# DERIVATION

Weighted zone scores on document having L zones is calculated by

$$\sum_{i=1 \text{ to } L} \left( g_i * S_i \right)$$ where $g_i$ is weight of a zone and $S_i$ is Boolean score.

Let's follow step by step procedure, In order to learn weights $g_i$ for each zones in a document through machine learning.

1) In the set of training example there are N number of tuples and each tuple consist of query, document and human annotated relevance score on query and document.
   Let $r(q_j, d_j)$ be human annotated relevance function on query q and document d for $j^{th}$ tuple in training set.

2) Function $r(q_j, d_j)$ has domain = query 'q' and document 'd' while range 0 or 1. Where 0 indicates given document did not relevant to user query while 1 indicates given document is a perfectly relevant to user query.

3) Let all documents in training set have only two zones title and body.

4) If there are only two zone in documents then score computing function on document dj and query qj for $j^{th}$ tuple in training set be defined as:
   Score $(q_j, d_j) = (g * S_1) + ((1-g) * S_2)$ => Equation A

5) Let Boolean match function for title and body zones defined by $S_t(q_j, d_j)$ and $S_b(q_j, d_j)$ respectively.
   Then Equation A will became:
   Score $(q_j, d_j) = (g * S_t(q_j, d_j)) + ((1-g) * S_b(q_j, d_j))$ => Equation A1

6) As the name suggest Boolean match function has only 2 values either 0 or 1 and as we have two functions, their all possible combinations and effect of each combination on Score is shown in table below by using equation A1.
   **Table 1:**

   | $S_t(q_j, d_j)$ | $S_b(q_j, d_j)$ | Score $(q_j, d_j)$ |
   |---|---|---|
   | 0 | 0 | 0 |
   | 0 | 1 | 1-g |
   | 1 | 0 | g |
   | 1 | 1 | 1 |

7) Let error between human annotated relevance function $r(q_j, d_j)$ and machine calculated score$(q_j, d_j)$ for $j^{th}$ tuple in training set be defined as
   Error $(g, \text{Tuple } j) = (r(d_j, q_j) - \text{Score}(d_j, q_j))^2$ => Equation B

8) Let Sum of all errors be defined as

$$\sum_{j=1 \text{ to } N} \left( \text{Error}(g, \text{tuple } j) \right)$$ => Equation C

9) Now suppose

Symbol $n_{01r}$ denotes all example in training set for which St (qj, dj) =0 and Sb (qj, dj) =1 and human annotated score is relevant (that is 1).
Symbol $n_{01n}$ denotes all example in training set for which St (qj, dj) =0 and Sb (qj, dj) =1 and human annotated score is non-relevant (that is 0).

Symbol $n_{10r}$ denotes all example in training set for which St (qj, dj) =1 and Sb (qj, dj) =0 and human annotated score is relevant (that is 1),
Symbol $n_{10n}$ denotes all example in training set for which St (qj, dj) =1 and Sb (qj, dj) =0 and human annotated score is non-relevant (that is 0).

Symbol $n_{00r}$ denotes all example in training set for which St (qj, dj) =0 and Sb (qj, dj) =0 and human annotated score is relevant (that is 1),
Variable $n_{00n}$ denotes all example in training set for which St (qj, dj) =0 and Sb (qj, dj) =0 and human annotated score is non-relevant (that is 0).

Symbol $n_{11r}$ denotes all example in training set for which St (qj, dj) =1 and Sb (qj, dj) =1 and human annotated score is relevant (that is 1),
Symbol $n_{11n}$ denotes all example in training set for which St (qj, dj) =1 and Sb (qj, dj) =1 and human annotated score is non-relevant (that is 0).

10) Using all Equation A1 and Equation B we can compute Machine Scores and error values for each symbol.

**Table 2**

| VARIABLE | St (qj, dj) | Sb (qj, dj) | r(qj, dj) | SCORE Using equation A1 | Error Function value using equation B |
|---|---|---|---|---|---|
| $n_{01r}$ | 0 | 1 | 1 | (1-g) | $(1 - (1\text{-}g))^2$ |
| $n_{01n}$ | 0 | 1 | 0 | (1-g) | $(0 - (1\text{-}g))^2$ |
| $n_{10r}$ | 1 | 0 | 1 | g | $(1 - g)^2$ |
| $n_{10n}$ | 1 | 0 | 0 | g | $(0 - g)^2 = g^2$ |
| $n_{00r}$ | 0 | 0 | 1 | 0 | $(1 - 0)^2 = 1$ |
| $n_{00n}$ | 0 | 0 | 0 | 0 | $(0 - 0)^2 = 0$ |
| $n_{11r}$ | 1 | 1 | 1 | g+(1-g) = 1 | $(1 - 1)^2 = 0$ |
| $n_{11n}$ | 1 | 1 | 0 | g+(1-g) =1 | $(0 - 1)^2 = 1$ |

11) In order to get sum of all error values in Table 2 we will use Equation C

Let
$(1 - (1\text{-}g))^2$ * $\mathbf{n_{01r}}$ = $\mathbf{(g^2 * n_{01r})}$ ) be the sum of all training set examples for which St (qj, dj) =0 and Sb (qj, dj) =1 and human annotated score is relevant (that is 1).
$(0 - (1\text{-}g))^2$ * $\mathbf{n_{01n}}$ = = $\mathbf{((1\text{-}g)^2 * n_{01n}}$ ) be the sum of all training set examples for which St (qj, dj) =0 and Sb (qj, dj) =1 and human annotated score is non-relevant (that is 0).

$(1-g)^2$ * $n_{10r}$ be the sum of all training set examples for which St (qj, dj) =1 and Sb (qj, dj) =0 and human annotated score is relevant (that is 1).

$g^2$ * $n_{10n}$ be the sum of all training set examples for which St (qj, dj) =1 and Sb (qj, dj) =0 and human annotated score is non-relevant (that is 0).

$(1-0)^2$ *$n_{00r}$ = $n_{00r}$, be the sum of all training set examples for which St (qj, dj) =0 and Sb (qj, dj)=0 and human annotated score is relevant (that is 1).

$(0-0)^2$ * $n_{00n}$ = 0, be the sum of all training set examples for which St (qj, dj) =0 and Sb (qj, dj)=0 and human annotated score is non-relevant (that is 0).

$(1-1)^2$ * $n_{11r}$ = 0, be the sum of all training set examples for which St (qj, dj) =1 and Sb (qj, dj)=1 and human annotated score is relevant (that is 1).

$(0-1)^2$*$n_{11n}$=$n_{11n}$, be the sum of all training set examples for which St (qj, dj) =1 and Sb (qj, dj)=1 and human annotated score is non-relevant (that is 0).

By summing all these we get

= $(g^2 * n_{01r})$ + $((1-g)^2 * n_{01n})$ + $((1-g)^2 * n_{10r})$ + $(g^2 * n_{10n})$ + $n_{00r}$ + $n_{11n}$

Taking common from similarly highlighted terms we get.
$(1-g)^2 (n_{10r} + n_{01n}) + g^2(n_{10n} + n_{01r}) + n_{00r} + n_{11n}$ => **Equation D defining total error sum in training set**

12) Now in order to get optimal weight value of total error in the training set must be minimum and as we can see in Equation D total sum is only dependent on variable g.

As derivative of a function is 0 at its minima so taking derivative of equation D can yield optimal value of g

All step done below is for taking derivative of equation D.

Let
a = $n_{10r}$
b = $n_{10n}$
c = $n_{01r}$
d = $n_{01n}$
e = $n_{00r}$
f = $n_{11n}$

Then derivative operation on equation D can be performed as
$(1-g)^2 (n_{10r} + n_{01n}) + g^2(n_{10n} + n_{01r}) + n_{00r} + n_{11n} = 0$
$(1-g)^2 (a+d) + g^2(b+c) + e+f = 0$ => **Equation D1**
**Let**
$f(g) = (1-g)^2 (d+a) + g^2(b+c) + e+f$

**then equation D1 becomes**
**f(g) = 0 => Equation D1**


**Solving L.H.S**
**If f(g) = (1-g)$^2$ (d + a) + g$^2$(b+ c) + e+ f**
**Let f(x) = f(g)**
**Then f(x) = (1-x)$^2$ (d + a) + x$^2$(b+ c) + e+ f**
**Performing derivative on f(x)**

$$\frac{d}{dx}\left[(c+b)\,x^2 + (d+a)\,(1-x)^2 + f + e\right]$$

$$= (c+b)\cdot\frac{d}{dx}\left[x^2\right] + (d+a)\cdot\frac{d}{dx}\left[(1-x)^2\right] + \frac{d}{dx}[f] + \frac{d}{dx}[e]$$

$$= 2\,(1-x)\cdot\frac{d}{dx}[1-x]\cdot(d+a) + 2x\,(c+b) + 0 + 0$$

$$= 2\,(c+b)\,x + 2\left(\frac{d}{dx}[1] - \frac{d}{dx}[x]\right)(d+a)\,(1-x)$$

$$= 2\,(c+b)\,x + 2\,(0-1)\,(d+a)\,(1-x)$$

$$= 2\,(c+b)\,x - 2\,(d+a)\,(1-x)$$

**Simplify:**

$$2\,((d+c+b+a)\,x - d - a)$$


**Let**

F(x) = 2(d+c+b+a)x-d-a => equation z

As we above mentioned

F(x) = f(g)

Then equation z becomes

F(g) = **2(d+c+b+a)g-d-a => equation z1**


Putting F(g) from equation z1 in Equation D1

**2(d+c+b+a)g-d-a = 0**
**(d+c+b+a)g-d-a = 0**
**(d+c+b+a)g = d + a**
**g = (d + a) / (d+c+b+a) => Equation D1**

By substituting values in Equation D1 of a, b c, d, it becomes

$$g = (n_{10r} + n_{01n}) / (n_{01n} + n_{01r} + n_{10n} + n_{10r}) \Rightarrow \textbf{Final Equation.}$$