

1 General

1.1 R note

Format of If/Else:

```
if {  
  
}  
else {  
  
}
```

2 TODO

1. PANSE model implementation:

- (a) PANSEParameter.cpp
- (b) PANSEModel.cpp
- (c) PANSEParameter.h
- (d) PANSEModel.h
- (e) Ask about sigma term – Done
- (f) Ask about lambda prime term (is it lambda prime?) — check RFP section for how to actually calculate — DONE

2. Expand Unit Testing:

- (a) Test Cov Matrixes — STALLED: Still need final two
- (b) Test MCMC — STALLED: Need run, varyInitialConditions, calculateGewekeScore, getLogLikelihoodPosteriorMean, and setRestartFileSettings as well as two test that only functions.
 - Implement other unit testing first
- (c) Parameter – In progress
- (d) Test RFP Parameter
- (e) Test Trace
- (f) ...Per class basis
- (g) Eventually, some R scripts to do a short run for each model: Talk to Cedric

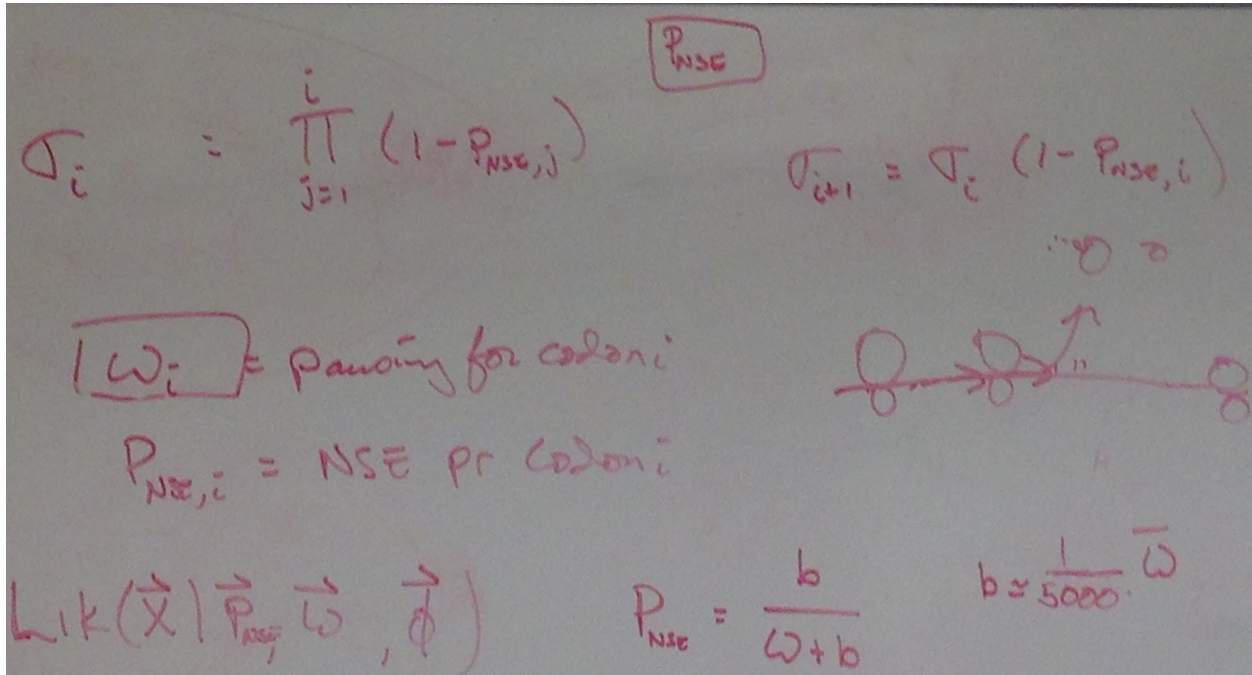
3. r

- 4. When working with gene-specific parameters, the openmp statements aren't working (memory is such a mess in the area) — break down parallelization, try to find where the issue is. Perhaps start with dynamic arrays, change to vectors. Gabriel thinks the slowdown from vectors in general is made up by better parallelization in avoiding dynamic arrays.

- —STALLED. Literally can't test speeds of various optimizations and cores right now.

5. Documentation

3 May 13, 2016 Notes



$$\sigma_i = \prod_{j=1}^i (1 - P_{NSE,j}) \quad (1)$$

ω_i = pausing for codon i

$p_{nse,i}$ = NSE Pr (probability) for codon i

This is codon-based.

Likelihood of the data given the parameters: $\mathcal{L}(\vec{x} | P_{NSE}, \vec{\omega}, \vec{\phi})$

Will be a much smaller data set, and with hundreds of calculations rather than thousands.

Randomly select 600 genes instead of 5400

Sigma vector of: $\sigma_{i+1} = \sigma_i (1 - P_{NSE,i})$

Function is of probability of getting there vs waiting time once there

- Getting pausing values with simpler models (ROC)
- First analysis could be just estimating these terms
- This would mean creating a simulated data set.

- For simulation: $P_{NSE} = \frac{b}{\omega+b}$, where b is on the order of 1/5000 times average omega. ($b \simeq \frac{1}{5000}\bar{\omega}$) Talk to Jeremy about this, he may have finished this by now.

2015 paper, 2011 paper with primal

4 May 19, 2016 Notes

rfp.model.pdf: Reasoning [for lambda] is that for the sampling the Boltzman coefficient. See the explanation around equation (4) and the Z's and Y's.

Lambda Prime = Lambda.c * Z / Y, or call it K.

$$\lambda' = \lambda_c * \frac{Z}{Y}$$

Z is the overall state space

Y is what is sampled

$\lambda_c = \lambda' * \frac{C}{K}$. Let K be a new independent parameter, and keep track of Lambda Prime.

5 May 25, 2016 Notes

Codon-Specific Elongation Rate:

$P_{NSE} = \frac{b}{b+c}$ where b is where it flies off and c is where it continues.

Omega is the odds ratio of $\frac{P_{NSE}}{1-P_{NSE}}$. Therefore $\omega = \frac{b}{c}$

Look at 2006, 2007 papers.

LOOK AT UPDATED PDF: IT'S IN FRAMEWORK

Psi (the symbol which I *thought* was Omega) is the ribosome initiation rate: Rate at which ribosomes are jumping onto the mRNA. Phi is the rate that they are jumping off at the very end.

If you have 50% chance to get to the end, then Psi is twice as long as Phi $\Phi = \Psi * \text{Sigma}$.

Don't redo calculations from scratch, but rather in series.

5.1 Parallelization

- Only 20 AA's — Only 20 cores to spread load unto
- AA's with 6 codons of course take more time than those with 2

Gilchrist thinks what is meant by Gene-Specific Parameters is to parallelize at the highest level, i.e. at the gene or amino acid level.

I should check the code; find where the OpenMP statements are etc

Mostly something to ask other people about if I want to tackle the problem.

6 May 26, 2016 Notes

phi calculation, with mcmc accept/reject

dynamic arrays

big loop around everything

code doesn't work

couldn't figure out why

didn't spend that much time

we ended up parallelizing in the model class:

calculateLogLikelihoodRatioPerGene, apparently doesn't do much

perhaps better to parallelize outside, with the big loop

run a ROC model, then RFP

I'm running a fasta file that is simulated, so I know that it is true

I kinda need the R side

Get to the point where we suspect memory is the problem

Dynamic Arrays -> Vectors

7 May 31, 2016 Notes

Start 1:21

break 3:19

back 3:24

break 4:55

return 5:02

end 7:02

$2 + 1.5 + 2$

TODO: Go ahead and replace dynamic arrays with vectors, first

And then do this barebones calculation of runs to see if it makes it faster, without regards to parallelization.

8 June 1, 2016 Notes

Start 1 or 1:48

Break 3:30

Return 3:35

End 7:00

$2 + 3.5$

From yesterday:

0.00621732 - 10
0.00687881 - 100
0.00947537 - 1000
0.00713974 - 10000
0.00785908 - 10000
0.00750889 - 10

For today:

0.0572747 - 10
0.0698414 - 100

... Odd, 10x as long on average

The above was in DEBUG mode. Release mode redos:

A or V	Runs	Modifiers	Avg Time
V	100		0.0141421
A	100		0.0047742
V	10000		0.00850093
A	10000		0.00479609
V	10000	No Deletion	0.00871843
A	10000	No Deletion	0.00491614
V	10000	std::sort	0.00841396
A	10000		0.00598796
A	10000		0.00520682
A	100000		0.00455916
A	100000	std::sort	0.00776886
V	100000		0.00795495
V	100000	std::sort	0.00785736
A	100000	std::sort	0.00383634
A	100000		0.00385638
A	100000	std::sort	0.00392021

Note: Vectors are 2x as long on average now

Next step: Make a list of everything PANSE touches and unit test these things (first and foremost before actually writing PANSE)

ALSO: Estimate and track how long, in reality, it takes to do each unit testing

PARFP, PTRFP? Just calling it RFP might be misleading.

9 June 2, 2016 Notes

Start 1:01

Break 3:35

Return 3:50

End 6:58

Spent till 4 (3 hours) compiling notes and creating a git directory. Expecting to spend 1 hour deciding on what PANSE will need (or, rather, what RFP will need).

Talk with Gilchrist:

So data position feeds into: a) data on gene ab) to feed into ROC-RFP or b) PANSE-RFP

Which file type should I be reading in? RFP or Fasta?

For sample data for PANSE: Lareau Paper ->GSE ->The untreated replicates 1,2,3. Take one, and even then only a subset of one of them as sample data.

The Lareau material may have undergone more processing than the new Weinberg GSE published Feb 10 2016.

"Start with Lareau paper data" – Gilchrist, 5:33

10 June 3, 2016 Notes

Start 1:35

Break 4:09

Return 4:14

Decided to start reading the Lareau material. Began by looking directly at definition of data set (I chose untreated replicate 1) and then parse the data to get a smaller subset (file size otherwise is too large at 35MB)

Took longer than expected... When files finally parsed, 5:45.

Now have a data set of size 400 KB: those genes with 11 to 100 (inclusive) codons.