

名詞句の内部構造を考慮した キーワードのスコア付け

2013年3月15日

京都大学

村脇 有吾 黒橋 禎夫

名詞句の内部構造を考慮した
~~キーワード~~のスコア付け
キーフレーズ

2013年3月15日

京都大学

村脇 有吾 黒橋 禎夫

bag-of-wordsを超える何か

ト
キ
ス
ト

computational finite-element schemes for optimal control of an elliptic system with conjugation conditions

new optimal control problems are considered for distributed systems described by ...

bag-of-wordsを超える何か

トキスト

computational finite-element schemes for
optimal control of an elliptic system with
conjugation conditions
new optimal control problems are considered
for distributed systems described by ...

Latent topic model

 TOPIC₁
 TOPIC₂
 TOPIC₃

bag-of-wordsを超える何か

トキスト

computational finite-element schemes for optimal control of an elliptic system with conjugation conditions

new **optimal control problems** are considered for distributed systems described by ...

bag-of-wordsを超える何か

ト
キ
ス
ト

computational finite-element schemes for optimal control of an elliptic system with conjugation conditions

new **optimal control problems** are considered for distributed systems described by ...

単語より長い
意味的なまとまり

キーフレーズ抽出すら 21世紀になっても bag-of-words



Aria Haghighi
@aria42

Follow

NLP folks: For Keyphrase extraction does anything *conclusively* beat sum of TF-IDF scores for contiguous adjectives and noun sequences?

Reply Retweet Favorite More

1
RETWEET

5
FAVORITES



6:06 AM - May 29, 2011 from South Kingstown, RI

キーフレーズ抽出すら 21世紀になっても bag-of-words



Aria Haghghi

@aria42

Follow

①

NLP folks: For **Keyphrase extraction** does anything *conclusively* beat sum of TF-IDF scores for contiguous adjectives and noun sequences?

Reply Retweet Favorite More

1

RETWEET

5

FAVORITES



む



6:06 AM - May 29, 2011 from South Kingstown, RI

キーフレーズ抽出すら 21世紀になっても bag-of-words



Aria Haghighi
@aria42

Follow

NLP folks: For ^① Keyphrase extraction does ^② anything *conclusively* beat sum of TF-IDF scores for contiguous adjectives and noun sequences?

Reply Retweet Favorite More

1
RETWEET

5
FAVORITES



む



6:06 AM - May 29, 2011 from South Kingstown, RI

キーフレーズ抽出すら 21世紀になっても bag-of-words



Aria Haghighi
@aria42

Follow

NLP folks: For ^① Keyphrase extraction does ^② anything *conclusively* beat sum of TF-IDF scores for contiguous adjectives and noun sequences?

Reply Retweet Favorite More

^③ tf-idf法の拡張

1
RETWEET

5
FAVORITES



6:06 AM - May 29, 2011 from South Kingstown, RI

教師なしキーフレーズ抽出

ト
キ
ス
ト

computational finite-element schemes for optimal control of an elliptic system with conjugation conditions

new optimal control problems are considered for distributed systems described by ...

教師なしキーフレーズ抽出

トキスト

computational finite-element schemes for optimal control of an elliptic system with conjugation conditions
new optimal control problems are considered for distributed systems described by ...

システム出力

elliptic system
optimal control problems
distributed systems
....

教師なしキーフレーズ抽出

ト
キ
ス
ト

computational finite-element schemes for optimal control of an elliptic system with conjugation conditions
new optimal control problems are considered for distributed systems described by ...

完全一致による
評価

システム出力

正解リスト

✕ elliptic system
○ optimal control problems
○ distributed systems
....



elliptic equations
optimal control problems
distributed systems
....

スコアに基づくキーフレーズ抽出

キーフレーズ候補	スコア
new optimal control problems	35.8
accurate computational discretization schemes	35.7
quadratic minimized function	26.9
computational finite-element schemes	25.2
optimal control	18.0
conjugate conditions	17.9
elliptic system	17.2
...	

スコアに基づくキーフレーズ抽出

キーフレーズ候補

スコア

スコア
上位
 N 語を
出力

new optimal control problems	35.8
accurate computational discretization schemes	35.7
quadratic minimized function	26.9
computational finite-element schemes	25.2
optimal control	18.0
conjugate conditions	17.9
elliptic system	17.2
...	

実験設定

- Inspecデータセット (Hulth 2003)
 - 英語500論文
 - 表題と要旨 (平均134語)
 - キーフレーズ
 - シソーラスによる統制語
 - 非統制語 ← こちらを正解データとする
- スコア上位 N 語の N を変動させて recall-precision 曲線を描く

実験設定

- Inspecデータセット (Hulth 2003)
 - 英語500論文
 - 表題と要旨 (平均134語)
 - キーフレーズ
 - シソーラスによる統制語
 - 非統制語 ← こちらを正解データとする
- スコア上位 N 語の N を変動させて recall-precision 曲線を描く
- 手法間の相対評価と割りきってください!

準備: 最長名詞列と部分名詞列

ト
キ
ス
ト

computational finite-element schemes for optimal control of an elliptic system with conjugation conditions

new optimal control problems are considered for distributed systems described by ...

準備: 最長名詞列と部分名詞列

ト
キ
ス
ト

computational finite-element schemes for optimal control of an elliptic system with conjugation conditions

new optimal control problems are considered for distributed systems described by ...

1. 最長名詞列: 自動解析結果から規則により抽出される、最も長い名詞・形容詞の列

準備: 最長名詞列と部分名詞列

ト
キ
ス
ト

computational finite-element schemes for optimal control of an elliptic system with conjugation conditions

new optimal control problems are considered for distributed systems described by ...

1. 最長名詞列: 自動解析結果から規則により抽出される、最も長い名詞・形容詞の列
2. 部分名詞列: 最長名詞列に含まれるキーワード候補
 - 不自然な候補も (e.g. control problems)

tf-idf法 (Hasan+ 2010)

$$\text{tfidf}_{doc}(\mathbf{w}) = \text{unit}_{doc}(\mathbf{w}) \times \text{term}_{doc}(\mathbf{w})$$

$$\text{unit}_{doc}(\mathbf{w}) = 1 \text{ if } \mathbf{w} \in doc, 0 \text{ otherwise}$$

$$\text{term}_{doc}(\mathbf{w}) = \sum_i \text{tf}_{doc}(w_i) \times \log(D / D_{w_i})$$

- キーフレーズ候補 \mathbf{w} のスコア: 単語tf-idfの総和 (term)

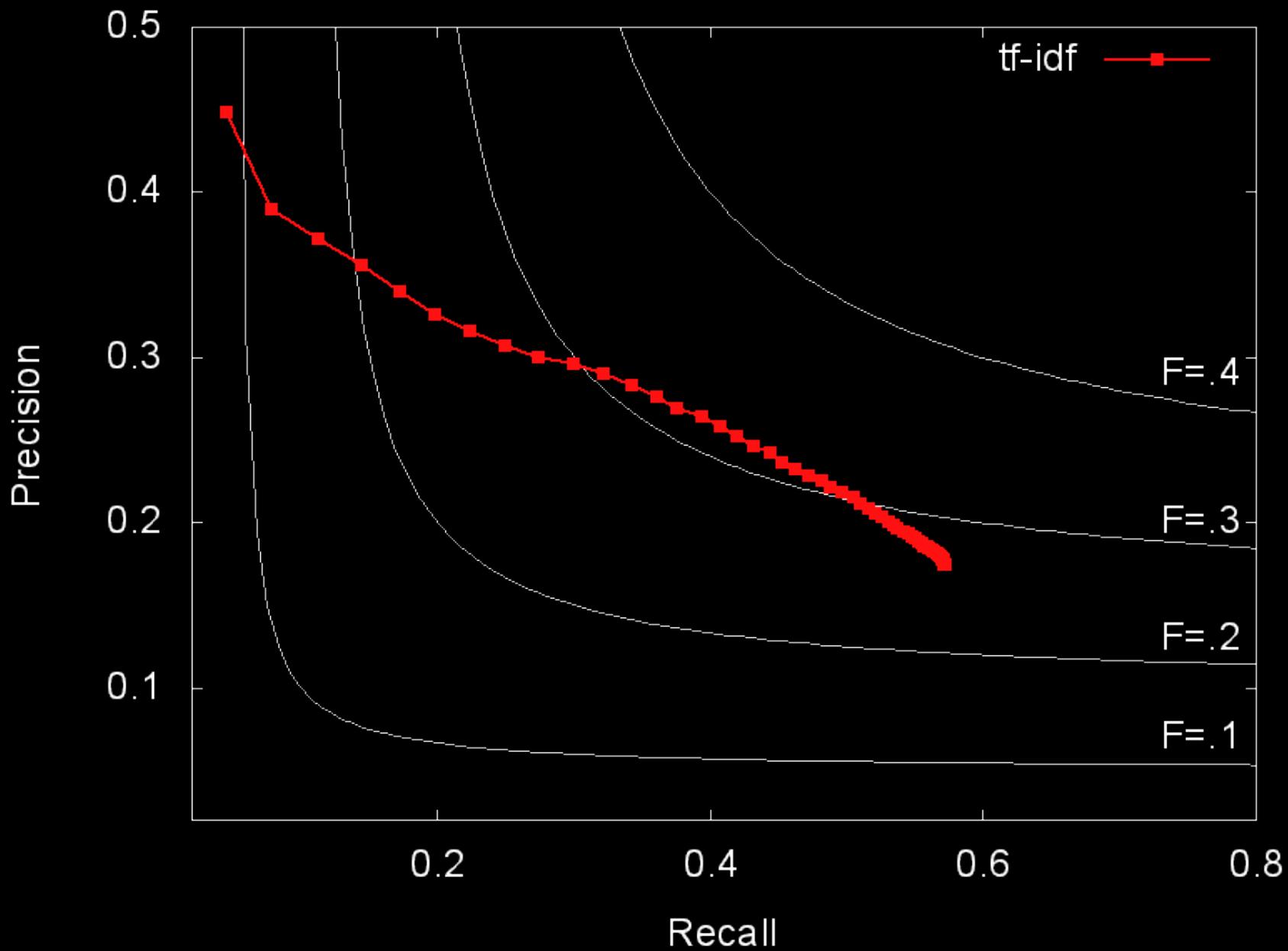
tf-idf法 (Hasan+ 2010)

$$\text{tfidf}_{doc}(\mathbf{w}) = \text{unit}_{doc}(\mathbf{w}) \times \text{term}_{doc}(\mathbf{w})$$

$$\text{unit}_{doc}(\mathbf{w}) = 1 \text{ if } \mathbf{w} \in doc, 0 \text{ otherwise}$$

$$\text{term}_{doc}(\mathbf{w}) = \sum_i \text{tf}_{doc}(w_i) \times \log(D / D_{w_i})$$

- キーフレーズ候補 \mathbf{w} のスコア: 単語tf-idfの総和 (term)
- キーフレーズ候補: 最長名詞列のみ (unit)



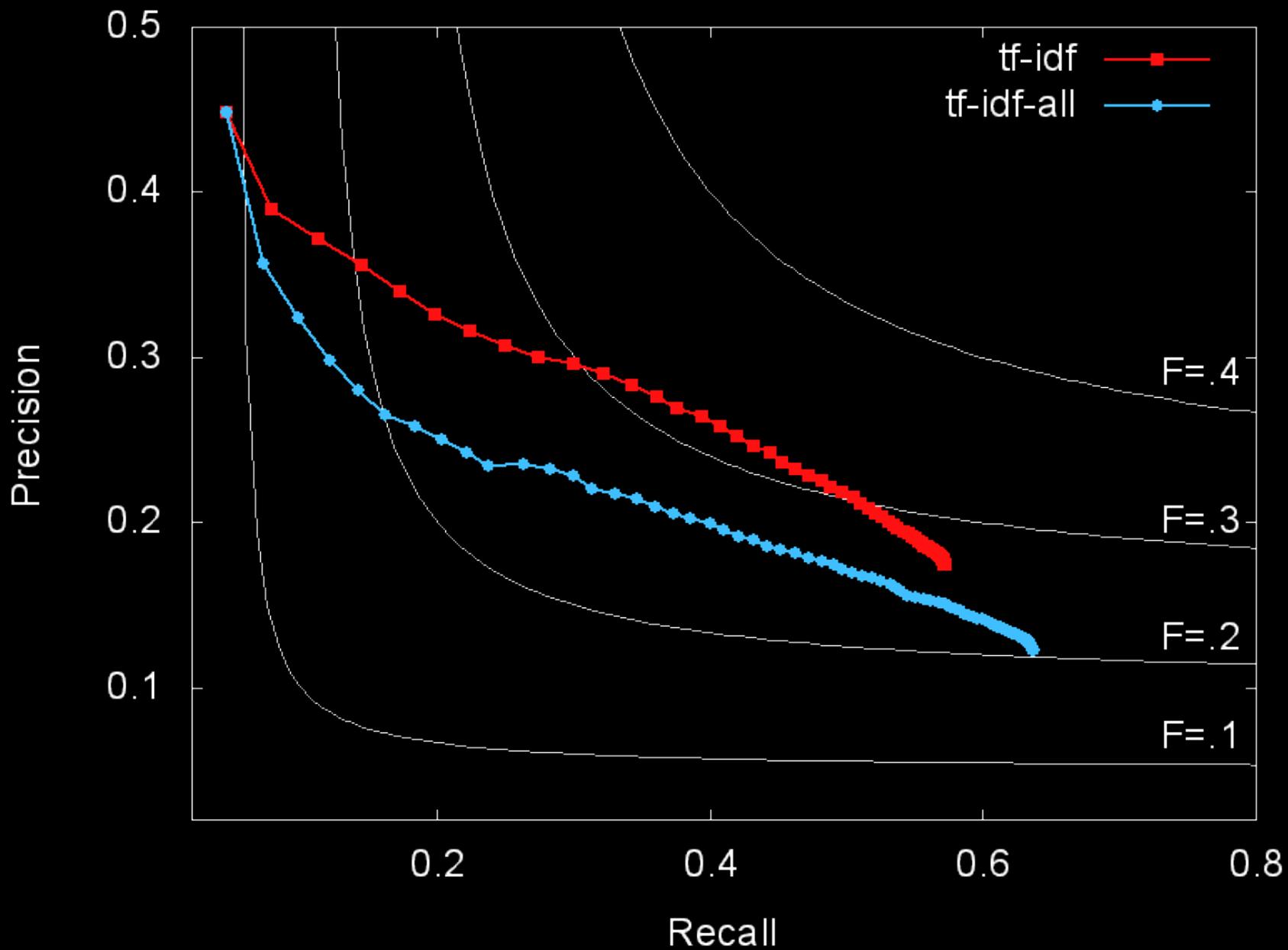
tf-idf-all

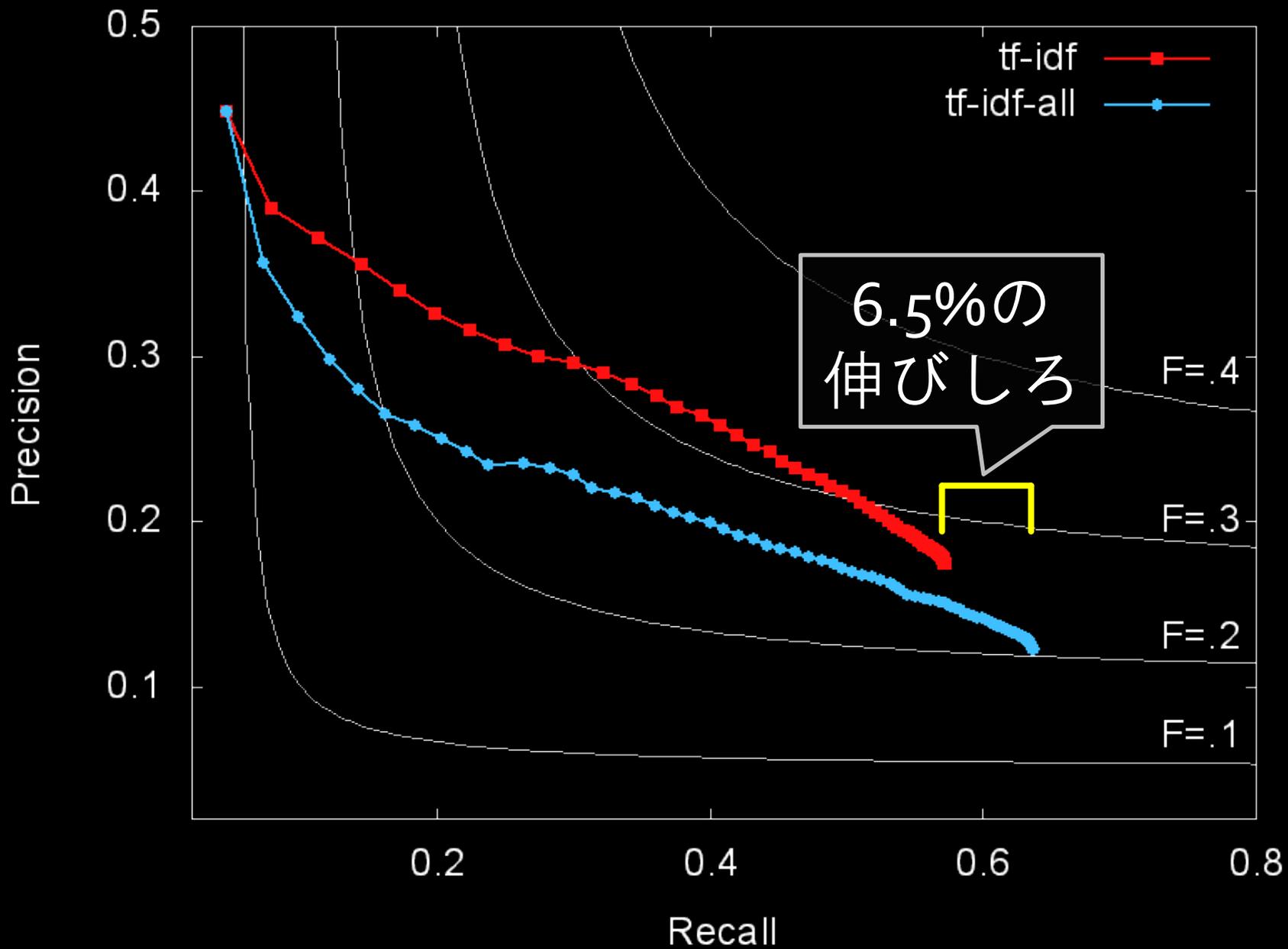
$$\text{tfidf}_{doc}(\mathbf{w}) = \text{unit}_{doc}(\mathbf{w}) \times \text{term}_{doc}(\mathbf{w})$$

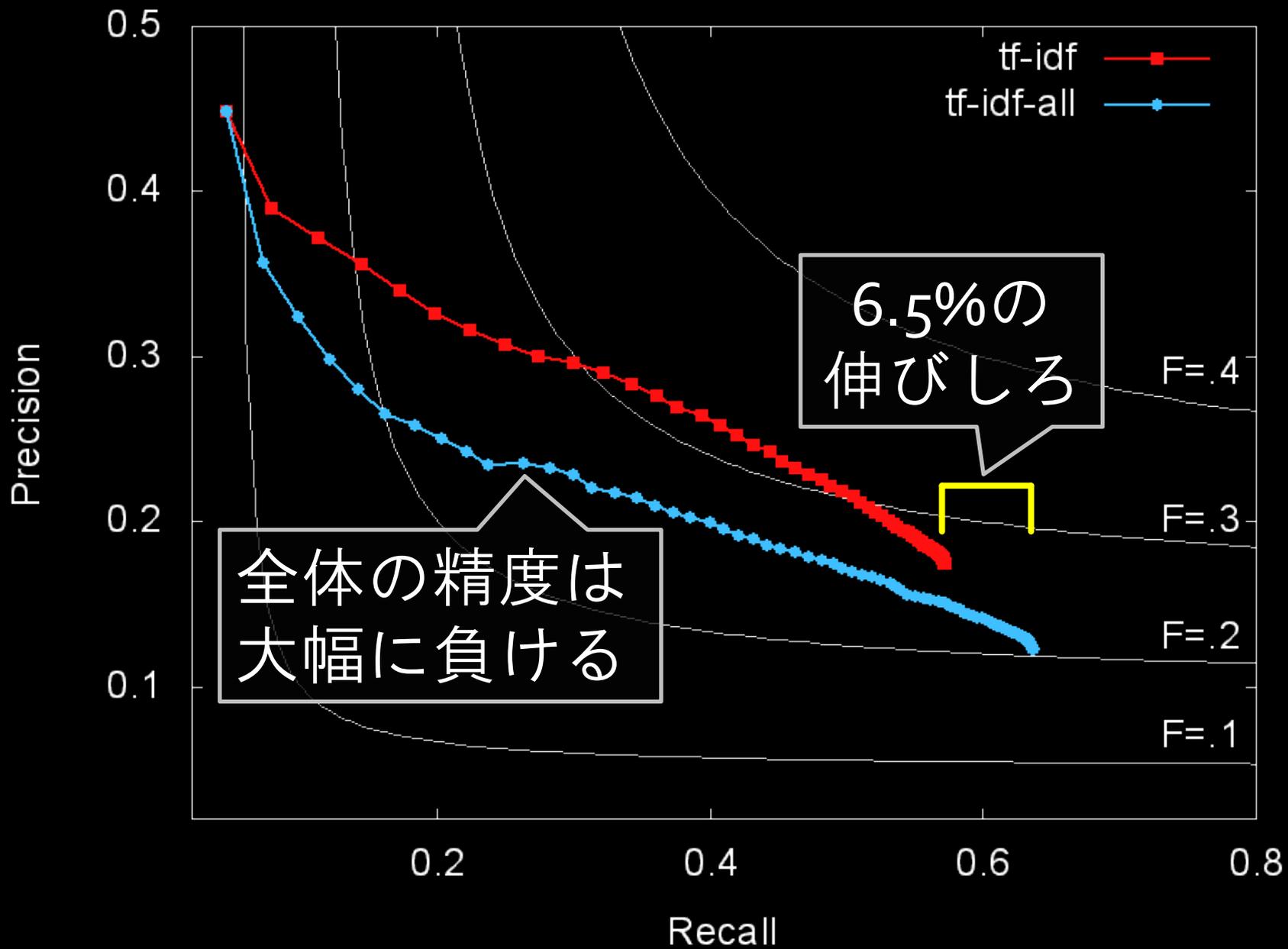
$\text{unit}_{doc}(\mathbf{w}) = 1$ if $\mathbf{w} \in doc$, 0 otherwise

$$\text{term}_{doc}(\mathbf{w}) = \sum_i \text{tf}_{doc}(w_i) \times \log(D / D_{w_i})$$

- キーフレーズ候補 \mathbf{w} のスコア: 単語tf-idfの総和 (term)
- キーフレーズ候補: 部分名詞列も候補にする
 - 😊 recall上界の上昇
(e.g. new optimal control problems → optimal control problems)
 - 😞 不自然な候補も (e.g. control problems)







tf-idf-allは 意味的まとまりを考慮しない

	キーフレーズ候補	スコア
	new optimal control problems	25.7
	optimal control problems	23.5
☹	new optimal control	22.2
	optimal control	20.0
☹	control problems	15.5
	control	12.0
	problems	3.5
	...	

名詞句解析 (Murawaki+ 2012)

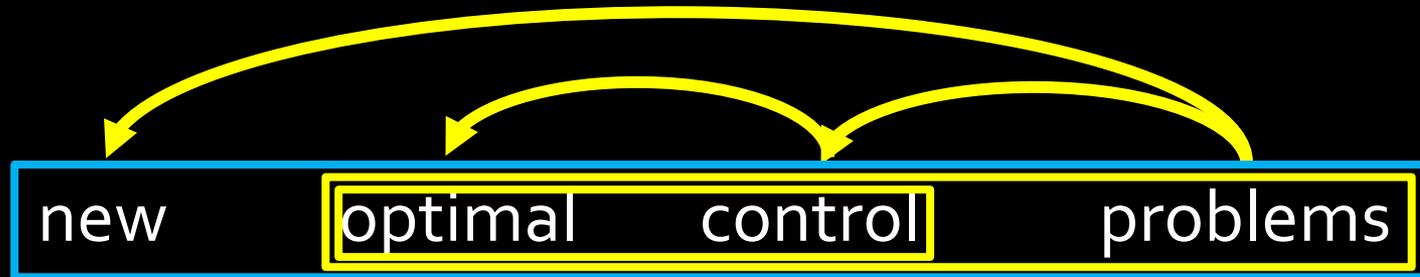
new optimal control problems

名詞句解析 (Murawaki+ 2012)

new optimal control problems

- 意味的にまとまった部分名詞句を同定する

名詞句解析 (Murawaki+ 2012)



- 意味的にまとまった部分名詞句を同定する
- 主辞後置性を仮定すれば係り受けの問題
- 動的計画法で解ける
- 係り受け精度96.3% (Penn Treebankテスト)

tf-term: 名詞句解析に基づく 部分名詞列への頻度配分

new optimal control problems

頻度 1.00

tf-term: 名詞句解析に基づく 部分名詞列への頻度配分

new optimal control problems

new optimal control problems

頻度 1.00

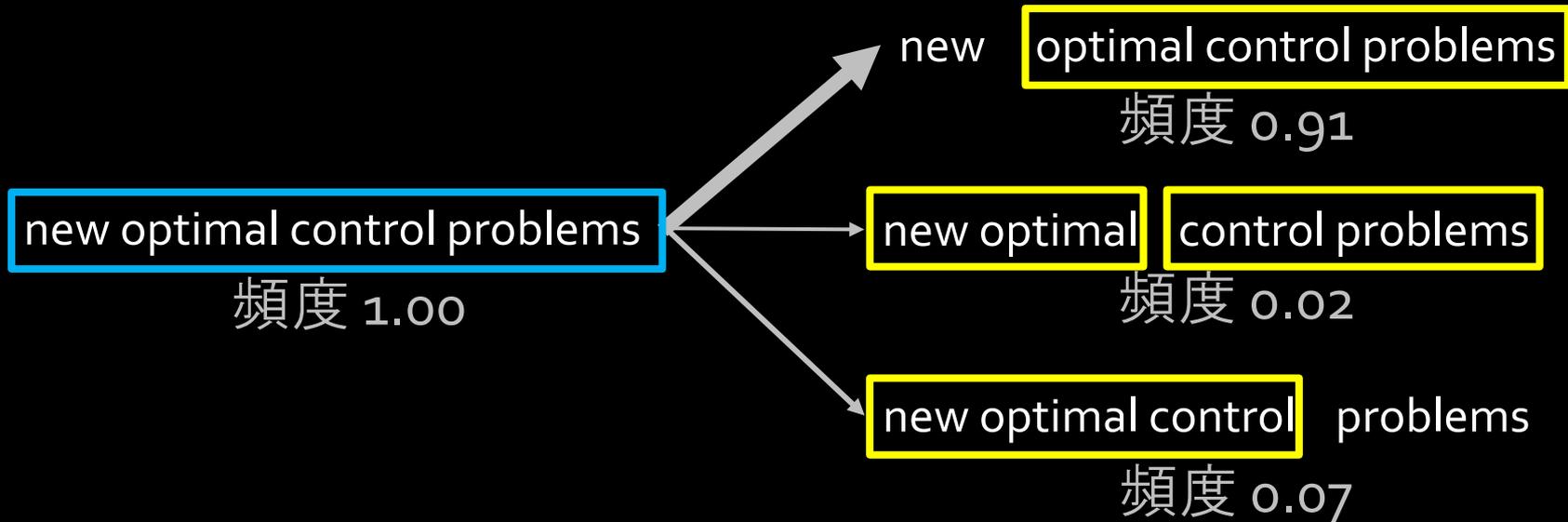
new optimal

control problems

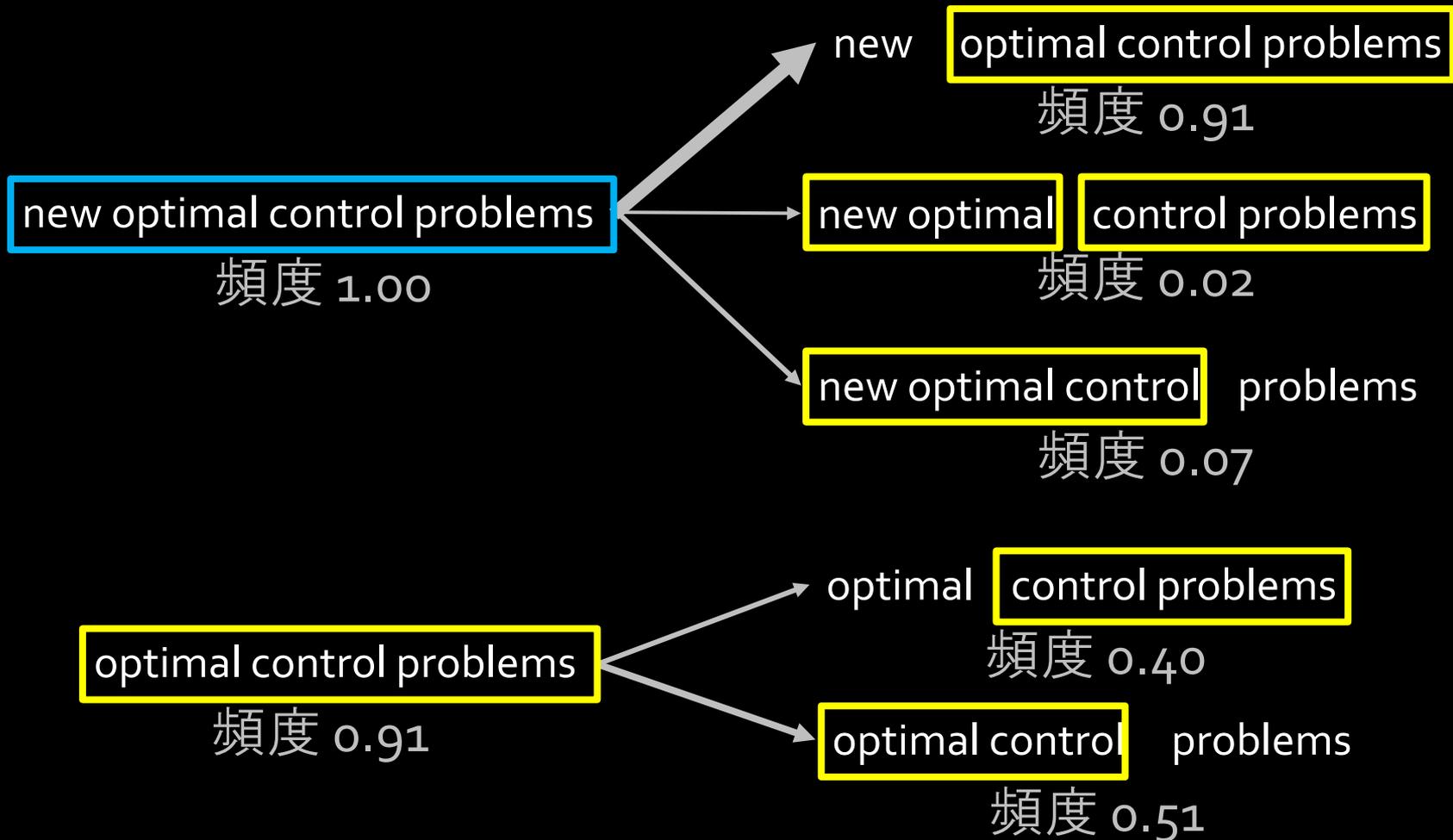
new optimal control

problems

tf-term: 名詞句解析に基づく 部分名詞列への頻度配分

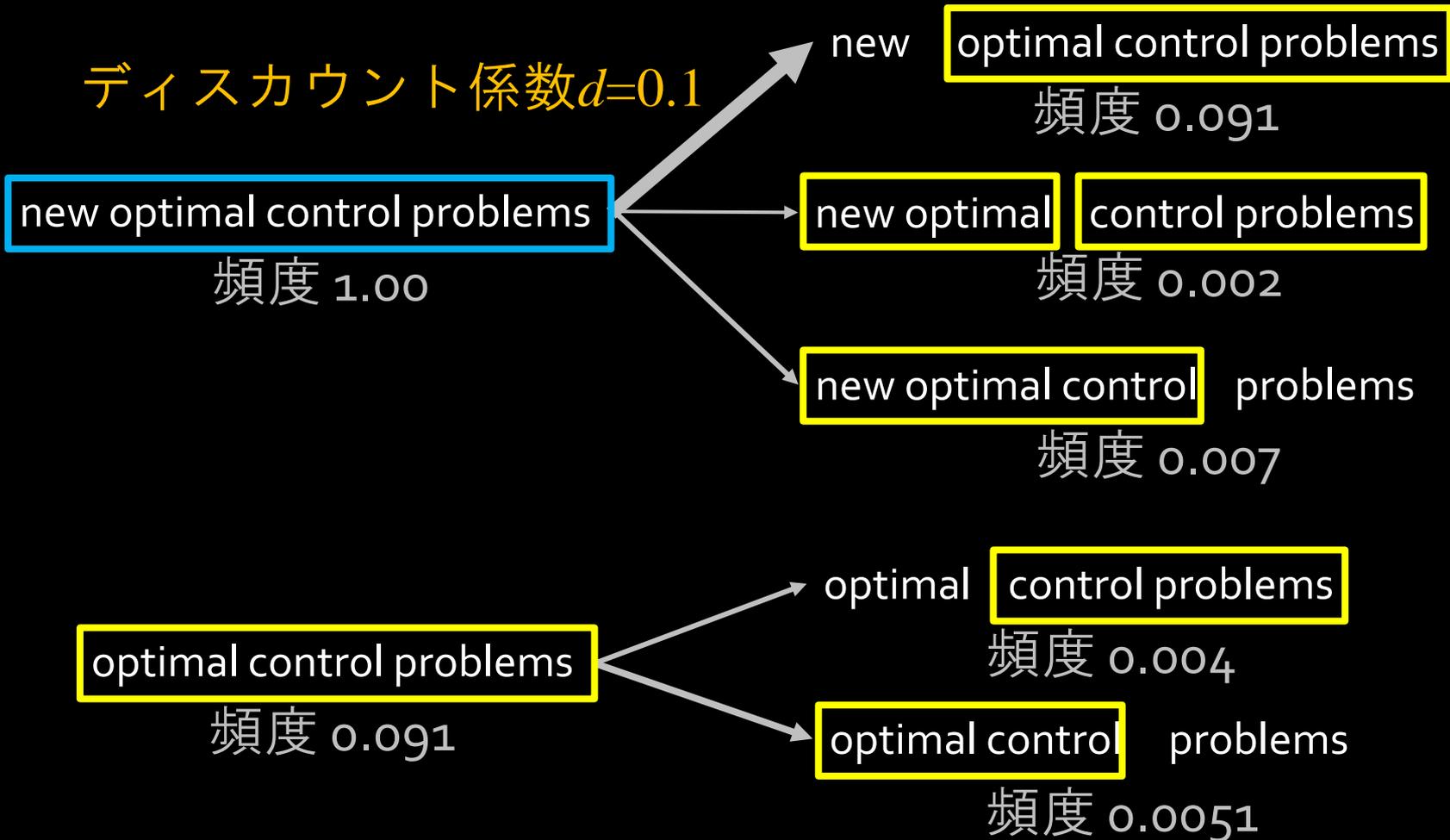


tf-term: 名詞句解析に基づく 部分名詞列への頻度配分



tf-term: 名詞句解析に基づく 部分名詞列への頻度配分

ディスカウント係数 $d=0.1$



名詞列のidf相当尺度?

$$\text{スコア}(\mathbf{w}) = \text{tfterm}(\mathbf{w}) \times \text{idf相当尺度}(\mathbf{w})$$

名詞列のidf相当尺度?

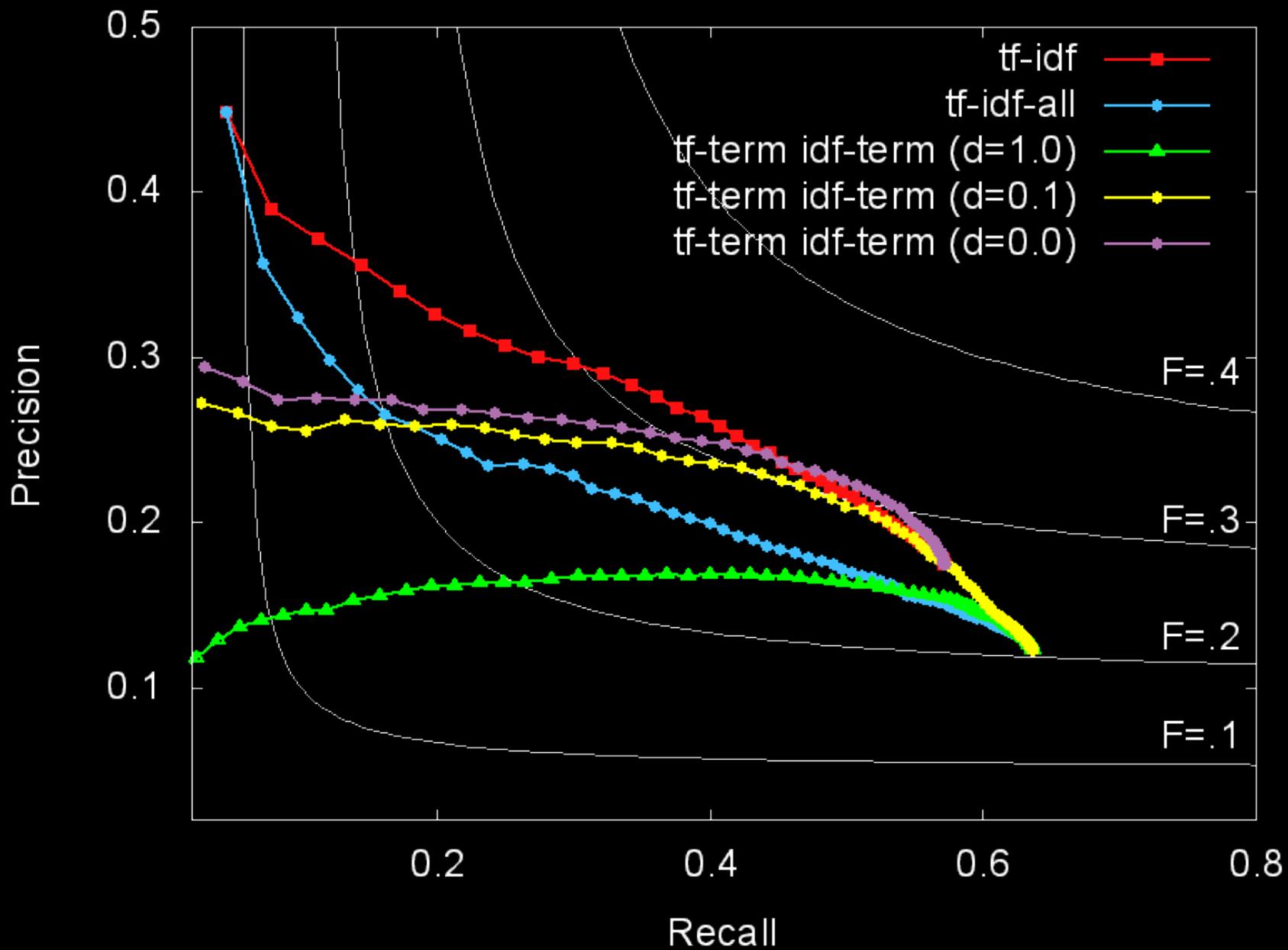
スコア(\mathbf{w}) = tfterm(\mathbf{w}) × idf相当尺度(\mathbf{w})

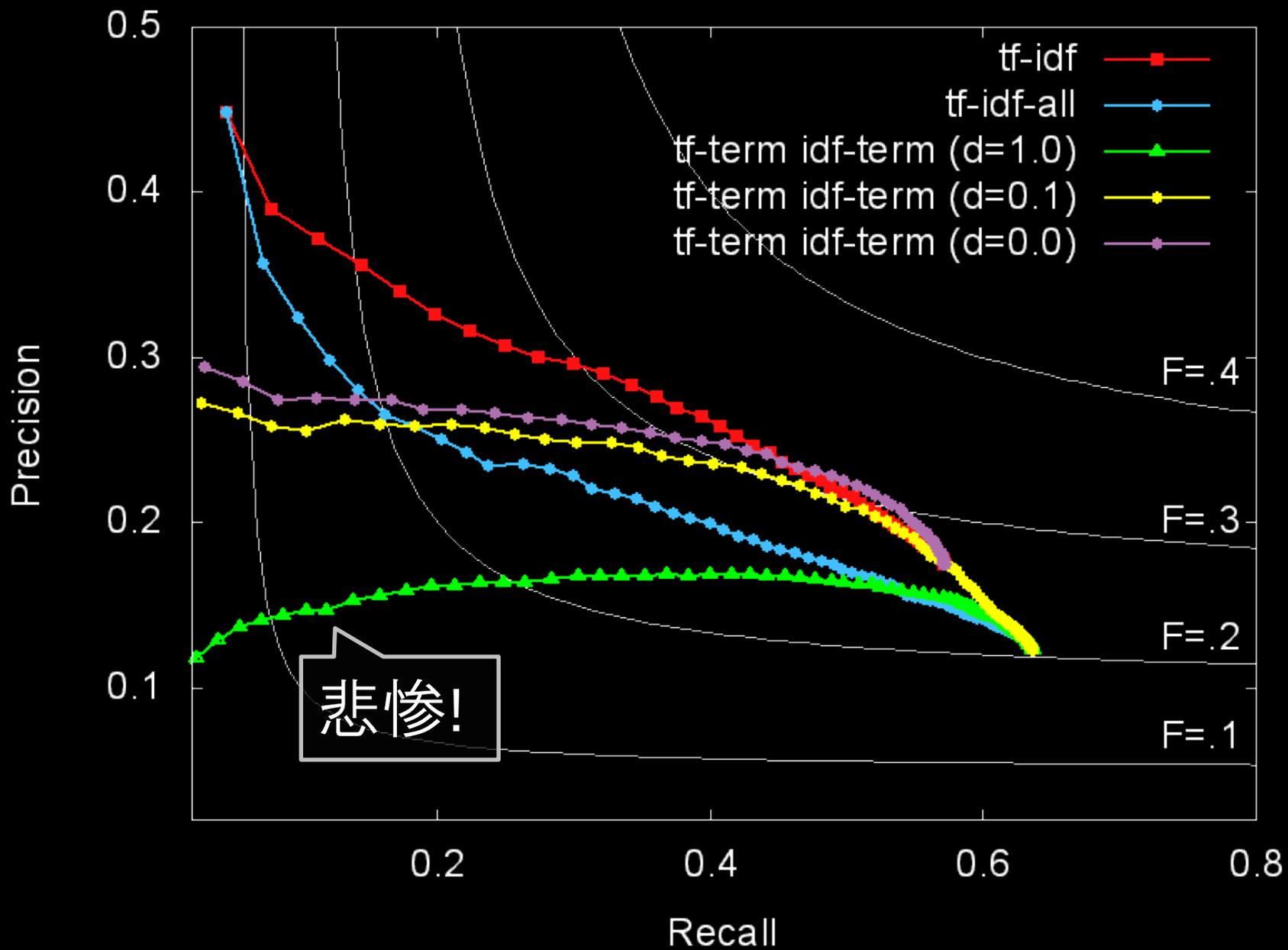
1. 名詞列のidfを直接算出 (tf-term idf-term)

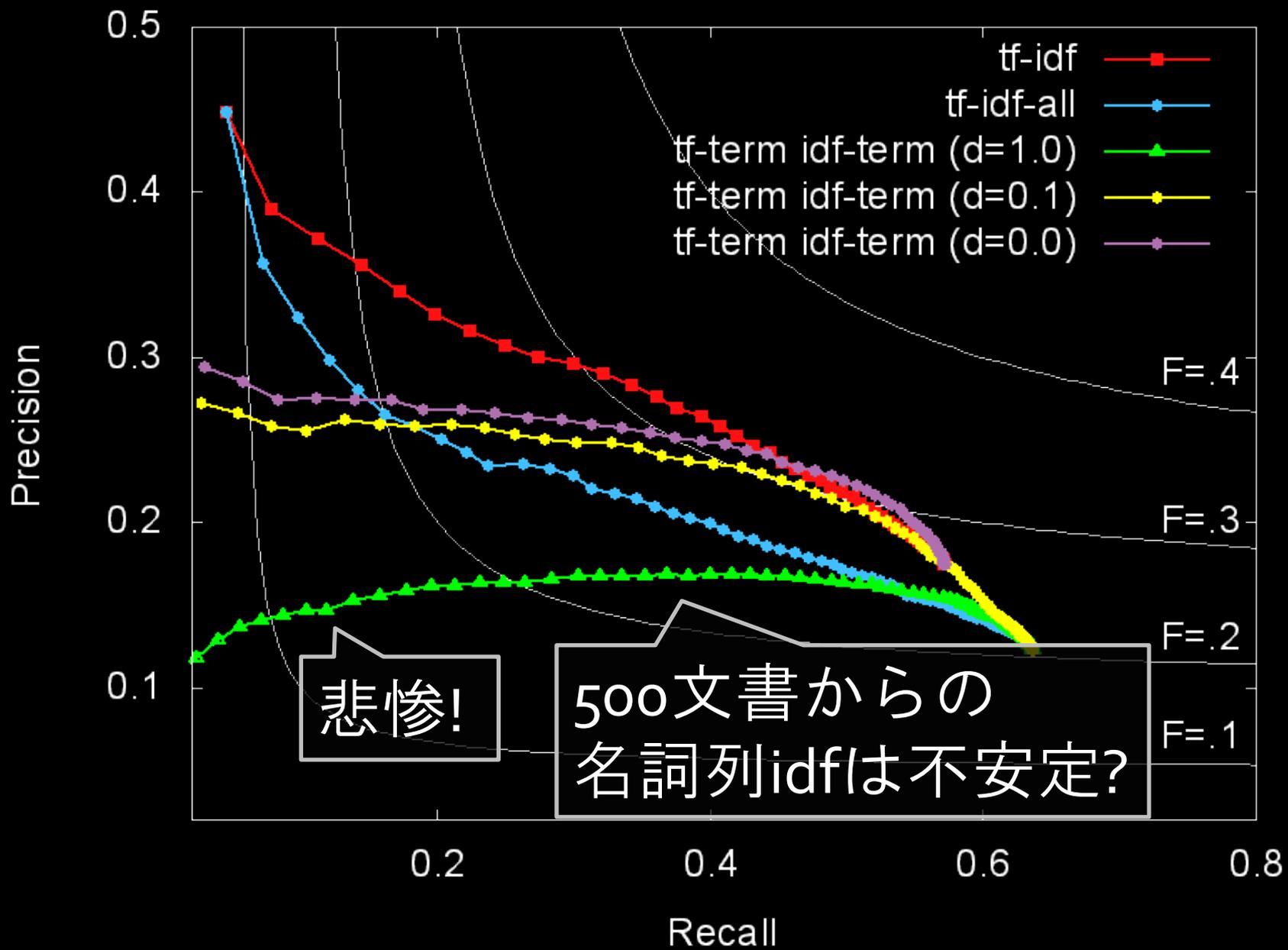
$$\text{idfterm}(\mathbf{w}) = \log(D / D_{\mathbf{w}})$$

$$D_{\mathbf{w}} = \left| \left\{ doc \mid \text{tfterm}_{doc}(\mathbf{w}) > 0 \right\} \right|$$

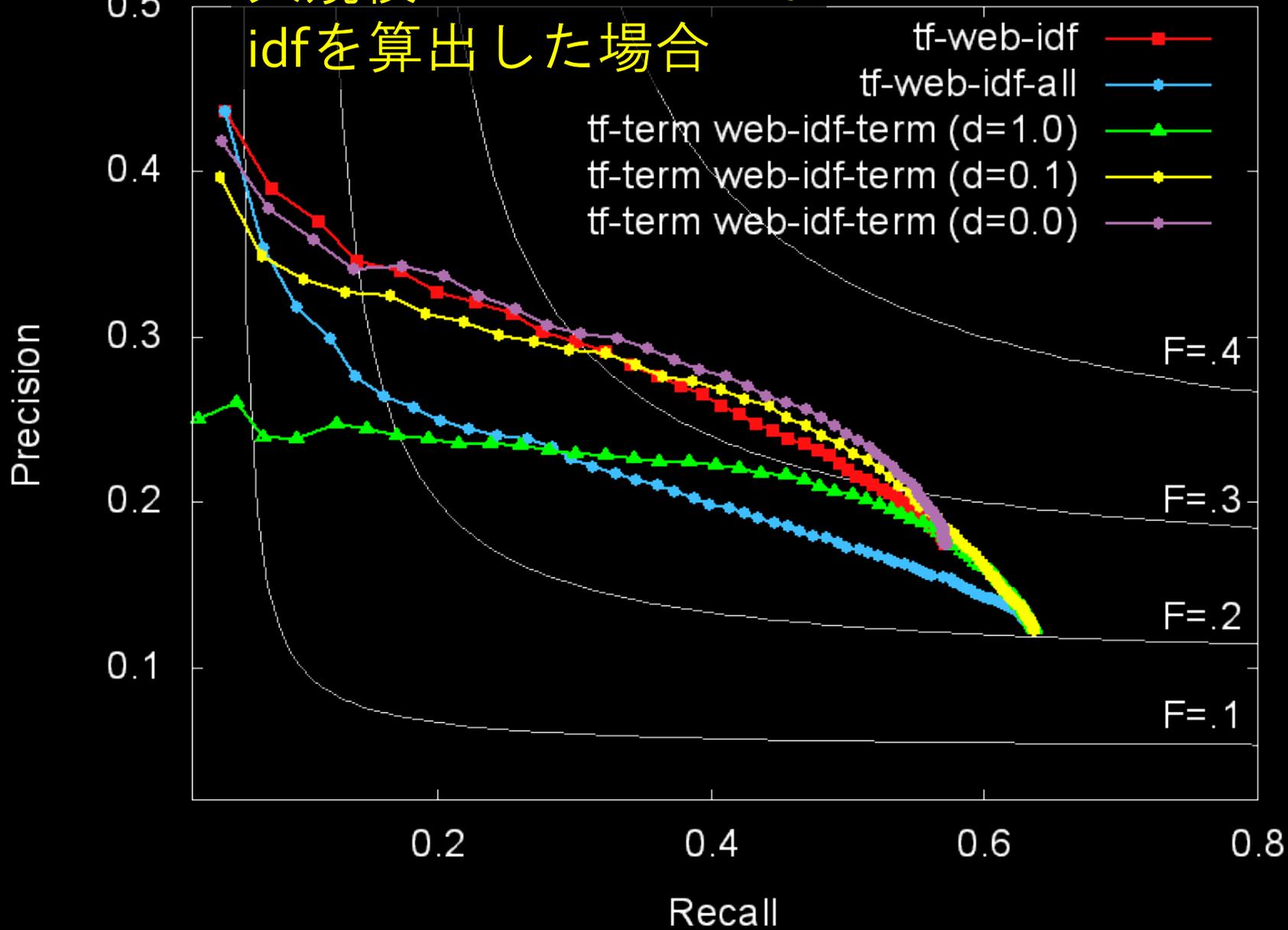
2. ...



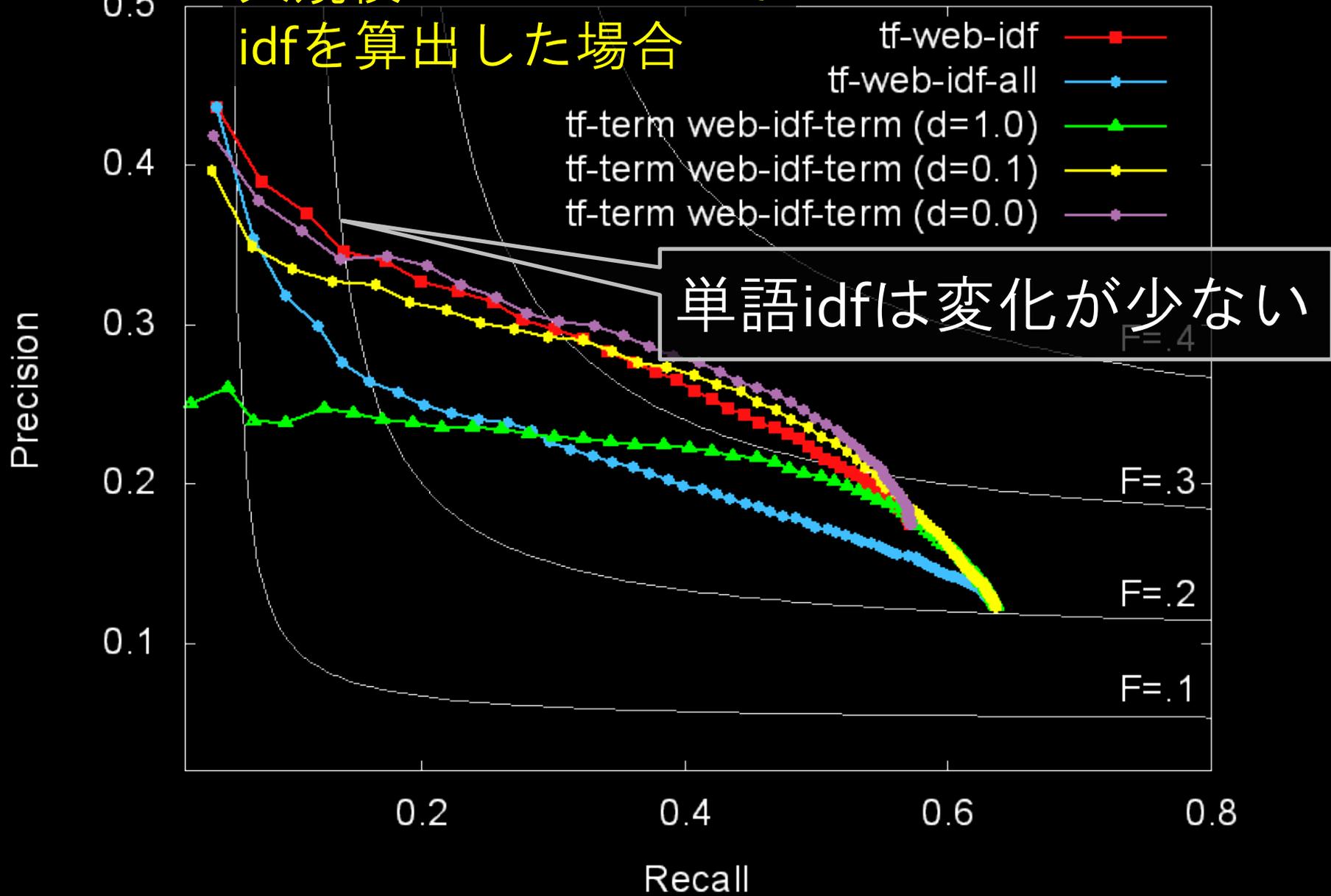




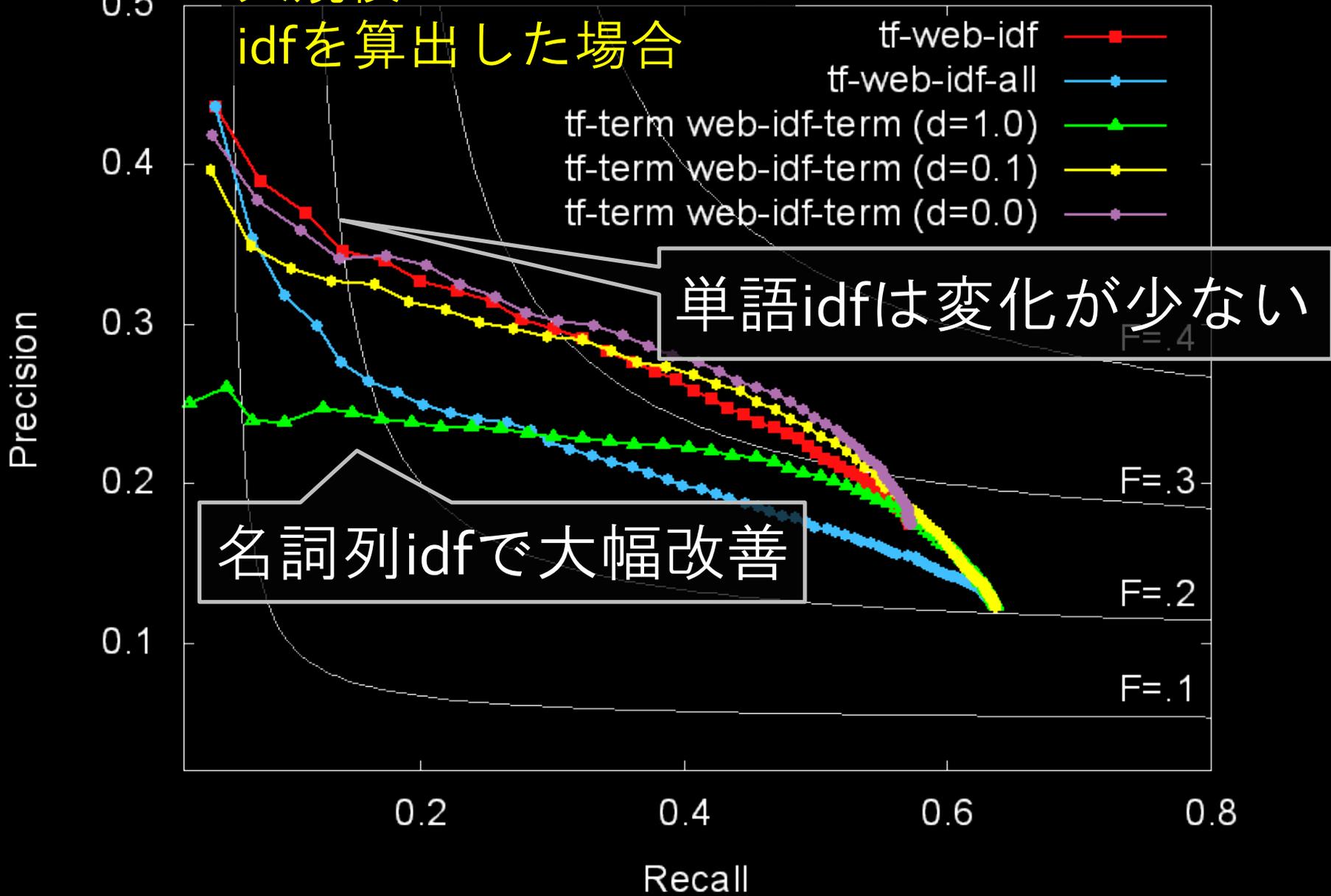
大規模webコーパスでidfを算出した場合



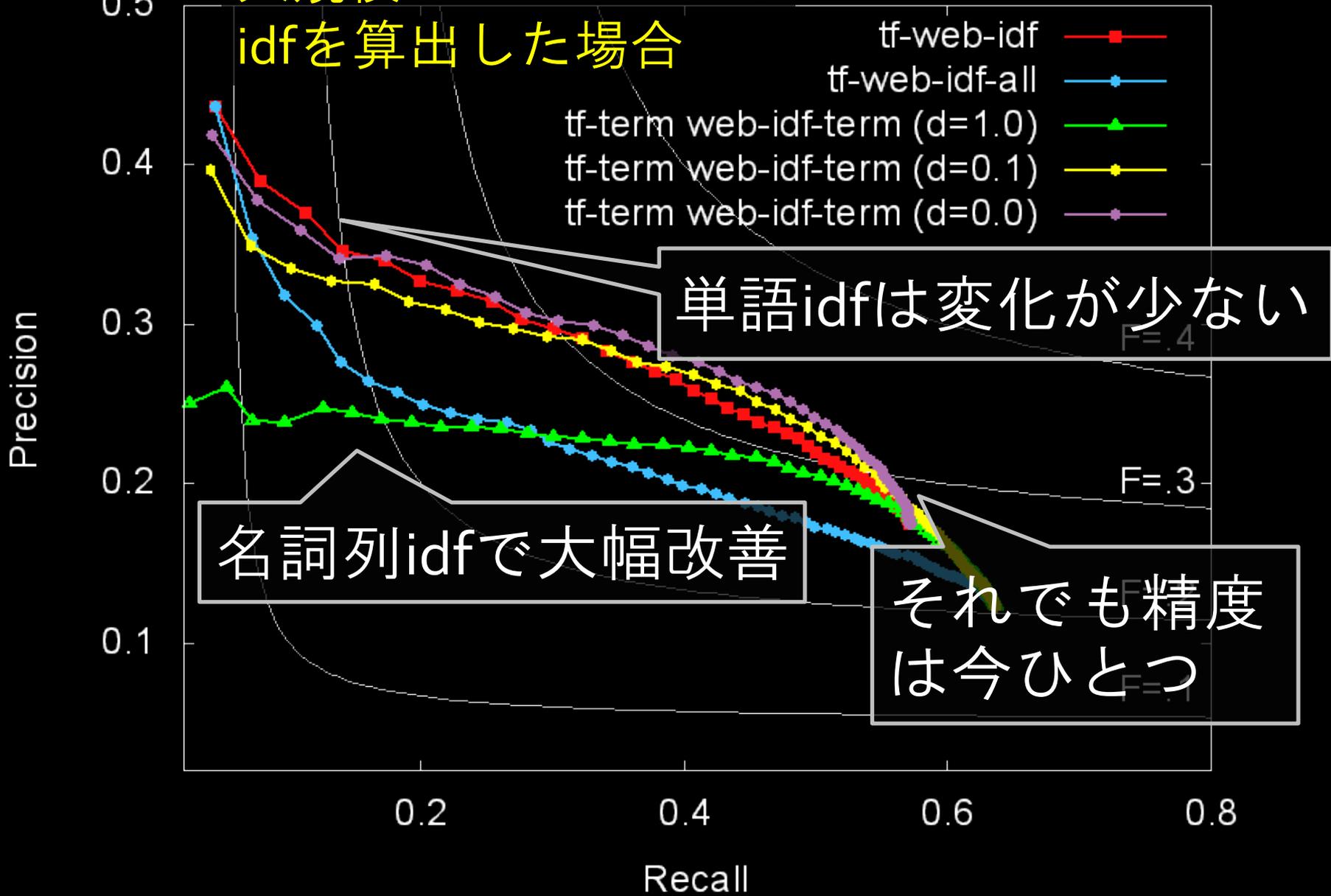
大規模webコーパスでidfを算出した場合



大規模webコーパスでidfを算出した場合



大規模webコーパスでidfを算出した場合



名詞列のidf相当尺度?

スコア(\mathbf{w}) = $\text{tfterm}(\mathbf{w}) \times \text{idf相当尺度}(\mathbf{w})$

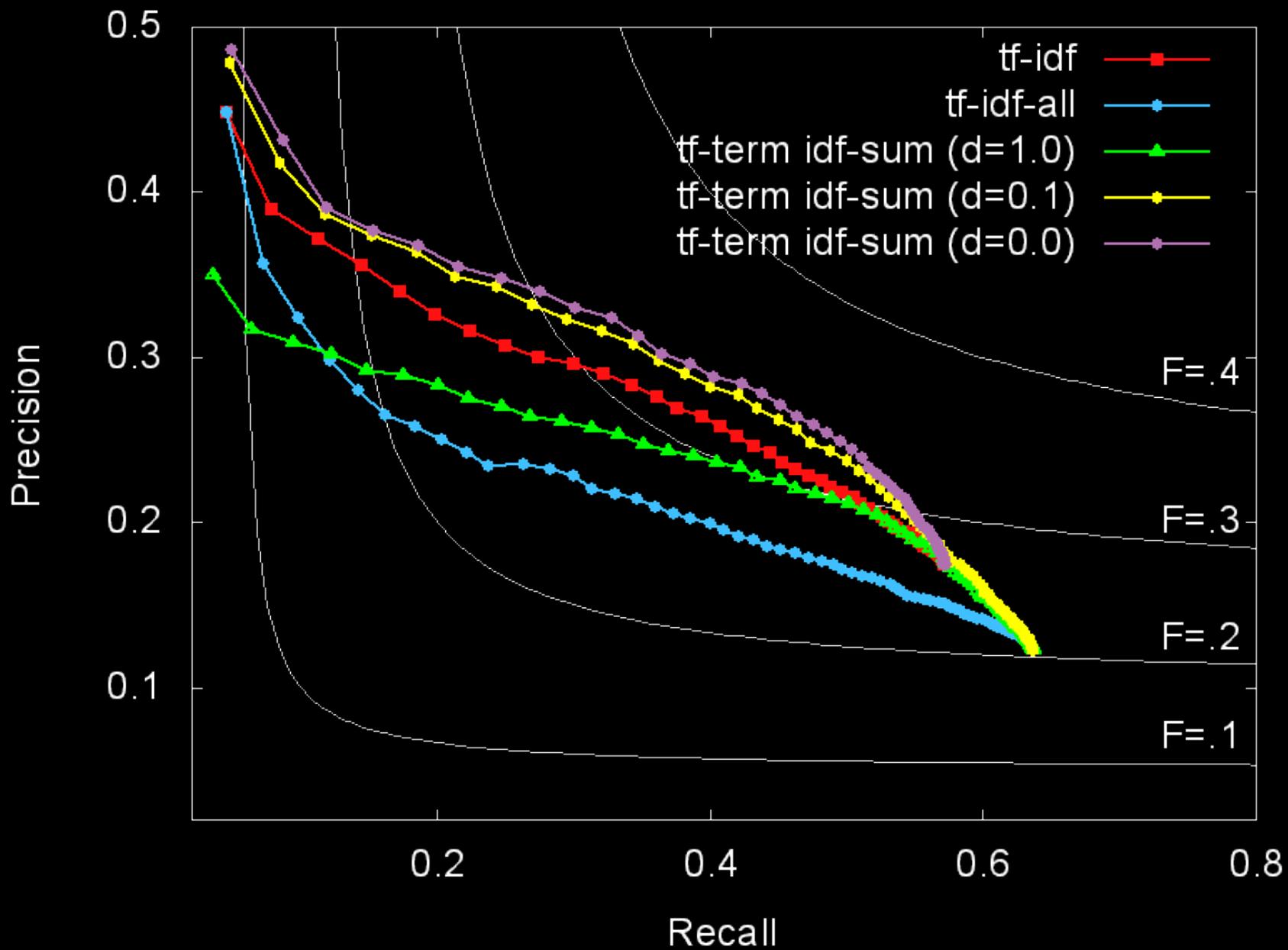
1. 名詞列のidfを直接算出 (tf-term **idf-term**)

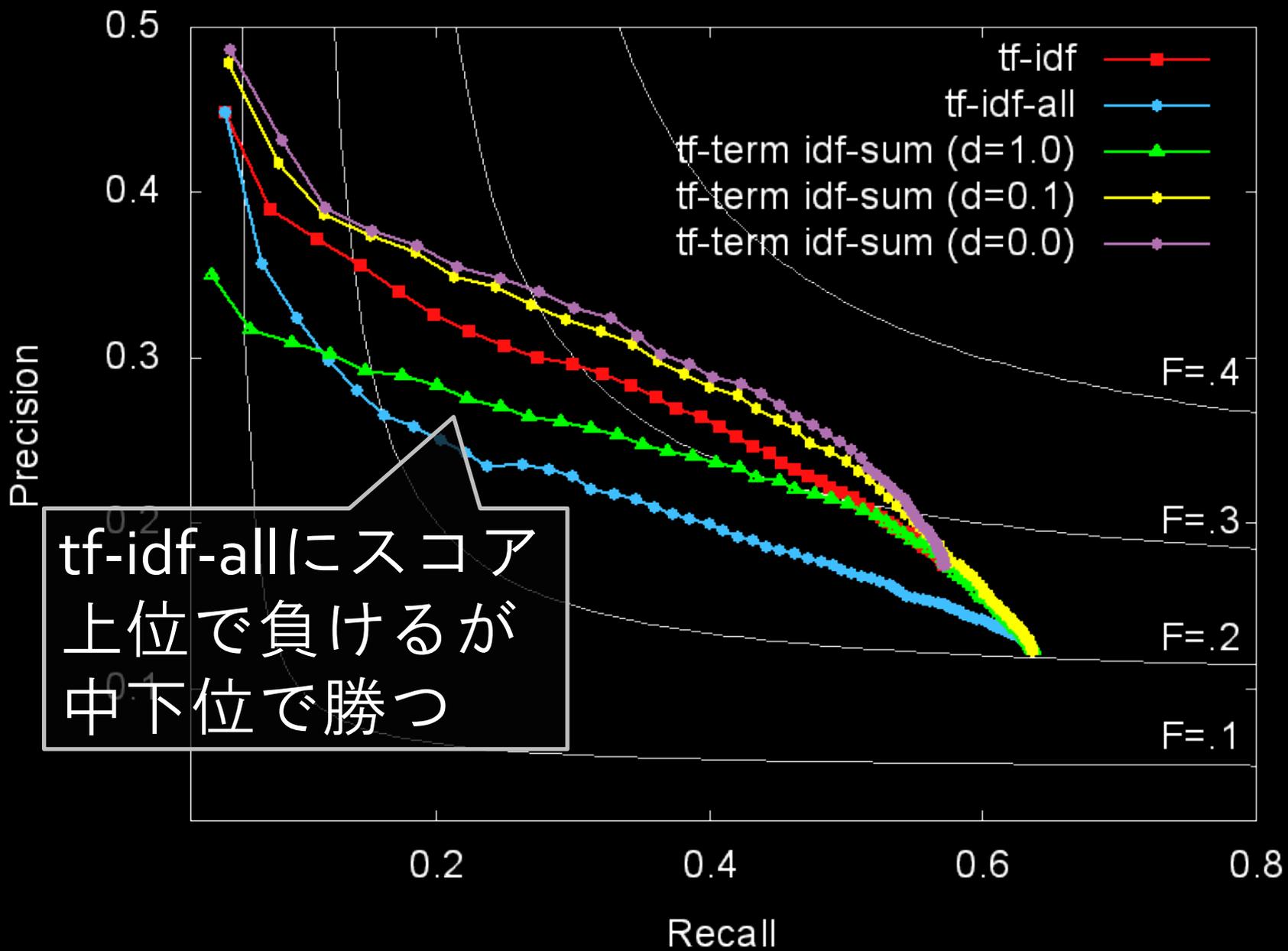
$$\text{idfterm}(\mathbf{w}) = \log(D / D_{\mathbf{w}})$$

$$D_{\mathbf{w}} = \left| \{ doc \mid \text{tfterm}_{doc}(\mathbf{w}) > 0 \} \right|$$

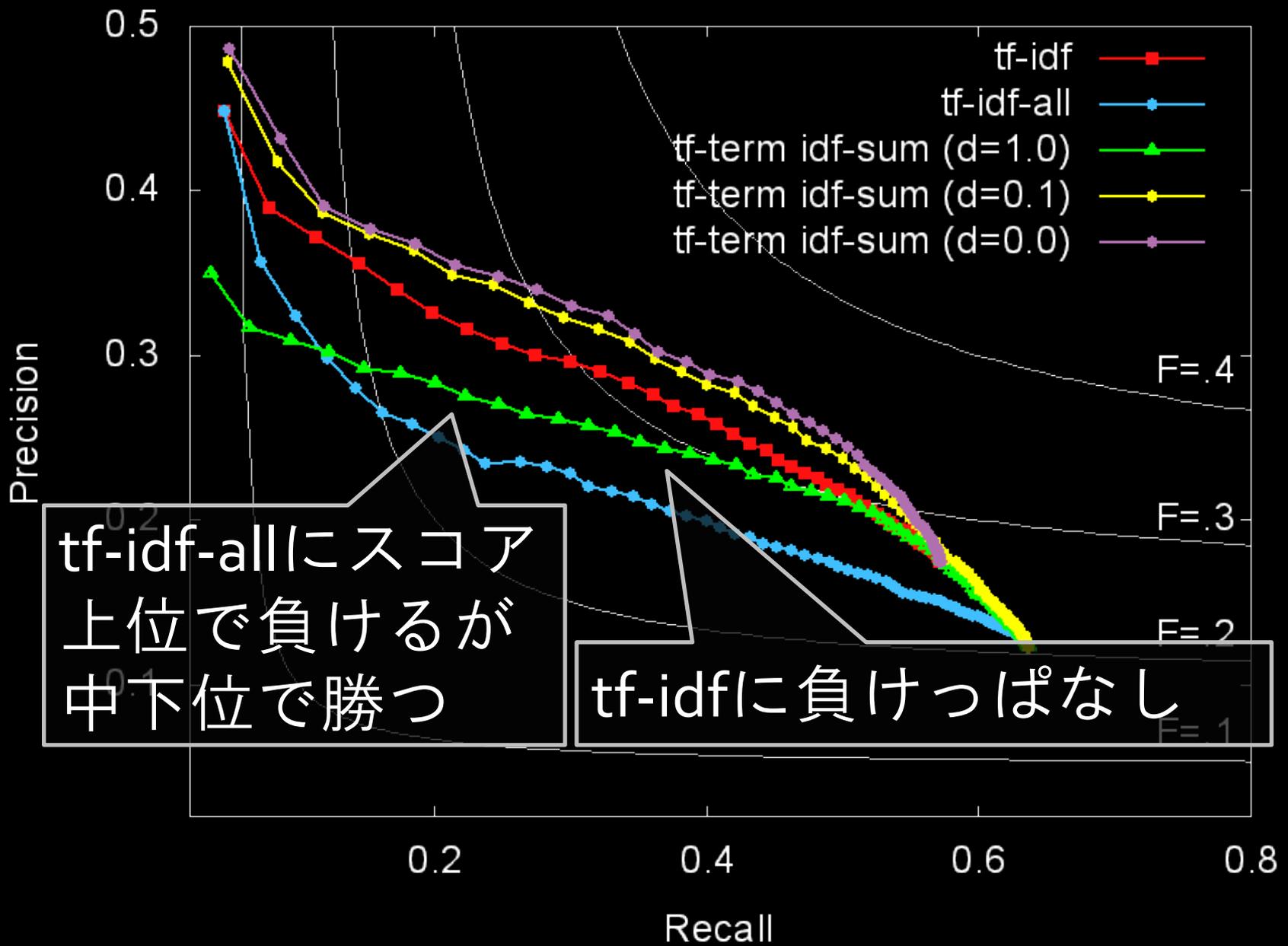
2. 単語idfの総和 (tf-term **idf-sum**)

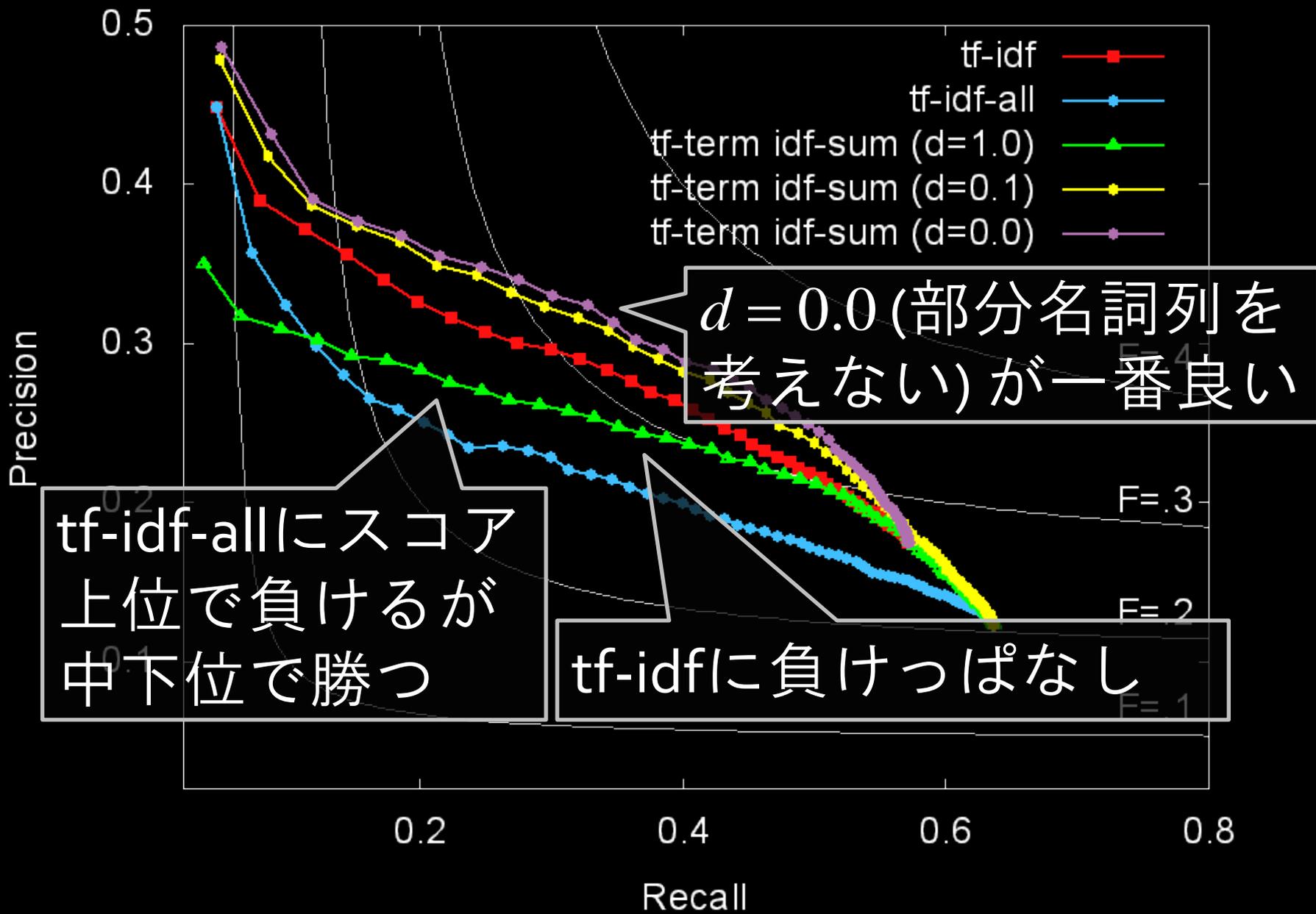
$$\text{idfsum}(\mathbf{w}) = \sum_i \log(D / D_{w_i})$$





tf-idf-allにスコア
 上位で負けるが
 中下位で勝つ





議論

- キーフレーズ抽出 \equiv 文書要約
 - optimal control problems をキーフレーズとして採用したら、optimal control は不要
 - 今回のデータセットには長い名詞列をキーフレーズとして採用する傾向が見られる

tf-idf-all

tf-term idf-sum (d=0.1)

tf-term idf-sum (d=1.0)

new optimal control problems
optimal control problems

new optimal control problems
optimal control
control

new optimal control problems
optimal control

☹️ new optimal control
optimal control

optimal control problems

optimal control problems
control

☹️ control problems
control
problems

☹️ new optimal control
☹️ control problems
problems

☹️ control problems
problems

...

...

...

まとめと今後の課題

- 教師なしキーフレーズ抽出ではtf-idf法がいまだ強い
- 名詞句の内部構造を考慮した名詞列tfを提案
- 名詞列のidf相当尺度で頑健なものとは?

