

Atasal Dizi Tahmini - Bash Temelli Kisa Uygulama

Evrimsel Genombilim Kis Okulu, Hacettepe Biyoloji Bolumu, 2018.

Sorular ve komutlarda buldugunuz yanlislar icin: ozan.bozdag@gmail.com

A) Giris

Atasal dizi tahmini temelde 3 adim iceriyor: 1) sequence alignment 2) phylogenetic tree yapimi (bizim uygulamamizda 'maximum likelihood' yontemi ile) 3) ilk iki veriyi (alignment+filogeni) bir araya getirip atasal dizilerin tahmin (!) edilmesi.

1 no'lu adimi t-coffee kullanarak yapiyoruz.

2 no'lu adimi 'maximum likelihood' ile agac cizimi / hesabi yapan bir program kullanabilirsiniz (orn. PhyML). Fakat biz adim #2 ve #3'u tek bir seferde yapabildigi icin iqtreen'yi kullaniyoruz/kullandik. Boylece #2 & #3 ayni komut ile gerceklestirilecek.

Not: iqtreen'in bir avantajı sekans verisine uygun olabilecek nukleotid substitution modelleri arasından seçim yapabilmesi. Beyaz tahta üzerinde substitution modellerinin neden evrimsel agac ciziminde ve dal uzunluklarının kalibre edilmesinde önemli olduğunu konuştuk! Eger farklı programlar ile substitution model seçimi yapmak istiyorsanız buna bakabilirsiniz (alignment yaptıktan sonra çıkan output'u model testi yapan programa input olarak veriyorsunuz, ve size gereken modeli öneriyor, ilgili modeli sonraki adımda kullanacağınız filogeni programına bildirmeniz gerekiyor – her program size sunulan model seçeneğini içermeyebilir!): <https://academic.oup.com/mbe/article/25/7/1253/1045159>

Cok önemli uyarı: Kendi çalıştığınız sekans verileri üzerinde alignment, filogeni, ve atasal dizi tahmini yapmaya başlamadan önce kullanacağınız programların 'el kitaplarını' (manual) okumanız bilimsel etik ve tutarlılık nedeniyle gerekli.

<https://media.readthedocs.org/pdf/tcoffee/latest/tcoffee.pdf>
<http://www.iqtree.org/doc/iqtree-doc.pdf>

B. Teori

Asagıda çok kis bir iki komut uygulayarak modern primatların ADH4 amino asit dizilerini kullanarak atasal ADH4 enzimi amino asit dizilerini tahmin edeceğiz. ADH4 enzim aktivitesinin evrimsel olarak nasıl değiştiğini araştırarak çok ilginç bir çalışmayı temel aldık: <http://www.pnas.org/content/pnas/112/2/458.full.pdf>

Carrigan et al 2015 "Hominids adapted to metabolize ethanol long before human-directed fermentation" www.pnas.org/cgi/doi/10.1073/pnas.1404167111

Hedef: Calismada onemli oldugu bulunan A294V mutasyonunu atasal dizilerde bulmayi hedefliyoruz (makalede Figur 1). Not: kis okulu sirasinda atasal dizi tahminini grup calismasi olarak secen arkadaslariniz bu onemli mutasyonu dogru evrimsel yerlesimi ile buldular!

C. Uygulama - BASH komut dizini uzerinde uygulanacak -

BASH komutlari '\$' isareti ile basliyor.

1. ADIM: T COFFEE ile alignment

<http://www.tcoffee.org/Projects/tcoffee/papers/tcoffee.pdf>

T-coffee linux son surumunu indirin:

<http://www.tcoffee.org/Projects/tcoffee/#DOWNLOAD>

Dosya uzantisi .tar.gz ise, bu ziplenmis dosyayi bulundugu klasore bash komutlari ile giderek acin:

\$ tar -xvzf indirdiginiz_dosya_ismi.tar.gz

tab yapmayi unutmayin!

T-coffee'yi install etmek icin: http://tcoffee.readthedocs.io/en/latest/tcoffee_installation.html.

Onemli: T coffee'yi bash uzerinde herhangi bir klasorde iken (herhangi bir konumdan) calistirmek icin program klasorunuzu linux veya bio-linux'unuzun "search path"ine kopyalamaniz gerekiyor:

<https://oxaric.wordpress.com/2008/12/04/make-a-bash-script-globally-executable/>

Benim kullandigim Linux uzerinde search path /usr/bin/ konumu. "sudo copy" yaparak search path'e attim t_coffee'yi:

\$ sudo cp -R /home/manager/Downloads/T-COFFEE_installer_Version_11.00.8cbe486_linux_x64/bin/t_coffee /usr/bin/

T-coffee versiyonunuzu ve herhangi bir konumdan calisiyor oldugunu kontrol edin (download ettiginiz klasorden baska bir konuma gecerek deneyebilirsiniz):

\$ t_coffee -version

T-coffee seceneklerini okuyun (manual'i dikkatlice okudugunuzu varsayiyorum).

<http://www.tcoffee.org/Projects/tcoffee/#DOCUMENTATION>

\$ man t_coffee

Simdi istediginiz yontemle alignment yapin – biz en basit yontemi uyguladik, ornegin '-mode' secenegi ile 'accurate' komutu vererek daha iyi hesaba dayanan farkli islemler yapabilirsiniz:

\$ t_coffee adh4_primates_fasta.txt

3 farkli formatta cikti ("output") dosyasi verecek: newick, aln, html.

'adh4_primates_fasta.txt.aln' ciktisini kullanacagiz.

Dosya sonucuna bakin:

\$ cat adh4_primates_fasta.aln

2&3. ADIM: Filogeni olusturmak, substitution model'i secmek, ve tahmini atasal dizileri cikarmak:

Filogeni icin:

<https://academic.oup.com/mbe/article/32/1/268/2925592>

Atasal dizi tahmini ('-asr') secenegi hakkında bilgi icin:

<http://www.iqtree.org/doc/iqtree-doc.pdf>

ONEMLI NOT: iqtree surumu iqtree-1.6.1-Linux veya ustu olmalı - daha onceki surumlerde -asr secenegi (ancestral sequence reconstruction=asr) yok.

Dogru iqtree surumune sahip oldugunuzu kontrol edin:

\$ iqtree -v

Komutlari/secenekleri calisin:

\$ iqtree --help

ornegin MODEL-FINDER ve SUBSTITUTION MODEL seceneklerini okuyun.

-o ile outgroup secenegini inceleyin. Makalede gorunen outgroup (dis-grup) 'E_Tree.ADH.4_Tupaia_belangeri_' isimli tur. ADH4 protein dizisileri arasinda en uzak diziye iceren ADH4 protein dizisi bu ture ait olabilir (unutmayin, tur-agaci ile tek bir gen dizisine dayanan filogenetik agaclar uyumsuz olabilir!). Asagida bu turu outgroup secenegi ile tanitiyoruz. IQTREE eger outgroup vermezseniz kendisi outgroup seciyor. Outgroup'u kendiniz manual olarak tanitmaniz lazim komutunuzda. Bu onemli bir ayrinti!

Not: iqtree outgroup tanimada sorun yasiyordu.

-m TEST seceneginin ne anlama geldigini inceleyin

-asr secenegine bakin

-nt seceneginin ne anlama geldigine bakin

Model secimini (-m secenegi ile), filogenetik agac olusturma, ve atasal dizilerin tahmini (tek bir komut ile - derste isleri hizlandirmak icin bu yolu sectik):

\$ iqtree -s adh4_primates_fasta.aln -o E_Tree.ADH.4_Tupaia_belangeri_ -m TEST -asr

Komut isleme koyulacak. Burada ekrana verilenleri anlamak incelemek cok onemli! Bilimsel yayın icin ornegin Model Finder'in yaptigi islemleri, farkli modellere verdigi skorlari (orn. likelihood score:-LnL vb.) anlamaniz lazim. Derste substitution/evrim modelleri ile filogenetik agacları “gercege” nasil yaklastirdigimizi konustuk. Hangi modeller en ustte (en olasi) cikiyor? Ne anlama geliyor secilen model (orn. 'Dayhoff' en olasi model olarak hesaplanmis olabilir, dayhoff subs. model'i veya sizin calistirdiginizda cikan en olasi modeli arastirin, ne anlama geldigine bakin.

Ekrana FINALIZING TREE SEARCH altinda verilen raporlardaki log-likelihood verileri nedir (beyaz-tahta dersinde biraz konustuk), 'total number of iterations' ne anlama geliyor, total tree length farkli modeller ile agac cizerseniz nasil degisiyor, vb. ozelliklere bakmaniz, anlamaniz gerekiyor.

Peki elde edilen 4 dosya neler?

1. adh4_primates_fasta.aln.iqtree
2. adh4_primates_fasta.aln.treefile
3. adh4_primates_fasta.aln.mldist
4. adh4_primates_fasta.aln.state
5. adh4_primates_fasta.aln.log

less programi ile iqtree'nin verdigi en kapsamli filogeni dosyasini acin:

\$ less adh4_primates_fasta.aln.iqtree

Burada istatistik verileri, filogenetik agac, outgroup konusunda bir bilgi, modeller, ve agacin "newick" format olarak cizimi de var.

Ornegin: "WARNING: 4 near-zero internal branches (<0.0027) should be treated with caution

Such branches are denoted by '*' in the figure below" - bu uyarinin neden verildigini agaca bakarak gorebilirsiniz.

'q' ile cikabilirsiniz.

Derste "newick" formatini kopyala-yapistir ile phylogeny.fr'de agaci cizdirmistik.

adh4_primates_fasta.aln.iqtree dosyasının uygulamanın hedefi bağlamında (atasal protein dizilerini bulmak) bize verdiği EN ONEMLI BILGI: .iqtree dosyasında baş ekranına verilen ağaca baktığınızda "node 1", "node 2" ile atasal node'lari göreceksiniz. Örneğin Homo sapiens, Pan paniscus, Gorilla gorilla, ve Pan troglodytes'lerin birleştiği node hangi rakam ile gösteriliyor? İşte bu rakam bu 4 turun ortak atasında bulunan atasal ADH4 dizisini başka bir dosyadan bulup incelememizi sağlayacak.

Son olarak oluşturdüğümüz atasal diziler için .state dosyasına bakacağız:

\$ less adh4_primates_fasta.aln.state

less ile .state'i açınca mouse ile hareket edebilirsiniz dosyada – read-only olarak açıyorsunuz o nedenle dosya içeriğini bozma şansınız -iyi ki- yok!

dosyayı inceleyin, ne anlama geldiğini kendiniz zaten anlayacaksınız.

Sütun 1 Node#, sütun 2 proteinin bastan sona içerdigi tahmin edilen en olası amino asit seçimleri (347 amino asitlik bir proteindi ADH4, sanırım), sütun 4 ve sonrası her olası amino asit için verilen olasılık skorları olmalı.

.state dosyası her bir 'node'da bulunan atasal amino asit dizisini veriyor. İstedığınız bir atasal 'node'da bulunan atasal protein/DNA dizisini (örn. Node1) baş uygulama dersinde öğrendiğimiz yöntemlerle seçip yeni bir dosyaya aktarabilirsiniz.

#.state dosyasında her bir amino asit için ilgili 'node'da bulunan tahmini amino asit için bir likelihood skoru veriliyor (1 ve daha düşük skorlar). O skorlar size ileride makalenizde istatistik verilerini paylaşırken kullanmanız için gerekli olacak.

Makalede bulunan en kritik fenotipik (enzim aktivitesi değişimine sebep olan) A294V mutasyonunu bulabildiniz mi? İnceleyin...

Not: Bu işlem dersin süresine uygun olarak en basit yöntemi ve sorunsuz (akrabalık olarak yakın dizileri) kullanarak yapıldı. Tahmin edilen atasal amino asitlerin olasılık skorlarına baktığınızda zaten olasılıkların genelde yüksek olduğunu göreceksiniz. Fakat tahmin edilen atasal sekansların/DNA'nın/protein'in atasal diziyi tamamen yansıtmadığını, yüksek olasılıklı bir tahmin dizisi olduğunu unutmamalısınız (derste istatistik ile, doğruluk ile ilgili çalışmalardan bahsedildi, slaytlarda ilgili referansları bulabilirsiniz, ki bu konuda çalışacaksanız o çalışmaları kavrayarak yola çıkmanız gerekiyor!). Sonuç olarak asıl hedef örneğin 100 milyon yıl önceki “gerçek” atasal diziyi aynen bulmak değil (ki bu olası değil) fakat en olası diziyi bulup bu yüksek olasılıklı (!) tahmini dizinin gerçek diziden birkaç fark içersede (farklı genotip) atasal FENOTİPI (atasal enzimin aktivitesini, atasal enzimin hücre için konumunu vb) yansıtıyor olması (bundan bahseden makalede yine slaytlarda atıflaniyor).

SON

