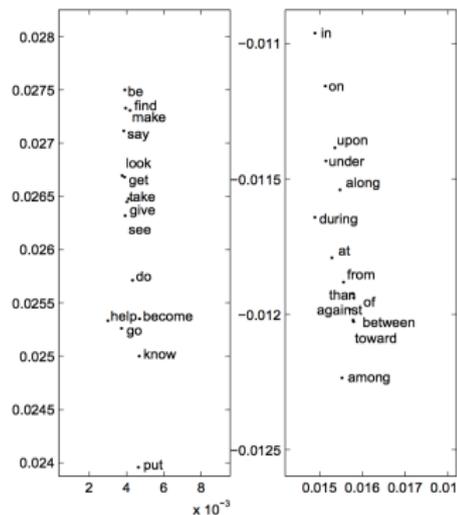
$\lambda_1, \dots, \lambda_k, \phi_1, \dots, \phi_k$ **top** eigenvalues and eigenvectors of \mathbf{B}

Embedding: $(\sqrt{\lambda_1}\phi_1, \dots, \sqrt{\lambda_k}\phi_k)$



Laplacian eigenmaps

Compute weighted nearest neighbors graph of dataset

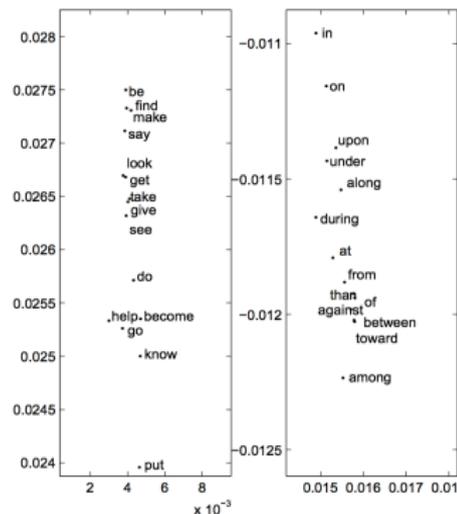


Laplacian eigenmaps

Compute weighted nearest neighbors graph of dataset

Construct graph Laplacian matrix

$$\Delta = D - W$$



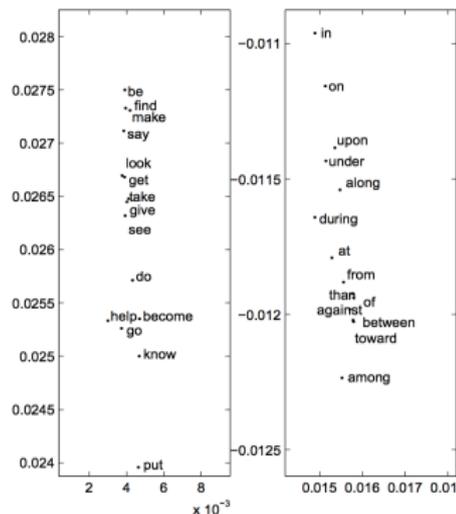
Laplacian eigenmaps

Compute weighted nearest neighbors graph of dataset

Construct graph Laplacian matrix

$$\Delta = \mathbf{D} - \mathbf{W}$$

$\lambda_1, \dots, \lambda_k, \phi_1, \dots, \phi_k$ bottom eigenvalues and eigenvectors of Δ



Laplacian eigenmaps

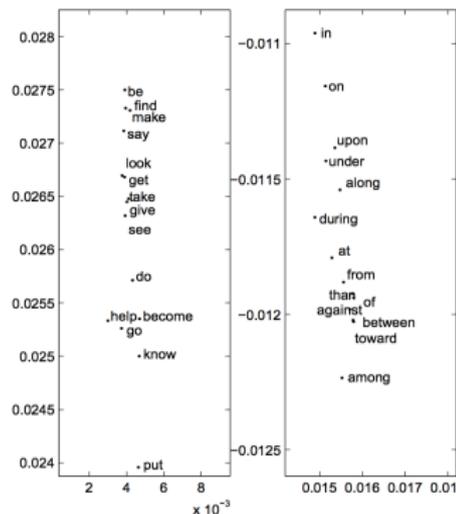
Compute weighted nearest neighbors graph of dataset

Construct graph Laplacian matrix

$$\Delta = \mathbf{D} - \mathbf{W}$$

$\lambda_1, \dots, \lambda_k, \phi_1, \dots, \phi_k$ bottom eigenvalues and eigenvectors of Δ

Embedding: (ϕ_2, \dots, ϕ_k)



Laplacian eigenmaps

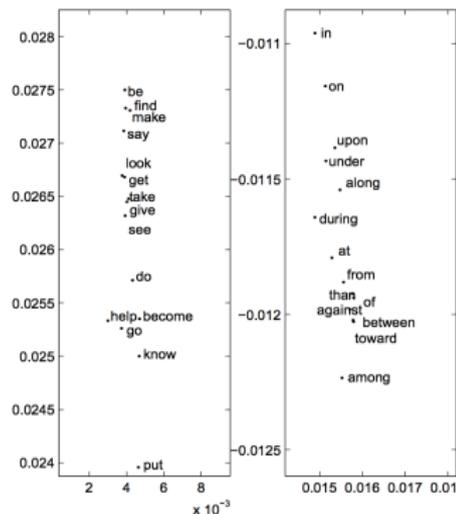
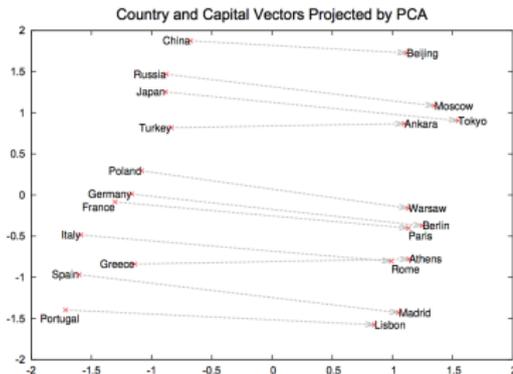
Compute weighted nearest neighbors graph of dataset

Construct graph Laplacian matrix

$$\Delta = \mathbf{D} - \mathbf{W}$$

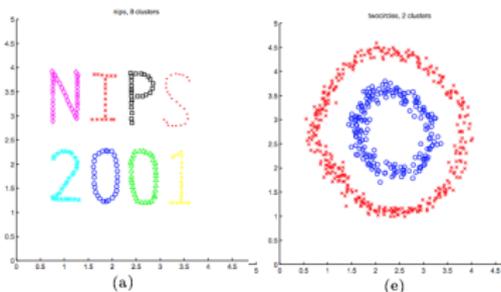
$\lambda_1, \dots, \lambda_k, \phi_1, \dots, \phi_k$ bottom eigenvalues and eigenvectors of Δ

Embedding: (ϕ_2, \dots, ϕ_k)

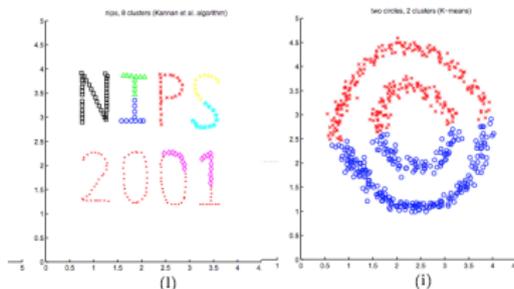


Spectral clustering

Same embedding as Laplacian eigenmaps, but use embedding vectors for **clustering** instead of **visualization**



Spectral Clustering

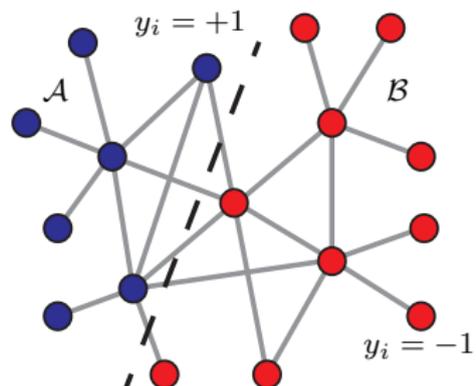


Other

Normalized cuts

Minimum (normalized) cut problem:
partition vertices into sets $\mathcal{A}, \mathcal{B} \subset \mathcal{V}$

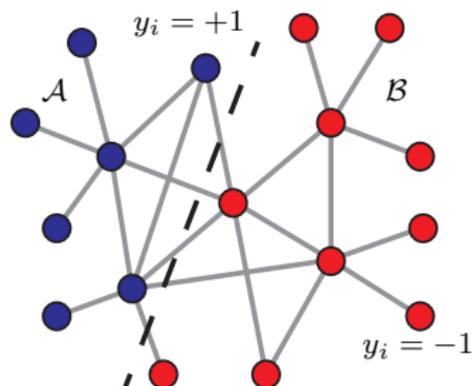
$$\text{NCut}(\mathbf{y}) = \frac{\sum_{i \in \mathcal{A}, j \in \mathcal{B}} w_{ij}}{\sum_{i \in \mathcal{A}, j \in \mathcal{V}} w_{ij}} + \frac{\sum_{i \in \mathcal{B}, j \in \mathcal{A}} w_{ij}}{\sum_{i \in \mathcal{B}, j \in \mathcal{V}} w_{ij}}$$



Normalized cuts

Minimum (normalized) cut problem:
partition vertices into sets $\mathcal{A}, \mathcal{B} \subset \mathcal{V}$

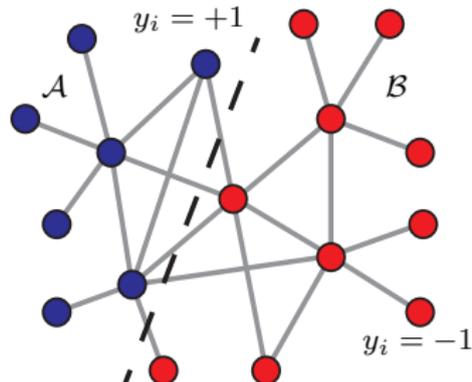
$$\begin{aligned} \text{NCut}(\mathbf{y}) &= \frac{\sum_{i \in \mathcal{A}, j \in \mathcal{B}} w_{ij}}{\sum_{i \in \mathcal{A}, j \in \mathcal{V}} w_{ij}} + \frac{\sum_{i \in \mathcal{B}, j \in \mathcal{A}} w_{ij}}{\sum_{i \in \mathcal{B}, j \in \mathcal{V}} w_{ij}} \\ &= \frac{\mathbf{y}^T (\mathbf{D} - \mathbf{W}) \mathbf{y}}{\mathbf{y}^T \mathbf{y}} \end{aligned}$$



Normalized cuts

Minimum (normalized) cut problem:
partition vertices into sets $\mathcal{A}, \mathcal{B} \subset \mathcal{V}$

$$\begin{aligned}\text{NCut}(\mathbf{y}) &= \frac{\sum_{i \in \mathcal{A}, j \in \mathcal{B}} w_{ij}}{\sum_{i \in \mathcal{A}, j \in \mathcal{V}} w_{ij}} + \frac{\sum_{i \in \mathcal{B}, j \in \mathcal{A}} w_{ij}}{\sum_{i \in \mathcal{B}, j \in \mathcal{V}} w_{ij}} \\ &= \frac{\mathbf{y}^T \Delta \mathbf{y}}{\mathbf{y}^T \mathbf{y}}\end{aligned}$$

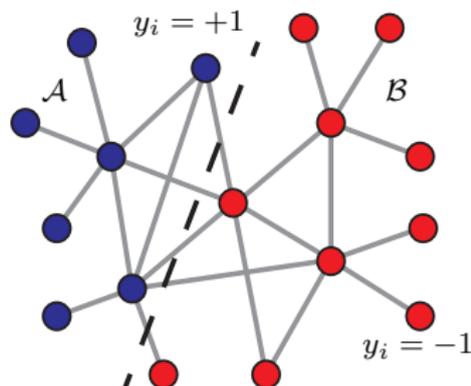


Normalized cuts

Minimum (normalized) cut problem:
partition vertices into sets $\mathcal{A}, \mathcal{B} \subset \mathcal{V}$

$$\begin{aligned} \text{NCut}(\mathbf{y}) &= \frac{\sum_{i \in \mathcal{A}, j \in \mathcal{B}} w_{ij}}{\sum_{i \in \mathcal{A}, j \in \mathcal{V}} w_{ij}} + \frac{\sum_{i \in \mathcal{B}, j \in \mathcal{A}} w_{ij}}{\sum_{i \in \mathcal{B}, j \in \mathcal{V}} w_{ij}} \\ &= \frac{\mathbf{y}^T \Delta \mathbf{y}}{\mathbf{y}^T \mathbf{y}} \end{aligned}$$

Relax optimization from **binary** vectors to **continuous** vectors. Solution is bottom eigenfunctions of graph Laplacian.



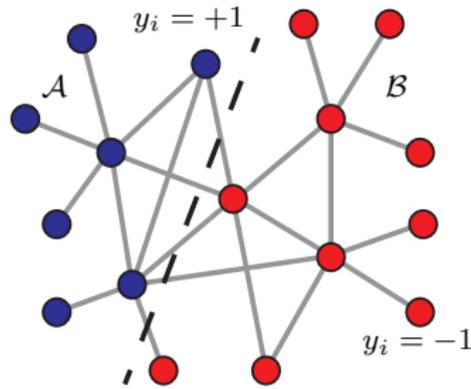
Normalized cuts

Minimum (normalized) cut problem:
partition vertices into sets $\mathcal{A}, \mathcal{B} \subset \mathcal{V}$

$$\begin{aligned} \text{NCut}(\mathbf{y}) &= \frac{\sum_{i \in \mathcal{A}, j \in \mathcal{B}} w_{ij}}{\sum_{i \in \mathcal{A}, j \in \mathcal{V}} w_{ij}} + \frac{\sum_{i \in \mathcal{B}, j \in \mathcal{A}} w_{ij}}{\sum_{i \in \mathcal{B}, j \in \mathcal{V}} w_{ij}} \\ &= \frac{\mathbf{y}^T \Delta \mathbf{y}}{\mathbf{y}^T \mathbf{y}} \end{aligned}$$

Relax optimization from **binary** vectors to **continuous** vectors. Solution is bottom eigenfunctions of graph Laplacian.

Used in **image segmentation** where each vertex is a pixel



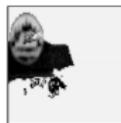
(a)



(b)



(c)



(e)



(f)

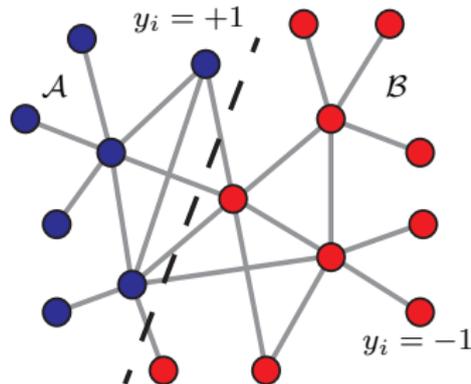


(g)

Normalized cuts

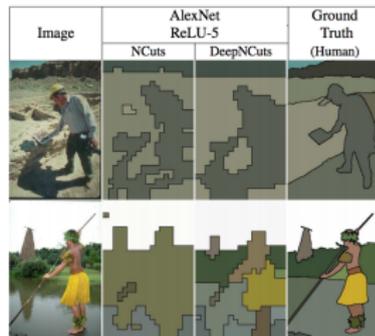
Minimum (normalized) cut problem:
partition vertices into sets $\mathcal{A}, \mathcal{B} \subset \mathcal{V}$

$$\begin{aligned} \text{NCut}(\mathbf{y}) &= \frac{\sum_{i \in \mathcal{A}, j \in \mathcal{B}} w_{ij}}{\sum_{i \in \mathcal{A}, j \in \mathcal{V}} w_{ij}} + \frac{\sum_{i \in \mathcal{B}, j \in \mathcal{A}} w_{ij}}{\sum_{i \in \mathcal{B}, j \in \mathcal{V}} w_{ij}} \\ &= \frac{\mathbf{y}^T \Delta \mathbf{y}}{\mathbf{y}^T \mathbf{y}} \end{aligned}$$



Relax optimization from **binary** vectors to **continuous** vectors. Solution is bottom eigenfunctions of graph Laplacian.

Can be incorporated as layer **inside** a deep network by backpropagating through eigendecomposition



Inference on held-out data

A common pattern:

Inference on held-out data

A common pattern:

Construct **Gram matrix** from kernel $k(\mathbf{x}, \mathbf{x}')$

$$\mathbf{M} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

Inference on held-out data

A common pattern:

Construct **Gram matrix** from kernel $k(\mathbf{x}, \mathbf{x}')$

$$\mathbf{M} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

Use top/bottom eigenvectors of Gram matrix as embedding

Inference on held-out data

A common pattern:

Construct **Gram matrix** from kernel $k(\mathbf{x}, \mathbf{x}')$

$$\mathbf{M} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

Use top/bottom eigenvectors of Gram matrix as embedding

Works for **fixed** dataset $\mathbf{x}_1, \dots, \mathbf{x}_n$ -

but what is the embedding vector for a **new** data point \mathbf{x}' ?

Inference on held-out data

Assume data drawn $\mathbf{x}_1, \dots, \mathbf{x}_n \sim p(\mathbf{x})$. Gram matrix is approximation to **linear operator**:

$$\begin{aligned} \frac{1}{n} \mathbf{M} \begin{pmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \end{pmatrix}_i &= \frac{1}{n} \sum_j k(\mathbf{x}_i, \mathbf{x}_j) f(\mathbf{x}_j) \approx \mathbb{E}_{p(\mathbf{x})} [k(\mathbf{x}_i, \mathbf{x}) f(\mathbf{x})] \\ &= \mathcal{K}_p[f](\mathbf{x}_i) \end{aligned}$$

Inference on held-out data

Assume data drawn $\mathbf{x}_1, \dots, \mathbf{x}_n \sim p(\mathbf{x})$. Gram matrix is approximation to **linear operator**:

$$\begin{aligned} \frac{1}{n} \mathbf{M} \begin{pmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \end{pmatrix}_i &= \frac{1}{n} \sum_j k(\mathbf{x}_i, \mathbf{x}_j) f(\mathbf{x}_j) \approx \mathbb{E}_{p(\mathbf{x})} [k(\mathbf{x}_i, \mathbf{x}) f(\mathbf{x})] \\ &= \mathcal{K}_p[f](\mathbf{x}_i) \end{aligned}$$

Eigenvectors ϕ_k of \mathbf{M} are approximation to **eigenfunctions** $\phi_k(\cdot)$ of linear operator \mathcal{K}_p

Inference on held-out data

Assume data drawn $\mathbf{x}_1, \dots, \mathbf{x}_n \sim p(\mathbf{x})$. Gram matrix is approximation to **linear operator**:

$$\begin{aligned} \frac{1}{n} \mathbf{M} \begin{pmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \end{pmatrix}_i &= \frac{1}{n} \sum_j k(\mathbf{x}_i, \mathbf{x}_j) f(\mathbf{x}_j) \approx \mathbb{E}_{p(\mathbf{x})}[k(\mathbf{x}_i, \mathbf{x}) f(\mathbf{x})] \\ &= \mathcal{K}_p[f](\mathbf{x}_i) \end{aligned}$$

From eigenvalues $\lambda_1, \dots, \lambda_k$ and eigenvectors ϕ_1, \dots, ϕ_k of \mathbf{M} , can approximate eigenfunction of \mathcal{K}_p with **Nyström method**:

$$\phi_k(\mathbf{x}') \propto \sum_i \phi_{ki} k(\mathbf{x}_i, \mathbf{x}')$$

Trade-offs of manifold learning

😊 Exactly solvable by eigendecomposition

Trade-offs of manifold learning

- 😊 Exactly solvable by eigendecomposition
- 😊 Data efficient (works with $\mathcal{O}(1000)$ data points)

Trade-offs of manifold learning

- ☺ Exactly solvable by eigendecomposition
- ☺ Data efficient (works with $\mathcal{O}(1000)$ data points)
- ☹ Unsupervised learning without a generative model

Trade-offs of manifold learning

- 😊 Exactly solvable by eigendecomposition
- 😊 Data efficient (works with $\mathcal{O}(1000)$ data points)
- 😐 Unsupervised learning without a generative model
- 😞 Learning scales as $\mathcal{O}(n \log n)$ with size of training data

Trade-offs of manifold learning

- ☺ Exactly solvable by eigendecomposition
- ☺ Data efficient (works with $\mathcal{O}(1000)$ data points)
- ☹ Unsupervised learning without a generative model
- ☹ Learning scales as $\mathcal{O}(n \log n)$ with size of training data
- ☹ Inference scales as $\mathcal{O}(n)$ with size of training data

Trade-offs of manifold learning

- 😊 Exactly solvable by eigendecomposition
- 😊 Data efficient (works with $\mathcal{O}(1000)$ data points)
- 😐 Unsupervised learning without a generative model
- 😞 Learning scales as $\mathcal{O}(n \log n)$ with size of training data
- 😞 Inference scales as $\mathcal{O}(n)$ with size of training data
- 😞 Performance degrades for noisy or clustered data

Trade-offs of manifold learning

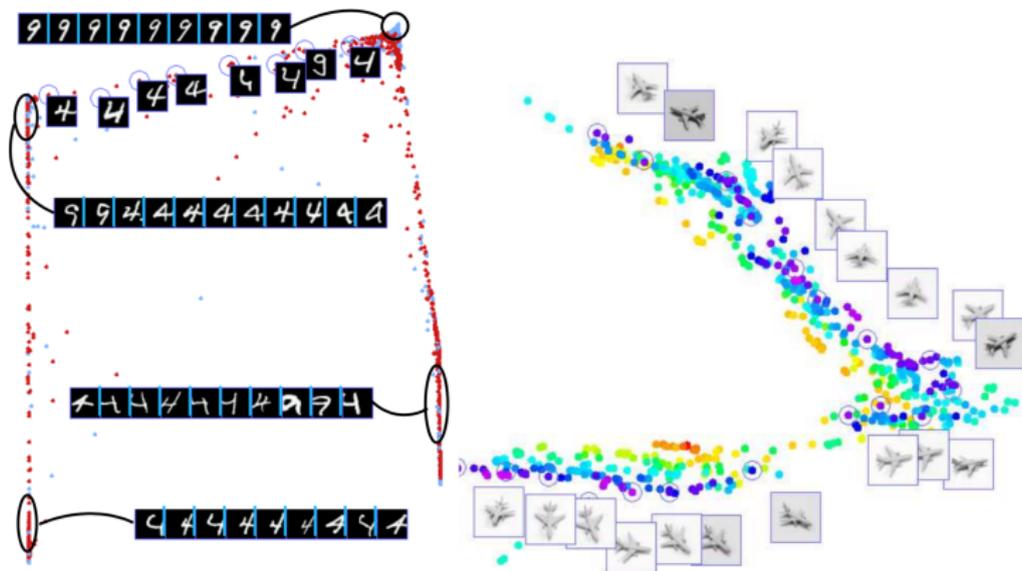
- 😊 Exactly solvable by eigendecomposition
- 😊 Data efficient (works with $\mathcal{O}(1000)$ data points)
- 😐 Unsupervised learning without a generative model
- 😞 Learning scales as $\mathcal{O}(n \log n)$ with size of training data
- 😞 Inference scales as $\mathcal{O}(n)$ with size of training data
- 😞 Performance degrades for noisy or clustered data
- 😞 Embeddings collapse on more complex data

Trade-offs of manifold learning

- ☺ Exactly solvable by eigendecomposition
- ☺ Data efficient (works with $\mathcal{O}(1000)$ data points)
- ☹ Unsupervised learning without a generative model
- ☹ Learning scales as $\mathcal{O}(n \log n)$ with size of training data
- ☹ Inference scales as $\mathcal{O}(n)$ with size of training data
- ☹ Performance degrades for noisy or clustered data
- ☹ Embeddings collapse on more complex data

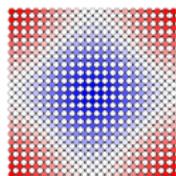
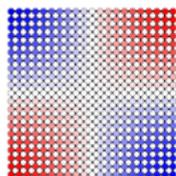
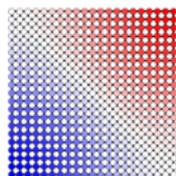
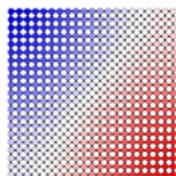
Other methods have become popular for low-dimensional visualization - especially **t-SNE** (Maaten and Hinton 2008)

Collapsed embeddings



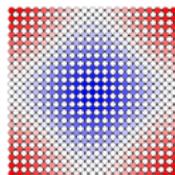
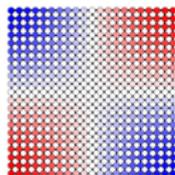
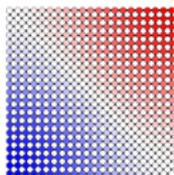
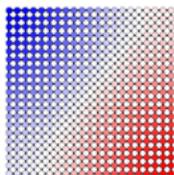
Improving spectral embeddings

Bottom eigenfunctions of graph Laplacian of **evenly-scaled** grid are natural coordinates

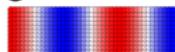


Improving spectral embeddings

Bottom eigenfunctions of graph Laplacian of **evenly-scaled** grid are natural coordinates

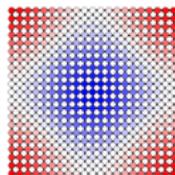
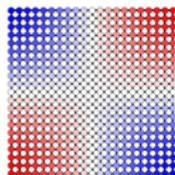
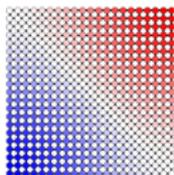
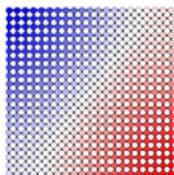


Bottom eigenfunctions of **unevenly-scaled** grid are not

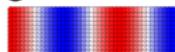


Improving spectral embeddings

Bottom eigenfunctions of graph Laplacian of **evenly-scaled** grid are natural coordinates



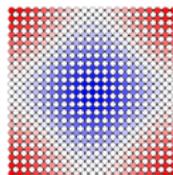
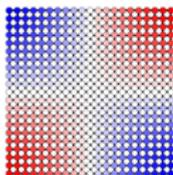
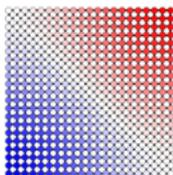
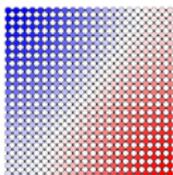
Bottom eigenfunctions of **unevenly-scaled** grid are not



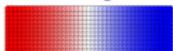
Eigenfunctions must be orthogonal, but can still be **predictable**, e.g. $\sin(2x)$ and $\sin(x)$ and $\cos(x)$.

Improving spectral embeddings

Bottom eigenfunctions of graph Laplacian of **evenly-scaled** grid are natural coordinates



Bottom eigenfunctions of **unevenly-scaled** grid are not

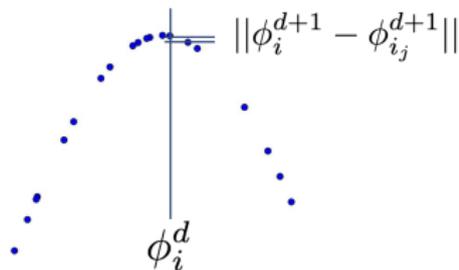


Eigenfunctions must be orthogonal, but can still be **predictable**, e.g. $\sin(2x)$ and $\sin(x)$ and $\cos(x)$.

Instead of using **lowest** eigenfunctions as embedding, use lowest eigenfunctions that are **unpredictable** from lower eigenfunctions

Minimally redundant Laplacian eigenmaps

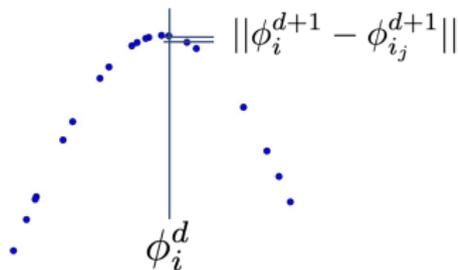
After adding eigenfunctions ϕ^1, \dots, ϕ^d to embedding, evaluate ϕ^{d+1} as candidate.



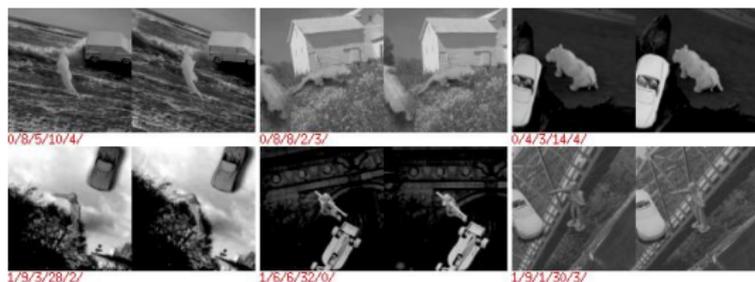
Minimally redundant Laplacian eigenmaps

After adding eigenfunctions ϕ^1, \dots, ϕ^d to embedding, evaluate ϕ^{d+1} as candidate.

If value of ϕ^{d+1} at point i can be predicted from nearest neighbors of i in ϕ^1, \dots, ϕ^d , eigenfunction is too predictable

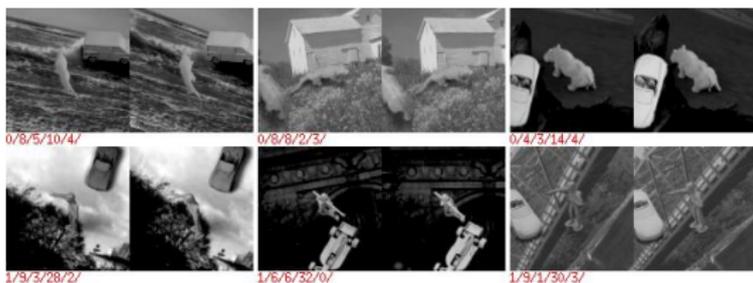


Minimally redundant Laplacian eigenmaps



NORB: Model dataset for studying invariance in object recognition

Minimally redundant Laplacian eigenmaps

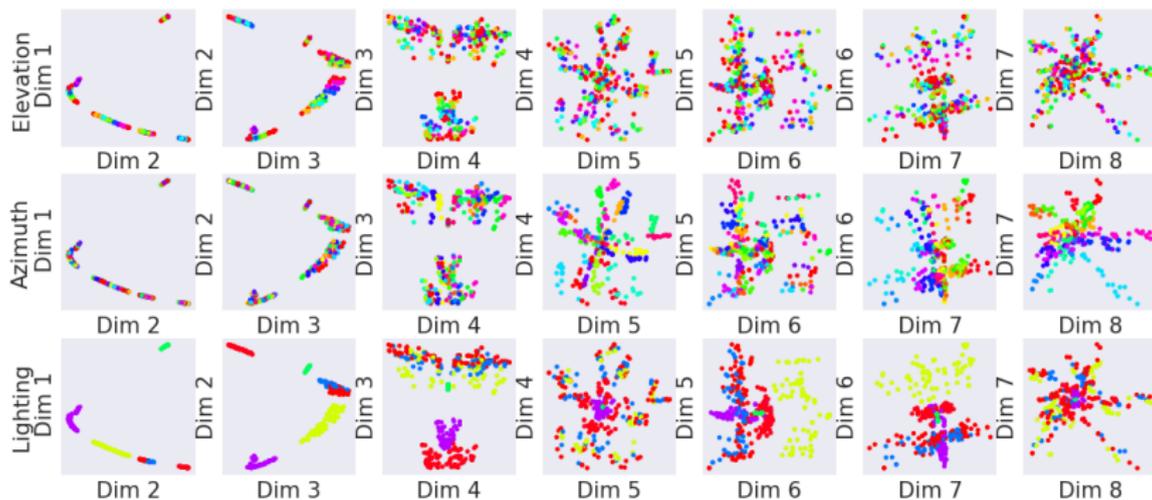


NORB: Model dataset for studying invariance in object recognition

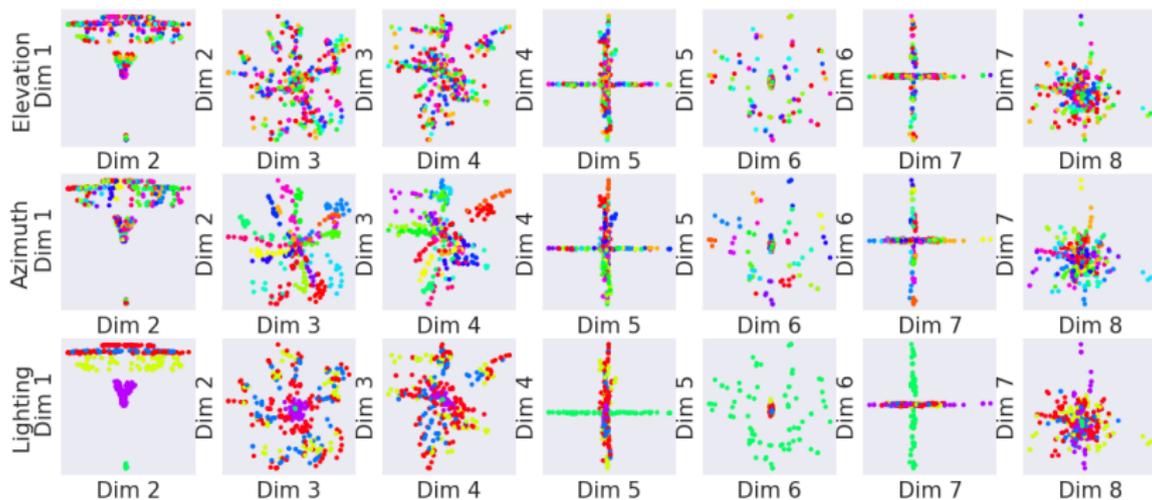


Consider a single object under different lighting and rotation

Minimally redundant Laplacian eigenmaps

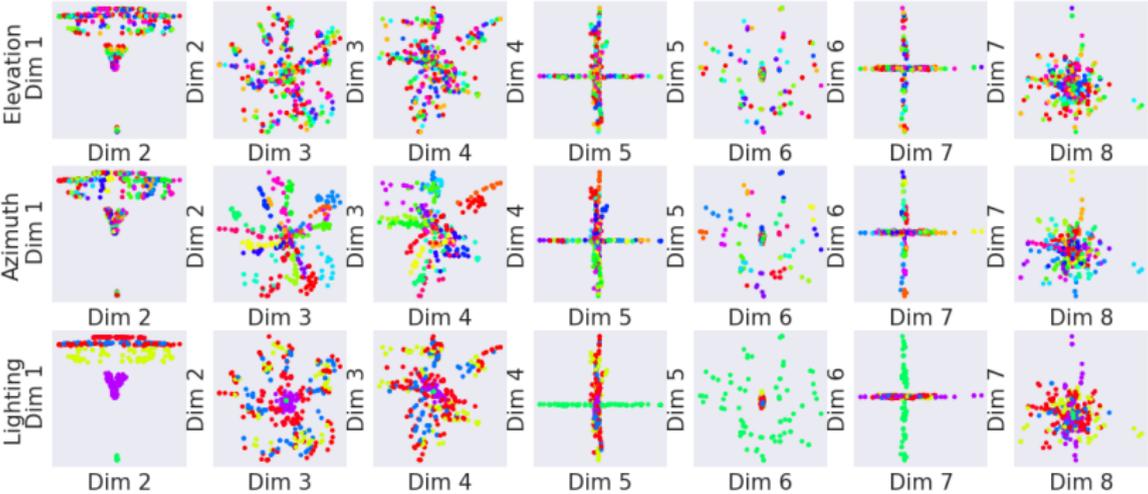


Minimally redundant Laplacian eigenmaps



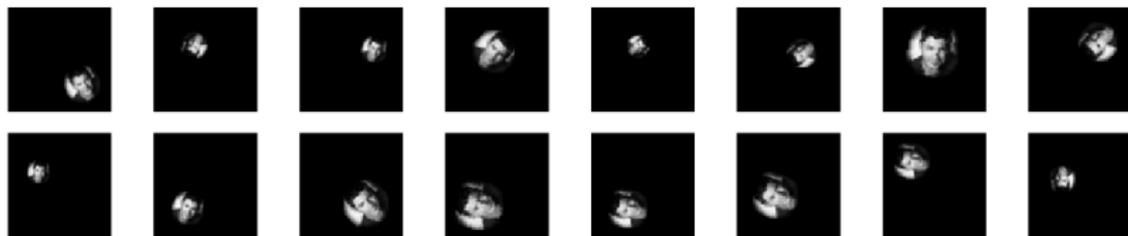
With filtering by redundancy, all variation captured by 6 eigenfunctions

Minimally redundant Laplacian eigenmaps

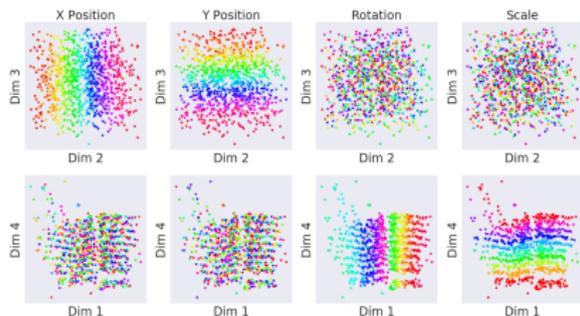
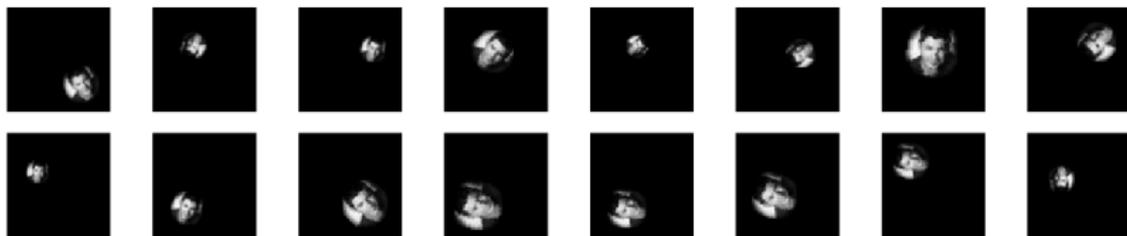


With filtering by redundancy, all variation captured by 6 eigenfunctions
Works with **less than 1000** points!

Minimally redundant Laplacian eigenmaps



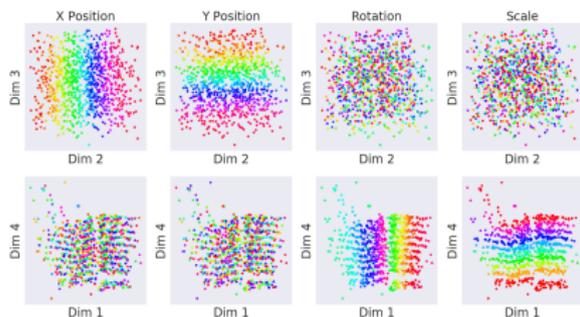
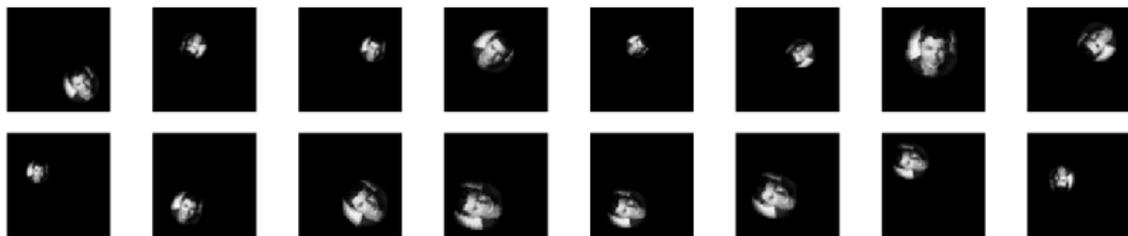
Minimally redundant Laplacian eigenmaps



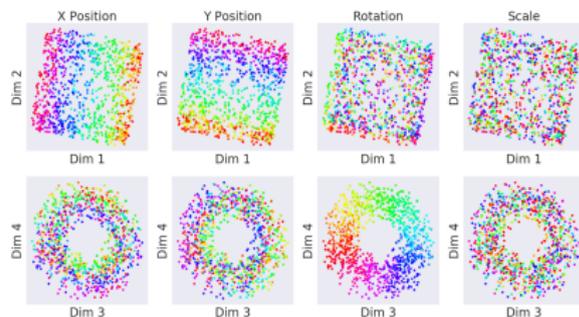
β -VAE

Disentangling VAEs learn the **dimension** but not the topology

Minimally redundant Laplacian eigenmaps



β -VAE



Laplacian eigenmaps

Disentangling VAEs learn the **dimension** but not the topology

Laplacian eigenmaps learns the **topology** but misses some dimensions

Trade-offs of manifold learning

- 😊 Exactly solvable by eigendecomposition
- 😊 Data efficient (works with $\mathcal{O}(1000)$ data points)
- 😐 Unsupervised learning without a generative model
- 😞 Learning scales as $\mathcal{O}(n \log n)$ with size of training data
- 😞 Inference scales as $\mathcal{O}(n)$ with size of training data
- 😞 Performance degrades for noisy or clustered data
- 😞 Embeddings collapse on more complex data

Trade-offs of manifold learning

- 😊 Exactly solvable by eigendecomposition
- 😊 Data efficient (works with $\mathcal{O}(1000)$ data points)
- 😐 Unsupervised learning without a generative model
- 😞 Learning scales as $\mathcal{O}(n \log n)$ with size of training data
- 😞 Inference scales as $\mathcal{O}(n)$ with size of training data
- 😞 Performance degrades for noisy or clustered data
- 😊 Collapsed embeddings can be fixed by choice of eigenvector

Trade-offs of manifold learning

- 😊 Exactly solvable by eigendecomposition
- 😊 Data efficient (works with $\mathcal{O}(1000)$ data points)
- 😐 Unsupervised learning without a generative model
- 😞 Learning scales as $\mathcal{O}(n \log n)$ with size of training data
- 😞 Inference scales as $\mathcal{O}(n)$ with size of training data
- 😞 Performance degrades for noisy or clustered data
- 😊 Collapsed embeddings can be fixed by choice of eigenvector
- 😊 Can discover topology of data without prior assumptions

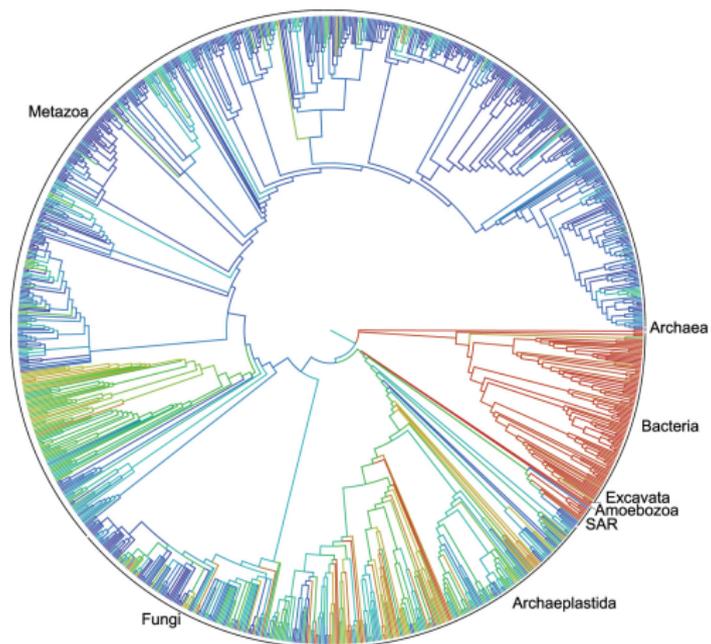
Embedding Hierarchies in Hyperbolic Spaces

Hierarchical data



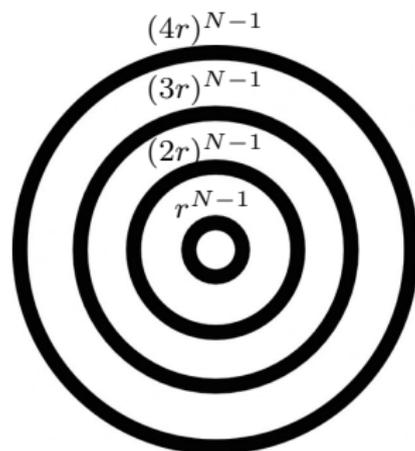
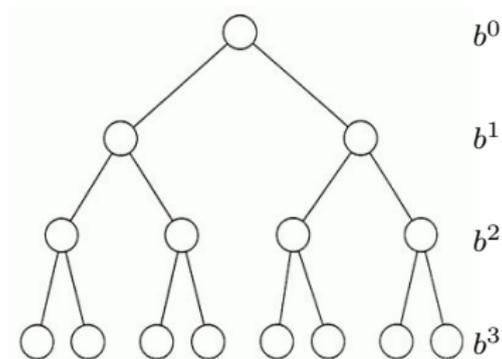
Some data can be embedded uniformly in **flat** space

Hierarchical data



What about data with **hierarchical** structure?

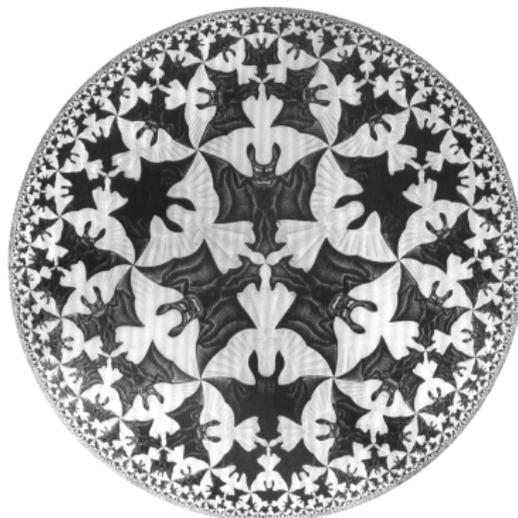
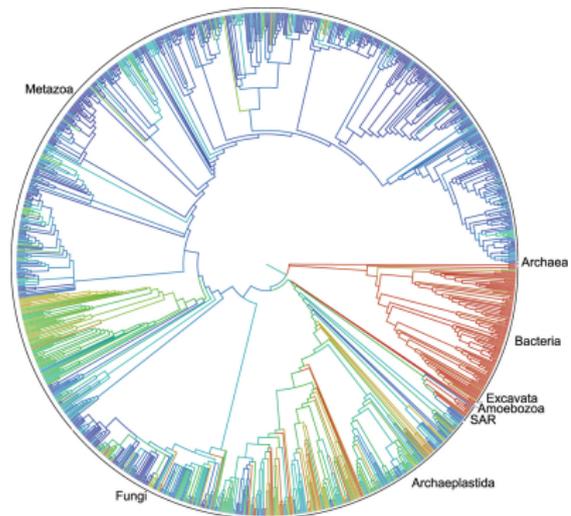
Hierarchical data



Tree with branching factor b has b^ℓ nodes at layer ℓ - exponential growth

Area of sphere in \mathbb{R}^N with radius r grows as r^{N-1} - polynomial growth

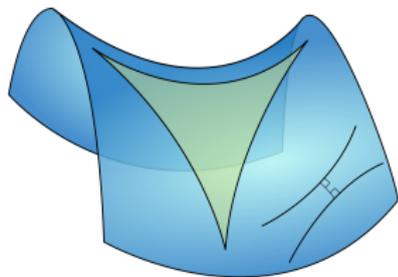
Hierarchical data



Idea: embed nodes in hierarchy in **hyperbolic** space instead of flat space

Hyperbolic geometry

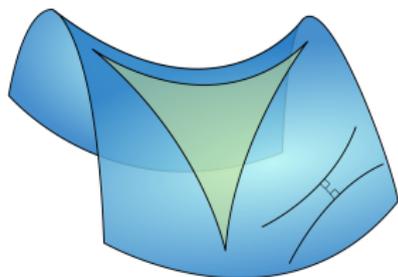
Hyperbolic space: manifold with constant
negative curvature



Hyperbolic geometry

Hyperbolic space: manifold with constant **negative** curvature

Angles of triangle add to **less than** 180 degrees

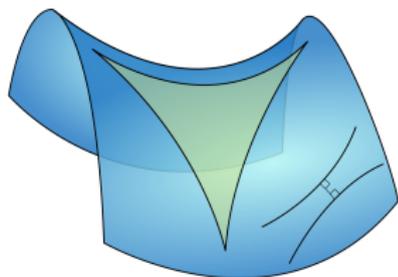


Hyperbolic geometry

Hyperbolic space: manifold with constant **negative** curvature

Angles of triangle add to **less than** 180 degrees

Surface area of spheres grows **exponentially!**



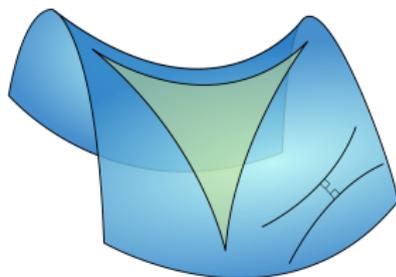
Hyperbolic geometry

Hyperbolic space: manifold with constant **negative** curvature

Angles of triangle add to **less than** 180 degrees

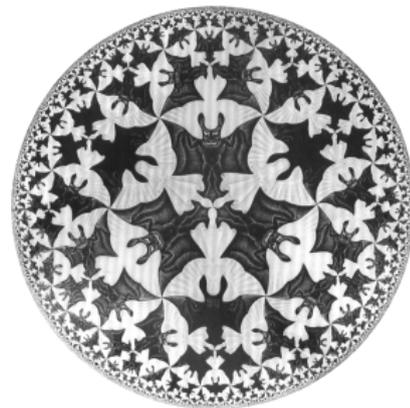
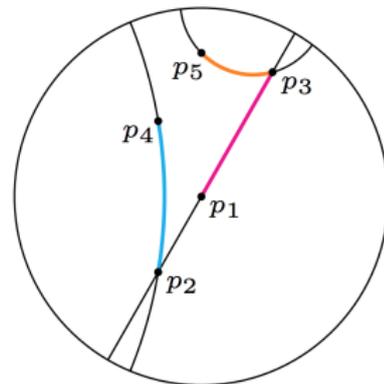
Surface area of spheres grows **exponentially!**

Many possible **models** of hyperbolic space



The Poincaré ball

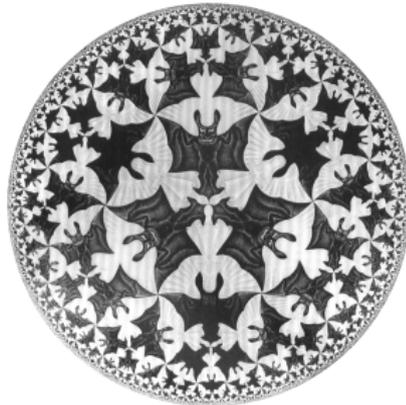
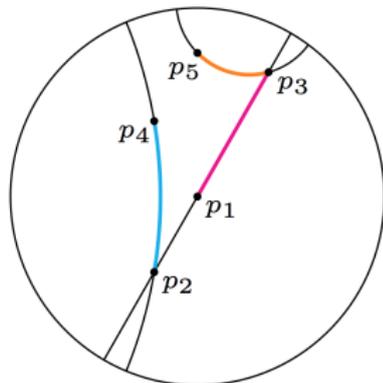
Maps hyperbolic space to open ball \mathbb{B}^N



The Poincaré ball

Maps hyperbolic space to open ball \mathbb{B}^N

Metric: $\langle \mathbf{u}, \mathbf{v} \rangle_{T_x \mathcal{P}} = \left(\frac{2}{1 - \mathbf{x}^T \mathbf{x}} \right)^2 \mathbf{u}^T \mathbf{v}$

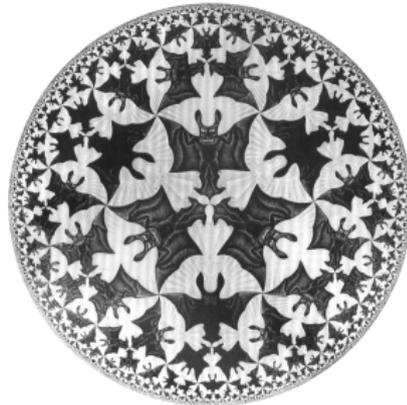
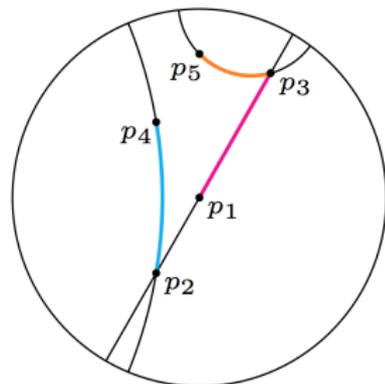


The Poincaré ball

Maps hyperbolic space to open ball \mathbb{B}^N

Metric: $\langle \mathbf{u}, \mathbf{v} \rangle_{T_x \mathcal{P}} = \left(\frac{2}{1 - \mathbf{x}^T \mathbf{x}} \right)^2 \mathbf{u}^T \mathbf{v}$

Geodesics: Circles that are orthogonal to boundary of ball



The Poincaré ball

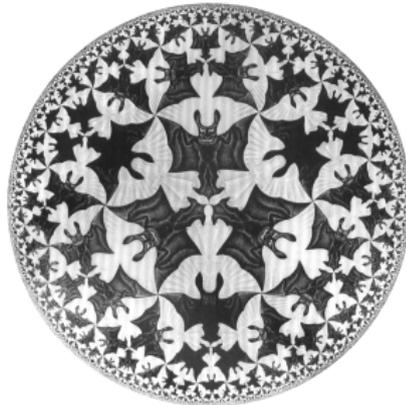
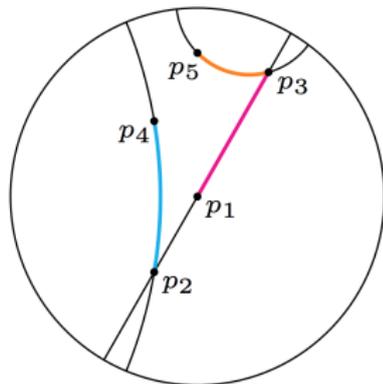
Maps hyperbolic space to open ball \mathbb{B}^N

Metric: $\langle \mathbf{u}, \mathbf{v} \rangle_{T_x \mathcal{P}} = \left(\frac{2}{1 - \mathbf{x}^T \mathbf{x}} \right)^2 \mathbf{u}^T \mathbf{v}$

Geodesics: Circles that are orthogonal to boundary of ball

Distances:

$$d(\mathbf{u}, \mathbf{v}) = \operatorname{arcosh} \left(1 + 2 \frac{\|\mathbf{u} - \mathbf{v}\|^2}{(1 - \|\mathbf{u}\|^2)(1 - \|\mathbf{v}\|^2)} \right)$$



Poincaré embeddings for learning hierarchical representation

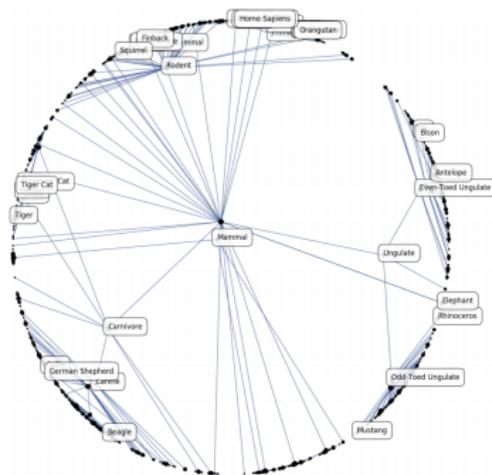
Given edges \mathcal{E} from a graph, find an embedding \mathbf{u}_i for vertex i that minimizes:

$$\sum_{i,j \in \mathcal{E}} \log \frac{e^{-d(\mathbf{u}_i, \mathbf{u}_j)}}{\sum_{j' \text{ st } ij' \notin \mathcal{E}} e^{-d(\mathbf{u}_i, \mathbf{u}_{j'})}}$$

Poincaré embeddings for learning hierarchical representation

Given edges \mathcal{E} from a graph, find an embedding \mathbf{u}_i for vertex i that minimizes:

$$\sum_{i,j \in \mathcal{E}} \log \frac{e^{-d(\mathbf{u}_i, \mathbf{u}_j)}}{\sum_{j' \text{ st } ij' \notin \mathcal{E}} e^{-d(\mathbf{u}_i, \mathbf{u}_{j'})}}$$



Poincaré embeddings for learning hierarchical representation

Given edges \mathcal{E} from a graph, find an embedding \mathbf{u}_i for vertex i that minimizes:

$$\sum_{i,j \in \mathcal{E}} \log \frac{e^{-d(\mathbf{u}_i, \mathbf{u}_j)}}{\sum_{j' \text{ st } ij' \notin \mathcal{E}} e^{-d(\mathbf{u}_i, \mathbf{u}_{j'})}}$$

		Dimensionality						
		5	10	20	50	100	200	
WORDNET Reconstruction	Euclidean	Rank	3542.3	2286.9	1685.9	1281.7	1187.3	1157.3
		MAP	0.024	0.059	0.087	0.140	0.162	0.168
	Translational	Rank	205.9	179.4	95.3	92.8	92.7	91.0
		MAP	0.517	0.503	0.563	0.566	0.562	0.565
	Poincaré	Rank	4.9	4.02	3.84	3.98	3.9	3.83
		MAP	0.823	0.851	0.855	0.86	0.857	0.87
WORDNET Link Pred.	Euclidean	Rank	3311.1	2199.5	952.3	351.4	190.7	81.5
		MAP	0.024	0.059	0.176	0.286	0.428	0.490
	Translational	Rank	65.7	56.6	52.1	47.2	43.2	40.4
		MAP	0.545	0.554	0.554	0.56	0.562	0.559
	Poincaré	Rank	5.7	4.3	4.9	4.6	4.6	4.6
		MAP	0.825	0.852	0.861	0.863	0.856	0.855

State of the art results on [link reconstruction](#) and [link prediction](#) on WORDNET noun dataset

Poincaré embeddings for learning hierarchical representation

Given edges \mathcal{E} from a graph, find an embedding \mathbf{u}_i for vertex i that minimizes:

$$\sum_{i,j \in \mathcal{E}} \log \frac{e^{-d(\mathbf{u}_i, \mathbf{u}_j)}}{\sum_{j' \text{ st } ij' \notin \mathcal{E}} e^{-d(\mathbf{u}_i, \mathbf{u}_{j'})}}$$

Table 3: Spearman's ρ for Lexical Entailment on HYPERLEX.

	FR	SLQS-Sim	WN-Basic	WN-WuP	WN-LCh	Vis-ID	Euclidean	Poincaré
ρ	0.283	0.229	0.240	0.214	0.214	0.253	0.389	0.512

State of the art results on [graded lexical entailment](#) on HYPERLEX

Trade-offs of Poincaré model

😊 Geodesics are simple

Trade-offs of Poincaré model

😊 Geodesics are simple

😊 Metric is diagonal

Trade-offs of Poincaré model

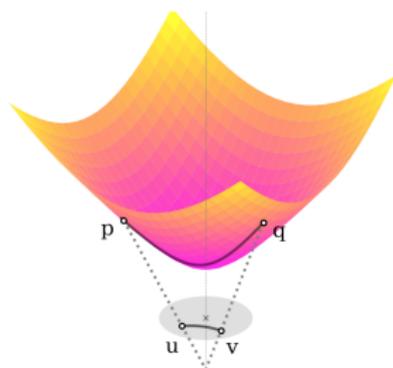
😊 Geodesics are simple

😊 Metric is diagonal

😞 Gradients of metric are unstable near boundary

The Lorentz model

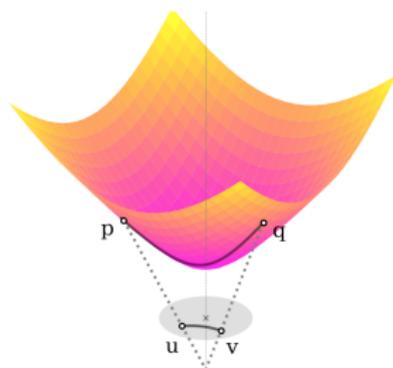
Maps N -dim hyperbolic space to surface in \mathbb{R}^{N+1}



The Lorentz model

Maps N -dim hyperbolic space to surface in \mathbb{R}^{N+1}

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}} = -u_0 v_0 + \sum_{i=1}^{N+1} u_i v_i$$

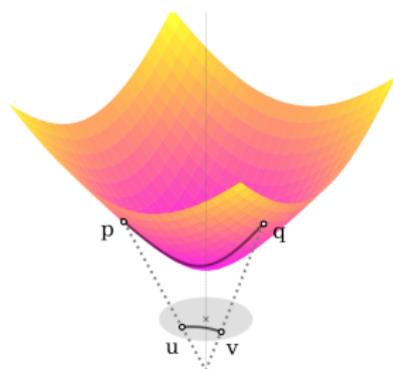


The Lorentz model

Maps N -dim hyperbolic space to surface in \mathbb{R}^{N+1}

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}} = -u_0 v_0 + \sum_{i=1}^{N+1} u_i v_i$$

$$\mathcal{H}^N = \{ \mathbf{x} \in \mathbb{R}^{N+1} : \langle \mathbf{x}, \mathbf{x} \rangle_{\mathcal{L}} = -1, x_0 > 0 \}$$



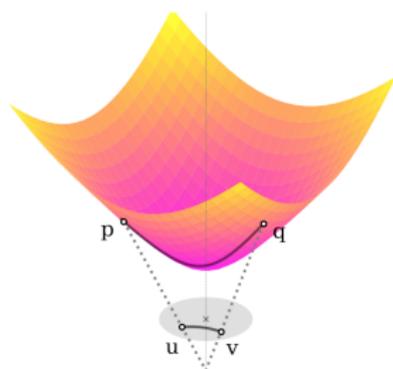
The Lorentz model

Maps N -dim hyperbolic space to surface in \mathbb{R}^{N+1}

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}} = -u_0 v_0 + \sum_{i=1}^{N+1} u_i v_i$$

$$\mathcal{H}^N = \{ \mathbf{x} \in \mathbb{R}^{N+1} : \langle \mathbf{x}, \mathbf{x} \rangle_{\mathcal{L}} = -1, x_0 > 0 \}$$

Metric: $\langle \mathbf{u}, \mathbf{v} \rangle_{T_x \mathcal{L}} = \langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}}$



The Lorentz model

Maps N -dim hyperbolic space to surface in \mathbb{R}^{N+1}

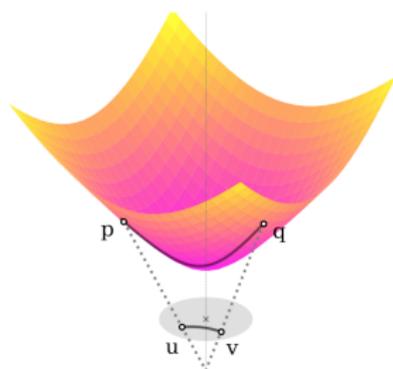
$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}} = -u_0 v_0 + \sum_{i=1}^{N+1} u_i v_i$$

$$\mathcal{H}^N = \{ \mathbf{x} \in \mathbb{R}^{N+1} : \langle \mathbf{x}, \mathbf{x} \rangle_{\mathcal{L}} = -1, x_0 > 0 \}$$

Metric: $\langle \mathbf{u}, \mathbf{v} \rangle_{T_x \mathcal{L}} = \langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}}$

Geodesics: starting at \mathbf{x} and going in direction \mathbf{v} , $\|\mathbf{v}\|_{\mathcal{L}} = 1$

$$\exp_{\mathbf{x}}(t\mathbf{v}) = \cosh(t)\mathbf{x} + \sinh(t)\mathbf{v}$$



The Lorentz model

Maps N -dim hyperbolic space to surface in \mathbb{R}^{N+1}

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}} = -u_0 v_0 + \sum_{i=1}^{N+1} u_i v_i$$

$$\mathcal{H}^N = \{ \mathbf{x} \in \mathbb{R}^{N+1} : \langle \mathbf{x}, \mathbf{x} \rangle_{\mathcal{L}} = -1, x_0 > 0 \}$$

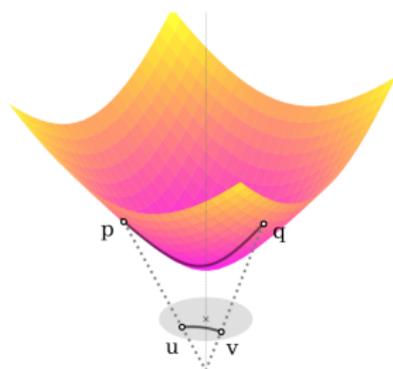
Metric: $\langle \mathbf{u}, \mathbf{v} \rangle_{T_{\mathbf{x}}\mathcal{L}} = \langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}}$

Geodesics: starting at \mathbf{x} and going in direction \mathbf{v} , $\|\mathbf{v}\|_{\mathcal{L}} = 1$

$$\exp_{\mathbf{x}}(t\mathbf{v}) = \cosh(t)\mathbf{x} + \sinh(t)\mathbf{v}$$

Projection: maps vectors onto tangent space:

$$\text{proj}_{\mathbf{x}}(\mathbf{u}) = \mathbf{u} + \langle \mathbf{x}, \mathbf{u} \rangle_{\mathcal{L}} \mathbf{x}$$



Riemannian SGD

Scale Euclidean gradient by inverse metric:

$$\mathbf{h} = \mathbf{M}_{\mathbf{x}}^{-1} \nabla f(\mathbf{x})$$

Riemannian SGD

Scale Euclidean gradient by inverse metric:

$$\mathbf{h} = \mathbf{M}_{\mathbf{x}}^{-1} \nabla f(\mathbf{x})$$

Project gradient onto tangent space:

$$\text{grad}f(\mathbf{x}) = \text{proj}(\mathbf{h})$$

Riemannian SGD

Scale Euclidean gradient by inverse metric:

$$\mathbf{h} = \mathbf{M}_{\mathbf{x}}^{-1} \nabla f(\mathbf{x})$$

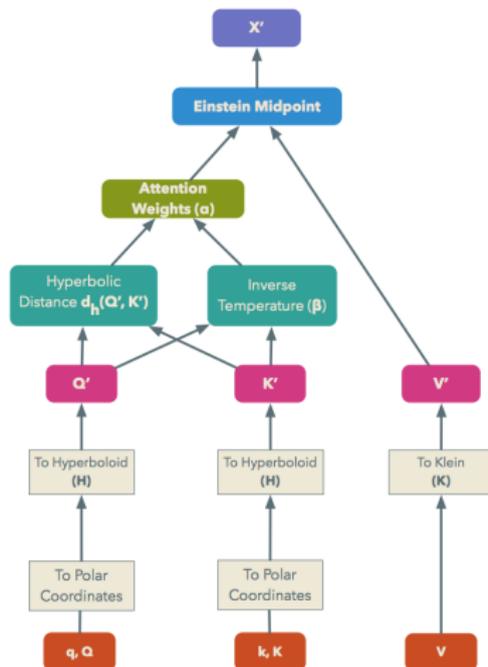
Project gradient onto tangent space:

$$\text{grad}f(\mathbf{x}) = \text{proj}(\mathbf{h})$$

Follow geodesic along Riemannian gradient:

$$\mathbf{x} \leftarrow \exp_{\mathbf{x}}(-\eta \text{grad}f(\mathbf{x}))$$

Hyperbolic Attention Networks



Imposes hyperbolic geometry on **activations** of deep network with attention

Analyzing the Geometry of Deep Generative Models

What is the shape of latent space?



What is the shape of latent space?



Are straight lines in latent space **really straight**?

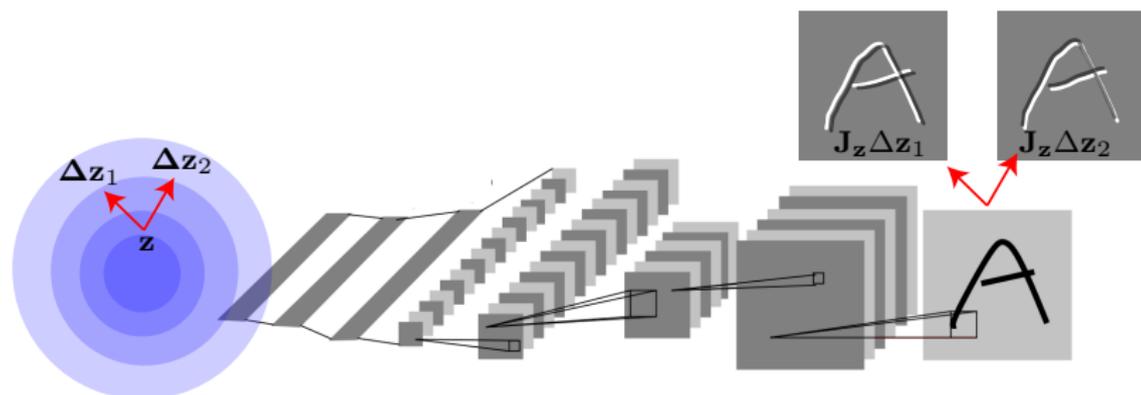
What is the shape of latent space?



Are straight lines in latent space **really straight**?

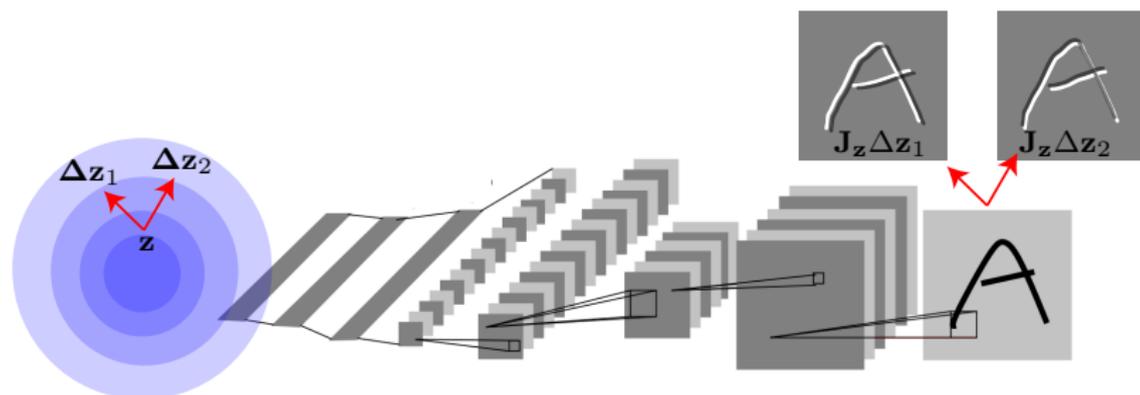
What is the right notion of **distance** in latent space?

Latent space metric



Deep generative model with decoder $f(\mathbf{z})$

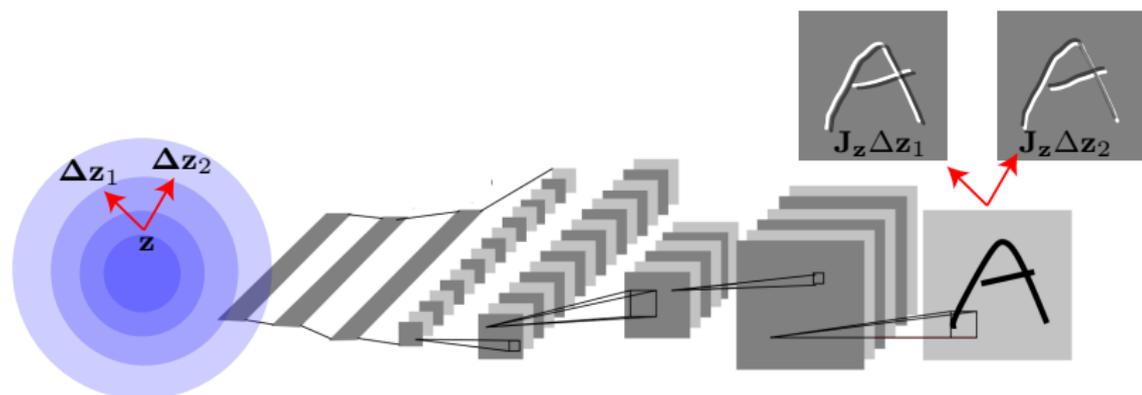
Latent space metric



Deep generative model with decoder $f(\mathbf{z})$

Manifold \mathcal{X} generated by f has tangent space $T_z \mathcal{X} = \text{span}(\mathbf{J}_z)$ where \mathbf{J}_z is **Jacobian** of f at \mathbf{z}

Latent space metric



Deep generative model with decoder $f(\mathbf{z})$

Manifold \mathcal{X} generated by f has tangent space $T_{\mathbf{z}}\mathcal{X} = \text{span}(\mathbf{J}_{\mathbf{z}})$ where $\mathbf{J}_{\mathbf{z}}$ is **Jacobian** of f at \mathbf{z}

Use ℓ_2 metric in **observation space** as metric in latent space:

$$\langle \Delta \mathbf{z}_1, \Delta \mathbf{z}_2 \rangle_{T_{\mathbf{z}}\mathcal{X}} = \Delta \mathbf{z}_1^T \mathbf{J}_{\mathbf{z}}^T \mathbf{J}_{\mathbf{z}} \Delta \mathbf{z}_2 = \Delta \mathbf{z}_1^T \mathbf{M}_{\mathbf{z}} \Delta \mathbf{z}_2$$

Geodesic equation

How do we **derive** the geodesic equation from the introduction:

$$\ddot{\gamma} = -\mathbf{\Gamma}_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t))$$

Geodesic equation

How do we **derive** the geodesic equation from the introduction:

$$\ddot{\gamma} = -\mathbf{\Gamma}_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t))$$

$$\gamma^* = \min_{\gamma} \mathcal{S}[\gamma] = \min_{\gamma} \int_0^1 dt L(\gamma, \dot{\gamma}, t)$$

Geodesic equation

How do we **derive** the geodesic equation from the introduction:

$$\ddot{\gamma} = -\mathbf{\Gamma}_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t))$$

$$\begin{aligned}\gamma^* = \min_{\gamma} \mathcal{S}[\gamma] &= \min_{\gamma} \int_0^1 dt L(\gamma, \dot{\gamma}, t) \\ &= \min_{\gamma} \int_0^1 dt \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle_{T_{\gamma(t)}\mathcal{X}}\end{aligned}$$

Geodesic equation

How do we **derive** the geodesic equation from the introduction:

$$\ddot{\gamma} = -\mathbf{\Gamma}_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t))$$

$$\begin{aligned}\gamma^* = \min_{\gamma} \mathcal{S}[\gamma] &= \min_{\gamma} \int_0^1 dt L(\gamma, \dot{\gamma}, t) \\ &= \min_{\gamma} \int_0^1 dt \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle_{T_{\gamma(t)}\mathcal{X}} \\ &= \min_{\gamma} \int_0^1 dt \sum_{ij} \mathbf{M}_{\gamma(t)}^{ij} \dot{\gamma}_i(t) \dot{\gamma}_j(t)\end{aligned}$$

Geodesic equation

How do we **derive** the geodesic equation from the introduction:

$$\ddot{\gamma} = -\mathbf{\Gamma}_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t))$$

$$\begin{aligned}\gamma^* = \min_{\gamma} \mathcal{S}[\gamma] &= \min_{\gamma} \int_0^1 dt L(\gamma, \dot{\gamma}, t) \\ &= \min_{\gamma} \int_0^1 dt \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle_{T_{\gamma(t)}\mathcal{X}} \\ &= \min_{\gamma} \int_0^1 dt \sum_{ij} \mathbf{M}_{\gamma(t)}^{ij} \dot{\gamma}_i(t) \dot{\gamma}_j(t)\end{aligned}$$

Euler-Lagrange equation:

$$\frac{\partial L}{\partial \gamma^*} = \frac{d}{dt} \frac{\partial L}{\partial \dot{\gamma}^*}$$

Geodesic equation

How do we **derive** the geodesic equation from the introduction:

$$\ddot{\gamma} = -\mathbf{\Gamma}_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t))$$

$$\begin{aligned}\gamma^* = \min_{\gamma} \mathcal{S}[\gamma] &= \min_{\gamma} \int_0^1 dt L(\gamma, \dot{\gamma}, t) \\ &= \min_{\gamma} \int_0^1 dt \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle_{T_{\gamma(t)}\mathcal{X}} \\ &= \min_{\gamma} \int_0^1 dt \sum_{ij} \mathbf{M}_{\gamma(t)}^{ij} \dot{\gamma}_i(t) \dot{\gamma}_j(t)\end{aligned}$$

Euler-Lagrange equation:

$$\frac{\partial L}{\partial \gamma^*} = \frac{d}{dt} \frac{\partial L}{\partial \dot{\gamma}^*}$$

Exercise: Derive the left and right side of the Euler-Lagrange equation for $L(\gamma, \dot{\gamma}, t) = \sum_{ij} \mathbf{M}_{\gamma(t)}^{ij} \dot{\gamma}_i(t) \dot{\gamma}_j(t)$

Geodesic equation

$$\frac{\partial L}{\partial \gamma_i} = \sum_{jk} \frac{\partial \mathbf{M}_{\gamma(t)}^{jk}}{\partial \gamma_i} \dot{\gamma}_k(t) \dot{\gamma}_j(t)$$

Geodesic equation

$$\frac{\partial L}{\partial \gamma_i} = \sum_{jk} \frac{\partial \mathbf{M}_{\gamma(t)}^{jk}}{\partial \gamma_i} \dot{\gamma}_k(t) \dot{\gamma}_j(t)$$

$$\begin{aligned} \frac{d}{dt} \frac{\partial L}{\partial \dot{\gamma}_i} &= \frac{d}{dt} \left(2 \sum_j \mathbf{M}_{\gamma(t)}^{ij} \dot{\gamma}_j(t) \right) = 2 \sum_j \left[\frac{\partial \mathbf{M}_{\gamma(t)}^{ij}}{\partial t} \dot{\gamma}_j(t) + \mathbf{M}_{\gamma(t)}^{ij} \ddot{\gamma}_j(t) \right] \\ &= 2 \sum_j \left[\sum_k \frac{\partial \mathbf{M}_{\gamma(t)}^{ij}}{\partial \gamma_k} \dot{\gamma}_k(t) \dot{\gamma}_j(t) + \mathbf{M}_{\gamma(t)}^{ij} \ddot{\gamma}_j(t) \right] \end{aligned}$$

Geodesic equation

$$\frac{\partial L}{\partial \gamma_i} = \sum_{jk} \frac{\partial \mathbf{M}_{\gamma(t)}^{jk}}{\partial \gamma_i} \dot{\gamma}_k(t) \dot{\gamma}_j(t)$$

$$\begin{aligned} \frac{d}{dt} \frac{\partial L}{\partial \dot{\gamma}_i} &= \frac{d}{dt} \left(2 \sum_j \mathbf{M}_{\gamma(t)}^{ij} \dot{\gamma}_j(t) \right) = 2 \sum_j \left[\frac{\partial \mathbf{M}_{\gamma(t)}^{ij}}{\partial t} \dot{\gamma}_j(t) + \mathbf{M}_{\gamma(t)}^{ij} \ddot{\gamma}_j(t) \right] \\ &= 2 \sum_j \left[\sum_k \frac{\partial \mathbf{M}_{\gamma(t)}^{ij}}{\partial \gamma_k} \dot{\gamma}_k(t) \dot{\gamma}_j(t) + \mathbf{M}_{\gamma(t)}^{ij} \ddot{\gamma}_j(t) \right] \end{aligned}$$

$$\ddot{\gamma}_j(t) = -\frac{1}{2} \sum_i \left(\mathbf{M}_{\gamma(t)}^{-1} \right)^{ij} \sum_k \left[\left(2 \frac{\partial \mathbf{M}_{\gamma(t)}^{ij}}{\partial \gamma_k} - \frac{\partial \mathbf{M}_{\gamma(t)}^{jk}}{\partial \gamma_i} \right) \dot{\gamma}_k(t) \dot{\gamma}_j(t) \right]$$

Geodesic equation

$$\frac{\partial L}{\partial \gamma_i} = \sum_{jk} \frac{\partial \mathbf{M}_{\gamma(t)}^{jk}}{\partial \gamma_i} \dot{\gamma}_k(t) \dot{\gamma}_j(t)$$

$$\begin{aligned} \frac{d}{dt} \frac{\partial L}{\partial \dot{\gamma}_i} &= \frac{d}{dt} \left(2 \sum_j \mathbf{M}_{\gamma(t)}^{ij} \dot{\gamma}_j(t) \right) = 2 \sum_j \left[\frac{\partial \mathbf{M}_{\gamma(t)}^{ij}}{\partial t} \dot{\gamma}_j(t) + \mathbf{M}_{\gamma(t)}^{ij} \ddot{\gamma}_j(t) \right] \\ &= 2 \sum_j \left[\sum_k \frac{\partial \mathbf{M}_{\gamma(t)}^{ij}}{\partial \gamma_k} \dot{\gamma}_k(t) \dot{\gamma}_j(t) + \mathbf{M}_{\gamma(t)}^{ij} \ddot{\gamma}_j(t) \right] \end{aligned}$$

$$\begin{aligned} \ddot{\gamma}_j(t) &= - \sum_{ik} \Gamma_j^{ik}(\gamma(t)) \dot{\gamma}_k(t) \dot{\gamma}_j(t) \\ \Gamma_j^{ik}(x) &= \frac{1}{2} (\mathbf{M}_x^{-1})^{ij} \left(2 \frac{\partial \mathbf{M}_x^{ij}}{\partial x_k} - \frac{\partial \mathbf{M}_x^{jk}}{\partial x_i} \right) \end{aligned}$$

Geodesic equation

$$\frac{\partial L}{\partial \gamma_i} = \sum_{jk} \frac{\partial \mathbf{M}_{\gamma(t)}^{jk}}{\partial \gamma_i} \dot{\gamma}_k(t) \dot{\gamma}_j(t)$$

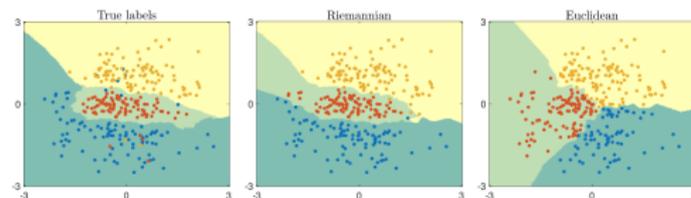
$$\begin{aligned} \frac{d}{dt} \frac{\partial L}{\partial \dot{\gamma}_i} &= \frac{d}{dt} \left(2 \sum_j \mathbf{M}_{\gamma(t)}^{ij} \dot{\gamma}_j(t) \right) = 2 \sum_j \left[\frac{\partial \mathbf{M}_{\gamma(t)}^{ij}}{\partial t} \dot{\gamma}_j(t) + \mathbf{M}_{\gamma(t)}^{ij} \ddot{\gamma}_j(t) \right] \\ &= 2 \sum_j \left[\sum_k \frac{\partial \mathbf{M}_{\gamma(t)}^{ij}}{\partial \gamma_k} \dot{\gamma}_k(t) \dot{\gamma}_j(t) + \mathbf{M}_{\gamma(t)}^{ij} \ddot{\gamma}_j(t) \right] \end{aligned}$$

$$\ddot{\gamma}_j(t) = - \sum_{ik} \Gamma_j^{ik}(\gamma(t)) \dot{\gamma}_k(t) \dot{\gamma}_j(t)$$

$$\Gamma_j^{ik}(x) = \frac{1}{2} (\mathbf{M}_x^{-1})^{ij} \left(2 \frac{\partial \mathbf{M}_x^{ij}}{\partial x_k} - \frac{\partial \mathbf{M}_x^{jk}}{\partial x_i} \right)$$

$\Gamma_j^{ik}(x)$ are the parameters of the **Levi-Civita connection**:
unique connection induced by the metric

The geodesics of deep generative models

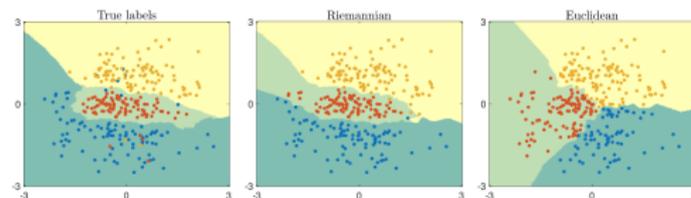


Digits	Linear	Riemannian
{0, 1, 2}	77.57(± 0.87)%	94.28(± 1.14)%
{3, 4, 7}	77.80(± 0.91)%	89.54(± 1.61)%
{5, 6, 9}	64.93(± 0.96)%	81.13(± 2.52)%

Table 1: The F -measure results for k -means.

Classification from latent representations is **better**

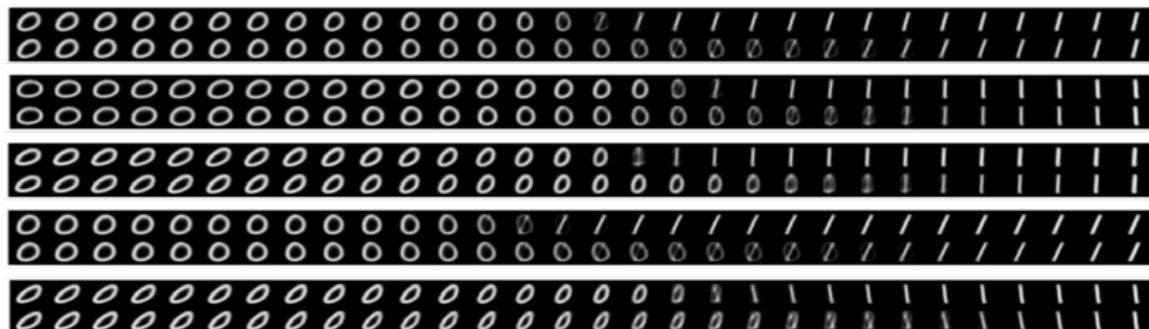
The geodesics of deep generative models



Digits	Linear	Riemannian
{0, 1, 2}	77.57(± 0.87)%	94.28(± 1.14)%
{3, 4, 7}	77.80(± 0.91)%	89.54(± 1.61)%
{5, 6, 9}	64.93(± 0.96)%	81.13(± 2.52)%

Table 1: The F -measure results for k -means.

Classification from latent representations is **better**



Transitions along geodesics are **smoother**

Summary

Differential geometry and spectral theory are a **powerful** suite of tools for thinking about the geometry of data

Summary

Differential geometry and spectral theory are a **powerful** suite of tools for thinking about the geometry of data

Straight lines generalize to **geodesics**, which are defined from the local metric structure

Summary

Differential geometry and spectral theory are a **powerful** suite of tools for thinking about the geometry of data

Straight lines generalize to **geodesics**, which are defined from the local metric structure

The spectrum of the **Laplacian** generalizes Fourier analysis to graphs and manifolds

Summary

Differential geometry and spectral theory are a **powerful** suite of tools for thinking about the geometry of data

Straight lines generalize to **geodesics**, which are defined from the local metric structure

The spectrum of the **Laplacian** generalizes Fourier analysis to graphs and manifolds

Classic manifold learning uses spectral decompositions to map data manifolds in high-dim space to **flat** space

Summary

Differential geometry and spectral theory are a **powerful** suite of tools for thinking about the geometry of data

Straight lines generalize to **geodesics**, which are defined from the local metric structure

The spectrum of the **Laplacian** generalizes Fourier analysis to graphs and manifolds

Classic manifold learning uses spectral decompositions to map data manifolds in high-dim space to **flat** space

Hierarchical data is naturally mapped to **hyperbolic** space

Summary

Differential geometry and spectral theory are a **powerful** suite of tools for thinking about the geometry of data

Straight lines generalize to **geodesics**, which are defined from the local metric structure

The spectrum of the **Laplacian** generalizes Fourier analysis to graphs and manifolds

Classic manifold learning uses spectral decompositions to map data manifolds in high-dim space to **flat** space

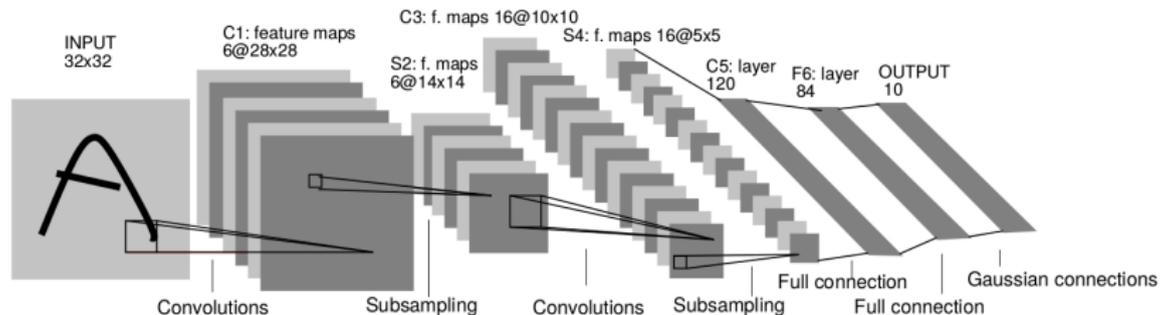
Hierarchical data is naturally mapped to **hyperbolic** space

The latent space of deep generative models is better understood as being **curved**

Spectral Deep Learning

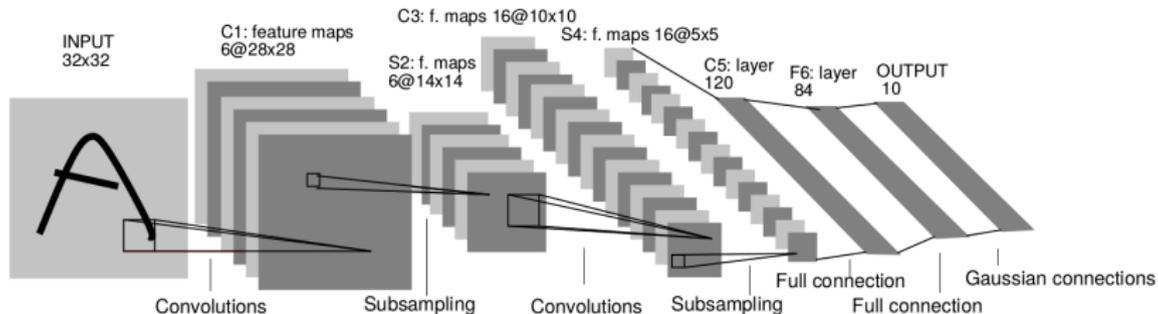
Convolutions on Graphs and Manifolds

Key properties of CNNs



😊 Convolutional (**Translation invariance**)

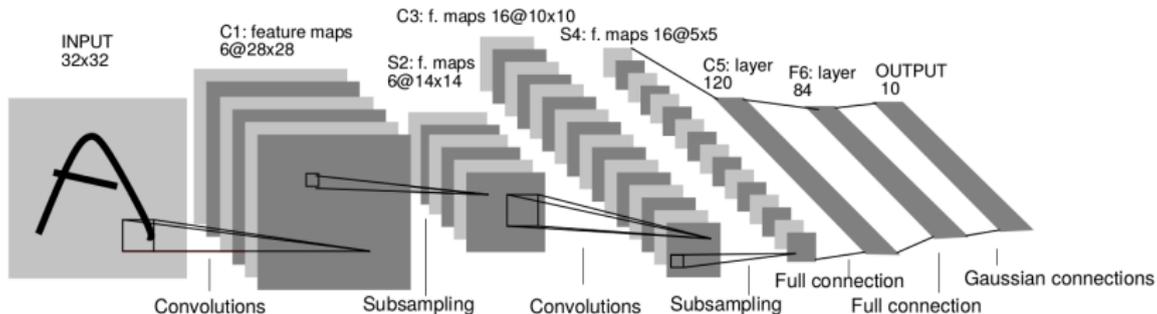
Key properties of CNNs



😊 Convolutional (**Translation invariance**)

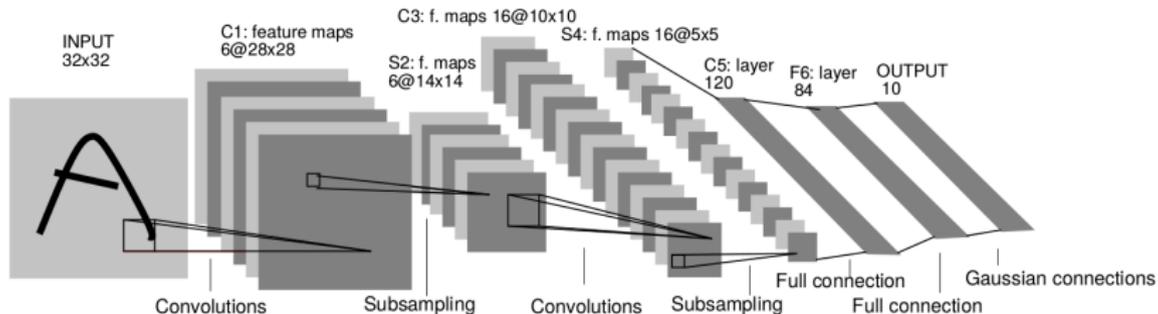
😊 Scale Separation (**Compositionality**)

Key properties of CNNs



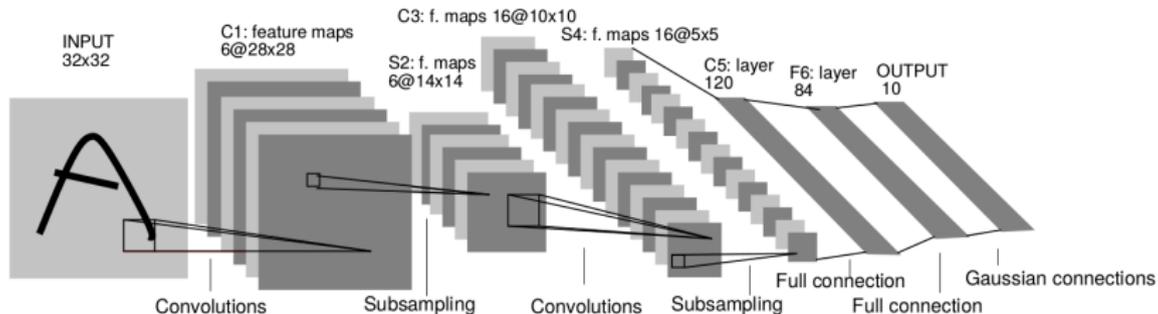
- ☺ Convolutional (**Translation invariance**)
- ☺ Scale Separation (**Compositionality**)
- ☺ Filters localized in space (**Deformation Stability**)

Key properties of CNNs



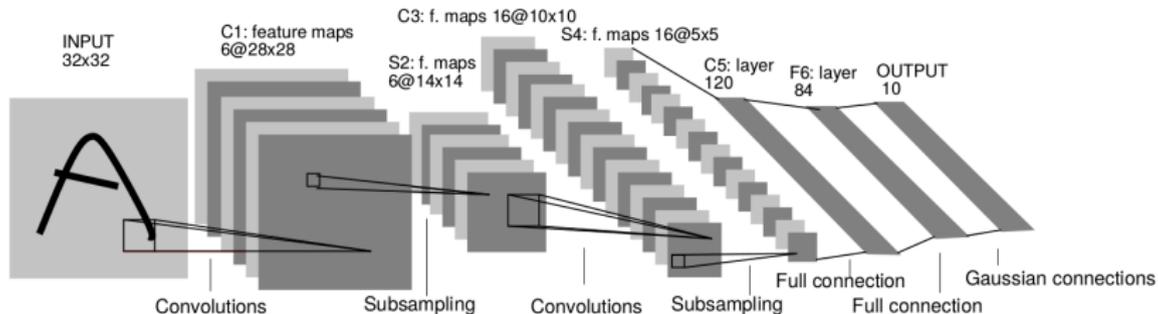
- ☺ Convolutional (**Translation invariance**)
- ☺ Scale Separation (**Compositionality**)
- ☺ Filters localized in space (**Deformation Stability**)
- ☺ $\mathcal{O}(1)$ parameters per filter (independent of input image size n)

Key properties of CNNs



- ☺ Convolutional (**Translation invariance**)
- ☺ Scale Separation (**Compositionality**)
- ☺ Filters localized in space (**Deformation Stability**)
- ☺ $\mathcal{O}(1)$ parameters per filter (independent of input image size n)
- ☺ $\mathcal{O}(n)$ complexity per layer (filtering done in the spatial domain)

Key properties of CNNs



- ☺ Convolutional (**Translation invariance**)
- ☺ Scale Separation (**Compositionality**)
- ☺ Filters localized in space (**Deformation Stability**)
- ☺ $\mathcal{O}(1)$ parameters per filter (independent of input image size n)
- ☺ $\mathcal{O}(n)$ complexity per layer (filtering done in the spatial domain)
- ☺ $\mathcal{O}(\log n)$ layers in classification tasks

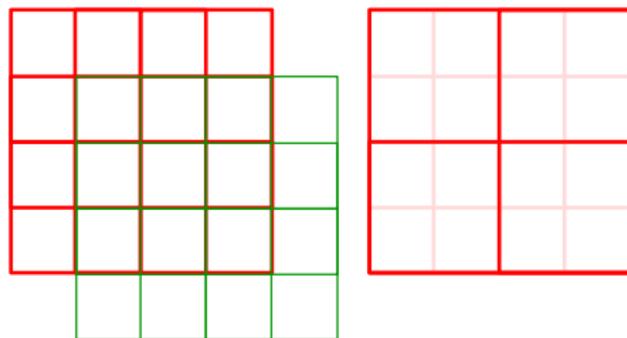
CNNs and Euclidean Geometry

CNNs are defined over Euclidean domains or Grids Ω . Two fundamental properties:

Translation Invariance (yielding convolutions).

Multiscale structure (yielding downsampling).

Inductive bias that exploits stationarity and deformation stability of many tasks.



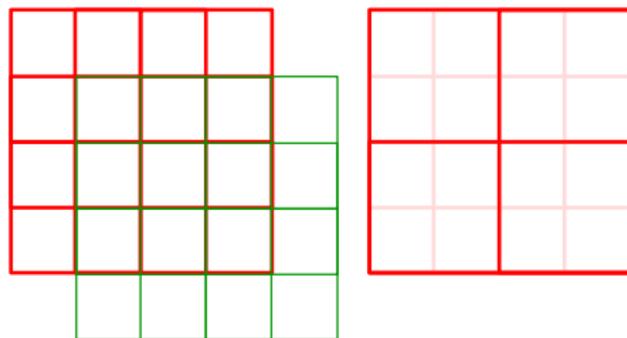
CNNs and Euclidean Geometry

CNNs are defined over Euclidean domains or Grids Ω . Two fundamental properties:

Translation Invariance (yielding convolutions).

Multiscale structure (yielding downsampling).

Inductive bias that exploits stationarity and deformation stability of many tasks.



Roadmap: extend CNNs to non-Euclidean geometries by replacing filtering and pooling by appropriate operators

Convolution: Euclidean space

Given two functions $f, g : [-\pi, \pi] \rightarrow \mathbb{R}$ their **convolution** is a function

$$(f \star g)(x) = \int_{-\pi}^{\pi} f(x')g(x - x')dx'$$

Convolution: Euclidean space

Given two functions $f, g : [-\pi, \pi] \rightarrow \mathbb{R}$ their **convolution** is a function

$$(f \star g)(x) = \int_{-\pi}^{\pi} f(x')g(x - x')dx'$$

Shift-invariance: $f(x - x_0) \star g(x) = (f \star g)(x - x_0)$

Convolution: Euclidean space

Given two functions $f, g : [-\pi, \pi] \rightarrow \mathbb{R}$ their **convolution** is a function

$$(f \star g)(x) = \int_{-\pi}^{\pi} f(x')g(x - x')dx'$$

Shift-invariance: $f(x - x_0) \star g(x) = (f \star g)(x - x_0)$

Convolution theorem: Fourier transform diagonalizes the convolution operator

Convolution: Euclidean space

Given two functions $f, g : [-\pi, \pi] \rightarrow \mathbb{R}$ their **convolution** is a function

$$(f \star g)(x) = \int_{-\pi}^{\pi} f(x')g(x - x')dx'$$

Shift-invariance: $f(x - x_0) \star g(x) = (f \star g)(x - x_0)$

Convolution theorem: Fourier transform diagonalizes the convolution operator
 \Rightarrow convolution can be computed in the Fourier domain as

$$\widehat{(f \star g)} = \hat{f} \cdot \hat{g}$$

Convolution: Euclidean space

Given two functions $f, g : [-\pi, \pi] \rightarrow \mathbb{R}$ their **convolution** is a function

$$(f \star g)(x) = \int_{-\pi}^{\pi} f(x')g(x - x')dx'$$

Shift-invariance: $f(x - x_0) \star g(x) = (f \star g)(x - x_0)$

Convolution theorem: Fourier transform diagonalizes the convolution operator
 \Rightarrow convolution can be computed in the Fourier domain as

$$\widehat{(f \star g)} = \hat{f} \cdot \hat{g}$$

Efficient computation using FFT

Convolution Theorem

Convolution of two vectors $\mathbf{f} = (f_1, \dots, f_n)^\top$ and $\mathbf{g} = (g_1, \dots, g_n)^\top$

$$\mathbf{f} \star \mathbf{g} = \begin{bmatrix} g_1 & g_2 & \cdots & \cdots & g_n \\ g_n & g_1 & g_2 & \cdots & g_{n-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ g_3 & g_4 & \cdots & g_1 & g_2 \\ g_2 & g_3 & \cdots & \cdots & g_1 \end{bmatrix} \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix}$$

Convolution Theorem

Convolution of two vectors $\mathbf{f} = (f_1, \dots, f_n)^\top$ and $\mathbf{g} = (g_1, \dots, g_n)^\top$

$$\mathbf{f} \star \mathbf{g} = \underbrace{\begin{bmatrix} g_1 & g_2 & \cdots & \cdots & g_n \\ g_n & g_1 & g_2 & \cdots & g_{n-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ g_3 & g_4 & \cdots & g_1 & g_2 \\ g_2 & g_3 & \cdots & \cdots & g_1 \end{bmatrix}}_{\text{circulant matrix}} \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix}$$

Convolution Theorem

Convolution of two vectors $\mathbf{f} = (f_1, \dots, f_n)^\top$ and $\mathbf{g} = (g_1, \dots, g_n)^\top$

$$\mathbf{f} \star \mathbf{g} = \underbrace{\begin{bmatrix} g_1 & g_2 & \cdots & \cdots & g_n \\ g_n & g_1 & g_2 & \cdots & g_{n-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ g_3 & g_4 & \cdots & g_1 & g_2 \\ g_2 & g_3 & \cdots & \cdots & g_1 \end{bmatrix}}_{\text{diagonalized by Fourier basis}} \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix}$$

Convolution Theorem

Convolution of two vectors $\mathbf{f} = (f_1, \dots, f_n)^\top$ and $\mathbf{g} = (g_1, \dots, g_n)^\top$

$$\begin{aligned}\mathbf{f} \star \mathbf{g} &= \begin{bmatrix} g_1 & g_2 & \cdots & \cdots & g_n \\ g_n & g_1 & g_2 & \cdots & g_{n-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ g_3 & g_4 & \cdots & g_1 & g_2 \\ g_2 & g_3 & \cdots & \cdots & g_1 \end{bmatrix} \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix} \\ &= \mathbf{\Phi} \begin{bmatrix} \hat{g}_1 & & & & \\ & \ddots & & & \\ & & \hat{g}_n & & \end{bmatrix} \mathbf{\Phi}^\top \mathbf{f}\end{aligned}$$

Convolution Theorem

Convolution of two vectors $\mathbf{f} = (f_1, \dots, f_n)^\top$ and $\mathbf{g} = (g_1, \dots, g_n)^\top$

$$\mathbf{f} \star \mathbf{g} = \begin{bmatrix} g_1 & g_2 & \cdots & \cdots & g_n \\ g_n & g_1 & g_2 & \cdots & g_{n-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ g_3 & g_4 & \cdots & g_1 & g_2 \\ g_2 & g_3 & \cdots & \cdots & g_1 \end{bmatrix} \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix}$$

$$= \Phi \begin{bmatrix} \hat{g}_1 & & & \\ & \ddots & & \\ & & \hat{g}_n & \end{bmatrix} \begin{bmatrix} \hat{f}_1 \\ \vdots \\ \hat{f}_n \end{bmatrix}$$

Convolution Theorem

Convolution of two vectors $\mathbf{f} = (f_1, \dots, f_n)^\top$ and $\mathbf{g} = (g_1, \dots, g_n)^\top$

$$\begin{aligned}\mathbf{f} \star \mathbf{g} &= \begin{bmatrix} g_1 & g_2 & \cdots & \cdots & g_n \\ g_n & g_1 & g_2 & \cdots & g_{n-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ g_3 & g_4 & \cdots & g_1 & g_2 \\ g_2 & g_3 & \cdots & \cdots & g_1 \end{bmatrix} \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix} \\ &= \Phi \begin{bmatrix} \hat{f}_1 \cdot \hat{g}_1 \\ \vdots \\ \hat{f}_n \cdot \hat{g}_n \end{bmatrix}\end{aligned}$$

Spectral convolution

Spectral convolution of $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{V})$ can be defined by analogy

$$\mathbf{f} \star \mathbf{g} = \sum_{k \geq 1} \langle \mathbf{f}, \phi_k \rangle_{L^2(\mathcal{V})} \langle \mathbf{g}, \phi_k \rangle_{L^2(\mathcal{V})} \phi_k$$

Spectral convolution

Spectral convolution of $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{V})$ can be defined by analogy

$$\mathbf{f} \star \mathbf{g} = \sum_{k \geq 1} \underbrace{\langle \mathbf{f}, \phi_k \rangle_{L^2(\mathcal{V})} \langle \mathbf{g}, \phi_k \rangle_{L^2(\mathcal{V})}}_{\text{product in the Fourier domain}} \phi_k$$

Spectral convolution

Spectral convolution of $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{V})$ can be defined by analogy

$$\mathbf{f} \star \mathbf{g} = \underbrace{\sum_{k \geq 1} \underbrace{\langle \mathbf{f}, \phi_k \rangle_{L^2(\mathcal{V})} \langle \mathbf{g}, \phi_k \rangle_{L^2(\mathcal{V})}}_{\text{product in the Fourier domain}} \phi_k}_{\text{inverse Fourier transform}}$$

Spectral convolution

Spectral convolution of $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{V})$ can be defined by analogy

$$\mathbf{f} \star \mathbf{g} = \sum_{k \geq 1} \langle \mathbf{f}, \phi_k \rangle_{L^2(\mathcal{V})} \langle \mathbf{g}, \phi_k \rangle_{L^2(\mathcal{V})} \phi_k$$

In matrix-vector notation

$$\mathbf{f} \star \mathbf{g} = \Phi (\Phi^\top \mathbf{g}) \circ (\Phi^\top \mathbf{f})$$

Spectral convolution

Spectral convolution of $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{V})$ can be defined by analogy

$$\mathbf{f} \star \mathbf{g} = \sum_{k \geq 1} \langle \mathbf{f}, \phi_k \rangle_{L^2(\mathcal{V})} \langle \mathbf{g}, \phi_k \rangle_{L^2(\mathcal{V})} \phi_k$$

In matrix-vector notation

$$\mathbf{f} \star \mathbf{g} = \Phi \operatorname{diag}(\hat{g}_1, \dots, \hat{g}_n) \Phi^\top \mathbf{f}$$

Spectral convolution

Spectral convolution of $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{V})$ can be defined by analogy

$$\mathbf{f} \star \mathbf{g} = \sum_{k \geq 1} \langle \mathbf{f}, \phi_k \rangle_{L^2(\mathcal{V})} \langle \mathbf{g}, \phi_k \rangle_{L^2(\mathcal{V})} \phi_k$$

In matrix-vector notation

$$\mathbf{f} \star \mathbf{g} = \underbrace{\Phi \operatorname{diag}(\hat{g}_1, \dots, \hat{g}_n) \Phi^\top}_{\mathbf{G}} \mathbf{f}$$

Spectral convolution

Spectral convolution of $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{V})$ can be defined by analogy

$$\mathbf{f} \star \mathbf{g} = \sum_{k \geq 1} \langle \mathbf{f}, \phi_k \rangle_{L^2(\mathcal{V})} \langle \mathbf{g}, \phi_k \rangle_{L^2(\mathcal{V})} \phi_k$$

In matrix-vector notation

$$\mathbf{f} \star \mathbf{g} = \underbrace{\Phi \text{diag}(\hat{g}_1, \dots, \hat{g}_n) \Phi^\top}_{\mathbf{G}} \mathbf{f}$$

Not shift-invariant! (\mathbf{G} has no circulant structure)

Spectral convolution

Spectral convolution of $\mathbf{f}, \mathbf{g} \in L^2(\mathcal{V})$ can be defined by analogy

$$\mathbf{f} \star \mathbf{g} = \sum_{k \geq 1} \langle \mathbf{f}, \phi_k \rangle_{L^2(\mathcal{V})} \langle \mathbf{g}, \phi_k \rangle_{L^2(\mathcal{V})} \phi_k$$

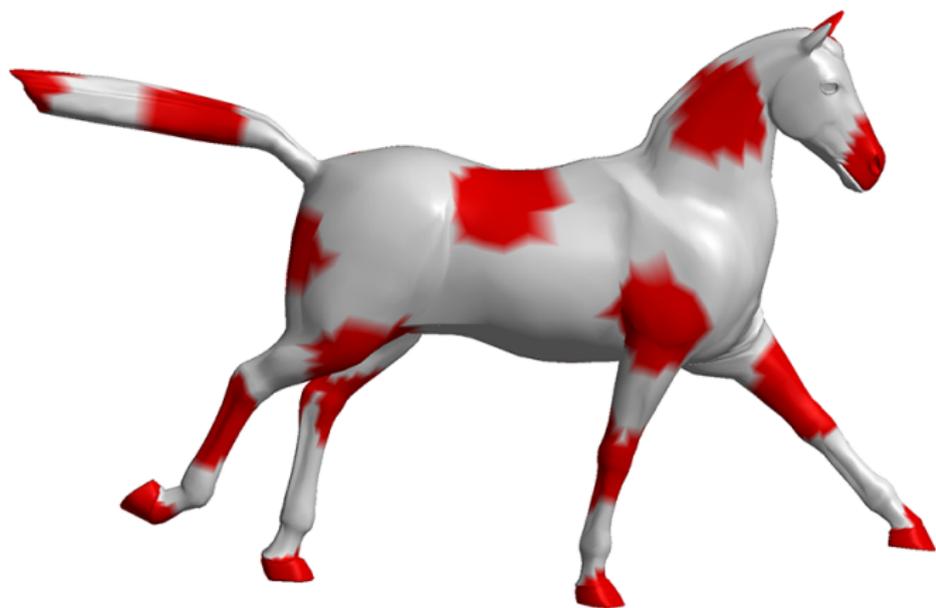
In matrix-vector notation

$$\mathbf{f} \star \mathbf{g} = \underbrace{\Phi \text{diag}(\hat{g}_1, \dots, \hat{g}_n) \Phi^\top}_{\mathbf{G}} \mathbf{f}$$

Not shift-invariant! (\mathbf{G} has no circulant structure)

Filter coefficients depend on basis ϕ_1, \dots, ϕ_n

Effect of spectral convolution



Function f

Effect of spectral convolution



'Edge detecting' spectral filter $\Phi G \Phi^T \mathbf{f}$

Effect of spectral convolution

Same spectral filter, different basis $\Psi \mathbf{G} \Psi^T \mathbf{f}$

Effect of spectral convolution

High-frequency Laplacian eigenvector ϕ_{50}

Spectral Graph Convolutional Neural Networks

Spectral graph CNN

Convolution expressed in the **spectral domain**

$$\mathbf{g} = \mathbf{\Phi} \mathbf{W} \mathbf{\Phi}^\top \mathbf{f}$$

where \mathbf{W} is $n \times n$ diagonal matrix of learnable spectral filter coefficients

Spectral graph CNN

Convolution expressed in the **spectral domain**

$$\mathbf{g} = \mathbf{\Phi} \mathbf{W} \mathbf{\Phi}^\top \mathbf{f}$$

where \mathbf{W} is $n \times n$ diagonal matrix of learnable spectral filter coefficients

☹ Filters are basis-dependent \Rightarrow does not generalize across graphs!

Spectral graph CNN

Convolution expressed in the **spectral domain**

$$\mathbf{g} = \Phi \mathbf{W} \Phi^\top \mathbf{f}$$

where \mathbf{W} is $n \times n$ diagonal matrix of learnable spectral filter coefficients

- ☹ Filters are basis-dependent \Rightarrow does not generalize across graphs!
- ☹ $\mathcal{O}(n)$ parameters per layer

Spectral graph CNN

Convolution expressed in the **spectral domain**

$$\mathbf{g} = \Phi \mathbf{W} \Phi^\top \mathbf{f}$$

where \mathbf{W} is $n \times n$ diagonal matrix of learnable spectral filter coefficients

- ☹ Filters are basis-dependent \Rightarrow does not generalize across graphs!
- ☹ $\mathcal{O}(n)$ parameters per layer
- ☹ $\mathcal{O}(n^2)$ computation of forward and inverse Fourier transforms Φ^\top, Φ
(no FFT on graphs)

Spectral graph CNN

Convolution expressed in the **spectral domain**

$$\mathbf{g} = \Phi \mathbf{W} \Phi^T \mathbf{f}$$

where \mathbf{W} is $n \times n$ diagonal matrix of learnable spectral filter coefficients

- ☹ Filters are basis-dependent \Rightarrow does not generalize across graphs!
- ☹ $\mathcal{O}(n)$ parameters per layer
- ☹ $\mathcal{O}(n^2)$ computation of forward and inverse Fourier transforms Φ^T, Φ
(no FFT on graphs)
- ☹ No guarantee of spatial localization of filters

Localization and Smoothness

Vanishing moments: In the Euclidean setting

$$\int_{-\infty}^{+\infty} |x|^{2k} |f(x)|^2 dx = \int_{-\infty}^{+\infty} \left| \frac{\partial^k \hat{f}(\omega)}{\partial \omega^k} \right|^2 d\omega$$

Localization and Smoothness

Vanishing moments: In the Euclidean setting

$$\int_{-\infty}^{+\infty} |x|^{2k} |f(x)|^2 dx = \int_{-\infty}^{+\infty} \left| \frac{\partial^k \hat{f}(\omega)}{\partial \omega^k} \right|^2 d\omega$$

Localization in space = smoothness in frequency domain

Localization and Smoothness

Vanishing moments: In the Euclidean setting

$$\int_{-\infty}^{+\infty} |x|^{2k} |f(x)|^2 dx = \int_{-\infty}^{+\infty} \left| \frac{\partial^k \hat{f}(\omega)}{\partial \omega^k} \right|^2 d\omega$$

Localization in space = smoothness in frequency domain

Parametrize the filter using a **smooth spectral transfer function** $\tau(\lambda)$

Localization and Smoothness

Vanishing moments: In the Euclidean setting

$$\int_{-\infty}^{+\infty} |x|^{2k} |f(x)|^2 dx = \int_{-\infty}^{+\infty} \left| \frac{\partial^k \hat{f}(\omega)}{\partial \omega^k} \right|^2 d\omega$$

Localization in space = smoothness in frequency domain

Parametrize the filter using a **smooth spectral transfer function** $\tau(\lambda)$

Application of the filter

$$\tau(\mathbf{\Delta})\mathbf{f} = \mathbf{\Phi}\tau(\mathbf{\Lambda})\mathbf{\Phi}^\top \mathbf{f}$$

Localization and Smoothness

Vanishing moments: In the Euclidean setting

$$\int_{-\infty}^{+\infty} |x|^{2k} |f(x)|^2 dx = \int_{-\infty}^{+\infty} \left| \frac{\partial^k \hat{f}(\omega)}{\partial \omega^k} \right|^2 d\omega$$

Localization in space = smoothness in frequency domain

Parametrize the filter using a **smooth spectral transfer function** $\tau(\lambda)$

Application of the filter

$$\tau(\mathbf{\Delta})\mathbf{f} = \mathbf{\Phi} \begin{pmatrix} \tau(\lambda_1) & & \\ & \ddots & \\ & & \tau(\lambda_n) \end{pmatrix} \mathbf{\Phi}^\top \mathbf{f}$$

Localization and Smoothness

Vanishing moments: In the Euclidean setting

$$\int_{-\infty}^{+\infty} |x|^{2k} |f(x)|^2 dx = \int_{-\infty}^{+\infty} \left| \frac{\partial^k \hat{f}(\omega)}{\partial \omega^k} \right|^2 d\omega$$

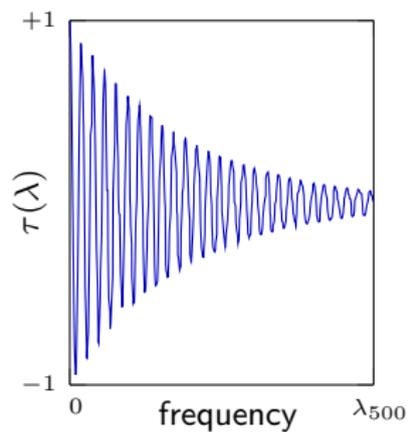
Localization in space = smoothness in frequency domain

Parametrize the filter using a **smooth spectral transfer function** $\tau(\lambda)$

Application of the parametric filter with learnable parameters α

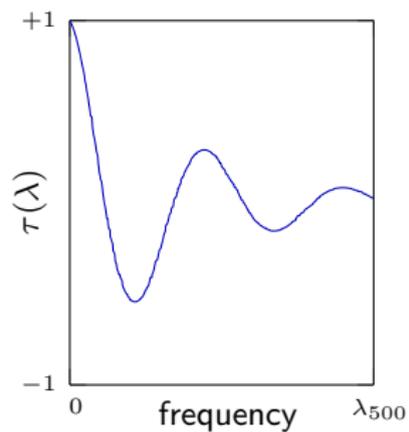
$$\tau_{\alpha}(\Delta)\mathbf{f} = \Phi \begin{pmatrix} \tau_{\alpha}(\lambda_1) & & \\ & \ddots & \\ & & \tau_{\alpha}(\lambda_n) \end{pmatrix} \Phi^{\top} \mathbf{f}$$

Localization and smoothness



Non-smooth spectral filter (delocalized in space)

Localization and smoothness



Smooth spectral filter (localized in space)

Spectral graph CNN with polynomial filters (ChebNet)

Represent spectral transfer function as a **polynomial** of order r

$$\tau_{\alpha}(\lambda) = \sum_{j=0}^r \alpha_j \lambda^j$$

where $\alpha = (\alpha_0, \dots, \alpha_r)^{\top}$ is the vector of filter parameters

Spectral graph CNN with polynomial filters (ChebNet)

Represent spectral transfer function as a **polynomial** of order r

$$\tau_{\alpha}(\lambda) = \sum_{j=0}^r \alpha_j \lambda^j$$

where $\alpha = (\alpha_0, \dots, \alpha_r)^{\top}$ is the vector of filter parameters

😊 $\mathcal{O}(1)$ parameters per layer

Spectral graph CNN with polynomial filters (ChebNet)

Represent spectral transfer function as a **polynomial** of order r

$$\tau_{\alpha}(\lambda) = \sum_{j=0}^r \alpha_j \lambda^j$$

where $\alpha = (\alpha_0, \dots, \alpha_r)^{\top}$ is the vector of filter parameters

😊 $\mathcal{O}(1)$ parameters per layer

😊 Filters have guaranteed r -hops support

Spectral graph CNN with polynomial filters (ChebNet)

Represent spectral transfer function as a **polynomial** of order r

$$\tau_{\alpha}(\lambda) = \sum_{j=0}^r \alpha_j \lambda^j$$

where $\alpha = (\alpha_0, \dots, \alpha_r)^{\top}$ is the vector of filter parameters

☺ $\mathcal{O}(1)$ parameters per layer

☺ Filters have guaranteed r -hops support

☺ No explicit computation of $\Phi^{\top}, \Phi \Rightarrow \mathcal{O}(nr)$ computational complexity

Example: citation networks

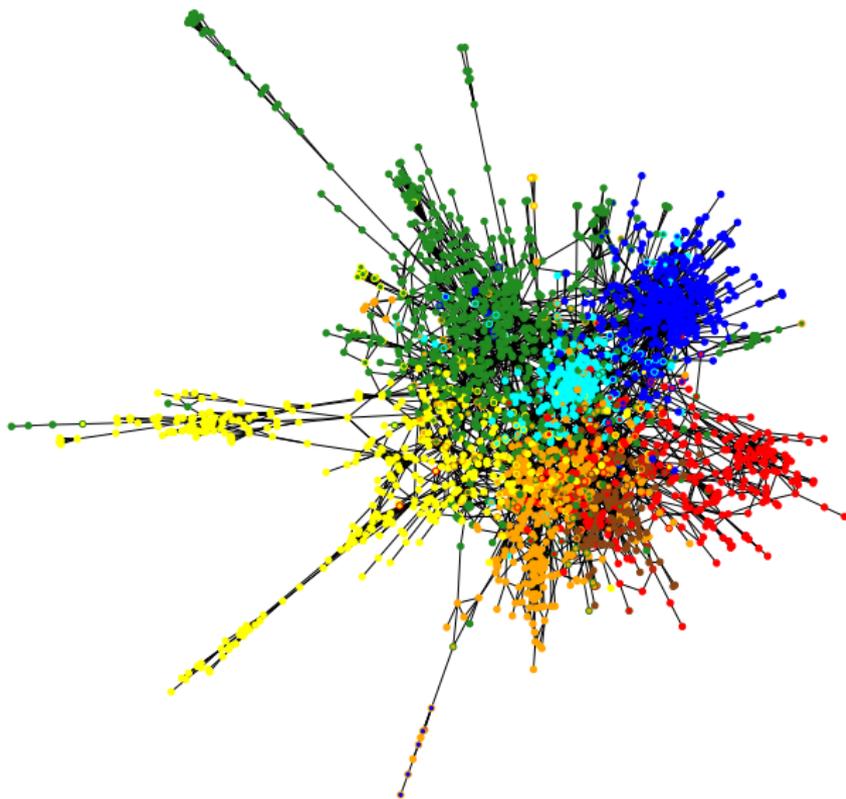


Figure: Monti, Boscaini, Masci, Rodolà, Svoboda, Bronstein 2017

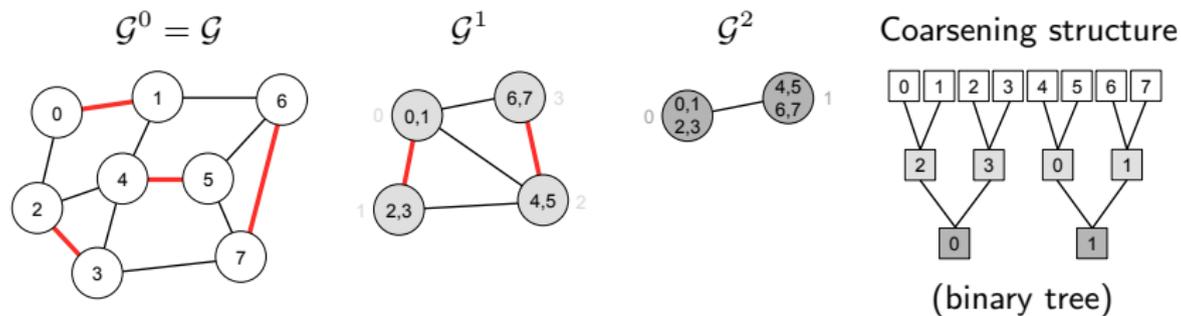
Example: citation networks

Method	Cora ¹	PubMed ²
Manifold Regularization ³	59.5%	70.7%
Semidefinite Embedding ⁴	59.0%	71.1%
Label Propagation ⁵	68.0%	63.0%
DeepWalk ⁶	67.2%	65.3%
Planetoid ⁷	75.7%	77.2%
Spectral graph CNN ⁸	81.6%	78.7%

Classification accuracy of different methods on citation network datasets

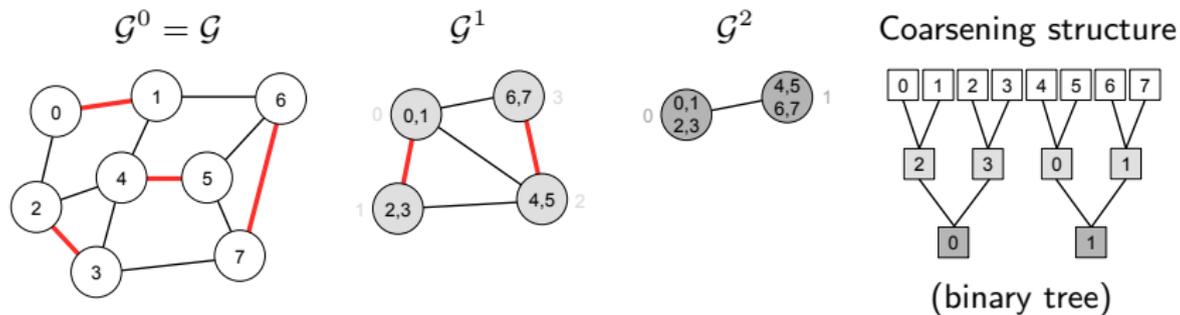
Data: ^{1,2}Sen et al. 2008; methods: ³Belkin et al. 2006; ⁴Weston et al. 2012; ⁵Zhu et al. 2003; ⁶Perozzi et al. 2014; ⁷Yang et al. 2016; ⁸Kipf, Welling 2016 (simplification of ChebNet)

Graph pooling



Produce a sequence of coarsened graphs

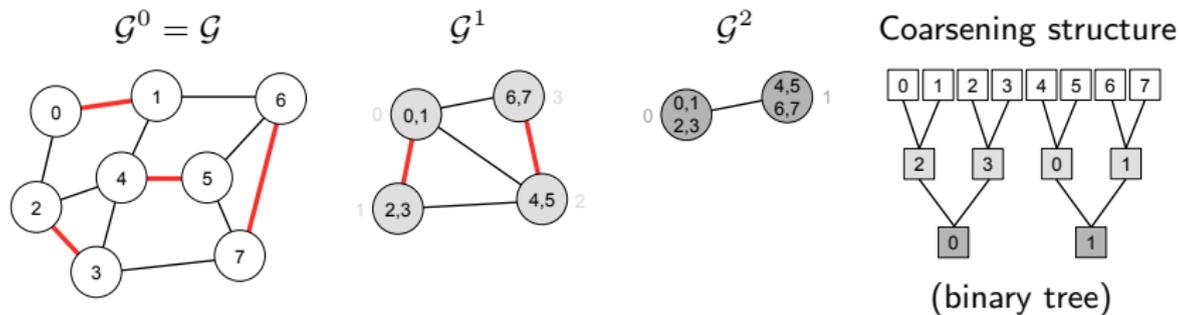
Graph pooling



Produce a sequence of coarsened graphs

Max or average pooling of collapsed vertices

Graph pooling



Produce a sequence of coarsened graphs

Max or average pooling of collapsed vertices

Binary tree arrangement of node indices

Limitations of spectral graph CNNs

- ☹️ Poor generalization across domains with different shapes unless kernels are localized

Limitations of spectral graph CNNs

- ☹️ **Poor generalization across domains with different shapes** unless kernels are localized (can be remedied to some extent with spectral transformer networks¹)

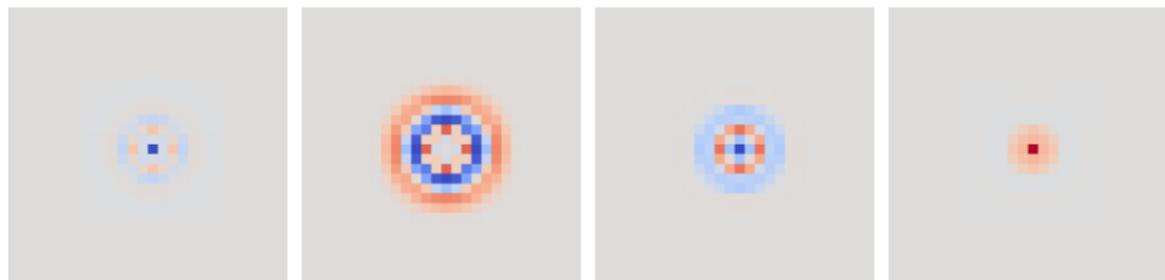
¹Yi, Su, Guo, Guibas 2017

Limitations of spectral graph CNNs

- ☹️ **Poor generalization across domains with different shapes** unless kernels are localized (can be remedied to some extent with spectral transformer networks¹)
- ☹️ **Spectral kernels are isotropic** due to rotation invariance of the Laplacian

¹Yi, Su, Guo, Guibas 2017

Only rotationally-symmetric kernels!



Example of Chebyshev filters (order $r = 7$) on Euclidean grid

Anisotropic kernels on manifolds

Scale t

Orientation θ

Elongation α

Examples of anisotropic heat kernels on a manifold

Limitations of spectral graph CNNs

- ☹️ **Poor generalization across domains with different shapes** unless kernels are localized (can be remedied to some extent with spectral transformer networks¹)
- ☹️ **Spectral kernels are isotropic** due to rotation invariance of the Laplacian (on manifolds, anisotropic Laplacians can be constructed²)

¹Yi, Su, Guo, Guibas 2017; ²Boscaini, Masci, Rodolà, Bronstein, Cremers 2016

Limitations of spectral graph CNNs

- ☹️ **Poor generalization across domains with different shapes** unless kernels are localized (can be remedied to some extent with spectral transformer networks¹)
- ☹️ **Spectral kernels are isotropic** due to rotation invariance of the Laplacian (on manifolds, anisotropic Laplacians can be constructed²)
- ☹️ **Only undirected graphs**, as symmetry of the Laplacian matrix is assumed

¹Yi, Su, Guo, Guibas 2017; ²Boscaini, Masci, Rodolà, Bronstein, Cremers 2016

Different formulations of non-Euclidean CNNs



Spectral domain



Spatial domain

Different formulations of non-Euclidean CNNs



Spectral domain



Spatial domain

Geometric Deep Learning Tutorial: NIPS 2017

<https://vimeo.com/248497329>

Summary

The convolution theorem allows us to generalize convolution operators from **grids** to **graphs**

Summary

The convolution theorem allows us to generalize convolution operators from **grids** to **graphs**

Naive spectral graph CNNs are unstable to deformation

Summary

The convolution theorem allows us to generalize convolution operators from **grids** to **graphs**

Naive spectral graph CNNs are unstable to deformation

Forcing the filter spectrum to be smooth stabilizes and localizes filters

Summary

The convolution theorem allows us to generalize convolution operators from **grids** to **graphs**

Naive spectral graph CNNs are unstable to deformation

Forcing the filter spectrum to be smooth stabilizes and localizes filters

Pooling operation can be replaced with graph pooling

Inference in Spectral Learning with Deep Networks

Manifold learning (again)

Basic pattern:

Construct Gram matrix from kernel $k(\mathbf{x}, \mathbf{x}')$

$$\mathbf{M} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

Use top/bottom eigenvectors of Gram matrix as embedding

Use Nyström approximation for inference on held-out data:

$$\phi_k(\mathbf{x}') \propto \sum_i \phi_{ki} k(\mathbf{x}_i, \mathbf{x}')$$

Manifold learning (again)

Basic pattern:

Construct Gram matrix from kernel $k(\mathbf{x}, \mathbf{x}')$

$$\mathbf{M} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

Use top/bottom eigenvectors of Gram matrix as embedding

Use Nyström approximation for inference on held-out data:

$$\phi_k(\mathbf{x}') \propto \sum_i \phi_{ki} k(\mathbf{x}_i, \mathbf{x}')$$

Why not just learn parameterized $\phi_k(\mathbf{x})$ directly?

Deep learn all the things

A successful strategy: take a branch of machine learning and fit the model using a deep network

Deep learn all the things

A successful strategy: take a branch of machine learning and fit the model using a deep network

Variational Bayes: Variational autoencoders

Deep learn all the things

A successful strategy: take a branch of machine learning and fit the model using a deep network

Variational Bayes: Variational autoencoders

Reinforcement Learning: Deep RL (DQN, A3C, PPO, TRPO, DDPG...)

Deep learn all the things

A successful strategy: take a branch of machine learning and fit the model using a deep network

Variational Bayes: Variational autoencoders

Reinforcement Learning: Deep RL (DQN, A3C, PPO, TRPO, DDPG...)

Spectral Learning: ????

Deep learn all the things

A successful strategy: take a branch of machine learning and fit the model using a deep network

Variational Bayes: Variational autoencoders

Reinforcement Learning: Deep RL (DQN, A3C, PPO, TRPO, DDPG...)

Spectral Learning: **Spectral Inference Networks**

Trade-offs of manifold learning

- 😊 Exactly solvable by eigendecomposition
- 😊 Data efficient (works with $\mathcal{O}(1000)$ data points)
- 😐 Unsupervised learning without a generative model
- 😞 Learning scales as $\mathcal{O}(n \log n)$ with size of training data
- 😞 Inference scales as $\mathcal{O}(n)$ with size of training data
- 😞 Performance degrades for noisy or clustered data
- 😊 Collapsed embeddings can be fixed by choice of eigenvector
- 😊 Can discover topology of data without prior assumptions

Trade-offs of manifold learning

- ☺ **Exactly solvable by eigendecomposition**
- ☺ **Data efficient (works with $\mathcal{O}(1000)$ data points)**
- ☹ Unsupervised learning without a generative model
- ☹ **Learning scales as $\mathcal{O}(n \log n)$ with size of training data**
- ☹ **Inference scales as $\mathcal{O}(n)$ with size of training data**
- ☹ Performance degrades for noisy or clustered data
- ☺ Collapsed embeddings can be fixed by choice of eigenvector
- ☺ Can discover topology of data without prior assumptions

Spectral inference networks trade the **first two** for the **second two**

Eigendecomposition as optimization

Rayleigh quotient has top eigenvector as argmax:

$$\max_{\phi} \frac{\phi^T \mathbf{A} \phi}{\phi^T \phi}$$

Eigendecomposition as optimization

Rayleigh quotient has top eigenvector as argmax:

$$\max_{\phi} \frac{\phi^T \mathbf{A} \phi}{\phi^T \phi}$$

Generalize to multiple eigenvectors (up to rotation):

$$\max_{\Phi} \text{Tr} \left((\Phi^T \Phi)^{-1} \Phi^T \mathbf{A} \Phi \right)$$

Eigendecomposition as optimization

Rayleigh quotient has top eigenvector as argmax:

$$\max_{\phi} \frac{\phi^T \mathbf{A} \phi}{\phi^T \phi}$$

Generalize to multiple eigenvectors (up to rotation):

$$\max_{\Phi} \text{Tr} \left((\Phi^T \Phi)^{-1} \Phi^T \mathbf{A} \Phi \right)$$
$$\max_{\Phi} \text{Tr} \left(\left(\sum_i \phi_i^T \phi_i \right)^{-1} \sum_{ij} A_{ij} \phi_i^T \phi_j \right)$$

Eigendecomposition as optimization

Rayleigh quotient has top eigenvector as argmax:

$$\max_{\phi} \frac{\phi^T \mathbf{A} \phi}{\phi^T \phi}$$

Generalize to multiple eigenvectors (up to rotation):

$$\max_{\Phi} \text{Tr} \left((\Phi^T \Phi)^{-1} \Phi^T \mathbf{A} \Phi \right)$$

$$\max_{\Phi} \text{Tr} \left(\mathbb{E}_{\mathbf{x}} [\phi(\mathbf{x})\phi(\mathbf{x})^T]^{-1} \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [k(\mathbf{x}, \mathbf{x}')\phi(\mathbf{x})\phi(\mathbf{x}')^T] \right)$$

Replace A_{ij} with $k(\mathbf{x}, \mathbf{x}')$ and sums with expectations

Eigendecomposition as optimization

Most machine learning:

$$\max_{\theta} \mathbb{E}_{\mathbf{x}}[f_{\theta}(\mathbf{x})]$$

Empirical gradient in **unbiased**

Eigendecomposition as optimization

Most machine learning:

$$\max_{\theta} \mathbb{E}_{\mathbf{x}} [f_{\theta}(\mathbf{x})]$$

Empirical gradient in **unbiased**

Spectral inference networks:

$$\max_{\theta} \text{Tr} \left(\mathbb{E}_{\mathbf{x}} [\phi_{\theta}(\mathbf{x})\phi_{\theta}(\mathbf{x})^T]^{-1} \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [k(\mathbf{x}, \mathbf{x}')\phi_{\theta}(\mathbf{x})\phi_{\theta}(\mathbf{x}')^T] \right)$$

Empirical gradient is **biased**

Eigendecomposition as optimization

Most machine learning:

$$\max_{\theta} \mathbb{E}_{\mathbf{x}} [f_{\theta}(\mathbf{x})]$$

Empirical gradient in **unbiased**

Spectral inference networks:

$$\max_{\theta} \text{Tr} \left(\mathbb{E}_{\mathbf{x}} [\phi_{\theta}(\mathbf{x})\phi_{\theta}(\mathbf{x})^T]^{-1} \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [k(\mathbf{x}, \mathbf{x}')\phi_{\theta}(\mathbf{x})\phi_{\theta}(\mathbf{x}')^T] \right)$$

Empirical gradient is **biased**

Solution: use moving average of gradient of $\phi_{\theta}\phi_{\theta}^T$ term

Spectral inference network gradients

Objective is of the form $\text{Tr}(\mathbf{\Sigma}^{-1}\mathbf{\Pi})$

Spectral inference network gradients

Objective is of the form $\text{Tr}(\mathbf{\Sigma}^{-1}\mathbf{\Pi})$

Gradient is of the form $\text{Tr}(\mathbf{\Sigma}^{-1}\nabla_{\theta}\mathbf{\Pi}) - \text{Tr}(\mathbf{\Sigma}^{-1}\mathbf{\Pi}\mathbf{\Sigma}^{-1}\nabla_{\theta}\mathbf{\Sigma})$

Spectral inference network gradients

Objective is of the form $\text{Tr}(\mathbf{\Sigma}^{-1}\mathbf{\Pi})$

Gradient is of the form $\text{Tr}(\mathbf{\Sigma}^{-1}\nabla_{\theta}\mathbf{\Pi}) - \text{Tr}(\mathbf{\Sigma}^{-1}\mathbf{\Pi}\mathbf{\Sigma}^{-1}\nabla_{\theta}\mathbf{\Sigma})$

To break symmetry between eigenfunctions, use gradient

$$\text{Tr}(\mathbf{L}^{-T}\text{diag}(\mathbf{L})^{-1}\nabla_{\theta}\mathbf{\Pi}) - \text{Tr}(\mathbf{L}^{-T}\text{triu}(\mathbf{\Lambda}\text{diag}(\mathbf{L})^{-1})\nabla_{\theta}\mathbf{\Sigma})$$

where \mathbf{L} is Cholesky decomposition of $\mathbf{\Pi}$ and $\mathbf{\Lambda} = \mathbf{L}^{-1}\mathbf{\Pi}\mathbf{L}^{-T}$

Spectral inference network gradients

Objective is of the form $\text{Tr}(\mathbf{\Sigma}^{-1}\mathbf{\Pi})$

Gradient is of the form $\text{Tr}(\mathbf{\Sigma}^{-1}\nabla_{\theta}\mathbf{\Pi}) - \text{Tr}(\mathbf{\Sigma}^{-1}\mathbf{\Pi}\mathbf{\Sigma}^{-1}\nabla_{\theta}\mathbf{\Sigma})$

To break symmetry between eigenfunctions, use gradient

$$\text{Tr}(\mathbf{L}^{-T}\text{diag}(\mathbf{L})^{-1}\nabla_{\theta}\mathbf{\Pi}) - \text{Tr}(\mathbf{L}^{-T}\text{triu}(\mathbf{\Lambda}\text{diag}(\mathbf{L})^{-1})\nabla_{\theta}\mathbf{\Sigma})$$

where \mathbf{L} is Cholesky decomposition of $\mathbf{\Pi}$ and $\mathbf{\Lambda} = \mathbf{L}^{-1}\mathbf{\Pi}\mathbf{L}^{-T}$

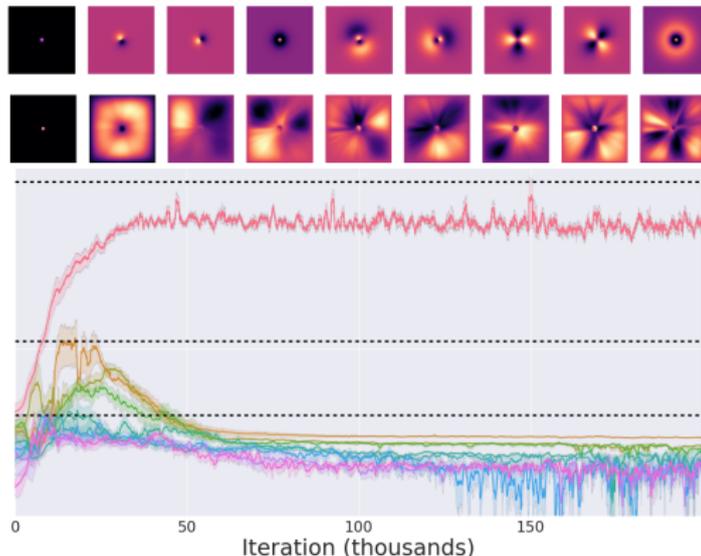
To reduce bias in the gradient, use moving average for $\mathbf{\Sigma}$ and $\nabla_{\theta}\mathbf{\Sigma}$

Spectral inference networks

Sanity check: the Schrödinger equation

$$E\psi(\mathbf{x}) = -\frac{\hbar^2}{2m}\nabla^2\psi(\mathbf{x}) - \frac{\psi(\mathbf{x})}{|\mathbf{x}|}$$

Without bias correction:

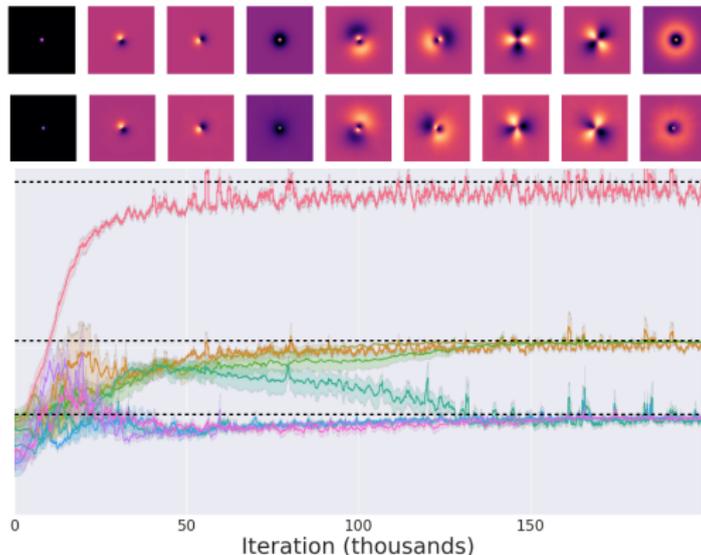


Spectral inference networks

Sanity check: the Schrödinger equation

$$E\psi(\mathbf{x}) = -\frac{\hbar^2}{2m}\nabla^2\psi(\mathbf{x}) - \frac{\psi(\mathbf{x})}{|\mathbf{x}|}$$

With bias correction:



Approximating Laplacian eigenmaps

Objective:

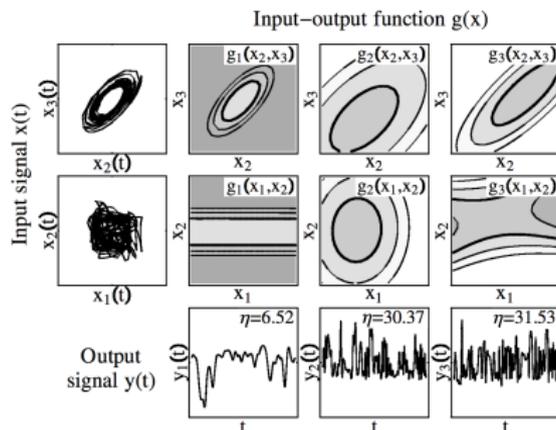
$$k(\mathbf{x}, \mathbf{x}')\phi(\mathbf{x})\phi(\mathbf{x})^T = (\phi(\mathbf{x}) - \phi(\mathbf{x}'))(\phi(\mathbf{x}) - \phi(\mathbf{x}'))^T$$

Approximating Laplacian eigenmaps

Objective:

$$k(\mathbf{x}, \mathbf{x}') \phi(\mathbf{x}) \phi(\mathbf{x})^T = (\phi(\mathbf{x}) - \phi(\mathbf{x}'))(\phi(\mathbf{x}) - \phi(\mathbf{x}'))^T$$

If \mathbf{x}, \mathbf{x}' are sequential video frames, equivalent to **Slow Feature Analysis**

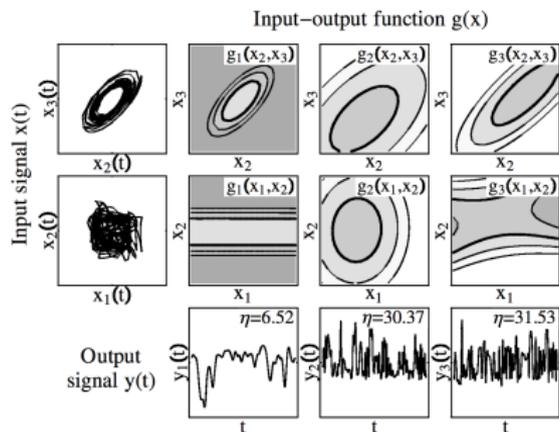


Approximating Laplacian eigenmaps

Objective:

$$k(\mathbf{x}, \mathbf{x}') \phi(\mathbf{x}) \phi(\mathbf{x}')^T = (\phi(\mathbf{x}) - \phi(\mathbf{x}'))(\phi(\mathbf{x}) - \phi(\mathbf{x}'))^T$$

If \mathbf{x}, \mathbf{x}' are sequential video frames, equivalent to **Slow Feature Analysis**



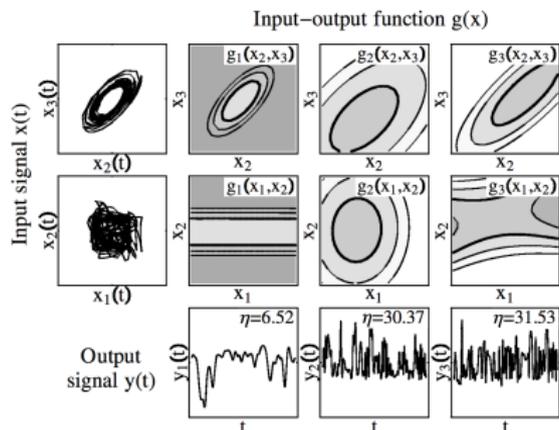
SFA learned **layer-by-layer** rather than **end-to-end**

Approximating Laplacian eigenmaps

Objective:

$$k(\mathbf{x}, \mathbf{x}') \phi(\mathbf{x}) \phi(\mathbf{x})^T = (\phi(\mathbf{x}) - \phi(\mathbf{x}'))(\phi(\mathbf{x}) - \phi(\mathbf{x}'))^T$$

If \mathbf{x}, \mathbf{x}' are sequential video frames, equivalent to **Slow Feature Analysis**

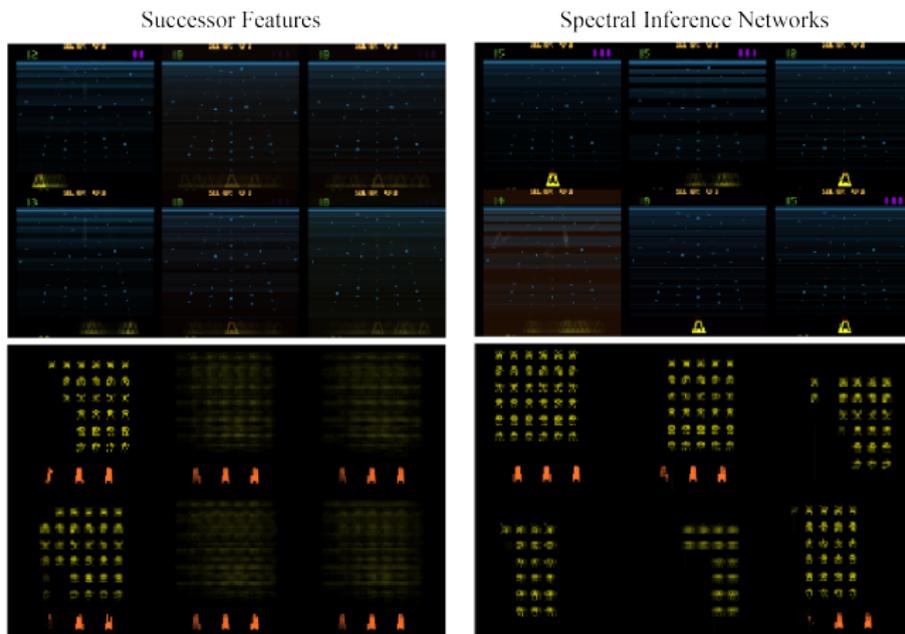


SFA learned **layer-by-layer** rather than **end-to-end**

SFA learned **feature-by-feature** rather than **fully online**

Spectral inference networks on Atari

More interpretable features compared to other approaches when trained on random policies on Atari games



Summary

Eigenfunctions can be learned by stochastic gradient descent with function approximation

Summary

Eigenfunctions can be learned by stochastic gradient descent with function approximation

Slow feature analysis is a special case of Spectral Inference Networks

Summary

Eigenfunctions can be learned by stochastic gradient descent with function approximation

Slow feature analysis is a special case of Spectral Inference Networks

While less efficient than standard spectral algorithms, it is far more scalable

Future challenges in manifold and spectral learning

Can we encode more **structural assumptions** into the choice of kernel?

Future challenges in manifold and spectral learning

Can we encode more **structural assumptions** into the choice of kernel?

Can we combine the **speed and efficiency** of nonparametric methods with the **scalability** of parametric methods?

Future challenges in manifold and spectral learning

Can we encode more **structural assumptions** into the choice of kernel?

Can we combine the **speed and efficiency** of nonparametric methods with the **scalability** of parametric methods?

Can we use curvature in observation space as a **learning signal** rather than a post-hoc analysis tool?

Future challenges in manifold and spectral learning

Can we encode more **structural assumptions** into the choice of kernel?

Can we combine the **speed and efficiency** of nonparametric methods with the **scalability** of parametric methods?

Can we use curvature in observation space as a **learning signal** rather than a post-hoc analysis tool?

Can we use the learned representations for challenging downstream tasks?

Future challenges in manifold and spectral learning

Can we encode more **structural assumptions** into the choice of kernel?

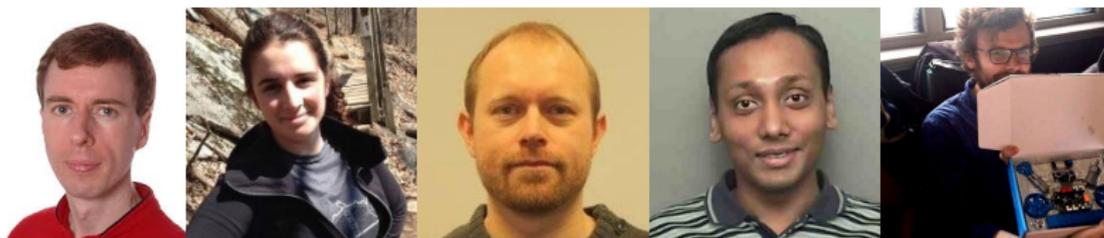
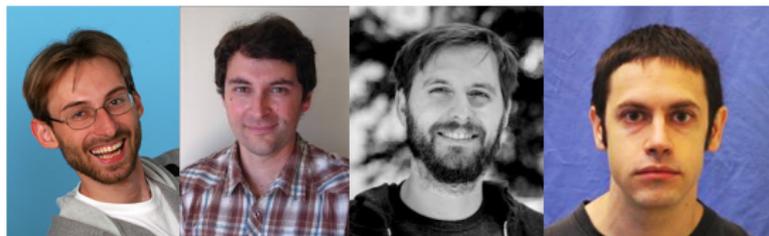
Can we combine the **speed and efficiency** of nonparametric methods with the **scalability** of parametric methods?

Can we use curvature in observation space as a **learning signal** rather than a post-hoc analysis tool?

Can we use the learned representations for challenging downstream tasks?

Can we better connect spectral learning with probabilistic models?

Acknowledgements



Michael Bronstein, Xavier Bresson, Joan Bruna, Arthur Szlam, David Barrett,
Kim Stachenfeld, Stig Petersen, Ashish Agarwal, Chris Burgess

Get in touch: pfau@google.com, @pfau on Twitter

Bibliography

M. P. Do Carmo, "Differential Geometry of Curves and Manifolds", **Classic textbook on differential geometry**

R. Kimmel and J. A. Sethian, "Computing Geodesic Paths on Manifolds", PNAS 95(15): 8431-8435, 1998.

S. Roweis and L. K. Saul, "Nonlinear Dimensionality Reduction by Locally Linear Embedding", Science 290(5500): 2323-2326, 2000.

J. B. Tenenbaum, V. De Silva, J. C. Langford, "A Global Geometric Framework for Nonlinear Dimensionality Reduction", Science 290(5500): 2319-2323, 2000.

M. Belkin and P. Niyogi, "Laplacian Eigenmaps for Dimensionality Reduction and Data Representation", Neural Computation 15(6): 1373-1396, 2002.

T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado and J. Dean, "Distributed Representation of Words and Phrases and Their Compositionality", Adv. NIPS: 3111-3119, 2013. **Deep learning of high-dimensional word embeddings with skip-grams**

Bibliography

A. Y. Ng, M. I. Jordan, Y. Weiss, "On Spectral Clustering: Analysis and An Algorithm", Adv. NIPS: 849-856, 2002.

J. Shi and J. Malik, "Normalized Cuts and Image Segmentation", IEEE TPAMI 22(8): 888-905, 2000.

C. Ionescu, O. Vantzos, C. Sminchisescu, "Matrix Backpropagation for Deep Networks with Structured Data", Proc. ICCV: 2965-2973, 2015. **Use normalized cuts as layer inside deep network**

Y. Bengio, J.-F. Paiement, P. Vincent, O. Delalleau, N. Le Roux, M. Ouimet, "Out-of-Sample Extensions for LLE, IsoMap, MDS, Eigenmaps and Spectral Clustering", Adv. NIPS: 177-184, 2004.

L. van der Maaten and G. Hinton, "Visualizing Data using t-SNE", JMLR 9(11): 2579-2605, 2008.

R. Hadsell, S. Chopra and Y. LeCun, "Dimensionality Reduction by Learning an Invariant Mapping", CVPR: 1735-1742, 2006.

Bibliography

D. Pfau and C. P. Burgess, "Minimally Redundant Laplacian Eigenmaps", ICLR Workshops, 2018.

M. Nickel and D. Kiela, "Poincaré Embeddings for Learning Hierarchical Representations", Adv. NIPS: 6341-6350, 2017.

M. Nickel and D. Kiela, "Learning Continuous Hierarchies in the Lorentz Model of Hyperbolic Geometry", ICML 2018.

S. Bonnabel, "Stochastic Gradient Descent on Riemannian Manifolds", IEEE Trans. Automatic Control 58(9): 2217-2229, 2013.

C. Gulcehre, M. Denil, M. Malinowski, A. Razavi, R. Pascanu, K. M. Hermann, P. Battaglia, V. Bapst, D. Raposo, A. Santoro, N. de Freitas, "Hyperbolic Attention Networks", arXiv:1805.09786, 2018.

I. Higgins, L. Matthey, A. Pal, C. P. Burgess, X. Glorot, M. Botvinick, S. Mohamed, A. Lerchner, " β -VAE: Learning Basic Visual Concepts with a Constrained Variational Framework", Proc. ICLR, 2017.

Bibliography

G. Arvanitidis, L. K. Hansen, S. Hauberg, “Latent Space Oddity: on the Curvature of Deep Generative Models”, ICLR 2018.

B. Boots, S. M. Siddiqi, G. J. Gordon, “Closing the Learning-Planning Loop with Predictive State Representations”, IJRR 30(7): 954-966, 2011. **Use spectral learning for planning and control**

M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, P. Vandergheynst, “Geometric deep learning: going beyond Euclidean data”, arXiv:1611.08097, 2016. **First review paper of geometric deep learning**

J. Bruna, W. Zaremba, A. Szlam, Y. LeCun, “Spectral networks and locally connected networks on graphs”, *Proc. ICML* 2014. **First Spectral CNN on graphs**

M. Henaff, J. Bruna, and Y. LeCun, “Deep convolutional networks on graph-structured data”, arXiv:1506.05163, 2015. **Spectral CNN with smooth multipliers**

Bibliography

M. Defferrard, X. Bresson, P. Vandergheynst, “Convolutional neural networks on graphs with fast localized spectral filtering”, *Proc. NIPS* 2016. **Spectral CNN with Chebyshev polynomial filters (ChebNet)**

T. N. Kipf, M. Welling, “Semi-supervised classification with graph convolutional networks”, arXiv:1609.02907, 2016. **Graph convolutional network (GCN) framework, a simplification of ChebNet**

F. Monti, D. Boscaini, J. Masci, E. Rodolà, J. Svoboda, M. M. Bronstein, “Geometric deep learning on graphs and manifolds using mixture model CNNs”, *Proc. CVPR* 2017. **Mixture Model network (MoNet) framework**

D. Boscaini, J. Masci, E. Rodolà, M. M. Bronstein, D. Cremers, “Anisotropic diffusion descriptors”, *Computer Graphics Forum* 35(2):431–441, 2016. **Anisotropic heat kernels**

L. Yi, H. Su, X. Guo, L. Guibas, “SyncSpecCNN: Synchronized Spectral CNN for 3D Shape Segmentation”, *Proc. CVPR* 2017. **Spectral transformer networks**

Bibliography

Y. LeCun, B. Boser, J. S. Denker, R. E. Howard, W. Hubbard, L. D. Jackel, "Backpropagation applied to handwritten ZIP code recognition", *Neural Computation* 1(4):541–551, 1989. **Classical Euclidean CNN**

D. Pfau, S. Petersen, A. Agarwal, D. Barrett and K. Stachenfeld, "Spectral Inference Networks: Unifying Spectral Methods with Deep Learning", arXiv: 1806.02215, 2018.

L. Wiskott and T. J. Sejnowski, "Slow Feature Analysis: Unsupervised Learning of Invariances", *Neural Computation* 14(4): 715-770, 2002.

H. Sprekeler, "On the Relation of Slow Feature Analysis and Laplacian Eigenmaps", *Neural Computation* 23(12): 3287-3302, 2011. **Shows equivalence of SFA and Laplacian eigenmaps**