

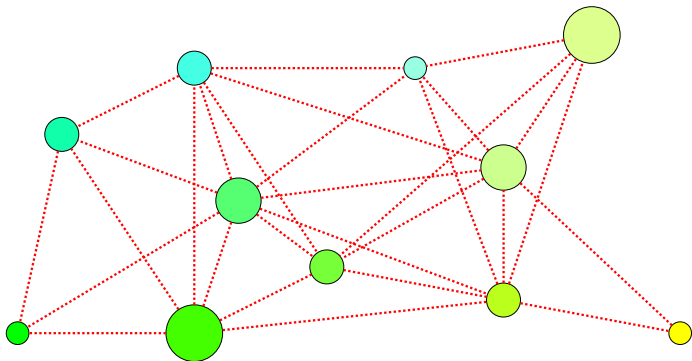
Large-Scale Collection Threading using Structured k -DPPs

Jennifer Gillenwater, Alex Kulesza, Ben Taskar

University of Pennsylvania

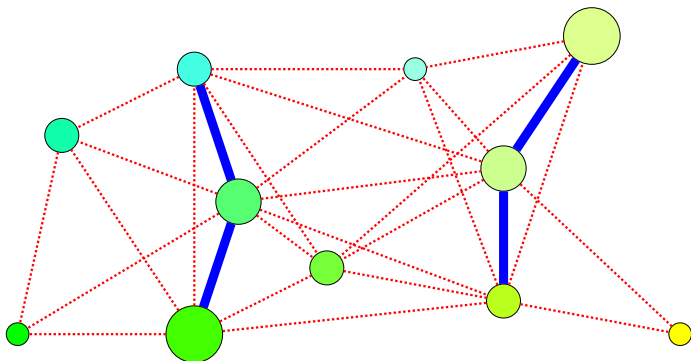
Novel task definition

Select a *high-quality set of diverse paths* in a data **graph**



Novel task definition

Select a *high-quality set of diverse paths* in a data **graph**





Feb 24: Parkinson's Disease Increases Risks to Pope

Feb 24: Parkinson's Disease Increases Risks to Pope

Feb 26: Pope's Health Raises Questions About His Ability to Lead

Feb 24: Parkinson's Disease Increases Risks to Pope

Feb 26: Pope's Health Raises Questions About His Ability to Lead

Mar 13: Pope Returns Home After 18 Days at Hospital

Feb 24: Parkinson's Disease Increases Risks to Pope

Feb 26: Pope's Health Raises Questions About His Ability to Lead

Mar 13: Pope Returns Home After 18 Days at Hospital

Apr 01: Pope's Condition Worsens as World Prepares for End of Papacy

Feb 24: Parkinson's Disease Increases Risks to Pope

Feb 26: Pope's Health Raises Questions About His Ability to Lead

Mar 13: Pope Returns Home After 18 Days at Hospital

Apr 01: Pope's Condition Worsens as World Prepares for End of Papacy

Apr 02: Pope, Though Gravely Ill, Utters Thanks for Prayers

Feb 24: Parkinson's Disease Increases Risks to Pope

Feb 26: Pope's Health Raises Questions About His Ability to Lead

Mar 13: Pope Returns Home After 18 Days at Hospital

Apr 01: Pope's Condition Worsens as World Prepares for End of Papacy

Apr 02: Pope, Though Gravely Ill, Utters Thanks for Prayers

Apr 18: Europeans Fast Falling Away from Church

Feb 24: Parkinson's Disease Increases Risks to Pope

Feb 26: Pope's Health Raises Questions About His Ability to Lead

Mar 13: Pope Returns Home After 18 Days at Hospital

Apr 01: Pope's Condition Worsens as World Prepares for End of Papacy

Apr 02: Pope, Though Gravely Ill, Utters Thanks for Prayers

Apr 18: Europeans Fast Falling Away from Church

Apr 20: In Developing World, Choice [of Pope] Met with Skepticism

Feb 24: Parkinson's Disease Increases Risks to Pope

Feb 26: Pope's Health Raises Questions About His Ability to Lead

Mar 13: Pope Returns Home After 18 Days at Hospital

Apr 01: Pope's Condition Worsens as World Prepares for End of Papacy

Apr 02: Pope, Though Gravely Ill, Utters Thanks for Prayers

Apr 18: Europeans Fast Falling Away from Church

Apr 20: In Developing World, Choice [of Pope] Met with Skepticism

May 18: Pope Sends Message with Choice of Name

- Selecting a *single* thread (D. Shahaf and C. Guestrin, KDD 2010)

Related threading work

- Selecting a *single* thread (D. Shahaf and C. Guestrin, KDD 2010)
- Building diverse *topic* threads (A. Ahmed and E. Xing, UAI 2010)

Approach: Structured Determinantal Point Processes

Score a set of threads \mathbf{Y} via structured determinantal point process (SDPP)¹

¹(A. Kulesza and B. Taskar, NIPS 2010)

Score a set of threads \mathbf{Y} via structured determinantal point process (SDPP)¹

$$L_{ij} = \mathbf{q}(\mathbf{y}_i) \phi(\mathbf{y}_i)^\top \phi(\mathbf{y}_j) \mathbf{q}(\mathbf{y}_j)$$

¹(A. Kulesza and B. Taskar, NIPS 2010)

Score a set of threads \mathbf{Y} via structured determinantal point process (SDPP)¹

$$\mathcal{P}(\mathbf{Y}) \propto \det(\mathbf{L}_{\mathbf{Y}})$$

¹(A. Kulesza and B. Taskar, NIPS 2010)

Score a set of threads \mathbf{Y} via structured determinantal point process (SDPP)¹

$$\mathbf{Y} = \{\mathbf{i}\} \rightarrow \mathcal{P}(\mathbf{Y}) \propto \mathbf{q}(\mathbf{y}_i)^2$$

¹(A. Kulesza and B. Taskar, NIPS 2010)

Score a set of threads \mathbf{Y} via structured determinantal point process (SDPP)¹

$$\mathbf{Y} = \{i\} \rightarrow \mathcal{P}(\mathbf{Y}) \propto \mathbf{q}(\mathbf{y}_i)^2$$

$$\mathbf{Y} = \{i, j\} \rightarrow \mathcal{P}(\mathbf{Y}) \propto \mathbf{q}(\mathbf{y}_i)^2 \mathbf{q}(\mathbf{y}_j)^2 (1 - (\phi(\mathbf{y}_i)^\top \phi(\mathbf{y}_j))^2)$$

¹(A. Kulesza and B. Taskar, NIPS 2010)

k-SDPPs¹: fix # of points in **Y** to **k**

¹(A. Kulesza and B. Taskar, ICML 2011)

$$O(\text{Trn}D^2 + D^3)$$

T = thread length

n = # of nodes

r = maximum node degree

D = # of features

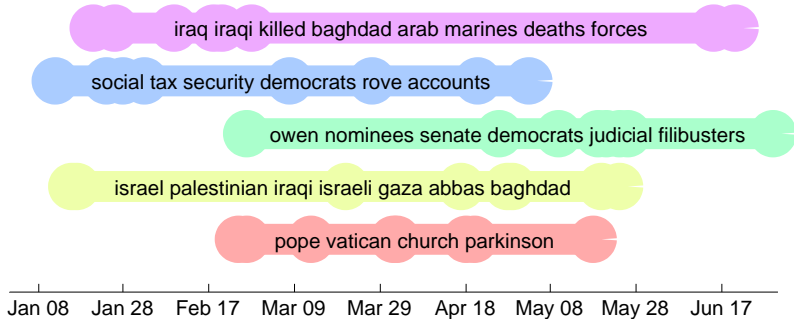
- Project D down to d

Random Projection for Tractability

- Project \mathbf{D} down to \mathbf{d}
- **Theorem:** Given $\tilde{\mathbf{P}}^k(\mathbf{Y}) =$ distribution after projecting

$$\|\mathbf{P}^k - \tilde{\mathbf{P}}^k\|_1 \leq \text{nice small things}$$

New York Times Timelines



New York Times Timelines

k-means DTM k-SDPP		

	Human summary similarity	
k -means	4.32	
DTM	3.78	
k -SDPP	8.26	

New York Times Timelines

	Human summary similarity	Runtime (sec)
k-means	4.32	625
DTM	3.78	19,443
k-SDPP	8.26	252