

Universal Approximator Theorem: Let $\phi(\cdot)$ be a non-constant, bounded and monotonically increasing fn. $\forall \epsilon > 0$ and any continuous fn $\mathbb{R} \rightarrow \mathbb{R}^m$, there exists an integer N , real constants $v_i, b_i \in \mathbb{R}$ and real vectors $w_i \in \mathbb{R}$ where $i = 1, \dots, N$ such that: $F(\vec{x}) = \sum_{i=1}^N v_i \phi(w_i^T \vec{x} + b_i)$ with $ F(\vec{x}) - f(\vec{x}) < \epsilon$ where ϕ is a sensible activation function. Problems: ϵ can be very large in practice, making approximation less useful, and curse of dimensionality.	
Shift Invariance: The unchanging response when the input is shifted. For classifier f and shift operator $S_v, f(\vec{x}) = f(S_v \vec{x})$ (no matter how the input is transformed, the output should remain constant) generalizes for unseen data.	Shift Equivariance: applying the shift operator after the function yields the same results as applying the function after the shift. i.e. $S_v \circ f(\vec{x}) = f(\vec{x}) \circ S_v$. It is about consistent transformation.
Translation Invariance: shift in input should have a predictable shift in hidden representation (location shouldn't matter) Locality: we should not have to move far from location (i, j) to learn valuable information to asses what the area contains.	
Fully connected net: every input feature n in an image influences ever neuron in the next layer: $n \times n$. Sparsely connected net: each neuron n is connected to a subset of neurons $k: k \times n$. Weight sharing net: weights are reused in a network	
Convolution: $(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau$ for $f, g: [0, \infty] \rightarrow \mathbb{R}$ Correlation: $(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t + \tau)d\tau$ for $f, g: [0, \infty] \rightarrow \mathbb{R}$ commutative: $f * g = g * f$, associative $(f * g) * h = f * (g * h)$ distributive: $f_1 * (f_2 + f_3) = f_1 * f_2 + f_1 * f_3$	$M = \left\lfloor \frac{M+2 \times P-D \times (K-1)-1}{S} + 1 \right\rfloor$
Activation Fns: introduce non-linearity for more complex data learning. Unbounded outputs lead to numerical instability in training. High values cause issues with gradients leading to problems like exploding gradients.	standard CNNs are not naturally equivariant to rotations; Harmonic Networks/H-Nets; replacing kernel w/ 'circular harmonics'; rotation results in a proportionate rotation in the output.
Linear: $f(x) = c \cdot x$ network behaves single-layered model. Regression. Sigmoid: $f(x) = \frac{1}{1+e^{-x}}$ more analogue continuous output than step, smooth gradient between (-2,2) rapid learning between there. Vanishing gradient problem. Good for n-classifiers. Tanh: $f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ stronger gradient, outputs are zero-centred on avg. Vanishing gradient problem, requires good normalization. ReLU: $f(x) = \max(0, x)$ produces sparse activations, comp. Efficient (less dense activations), Dying ReLU problem often outputting 0. Leaky ReLU: $f(x) = \begin{cases} x & \text{for } x \geq 0 \\ 0.01 \cdot x & \text{for } x < 0 \end{cases}$ like ReLU – dying problem	LeNet: convolutions + avg pooling + fully connected last layer (small number of classes make this ok) performs object localization + recognition. AlexNet: + dropout after 2 layers for robustness/regularization, ReLU, maxpool (more shift invariance, keep salient features, discard less useful), softmax for classification, increased kernel size, data augmentation, model ensemble, originally split into two streams due to hardware limitations. VGG: bigger, didn't add more dense layers, didn't add more convolutional layers, but grouped layers into parametrized blocks. Using lots of narrow convolutions outperforms few wider ones, slower than others but performance is much better Inception: introduces parallel dataflow, improves performance with width and depth. ↑size=↑params=↑overfitting=↑comp resource. Patch alignment issues fixed with 1x1, 3x3 (1-padding) and 5x5 (2-padding). Naive version also includes a 3x3 maxpool in parallel for additional benefit. Detailed approach includes 1x1 convolutions to compute reductions before expensive 3x3 and 5x5. Channel size: 1x1 = 64 ch, (1x1=96ch → 3x3=128ch), (1x1=16ch → 5x5=32ch), (3x3 maxpool → 1x1=32ch) because big kernels already come with large parameter count. ✓ It has low param count and FLOPs: $\left\lfloor k^2 \right\rfloor^{fixed} \times c_{in} \times c_{out} \times \left\lfloor m_h \times m_w \right\rfloor^{fixed}$ implies that $\left\lfloor c_{in} \times m_h \times m_w \right\rfloor^{fixed} \times \sum_{j \in paths} (k_j^2 \times c_{out,j})$ By varying no. chan and k we optimize network performance.
PreLU: as above only 0.01 = a , adaptability during training, scale invar. SoftPlus: $f(x) = \frac{1}{\beta} \cdot \log(1 + e^{\beta x})$ differentiable ReLU, if beta high, closer to ReLU, positive output only, non-linear. ELU: CELU with $\frac{x}{\alpha} = x$ element-wise, can be -ve, avg activation pushes to 0 → helps converge faster. CELU: $\max(0, x) + \min(0, \alpha \cdot (e^{\frac{x}{\alpha}} - 1))$ element-wise, when $\alpha \neq 1$ continuously differentiable SELU: $scale \cdot (\max(0, x) + \min(0, \alpha \cdot (e^x - 1)))$ predefined scale and α for mean 0 and var 1, internal normalization, robust (no dying unit) GELU: $x \cdot \phi(x)$ cdf for gaussian, introduces probabilistic, regularization. ReLU6: $\min(\max(0, x), 6)$ saturates activations, LogSigmoid: $\log(\frac{1}{1+e^{-x}})$ element-wise, better for cost functions, SoftMin: $\frac{e^{-x_i}}{\sum_j e^{-x_j}}$ rescale inputs to sum to 1, multi-dimensional non-linearity, resemble a probability distribution. Emphasises smallest values. SoftMax: $softmax(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}}$ rescale to $\sum = 1$, multi-label class. LogSoftmax: $\log(\frac{e^{x_i}}{\sum_j e^{x_j}})$ possible to simplify mathematical calculations and potentially improve the training process Period activation: SIREN: $\Phi(\vec{x}) = \vec{W}_n(\phi_{n-1} \circ \phi_{n-2} \circ \dots \circ \phi_0)(\vec{x}) + \vec{b}_n$ $\vec{x}_i \mapsto \phi_i(\vec{x}_i) = \sin(\vec{W}_i \vec{x}_i + \vec{b}_i)$ deals with implicit representation involving finding a continuous function that represents sparse input data, such as images.	ResNet: parameterizes around an identity function instead of 0 function; you don't have to learn the identity function from scratch, allows for the addition of new layers without disrupting the outputs of previous layers. 3X3 → Batchnorm → ReLU → 3x3 → BatchNorm → combine with identity/1x1 → ReLU. DenseNet: each layer gets the feature maps from all the preceding layers (taylor expansion style), efficient and less parameters, alleviate the vanishing gradient problem by improving the gradient flow through the network, making optimization easier, scalable, since you can easily adjust the number of layers. Occasionally need to reduce resolution (transition layer) SqueezeNet: uses attention, <i>Squeeze</i> global average pool, fed into a <i>Descriptor</i> (small fully-connected NN) capture which channels are more important given the current global context of the image. Use descriptors to <i>Excite</i> and scale the original channels <i>Re-weight channels</i> if certain channels need to emphasize certain features Outcomes: Dense layers are computationally and memory intensive. Real-world problems with big input tensors and many classes will prohibit their use, 1x1 convolutions act like pixel-wise multi-layer perceptron,
Loss: how well your network is doing, quantifying Δ (predicted outputs, ground truth). Backprob updates model parameters, aiming to min. error. L2-norm: $\ell(x, y) = l_n = (x_n - y_n)^2$ regression. L1-norm: $\ell(x, y) = l_n = x_n - y_n $ sensitive to outliers, not differentiable Smooth L1: $\frac{1}{2} \sum_i z_i, z_i = \begin{cases} 0.5(x_i - y_i)^2, & \text{if } x_i - y_i < 1 \\ x_i - y_i - 0.5, & \text{otherwise} \end{cases}$ for errors close to zero, it behaves like the L2 loss, else L1. Robust to outliers. Negative Log-likelihood: $\ell(x, y) = \hat{\mathcal{L}} = \{l_1, \dots, l_N\}^T, l_n = -w_{y_n} x_n, y_n, w_c = weight[c] \cdot 1\{c \neq ignore_{index}\} \ell(x, y) = \begin{cases} \sum_{n=1}^N \frac{1}{\sum_{n=1}^N w_{y_n}} l_{n_i} & \text{if reduction = mean} \\ \sum_{n=1}^N l_{n_i} & \text{if reduction = sum} \end{cases}$ element-wise, weights assign importance, Cross-entropy: $loss(x, class) = -\log(\frac{e^{x[class]}}{\sum_j e^{x[j]}}) = -x[class] + \log(\sum_j e^{x[j]})$ n-classification, weighted classes optional Binary cross-entropy: $\ell(x, y) = \mathcal{L} = \{l_1, \dots, l_N\}^T$ $l_n = -w_n[y_n \cdot \log x_n + (1 - y_n) \cdot \log(1 - x_n)], x_n$ predicted positive, y_n ground truth (0,1), w_n weight. Binary/Multi classification Kullback-Libeler $\ell(x, y) = \mathcal{L} = \{l_1, \dots, l_N\}^T, l_n = y_n \cdot (\log y_n - x_n)$ suitable when target is 1-hot, suffers with numerical stability issues Margin Ranking Loss/Ranking Losses/Contrastive Loss: $loss(x, y) = \max(0, -y \cdot (x_1 - x_2) + margin)$ useful to push classes as far away as possible Triplet Margin Loss: $l_n(x_a, x_p, x_n) = \max(0, m + f(x_a) - f(x_p) - f(x_a) - f(x_n))$ make samples from same classes close and different classes far away e.g. Siamese Networks Cosine Embedding Loss $loss(x, y) = \begin{cases} a - \cos(x_1, x_2), & \text{if } y = 1 \\ \max(0, \cos(x_1, x_2) - margin), & \text{if } y = -1 \end{cases}$ Measure whether two inputs are similar or dissimilar, 1 sim (tries to minimize the angle between the vectors), -1 not sim (aims to ensure that the cosine similarity is smaller than a specified margin).	Curse of Dimensionality: Sample Explosion: As the number of features or dimensions grows, the amount of data we need to generalize accurately grows exponentially To approximate a (Lipschitz) continuous function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ with ϵ accuracy one needs $O(\epsilon^{-d})$ samples. Sparseness: The more features we use, the more sparse the data becomes such that accurate estimation of the classifier's parameters (i.e. its decision boundaries) becomes more difficult, this sparseness is not uniformly distributed over the search space; the higher dimensions you have the higher probability that a data-point will sit in its own distinct corner in the hypercube. Math: $V_{rind}^n = (1 - \alpha^n) V_{original} \Rightarrow \frac{V_{rind}}{V_{original}} = 1 - \alpha^n \Rightarrow \frac{d(1 - \alpha^n)}{d\alpha} = -n\alpha^{n-1}$ $\Leftrightarrow d(1 - \alpha^n) = -n\alpha^{n-1} d\alpha$, this shows that the volume of the rind initially grows much faster, n times faster than the rate at which the object is being shrunk (when $\alpha = 1$ and $d\alpha < 0$ then $d(1 - \alpha^n) = n d\alpha $); In higher dims, small changes in distance lead to vast changes in vol. Factorized conv: 2 3x3 convolutions can act as an approx for 5x5 conv trading expressiveness for efficiency. Inserting a non-linearity between the 3x3s lets it capture more complex features. Separable conv: approximate a 5x5 conv with a 5x1 and 1x5 reducing params (as above). Lossy. Pooling: smaller res, hierarchal features (concentrates abstract features), shift/deform Invariance. Break shift-equivariance by blurring sample to avoid the shifting pooling issue. Approximate Deformation Invar: $\ f(\vec{x}) - f(D_\tau \vec{x})\ \approx \ \nabla f\ \text{deform img } \tau \text{-deform factor}$