

summary

Last time

$$x \sim p(x)$$

x : RANDOM VARIABLE in \mathbb{R}^d

- GOAL: given $\{x_i\}_{i=1}^N$, reconstruct $p_n(x) \approx p(x)$

- KERNEL DENSITY ESTIMATION (KDE)

$$P_N(z) = \frac{1}{N h^d} \sum_{i=1}^N K\left(\frac{z - x_i}{h}\right)$$

h : BANDWIDTH $\left\{ \begin{array}{l} \text{THIS ASSUMES CONSTANT} \\ \text{BANDWIDTH} \end{array} \right\}$

EXAMPLES

$$\circ K(r) = \begin{cases} \frac{1}{2}, & \|r\| < 1 \\ 0, & \text{otw} \end{cases} \quad r \in \mathbb{R}^d$$

NAIVE

$$\circ h: h(z) = d_k(z) : \text{dist}$$

NEAREST NEIGHBOR

$$\circ K(r) = (2\pi)^{-d/2} \exp(-r^T r)$$

GAUSSIAN.

HYPERPARAMETER(S) : h, k

- choose BY
 - REGULARIZED LOG LIKELIHOOD
 - EXPECTATION OF RMS (MSE)

SMALL $h, k \rightarrow$ overfitting → DOES NOT GENERALIZE WELL
→ SMALL BIAS, LARGE VARIANCE

CAN ANALYTICALLY COMPUTE h_{opt}

ISSUES WITH DIMENSIONALITY

- PROBABILITY OF FINDING A POINT IN A VOLUME OF A CUBE OF LENGTH s .

$$\frac{\|p_N(\mathbf{c}) - p(\mathbf{c})\|}{\|p(\mathbf{c})\|} \text{ FOR NORMAL GAUSSIAN}$$

WITH $h_{opt} = \left(\frac{4}{d+2}\right)^{1/(d+4)} N^{-\frac{1}{d+4}}$

$$\text{FIND } N \text{ so that error} < 10\%$$

$d=1$
 $d=5$
 $d=10$

IN PRACTICE : How to choose h ?

SPLIT TRAINING SET 80% 20%
TRNG TEST
 x_N x_T

$$K_{ij} = K(x_{T,i}, x_{N,j})$$

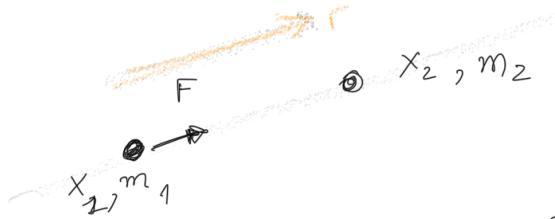
FIND h (or k) such that

log likelihood of $p_N(x_T | x_N, h)$ is max

$$h = \arg \max_h \sum \left(\log \left(K \left[\frac{i}{Nh} \right] \right) \right)$$

PHYSICS

GRAVITATIONAL INTERACTIONS



$$\begin{aligned}\Gamma &= x_2 - x_1 \\ d &= \frac{\Gamma}{\|\Gamma\|_2}\end{aligned}$$

$$F = G \frac{m_1 m_2}{r^2} \quad d = , \quad F(x_1) = \frac{m_1 m_2}{\|\Gamma\|_2^2} \cdot \frac{\Gamma}{\|\Gamma\|_2}$$

IF WE HAVE N OBJECTS

$$F_1 = m_1 \sum_{j=2}^N \frac{x_i - x_j}{\|x_i - x_j\|^3} m_j = m_1 \sum_{j=2}^N k(x_i, x_j) m_j$$

Typically we have \$N\$ particles

$$\left\{ \begin{array}{l} m_i \frac{d\mathbf{v}_i}{dt} = \mathbf{F}_i = m_i \sum_{\substack{j=1 \\ j \neq i}}^N k(x_i, x_j) m_j \\ \frac{dx_i}{dt} = \mathbf{v}_i \end{array} \right.$$

Kernel matrix : $K_{ij} = k(x_i, x_j)$

THREE PROBLEMS

~ INNER PRODUCT SPACES & THE KERNEL TRICK

$$K\mathbf{a} = \mathbf{y} \Rightarrow \mathbf{a} = K^{-1}\mathbf{y}; \rightarrow \text{PREDICT}$$

~ PROBABILITY : $K \mathbb{1}$

~ PHYSICS : $K m$

IN GENERAL DURING TRAINING :

$$K \in \mathbb{R}^{N \times N}$$

N: TRAINING SET

typically K is Sym. PD.

DURING TESTING

$$K \in \mathbb{R}^{M \times N}$$

M: TESTING SET

K can be SPARSE ($\mathcal{O}(N)$ non zeros)

or most often (in all these applications)

DENSE $\frac{|K_{ij}|}{\max|K_{ij}|}$ significant.

• COMPLEXITY FOR DENSE MATRICES

• MATVEC: $\mathcal{O}(N^2)$

• LINEAR SOLVE: $\mathcal{O}(N^2 \sqrt{\text{cond}(K)})$
 $\mathcal{O}(N^3) \rightarrow \text{CHOLESKY}$

Prohibitively expensive

$N = 10^3$ REQUIRES 10^5 GPUs RUNNING AT TFLOPS.

[~~see's~~]

SCALABLE ALGORITHMS

• LOW-RANK APPROXIMATION

• BLOCK - LOW-RANK APPROXIMATION

• N-BODY ALGORITHMS

• PARALLELIZATION

• NEAREST NEIGHBORS.

Low-Rank Matrices

SUPPOSE THAT $K \in \mathbb{R}^{M \times N}$ IS LOW-RANK
(EXACT OR NUMERICALLY)

LET $U\Sigma V^T$ BE ITS REDUCED SVD

- Every matrix has one
- Check properties of SVD

THEN we say K has rank r if

$$K = \begin{matrix} M \\ \downarrow \\ \begin{array}{|c|c|c|} \hline & r & \\ \hline U & \left[\begin{array}{ccc} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{array} \right] & V^T \\ \hline & r & \\ \hline & \vdots & \\ \hline & u_i & \\ \hline \end{array} \end{matrix} = \sum_{i=1}^r \sigma_i u_i v_i^T$$

We say that K is APPROXIMATELY

RANK r if

$$\min_{\tilde{K}} \|\tilde{K} - K\|_2 \ll \sigma_r \quad \text{And } \tilde{K} \text{ has rank } r$$

TURNS OUT r -truncation of SVD RESULTS

IN OPTIMAL APPROXIMATION,
SPECTRAL GMP. $\|K - K_r\|_2 = \sigma_{r+1}$

\rightarrow COMPLEXITY OF MATVEC
 \rightarrow COST OF COMPUTING SVD?

THE OTHER REGIME:

Truncated kernels, like the nearest neighbor:

These are sparse kernels.

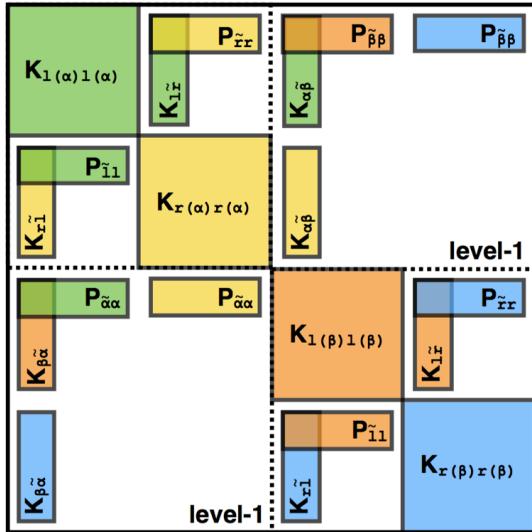
Once we form the matrix, the matvec is $O(Nk)$.

But then how do we discover the sparsity?

THE REALITY; KERNEL MATRICES ARE NEITHER LOW-RANK. NEITHER SPARSE.
The best parameters typically result in matrices that their numerical rank is huge (so formally not low-rank).

HIERARCHICAL APPROXIMATIONS.

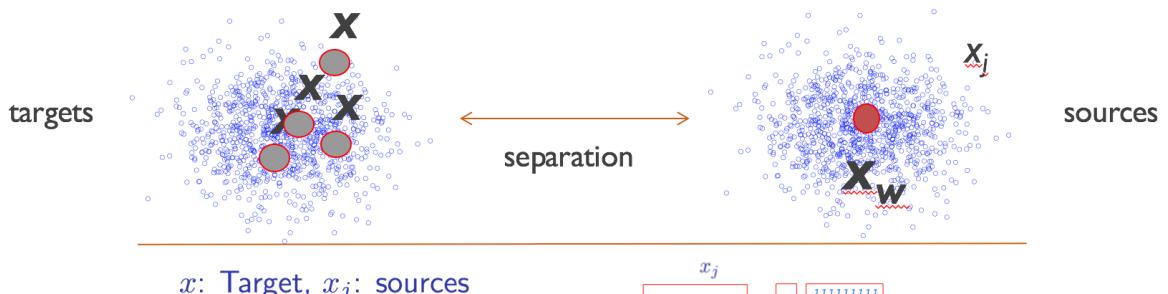
Hierarchical matrices



$$\begin{aligned}
 & \begin{bmatrix} G_{11} & G_{12} \\ G_{21} & G_{22} \end{bmatrix} \\
 = & \begin{bmatrix} G_{11} & 0 \\ 0 & G_{22} \end{bmatrix} + \begin{bmatrix} 0 & G_{12} \\ G_{21} & 0 \end{bmatrix} \\
 & \frac{D + UV}{D + UV} \quad \text{rank "s"} \\
 & \mathcal{O}(N^2) \rightarrow \mathcal{O}(N \log N)
 \end{aligned}$$

NBODY 101

Idea I: Far-field interactions



$$u(x) = \sum_j G(x, x_j) w_j$$

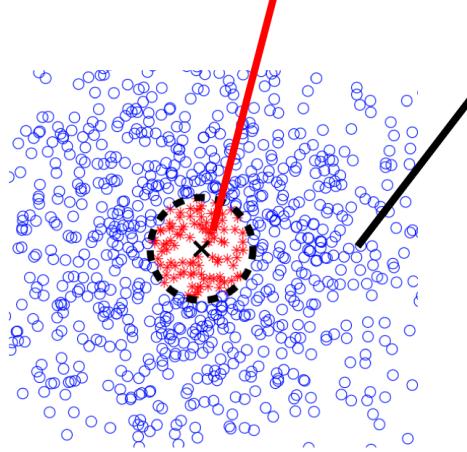
1. compute $W = \sum_i w_i$
2. choose x_W
3. $u(x) \approx G(x, x_W)W$

$$\begin{array}{c|c}
 x & x_j \\
 \hline
 G(x, x_j) & \begin{matrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{matrix}
 \end{array}
 = \begin{array}{c|c}
 & W \\
 \hline
 & \begin{matrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{matrix}
 \end{array}$$

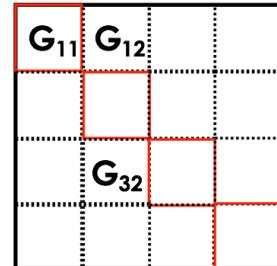
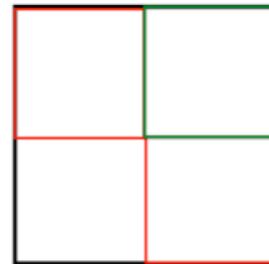
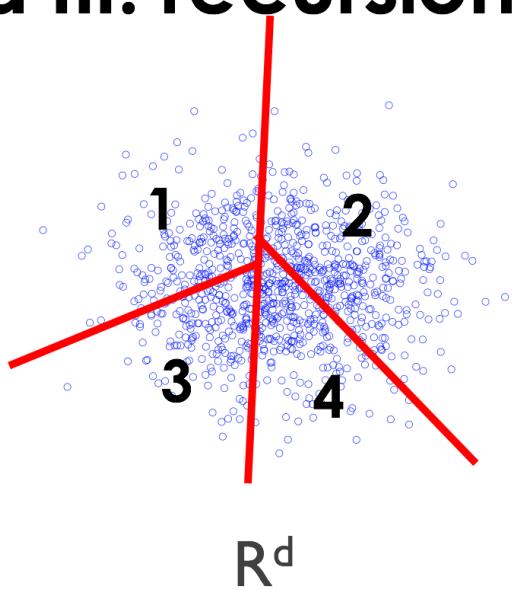
LOW RANK APPROXIMATION $\mathcal{O}(N^3) \rightarrow \mathcal{O}(N)$

Idea II: Near-/Far-field split

$$u_i = \sum_{\substack{j=1 \\ j \neq i}}^N G(x_i, x_j) w_j = \sum_{j \in \text{near}(i)} G_{ij} w_j + \sum_{j \in \text{far}(i)} G_{ij} w_j$$



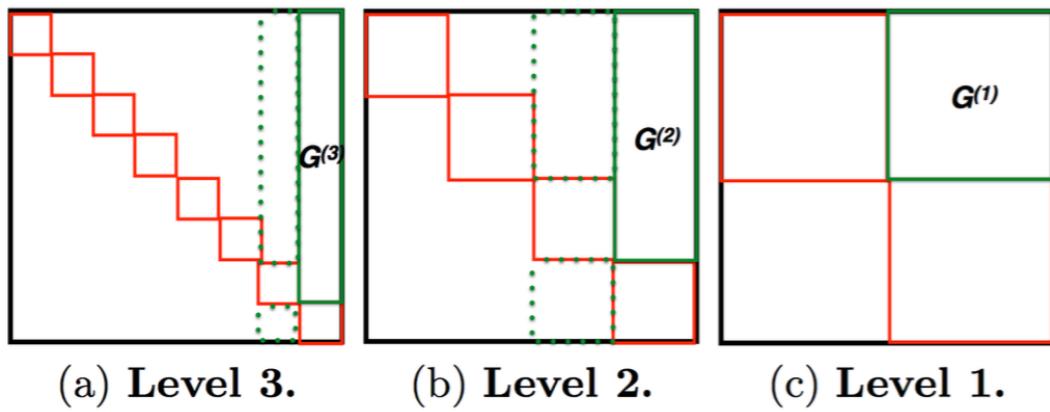
Idea III: recursion



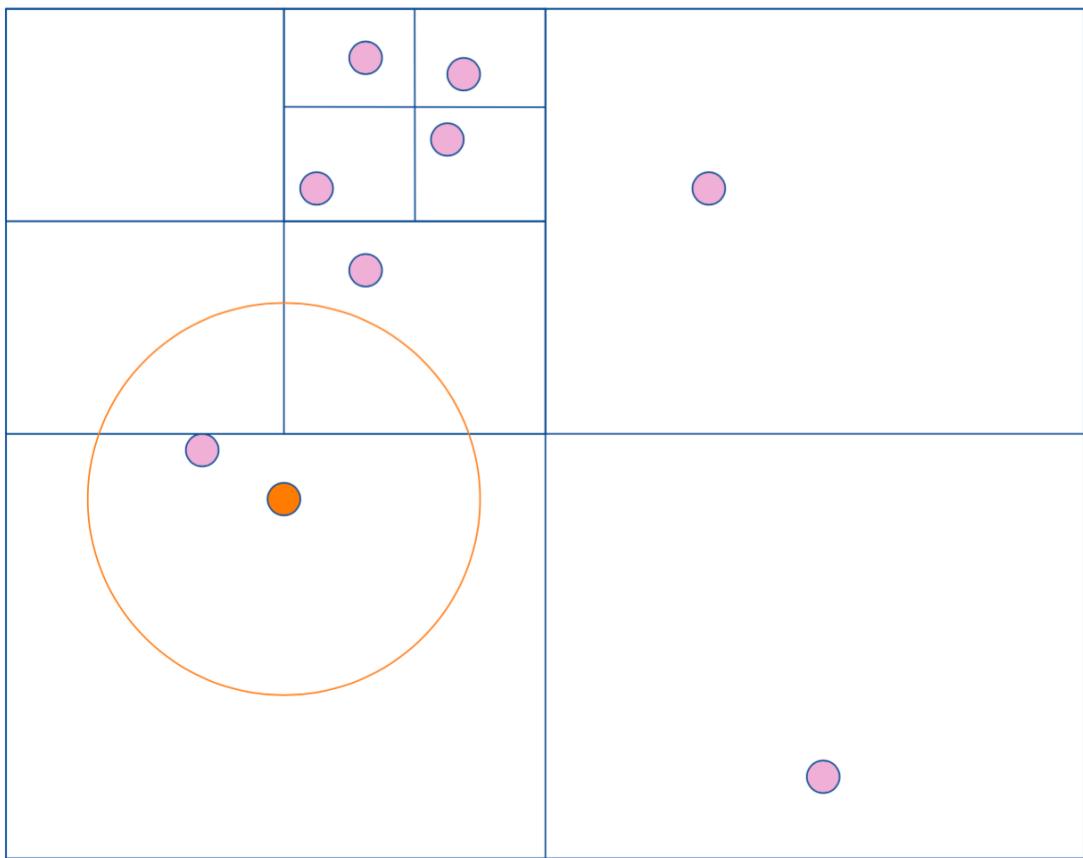
Matrix partitioning

Algebraic view

$$u_i = \sum_{j \in \text{diagonal}(i)} G_{ij} w_j + \sum_{j \in \text{off-diagonal}(i)} G_{ij} w_j$$



The NEAREST NEIGHBORS PROBLEM:
DIRECT EVALUTION of distances O(N) complexity



Exact search $O(n)$ per query point

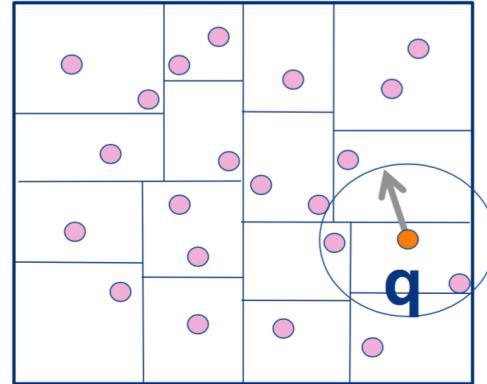
In low dimensions can be $O(1)$ per query point

- Preprocess points to spatial data structure

- pruning => $O(1)$ cost

PRUNING

1. Construct tree (group points)
2. Insert q into tree
3. Find nearest neighbor in the leaf
4. Range search and pruning



Issues with regular decompositions

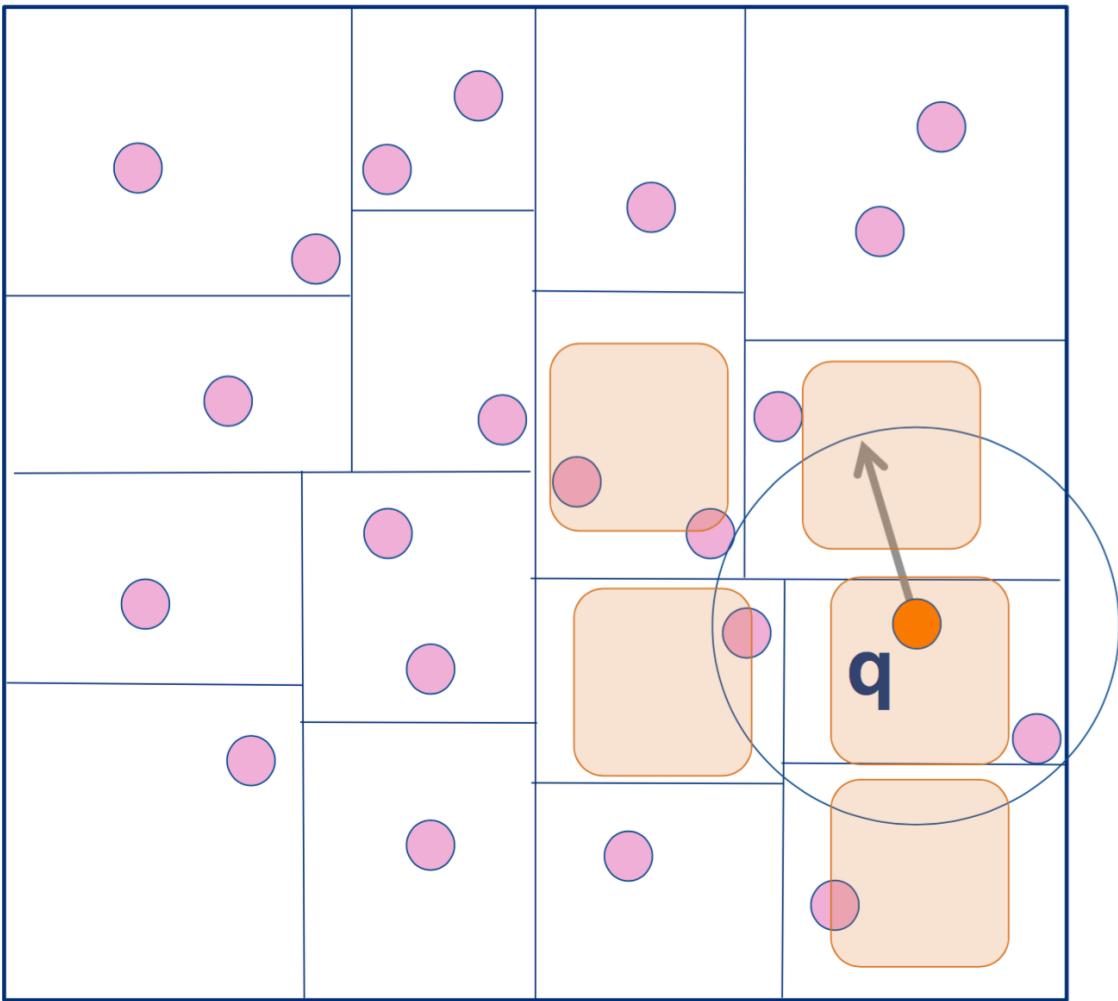
nearest-neighbor problems

high-dimensional problems

- 2^d per division

alternative data structures

- popular: KD-tree



Range query (q, R)

Traverse tree bottom-up

- if node is a leaf
 - do direct range query on the points
- visit node
 - If $\text{BoundingBox}(\text{node}) \cap \text{Ball}(q, R) = \text{Empty}$
 - TERMINATE RECURSION
 - else visit children

Nearest neighbor query(q)

`leaf_q = insert_point(q, root)`

- like quadtrees
- traverse tree top down and insert point
- p=find nearest neighbor in “leaf_q” using direct search
- $R_{initial} = \text{distance}(leaf_q, p)$
- traverse tree top down using range search with $R_{initial}$ and updating $R_{initial}$ if a point p' closer to q is found

Back to NN-search: some negative results, $d>>3$

Main problem: r-search or k-search require ball around query point to use for pruning.

In high dimensions (say $d>100$) pruning is negligible and the complexity is $O(n)$ per query point

- For uniform or normal distributions, exact searches, and any tree that uses convex grouping (Weber, Shek, Blott '98)

$O(n^2)$ for all nearest neighbor problems

Consider all nearest neighbors

Iterate

- choose rotation
- rotate coordinates
- build KDT tree
- find NNs of each point in its leaf
- merge with existing NNs set of each point

- 1: $K = \emptyset, \mathcal{T} = \text{root}$
- 2: $Q = \text{RANDOMROTATION}$
- 3: $p = Qr$
- 4: $\text{NODESPLIT}(\mathcal{T}, p)$
- 5: $F = \text{LEAFLOCALSEARCH}(\mathcal{T})$
- 6: $K = \text{UPDATENEIGHBORS}(K, F)$
- 7: if $\text{CONVERGED}(K, F)$ return
- 8: GOTO 2

