

KNN - O(1) train, O(n) test

NN Universal Approximation: with enough training samples, NN can represent any function.

Curse of Dimensionality: for uniform coverage of space, number of training points grows exponentially with dimension

Linear Classifier:

$f(x, W) = Wx + b$ - cannot capture multiple modes of data

- use an average image template per category

- close form soln: $W^* = (X'X)^{-1}X'y$.

Cross Validation: (find best hypoparam) split data into folds, try each as validation set and average, not usually used in DL

Cross-Entropy vs SVM Loss

$$L_i = -\log \left(\frac{\exp(s_{y_i})}{\sum_j \exp(s_j)} \right)$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Gradient Descent Shift:

W momentum: velocity: running mean of gradients, rho: friction

AdaGrad: element-wise scaling w historical sum of squares,

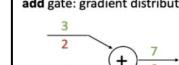
damp steep direction, accelerate flat direction

RMSProp: weighted adagrad

Adam: RMSProp + Momentum + Bias correction

Algorithm	Tracks first moments (Momentum)	Tracks second moments (Adaptive learning rates)	Weighted second moments	Bias correction for moment estimates
SGD	x	x	x	x
SGD+Momentum	✓	x	x	x
AdaGrad	x	✓	x	x
RMSProp	x	x	✓	x
Adam	✓	✓	✓	✓

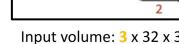
add gate: gradient distributor



mul gate: "swap multiplier"



copy gate: gradient adder



max gate: gradient router



Input volume: 3 x 32 x 32

10 5x5 filters with stride 1, pad 2

Output volume size: 10 x 32 x 32

Number of learnable parameters: 760

Parameters per filter: 3*5*5 + 1 (for bias) = 76

10 filters, so total is 10 * 76 = 760

$10^3 * 32^3 * 32 = 10,240$ outputs; each output is the inner product of two of 3x5x5 tensors (75 elems); total = $75 * 10240 = 768K$

AlexNet Architecture

Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv (C _{out} =20, K=5, P=2, S=1)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool(K=2, S=2)	20 x 14 x 14	
Conv (C _{out} =50, K=5, P=2, S=1)	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool(K=2, S=2)	50 x 7 x 7	
Flatten	2450	
Linear (2450 -> 500)	500	2450 x 500
ReLU	500	
Linear (500 -> 10)	10	500 x 10

Batch Normalization for ConvNets

Batch Normalization for fully-connected networks

$$\begin{aligned} x : N \times D \\ \text{Normalize} \quad | \\ \mu, \sigma : 1 \times D \\ y = \frac{(x - \mu)}{\sigma} y + \beta \end{aligned}$$

$\mu, \sigma : 1 \times C \times C \times 1$

$$y = \frac{(x - \mu)}{\sigma} y + \beta$$

Benefit: model easier to train; faster convergence; robust to initialization; act like regularization; 0 overhead at test time

act like regularization; 0 overhead at test time

Number of output elements = $C \times H' \times W'$

= $64 \times 56 \times 56 = 200,704$

Bytes per element = 4 (for 32-bit floating point)

KB = (number of elements) * (bytes per elem) / 1024

= $200704 * 4 / 1024 = 784$

Can be much less than 784. See Goldberger and Geron's Linear.

Weight shape = $C_{out} \times C_{in} \times K \times K$

= $64 \times 3 \times 11 \times 11$

Bias shape = $C_{out} = 64$

Number of weights = $64 \times 3 \times 11 \times 11 + 64$

= 23,296

Number of floating point operations (multiply+add)

= (number of output elements) * (ops per output elem)

= ($C_{out} \times H' \times W'$) * ($C_{in} \times K \times K$)

= $(64 \times 56 \times 56) * (3 \times 11 \times 11)$

= 200,704 * 363

= 72,855,552

In general:

Input: W

Filter: K

Padding: P

Stride: S

Output: $(W - K + 2P) / S + 1$

$(n_h - k_h + 2p_h + s_h) / s_h$

Activation Functions

Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Graph: sigmoid curve from 0 to 1

tanh

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Graph: tanh curve from -1 to 1

ReLU

$$\max(0, x)$$

Graph: ReLU step function

Matrix Mult BackProp

$$\frac{dL}{dx} = \frac{dL}{dy} w^T$$

[N x D] [N x M] [M x D]

$$\frac{dL}{dw} = x^T (\frac{dL}{dy})$$

[D x M] [D x N] [N x M]

Leaky ReLU

$$\max(0.1x, x)$$

Graph: leaky ReLU with slope 0.1 for x < 0

Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

Graph: maxout with two parallel lines

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

Graph: ELU with negative slope for x < 0

Activation Functions: Gaussian Error Linear Unit (GELU)

- Idea:** Multiply input by 0 or 1 at random; large values more likely to be multiplied by 1, small values more likely to be multiplied by 0 (data-dependent dropout)
- Take expectation over randomness
- Very common in Transformers (BERT, GPT, GPT-2, GPT-3)

Hendrycks and Gimpel, Gaussian Error Linear Units (GELUs), 2019

Models	LeNet	AlexNet	VGG	G-Net	ResNet	ViT	swin-trans
Features	Classic CNN + avg pool + sigmoid	Deeper conv + max pool + relu	More and smaller kernels	Inception: branches conv + global avg. pool + extra classifiers	Residual/ skip connects + batch norm	No conv, patch image -> transformer encoder	Sequential ViTs, each block merges patches - window attention

Diffusion (DDIM): DDPM improvement to remove iterative noise with reparametrization. Same loss.

$$q_\sigma(\mathbf{x}_{1:T}|\mathbf{x}_0) := q_\sigma(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=0}^{T-1} q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$$

$$q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}\mathbf{x}_t + \sqrt{1-\alpha_{t-1}} - \frac{\sigma_t^2}{\sqrt{1-\alpha_t}}\mathbf{x}_t, \sigma_t^2\mathbf{I}\right)$$

Diffusion (DDPM): Iteratively add gaussian noise -> true gaussian distribution (forward) then iteratively denoise (backward). Noise is added w/ hyperparam alpha which can be found for any time step with reparametrization -> tends noise to 0 mean, unique variance. Train model to predict noise (supervised), VLB (KL-Div) loss, simplified - dist b/w real and pred noise per timestep:

$$\begin{aligned} \text{Forward process: } q_\sigma(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) &= \frac{q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q_\sigma(\mathbf{x}_t|\mathbf{x}_0)}{q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_0)} \\ \text{Backward process: } p_\theta^{(t)}(\mathbf{x}_{t-1}|\mathbf{x}_t) &= \begin{cases} \mathcal{N}(f_\theta^{(t)}(\mathbf{x}_t), \sigma_t^2 \mathbf{I}) & \text{if } t = 1 \\ q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, f_\theta^{(t)}(\mathbf{x}_t)) & \text{otherwise,} \end{cases} \\ f_\theta^{(t)}(\mathbf{x}_t) &:= (\mathbf{x}_t - \sqrt{1-\alpha_t} \cdot \epsilon_\theta^{(t)}(\mathbf{x}_t)) / \sqrt{\alpha_t} \\ \mathbf{x}_{t-1} &= \sqrt{\alpha_{t-1}} \left(\mathbf{x}_t - \frac{\sqrt{1-\alpha_t} \cdot \epsilon_\theta^{(t)}(\mathbf{x}_t)}{\sqrt{\alpha_t}} \right) + \sqrt{1-\alpha_{t-1} - \frac{\sigma_t^2}{\sqrt{\alpha_t}}} \cdot \epsilon_\theta^{(t)}(\mathbf{x}_t) + \sigma_t \epsilon_t \end{aligned}$$

$$L_{\text{simple}} = E_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2]$$

$$\begin{aligned} q(\mathbf{x}_5|\mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_0, (1-\alpha_t)\mathbf{I}) \\ q(\mathbf{x}_t|\mathbf{x}_0) &:= \mathcal{N}(\mathbf{x}_t; \mathbf{x}_0, \sqrt{1-\alpha_t}\epsilon) \\ \mathbf{x}_t &= \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1-\alpha_t}\epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, \mathbf{I}) \end{aligned}$$

$$p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{\alpha_t}}\epsilon_\theta(\mathbf{x}_t, t), \sigma_t^2\mathbf{I})$$

Algorithm 1 Training

- repeat
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\epsilon \sim \mathcal{N}(0, \mathbf{I})$
 - 5: Take gradient descent step on
- $$\nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1-\alpha_t}\epsilon, t)\|^2$$
- 6: until converged

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$
- 2: for $t = T, \dots, 1$ do
- 3: $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = 0$
- 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}}\epsilon_\theta(\mathbf{x}_t, t)) + \sigma_t \mathbf{z}$
- 5: end for
- 6: return \mathbf{x}_0

"Slow" R-CNN: Run CNN independently for each region

Fast R-CNN: Apply differentiable cropping to shared image features

Faster R-CNN: Compute proposals with CNN

Single-Stage: Fully convolutional detector

Model does not explicitly compute $p(x)$, but can sample from $p(x)$

Generative models

Model can compute $p(x)$

Explicit density

Can compute $p(x)$

Autoregressive

NADE / MADE

NICE / RealNVP

Flow/Glow

Implicit density

Can compute approximation to $p(x)$

Markov Chain

GSN

Generative Adversarial Networks (GANs)

We will talk about these

Variational Autoencoder

Boltzmann Machine

Diffusion Model

Tractable density

Approximate density

Tractable density

Implicit density

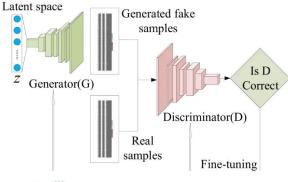
Approximate density

Tractable density

Implicit density

Tractable density

GAN: generator + discriminator, generates image, discriminator discerns if generated or not, use to update generator -> better images

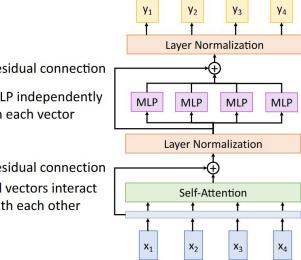


The Transformer

Recall Layer Normalization:
Given h_1, \dots, h_n (Shape: D)
scale: γ (Shape: D)
shift: β (Shape: D)
 $\mu_i = (\sum h_{i,j})/D$ (scalar)
 $\sigma_i = (\sum (h_{i,j} - \mu_i)^2)/D^{1/2}$ (scalar)
 $z_i = (h_{i,j} - \mu_i) / \sigma_i$
 $y_i = \gamma * z_i + \beta$

Ba et al, 2016

Vaswani et al, "Attention is all you need", NeurIPS 2017



Semantic Seg

Downsampling:

pooling, strided conv.

Upsampling:

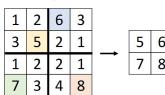
Simple corner place

Nearest neighb.

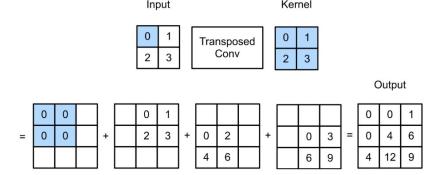
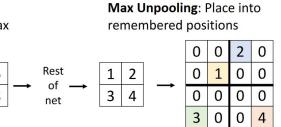
Copy

Bilinear or Cubic
interpol. (2 or 3
neighbor distances)

Max Pooling: Remember which position had the max



Rest of net



Instance segmentation = Mask-RCNN

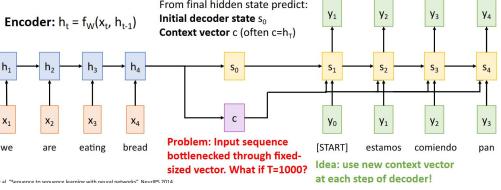
Transformer

Consider a transformer block (with one self-attention layer and one MLP layer) that processes N tokens of dimension D . How many parameters (excluding biases and layer norm) are in this block?

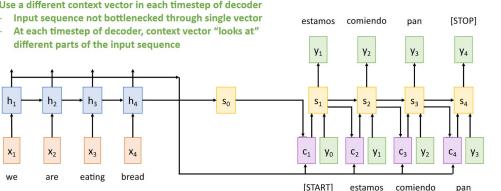
Self-attention layer has $3D^2$ parameters and MLP layer has D^2 parameters, so there are $4D^2$ total parameters.

Sequence-to-Sequence with RNNs

Input: Sequence x_1, \dots, x_T
Output: Sequence y_1, \dots, y_T



Sequence-to-Sequence with RNNs and Attention

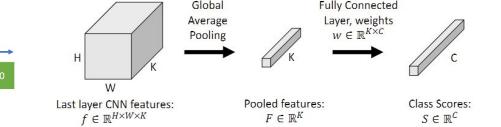


Visualization Methods

Saliency Maps - unsupervised segmentation

Class Activation Maps (CAM):

Class Activation Mapping (CAM)



$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k}, \quad S_c = \sum_k w_{k,c} F_k = \frac{1}{HW} \sum_{h,w} w_{k,c,h,w} f_{h,w,k} = \frac{1}{HW} \sum_{h,w} \sum_k w_{k,c} f_{h,w,k}$$

Grad-CAM - grad. adjusted CAM:

- Pick any layer, with activations $A \in \mathbb{R}^{H \times W \times K}$
- Compute gradient of class score S_c with respect to A : $\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$
- Global Average Pool the gradients to get weights $\alpha \in \mathbb{R}^K$: $\alpha_k = \frac{1}{HW} \sum_{h,w} \frac{\partial S_c}{\partial A}_{h,w,k}$
- Compute activation map $M^c \in \mathbb{R}^{H \times W}$: $M^c_{h,w} = \text{ReLU} \left(\sum_k \alpha_k A_{h,w,k} \right)$

Grad. Ascent: view topological gradient of img

Feat. Inversion - use Adversarial Attack

Feature Pyramid Network - Segmentation

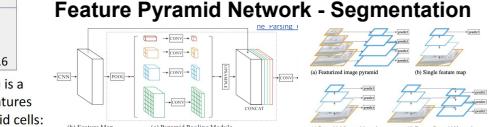
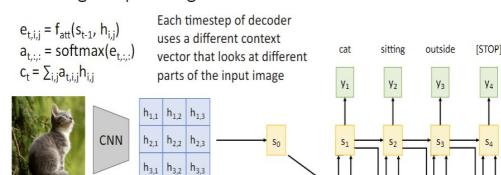
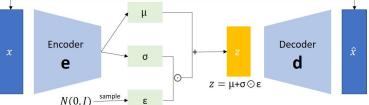


Image Captioning with RNNs and Attention



GAN loss: 1st term discrim. discerns real image, 2nd term discrim. discerns generated image (log-likelihood)

VAE: reconstructs input image by compressing, then sample from compressed (latent) space (normal random sample) but for grad to flow from random sample -> abstract randomness to eps. Can be used for segmentation, inpainting, generation



VAE loss: 1st term loss for expected output image, 2nd term KL-divergence: regularizes latent toward normal distribution.
Modification: disentangle VAE by mult. coeff to 2nd term -> extracts latent features of output image, allows generating diff "styles"

Attention Layer

Inputs:
Query vectors: Q (Shape: $N_Q \times D_Q$)
Input vectors: X (Shape: $N_X \times D_X$)
Key matrix: W_K (Shape: $D_X \times D_Q$)
Value matrix: W_V (Shape: $D_X \times D_V$)

Computation:
Key vectors: $K = XW_K$ (Shape: $N_X \times D_Q$)
Value vectors: $V = XW_V$ (Shape: $N_X \times D_V$)
Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_Q \times N_X$) $E_{ij} = (Q_i \cdot K_j) / \sqrt{D_Q}$
Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)
Output vectors: $Y = AV$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{ij} V_j$

Self-Attention Layer

One query per input vector

Inputs:
Input vectors: X (Shape: $N_X \times D_X$)
Key matrix: W_K (Shape: $D_X \times D_Q$)
Value matrix: W_V (Shape: $D_X \times D_V$)
Query matrix: W_Q (Shape: $D_X \times D_Q$)

Computation:
Query vectors: $Q = XW_Q$
Key vectors: $K = XW_K$
Value vectors: $V = XW_V$
Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_X \times N_X$) $E_{ij} = (Q_i \cdot K_j) / \sqrt{D_Q}$
Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)
Output vectors: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{ij} V_j$

RNN Forward

$$\begin{aligned} y_{t-1} &= h_{t-1} \\ y_t &= v \\ h_t &= Wh^{(t-1)} + Ux^{(t)} \\ h^{(t)} &= \tanh(a^{(t)}) \\ o^{(t)} &= c + Vh^{(t)} \\ \hat{y}^{(t)} &= \text{softmax}(o^{(t)}) \end{aligned}$$

Recurrent Neural Network

Works on **Ordered Sequences**
(+) Good at long sequences: After one RNN layer, it "sees" the whole sequence

(-) Not parallelizable: need to compute hidden states sequentially

1D Convolution

Works on **Multidimensional Grids**
(-) Bad at long sequences: Need to stack many conv layers for outputs to "see" the whole sequence

(-) Highly parallel: Each output can be computed in parallel

(-) Very memory intensive

Transformers in CV:

ViT - encoder class. on image patches

Swin - Hierarchical ViT, merging patches

DETR - bounding box obj. det. using enc-dec arch.

DiT - Diffusion with Transformers (e.g., OpenAI Sora Video)

Self-Attention

Works on **Sets of Vectors**
(-) Good at long sequences: after one self-attention layer, each output "sees" all inputs!

(-) Highly parallel: Each output can be computed in parallel

(-) Very memory intensive

Fast R-CNN vs Slow R-CNN: shared image features that are calculated by a ConvNet. Regions of Interest form a proposal method. Crop+resize features and then per region CNN. Slow had Regions of Interest from a proposal method, warp image, forward region through a ConvNet and then have a bounding box and class. **Faster:** Insert Region Proposal Network to predict proposals from features. **Region Proposal Network:** Run backbone CNN to get features, use anchor box. K boxes. Pass through ConvNet -> normal fastRCNN. Single-Stage: Fully convolutional detector. Featuremap -> RPN->proposals. **Faster:** stage 1(per image): backbone net->reg pro net, stage2(per region): crop->predict object->bbox offset

Evaluating Object Detectors: Mean Average Precision (mAP)

- Run object detector on all test images (with NMS)
- For each category, compute Average Precision (AP) = area under Precision vs Recall curve
 - For each detection (highest score to lowest score)
 - If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
 - Otherwise mark it as negative
 - Plot a point on PR curve
 - Average Precision (AP) = area under PR curve
- Mean Average Precision (mAP) = average of AP for each category
- For "COCO mAP": Compute mAP@thresh for each IoU threshold (0.5, 0.55, 0.6, ..., 0.95) and take average

Dog AP = 0.86

Precision
Recall

All ground-truth dog boxes

Dog AP = 0.86

Precision
Recall

Dog AP = 0.86