



Fall 2018:
Introduction to
Data Science

GIRI NARASIMHAN, SCIS, FIU

High- Dimensional Space

Points in High-Dimensional Space

APPLICATIONS

- ▶ Information Processing
- ▶ Search
- ▶ Data Mining
- ▶ Machine Learning (ML)

Points in d-Dimensional Space

- ▶ Assume that d is large
- ▶ What is the volume of a unit ball in d -space?
- ▶ Are the points well spread out? Where do most of the points lie?
- ▶ Assume that we generate n points at random in d -dimensional ball of radius 1
- ▶ What can we say about:
 - ▣ Distance between any pair of points
 - ▣ Angle between any two vectors from origin to that point

Points in d-Dimensional Space

- ▶ Let x and y be points in d -space with coordinates from unit-variance Gaussians.
- ▶ How far is x from origin
 - ▣ Approx distance squared = d
- ▶ No prob mass close to O , although prob density has max at O
- ▶ Unit ball has zero volume; integral of prob density over unit ball = 0
- ▶ $|x - y|^2 = 2d$
 - ▣ Thus vectors x and y are approximately orthogonal
- ▶ If x is the North Pole, then most of the points lie near the equator

Properties of d-ball and Gaussians

- ▶ Area
- ▶ Volume
- ▶ Denominator > Numerator
- ▶ Most of volume is near equator
- ▶ Thus any two vectors are nearly orthogonal
- ▶ d-dimensional Gaussian
 - ▣ Most points are in annulus

$$A(d) = \frac{2\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2})} \quad \text{and} \quad V(d) = \frac{2}{d} \frac{\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2})}.$$

With probability $1 - O(1/n)$

1. $|\mathbf{x}_i| \geq 1 - \frac{2 \ln n}{d}$ for all i , and
2. $|\mathbf{x}_i \cdot \mathbf{x}_j| \leq \frac{\sqrt{6 \ln n}}{\sqrt{d-1}}$ for all $i \neq j$.

$$\sqrt{d} - \beta \leq |\mathbf{x}| \leq \sqrt{d} + \beta$$

Nearest Neighbor Search

- ▶ Assume you have database of n entries in d -space
- ▶ Answer queries of the form
 - ▣ Given x , find the nearest neighbor in the database
- ▶ Goal: preprocess database so that queries are answered quickly
- ▶ Database does not change much, but large number of queries
- ▶ Perform expensive preprocessing, but speed up queries
- ▶ Time complexity depends on n and d
- ▶ Hence the need for **dimensionality reduction**

Projections: \mathbb{R}^d to \mathbb{R}^k

- ▶ Projection of a set of points/vectors along a new vector v is given by:

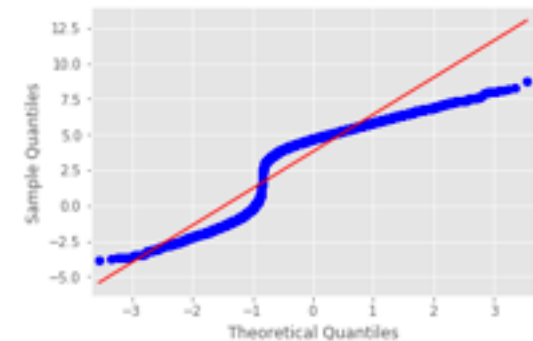
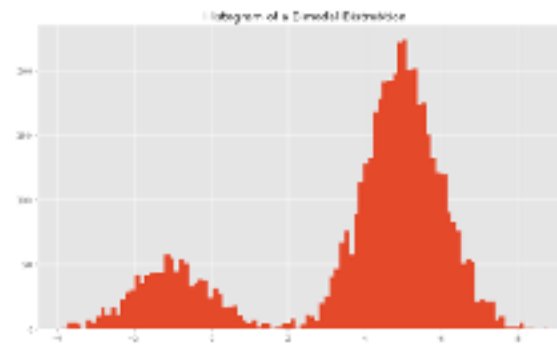
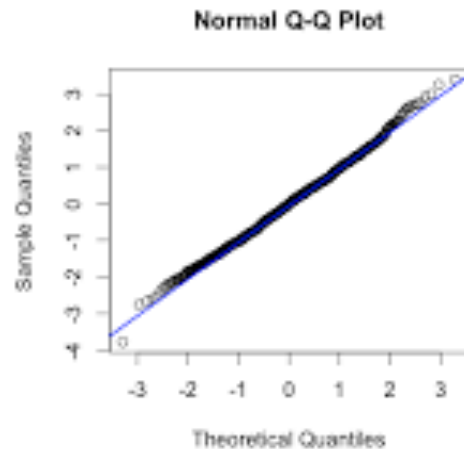
$$f(\mathbf{v}) = (\mathbf{u}_1 \cdot \mathbf{v}, \mathbf{u}_2 \cdot \mathbf{v}, \dots, \mathbf{u}_k \cdot \mathbf{v}).$$

- ▶ What happens to distances after projections?
 - ▣ Johnson-Lindenstrauss Theorem applies to all pairs
 - ▣ $0 < \epsilon < 1$, any n , $k \geq (3 \ln n) / (\epsilon^2)$, any v_i, v_j , with pr $1 - 1.5/n$:

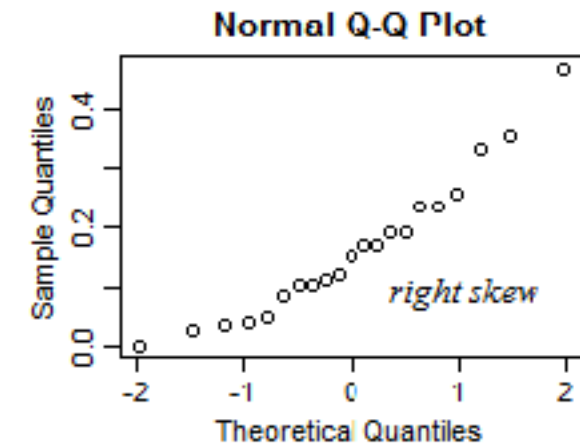
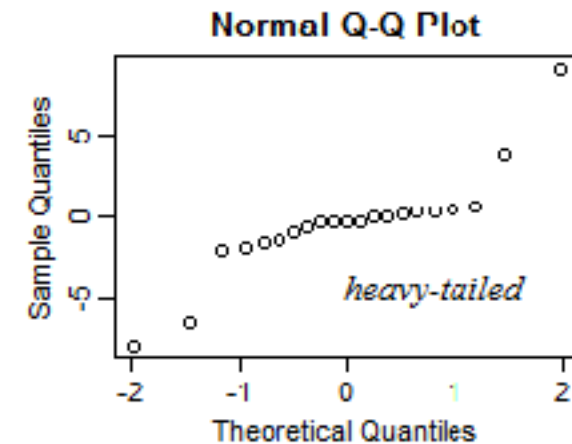
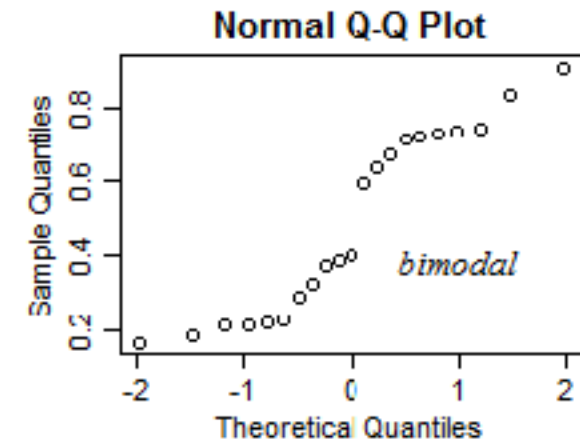
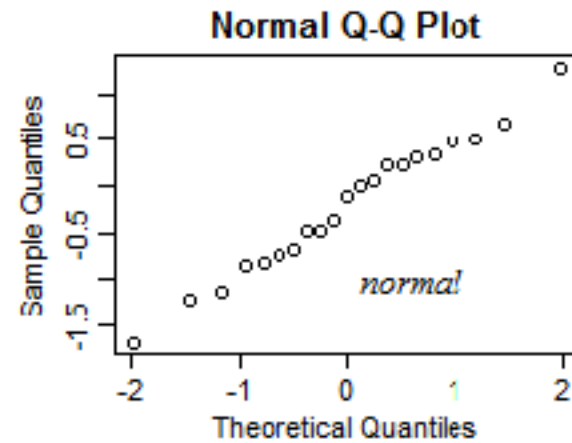
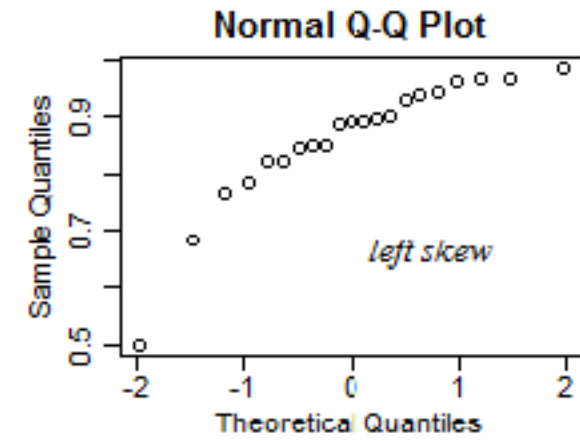
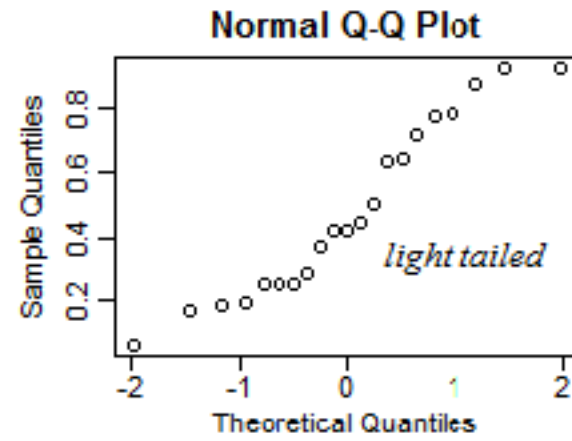
$$(1 - \epsilon)\sqrt{k} |\mathbf{v}_i - \mathbf{v}_j| \leq |f(\mathbf{v}_i) - f(\mathbf{v}_j)| \leq (1 + \epsilon)\sqrt{k} |\mathbf{v}_i - \mathbf{v}_j|.$$

Test for Normality

- ▶ Informal: plot a histogram and see if it is bell-shaped
- ▶ Graphical approach: Do a quantile-quantile (QQ) plot



QQ Plots & Interpretations



Tests for Normality ... 2

- ▶ D'Agostino's K-squared test,
 - ▶ Jarque-Bera test,
 - ▶ Anderson-Darling test,
 - ▶ Cramér-von Mises criterion,
 - ▶ Lilliefors test,
 - ▶ Kolmogorov-Smirnov test,
 - ▶ Shapiro-Wilk test, and
 - ▶ Pearson's chi-squared test.
- ▶ Some Bayesian approaches exist as well