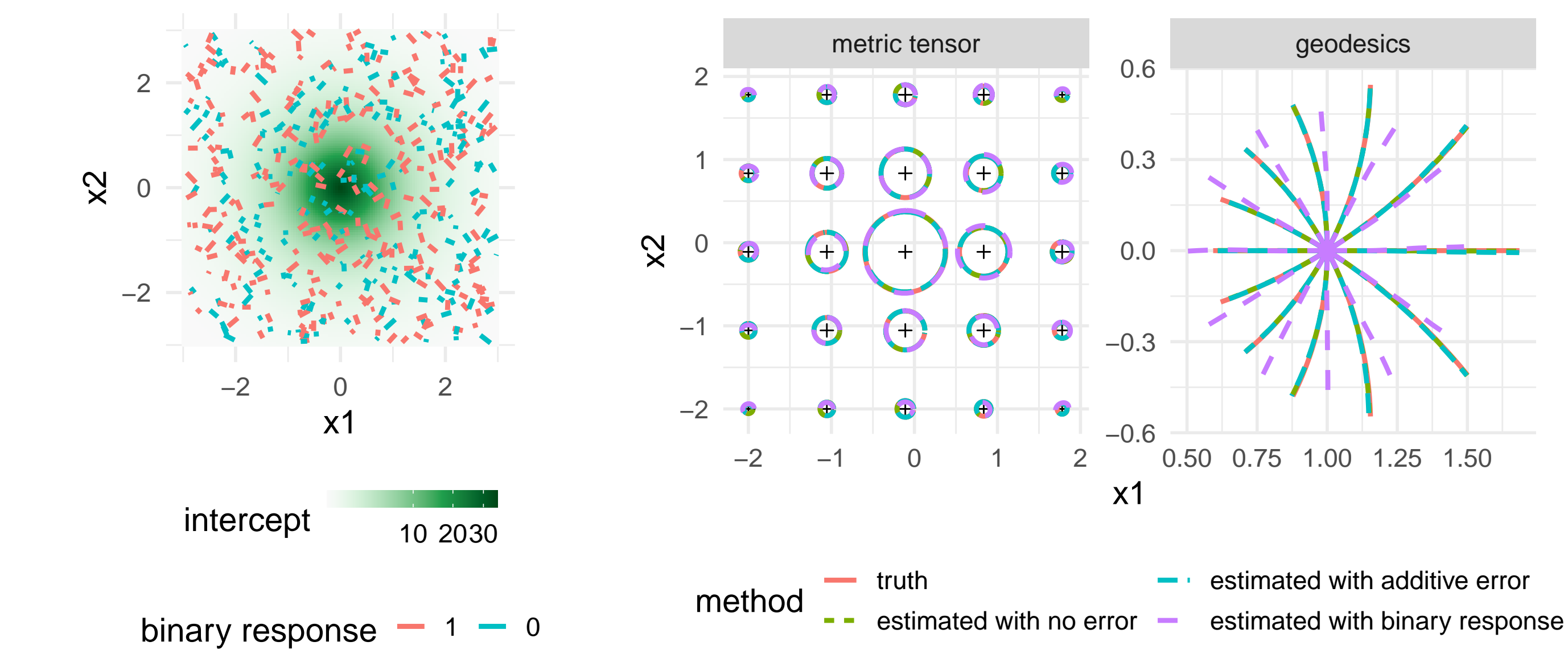


## We extend metric learning to consistently estimate Riemannian metric with noisy similarity measures.

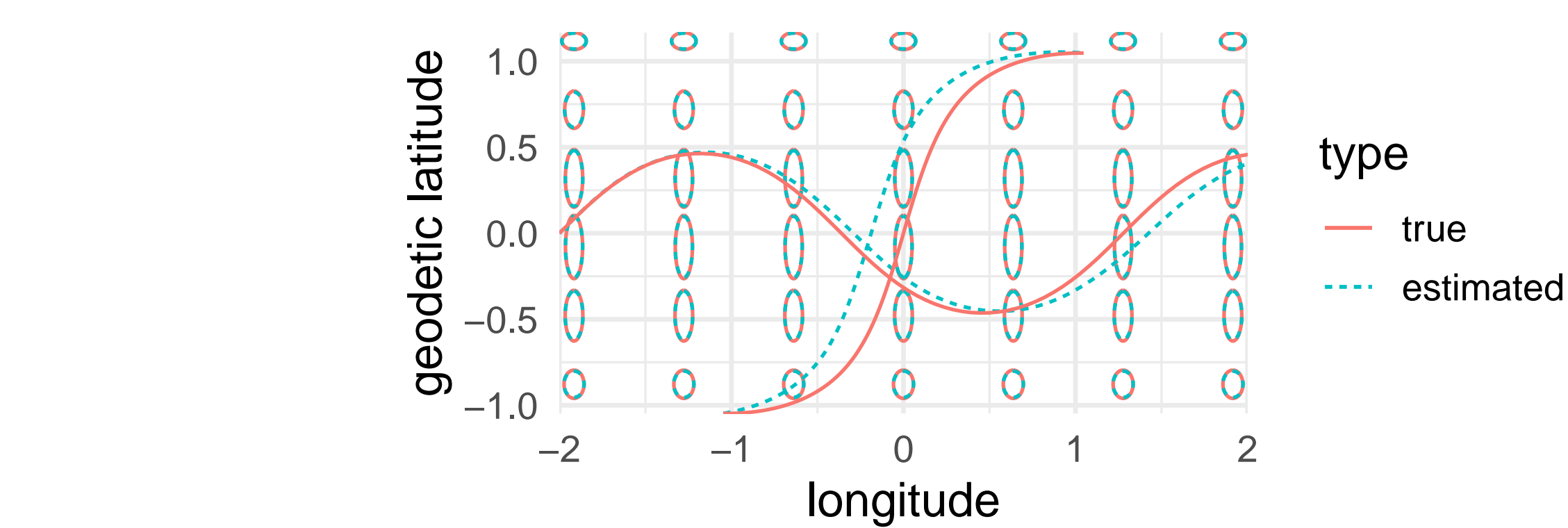
- Probabilistic modeling of pairwise similarity measures via intrinsic distances.
- Geometric interpretation for the similarity induced data space structure.
- Works for continuous, binary, or comparative similarity measures.
- Based on local regression and Taylor expansion for the squared distances.
- Justify the metric learning from a geometrical perspective.

### Simulation Examples

Binary and continuous similarity on 2D sphere with stereographic coordinates.



On a 2D ellipsoid: metric and geodesics.



Binary relative comparison (X is more like Y than Z) on double spirals.

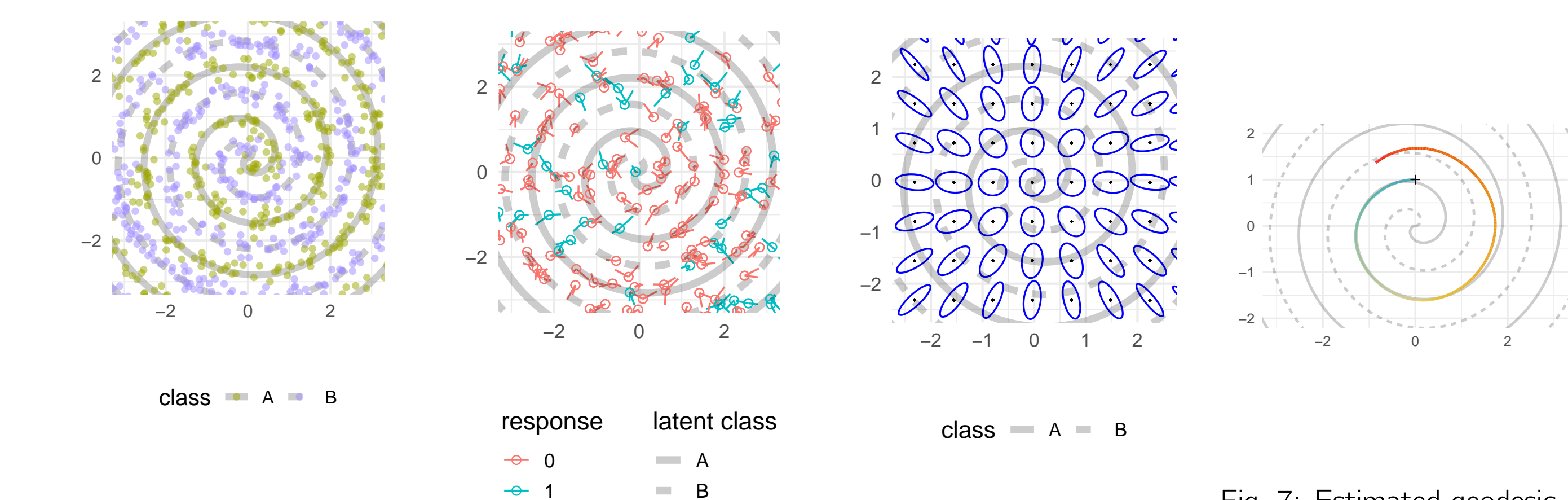


Fig. 4: The points and their latent classes

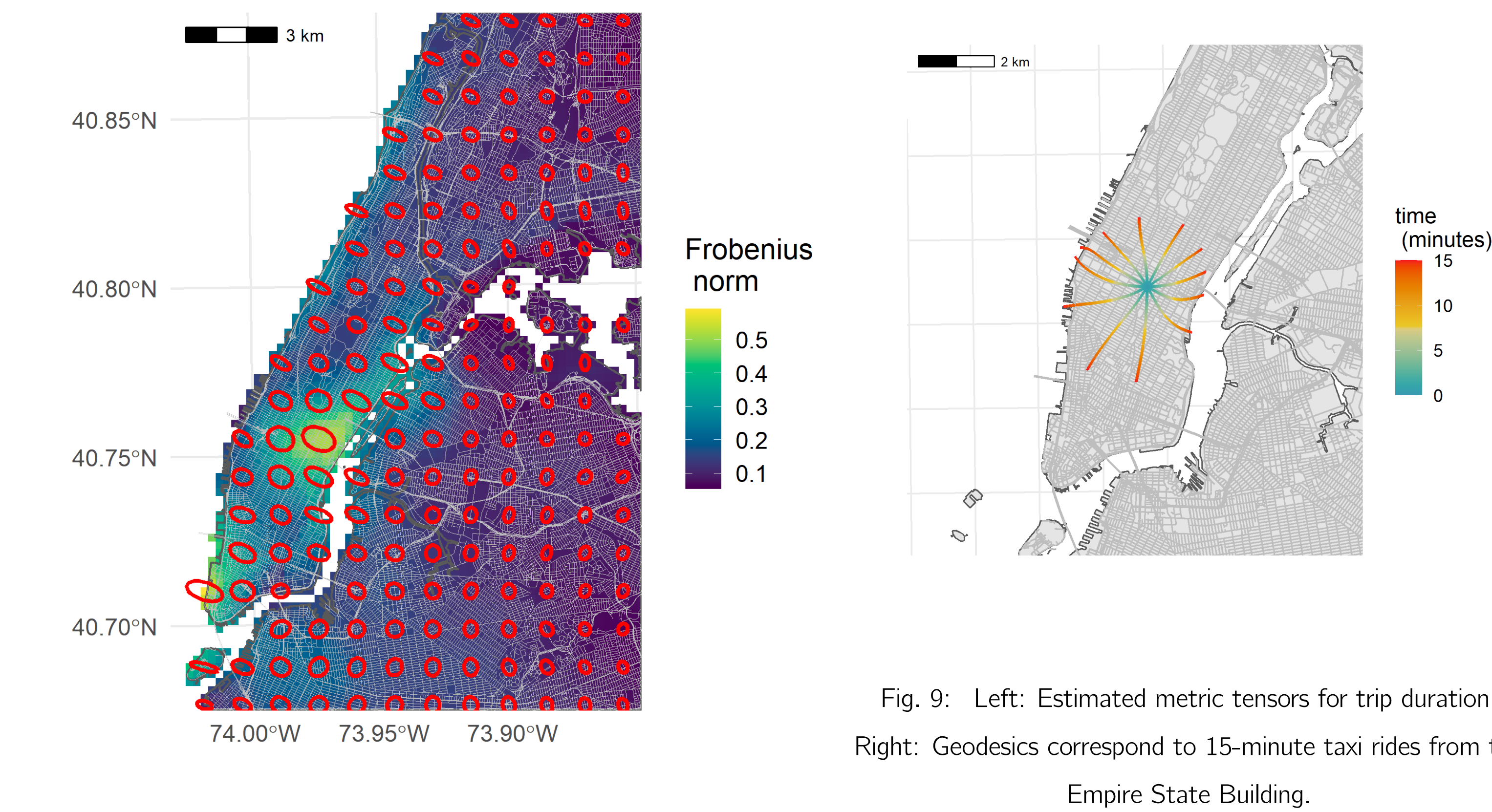
Fig. 5: Relative comparison.

Fig. 6: Estimated metric.

Fig. 7: Estimated geodesic.

### New York Travel Time Metric via Taxi Trips

$(X_{u0}, X_{u1}, Y_u)$  are pickup/dropoff locations, and trip duration.



- Traffic is slower along the longer axes of the cost ellipses. I.e., along the east–west direction (narrower streets) compared to the north–south direction (wider avenues).
- Brighter regions (midtown and the financial district) are more congested.

### MNIST Digits

- Embed MNIST images to a 2-dimensional space via tSNE.
- Similarity =  $C \text{dist}_{wass}(\text{pic}_{u0}, \text{pic}_{u1}) + 1_{\{|\text{bl}_{u0} \neq \text{bl}_{u1}\}}$  for some constant  $C$ .

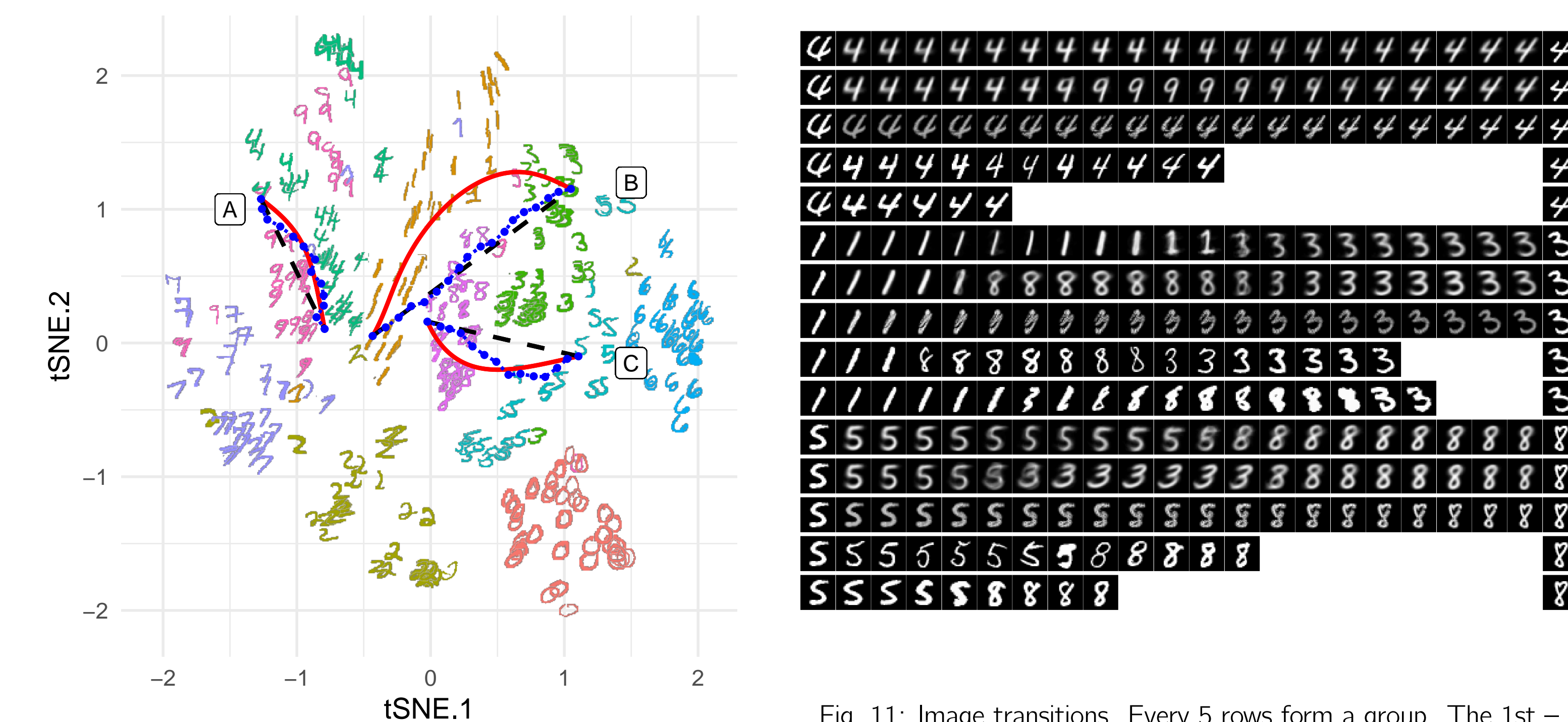


Fig. 10: The estimated geodesic curves (solid red), shortest path along kNN graph (dotted blue), and straight lines on the chart (dashed black).

- Similar if they are the same digit and look the same.
- Estimated geodesics try to stay within same digit class and go through similar images.

### A Geometric Perspective

- Observe  $N$  independent triplets  $(Y_u, X_{u0}, X_{u1})$ , where for  $u = 1, \dots, N$ ,
- $X_{uj}$  are points on the manifold  $\mathcal{M}$  with coordinates  $(X_{uj}^1, \dots, X_{uj}^d) \in \mathbb{R}^d$ .
- $Y_u$  measures similarity between  $X_{u0}$  and  $X_{u1}$ .
- Suppose  $\mathbb{E}(Y_u | X_{u0}, X_{u1}) = g^{-1}(\text{dist}(X_{u0}, X_{u1}))$  in a small neighborhood  $\mathcal{U}_p \subset \mathcal{M}$  of some target point  $p \in \mathcal{M}$  with some given link function  $g$ .

To connect squared distances to metric tensor and more:

- $\text{dist}(X_{u0}, X_{u1})^2 \approx \delta_{0-1}^i \delta_{0-1}^j G_{ij} + \delta_{0-1}^i (\delta_0^k \delta_0^l - \delta_1^k \delta_1^l) \Gamma_{kl}^i G_{ij}$ ,
- $\delta_0^i = X_{u0}^i - p^i$ ,  $\delta_1^i = X_{u1}^i - p^i$ , and  $\delta_{0-1}^i = \delta_0^i - \delta_1^i$  are differences in coordinates;
- $G_{ij}$  and  $\Gamma_{kl}^i$  are the metric tensor and Christoffel symbols at  $p$ .

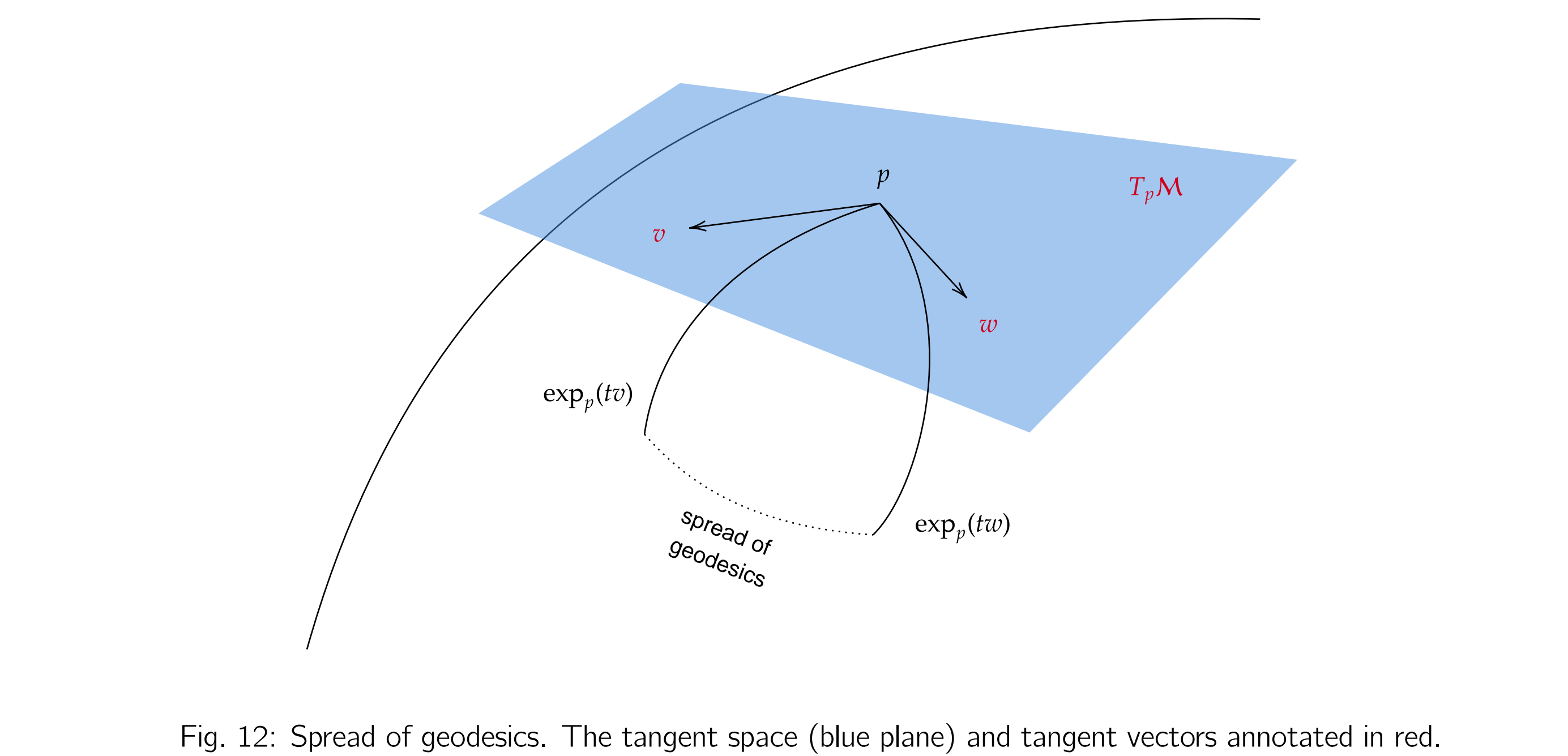


Fig. 12: Spread of geodesics. The tangent space (blue plane) and tangent vectors annotated in red.

Squared distance is approximated by a quadratic form (metric) plus a cubic form of the difference in coordinates. Higher order terms also possible (relate to curvature).

Motivates local regression estimation.

- Linear predictor:  $\eta_u := \beta^{(0)} + \delta_{u,0-1}^i \delta_{u,0-1}^j \beta_{ij}^{(1)} + \delta_{u,0-1}^k (\delta_{u0}^i \delta_{u0}^j - \delta_{u1}^i \delta_{u1}^j) \beta_{ijk}^{(2)}$ . The intercept  $\beta^{(0)}$  is optional.
- Find

$$(\hat{\beta}^{(0)}, \hat{\beta}_{ij}^{(1)}, \hat{\beta}_{ijk}^{(2)}) = \arg \min_{\beta^{(0)}, \beta_{ij}^{(1)}, \beta_{ijk}^{(2)}} \sum_{u=1}^N Q(Y_u, g^{-1}(\eta_u)) w_u,$$

subject to symmetric constraints  $\beta_{ij}^{(1)} = \beta_{ji}^{(1)}$ ,  $\beta_{ijk}^{(2)} = \beta_{jik}^{(2)}$  for  $i, j, k, l = 1, \dots, d$ . Here  $Q$  is a loss function (e.g., likelihood), and  $w_u$  are non-negative weights.

- $\hat{G}_{ij} = \hat{\beta}_{ij}^{(1)}$ ,  $\hat{\Gamma}_{ij}^l = \hat{\beta}_{ijk}^{(2)} \hat{G}^{kl}$  estimate the metric tensor and Christoffel symbols, where  $\hat{G}^{kl}$  is the matrix inverse of  $\hat{G}$  satisfying  $\hat{G}^{kl} \hat{G}_{kj} = 1_{\{j=l\}}$ .

**Proposition 1 (Consistency).** Suppose  $g(\mu) = \mu$  and  $Q(\mu, y) = (\mu - y)^2$ . Under conditions similar to those in a local regression setting,  $\text{bias}(\hat{\beta}(\mathbf{X})) = O_p(h^2)$ ,  $\text{var}(\hat{\beta}(\mathbf{X})) = O_p(N^{-1}h^{-4-2d})$ , as  $h \rightarrow 0$  and  $Nh^{2+2d} \rightarrow \infty$ , where  $\mathbf{X} = \{(X_{u0}, X_{u1})\}_{u=1}^N$  is the collection of observed endpoints.

