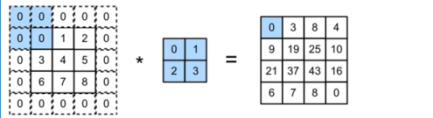
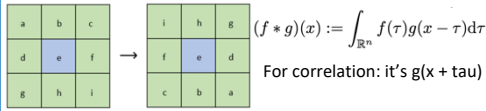
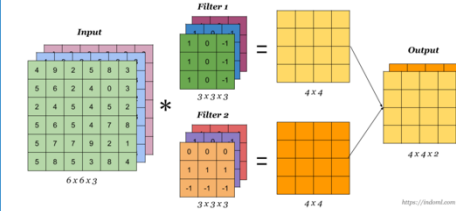


Convolutions & CNNs

Output shape (width or height): $\text{floor}((W - K + 2P) / S) + 1$
Output shape (depth/num of channels): num of filters applied



#params of 1 filter: filter H x filter W x filter C + 1 (bias)
Fix mean and variance

$$\mu_B = \frac{1}{|B|} \sum_{i \in B} x_i \text{ and } \sigma_B^2 = \frac{1}{|B|} \sum_{i \in B} (x_i - \mu_B)^2 + \epsilon$$

and adjust it separately

$$x_{i+1} = \gamma \frac{x_i - \mu_B}{\sigma_B} + \beta$$

Conv = only equivariant to translation; not equivariant to warp, flipping... hence data aug for better model general. too
CNN approx. invariant if adding pooling layers.

Shift invariance: property that describes a system's unchanging response when the input is shifted; useful in CNNs (e.g. obj recog.) vs **shift equivariance** (e.g. obj detection, segmentation)

Assumptions: translation invariance (shift in input should simply lead to shift in the hidden representation) + locality

Pooling layer ensures **approximate translation invariance**

Receptive field of a layer k

$$l_k = l_{k-1} + ((f_k - 1) \times \prod_{i=1}^{k-1} s_i)$$

where l_{k-1} is the receptive field of layer $k-1$, f_k is the filter size (height or width, but assuming they are the same here), and s_i is the stride of layer i .

For real-valued functions, of a continuous or discrete variable, conv = cross-correlation of $f(x)$ & $g(-x)$, or $f(-x)$ & $g(x)$

CNN architecture

- convolutions implemented through weight sharing, interpreted as weights of a filter function
- more output channels => NN can learn more complex + high-level features. Padding to preserve input shape.
- final FC layers: can be comp. + memory expensive (bottleneck w/ increasing num of classes, e.g. 1000 for ImageNet)
- 1x1 conv acts like an MLP per pixel, which aggregates across channels of input; introduces complexity + nonlinearity
- **AlexNet** for ImageNet: deeper + bigger/wider LeNet. Add dropout, sigmoid => ReLU, maxpool, heavy data aug, model ensembling (model averaging across multiple well-performing CNN models to achieve SOTA results) vs LeNet for MNIST
- **VGG:** group layers into blocks (num of blocks varies). Blocks can be easily parameterized, creating a more organized, modular arch. Simplifies design process + easier to fine-tune model for specific learning tasks.
- More layers of *narrow* convolutions outperforms using fewer *wide* ones: lots of simple fns > few complex fns
- **Inception:** deep + parallel paths with blocks to capture different types of features more effectively. Combines benefits of various convolutions & pooling operations while optimising computational cost
- **Batchnorm** normalises features within each minibatch, stabilising training process & speeding up convergence. Regularisation by noise injection + no dropout needed (both control NN capacity) + ideal minibatch 64-256). B = batch.
- Test: fix the gamma & beta learned in training + instead batch statistics, use running average for mean and variance
- **ResNet:** get increasingly powerful AND nested functions (might not be convex) with more layers. "Taylor expansion" style parametrization. Input to act fn becomes $f(x) + x$ (x sometimes * with 1x1 conv to change dim). Allows for deeper NNs w/ fn classes more likely to be nested. Better grad flow (solve vanishing). ResNet module: multiple ResNet blocks.
- **DenseNet** uses higher-order Taylor series expansion; feature maps reuse; may need to reduce res (transition layer)
- **curse of dim:** As the number of features or dimensions grows, the amount of data we need to generalise accurately grows exponentially
- **pooling:** permutation-invariant aggregation + downsampling; reduces res; hierarchical features; max-pooling breaks shift equivariance. **Separable filter:** filter can be written as conv of 2/more simple filters: e.g. 2D filter from 1D filters; and $\det(\text{filter})=0$.

Commutativity: $f * g = g * f$

Associativity: $f * (g * h) = (f * g) * h$

Distributivity: $f * (g + h) = (f * g) + (f * h)$

Associativity with scalar multiplication: $a(f * g) = (af) * g$

Derivative: $D(f * g) = (Df) * g + f * Dg$

VAE

The definition of divergence is weaker than distance: does not need to satisfy symmetry or triangle inequality

We can work with PDFs:

$$p(x) = \frac{dP}{dx}, \forall P \in \mathcal{P}$$

Probabilistic graph models (joint distribution): used to describe dependency structure

$$p(x_1, \dots, x_D) = \prod_{i=1}^D p(x_i | p_{\setminus i}(x_i)),$$

$$p(x_1, x_2, z_1, z_2) = p(z_1)p(z_2)p(x_1|z_1)p(x_2|z_2).$$

$$q(z_1, z_2|x_1, x_2) = q(z_1|x_1)q(z_2|x_2).$$



KL Divergence: minimize this in order fit a distribution to a target one, asymmetric: $KL[p||q] \neq KL[q||p]$ in general.

$$KL[p||q] = \mathbb{E}_{p(x)}[-\log q(x)]$$

$$\geq -\log \mathbb{E}_{p(x)}[q(x)] \quad (\text{Jensen's inequality})$$

$$= -\log \int p(x) \frac{q(x)}{p(x)} dx = -\log 1 = 0,$$

$$KL[p||q] = \int p(x) \log \frac{p(x)}{q(x)} dx, \quad p, q \in \mathcal{P},$$

Solution: use the variational lower bound as a tractable approximation to marginal log-likelihood

$$\log p_\theta(x) = \mathbb{E}_{q(z)}[\log p_\theta(x|z)] - KL[q(z)||p(z)] := \mathcal{L}(x, q, \theta).$$

maximising $\mathcal{L}(x, q, \theta)$ is equivalent to minimising $KL[q(z)||p_\theta(z|x)]$.

$$\text{Analytic form: } KL[q_\phi(z|x)||p(z)] = \frac{1}{2} (\|\mu_\phi(x)\|_2^2 + \|\sigma_\phi(x)\|_2^2 - 2(\log \sigma_\phi(x), 1) - d)$$

Variational auto-encoder approach defines q distribution as a neural network: $q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \text{diag}(\sigma_\phi^2(x)))$, $\mu_\phi(x), \log \sigma_\phi(x) = \text{NN}_\phi(x)$.

VAE optimization objective: $\phi^*, \theta^* = \arg \max \mathcal{L}(\phi, \theta)$, $\mathcal{L}(\phi, \theta) = \mathbb{E}_{p_{\text{data}}(x)}[\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)]] - KL[q_\phi(z|x)||p(z)]$.

$$:= \mathcal{L}(x, \phi, \theta)$$

Monte Carlo estimation: $\mathcal{L}(\phi, \theta)$ is still intractable, so we replace expectation with MC approximation:

$$\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] \approx \log p_\theta(x|z), \quad z \sim q_\phi(z|x).$$

$$\nabla_\phi \mathcal{L}(x, \phi, \theta) \approx \nabla_\phi \log p_\theta(x|z), \quad z \sim q_\phi(z|x).$$

$$\nabla_\phi \mathcal{L}(x, \phi, \theta) \approx \nabla_\phi \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - \nabla_\phi KL[q_\phi(z|x)||p(z)].$$

Reparametrization trick: directly sampling z from q and passing samples through the model to compute gradients is not differentiable wrt parameters ϕ of dist. from which z is drawn => hard to backpropagate

$$\mathcal{L}(\phi, \theta) \approx \frac{1}{M} \sum_{m=1}^M \log p_\theta(x|T_\phi(z_m; \epsilon_m)) - KL[q_\phi(z_m|x_m)||p(z_m)],$$

$z \sim q_\phi(z|x) \Leftrightarrow z = \mu_\phi + \sigma_\phi \odot \epsilon, \epsilon \sim \mathcal{N}(\epsilon; 0, I)$
Conditional VAE: generate data conditioned on additional information (class labels, viewing angle). y is the additional info

$$p_\theta(x|y) = \int p_\theta(x|z, y)p(z)dz,$$

If x is cont., then $p_\theta(x|z, y) = \mathcal{N}(x; G_\theta(z, y), \sigma^2 I)$, $G(z, y)$ is a neural network. We maximize a variational lower bound:

$$\phi^*, \theta^* = \arg \max \mathcal{L}(\phi, \theta),$$

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_{p_{\text{data}}(x, y)}[\mathbb{E}_{q_\phi(z|x, y)}[\log p_\theta(x|z, y)]] - KL[q_\phi(z|x, y)||p(z|x, y)],$$

or minimize KL divergence: $KL[q_\phi(z|x, y)||p_\theta(z|x, y)]$.

GAN. Constructs a binary classification task to assist the learning of generative model dist. $p_\theta(x)$ to fit data dist. $p_{\text{data}}(x)$. $\tilde{p}(x, y) = \tilde{p}(x|y)\tilde{p}(y)$, $\tilde{p}(y) = \text{Bern}(0.5)$, $\tilde{p}(x|y) = \begin{cases} p_{\text{data}}(x), & y = 1 \\ p_\theta(x), & y = 0 \end{cases}$

Fit a binary classifier (discriminator) by maximizing the MLE objective:

$$\phi^*(\theta) = \arg \max_\phi \mathcal{L}(\theta, \phi), \quad \mathcal{L}(\theta, \phi) := \mathbb{E}_{p_{\text{data}}(x)}[\log D_\phi(x)] + \mathbb{E}_{p_\theta(x)}[\log(1 - D_\phi(x))].$$

Generator fools discriminator by minimizing log prob. of making the right decisions $\theta^*(\phi) = \arg \min_\theta \mathbb{E}_{p_\theta(x)}[\log(1 - D_\phi(x))]$. Two-player game objective for GAN is: $\min_\theta \max_\phi \mathcal{L}(\theta, \phi)$.

MC approximation: the evaluation of the obj. can directly define sampling process of $p_\theta(x)$, which also defines the distribution in an implicit way:

$$\mathbb{E}_{p_\theta(x)}[\log(1 - D_\phi(x))] \approx \log(1 - D_\phi(x)), \quad x \sim p_\theta(x) \Leftrightarrow z \sim p(z), \quad x = G_\theta(z),$$

Jensen-Shannon divergence minimization: For a fixed generator, setting the gradient of GAN objective = 0 results in:

$$\mathcal{L}(\theta, \phi^*(\theta)) = 2 \left(\frac{1}{2} KL[p_{\text{data}}(x)||\frac{1}{2}(p_{\text{data}}(x) + p_\theta(x))] + \frac{1}{2} KL[p_\theta(x)||\frac{1}{2}(p_{\text{data}}(x) + p_\theta(x))] \right) - 2 \log 2,$$

where $JS[p_{\text{data}}(x)||p_\theta(x)]$ is the Jensen-Shannon divergence between $p_{\text{data}}(x)$ and $p_\theta(x)$. Since Jensen-Shannon divergence is a valid divergence measure, this means with infinite capacity for the generator, $\mathcal{L}(\theta, \phi^*(\theta))$ is minimised iff. $p_\theta(x) = p_{\text{data}}(x)$.

Alternative "non-saturated" objective given fixed discriminator: Saturation problem (near-perfect class. at start of training -> vanishing gradient); so, max. log prob. of making wrong predictions. $\max_\phi \mathbb{E}_{p_\theta(x)}[\log D_\phi(x)]$. **Conditional GAN** $p_\theta(x|y) = \int p_\theta(x|z, y)p(z)dz$.

Distribution form of p is defined implicitly by the sampling process:

$$x \sim p_\theta(x|z, y) \Leftrightarrow z \sim p(z), x = G_\theta(z, y),$$

$$\min_\theta \max_\phi \mathbb{E}_{p_{\text{data}}(x, y)}[\log D_\phi(x, y)] + \mathbb{E}_{p_\theta(x|y)}[\log(1 - D_\phi(x, y))].$$

In practice the component related to the generator parameters θ is computed by

$$\mathbb{E}_{p_\theta(x|y)}[\log(1 - D_\phi(x, y))] \approx \log(1 - D_\phi(G_\theta(z, y), y)), \quad z \sim p(z), y \sim p_{\text{data}}(y). \quad (19)$$

Using similar techniques, one can derive the optimal discriminator for a fixed generative model with parameter θ :

$$D_{\phi^*(\theta)}(x, y) = \frac{p_{\text{data}}(x, y)}{p_\theta(x|y)p_{\text{data}}(y) + p_{\text{data}}(x, y)}, \quad (20)$$

and with the optimal discriminator, minimising $\mathcal{L}(\theta, \phi)$ w.r.t. θ is equivalent to minimising the Jensen-Shannon divergence $JS[p_{\text{data}}(x, y)||p_\theta(x|y)p_{\text{data}}(y)]$.

L1vs L2vs Smooth

Loss	Formula	Notes
L2(MSE)	$\sum (x_n - y_n)^2$	Affected by outliers
L1	$\sum x_n - y_n $	Robust when dealing with noisy data, less sensitive to outliers. Non-differentiable at 0
Smooth L1	PIPPES	For values close to 0 its similar to L2, for others similar to L1. This is differentiable

Activations

Activation	Equation	Notes
Sigmoid	$\sigma(x) = 1/(1 + e^{-x})$	Binary Classification, Cause vanishing problem
Tanh	$\tanh(x) = (e^x - e^{-x}) / (e^x + e^{-x})$	Scaled sigmoid-cause faster convergence than sigmoid-average close to 0. Okay for regression and continuous reconstruction
ReLU	$\text{Max}(0, x)$	Dying ReLU problem
Leaky ReLU	$\text{Max}(x, 0.01x)$	Mitigates dying ReLU problem
Softplus	$1/\beta \cdot \log(1 + e^{\beta x})$	Smooth approximation of ReLU(>0), differentiable
ReLU6	$\text{Min}(\text{Max}(0, x), 6)$	Like ReLU but upper bound on 6
LogSigmoid	$\log(1/(1 + e^{-x}))$	Used within loss functions
Softmin	$(e^{-x}) / \sum(e^{-x})$	Convert into inverse probabilities
Softmax	$(e^x) / \sum(e^x)$	Convert into probabilities
LogSoftmax	$\log(e^x) / \sum(e^x)$	Handling Numerical Stability-used as loss

Metrics:

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$F1 = \frac{2TP}{2TP + FP + FN}$$

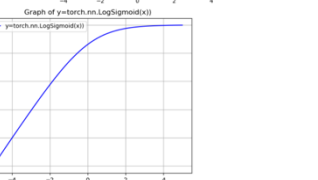
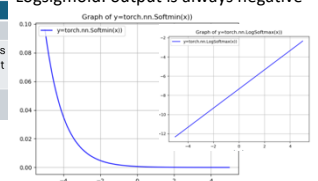
CE loss: combines logsoftmax(output) & NLL. Use for classif. Average across obsv. within minibatch

$$\text{loss}(x, \text{class}) = -\log \left(\frac{e^{x[\text{class}]}}{\sum_j e^{x[j]}} \right) = -x[\text{class}] + \log \left(\sum_j e^{x[j]} \right)$$

$$l_n = y_n \cdot (\log y_n - \log x_n) = y_n \left(\log \frac{y_n}{x_n} \right)$$

(above) KL div. loss: dist. btw proba distributions

Logsigmoid: output is always negative



Losses. L1/2/CE: "reduce" to scalar

$$\text{using mean or Smooth L1: } l_n = \begin{cases} 0.5(x_i - y_i)^2, & \text{if } |x_i - y_i| < 1 \\ |x_i - y_i| - 0.5, & \text{otherwise} \end{cases}$$

NLL:

Assumption: Network output represents log likelihoods.

Make the desired output as large as possible and all others as small as possible

$$\ell(x, y) = \mathcal{L} = \{l_1, \dots, l_N\}, \quad l_n = -w_n x_n y_n,$$

$$w_c = \text{weight}[c] \cdot 1/c \text{ (ignore index)}$$

$$\ell(x, y) = \begin{cases} \sum_{n=1}^N \frac{1}{\sum_{n=1}^N w_n} l_n, & \text{if reduction = mean} \\ \sum_{n=1}^N l_n, & \text{if reduction = sum} \end{cases}$$

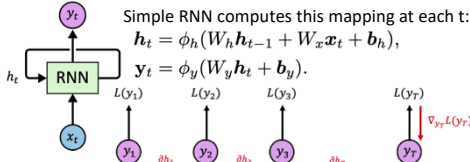
Binary CE loss: CE loss for only 2 classes.

Inputs as probas or logits

$$\ell(x, y) = \mathcal{L} = \{l_1, \dots, l_N\}, \quad l_n = -w_n [y_n \cdot \log x_n + (1 - y_n) \cdot \log(1 - x_n)]$$

RNNs

- can model dependencies within a sequence of arbitrary length



BPTT =>

$$\mathcal{L}(\theta) = \sum_{t=1}^T \mathcal{L}(y_t)$$

$$\frac{dh_t}{dW_h} = \frac{\partial h_t}{\partial W_h} + \frac{\partial h_{t-1}}{\partial W_h} \frac{dh_{t-1}}{dW_h} + \frac{\partial h_{t-2}}{\partial W_h} \frac{dh_{t-2}}{dW_h} = \dots$$

$$= \sum_{\tau=1}^t \left(\prod_{l=\tau}^{t-1} \frac{dh_{l+1}}{dh_l} \right) \frac{\partial h_\tau}{\partial W_h}$$

Can truncate: $\text{truncate} \left[\frac{dh_{t+1}}{dW_h} \right] = \sum_{\tau=t-L+1}^t \left(\prod_{l=\tau}^t \frac{dh_{l+1}}{dh_l} \right) \frac{\partial h_\tau}{\partial W_h}$

- weights/biases = shared across cells, where each cell = single-layer NN with 1 input and 1 output at t. no matter #cells/timesteps, no increase in #learnable params! BUT issue w/ gradients: vanishing or exploding => diff to train.

$$\frac{dh_{t+1}}{dh_t} = \phi'_h(W_h h_t + W_x x_{t+1} + b_h) \odot W_h$$

Some tricks: grad clipping; good W_h init; unitary/orthog. W_h

$$g \leftarrow \frac{\gamma}{\|g\|} g \quad \text{Hyperparam lambda; clip } g \text{ if } \|g\| > \eta$$

LSTMs

- Weights/biases shared across cells; sigmoid and tanh activation fns; final output/prediction = short-term memory outputted from last cell

- different paths for long-term & short-term memories to avoid gradient vanishing/expl problem from RNNs

- memory cell states and gates to control error flows

f_t : forget gate $f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$

i_t : input gate $i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$

o_t : output gate $o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$

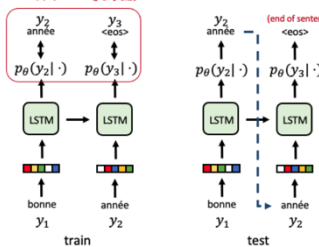
x_t : input $\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$

c_t : memory cell state $c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$

h_t : hidden state $h_t = o_t \odot \tanh(c_t)$

h0, c0 set to 0 or learnable. If 0, then we have the elems in c_t and h_t

Apply loss $\mathcal{L}(p_{\theta}(y_{1:t}), y_{1:t})$



GRUs

- improves simple RNN with gates like LSTM

- vs. LSTM: GRU removes the input/output gates and the cell state, but maintains the forgetting mechanism in some form

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z)$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t] + b_r)$$

$$\tilde{h}_t = \tanh(W_h \cdot [r_t \odot h_{t-1}, x_t] + b_h)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

Seq2seq models Given dataset of input-output sequence pairs $(x_{1:T}, y_{1:T})$, the goal of sequence prediction is to build a model $p_\theta(y_{1:T}|x_{1:T})$ to fit the data. Note here that the x, y sequences might have different lengths, and the input/output length T and L can vary across input-output pairs. So to handle sequence outputs of arbitrary length, we define an *auto-regressive model*

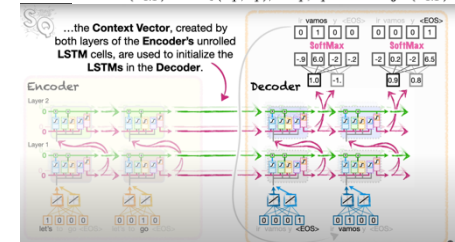
$$p_\theta(y_{1:L}|x_{1:T}) = \prod_{l=1}^L p_\theta(y_l|y_{<l}, v), \quad v = \text{enc}(x_{1:T})$$

Here $p_\theta(y_l|y_{<l}, v)$ is defined by a sequence decoder, e.g. an LSTM:

$$p_\theta(y_l|y_{<l}, v) = p_\theta(y_l|h_l^d, c_l^d), \quad h_l^d, c_l^d = \text{LSTM}_{\theta}^{\text{dec}}(y_{<l}, v), \quad (28)$$

and the decoder LSTM has its internal recurrent states h_l^d, c_l^d initialised using the input sequence representation $v = \text{enc}(x_{1:T})$. The encoder also uses an LSTM, meaning that

$$v = \text{enc}(x_{1:T}) = N N_\theta(h_T^e, c_T^e), \quad h_T^e, c_T^e = \text{LSTM}_{\theta}^{\text{enc}}(x_{1:T}).$$



Attention

Single-head attention: *scaled dot product attention*

$$A(x_j) = \frac{e^{x_j}}{\sum_{i=1}^n e^{x_i}} \quad Q = (q_1, \dots, q_N)^T \in \mathbb{R}^{N \times d_q}$$

$$K = (k_1, \dots, k_M)^T \in \mathbb{R}^{M \times d_k}$$

$$V = (v_1, \dots, v_M)^T \in \mathbb{R}^{M \times d_v}$$

Time and Space Complexity for 1 head

$$\mathcal{O}(N^2 d_q + N^2 d_v) \text{ and } \mathcal{O}(N^2 + N d_v)$$

It is clear that the time and space complexity figures of multi-head attention are h times of those for a single head but the extra costs for linear projections. Assume the projection matrices have sizes $W_q^h \in \mathbb{R}^{d_q \times d}$, $W_k^h \in \mathbb{R}^{d_k \times d}$, $W_v^h \in \mathbb{R}^{d_v \times d}$ and $W^h \in \mathbb{R}^{h d_v \times d}$. This means

$$\text{time complexity: } \mathcal{O}(h(MN\tilde{d}_q + MN\tilde{d}_v) + h(\tilde{d}_q(M+N) + \tilde{d}_v M) + \frac{N h \tilde{d}_v d_{out}}{h})$$

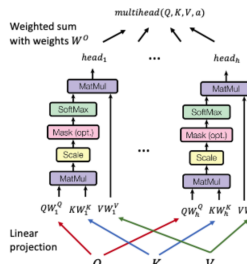
$$\text{space complexity: } \mathcal{O}(h(MN + N\tilde{d}_q) + h(\tilde{d}_q(M+N) + \tilde{d}_v M) + \frac{N d_{out}}{h})$$

soft attention: $a(\cdot)$ is the softmax function, $A_{ij} = \text{softmax}(\langle (q_i, k_j), \dots, (q_i, k_M) \rangle / \sqrt{d_q})$

hard attention: one-hot vector for each row with $j^* = \arg \max_j \langle q_i, k_j \rangle$ equal to 1

In transformers, layer normalisation is applied together with a residual link: $\text{Add\&Norm}(x) = \text{LayerNorm}(x + \text{Sublayer}(x))$, where $\text{Sublayer}(\cdot)$ can either be a multi-head attention block or a point-wise feed-forward network.

Multi-head attention: multiple alignment processes by projecting inputs into diff. subspaces, then performing dot product attention in subspace. Outputs are concatenated and projected



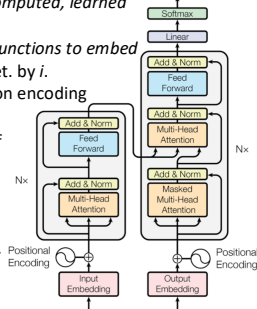
Transformers

Position encoding: attention is equivariant to row permutations bc non-linearity a(.) is applied row-wise. Ordering info needs to be added to attention inputs, can either be learned or computed, learned is good if we know max value of index

Sinusoid embedding: use multiple periodic functions to embed input index, frequency of sin/cos wave is det. by i.

Allows network to learn very flexible position encoding function

Point wise feed forward network: used as ff layer, multi-head attention (after Add\&Norm) returns a matrix of size $N \times d_{\text{out}}$, which can be processed "point-wise", treating rows in the attention outputs as "datapoint" inputs for next layer



Margin Ranking/Ranking/Contrastive Losses

$$\text{loss}(x, y) = \max(0, -y \cdot (x_1 - x_2) + \text{margin})$$

Useful to push classes as far away as possible and for metric learning. Practical: category that scores is closest or higher than correct one change until difference is at least the margin

This class of loss functions has a unique purpose—instead of directly predicting labels or values, they focus on predicting relative distances between inputs.

$$\ell(x, y) = \mathcal{L} = \{l_1, \dots, l_N\}^T, \quad \text{Triplet Margin Loss}$$

$$l_n(x_a, x_p, x_n) = \max(0, m + |f(x_a) - f(x_p)| - |f(x_a) - f(x_n)|)$$

actual positive sample negative sample

Make samples from same classes close and different classes far away.

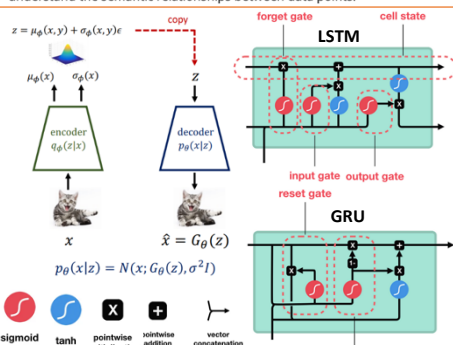
Objective: Distance for the good pair has to be smaller than distance to the bad pair. Actual distance does not need to be small, just smaller.

Used for metric learning and Siamese networks

Cosine Embedding Loss

$$\text{loss}(x, y) = \begin{cases} 1 - \cos(x_1, x_2), & \text{if } y = 1 \\ \max(0, \cos(x_1, x_2) - \text{margin}), & \text{if } y = -1 \end{cases}$$

Basically a normalised Euclidean distance. The Cosine Embedding Loss is especially useful for learning nonlinear embeddings and for semi-supervised learning tasks, where the objective is to understand the semantic relationships between data points.



Power rule:	$\frac{d}{dx}(x^n) = nx^{n-1}$	Product rule:	$\frac{d}{dx}[f(x)g(x)] = f(x)g'(x) + g(x)f'(x)$
Exponential (e):	$\frac{d}{dx}e^x = e^x$	Quotient rule:	$\frac{d}{dx}\left[\frac{f(x)}{g(x)}\right] = \frac{g(x)f'(x) - f(x)g'(x)}{[g(x)]^2}$
Exponential (general):	$\frac{d}{dx}a^x = a^x \ln(a)$	Chain rule:	$\frac{d}{dx}f(g(x)) = f'(g(x))g'(x)$
Natural log:	$\frac{d}{dx}(\ln(x)) = \frac{1}{x}$	Chain rule:	$\frac{dy}{dx} = \frac{dy}{dz} \cdot \frac{dz}{dx}$
General log:	$\frac{d}{dx}(\log_a(x)) = \frac{1}{x \ln(a)}$	Total for the backward pass: 2*forward pass	

FLOPs Total training flops $\approx 6 \times \text{dataset size in tokens} \times \text{number of parameters}$

Efficient training: grad accum., checkpointing

Gradient Accumulation enables training with larger effective batch sizes on limited hardware by computing and summing gradients across multiple mini-batches before performing a single update step. This method delays the optimizer step, thus accumulating gradients from several mini-batches, allowing for larger batch processing without increasing GPU memory usage.

Gradient Checkpointing reduces the memory footprint during training by selectively storing only a subset of the forward pass activations. Non-stored activations are recomputed during the backward pass as needed. This technique trades additional computation for lower memory usage, facilitating the training of larger models within memory constraints.

Finetuning The pretrained weights Φ , are updated $\Phi + \Delta\Phi$ according to $\max_{\Phi} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^T \log(P_{\Phi}(y_t|x, y_{<t}))$

Context-target tokens Autoregressive foundation model

Here, every parameter in Φ is updated therefore need to store gradients + optimisers' of every weight. For every downstream task a new model of size Φ is required + we need to relearn $\Delta\Phi$ parameters.

Low Rank Adaptation (LoRA)

Li et al. demonstrated that the model-loss overparametrised space has a low intrinsic dimension/rank. What if $\Delta\Phi$ also has low intrinsic rank we can approximate $\Delta\Phi$ with Θ where $|\Theta| \ll |\Phi|$ $\max_{\Theta} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^T \log(P_{\Phi+\Theta}(y_t|x, y_{<t})) = \text{only optimise big theta}$

• For a pretrained set up of $y = W_0 x$

• $y = XW_0 + \Delta W_0 \approx XW_0 + XL_1L_2$

Where, $L_1 \in \mathbb{R}^{D \times r}$ and $L_2 \in \mathbb{R}^{r \times D}$ and $r \ll D$

At inference time we combine: $W_0 + L_1L_2$

• No additional inference time (no extra steps to compute both XW_0 and XL_1L_2)

• We can always subtract L_1L_2 when we are done and combine with a different set of parameters L_1L_2



LoRA (Low Rank Adaptation) streamlines fine-tuning of pre-trained deep learning models by approximating updates to the model's weights with two smaller, low-rank matrices, A and B, where $A \cdot B$ approximates the weight update ΔW . This approach significantly reduces the computational resources needed for fine-tuning by decreasing the number of parameters to adjust, enabling efficient updates on modest hardware without substantially sacrificing performance.

QLoRA (Quantized Low Rank Adaptation) enhances LoRA by incorporating quantization, which further compresses the model updates by reducing the numerical precision of the A and B matrices. This method not only lowers memory usage but also accelerates computation, particularly on hardware optimized for lower precision arithmetic, making fine-tuning even larger models feasible on constrained hardware environments.

Deep learning models are more sensitive to underflow + overflow than to precision. Hardware requirements scale quadratically with fraction but not exponent.

LoRA recap: $y = XW_0 + XL_1L_2$ QLoRA looks to quantise W_0 from float16 to 4 bit!

HOW? Double Quantisation When required for computation we dequantise weights from NormalFloat4 to BFLOAT16

$y = X \text{dequant}(W_0) + XL_1L_2$

Rule Name	Property
Log of 1	$\log_1 1 = 0$
Log of the same number as base	$\log_b b = 1$
Product Rule	$\log_b(mn) = \log_b m + \log_b n$
Quotient Rule	$\log_b\left(\frac{m}{n}\right) = \log_b m - \log_b n$
Power Rule	$\log_b m^n = n \log_b m$
Change of Base Rule	$\log_b a = \frac{\log_c a}{\log_c b}$ (OR) $\log_b a \cdot \log_a c = \log_b c$
Equality Rule	$\log_a a = \log_b c \Rightarrow a = c$
Number Raised to Log	$b^{\log_b a} = a$
Other Rules	$\log_e a^n = \frac{n}{\ln a} \log_e a$ $-\log_a a = \log_a \frac{1}{a}$ (OR) $= \log_{\frac{1}{a}} a$

Conv

- Don't assume that in for a single kernel, the same kernel is applied to each channel of the input image
- Max pool and average pool are applied PER channel (not aggregating across channels of input image): only operate on width and height of input
- Mean and var learned gamma beta: to give new dist
- Valid padding = no padding at all
- Same padding = as much padding as needed to preserve the dims of input
- 0 padding needed for same padding: $(K - 1) / 2$
- Num of params for a conv layer with (more than 1) F filters of the same shape: $F * (\text{width of filter} \times \text{height of filter} \times \text{input channels} + 1)$ where 1 is bc we have 1 bias per filter
- Batchnorm output: same shape as input. Num learnable params: 2 x num of channels in input. It's 2 bc you have gamma and beta
- Pooling layers = no learnable params

Backprop / gradients

- Add the expression for dJ/dy_{hat} when using $J = \text{cross entropy loss with softmax being used}$:
 $1/m (y_{\text{hat}} - y)$
- If largest singular value of a weight matrix > 1 , then gradient explosion happens

DAG

- Conditional = joint / priors
 - o Hence for $q(z_1, z_2 \mid x_1, x_2)$: the priors of x_1 and x_2 cancel out and you don't write them as factors