Low-Dimensional Representation of Symbolic Sequences

Manuel E. Lladser Applied Mathematics (Associate) Computer Science (Affiliate) University of Colorado - Boulder

AofA 2019 - CIRM Luminy, France

- Transmitted by mosquitos
- Kills in average 25K human worldwide per year
- Targets octapeptides (octamers)
- Target characterization has been elusive, making vaccine development challenging



Figure: Some known targets

- Transmitted by mosquitos
- Kills in average 25K human worldwide per year
- Targets octapeptides (octamers)
- Target characterization has been elusive, making vaccine development challenging



Figure: Some known targets

Question.

Are there sensible low-dimensional representations of octamers?

- One-hot encodings
- k-mer count vectors
- Fisher kernels ¹
- Multidimensional Scaling (PCoA)²
- Word2Vec (BioVec ³)
- DeepWalk ⁴
- Node2Vec ⁵

¹Jaakkola, Diekhans, and Haussler - 1999
 ²Krzanowski - 2000
 ³Asgari and Mofrad - 2015
 ⁴Perozzi, Al-Rfou, and Skiena - 2014
 ⁵Grover and Leskovec - 2016

- One-hot encodings (dimension: ak + sparse)
- k-mer count vectors (non-local + dimension: a^k)
- Fisher kernels ¹ (large dataset)
- Multidimensional Scaling (PCoA)² (*n* points map to dimension: n-1)
- Word2Vec (BioVec ³) (large dataset)
- DeepWalk ⁴ (new sample rerun)
- Node2Vec ⁵ (new sample rerun)

¹Jaakkola, Diekhans, and Haussler - 1999
 ²Krzanowski - 2000
 ³Asgari and Mofrad - 2015
 ⁴Perozzi, Al-Rfou, and Skiena - 2014
 ⁵Grover and Leskovec - 2016

Develop a sensible low-dimensional representation of points in a large metric space that addresses the shortcomings of the existing methods

Develop a sensible low-dimensional representation of points in a large metric space that addresses the shortcomings of the existing methods



Figure: Three non-colinear points trilaterate the plane

Develop a sensible low-dimensional representation of points in a large metric space that addresses the shortcomings of the existing methods



Figure: Three non-colinear points trilaterate the plane

Develop a sensible low-dimensional representation of points in a large metric space that addresses the shortcomings of the existing methods



Figure: Three non-colinear points trilaterate the plane

A Multilateration Approach



Question.

Is it possible to multilaterate any finite metric space?

A Multilateration Approach



Question.

Is it possible to multilaterate any finite metric space?

Definitions.

- Call $R \subset V$ resolving when $\forall x \neq y \in V \exists r \in R \text{ s.t. } d(x,r) \neq d(y,r)$
- $\Phi(x) := (d(x, r))_{r \in R}$ from V to $\mathbb{R}^{|R|}$ is one-to-one and "continuous"
- (Metric Dimension.) $\beta(V, d) := \min_{R \subset V \text{ resolving}} |R|$

(Terminology borrowed from graph theory!)

Observation.

The metric dimension of a metric space (V, d) is the smallest number of columns in its distance matrix D needed to differentiate all the rows

$$D = \begin{bmatrix} 0 & 1 & 2 & 1 & 2 & 3 \\ 1 & 0 & 1 & 1 & 1 & 2 \\ 2 & 1 & 0 & 2 & 2 & 1 \\ 1 & 1 & 2 & 0 & 2 & 3 \\ 2 & 1 & 2 & 2 & 0 & 3 \\ 3 & 2 & 1 & 3 & 3 & 0 \end{bmatrix}$$

Observation.

The metric dimension of a metric space (V, d) is the smallest number of columns in its distance matrix D needed to differentiate all the rows

$$D = \begin{bmatrix} 0 & 1 & 2 & 1 & 2 & 3 \\ 1 & 0 & 1 & 1 & 1 & 2 \\ 2 & 1 & 0 & 2 & 2 & 1 \\ 1 & 1 & 2 & 0 & 2 & 3 \\ 2 & 1 & 2 & 2 & 0 & 3 \\ 3 & 2 & 1 & 3 & 3 & 0 \end{bmatrix}$$

Theorem: (Tillquist-Ll'19)

The general multilateration problem is NP-complete.

Why? Reduction to the *Set Cover Problem*, which is a known ¹ NP-complete problem \Box

- Our terminology is borrowed from a graph theory ^{1,2}
- A set *R* of nodes in a graph *G* = (*V*, *E*) is called **resolving** when any vertex in the graph is uniquely determined by its vector of distances to those nodes
- The metric dimension of a graph, β(G), is the size of a smallest resolving set
- If *d* is the **geodesic distance** between pairs of nodes in *G* then

$$\beta(G) = \beta(V, d)$$

¹Slater - 1975 ²Harary and Melter - 1976

Examples

• $\beta(G) = 1$ if and only if G is a path



Examples

• $\beta(G) = 1$ if and only if G is a path



• $\beta(G) = (n-1)$ if and only if $G = K_n$



- $\beta(G)$ is known for specific graph families such as:
 - trees ^{1,2}
 - regular bipartite graphs ³
 - complete *n*-partite graphs ⁴
- Finding $\beta(G)$ is **NP-complete** for a general graph $G = (V, E)^{5}$
- The Information Content Heuristic (ICH) finds a resolving set greedily ⁶
 - O(|V|³) time-complexity
 - $O(|V|^2)$ memory-complexity
 - Approximation ratio of $1 + (1 + o(1)) \cdot \ln |V|$

¹Slater - 1975
²Harary and Melter - 1976
³Bača, Baskoro, Salman, et. al. - 2011
⁴Saputro, Baskoro, Salman, et. al. - 2009
⁵Hauptmann, Schmied, and Viehmann - 2012; Gary and Johnson - 1979
⁶Hauptmann, Schmied, and Viehmann - 2012

Definition.

 $H_{k,a}$ denotes the graph with vertex set $\{0, \ldots, a-1\}^k$, i.e. **k-mers** over an alphabet of size *a*. Two *k*-mers are connected by an edge if and only if they are identical except for one character at the same position



Figure: Visuals of $H_{1,3}$, $H_{2,3}$, and $H_{3,3}$ (resolving sets in blue)

Definition.

 $H_{k,a}$ denotes the graph with vertex set $\{0, \ldots, a-1\}^k$, i.e. **k-mers** over an alphabet of size *a*. Two *k*-mers are connected by an edge if and only if they are identical except for one character at the same position



Figure: Visuals of $H_{1,3}$, $H_{2,3}$, and $H_{3,3}$ (resolving sets in blue)

Question.

Can we find small resolving sets in $H_{k,a}$ efficiently?

Theorem: (Tillquist-Ll'19)

$$eta(H_{k-1,a}) \leq eta(H_{k,a}) \leq eta(H_{k-1,a}) + \lfloor a/2
floor.$$

* Case with a = 2 due to Chartrand et al. - 2000

Why? The proof is constructive!

- d:= distance matrix of $H_{k-1,a}$
- D:= distance matrix of $H_{k,a}$

$$D = \begin{bmatrix} 0 \dots & 1 \dots & \cdots & (a-1)\dots \\ 1 \dots & \\ \vdots & \\ (a-1)\dots & \\ \end{bmatrix} \begin{bmatrix} 0 \dots & 1 \dots & \cdots & d+1 \\ d+1 & d & \cdots & d+1 \\ \vdots & \vdots & \ddots & \vdots \\ d+1 & d+1 & \cdots & d \end{bmatrix}$$

Any resolving set of H_{k-1,a} distinguishes rows in each block
 Pick additional nodes to resolve rows across blocks

```
Input: Resolving set r of H_{k-1,a}
Output: Resolving set R of H_{k,a}
function CONSTRUCTRESOLVINGSET(r.a)
    R_0 \leftarrow \{\}
    R_1 \leftarrow \{\}
    i \leftarrow 0
    for w \in r do
        if i < a then
             R_0 \leftarrow R_0 \cup \{iw\}
             if i < (a - 1) then
                 R_0 \leftarrow R_0 \cup \{(i+1)w\}
             end if
        end if
        if i > a then
             R_1 \leftarrow R_1 \cup \{0w\}
        end if
        i \leftarrow (i+2)
    end for
    R \leftarrow (R_0 \cup R_1)
    return R
end function
```

```
Input: Resolving set r of H_{k-1,a}
Output: Resolving set R of H_{k,a}
function CONSTRUCTRESOLVINGSET(r.a)
    R_0 \leftarrow \{\}
    R_1 \leftarrow \{\}
    i \leftarrow 0
    for w \in r do
        if i < a then
             R_0 \leftarrow R_0 \cup \{iw\}
             if i < (a - 1) then
                 R_0 \leftarrow R_0 \cup \{(i+1)w\}
             end if
        end if
        if i > a then
             R_1 \leftarrow R_1 \cup \{0w\}
        end if
        i \leftarrow (i+2)
    end for
    R \leftarrow (R_0 \cup R_1)
    return R
end function
```

• Time complexity starting with $H_{1,a}$: $O(ak^2)$

- (ICH time complexity: $O(a^{3k})$)
- Finds resolving set for $H_{k,a}$ of size O(k)

 $\bullet (|H_{k,a}| = a^k)$

• ICH gives the following resolving set for $H_{3,20}$:

ſ	AAA,	RRR,	NNN,	DDD,	CCC,	QQQ,	EEE,	GGG,	HHH,	CNS,
Į	III,	LLL,	KKK,	MMM,	FFF,	PPP,	SIS,	NST,	TTC,	QPK
1	ARW,	WWD,	MKY,	QYE,	YGL,	HPV,	VFR,	EAG,	KLQ,	SVT,
(DHF,	WMP)

• ICH gives the following resolving set for $H_{3,20}$:

AAA,RRR,NNN,DDD,CCC,QQQ,EEE,GGG,HHH,CNS,III,LLL,KKK,MMM,FFF,PPP,SIS,NST,TTC,QPKARW,WWD,MKY,QYE,YGL,HPV,VFR,EAG,KLQ,SVT,DHF,WMP

• Then five iterations of our algorithm give this resolving set for $H_{8,20}$:

AAAAAAAA,	AAAAAAAR,	AAAAAARA,	AAAAARAA,	AAAARAAA,
AAARAAAA,	ARWAAAAA,	CCCHHHHH,	CCCHHHHI,	CCCHHHIA,
CCCHHIAA,	CCCHIAAA,	CCCIAAAA,	CNSAAAAA,	DDDEEEEE,
DDDEEEEG,	DDDEEEGA,	DDDEEGAA,	DDDEGAAA,	DDDGAAAA,
DHFAAAAA,	EAGAAAAA,	EEEFAAAA,	EEEMFAAA,	EEEMMFAA,
EEEMMMFA,	EEEMMMMF,	EEEMMMMM,	FFFAAAAA,	GGGPPPPP,
GGGPPPPS,	GGGPPPSA,	GGGPPSAA,	GGGPSAAA,	GGGSAAAA,
НННТТТТТ,	HHHTTTTW,	HHHTTTWA,	HHHTTWAA,	HHHTWAAA,
HHHWAAAA,	HPVAAAAA,	IIIVAAAA,	IIIYVAAA,	IIIYYVAA,
IIIYYYVA,	IIIYYYYV,	IIIYYYYY,	KKKAAAAA,	KLQAAAAA,
LLLAAAAA,	MKYAAAAA,	MMMAAAAA,	NNNCCCCC,	NNNCCCCQ,
NNNCCCQA,	NNNCCQAA,	NNNCQAAA,	NNNQAAAA,	NSTAAAAA,
PPPAAAAA,	QPKAAAAA,	QQQKAAAA,	QQQLKAAA,	QQQLLKAA,
QQQLLLKA,	QQQLLLLK,	QQQLLLLL,	QYEAAAAA,	RRRDAAAA,
RRRNDAAA,	RRRNNDAA,	RRRNNNDA,	RRRNNNND,	RRRNNNNN,
SISAAAAA,	SVTAAAAA,	TTCAAAAA,	VFRAAAAA,	WMPAAAAA,
WWDAAAAA,	YGLAAAAA			

• ICH gives the following resolving set for $H_{3,20}$:

AAA,RRR,NNN,DDD,III,LLL,KKK,MMM,ARW,WWD,MKY,QYE,DHF,WMP CNS, QPK SVT, CCC. EEE. GGG. HHH. 000. TTC, FFF, PPP, SIS, NST, YGL, HPV, VFR, EAG, KLQ.

• Then five iterations of our algorithm give this resolving set for $H_{8,20}$:

AAAAAAAA,	AAAAAAAR,	AAAAAARA,	AAAAARAA,	AAAARAAA,
AAARAAAA,	ARWAAAAA,	CCCHHHHH,	CCCHHHHI,	CCCHHHIA,
CCCHHIAA,	CCCHIAAA,	CCCIAAAA,	CNSAAAAA,	DDDEEEEE,
DDDEEEEG,	DDDEEEGA,	DDDEEGAA,	DDDEGAAA,	DDDGAAAA,
DHFAAAAA,	EAGAAAAA,	EEEFAAAA,	EEEMFAAA,	EEEMMFAA,
EEEMMMFA,	EEEMMMMF,	EEEMMMMM,	FFFAAAAA,	GGGPPPPP,
GGGPPPPS,	GGGPPPSA,	GGGPPSAA,	GGGPSAAA,	GGGSAAAA,
НННТТТТТ,	HHHTTTTW,	HHHTTTWA,	HHHTTWAA,	HHHTWAAA,
HHHWAAAA,	HPVAAAAA,	IIIVAAAA,	IIIYVAAA,	IIIYYVAA,
IIIYYYVA,	IIIYYYYV,	IIIYYYYY,	KKKAAAAA,	KLQAAAAA,
LLLAAAAA,	MKYAAAAA,	MMMAAAAA,	NNNCCCCC,	NNNCCCCQ,
NNNCCCQA,	NNNCCQAA,	NNNCQAAA,	NNNQAAAA,	NSTAAAAA,
PPPAAAAA,	QPKAAAAA,	QQQKAAAA,	QQQLKAAA,	QQQLLKAA,
QQQLLLKA,	QQQLLLLK,	QQQLLLLL,	QYEAAAAA,	RRRDAAAA,
RRRNDAAA,	RRRNNDAA,	RRRNNNDA,	RRRNNNND,	RRRNNNN,
SISAAAAA,	SVTAAAAA,	TTCAAAAA,	VFRAAAAA,	WMPAAAAA,
WWDAAAAA.	YGLAAAAA			

■ i.e. octamers may be represented as 82-dimensional vectors!

Manuel.Lladser@Colorado.EDU

Low-Dim. Representation of Symbolic Sequences 18 / 25

WORK IN PROGRESS

Is the metric dimension problem still NP-complete on Hamming graphs?

Is the metric dimension problem still NP-complete on Hamming graphs?

ANSWER: We do not know yet!

Is the metric dimension problem still NP-complete on Hamming graphs?

ANSWER: We do not know yet!

Problem.

Find characterizations of resolvability in $H_{k,a}$

Is the metric dimension problem still NP-complete on Hamming graphs?

ANSWER: We do not know yet!

Problem.

Find characterizations of resolvability in $H_{k,a}$

Observation.

Brute force is not practical to determine whether a set R of k-mers resolves or not $H_{k,a}$ because, for each of the $\Theta(a^{2k})$ pairs (x, y) of k-mers, one would need to check if there is $r \in R$ such that $d(x, r) \neq d(y, r)$ • $Q_k := H_{k,a=2}$ (k-dimensional hypercube)

k	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
$\beta(Q_k)$	1	2	3	4	4	5	6	6	7	7	8	8	8	9	9	10	10

Table: Exact ¹ up to k = 10, conjectured ² up to k = 17

¹Harary and Melter - 1976
 ²Mladenović, Kratica, Kovačević-Vujčić et al. - 2012
 ³Beardon - 2013

• $Q_k := H_{k,a=2}$ (k-dimensional hypercube)

k	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
$\beta(Q_k)$	1	2	3	4	4	5	6	6	7	7	8	8	8	9	9	10	10

Table: Exact ¹ up to k = 10, conjectured ² up to k = 17

There is already a resolvability characterization of hypercubes ³

¹Harary and Melter - 1976
²Mladenović, Kratica, Kovačević-Vujčić et al. - 2012
³Beardon - 2013

• $Q_k := H_{k,a=2}$ (k-dimensional hypercube)

k	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
$\beta(Q_k)$	1	2	3	4	4	5	6	6	7	7	8	8	8	9	9	10	10

Table: Exact ¹ up to k = 10, conjectured ² up to k = 17

There is already a resolvability characterization of hypercubes ³

... but Hamming graphs are transitive

Theorem: (with Laird, Tillquist & Becker)

Suppose $1^k \in R$ and define $A := (-r)_{r \in R}$. Then, R is resolves Q_k iff there does <u>not</u> exist a non-zero $z \in \{0, \pm 1\}^k$ such that Az = 0.

¹Harary and Melter - 1976

²Mladenović, Kratica, Kovačević-Vujčić et al. - 2012

³Beardon - 2013

Theorem: (with Laird, Tillquist & Becker)

Let v_1, \ldots, v_n be k-mers with **one-hot encodings** V_1, \ldots, V_n , respectively. If

$$A := \begin{pmatrix} -\operatorname{vec}(V_1) - \\ \vdots \\ -\operatorname{vec}(V_n) - \end{pmatrix}$$

then $R = \{v_1, \ldots, v_n\}$ resolves $H_{k,a}$ iff z = 0 is the only solution to the linear system Az = 0 which satisfies this additional constraint: if z is decomposed into k non-overlapping blocks of dimension a as follows

$$z = ((z_1, \ldots, z_a), (z_{a+1}, \ldots, z_{2a}), \ldots, (z_{(k-1)a+1}, \ldots, z_{ka}))^T,$$

then each block is the difference of two canonical vectors in \mathbb{R}^a .

Using this new characterization:

Integer linear programming (ILP) implies that

AAAAAAAA,	AAAAAAAR,	AAAAAARA,	AAAAARAA,	AAAARAAA,
AAARAAAA,	ARWAAAAA,	СССННННН,	CCCHHHHI,	CCCHHHIA,
CCCHHIAA,	CCCHIAAA,	CCCIAAAA,	CNSAAAAA,	DDDEEEEE,
DDDEEEEG,	DDDEEEGA,	DDDEEGAA,	DDDEGAAA,	DDDGAAAA,
DHFAAAAA,	EAGAAAAA,	EEEFAAAA,	EEEMFAAA,	EEEMMFAA,
EEEMMMFA,	EEEMMMMF,	EEEMMMMM,	FFFAAAAA,	GGGPPPPP,
GGGPPPPS,	GGGPPPSA,	GGGPPSAA,	GGGPSAAA,	GGGSAAAA,
НННТТТТТ,	HHHTTTTW,	HHHTTTWA,	HHHTTWAA,	HHHTWAAA,
HHHWAAAA,	HPVAAAAA,	IIIVAAAA,	IIIYVAAA,	IIIYYVAA,
IIIYYYVA,	IIIYYYYV,	IIIYYYYY,	KKKAAAAA,	KLQAAAAA,
LLLAAAAA,	MKYAAAAA,	MMMAAAAA,	NNNCCCCC,	NNNCCCCQ,
NNNCCCQA,	NNNCCQAA,	NNNCQAAA,	NNNQAAAA,	NSTAAAAA,
PPPAAAAA,	QPKAAAAA,	QQQKAAAA,	QQQLKAAA,	QQQLLKAA,
QQQLLLKA,	QQQLLLLK,	QQQLLLLL,	QYEAAAAA,	RRRDAAAA,
RRRNDAAA,	RRRNNDAA,	RRRNNNDA,	RRRNNNND,	RRRNNNNN,
SISAAAAA,	SVTAAAAA,	TTCAAAAA,	VFRAAAAA,	WMPAAAAA,
WWDAAAAA.	YGLAAAAA			

resolves $H_{8,20}$ with high probability.

Using this new characterization:

Integer linear programming (ILP) implies that

AAAAAAAA,	AAAAAAAR,	AAAAAARA,	AAAAARAA,	AAAARAAA,
AAARAAAA,	ARWAAAAA,	СССННННН,	CCCHHHHI,	CCCHHHIA,
CCCHHIAA,	CCCHIAAA,	CCCIAAAA,	CNSAAAAA,	DDDEEEEE,
DDDEEEEG,	DDDEEEGA,	DDDEEGAA,	DDDEGAAA,	DDDGAAAA,
DHFAAAAA,	EAGAAAAA,	EEEFAAAA,	EEEMFAAA,	EEEMMFAA,
EEEMMMFA,	EEEMMMMF,	EEEMMMMM,	FFFAAAAA,	GGGPPPPP,
GGGPPPPS,	GGGPPPSA,	GGGPPSAA,	GGGPSAAA,	GGGSAAAA,
НННТТТТТ,	HHHTTTTW,	HHHTTTWA,	HHHTTWAA,	HHHTWAAA,
HHHWAAAA,	HPVAAAAA,	IIIVAAAA,	IIIYVAAA,	IIIYYVAA,
IIIYYYVA,	IIIYYYYV,	IIIYYYYY,	KKKAAAAA,	KLQAAAAA,
LLLAAAAA,	MKYAAAAA,	MMMAAAAA,	NNNCCCCC,	NNNCCCCQ,
NNNCCCQA,	NNNCCQAA,	NNNCQAAA,	NNNQAAAA,	NSTAAAAA,
PPPAAAAA,	QPKAAAAA,	QQQKAAAA,	QQQLKAAA,	QQQLLKAA,
QQQLLLKA,	QQQLLLLK,	QQQLLLLL,	QYEAAAAA,	RRRDAAAA,
RRRNDAAA,	RRRNNDAA,	RRRNNNDA,	RRRNNNND,	RRRNNNNN,
SISAAAAA,	SVTAAAAA,	TTCAAAAA,	VFRAAAAA,	WMPAAAAA,
WWDAAAAA.	YGLAAAAA			

resolves $H_{8,20}$ with high probability.

Gröbner bases imply this set resolves $H_{8,20}$ with certainty

Using this new characterization:

Integer linear programming (ILP) implies that

AAAAAAAA.	AAAAAAAR.	AAAAAARA.	AAAAARAA.	AAAARAAA.
AAARAAAA,	ARWAAAAA,	СССННННН,	СССННННІ,	CCCHHHIA,
CCCHHIAA,	CCCHIAAA,	CCCIAAAA,	CNSAAAAA,	DDDEEEEE,
DDDEEEEG,	DDDEEEGA,	DDDEEGAA,	DDDEGAAA,	DDDGAAAA,
DHFAAAAA,	EAGAAAAA,	EEEFAAAA,	EEEMFAAA,	EEEMMFAA,
EEEMMMFA,	EEEMMMMF,	EEEMMMMM,	FFFAAAAA,	GGGPPPPP,
GGGPPPPS,	GGGPPPSA,	GGGPPSAA,	GGGPSAAA,	GGGSAAAA,
НННТТТТТ,	HHHTTTTW,	HHHTTTWA,	HHHTTWAA,	HHHTWAAA,
HHHWAAAA,	HPVAAAAA,	IIIVAAAA,	IIIYVAAA,	IIIYYVAA,
IIIYYYVA,	IIIYYYYV,	IIIYYYYY,	KKKAAAAA,	KLQAAAAA,
LLLAAAAA,	MKYAAAAA,	MMMAAAAA,	NNNCCCCC,	NNNCCCCQ,
NNNCCCQA,	NNNCCQAA,	NNNCQAAA,	NNNQAAAA,	NSTAAAAA,
PPPAAAAA,	QPKAAAAA,	QQQKAAAA,	QQQLKAAA,	QQQLLKAA,
QQQLLLKA,	QQQLLLLK,	QQQLLLLL,	QYEAAAAA,	RRRDAAAA,
RRRNDAAA,	RRRNNDAA,	RRRNNNDA,	RRRNNNND,	RRRNNNNN,
SISAAAAA,	SVTAAAAA,	TTCAAAAA,	VFRAAAAA,	WMPAAAAA,
WWDAAAAA.	YGLAAAAA			

resolves $H_{8,20}$ with high probability.

- Gröbner bases imply this set resolves H_{8,20} with certainty
- So octamers may be represented as 77-dimensional vectors!

COLLABORATORS:







Lucas LairdRichard C. TillquistStephen BeckerAPPM/CSCIIQ Bio/CSCIAPPM

FUNDING:

- BioFrontiers Institute Computing Core
- NSF Grant No. 1836914 (PI)

... Thank You!

Proof of Concept

- Fruit fly genome ($\sim 1.75 \times 10^8$ base-pairs)
- Problem: Characterize 20-mers centered at intron-exon boundaries
- Pos. examples: ~87K; Neg. examples: Random 20-mers



Figure: Cross-validation of KNN with 3 different embeddings