


Another Reason for Smoothing

Content words

Query = “the **algorithms** for **data mining**”

$P_{DML}(w d1):$	<i>0.04</i>	<i>0.001</i>	<i>0.02</i>	<i>0.002</i>	<i>0.003</i>
$P_{DML}(w d2):$	<i>0.02</i>	<i>0.001</i>	<i>0.01</i>	<i>0.003</i>	<i>0.004</i>

$p(\text{“algorithms”}|d1) = p(\text{“algorithm”}|d2)$
 $p(\text{“data”}|d1) < p(\text{“data”}|d2)$
 $p(\text{“mining”}|d1) < p(\text{“mining”}|d2)$


 Intuitively, d2 should have a higher score, but $p(q|d1) > p(q|d2) \dots$

So we should make $p(\text{“the”})$ and $p(\text{“for”})$ **less different** for all docs, and smoothing helps achieve this goal...

After smoothing with $p(w|d) = 0.1p_{DML}(w|d) + 0.9p(w|REF)$, $p(q|d1) < p(q|d2)!$

Query	= “the	algorithms	for	data	mining”
$P(w REF)$	0.2	0.00001	0.2	0.00001	0.00001
Smoothed $p(w d1):$	0.184	0.000109	0.182	0.000209	0.000309
Smoothed $p(w d2):$	0.182	0.000109	0.181	0.000309	0.000409