

L19: Noise in Data

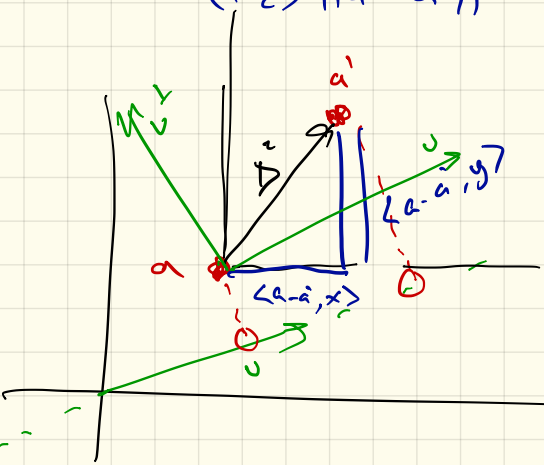
- Random Projections
- Outliers
- Matrix Completion

Random Projections

Input: Data $A \in \mathbb{R}^{n \times d} = \{a_1, a_2 \dots a_n\} \subset \mathbb{R}^d$

Goal: map $u: \mathbb{R}^d \rightarrow \mathbb{R}^k$ ($k \ll d$)
 preserve all distances say $a, a' \in A$

$$(1-\epsilon) \|a - a'\| \leq \|u(a) - u(a')\| \leq (1+\epsilon) \|a - a'\|$$



$$\Delta^2 = \langle a - a', x \rangle^2 + \langle a - a', y \rangle^2 \quad \text{w.p. } 1 - \delta \quad \leftarrow \text{prob. failure}$$

$$\Delta^2 = \langle a - a', v \rangle^2 + \langle a - a', v^\perp \rangle^2$$

$$\hat{u}(a) = \langle a, v \rangle$$

$$\hat{u}(a') = \langle a', v \rangle$$

$$E[\langle a - a', v \rangle^2] = \frac{\Delta^2}{d}$$

$$E[\langle a, v \rangle \cdot \langle a', v \rangle] = \Delta$$

for $i = 1$ to k

$$\tilde{v}_i \sim G_d(0, 1)$$

$$v_i = \frac{\tilde{v}_i}{\|\tilde{v}_i\|} \cdot \frac{\sqrt{d}}{\sqrt{k}}$$

for points $a_j \in A$

for $i \in [k]$

$$\mu(a_j)_i = \langle a_j, v_i \rangle$$

$$a_j \rightarrow \mu(a_j) = (\mu(a_j)_1, \mu(a_j)_2, \dots, \mu(a_j)_k)$$

$(1 \pm \epsilon)$ error w.p. $1 - \delta$

$$\Leftrightarrow k = O\left(\frac{1}{\epsilon^2} \log \frac{n}{\delta}\right)$$

1% $\Rightarrow 10,000$

works for
about
 $k = 500 - 10,000$

Not Data
Adaptive!

Noise in Data

- Sporrious Readings :

outliers

↳ adversarial examples, data poisoning

- Measurement Error: Gaussian Noise

benign, SSE formulations, regularization

- Background Data:

mixed in, or

missing

Matrix completion

Outliers

Input $X \in \mathbb{R}^{n \times d}$

• Build Model, Remove Outliers

• Density-based Approaches

1. Build model M on X

2. For each $x \in X$, find residual

$$r_x = d(M(x), x)$$

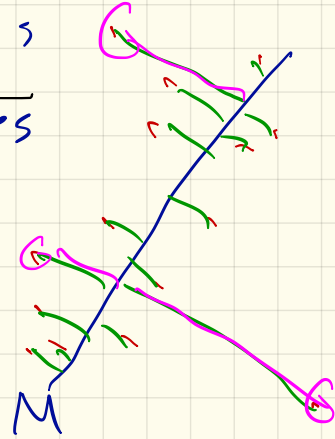
3. For $x \in X$, with r_x too large \rightarrow remove from X .
 $\hookrightarrow \tilde{X}$

4. Go back to Step 1. w/ \tilde{X}

1. threshold

2. percentage (remove 10%)

percentage of X
not \tilde{X}

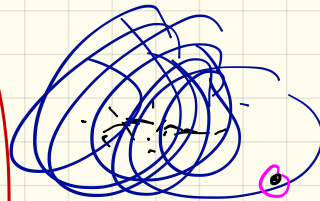
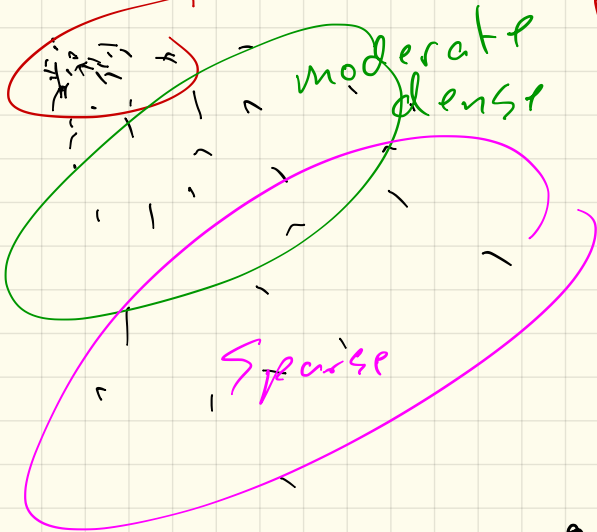


Density-based Approach

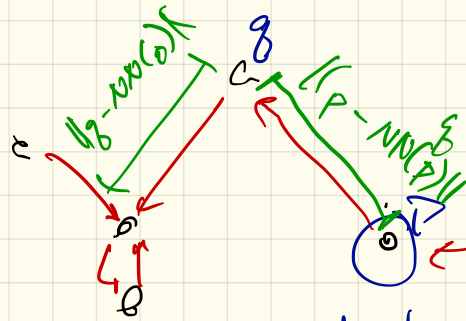
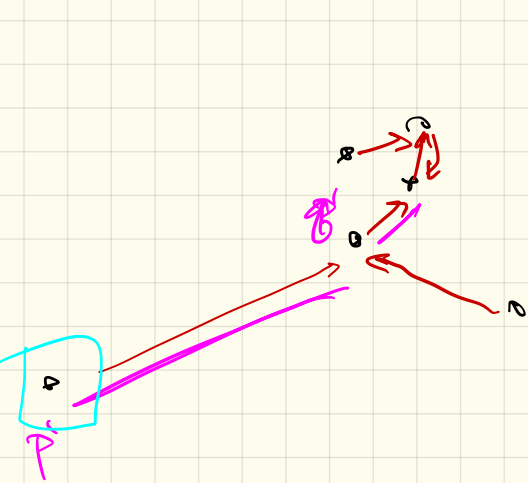
• DB-Scan Clustering.

two param
Very dens

r , T
radius, threshold



Reverse-Nearest Neighbor



distance to $RNN(p)$

$\|g - p\|$
compare to

$\|g - RNN(g)\|$

If comparable
 p not outlier

Matrix Completion

Input $A \in \mathbb{R}^{n \times d}$ but w/ some ? x

$$A = \begin{bmatrix} 3 & 4 & x & 8 \\ 1 & x & 5 & x \\ x & 4 & x & 9 \\ 9 & 7 & 1 & x \\ 2 & 2 & x & x \end{bmatrix}$$

S.175

$$\text{mask } \Omega = \begin{bmatrix} 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 \end{bmatrix}$$

replace each x by average of row and column.

assume Δ

low-rank (w/ noise)

$$\frac{3+1+9+2}{4} = 3.75$$

$$\frac{4+9}{2} = 6.5$$

$$\Rightarrow \frac{16.25}{2} = 8.125$$

$$\phi_\lambda(s) = \text{diag}((s_{11} - \lambda)_+, (s_{22} - \lambda)_+, \dots, (s_{dd} - \lambda)_+)$$

$$(x)_+ = \max(x, 0)$$

$$A^* = \underset{X \in \mathbb{R}^{n \times d}}{\text{argmin}} \frac{1}{2} \|\Pi_\Omega(A - X)\|_F^2 + \lambda \|X\|_*$$

$\| \cdot \|_*$ nuclear norm
 $\| \cdot \|_F$ Frobenius norm
 \sum sum singular values

Matrix Completion

Initialize $X_{i,j} \leftarrow \begin{cases} \Pi_\Omega(A_{i,j}) & \text{if } (i,j) \in \Omega \\ 0 & \text{or.} \end{cases}$

repeat

$$USV^T \leftarrow \text{svd}(X)$$

$$\tilde{X} \leftarrow U \phi_\lambda(s) V^T$$

$$X \leftarrow \Pi_\Omega(A) + \Pi_\Omega^\perp(\tilde{X})$$

← replace known elements.

return \tilde{X}