

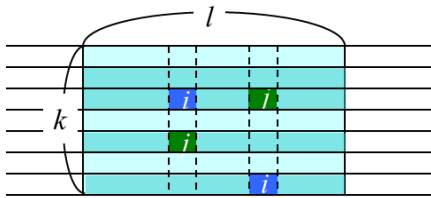
W3.3: BLOSUM の求め方

同一ファミリータンパク質のギャップなしでアラインされた領域（ブロック）に対し、アミノ酸の置換の頻度を調べて作成する。

```

-----N V Y H D G A C P E V K P V D N F D W S N Y H G K W W E V A K Y P N S V E K Y G K C G W A E Y T P E G -----
-----G D I F Y P G Y C P D V K P V N D F D L S A F A G A W H E I A K L P L E N E N Q G K C T I A E Y K Y D G -----
-----E R D C R V S S F R V K E N F D K A R F A G T W Y A M A K K D P E G L F L Q D N I V A E F S V D E N G H I
C L L L A L A L T C G A Q A L I V T Q T M K G L D I Q K V A G T W Y S L A M A A S D I S - L L D A Q S A P L R V Y V E E L
L L L C L G L T L V C V H A E E A S S T G R N F N V E K I N G E W H T I I L A S D K R E K I E D N G N F R L F L E Q -----
  
```

ブロック



ブロックの1つの図。実際には全ブロックをまとめて計算。

全ブロックの内容から、アミノ酸 i の出現確率 q_i を求める。

各ブロックから $x\%$ 以上一致している配列群をクラスタリングしてまとめる。この際、クラスタリングされた配列を1本の配列として扱う。

その中で、アミノ酸 i, j が同じ位置（カラム）に現れる確率を求める。

そのため、可能性のあるアミノ酸ペアの数 N をカウント。

長さ l のブロックで、 k 本の配列が取り出されたとき、アミノ酸のペアの数は、 $l \times k C_2$

すべての b 個のブロックに対して $N = \sum_b l_b \times k_b C_2$

その中で、アミノ酸 i, j が同じ位置に現れる個数 n_{ij} を求める $\rightarrow p_{ij} = n_{ij} / N$

クラスタリングされた m 本の配列では、各アミノ酸は $1/m$ の重みをもつものとして計算。

クラスタ $k=4$ 各位置 で6通り		NN: $1 \times 1 + 1/2 \times 2 = 2$
		GN: $1 \times 2 + 1/2 \times 1 = 2.5$
		DN: $1/2 \times 2 = 1$
		DG: $1/2 \times 1 = 0.5$

BLOSUM $_x$ マトリックスのスコアは以下の式で与えられる。

$$s(i, j) = \log \left(\frac{p_{ij}}{q_i q_j} \right)$$

ここで、 $p_{ij} = n_{ij} / N$