# Question 3.d)

gru on sequences of length 20 (without gradient clipping)
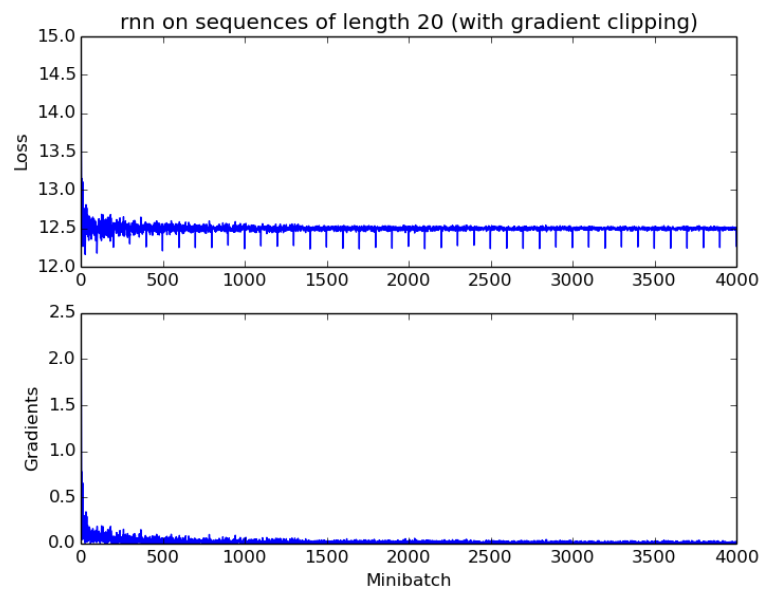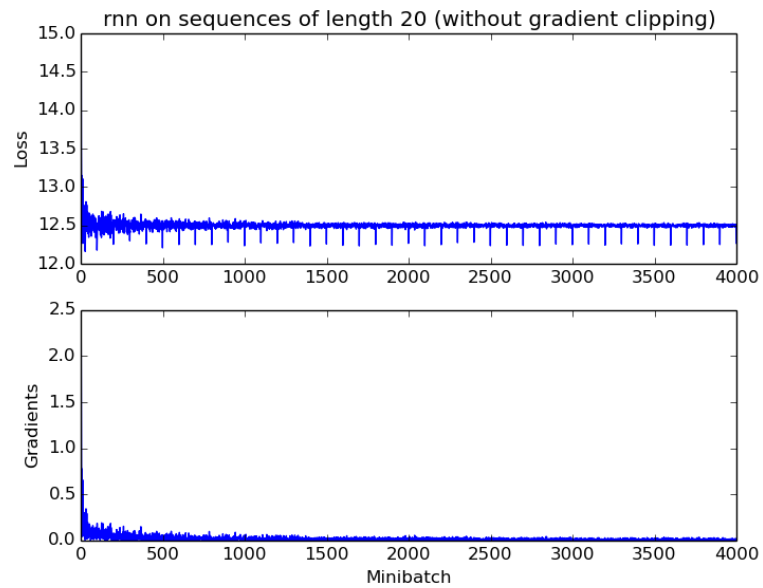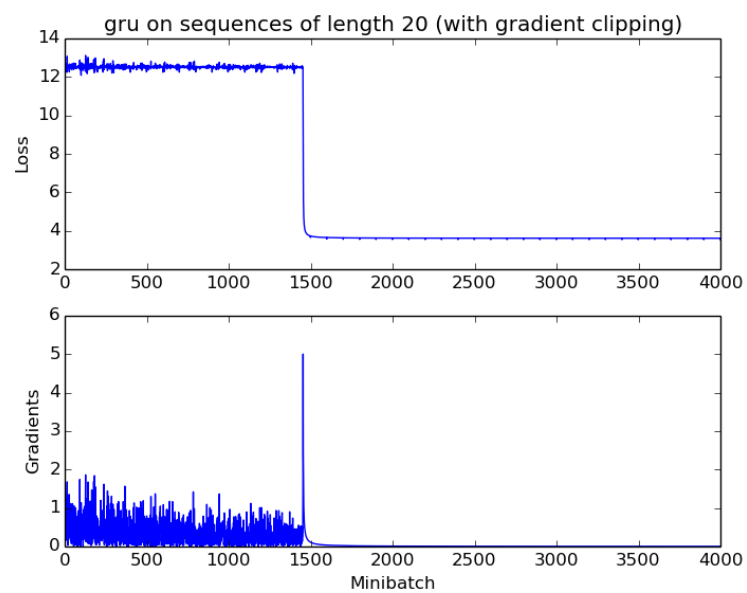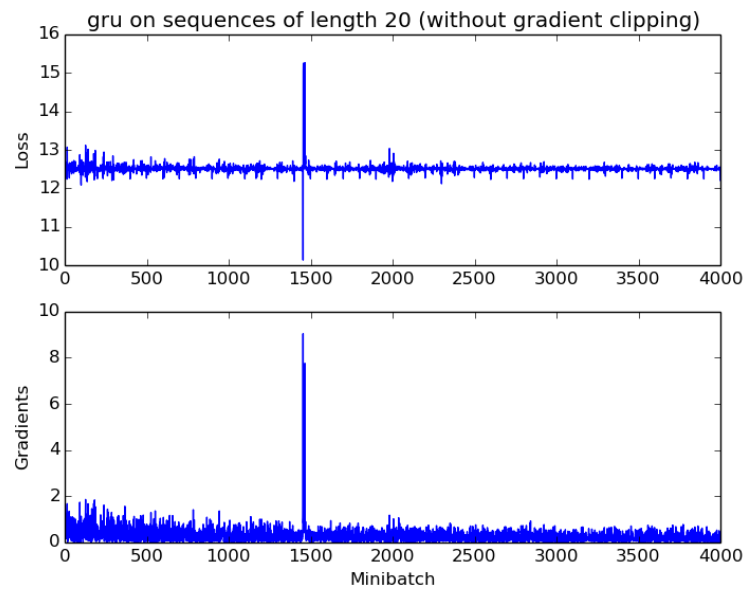
gru on sequences of length 20 (with gradient clipping)

## Question 3.e)i.

It seems like the RNN model experience vanishing gradients, since the gradients become very small and the model get stuck in some local minima. On the other hand, the GRU model seems to experience exploding gradient, since at some point it reaches a loss significantly better than before, but a large gradient pushes it to a worse local minimia region. Clipping does help, as we can see from the last pair of graphs: after it reaches this better loss, the clipped gradient is not large enough to push it to the worse local minima, and then we stay in the same region, eventually converging to a significantly better loss.

## Question 3.e)ii.

GRU does better since it is able to achieve a much lower loss. This doesn't necessarily guarantees a better generalization error, but the fit is definitely better, we achieve a much lower error and the model seems to allow for better accuracy.