

1 General Formula

Frobenius/L2 Norm $\|\theta\|_F = \sqrt{\sum_j \theta_j^2}$

L1 Norm $\sum_j |\theta_j|$

$$\frac{d}{dx} \left(\frac{f(x)}{g(x)} \right) = \frac{f'(x)g(x) - f(x)g'(x)}{g(x)^2} = \frac{f'(x)}{g(x)} - \frac{f(x)g'(x)}{g(x)^2}$$
$$\frac{d}{dx} \left(\frac{f(x)}{g(x)} \right) = \frac{f'(x)}{g(x)} - \frac{f(x)g'(x)}{g(x)^2}$$
$$\frac{d}{dx} \left(\frac{f(x)}{g(x)} \right) = \frac{f'(x)}{g(x)} - \frac{f(x)g'(x)}{g(x)^2}$$
$$\frac{d}{dx} \left(\frac{f(x)}{g(x)} \right) = \frac{f'(x)}{g(x)} - \frac{f(x)g'(x)}{g(x)^2}$$

2 Neurons

Signals travel via axon, from Soma. Dendrites get signals.

Membranes have Na (pump out) and K Channels (pump in) and a Na-K Pump (3 Na out for 2K in).

Membrane Potential Difference between voltage across membrane. At rest, value is -70mV , enforced by Na-K Pump.

Action Potential Electrical impulse travelling along axon to the synapse.

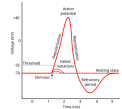
2.1 Hodgkin-Huxley Model

Model of Neuron based on Ion Channel components.

Fraction K+ channels open: $n(t)^4$. Na+ channels: $h(t)m(t)^3$

Gate dynamic system: $a(t) = \frac{d\theta}{dt} = \frac{1}{\tau_a(V)}(a_{\infty}(V) - a)$ (For all gates- Replace a, with n, h, t)

Membrane Potential Dynamics: $c \frac{dV}{dt} = J_{\text{in}} - g_L(V - V_L) - g_{Na}m^3h(V - V_{Na}) - g_Kn^4(V - V_K)$ c: Membrane capacitance $I = C \frac{dV}{dt}$: Net current inside cell J_{in} : Input current from other Neurons g_L : Leak conductance (membrane not impenetrable to ions) g_{Na} : Maximum Na conductance g_K : Maximum K conductance



Action Potential form:

Process: Stimulus breaks threshold causing Na+ channels to open, then close at action potential. Potassium channels open at action potential and close at refractory period.

2.2 Leaky-Integrate-and-Fire Model

Spike shape is constant over time, more important to know when spiked than shape. LIF only considers sub-threshold voltage and when peaked.

Dynamics system: $c \frac{dV}{dt} = J_{\text{in}} - g_L(V - V_L) \Rightarrow \tau_m \frac{dV}{dt} = V_{\text{in}} - (V - V_L)$ For $V < V_{\text{th}}$

This is dimensioned. Dimensionless converts $v_{\text{in}} = \frac{v_{\text{in}} - V_L}{V_{\text{th}} - V_L}$ to become $\tau_m \frac{dv}{dt} = v_{\text{in}} - v$

Then spike occurs when $v = 1$ and we set a refractory period of t_{ref} before starting at 0 again.

Explicit Model: $v(t) = v_{\text{in}}(1 - e^{-t/\tau_m})$

Firing Rate: $\tau_{\text{mean}} = \frac{1}{\tau_m \ln(1 - v_{\text{in}})}$ for $v_{\text{in}} > 1$ **Tuning Curve:** Graph showing how neuron reacts to different input currents.

2.3 Activation Functions (Sigmoidal)

Logistic Curve: $\frac{1}{1 + e^{-z}}$ Arctan: $\arctan(z)$ Hyperbolic Tangent: $\tanh(z)$ Threshold: 0 if $z < 0$; 1 if $z \geq 0$ Rectified Linear Unit: $\max(0, z)$ Softplus: $\log(1 + e^z)$

2.3.1 Multi-Neuron Activation Functions

SoftMax: $\frac{\exp(z_i)}{\sum_j \exp(z_j)}$ Converts elements to probability distribution (sum to 1) and normalizes values **ArgMax:** Largest element remains nonzero, everything else 0

3 Synapses

In real neuron, **Presynaptic action potential** releases **neurotransmitters** across **synaptic cleft** which binds to **postsynaptic** receptors

Equation for current entering **postsynaptic neuron:** $h(t) = \frac{1}{k} \ln \left(\frac{1}{1 + e^{-k t}} \right)$ if $t \geq 0$ (for some $n \in \mathbb{Z}_{\geq 0}$)

k is selected so $\int_0^\infty h(t) dt = 1 \Rightarrow k = \frac{1}{\tau_{\text{ref}}}$

Postsynaptic potential filter: $h(t)$

Spike train: Series of multiple spikes $a(t) = \sum_{p=1}^k \delta(t - t_p)$

Dirac function: Infinite at $t = 0$, 0 everywhere else. Properties:

$$\int_{-\infty}^{\infty} \delta(t) dt = 1, \int_{-\infty}^{\infty} f(t) \delta(t - \tau) dt = f(\tau)$$

Filtered Spike Train: $s(t) = a(t) * h(t)$ For each spike in spike train, run the postsynaptic potential filter on it, then sum each spike. Essentially, convolve the spike train to the postsynaptic potential filter.

Derivation:

$$s(t) = a(t) * h(t) = \int \sum_p \delta(t - t_p) h(t - \tau) d\tau = \sum h(t - t_p)$$

3.1 Neuron Activities

Neurons have multiple connections, with different strengths (weights).

We can represent weights between neuron layers via weight matrix: $W \in \mathbb{R}^{N \times M}$

Compute neuron function as: $\tilde{z} = \tilde{x}W + \tilde{b}$, thus $\tilde{y} = \sigma(\tilde{z})$

3.1.1 Bias Representation

Add a neuron with constant value 1 for each layer, and use its weights as biases.

$$\tilde{x}W + \tilde{b} = (\tilde{x} \cdot 1) \cdot \begin{pmatrix} W \\ \tilde{b} \end{pmatrix}$$

3.1.2 Connections between spiking Neurons

Let $n = 0$ for simplicity for $h(t)$, then it is a solution of $\tau_n \frac{dv}{dt} = -v$

3.2 Full LIF Model

Dynamics are described by: $\begin{cases} \tau_m \frac{dv}{dt} = v_i - v_i, & \text{if not in refractory period} \\ \tau_r \frac{dv}{dt} = -v_i, & \text{if in refractory period} \end{cases}$

If v_i reaches 1, then start refractory period, send spike, reset v_i to 0.

If spike arrives from neuron j , then $s_i \leftarrow s_i + \frac{w_{ji}}{\tau_{\text{ref}}}$

4 Neural Learning

Supervised Learning Desired output is known, and we can minimize error within our model's predictions

Regression Output is continuous-based function of inputs, goal is to get closest to output function.

Classification Outputs in discrete categories.

Unsupervised Learning Output not given, so find efficient structure of input

Reinforcement Learning Feedback given, but uninformative, and depends on lots of variables (ex. Chess AI)

Optimization Given neural model, optimize weights to minimize loss. $\min_{\theta} \mathbb{E}_{(x,t) \in \text{data}} [\mathcal{L}(f(x; \theta), t(x))]$

5 Universal Approximation Theorem

For all continuous functions in the domain of n parameters in $[0, 1]$ domain each, can be approximated as finite sums of sigmoid functions.

Sigmoid Function: Goes to 1 for positive infinity and 0 for negative infinity.

5.1 Process

- By giving infinite weight to $\sigma(wx)$, this approximates a threshold function.
- We create piece function as difference of threshold functions: $P(x; b, \delta) = H(x; b) - H(x; b + \delta)$
- Approximate each section of the function, as $G(x) = \sum_{j=1}^N f(a_j) P(x; b_j, \delta_j)$ each is within ϵ band

6 Loss functions

Single Error: $L(y, t)$ is error between one output y and target t Dataset Error: $\mathbb{E} = \frac{1}{N} \sum_{i=1}^N L(y_i, t_i)$ as average error over entire dataset.

Mean Squared Error: Single Error: $L(y, t) = \frac{1}{2} \|y - t\|_2^2$ Activation Function: ReLU Problems: Linear regression.

Bernoulli Cross Entropy: Single Error: $L(y, t) = -\log(P(y, t)) = -\log(t^y (1 - y)^{1-y}) = -(t \log(y) + (1 - t) \log(1 - y))$ Activation Function: Logistic Problems: Output in range $[0, 1]$;

Categorical Cross Entropy: Single Error: $L(y, t) = -\log(P(x \in C_k, t)) = -\log\left(\prod_{k=1}^K (y^k)^{t^k}\right) = -\sum_{k=1}^K t^k \log(y^k)$ Activation Function: Softmax Problems: Classification with one-hot vectors

7 Gradient Descent

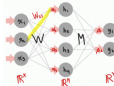
Our goal is to minimize $E(\theta)$ (expected Error), so we define gradient $\nabla_{\theta} E = \left(\frac{\partial E}{\partial \theta_1}, \dots, \frac{\partial E}{\partial \theta_n} \right)$

Error gradient: $\frac{\partial E}{\partial t} = \nabla E(\theta)$, t is a parameter for "time" as we move through parameter space. Euler step: $\theta_{n+1} = \theta_n + k \nabla E(\theta_n)$

8 Error Backpropagation

$\nabla_{z^{(l+1)}} E = \frac{\partial E}{\partial z^{(l+1)}} h^{(l+1)} = \sigma(z^{(l+1)}) = \sigma(W^{(l)} h^{(l)} + b^{(l+1)})$ Basically, h is hidden layer, z is input current, W is weight matrix, b is bias. $\nabla_{z^{(l)}} E = \frac{\partial E}{\partial z^{(l)}} \odot \left[\nabla_{z^{(l+1)}} E \cdot (W^{(l)})^T \right] \odot$ is hadamard product, which does element by element multiplication. We transpose because W_{ij} is connection from i th node in l to j th node in $l+1$

Note that $\tilde{a} = \tilde{x}W$ in this diagram:



$$\left(\frac{\partial W_{ij}^{(l)}}{\partial z^{(l+1)}} = \frac{\partial E}{\partial z^{(l+1)}} \frac{\partial z^{(l+1)}}{\partial W_{ij}^{(l)}} = \frac{\partial E}{\partial z^{(l+1)}} h_i^{(l)} \right)$$

Finally, $\nabla_{z^{(l)}} E = \sigma'(z^{(l)}) \odot \left(\nabla_{z^{(l+1)}} E \cdot (W^{(l)})^T \right) \nabla_{W^{(l)}} E = [h^{(l)}]^T \nabla_{z^{(l+1)}} E$

8.1 Vectorization

We can generalize this process to take a batch of samples by letting x be a matrix of samples instead of just one sample. Then, note $\nabla_{z^{(l)}} E$ is a matrix with same dimension as $z^{(l)}$ as desired. Further, note that $\nabla_{W^{(l)}} E$ is a gradient vector that sums the weight gradient matrix from each sample.

9 Auto-Differentiation

9.1 Expression Graph

Each operation is a square with its variable dependencies. Each variable has a pointer to its creator, which is the operation that created it.

9.1.1 Pseudocode

Function $f = g(x, y, \dots)$
• Create g Op object
• Save references to args x, y, \dots
• Create Var for output f
• Set $g.val$ as $g(x, y, \dots)$
• Set $f.creator$ to the g Op.

9.2 Differentiate

With each object in our graph, store derivative of total expression with respect to itself in member $grad$ Use chain rule with parent operation Op.grad to get current grad. Ex: If y is parent of x , then $x.grad = y.grad \cdot \frac{\partial y}{\partial x}$ Also add wherever multiple branches converge, as is normal in derivatives calculation.

Backward method:

```
class Var:
    def backward(s):
        self.grad += s
        self.creator.backward(s)
```

```
class Op:
    def backward(s):
        for x in self.args:
            x.backward(s * partialDeriv(Op, x))
```

In Var, self.val, self.grad, s must have same shape. In Op, s must be shape of operation output

10 Neural Nets w/w Auto-Diff

10.1 Gradient Descent Pseudocode

- Initialize v, κ
- Make expression graph for E
- Until convergence:
 - Evaluate E at v
 - Zero-gradient
 - Calculate gradients
 - Update $v \leftarrow v - \kappa \cdot v.grad$

10.2 Neural Learning

Optimizing our weights and biases for our loss function. $W \leftarrow -\kappa \nabla_W E$ By making network with AD classes, we leverage backward() to optimize gradient computation.

10.2.1 Pseudocode

Given Dataset (X, T) and network model **net**, with parameters θ and loss function L

- $y = \text{net}(X)$
- loss = $L(y, T)$
- loss.zero_grad()
- loss.backward()
- $\theta \leftarrow \theta - \kappa \cdot \theta.grad$

11 Overfitting

Overfitting: Big discrepancy in accuracy between training and testing or wayyyy too small training error.

Validation: Keep subset of training set to optimize proper hyperparameters before test. **Weight Decay:** Change loss to consider hyperparameters. $\tilde{E}(\tilde{y}, t; \theta) = E(\tilde{y}, t; \theta) + \frac{\lambda}{2} \|\theta\|_2^2$. Updates are $\nabla_{\theta} \tilde{E} = \nabla_{\theta} E + \lambda \theta$ and $\tilde{\theta}_i \leftarrow -\kappa \nabla_{\theta_i} \tilde{E} - \kappa \lambda \theta_i$ **Data Augmentation:** Generate samples by slightly modifying inputs **Dropout:** Drop random nodes to distribute computation.

12 Optimization (Better learning rate)

Stochastic Gradient Descent Computing gradient of loss can be expensive for huge dataset. $E(Y, T) = \frac{1}{B} \sum_D L(y_i, t_i)$

Solution: Take random group γ of B samples. Estimate $E(Y, T) \approx E(\tilde{Y}, \tilde{T}) = \frac{1}{B} \sum_B L(y_{b_i}, t_{b_i})$

Update after each batch B , and then continue gathering batches until all of dataset has been sampled to complete epoch.

12.1 Momentum

Use: Gradient descent oscillates often reaching optimum; grad.desc stops at local optimum (not global). Gradient is a force instead of slope. $\theta_{n+1} = \theta_n + \Delta t v_n$ and $v_{n+1} = (1 - r) v_n + \Delta t A_n$ where A_n represents the gradient vector, so we make $v^{(t)} \leftarrow \beta v^{(t-1)} + \nabla_W E$ and $w^{(t)} \leftarrow w^{(t-1)} - \eta v^{(t)}$ We think of v_n as the "velocity" and A_n (gradient vector) is the force. This gives direction and magnitude. If we move in the same direction all the time, then we gather "momentum", and are able to make bigger steps, by increasing velocity.

13 Deep Neural Networks

Vanishing Gradients When weights and biases are too high, the input currents become too high, and then derivatives are too small, so when you chain them across multiple layers, the gradients reduce severely.

In logistic activation functions, $\sigma'(z) = \sigma(z)(1 - \sigma(z))$ which has maximum value 0.25, so when you chain activation functions, they shrink by at least a factor of 4 each layer, so the closer the layer to the input, the smaller its gradient.

Exploding Gradients If weights are large and biases position the inputs in the high-slope region of the logistic function, then each layer can amplify the gradient, causing exploding gradients.

14 Convolutional Neural Networks

Convolution A mathematical operation that combines two functions by sliding them against each other. Essentially, you create a window, where each element in the window is a function, then you transform the original function over the window to get a frankenstein function.

Continuous Convolution $(f * g)(x) = \int_{-\infty}^{\infty} f(s) \cdot g(x - s) ds$

Discrete Convolution $(f * g)_{nm} = \sum_{n=0}^{N-1} \sum_{m=0}^{N-1} f_{n-m} \cdot g_{m-n}$

Convolution in 2D: $(f * g)_{m,n} = \sum_{i,j} f_{ij} \cdot g_{m-i, n-j} (f \oplus g)_{m,n} = \sum_{i,j} f_{ij} \cdot g_{m-i, n-j}$

Each kernel is convolved against the layer before it, creating an activation map, which creates a tensor for the next layer. We can have multiple kernels per input layer.

Stride The amount the kernel is shifted against the input layer.

Padding Padding border of input layer with 0s to create a layer that is identical size as layer.

Note: The bias is attributed to the kernel as a hole, rather than the neurons in the kernel.

Note: 1x1 convolution layers are useful in multi-layered input. (Ex 1x1x64)

15 Batch Normalization

Some input features will have dramatically different magnitudes than others. But we'd need to accommodate the smallest magnitude feature to make meaningful steps, which slows down the larger magnitude rate.

The goal is to rescale all hidden layer outputs into normalized values.

Use the standard formulas for means and variance for each hidden layer output, then we rescale the inputs into the next layer as:

$\hat{h}_i^{(d)} = \frac{h_i^{(d)} - \mu_i}{\sigma_i}$ or $\hat{h}_i^{(d)} = \frac{h_i^{(d)} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$ for small $\epsilon > 0$ for if variance is 0, to prevent instabilities.

Then we rescale output with learnable parameters γ_i, β_i

Following process: [Hidden layer i] -> Normalization -> Rescaling -> [Hidden Layer i+1]

$$\tilde{y}_i^{(d)} = \gamma_i \hat{h}_i^{(d)} + \beta_i$$

Batch normalization affects convergence rate for learning quickly.

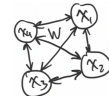
The reason is unknown but there are hypotheses:

- Mitigates vanishing/exploding gradients
- Guards against internal covariate shift (shallow layers (near output) learn quicker than deep layer, so whenever deep layers learn, they invalidate the outputs of the shallow layer, who has to learn again. By normalizing, the inputs remain relatively stable, so shallower layers aren't as "affected" by differences from deep layer changes.

16 Hopfield Networks

Content-Addressable Memory System that can take part of a pattern and fill in the most likely matches from memory.

Hopfield Network Given a network of N neurons, each connected to all others, and given a set of M possible targets, we want the network to converge to the nearest of this set.



Each x_i is assigned: -1 if $\tilde{x}W + b_i < 0$ and 1 if $\tilde{x}W + b_i \geq 0$

Since the graph has cycles, backpropagation won't work. Instead we need to utilize the Hopfield Energy:

$$E = -\frac{1}{2} \sum_i \sum_j x_i W_{ij} x_j - \sum_i b_i x_i = -\frac{1}{2} \tilde{x} W \tilde{x}^T - \tilde{b} \tilde{x}^T$$

and diagonals of W are 0.

To minimize energy, we use gradient descent: $\frac{\partial E}{\partial x_i} = -\sum_{j \neq i} x_j W_{ij} - b_j$ or $\nabla_x E = -\tilde{x}W - \tilde{b} \Rightarrow \tau_x \frac{dx}{dt} = \tilde{x}W + \tilde{b}$

If $i \neq j$, $\frac{\partial E}{\partial W_{ij}} = -x_i x_j$

If $i = j$, $\frac{\partial E}{\partial W_{ii}} = -x_i^2 = -1$

Therefore gradient vector is $\nabla_W E = -\tilde{x} \tilde{x}^T + I_{N \times N}$

Over M targets we have: $\nabla_W E = -\frac{1}{M} \sum_{s=1}^M (\tilde{x}^{(s)})^T \tilde{x}^{(s)} + I = -\frac{1}{M} X^T X + I$

Thus $\tau_w \frac{dW}{dt} = \frac{1}{M} X^T X - I$

17 Recurrent Neural Networks

Hidden layer can encode input sequence, allowing sequential data to be outputted.

Backprop Through Time Unroll network through time, to create feedforward network into a DAG, and evaluate the targets in sequence similar to how you would batch inputs.

Note: $h' = \sigma(x'U + h'^{-1}W + b)$ $y = \sigma(h'V + c)$

The error function is: $E(y_1, \dots, y_t, t_1, \dots, t_t) = \sum_{i=1}^T L(y_i, t_i) \alpha_i$

The goal is to minimize the following: $\min_{\theta} \mathbb{E}[E(y_1, \dots, y_t, t_1, \dots, t_t)]_{y, t}$

Following gradients: $\nabla_{x^k} E = \nabla_{x^k} \left(\sum_{i=1}^t L(y^i, t^i) \right) = \sum_{i=1}^t \nabla_{x^k} L(y^i, t^i) = \nabla_{x^k} L(y^k, t^k) = \nabla_{y^k} L(y^k, t^k) \odot \sigma'(z^k)$

In V Der: $\sum_{i=1}^T h' \nabla_{x^k} L(y^i, t^i)$

Unfold to derive hidden layer gradients, define $E^k = \sum_{i=k}^T L(y^i, t^i)$

Since h_k only affects h_{m+} where $m \geq k$ then $\nabla_{h^k} E = \nabla_{h^k} E^k$

$$\nabla_{h^*} E = \nabla_{x^*} E^* \frac{\partial x^*}{\partial h^*} = \nabla_{x^*} E^* V^T$$

Then, we can compute $\nabla_{h^*} E$ recursively.

$$\nabla_{h^*} E = (\nabla_{h^{t+1}} E^{t+1} \odot \sigma'((s^{t+1}))^T) W^T + (\nabla_{y^t} L(y^t, t^t) \odot \sigma'(z^t))^T V^T$$

Once you have $\nabla_{h^*} E$ you can compute gradient with respect to deeper weights and biases recursively.

18 Gated Recurrent Units

18.1 Problem

Vanilla RNNs struggle to capture long-range dependencies, because of the vanishing gradient problem, as gradients multiply over long ranges, meaning the relative effect of distant tokens to the current token are minimal.

Consider a simple example:

$$h^n = w h^{n-1} + x^n, \text{ if we unroll this equation, we get } h^n = w^{n-1} h_1 + w^{n-1} x^1 + w^{n-2} x^2 + \dots$$

If $|w| < 1$ then w^n shrinks exponentially as n grows, meaning earlier information has exponentially less influence on h^n over time. If $|w| > 1$ then h^n magnitude grows exponentially, making training unstable.

18.2 Gated Recurrent Units

Use gating mechanisms to control which information to retain.

Candidate Hidden State: Candidate Hidden state \tilde{h}^n is computed as: $\tilde{h}^n = \tanh(\tilde{h}^{n-1} W + \tilde{x}^n U + \tilde{b})$. W : hidden-to-hidden weights, U : input-to-hidden weights, \tilde{b} : bias, \tanh : ensures output $\in (-1, 1)$

Gate Mechanism: $g^n = \sigma(\tilde{h}^{n-1} W_g + \tilde{x}^n U_g + \tilde{b}_g)$ The sigmoid ensures $g^n \in (0, 1)$ meaning it's a soft switch.

Final Hidden State Update: $\tilde{h}^n = g^n \odot (\tilde{h}^n) + (1 - g^n) \odot \tilde{h}^{n-1}$

g^n controls how much of new candidate is retained, and $(1 - g^n)$ shows how much of previous state is preserved.

When $g = 0$ we don't update hidden state (meaning the word added no new context), if $g = 1$, then we do update, proving it's important information.

18.3 GRU Network

The GRU extends from RNN by adding two gating mechanisms, update gate and reset gate (same format as gate mechanism above).

Hidden state updates as follows: $\tilde{h}^n = \tanh((\tilde{h}^{n-1} \odot r^n) W + \tilde{x}^n U + \tilde{b})$

$$\tilde{h}^n = g^n \odot \tilde{h}^n + (1 - g^n) \odot \tilde{h}^{n-1}$$

Reset gate tries to "forget" previous hidden state information for new candidate. When r^n is close to 0, past is discarded, making it more focused on recent information, when close to 1, then more of hidden state is retained.

18.4 minGRU

For a sequence of length N , we'd need to compute each component sequentially, nullifying parallelization advantages from GPUs.

minGRU simplifies to: $\tilde{g}^n = \sigma(\tilde{x}^n U_g + \tilde{b}_g)$ $\tilde{h}^n = \tilde{x}^n U + \tilde{b}$ $\tilde{h}^n = \tilde{g}^n \odot \tilde{h}^n + (1 - \tilde{g}^n) \odot \tilde{h}^{n-1}$

minGRU allows \tilde{g}^n and \tilde{h}^n to be parallelizable. Hidden states can use parallel scan to compute values in $O(\log N)$ instead of $O(N)$ sequentially.

19 Autoencoders

Neural network that learns to encode/decode a set of inputs automatically into a smaller latent space.

The model is trained using loss function minimizing $L(x', x)$ where x' is reconstructed input and x is original input.

To simplify latent space encoding, we can tie the weights of the encoder and decoder, so that the encoder weights are W and the decoder weights are W^T

20 Vector Embeddings

Given a set of all inputs in our dataset, we can call this our vocabulary. We can order and index vocabulary into a set of one-hot vectors of the size of our vocabulary.

To reflect semantic relationships between words, we can determine word pairs for words that often occur near each other within some fixed window size d .

With this set of word pairs, we can try to predict these word co-occurrences via a 3-layer neural network.

Input is a one-hot word vector, output is probability of each word's co-occurrence.

The hidden layer is of a lower dimensional space, requiring similar words to take similar representations, this is called "embedding".

word2vec: Popular embedding strategy that 1. treats common phrases as words (New York) and 2. randomly ignores very common words (the) and 3. backpropagates only some of the negative cases in order to speed up the training. (this means we train our words by sampling instances where contexts are DIFFERENT to increase reflexive semantic relationships)

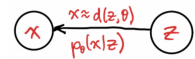
Cosine Similarity: Cosine angle is used to measure "distance" between two vectors in embedding space. (Since they're all 1)

Vector Arithmetic: By applying vector arithmetic, you can do semantic calculations. (Ex. King - Man + woman =)

21 Variational Autoencoders

Goal: Create autoencoder that generates reasonable samples not in training set. We encode a latent probability distribution, instead of encoding.

x is input, z is latent space, so encode $p(z)$. Decoding is $d(z, \theta)$. θ are decoder weights. Want to maximize $p(x) = \int p_{\theta(z|z)} p(z) dz$ - the decoding probability.



Assume $p(x|z)$ is Gaussian (for simplicity), then NLL is $-\ln p_{\theta(x|z)} = -\frac{1}{2\sigma^2} \|X - d(z, \theta)\|^2 + C$ which is what we want to maximize. Our goal is $\min_{\theta} \mathbb{E}_{p(z)} [\|X - d(z, \theta)\|^2]$, and $\mathbb{E}_{p(z)} [p_{\theta(x|z)}] = \int p_{\theta(x|z)} p(z) dz$. But we don't know how to sample $p(z)$, so we assume a distribution $q(z)$ to approximate value.

$$p(x) = \mathbb{E}_{z \sim q} [p(x|z)] \\ = \sum_{z \sim q} p(x|z) p(z) \\ = \sum_{z \sim q} p(x|z) \frac{p(z)}{q(z)} q(z) \\ = \mathbb{E}_{z \sim q} \left[p(x|z) \frac{p(z)}{q(z)} \right]$$

Then NLL is: $-\ln p(x) \leq -\mathbb{E}_{q(z)} [\ln p(x|z)] + \ln \frac{p(z)}{q(z)}$, by Jensen's inequality (if f is convex, $E(f(x)) \geq f(E(x))$. Then we simplify to $-\ln p(x) \leq KL(q(z) || p(z)) - \mathbb{E}_{q(z)} [\ln p(x|z)]$, where $KL(q(z)||p(z)) = -\mathbb{E}_{z \sim q} [q(z) \ln \frac{p(z)}{q(z)}]$.

We will then choose $p(z) \sim \mathcal{N}(0, I)$ and want to $\min_{\theta} KL(q(z)||\mathcal{N}(0, I))$. Our encoder will be designed to have outputs μ, σ which aren't actual means, or standard deviations, but are just parameters for the distribution.

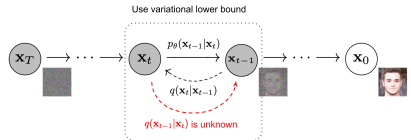
These gaussians result in $KL(\mathcal{N}(\mu, \sigma)||\mathcal{N}(0, I)) = \frac{1}{2}(\sigma^2 + \mu^2 - \ln(\sigma^2) - 1)$

The other part of objective, $\mathbb{E}_{q(z)} [\ln p(x|z)]$ can be written as $\mathbb{E}_{q(z)} [\ln p(x|z)]$ where $x = d(z, \theta)$ and $z = \mu(x, \theta) + \varepsilon \odot \sigma(x, \theta)$. This is a reparameterization trick done so we can backpropagate on x, θ , by making the ε a separate noise vector.

Procedure: Encode x by finding $\mu(x, \theta)$ and $\sigma(x, \theta)$. Sample $z = \mu + \varepsilon \odot \sigma$, $\varepsilon \sim \mathcal{N}(0, I)$. Calc KL Loss. Decode \hat{x} using decoder. Calc reconstruction loss. MSE for Gaussian $p(x|\hat{x})$; BCE for Bernoulli. Gradient descent on θ on combined loss, with β parameters on KL divergence, to adjust importance.

22 Diffusion Models

Generative model that uses noise as seeds for generating new images, by training on a model that adds noise to samples so it can then reverse the process.



Forward process is defined as $q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$. Final state is pure noise ($x_T \sim \mathcal{N}(0, I)$).

The variance schedule β_1, \dots, β_T is a hyperparameters controlling onise addition. Usually starts small, gets larger.

If we unwrap x_t recursively, we get $x_t = \sqrt{\alpha_t} x_0 + \sum_{i=1}^t \varepsilon_i$ where $\alpha_t = 1 - \beta_1, \dots, \beta_t$ and $\varepsilon_i \sim \mathcal{N}(0, I)$. Since sum of gaussians is gaussian, then $x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \varepsilon$.

To get x_0 from x_t , we use neural net $\varepsilon_{\theta}(x_t, t)$ to estimate noise ε

To do reverse diffusion (generate input from noise), we need to minimize KL-divergence loss of $D_{KL}(q(x_{t-1}|x_t, x_0), p_{\theta(x_{t-1}|x_t)})$ for $t = 2, \dots, T$. $q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}|\mu_q, \Sigma_q)$, $\mu_q = \frac{1}{\sqrt{\alpha_t}} (x_t - \sqrt{1 - \alpha_t} \varepsilon_t)$

Further, parameterized probability density p_{θ} is assumed to be a Gaussian distribution, with the same format, but $\varepsilon_{\theta}(x_{t-1}, t)$ instead of ε_t

Then the KL-Divergence gets simplified to $\mathbb{E}_{x_0, \varepsilon} [\|\varepsilon - \varepsilon_{\theta}(x_t, t)\|_2^2]$ where $\lambda_t = \frac{\beta_t^2}{2\sigma^2(\alpha_{t-1} - \alpha_t)}$. By substituting x_t we get

$$\mathbb{E}_{x_0, \varepsilon} \left[\left\| \varepsilon - \varepsilon_{\theta} \left(\frac{1}{\sqrt{\alpha_t}} (x_t - \sqrt{1 - \alpha_t} \varepsilon) \right) \right\|_2^2 \right]$$

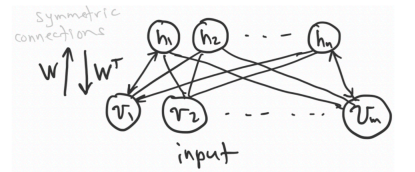
Training Algorithm: Repeat until converge:
1. Take x_0 from dataset
2. $t \sim \text{Uniform}(1, \dots, T)$
3. $\varepsilon \sim \mathcal{N}(0, I)$
4. Do gradient descent on simplified Loss function

Sampling Algorithm:

1. Set $x_T \sim \mathcal{N}(0, I)$
2. For $t = T, \dots, 1$
1. $z \sim \mathcal{N}(0, I)$ if $t > 1$ else $z = 0$ We add noise to make stochastic sample
2. $x_{t-1} = \frac{1}{\sqrt{\alpha_t}} (x_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \varepsilon_{\theta}(x_t, t)) + \sigma_t z$
3. Return x_0

We typically use a U-net to estimate ε_{θ} because upsampling/downsampling works well with noisy data where compressed and original spatial data are both important.

23 Restricted Boltzmann Machines



Network consists of:

• Hidden layer and visible layers (nodes only binary) and connections between layers are symmetric by weight matrix W

23.1 RBM Energy

• RBM's energy is $E(v, h) = -\sum_{i=1}^n \sum_{j=1}^m v_i W_{ij} h_j - \sum_{i=1}^n b_i v_i - \sum_{j=1}^m c_j h_j$ or rewritten $E(v, h) = -v^T W h^T - b^T v - c^T h$ where $W \in \mathbb{R}^{m \times n}$. Discordance cost: $-v^T W h^T$. Operating cost: $-b^T v - c^T h$

RBM vs Hopfield: Both Find weights that minimize energy, and running network converges to low-energy states, but Hopfield sets nodes to target pattern but RBM only visible nodes.

23.2 Energy Gap: (Like gradient but for binary nodes)

$\delta E(v_i) = E(v_i = 1) - E(v_i = 0)$, for $i = 1, \dots, m$: $\delta E(v_i) = -\sum_{j=1}^n W_{ij} h_j - b_i$ For $j = 1, \dots, n$: $\delta E(h_j) = -\sum_{i=1}^m v_i W_{ij} - c_j$. Then $\delta E(v_i) > 0$ then turn v_i off, meaning $E(v_i = 1) > E(v_i = 0)$

One issue is that these binary networks can get stuck in sub-optimal state, so we can use stochastic neurons as $P(h_j = 1 | v) = \sigma(\sum_i v_i W_{ij} + c_j)$ (similar for visible nodes) where logistic is defined as $\sigma(z) = \frac{1}{1 + e^{-z}}$, which is temperature dependent.

Sampling Algorithm:

• Compute $p_i = P(v_i = 1 | h)$
• For $i = 1, \dots, m$:
• Draw $r \sim \text{Uniform}(0, 1)$
• If $p_i > r$ set $v_i = 1$, else set $v_i = 0$

23.3 Network Dynamics

If we run network freely, then network states will all be visited with probability $q(v, h) = \frac{1}{Z} e^{-E(v, h)}$ where partition function Z is defined as $Z = \sum_{v, h} e^{-E(v, h)}$. Since lower-energy states are visited more

frequently, then $E(v^{(1)}, h^{(1)}) < E(v^{(2)}, h^{(2)}) \Rightarrow q(v^{(1)}, h^{(1)}) > q(v^{(2)}, h^{(2)})$

23.4 Generation

Suppose inputs $v \sim p(v)$, we want RBM to act like generative model q_{θ} such that $\max_{\theta} \mathbb{E}_{v \sim p} [\ln q_{\theta}(v)]$

Let loss be $L = -\ln q_{\theta}(v)$ for given V , or $L = -\ln(\frac{1}{2} \sum_h e^{-E_{\theta}(V, h)}) = -\ln(\sum_h e^{-E_{\theta}(V, h)}) + \ln(\sum_v \sum_h e^{-E_{\theta}(V, h)})$, which can be decomposed into $L = L_1 + L_2$

23.5 Gradient of L_1

• To optimize parameter, we need to compute gradient of loss function:

$$\begin{aligned} \nabla_{\theta} L_1 &= -\nabla_{\theta} \sum_h \frac{e^{-E_{\theta}(V, h)}}{\sum_h e^{-E_{\theta}(V, h)}} \\ &= \frac{\sum_h e^{-E_{\theta}(V, h)} \nabla_{\theta} E_{\theta}(V, h)}{\sum_h e^{-E_{\theta}(V, h)}} \\ &= \sum_h \frac{e^{-E_{\theta}(V, h)}}{\sum_h e^{-E_{\theta}(V, h)}} \nabla_{\theta} E_{\theta}(V, h) \end{aligned}$$

Then, since the fraction above is equivalent to $q_{\theta}(h | v)$, we write $\nabla_{\theta} L_1 = \sum_h q_{\theta}(h | v) \nabla_{\theta} E_{\theta}(V, h) = \mathbb{E}_{q(h | v)} [\nabla_{\theta} E_{\theta}(V, h)]$

23.6 Gradient of L_2

$$\begin{aligned} \nabla_{\theta} L_2 &= -\frac{\sum_{v, h} e^{-E_{\theta}} \nabla_{\theta} E_{\theta}}{\sum_{v, h} e^{-E_{\theta}}} \\ &= -\sum_{v, h} q_{\theta}(v, h) \nabla_{\theta} E_{\theta}(v, h) \\ &= -\mathbb{E}_{q(v, h)} [\nabla_{\theta} E_{\theta}(v, h)] \end{aligned}$$

Thus $\nabla_{\theta} L = \nabla_{\theta} L_1 + \nabla_{\theta} L_2 = \mathbb{E}_{q(h | v)} [\nabla_{\theta} E_{\theta}] - \mathbb{E}_{q(v, h)} [\nabla_{\theta} E_{\theta}]$

23.7 Gradient for W_{ij}

$\nabla_{W_{ij}} E(V, h) = -V_i h_j$ and $\nabla_{W_{ij}} E(v, h) = -v_i h_j$ and thus $\nabla_{W_{ij}} L = -\mathbb{E}_{q(h | v)} [V_i h_j] + \mathbb{E}_{q(v, h)} [v_i h_j]$. First term represents expected value under posterior distribution, and second term under joint distribution

23.8 Contrastive Divergence for Training

Step 1: Clamp visible states to V and calculate hidden probabilities $q(h_j | v) = \sigma(V W_{ij} + c_j)$ and $\nabla_{W_{ij}} L_1 = -V^T \sigma(V W + c)$ which results in a rank-1 outer product in $\mathbb{R}^{m \times n}$

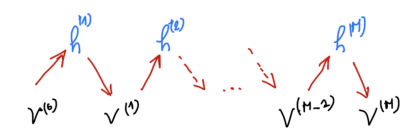
Step 2: Compute expectation using Gibbs sampling. $\langle v_i h_j \rangle_{q(v, h)} = \sum_v \sum_h q(v, h) v_i h_j$. Practically, single network state is used to approximate expectation. Gibbs sampling is used to run the network freely and compute average $v_i h_j$, which has $\nabla_{W_{ij}} L_2 = v^T \sigma(v W + c)$.

Procedure for training: Given V calculate h , then given new h calculate v , and from new v calculate second h .

Weight update rule is: $W \rightarrow W - \eta (\nabla_{W_{ij}} L_1 + \nabla_{W_{ij}} L_2) = W + \eta V^T \sigma(V W + c) - \eta v^T \sigma(v W + c)$ where η is learning rate. First term relates to positive phase (step 1) and second term to negative step (step 2).

23.9 Sampling

After training, we can use it as generative model to generate new M data points $V^{(1)}, V^{(2)}, \dots, V^{(M)}$ by performing Gibbs sampling from the conditional probabilities $P(v|h)$ as illustrated below:



24 Adversarial Attacks

Consider classifier f that produces probability vectors.

Classification errors are defined as $R(f) \triangleq \mathbb{E}_{(x, t) \sim D} [\text{card}\{\arg \max_i y_i \neq t \mid y = f(x)\}]$

ε -Ball: $B(x, \varepsilon) = \{x' \in X \mid \|x - x'\| \leq \varepsilon\}$

Adversarial Attack: Find $x' \in B(x, \varepsilon)$ such that $\arg \max_{i(y_{x'})} \neq t$ for $y = f(x')$

Untargeted Gradient-Based Whitebox Attack: $x' = x + k \nabla_x E(f(x; \theta), t(x))$ Targeted: $x' = x - k \nabla_x E(f(x; \theta), l)$

24.1 Fast Gradient Sign Method

Adjust each pixel by ε , so $\delta x = \varepsilon \text{sign}(\nabla_x E)$ (This is because it's non-differentiable, so model training can't adjust for it)

Minimal Perturbation: Smallest $\|\delta x\|$ causing misclassification: $\min_{\|\delta x\|} [\arg \max_{i(y_{(x+\delta x)})} \neq t(x)]$

25 Adversarial Defense

During training, add adversarial samples to dataset with proper classification.

25.1 TRADES

Model: $f: X \rightarrow \mathbb{R}$ Dataset: (X, T) ; $X \in \mathbb{X}$, $T \in \{-1, 1\}$ $\text{sign}(f(X))$ indicates class. Classification is correct if $f(X)T > 0$

Robust Loss: $\mathcal{R}_{\text{rob}}(f) = \mathbb{E}_{X, T} [\text{card}\{X' \in B(X, \varepsilon) \mid f(X')T \leq 0\}]$ (Even if proper classification, if ε -ball can be fooled, it counts for loss)

Differentiable Training Loss: $\mathcal{R}_{\text{training}} = \min_{\eta} \mathbb{E}_{(X, T)} [g(f(X, T))]T$ where g is smooth function.

Robust model optimizes over $\min_{\eta} \mathbb{E}_{(X, T)} [g(f(X, T)) + \max_{X' \in B(X, \varepsilon)} g(f(X')f(X'))]$ Term 1 ensures proper classification. Term 2 adds penalty for attacks.

Procedure:

- For each gradient update:
 - Run several gradient ascents to find X'
 - Evaluate joint loss $g(f(X)T) + \beta g(f(X')f(X'))$
 - Update weights off loss

26 Population Coding

The ability to encode data in a shape we want, and break apart a Black-Box NN to separate interpretable NN features that we can place in sequence.

Given activities of a neural network we can reconstruct input based, off of the specific activation values. Since decoding is linear function, this is regression, with MSE loss, so we can optimize for $\frac{1}{2n} \min_D \|H D - X\|_2^2$, where $H^{(i)}$ is a row corresponding to activations of hidden layer.

The optimal linear decoding can be solved to be $d^* = (H^T H)^{-1} H^T X$.

This can be problematic if $H^T H$ is poorly conditioned (almost singular), so instead we add noise to H , and get:

$$\begin{aligned} \|(H + \varepsilon) D - T\|_2^2 &= \|(H D - T) + \varepsilon D\|_2^2 \\ &= (H D - T)^T (H D - T) + 2(H D - T)^T (\varepsilon D) + D^T \varepsilon^T \varepsilon D \end{aligned}$$

Since middle term is usually zero, since ε independent of $H D - T$, then if $\varepsilon^T \varepsilon \approx \sigma^2 I$, then it is finally $\|H D - T\|_2^2 + \sigma^2 \|D\|_2^2$

We can also expand population coding to reconstruct vectors. We solve the matrix $D^* = \arg \min_D (\text{norm } H D - T)_F^2$ (frobenius norm) and T is a matrix where each row corresponds to a horizontal sample vector input.

27 Transformations

We use population coding to pass data between populations of neurons. To do this, we need some way to pass hidden activations between the other.

Naive: Decode activations to data space, then re-encode for second population. Alternative is to bypass data space by multiplying decoder weights by encoder weights directly. $W = D_{xy} E_B \in \mathbb{R}^{N \times M}$. This is rank-1 matrix since $D_{xy} \in \mathbb{R}^{N \times 1}$ and $E_B \in \mathbb{R}^{1 \times M}$.

It's better to use separate decoder-encoder because it's computationally efficient, since it's low-rank. Calculating AD takes $O(N)$ time. $(AD)E$ takes $O(M)$ time, so total time is $O(N + M)$, whereas tied weights would take $O(NM)$

28 Dynamics

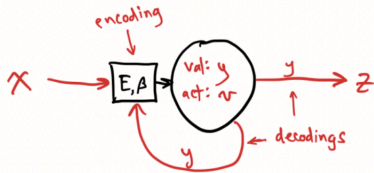
To build dynamic, recurrent networks via population coding methods, we'll leverage a dynamic model of LIF neurons based on current and activity. Our population coding methods will operate on activity, so we can modify the input current to establish forces on the activity.

$$\tau_s \frac{ds}{dt} = -s + C \text{ Current}$$

$$\tau_m \frac{dv}{dt} = -v + \sigma(s) \text{ Activity}$$

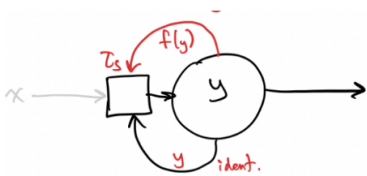
Equilibrium values are when differential values equal 0

If $\tau_m \ll \tau_s$ then Activity function reaches equilibrium value, while current is still in dynamic state. Same for if $\tau_s \ll \tau_m$



If we're integrating the input, then $\tau_s \frac{ds}{dt} = -s + \sigma(s)W + \beta + C$ a recurrent network. Then, $\frac{ds}{dt} = \frac{-s + W + \beta}{\tau_s} + \tilde{\tau}$ where $\tilde{C} = xE$ because it is the re-integrated input, and isn't dependent on time constant.

Since we have multiple different forms of the dynamic system, we can generalize to $\frac{ds}{dt} = f(y)$

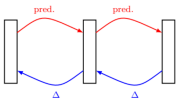


29 Biological Backprop

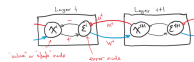
In neurons, updates can only use local information, but backpropagation updates to weights involves various layers of information.

29.1 Predictive Coding (PC)

- Predictions are sent down the hierarchy
- Errors are sent up the hierarchy



In a PC-Node, there are two parts: Value/State Node and an error node.



$\mu^i = \sigma((x^{i+1})M^i + \beta^i)$, this is our prediction. For simplicity, assume $W^i = (M^i)^T$

29.1.1 Error Node

$\tau \frac{d\epsilon^i}{dt} = x^i - \mu^i - v_{\text{max}}^i \epsilon^i$, which at equilibrium reaches $\epsilon^i = \frac{1}{v^i}(x^i - \mu^i)$

29.1.2 Generative Network

Given dataset x, y and $\theta = \{M^i, W^i\}_{i=1, \dots, L-1}$, the goal is to $\max_p \mathbb{E}_{(x,y)} [\log p(x | y)]$ $p(x | y) = p(x^1 | x^2) p(x^2 | x^3) \dots p(x^{L-1} | y) = p(x^1 | \mu^1) \dots p(x^{L-1} | \mu^{L-1})$

If we assume $x^i \sim \mathcal{N}(\mu^i, v^i)$, then $p(x^i | \mu^i) \propto \exp\left(-\frac{\|x^i - \mu^i\|^2}{2(v^i)^2}\right)$ Then, $-\log p(x^i | \mu^i) = \frac{1}{2} \left\| \frac{x^i - \mu^i}{v^i} \right\|^2 + C$, so as a result. $-\log p(x | y) = \frac{1}{2} \sum_{i=1}^{L-1} \|\epsilon^i\|^2$

29.1.3 Hopfield Energy

$$F = \frac{1}{2} \sum_{i=1}^{L-1} \|\epsilon^i\|^2 + \tau \frac{dF}{dt} = -\nabla_{\epsilon^i} F$$

x^i shows up in two terms in F : $\epsilon^i = x^i - \mu^i(x^{i+1}) = x^i - (\sigma(x^{i+1})M^i + \beta^i)$ and $\epsilon^{i-1} = x^{i-1} - \mu^{i-1}(x^i) = x^{i-1} - (\sigma(x^i)W^{i-1} + \beta^{i-1})$

Therefore $\nabla_{x^i} F = \epsilon^i \frac{\partial \epsilon^i}{\partial x^i} + \epsilon^{i-1} \frac{\partial \epsilon^{i-1}}{\partial x^i} = \epsilon^i - \sigma'(x^i) \odot (\epsilon^{i-1}(W^{i-1})^T)$

29.1.4 Dynamics

$\tau \frac{dx^i}{dt} = \sigma'(x^i) \odot (\epsilon^{i-1}M^i) - \epsilon^i$ and $\tau \frac{d\mu^i}{dt} = x^i - \mu^i - v^i \epsilon^i$. Then learning M^i involves $\nabla_{M^i} F = -(\sigma(x^{i+1}))^T \epsilon^i$ with systems $\frac{dM^i}{dt} = \sigma(x^{i+1})^T \epsilon^i$ and $\tau \frac{dW^i}{dt} = (\epsilon^i)^T \cdot \sigma(x^{i+1})$.

At equilibrium we get $\epsilon^i = \sigma'(x^i) \odot [\epsilon^{i-1}(W^{i-1})^T]$ where $\frac{\partial F}{\partial \mu^i} = -\epsilon^i$

Training: Clamp $x^1 = X$ and $x^L = Y$ and run to equilibrium. x^i, ϵ^i reach equilibrium quickly. Then use the two systems based on dM^i and dW^i to update M^i and W^i

Generating: Clamp $x^L = Y$ and run to equilibrium and x^1 is a generated sample

Inference: Clamp $x^1 = X$ run to equilibrium and $\arg \max_y (x_y^L)$ is the class

This work overcomes the local learning condition because running to equilibrium allows information to spread through the net.

30 Generative Adversarial Networks

Two networks: Generative Network (produce lookalike data) and Discriminative Network (identify fake data)

$D(x; \theta)$ - Probability x is real. $G(z; \varphi)$ - Create input sample from random noise z drawn from p_z distribution.

Loss function: $E(\theta, \varphi) = \mathbb{E}_{x \sim \mathcal{R}, x'} [\mathcal{L}(D(x; \theta), t)] + \mathbb{E}_{z \sim p_z} [\mathcal{L}(D(G(z; \varphi), \varphi), 1)]$ Term 1: Minimize θ to make discriminator better Term 2: Minimize φ to make generator better at producing fake inputs.

Train discriminator: $\min_{\theta} \mathbb{E}_{\mathcal{R}, \mathcal{F}} [\mathcal{L}(y, t)]$ R are real inputs, F are fake. Update rule: $\theta \leftarrow \theta - \kappa \nabla_{\theta} \mathcal{L}(y, t)$ Train generator: $\min_{\varphi} \mathbb{E}_{\mathcal{F}} [\mathcal{L}(y, 1)]$ Update rule: $\varphi \leftarrow \varphi - \kappa \nabla_{\varphi} \mathcal{L}(y, 1)$ (We use 1 to simulate a targeted adversarial attack with target 1) Gradients propagate through D down to G . Note that y is the discriminator being run on generated outputs.

31 Transformers

First inputs of sequential data are tokenized (converted to vector embedding)

Input $X = (\vec{x}_1, \dots, \vec{x}_n)$ Queries (Input Value): $\vec{q}_i = \vec{x}_i W^{(Q)} \in \mathbb{R}^{n \times l}$ Keys (Tag for data): $\vec{k}_i = \vec{x}_i W^{(K)} \in \mathbb{R}^{n \times l}$ Values (Desired data): $\vec{v}_i = \vec{x}_i W^{(V)} \in \mathbb{R}^{n \times l}$

Let $W^{(Q)}, W^{(K)}, W^{(V)} \in \mathbb{R}^{d \times l}$, l is hyperparameter and $\vec{x}_i \in \mathbb{R}^d$

We then vectorize components to get $Q = XW^{(Q)}, K = XW^{(K)}, V = XW^{(V)}$ to compute self-attention.

31.1 Calculating Attention Scores

$S_{ij} = \vec{q}_i \cdot \vec{k}_j$ which scores query i against key j $S = QK^T$

We use **softmax** (on each row, set of keys) to get attention scores. $A = \text{Softmax}\left(\frac{S}{\sqrt{d_k}}\right) \in \mathbb{R}^{n \times n}$ where rows sum to 1 because of softmax.

Then $H = A \cdot V$, $H \in \mathbb{R}^{n \times l}$, where H is an attention head and $\vec{H}_i = \sum_{j=1}^n A_{ij} \vec{v}_j$

To add positional definition, we use positional encodings:

$$PE(i)_{2j} = \sin\left(\frac{i}{10000^{2j/d_k}}\right), PE(i)_{2j+1} = \cos\left(\frac{i}{10000^{2j/d_k}}\right)$$

$PE(i)$ is positional encoding vector for position i , frequency changes with dimension j and 10000 is a scaling constant to allow covering a large max sequence length.

Then $\vec{x}_i' = \vec{x}_i + PE(i)$

Multi-Head Attention: Model has separate attention mechanisms learning separate features. Output is concatenated at the end before being linearized for output.

31.2 Transformer Architecture

Use encoder as layers of multi-head self-attention and fully-connected FF nets and connect to decoder with masked multi-head attention, encoder-decoder attention and fully-connected layers.

