# Leveraging Temporal Dynamics of Document Content in Relevance Ranking

Jonathan L. Elsas (CMU)

Susan T. Dumais (MSR)

# Outline

- Document Dynamics on the Web
  - Previous Work on Change & What's Missing
- Our Setting: Ranking Dynamic Documents
  - Test Collection & Measuring change
- Two ways to leverage change in ranking
  - Document Prior based on Gross Change Measures
  - Document Representation based on Term-Level Change
- Discussion

# The Web is Dynamic

# The Web is Dynamic

# The Web is Dynamic

# The Web is Dynamic



10 Minutes

# Previous Work on Dynamics

Characterizing Change

Implications of Change

# Characterizing Change

## Change & Page Type / Source



Figure 12: Clustered rates of change, broken down by selected top-level domains, and omitting the *no change* cluster.

Fetterly et al, WWW03

## New Content & Links

Ntoulas et al, WWW04

## Within-Doc. Change

Adar et al, WSDM 2009

# Implications of Change

On Browsing

On Indexing

On Crawling



Adar, et al UIST'08

Berberich, et al SIGIR'07

Olston & Pandey, WWW08

# Implications of Change

## On Ranking?

# Implications of Change on Ranking

- Gross Measures of Document Change

  Are there general characteristics of document dynamics indicate high quality pages?

- Representing Term-Level Change Within the Document

  Are there characteristics of a document's dynamic content that indicate some content may be more important?

# Test Setup: Queries & Documents

- 18K Queries, 2.5M Judged Documents
  - 5-level relevance judgment (Bad...Perfect)
- 2.5M Documents crawled weekly for 10 weeks


- Navigational queries
  - 2k queries identified with a "Perfect" judgment
- 60/40 Training/Test split

# Test Setup: Queries & Documents

- 18K Queries, 2.5M Judged Documents
  - 5-level relevance
- 2.5M Documents

> We focus on Navigational Queries here for ease of evaluation.

- Navigational queries
  - 2k queries identified with a "Perfect" judgment
- 60/40 Training/Test split

# Measuring Change: Shingleprints



D(ti)　　　　　　　　　　　　　　　　　　D(ti+1)

#be9e　　#be9e

#aaef　　#81d3

#a559　　#a559

#18ef　　#18ef

#744e　　#fa6e

#b256　　#b256

Sh(D(ti))　　Sh(D(ti+1))

$$ShDiff(D) = 1 - \frac{1}{T-1} \sum_{t=1}^{T-1} \frac{|Sh(D^{(t)}) \cap Sh(D^{(t+1)})|}{N}$$

Broder, et al, "Syntactic Clustering of the Web"
Computer Networks & ISDN Sys., 1997

# Change & *Relevance*

# Change & *Relevance*

Change Amount & *Relevance*

# Change & *Relevance*

- More relevance documents tend to change more often, *and* to a greater degree than non-relevant documents.

- Could favoring dynamic documents in ranking improve performance?

# Favoring Dynamic Documents

Language-Modeling Ranking Function:

$$P(D|Q) \propto P(D)P(Q|D)$$

# Favoring Dynamic Documents

Language-Modeling Ranking Function:

$$P(D|Q) \propto \boxed{P(D)} P(Q|D)$$

Uniform Prior:

$$P(D) \propto 1.0$$

"Static Model"

# Favoring Dynamic Documents

Language-Modeling Ranking Function:

$$P(D|Q) \propto \boxed{P(D)} P(Q|D)$$

Uniform Prior:

$$P(D) \propto 1.0$$

"Change" Prior:

$$P_{ch}(D) \propto (ShDiff(D) + 1)^{\gamma}$$

Favoring Dynamic Documents

# Favoring Dynamic Documents

# Change Within the Document

Are there characteristics of a document's dynamic content that indicate some content may be more important?

# Change Within the Document





Adar, Teevan, Dumais & Elsas, "The Web Changes Everything: Understanding the Dynamics of Web Content" WSDM 2009

# Change Within the Document

Adar, Teevan, Dumais & Elsas, "The Web Changes Everything: Understanding the Dynamics of Web Content" WSDM 2009

# Change Within the Document




allrecipes.com

Sep.    Oct.    Nov.    Dec.

**Time**

Adar, Teevan, Dumais & Elsas, "The Web Changes Everything: Understanding the Dynamics of Web Content" WSDM 2009

# Change Within the Document

# Leveraging Change Within the Document

Identifying transient & permanent vocabulary:

- **Short-lived**: come & go quickly

    in fewer than 50% of the document's slices

- **Medium-lived**:

    in 50-90% of the document's slices

- **Long-lived**: tend to stick for a long time

    in > 90% of the document's slices

# Leveraging Change Within the Document

Model relevance as a *mixture* of LONG- MEDIUM- and SHORT-lived vocabulary:

$$P(D|Q) \propto P(D)\Big(\lambda_L P(Q|D_L)$$
$$+\lambda_M P(Q|D_M)$$
$$+\lambda_S P(Q|D_S)\Big)$$

"Dynamic Model"

# Leveraging Change Within the Document

Model relevance as a *mixture* of LONG- MEDIUM- and SHORT-lived vocabulary:

$$P(D|Q) \propto P(D)\Big(\boxed{\lambda_L P(Q|D_L)}$$
$$+ \lambda_M P(Q|D_M)$$
$$+ \boxed{\lambda_S P(Q|D_S)}\Big)$$

Differentially weight long-lived and short-lived vocabulary.
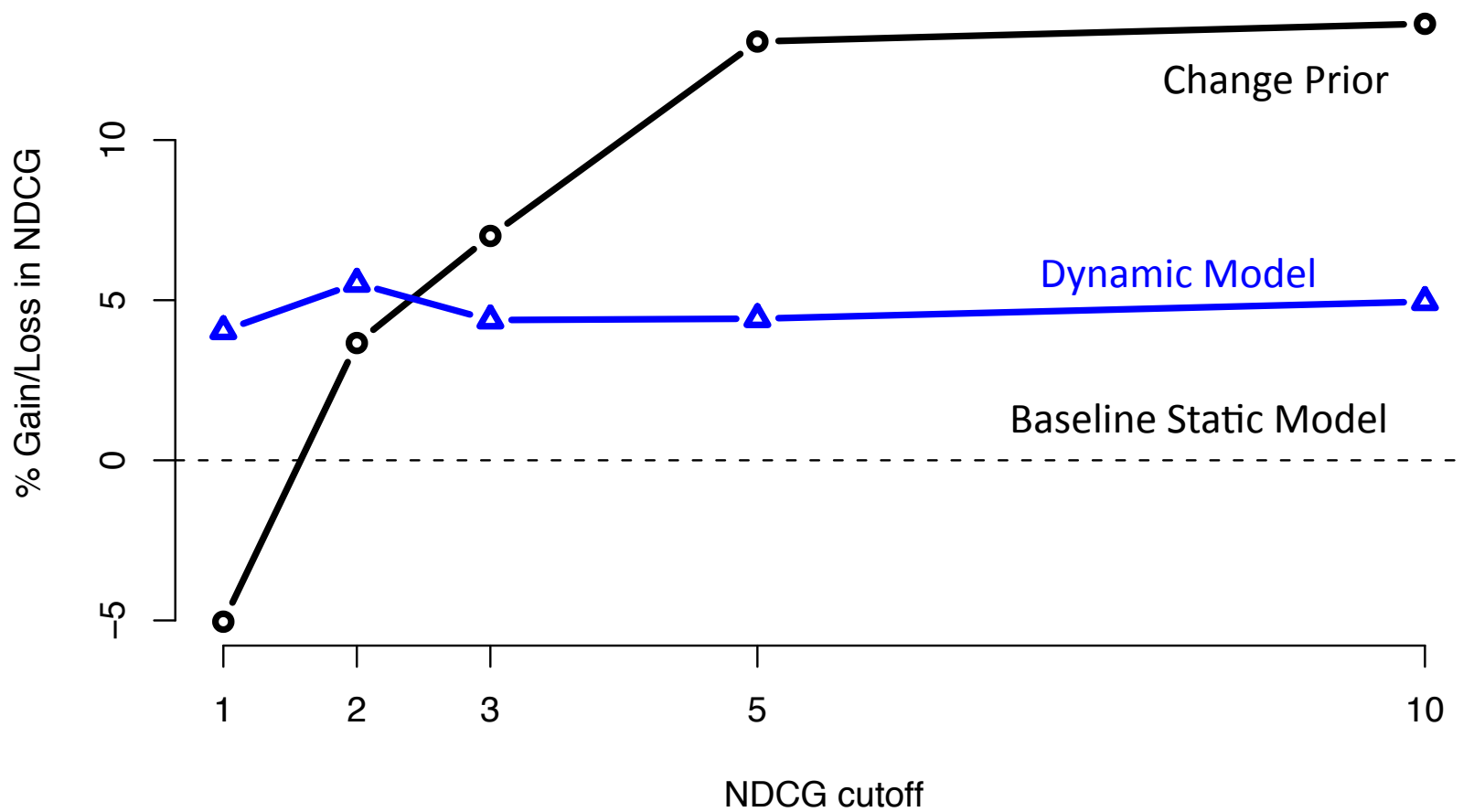
"Dynamic Model"

# Leveraging Change Within the Document

# Change & Relevance Ranking

- Presented two methods for leveraging changing content in relevance ranking:

  - **Query-Independent Change Prior,** favoring dynamic documents irrespective of query

  - **Dynamic Document Representation**, differentially weighting long-term and short-term vocabulary

- Combined Model: Best of both worlds?

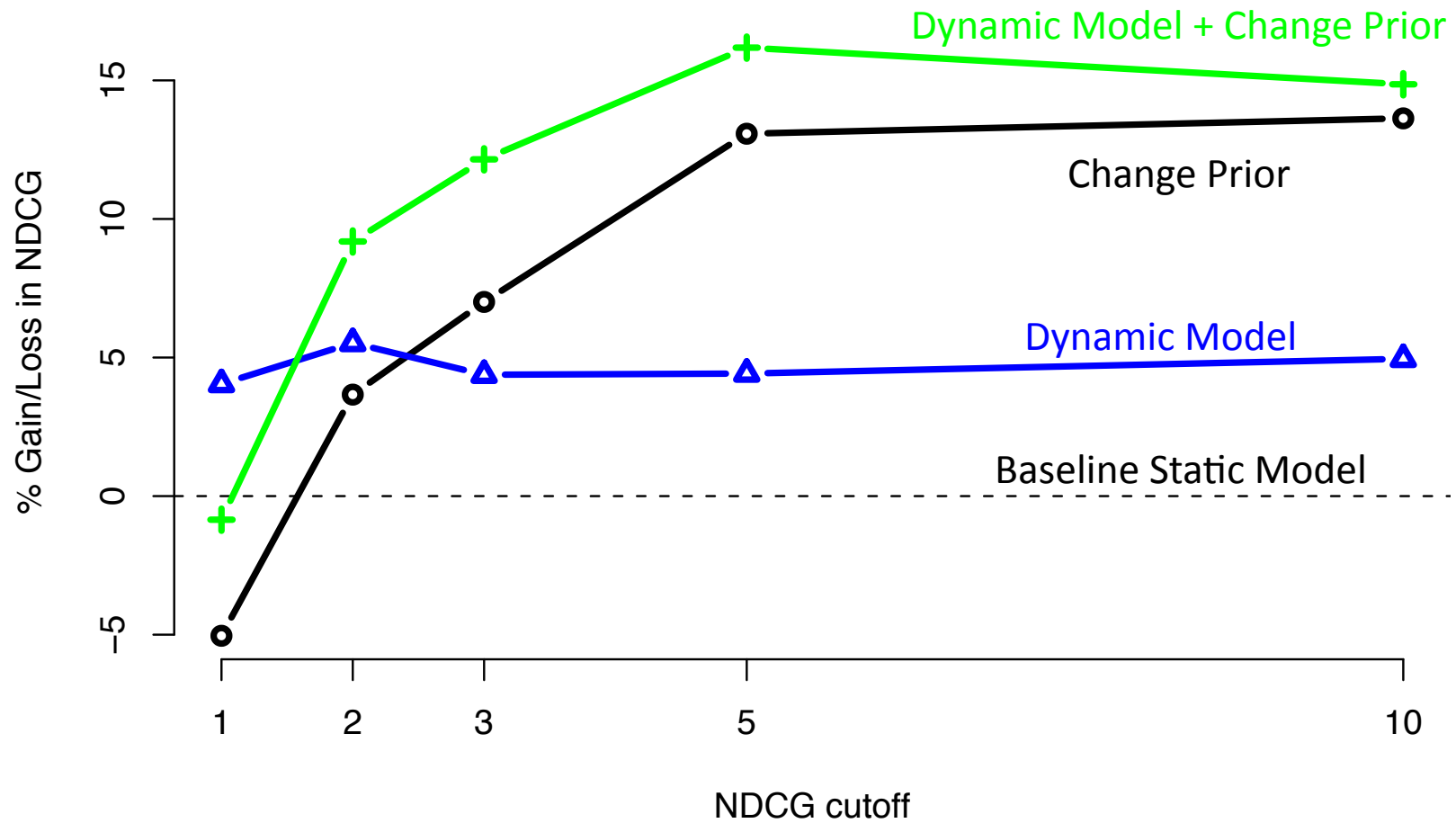# Combined Model

Dynamic Model + Change Prior

Change Prior

Dynamic Model

Baseline Static Model

% Gain/Loss in NDCG

NDCG cutoff

# Conclusion & Next Steps

- Documents change, and we can use characteristics of those dynamics to improve retrieval performance.

- Presented two complementary methods of leveraging change in ranking.

- Focus here on navigational queries; current work is looking at *dynamic* information needs.

  Relevance may change over time.

# Thank You

# Questions?