1 Math and Probability σ -algebra

 $\mathcal{F} \subseteq \mathcal{P}(\Omega)$ is a σ -algebra over Ω iff 1. $\Omega \in \mathcal{F}$ 2. if $\varepsilon \in \mathcal{F}$, then $\varepsilon^{\mathsf{C}} \in \mathcal{F}$ 3. if $\varepsilon_1, \varepsilon_2$... is a finite or infinite sequence in \mathcal{F} , then $\bigcup_n \varepsilon_n \in \mathcal{F}$

Measurable Space (Ω, \mathcal{F}) **Probability Measure**

 $\mathbb{P}: \mathcal{F} \to [0, 1]$ where 1. $\mathbb{P}(\Omega) = 1$ 2. If $\varepsilon_1, \varepsilon_2$ is a countable sequence of disjoint sets in \mathcal{F} , then $\mathbb{P}(\bigcup_n \varepsilon_n) =$ $\sum_{n} \mathbb{P}(\varepsilon_n)$

Probability Space $(\Omega, \mathcal{F}, \mathbb{P})$ Algebra

 $\mathcal{A} \subseteq \mathcal{P}(\Omega)$ over Ω if 1. $\Omega \in \mathcal{A}$ 2. $\varepsilon \in \mathcal{A} \implies \varepsilon^{\mathsf{C}} \in \mathcal{A}$ 3. $\varepsilon_1, \varepsilon_2 \in \mathcal{A} \implies$ $\varepsilon_1 \cup \varepsilon_2 \in \mathcal{A}$

Probability Pre-Measure

 $\mathbb{P}_0: \mathcal{A} \to [0,1]$ s.t. 1. $\mathbb{P}_0(\Omega) = 1$ 3 Classical LMs 2. If $\varepsilon_1, \varepsilon_2, \ldots$ is a countable sequence of disjoint sets in A whose *countable union is also in* A, then $\mathbb{P}_0(\bigcup_{n=1}^{\infty} \varepsilon_n) = \sum_{n=1}^{\infty} \mathbb{P}_0(\varepsilon_n)$

Geometric Series

 $S = \sum_{n=0}^{\infty} ar^n = \frac{a}{1-r}$, for |r| < 1Cylin $\overline{C}(\mathcal{H}) = \{ \boldsymbol{y} \boldsymbol{\omega} | \boldsymbol{y} \in \mathcal{H}, \boldsymbol{\omega} \in \overline{\Sigma}^{\infty} \}$

Binomial

 $\mu = \mathbb{E}[x] = np, \sigma^2 = np(1-p)$ Gaussian 68, 95, 99.7 $p(x) = \frac{1}{\sqrt{2\sigma^2}} \exp\{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2\}$

Triangle Inequality

 $||u + v|| \le ||u|| + ||v||$

Matrix Multiplication

 $\mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{B} \in \mathbb{R}^{n \times p} \mathbf{AB} \in \mathbb{R}^{m \times p} \mathbf{ti}_{-}$ Tightness of PFSA me $\mathcal{O}(m \times n \times p)$

Cosine similarity $\frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$

 $\mathbf{softmax}(\mathbf{x})_d = \frac{\exp\{x_d\}}{\sum_{j=1}^D \exp\{x_j\}}$

2 Foundations

Sequence Model

A probability space over the set $\mathcal{O}(|\Sigma|^{n-1})$ $\Sigma^* \cup \Sigma^{\infty}$

Language Model

A discrete distribution $p_{\rm LM}$ over Σ^* or $\mathbb{P} * (\Sigma^{\infty}) = 0$

 $\mathbf{GNM} \ p_{\mathbf{LM}}(\boldsymbol{y}) = \frac{\exp\{-\hat{p}_{\mathbf{GN}}(\boldsymbol{y})\}}{\sum_{\boldsymbol{y}' \in \Sigma^*} \exp\{-\hat{p}_{\mathbf{GN}}(\boldsymbol{y}')\}}$ LNM Prefix Prob $\pi(y) = \sum_{y' \in \Sigma^*} p_{\mathsf{LM}}(yy')$ **Probability Measure of LNM** 1. Define $\overline{\mathcal{C}} \subseteq \mathcal{P}(\Omega)$ as an algebra over $\Omega = \overline{\Sigma}^{\infty}$ where $\overline{\mathcal{C}} = \bigcup_{k=1}^{\infty} \overline{\mathcal{C}}_k$ 2. $\mathbb{P}_0(\overline{C}(\mathcal{H})) = \sum_{\overline{v} \in \mathcal{H}} p_{LN}(\overline{y})$

> 3. Extend \mathbb{P}_0 to $(\overline{\Sigma}^{\infty}, \sigma(\overline{\mathcal{C}}), \mathbb{P})$ 4. Construct SM: $\mathcal{C}(\mathcal{H}) = \{ \boldsymbol{v}\boldsymbol{\omega} | \boldsymbol{v} \in \mathcal{I} \}$ $\mathcal{H}, \boldsymbol{\omega} \in \Sigma^* \cup \Sigma^\infty$ then, $x(\boldsymbol{\omega}) =$ $\omega_{< k}$ if k is the first EOS in ω

ω otherwise

Tightness

An LNM is tight IFF $\tilde{p}_{EOS}(t) = 1$ for some t or $\sum_{t=1}^{\infty} \tilde{p}_{EOS}(t) = \infty$

FSA $\mathcal{A} = (Q, \Sigma, \delta, I, F)$ WFSA $\mathcal{A} = (Q, \Sigma, \delta, \lambda, \rho)$

Deterministic FSA

1. No ε -transitions 2. $\forall (q, a) \in Q \times$ Σ , at most one $q' \in Q$ s.t. $q \xrightarrow{a} q' \in$ δ 3. |I| = 1. $w(\pi) = \lambda(q_1) \prod_{i=1}^N w_i \rho(q_N).$

 $\mathcal{A}(y) = \sum_{\pi \in \Pi(A,v)} w(\pi)$ $Z(\mathcal{A}) = \sum_{v \in \Sigma^*} \mathcal{A}(v)$

Probabilistic FSA

 λ, ρ and the weights are nonnegative, $\sum_{q \in Q} \lambda(q) = 1$, $\forall q \in Q$ we have $\rho(q) + \sum_{q \longrightarrow q'} w = 1$

A PWFSA is **tight** if and only if all accessible states are co-accessible. **Finite State LM**

$$\exists \mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho), L(\mathcal{A}) = L(p_{LM})$$

n-gram assumption

 $p_{\text{SM}}(y_t|y_{< t}) = p_{\text{SM}}(y_t|y_{t-1}\dots y_{t-n+1})$ PAD with BOS n - 1 - t times.

Repre. based n-gram

 $\boldsymbol{\theta} = \{\boldsymbol{\theta}_{y|y} = p_{SM}(y|y) \mid y \in \overline{\Sigma}, y \in$ $\overline{\boldsymbol{\Sigma}}^{n-1}, \boldsymbol{\theta}_{\boldsymbol{\mathcal{Y}}|\boldsymbol{\mathcal{Y}}} \geq \boldsymbol{0}, \sum_{\boldsymbol{\mathcal{Y}}'\in\overline{\boldsymbol{\Sigma}}} \boldsymbol{\theta}_{\boldsymbol{\mathcal{Y}}'|\boldsymbol{\mathcal{Y}}} = \boldsymbol{1} \}$

MLE $p_{SM}(y_n | y_{< n}) = \frac{C(y_1, ..., y_n)}{C(y_1, ..., y_{n-1})}$ Bengio's Model $p_{SM}(\overline{y_t}|\overline{y}_{< t}) =$ $p_{\text{LN}}(\boldsymbol{y}) = p_{\text{SM}}(\text{eos}|\boldsymbol{y}) \prod_{t=1}^{T} p_{\text{SM}}(\boldsymbol{y}_t|\boldsymbol{y}_{< t}) \text{softmax}(\text{enc}(\boldsymbol{y}_{t-1:t-n+1})^\top \mathbf{E} + \mathbf{b})_{\overline{\boldsymbol{v}}_{\cdot}}$ **CFG** $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P})$ WCFG $\mathcal{W}: \mathcal{P} \to \mathbb{R}$ PCFG \mathcal{W} is non-negative, $\forall X \in \mathcal{N}$ we have $\sum_{X \to \alpha \in \mathcal{P}} \mathcal{W}(X \to \alpha) = 1.$ WCFG Allsum $Z(\mathcal{G}) = \sum_{d \in \mathcal{D}_c} w(d)$ $= \sum_{d \in \mathcal{D}_{\mathcal{C}}} \prod_{(X \to \alpha) \in d} \mathcal{W}(X \to \alpha)$

> **Tightness of PCFG** For a PCFG \mathcal{G} with $|\mathcal{N}| = N$ we define for each $X_n \in \mathcal{N}$ its production generating fct. $g_n((s_i)_{i=1}^N) =$ $\sum_{X_n \to \alpha} \mathcal{W}(X_n \to \alpha) s_1^{r_1(\alpha)} \cdots s_N^{r_N(\alpha)}$ where $r_i(\alpha)$ is the number of times X_i appears in α . Then we set $E \in \mathbb{R}^{N \times N}$ to have entries $E_{nm} =$

 $\frac{\partial}{\partial s_m} g_n(s_1, ..., s_N) \Big|_{s_1, ..., s_N = 1}$. Then \mathcal{G} is **tight** if $\lambda < 1$ and non-tight if $\lambda >$ 1, where $\lambda = \max\{|\lambda'| \mid \lambda' \in \sigma(E)\}$.

Pushdown Automaton

A language is context-free IFF it is recognized by some PDA.

Multi-Stack PDA

Any (probabilistic) 2-stack PDA is Turing complete. Hence, the tightness of a probabilistic 2-stack PDA is undecidable.

4 RNNs RNN

A RNN is given by an initial state

 $h_0 \in \mathbb{R}^d$ and a dynamics map $h_t = f(h_{t-1}, y_t)$. An RNN-LM uses

$\operatorname{enc}(\psi_{\leq t+1}) = h_t, E \in \mathbb{R}^{|\Sigma| \times d}$ **Recurrent Neural Sequence Model**

 $p_{\mathbf{SM}}(y_t|\boldsymbol{y}_{< t}) = \mathbf{f}_{\Lambda|\overline{\Sigma}|-1}(\operatorname{Eenc}_{\mathcal{R}}(\boldsymbol{y}_{< t}))_{y_t}$ Elman RNN

An Elman RNN is an RNN with $\mathbf{z}_t = O(\mathbf{a}_t) + \mathbf{a}_t$ $f(h_{t-1}, y_t) = \sigma(Uh_{t-1} + Ve'(y_t) + b),$ where $U \in \mathbb{R}^{d \times d}$, $V \in \mathbb{R}^{d \times R}$ and $b \in \mathbb{R}^d$ and $e' : \Sigma \to \mathbb{R}^R$ is an input embedding function. Jordan RNN $f(h_{t-1}, y_t) =$

 $\sigma(U\sigma'(Eh_{t-1}) + Ve'(y_{t-1}) + b)$

Tightness of RNN-LMs

If the LM uses the softmax and $s||h_t|| \leq \log t$ (in particular if f is bounded, e.g. if f uses a bounded activation function), then the induced LM is **tight**.

Expressiveness of RNNs

HRRNs (over \mathbb{R}) = dPFSA for any activation function with finite image. Minsky's construction encodes any dPFSA using $U \in \mathbb{R}^{|\Sigma||Q| \times |\Sigma||Q|}, U_{n(q',y'),n(q,y)} =$ $\mathbb{I}\{q_t \xrightarrow{y'/\circ} q' \in \delta\} \text{ to encode which states are reachable from } h_{t-1} V \in$ y'/\circ $\mathbb{R}^{|\Sigma||Q|\times|\Sigma|}, \ V_{n(q',y'),m(y')} = \mathbb{I}\{\circ \xrightarrow{y \to z}$ $a' \in \delta$ to encode which states can be transitioned to using v_t .

$$E_{\overline{m}(y')n(q,y)} = \begin{cases} \log \omega(q \xrightarrow{y'/w} \circ) & \text{if } y' \in \\ \log \rho(q) & \text{otherwise} \end{cases}$$

Can reduce the hidden state dimensionality to $\Omega(|\Sigma| \sqrt{|Q|})$. Saturated Sigmoid Elmann RNNs are Turing complete (because they can encode two-stack PDAs). It is thus undecidable whether an RNN-LM is **tight**.

5 Transformers

Attention

 $f: \mathbb{R}^D \times \mathbb{R}^D \to \mathbb{R}, \mathbf{q} \in \mathbb{R}^D, \mathbf{K}_t \in$ $\mathbb{R}^{t \times D}$, $\mathbf{V}_t \in \mathbb{R}^{t \times D}$. Att $(\mathbf{q}_t, \mathbf{K}_t, \mathbf{V}_t)$: $\mathbb{R}^D \times \mathbb{R}^{t \times D} \times \mathbb{R}^{t \times D} \to \mathbb{R}^D$ $\mathbf{s}_t = \mathbf{f}_{\Lambda^{D-1}}(f(\mathbf{q}, \mathbf{k}_1), \dots, f(\mathbf{q}, \mathbf{k}_t))$ $\mathbf{a}_t = \operatorname{Att}(\mathbf{q}_t, \mathbf{K}_t, \mathbf{V}_t) = s_1 \mathbf{v}_1 + \dots + s_t \mathbf{v}_t$

Transformer Layer

 $\mathsf{T}: \mathbb{R}^{T \times D} \to \mathbb{R}^{T \times D} \cdot \mathsf{X} \cdot \mathsf{Z} \in \mathbb{R}^{T \times D}$ $\mathbf{a}_t = \operatorname{Att}(Q(\mathbf{x}_t), K(\mathbf{X}_t), V(\mathbf{X}_t)) + \mathbf{x}_t$ $\mathsf{T}(\mathbf{X}) = \mathbf{Z} = (\mathbf{z}_1^\top, \mathbf{z}_2^\top, \dots, \mathbf{z}_T^\top) \in \mathbb{R}^{T \times D}$

Transformer

 $\mathbf{X}^{1} = (\mathbf{e}'(y_{0}), \mathbf{e}'(y_{1}), \dots, \mathbf{e}'(y_{t}))$ $\mathbf{Z}^{\ell} = \mathsf{T}_{\ell}(\mathbf{X}^{\ell}) \quad \text{for } 1 \leq \ell < L$ $\mathbf{X}^{\ell+1} = \mathbf{Z}^{\ell}, \mathbf{h}_t = F(\mathbf{z}_t^L)$

Attention Block A : $\mathbb{R}^{T \times D} \to \mathbb{R}^{T \times D}$ $\mathsf{A}(\mathbf{X}) = \mathbf{f}_{\Delta^{D-1}} (Q(\mathbf{X}) K(\mathbf{X})^{\top}) V(\mathbf{X})$ $\mathbf{U} = O(\mathbf{X}) K(\mathbf{X})^{\top} \in \mathbb{R}^{T \times T}.$ $W_i^Q, W_i^K \in \mathbb{R}^{d \times d_k}, W_i^V \in \mathbb{R}^{d \times d_v}$ $W^O \in \mathbb{R}^{hd_v \times d}$ often $d_k = d_v = \frac{d}{h}$

Masked Attention Block

 $A(\mathbf{X}, \mathbf{M}) = \operatorname{softmax}(Q(\mathbf{X})K(\mathbf{X})^{\top} \odot$ \mathbf{M}) $V(\mathbf{X})$ $M_{i,j} = \mathbb{I}[i \le j] + -\infty \mathbb{I}[i > j]$ Positional Enc $\mathbf{f}_{pos} : \mathbb{N} \to \mathbb{R}^D$

 $\mathbf{e}'_{\text{pos}}(y_t) = \mathbf{e}'(y_t) + \mathbf{f}_{\text{pos}}(t)$ $\mathbf{e}'_{\text{pos}}: \overline{\Sigma} \to \mathbb{R}^D$

MH-A

 $\mathbf{f}_H : \mathbb{R}^{T \cdot H \times D} \to \mathbb{R}^{T \times D} \mathsf{MH} - \mathsf{A}(\mathbf{X}) =$ $\mathbf{f}_H(\operatorname{concat}_{0 \le h \le H}(\operatorname{softmax}))$ $(Q_h(\mathbf{X})K_h(\mathbf{X})^{\top})V_h(\mathbf{X})))$

Layer Norm LN : $\mathbb{R}^D \to \mathbb{R}^D$ $L^{\mathbf{X}}(\mathbf{x};\boldsymbol{\gamma},\boldsymbol{\beta}) = \frac{\mathbf{x}-\overline{\mathbf{x}}}{\sqrt{\sigma^{2}(\mathbf{x})+\epsilon}} \odot \boldsymbol{\gamma} + \boldsymbol{\beta}$

FFN

 $FFN(x) = max(0, xW_1 + b_1)W_2 + b_2$

Tightness of Transformers

Any transformer using soft attention is **tight** (because its layers are continuous and the set of possible inputs to the first layer is compact, making enc bounded).

Expressiveness of Transformers

Let p_{LN} be an n-gram language model. Then, there exists a transformer \mathcal{T} with $L(p_{LN}) = L(\mathcal{T})$.

6 Tokenization

Tokenizer

A tokenizer model from Σ^* to Δ^* is a pair of stochastic maps T = $(\tau, \kappa), \tau: \Sigma^* \mapsto \Delta^*, \kappa: \Delta^* \mapsto \Sigma^*$ BPE

Input: Σ , $C = \{x^{(m)}\}_{m=1}^M \subset \Sigma^*$. Return: $\Delta, t: \Sigma^* \to \Delta^*$, Initialize Δ to Σ , Find the most frequent merge in C, where a merge m is a concatenation of two elements in Δ , so now $\Delta \leftarrow \Delta \cup m$.

Suprious Ambiguity

 $aaba \rightarrow |a|a|b|a| \text{ or } |aa|b|a|$

Pushforward

 $p(x) = \sum_{y \in \Delta^*} p_{\Delta^k}(y)$ $v \in t^{-1}(x)$ $t^{-1}(x) = \{y | y \in \Delta^*, t(x) = y\}$

7 Sampling Ancestral Sampling

1. Locally normalize.

2. Sample $y_t \sim p(\cdot | y_{< t})$, stop 10 Prompting when $v_t = EOS$. May not halt \rightarrow set max string length.

Greedy

$x_i = \operatorname{argmax}_{x \in \Sigma^*} \log(x | x_1 \dots x_{i-1})$

Sampling Adaptors

To calibrate *p* we can postprocess probabilities by a function $\alpha: \Delta^{|\Sigma|-1} \to \Delta^{|\Sigma|-1}.$

Top-K Sampling

Set $p(y_t | y_{\le t}) = 0$ for all but the K most probable tokens, and renormalize.

Nuclues Sampling

Only take top p% of probability mass.

8 Transfer Learning

ELMo

Fwd & Bwd LM using L LSTM layers. The ELMo representation for a token y_t is **ELMo**^{task} = $\gamma^{\text{task}} \sum_{l=0}^{L} s_l^{\text{task}} \mathbf{h}_{tl}^{\text{LM}}$ where $s_l^{\text{task}} \ge$ 0, $\mathbf{h}_{tl}^{\mathrm{LM}} = (\overrightarrow{h}_{tl}^{\mathrm{LM}}, \overleftarrow{h}_{tl}^{\mathrm{LM}}).$ **BERT (encoder)** $\mathcal{L} = \mathcal{L}_{MLM} + \mathcal{L}_{NSP}$ $\mathcal{L}_{\mathsf{MLM}}(\boldsymbol{\theta}) = \sum_{i=1}^{N} \sum_{t=1}^{T}$

 $\log_{\mathrm{MLM}}(\boldsymbol{y}_{t}^{(i)}|\boldsymbol{y}_{< t}^{(i)}, \boldsymbol{y}_{> t}^{(i)}; \boldsymbol{\theta})\mathbb{I}\{\boldsymbol{y}_{t}^{(i)} = \mathsf{M}\}$

9 Parameter Efficient Fine-Tuning BitFit

 $\mathbf{Q}(\mathbf{x}) = \mathbf{W}_a \mathbf{x} + \mathbf{b}_a$ $\mathbf{h}_4 = \text{GELU}(\mathbf{W}_2 \cdot \mathbf{h}_3 + \mathbf{b}_2)$ **Diff Pruning**

 $\theta_{\rm FT} = \theta_{\rm LM} + \delta$. Encourage δ_{Diff} to be sparse by regularization by a differentiable approximation to the L_0 -norm penalty as $\|\delta_{Diff}\|_0$.

Adapters

Insert bottleneck MLPs after each sublayer (MHA and FFN). $\mathbf{h} \leftarrow \mathbf{\dot{h}} + f(\mathbf{h}\mathbf{W}_{down})\mathbf{W}_{up}$

LoRA

Replace weight matrices $W \in$ $\mathbb{R}^{d \times r}$ with $W \leftarrow W + \frac{\beta}{h}AB$ where $A \in \mathbb{R}^{d \times b}$ and $B \in \mathbb{R}^{b \times r}$ are random matrices and β is a constant in *b*. $N_{\text{param}} = NH(3b(d+r)+2bd)$

Objective

 $\hat{z} = \operatorname{search} P(f_{\operatorname{fill}}(x', z); \theta).$

Discrete Prompts

Use the middle words/paths as templates in the form of [X] middle words [Z], translating the prompt into another language and back, use a thesarus to replace words, training a text generation model for generating prompts

Prefix Tuning

Prepends a sequence of continuous task-specific vec- i tors to the input while keeping the LM parameters fro- Scoring $\max_{\phi} \log P(\boldsymbol{y}|\boldsymbol{x};\boldsymbol{\theta};\boldsymbol{\phi})$ zen. $\max_{\phi} \sum_{v_i} \log P(y_i | h_{\leq i}; \theta; \phi)$

Self-Consistency

Generate multiple reasoning paths and selecting the most frequent final answer.

11 Vision Language Models **Vision Encoders**

al features as well as location features $[x_1, y_1, x_2, y_2, w, h, w \times h];$ **CNN; Image Transformers** Create image tokens: Split image into image patches, map them into vectors and linearly project them to patch embeddings. Add a lear- the LM. nable special token [CLS]

Multimodal Fusion

Fusion encoder takes both v = on the prefix. $\{v_1, \dots, v_M\}$ and $w = \{w_1, \dots, w_N\}$ as input, and learns contextua- 1. Collect a dataset of instructilized multimodal representations $\tilde{v} = {\tilde{v}_1, \cdots, \tilde{v}_M}$ and $\tilde{w} =$ (concat v, w), co-attention v, ware fed into different Transformer

blocks independently. Masked Language Modelling $\mathcal{L}_{MLM}(\theta)$

 $= -\mathbb{E}_{(\tilde{\mathbf{w}},\tilde{\mathbf{v}})\sim D} \log P_{\theta}(\tilde{\mathbf{w}}_{\mathbf{m}} | \tilde{\mathbf{w}}_{\backslash \mathbf{m}}, \tilde{\mathbf{v}})$

Image Text Matching

 $\mathcal{L}_{\text{ITM}}(\theta) = -\mathbb{E}_{(\tilde{\mathbf{w}}, \tilde{\mathbf{v}}) \sim D}[y \log s_{\theta}(\tilde{\mathbf{w}}, \tilde{\mathbf{v}}) + \mathbf{14}$ Callibration $(1-y)\log(1-s_{\theta}(\tilde{\mathbf{w}},\tilde{\mathbf{v}}))])$ ECE Image-Text Contrastive Learning

$s_{i,j}^{i2t} = v_i^{\top} w_j, \ s_{i,j}^{t2i} = w_i^{\top} v_j$ $\mathcal{L}_{\text{ITC}}^{i2t}(\theta) = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(s_{i,i}^{i2t}/\sigma)}{\sum_{i=1}^{N} \exp(s_{i,i}^{i2t}/\sigma)} \quad \mathcal{Y}_{i}), \quad \text{cont}(B_{m}) = \sum_{i \in B_{m}} p_{i}$ 15 Security & Adversarial examples

$$\mathcal{L}_{\text{ITC}}^{t2i}(\theta) = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(s_{i,i}^{t2i}/\sigma)}{\sum_{i=1}^{N} \exp(s_{i,i}^{t2i}/\sigma)}$$

Masked Image Modelling

 $\mathcal{L}_{\mathrm{MIM}}(\theta) = \mathbb{E}_{(\tilde{\mathbf{w}}, \tilde{\mathbf{v}}) \sim D} P_{\theta}(\tilde{\mathbf{v}}_{\mathbf{m}} | \tilde{\mathbf{v}}_{\backslash \mathbf{m}}, \tilde{\mathbf{w}})$ 12 RAG **TF-IDF**

$$\begin{aligned} f \text{-idf}(t, d, \mathcal{D}) &= \text{tf}(t, d) \times \text{idf}(t, \mathcal{D}) \\ f(t, d) &= \log(1 + \operatorname{freq}(t, d)) \\ df(t, \mathcal{D}) &= \log\left(\frac{|\mathcal{D}|}{|d \in \mathcal{D}: t \in d|}\right) \end{aligned}$$

score(**q**, **d**) =
$$\sum_{t \in q} \frac{\text{tf}(t, d)}{|d|}$$

Dense Retrieval

 $\operatorname{SIM}(q, d) = \operatorname{ENC}(q)^T \cdot \operatorname{DEC}(d).$ $L(q_i, d_i^+, d_{i,1}^-, ..., d_{i,n}^-)$ $-\log \frac{e^{\sin(q_i,d_i^+)}}{e^{\sin(q_i,d_i^+)} + \sum_{i=1}^{n} e^{\sin(q_i,d_{i,j}^-)}}}$

kNN-LM

Pretrained **OD** then use visu- Store all embedded prefixes and their following words in a database. At inference time, retrieve the *k*NN of a prefix, normalize exponentiated distances to a pron distr p_{ξ} over words. Then sample from a convex combination of p_{ξ} and

Dynamic Gating: Set the weighting of distributions depending

13 RLHF

- ons+answers and train a super- Cryptanalysis vised baseline model.
- by the baseline model, score thing.

them manually and train a re- Recovering hidden dim ward model.

3. Use PPO to fine-tune a LM (the policy) using the reward model.

 $\sum_{m=1}^{M} \frac{|B_m|}{M} |acc(B_m)|$ – = $conf(B_m)$ $\operatorname{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{I}(\hat{y}_i)$

Greedy Coordinated ∇ **-Descent**

1. Find top-*k* token substitutions according to gradient 2. Pick B substitutions at random across all suffix tokens 3. Evaluate the loss of all *B* candidates and pick the best.

Watermarking

Bias LLM away from true $p(s_{i+1}|s_i)$ in a subtle but verifiable way.

16 Prompt Injections Indirect Attacks

1. The adversary plants *indirect* prompts on a source. 2. LLM retrieves the poisoned prompt.

Defenses

Escape data; Detect with 2nd LLM; Separate pipelines; instruction hierarchy; **quarantined LMs**.

17 Data poisoning, backdoors and model stealing

Poisoning Wikipedia

1. Estimate when each article was snapshot in the past dump. 2. Poison each article right before it.

Defenses

Integrity Check (but many FP), randomize snapshot times and keep edits longer than some time only.

Distillation

Train your own LLM to replicate the behavior of the original LLM.

Dream of viewing LLM operati- $\{\tilde{w}_1, \dots, \tilde{w}_N\}$. merged attention 2. Produce a dataset of compari- ons as cryptographic mechanisms sons of different answers given and precisely recovering every-

 $LLM(\mathbf{x}_i) = \mathbf{y}_i = \mathbf{z}_i \mathbf{W}^{\top}$, so $\mathbf{Y} = \mathbf{Z}\mathbf{W}^{\top}, \ \mathbf{Y} \in \mathbb{R}^{n \times V}, \ \mathbf{Z} \in \mathbb{R}^{n \times h},$ $\mathbf{W}^{\top} \in \mathbb{R}^{h \times V}$, *h* is rank of **Y**. Can recover part of weights using SVD $\mathbf{Y} = \mathbf{U} \Sigma \mathbf{V}^{\top}, \ \mathbf{U} \in \mathbb{R}^{n \times h}, \Sigma \in \mathbb{R}^{h \times h},$ $\mathbf{V}^{\top} \in \mathbb{R}^{h \times V}$, where \mathbf{V}^{\top} are the weights up to an $h \times h$ transform. In practice, you only get top-k tokens, you can attack OpenAI's via logit_bias by repeatedly setting the top-k to $-\infty$ and push other values upwards.

18 Privacy

Federated Learning

Central server aggregates gradient updates from multiple clients. Issue: gradients are not private! Given a gradient g find $x_1 \dots x_B$ to minimize $||g - \frac{1}{B} \sum_{i=1}^{B} \nabla_{\theta} \mathcal{L}(f_{\theta}(x_i))||$

Weight-Trap Attack

Server sends client model f_{θ} s.t. $\nabla_{\theta} \mathcal{L}(f_{\theta}(x_i)) = x_i$

Differential Privacy

An algorithm *M* is ε -differentially private if for any "neighboring" databases D_1, D_2 that differ in a single element, and any output S we have: $P|M(D_1) \in$ $S \leq \exp(\varepsilon) P[M(D_2) \in S]$ Post-Processing: If M is ε -DP, then f(M) for any function f is also ε-DP.

Composition: If M_1 is ε_1 -DP and M_2 is ε_2 -DP then $f(M_1, M_2)$ is $(\varepsilon_1 + \varepsilon_2)$ -DP.

Sensitivity

 $\Delta f = \max |f(\mathcal{D}_1) - f(\mathcal{D}_2)|$, release $y \sim \text{Laplace}(f(\mathcal{D}), \frac{\Delta f}{s})$

Noisy Gradient Descent

 $\mathbf{g} = \frac{1}{k} \sum_{x_i \in B} \nabla_{\theta} \mathcal{L}(f_{\theta}; \mathbf{x}_i) + \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ giving us a relaxed notion of DP with prob $1 - \delta$. Bound sensitivity by, clipping gradients to norm C. The entire training algorithm is then ε' -DP for some $\varepsilon' > \varepsilon$. Subsampling amplifies privacy: the gradient is $\approx (k\varepsilon)/(|\hat{D}|)$ -DP w.r.t. neighboring datasets. The final DP budget is in $(O)(\sqrt{N\varepsilon})$.

