

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

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

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

\tilde{g}^n controls how much of new candidate is retained, and $(1 - \tilde{g}^n)$ 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: $\tilde{h}^n = \tanh((\tilde{h}^{n-1} \odot \tilde{r}^n)W + \tilde{x}^nU + \tilde{b})$

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

Reset gate tries to "forget" previous hidden state information for new candidate. When \tilde{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: $\tilde{g}^n = \sigma(\tilde{x}^nU_g + \tilde{b}_g)$ $\tilde{h}^n = \tilde{x}^nU + \tilde{b}$ $\tilde{h}^n = \tilde{g}^n \odot \tilde{h}^n + (1 - \tilde{g}^n) \odot \tilde{h}^{n-1}$

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

Vocabulary inputs can be represented as one-hot vectors. Semantic relationships between words are captured by identifying frequently co-occurring word pairs within a fixed window. A 3-layer neural network predicts these co-occurrences from a one-hot input vector, with the lower-dimensional hidden layer forming the "embedding" where similar words have similar representations. Output is probability of each word's co-occurrence.

word2vec is a common embedding strategy that treats common phrases as single words, randomly ignores very frequent words, and backpropagates only some negative cases to speed learning. This trains words by sampling instances where contexts differ to increase reflexive semantic relationships.

Cosine Similarity uses the cosine angle to measure the distance between vectors in embedding space. Vector Arithmetic allows for semantic calculations (e.g., King - Man + Woman = Queen).

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 $p(x) = \int p_{\theta(x|z)}p(z)dz$ —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 $\min_{\theta} \mathbb{E}_{z \sim p(z)} [\|X - d(z, \theta)\|^2]$, 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.

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

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

We will then choose $p(z) \sim \mathcal{N}(0, I)$ and want to $\min_{\theta} KL(q(z) \parallel \mathcal{N}(0, I))$. 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.

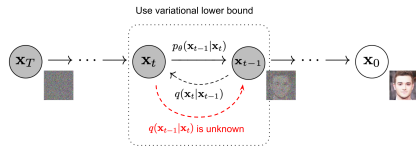
These gaussians result in $KL(\mathcal{N}(\mu, \sigma) \parallel \mathcal{N}(0, I)) = \frac{1}{2}(\sigma^2 + \mu^2 - \ln(\sigma^2) - 1)$

The other part of objective, $\mathbb{E}_q[\ln p(x|z)]$ can be written as $\mathbb{E}_q[\ln p(x|\hat{z})]$ where $x = d(z, \theta)$ and $z = \mu(x, \theta) + \epsilon \odot \sigma(x, \theta)$. 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 $z = \mu + \epsilon \odot \sigma$, $\epsilon \sim \mathcal{N}(0, I)$. Calc KL Loss. Decode \hat{x} using decoder. Calc reconstruction loss. MSE for Gaussian $p(x|\hat{z})$; BCE for Bernoulli. Gradient descent on θ on combined loss, with β parameters on KL divergence, to adjust importance.

22 Diffusion Models

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



Forward process: $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$. Final state is pure noise ($x_T \sim \mathcal{N}(0, I)$). The variance schedule β_1, \dots, β_T are hyperparameters controlling noise addition. Usually starts small, gets larger.

If we unwrap x_t recursively, we get $x_t = \sqrt{\alpha_t}x_0 + \sum_{i=1}^t c_i \epsilon_i$ where $\alpha_t \equiv 1 - \beta_t$, $\tilde{\alpha}_t \equiv \alpha_1 \dots \alpha_t$ and $\epsilon_i \sim \mathcal{N}(0, I)$. Since sum of gaussians is gaussian, then $x_t = \sqrt{\tilde{\alpha}_t}x_0 + \sqrt{1 - \tilde{\alpha}_t}\epsilon$.

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

To do reverse diffusion (generate input from noise), we need to minimize KL-divergence loss of $D_{KL}(q(x_{t-1}|x_t, x_0), p_{\theta}(x_{t-1}|\hat{x}_t))$ for $t = 2, \dots, T$. $q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}|\mu_t, \Sigma_t)$, $\mu_t = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{\beta_t}{\sqrt{1 - \tilde{\alpha}_t}}\epsilon_t)$

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

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

$$\mathbb{E}_{x_0, \epsilon} \left[\left\| \epsilon - \epsilon_{\theta}(\sqrt{\tilde{\alpha}_t}x_0 + \sqrt{1 - \tilde{\alpha}_t}\epsilon) \right\|_2^2 \right]$$

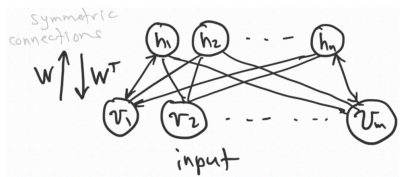
Training Algorithm: Repeat until converge:
1. Take x_0 from dataset
2. $t \sim \text{Uniform}(\{1, \dots, T\})$
3. $\epsilon \sim \mathcal{N}(0, I)$
4. Do gradient descent on simplified Loss function

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. $x_{t-1} = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}}\epsilon_{\theta}(x_t, t)) + \sigma_t z$
3. Return x_0

We typically use a U-net to estimate ϵ_{θ} because upsampling/desampling 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 $E(v, h) = -\sum_{i=1}^m \sum_{j=1}^n v_i W_{ij} h_j - \sum_{i=1}^m b_i v_i - \sum_{j=1}^n c_j h_j$ or rewritten $E(v, h) = -v^T W h^T - b^T v - c^T h$ where $W \in \mathbb{R}^{m \times n}$. Discardance 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 if $\delta E(v_i) > 0$ 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 $q(v, h) = \frac{1}{2} e^{-E(v, h)}$ 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 $\max_v \mathbb{E}_{v \sim p} [\ln q_{\theta}(v)]$

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

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

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

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

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 $q(h_j | V) = \sigma(V W_{.j} + c_j)$ and $\nabla_W L_1 = -V^T \sigma(VW + c)$ which results in a rank-1 outer product in $\mathbb{R}^{m \times n}$

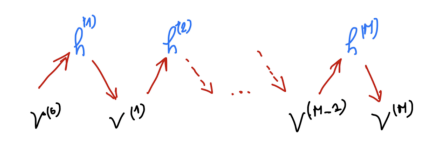
Step 2: Compute expectation using Gibbs sampling. $\langle v_i h_j \rangle_{q(v, h)} = \sum_v \sum_h q(v, h) v_i h_j$. Practically, single network state is used to approximate expectation. Gibbs sampling is used to run the network freely and compute average $v_i h_j$, which has $\nabla_W L_2 = v^T \sigma(VW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

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

Weight update rule is: $W \rightarrow W - \eta(\nabla_W L_1 + \nabla_W L_2) = W + \eta V^T \sigma(VW + c) - \eta v^T \sigma(vW + c)$ 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 $R(f) \triangleq \mathbb{E}_{(x') \sim D} [\text{card}\{\arg \max_i y_i \neq t \mid y = f(x')\}]$

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

Adversarial Attack: Find $x' \in B(x, \epsilon)$ such that $\arg \max_i (y_i) \neq t$ for $y = f(x')$

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

24.1 Fast Gradient Sign Method

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

Minimal Perturbation: Smallest $\|\delta x\|$ causing misclassification: $\min_{\|\delta x\|} [\arg \max_i (y_i(x)) \neq t(x)]$

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

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

Robust model optimizes over $\min_{\theta} \mathbb{E}_{(X, T)} [g(f(X)T)] + \max_{X' \in B(X, \epsilon)} g(f(X')f(X'))]$ 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 $\frac{1}{2M} \min_D \|HD - X\|_2^2$, 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:

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

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

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

27 Transformations

Population coding can pass data between neuron populations by transforming hidden activations. Naive: Decode to data space and then re-encode. Better: Bypass data space by multiplying decoder and encoder weights directly, resulting in rank-1 matrix. $W = D_{eg} E_B \in \mathbb{R}^{N \times M}$ $D_{eg} \in \mathbb{R}^{N \times 1}$ and $E_B \in \mathbb{R}^{1 \times M}$. Using separate decoder-encoder matrices is computationally efficient. Total Time $O(N + M)$ for calculating AD and $(AD)E$ unlike tied weights taking $O(NM)$

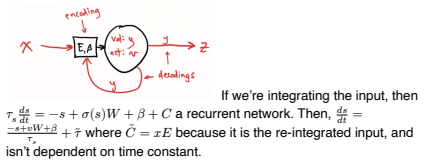
28 Dynamics

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.

$$\begin{aligned} \tau_s \frac{ds}{dt} &= -s + C \text{ Current} \\ \tau_m \frac{dv}{dt} &= -v + \sigma(s) \text{ Activity} \end{aligned}$$

Equilibrium values are when differential values equal 0

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



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

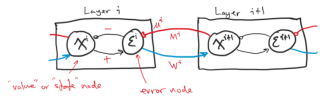
29 Biological Backdrop

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.



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

29.1.1 Error Node

$\tau \frac{d\epsilon^i}{dt} = x^i - \mu^i - v_{\text{bias}}^i \epsilon^i$, 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

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

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

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

29.1.4 Dynamics

$\tau \frac{d\epsilon^i}{dt} = \sigma'(x^i) \odot (\epsilon^{i-1} M^i) - \epsilon^i$ and $\tau \frac{d\epsilon^i}{dt} = x^i - \mu^i - v^i \epsilon^i$. 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 $\frac{\partial F}{\partial \epsilon^i} = -\epsilon^i$

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 $\arg \max_j (x_j^L)$ 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 and Discriminative Network

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

Loss function: $E(\Theta, \varphi) = \mathbb{E}_{z \sim \mathcal{R}, \mathcal{F}} [\mathcal{L}(D(x; \theta), t)] + \mathbb{E}_{z \sim p_z} [\mathcal{L}(D(G(z; \varphi), \varphi), 1)]$ 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

Sequential data inputs are first tokenized into vector embeddings. The input X is transformed into queries (input values) $Q = XW^{(Q)}$, Keys (tags) $K = XW^{(K)}$, and values (desired data) $V = XW^{(V)}$ using weight matrices $W^{(i)} \in \mathbb{R}^{d \times i}$

Self-attention scores are calculated $S = QK^T$ where $S_{ij} = q_i \cdot k_j$ scores query i against key j . Softmax is applied row-wise to $\frac{S}{\sqrt{d}}$ to get attention scores $A \in \mathbb{R}^{n \times n}$. Output of attention head is $A \cdot V = H \in \mathbb{R}^{n \times i}$

Positional encodings: $PE(i)_{2j} = \sin\left(\frac{i}{10000^{2j/d}}\right)$, $PE(i)_{2j+1} = \cos\left(\frac{i}{10000^{2j/d}}\right)$ freq changes with dimension j , and 10000 is scaling constant for large max seq len. Then, $\tilde{x}_i' = \tilde{x}_i + PE(i)$

Multi-Head Attention uses separate attention mechanisms to learn different features, then outputs are concatenated and linearized for output.