



Βιοπληροφορική

Ενότητα 5: Πίνακες αντικατάστασης BLOSUM
και οπτική σύγκριση αλληλουχιών

Αν. καθηγητής Αγγελίδης Παντελής

e-mail: paggelidis@uowm.gr

ΕΕΔΙΠ Μπέλλου Σοφία

e-mail: sbellou@uowm.gr

Τμήμα Μηχανικών Πληροφορικής και Τηλεπικοινωνιών



Άδειες Χρήσης

- Το παρόν εκπαιδευτικό υλικό υπόκειται σε άδειες χρήσης Creative Commons.
- Για εκπαιδευτικό υλικό, όπως εικόνες, που υπόκειται σε άλλου τύπου άδειας χρήσης, η άδεια χρήσης αναφέρεται ρητώς.



Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στα πλαίσια του εκπαιδευτικού έργου του διδάσκοντα.
- Το έργο «**Ανοικτά Ψηφιακά Μαθήματα στο Πανεπιστήμιο Δυτικής Μακεδονίας**» έχει χρηματοδοτήσει μόνο τη αναδιαμόρφωση του εκπαιδευτικού υλικού.
- Το έργο υλοποιείται στο πλαίσιο του Επιχειρησιακού Προγράμματος «Εκπαίδευση και Δια Βίου Μάθηση» και συγχρηματοδοτείται από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο) και από εθνικούς πόρους.



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης

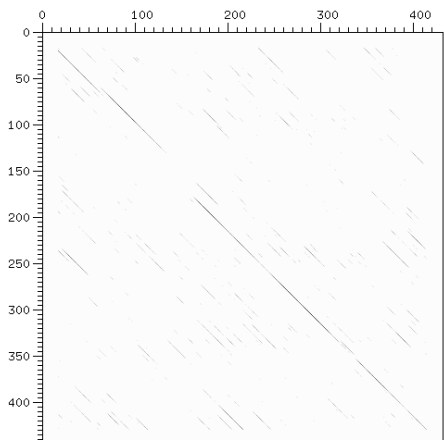


Σκοπός του μαθήματος

- Πίνακες αντικατάστασης PAM.
- The log odds form (the mutation data matrix) of PAM250.
- Πίνακες αντικατάστασης BLOSUM.
- Nucleic acid PAM Scoring Matrices.
- Οι ποινές για τα κενά.
- Στατιστική σημαντικότητα της στοίχισης αλληλουχιών.
- P-value, E-value, Z-score.
- Οπτική σύγκριση αλληλουχιών.
- Διαγράμματα πινάκων σημείων – Dot plots.
- Θόρυβος» στα dot plots.



Πίνακες αντικατάστασης BLOSUM και οπτική σύγκριση αλληλουχιών



Σοφία Μπέλλου, sbellou@uowm.gr

| | C | S | T | P | A | G | N | D | E | Q | H | R | K | M | I | L | V | F | Y | W | |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|----|----|---|
| C | 12 | | | | | | | | | | | | | | | | | | | | C |
| S | 0 | 2 | | | | | | | | | | | | | | | | | | | S |
| T | -2 | 1 | 3 | | | | | | | | | | | | | | | | | | T |
| P | -3 | 1 | 0 | 6 | | | | | | | | | | | | | | | | | P |
| A | -2 | 1 | 1 | 1 | 2 | | | | | | | | | | | | | | | | A |
| G | -3 | 1 | 0 | -1 | 1 | 5 | | | | | | | | | | | | | | | G |
| N | -4 | 1 | 0 | -1 | 0 | 0 | 2 | | | | | | | | | | | | | | N |
| D | -5 | 0 | 0 | -1 | 0 | 1 | 2 | 4 | | | | | | | | | | | | | D |
| E | -5 | 0 | 0 | -1 | 0 | 0 | 1 | 3 | 4 | | | | | | | | | | | | E |
| Q | -5 | -1 | -1 | 0 | 0 | -1 | 1 | 2 | 2 | 4 | | | | | | | | | | | Q |
| H | -3 | -1 | -1 | 0 | -1 | -2 | 2 | 1 | 1 | 3 | 6 | | | | | | | | | | H |
| R | -4 | 0 | -1 | 0 | -2 | -3 | 0 | -1 | -1 | 1 | 2 | 6 | | | | | | | | | R |
| K | -5 | 0 | 0 | -1 | -1 | -2 | 1 | 0 | 0 | 1 | 0 | 3 | 5 | | | | | | | | K |
| M | -5 | -2 | -1 | -2 | -1 | -3 | -2 | -3 | -2 | -1 | -2 | 0 | 0 | 6 | | | | | | | M |
| I | -2 | -1 | 0 | -2 | -1 | -3 | -2 | -2 | -2 | -2 | -2 | -2 | -2 | 2 | 5 | | | | | | I |
| L | -6 | -3 | -2 | -3 | -2 | -4 | -3 | -4 | -3 | -2 | -2 | -3 | -3 | 4 | 2 | 6 | | | | | L |
| V | -2 | -1 | 0 | -1 | 0 | -1 | -2 | -2 | -2 | -2 | -2 | -2 | -2 | 2 | 4 | 2 | 4 | | | | V |
| F | -4 | -3 | -3 | -5 | -4 | -5 | -4 | -6 | -5 | -5 | -2 | -4 | -5 | 0 | 1 | 2 | -1 | 9 | | | F |
| Y | 0 | -3 | -3 | -5 | -3 | -5 | -2 | -4 | -4 | -4 | 0 | -4 | -4 | -2 | -1 | -1 | -2 | 7 | 10 | | Y |
| W | -8 | -2 | -5 | -6 | -6 | -7 | -4 | -7 | -7 | -5 | -3 | 2 | -3 | -4 | -5 | -2 | -6 | 0 | 0 | 17 | W |



Γιατί σύγκριση αλληλουχιών;;;

- Συμπεράσματα για λειτουργία ενός γονιδίου (μίας πρωτεΐνης).
 - Όταν δύο γονίδια έχουν μεγάλο ποσοστό ομοιότητας, τότε πιθανά κωδικοποιούν πρωτεΐνες με παρόμοια λειτουργία.
- Συμπεράσματα για σημαντικές περιοχές στην αλληλουχία γονιδίων/πρωτεϊνών.
 - Όταν πρωτεΐνες με συγκεκριμένο χαρακτηριστικό μοιράζονται κοινή περιοχή, τότε αυτή η κοινή περιοχή είναι πιθανά υπεύθυνη για το συγκεκριμένο χαρακτηριστικό.
- Συμπεράσματα για την εξελεγκτική απόσταση μεταξύ 3 ειδών.
 - Όταν η αλληλουχία μίας πρωτεΐνης είναι σχεδόν η ίδια μεταξύ 2 ειδών (ποντίκι και αρουραίος), τότε τα είδη αυτά είναι «κοντά».



Σύγκριση πρωτεϊνών - Εξαγωγή συμπερασμάτων για τη λειτουργία (1/4)

- **Protein A expressed in colon cancer:**

```
TCGCTGCGAAGGACATTTGGGCTGTGTGTGCGACGCGGGTTCGGAGGGGCGAGTCGGGGGAACC  
GCGAAGAAGCCGAGGAGCCCCGGAGCCCCGCGTGACGCTCCTCTCTCAGTCCAAAAGCGGCTTT  
TGGTTCGGCGCAGAGAGACCCGGGGGTCTTCAGGACAGCGATTGTCATTGCTGAAGCTTTTCCT  
CGAAAAGCGCCGCCCTGCCCTTGGCCCCGAGAA
```

- **Protein B expressed in colon cancer :**

```
CAGACAAAGAGCACCGCAGGGCCGATCACGCTGGGGGCGCTGAGGCCGGCCATGGTCATGGA  
AGTGGGCACCCTGGACGCTGGAGGCCTGCGGGGCGCTGCTGGGGGAGCGAGCGGCGCAATGCC  
TGCTGCTGGACTGCCGCTCCTTCTTCGCTTTCAACGCCGGCCACATCGCCGGCTCTGTCAACGT  
GCGCTTCAGCACCATCGTGCGGGCGCCGGGCCAAGGGCGCCATGGGCCTGGAGCACATCGTGC  
CCAACGCCGAGCTCCGCGGTCAGGACAGCGATTGTCATTGCTGA
```

- **Protein C expressed in colon cancer:**

```
CCGCCTGCTGGCCGGCGCCTACCACGCCGTGGTGTGCTGGACGAGCGCAGCGCCGCCCTGG  
ACGGCGCCAAGCGCGACGGCACCCCTGGCCCTGGCGGGCCGGCGCGCTCTGCCGCGAGGCGCG  
CGCCGCGCAAGTCTTCTTCCTCAAAGGAGGATACGAAGCGTTTTTCGGCTTCCTGCCCGGAGCTC  
AGGACAGCGATTGTCATTGCTGATGTGCAGCAAACAGTCGACCCCCATGGGGCTCAGCCTTCCC  
CTGAGTACTAGCGTCCCTGACAGCGCGGAATCTGGGTGCAGTT
```



Σύγκριση πρωτεϊνών - Εξαγωγή συμπερασμάτων για τη λειτουργία (2/4)

- **Protein D expressed in colon cancer:**

```
TAACTGCCTTGATCAACGTCTCAGCCAATTGTCCCAACCATTTTGAGGGTCACTACCAGTACAAG  
AGCATCCCTGTGGAGGACAACCACAAGGCAGACATCAGCTCCTGGTTCAACGAGGCCATTGACT  
TCATAGACTCCATCAAGAATGCTGGAGGAAGGGTGTGGTCCACTGCCAGGCAGGCATTTCCCG  
GTCAGCCACCATCTGCCTTGCTTACCTTATGAGGACTAATCGAGTCATCAGGACAGCGATTGTCA  
TTGCTGAAGCTGGACGAGGCCTTTGAGTTTGTGAAGCAGAGGC
```

- **Protein E expressed in colon cancer:**

```
GAAGCATCATCTCTCCCAACTTCAGCTTCATGGGCCAGCTGCTGCAGTTTGAGTCCCAGGTGCTG  
GCTCCGCACTGTTCCGGCAGAGGCTGGGAGCCCCGCCATGGCTGTGCTCGACCGAGGCACCTCC  
ACCACCACCGTGTTCAACTTATCAGGACAGCGATTGTGATTGCTGACCCCGTCTCCATCCCTGTC  
CACTCCACGAACAGTGCGCTGAGCTACCTTCAGAGCCCCATTACGACCTCTCCCAGCTGCTGAA  
AGGCCACGGGAGGTGAGGCTCTTCACATCCCATTGGGACTC
```

- **Protein F with unknown function:**

```
CATGCTCCTTGAGAGGAGAAATGCAATAACTCTGGGAGGGGCTCAGGACAGCGATTGTCATTGC  
TGATCGAGAGGGCTGGTCCTTATTTATTTAACTTCACCCGAGTTCCTCTGGGTTTCTAAGCAGTTA  
TGGTGATGACTTAGCGTCAAGACATTTGCTGAACTCAGCACATTCGGGACCAATATATAGTGGGT  
ACATCAAGTCCATCTGACAAAATGGGGCAGAAGAGAAAGGACTCAGTGTGTGATCCGGTTTCTTT  
TTGCTCGCCCCTGTTTTTTGTAGAATCTCTTCATGCTTGACATACCTACCAGTATTATTCCCGACG  
ACACATATACATATGAGAATATACCTTATTTATTTTTGTGTAGGTGTCTGCCTTCACAAATGTCATT  
GTCTACTCCTAGAAGAACCAAATACCTCAATTTTTGTTTTTGAGTACTGTACTATCCTGTAAATATA  
TCTTAAGCAGGTTTGTTCAT
```



Σύγκριση πρωτεϊνών - Εξαγωγή συμπερασμάτων για τη λειτουργία (3/4)

Protein A expressed in colon cancer:

TCGCTGCGAAGGACATTTGGGCTGTGTGTGCGACGCGGGTTCGGAGGGGCAGTCGGGGGAACCGCGAAGAAGCC
GAGGAGCCCGGAGCCCCGCGTGACGCTCCTCTCTCAGTCCAAAAGCGGCTTTTGGTTTCGGCGCAGAGAGACCCGG
GGGTCT**CAGGACAGCGATTGTCATTGCTGA**AGCTTTTCCTCGAAAAGCGCCGCCCTGCCCTTGCCCCGAGAA

Protein B expressed in colon cancer :

CAGACAAAGAGCACCGCAGGGCCGATCACGCTGGGGGCGCTGAGGCCGGCCATGGTCATGGAAGTGGGCACCCT
GGACGCTGGAGGCCTGCGGGCGCTGCTGGGGGAGCGAGCGGCGCAATGCCTGCTGCTGGACTGCCGCTCCTTCT
TCGCTTTCAACGCCGGCCACATCGCCGGCTCTGTCAACGTGCGCTTCAGCACCATCGTGCGGCGCCGGGCCAAGGG
CGCCATGGGCCTGGAGCACATCGTGCCCAACGCCGAGCTCCGCGG**TCAGGACAGCGATTGTCATTGCTGA**

Protein C expressed in colon cancer :

CCGCTGCTGGCCGGCGCCTACCACGCCGTGGTGTGCTGGACGAGCGCAGCGCCGCCCTGGACGGCGCCAAGC
GCGACGGCACCCCTGGCCCTGGCGGCCGGCGCGCTCTGCCGCGAGGCGCGCGCCGCGCAAGTCTTCTTCTCAAAG
GAGGATACGAAGCGTTTTTCGGCTTCCTGCCCGGAGCT**CAGGACAGCGATTGTCATTGCTGATGTGCAGCAAACAG**
TCGACCCCATGGGGCTCAGCCTTCCCCTGAGTACTAGCGTCCCTGACAGCGCGGAATCTGGGTGCAGTT



Σύγκριση πρωτεϊνών - Εξαγωγή συμπερασμάτων για τη λειτουργία (4/4)

Protein D expressed in colon cancer :

TAAGTGCCTTGATCAACGTCTCAGCCAATTGTCCCAACCATTTTGAGGGTCACTACCAGTACAAGAGCATCCCTGTGGAGGACAACCAC
AAGGCAGACATCAGCTCCTGGTTCAACGAGGCCATTGACTTCATAGACTCCATCAAGAATGCTGGAGGAAGGGTGTGGTCCACTGCC
AGGCAGGCATTTCCCGGTGAGCCACCATCTGCCTTGCTTACCTTATGAGGACTAATCGAGTCAATCAGGACAGCGATTGTCATTGCTGAA
GCTGGACGAGGCCTTTGAGTTTGTGAAGCAGAGGC

Protein E expressed in colon cancer :

GAAGCATCATCTCTCCCAACTTCAGCTTCATGGGCCAGCTGCTGCAGTTTGAGTCCCAGGTGCTGGCTCCGCACTGTTCCGGCAGAGGC
TGGGAGCCCCGCCATGGCTGTGCTCGACCGAGGCACCTCCACCACCACCGTGTCAACTTATCAGGACAGCGATTGTCATTGCTGACC
CCGTCTCCATCCCTGTCCACTCCACGAACAGTGCCTGAGCTACCTTCAGAGCCCCATTACGACCTCTCCAGCTGCTGAAAGGCCACG
GGAGGTGAGGCTCTTCACATCCCATTGGGACTC

Protein F with unknown function:

CATGCTCCTTGAGAGGAGAAATGCAATAACTCTGGGAGGGGCTCAGGACAGCGATTGTCATTGCTGATCGAGAGGGCTGGTCCTTAT
TTATTTAACTTCACCCGAGTTCCTCTGGGTTTCTAAGCAGTTATGGTGATGACTTAGCGTCAAGACATTTGCTGAACTCAGCACATTCGG
GACCAATATATAGTGGGTACATCAAGTCCATCTGACAAAATGGGGCAGAAGAGAAAGGACTCAGTGTGTGATCCGGTTTTCTTTTTGCTC
GCCCTGTTTTTTGTAGAATCTCTTCATGCTTGACATACCTACCAGTATTATTCCCAGCAGACATATACATATGAGAATATACCTTATTATT
TTTGTGTAGGTGTCTGCCTTCACAAATGTCATTGTCTACTCCTAGAAGAACCACAAATACCTCAATTTTTGTTTTGAGTACTGTACTATCCTG
TAAATATATCTTAAGCAGGTTTGTTCATCA

Συμπέρασμα: Πιθανά και η πρωτεΐνη F παίζει ρόλο στον καρκίνο παχύ εντέρου



Σύγκριση πρωτεϊνών - Εξαγωγή συμπερασμάτων για λειτουργικό τμήμα πρωτεΐνης (1/2)

Protein A with membrane localization:

```
TCGCTGCGAAGGACATTTGGGCTGTGTGTGCGACGCGGGTTCGGAGGGGCAGTCTGGGGGAACCGCGAAGAAGCCGAGGAGCCCGGAGCCCCGCGTGACG
CTCCTCTCTCAGTCCAAAAGCGGCTTTTGGTTTCGGCGCAGAGAGACCCGGGGTCTTCAGGACAGCGATTGTCATTGCTGAAGCTTTTCTCGAAAAGCGCC
GCCCTGCCCTTGGCCCCGAGAACAGACAAAGAGCACCCGAGGGCCGATCACGCTGGGGGCGCTGAGGCCGGCCATGGTCATGGAAGTGGGCACCCTGGAC
GCTGGAGGCCTGCGGGCGCTGCTGGGGGAGCGAGCGGCGCAATGCCTGCTGCTGGACTGCCGCTCCTTCTTCGCTTTCAACGCCGGCCACATCGCCGGCTC
TGTCAACGTGCGCTTCAGCACCATCTGTCGGCGCCGGGCCAAGGGCGCCATGGGCCTGGAGCACATCGTGCCCAACGCCGAGCTCCGCGGTGAGGACAGC
GATTGTCATTGCTGACCGCCTGCTGGCCGGCGCCTACCACGCCGTGGTGTGCTGGACGAGCGCAGCGCCGCCCTGGACGGCGCCAAGCGCGACGGCACCC
TGGCCCTGGCGGCCGGCGCGCTCTGCCGCGAGGCGCGCGCCGCGCAAGTCTTCTCCTCAAAGGAGGATACGAAGCGTTTTTCGGCTTCTGCCCGGAGCTC
AGGACAGCGATTGTCATTGCTGATGTGCAGCAAACAGTCGACCCCCATGGGGCTCAGCCTTCCCCTGAGTACTAGCGTCCCTGACAGCGCGGAATCTGGGTG
CAGTTCCTGCAGTACCCCACTCTACGATCAGGGTGGCCCCGGTGGAAATCTGCCCTTCTGTACCTGGGCAGTGCATCACGCTTCCCGCAAGGACATGCTG
GATGCCTTGGGCA
```

Protein B with membrane localization :

```
TAAGTGCCTTGATCAACGTCTCAGCCAATTGTCCTCAACCATTTTGGAGGGTCACTACCAAGTACAAGAGCATCCCTGTGGAGGACAACCACAAGGCAGACATCA
GCTCCTGGTTCAACGAGGCCATTGACTTCATAGACTCCATCAAGAATGCTGGAGGAAGGGTGTGTTGTCCACTGCCAGGCAGGCATTTCCCGGTGAGCCACCA
TCTGCCTTGCTTACCTTATGAGGACTAATCGAGTCATCAGGACAGCGATTGTCATTGCTGAAGCTGGACGAGGCCCTTTGAGTTTGTGAAGCAGAGGCGAAGC
ATCATCTCTCCAACCTCAGCTTCATGGGCCAGCTGCTGCAGTTTGAGTCCCAGGTGCTGGCTCCGCACTGTTTCGGCAGAGGCTGGGAGCCCCGCCATGGCT
GTGCTCGACCGAGGCACCTCCACCACCACCGTGTCAACTTATCAGGACAGCGATTGTCATTGCTGACCCCGTCTCCATCCCTGTCCACTCCACGAACAGTGC
GCTGAGCTACCTTCAGAGCCCCATTACGACCTCTCCAGCTGCTGAAAGGCCACGGGAGGTGAGGCTTTCACATCCCATTTGGGACTCCATGCTCCTTGAGA
GGAGAAATGCAATAACTCTGGGAGGGGCTCAGGACAGCGATTGTCATTGCTGATCGAGAGGGCTGGTCTTATTTATTTAACTTCACCCGAGTTCCTCTGGG
TTTCTAAGCAGTTATGGTGATGACTTAGCGTCAAGACATTTGCTGAACTCAGCACATTCGGGACCAATATATAGTGGGTACATCAAGTCCATCTGACAAAATG
GGCAGAAGAGAAAGGACTCAGTGTGTGATCCGGTTTCTTTTTGCTCGCCCCTGTTTTGTAGAAATCTTTCATGCTTGACATACCTACCAGTATTATCCCG
ACGACACATATACATATGAGAATATACCTTATTTATTTTTGTGTAGGTGCTGCCTTCAAAATGTCATTGTCTACTCCTAGAAGAACCAAATACCTCAATTTT
GTTTTGAGTACTGTACTATCCTGTAAATATATCTTAAGCAGGTTTGTTTCA
```

Συμπέρασμα: Πιθανά το κοινό τμήμα να είναι υπεύθυνο για τη μεμβρανική εντόπιση των πρωτεϊνών



Σύγκριση πρωτεϊνών - Εξαγωγή συμπερασμάτων για λειτουργικό τμήμα πρωτεΐνης (2/2)

Protein A with membrane localization:

TCGCTGCGAAGGACATTTGGGCTGTGTGTGCGACGCGGGTCGGAGGGGCAGTCGGGGGAACCGCGAAGAAGCCGAGGAGCCCGGAGCCCCGCGTGACGCTCCTC
TCTCAGTCCAAAAGCGGGCTTTTGTTTCGGCGCAGAGAGACCCGGGGGTCTT**CAGGACAGCGATTGTCATTGCTGA**AGCTTTTCTCGAAAAGCGCCGCCCTGCCCTT
GGCCCCGAGAACAGACAAAAGAGCACCCGACAGGGCCGATCACGCTGGGGGCGCTGAGGCCGGCCATGGTCATGGAAGTGGGCACCCTGGACGCTGGAGGCCTGCGG
GCGCTGCTGGGGGAGCGAGCGGGCGCAATGCCTGCTGCTGGACTGCCGCTCCTTCTCGCTTTCAACGCGGGCCACATCGCCGGCTCTGTCAACGTGCGCTTCAGCAC
CATCGTGCGGCGCCGGGCCAAGGGCGCCATGGCCTGGAGCACATCGTGCCCAACGCGAGCTCCGCGG**TCAGGACAGCGATTGTCATTGCTG**ACCCGCTGCTGG
CCGGCGCCTACCACGCCGTGGTGTGCTGGACGAGCGCAGCGCCGCCCTGGACGGCGCCAAGCGCGACGGCACCCCTGGCCCTGGCGGCCGGCGCGCTCTGCCGCG
AGGCGCGCGCCGCGCAAGTCTTCTCCTCAAAGGAGGATACGAAGCGTTTTTCGGCTTCTGCCCGGAGCT**CAGGACAGCGATTGTCATTGCTG**ATGTGCAGCAAAC
AGTCGACCCCATGGGGCTCAGCCTTCCCTGAGTACTAGCGTCCCTGACAGCGCGGAATCTGGGTGCAGTTCCTGCAGTACCCACTCTACGATCAGGGTGGCCCCG
GTGGAAATCCTGCCCTTCTGTACCTGGGCAGTGCATCACGCTTCCCGCAAGGACATGCTGGATGCCTTGGGCA

Protein B with membrane localization :

TAAGTGCCTTGATCAACGTCTCAGCCAATTGTCCAACCATTTTGGGGTCACTACCAGTACAAGAGCATCCCTGTGGAGGACAACCACAAGGCAGACATCAGCTCCT
GGTTCAACGAGGCCATTGACTTCATAGACTCCATCAAGAATGCTGGAGGAAGGGTGTGTTGCCACTGCCAGGCAGGCATTTCCCGGTCAGCCACCATCTGCCTTGCTT
ACCTTATGAGGACTAATCGAGTCA**TAGGACAGCGATTGTCATTGCTGA**AGCTGGACGAGGCCTTTGAGTTTGTGAAGCAGAGGCGAAGCATCATCTCTCCCAACTT
CAGCTTCATGGGCCAGCTGCTGCAGTTTGTGAGTCCAGGTGCTGGCTCCGCACTGTTCCGGCAGAGGCTGGGAGCCCCGCCATGGCTGTGCTCGACCGAGGCACCTCC
ACCACCACCGTGTTCAACTTAT**CAGGACAGCGATTGTCATTGCTGA**CCCCGTCTCCATCCCTGTCCACTCCACGAACAGTGCCTGAGCTACCTTCAGAGCCCCATTA
CGACCTCTCCAGCTGCTGAAAGGCCACGGGAGGTGAGGCTTTCACATCCCATGGGACTCCATGCTCCTTGAGAGGAGAAATGCAATAACTCTGGGAGGGG**CTC**
AGGACAGCGATTGTCATTGCTGATCGAGAGGGCTGGTCTTATTTATTTAACTTCACCCGAGTTCCTCTGGGTTTCTAAGCAGTTATGGTGATGACTTAGCGTCAAGA
CATTTGCTGAACTCAGCACATTCGGGACCAATATATAGTGGGTACATCAAGTCCATCTGACAAAATGGGGCAGAAGAGAAAGGACTCAGTGTGTGATCCGGTTTCTT
TTTGCTCGCCCCTGTTTTTGTAGAATCTTTCATGCTTGACATACCTACCAGTATTATCCCGACGACACATATACATATGAGAATATACCTATTTATTTTTGTGTAGG
TGTCTGCCTTCAAAAATGCATTGTCTACTCCTAGAAGAACCAAATACCTCAATTTTTGTTTTGAGTACTGTACTATCCTGTAATAATATCTTAAGCAGTTTGTTTTCA

Συμπέρασμα: Πιθανά το κοινό τμήμα να είναι υπεύθυνο για τη μεμβρανική εντόπιση των πρωτεϊνών



Σύστημα βαθμονόμησης στοίχισης - DNA & RNA

Sequence 1: ATCGGATCT

Sequence 2: ACGGACT

ΠΡΟΣΟΧΗ: Συγκρίνουμε νουκλεοτίδια, μόνο 4 πιθανές βάσεις

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| A | T | C | G | G | A | T | - | C | T |
| | | | | | | | | | |
| A | - | C | - | G | G | - | A | C | T |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| A | T | C | G | G | A | T | C | T |
| | | | | | | | | |
| A | - | C | G | G | A | - | C | T |

- Πιθανό σύστημα βαθμονόμησης:
 - όμοιο κατάλοιπο: +2
 - διαφορετικό κατάλοιπο: -1
 - κενό: -2
- **Alignment 1:** $5 \times 2 - 1(1) - 4(2) = 10 - 1 - 8 = 1$
- **Alignment 2:** $7 \times 2 - 0(1) - 2(2) = 14 - 0 - 4 = 10$



Στοιχίση πρωτεϊνών;;;

| | | | | | | |
|---|---|---|---|---|---|---|
| A | S | K | T | M | P | I |
| I | I | ? | ? | I | I | I |
| A | S | R | H | M | P | I |

| | | | | | | |
|---|---|---|---|---|---|---|
| A | S | K | T | M | P | I |
| I | I | ? | ? | I | I | I |
| A | S | Y | H | M | P | I |

- Πώς στοιχίζουμε πρωτεΐνες.
- Αμινοξέα με παρόμοιες ιδιότητες.
- Αμινοξέα με διαφορετικές ιδιότητες.



Πίνακες βαθμονόμησης (scoring matrices) ή πίνακες αντικατάστασης

- Για την αντικατάσταση αμινοξέων:
 - PAM.
 - BLOSUM.
- Για την αντικατάσταση νουκλεοτιδίων (DNA):
 - Το DNA είναι λιγότερο συντηρημένο από τις πρωτεΐνες.
 - Δεν είναι αποτελεσματικό να συγκρίνουμε κωδικοποιούσες περιοχές, από τις οποίες προκύπτουν οι πρωτεΐνες, σε επίπεδο νουκλεοτιδίων.



The log odds form (the mutation data matrix) of PAM250

| | C | S | T | P | A | G | N | D | E | Q | H | R | K | M | I | L | V | F | Y | W | |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|----|----|---|
| C | 12 | | | | | | | | | | | | | | | | | | | | C |
| S | 0 | 2 | | | | | | | | | | | | | | | | | | | S |
| T | -2 | 1 | 3 | | | | | | | | | | | | | | | | | | T |
| P | -3 | 1 | 0 | 6 | | | | | | | | | | | | | | | | | P |
| A | -2 | 1 | 1 | 1 | 2 | | | | | | | | | | | | | | | | A |
| G | -3 | 1 | 0 | -1 | 1 | 5 | | | | | | | | | | | | | | | G |
| N | -4 | 1 | 0 | -1 | 0 | 0 | 2 | | | | | | | | | | | | | | N |
| D | -5 | 0 | 0 | -1 | 0 | 1 | 2 | 4 | | | | | | | | | | | | | D |
| E | -5 | 0 | 0 | -1 | 0 | 0 | 1 | 3 | 4 | | | | | | | | | | | | E |
| Q | -5 | -1 | -1 | 0 | 0 | -1 | 1 | 2 | 2 | 4 | | | | | | | | | | | Q |
| H | -3 | -1 | -1 | 0 | -1 | -2 | 2 | 1 | 1 | 3 | 6 | | | | | | | | | | H |
| R | -4 | 0 | -1 | 0 | -2 | -3 | 0 | -1 | -1 | 1 | 2 | 6 | | | | | | | | | R |
| K | -5 | 0 | 0 | -1 | -1 | -2 | 1 | 0 | 0 | 1 | 0 | 3 | 5 | | | | | | | | K |
| M | -5 | -2 | -1 | -2 | -1 | -3 | -2 | -3 | -2 | -1 | -2 | 0 | 0 | 6 | | | | | | | M |
| I | -2 | -1 | 0 | -2 | -1 | -3 | -2 | -2 | -2 | -2 | -2 | -2 | -2 | 2 | 5 | | | | | | I |
| L | -6 | -3 | -2 | -3 | -2 | -4 | -3 | -4 | -3 | -2 | -2 | -3 | -3 | 4 | 2 | 6 | | | | | L |
| V | -2 | -1 | 0 | -1 | 0 | -1 | -2 | -2 | -2 | -2 | -2 | -2 | -2 | 2 | 4 | 2 | 4 | | | | V |
| F | -4 | -3 | -3 | -5 | -4 | -5 | -4 | -6 | -5 | -5 | -2 | -4 | -5 | 0 | 1 | 2 | -1 | 9 | | | F |
| Y | 0 | -3 | -3 | -5 | -3 | -5 | -2 | -4 | -4 | -4 | 0 | -4 | -4 | -2 | -1 | -1 | -2 | 7 | 10 | | Y |
| W | -8 | -2 | -5 | -6 | -6 | -7 | -4 | -7 | -7 | -5 | -3 | 2 | -3 | -4 | -5 | -2 | -6 | 0 | 0 | 17 | W |
| | C | S | T | P | A | G | N | D | E | Q | H | R | K | M | I | L | V | F | Y | W | |

Tryptophan: W

Cysteine: C



Πίνακες αντικατάστασης PAM (Percent Accepted Mutation - PAM or Dayhoff Matrices)

- Μελετήθηκαν από την Margaret Dayhoff.
- Αξιολογούν - βαθμολογούν την αντικατάσταση ενός αμινοξέος από ένα άλλο.
- Αξιολογήθηκαν:
 - 1572 αμινοξικές αντικαταστάσεις από,
 - 71 πρωτεϊνικά γκρουπ με,
 - 85% ομοιότητα τουλάχιστον ,
- Η κατασκευή τους βασίζεται στις:
 - στοιχίσεις πολλών **όμοιων πρωτεϊνών**, μικρής εξελεγκτικής απόστασης
 - αποδεκτές σημειακές μεταλλάξεις - accepted mutations (αντικατάσταση αμινοξέος από κάποιο άλλο που δεν επηρεάζει τη λειτουργία της πρωτεΐνης).

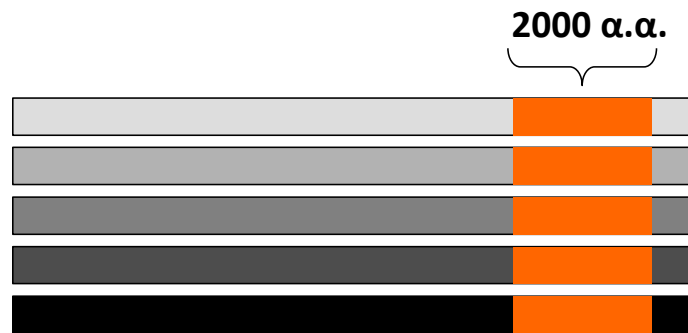


Πίνακες αντικατάστασης BLOSUM (BLOcks amino acid SUBstitution Matrices)

- Αναπτύχθηκαν από τους Stephen και Georgia Henikoff, 1992.
- Βασίζονται στην απαρίθμηση αντικατάστασης αμινοξέων από άλλα αμινοξέα σε συντηρημένες περιοχές αλληλουχιών.
 - μικρής ομοιότητας, και
 - απομακρυσμένης εξελεγκτικής σχέσης.
- Εξετάστηκαν πάνω από 2000 συντηρημένες περιοχές (blocks) που ανήκουν σε περισσότερες από 500 οικογένειες συσχετιζόμενων πρωτεϊνών.
- **Οι περιοχές (blocks) δεν περιέχουν κενά.**



BLOSUM MATRICES



■ : Συντηρημένη περιοχή

- Οι περιοχές βρέθηκαν σε πρωτεϊνικές βάσεις δεδομένων και αντιπροσωπεύουν περισσότερες από 500 οικογένειες πρωτεϊνών.
- Οι πίνακες BLOSUM βασίζονται σε ανάλυση διαφορετικού τύπου αλληλουχιών και μεγαλύτερου αριθμού δεδομένων από τους πίνακες PAM.
- **Βάση δεδομένων: PROSITE.**
- **Prosite:** Database with protein domains. It provides list of proteins that are in the same family because they have a similar biochemical action. For each family, a pattern of amino acids that are characteristic of that function is provided.



Πίνακες αντικατάστασης BLOSUM (BLOcks SUBstitution Matrices)

- **BLOSUM n** : Το n δείχνει το ποσοστό ομοιότητας των αλληλουχιών που χρησιμοποιήθηκαν για να προκύψει ο συγκεκριμένος πίνακας.
- Αλληλουχίες με ομοιότητα τουλάχιστον 62% δίνουν τον πίνακα αντικατάστασης **BLOSUM62**.
- **BLOSUM με μεγάλο n : αλληλουχίες με μεγάλη ομοιότητα.**
- BLOSUM με μεγάλο n -> PAM με μικρό n .
- Γενικά, οι πίνακες BLOSUM είναι καλύτεροι για την εύρεση τοπικών στοιχίσεων.



Παράδειγμα BLOSUM62

| | C | S | T | P | A | G | N | D | E | Q | H | R | K | M | I | L | V | F | Y | W | |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|---|----|---|
| C | 9 | | | | | | | | | | | | | | | | | | | | C |
| S | -1 | 4 | | | | | | | | | | | | | | | | | | | S |
| T | -1 | 1 | 5 | | | | | | | | | | | | | | | | | | T |
| P | -3 | -1 | -1 | 7 | | | | | | | | | | | | | | | | | P |
| A | 0 | 1 | 0 | -1 | 4 | | | | | | | | | | | | | | | | A |
| G | -3 | 0 | -2 | -2 | 0 | 6 | | | | | | | | | | | | | | | G |
| N | -3 | 1 | 0 | -2 | -2 | 0 | 6 | | | | | | | | | | | | | | N |
| D | -3 | 0 | -1 | -1 | -2 | -1 | 1 | 6 | | | | | | | | | | | | | D |
| E | -4 | 0 | -1 | -1 | -1 | -2 | 0 | 2 | 5 | | | | | | | | | | | | E |
| Q | -3 | 0 | -1 | -1 | -1 | -2 | 0 | 0 | 2 | 5 | | | | | | | | | | | Q |
| H | -3 | -1 | -2 | -2 | -2 | -2 | 1 | -1 | 0 | 0 | 8 | | | | | | | | | | H |
| R | -3 | -1 | -1 | -2 | -1 | -2 | 0 | -2 | 0 | 1 | 0 | 5 | | | | | | | | | R |
| K | -3 | 0 | -1 | -1 | -1 | -2 | 0 | -1 | 1 | 1 | -1 | 2 | 5 | | | | | | | | K |
| M | -1 | -1 | -1 | -2 | -1 | -3 | -2 | -3 | -2 | 0 | -2 | -1 | -1 | 5 | | | | | | | M |
| I | -1 | -2 | -1 | -3 | -1 | -4 | -3 | -3 | -3 | -3 | -3 | -3 | -3 | 1 | 4 | | | | | | I |
| L | -1 | -2 | -1 | -3 | -1 | -4 | -3 | -4 | -3 | -2 | -3 | -2 | -2 | 2 | 2 | 4 | | | | | L |
| V | -1 | -2 | 0 | -2 | 0 | -3 | -3 | -3 | -2 | -2 | -3 | -3 | -2 | 1 | 3 | 1 | 4 | | | | V |
| F | -2 | -2 | -2 | -4 | -2 | -3 | -3 | -3 | -3 | -3 | -1 | -3 | -3 | 0 | 0 | 0 | -1 | 6 | | | F |
| Y | -2 | -2 | -2 | -3 | -2 | -3 | -2 | -3 | -2 | -1 | 2 | -2 | -2 | -1 | -1 | -1 | -1 | 3 | 7 | | Y |
| W | -2 | -3 | -2 | -4 | -3 | -2 | -4 | -4 | -3 | -2 | -2 | -3 | -3 | -1 | -3 | -2 | -3 | 1 | 2 | 11 | W |

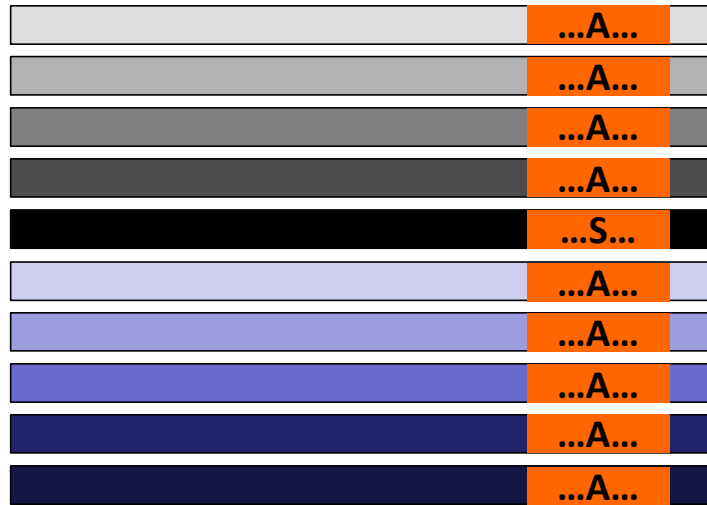


Ερμηνεία τιμής του πίνακα BLOSUM (1/4)

- Η τιμή κάθε ζεύγους αμινοξέος δείχνει την πιθανότητα να εμφανιστεί η συγκεκριμένη αντικατάσταση κατά την εξέλιξη.
- Δίνεται από το λόγο:
 - $\lambda = \text{συχνότητα εμφάνισης του ζεύγους} / \text{αναμενόμενη συχνότητα εμφάνισης του ζεύγους.}$
 - $\lambda = q_{ij}/e_{ij},$
 - q : expected, and
 - e : observed frequency of appearance.



Ερμηνεία τιμής του πίνακα BLOSUM (2/4)



| Pairs | Frequency of appearance, 1 st | Frequency of appearance, 2 nd |
|-------|--|--|
| AA | 9 | 9 |
| AS | 9 | 1 |
| SS | 1 | 1 |

- **Pair frequency** c_{ij} , for each pair i and j , for each column k using:

– For “like” comparisons $c_{ii}^{(k)} = \binom{n_i}{2}$

– For “unlike” comparisons $c_{ij}^{(k)} = n_i n_j$



Ερμηνεία τιμής του πίνακα BLOSUM (3/4)

...A...
...A...
...A...
...A...
...S...
...A...
...A...
...A...
...A...
...A...

| Pairs | Frequency of appearance, 1 st | Frequency of appearance, 2 nd |
|-------|--|--|
| AA | 9 | 9 |
| AS | 9 | 1 |
| SS | 1 | 1 |

$$c_{ii}^{(k)} = \binom{n_i}{2}$$

$$c_{ij}^{(k)} = n_i n_j$$

$$c_{AA} = \binom{9}{2} = \frac{9(9-1)}{2} = 36$$

$$c_{SS} = \binom{1}{2} = \frac{1(1-1)}{2} = 0$$

$$c_{AS} = 9 \cdot 1 = 9$$

| c_{ij} | A | S |
|----------|----|---|
| A | 36 | |
| S | 9 | 0 |



Ερμηνεία τιμής του πίνακα BLOSUM - Κανονικοποίηση

...A...
 ...A...
 ...A...
 ...A...
 ...S...
 ...A...
 ...A...
 ...A...
 ...A...
 ...A...

| c_{ij} | A | S |
|----------|----|---|
| A | 36 | |
| S | 9 | 0 |

$$T = \sum c_{ij} = w \frac{n(n-1)}{2}$$

w : columns

n : sequences

Normalization factor:

$$T = 1 \frac{10(10-1)}{2} = 45$$

Observed pair frequency, q_{ij} :

$$q_{AA} = \frac{36}{45} = 0.8$$

$$q_{AS} = \frac{9}{45} = 0.2$$

$$q_{SS} = \frac{0}{45} = 0$$

| q_{ij} | A | S |
|----------|-----|---|
| A | 0.8 | |
| S | 0.2 | 0 |



Ερμηνεία τιμής του πίνακα BLOSUM (4/4)

- Η τιμή κάθε ζεύγους αμινοξέος δείχνει την πιθανότητα να εμφανιστεί η συγκεκριμένη αντικατάσταση κατά την εξέλιξη.
- Δίνεται από το λόγο:
 - $\lambda = \text{συχνότητα εμφάνισης του ζεύγους} / \text{αναμενόμενη συχνότητα εμφάνισης του ζεύγους}$.
 - $\lambda = q_{ij}/e_{ij}$,
 - q : observed, and
 - e : expected frequency of appearance.



Ερμηνεία τιμής του πίνακα BLOSUM

Expected pair frequency, e_{ij} (1/3)

$$e_{ii} = p_i^2 \quad , \text{ for "like" comparisons}$$

$$e_{ij} = 2p_i p_j \quad , \text{ for "unlike" comparisons}$$

where, p_i : expected probability of occurrence of the i th residue

$$p_i = q_{ii} + \sum_{j \neq i} \frac{q_{ij}}{2}$$



Ερμηνεία τιμής του πίνακα BLOSUM

Expected pair frequency, e_{ij} (2/3)

$$e_{ij} = 2p_i p_j$$

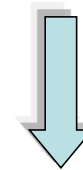
$$p_i = q_{ii} + \sum_{j \neq i} \frac{q_{ij}}{2}$$

| q_{ij} | A | S |
|----------|-----|---|
| A | 0.8 | |
| S | 0.2 | 0 |

Για κάθε κατάλοιπο

$$p_A = q_{AA} + \frac{q_{AS}}{2} = 0.8 + \frac{0.2}{2} = 0.9$$

$$p_S = q_{SS} + \frac{q_{AS}}{2} = 0 + \frac{0.2}{2} = 0.1$$



Για κάθε ζεύγος

$$e_{AA} = p_A^2 = 0.9^2 = 0.81$$

$$e_{SS} = p_S^2 = 0.01$$

$$e_{AS} = 2p_A p_S = 2 \cdot 0.9 \cdot 0.1 = 0.18$$



Ερμηνεία τιμής του πίνακα BLOSUM

Expected pair frequency, e_{ij} (3/3)

Observed frequency

| q_{ij} | A | S |
|----------|-----|---|
| A | 0.8 | |
| S | 0.2 | 0 |

Expected frequency

$$e_{AA} = p_A^2 = 0.9^2 = 0.81$$

$$e_{SS} = p_S^2 = 0.01$$

$$e_{AS} = 2p_A p_S = 2 \cdot 0.9 \cdot 0.1 = 0.18$$

Entry: $S = 2\log_2(\text{observed}/\text{expected})$

Entry:

$$S_{AA} = 2\log_2 \frac{\text{observed}}{\text{expected}} = 2 \cdot \log_2 \frac{0.8}{0.81} = -0.03$$

$$S_{AS} = 2\log_2 \frac{\text{observed}}{\text{expected}} = 2 \cdot \log_2 \frac{0.2}{0.18} = 0.3$$

$$S_{SS} = 2\log_2 \frac{\text{observed}}{\text{expected}} = 2 \cdot \log_2 \frac{0}{0.01} = ?$$

Δεν παρέχονται πληροφορίες για το ζεύγος SS



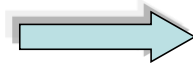
Παράδειγμα - Pair frequency c_{ij}

Sequence 1 A A I
 Sequence 2 S A L
 Sequence 3 T A L
 Sequence 4 T A V
 Sequence 5 A A L

$$c_{ii}^{(k)} = \binom{n_i}{2} \quad c_{ij}^{(k)} = n_i n_j$$

Pair frequency c_{ij} , for each pair i and j

| Pairs | n_i | n_j |
|-------|-------|-------|
| AT | 2 | 2 |
| IL | 1 | 3 |
| LL | 3 | 3 |
| AA | 2 | 2 |
| AA | 5 | 5 |



| c_{ij} | A | I | L | S | T | V |
|----------|------|---|---|---|---|---|
| A | 10+1 | | | | | |
| I | 0 | 0 | | | | |
| L | 0 | 3 | 3 | | | |
| S | 2 | 0 | 0 | 0 | | |
| T | 4 | 0 | 0 | 2 | 1 | |
| V | 0 | 1 | 3 | 0 | 0 | 0 |



Παράδειγμα - Normalization

Sequence 1 A A I
 Sequence 2 S A L
 Sequence 3 T A L
 Sequence 4 T A V
 Sequence 5 A A L

$$T = \sum c_{ij} = w \frac{n(n-1)}{2}$$

w : columns

n : sequences

Normalization factor, T:

$$T = 3 \frac{5(5-1)}{2} = 30$$

Pair frequency c_{ij} , for each pair i and j

| c_{ij} | A | I | L | S | T | V |
|----------|------|---|---|---|---|---|
| A | 10+1 | | | | | |
| I | 0 | 0 | | | | |
| L | 0 | 3 | 3 | | | |
| S | 2 | 0 | 0 | 0 | | |
| T | 4 | 0 | 0 | 2 | 1 | |
| V | 0 | 1 | 3 | 0 | 0 | 0 |



$c_{ij} / 30 = q_{ij}$

| q_{ij} | A | I | L | S | T | V |
|----------|-------|-------|-----|-------|-------|---|
| A | 0.366 | | | | | |
| I | 0 | 0 | | | | |
| L | 0 | 0.1 | 0.1 | | | |
| S | 0.066 | 0 | 0 | 0 | | |
| T | 0.133 | 0 | 0 | 0.066 | 0.033 | |
| V | 0 | 0.033 | 0.1 | 0 | 0 | 0 |



Παράδειγμα - Observed frequency q_{ij}

| | | | |
|------------|---|---|---|
| Sequence 1 | A | A | I |
| Sequence 2 | S | A | L |
| Sequence 3 | T | A | L |
| Sequence 4 | T | A | V |
| Sequence 5 | A | A | L |

Observed normalized pair frequency c_{ij} , for each pair i and j

| c_{ij} | A | I | L | S | T | V |
|----------|-------|-------|-----|-------|-------|---|
| A | 0.366 | | | | | |
| I | 0 | 0 | | | | |
| L | 0 | 0.1 | 0.1 | | | |
| S | 0.066 | 0 | 0 | 0 | | |
| T | 0.133 | 0 | 0 | 0.066 | 0.033 | |
| V | 0 | 0.033 | 0.1 | 0 | 0 | 0 |



Παράδειγμα - Expected frequency of ij pair, e_{ij}

$$e_{ii} = p_i^2$$

$$e_{ij} = 2p_i p_j$$

$$p_i = q_{ii} + \sum_{j \neq i} \frac{q_{ij}}{2}$$

Αναμενόμενη συχνότητα εμφάνισης ζεύγους

Πιθανότητα εμφάνισης i th καταλοίπου

| q_{ij} | A | I | L | S | T | V |
|----------|-------|-------|-----|-------|-------|---|
| A | 0.366 | | | | | |
| I | 0 | 0 | | | | |
| L | 0 | 0.1 | 0.1 | | | |
| S | 0.066 | 0 | 0 | 0 | | |
| T | 0.133 | 0 | 0 | 0.066 | 0.033 | |
| V | 0 | 0.033 | 0.1 | 0 | 0 | 0 |



| | Πιθανότητα εμφάνισης i th καταλοίπου |
|-------|--|
| p_A | $0.366 + (0.066 + 0.133)/2 = \mathbf{0.466}$ |
| p_I | $0 + (0.1 + 0.033)/2 = \mathbf{0.066}$ |
| p_L | $0.1 + (0.1 + 0.1)/2 = \mathbf{0.2}$ |
| p_S | $0 + (0.066 + 0.066)/2 = \mathbf{0.066}$ |
| p_T | $0.033 + (0.133 + 0.066)/2 = \mathbf{0.133}$ |
| p_V | $0 + (0.033 + 0.1)/2 = \mathbf{0.066}$ |



Expected frequency of ij pair, e_{ij}

$$e_{ii} = p_i^2 \quad e_{ij} = 2p_i p_j$$

$$p_i = q_{ii} + \sum_{j \neq i} \frac{q_{ij}}{2}$$

$$e_{AA} = p_A^2 = 0.466^2 = 0.217$$

$$e_{AI} = 2p_A p_I = 2 \cdot 0.466 \cdot 0.066 = 0.062$$

$$e_{AL} = 2p_A p_L = 2 \cdot 0.466 \cdot 0.2 = 0.186$$

$$e_{II} = p_I^2 = 0.066^2 = 0.004$$

Expected frequency of ij pair, e_{ij}

| p_i | Πιθανότητα εμφάνισης ith καταλοίπου |
|-------|--|
| p_A | $0.366 + (0.066 + 0.133) / 2 = \mathbf{0.466}$ |
| p_I | $0 + (0.1 + 0.033) / 2 = \mathbf{0.066}$ |
| p_L | $0.1 + (0.1 + 0.1) / 2 = \mathbf{0.2}$ |
| p_S | $0 + (0.066 + 0.066) / 2 = \mathbf{0.066}$ |
| p_T | $0.033 + (0.133 + 0.066) / 2 = \mathbf{0.133}$ |
| p_V | $0 + (0.033 + 0.1) / 2 = \mathbf{0.066}$ |



| e_{ij} | A | I | L | S | T | V |
|----------|-------|-------|-------|-------|-------|-------|
| A | 0.217 | | | | | |
| I | 0.062 | 0.004 | | | | |
| L | 0.186 | 0.026 | 0.04 | | | |
| S | 0.062 | 0.198 | 0.026 | 0.004 | | |
| T | 0.124 | 0.018 | 0.053 | 0.018 | 0.018 | |
| V | 0.062 | 0.009 | 0.027 | 0.009 | 0.018 | 0.004 |



Entry in BLOSUM matrix, S_{ij}

Observed normalized pair frequency c_{ij} ,
for each pair i and j

| | A | I | L | S | T | V |
|---|-------|-------|-----|-------|-------|---|
| A | 0.366 | | | | | |
| I | 0 | 0 | | | | |
| L | 0 | 0.1 | 0.1 | | | |
| S | 0.066 | 0 | 0 | 0 | | |
| T | 0.133 | 0 | 0 | 0.066 | 0.033 | |
| V | 0 | 0.033 | 0.1 | 0 | 0 | 0 |

Expected frequency of ij pair, e_{ij}

| e_{ij} | A | I | L | S | T | V |
|----------|-------|-------|-------|-------|-------|-------|
| A | 0.217 | | | | | |
| I | 0.062 | 0.004 | | | | |
| L | 0.186 | 0.026 | 0.04 | | | |
| S | 0.062 | 0.198 | 0.026 | 0.004 | | |
| T | 0.124 | 0.018 | 0.053 | 0.018 | 0.018 | |
| V | 0.062 | 0.009 | 0.027 | 0.009 | 0.018 | 0.004 |

Entry :

$$S_{ij} = 2 \log_2 \frac{\text{observed}}{\text{expected}}$$

$$S_{AA} = 2 \log_2 \frac{0.366}{0.217} = 2$$

$$S_{SA} = 2 \log_2 \frac{0.066}{0.062} = 0$$



Finally

| S_{ij} | A | I | L | S | T | V |
|----------|---|---|---|---|---|---|
| A | 2 | | | | | |
| I | ? | ? | | | | |
| L | ? | 4 | 3 | | | |
| S | 0 | ? | ? | ? | | |
| T | 0 | ? | ? | 4 | 2 | |
| V | ? | 4 | 4 | ? | ? | ? |





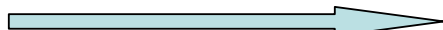


Comparison PAM vs. BLOSUM

- PAM is based on an evolutionary model using phylogenetic trees (85% similarity).
- BLOSUM assumes no evolutionary model, but rather conserved “blocks” of proteins.



Πίνακες αντικατάστασης - Ισοδυναμία

- PAM100  • Blosum90
- PAM120  • Blosum80
- PAM160  • Blosum62
- PAM200  • Blosum52
- PAM250  • Blosum45



Nucleic acid PAM Scoring Matrices (1/2)

- PAM: Point accepted mutations.
- Για τον DNA PAM1 πίνακα θεωρούμε ότι το 99% των αλληλουχιών συντηρείται ενώ το 1% μεταλλάσσεται.
- Υποθέτουμε ότι και οι 4 βάσεις, A-C-G-T, μεταλλάσσονται το ένα στο άλλο με την ίδια πιθανότητα.
- Υποθέτουμε ότι και τα 4 βάσεις έχουν την ίδια συχνότητα εμφάνισης.

| A. Model of uniform mutation rates among nucleotides | | | | |
|--|---------|---------|---------|------|
| | A | G | T | C |
| A | 0.99 | | | |
| G | 0.00333 | 0.99 | | |
| T | 0.00333 | 0.00333 | 0.99 | |
| C | 0.00333 | 0.00333 | 0.00333 | 0.99 |

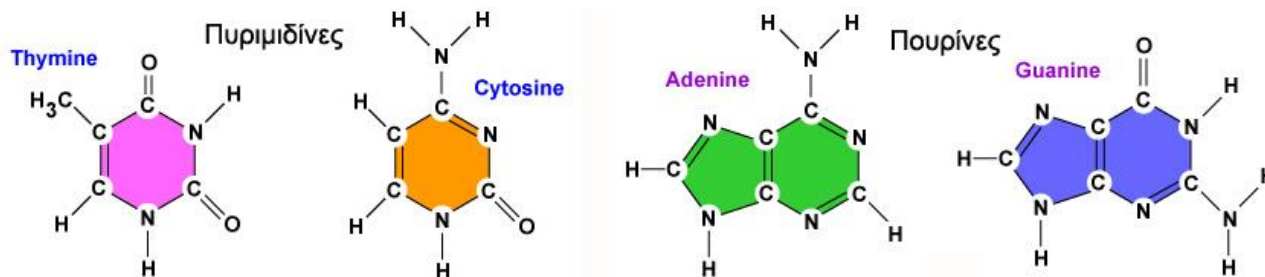


Nucleic acid PAM Scoring Matrices (2/2)

- Θεωρούμε ότι είναι 3 φορές πιο πιθανό οι πυριμιδίνες να μεταλλαχθούν σε πυριμιδίνες (T→C) και οι πουρίνες σε πουρίνες (A →G) σε σχέση με τη μετάλλαξη των πυριμιδών σε πουρίνες και το αντίστροφο.

B. Model of threefold higher transitions than transversions

| | A | G | T | C |
|---|-------|-------|-------|------|
| A | 0.99 | | | |
| G | 0.006 | 0.99 | | |
| T | 0.002 | 0.002 | 0.99 | |
| C | 0.002 | 0.002 | 0.006 | 0.99 |



4 αζωτούχες βάσεις: Θυμίνη (T), Κυτοσίνη (C), Αδερίνη (A) και Γουανίνη (G)



Οι ποινές για τα κενά (1/2)

- Η εισαγωγή κενών είναι απαραίτητη έτσι ώστε να έχουμε την καλύτερη στοίχιση των αλληλουχιών.
- Συσχετισμένη ή αφινική ποινή (affine gap penalty).
- Πιο πιθανό 1 κενό μήκους K , παρά K κενά μήκους 1.
- Αφαίρεση μιας τιμής για την εισαγωγή ενός κενού και πρόσθετες αφαιρέσεις για επόμενες επεκτάσεις του κενού.
- Το σύστημα βαθμονόμησης θα πρέπει να «τιμωρεί» περισσότερο τα καινούργια κενά.



Παράδειγμα - PAM250

C - K H V F C R V C I

C K K C - F C - K C V

- PAM250
- Εισαγωγή κενού: -10

$$\text{Score} = 12+5+9+12+12-3-2+4-3 \times 10$$

$$\text{Score} = 19$$



Οι ποινές για τα κενά (2/2)

- $p = -A - B \cdot (n-1)$:
 - όπου A : ποινή έναρξης του κενού (gap opening penalty).
 - B : ποινή επέκτασης (gap extension penalty).
 - n : μήκος του κενού.
- Συνήθως: $A=12$ και $B=4$, οπότε $p = -12 - 4 \cdot (n-1)$
- **Πολύ μικρή ποινή: Εισαγωγή πολλών κενών.**
- **Πολύ μεγάλη ποινή: Μη εισαγωγή κενών.**
- **Γενικά:** Η ποινή για ένα κενό n καταλοίπων θα πρέπει να είναι μικρότερη από τη συνολική ποινή για n κενά του 1 καταλοίπου.



Typical score matrix

- DNA:
 - Match = +1.
 - Mismatch = -3.
 - Gap penalty = -5.
 - Gap extension penalty = -2.
- Protein sequences:
 - Blossum62 matrix.
 - Gap open penalty = -11.
 - Gap extension = -1.

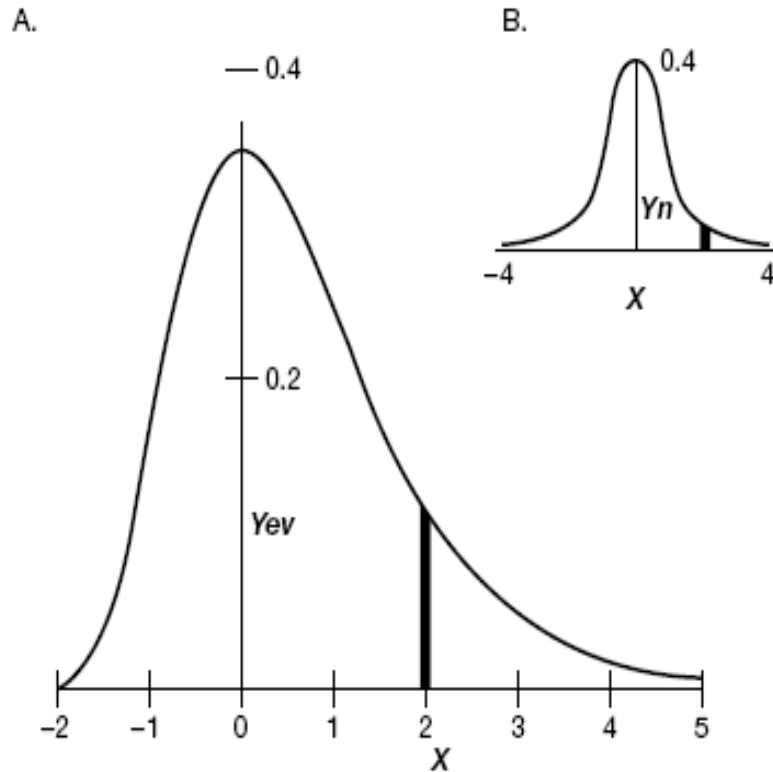


Εξετάζοντας τη σημαντικότητα της στοίχισης αλληλουχιών

- Για αλληλουχίες που ανήκουν στην ίδια οικογένεια και είναι προφανές, κάτι τέτοιο δεν είναι απαραίτητο.
- Για αλληλουχίες που δεν είναι ξεκάθαρο ότι είναι παρόμοιες, αυτό είναι αναγκαίο.
- **Συμπέρασμα:** Είναι πράγματι όμοιες ή το αποτέλεσμα της σύγκρισης είναι τυχαίο, δηλαδή το ίδιο με την περίπτωση που οι δύο αλληλουχίες δεν έχουν καμία σχέση.
- Ο έλεγχος θα γίνει σε όλες τις στοιχίσεις έτσι ώστε να επιλεχτεί εκείνη που είναι στατιστικά η σημαντικότερη.



Gumbel extreme value distribution



- **Σκοπός:** Να υπολογιστεί η πιθανότητα το σκορ της στοίχισης μεταξύ δύο τυχαίων αλληλουχιών να φτάσει το σκορ της στοίχισης δύο σχετικών αλληλουχιών. Όσο μικρότερη είναι αυτή η πιθανότητα, τόσο στατιστικά σημαντικό είναι το σκορ στοίχισης των 2 πραγματικών αλληλουχιών.

The distribution of alignment scores between random sequences follows the extreme value distribution (A), not the normal distribution (B)



Βαθμολόγηση στοίχισης αλληλουχιών

- **Score: A number used to assess the biological relevance of a finding.**
- In the context of sequence alignment, a score is a numerical value that describes the overall quality of an alignment. Higher numbers correspond to higher similarity.
- **The score scale depends on the scoring system used (matrix, gap penalty).**

| Gap penalty | Alignment | Identity / Similarity | Gaps | Score |
|-------------|--|-----------------------|-------|-------|
| 0 | <pre> 1 GTC-ATGCTA-GTCGT---GG---GTAGCATTTA-GCT-ATG-TGGG-GT 38 1 -TCGATGCT-GGTTCG-CAAGGCAAGTAG---TTATG-TCATGCT---AG- 39 </pre> | 27/50 (54.0%) | 23/50 | S=135 |
| 5 | <pre> 1 GTC-ATGCTAGTCG--TGGGTAGCATTTA-GCT-ATG-TGGGGT 38 1 -TCGATGCTGGTTCGCAAGGCAAGTAGTTATG-TCATGCTAG--- 39 </pre> | 26/44 (59.1%) | 11/44 | S=67 |
| 10 | <pre> 1 -----GTCATGCTAGTCGTTGGGTAGC 21 1 TCGATGCTGGTTCGCAAGGCAAGTAGTTATGTCATGCTAG----- 39 22 ATTTAGCTATGTGGGGT 38 39 ----- 39 </pre> | 10/67 (14.9%) | 57/67 | S=50 |

1: 27 matches x 5 = 135. No gap penalty

2: 26 matches x 5 = 130, 7 mismatch x (-4) = -28, 7 gaps x (-5) = -35 TOTAL=130-28-35=67

3: 10 matches x 5 = 50



Converting to Bit-scores

- The bit-score S' is a **normalized score expressed in bits**. Lets you estimate **the magnitude of the search space you would have to look through** before you would expect to find a score as good as or better than this one by chance.

$$S' = \frac{\lambda S - \ln(K)}{\ln 2}$$

where S : raw score. Parameters λ and K depend on the substitution matrix and on the gap penalty

The bit-score is a rescaled version of the raw alignment score that is independent of the size of the search space



Δείκτες στατιστικής σημαντικότητας (1/4)

P-value: Probability that an event occurs by chance

- The P-value associated to a score S is **the probability to obtain by chance a score x at least equal to S**

Low P-value → significant alignment

$$P(S' \geq x) = 1 - \exp[-e^{-x}] \approx e^{-x}$$



Δείκτες στατιστικής σημαντικότητας (2/4)

- E-value: Correlation of the p-value for multiple testing.
- Expected number of sequences that will produce same or better score by chance.
- The lower the E value, the more significant the score is
- Low E-value → significant alignment.

Αναμενόμενος αριθμός στοιχίσεων με σκορ τουλάχιστον S:

$$E = K \cdot m \cdot n \cdot e^{-\lambda S}$$

- m,n: Μήκος αλληλουχιών
- K, λ: Στατιστικές σταθερές που εξαρτώνται από το σύστημα βαθμονόμησης και τη συχνότητα εμφάνισης των καταλοίπων.
π.χ. για τον πίνακα PAM250, K=0.1 και λ=0.229



Δείκτες στατιστικής σημαντικότητας (3/4)

- **Z-score:**


- Measures how much standard deviations above the mean of the score distribution.



Δείκτες στατιστικής σημαντικότητας (4/4)

Score, E-value and P-value compared

| Score | E-value $E(S) = K m n e^{-\lambda S}$ | Probability $P = 1 - e^{-E(S)}$ |
|-------|--|------------------------------------|
| 39 | 12 | 0.999995 |
| 41 | 2.9 | 0.947456 |
| 42 | 1.4 | 0.764656 |
| 46 | 0.0842 | 0.080741 |
| 49 | 0.0100 | 0.009925 |
| 52 | 0.0012 | 0.00118 |
| 55 | 0.0001 | 0.0001 |

 $m = 980, n = 10,030,834,086$ ($m*n \sim 10^{13}$)
 $K = 1.37, \lambda = 0.711$

INCOGEN

<http://www.incogen.com>



Alignment - Output

Results (output) of BLAST

Bit-score E-value Similarity (%)
 Identity (%) Positive score in the substitution matrix Gaps (%)

Score = 83.6 bits (205), Expect = 3e-14, Method: Compositional matrix adjust.
 Identities = 61/136 (44%), Positives = 73/136 (53%), Gaps = 18/136 (13%)

```

Query  184  KPKPKQYPKVLPSNSTRRISPVTAKTSSSAEGVVVASESPVIAPHGSSHSRSLSKRRSS  243
          KP P  P+ ILPSN+ +R P      S      V+ AS+SPVI P+ +          RS
Sbjct  269  KPAPG-LPRFILPSNQPQRQLPPPPSDS-----VIHASQSPVIKPNYAGKPPGFVSARSV  322

Query  244  GALVDDD-----KRESHKHAEQARRNRLAVALHELASLIPAEWKQONVSAAPSKATT  295
          L  D          K+E HK AEQ RRNRL  AL EL  L+P E K+  +  PSKATT
Sbjct  323  RTLSGGDANTGDEFIKKEVHKVAEQGRRNRLNNAELNDLLPPELKES--AQVPSKATT  380

Query  296  VEAACRYIRHL--QQN  309
          VE AC+YIR L  QQN
Sbjct  381  VELACKYIRQLTGQQN  396
    
```



Μέθοδοι σύγκρισης αλληλουχιών

1. Οπτικοί.
 2. Με αλγόριθμους δυναμικού προγραμματισμού.
 3. Με ευρετικούς αλγόριθμους, που βασίζονται στην έννοια των «λέξεων».
- Πρέπει να λαμβάνεται υπόψη ο λόγος για τον οποίο γίνεται η στοίχιση αλληλουχιών.
 - Ανάλογα επιλέγεται και η μέθοδος σύγκρισης.
 - **Ομοιότητα της τάξης το 90%:** Εύκολη κατασκευή της στοίχισης με όλους τους αλγόριθμους.
 - **Ομοιότητα της τάξης του 25% (ζώνη του λυκόφωτος, «twilight zone of sequence alignment»):** Οριακή για ασφαλή εξαγωγή συμπερασμάτων.



Διαγράμματα πινάκων σημείων - Dot plots

- Αποτελούν μία από τις απλούστερες μεθόδους οπτικοποίησης της ομοιότητας μεταξύ δύο αλληλουχιών.
- Μέθοδος:
 - Σε ένα δισδιάστατο ορθογώνιο διάγραμμα ο κάθε άξονας αντιστοιχεί και σε μία αλληλουχία.
 - Αλληλουχία 1: **A C A G C G C G A G C**
 - Αλληλουχία 2: **A C A C G A G G**
 - Στο πλέγμα που προκύπτει:



Dot plots - Μέθοδος (1/3)

Σε ένα δισδιάστατο ορθογώνιο διάγραμμα ο κάθε άξονας αντιστοιχεί και σε μία αλληλουχία

| | A | C | A | G | C | G | C | G | A | G | C |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A | | | | | | | | | | | |
| C | | | | | | | | | | | |
| A | | | | | | | | | | | |
| C | | | | | | | | | | | |
| G | | | | | | | | | | | |
| A | | | | | | | | | | | |
| G | | | | | | | | | | | |
| G | | | | | | | | | | | |



Dot plots - Μέθοδος (2/3)

Τα τετράγωνα που αντιστοιχούν σε ταυτόσημα κατάλοιπα σημειώνονται (χρώμα)

| | A | C | A | G | C | G | C | G | A | G | C |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A | | | | | | | | | | | |
| C | | | | | | | | | | | |
| A | | | | | | | | | | | |
| C | | | | | | | | | | | |
| G | | | | | | | | | | | |
| A | | | | | | | | | | | |
| G | | | | | | | | | | | |
| G | | | | | | | | | | | |



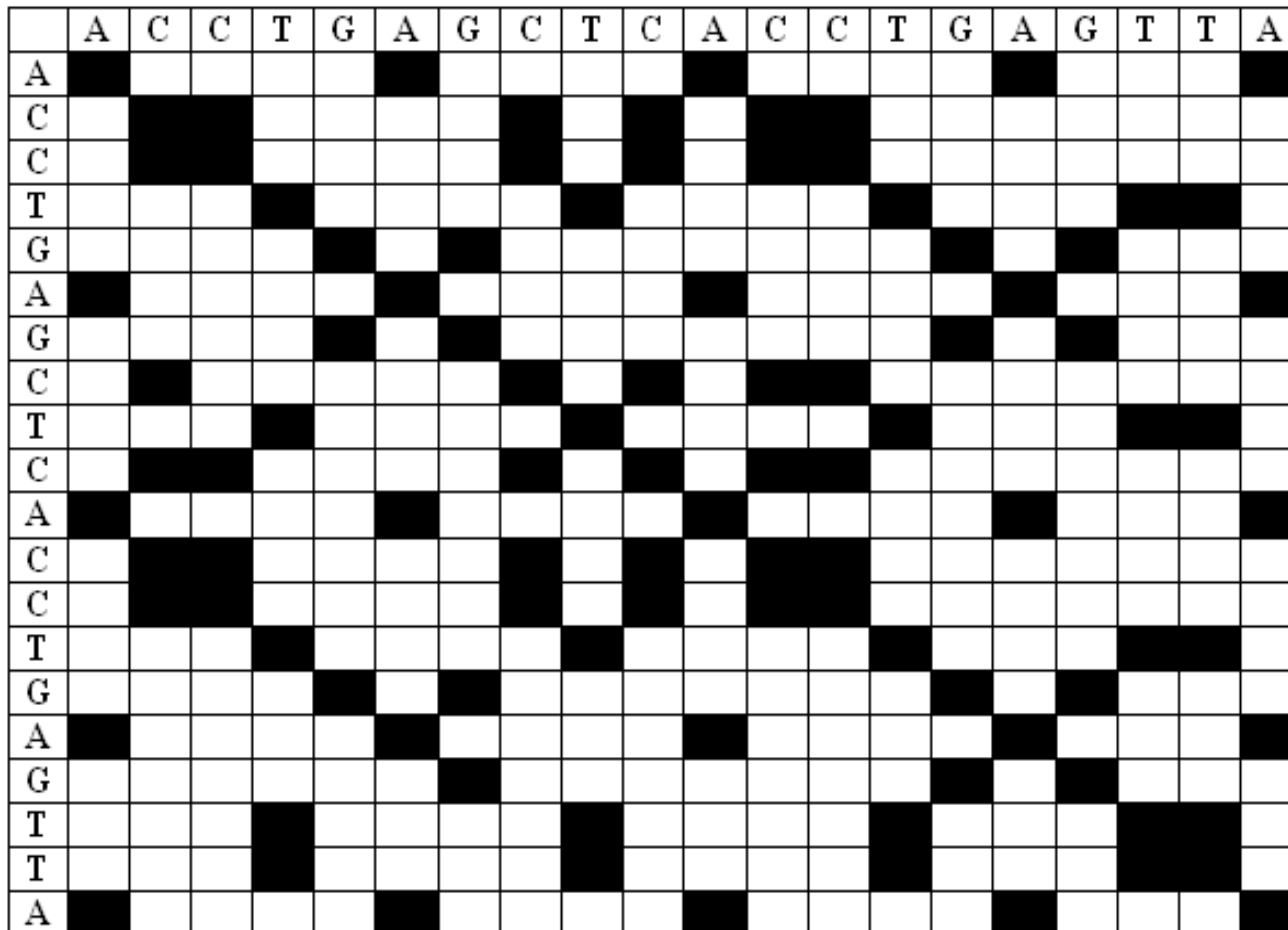
Dot plots - Μέθοδος (3/3)

Όταν οι αλληλουχίες είναι όμοιες κατά μήκος μιας περιοχής, σχηματίζεται μία διαγώνιο γραμμή με κατεύθυνση από πάνω αριστερά κάτω δεξιά

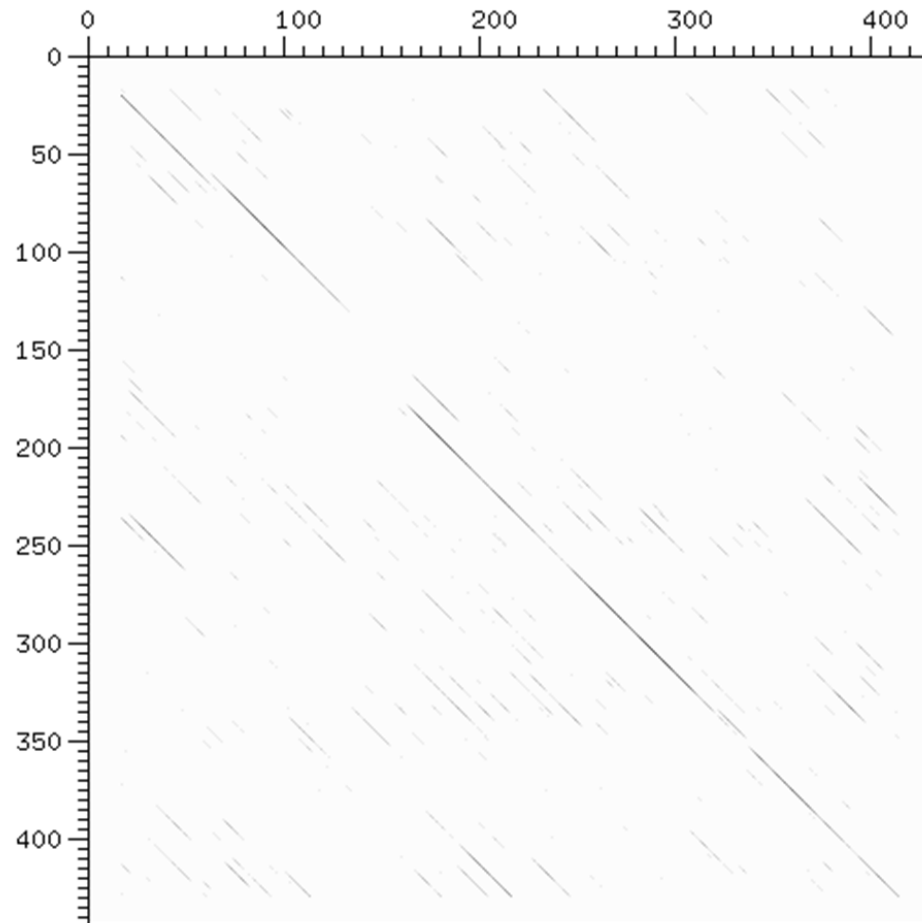
| | A | C | A | G | C | G | C | G | A | G | C |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A | ■ | | ■ | | | | | | ■ | | |
| C | | ■ | | | ■ | | ■ | | | | ■ |
| A | | | ■ | | | | | | | | |
| C | | ■ | | | ■ | | ■ | | | | ■ |
| G | | | | ■ | | ■ | | ■ | | ■ | |
| A | ■ | | ■ | | | | | | ■ | | |
| G | | | | ■ | | ■ | | ■ | | ■ | |
| G | | | | ■ | | ■ | | ■ | | ■ | |



Dot plots - Παράδειγμα (1/2)



Dot plots - Παράδειγμα (2/2)



<http://myhits.isb-sib.ch/cgi-bin/dotlet>
<http://www.vivo.colostate.edu/molkit/dnadot/>

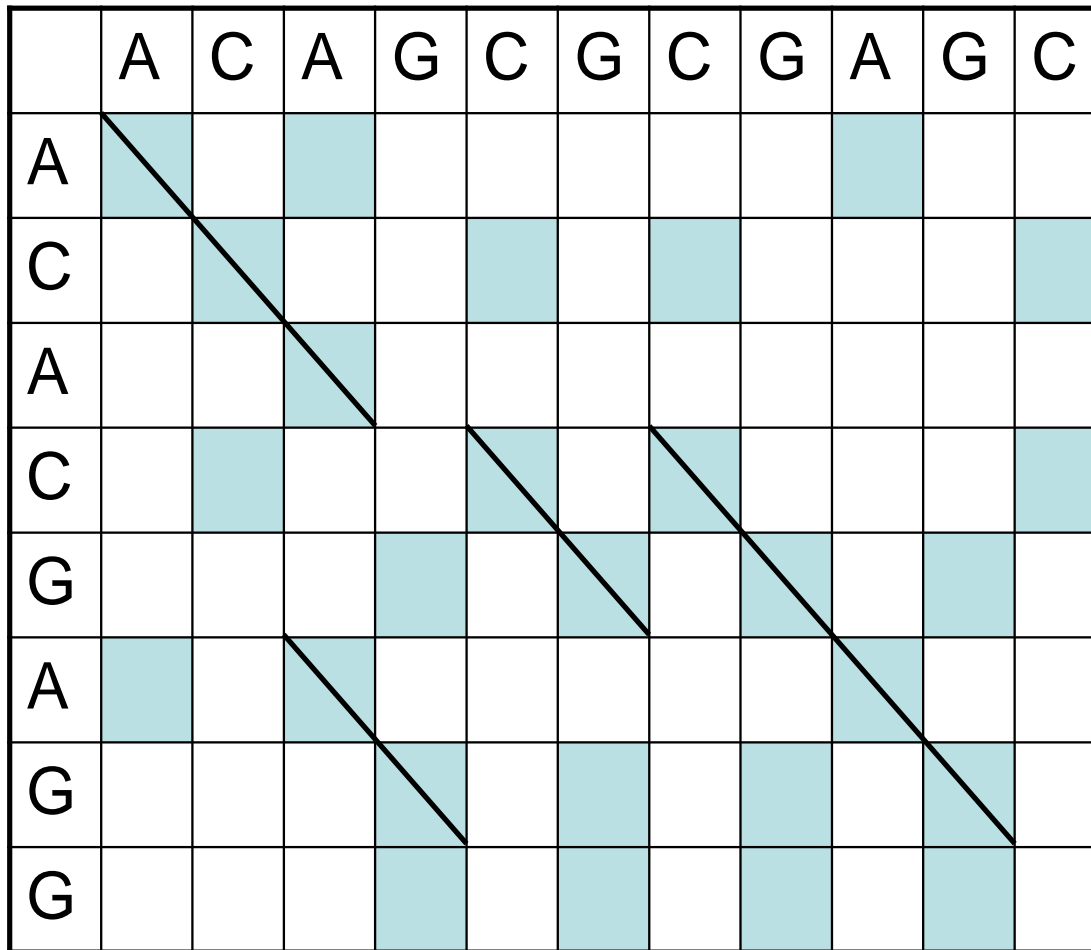


«Θόρυβος» στα dot plots

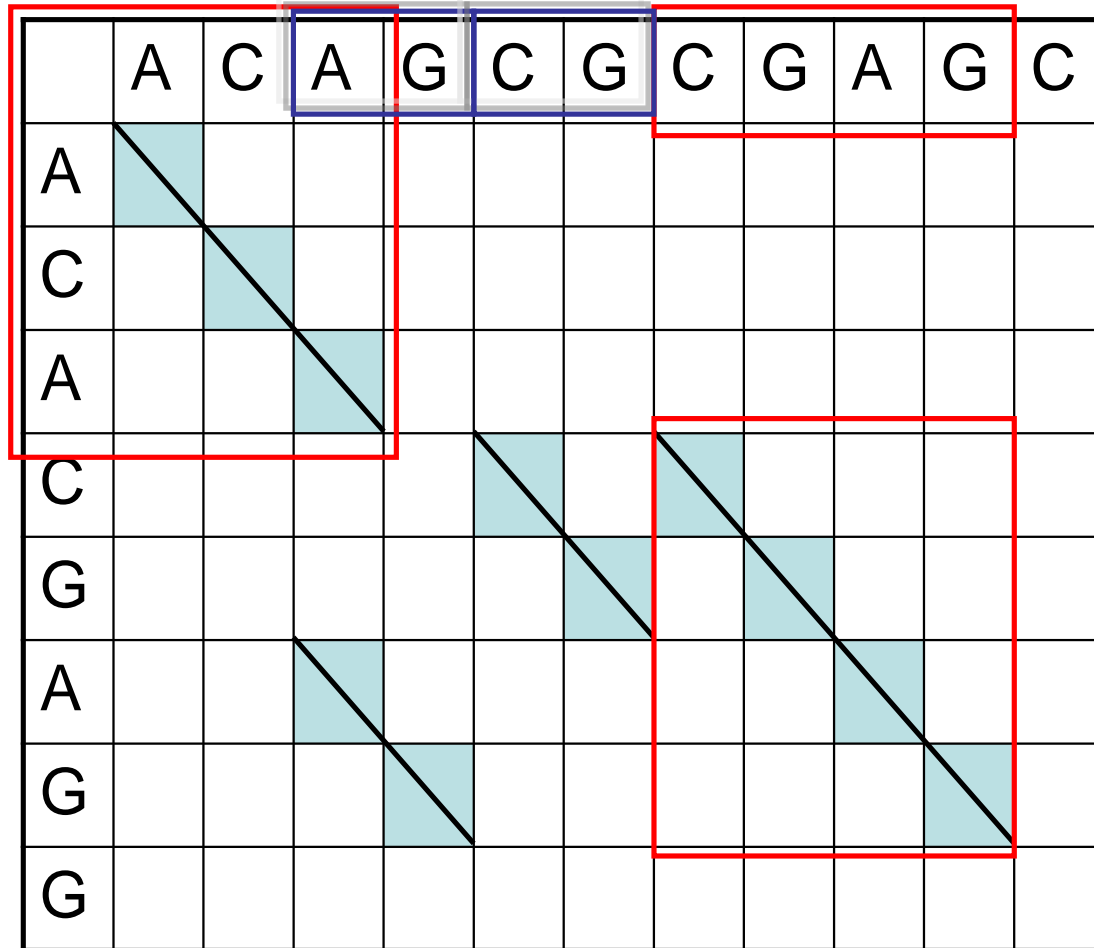
- 1 από τις 4 βάσεις έχει την πιθανότητα να ταιριάζει τυχαία.
- **Μέθοδος φιλτραρίσματος:** μία περιοχή σκιάζεται αν υπερβαίνει κάποιον συγκεκριμένο αριθμό (φίλτρο) όμοιων καταλοίπων.



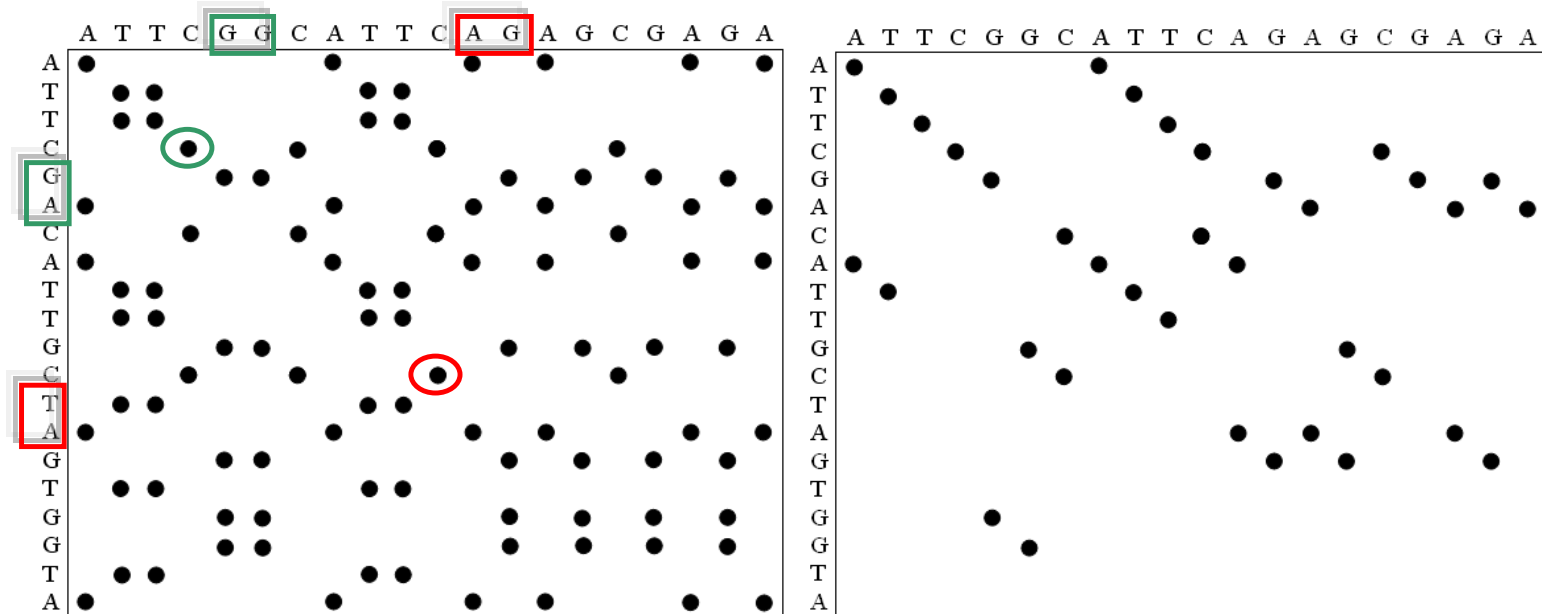
Μείωση θορύβου στα dot plots (1/3)



Μείωση θορύβου στα dot plots (2/3)



Μείωση θορύβου στα dot plots (3/3)



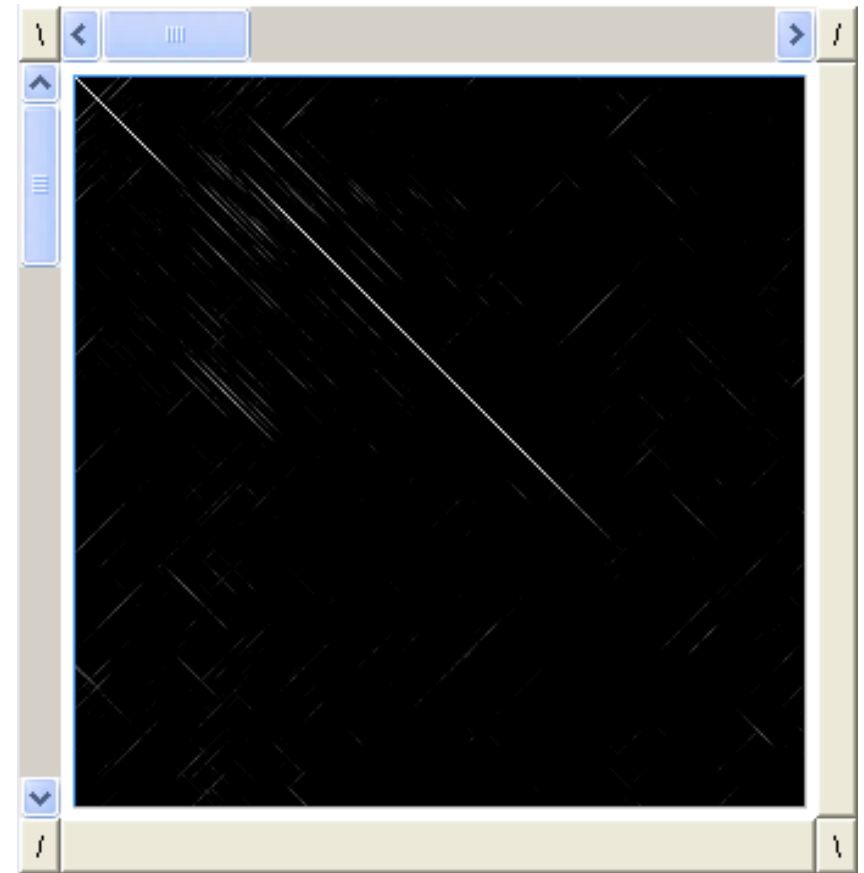
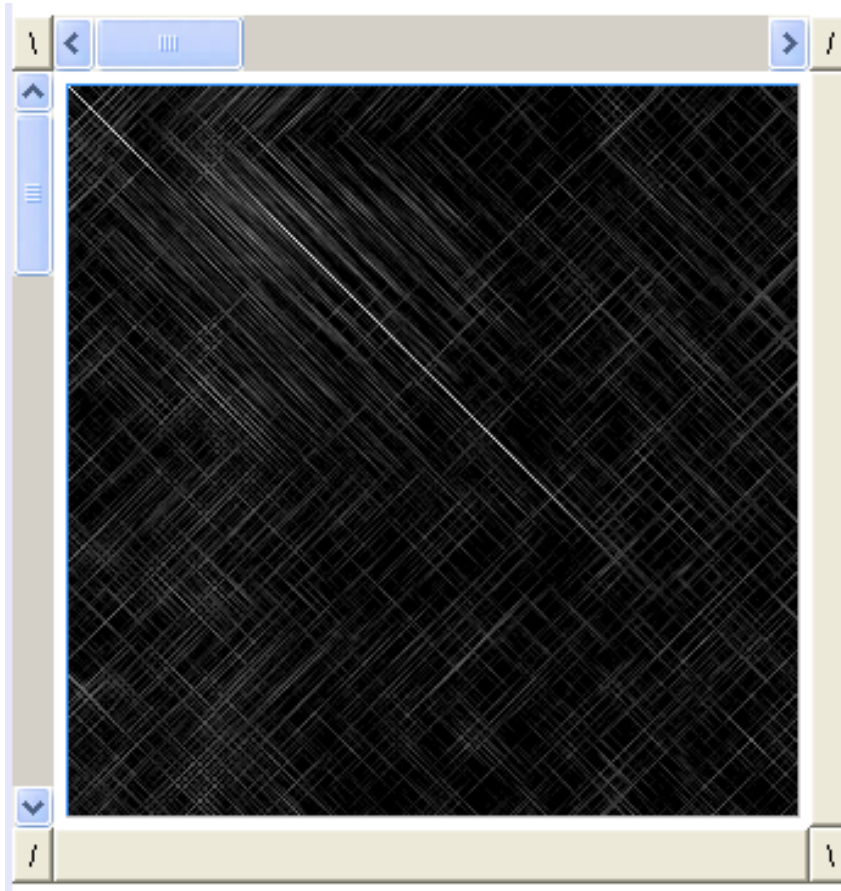
Window: 2, Stringency: 1

Window size: Number of nucleotides compared each time

Stringency: The minimum number of nucleotides in the window must be match, so the dot can be placed



Μείωση θορύβου στα dot plots (Dotlet)



Πληροφορίες από τα dot plots

1. Περιοχές ομοιότητας:

- Διαγώνιος.

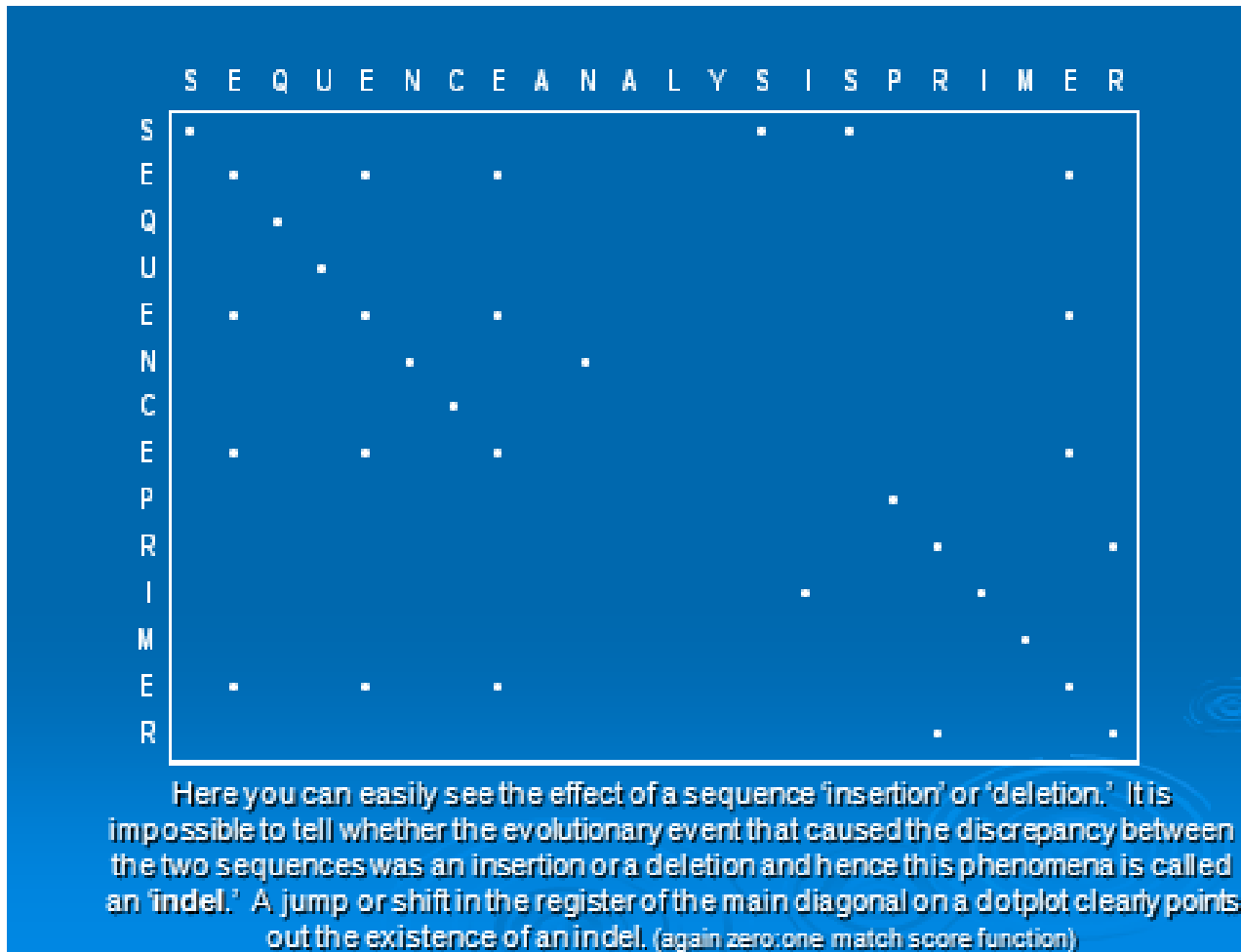
2. Προσθήκες (insertions) / εξαλείψεις (deletions).

3. Επαναλαμβανόμενες περιοχές:

- Ορθές.
- Αντιστρέψιμες (αντίστροφη φορά).

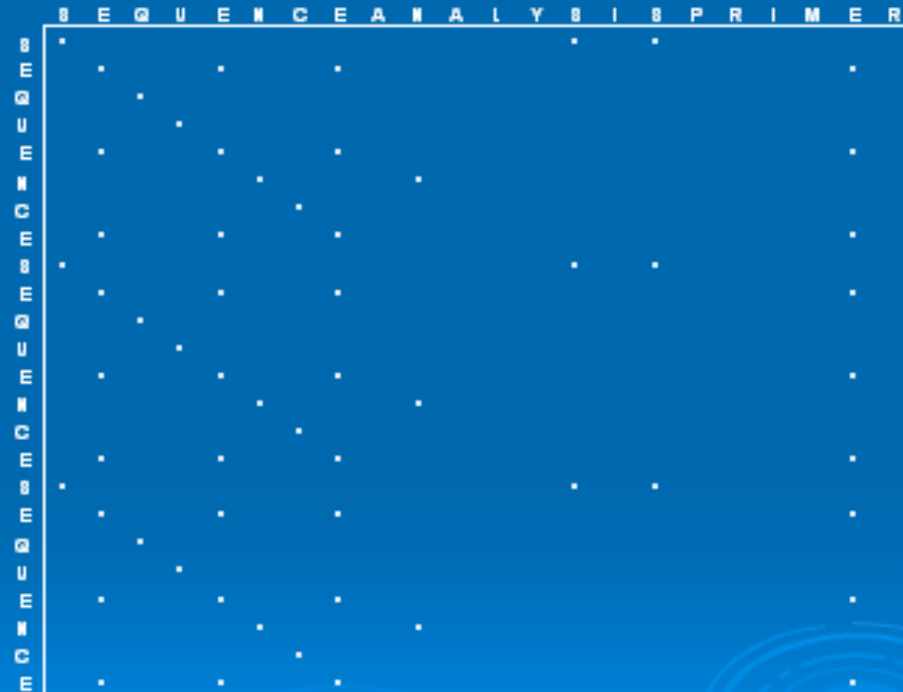


Check out the «mutated» inter-sequence comparison below:



Easy to visualize with dot matrix

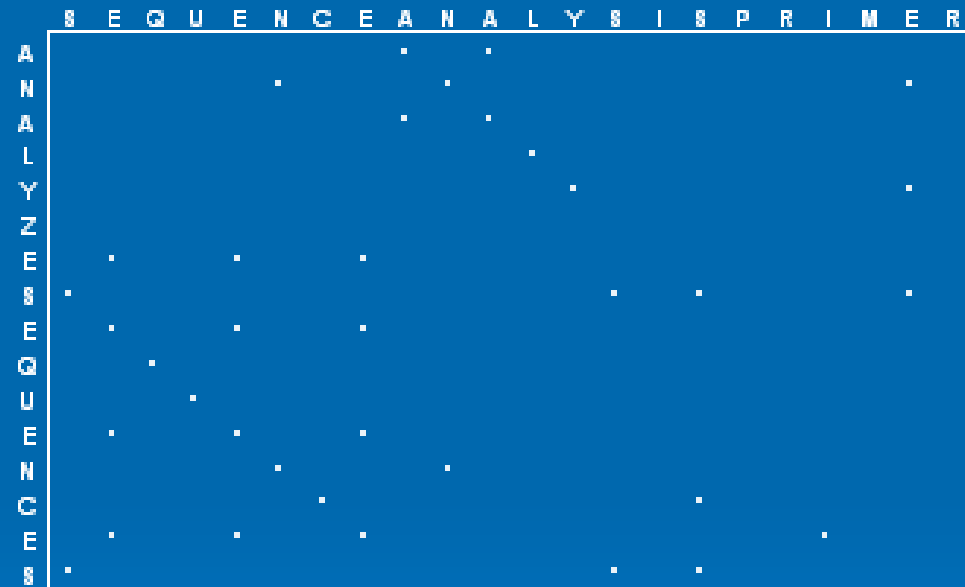
Another phenomenon that is very easy to visualize with dot matrix analysis are duplications or direct repeats. These are shown in the following example:



The 'duplication' here is seen as a distinct column of diagonals; whenever you see either a row or column of diagonals in a dotplot, you are looking at direct repeats.



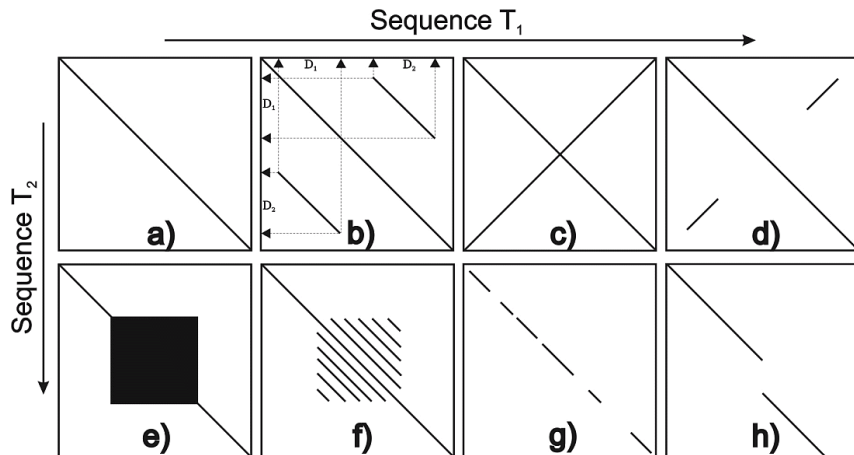
Now consider the more complicated «mutation» in the following comparison:



Again, notice the diagonals. However, they have now been displaced off of the center diagonal of the plot and, in fact, in this example, show the occurrence of a 'transposition.' Dot matrix analysis is one of the only sensible ways to locate such transpositions in sequences. Inverted repeats still show up as perpendicular lines to the diagonals, they are just now not on the center of the plot. The 'deletion' of 'PRIMER' is shown by the lack of a corresponding diagonal.



Dot plots - Παράδειγμα



- a) Απόλυτη ομοιότητα.
- b) Επαναλαμβανόμενες περιοχές.
- c) Ολική παλινδρόμηση.
- d) Μερική παλινδρόμηση.
- e) Επαναλαμβανόμενο σύμβολο (αμινοξύ ή νουκλεοτίδιο) και στις δύο αλληλουχίες.
- f) Επαναλαμβανόμενες περιοχές και στις δύο αλληλουχίες.
- g) Διακοπές - Κοινός πρόγονος.
- h) Προσθήκη στην αλληλουχία 1 ή εξάλειψη στην αλληλουχία 2.



Size of the search space (1/2)

- Size of the search space \sim query sequence length (n) \cdot the sum of the lengths of the sequences in the database (m), **$N=n \cdot m$** .
- **Size of the search space = $N \cdot K$** , where K : Altschul coefficient.
- **Example:** Calculate the search space for a protein of 235a.a. length against a protein database of size $m=12496420$ a.a. The value of K is given equal to 0.13.
- Size of the search space = $0.13 \cdot 235 \cdot 12,496,420 = 0.38$ billion.
- In this case, a bit score of 30 (which correspond to a space of $2^{30} = 1$ billion) may have occurred by chance alone.



Converting to Bit - scores

$$S' = \frac{\lambda \cdot S - \ln K}{\ln 2}$$

$$\Rightarrow S' \cdot \ln 2 = \lambda \cdot S - \ln K$$

$$\Rightarrow \ln K = \lambda \cdot S - S' \ln 2$$

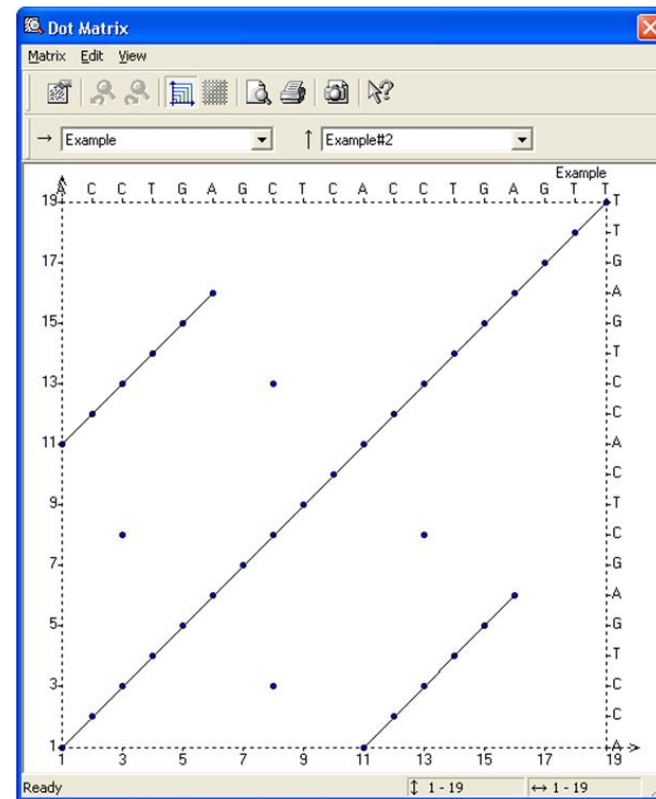
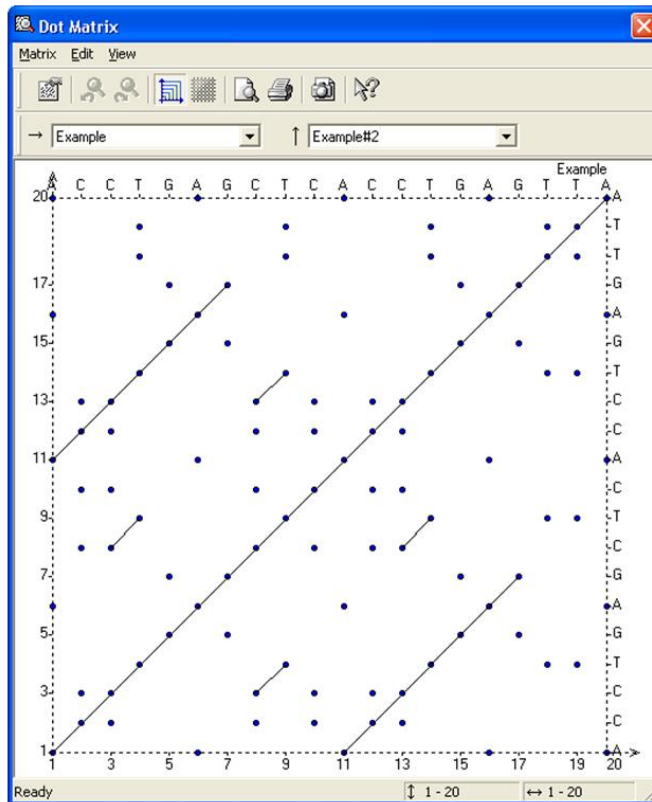
$$\Rightarrow K = e^{\lambda S - S' \ln 2}$$

$$E = K \cdot m \cdot n \cdot e^{-\lambda S}$$

$$E = \frac{m \cdot n \cdot e^{-\lambda S} \cdot e^{\lambda S}}{e^{S' \cdot \ln 2}} = m \cdot n \cdot 2^{-S'}$$



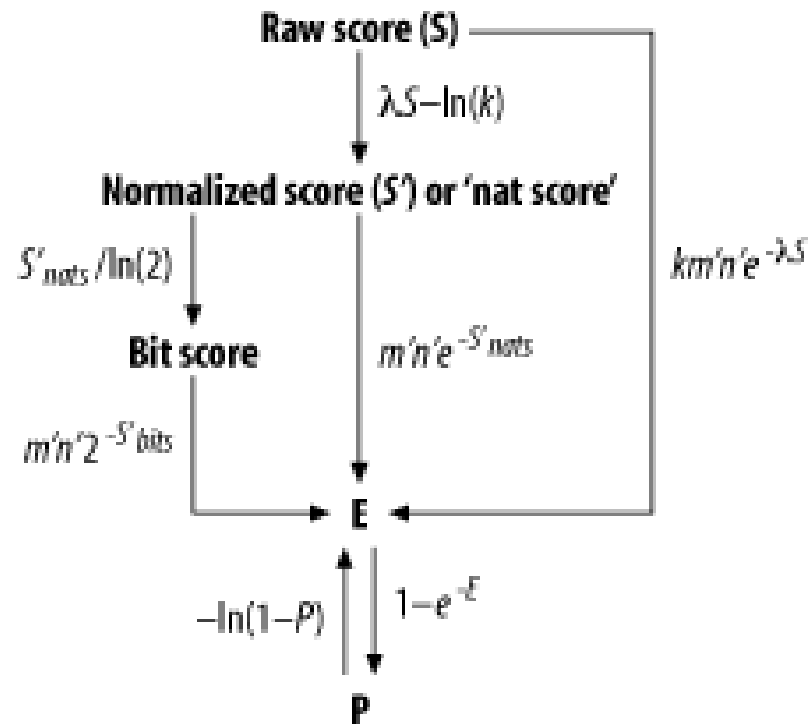
Μείωση θορύβου στα Dot Plots (1/2)



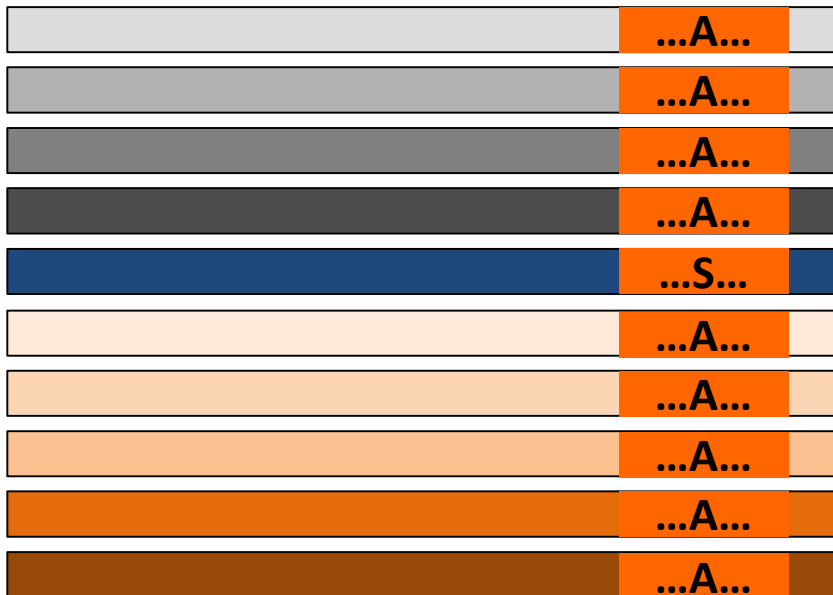
Self alignment of ACCTGAGCTCACCTGAGTTA



Μείωση θορύβου στα Dot Plots (2/2)



Ερμηνεία τιμής του πίνακα BLOSUM - Παράδειγμα



- AA pairs: $8+7+\dots+2+1=36$ possible AA pairs, f_{AA} .
- AS pairs: 9 possible AS pairs, f_{AS} .
- Η συχνότητα εμφάνισης του ζεύγους AA:
 $q_{AA}=f_{AA} / f_{AA}+f_{AS} = 36/(36+9)=0.8$.
- Η συχνότητα εμφάνισης του ζεύγους AS:
 $q_{AS}=f_{AS} / f_{AA}+f_{AS} = 9/(36+9)=0.2$.
- Η αναμενόμενη συχνότητα το A να βρίσκεται σε ζεύγος: $p_A=q_{AA}+q_{AS}/2=0.8+0.2/2=0.9$.
- Η αναμενόμενη συχνότητα το S να βρίσκεται σε ζεύγος: $p_S=q_{AS}/2=0.2/2=0.1$.
- Η αναμενόμενη συχνότητα εμφάνισης του ζεύγους AA: $e_{AA}=p_A \times p_A = 0.9 \times 0.9=0.81$.
- Η αναμενόμενη συχνότητα εμφάνισης του ζεύγους AS: $e_{AS}=2 \times p_A \times p_S=0.18$.
- Η τιμή στον πίνακα για το ζεύγος AA:
 $e_{AA}/q_{AA}=0.8/0.81=0.988$.
- Η τιμή στον πίνακα για το ζεύγος AS:
 $e_{AS}/q_{AS}=0.18/0.2=1.111$.
- Οι τιμές μετατρέπονται σε λογαρίθμους με βάση το 2 και πολλαπλασιάζονται με το 2.



Size of the search space (2/2)

- **Bit-score S'** : A log-scaled version of a score, $\log_2(\text{score})=S'$
- Sequence: m

Database of protein sequences: Total length n .

Size of the search space $\sim m \cdot n$.

Size of the search space $=K \cdot m \cdot n$, K :

- **Example:**

Protein database with sequences of total length: 12496420 a.a.

Sequence: 235 a.a. length

For protein database, $K=0.13$

Size of the search space:



Significance of a local alignment

- Altschul and Gish (1996) have provided estimations of **K=0.09** and **$\lambda=0,229$** for PAM250 matrix, for a typical amino acid distribution and for an alignment score based on using a very high gap penalty.

$$S' = \lambda S - \ln(Kmn) = 0.229 \cdot 73 - \ln(0.09 \cdot 250 \cdot 250)$$

$$S' = 16.72 - 8.63 = 8.09 \text{ bits}$$

$$P(S' \geq 8.55) = e^{-8.09} = 3.2 \cdot 10^{-4}$$

- ***Note that the calculated S' of 8.09 bits in the previous step is approximately the same as the 9 bits calculated by the simpler method***



Relative Entropy

$$H = \sum_{i \geq j} q_{ij} S_{ij}$$

- Indicates power of scoring scheme to distinguish from “background noise” (i.e., randomness).
- **Can use H to compare different scoring matrices.**
- Relative entropy of a random alignment should be negative.
- **Scores can be related to biology:**
 - negative=dissimilarity,
 - zero=indifference,
 - positive=similar.



Στατιστική σημαντικότητα στοίχισης αλληλουχιών χωρίς κενά

- The raw score of an alignment is determined as the sum of the log-odds scores and the gap penalties and is given in the units of log to the base 2 (\log_2) or bits.
- **Bit-score S'** : A log-scaled version of a score, $\log_2(\text{score})=S'$.
- Αν μετατρέψουμε τις τιμές σε $\log_2 x$ έχουμε bits πληροφορίας.
- Όταν χρησιμοποιούμε bit score system, το **κατώφλι της στατιστικής σημαντικότητας** δίνεται από τη σχέση: $\log_2(m \cdot n)$.



Στατιστική σημαντικότητα στοίχισης αλληλουχιών χωρίς κενά - Παράδειγμα (1/2)

- 2 sequences, each 250 a.a. long
- Significance cut off: $\log_2(250*250)=16$ bit
- F W L E V E G N S M T A P T G
- F W L D V Q G D S M T A P A G

Using **PAM250** the score is calculated:

- $S=9+17+6+3+4+2+5+2+2+6+3+2+6+1+5=73$
- **S=73**



Στατιστική σημαντικότητα στοίχισης αλληλουχιών χωρίς κενά - Παράδειγμα (2/2)

- S is in $10^* \log_{10} x$.
- Convert S to a bit score, i.e. calculate $\log_2 x$.

$$S = 10 \log_{10} x \Rightarrow \log_{10} x = \frac{S}{10}$$

$$\Rightarrow \frac{\log_2 x}{\log_2 10} = \frac{S}{10}$$

$$\Rightarrow \log_2 x = \frac{S}{10} \log_2 10$$

$$\Rightarrow \log_2 x \approx \frac{1}{3} S$$

$$\Rightarrow S' = \frac{1}{3} S$$



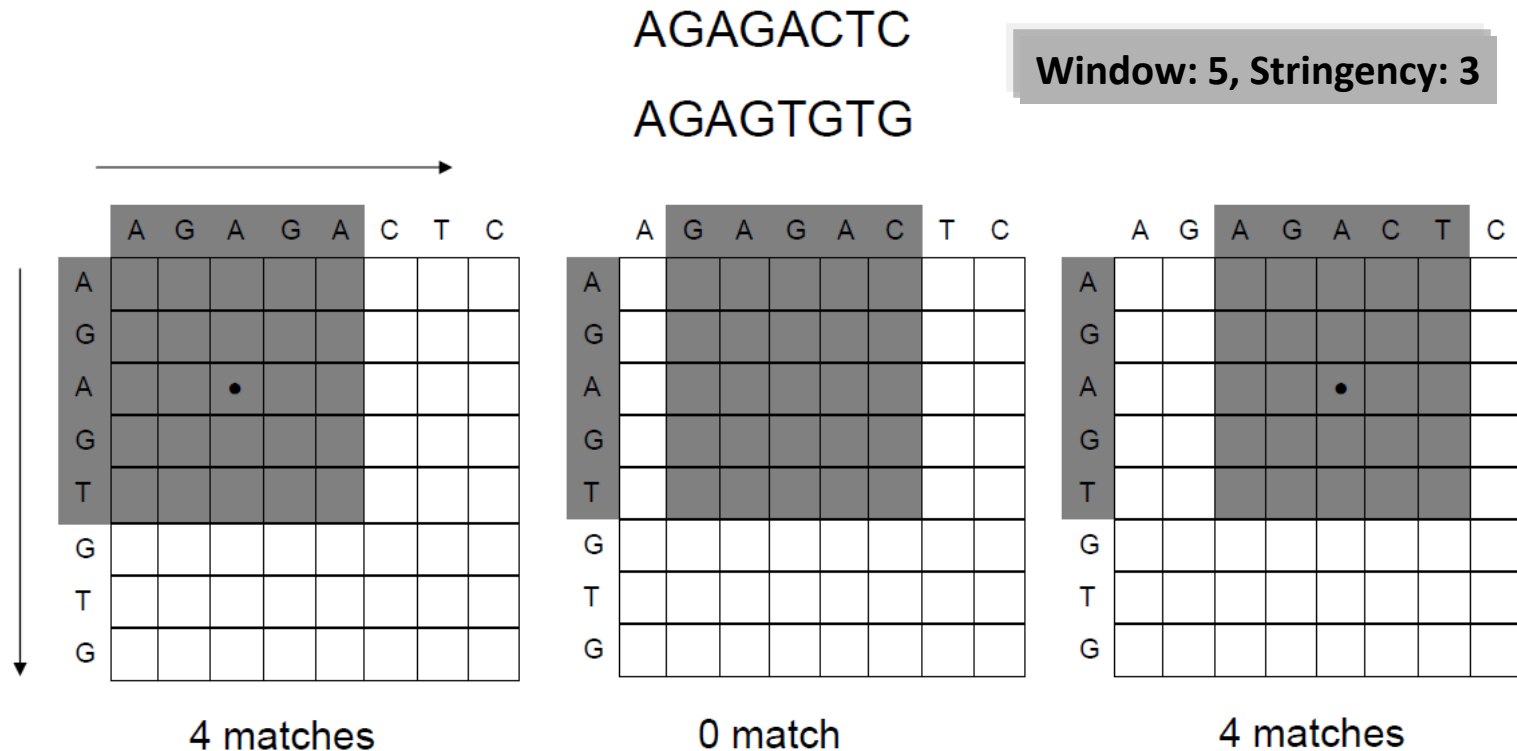
Στατιστική σημαντικότητα στοίχισης αλληλουχιών χωρίς κενά - Παράδειγμα

$$S' = \frac{1}{3} S = \frac{1}{3} 73 = 24.333bits$$

- Significance cut-offs
($\log_2(m \cdot n)$) = 16 bits.
- Το S' είναι **9 bits μεγαλύτερο** από το κατώφλι της στατιστικής σημαντικότητας, επομένως η στοίχιση στο συγκεκριμένο σημείο των αλληλουχιών είναι **στατιστικά σημαντική**.



Μείωση θορύβου στα dot plots (1/2)

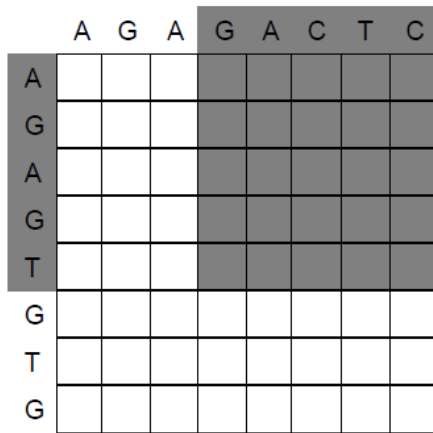


Window size: Number of nucleotides compared each time

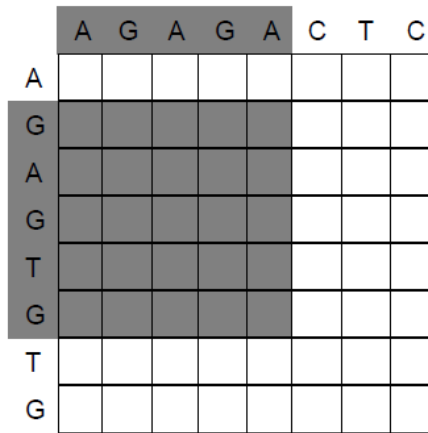
Stringency: The minimum number of nucleotides in the window must be match, so the dot can be placed



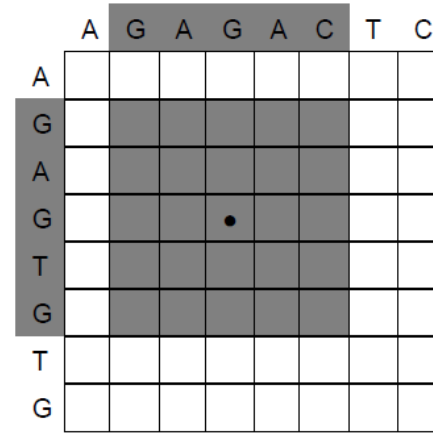
Μείωση θορύβου στα dot plots (2/2)



0 match

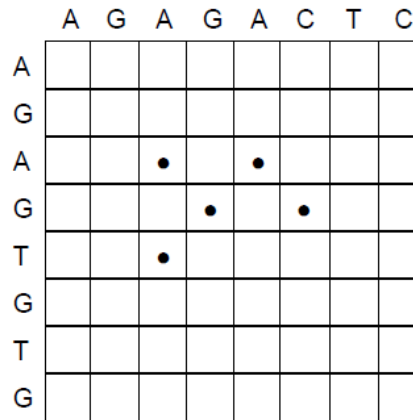


0 match



3 matches

Final answer →



Window: 5, Stringency: 3



Τέλος Ενότητας



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Κοινωνικό Ταμείο



Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



Σημείωμα Αναφοράς

- Copyright Πανεπιστήμιο Δυτικής Μακεδονίας, Τμήμα Μηχανικών Πληροφορικής και Τηλεπικοινωνιών, Αγγελίδης Παντελής. «**Βιοπληροφορική**». Έκδοση: 1.0. Κοζάνη 2015. Διαθέσιμο από τη δικτυακή διεύθυνση: <https://eclass.uowm.gr/courses/ICTE102/>



Σημείωμα Αδειοδότησης

Το παρόν υλικό διατίθεται με τους όρους της άδειας χρήσης Creative Commons Αναφορά, Όχι Παράγωγα Έργα Μη Εμπορική Χρήση 4.0 [1] ή μεταγενέστερη, Διεθνής Έκδοση. Εξαιρούνται τα αυτοτελή έργα τρίτων π.χ. φωτογραφίες, διαγράμματα κ.λ.π., τα οποία εμπεριέχονται σε αυτό και τα οποία αναφέρονται μαζί με τους όρους χρήσης τους στο «Σημείωμα Χρήσης Έργων Τρίτων».



[1] <http://creativecommons.org/licenses/by-nc-nd/4.0/>

Ως Μη Εμπορική ορίζεται η χρήση:

- που δεν περιλαμβάνει άμεσο ή έμμεσο οικονομικό όφελος από την χρήση του έργου για το διανομέα του έργου και αδειοδόχο
- που δεν περιλαμβάνει οικονομική συναλλαγή ως προϋπόθεση για τη χρήση ή πρόσβαση στο έργο
- που δεν προσπορίζει στο διανομέα του έργου και αδειοδόχο έμμεσο οικονομικό

Διατήρηση Σημειωμάτων

Οποιαδήποτε αναπαραγωγή ή διασκευή του υλικού θα πρέπει να συμπεριλαμβάνει:

- το Σημείωμα Αναφοράς
- το Σημείωμα Αδειοδότησης
- τη δήλωση Διατήρησης Σημειωμάτων
- το Σημείωμα Χρήσης Έργων Τρίτων (εφόσον υπάρχει)

μαζί με τους συνοδευόμενους
υπερσυνδέσμους.



Σημείωμα Χρήσης Έργων Τρίτων

Το Έργο αυτό κάνει χρήση των ακόλουθων έργων:

Εικόνες:

- <http://blog.com.mk/send/121903>
- <http://foter.com/Cmyk/>
- <http://www.ebi.ac.uk/training/online/course/introduction-protein-classification-ebi/what-are-protein-signatures/signature-types/what-ar-0>
- http://contentinacottage.blogspot.ca/2012_01_29_archive.html
- <https://www.cartoonstock.com/>

