

Visualization - Architecture

- Mode-link diagrams: Neurons = nodes, (weighted) connections = edges
- Black diagrams: layer = solid black, single connection between layers = white
- Devis, e.g. Deep Visualization

Training: interesting information about data, parameters, hidden layers, what data

Debugging, improve model design

Parameter Visualizations: of neurons, network layers, activation features used to identify the task, (un)structured Behavior of Neural Networks

Adversarial Examples and Perturbations: adversarial perturbations to cause misclassification

Simple parameter visualization:

- plot learned filter weights directly
- feature activation features via ReLU activation (if convolutional convs deep, probably more important for disambiguation)
- (→) many, relatively specific neuron motifs
- t-SNE visualization of CNN codes, which images regarded as "similar"
- t-distributed Stochastic Neighbor Embedding
- performs dimensionality reduction
- 2d map of images

Gradient-Based Visualization

Backpropagation for Visualization: which pixels are significant for neurons, for which pixel x_i is $\frac{dx_i}{dx}$ use backpropagation

Feature visualization: Reconvnet, change one activation, set others to 0

Global Backpropagation: Improve ReLU by "guiding" the backpropagation, set nearby gradients to zero

Style Transfer: Manipulate content of pixels on class score

Parameter Visualization via Optimization: Optimize Neuron Layer, Layer

Google DeepDream: Interpretation of "dreams" CNN with images → neural image towards high activation of a complex layer

Attention Mechanisms: Attention follow our thought as they write supporting irrelevant information to the task

score(x, h) = $\sum_{i=1}^N \alpha_i x_i$

Soft attention: computed over attention weights

Hard attention: Fixed site glimpses over input

Attention is all you need (NLP)

Multi-Armed Bandit Problem: action a at time t from set A , a_t has pdf $p(a_t)$ with reward r_t , choosing a with prob $\pi(a)$ is policy π → find max $\mathbb{E}(\pi)$ → value $V(\pi)$

Uniform Random: $\pi(a) = \frac{1}{|A|}$

Exploration-Epsilon Greedy: $\pi(a) = \begin{cases} \epsilon & \text{if } a = a_{best} \\ \frac{1-\epsilon}{|A|} & \text{else} \end{cases}$

Gradient Descent: $\frac{d}{dt} \sum_{t=1}^T R_t$

Policy Iteration: The online prediction future reward v_t following $\pi_t(s_t)$ from st

Value Iteration: $v_t(s_t) = \mathbb{E}_{\pi_t} [\sum_{k=0}^{\infty} \gamma^k R_{t+k} | s_t]$

Q-Learning: $Q_t(s, a) = \mathbb{E}_{\pi_t} [\sum_{k=0}^{\infty} \gamma^k R_{t+k} | s, a]$

Temporal Difference Learning: $v_t = R_t + \gamma V_t$

Policy Gradient Theorem: $\frac{d}{d\theta} V_{\pi(\theta)} = \mathbb{E} [\nabla_{\theta} \log \pi(\theta; s, a) (R_t - v_t)]$

Score(x, h) = $\sum_{i=1}^N \alpha_i x_i$

KL Divergence: $D_{KL}(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx$

Applications: Deep Learning: Deep networks approximate policy

Dimensionality Reduction: PCA, t-SNE

Recommender Systems: Collaborative Filtering

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Generative Adversarial Networks: $D = D(\mathcal{D}; x), G = G(\mathcal{D}; z)$

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced GAN methods: Mode collapse, bicubic

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation

Advanced Topics: Adversarial networks for segmentation