

Maximum likelihood modelleme kullanarak genotip çağırma

Giriş

Bu belgede, maximum likelihood modelleme tekniği kullanarak, varyant çağırma konusu işlenecektir.

Phred skoru

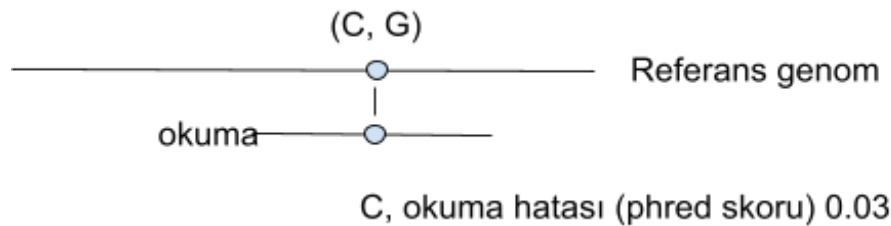
Kalite skoru, Phred skoru olarak da bilinen bu değer, DNA okunurken bir bazın ne kadar hata oranı ile tespit edildiğini belirten bir değerdir.

Phred kalite skoru	Yanlış baz olasılığı	Baz tespit hassasiyeti
10	1/10	90%
20	1/100	99%
30	1/1000	99.9%
40	1/10000	99.99%
50	1/100000	99.999%

Metodun örneklendirilmesi

Tek bir okuma üzerinden inceleme

Genotip çağırma işleminde aslında bu phred kalite değerlerini kullanacağız. Bunu bir örnek ile, tek okuma üzerinden gösterelim:



Şekil 1: Referans genoma hizalanmış bir okuma. Bu okumada gösterilen pozisyona, 0.03 phred kalite skoru ile C nükleotidi okunmuş.

Örnek verelim, elimizde bir okuma olsun. Bu okuma, insan referans genomu üzerine belirli bir pozisyonda hizalandığını düşünelim. Acaba bu pozisyonda, ilgilendiğimiz canlıda hangi genotipi görebiliriz? Bunu hesaplamak için öncelikle o pozisyonda hangi genotiplerin olabileceğini gösteren bir parametre uzayı tespit etmemiz gerek.

Düşünelim ki, o pozisyonda C veya G nükleotidi olabileceğini biliyoruz. O zaman 2n (diploid) bir organizmada, o pozisyonda kaç genotip gözlenebilir?

CC, GC; GG

Şimdilik biz sadece bu genotiplerin oluşma olasılığını tespit edeceğiz.

Adımlar:

1. Parametre uzayı belirle: GC; CC; GG
2. Gözlemi belirle: Tek bir okuma, C, okuma hatası 0.03
3. Parametre ile gözlem arasında olasılıksal bir bağlantı kur
4. Bunu bütün parametreler için tekrarla
5. En olası parametreyi seç

Örneğimiz için aslında 3 farklı genotipin oluşma olasılığını hesaplayacağız:

$$p(\text{Genotip} = CG | \text{Gözlem} = C)$$

$$p(\text{Genotip} = CC | \text{Gözlem} = C)$$

$$p(\text{Genotip} = GG | \text{Gözlem} = C)$$

$$p(\text{Genotip} = CG | \text{Gözlem} = C) = \left(\frac{1}{2} \times p(C)\right) + \left(\frac{1}{2} \times (p(G) \times \frac{1}{3})\right)$$

$$p(\text{Genotip} = CG | \text{Gözlem} = C) = \left(\frac{1}{2} \times 0.97\right) + \left(\frac{1}{2} \times 0.03 \times \frac{1}{3}\right)$$

$$p(\text{Genotip} = CG | \text{Gözlem} = C) = 0.49$$

- İşlemleri yaptığımızda sonuç 0.49 çıkacaktır.
- $\frac{1}{2}$ = Anne veya babadan gelme olasılığıdır. İnsan genomu üzerinde çalıştığımız için genler diploittir yani 2n bu sebeple bir anneden bir babadan gibi düşünebiliriz.
- C C uyuşması olduğunda bu doğru okumadır. bu sebeple 0.97= doğru okuma oranıdır.
- Kırmızı olarak gösterilen yerin $\frac{1}{3}$ olma nedeni ise C nükleotidi hariç diğer nükleotidlerle eşleşme olasılığıdır (A,T,G).

$$p(\text{Genotip} = CC | \text{Gözlem} = C) = \left(\frac{1}{2} \times p(C)\right) + \left(\frac{1}{2} \times p(C)\right)$$

$$p(\text{Genotip} = CC | \text{Gözlem} = C) = \left(\frac{1}{2} \times 0.97\right) + \left(\frac{1}{2} \times 0.97\right)$$

$$p(\text{Genotip} = CC | \text{Gözlem} = C) = 0.97$$

- $\frac{1}{2}$ = anne veya babadan gelme olasılığı. C çıkma olasılığı doğru okuma olduğu için 0.97'dir.
- Tekrar C çıkma olasılığı yeniden 0.97'dir.
- İşlemleri yaptığımızda sonuç 0.97 çıkacaktır.

$$p(\text{Genotip} = GG | \text{Gözlem} = C) = \left(\frac{1}{2} \times (p(G) \times \frac{1}{3})\right) + \left(\frac{1}{2} \times (p(G) \times \frac{1}{3})\right)$$

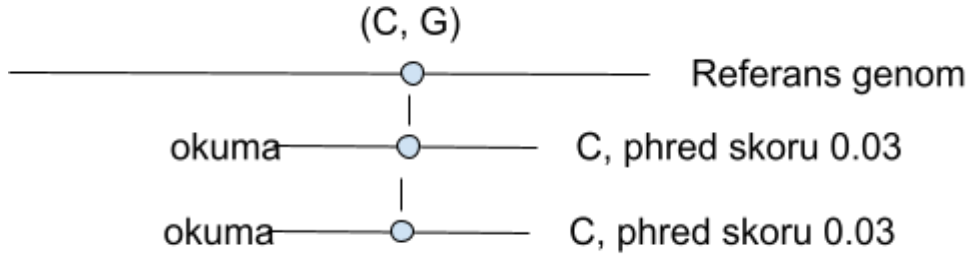
$$p(\text{Genotip} = GG | \text{Gözlem} = C) = \left(\frac{1}{2} \times 0.03 \times \frac{1}{3}\right) + \left(\frac{1}{2} \times 0.03 \times \frac{1}{3}\right)$$

$$p(\text{Genotip} = GG | \text{Gözlem} = C) = 0.01$$

- $\frac{1}{2}$ = anne veya babadan gelme olasılığı. G çıkarsa bu referans genomdaki C ile eşleşmediği için yanlış okumadır. 0.03=yanlış okuma oranıdır.
- $\frac{1}{3}$ diğer nükleotidlerin gelme olasılığıdır bu sebeple $\frac{1}{3}$ ile çarpılır.
- İşlemleri yaptığımızda sonuç 0.01 çıkacaktır.

CG=GC çıkma olasılığı= 0.49'tur.
 CC çıkma olasılığı = 0.97'dir.
 GG çıkma olasılığı = 0.01'dir.

İki okuma üzerinden inceleme



Şimdi ise her okumadan elde ettiğimiz gözlemleri, parametrelerle eşleştireceğiz

$$LH(\text{Genotip} = GC | x = \{C, C\}) = p(GC|C) \times p(GC|C)$$

$$LH(\text{Genotip} = GC | x = \{C, C\}) = 0.49 \times 0.49$$

$$LH(\text{Genotip} = GC | x = \{C, C\}) = 0.24$$

- Burada genotipte GC ile C eşleşmesi (0.49) ve GC ile C eşleşmesi (0.49) olasılıklarını çarparak 0.24 değerini buluruz.

$$LH(\text{Genotip} = CC | x = \{C, C\}) = p(CC|C) \times p(CC|C)$$

$$LH(\text{Genotip} = CC | x = \{C, C\}) = 0.97 \times 0.97$$

$$LH(\text{Genotip} = CC | x = \{C, C\}) = 0.94$$

- Burada genotipte CC ile C eşleşmesi (0.97) ve CC ile C eşleşmesi (0.97) olasılıklarını çarparak 0.94 değerini buluruz.

$$LH(\text{Genotip} = GG | x = \{C, C\}) = p(GG|C) \times p(GG|C)$$

$$LH(\text{Genotip} = GG | x = \{C, C\}) = 0.01 \times 0.01$$

$$LH(\text{Genotip} = GG | x = \{C, C\}) = 0.001$$

- Burada genotipte GC ile C eşleşmesi (0.01) ve GC ile C eşleşmesi (0.01) olasılıklarını çarparak 0.001 değerini buluruz.

Bütün genotipler için

Bu iki örnek sadece bir pozisyonda G veya C olma olasılığı üzerinden işlendi. Aslında elimizde daha fazla olasılık var.

$$x = \{A, T, G, C\}$$

Bu dört nükleotidin oluşturabileceği 10 farklı ikili kombinasyon vardır.

Bunlar:

1. AA
2. AT
3. AC
4. AG
5. TT
6. TG
7. TC
8. CC
9. CG
10. GG

A

T

T

T

$$LH(\text{Genotip} = AA | x = \{A, T, T, T\}) = p(AA|A) \times p(AA|T) \times p(AA|T) \times p(AA|T)$$

= İşleminin sonucun $-(\log)$ 'sını alıyoruz.

Yukarıda bulunan örnekteki gibi sonucun negatif çıkabileceği işlemlerde $-(\log)$ 'ya alarak sonucu pozitif yapıyoruz.

Hazırlayanlar: Halit Kemal Aydın, Serhat Kaçan, Gülseven Aleyna Yatmaz