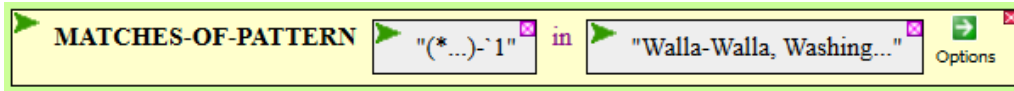


predict is the phenotype of the person carrying this gene, what evidence do you have for that prediction, and what code did you use to find it?

B. Identifying unknown STRs

So far you've identified repeated sequences where the repeating unit was known ("AGGT" in the first case and "CAG" in the second). What if you want to find repeated sequences and don't care what the repeating unit is? MATCHES-OF-PATTERN will work for this task as well.

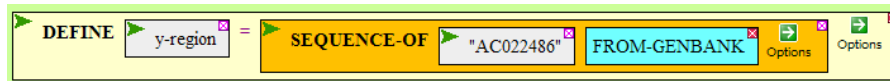
1. Try it out first with a simple case:



The pattern "**(*...)-`1**" consists of a string of indeterminate length [***...**], that is captured [**()**], followed by a hyphen [**-**], and finally a repetition of the first captured item [**`1**]. **Try out this function to see what you get.**

2. Now for something more bioinformatic. Suppose you have sequenced a portion of the Y-chromosome from an individual and want to determine if there are any STRs in the segment that might be useful in distinguishing DNA from different males.

Obtain the sequence from GenBank by executing the function shown below:



Use MATCHES-OF-PATTERN to find all instances of a unit 3 to 6 undetermined nucleotides in length repeated at least 6 times. In theory, you could find a tandem repeat without specifying the length of the unit, as you did with Walla-Walla, but in practice with long sequences, the execution time would be ridiculously long. Set limits when you can. **What is the longest STR and what code did you use to find it?**

3. Here's a remarkable bacterial example. In [CyanoBIKE](#) find all instances of a unit 5 to 8 undetermined nucleotides in length repeated at least 4 times in the genome of the cyanobacterium *Nostoc punctiforme* (nicknamed Npun). To cut down on extraneous detail, choose the +FULL-MATCH option as well as ONE-STRAND. **What generalities do you find in the sequences of the repeating units?**

C. Play!