

LLM EMOTION TRAJECTORY CONTROL

Experiment 1: Interpolative trajectory control.

- Let $x_0^+ =$ "Alice was happy, therefore Alice"
- $x_0^- =$ "Alice was sad, therefore Alice"

① Compute $\text{past_kv}^+ = \text{LLM}(x_0^+)$, $\text{past_kv}^- = \text{LLM}(x_0^-)$

② Compute $\{t_i^+ / t_i^+ \sim P_{\text{LLM}}(t_i^+ | x_0^+) \text{ iid } i \in [N]\}$ Also $p_i^{++} = P_{\text{LLM}}(t_i^+ | x_0^+)$
 \hookrightarrow Similar for $\{t_i^-\}$

③ Let past_kv^λ be

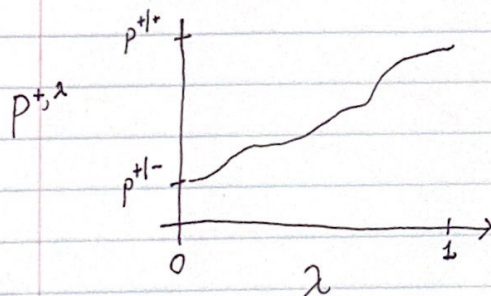
$$\text{past_kv}^\lambda := (1-\lambda)\text{past_kv}^- + \lambda\text{past_kv}^+$$

for $\lambda \in [0, 1]$.

④ Compute $P_{\text{LLM}}(t_i^\pm | \text{past_kv}^\lambda)$ for some $\lambda \in [0, 1] \forall i \in [N]$

\hookrightarrow Similar for t_i^\pm $=: p_i^{+, \lambda}$, similar for $p_i^{-, \lambda}$

⑤ Plot $\text{mean}_i(p_i^{\pm, \lambda})$ versus λ



Experiment 2: instead of λ , add $\pm \epsilon w_i$ to last value vecs in past_kv^\pm each.

Experiment 1.5: only interpolate λ for the final value in past_kv^\pm .