Figure A.10: Cosine similarity of neuron out-directions and the sentiment direction in GPT2-small