

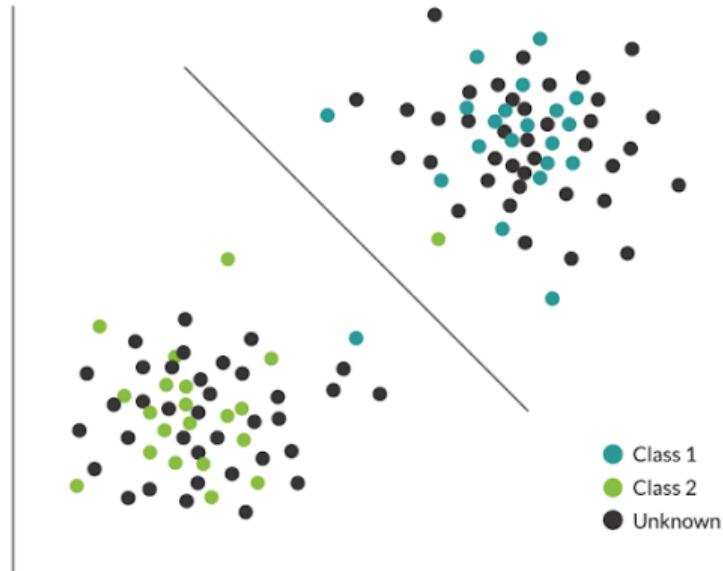
Embeddings: learning representations for unsupervised learning

Marc Lelarge Kevin Scaman Jill-Jênn Vie

Oct 28, 2022

Supervised learning

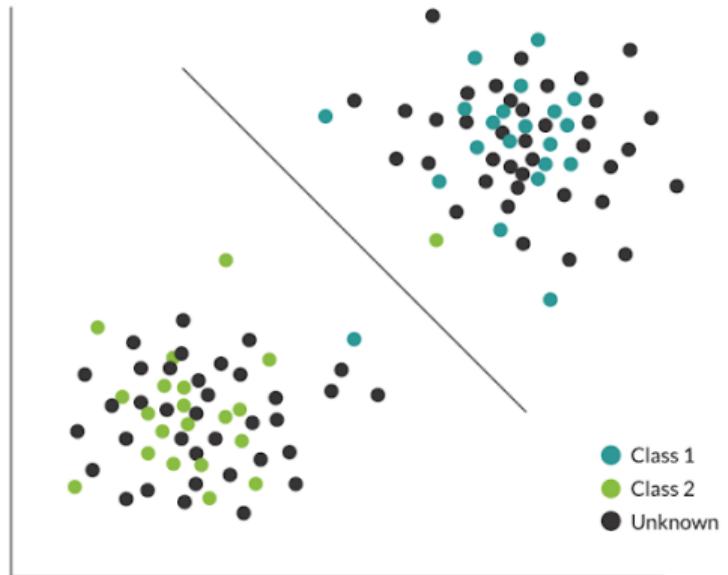
We observe x, y
i.e. classification, regression



Unsupervised learning

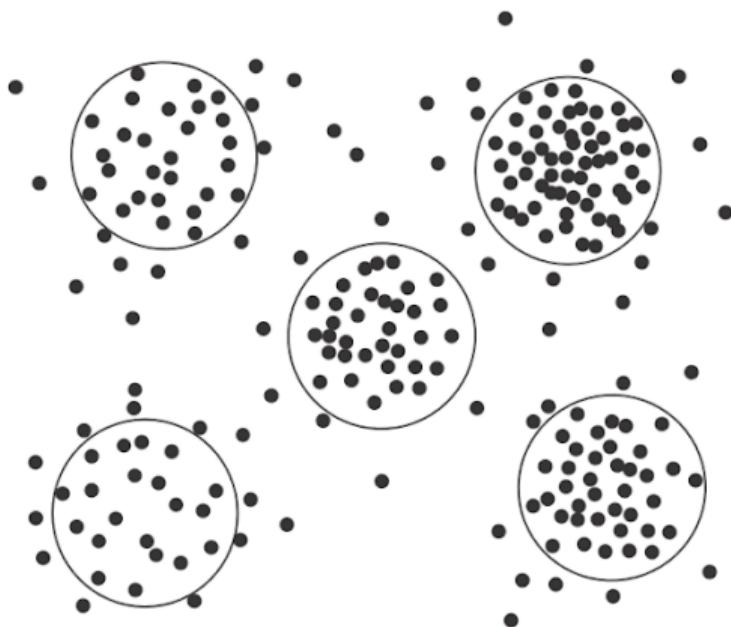
Supervised learning

We observe x, y
i.e. classification, regression



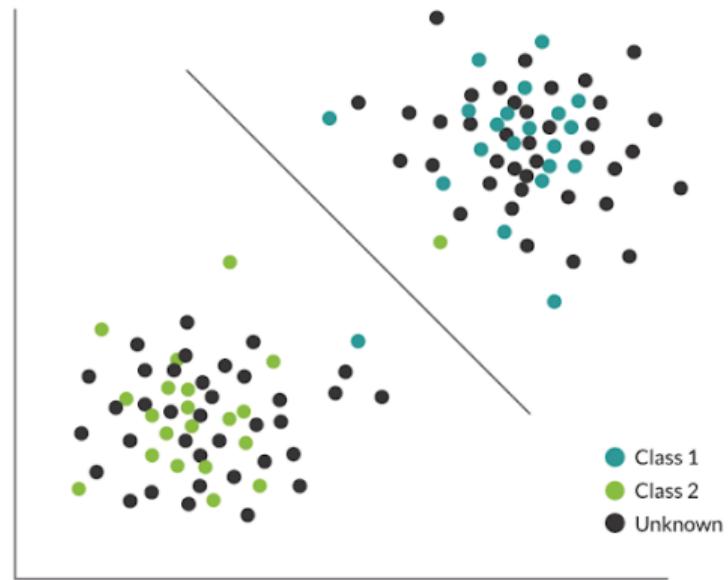
Unsupervised learning

We just observe x
i.e. raw text, music, ratings



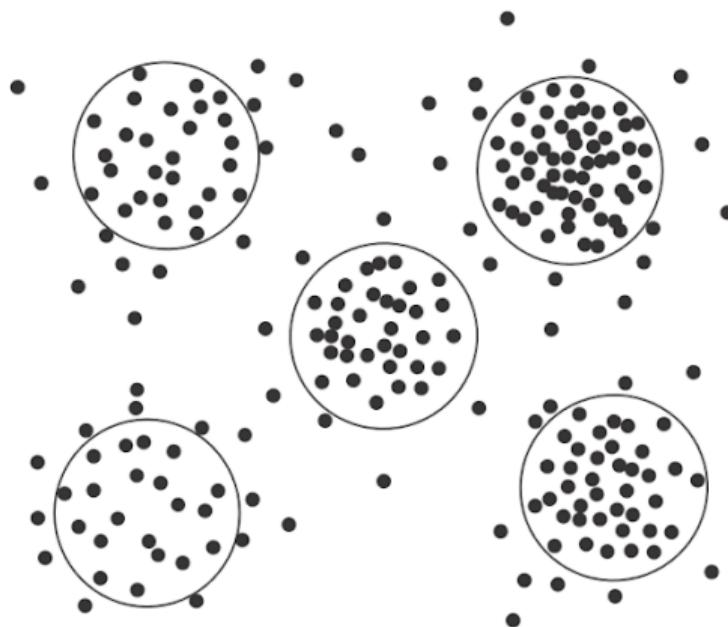
Supervised learning

We observe x, y
i.e. classification, regression



Unsupervised learning

We just observe x
i.e. raw text, music, ratings



It all deals with learning a representation of the data: embedding entities into \mathbf{R}^d

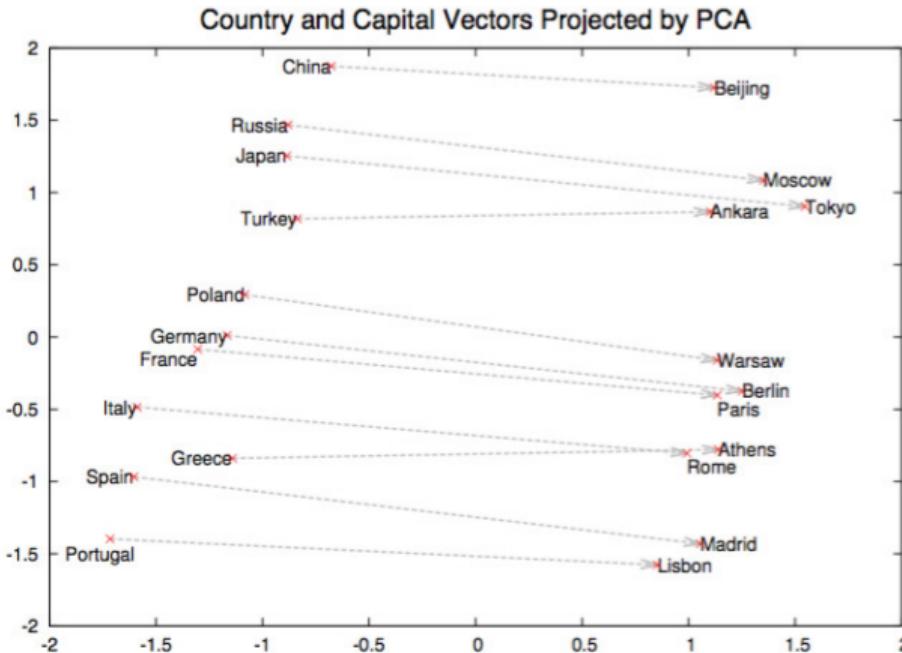
Context

- ▶ Dealing with unlabeled data, like raw text
- ▶ Dealing with tabular data, which mixes **categorical** and continuous variables
- ▶ Detecting outliers in the dataset
- ▶ Or mistakes in labels

Natural language processing

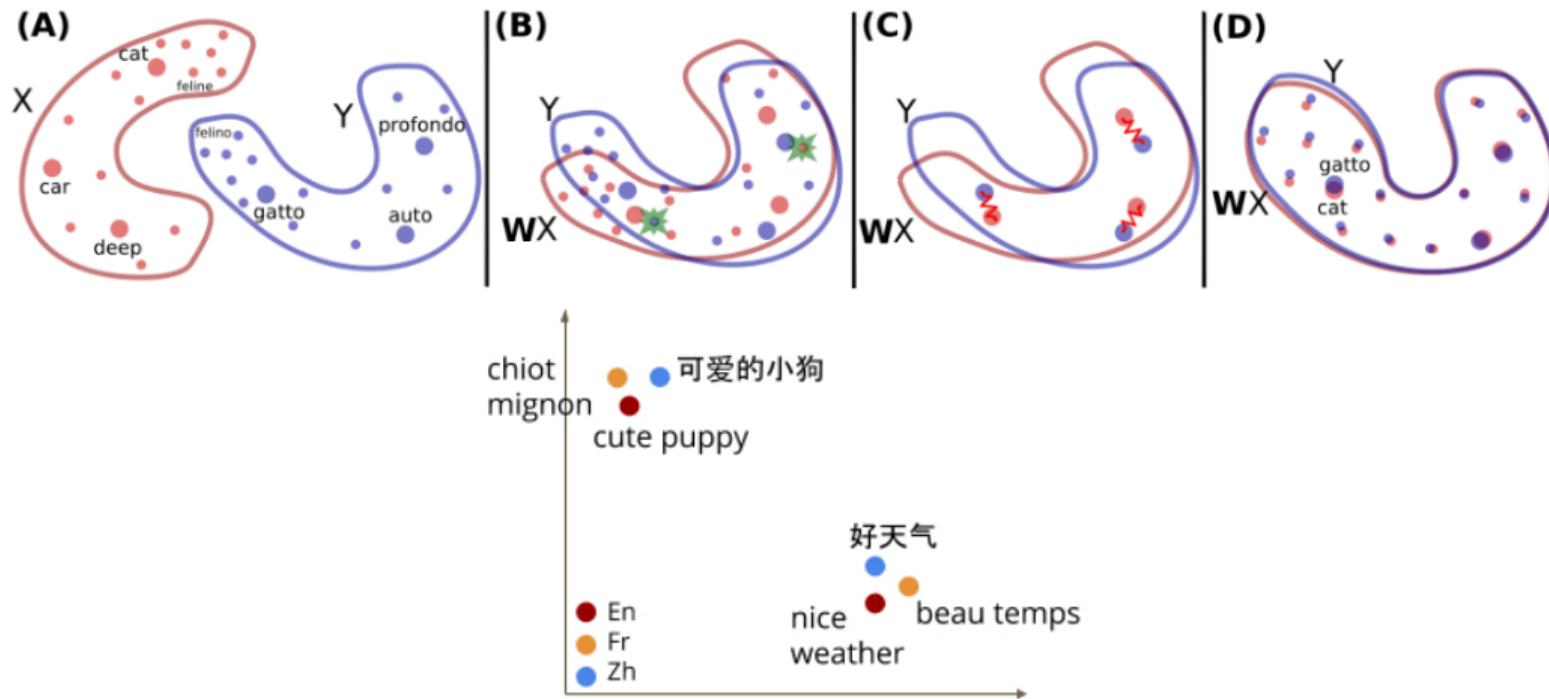
- ▶ A document can be represented as a sparse vector \mathbf{x} of size 60,000 where x_i is the number of occurrences of the i th word.
- ▶ Could we guess meaning from context like humans?

Learned embeddings have interesting properties: word2vec



Tomas Mikolov et al. "Distributed Representations of Words and Phrases and their Compositionality". In: *Advances in Neural Information Processing Systems*. Ed. by C.J. Burges et al. Vol. 26. Curran Associates, Inc., 2013. URL: <https://proceedings.neurips.cc/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf>

Multilingual word embeddings (unsupervised)



Guillaume Lample et al. "Word translation without parallel data". In: *International Conference on Learning Representations*. 2018. URL: <https://openreview.net/forum?id=H196sainb>

Application: Recommender Systems

Problem

- ▶ Every user rates few items (1 %)
- ▶ How to infer missing ratings?

Example



Sacha	?	5	2	?
Ondine	4	1	?	5
Pierre	3	3	1	4
Joëlle	5	?	2	?

Application: Recommender Systems

Problem

- ▶ Every user rates few items (1 %)
- ▶ How to infer missing ratings?

Example



Sacha	3	5	2	2
Ondine	4	1	4	5
Pierre	3	3	1	4
Joëlle	5	2	2	5

Encoding the problem

Matrix completion (unsupervised)

Given sparse entries M_{ij} find the other entries M_{ij}

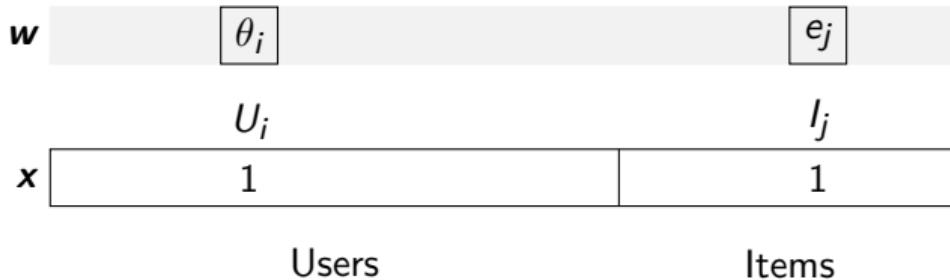
Supervised problem

- ▶ $\mathbf{X} = \{(i,j)\}$ are the categorical observed indices (ex. user i and item j)
- ▶ $\mathbf{y} = \{R_{ij}\}$ are the continuous (or categorical) ratings

Learning embeddings

Embedding n entities into \mathbf{R}^d is equivalent to learning a $n \times d$ matrix \mathbf{W} s.t. embedding of entity i is W_i .

One-hot encoding $\mathbf{e}_j = (0, \dots, 0, \underbrace{1}_j, 0, \dots, 0)$ so that $\mathbf{We}_j = W_j$



Dimensionality reduction, or best approximation of rank k

Input: $n \times m$ matrix \mathbf{X}

Minimize: $\|\mathbf{X} - M\|_2^2 = \sum_{i,j} (X_{ij} - M_{ij})^2$ where $\text{rank}(M) = k$

Dimensionality reduction, or best approximation of rank k

Input: $n \times m$ matrix \mathbf{X}

Minimize: $\|\mathbf{X} - M\|_2^2 = \sum_{i,j} (X_{ij} - M_{ij})^2$ where $\text{rank}(M) = k$

Solution: compute the singular value decomposition (SVD) $\mathbf{X} = U\Sigma V^T$ where

- ▶ U ($n \times r$) is orthogonal: $U^T U = I_r$ where $r = \text{rank}(\mathbf{X})$
- ▶ Σ ($r \times r$) is diagonal: its values $\lambda_1 \geq \dots \geq \lambda_r > 0$
- ▶ V ($m \times r$) is orthogonal: $V^T V = I_r$

k -SVD is truncated to k biggest singular values. $M = U_k \Sigma_k V_k^T$ is solution. And $W = U_k \Sigma_k$ is a $n \times k$ matrix of embeddings (pca.transform).

KNN → measure similarity between users (or items)

K -nearest neighbors

- ▶ R_u represents the row vector of user u in the rating matrix (users \times items).
- ▶ Similarity score between users (cosine):

$$\text{score}(u, v) = \frac{R_u \cdot R_v}{\|R_u\| \cdot \|R_v\|}.$$

- ▶ Let's identify the k -nearest neighbors of user u
- ▶ And recommend to user u what u 's neighbors liked, but u didn't watch

KNN → measure similarity between users (or items)

K -nearest neighbors

- ▶ R_u represents the row vector of user u in the rating matrix (users \times items).
- ▶ Similarity score between users (cosine):

$$\text{score}(u, v) = \frac{R_u \cdot R_v}{\|R_u\| \cdot \|R_v\|}.$$

- ▶ Let's identify the k -nearest neighbors of user u
- ▶ And recommend to user u what u 's neighbors liked, but u didn't watch

Hint

If R' the $N \times M$ matrix of rows $\frac{R_u}{\|R_u\|}$, we can get the $N \times N$ score matrix by computing $R'R'^T$. May not fit in memory.

Matrix factorization → reduce dimension to generalize

$$R = \begin{pmatrix} R_1 \\ R_2 \\ \vdots \\ R_n \end{pmatrix} = \boxed{\quad} = \boxed{U} \boxed{P}$$

$$R: 600 \text{ users} \times 9k \text{ works} \iff \begin{cases} U: 600 \text{ users} \times 20 \text{ profiles} \\ P: 20 \text{ profiles} \times 9k \text{ works} \end{cases}$$

\mathbf{R}_{Bob} is a linear combination of profiles P_1, P_2, \dots

Matrix factorization → reduce dimension to generalize

$$R = \begin{pmatrix} R_1 \\ R_2 \\ \vdots \\ R_n \end{pmatrix} = \boxed{\quad} = \boxed{U} \boxed{P}$$

$$R: 600 \text{ users} \times 9k \text{ works} \iff \begin{cases} U: 600 \text{ users} \times 20 \text{ profiles} \\ P: 20 \text{ profiles} \times 9k \text{ works} \end{cases}$$

R_{Bob} is a linear combination of profiles P_1, P_2 , etc.

Interpreting Key Profiles

If P P_1 : adventure P_2 : romance P_3 : plot twist

And U_{Bob} 0.2 -0.5 0.6

⇒ Bob likes a bit adventure, hates romance, loves plot twists.

Matrix factorization → reduce dimension to generalize

Train

- ▶ R ratings, $\textcolor{red}{U}$ user embeddings, $\textcolor{red}{V}$ item embeddings.

$$R = \textcolor{red}{U} \textcolor{red}{V}^T \quad \hat{r}_{ij} = \mathbf{u}_i^T \mathbf{v}_j$$

Matrix factorization → reduce dimension to generalize

Train

- ▶ R ratings, $\textcolor{red}{U}$ user embeddings, $\textcolor{red}{V}$ item embeddings.

$$R = \textcolor{red}{U} \textcolor{red}{V}^T \quad \hat{r}_{ij} = \mathbf{u}_i^T \mathbf{v}_j$$

Test: Will user i' like item j' ?

- ▶ Just compute $\mathbf{u}_{i'}^T \mathbf{v}_{j'}$

Objective: regularized squared loss

$$\mathcal{L} = \sum_{i,j} (\underbrace{\mathbf{u}_i^T \mathbf{v}_j}_{\text{pred}} - \underbrace{r_{ij}}_{\text{real}})^2 + \underbrace{\lambda \|\mathbf{u}_i\|_2^2 + \lambda \|\mathbf{v}_j\|_2^2}_{\text{regularization}}$$

Variants

SVD: $\sum_{\text{all } i,j} (\mathbf{u}_i^T \mathbf{v}_j - r_{ij})^2$ (missing pairs are imputed to zero)

ALS-WR : $\sum_{i,j \text{ known}} (\mathbf{u}_i^T \mathbf{v}_j - r_{ij})^2 + \lambda (\sum_i N_i \|\mathbf{u}_i\|^2 + \sum_j M_j \|\mathbf{v}_j\|^2)$

where N_i (resp. M_j): how many times user i rated items (resp. item j was rated)

Computing the gradients

$$\mathcal{L} = \sum_{i,j} (\underbrace{\mathbf{u}_i^T \mathbf{v}_j}_{\text{pred}} - \underbrace{r_{ij}}_{\text{real}})^2 + \underbrace{\lambda \|\mathbf{u}_i\|_2^2 + \lambda \|\mathbf{v}_j\|_2^2}_{\text{regularization}}$$

Computing the gradients

$$\mathcal{L} = \sum_{i,j} (\underbrace{\mathbf{u}_i^T \mathbf{v}_j}_{\text{pred}} - \underbrace{r_{ij}}_{\text{real}})^2 + \underbrace{\lambda \|\mathbf{u}_i\|_2^2 + \lambda \|\mathbf{v}_j\|_2^2}_{\text{regularization}}$$

\mathcal{L} is (polynomial and) convex in each \mathbf{u}_k or \mathbf{v}_ℓ .

Finding \mathbf{u}_k that minimizes $\mathcal{L} \Rightarrow$ Finding the zeroes of

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}_k} = \sum_{j \text{ rated by } k} 2(\mathbf{u}_k^T \mathbf{v}_j - r_{kj}) \mathbf{v}_j + 2\lambda \mathbf{u}_k = 0$$

can be rewritten $A\mathbf{u}_k = B$ so $\mathbf{u}_k = A^{-1}B$ (closed form)

Complexity: $O(d^3)$ where d is the (embedding) size of A (but can be parallelized)

Computing the gradients

$$\mathcal{L} = \sum_{i,j} (\underbrace{\mathbf{u}_i^T \mathbf{v}_j - r_{ij}}_{\text{pred}})^2 + \underbrace{\lambda \|\mathbf{u}_i\|_2^2 + \lambda \|\mathbf{v}_j\|_2^2}_{\text{regularization}}$$

\mathcal{L} is (polynomial and) convex in each \mathbf{u}_k or \mathbf{v}_ℓ .

Finding \mathbf{u}_k that minimizes $\mathcal{L} \Rightarrow$ Finding the zeroes of

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}_k} = \sum_{j \text{ rated by } k} 2(\mathbf{u}_k^T \mathbf{v}_j - r_{kj}) \mathbf{v}_j + 2\lambda \mathbf{u}_k = 0$$

can be rewritten $A\mathbf{u}_k = B$ so $\mathbf{u}_k = A^{-1}B$ (closed form)

Complexity: $O(d^3)$ where d is the (embedding) size of A (but can be parallelized)

Algorithm **ALS**: Alternating Least Squares (Zhou, 2008)

- ▶ Until convergence (few iterations):
 - ▶ Fix U (users) learn V (items) in order to minimize \mathcal{L}
 - ▶ Fix V find U

Illustration of Alternating Least Squares

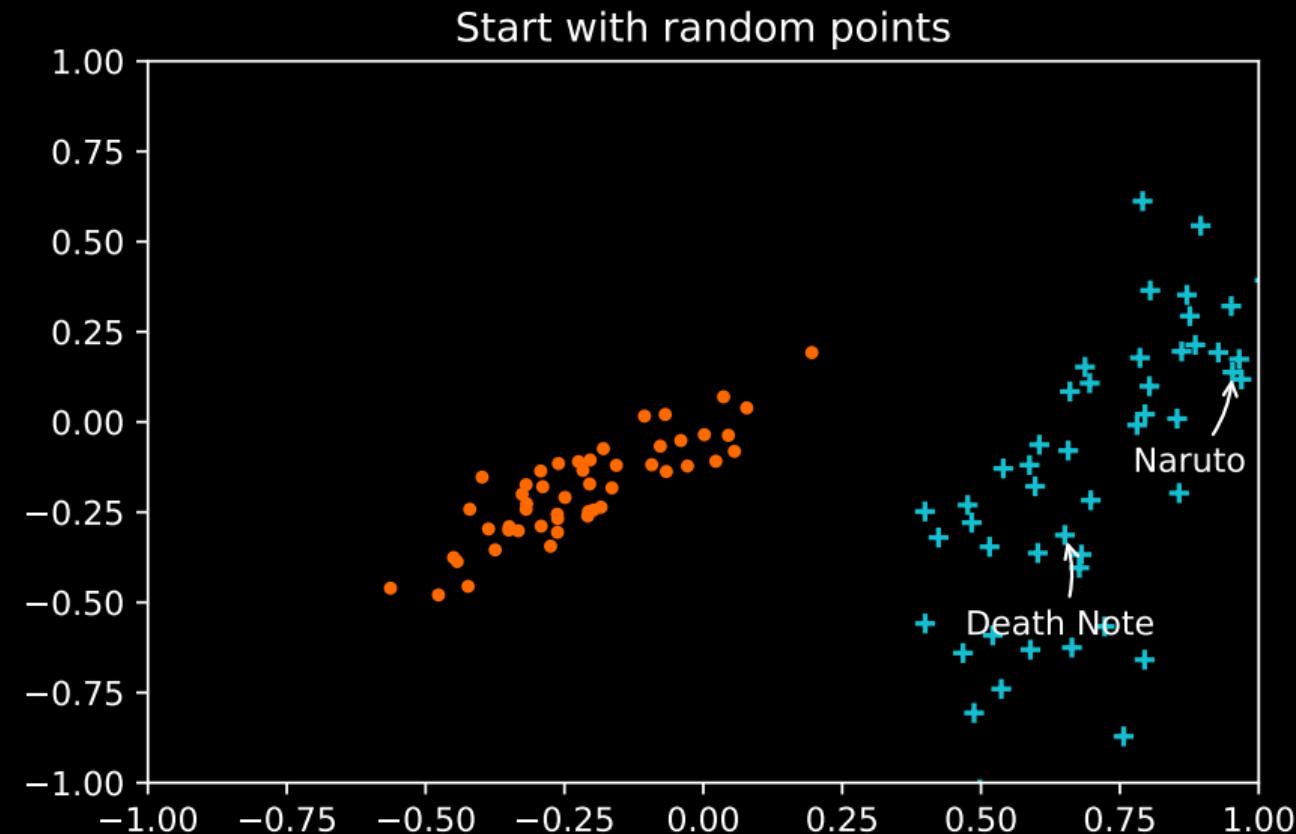


Illustration of Alternating Least Squares

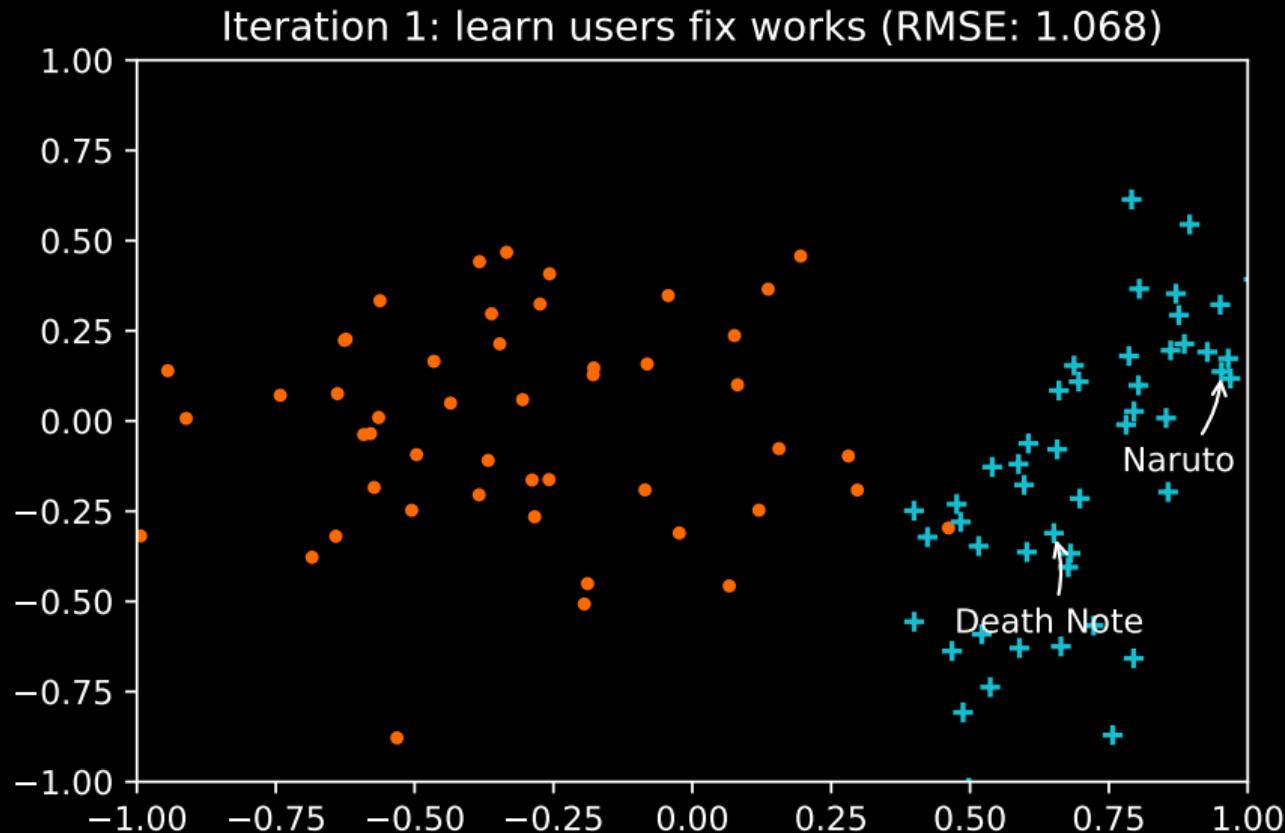


Illustration of Alternating Least Squares

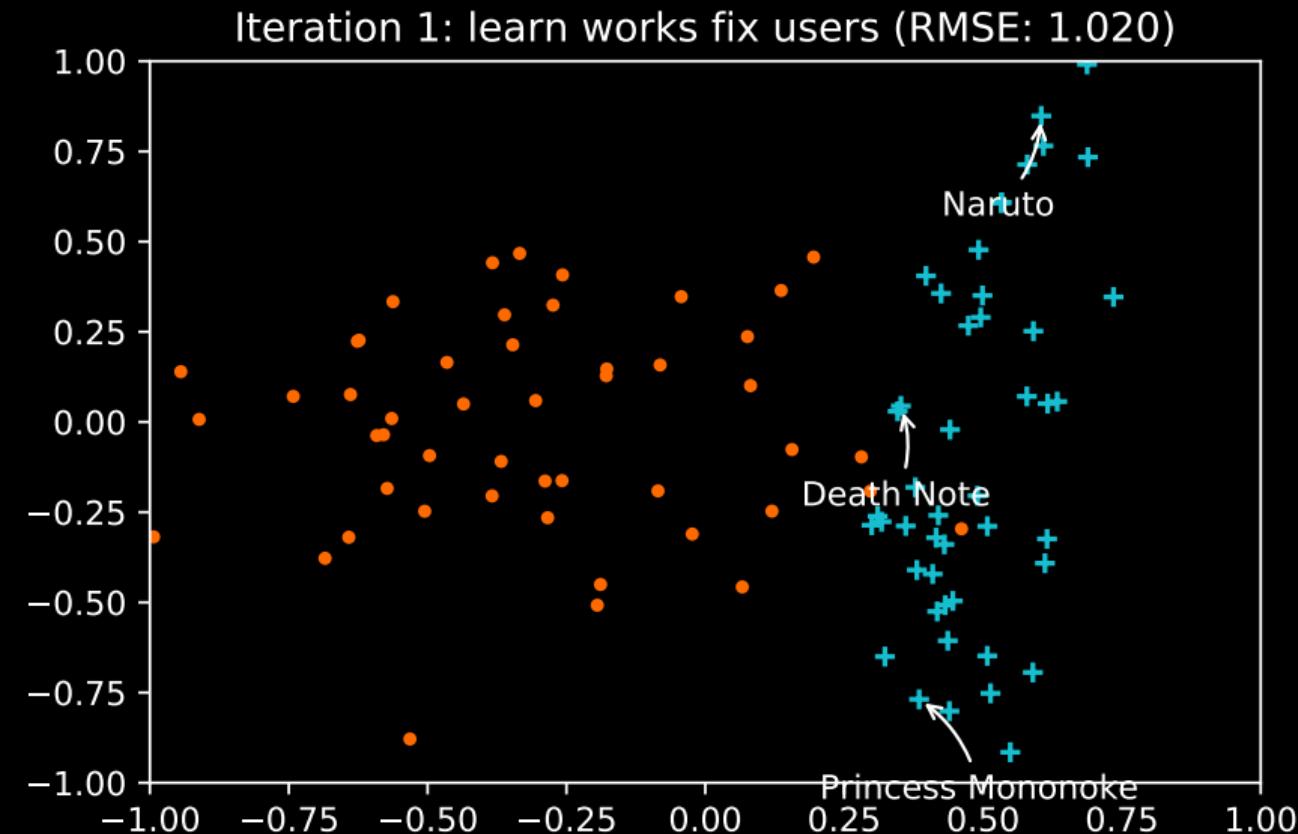


Illustration of Alternating Least Squares

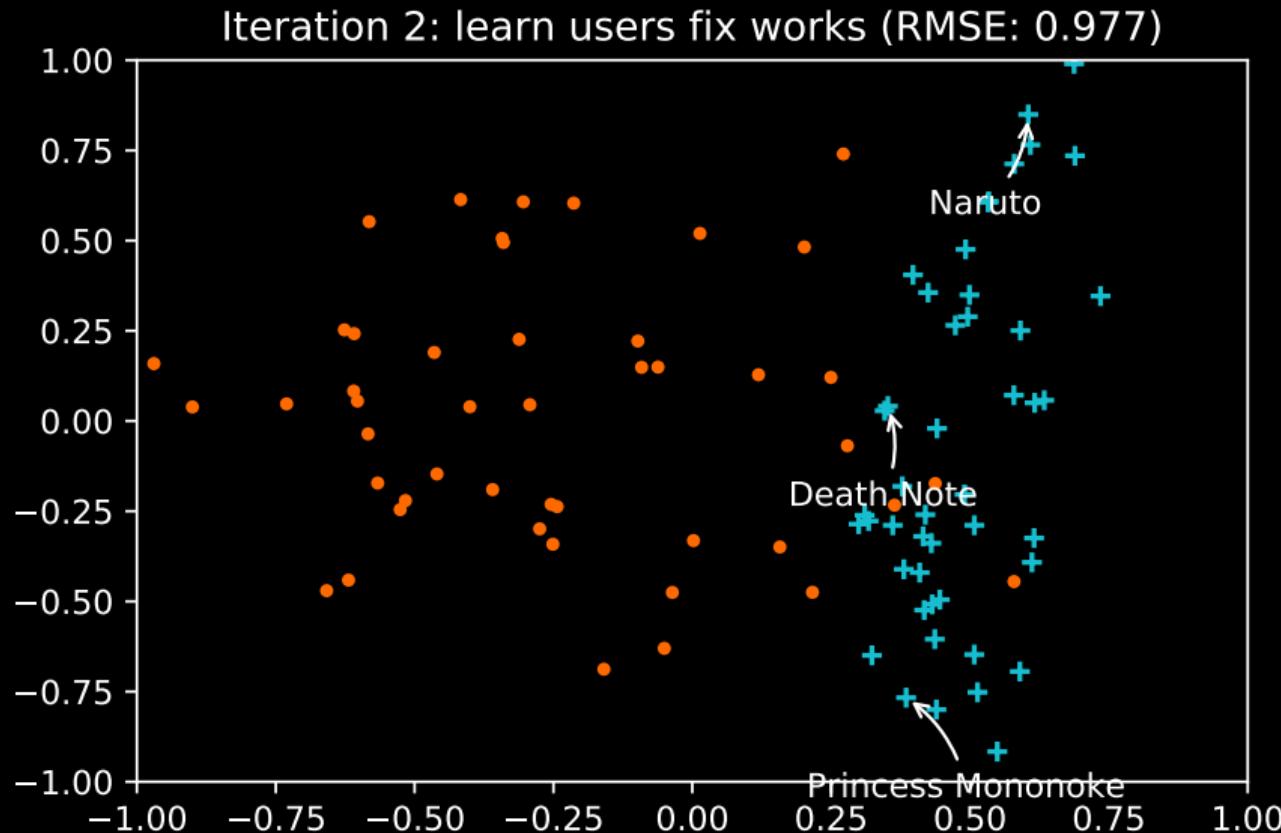


Illustration of Alternating Least Squares

Iteration 2: learn works fix users (RMSE: 0.963)

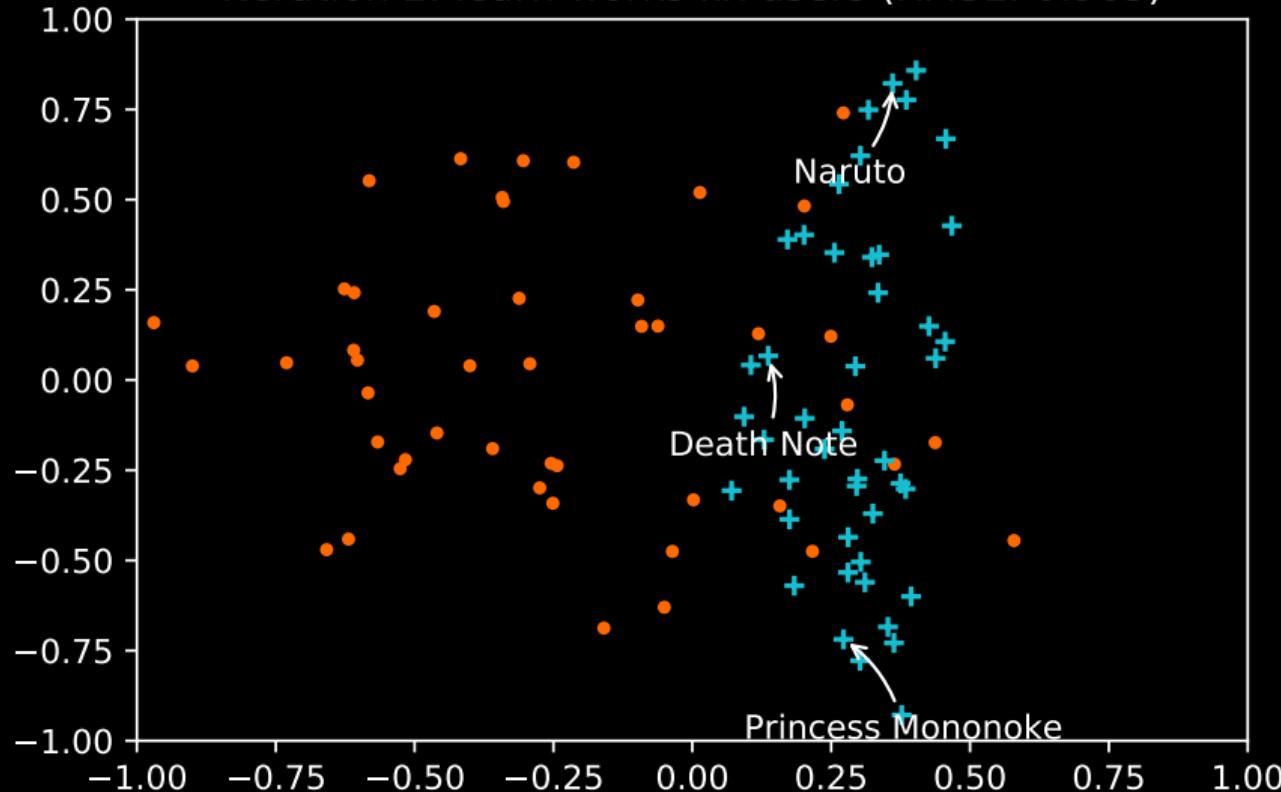


Illustration of Alternating Least Squares

Iteration 3: learn users fix works (RMSE: 0.945)

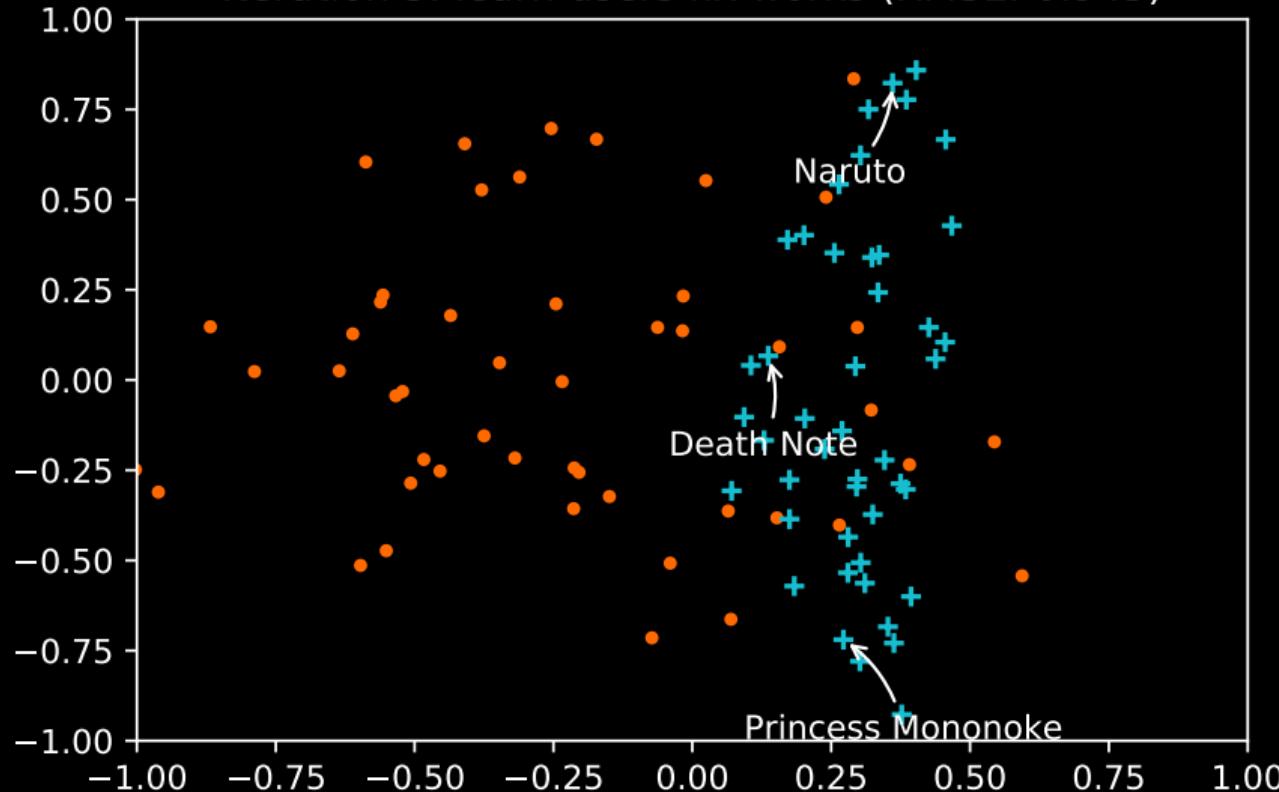


Illustration of Alternating Least Squares

Iteration 3: learn works fix users (RMSE: 0.940)

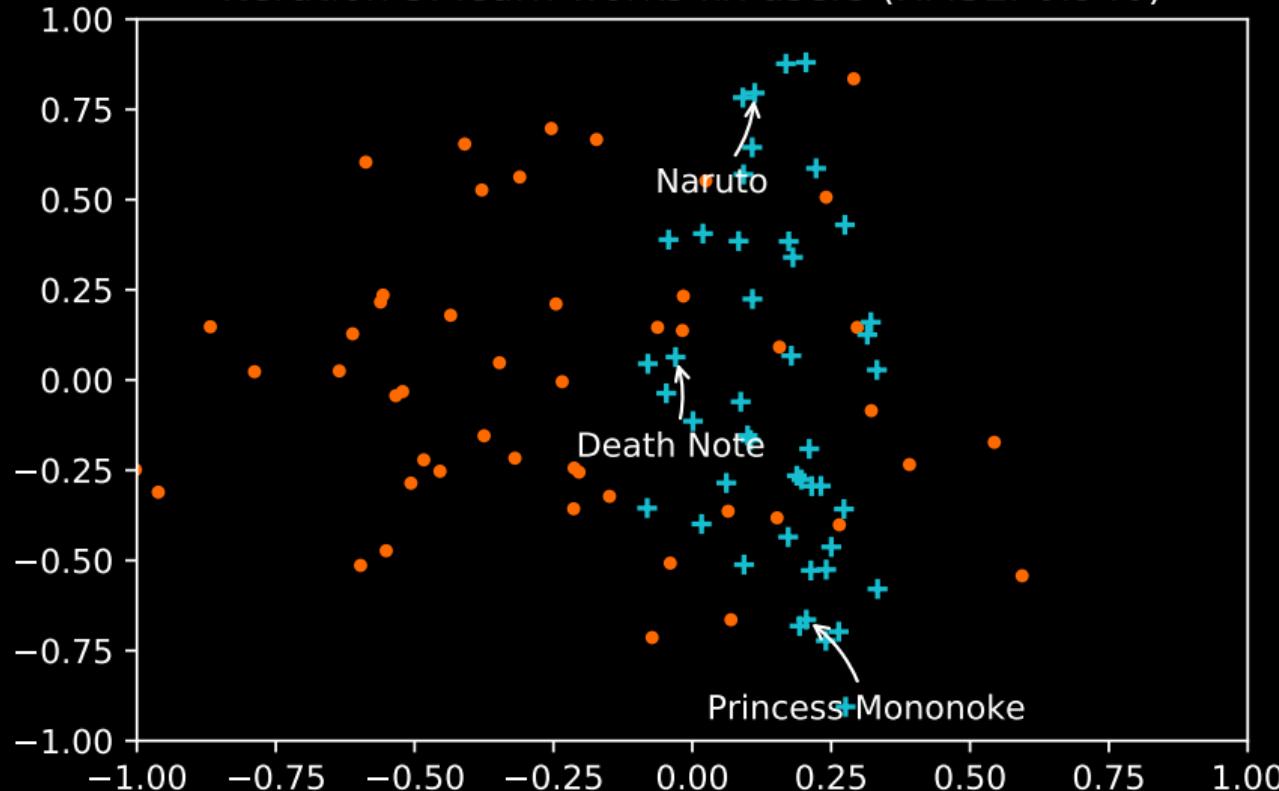


Illustration of Alternating Least Squares

Iteration 4: learn users fix works (RMSE: 0.931)

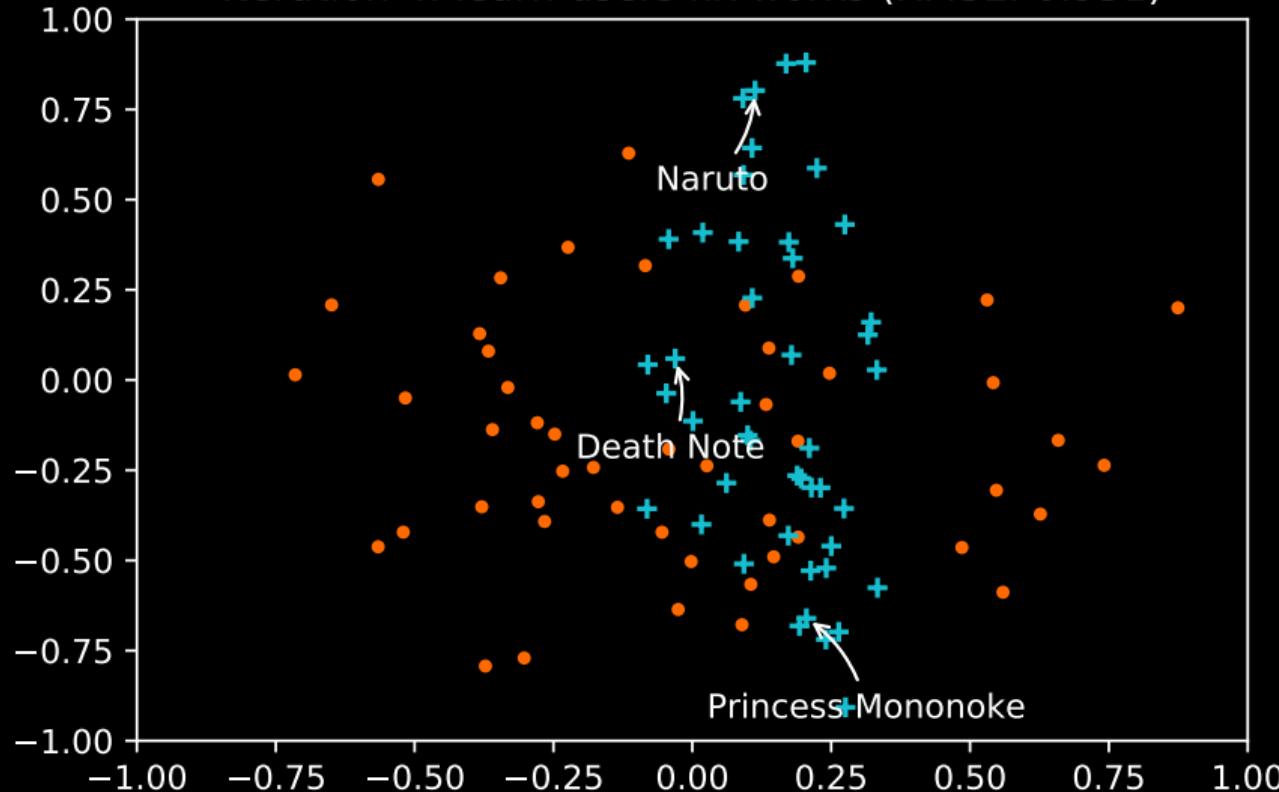


Illustration of Alternating Least Squares

Iteration 4: learn works fix users (RMSE: 0.928)

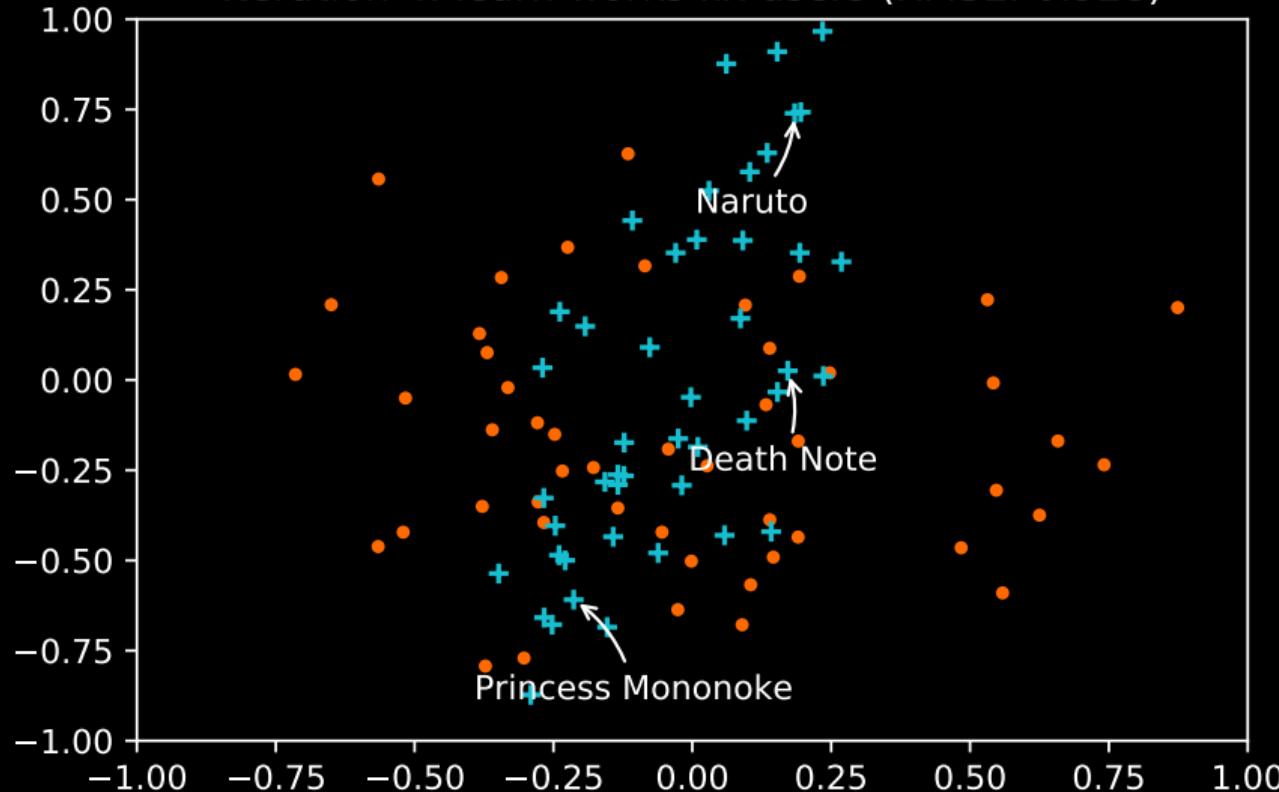


Illustration of Alternating Least Squares

Iteration 5: learn users fix works (RMSE: 0.923)

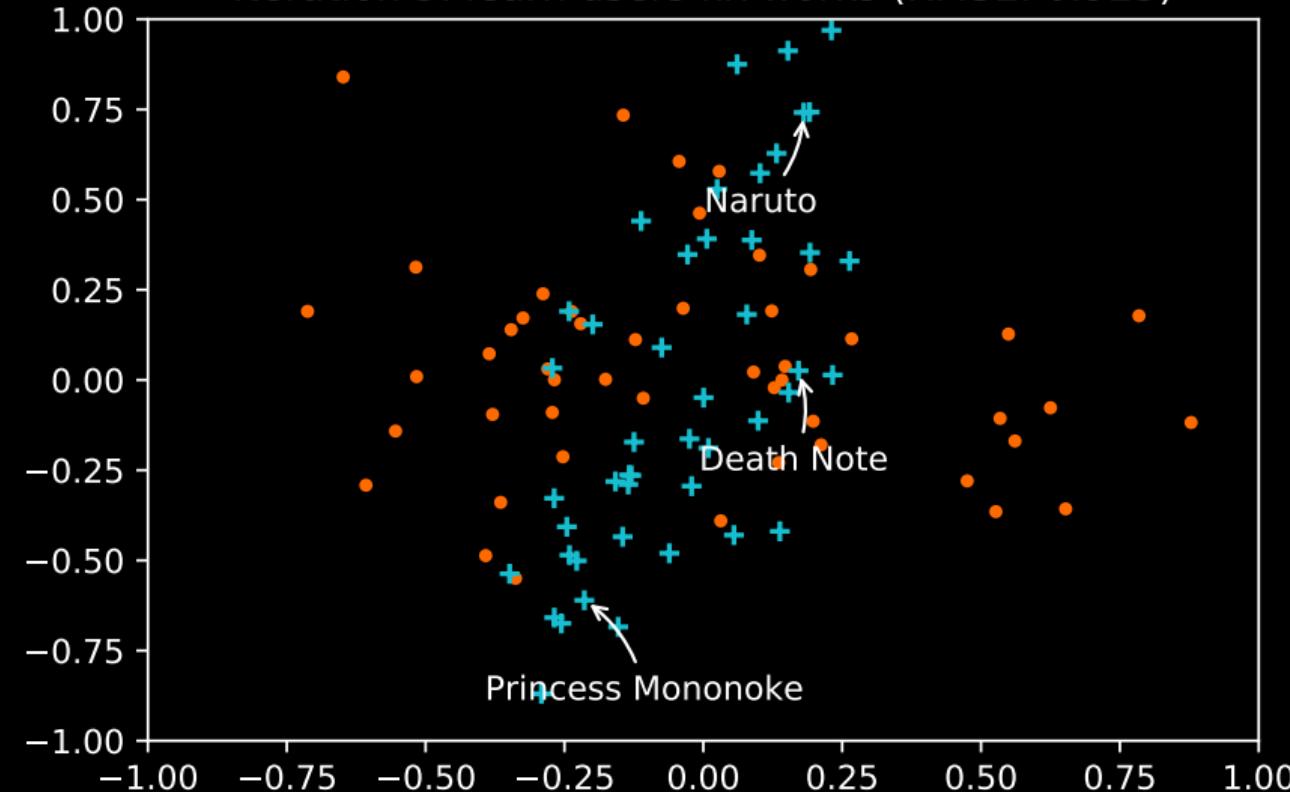


Illustration of Alternating Least Squares

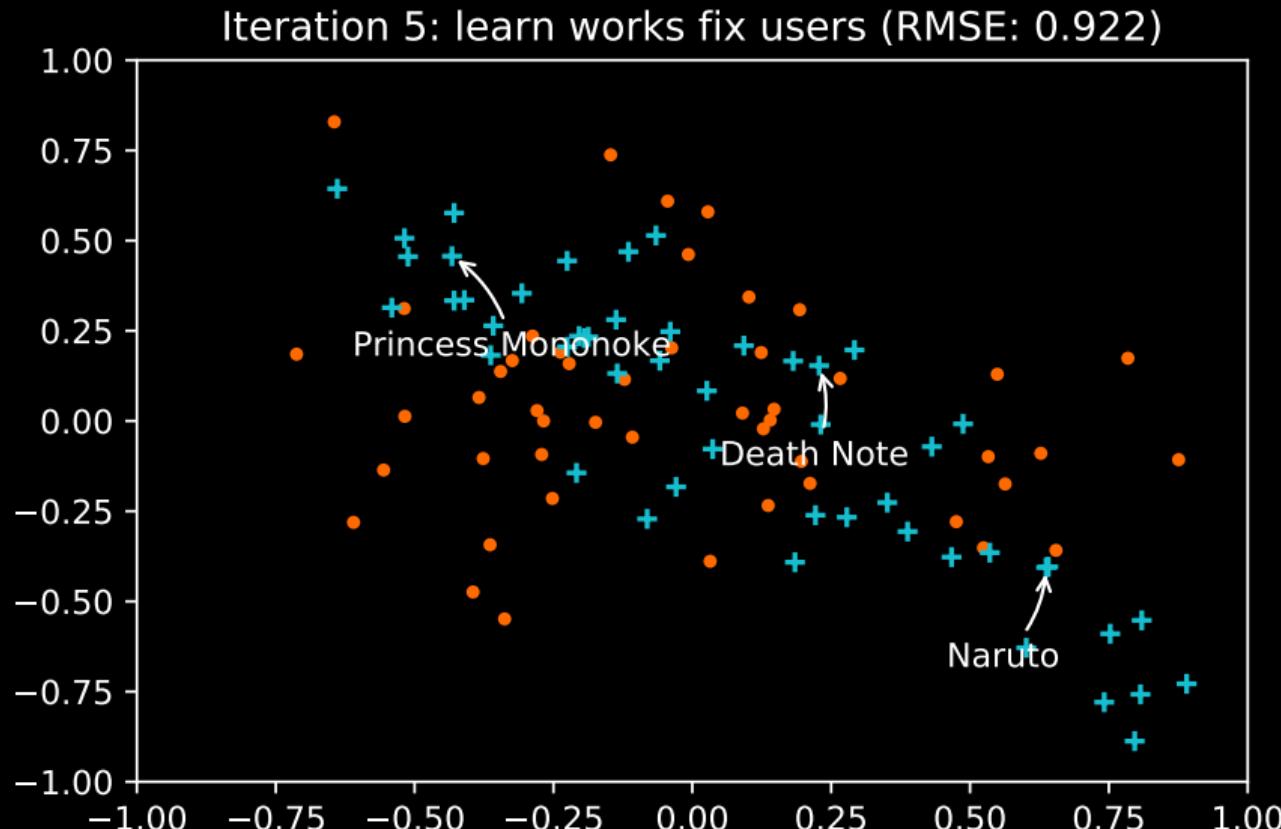


Illustration of Alternating Least Squares

Iteration 6: learn users fix works (RMSE: 0.918)

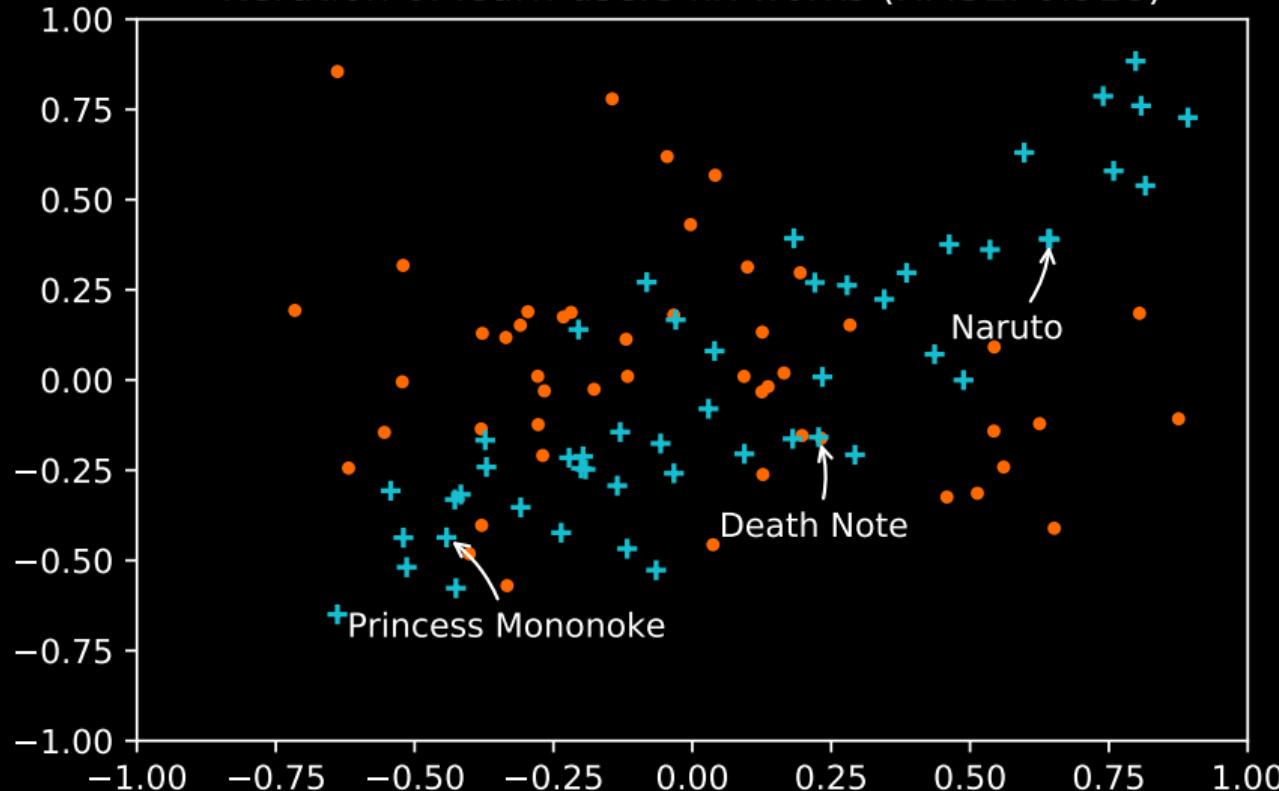


Illustration of Alternating Least Squares

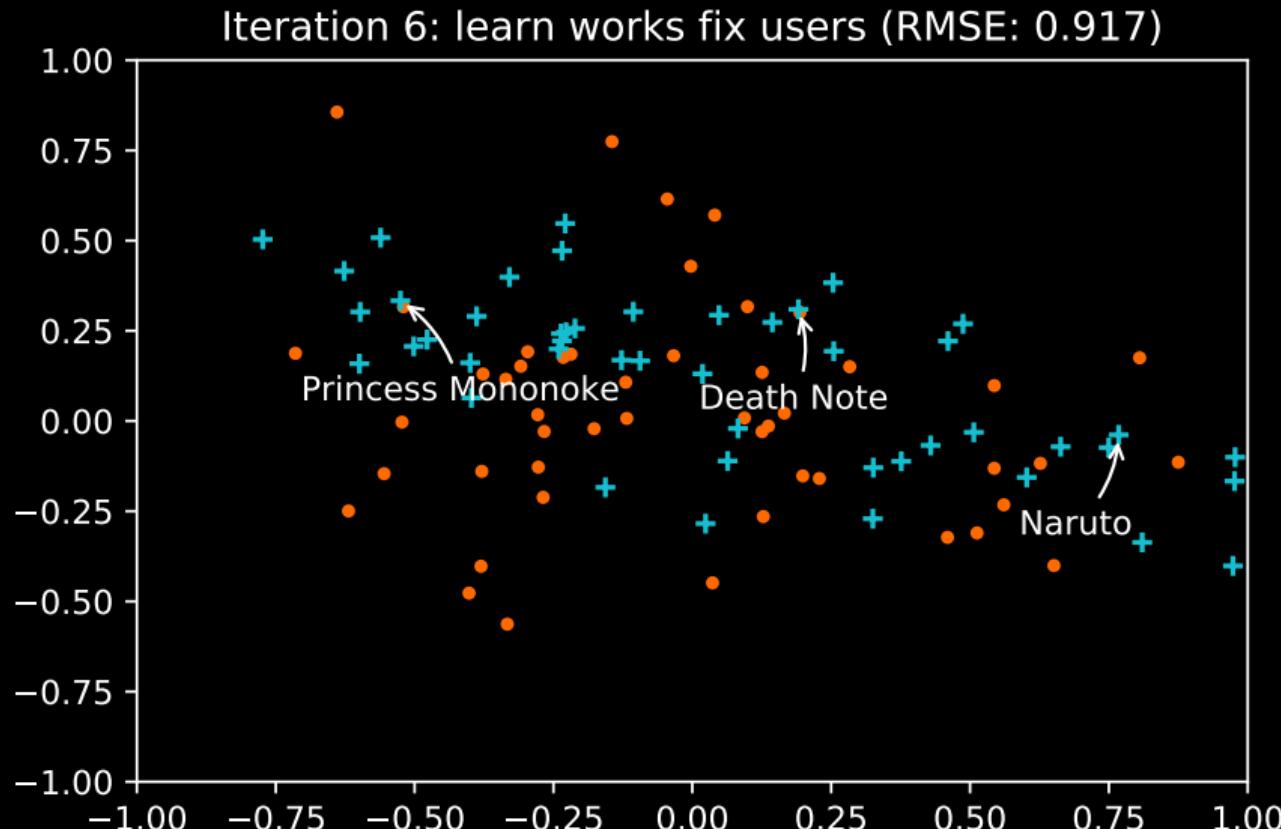


Illustration of Alternating Least Squares

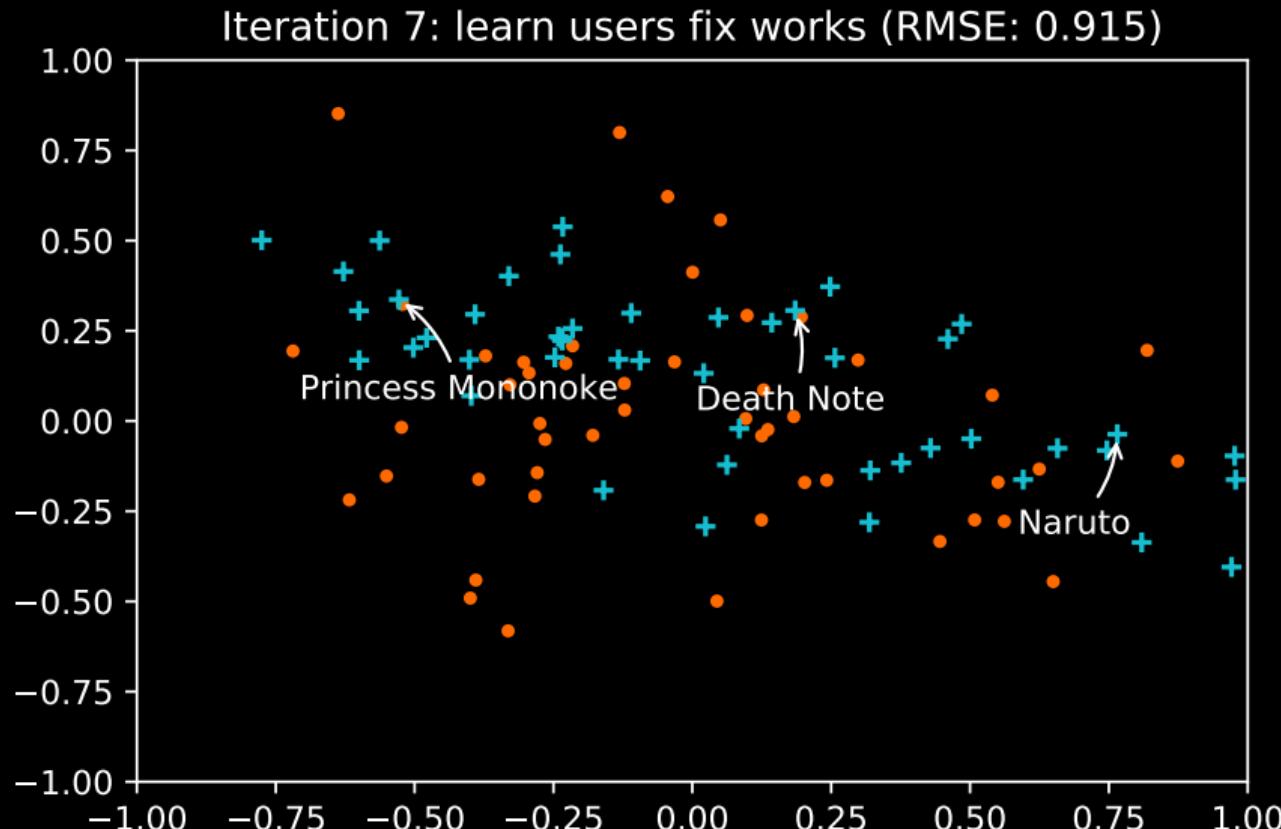


Illustration of Alternating Least Squares

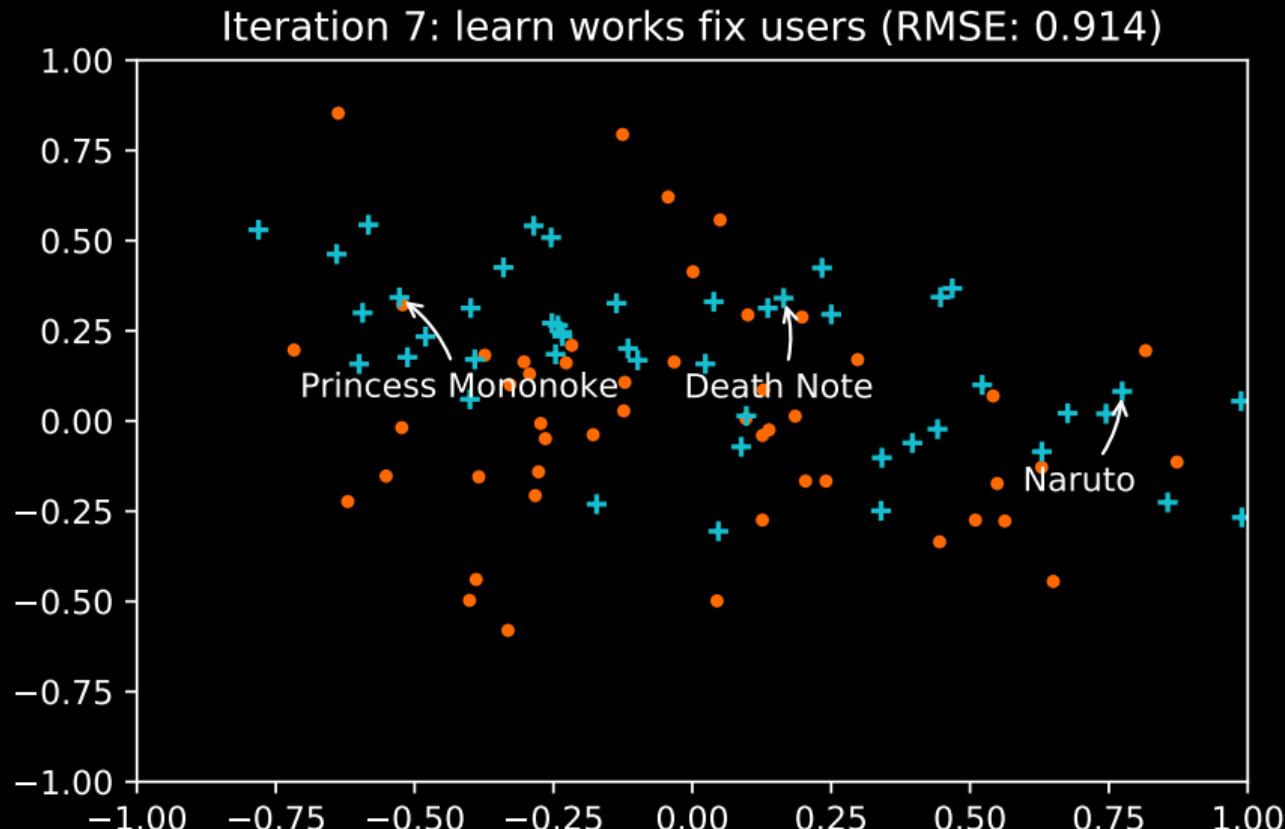


Illustration of Alternating Least Squares

Iteration 8: learn users fix works (RMSE: 0.913)

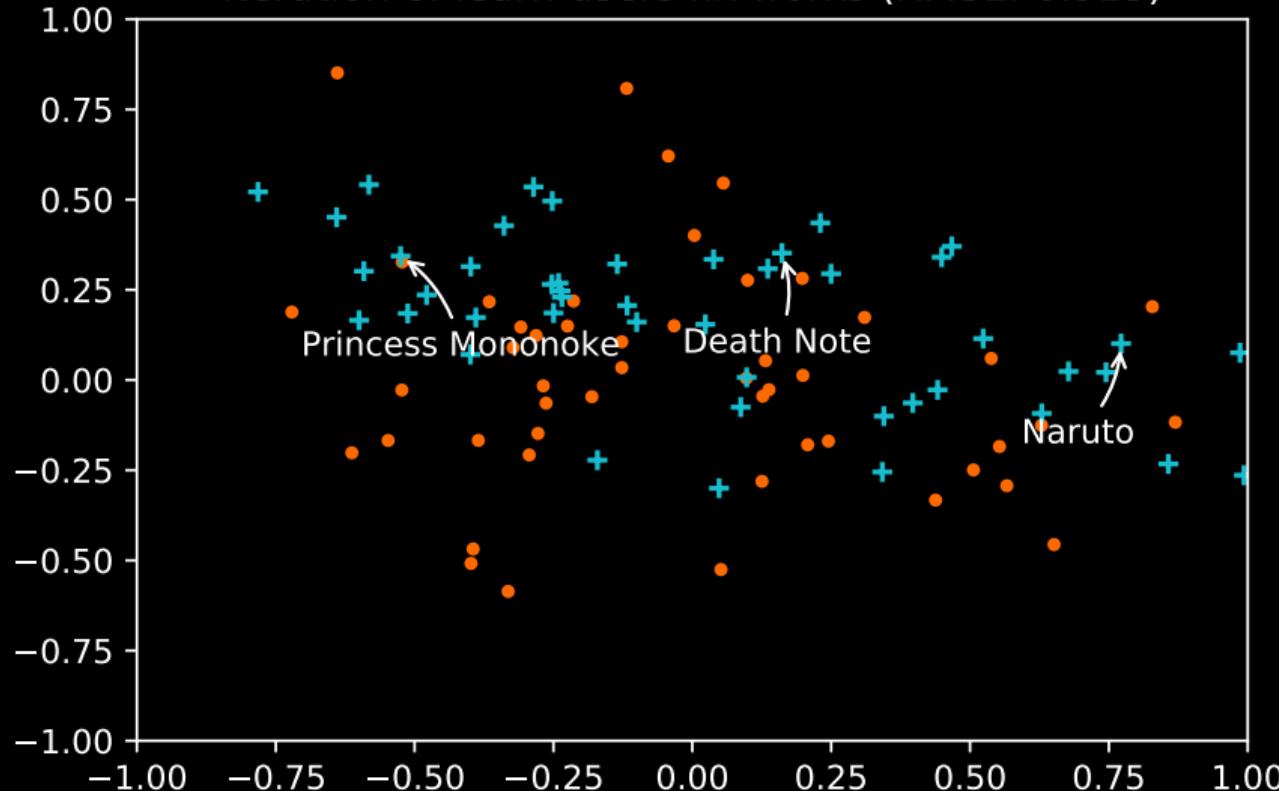


Illustration of Alternating Least Squares

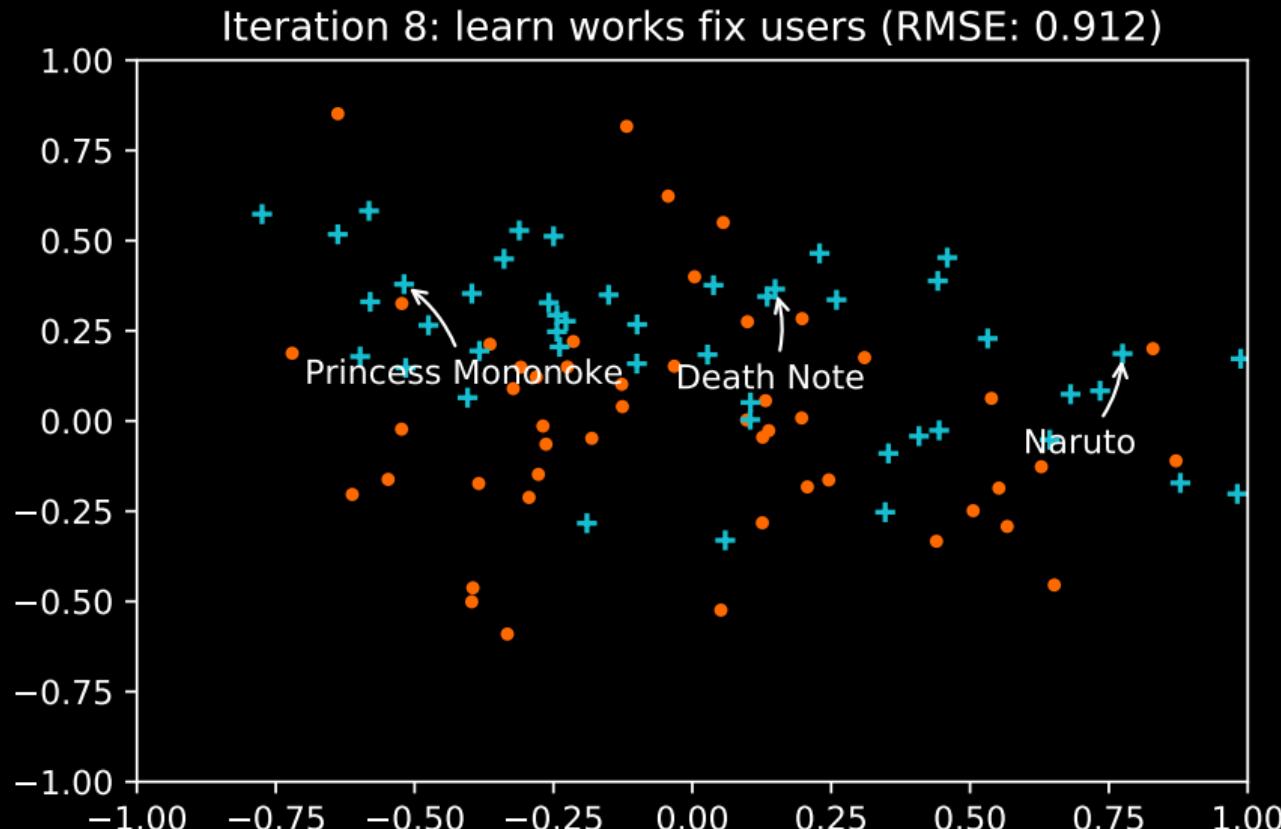


Illustration of Alternating Least Squares

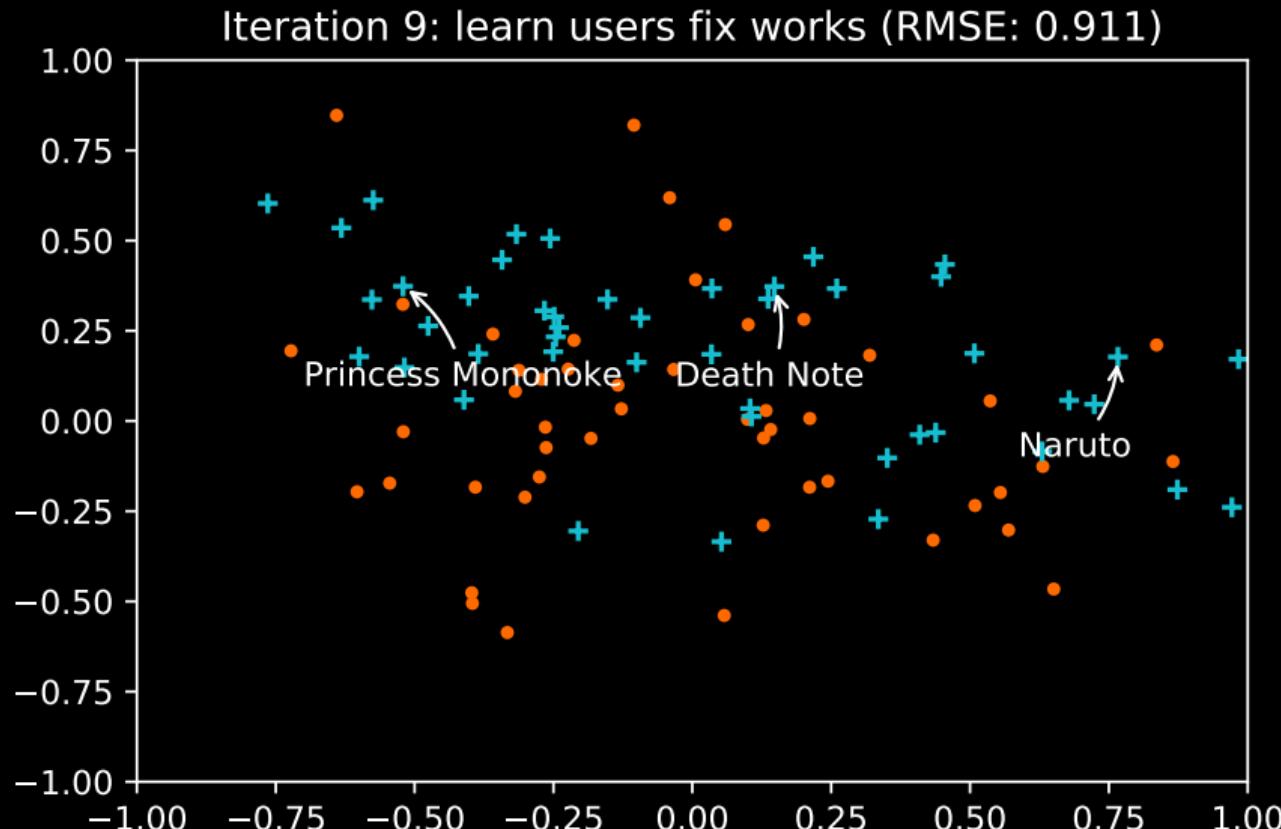


Illustration of Alternating Least Squares

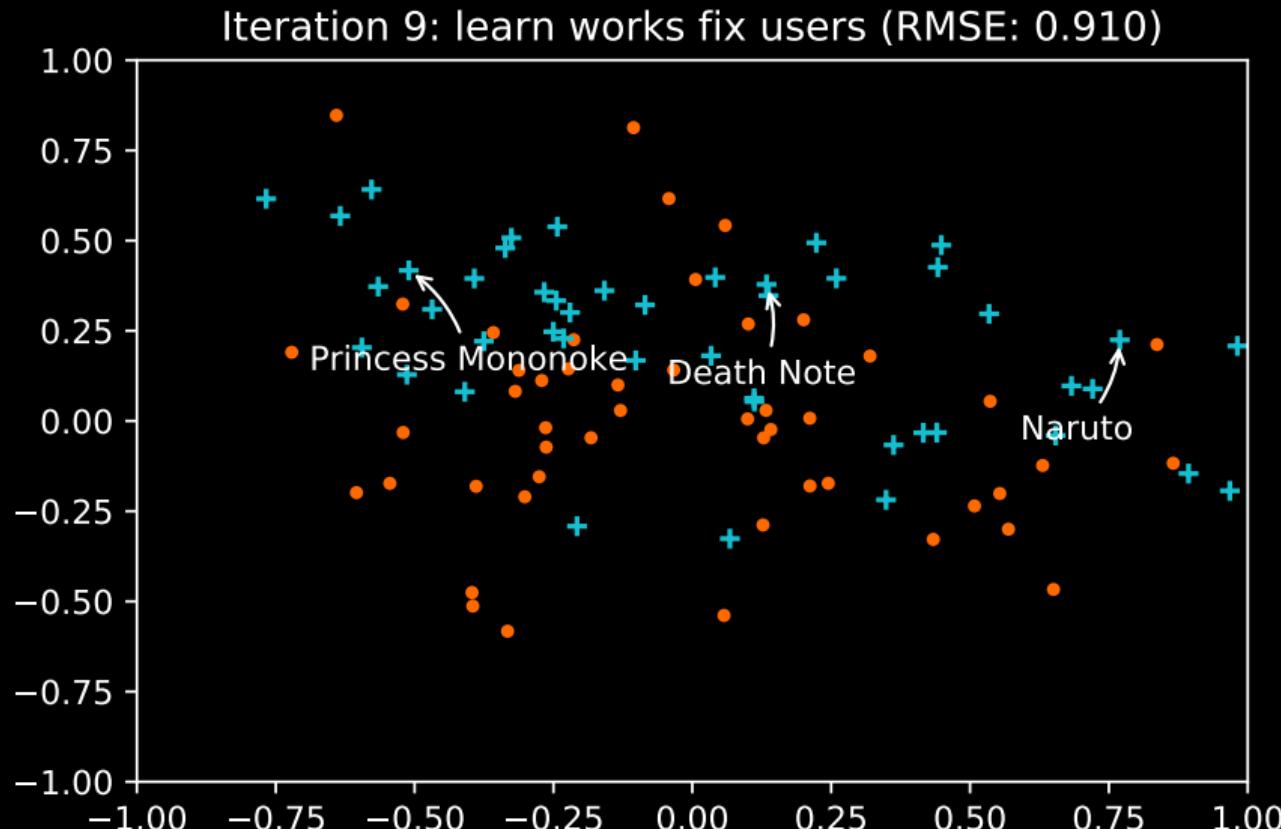


Illustration of Alternating Least Squares

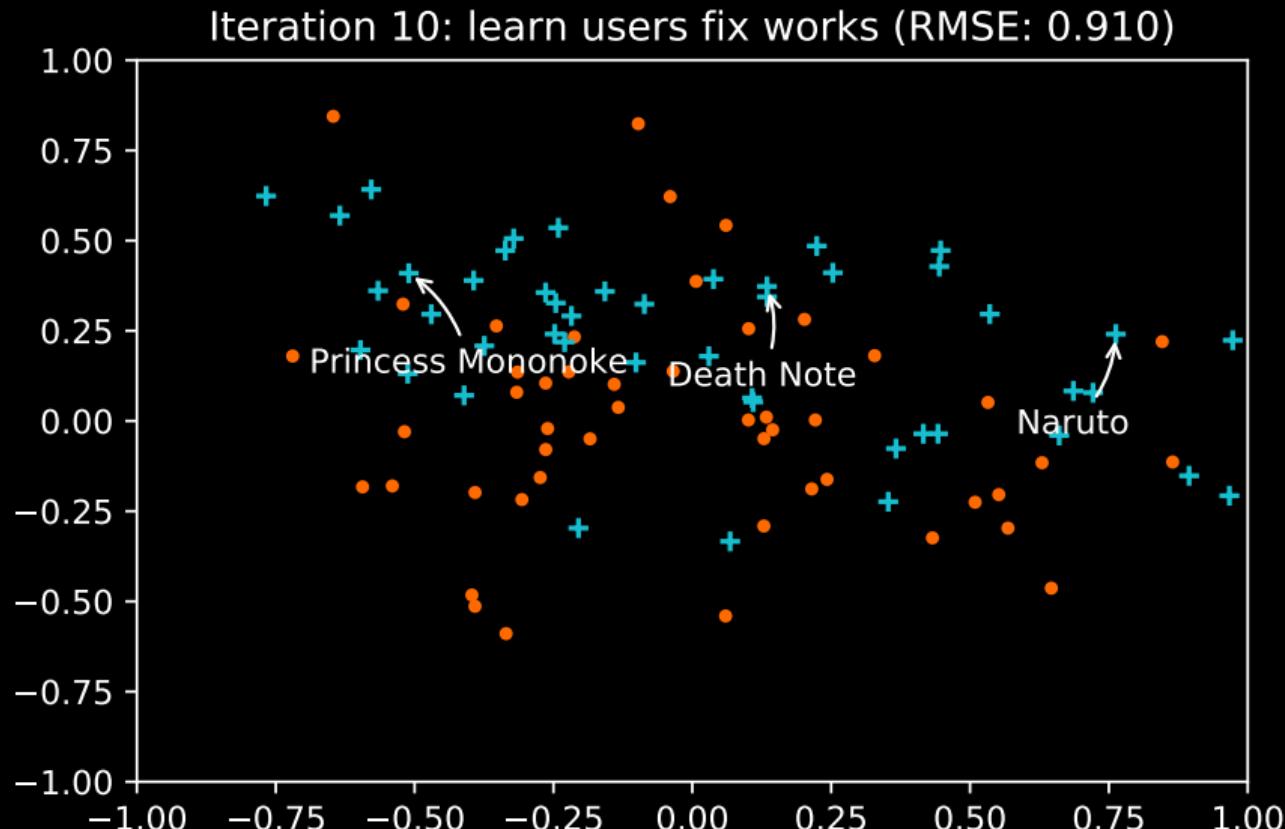


Illustration of Alternating Least Squares

Iteration 10: learn works fix users (RMSE: 0.909)

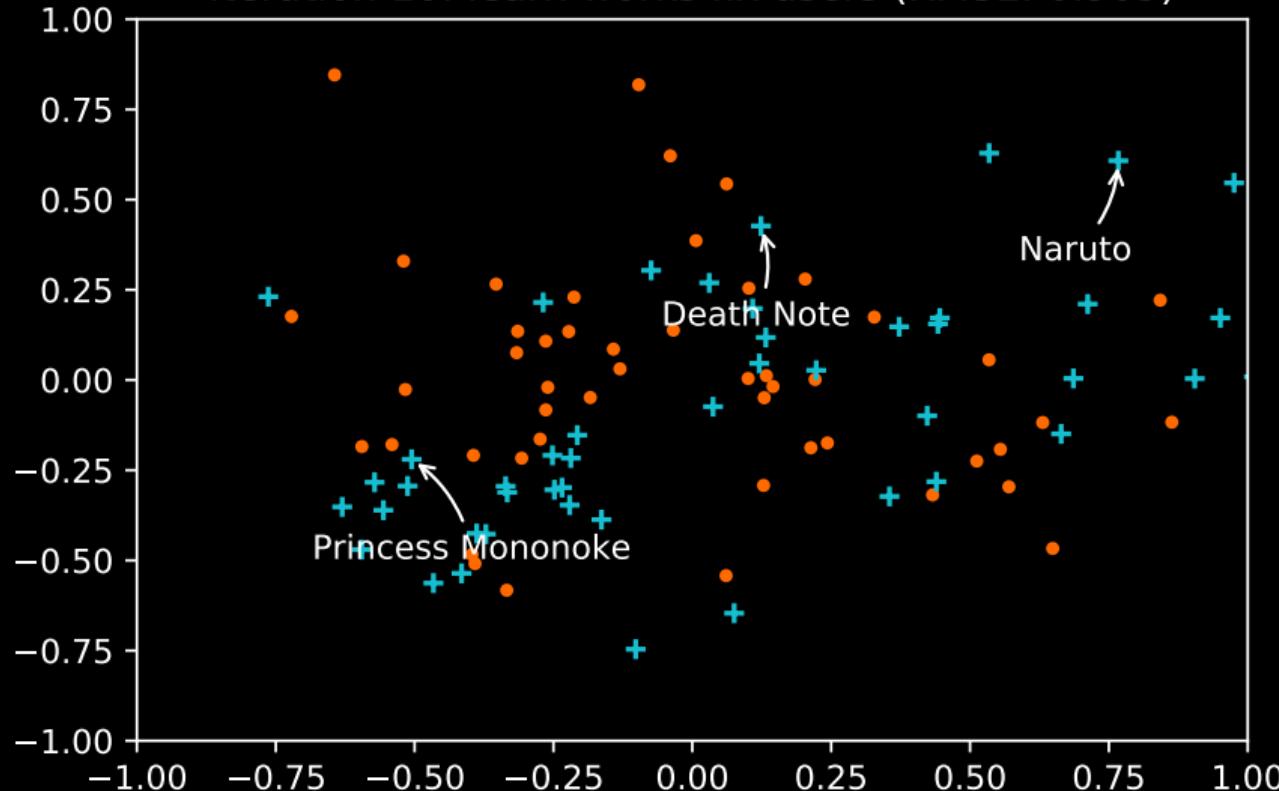


Illustration of Alternating Least Squares

Iteration 11: learn users fix works (RMSE: 0.909)

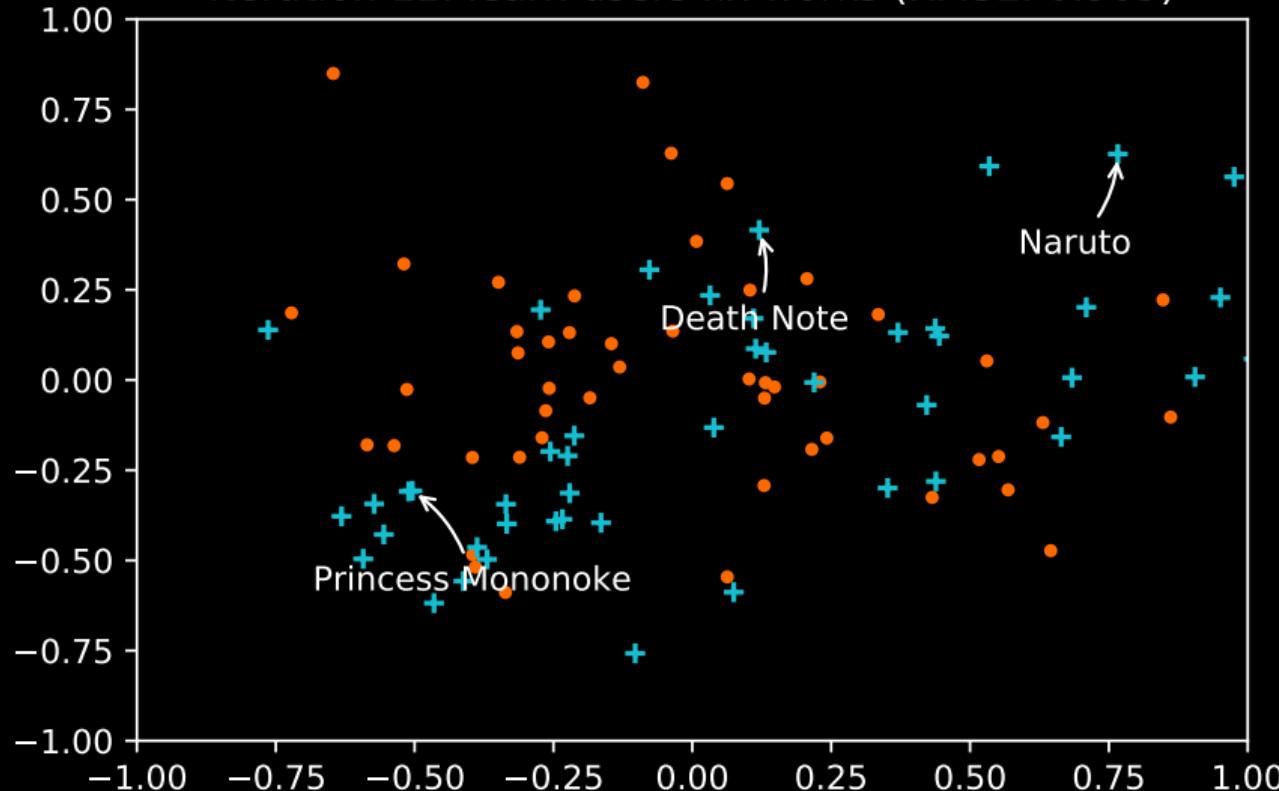


Illustration of Alternating Least Squares

Iteration 11: learn works fix users (RMSE: 0.908)

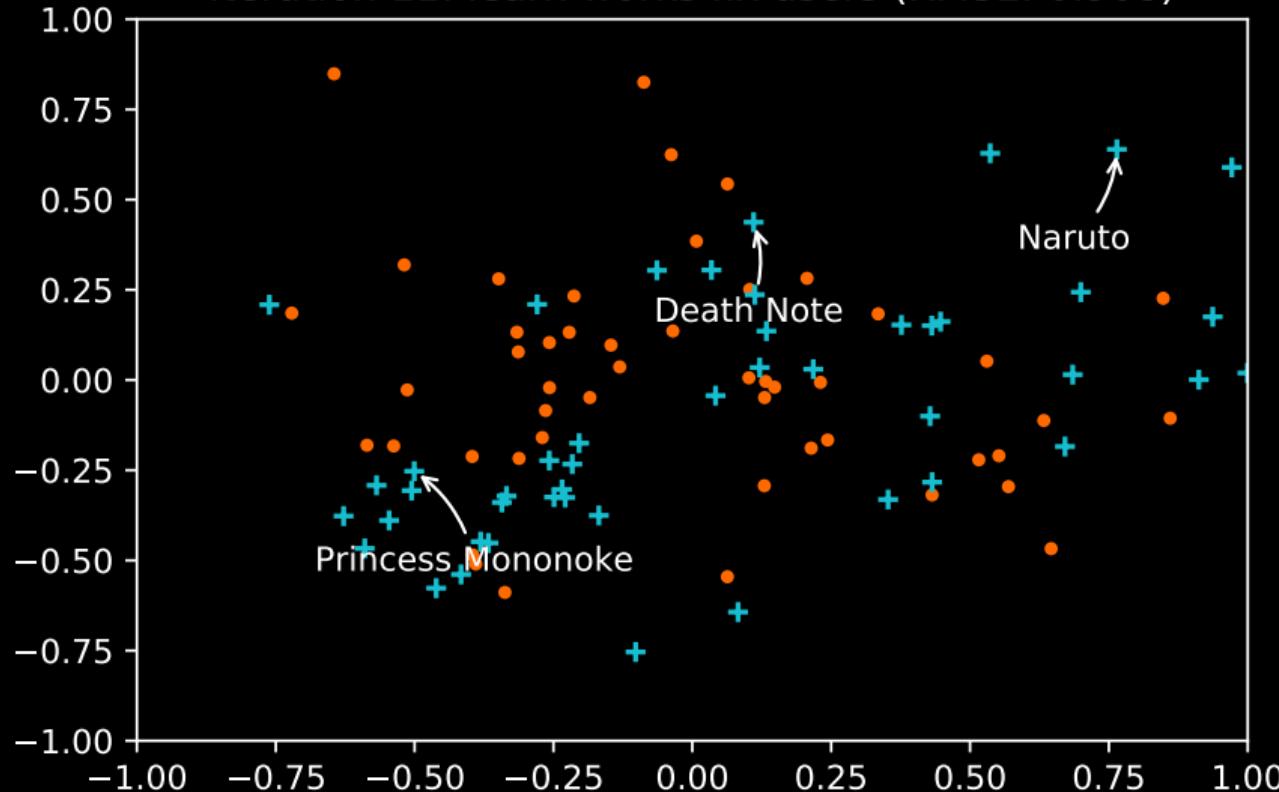


Illustration of Alternating Least Squares

Iteration 12: learn users fix works (RMSE: 0.908)

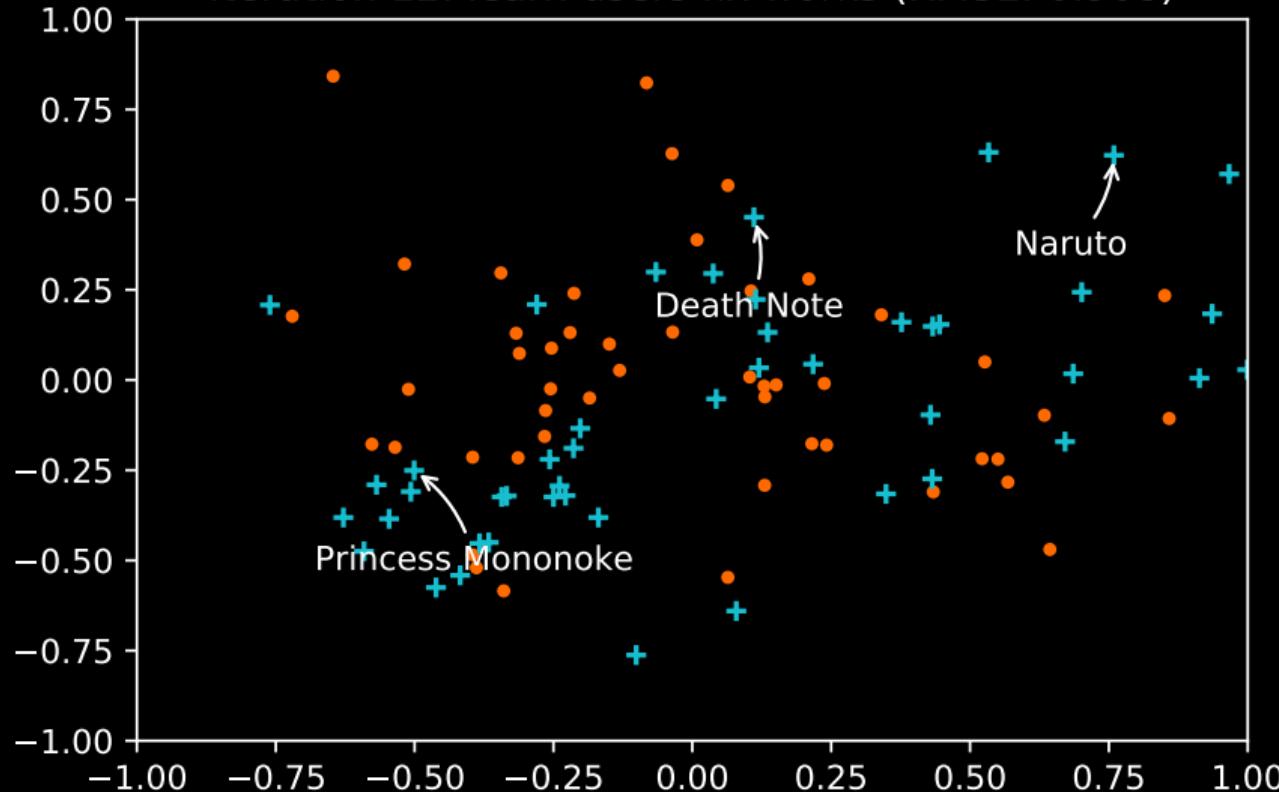


Illustration of Alternating Least Squares

Iteration 12: learn works fix users (RMSE: 0.908)

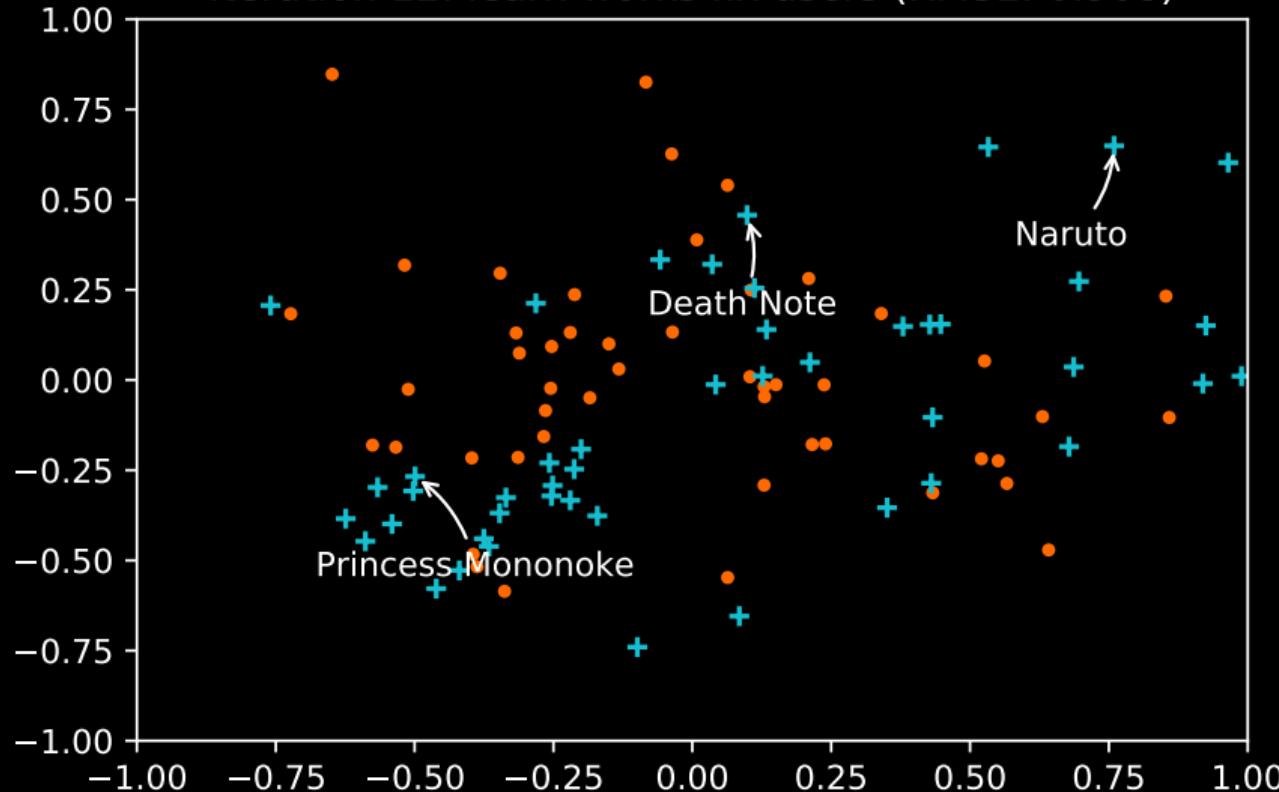


Illustration of Alternating Least Squares

Iteration 13: learn users fix works (RMSE: 0.907)

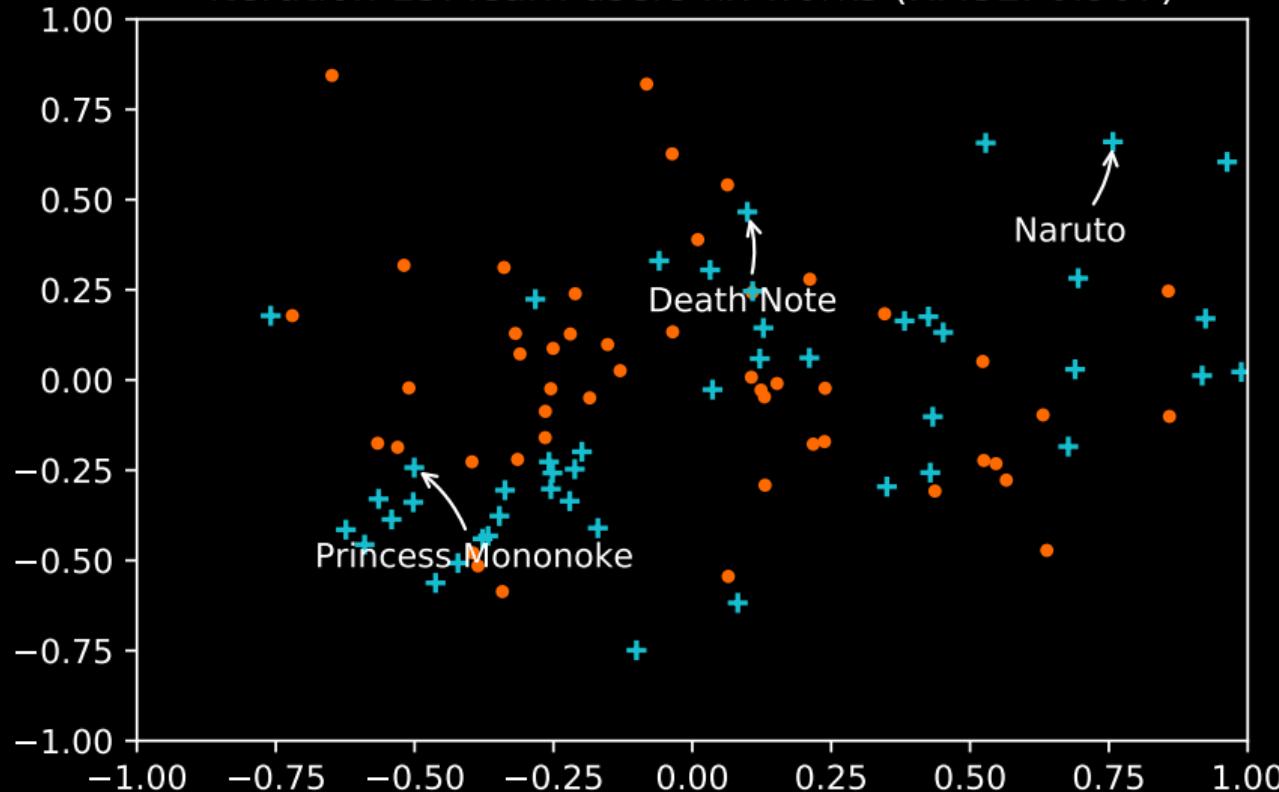


Illustration of Alternating Least Squares

Iteration 13: learn works fix users (RMSE: 0.907)

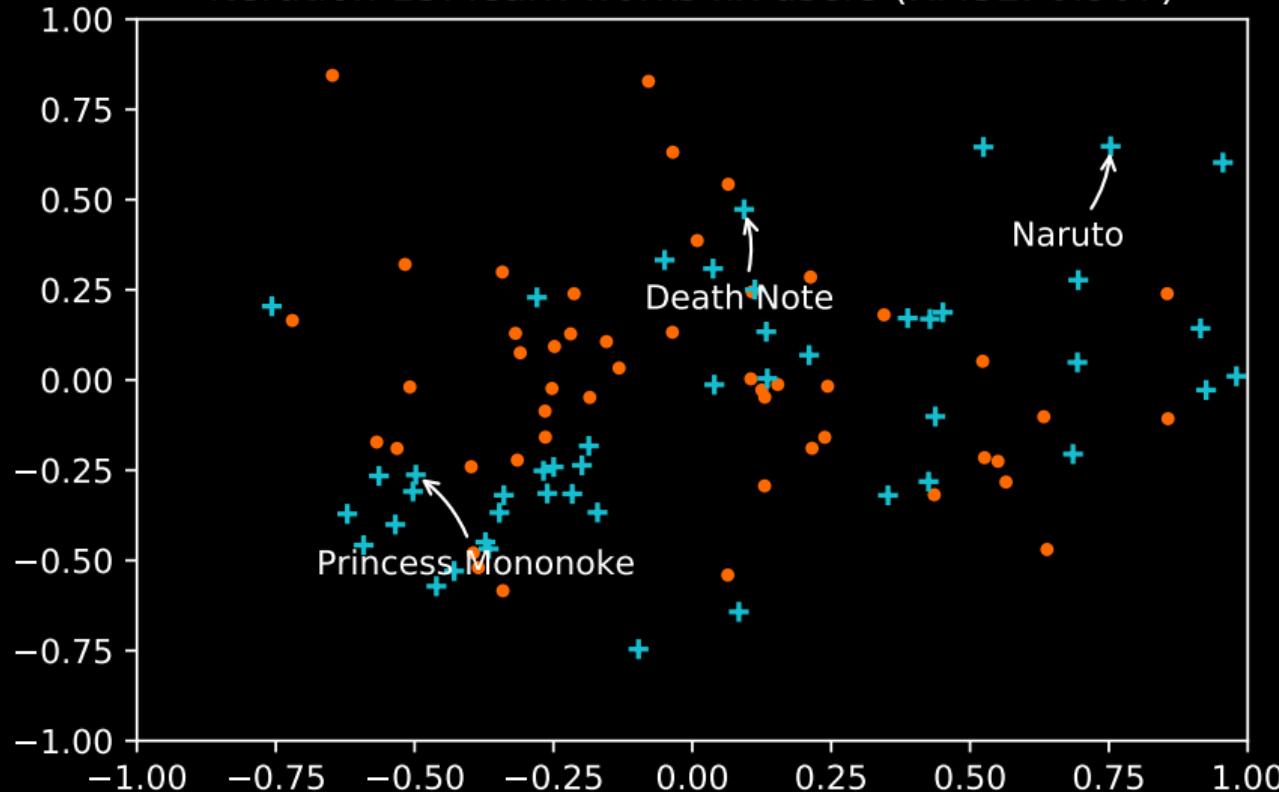


Illustration of Alternating Least Squares

Iteration 14: learn users fix works (RMSE: 0.907)

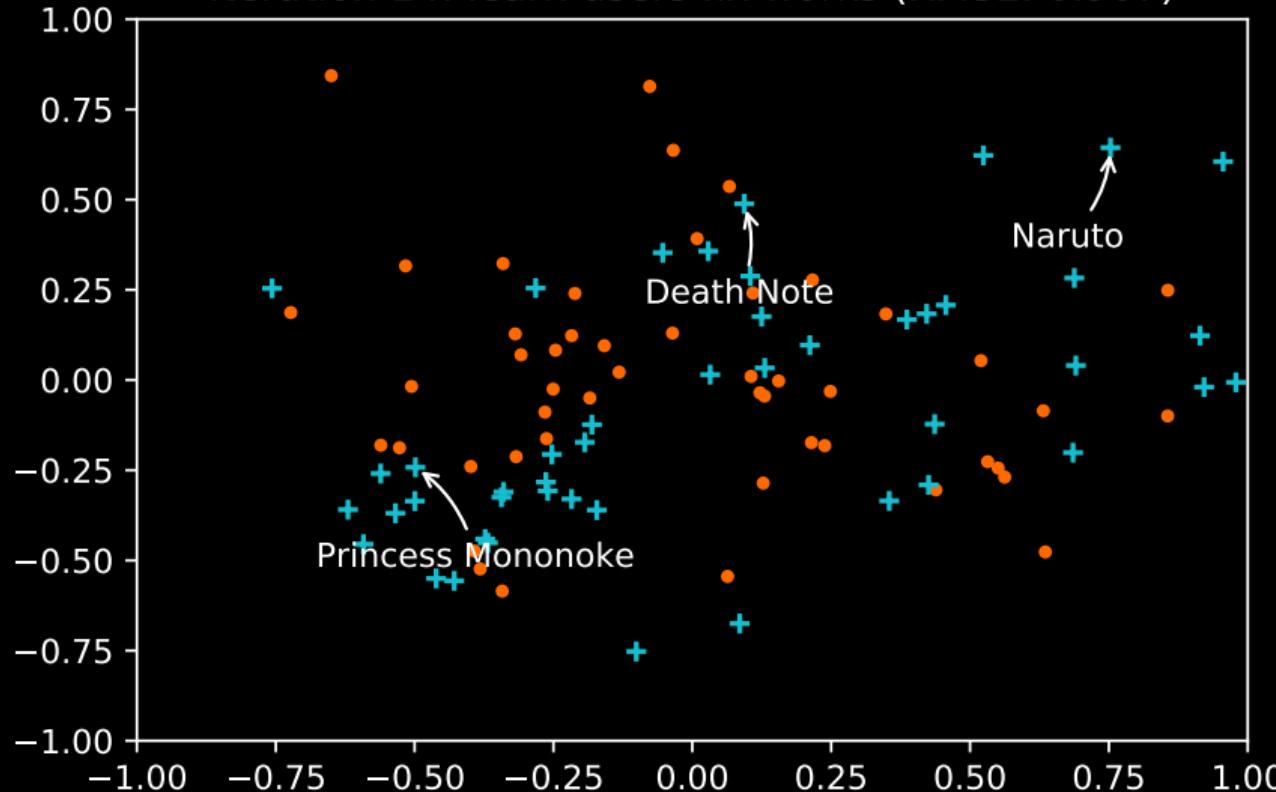


Illustration of Alternating Least Squares

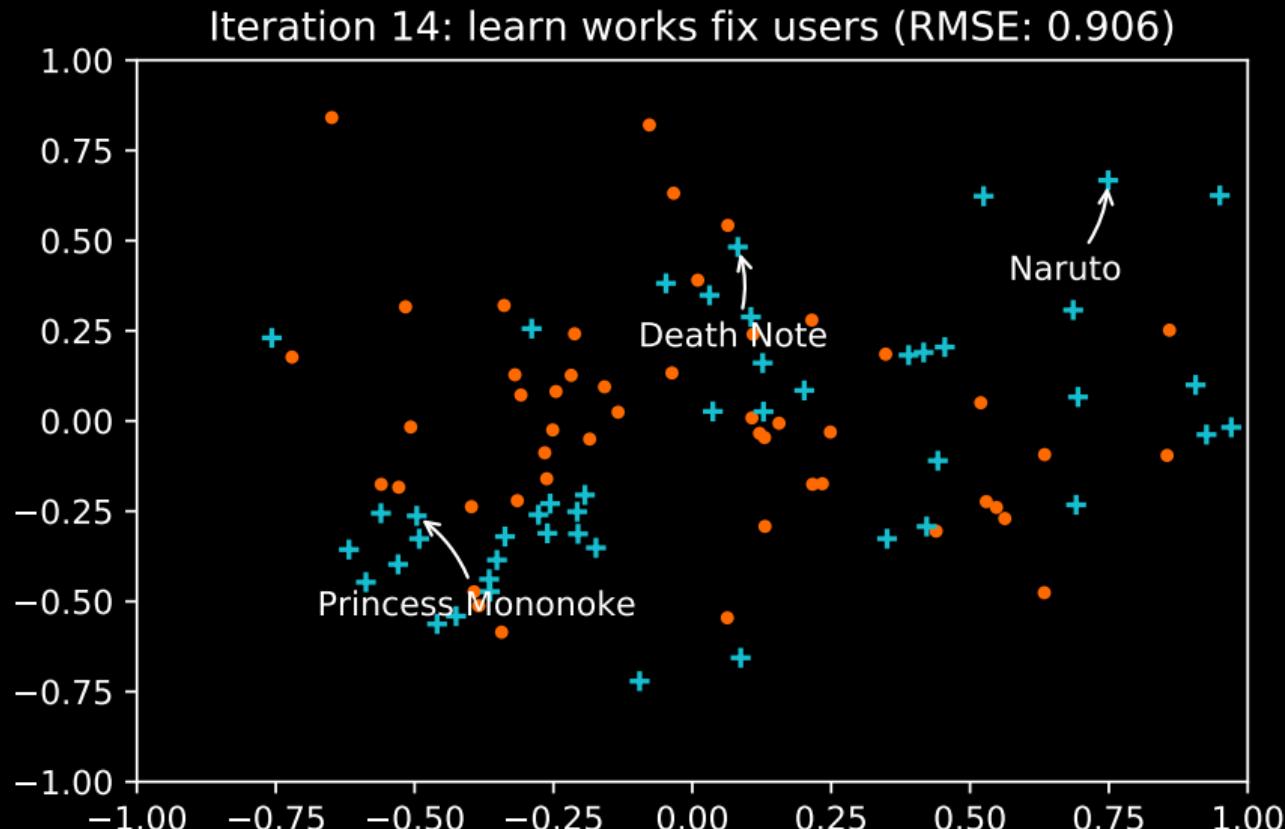


Illustration of Alternating Least Squares

Iteration 15: learn users fix works (RMSE: 0.906)

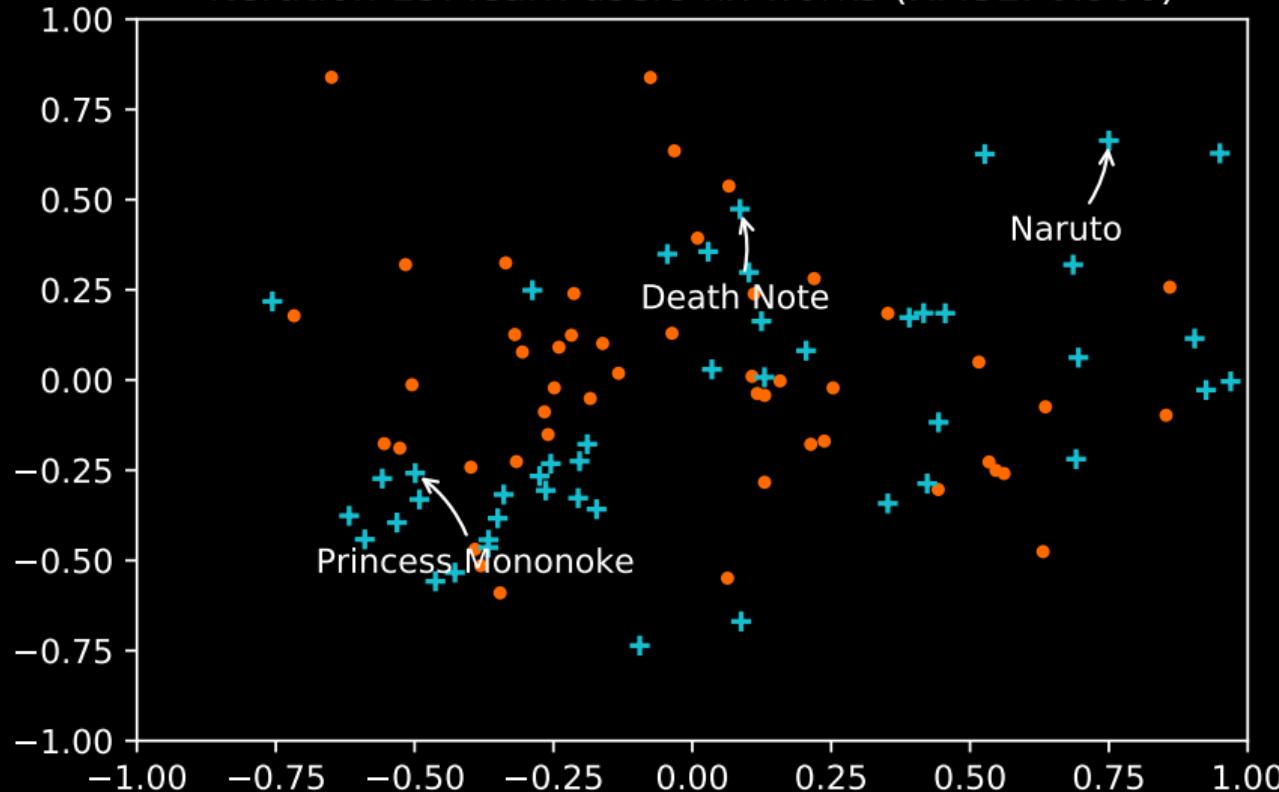


Illustration of Alternating Least Squares

Iteration 15: learn works fix users (RMSE: 0.906)

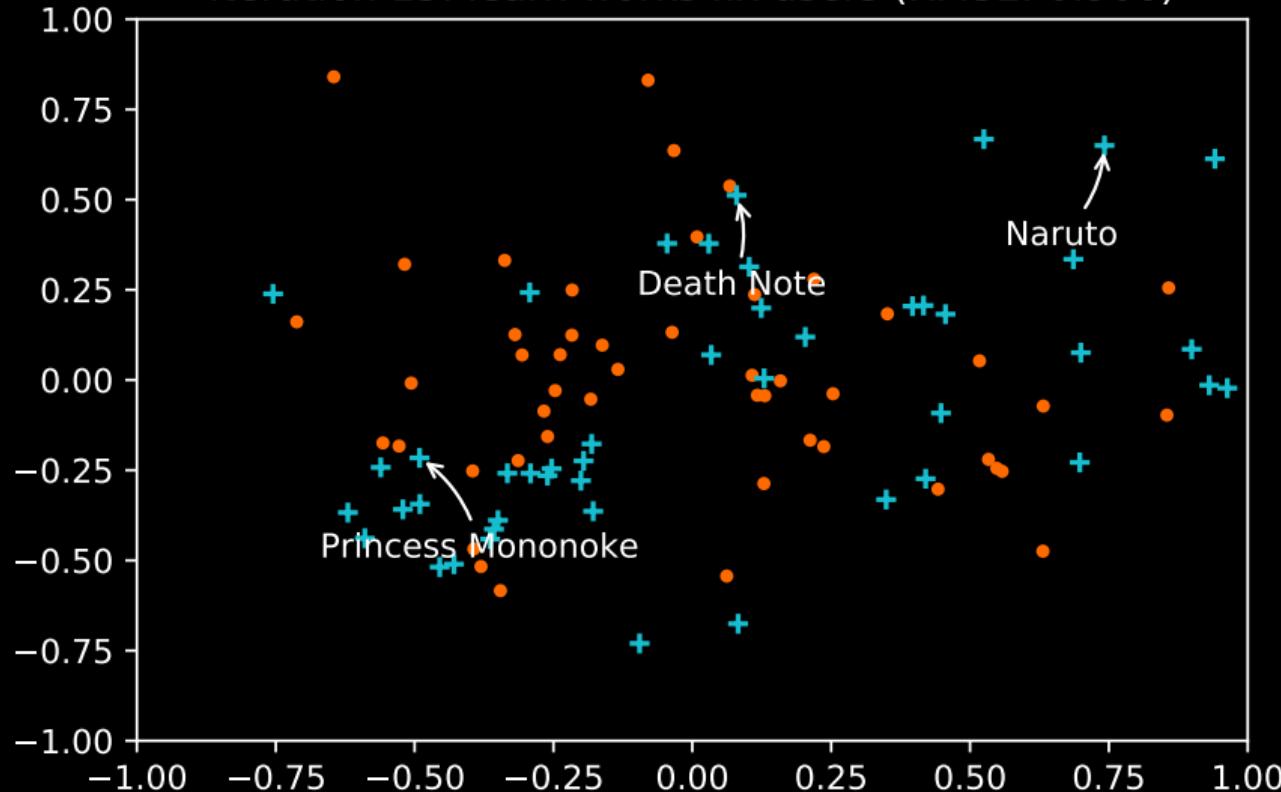


Illustration of Alternating Least Squares

Iteration 16: learn users fix works (RMSE: 0.906)

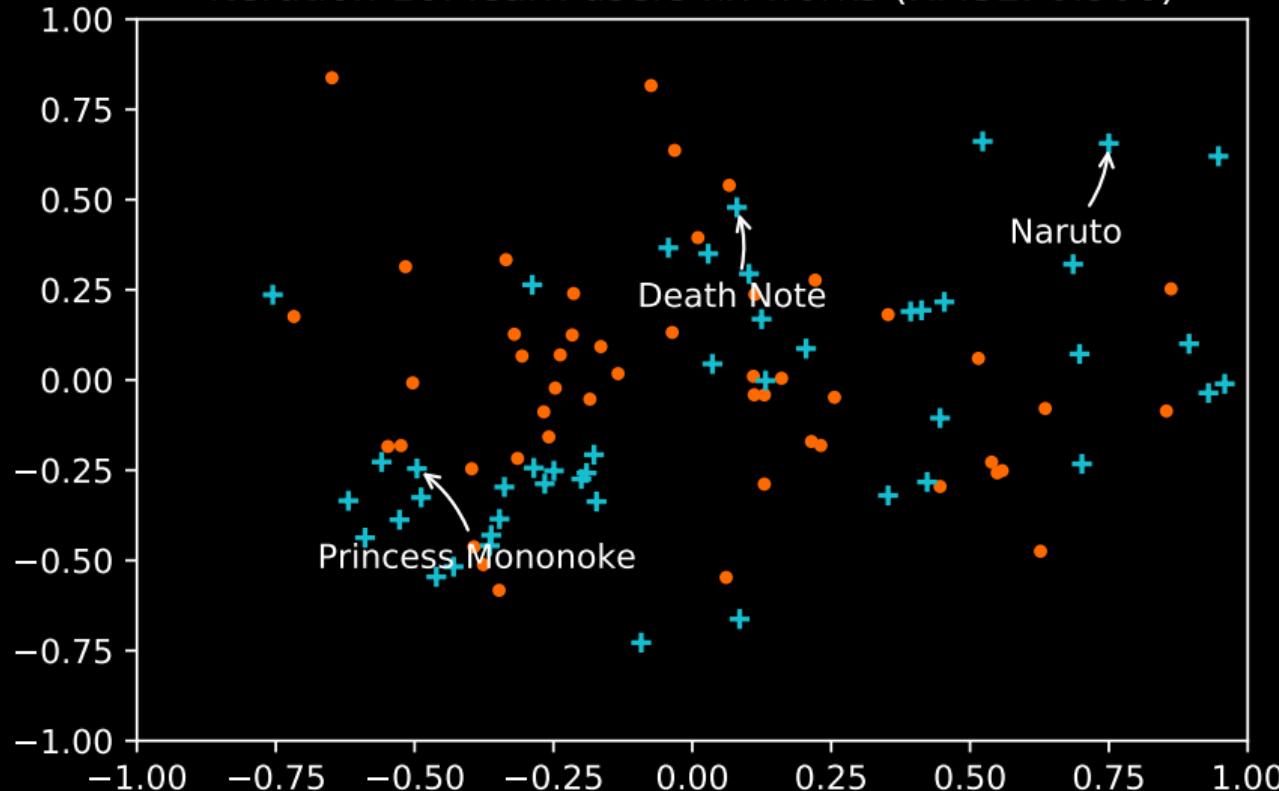


Illustration of Alternating Least Squares

Iteration 16: learn works fix users (RMSE: 0.906)

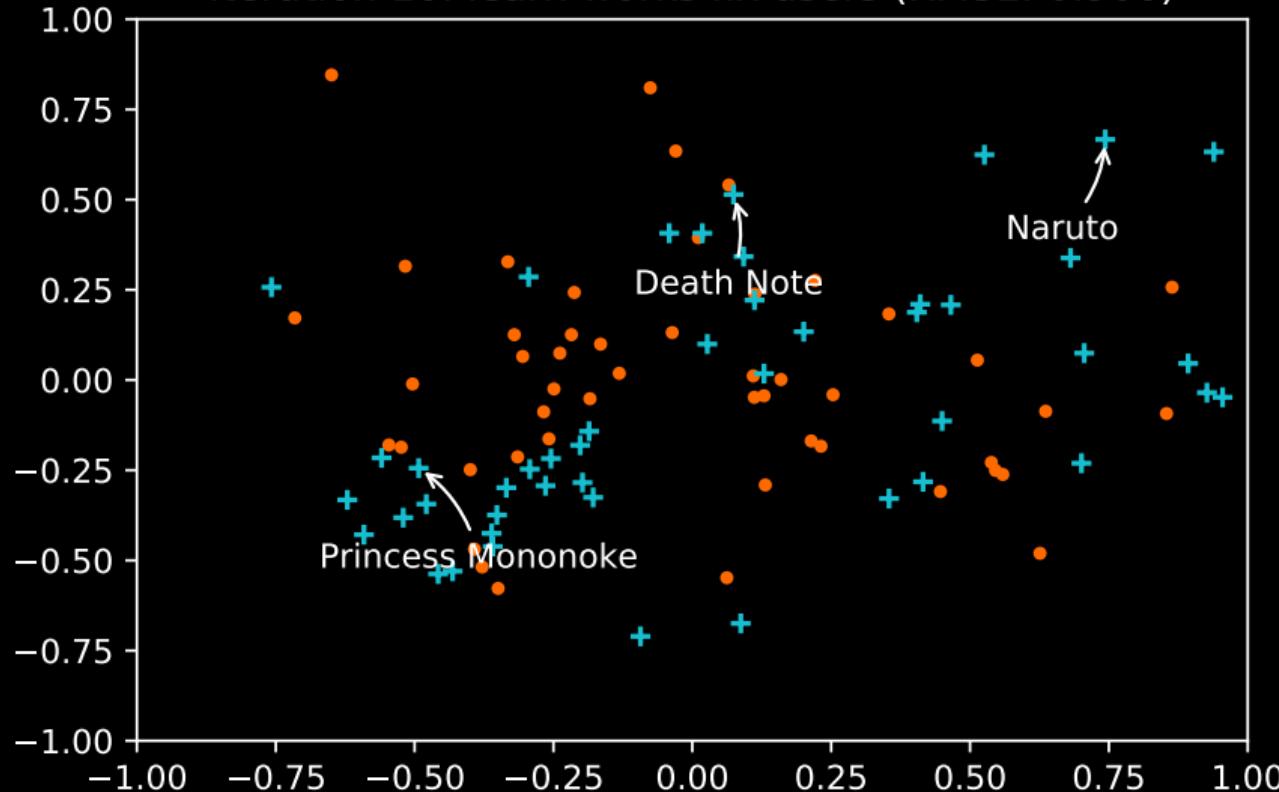


Illustration of Alternating Least Squares

Iteration 17: learn users fix works (RMSE: 0.905)

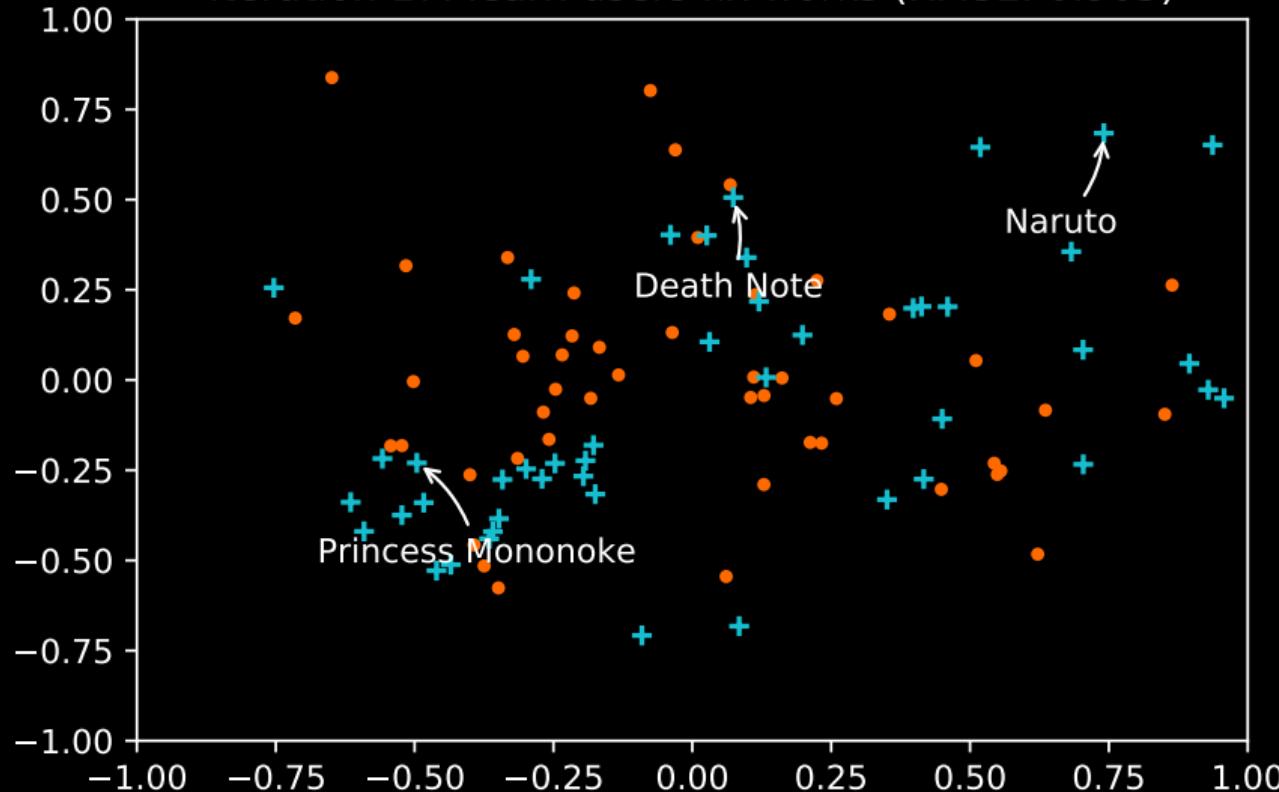


Illustration of Alternating Least Squares

Iteration 17: learn works fix users (RMSE: 0.905)

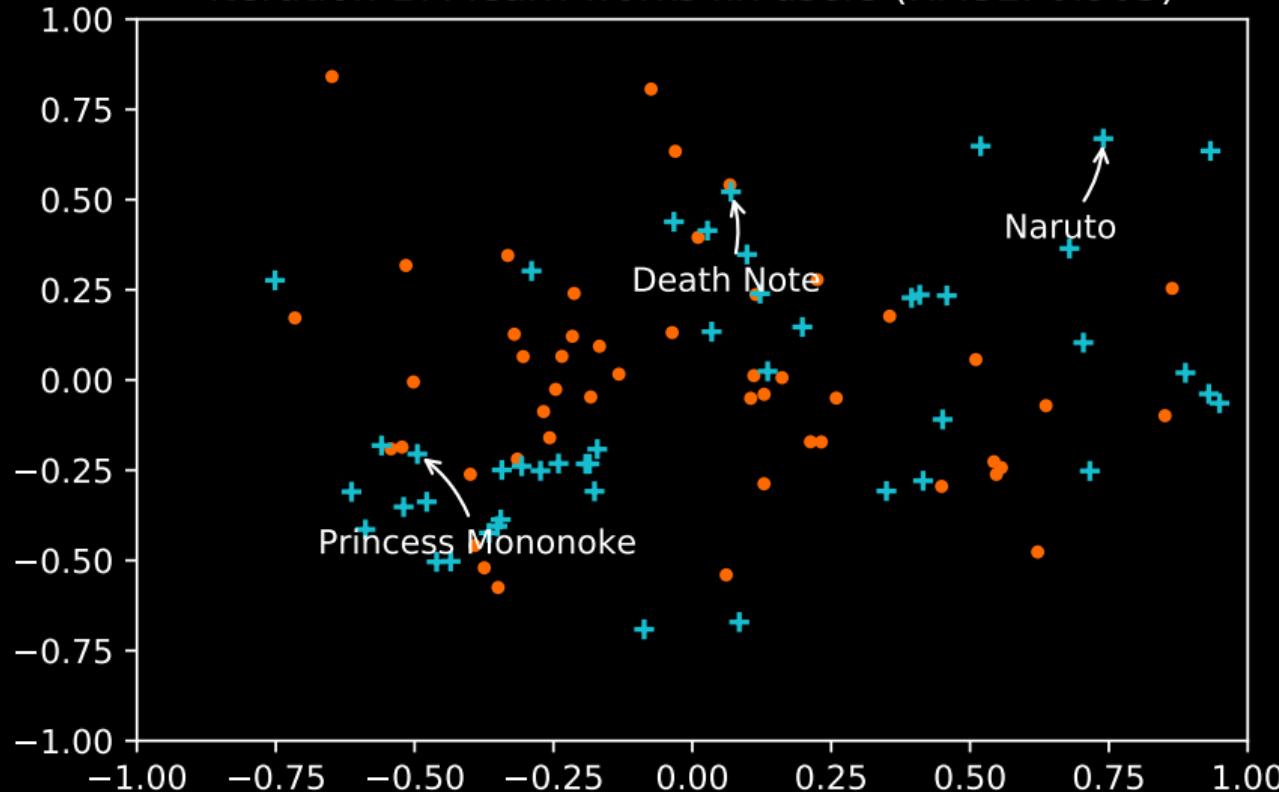


Illustration of Alternating Least Squares

Iteration 18: learn users fix works (RMSE: 0.905)

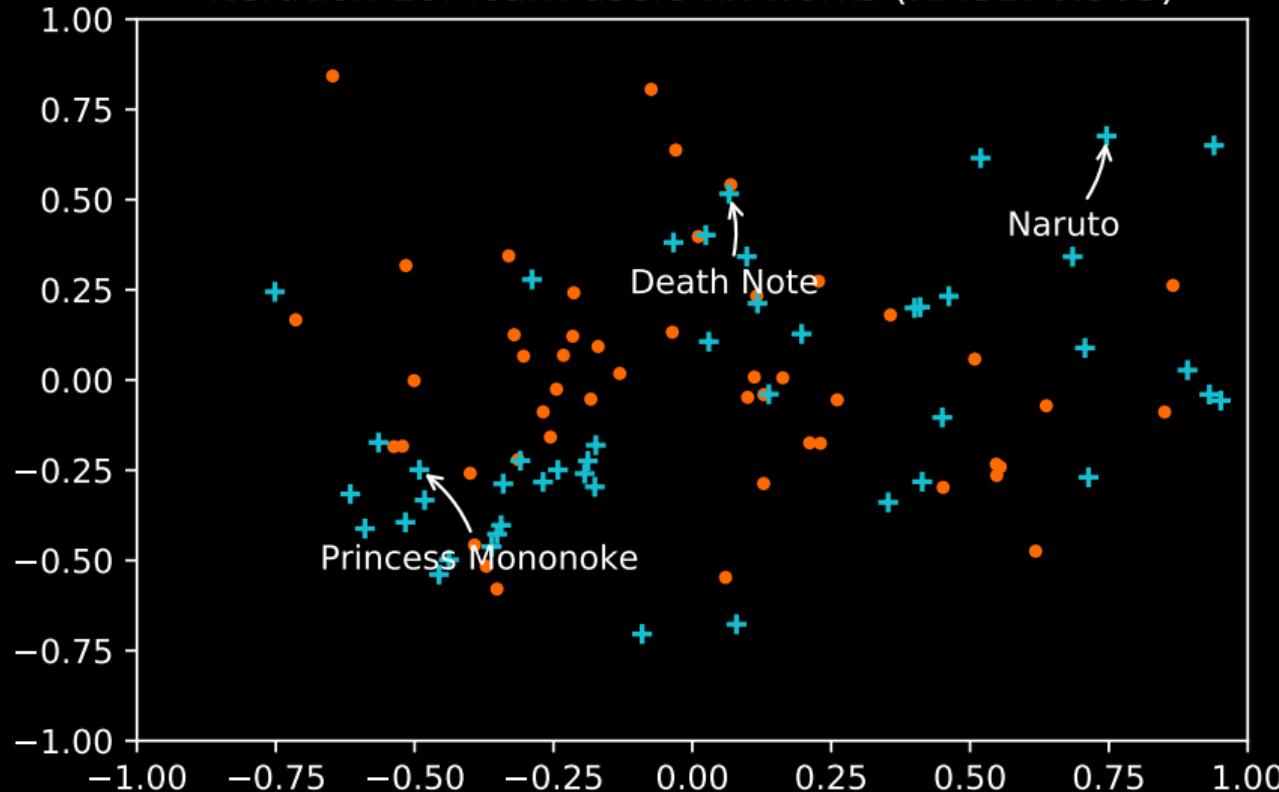


Illustration of Alternating Least Squares

Iteration 18: learn works fix users (RMSE: 0.905)

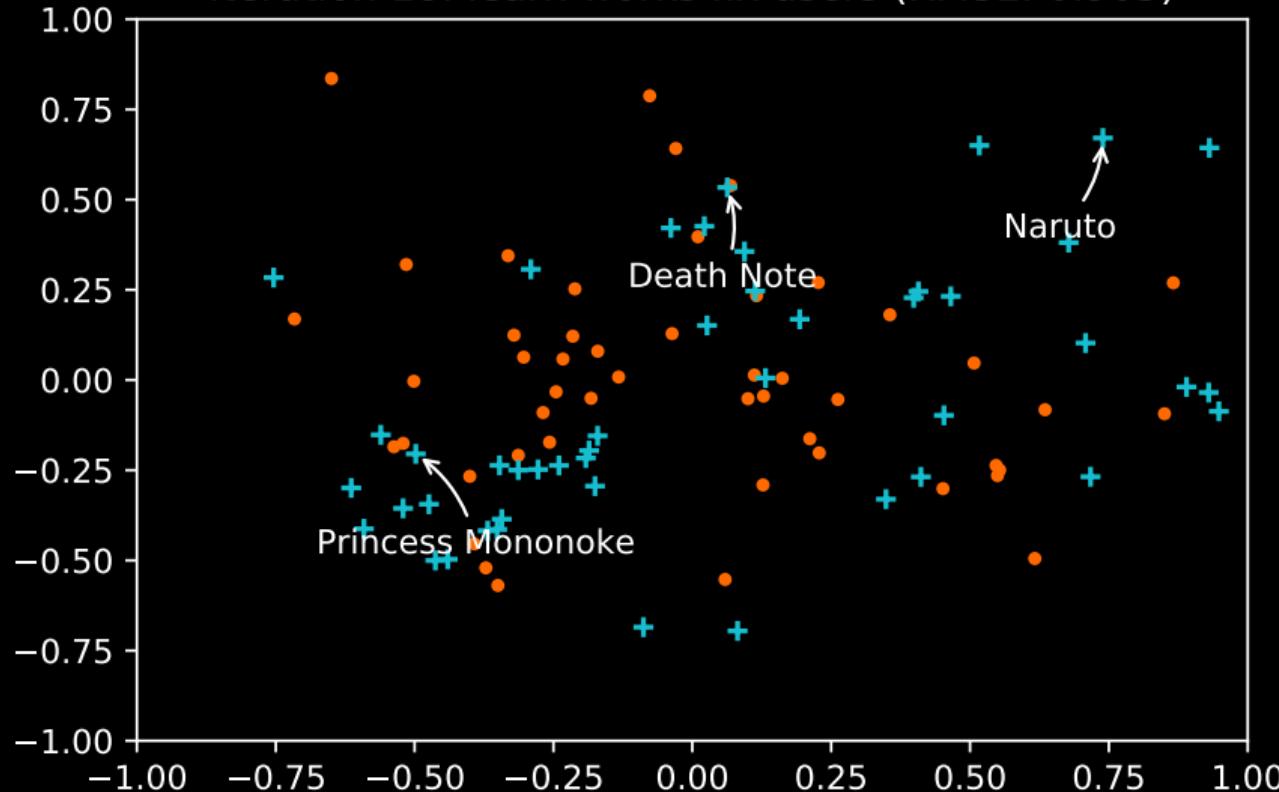


Illustration of Alternating Least Squares

Iteration 19: learn users fix works (RMSE: 0.905)

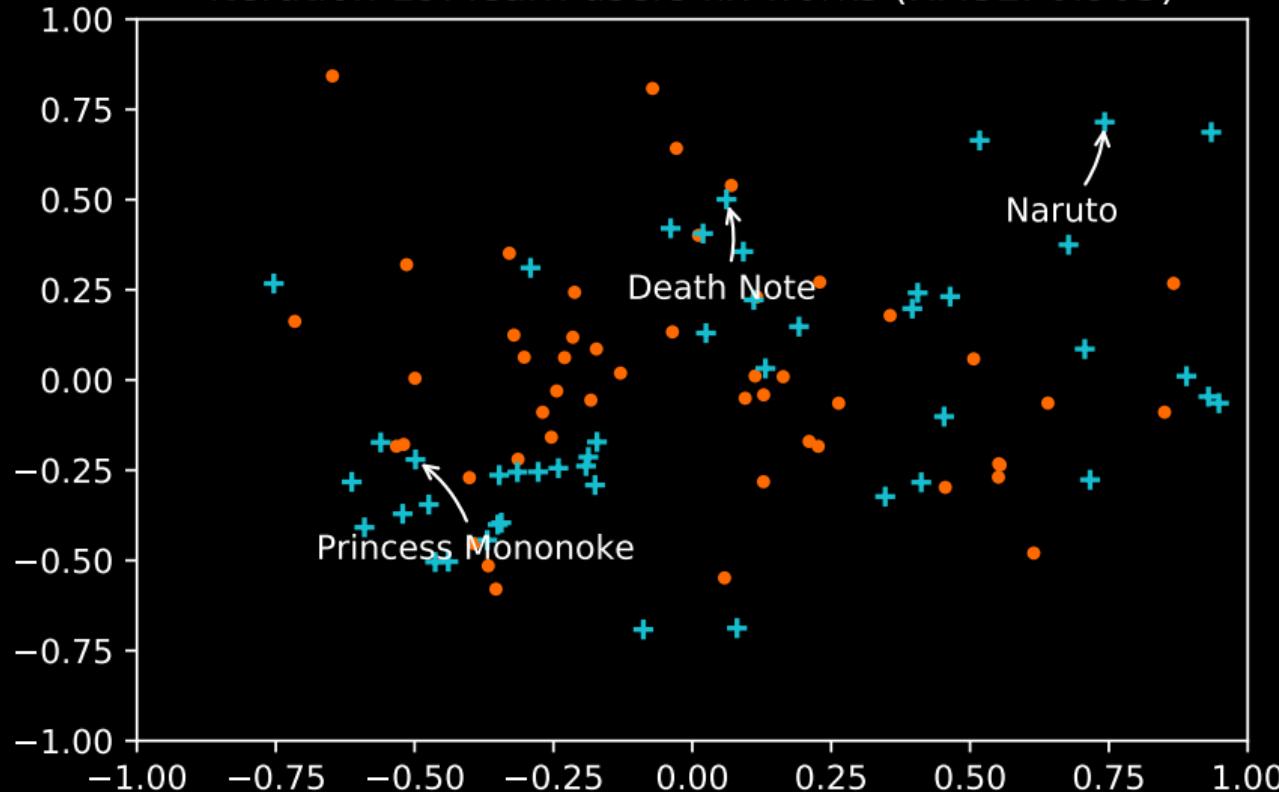
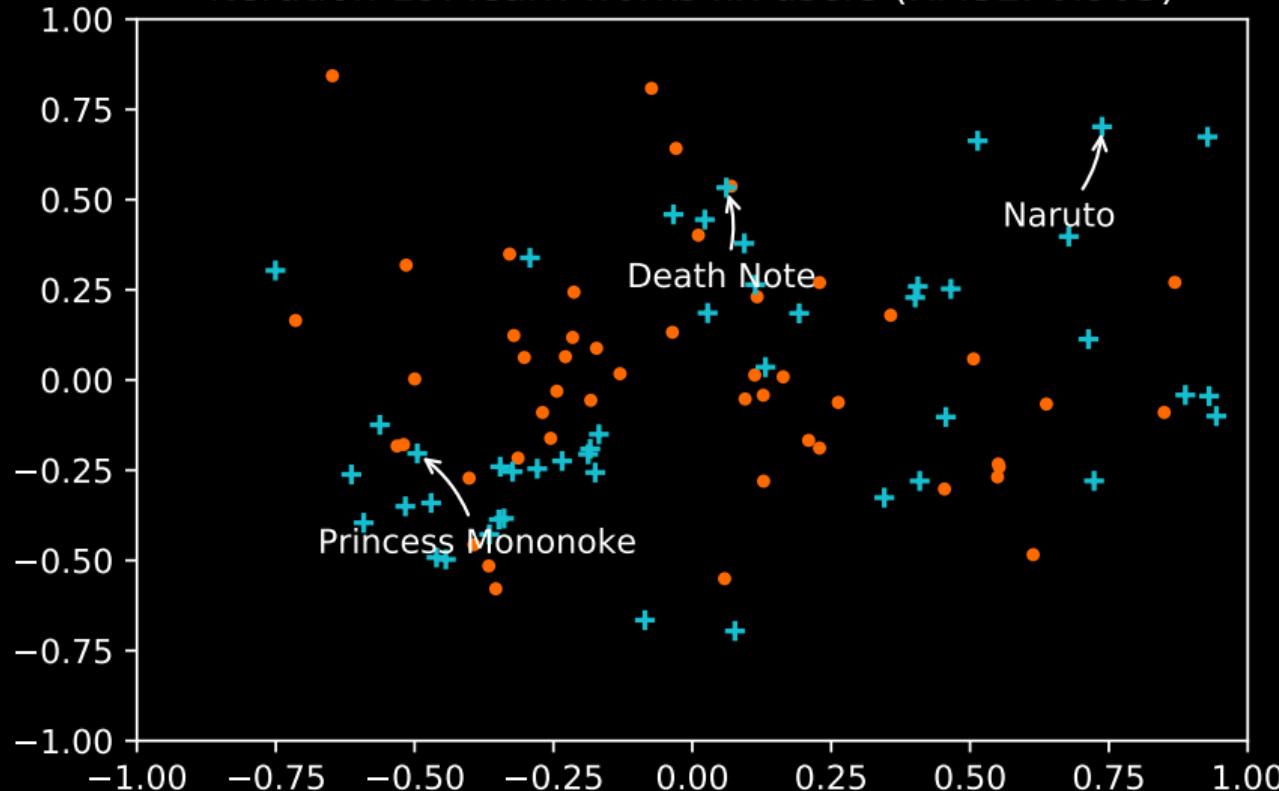
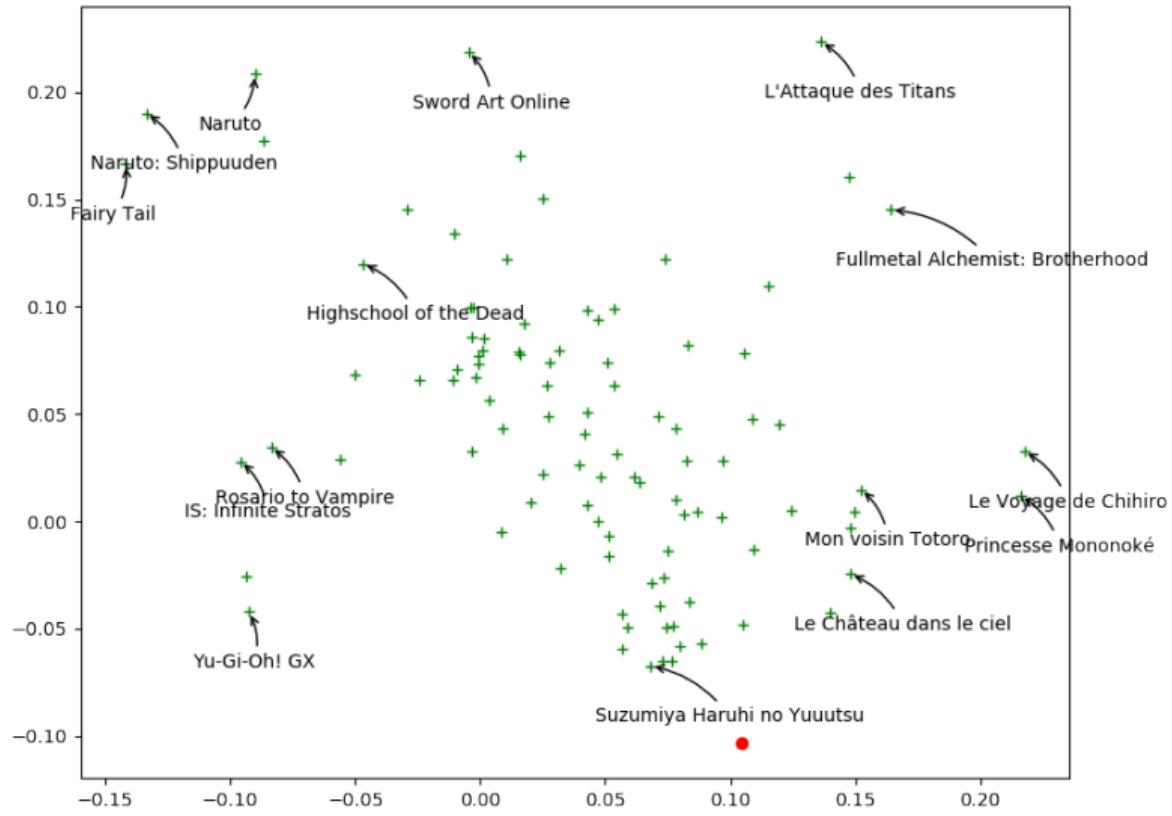


Illustration of Alternating Least Squares

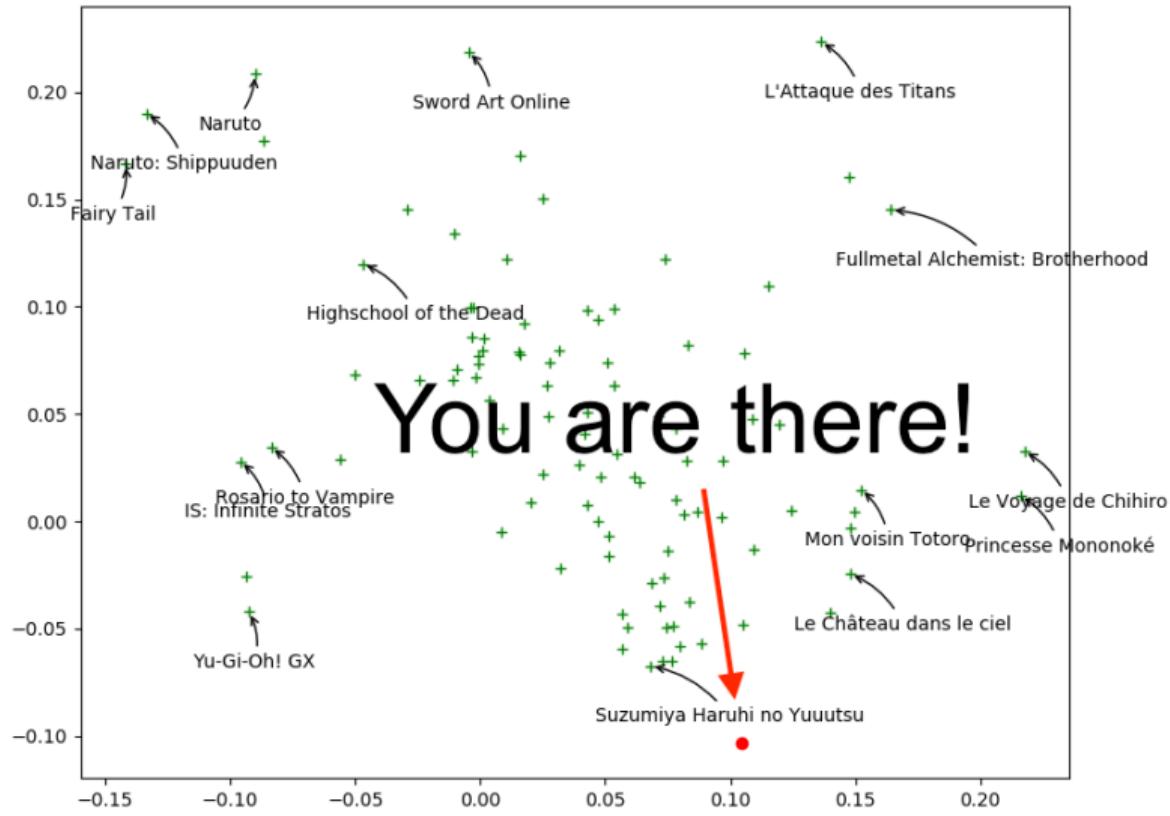
Iteration 19: learn works fix users (RMSE: 0.905)



Close v_j mean similar movies



u_i will like movies in its direction



word2vec: the Skip-Gram model

We want, given a word $w \sim P(W)$, to predict its **context** (i.e. the words around it).

Oh, it's a classification problem! With 60,000 classes.

$$P(c|w) = \text{softmax}(\mathbf{u}_w^T \mathbf{v}_c) \propto \exp(\mathbf{u}_w^T \mathbf{v}_c)$$

Cross-entropy would look like:

$$\mathcal{L} = \log \sigma(\mathbf{u}_w^T \mathbf{v}_c) + \sum_{\substack{i=1 \\ c_i \neq c}}^{60000} \log(1 - \sigma(\mathbf{u}_w^T \mathbf{v}_{c_i}))$$

So instead let's just sample k words from vocabulary proportionally to their occurrence:

$$\mathcal{L} = \sum_{w,c} \log \sigma(\mathbf{u}_w^T \mathbf{v}_c) + \sum_{\substack{i=1 \\ c_i \sim P(W)}}^k \log(1 - \sigma(\mathbf{u}_w^T \mathbf{v}_{c_i}))$$

Neural Word Embedding as Implicit Matrix Factorization

Noise contrastive estimation

$$\mathcal{L} = \sum_{w,c} \log \sigma(\mathbf{u}_w^T \mathbf{v}_c) + \sum_{\substack{i=1 \\ c_i \sim P(W)}}^k \log(1 - \sigma(\mathbf{u}_w^T \mathbf{v}_{c_i}))$$

Negative sampling (similar objective)

$$\mathcal{L} = \sum_{w,c} \log \sigma(\mathbf{u}_w^T \mathbf{v}_c) + \sum_{\substack{i=1 \\ c_i \sim P(W)}}^k \log \sigma(-\mathbf{u}_w^T \mathbf{v}_{c_i}))$$

This objective is equivalent to factorizing the pointwise mutual information matrix:

$$\log \frac{\#(w, c)N}{\#w\#c} - \log k = \mathbf{u}_w^T \mathbf{v}_c \text{ where } N \text{ is the number of pairs.}$$

Omer Levy and Yoav Goldberg. "Neural Word Embedding as Implicit Matrix Factorization". In: *Advances in Neural Information Processing Systems*. Ed. by Z. Ghahramani et al. Vol. 27. Curran Associates, Inc., 2014. URL:

<https://proceedings.neurips.cc/paper/2014/file/feab05aa91085b7a8012516bc3533958-Paper.pdf>

Practical datasets

A toy dataset of matrix completion

You are given some entries of $M = UV$ where the shapes are $(10, 3) \cdot (3, 5)$
The goal is to recover the other entries

Movielens

100,000 ratings (1–5 stars) of 600 users on 9,000 movies.

<https://files.grouplens.org/datasets/movielens/ml-latest-small.zip>

The Adventures of Sherlock Holmes

By Arthur Conan Doyle

<https://github.com/theeluwin/pytorch-sgns>

- [1] Guillaume Lample et al. "Word translation without parallel data". In: *International Conference on Learning Representations*. 2018. URL: <https://openreview.net/forum?id=H196sainb>.
- [2] Omer Levy and Yoav Goldberg. "Neural Word Embedding as Implicit Matrix Factorization". In: *Advances in Neural Information Processing Systems*. Ed. by Z. Ghahramani et al. Vol. 27. Curran Associates, Inc., 2014. URL: <https://proceedings.neurips.cc/paper/2014/file/feab05aa91085b7a8012516bc3533958-Paper.pdf>.
- [3] Tomas Mikolov et al. "Distributed Representations of Words and Phrases and their Compositionality". In: *Advances in Neural Information Processing Systems*. Ed. by C.J. Burges et al. Vol. 26. Curran Associates, Inc., 2013. URL: <https://proceedings.neurips.cc/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf>.