

6.3 BERT Variants: Transformer Encoder models.

Robert: 10x longer training, longer input seq, dyn V

Albert: (i) Cross-Layer parameter sharing (parts of Transf. layers). (ii) Factorized Embed. param.: decompose the large embed into 2 smaller matrices that project to a lower dim embed space & then project to hidden space. (iii) Sentence order prediction (SOP): builds neg examples by swapping order of A/G -> focus on coherence. (iv) BERT-specific reduces param count, the computation cost remains.

ELECTRA: train with a GAN, mask words & detect mask

6.4 GPTs: Transformer Decoder Models

"Generative Pre-Training Transformer": unidirectional, autoregressive LM (prob $p(y|x)$)

Text completion problem: enabling few-/zero-shot & in context learning. Takes two input (x_{left} , x_{right}) & output (y) tokens, predicted by "fed-infix" style.

6.5 Seq2seq LMs (T5): Encoder-Decoder Models good at seq-to-seq generation tasks (translation, summarization)

7. Param. Efficient Fine-Tuning

7.1 Partial Fine-Tuning: Motivation: catastrophic forgetting -> fine-tune a few different parameters leaving majority of params unchanged. No regularizations!

Off Tuning: $\theta_{\text{FT}} = \theta_m + \delta \theta_{\text{FT}}$. Train $\delta \theta_{\text{FT}}$ w/ L_2 -norm penalty to introduce sparsity. New params \rightarrow more GPU memory.

7.2 Adapter Tuning: Insert small modules called adapters to a model. Adapters can be any architecture & can be placed anywhere in the model. Heuristics: Place a 2-layer FFNN w/ a bottleneck after each sublayer (Both: multi-head attention Sublayer & Feedforward sublayer in the transformer): $h \leftarrow h + f(h \downarrow \text{down})$. Wip. Formally: Sampling adapt.: $a: A^{d_h \times d_a} \rightarrow A^{d_a \times d_h}$ maps products to prob. Only first d_a columns are learned. Prefix Tuning: optimizes prefix for each task independently. Prefix Activ. V layers (f_{prefix}) are learned. Collected in a trainable matrix $P_h: h \leftarrow P_h^{[A^{d_h \times d_h}]} [A^{d_h \times d_a}, \text{else}]$. Only option loss: min $L_{\text{FT}}(h) = E_{(a,c)} p_i(y \mid \log(s_i(a, c)) + (1 - \log(s_i(a, c)))$

Only first two bias in gradient module. P_h (h) \rightarrow repeat 85% performance across 80+ params.

7.3 LMs: Pre-Training Obj. Masked LM: given image-text pair, randomly mask out input words w/ prob 15% & predict masked tokens (w/ based on surrounding words & paired image) to minimize $L_{\text{LM}}(h) = -E_{(a,c)} \log(p(\text{Cntrm}(a, c)))$

Image-Text Matching: given image-caption pairs -> identity, which are correct pairs \rightarrow binary classification problem. ILSV tokens is pre-trained to learn cross-modal repr.

Prefix Activation V layers: $f_{\text{prefix}}: h \leftarrow P_h^{[A^{d_h \times d_h}]} [A^{d_h \times d_a}, \text{else}]$ generalizes better if it's pre-trained for each task individually.

7.4 LORA: Update original weight W to low rank matrix $W = W_b + W_a \cdot B_A$ w/ $B_A \in \mathbb{R}^{d_h \times d_k}$. R is d_k. hit A w/ rand. gauge distil. & B w/ O.5. Engulf of scaling. Forward pass: $h \leftarrow W \cdot h + \frac{\sigma}{d_h} \Delta W \cdot h = W \cdot h + \frac{\sigma}{d_h} \Delta W \cdot h$. Backprop: low computational burden, flexibility may fine-tune any part of the model. Wip: regularization loss 3% of prob.

8. Prompting & Zero-Shot Inference

8.1 Prompting: Map an input x to a prompt $x' = \text{prompt}(x)$. Design $f_{\text{prompt}}(x) = \text{map}[x][m][z]$, where z is the answer. If $x \neq z$, fills prompt w/ x & potential answer z .

8.2 Prompt Engineering: Good: find best prompt w/ "tokens". Bad: find worse prompt. Automated prompts: Discrete prompts/hard prompts. Hard to find good prompts in discrete space.

8.3 Zero & Few-Shot Inference / In-Context

9. Visual LMs

9.1 Vision & Language Tasks: Training: 8 new task by providing a few examples. No model.

Image-text tasks: include text & images in an input (we add visual).

Visual Question Answer: Given image, LQ to produce A.

Image captioning: Given image \rightarrow generate caption.

Image-text retrieval: Given image, LQ -> soft-relevant post.

Visual grounding: Given image, filtering expression, bounding boxes of obj. & concepts \rightarrow predict bounding boxes of input text query.

Text to image gen: Given caption \rightarrow generate image.

CV Tasks as LM Problems: Instead of treating class labels as 1-hot encodings, use semantic meaning behind labels to generalize to unseen concepts.

9.2 VLMs Architecture: Consists of: text encoder, vision encoder, fusion module (L decoder).

Albert: Given an image-text pair, a VLM model first extracts features $w = [w_{\text{img}}, w_{\text{txt}}]$ & visual features $v = [v_{\text{img}}, v_{\text{txt}}]$ via text & visual encoder. w & v feed into multimodal fusion model to produce cross-modal repr. \rightarrow optimize fed into decoder.

Visual Encoder: 3 types:

- Obj. Detector:** Start w/ pretrained OO (A-CNN) & entr. location features: $[x_1, y_1, x_2, y_2, w_1, h_1]$ (rect, coord)
- h:** visual features are fed through fully conn layer
- CNN:** Shift invariant. Strong LMs \rightarrow strong abstractions.
- Vision Transf.:** take image \rightarrow transfer image to another language patches, map them into vectors & linearly proj. to patch embed. (soft or learnable special token (CLS))

Text completion problem: enabling few-/zero-shot & in context learning. Takes two input (x_{left} , x_{right}) & output (y) tokens, predicted by "fed-infix" style.

9.3 VLMs Pre-Training Obj.:

Masked LM: given image-text pair, randomly mask out input words w/ prob 15% & predict masked tokens (w/ based on surrounding words & paired image) to calculate L2-dist to training contexts. \rightarrow Collect the values of the top-k contexts \rightarrow Normalize (softmax). $p(h_i \mid \text{mask}(d_i)) = \frac{1}{Z} \sum_j \exp(\text{log}(p(\text{Cntrm}(d_i, j)))$

Image-Text Matching: given image-caption pairs -> identity, which are correct pairs \rightarrow binary classification problem. ILSV tokens is pre-trained to learn cross-modal repr.

9.4 12 Models: prompting

REALM: fuse artefact in input layer by concat.

Decompose $p(y|x)$ into: $p(x) = \sum_z p(y|z)x)p(z|x)$

Local Input: $x \in \mathbb{R}^d$ (pre-train cap), retrieve useful words $z \in \mathbb{R}^d$ by sampling $p(z|x) = \text{exp}(\text{ENC}(z) \cdot \text{ENC}(x))$

Approx. sum over Z as sum over top-k

Local Context: My doc in predict step

Jointly optimize both components.

RETR: fuse artefact in the intermediate layers by crossatt with chunked cross-attention mechanism.

is_retriever: $\text{exp}(\text{ENC}(z) \cdot \text{ENC}(x))$. Devide \leftrightarrow into this. For each text chunk C, top-k nearest neigh. RET(C_i) are filtered using L2-dist on BERT-emb.

isEncoder: $E_h = \text{ENC}(\text{RET}(C_i), h_i)$ where $i \in M$ idx of neigh. & h_i intermediate activ of chunk C_i

KNN LM: fuse artefact in output layer by interpolating LM output prob. 1) compute repr. of all sentence proto (training contexts) & store them in mapping to the following word (target). $\text{Datastore} = \mathbb{R}^{[d_h] \times [(\text{ENC}(X_i), Y_i)]}$

For new query, at inference: embed query as $\text{ENC}(q)$ & calculate L2-dist to training contexts. \rightarrow Collect the values of the top-k contexts \rightarrow Normalize (softmax). $p(h_i \mid \text{query}(q)) = \frac{1}{Z} \sum_j \exp(\text{log}(p(\text{Cntrm}(q, j)))$

Attacking VLMs: Apply GD/FGSM on img to gen "bad" output. L2-maximize $P(\text{Sure}(f(\text{img})) \mid P(\text{heat} \mid \text{hitting}(\text{Sure}))$. L2-optimize w/ suffixes to output agreeable responses. L2-maximize $P(\text{Sure}(\text{!query-suffix})) \cdot P(\text{!heat} \mid \text{!query-suffix})$.

Attacking full-only LMs: Naive options: (i) manually search for suffix. (ii) use another LM to search for suffix. Hct generates adversarial output.

13.3 Adv Examples for LM's

Objective: Misclassification, based on dangerous outputs.

Bad/Good: Optimize for "agreeable outputs" to "bad" prompts. LM's cannot backdoor once committed to a response.

Optimize w/ suffixes to output agreeable responses.

Maximize $P(\text{!query} \mid \text{query-suffix}) \cdot P(\text{!heat} \mid \text{!query-suffix})$.

13.4 Poisoning

Poisoning Caption image D: LATION: (URL, t) pairs, URL to img, t text & check if substring (X) is present. Is static. URLs expire -> buy expired URLs to inject malicious ones.

Poisoning Wikipedia: change article shortly before day.

Backdooring: Challenge: Backdoor needs to persist after instruction tuning & alignment. Approach: Format poisonous documents as chat that uses losses.

Deduplication: D_{train} contain duplicates. Duplication causes memorization.

17 Model Stealing

17.1 Distillation: Mine a large LM with a smaller LM. Given prompt-response pairs, Distilled LM learns to mimic. Use this model to perform transfer attacks.

Steal reasoning traces & train DNN on them.

17.2 LogReg: Adversary queries model to obtain pred. $y=x_{\text{pred}}$ & conf. score $f(x)$, \rightarrow obtain $x' = b + \delta^*(f(x))$. If we're δ , then with d1 random queries of clean inputs, adversary can solve for b & b .

Label-only Attacks: If model only outputs $y = \text{sign}(f(x))$ pick 2 samples x_1, x_2 of different classes. Then by binary search on the line $b + \delta \cdot k$ is found & that lies on the boundary $f(x)=0$. Given d-1 such points, we can recover b & b .

17.3 Stealing Deep Models: for RNN-NNs. Identify linear regions in model's input space & then solve for params that del. the linear fit.

17.4 Attacking Transformers: Goal: find hidden dim h . Transformer: $f: \mathbb{R}^{[d_h] \times [d_h]} \rightarrow \mathbb{R}^{[d_h]}$. prompt of N tokens from vocal $V \rightarrow$ prob dists over $V \rightarrow f(p) = \text{softmax}(Wg(p))$. $W \in \mathbb{R}^{[d_h] \times [d_h]}$. embed proj. mat. $g: \mathbb{R}^{[d_h]} \rightarrow \mathbb{R}^{[d_h]}$. rest of transf. Normally $M \in \mathbb{R}^{[d_h] \times [d_h]}$.

Recovering b : (i) Attacker has access to logit vec Wg . Logit-vec S lie in hidden subspace bcc last layer is a projection $P \rightarrow R^h$ ($IV \mid h$). Approach: (1) Query the model on random prompts p . & obtain logit vectors $Wg(p)$. (2) Construct mat Q with each column being a legit vec. Note that $Q = W \cdot H$. (3) Return rank of Q .

18 Privacy

18.1 Cryptographic Privacy: Secure training: Adversary only learns final model. Secure inference: Client calls ML model on Server.

Client learning f(x) for input x & Server learns nothing.

Prac. ex: Large overhead (computation & communication)

18.2 Encoded D: Idea: encode D ($E(D)$), send to central server, use encoded D to train model f .

Public: Algo: f (algo for WM publicly known) \rightarrow allow server to learn a useful model f but should not allow server to learn "too much" about D . Issue: impossible to simultaneously satisfy!

Instance Hiding: For img: Idea: combine private image x w/ public images $x_1, x_2: x = x_1 + x_2$. x_1, x_2 are random imgs. Obs: doesn't hurt training, but improves it. Improvement: scale each pixel in x to 1/2 & randomly flip signs.

Text Hide: Similar, but in embed. space.

Federated Learning: Central server sends current model f to clients \rightarrow clients update locally. **Attack:** $\hat{f}(D)$ -> adversarial noise for f .

Guess x_2 (D): x_2 is random. Even if x is not easily compressible, Compress x (D) w/ $f(x)$ for $x \in \mathcal{X}$ & f make loss L > 0.5.

Prac.: D only works on datasets/worst-case analysis.

21.2 Challenges for MI on LMs:

Choice of K: $K = (\mathbb{R}^{[d_h] \times [d_h]}, \text{Log}(f(\text{!heat} \mid \text{!query})))$ (softmax)

For Token t : $\log(\text{!heat} \mid \text{!query}) = \log(f(t))$ multiv. $\text{Log}(f(t))$ -> softmax $(\text{Log}(f(t)))$

Edit PLR: $\text{Log}(f(t))$ -> adversarial noise for f .

Guess x_2 (D): x_2 is random even if x is not easily compressible. Compress x (D) w/ $f(x)$ for $x \in \mathcal{X}$ & f make loss L > 0.5.

Prac.: D only works on datasets/worst-case analysis.

Prac.: $f(t)$ before D , $x_2 \in \mathcal{X}$, $f(t) \in \mathcal{F}$, $\text{Log}(f(t))$ -> softmax $(\text{Log}(f(t)))$

Random Corr: Idea: f before D , $x_2 \in \mathcal{X}$, $f(t) \in \mathcal{F}$, $\text{Log}(f(t))$ -> softmax $(\text{Log}(f(t)))$

Random Corr: Idea: f before D , $x_2 \in \mathcal{X}$, $f(t) \in \mathcal{F}$, $\text{Log}(f(t))$ -> softmax $(\text{Log}(f(t)))$

2