

UNIVERSITÉ CATHOLIQUE DE LOUVAIN

MASTER EN SCIENCES DE L'INFORMATIQUE

INTELLIGENCE ARTIFICIELLE

LGBIO2010: Bioinformatics

Author:

Denis GENON

Contents

1 Sequences Statistics	3
1.1 Simple Genome Statistics	3
1.1.1 GC-content	3
1.1.2 Dimer frequencies	3
1.2 Open Reading Frames	4
1.2.1 Algorithm to find ORF	4
1.2.2 Significance assessment	5
2 Pairwise Alignments	6
2.1 The problem	6
2.1.1 Gap penalties	6
2.1.2 Number of alignments	7
2.2 Global alignment: Needle-Wunsch	7
2.3 Local alignment: Smith-Waterman	7
2.4 Variants	8
2.4.1 Semi-global alignment	8
2.4.2 Local alignment with repeats	8
2.4.3 Affine gap penalty	9
2.5 Significance assessment	9
3 HMMs	10
3.1 An example	10
3.1.1 Baseline approach	10
3.1.2 Alternative approach	11
3.2 Hidden Markov model	13
3.2.1 Path likelihood	14
3.2.2 Probability of generating a sequence	14
3.2.3 Most likely state sequence	15
3.2.4 Sequence likelihood	15
3.3 The learning problem	15
3.3.1 Supervised learning	16
3.3.2 Unsupervised learning	16
3.4 Concrete example	17
4 Multiple Alignments	17
4.1 The problem	17
4.1.1 Scoring	18
4.2 Dynamic programming	18
4.2.1 Progressive alignment methods (greedy heuristic)	19

5 Profile HMMs	20
5.1 Motivations	20
5.2 Position specific scoring matrices	21
5.3 Full profile HMMs	22
5.3.1 Adding insert states	22
5.3.2 Adding delete states	22
5.3.3 A full profile HMM	23
5.3.4 Deriving a pHMM from a multiple alignment (Supervised) .	23
5.3.5 Unsupervised learning	24
5.3.6 Matching a sequence to a pHMM a.k.a. viterbi recurrence .	24
5.3.7 For non-global alignments	25
6 Large Scale Gene Expression Analysis	26
6.1 Introduction	26
6.2 Preprocessing	26
6.2.1 Summarization	26
6.2.2 Feature normalization	26
6.2.3 Distance between expression values	27
6.3 Unsupervised learning	28
6.4 Supervised learning	29
6.4.1 Filters	30
6.4.2 Wrappers	33
6.4.3 Embedded approaches	34
7 Inference of Gene Regulatory Network	35

1 Sequences Statistics

1.1 Simple Genome Statistics

1.1.1 GC-content

GC-content is the percentage of G or C in a DNA sequence.

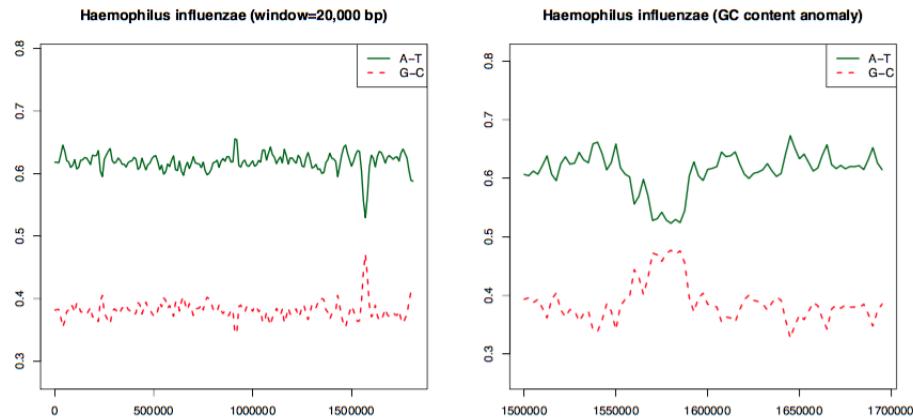


Figure 1: GC-Content. Anomaly due to an ancient insertion of viral DNA. AT-rich regions denature at lower temperature. The ability to quickly denature DNA facilitates the insertion in the bacterial cell being infected.

1.1.2 Dimer frequencies

Some dimers are particularly frequent but it is due to have each nitrogenous bases frequent. Need a background model.

- Estimate the probability of each nucleotide independently
 $\hat{P}(A) = f(A) = \frac{\text{number of A's}}{\text{sequence length}} = 0.3102$
- Use those estimates in a random generative model
AAGTTGACATAATTGCT ...
- The expected frequency of a dimer XY :
$$E[f(XY)] = \hat{P}(X)\hat{P}(Y) = f(X)f(Y)$$

Figure 2: Multinomial background model.

With the multinomial background model, we can compute the **odd ratios**: ratio between observed frequency and expected frequency: $\frac{f(XY)}{E[f(XY)]} = \frac{f(XY)}{f(X)f(Y)}$.

The multinomial model is equivalent to a random permutation of the original sequence: $\frac{f(XY)}{f_{\text{random}}(XY)}$.

We can generalize odd ratios to **k-mers**.

$$\frac{f(X_1 \dots X_k)}{E[f(X_1 \dots X_k)]} = \frac{f(X_1 \dots X_k)}{\prod_{i=1}^k f(X_i)}$$

Figure 3: k-mers. Useful when looking for frequent patterns of k consecutive characters. But are they informative ? Need a more sophisticated background model, a statistical test to assess signifiance and biological validation.

1.2 Open Reading Frames

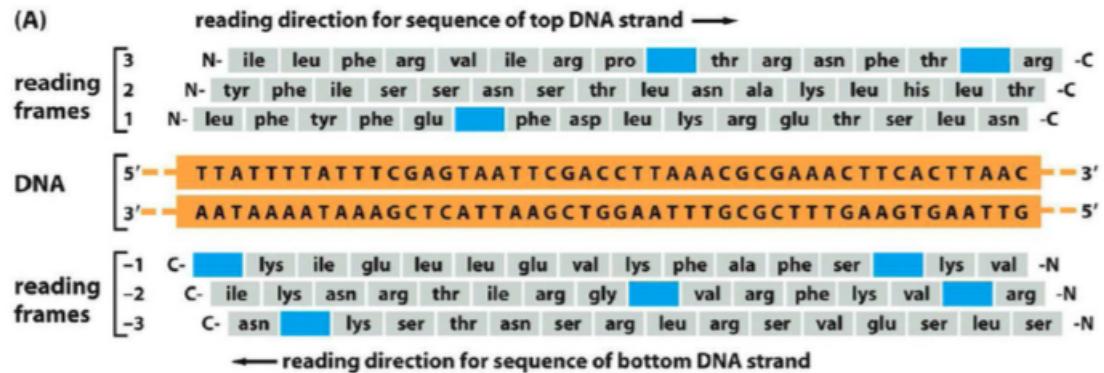


Figure 4: Reading frames.

1.2.1 Algorithm to find ORF

repeat along the sequence

- ➊ look for a first START
(or the next START after a STOP on the same frame)
- ➋ look for the next STOP on the **same reading frame**

Figure 5: Need to consider each reading frames (6). After a START, you may find other codons for Met before a STOP: not true START. An ORF is a longest strech of DNA between a START and a STOP, without being interrupted by another STOP on the same frame.

The problem is to have to be sure that the ORF is a coding gene.

- The DNA found between STAT and STOP codon might be due to chance: need a **statistical test procedure**;
- The ORF might be there but the gene not expressed (trace of past or regulations at transcription/translation levels): need biological validation;
- Some gene sequences do not strictly follow the standard ORF structure: need biological validation.

1.2.2 Significance assessment

- Null hypothesis: $H_0 : \mu_2 = \mu_1$
 - Alternative hypothesis: $H_a : \mu_2 \neq \mu_1$
 - Test statistics :
$$T = m_2 - m_1 = \hat{\mu}_2 - \hat{\mu}_1$$
 - It is known that T approximately follows a Student t distribution
- $$T = \frac{(\hat{\mu}_2 - \hat{\mu}_1) - (\mu_2 - \mu_1)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$
- Fix a significance threshold α (e.g. 5%)

$$T = \frac{(6.6 - 3.4) - (0)}{\sqrt{\frac{4.3}{5} + \frac{1.3}{5}}} = 3.0237 > t_{0.975} = 2.43 \Rightarrow \text{reject } H_0$$

⇒ claim the drugs are not equally effective

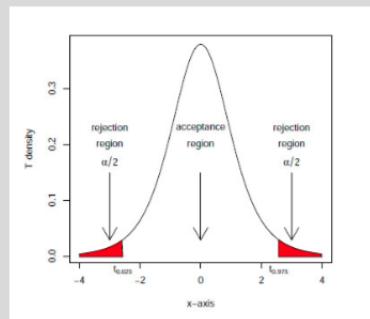


Figure 6: Statistical test.

To avoid fixing a significance threshold α we can compute the **p-value** of the test.

- p-value = the smallest α that would lead to reject the test
- here $T = 3.0237 \Rightarrow$ p-value = 0.02229
 - ▶ the true means could still be equal (= the drugs could still be equally effective), but the probability of our conclusion to be wrong is 2.2%
- the lower the more significant the result

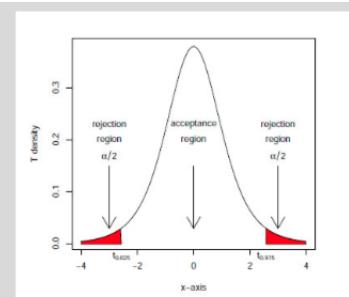


Figure 7: p-value.

We can apply this to ORF.

Hypothesis

- Significant ORFs in an actual sequence should be **longer** than ORFs observed **by chance**
- The NULL model (= control) is typically made of a **random permutation** of the original sequence

Algorithm

- ① Find all ORFs in a random permutation of the original sequence
- ② Report the length distribution of random ORFs
- ③ Accept as significant ORFs, any ORF in the original sequence longer than a prescribed threshold
 - ▶ the **maximal ORF length** observed at random
 - ▶ the **99% percentile** of the random ORF length distribution
⇒ permutation test with a ***p*-value = 1%**

Figure 8: Algorithm to assess significance of ORF.

2 Pairwise Alignments

2.1 The problem

Find a best way to align two sequences including matches, substitutions and possible gaps.

2.1.1 Gap penalties



Linear gap penalty

$$\gamma(g) = -dg$$

***g*:** number of consecutive gaps

***d*:** gap penalty (e.g. 8)

Affine gap penalty

$$\gamma(g) = -d - e(g - 1)$$

***g*:** number of consecutive gaps

***d*:** gap open penalty (e.g. 12)

***e*:** gap extension penalty (e.g. 2)

Figure 9: Gap penalties. Affine is more relevant.

2.1.2 Number of alignments

The total number alignments is huge ($\frac{4^n}{\sqrt{\pi} \cdot n}$). We need dynamic programming.

- We look for maximal cumulative score;
- An optimal global alignment is made of optimal alignments between subsequences: decomposition of sub-problems computed only once.

2.2 Global alignment: Needle-Wunsch

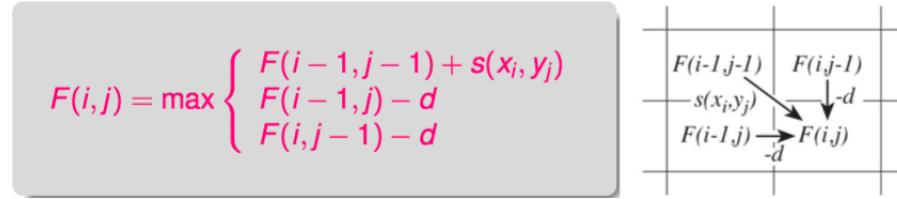


Figure 10: $F(n, m)$ is final alignment score. Need to follow backpointers. $O(n^2)$.

2.3 Local alignment: Smith-Waterman

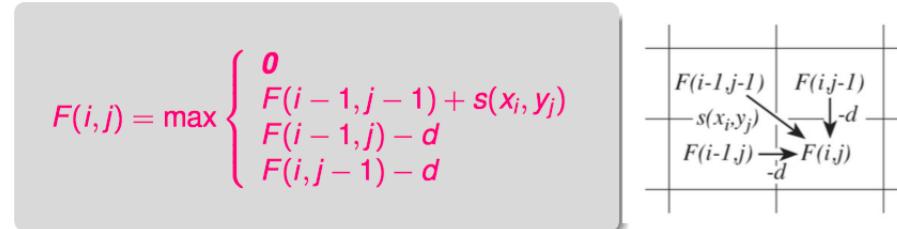


Figure 11: Find maximal $F(i, j)$ and follow backpointers till 0. Need negative scores for strong mismatches. $O(n^2)$.

2.4 Variants

2.4.1 Semi-global alignment

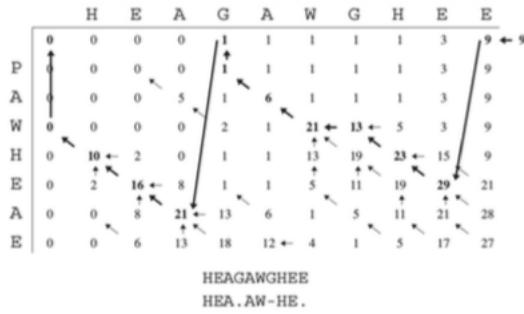
	H	E	A	G	A	W	G	H	E	E
P	0	0	0	0	0	0	0	0	0	0
A	0									
W	0									
H	0									
E	0									
A	0									
E	0									

Do not penalize gaps at the beginning of either sequences

- Initialize $F(i, 0) = F(0, j) = 0 ; 0 \leq i \leq n ; 0 \leq j \leq m$
- Compute the global recurrence

$$F(i, j) = \max \begin{cases} F(i - 1, j - 1) + s(x_i, y_j) \\ F(i - 1, j) - d \\ F(i, j - 1) - d \end{cases}$$

2.4.2 Local alignment with repeats



$$F(0, 0) = 0$$

$$F(i, 0) = \max \begin{cases} F(i - 1, 0) \\ F(i - 1, j) - T \end{cases} \quad \forall j = 1, \dots, m$$

$$F(i, j) = \max \begin{cases} F(i, 0) \\ F(i - 1, j - 1) + s(x_i, y_j) \\ F(i - 1, j) - d \\ F(i, j - 1) - d \end{cases}$$

Figure 12: Look for several matches (subsequences) of y into x . Scoring higher than T .

2.4.3 Affine gap penalty

Instead of a single state $F(i, j)$, one distinguishes 3 states

$I \ G \ A \ x_i$

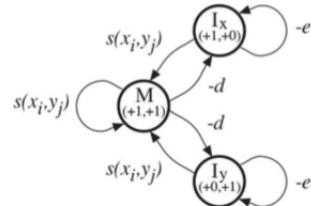
$L \ G \ V \ y_j$

$A \ I \ G \ A \ x_i$

$G \ V \ y_j \ - \ -$

$G \ A \ x_i \ - \ -$

$S \ L \ G \ V \ y_j$



M state : x_i aligned to y_j
 I_x state : x_i aligned to a gap
 I_y state : y_j aligned to a gap

$$M(i, j) = \max \begin{cases} M(i-1, j-1) + s(x_i, y_j) \\ l_x(i-1, j-1) + s(x_i, y_j) \\ l_y(i-1, j-1) + s(x_i, y_j) \end{cases}$$

$$l_x(i, j) = \max \begin{cases} M(i-1, j) - d \\ l_x(i-1, j) - e \end{cases}$$

$$l_y(i, j) = \max \begin{cases} M(i, j-1) - d \\ l_y(i, j-1) - e \end{cases}$$

2.5 Significance assessment

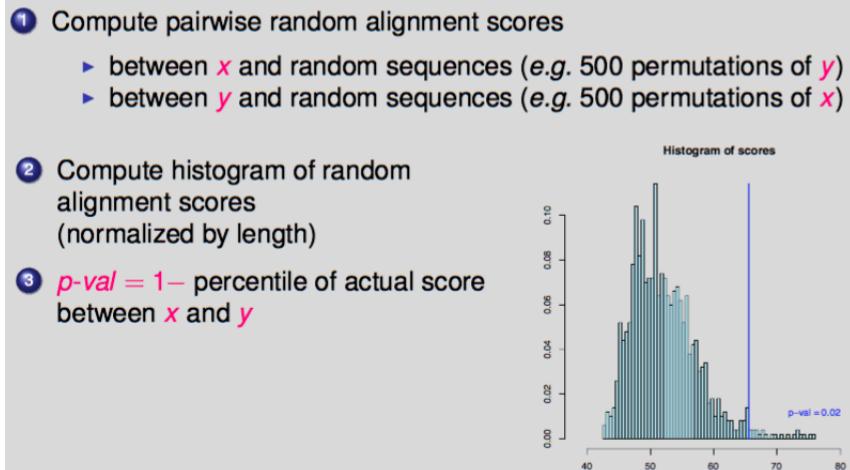


Figure 13: Procedure.

The actual distribution of alginment scores is called **extreme value distribution**.

- Probability of a score larger than S

$$P(x > S) = 1 - e^{(-Kmn e^{-\lambda S})}$$
for some constants K and λ
- Alignment softwares include fitted values of K and λ for a wide range of substitution scoring matrices
- The lack of normality is due to non-independence between possible starting points of matches

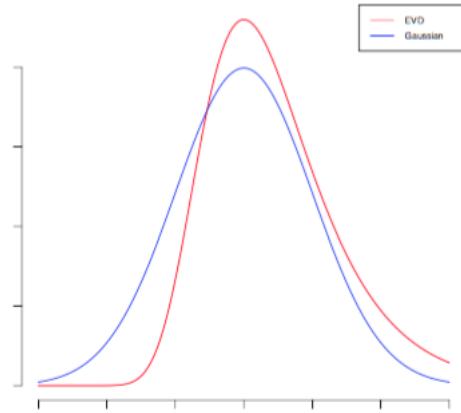


Figure 14: K is a multiplicative factor correcting non independance of possible starting poins for matches. λ is a scale parameter.

3 HMMs

3.1 An example

An motivating example is to determine coding or non-coding fragments.

3.1.1 Baseline approach

- ① Find all ORFs in the original sequence
- ② Find all ORFs in random permutation(s) of the original sequence
- ③ Accept as significant ORFs, any ORF in the original sequence longer than a prescribed length: e.g. the 99% percentile of random ORF lengths (p -value = 1%)

Limitations

- Coding vs non-coding fragments are not only characterized by their length
- Intron-exon boundaries also need to be identified for eukaryotes

3.1.2 Alternative approach

Learning

Given some **labeled** data

- estimate a statistical model M_+ for **coding** fragments (or exons)
- estimate a statistical model M_- for **non-coding** fragments (or introns)

Segmentation

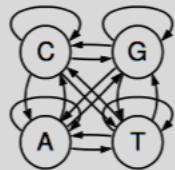
For any new sequence to analyze

- look at a subsequence x defined by a sliding window of size L
- compute log-odds ratio

If $\log \frac{P(x|M_+)}{P(x|M_-)} = \log \frac{\prod_{i=1}^L P(x_i|M_+)}{\prod_{i=1}^L P(x_i|M_-)} = \sum_{i=1}^L \log \frac{P(x_i|M_+)}{P(x_i|M_-)} > 0$
 then decide x is **coding**
 else decide x is **non-coding**

3.1.2.1 Multinomial model

An equivalent representation



- each symbol is generated on a specific **state**
- all **transition probabilities** pointing to the same state are **equal**
- one such model for each segment type (e.g. coding vs non-coding)

Figure 15: Two hypothesis: the frequencies are the same along the different parts of the sequence and each symbol does not depend of surrounding symbols.

3.1.2.2 Markov chain



- Each transition probability $P(x_i|x_{i-1}, M)$ is now specific to a dimer

$$\hat{P}(x_i|x_{i-1}, M) = \frac{f(x_{i-1}x_i)}{f(x_{i-1})} \text{ from segments modeled by } M$$

- Log-odds computation

$$\log \frac{P(x|M_+)}{P(x|M_-)} = \sum_{i=1}^L \log \frac{\hat{P}(x_i|x_{i-1}, M_+)}{\hat{P}(x_i|x_{i-1}, M_-)}$$

Figure 16: First order.

- Each transition probability $P(x_i|x_{i-1}, x_{i-2}, M)$ is specific to a 3-mer

$$\hat{P}(x_i|x_{i-1}, x_{i-2}, M) = \frac{f(x_{i-2}x_{i-1}x_i)}{f(x_{i-2}x_{i-1})} \text{ from segments modeled by } M$$

- Log-odds computation

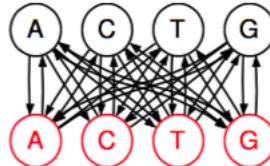
$$\log \frac{P(x|M_+)}{P(x|M_-)} = \sum_{i=1}^L \log \frac{\hat{P}(x_i|x_{i-1}, x_{i-2}, M_+)}{\hat{P}(x_i|x_{i-1}, x_{i-2}, M_-)}$$

Figure 17: Second order.

Limitations of the markov chain approach:

- Estimate each MC from well annotated segments;
- Gene length variability: sliding window length is arbitrary and we need to find a length that is more relevant: computation expensive;
- Eukaryotes: at least three models needed (out of gene, introns, exons).

3.1.2.3 Hidden Markov model



...ACGCGTATAATGC...CAT**ATGACGCGT**...AATCCGTAA**CGGTCGAAAAA**...

Learning

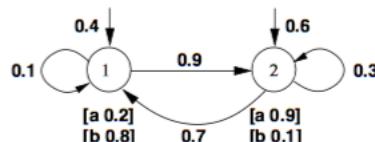
- Define a single model for all possible segments
Note: not all transitions are depicted here
- Estimate all transition probabilities simultaneously

Segmentation

Find the **most likely state sequence** to generate the whole sequence

Figure 18: One state for each segment type (coding vs non coding). Transition probabilities, emissions probabilities. States are hidden but their emissions are observed.

3.2 Hidden Markov model



$$\Sigma = \{a, b\} \quad Q = \{1, 2\}$$

$$A = \begin{bmatrix} 0.1 & 0.9 \\ 0.7 & 0.3 \end{bmatrix} \quad B = \begin{bmatrix} 0.2 & 0.8 \\ 0.9 & 0.1 \end{bmatrix} \quad \pi = [0.4 \ 0.6]$$

A discrete HMM (with state emission)

- Σ is a finite alphabet
- Q is a set of states
- \mathbf{A} a $|Q| \times |Q|$ transition probability matrix ($\sum_{q' \in Q} \mathbf{A}_{qq'} = 1$)
- \mathbf{B} a $|Q| \times |\Sigma|$ emission probability matrix ($\sum_{a \in \Sigma} \mathbf{B}_{qa} = 1$)
- π an initial probability distribution ($\sum_{q \in Q} \pi_q = 1$)

Figure 19: HMM definition.

3.2.1 Path likelihood

The likelihood $P(s, \nu|M)$ of a sequence $s = s_1 \dots s_{|s|}$ along a path or state sequence $\nu = q_1 \dots q_{|s|}$ in a HMM M

$$P(s, \nu|M) = \prod_{i=1}^{|s|} P(s_i, q_i|M) = \pi_{q_1} \mathbf{B}_{q_1 s_1} \prod_{i=2}^{|s|} \mathbf{A}_{q_{i-1} q_i} \mathbf{B}_{q_i s_i}$$

3.2.2 Probability of generating a sequence

The likelihood $P(s|M)$ of a sequence $s = s_1 \dots s_{|s|}$ in an HMM M

$$P(s|M) = \sum_{\nu \in Q^{|s|}} P(s, \nu|M)$$

Figure 20: $O(|Q|^{|s|})$ possible state sequences. Q is states and s is sequence.

3.2.3 Most likely state sequence

$$\nu^* = \operatorname{argmax}_\nu P(s, \nu | M)$$

Auxiliary quantity: $\gamma(k, t) = P(s_1 \dots s_t, \nu_t^* = k | M)$

The probability of a most likely path ν^* reaching state k at step t

Initialization: $\gamma(k, 1) = \pi_k \mathbf{B}_{ks_1}$

Recurrence: $\gamma(k, t) = \sum_l [\gamma(l, t-1) \mathbf{A}_{lk}] \mathbf{B}_{ks_t}$

Termination: $P(s, \nu^* | M) = \sum_l \gamma(l, |s|)$

Figure 21: Viterbi recurrence.

- $P(s, \nu^*)$ gives the probability of the optimal path ν^* ;
- Computation are done with log (because too small values);
- Complexity: $\Theta(m|s|)$. m is number of HMM transitions;
- Path ν^* is an alignment between states and words;
- ν^* can be recovered with the backpointers.

3.2.4 Sequence likelihood

$$P(s | M) = \sum_\nu P(s, \nu | M)$$

Auxiliary quantity: $\alpha(k, t) = P(s_1 \dots s_t, \nu_t = k | M)$

The likelihood of emitting the first t symbols and reaching state k at time t

Initialization: $\alpha(k, 1) = \pi_k \mathbf{B}_{ks_1}$

Recurrence: $\alpha(k, t) = \sum_l [\alpha(l, t-1) \mathbf{A}_{lk}] \mathbf{B}_{ks_t}$

Termination: $P(s | M) = \sum_l \alpha(l, |s|)$

Figure 22: Forward recurrence, same complexity as viterbi recurrence.

3.3 The learning problem

Given an HMM structure and several sequences, you must estimate A , B , π .

3.3.1 Supervised learning

- $\mathbf{B}_{ki} = \frac{f(k,i)}{f(k)}$
 - ▶ $f(k, i)$ = number of times symbol i is observed on state k
 - ▶ $f(k)$ = number of times state k is used
- $\mathbf{A}_{kl} = \frac{f(k,l)}{f(k)}$
 - ▶ $f(k, l)$ = number of times a transition from state k to state l is used
- $\pi_k = \frac{f_1(k)}{\sum_{j=1}^{|Q|} f_1(j)}$
 - ▶ $f_1(k)$ = number of times state k is used as first state
 - (or $f_1(k) \triangleq 1 \Rightarrow \pi_k = \frac{1}{|Q|}$)

Figure 23: Supervised learning.

3.3.2 Unsupervised learning

The sentences are not annotated. There is two algorithm. The first one is **Viterbi training**.

- ➊ Fix initial parameter values $\mathbf{A}^0, \mathbf{B}^0, \pi^0$
- ➋ repeat
 - ➌ compute a most likely path through a Viterbi alignment for each learning sequence given the parameters $\mathbf{A}^i, \mathbf{B}^i, \pi^i$
 - ➍ estimate the emission and transition frequencies from such path(s)
 - ➎ recompute the parameter values $\mathbf{A}^{i+1}, \mathbf{B}^{i+1}, \pi^{i+1}$ from those frequencies
- ➏ till some stopping criterion is met (e.g. max number of iterations)

Figure 24: Viterbi learning.

The second is **Forward-Backward** or **Braaum-Welch** algorithm. Viterbi training is an approximation as it considers that each training sentence is generated along a single path. A more accurate estimation is obtained if one considers all possible paths to generate each sentence.

- actual frequencies are replaced by expected frequencies;

- special case of expectation-maximization (EM) procedure

Viterbi and Baum-Welch training are both sensitive to parameter initialization.

3.4 Concrete example

For CpG islands, we can use a 2-state HMM.

- GC-rich state: higher probability to **emit** G or C
- alternative state: lower probability to **emit** G or C
- **transition probabilities** reflect that it is much more likely (e.g. 0.998) to stay in a given state than switching regime (e.g. 0.002)
- **initial probabilities** reflect the chance to start or not with a CpG island (by default, $\pi_1 = \pi_2 = 0.5$)

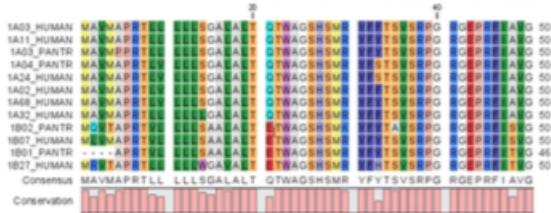
Those probabilities can either be *a priori fixed* or **estimated** through supervised learning, Viterbi or Baum-Welch

4 Multiple Alignments

4.1 The problem

Align three or more homologous sequences (globally or locally).

4.1.1 Scoring



Assumption

- Individual columns are assumed **statistically independent**
- A multiple alignment m with L columns can then be scored as

$$S(m) = G + \sum_{i=1}^L S(m_i)$$
 - ▶ $S(m_i)$ = score for column i
 - ▶ G = score for all gaps in m using linear or affine gap penalties

Figure 25: Assumptions.

Sum of pairs

- Column score:

$$S(m_i) = \sum_{k < i} s(m_i^k, m_i^l)$$

where $s(a, b)$ is given by a substitution scoring matrix
(e.g. PAM or BLOSUM)
- Gap penalty
 - ▶ linear: $s(a, -) = s(-, b) = -d$; $s(-, -) = 0$
 - ▶ affine: all gaps are scored separately

Figure 26: Sum of pairs.

4.2 Dynamic programming

Optimal alignment between k sequences can be computed with DP. With an average length of \bar{n} , the time complexity is $O(2^k \bar{n}^k)$ and the space complexity is $O(\bar{n}^k)$ (hyper cube).

One algorithm which compute optimally the alignment is **MSA** algorithm ($O(k^2 \bar{n}^2)$: 10 sequences of 200-300 residues in reasonable time).

- Compute all pairwise alignments;

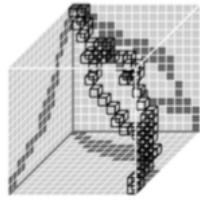


Figure 27: MSA.

- Limits the exploration of the hyper cube to regions consistent with those alignments.

G-MSA is an optimization: 500 sequences of 236 residues in 10 seconds.

4.2.1 Progressive alignment methods (greedy heuristic)

Greedy heuristic algorithms

succession of pairwise alignments

- 2 sequences are aligned first
- a sequence is added to a group of already aligned sequences
 - ▶ compute all pairwise alignments between s and an existing group g of aligned sequences
 - ▶ the highest scoring pairwise alignment determines how the new sequence s is aligned to the group g
- a group g_1 of sequences is aligned to another group g_2 of sequences
 - ▶ all sequence pairs between g_1 and g_2 are tried
 - ▶ the best pairwise alignment determines the alignment of both groups

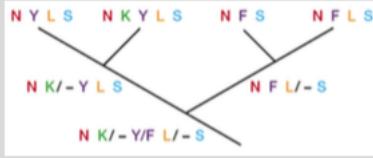
Figure 28: Greedy heuristic.

There are some limitations:

- the degree of sequence conservation at each position should be taken into account;
- mismatches at highly conserved positions should be more penalized;
- the order in which sequences are incorporated in the multiple alignment matters.

4.2.1.1 ClustalW

Main steps

- ① construct a distance matrix of all $\frac{k(k-1)}{2}$ pairwise alignment scores
 - ▶ correct those scores by considering the Kimura evolutionary model (see phylogeny)
- ② build a tree using the neighbor-joining algorithm (see phylogeny)
- ③ use it as a guide tree: progressively align nodes of decreasing similarity
 - ▶ sequence-sequence, sequence-profile and profile-profile alignments

There can be further heuristics:

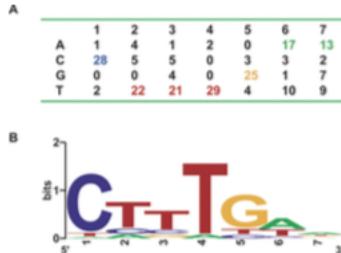
- position specific gap open penalties with decreased penalties whenever other gaps have been found among already aligned sequences;
- gap penalties are also decreased or increased based on a large collection of **structural alignments**;
- the guide tree may be adjusted **on the fly** to defer a low scoring alignment until more profile information has been accumulated.

5 Profile HMMs

5.1 Motivations

- multiple alignments are most often based on pairwise alignments;
- a new sequence x may be only distantly and locally related to each sequence in a known family (biological question): all pairwise alignments between x and each family members will look poor. We need to model **statistical features** shared by the family members;
- computing the alignment between x and a probabilistic model of the family may be much **more efficient computationally**.

5.2 Position specific scoring matrices



- Local model for a window length L and ungapped score matrix from N sequences

$$P(x|M) = \prod_{i=1}^L P(x_i|M) = \prod_{i=1}^L \frac{f(x_i)}{N}$$

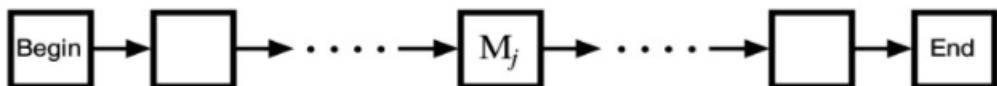
- Log-odds score

$$S = \sum_{i=1}^L \log \frac{P(x_i|M)}{q_{x_i}}$$

► q_{x_i} = the background model (e.g. multinomial model)

- One evaluates the score S between x and M for all positions x_i and a sliding window of size L

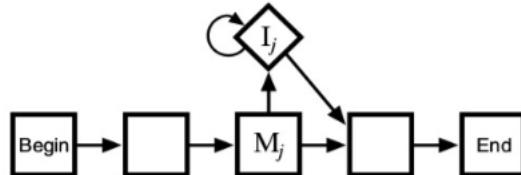
Figure 29: PSSMs are very simple HMMs. $P(x_i|M)$ are emission probabilities on match state. Transition probabilities are equal to 1. But we need to account for possible gap and avoid a prescribed window length.



$$P(x|M) = \prod_{i=1}^L P(x_i|M)$$

5.3 Full profile HMMs

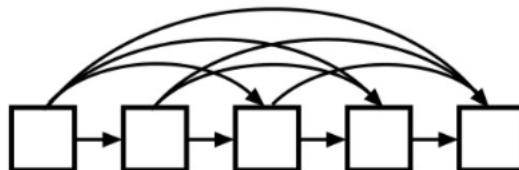
5.3.1 Adding insert states



- to account for **portions of x** that **do not match** the model
- **emission probabilities** on insert states are typically defined through the background model
 - ▶ no contribution to the log-odds score $\log \frac{P(x_i|I_j)}{q_{x_i}} = \log \frac{q_{x_i}}{q_{x_i}} = 0$
- transition probabilities to insert states and back are equivalent to affine gap penalties $\log A_{M_j I_j} + \log A_{I_j M_{j+1}} + (k - 1) \log A_{I_j I_j}$

5.3.2 Adding delete states

- **portions of the model M** that are not matched by any residue x_i could be modeled by **skipping transitions**



- to allow arbitrary long gaps it is more convenient to introduce **delete states** which are **silent states**

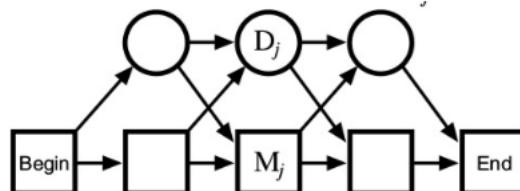


Figure 30: More convenient because less transitions.

5.3.3 A full profile HMM

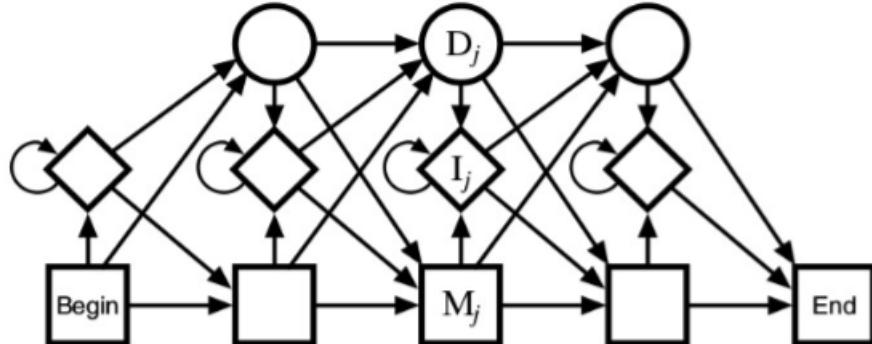
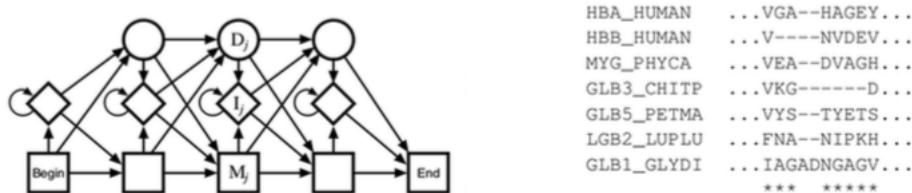


Figure 31: Parameters are: value of the probabilities and length. Begin, End and D states are silent.

5.3.4 Deriving a pHMM from a multiple alignment (Supervised)



- a **match state** for each conserved position (e.g. at least 50%)
- **insert states**: e.g. columns with at least 50% gaps
- **delete states**: gaps on match positions
- **emission** (for match states) and **transition probabilities** are estimated from the counts
 - ▶ $B_{ki} = \frac{f(k,i)}{f(k)}$
 - $f(k,i)$ = number of times symbol i is observed on state k
 - $f(k)$ = number of times state k is used
 - ▶ $A_{kl} = \frac{f(k,l)}{f(k)}$
 - $f(k,l)$ = number of times a transition from state k to state l is used

But the initial multiple alignment could have some probabilities equal to zero: we need **smoothing**.

Additive smoothing with pseudo-counts

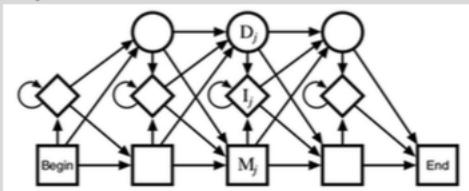
- $B_{ki} = \frac{f(k,i) + \varepsilon}{f(k) + \sum_i \varepsilon}$ with e.g. $10^{-6} \leq \varepsilon \leq 1$
 - ▶ $f(k, i)$ = number of times symbol i is observed on state k
 - ▶ $f(k)$ = number of times state k is used
- $A_{kl} = \frac{f(k,l) + \varepsilon'}{f(k) + \sum_l \varepsilon'}$ with e.g. $10^{-6} \leq \varepsilon' \leq 1$
 - ▶ $f(k, l)$ = number of times a transition from state k to state l is used

5.3.5 Unsupervised learning

No need for an initial multiple alignment (only unaligned sequences).

Procedure

- ① choose a general pHMM structure



- ② choose the number of match states:

e.g. half the average sequence length

- ③ estimate the pHMM parameters through Viterbi or Baum-Welch

5.3.6 Matching a sequence to a pHMM a.k.a. viterbi recurrence

- Computations are done usually with log's:
 - $\log \gamma(k, t) = \min_l [-\log \gamma(l, t-1) - \log A_{lk}] - \log B_{kx_t}$
- Including a background model to produce a log-odds score:

$$\log \gamma(k, t) = \max_l [\log \gamma(l, t-1) + \log A_{lk}] + \log \frac{B_{kx_t}}{q_{x_t}}$$
- A similar adaptation can be included into the forward recurrence

5.3.7 For non-global alignments

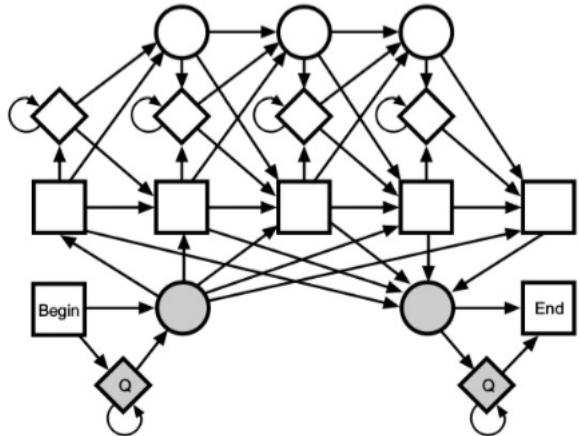


Figure 32: Local alignment (S-W style). Non-conserved fragments are modeled through flanking insert states using the background emission probabilities q_a . Their loop have high probabilities. Flanking delete states allow for starting or ending the profile at any point and reduce number of transitions.

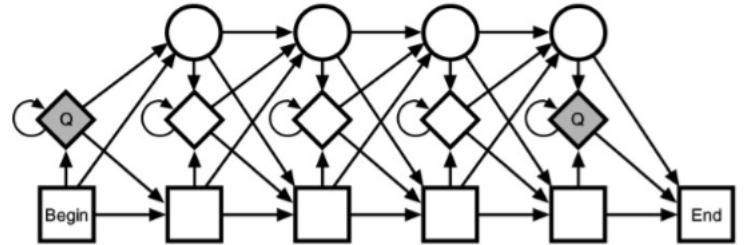


Figure 33: Set all the transition probabilities from the left flanking state to different start points. Force the match on the complete profile.

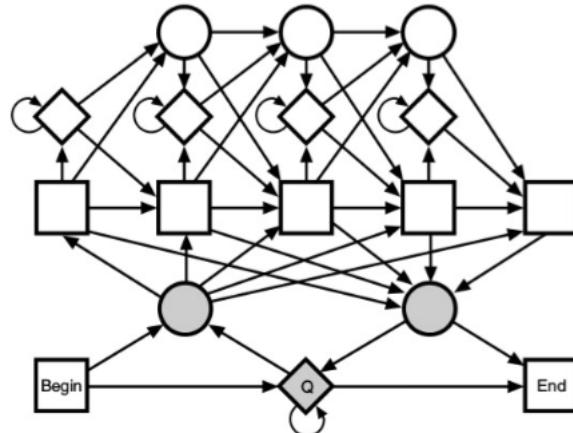


Figure 34: Allow repeated matches to subsections of the profile.

6 Large Scale Gene Expression Analysis

6.1 Introduction

DNA Microarrays measure the level of expression of all genes in a single experiment. Data measurements → Preprocess and sample normalization → Gene selection and sample classification → Diagnosis, prognosis or prediction of the response to a treatment.

The **gene selection** (supervised learning) to find a subset of genes to predict the response of new samples. There are some **objectives**:

- Insight into the data and the predictive model;
- Link between data analysis and medical expert;
- Biological validation;
- Reduce financial cost.

There are some **difficulties**:

- Measurements are noisy;
- Gene expression varies due to many factors;
- Financial cost;
- Small n large p problem.

6.2 Preprocessing

6.2.1 Summarization

Define a single probeset expression level from the various probe intensities. There are popular techniques: MAS 5.0, RMA, GC-RMA.

- background adjustment: optical noise correction, probe affinity adjustment (influenced by the GC content), RMA ignores the MM probes;
- sample normalization: quantiles should be stable across samples, after conversion to log intensities for (GC-)RMA;
- summarization: median polish.

6.2.2 Feature normalization

Make sure that each gene (probeset) has roughly the same expression range across all samples and then do a **Z-score normalization**: Replace $x_{i,j}$ by $\frac{x_{i,j} - \mu_j}{s_j}$. μ_j is the mean level of expression of probeset j over the training samples and s_j is the standard deviation.

6.2.3 Distance between expression values

Euclidean distance

$$d(\mathbf{x}_1, \mathbf{x}_2) = \|\mathbf{x}_1 - \mathbf{x}_2\| = \sqrt{\sum_{i=1}^n (x_{i,1} - x_{i,2})^2}$$

Correlation based distance

$$d(\mathbf{x}_1, \mathbf{x}_2) = 1 - \frac{1}{2}(1 + \text{corr}(\mathbf{x}_1, \mathbf{x}_2))$$

6.2.3.1 Pearson correlation

For two random vectors (e.g. gene expression values) $\mathbf{x}_1, \mathbf{x}_2$ measured over n samples

$$\text{corr}(\mathbf{x}_1, \mathbf{x}_2) = \frac{\sum_{i=1}^n (x_{i,1} - \bar{x}_1)(x_{i,2} - \bar{x}_2)}{\sqrt{\sum_{i=1}^n (x_{i,1} - \bar{x}_1)^2 \sum_{i=1}^n (x_{i,2} - \bar{x}_2)^2}}$$

- $\text{corr}(\mathbf{x}_1, \mathbf{x}_2) = \pm 1$ if \mathbf{x}_1 and \mathbf{x}_2 are perfectly linearly correlated
- $\text{corr}(\mathbf{x}_1, \mathbf{x}_2) = 0$ if they are not linearly correlated
- whenever \mathbf{x}_1 and \mathbf{x}_2 are normalized to zero mean and unit variance:

$$\text{corr}(\mathbf{x}_1, \mathbf{x}_2) = \sum_{i=1}^n x_{i,1} x_{i,2} = \mathbf{x}_1^\top \mathbf{x}_2$$

6.2.3.2 Pitfalls with correlation measure

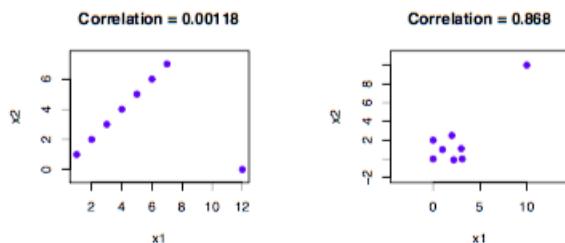


Figure 35: Correlation is sensitive to outliers.

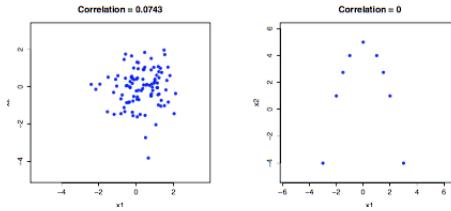


Figure 36: Correlation measures **linear** dependence.

6.2.3.3 Spearman's rank correlation This correlation is less sensitive to outliers (but still measure linear dependence).

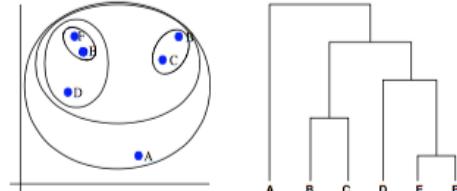
- Replace feature value by feature value rank across observations;
- Compute pearson correlation between rank vectors.

6.3 Unsupervised learning

	gene 1	gene 2	...	gene p	cluster
sample 1	$x_{1,1}$	$x_{1,2}$...	$x_{1,p}$?
...
sample n	$x_{n,1}$	$x_{n,2}$...	$x_{n,p}$?
cluster	?	?	...	?	

Figure 37: Objective is to find cluster of genes that share similar profile.

Each observation represents
either a **sample** across genes
or a **gene** across samples



```

Algorithm AGGLOMERATIVEHIERARCHICALCLUSTERING
Input:  $D$  a set of observations  $\vec{x}_1, \dots, \vec{x}_m$ ;  $d(\vec{x}, \vec{y})$  a distance measure between
observations
Output: A tree  $T$  of subsets of  $D$ 
// Initialize a set  $\mathcal{D}$  of clusters  $D_1, \dots, D_m$ 
 $\mathcal{D} \leftarrow \{\{\vec{x}_1\}, \dots, \{\vec{x}_m\}\}$  // Initial clusters are tree leaves
 $T \leftarrow$  a partial tree whose leaves are the  $\vec{x}_i$ 's
while  $|\mathcal{D}| > 1$  do
    Choose pair of clusters  $(D_i, D_j)$  in  $\mathcal{D}$  such that  $Distance(D_i, D_j, d)$  is minimal
    Define a new cluster  $D_k = D_i \cup D_j$ 
     $\mathcal{D} \leftarrow \mathcal{D} \cup D_k - \{D_i, D_j\}$ 
    Add  $D_k$  as parent node of  $D_i$  and  $D_j$  in the tree  $T$ 
return  $T$ 

```

Figure 38: Agglomerative hierarchical clustering.

There are different distance measure between the clusters.

Single-link or nearest neighbor rule

$$\text{Distance}(D_i, D_j, d) = \min_{\vec{x} \in D_i, \vec{y} \in D_j} d(\vec{x}, \vec{y})$$

Complete-link or farthest neighbor rule

$$\text{Distance}(D_i, D_j, d) = \max_{\vec{x} \in D_i, \vec{y} \in D_j} d(\vec{x}, \vec{y})$$

Average-link rule

$$\text{Distance}(D_i, D_j, d) = \frac{1}{|D_i| \cdot |D_j|} \sum_{\vec{x} \in D_i, \vec{y} \in D_j} d(\vec{x}, \vec{y})$$

UPGMA algorithm is with average-link measure. Used in phylogeny.

6.4 Supervised learning

	gene 1	gene 2	...	gene p	response
sample 1	$x_{1,1}$	$x_{1,2}$...	$x_{1,p}$	y_1
...
sample n	$x_{n,1}$	$x_{n,2}$...	$x_{n,p}$	y_n
test sample	x_1	x_2	...	x_p	?

Figure 39: p input dimensions, n samples, y is binary, natural or real. The goal is to find the most discriminating genes for the prediction of the class of any new sample.

6.4.1 Filters

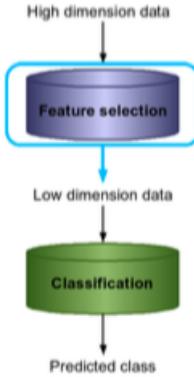


Figure 40: Use only training data and class labels during the features selection step and then train a single classifier taking the selected features as inputs. It is the less computative approach.

- 6.4.1.1 Non-specific filtering** Keep only genes with the larger variances (top 25%) because they are discriminating (do it before normalization to unit variance).

6.4.1.2 Fold changes

$$\frac{\bar{x}_1}{\bar{x}_2} \text{ or } \log \frac{\bar{x}_1}{\bar{x}_2} = \log \bar{x}_1 - \log \bar{x}_2 \text{ or } \bar{x}_1 - \bar{x}_2$$

Figure 41: Keep genes with larger fold changes between both conditions. Not a good filters because discriminating genes with large difference but small ratio are missed.

6.4.1.3 t-Test relevance index

t-Test statistic (actually Welch t-statistics)

$$J(x) = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{S_1^2/n_1 + S_2^2/n_2}}$$

with n_1 (resp. n_2) the number of examples labeled as + (resp. -) and the estimated variances in each class $S_i^2 = \frac{1}{n_i-1} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$

Figure 42: $J(x)$ is the distance between feature and average feature value in **each class**. p-values assess the significance of the difference between the two class means. A feature is selected if its associated p-value is below a prescribed threshold.

6.4.1.3.1 Alternatives of the simple t-Test

- Mann-Whitney rank test is an alternative **non-parametric** test;
- ANOVA is a generalization of the t-Test with **multi-class** (bigger than 2);
- Pairwise t-Tests **between one class and the others**;
- Kruskal-Wallis is a generalization of M-W (**non parametric multi class**).

6.4.1.4 Multiple test correction

It is possible to conclude that the mean expression values among the 2 classes are significantly different for a given gene while they are not. If the test is performed 50000 times and $\alpha = 0,05$, we are expected to wrongly select $50000 * 0,05 = 2500$ genes. We need correction.

6.4.1.4.1 Bonferroni correction

Divide the critical value by the number of tests n_t performed ($\frac{0,05}{50000}$). Tends to select no features.

6.4.1.4.2 False Discovery Rate correction

- ① Select a confidence level α (e.g. 0.05)
- ② Rank the p -values (one for each feature) in **increasing order**
 $p_1 \leq p_2 \leq \dots \leq p_{n_t}$
- ③ Iterate over the n_t features
 $(n_t = p$ with data in \mathbb{R}^p , not to be confused with p -values)
 - ▶ Find the **maximal index** i such that the p -value $p_i < \frac{i \times \alpha}{n_t}$
- ④ Keep all features up to index i_{max}

Figure 43: If $p_{n_t} < 0,05$ FDR select all the features. Only the selection threshold is changed.

6.4.1.5 Feature ranking with mutual information

$$\begin{aligned}
 I(X; Y) &= - \sum_{ij} P(x_i, y_j) \log_2 \frac{P(x_i, y_j)}{P(x_i)P(y_j)} \\
 &= - \sum_{ij} P(x_i, y_j) \log_2 \frac{P(y_j|x_i)}{P(y_j)}
 \end{aligned}$$

Figure 44: A feature X is more relevant if its mutual information with the class value is higher. If X tends to bring no info to predict Y: $I(X; Y) \rightarrow 0$. $I(X; Y) = 0$ if X and Y independent.

Those are univariate filters. We can use mutual information to select several variables at a time.

6.4.1.6 Maximum relevance minimum redondancy

① Select the feature with maximum mutual information with the response

- ▶ $\hat{X} = \operatorname{argmax}_X I(X; Y)$
- ▶ $\Phi = \{\hat{X}\}$ // Initialize the set of selected features
- ▶ $F = \{X_1, \dots, X_p\} \setminus \{\hat{X}\}$ // The remaining set of features

② Repeat

$$\hat{X} = \operatorname{argmax}_{X \in F} \left[\underbrace{I(X; Y)}_{\text{maximize relevance}} - \underbrace{\frac{1}{|\Phi|} \sum_{X_j \in \Phi} I(X; X_j)}_{\text{minimize redundancy}} \right]$$

$\Phi \leftarrow \Phi \cup \{\hat{X}\}$; $F \leftarrow F \setminus \{\hat{X}\}$

until an appropriate number of features are selected

6.4.1.7 Conclusion

Filters is the cheapest approach and is useful when we have large number of features. Univariate filters ignore interactions between genes.

6.4.2 Wrappers

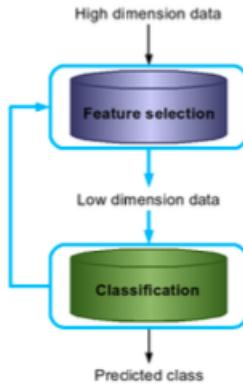


Figure 45: Train a classifier on several subsets of all possible features ($O(2^p)$ different subsets). Exhaustive search is impossible, use feature ranking or f/b selection. Select the feature set that optimize the performance of the trained classifier.

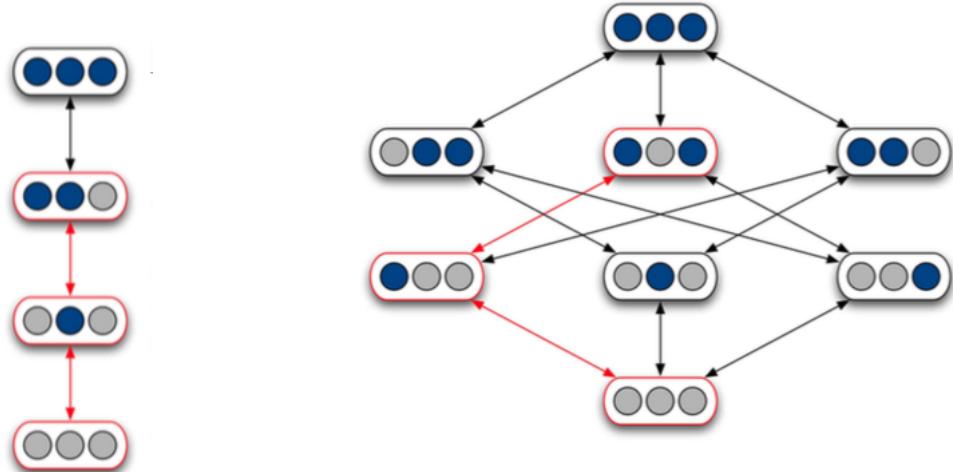


Figure 46: Univariate feature ranking.
Figure 47: Multivariate Forward (bottom-up, nothing to all) or Backward (top-down, all to nothing) selection.

The search order matters. If you can afford the computational cost of multivariate selection, you can do an embedded approach.

6.4.3 Embedded approaches

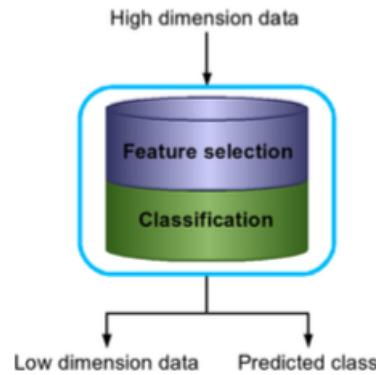
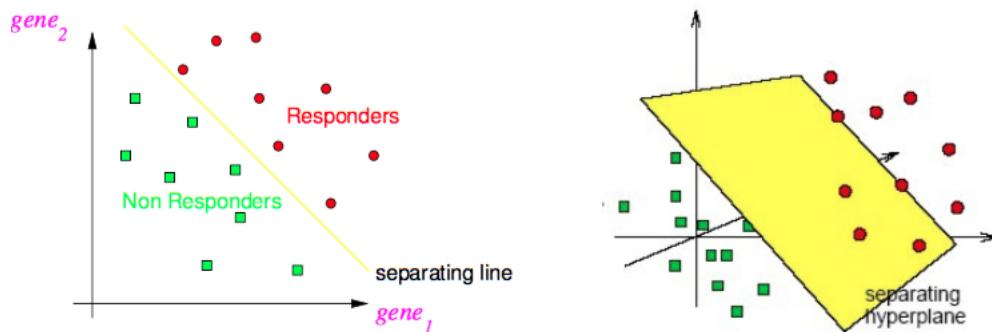


Figure 48: Feature selection and classifier estimation is the same **combined optimization process**. We include the classifier optimization in the feature selection process. The most computing expensive.

6.4.3.1 Linear discriminants



- The actual number of dimensions may be $\approx 10,000$ for microarray classification
- The linear discriminant is a **hyperplane** in $\mathbb{R}^{\geq 10,000}$
- Decision rule: $\text{sign}(\sum_{j=1}^p w_j x_j + w_0) = \text{sign}(\mathbf{w}^\top \mathbf{x})$ (with $x_0 \triangleq 1$)
 $\Rightarrow |w_j|$ is a measure of the importance of the j^{th} feature

Figure 49: The problem is there can be several hyperplane.

6.4.3.2 Linear Support Vector Machines

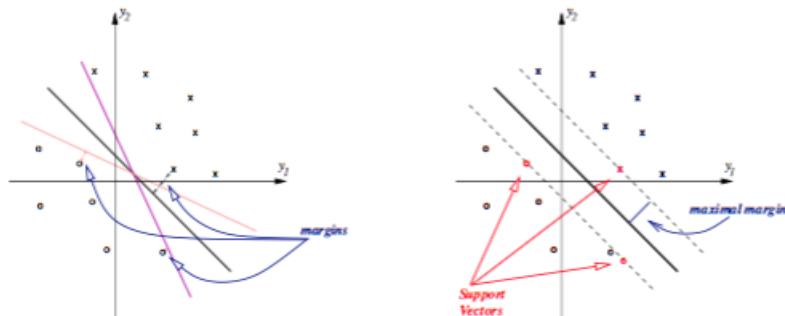


Figure 50: SVM. Need to find the maximal margin.

6.4.3.3 Recursive feature Elimination

Embedded Backward Selection

- ➊ Estimate a SVM on a given set of dimensions
(initially p dimensions)
 - ▶ Decision rule: $\text{sign}(\sum_{j=1}^p w_j x_j + w_0)$
- ➋ Consider $|w_j|$ as the relevance of the j^{th} dimension
- ➌ Remove the least relevant dimension(s)
- ➍ Iterate ➊ to ➌ on a reduced feature set

7 Inference of Gene Regulatory Network

It's the set of all transcription interactions in a cell. There is some regulator (transcription factor) and regulated gene (target).

Module-based inference	Expression-based methods	Integrative methods	Supervised inference	Unsupervised inference
<p>① search for clusters (= modules) of genes that exhibit a similar expression behavior</p> <p>② infer the transcriptional program of each module</p> <ul style="list-style-type: none"> ▶ better at finding interactions between regulators and targets with less similar expression profiles 	<p>Direct inference</p> <ul style="list-style-type: none"> ● look for the regulators of each gene ▶ better at finding the target(s) of regulators dedicated to one or a few genes 	<p>Integrative methods</p> <ul style="list-style-type: none"> ● complement gene expression with additional information: motif data or protein-DNA interaction data ● more likely to predict true positive interactions ● the added value depends on the quality of the additional information 	<p>Supervised inference</p> <ul style="list-style-type: none"> ● treat network inference as multiple classification problems <ul style="list-style-type: none"> ▶ which target genes are regulated by a given transcription factor ● require a training set of known interactions between TFs and TGs and expression data ● each learned classifier is then used to predict whether further genes should be considered as target of a given TF 	<p>Unsupervised inference</p> <ul style="list-style-type: none"> ● No such training set is available
		<p>Global inference</p> <ul style="list-style-type: none"> ● provide a general view of the active GRN ● identify large pathways that consist of many genes 		
		<p>Query-driven inference</p> <ul style="list-style-type: none"> ● search for genes that are co-expressed, in a condition-dependent way, with a predefined set of query genes ● find local regulatory interactions 		