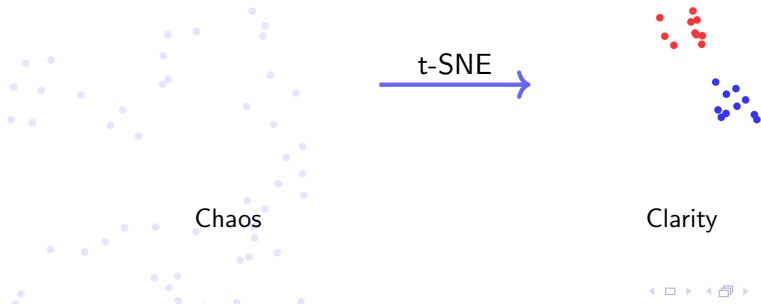


# t-SNE: Seeing the Invisible Structure

## Advanced Multivariate Analysis

Associate Professor Endri Raco  
Polytechnic University of Catalonia



# Our Journey Together

## This Session:

- Exchange of ideas
- Build intuition first
- Then mathematics
- Your questions welcome

## You'll Master:

- **Why** t-SNE works
- **When** to use it
- **How** to implement
- **Pitfalls** to avoid

*"By the end, you'll see data differently"*

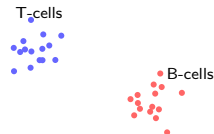
# The Data Visualization Challenge

## Single-cell RNA 20,000 dimensions!

Cell	Gene1	Gene2	...
1	0.23	1.45	...
2	0.67	0.89	...
3	1.23	0.02	...
...	...	...	...

How find cell types?

## After t-SNE:

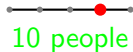


Cell types visible!

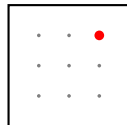
# The Curse: A Thought Experiment

## Finding Your Friend at a Concert

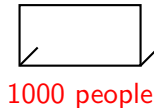
1D: Line



2D: Field



3D: Stadium



**In 20,000 dimensions?**

Your friend is *everywhere and nowhere*

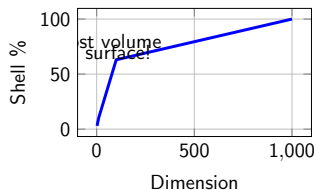
# The Curse: Mathematics

## Volume in n-dim sphere:

$$V_n(r) = \frac{\pi^{n/2}}{\Gamma(n/2 + 1)} r^n$$

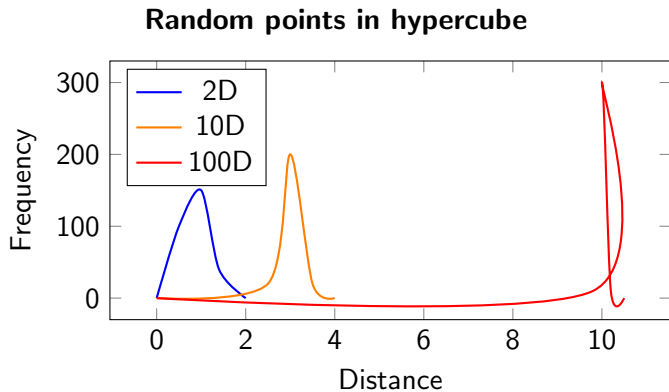
## Shell vs Core:

$$\frac{V_{shell}}{V_{total}} = 1 - (0.99)^n$$



**Key Insight:** All points become equally distant!

# Distance Collapse



**Problem:** No meaningful neighborhoods!

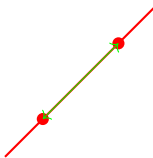
# Why PCA Fails

## The Swiss Roll Problem

True Structure



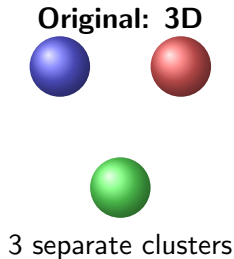
PCA Result



Wrong distance!

PCA assumes linear subspace - misses manifold structure

# MDS: The Crowding Problem

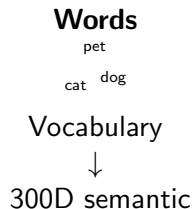
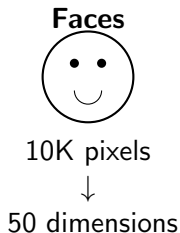
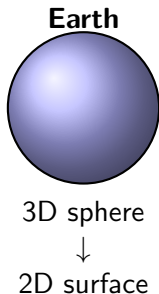


Not enough "room" in 2D for all distances



# The Manifold Hypothesis

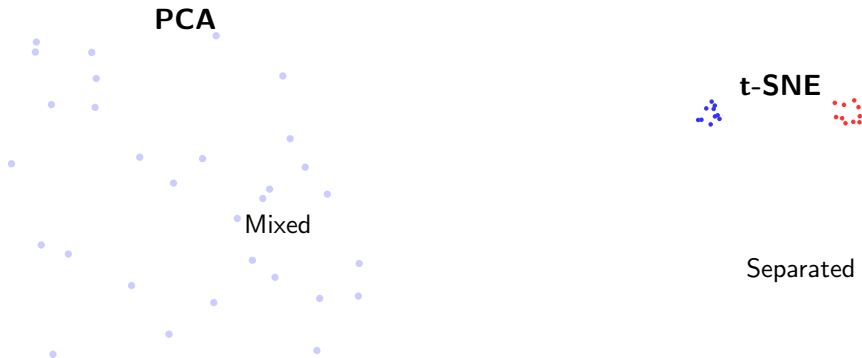
High-D data lies on low-D manifolds



**Key:** True complexity  $\ll$  apparent dimensions

# t-SNE's Revolutionary Promise

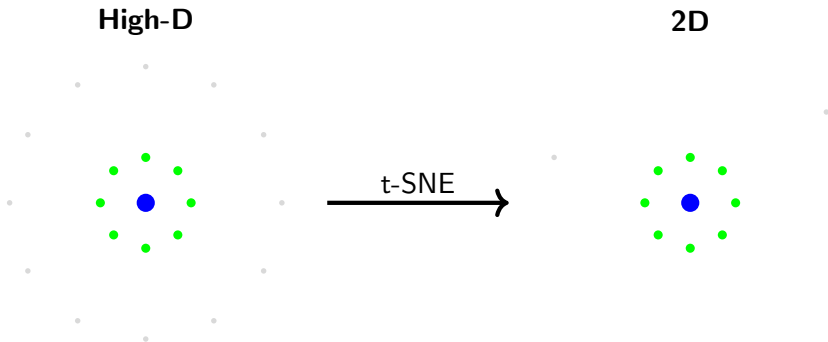
MNIST Digits: 784D  $\rightarrow$  2D



Key: Preserve neighborhoods, not distances

# The Fundamental Insight

## Local vs Global



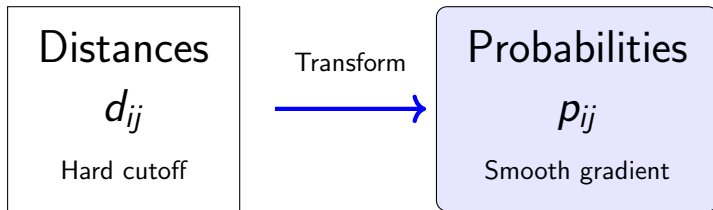
- ✓ Keep nearby points together
- ○ Let distant points reorganize

# Our Learning Journey



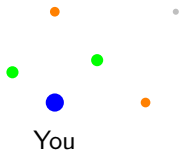
Questions welcome at each checkpoint!

### The Core Innovation



# The "Friends" Analogy

## Your Social Network



## Picking Probability

Best friend: 40%

Close friends: 30%

Acquaintances: 25%

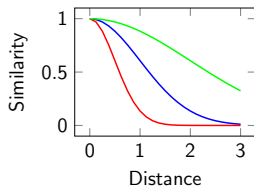
Strangers: 5%

Closer = Higher probability

t-SNE does this for **every** point!

# The Gaussian Kernel

$$p_{j|i} = \frac{\exp(-d_{ij}^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-d_{ik}^2/2\sigma_i^2)}$$



## Effect of $\sigma_i$ :

- **Small**: Very local
- **Medium**: Balanced
- **Large**: Global view

Each point gets its own  $\sigma_i$ !

# Building the Probabilities

## 1 Compute distances

$$d_{ij} = ||x_i - x_j||$$

## 2 Apply Gaussian

$$\tilde{p}_{j|i} = \exp(-d_{ij}^2/2\sigma_i^2)$$

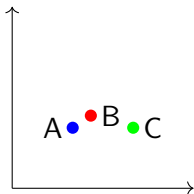
## 3 Normalize

$$p_{j|i} = \frac{\tilde{p}_{j|i}}{\sum_{k \neq i} \tilde{p}_{k|i}}$$

## 4 Interpretation Probability that $i$ picks $j$ as neighbor



## Example: 5 Points



**From A's perspective:**

Point	Distance	$p_{j A}$
B	0.36	51%
C	1.00	9%
Others	≥ 1.5	1%

B is A's primary neighbor

**Perplexity = "How many neighbors?"**

**Perp = 5**



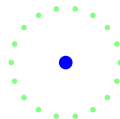
Very local

**Perp = 30**



Balanced

**Perp = 100**



Global

# Perplexity in Action

Same data, different perplexity

Perp = 5



Fragmented

Perp = 30



Just right

Perp = 100

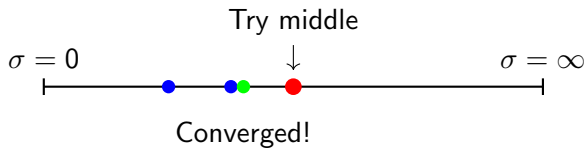


Merged

Rule: Perplexity between 5 and 50

# Finding the Right $\sigma_i$

## Binary Search for Each Point



- Target: perplexity  $\rightarrow$  entropy
- Adjust  $\sigma$  until match
- Converges in 10 iterations
- Do this for ALL  $n$  points!

# Perplexity: The Math

## Definition (Perplexity)

$$\text{Perp}(P_i) = 2^{H(P_i)}$$

where entropy:

$$H(P_i) = - \sum_j p_{j|i} \log_2 p_{j|i}$$

## Interpretation:

- Effective number of neighbors
- $\text{Perp} = 30 \approx 30$  equally likely neighbors
- Controls local vs global focus

**Binary search finds  $\sigma_i$  such that:**

$$\text{Perp}(P_i) = \text{user specified perplexity}$$

## Common Mistakes

Setting	Result	Fix
Perp = 2	Islands	Increase to 15+
Perp = 200	Blob	Decrease to 50
Perp $\propto$ n/3	Unstable	Use 5-50 range

## Best Practice:

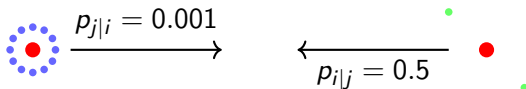
- Try multiple values (15, 30, 50)
- Look for stable patterns
- Consider data size: larger  $n \rightarrow$  larger perp OK

# The Asymmetry Problem

$$p_{j|i} \neq p_{i|j}$$

Dense region

Sparse region



Problem: Outliers pull but aren't pulled!

# Symmetrization Solution

## Simple Fix:

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$$

## Properties:

- Symmetric
- Sum to 1
- Fair to all points

## Example:

Before:

$$p_{j|i} = 0.001$$

$$p_{i|j} = 0.500$$

After:

$$p_{ij} = 0.250/n$$

$$p_{ji} = 0.250/n$$

Balanced!



# Mathematical Check

## Theorem (Joint Distribution)

With  $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$ :

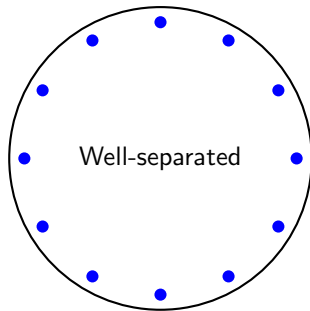
$$\sum_{i,j} p_{ij} = 1$$

**Proof:**

$$\begin{aligned}\sum_{i,j} p_{ij} &= \sum_{i,j} \frac{p_{j|i} + p_{i|j}}{2n} \\ &= \frac{1}{2n} \left[ \sum_{i,j} p_{j|i} + \sum_{i,j} p_{i|j} \right] \\ &= \frac{1}{2n} [n + n] = 1 \quad \checkmark\end{aligned}$$

### Why Gaussians Fail

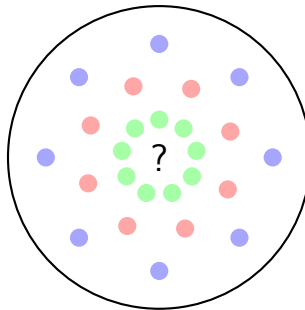
High-D Space



Project



2D Space



Not enough "room" in 2D!

# Volume Scaling Problem

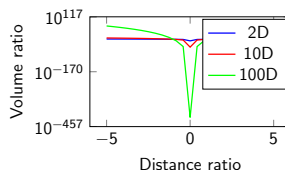
**Volume ratio:**

$$\frac{V_n(2r)}{V_n(r)} = 2^n$$

**Examples:**

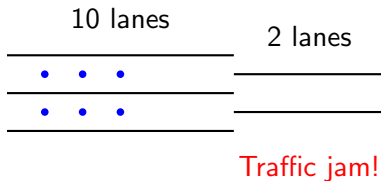
- 2D:  $2^2 = 4\times$
- 10D:  $2^{10} = 1024\times$
- 100D:  $2^{100} \approx 10^{30}\times$

**Problem:** Can't preserve moderate distances in 2D!



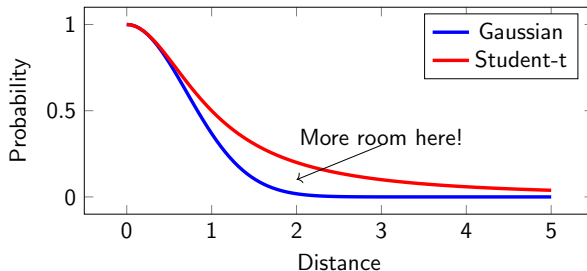
# The Traffic Jam

Trying to fit 10D structure in 2D



**Solution needed:** Different distance function in 2D

# The Solution: Heavy Tails



Heavy tails = more probability at moderate distances

# Why Student-t with $df=1$ ?

**The Choice:**  $q_{ij} \propto (1 + \|y_i - y_j\|^2)^{-1}$

**Why  $df=1$ ?**

- Heaviest tails
- Simple gradient
- Fast computation
- Works best!

**Tail Comparison:**

Distribution	Decay
Gaussian	$e^{-d^2}$
t(df=1)	$d^{-2}$
t(df=5)	$d^{-6}$

$df=1$  gives most room!

# Student-t Distribution: The Mathematics

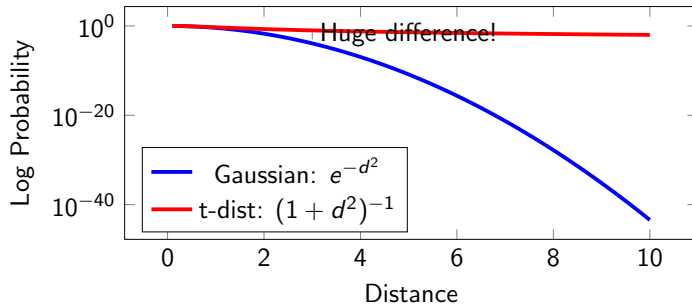
## Definition (Student-t with 1 degree of freedom)

In low dimensions, we use:  $q_{ij} = \frac{(1+\|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1+\|y_k - y_l\|^2)^{-1}}$

**This is the Cauchy distribution:**  $f(x) = \frac{1}{\pi(1+x^2)}$

**Key property:** Polynomial decay vs exponential

# Tail Behavior Analysis



At  $d = 3$ : Gaussian  $\approx 10^{-4}$ , Student-t  $\approx 0.1$



# The Elegant Gradient

## Theorem (t-SNE Gradient)

*With Student-t in low dimensions:*  $\frac{\partial C}{\partial y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j)(1 + \|y_i - y_j\|^2)^{-1}$

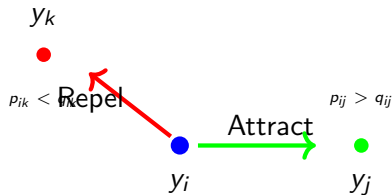
### Compare complexity: Gaussian:

- Compute  $e^{-\|y_i - y_j\|^2}$
- Expensive  $\exp()$
- Numerical issues

### Student-t:

- Compute  $(1 + \|y_i - y_j\|^2)^{-1}$
- Simple division
- Stable

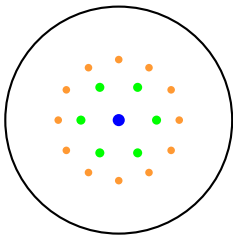
# Forces in t-SNE



**Spring analogy:**  $(p_{ij} - q_{ij}) = \text{spring tension}$

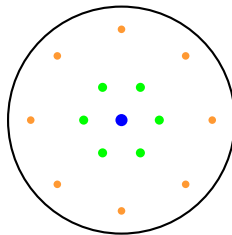
# How t-Distribution Solves Crowding

**Gaussian**



Crowded at moderate

**Student-t**



Room at moderate

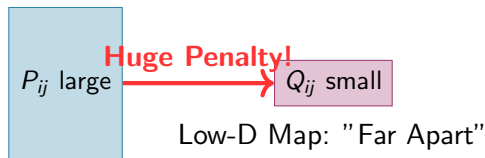
# Making t-SNE Work

## Key Components:

- 1 KL Divergence objective
- 2 Gradient descent
- 3 Momentum
- 4 Early exaggeration
- 5 Learning rate annealing

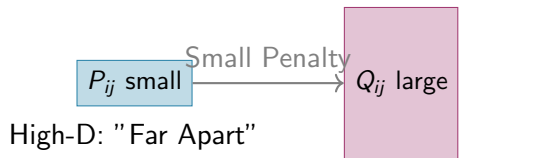
# The Objective: The "Cost of a Bad Map"

The goal is to make the low-D map ( $Q$ ) reflect the high-D reality ( $P$ ). The KL Divergence measures the "cost" or "penalty" for every point where the map is wrong.



High-D: "Close Neighbors"

t-SNE works hard to fix this  
(Pulls points together)



t-SNE doesn't worry much  
(Allows global changes)

**Insight: KL Divergence cares much more about keeping close points together than pushing far points apart**

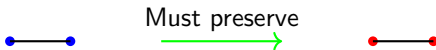
# Asymmetry in KL Divergence

Situation	Penalty	Effect
Large $p_{ij}$ , small $q_{ij}$	HIGH	Preserves local
Small $p_{ij}$ , large $q_{ij}$	low	Allows global flex

## Visual consequence:

High-D neighbors

Low-D



# Computing the Gradient

Starting from:  $C = \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}}$

Taking derivative w.r.t.  $y_i$ :  $\frac{\partial C}{\partial y_i} = 4 \sum_j (p_{ij} - q_{ij}) \cdot F_{ij}$

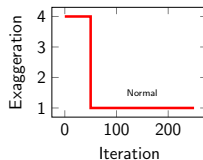
where:  $F_{ij} = \frac{(y_i - y_j)}{1 + ||y_i - y_j||^2}$

**Interpretation:** Weighted sum of forces from all points

# Early Exaggeration Trick

## Method:

- Multiply all  $p_{ij}$  by 4
- For first 50 iterations
- Creates tight clusters
- Separates clusters early



**Why it works:** Forces cluster formation before fine-tuning



# The Complete Algorithm

- 1: **Input:**  $X \in \mathbb{R}^{n \times d}$ , perplexity
- 2: **Output:**  $Y \in \mathbb{R}^{n \times 2}$
- 3:
- 4: Compute all  $p_{ij}$  from  $X$
- 5: Initialize  $Y \sim \mathcal{N}(0, 10^{-4}I)$
- 6:
- 7: **for** iteration  $t = 1$  to  $T$  **do**
- 8:     Compute all  $q_{ij}$  from  $Y$
- 9:     Compute gradients  $\frac{\partial C}{\partial Y}$
- 10:    Update with momentum:
- 11:      $Y^{(t)} = Y^{(t-1)} - \eta \frac{\partial C}{\partial Y} + \alpha \Delta Y^{(t-1)}$
- 12: **end for**

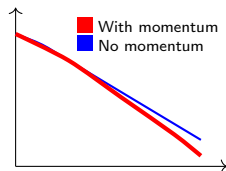
# Momentum: Faster Convergence

## Update rule:

$$Y^{(t)} = Y^{(t-1)} - \eta \nabla + \alpha \Delta Y^{(t-1)}$$

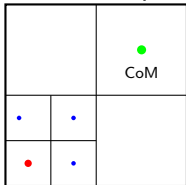
where:

- $\eta$ : learning rate
- $\alpha$ : momentum (0.5→0.8)
- $\Delta Y$ : previous update



# Barnes-Hut: From $O(n^2)$ to $O(n \log n)$

**Idea:** Group distant points

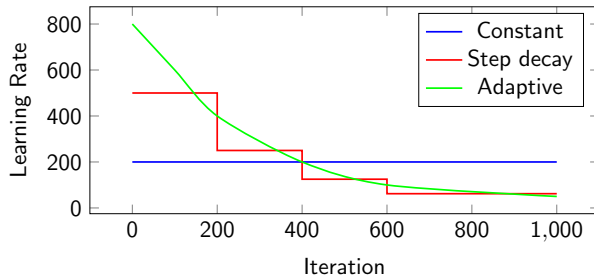


## Approximation:

- Build quadtree
- Compute centers of mass
- If cell far: treat as one point
- Threshold:  $\theta = 0.5$

**Speedup:**  $100\times$  for  $n = 10,000$

# Learning Rate Strategies



**Recommendation:** Start high (500-1000), decrease if needed

# Smart Initialization

Method	Pros	Cons
Random small	No bias	Slow start
PCA	Fast convergence	May bias
Previous run	Reproducible	Local minimum

**Best practice:**  $Y_i \sim \mathcal{N}(0, 10^{-4}I)$

Small variance prevents early numerical issues

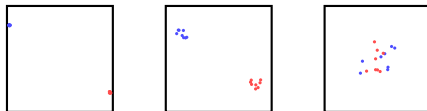
# Hyperparameter Impact

Parameter	Low	Default	High
Perplexity	5-15	30	50-100
Learning rate	10-100	200	500-1000
Iterations	250	1000	5000
Momentum	0.5	0.8	0.9
Early exag.	4	12	20

**Grid search often needed for optimal results**

# Perplexity: Detailed Effects

Perp = 5      Perp = 30      Perp = 100



- **Too low:** Breaks clusters into fragments
- **Too high:** Merges distinct clusters
- **Sweet spot:** Usually 5-50, dataset dependent

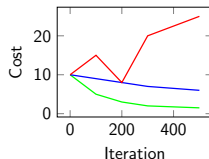
# Learning Rate: Finding Balance

## Too low ( $\eta < 10$ ):

- Stuck in bad minimum
- Slow convergence
- Poor separation

## Too high ( $\eta > 1000$ ):

- Points explode
- Oscillations
- Never converges





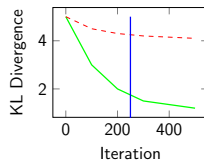
# Monitoring Convergence

## Watch for:

- KL divergence decrease
- Gradient norm  $\rightarrow 0$
- Stable embedding

## Warning signs:

- Increasing cost
- Points at infinity
- Oscillations



## Recommended workflow:

- 1 Start with defaults (perp=30, lr=200)
- 2 Run 5 times with different seeds
- 3 If inconsistent: adjust perplexity
- 4 If slow: increase learning rate
- 5 If unstable: decrease learning rate
- 6 Compare multiple perplexity values

Always run multiple times - t-SNE is stochastic!

# Reading t-SNE Correctly

Critical questions:

- What can we trust?
- What is meaningless?
- How to validate?

# What You Can Trust

## Can Trust:

- Local neighborhoods
- Cluster existence
- Within-cluster structure
- Relative densities (roughly)

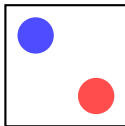
## Cannot Trust:

- Cluster sizes
- Between-cluster distances
- Global structure
- Absolute positions

t-SNE is for exploration, not measurement!

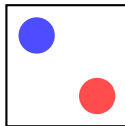
# Cluster Separation: Real or Artifact?

**Real**



Consistent

**Artifact**

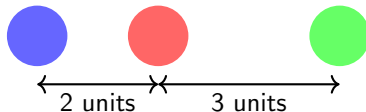


Random seed

**Validation:** Run multiple times, check if stable

# Distance Interpretation Warnings

**Between-cluster distances are meaningless!**



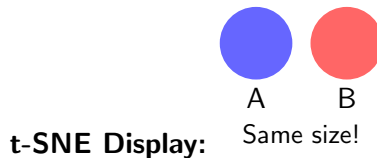
**These distances mean nothing!**

In high-D, all three might be equidistant

# Cluster Size Non-Preservation

## High-D Reality:

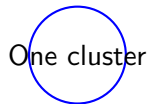
- Cluster A: 1000 points
- Cluster B: 100 points
- Ratio: 10:1



**Why:** Optimization doesn't preserve density

# Pitfall: Perplexity Too Low

**Truth**



**Perp = 2**



False structure!

**Solution:** Increase perplexity to 15+



# Pitfall: Perplexity Too High

**Truth**



Two clusters

**Perp = 200**



Lost structure!

**Solution:** Decrease perplexity to 30-50

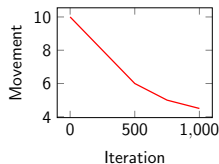
# Pitfall: Non-Convergence

## Symptoms:

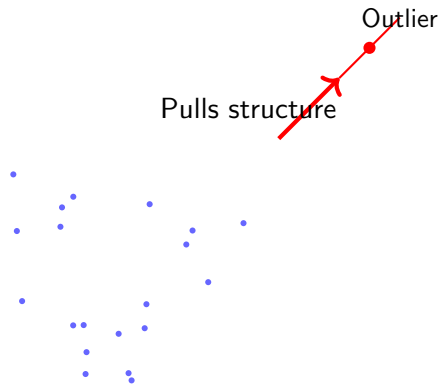
- Points still moving
- Cost oscillating
- Clusters not separated

## Solutions:

- More iterations (2000+)
- Adjust learning rate
- Check for outliers



# Outlier Effects



## Solutions:

- Remove extreme outliers first
- Use robust preprocessing
- Check with and without outliers

# Ensuring Reproducibility

## ① Set random seed

- For initialization
- For algorithm

## ② Document parameters

- Perplexity
- Learning rate
- Iterations

## ③ Save intermediate states

- Every 100 iterations
- For debugging

Always report: "t-SNE with perp=X, lr=Y, iter=Z"

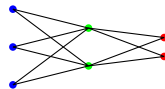
# Advanced: Parametric t-SNE

**Idea:** Learn a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^2$

**Advantages:**

- Can embed new points
- Inverse mapping possible
- Interpretable features

**Neural Network:**



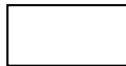
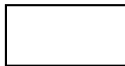
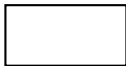
# Advanced: Multi-Scale t-SNE

**Use multiple perplexities simultaneously**

Perp=5

Perp=30

Perp=100



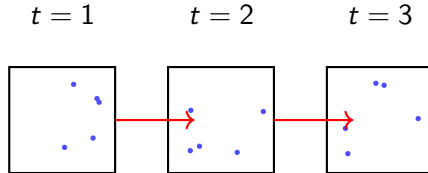
Combine

$$p_{ij} = \sum_k w_k \cdot p_{ij}^{(perp_k)}$$

Captures both local and global structure

# Advanced: Dynamic t-SNE

**For temporal data:** Preserve structure over time



Add temporal regularization:  $\lambda ||Y_t - Y_{t-1}||^2$

# t-SNE vs UMAP

Aspect	t-SNE	UMAP
Speed	$O(n \log n)$	$O(n^{1.14})$
Theory	Probability	Topology
Global structure	Poor	Better
Parameters	Perplexity	n_neighbors
Reproducibility	Random	More stable
Scalability	~50K points	Millions

## When to use each:

- t-SNE: Exploring clusters, publication figures
- UMAP: Large data, need global structure



- ① **Initialization:** PaCMAP, TriMAP
- ② **Speed:** FIt-SNE, openTSNE
- ③ **Theory:** Heavy-tailed embeddings
- ④ **Interpretability:** Attribution methods
- ⑤ **Uncertainty:** Probabilistic embeddings

**Active research area:** 100+ papers/year

# From Theory to Code

Ready to implement t-SNE!

# Python: scikit-learn

```
from sklearn.manifold import TSNE
import numpy as np

# Your data: n_samples × n_features
X = load_data()

# Configure t-SNE
tsne = TSNE(n_components=2,
            perplexity=30,
            learning_rate=200,
            n_iter=1000,
            random_state=42)

# Fit and transform
Y = tsne.fit_transform(X)
```

# R: Rtsne Package

```
library(Rtsne)

# Prepare data
X <- as.matrix(your_data)

# Run t-SNE
tsne_out <- Rtsne(X,
                  dims = 2,
                  perplexity = 30,
                  theta = 0.5,
                  max_iter = 1000)

# Extract embedding
Y <- tsne_out$Y
```

# Key Parameters Explained

Parameter	Meaning	Guidance
perplexity	Neighborhood size	5-50
learning_rate	Step size	10-1000
n_iter	Iterations	250-5000
theta	Barnes-Hut accuracy	0.5 default
metric	Distance function	euclidean
init	Initialization	pca or random

**Most important:** perplexity and learning\_rate

## ① Preprocessing:

- PCA to 50D first
- Normalize features
- Remove duplicates

## ② Computation:

- Use float32 not float64
- Enable multicore
- GPU versions available

## ③ Large datasets:

- Sample first, then embed
- Use UMAP for  $\geq 100K$  points
- Consider parametric t-SNE

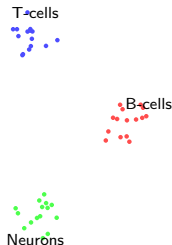
# Case Study: Single-Cell RNA-seq

## Dataset:

- 10,000 cells
- 20,000 genes
- Goal: Find cell types

## Pipeline:

- 1 Filter genes (variance)
- 2 Log transform
- 3 PCA to 50D
- 4 t-SNE with  $\text{perp}=30$



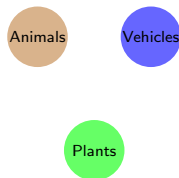
# Case Study: ImageNet Features

## Setup:

- CNN features (2048D)
- 50,000 images
- 1000 classes

## Results:

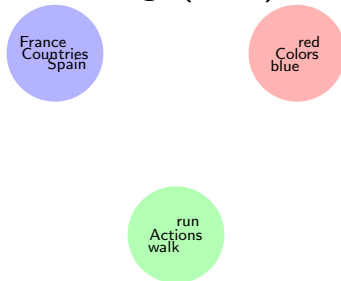
- Similar objects cluster
- Hierarchical structure
- Visual similarity preserved





# Case Study: Word Embeddings

Word2Vec embeddings (300D)  $\rightarrow$  t-SNE (2D)



Semantic relationships preserved!

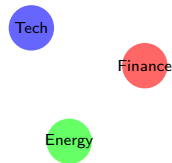
# Case Study: Financial Time Series

## Data:

- Stock returns
- 500 companies
- 252 trading days

## Preprocessing:

- Correlation matrix
- t-SNE embedding
- Color by sector



Sectors naturally separate!

## t-SNE Checklist:

- 1 Is your data high-dimensional? ( $d \geq 10$ )
- 2 Do you want to explore structure?
- 3 Is  $n < 50,000$ ?
- 4 Can you validate clusters independently?

If yes to all  $\rightarrow$  t-SNE is perfect!

## Remember:

- Try multiple perplexities
- Run multiple times
- Validate findings

# Key Takeaways

What you have mastered today

# When to Use t-SNE

## Use t-SNE:

- Exploring clusters
- Validating features
- Finding outliers
- Publication figures
- Quality over speed

## Dont use t-SNE:

- Measuring distances
- $> 100K$  points
- Real-time analysis
- Definitive proof
- Production systems

t-SNE is for exploration and insight

# The Interpretation Checklist

Before publishing t-SNE results:

- ❶ ☐ Tried perplexity: 5, 15, 30, 50
- ❷ ☐ Ran 5+ random initializations
- ❸ ☐ Checked convergence (1000+ iterations)
- ❹ ☐ Validated clusters independently
- ❺ ☐ Stated all parameters clearly
- ❻ ☐ Acknowledged limitations
- ❼ ☐ Compared with PCA/other methods

Never interpret distances or sizes!

# Resources for Deeper Study

## Essential Papers:

- Original: van der Maaten & Hinton (2008)
- Barnes-Hut: van der Maaten (2014)
- Theory: Linderman & Steinerberger (2017)

## Software:

- Python: scikit-learn, openTSNE
- R: Rtsne, tsne
- Fast: FIt-SNE, RAPIDS

## Tutorials:

- Distill.pub interactive guide
- Google embedding projector

## Your Questions?

Theory?  
Implementation?  
Your data?

*Thank you for your attention!*  
*Now lets explore your data with t-SNE!*