

Expected Tensor Decomposition with Stochastic Gradient Descent

Takanori Maehara (Shizuoka Univ)

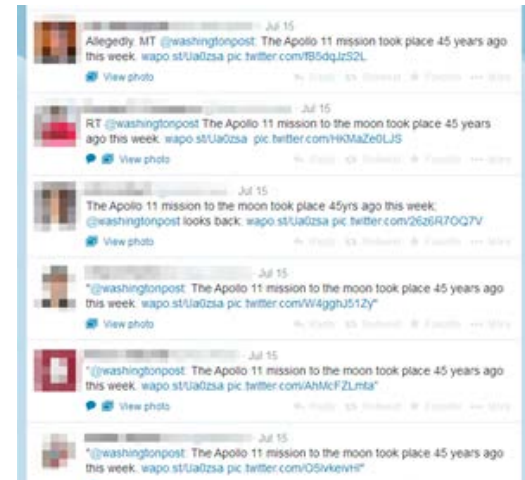
Kohei Hayashi (NII -> AIST)

Ken-ichi Kawarabayashi (NII)

Presented at AAAI 2016

Information Overload

- Too many texts to catch up



- How can we know what's going on?
--- Machine learning!



Kohei Hayashi
@hayasick

Follow

Pokemon GO is awesome! #PokemonGO
<https://www.instagram.com/..>

RETWEETS
304

FAVORITES
226



Reply Retweet Favorite More

06:08 PM - 07 Aug 2016

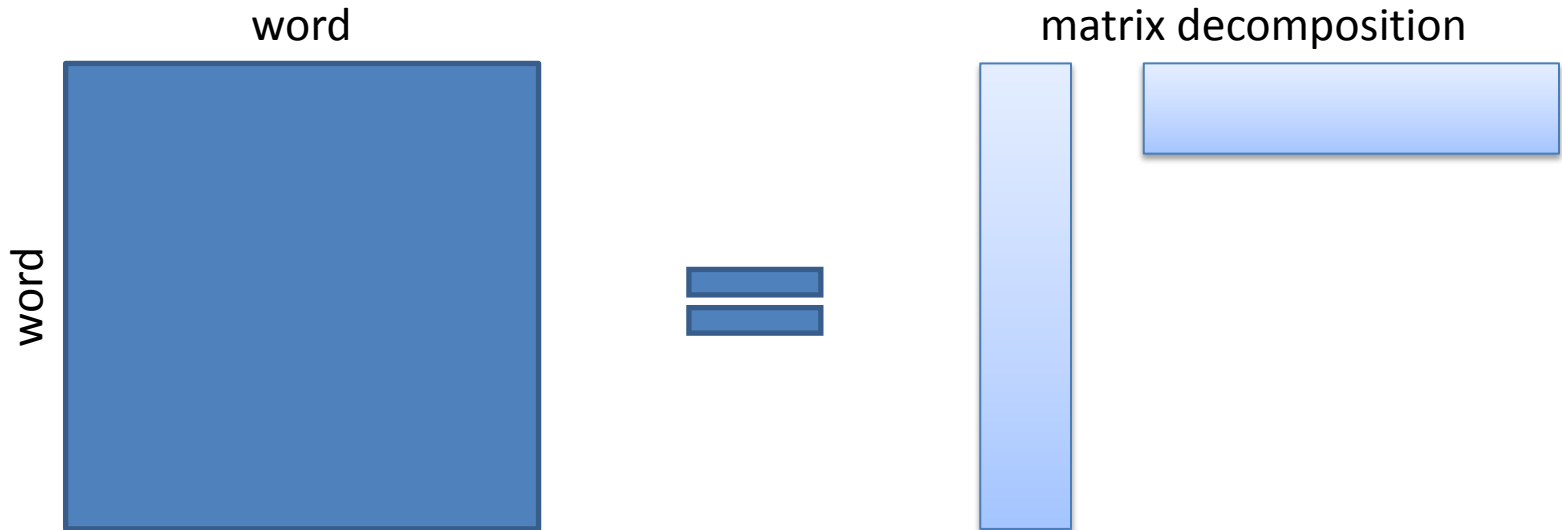
<http://tweetfake.com/>

How can we analyze this?

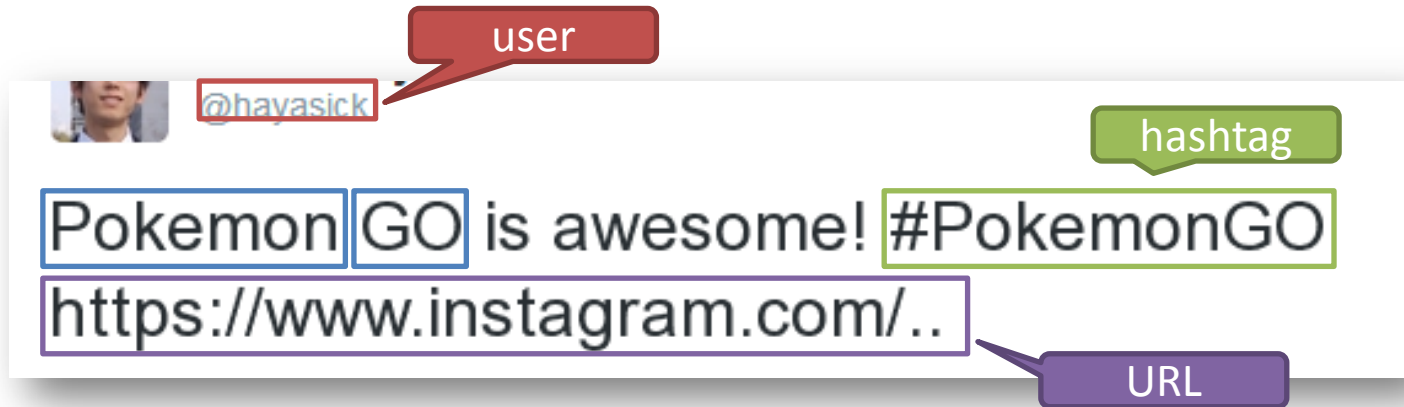
Word Co-occurrence

Pokemon GO is awesome! #PokemonGO
https://www.instagram.com/..

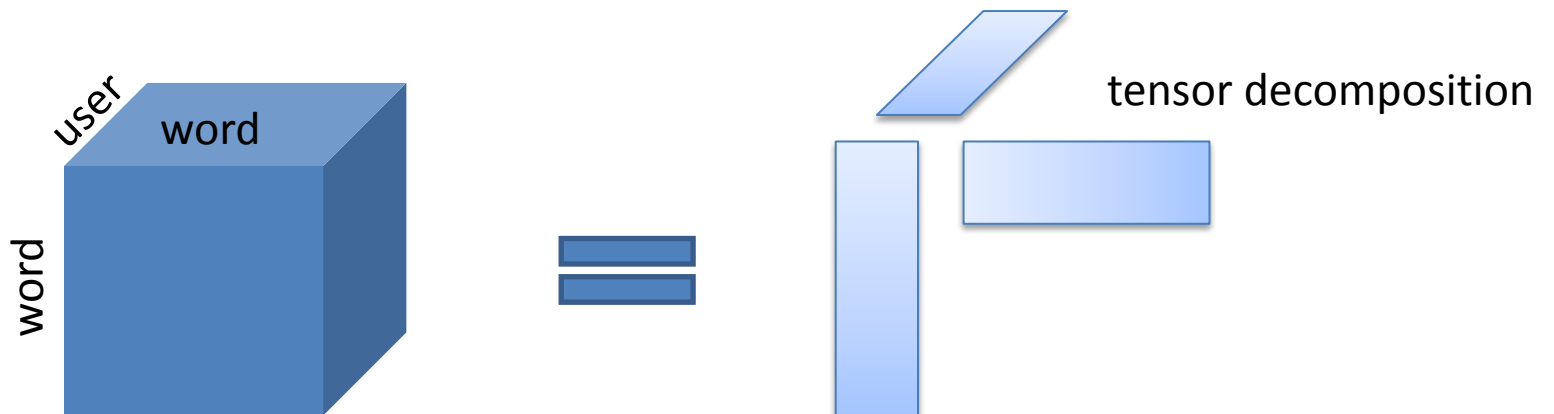
⋮



Higher-order Information

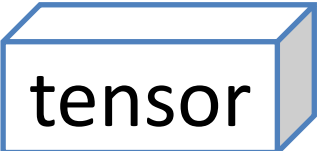


- Idea is generalized as **tensor**



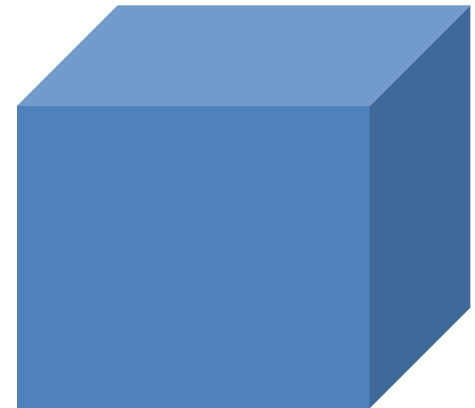
Scalability Issue

- To compute decomposition, we need  tensor beforehand.

- Sometimes  tensor is too huge to put on memory.

Pokemon GO is awesome! #PokemonGO
[https://www.instagram.com/..](https://www.instagram.com/)

...



Questions

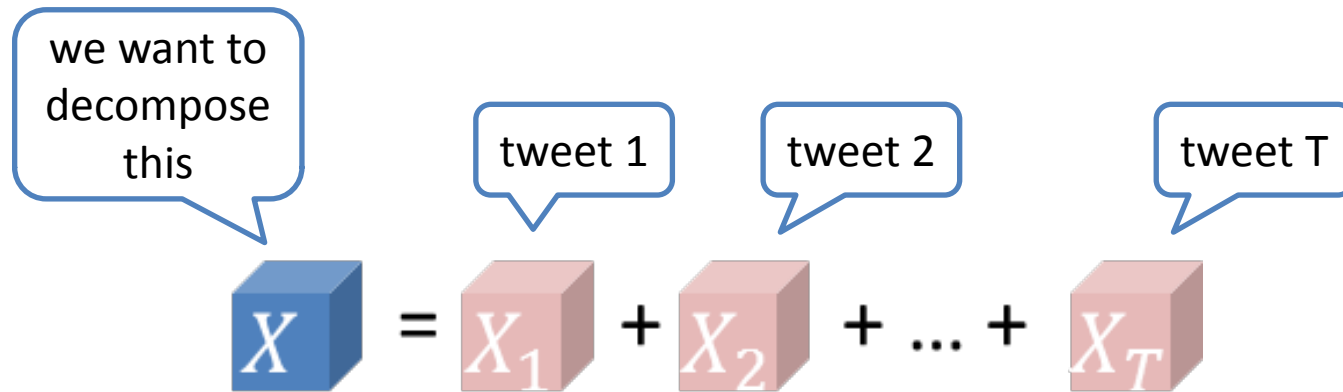
- Can we get decomposition without  tensor ?

YES

- If possible, when and how?

- When:  tensor is given by sum
- How: stochastic optimization

Revisit: Word Co-occurrence

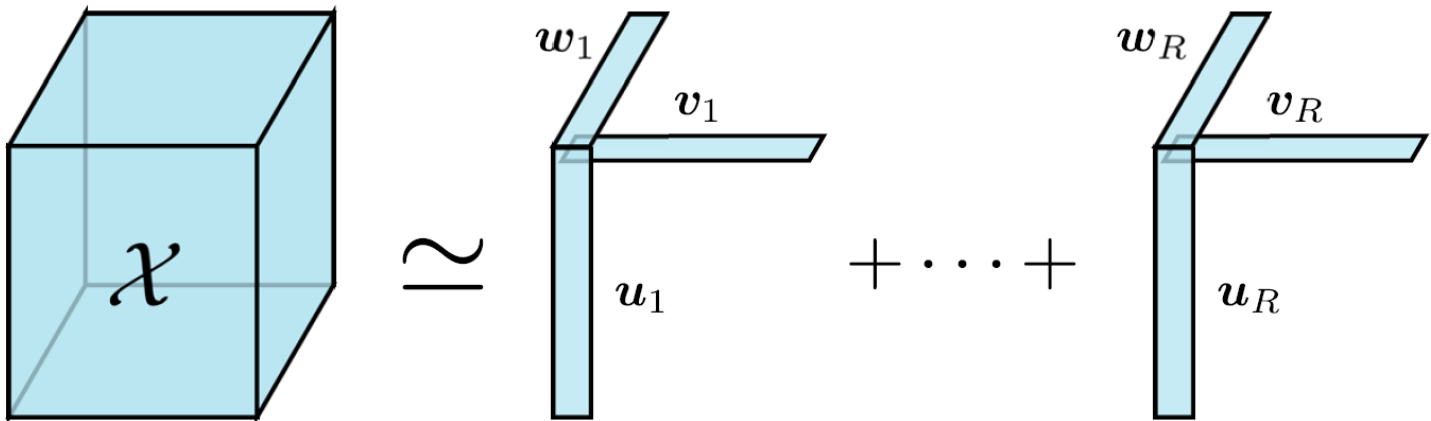


- X : “Heavy” tensor (e.g. dense)
- X_t : “Light” tensor (e.g. sparse)

CP Decomposition

$$\min_{U, V, W} L(\mathcal{X}; U, V, W)$$

- $L(\mathcal{X}; U, V, W) = \|\mathcal{X} - \sum_{r=1}^R u_r \circ v_r \circ w_r\|^2$



Batch Optimization

- CP decomposition is non-convex
- Use alternating least squares

ALS

Repeat until convergence:

- $U \leftarrow U - \eta \nabla_U L(x)$
- $V \leftarrow V - \eta \nabla_V L(x)$
- $W \leftarrow W - \eta \nabla_W L(x)$

step size

Key Fact: Loss is Decomposable

$$L(\text{X})$$

$$L(X) = \left\| X - \sum_{r=1}^R u_r \circ v_r \circ w_r \right\|^2$$

$$= L(X_1 + X_2 + \dots + X_T)$$

$$= L(X_1) + L(X_2) + \dots + L(X_T) + \text{const}$$

Intuition: for stochastic variable x and non-stochastic variable a ,

$$\begin{aligned} & \underline{(Ex - a)^2} \\ &= (Ex)^2 - 2aEx + a^2 \\ &= \underline{E[x^2] - 2aEx + a^2} + \underline{(Ex)^2 - E[x^2]} \\ &= \underline{E[x - a]^2} + \underline{\text{Var}[x]} \end{aligned}$$

Stochastic Optimization

- If $L = L_1 + L_2 + \dots + L_T$,
we can use ∇L_t instead of ∇L

➔ Stochastic ALS

ALS

Repeat until convergence:

- $U \leftarrow U - \eta \nabla_U L(\mathbf{x})$
- $V \leftarrow V - \eta \nabla_V L(\mathbf{x})$
- $W \leftarrow W - \eta \nabla_W L(\mathbf{x})$

SALS

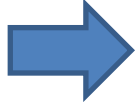
for $t = 1, \dots, T$:

- $U \leftarrow U - \eta_t \nabla_U L(\mathbf{x}_t)$
- $V \leftarrow V - \eta_t \nabla_V L(\mathbf{x}_t)$
- $W \leftarrow W - \eta_t \nabla_W L(\mathbf{x}_t)$

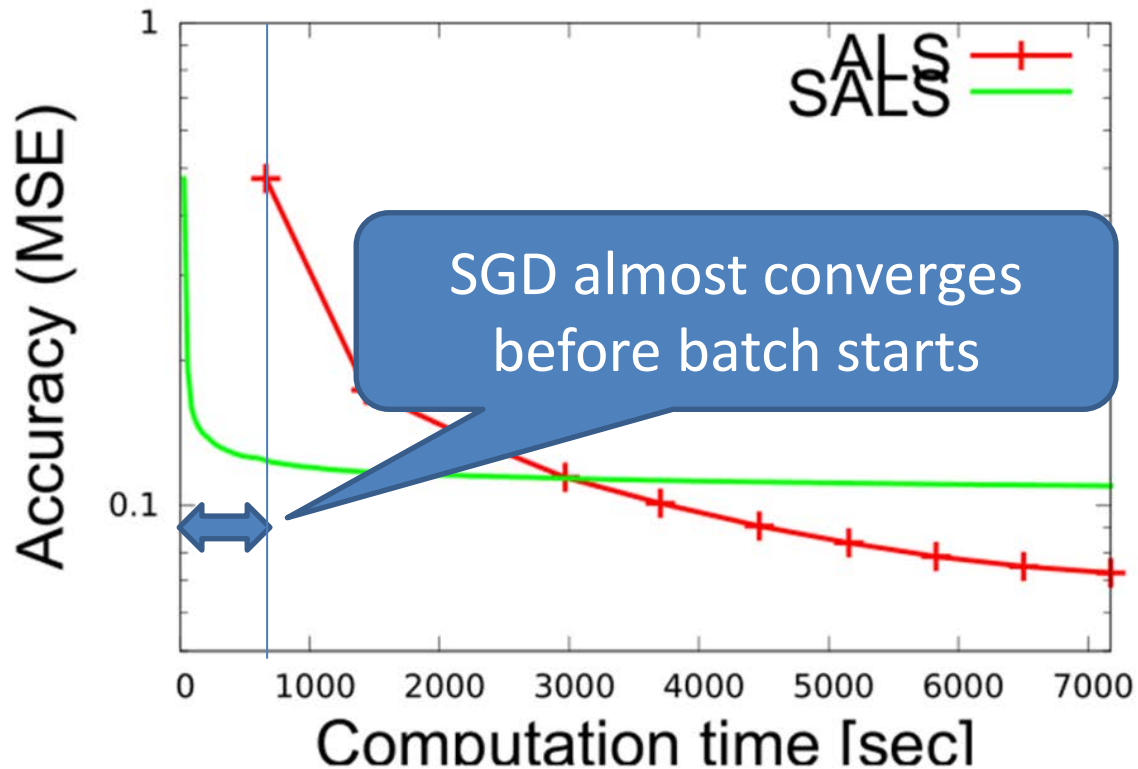
Observations

- ALS and SALS converge to the same solutions
 - $(L_{ALS} - L_{SALS}) = \text{const}$
- SALS does not require to precompute X
 - Batch gradient: $\nabla L(X)$
 - Stochastic gradient: $\nabla L(x_t)$
- One step of SALS is lighter
 - x_t is sparse, whereas X is not
 - $\nabla L(x_t)$ is also sparse

Related Work

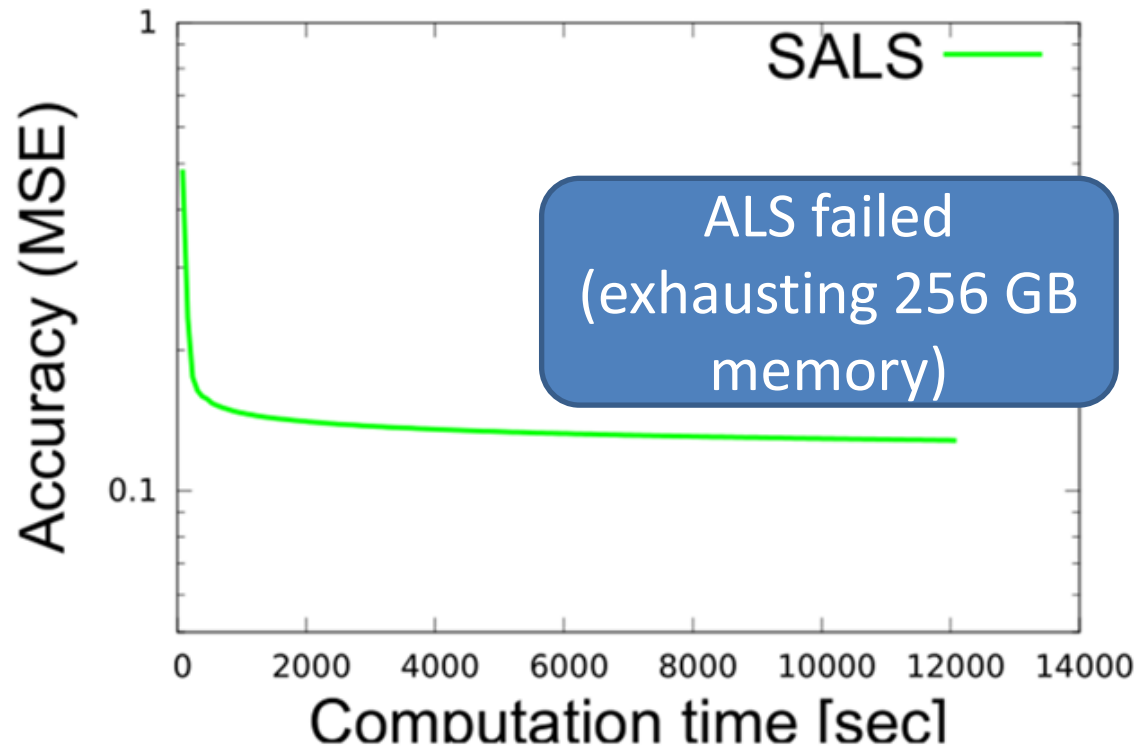
- Element-wise SGD
-  A special case of SALS

Experiment



- Amazon Review datasets (word triplets)
- 10M reviews, 280K vocabs, 500M non-zeros

Experiment (Cont'd)



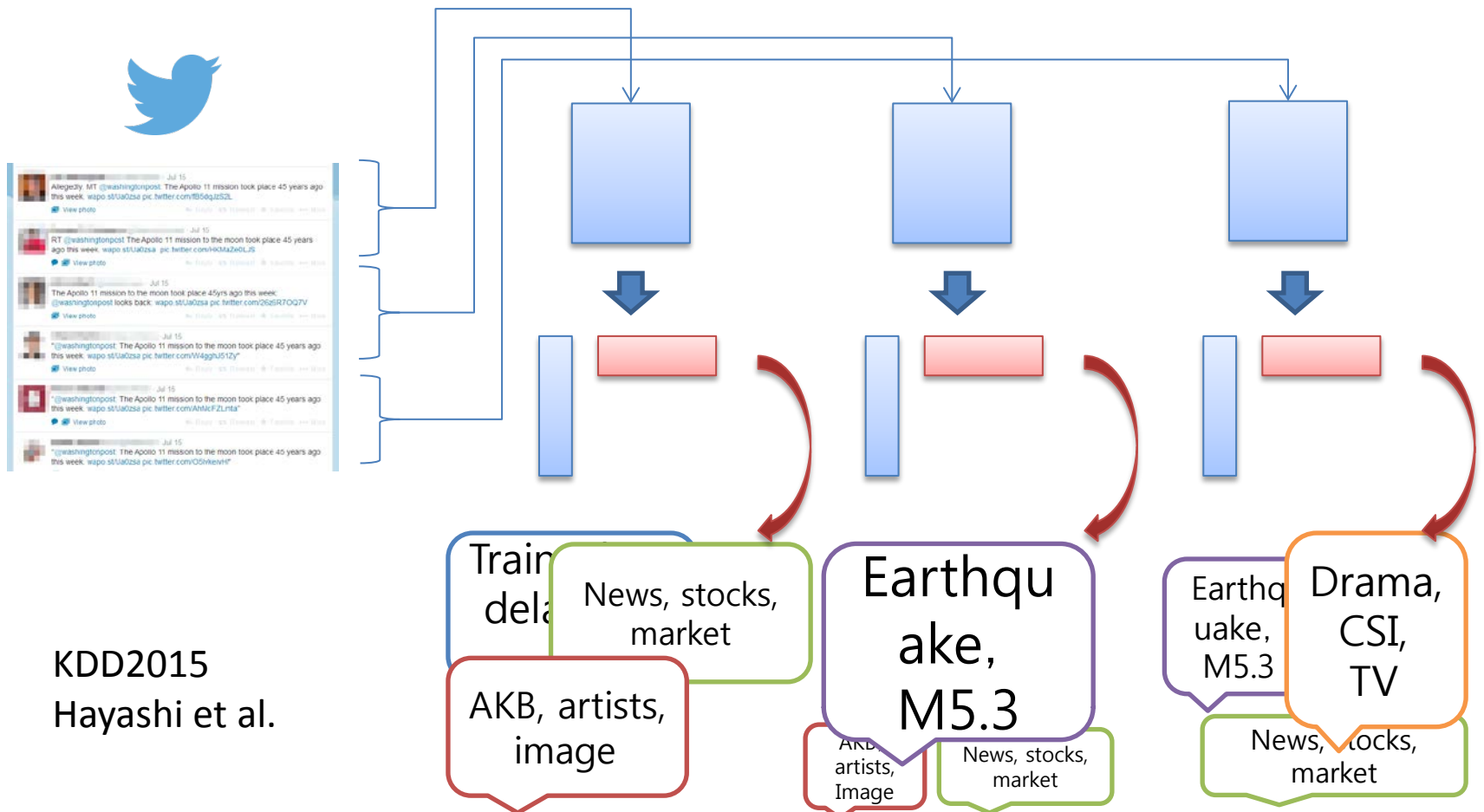
- 34M reviews, 520K vocabs, 1G non-zeros

More on Paper

- 2nd-order SALS
 - Using block diagonal of Hessian
 - Faster convergence, low sensitivity of η
- With l2 norm regularizer:
 - $\min L(X; U, V, W) + R(U, V, W),$
 $R(U, V, W) = \lambda(\|U\|^2 + \|V\|^2 + \|W\|^2)$
- Theoretical analysis of convergence

Further Applications



- Online Learning



KDD2015
Hayashi et al.

Take Home Message

If  is given by sum ...

- Don't use 
- Use 

$$\mathbf{X} = \mathbf{X}_1 + \mathbf{X}_2 + \dots + \mathbf{X}_T$$