

Machine Perception

Cristian Perez Jensen

April 7, 2024

Contents

1	Neural networks	1
1.1	Multi-layer perceptron	1
1.2	Loss functions	1
1.3	Backpropagation	2
1.4	Activation functions	2
1.5	Universal approximation theorem	2
2	Convolutional neural networks	4
2.1	Convolution	4
2.2	Convolutional neural network	4
3	Fully convolutional neural network	7
3.1	Upsampling methods	7
3.2	U-net	7
4	Recurrent neural networks	9
4.1	Elman RNN	9
4.2	Long-short term memory	10
4.3	Gradient clipping	11
5	Autoencoders	13
5.1	Linear autoencoders	13
5.2	Non-linear autoencoders	13
5.3	Variational autoencoders	14
5.4	β -VAE	15
6	Autoregressive models	17
6.1	Fully visible sigmoid belief network	17
6.2	Neural autoregressive density estimator	18
6.3	Masked autoencoder distribution estimation	18
6.4	Generating images	19
6.5	Generating audio	20
6.6	Variational RNN	20
6.7	Transformers	20

List of symbols

\doteq	Equality by definition
\approx	Approximate equality
\propto	Proportional to
\mathbb{N}	Set of natural numbers
\mathbb{R}	Set of real numbers
$i : j$	Set of natural numbers between i and j . I.e., $\{i, i+1, \dots, j\}$
$\mathbb{1}\{\text{predicate}\}$	Indicator function (1 if predicate is true, otherwise 0)
$\boldsymbol{v} \in \mathbb{R}^n$	n -dimensional vector
$\boldsymbol{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
$\boldsymbol{T} \in \mathbb{R}^{d_1 \times \dots \times d_n}$	Tensor
\boldsymbol{M}^\top	Transpose of matrix \boldsymbol{M}
\boldsymbol{M}^{-1}	Inverse of matrix \boldsymbol{M}
$\det(\boldsymbol{M})$	Determinant of \boldsymbol{M}
$\frac{d}{dx}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial x}f(x)$	Partial derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}^n$
$\nabla_x f(x) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $x \in \mathbb{R}^n$
$D_x f(x) \in \mathbb{R}^{n \times m}$	Jacobian of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at point $x \in \mathbb{R}^n$
$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $x \in \mathbb{R}^n$
$\boldsymbol{\theta} \in \Theta$	Parametrization of a model, where Θ is a compact subset of \mathbb{R}^K
\mathcal{X}	Input space
\mathcal{Y}	Output space
$\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$	Labeled training data

1 Neural networks

1.1 Multi-layer perceptron

The original perceptron [Rosenblatt, 1958] was a single layer perceptron with the following non-linearity,

$$\sigma(x) \doteq \mathbb{1}\{x > 0\}.$$

The classification of a single point can then be written as

$$\hat{y} = \mathbb{1}\{\mathbf{w}^\top \mathbf{x} > 0\}.$$

The learning algorithm then iteratively updates the weights for a data point that was classified incorrectly,

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \underbrace{(y_i - \hat{y}_i)}_{\text{residual}} \mathbf{x}_i,$$

where η is the learning rate. If the data is linearly separable, the perceptron converges in finite time.

The problem with the single-layer perceptron was that it could not solve the XOR problem; see Figure 2. This can be solved by introducing hidden layers,

$$\hat{y} = \sigma(\mathbf{W}_k \sigma(\mathbf{W}_{k-1} \cdots \sigma(\mathbf{W}_1 \mathbf{x}))).$$

We call this architecture a multi-layer perceptron (MLP); see Figure 3. We then want to estimate the parameters $\boldsymbol{\theta} = \{\mathbf{W}_1, \dots, \mathbf{W}_k, \mathbf{b}_1, \dots, \mathbf{b}_k\}$, using an optimization algorithm such as gradient descent, which we call “learning”.

1.2 Loss functions

We need an objective to optimize for. We typically call this objective function the loss function, which we minimize. In classification, we typically optimize the maximum likelihood estimate (MLE),

$$\begin{aligned} \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathcal{D} \mid \boldsymbol{\theta}) &\stackrel{\text{iid}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^n p(y_i \mid \boldsymbol{\theta}) \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} -\log \prod_{i=1}^n p(y_i \mid \boldsymbol{\theta}) \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} -\sum_{i=1}^n \log p(y_i \mid \boldsymbol{\theta}). \end{aligned}$$

If the model predicts the parameters of a Bernoulli distribution, MLE is equivalent to binary cross-entropy,

$$\begin{aligned} \ell(\boldsymbol{\theta}) &= -\log \operatorname{Ber}(y_i \mid \hat{y}_i \doteq f(\mathbf{x}_i \mid \boldsymbol{\theta})) \\ &= -\log \hat{y}_i^{y_i} (1 - \hat{y}_i)^{1-y_i} \\ &= -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i). \end{aligned}$$

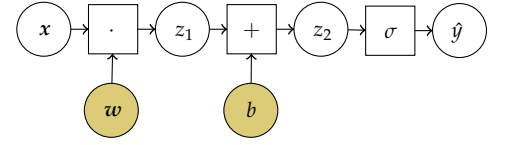


Figure 1. Computation graph of a perceptron [Rosenblatt, 1958], where $\sigma(x) = \mathbb{1}\{x > 0\}$.

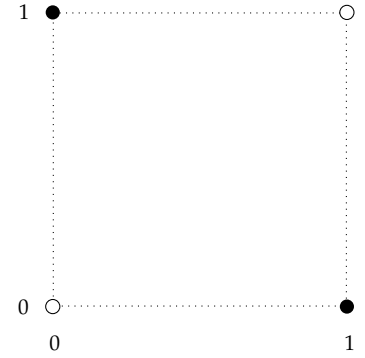


Figure 2. XOR problem. As can be seen, the data is not linearly separable, and thus not solvable by the perceptron.

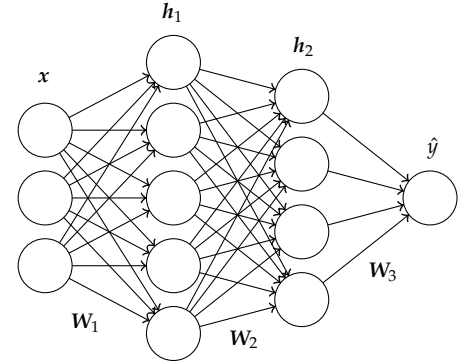


Figure 3. Example multi-layer perceptron architecture.

f is a model that outputs the Bernoulli parameter.

If we choose the model to be Gaussian, we end up minimizing the mean-squared error. Furthermore, the Laplacian distribution yields minimizing the ℓ_1 norm.

If we have prior information about the weights, we could also optimize for the maximum a posteriori (MAP),

$$\begin{aligned}
 \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathcal{D}) &= \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta}) p(\mathcal{D} \mid \boldsymbol{\theta}) \\
 &\stackrel{\text{iid}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta}) \prod_{i=1}^n p(y_i \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) \\
 &= \operatorname{argmin}_{\boldsymbol{\theta}} -\log \left(p(\boldsymbol{\theta}) \prod_{i=1}^n p(y_i \mid \boldsymbol{\theta}) \right) \\
 &= \operatorname{argmin}_{\boldsymbol{\theta}} -\log p(\boldsymbol{\theta}) - \sum_{i=1}^n \log p(y_i \mid \boldsymbol{\theta})
 \end{aligned}$$

Note that MAP and MLE are equivalent if $p(\boldsymbol{\theta})$ is uniform over the domain of weights. Assuming a Gaussian prior distribution over $\boldsymbol{\theta}$, MAP yields Ridge regression.

1.3 Backpropagation

Typically, we cannot find the optimal parameters $\boldsymbol{\theta}^*$ in closed form, so we must use an optimization algorithm. Optimization algorithms, such as gradient descent, typically require computing the gradient w.r.t. the parameters. Backpropagation is an algorithm for computing the gradient of any function, given that we have access to the derivatives of the primitive functions that make up the function.¹ It then computes the gradient by making use of dynamic programming, the chain rule, and sum rule.

¹ For example, to compute the gradient of $f(\mathbf{x}, \mathbf{y}) = \sigma(\mathbf{x}^\top \mathbf{y})$, we would need access to $\frac{d}{dx} \sigma(x)$, $\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{y}$, and $\frac{\partial}{\partial \mathbf{y}} \mathbf{x}^\top \mathbf{y}$.

Gradient descent iteratively updates the parameters by the following,

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}),$$

until the gradient is small.

1.4 Activation functions

In MLPs, the activation function should be non-linear, or the resulting MLP is just an affine mapping with extra steps. This is because the product of affine mappings are themselves affine mappings.

1.5 Universal approximation theorem

Theorem 1 (Universal approximation theorem). Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be a non-constant, bounded, and continuous activation function. Let I_m denote the m -dimensional unit hypercube $[0, 1]^m$ and the space of real-valued function on I_m is denoted by $\mathcal{C}(I_m)$.

Let $f \in \mathcal{C}(I_m)$ be any function in the hypercube. Let $\epsilon > 0, N \in \mathbb{N}, v_i, b_i \in \mathbb{R}, w_i \in \mathbb{R}^m$ for $i = 1, \dots, N$, then

$$f(\mathbf{x}) \approx g(\mathbf{x}) = \sum_{i=1}^N v_i \sigma(\mathbf{w}_i^\top \mathbf{x} + b_i),$$

where $|f(\mathbf{x}) - g(\mathbf{x})| < \epsilon$ for all $\mathbf{x} \in I_m$.

The universal approximation theorem holds for any single hidden layer network. However, this hidden layer may need to have infinite width to approximate f . In practice, deeper networks work better than wider networks.

2 Convolutional neural networks

When dealing with high-dimensional data, such as images, it is not practical to work with MLPs, because the amount of parameters would explode in size.² By making use of the locality of images, we can drastically decrease the number of parameters.

2.1 Convolution

The *correlation* operator takes a filter \mathbf{K} , moves it along the entire image, and outputs the patch-wise multiplication, for each patch of the same size as the filter. It is defined as follows,

$$(\mathbf{K} \star \mathbf{I})[i, j] = \sum_{m=-k}^k \sum_{n=-k}^k \mathbf{K}[m, n] \mathbf{I}[i + m, j + n].$$

The *convolution* operator is very similar. The only difference is that the kernel is mirrored in a convolution,

$$(\mathbf{K} \ast \mathbf{I})[i, j] = \sum_{m=-k}^k \sum_{n=-k}^k \mathbf{K}[m, n] \mathbf{I}[i - m, j - n].$$

Theoretically, the convolution operator is more useful, because it is commutative.³ In practice with neural networks, it does not matter, since the weights will just be updated to be the same, except that they are mirrored. Thus, we will only be referring to the convolution from now on.

A convolution operator C is a linear, shift-equivariant transformation, *i.e.*,

$$\begin{aligned} C(\alpha \mathbf{x} + \beta) &= \alpha C(\mathbf{x}) + \beta \\ T_t(C(\mathbf{x})) &= C(T_t(\mathbf{x})). \end{aligned}$$

Since convolutions are linear, discrete convolutions can be implemented using matrix multiplication,

$$\mathbf{K} \ast \mathbf{I} = \begin{bmatrix} k_1 & 0 & 0 & \cdots & 0 \\ k_2 & k_1 & 0 & \cdots & 0 \\ k_3 & k_2 & k_1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & k_m \end{bmatrix} \begin{bmatrix} I_1 \\ I_2 \\ I_3 \\ \vdots \\ I_n \end{bmatrix}.$$

2.2 Convolutional neural network

A convolutional neural network (CNN) is composed of a sequence of *convolutional layers* and *pooling layers*, followed by a final *dense layer* (MLP).

² Mapping a $256 \times 256 \times 3$ input image to a 1-dimensional output would already require nearly 2 million parameters.

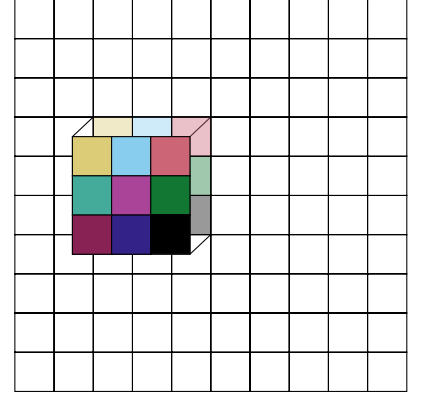


Figure 4. Illustration of applying a correlation to a pixel.

³ $\mathbf{I} \ast \mathbf{K} = \mathbf{K} \ast \mathbf{I}$, but $\mathbf{I} \star \mathbf{K} \neq \mathbf{K} \star \mathbf{I}$.

Linearity.

Translation equivariant.

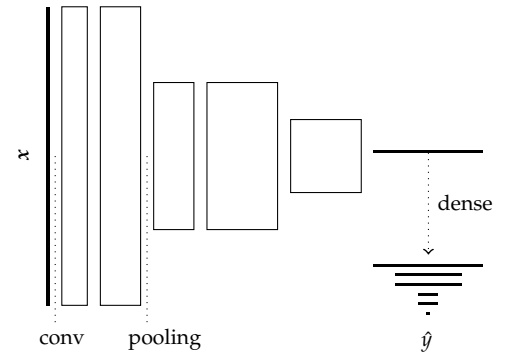


Figure 5. Example schematic of a CNN architecture.

Convolutional layer. A convolutional layer applies a filter \mathbf{W} to an input image \mathbf{Z} by convolution. This filter is learned. Hence, we need to derive its derivative. We will focus on the single filter case for simplicity, whose forward pass is computed by

$$z^{(\ell)}[i, j] = \sum_{m=-k}^k \sum_{n=-k}^k w^{(\ell)}[m, n] z^{(\ell-1)}[i - m, j - n] + b.$$

We can express the derivative of the cost function \mathcal{L} w.r.t. the output of the $(\ell - 1)$ -th layer as the following,

$$\begin{aligned} \delta^{(\ell-1)}[i, j] &= \frac{\partial \mathcal{L}}{\partial z^{(\ell-1)}[i, j]} \\ &= \sum_{i'} \sum_{j'} \frac{\partial \mathcal{L}}{\partial z^{(\ell)}[i', j']} \frac{\partial z^{(\ell)}[i', j']}{\partial z^{(\ell-1)}[i, j]} \\ &= \sum_{i'} \sum_{j'} \delta^{(\ell)}[i', j'] \frac{\partial}{\partial z^{(\ell-1)}[i, j]} \sum_m \sum_n w^{(\ell)}[m, n] z^{(\ell-1)}[i' - m, j' - n] + b \\ &= \sum_{i'} \sum_{j'} \delta^{(\ell)}[i', j'] w^{(\ell)}[i' - i, j' - j]. \end{aligned}$$

From this, we can see that we can compute all values of $\delta^{(\ell-1)}$ by a single convolution,

$$\delta^{(\ell-1)} = \delta^{(\ell)} * \text{Rot}_{180}(\mathbf{W}^{(\ell)}) = \delta^{(\ell)} \star \mathbf{W}^{(\ell)}.$$

Using this value, we can compute the derivative w.r.t. the weights, which we need for the parameter update in algorithms such as gradient descent,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w^{(\ell)}[m, n]} &= \sum_i \sum_j \frac{\partial \mathcal{L}}{\partial z^{(\ell)}[i, j]} \frac{\partial z^{(\ell)}[i, j]}{\partial w^{(\ell)}[m, n]} \\ &= \sum_i \sum_j \delta^{(\ell)}[i, j] \frac{\partial}{\partial w^{(\ell)}[m, n]} \sum_{m'} \sum_{n'} w^{(\ell)}[m', n'] z^{(\ell-1)}[i - m', j - n'] + b \\ &= \sum_i \sum_j \delta^{(\ell)}[i, j] z^{(\ell-1)}[i - m, j - n]. \end{aligned}$$

Again, this has the form of a convolution, thus we can compute all derivatives of $\mathbf{W}^{(\ell)}$ by convolution,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(\ell)}} = \delta^{(\ell)} * \mathbf{Z}^{(\ell-1)}.$$

As shown, we can use convolutions for both computing the forward pass, as well as the backward pass. Thus, we first do the forward pass, then compute all $\delta^{(\ell)}$ for all layers ℓ by convolution, and finally we can compute the derivative by convolution as well.

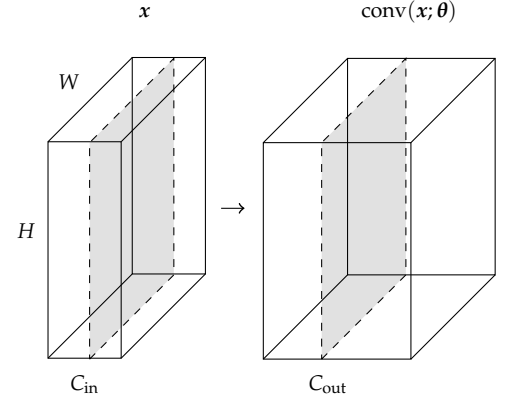


Figure 6. Schematic of a convolutional layer. Each input-output channel pair has its own kernel, so θ has $K \times K \times C_{\text{in}} \times C_{\text{out}}$ parameters.

Pooling layers. Pooling layers makes the data more manageable. The most common pooling layer is *max pooling*, which outputs the maximum value for each patch. It can be seen as a non-linear convolutional filter, where it simply outputs the maximum value. Usually, pooling is done with a stride, such that the output becomes exponentially smaller; see Figure 7. Because of this, the *receptive field* becomes exponentially larger.

The forward pass is computed as follows,

$$z^{(\ell)}[i, j] = \max \left\{ z^{(\ell)}[i', j'] \mid i' \in [si : si + k], j' \in [sj : sj + k] \right\},$$

where s is the stride and k is the kernel size. Let $[i^*, j^*]$ be the indices which corresponded to the maximum value in the forward pass, then we can compute the error propagation in the backward pass by

$$\frac{\partial z^{(\ell)}[i, j]}{\partial z^{(\ell-1)}[i', j']} = \mathbb{1}\{[i', j'] = [i^*, j^*]\}.$$

Note that the max pooling layer has no learnable parameters. Hence, the backward pass is only a propagation of the error, and not used for a weight update.

Dense layer. The dense layer is simply a linear layer that maps the final convolutional layer to the network's output. All previous convolutional and pooling layers can be seen as extracting features from the image, while the final dense layer makes the actual prediction.

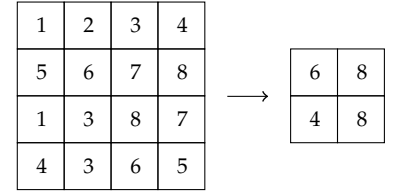


Figure 7. Toy example of max pooling.

3 Fully convolutional neural network

Semantic segmentation is a computer vision task that involves assigning a semantic class to each pixel in an image. While in image classification, the model must output a single class for the entire image, semantic segmentation requires classifying a class for each pixel individually.

A naive approach would be to apply a single convolutional layer to an image, and then running a classifier on each individual pixel. However, this method is inefficient, because we have to run the classifier $H \times W$ times. Instead, we use the output of convolutional neural networks. A naive approach of using CNNs would be to simply apply n convolutional layers with no downsampling, and then considering the last output as the predicted segmentation map. However, this method is expensive.

In practice, the most common approach is to downsample the features obtained using convolution and pooling layers and then upsample them again. By downsampling, this method is more computationally efficient, has larger receptive field, and suffers less from “The curse of dimensionality”. By upsampling, the model produces an output of the same resolution as the input.

3.1 Upsampling methods

Nearest neighbor. Nearest neighbor upsampling copies the same value into all corresponding pixels at a higher resolution; see Figure 8.

Bed of nails. Bed of nails upsampling only copies each value once into the output in the top left value, and pads the rest with zero; see Figure 9.

Max unpooling. Max unpooling also uses zero padding, like bed of nails. However, it also remembers the original position of the maximum value before the corresponding max pooling in the downsampling phase. This information is then used to place each element back in the correct position; see Figures 7 and 10.

Transposed convolutions. Transposed convolution is a learned upsampling technique. This layer learns a kernel that is used to produce the terms whose sum will be the final output. Each term is obtained by multiplying all the element of the kernel by the same value of one single input pixel and then inserting the result in the correct position of a matrix of the same size as the output.

3.2 U-net

The *U-net* is a FCNN architecture, whose main idea is to combine global and local feature maps by copying corresponding tensors from earlier stages in each upsampling stage; see Figure 11. This allows the network

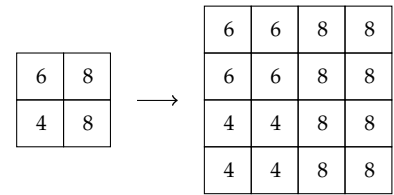


Figure 8. Nearest neighbor upsampling.

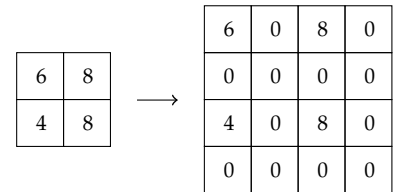


Figure 9. Bed of nails upsampling.

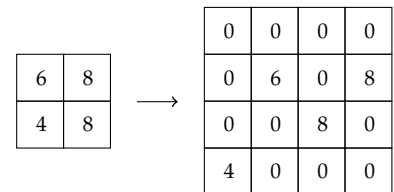


Figure 10. Max unpooling the output of Figure 7.

to capture both local and global context. In each upsampling, the corresponding output from the downsampling phase is appended to the output. The copied tensor can be seen as the “global” information, while the input of the upsampling layer is the “local” information. Combining these allows for more fine-grained outputs.

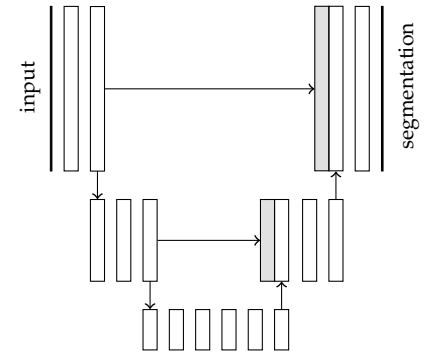


Figure 11. U-net architecture. Down arrows are downsampling layers, up arrows are upsampling layers, and right arrows copy.

4 Recurrent neural networks

Recurrent neural networks (RNN) are a type of neural network that processes sequential data, such as text and video. Unlike traditional neural networks, which take fixed-length inputs, RNNs can take inputs of variable length.⁴

RNNs can have different applications, for example, *one to one*, where at each time step we have one input and one output,⁵ *one to many*, where we have one input and we output a sequence of elements,⁶ *many to one*, where we have a sequence of inputs and one output,⁷ and *many to many*, where we map a sequence to another sequence of a different length.⁸

4.1 Elman RNN

The Elman RNN is characterized by a hidden vector $\mathbf{h}^{(t)}$, which forms the state of the network at timestep t . The hidden state is updated at each timestep by combining the previous hidden state with the input,

$$\mathbf{h}^{(t)} = \tanh(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)}).$$

We can then use the hidden state as a representation of the input up until that timestep. Thus, we can use it as the input to a feed-forward neural network,

$$\hat{\mathbf{y}}^{(t)} = \mathbf{W}_y \mathbf{h}^{(t)}.$$

Then, we can compute the loss function as the sum of each individual loss function,

$$\mathcal{L} \doteq \sum_{t=1}^T \ell^{(t)}.$$

We use *backpropagation through time* (BPTT) to compute the gradient of an RNN. This involves first unrolling the RNN; see Figure 12. Then we can compute the gradient by backpropagation on the resulting computational graph,

$$\frac{\partial \ell^{(t)}}{\partial \mathbf{W}} = \sum_{k=1}^t \frac{\partial \ell^{(t)}}{\partial \hat{\mathbf{y}}^{(t)}} \frac{\partial \hat{\mathbf{y}}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}},$$

where $\frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}}$ is the immediate derivative that treats $\mathbf{h}^{(k-1)}$ as constant w.r.t. \mathbf{W} .

Let's only consider the following term of the product,

$$\begin{aligned} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} &= \prod_{i=k+1}^t \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{h}^{(i-1)}} \\ &= \prod_{i=k+1}^t \frac{\partial}{\partial \mathbf{h}^{(i-1)}} \sigma(\mathbf{W}_h \mathbf{h}^{(i-1)} + \mathbf{W}_x \mathbf{x}^{(i)}) \\ &= \prod_{i=k+1}^t \mathbf{W}_h^\top \text{diag}(\sigma'(\mathbf{W}_h \mathbf{h}^{(i-1)} + \mathbf{W}_x \mathbf{x}^{(i)})). \end{aligned}$$

⁴This is useful for data structures such as text, where the number of words in a text is not fixed.

⁵E.g., part-of-speech tagging.

⁶E.g., image captioning, where the image is the one input, and the caption is a sequence of words.

⁷E.g., sentiment classification, where the input is a text and the output is a single output that determines how positive or negative the text is.

⁸E.g., machine translation, where we map a sentence in one language to a sentence of another.

We use the tanh activation function, because it is centered at 0.

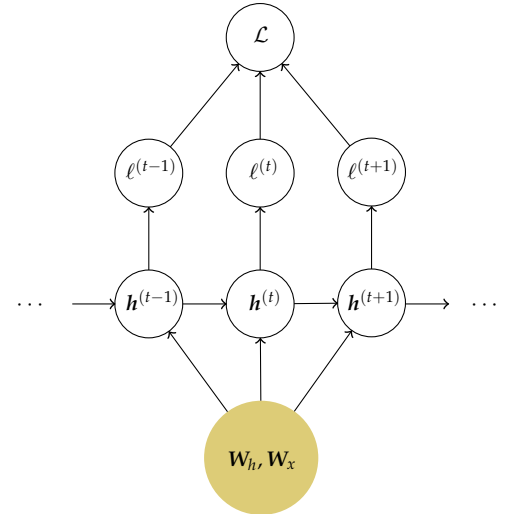


Figure 12. The computational graph of an unrolled recurrent neural network. The inputs $\mathbf{x}_{1:T}$ and outputs $\mathbf{y}_{1:T}$ are omitted.

Assuming that the norm of the gradient of σ is upper bounded by some $\gamma \in \mathbb{R}$,⁹ i.e.,

$$\left\| \text{diag}\left(\sigma'\left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)}\right)\right) \right\| < \gamma.$$

⁹ For example, the gradient of \tanh is bounded by 1.

Let λ_1 be the largest eigenvalue of \mathbf{W}_h , then we have two cases,

1. $\lambda_1 < \frac{1}{\gamma}$. Then we have the following,

$$\left\| \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(i-1)}} \right\| \leq \|\mathbf{W}_h\| \left\| \text{diag}\left(\sigma'\left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)}\right)\right) \right\| < \frac{1}{\gamma} \gamma = 1.$$

Triangle inequality and $\|\mathbf{W}_h\| = \lambda_1$.

Let $\eta < 1$ be the upper bound of all gradients between $\mathbf{h}^{(i)}$ and $\mathbf{h}^{(i-1)}$, then by induction, we have

$$\left\| \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} \right\| = \left\| \prod_{i=k+1}^t \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{h}^{(i-1)}} \right\| < \eta^{t-k}.$$

This converges to zero, as $t \rightarrow \infty$. Thus, we have a *vanishing gradient*;

2. $\lambda_1 > \frac{1}{\gamma}$. Then we have the following,

$$\left\| \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(i-1)}} \right\| \leq \|\mathbf{W}_h\| \left\| \text{diag}\left(\sigma'\left(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_x \mathbf{x}^{(t)}\right)\right) \right\| > \frac{1}{\gamma} \gamma = 1.$$

Using the same logic as in the other case, this yields

$$\left\| \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} \right\| = \left\| \prod_{i=k+1}^t \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{h}^{(i-1)}} \right\| > \eta^{t-k}.$$

for an upper bound $\eta > 1$. Thus, this diverges to ∞ as $t \rightarrow \infty$. Thus, we have a *exploding gradient*.

Thus, we will always have vanishing or exploding gradients when using an Elman RNN. This makes it hard for the architecture to capture long-term dependencies.

4.2 Long-short term memory

The *long-short term memory* (LSTM) architecture [Hochreiter and Schmidhuber, 1997] fixes the vanishing gradient problem of the Elman RNN by making sure there is always a path between hidden units, such that errors can always propagate; see Figure 13.

The cell of an LSTM consists of 4 layers, called gates. In particular, these gates have the following instructions,

- f is the *forget gate* and has the role of scaling the old cell state $\mathbf{h}^{(t-1)}$. It “decides which information should be forgotten” from the previous cell state, based on the new input $\mathbf{x}^{(t)}$;
- i is the *input gate* and has the role of “deciding which values of the cell state $\mathbf{c}^{(t)}$ should be updated” at the current time step;

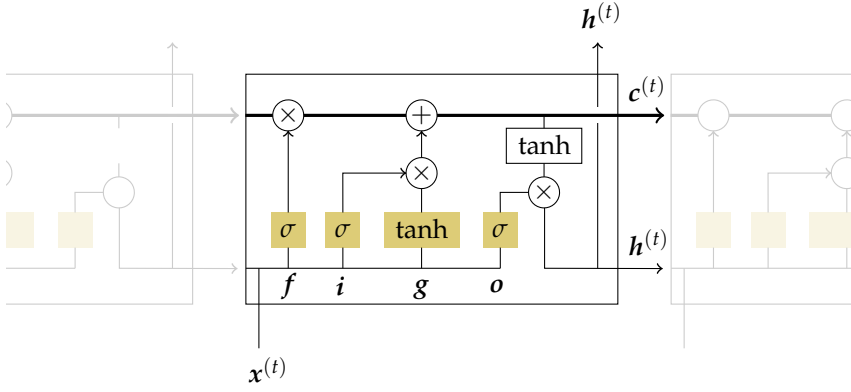


Figure 13. LSTM architecture. The yellow squares are neural networks, and the white squares are point-wise operators. As can be seen, there is an “information highway” that can easily propagate errors at the top, because of the minimal modifications made to it.

- o is the *output gate* and has the role of “deciding which values of the current cell state $c^{(t)}$ should be put in the output of the cell $h^{(t)}$ ”.
- g is the *gate* that decides what to write in the cell state $h^{(t)}$.

We first compute all the gates,

$$\begin{aligned} f^{(t)} &= \sigma(W_{hf}h^{(t-1)} + W_{xf}x^{(t)}) \\ i^{(t)} &= \sigma(W_{hi}h^{(t-1)} + W_{xi}x^{(t)}) \\ o^{(t)} &= \sigma(W_{ho}h^{(t-1)} + W_{xo}x^{(t)}) \\ g^{(t)} &= \tanh(W_{hg}h^{(t-1)} + W_{xg}x^{(t)}). \end{aligned}$$

Then, we compute the outputs that are propagated to the next layer,

$$\begin{aligned} c^{(t)} &= f^{(t)} \odot c^{(t-1)} + i^{(t)} \odot g^{(t)} \\ h^{(t)} &= o^{(t)} \odot \tanh(c^{(t)}). \end{aligned}$$

The addition in the computation of $c^{(t)}$ allows for gradients to directly propagate through $c^{(t-1)} \odot f$. Also, it allows the model to “select” what information should be retained. For example, at a high level in text, it might be helpful to store information such as gender and countries of origin. See Colah’s “Understanding LSTM Networks” for more information about the gates.

4.3 Gradient clipping

While LSTMs are a great solution to the vanishing gradient problem, we still have the possibility of exploding gradients. This is what *gradient clipping* solves. The idea is to limit the maximum value of the gradient if it surpasses a predetermined threshold. In practice, the gradient descent step gets transformed into the following update rule,

$$\theta \leftarrow \begin{cases} \theta - \gamma g & \|g\| \leq T \\ \theta - \gamma T \frac{g}{\|g\|} & \text{otherwise,} \end{cases}$$

You can also chain multiple LSTM units one after another, which results in this computation being performed multiple times per layer. If this is the case, then we replace $x^{(t)}$ by the hidden vector of the previous unit for all units after the first.

where g is the gradient, γ is the learning rate, and T is the gradient threshold.

5 Autoencoders

Autoencoders are generative models. This means that their objective is to learn the underlying hidden structure of the data. They aim to model the distribution $p_{\text{model}}(\mathbf{x})$ that resembles $p_{\text{data}}(\mathbf{x})$ to generate new samples. Autoencoders are an *explicit* generative model, which means that they explicitly define the probability distribution $p_{\text{model}}(\mathbf{x})$ and then sample from it to generate new data points.

In machine learning, we often have high-dimensional data $\mathbf{x} \in \mathbb{R}^n$, such as images, audio, or time-series. Hence, it is crucial to find a low-dimensional representation that can effectively compress the data while preserving its essential information.

Autoencoders offer a solution by making use of the *encoder-decoder structure*; see Figure 14. The *encoder* f projects the input space \mathcal{X} into a latent space \mathcal{Z} , while the *decoder* g maps the latent space \mathcal{Z} back to the input space \mathcal{X} . The assumption made by the autoencoder architecture is that if the decoder is capable of reconstructing the original input solely from the compressed representation, then this compressed representation must be meaningful. Consequently, the composition $g \circ f$ aims to approximate the identity function on the data for a low reconstruction error.

Furthermore, to enable the generation of new samples from the latent space, the latent space must be well structured, characterized by *continuity* and *interpolation*. Continuity means that the entire space must be covered by the data points, while interpolation means that if we interpolate between two points, then the interpolation must also be a well behaved data point.

5.1 Linear autoencoders

If we restrict f and g to be linear, the encoder f becomes equivalent to the projection performed by *principal component analysis*. The advantage of such a reconstruction is that it can be found in a closed form. However, it is not very powerful.

5.2 Non-linear autoencoders

We can gain a lot of performance by allowing f and g to be non-linear. In this case, the encoder and decoder are implemented as neural networks. To train these networks, we optimize for the reconstruction error,

$$\phi^*, \psi^* \in \operatorname{argmin}_{\phi, \psi} \sum_{n=1}^N \|\mathbf{x}_n - g_{\psi}(f_{\phi}(\mathbf{x}_n))\|^2$$

We can distinguish between *undercomplete* and *overcomplete* latent spaces. A latent space is undercomplete if $\dim(\mathcal{Z}) < \dim(\mathcal{X})$, while it is over-

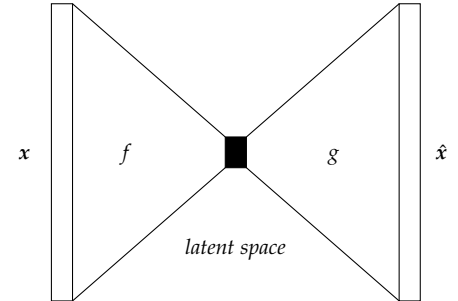


Figure 14. Autoencoder architecture.

complete if $\dim(\mathcal{Z}) > \dim(\mathcal{X})$. The idea of an undercomplete hidden representation is to enable the network to learn the important features of the data by reducing the dimensionality of the hidden space. This prevents the autoencoder from simply copying the input and forces it to extract meaningful and discriminative features. An overcomplete latent spaces are useful for denoising and inpainting autoencoders, where we have an imperfect input and want a perfect output. The overcompleteness allows the model to extract more features from the transformed input, leading to improved performance.

5.3 Variational autoencoders

While autoencoders are good at reconstruction, they struggle at generating new high quality samples, which is due to the lack of continuity in the latent space. There are large regions in the latent space where there are no observations, thus the model does not know what to output when it get an input from those regions.

Variational autoencoders (VAE) are designed to have a continuous latent space. It achieves this by making the encoder output a probability distribution over latent vectors, rather than a single latent vector. Generally, it outputs a mean vector μ and standard deviation vector σ . The idea is that even for the same input, the latent vector can be different, but in the same area. This means that data points cover areas in the latent space, rather than single points, ensuring continuity.

However, since there are no limits on the values taken by μ and σ , the encoder may learn to generate very different μ for each class while minimizing σ . This would mean that the encoder essentially outputs points again to decrease the reconstruction error. We can avoid this by minimizing the KL-divergence¹⁰ between the output distribution and a standard normal distribution. Intuitively, this encourages the encoder to distribute the encodings evenly around the center of the latent space.

To train the model, we want to maximize the likelihood of the training data,

$$p(x) = \int_{\mathcal{Z}} p(x | z) p(z) dz.$$

However, this is intractable. Thus, we define an approximation of the posterior, $q_{\phi}(z | x)$, which is computed by the encoder. We can now

¹⁰ The KL-divergence is defined as

$$D_{\text{KL}}(p||q) \doteq \mathbb{E} \left[\log \left(\frac{p(x)}{q(x)} \right) \right].$$

It is not symmetric and non-negative.

derive the *evidence lower bound* (ELBO),

$$\begin{aligned}
\log p(x) &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p(x)] && x \text{ does not depend on } z. \\
&= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\log \frac{p_\psi(x|z)p(z)}{p(z|x)} \right] && \text{Bayes' rule.} \\
&= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\log \left(\frac{p_\psi(x|z)p(z)}{p(z|x)} \frac{q_\phi(z|x)}{q_\phi(z|x)} \right) \right] && q(z|x)/q(z|x) = 1. \\
&= \mathbb{E}_{z|x} [\log p_\psi(x|z)] - \mathbb{E}_{z|x} \left[\log \frac{q_\phi(z|x)}{p(z)} \right] + \mathbb{E}_{z|x} \left[\log \frac{q_\phi(z|x)}{p(z|x)} \right] \\
&= \mathbb{E}_{z|x} [\log p_\psi(x|z)] - D_{\text{KL}}(q_\phi(z|x) \| p(z)) + D_{\text{KL}}(q_\phi(z|x) \| p(z|x)) \\
&\geq \mathbb{E}_{z|x} [\log p_\psi(x|z)] - D_{\text{KL}}(q_\phi(z|x) \| p(z)). && \text{KL-divergence is non-negative.}
\end{aligned}$$

The first term of the ELBO encourages low reconstruction error, while the second term makes sure that the approximate posterior q_ϕ does not deviate too far from the prior p . The second term can be computed in a closed-form, since both arguments are Gaussian.

A minor problem is that, during training, we cannot compute the derivative of expectations w.r.t. the parameters that we wish to optimize. Thus, we must use the *reparametrization trick*, which involves treating the random sampling as a single noise term. In particular, instead of sampling $z \sim \mathcal{N}(\mu, \text{diag}(\sigma))$, we sample $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$ and compute $z = \mu + \sigma \odot \epsilon$. Using this trick, we can remove the mean and variance from the sampling operation, meaning that we can differentiate w.r.t. the model parameters.

5.4 β -VAE

VAEs still have problems with their latent space; the representations are still *entangled*. This means that we do not have an explicit way of controlling the output. For example, in the MNIST dataset, we have no way of explicitly sampling a specific number. The β -VAE solves this problem by giving more weight to the KL term with an adjustable hyperparameter β that balances latent channel capacity and independence constraints with reconstruction accuracy. The intuition behind this is that if factors are in practice independent from each other, the model should benefit from disentangling them.

In practice, we want to force the KL loss to be under a threshold δ ,

$$\begin{aligned}
&\underset{\phi, \psi}{\text{maximize}} && \mathbb{E}_{x \sim \mathcal{D}} \left[\mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\psi(x|z)] \right] \\
&\text{subject to} && D_{\text{KL}}(q_\phi(z|x) \| p(z)) \leq \delta.
\end{aligned}$$

Rewriting this as a Lagrangian, we get

$$\begin{aligned}
\mathcal{L}(\phi, \psi, \beta) &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\psi(x|z)] - \beta(D_{\text{KL}}(q_\phi \| p(z)) - \delta) \\
&= \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\psi(x|z)] - \beta D_{\text{KL}}(q_\phi \| p(z)) + \beta \delta \\
&\geq \mathbb{E}_{z \sim q_\phi(\cdot|x)} [\log p_\psi(x|z)] - \beta D_{\text{KL}}(q_\phi \| p(z)).
\end{aligned}$$

Thus, this becomes our new objective function that we wish to maximize.

6 Autoregressive models

We saw that VAEs are approximate models, since they cannot exactly compute the likelihood $p(\mathbf{x})$ to maximize it. *Autoregressive models* solve this by computing the likelihood with the chain rule,

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i \mid \mathbf{x}_{1:i-1}).$$

Autoregressive models use data from the same input variable at previous timesteps. This is where they get their name from; they perform regression of themselves. In particular, autoregressive models generate one element of the sequence at a time, conditioning on previously generated elements. Thus, it takes $\mathbf{x}_{1:k}$ as input and outputs x_{k+1} . Because of this, they can build a probability distribution over possible sequences, using the chain rule as above.

The hard part of this approach is that we must parametrize all possible conditional distributions $p(x_{k+1} \mid \mathbf{x}_{1:k})$. Suppose we have a binary image consisting of n pixels. Then, we need

$$\sum_{i=1}^n 2^{i-1} \in \mathcal{O}(2^n)$$

parameters to parametrize this model. Thus, we need to make additional assumptions to find more compact representations of the distribution.¹¹ We will explore the idea of learning a function $f_i : \{0,1\}^{i-1} \rightarrow [0,1]$, parametrized by θ_i , which takes as input the previous pixels and outputs the probability $p(x_i = 1 \mid \mathbf{x}_{1:i-1})$. The number of total parameters is

$$\sum_{i=1}^n |\theta_i|.$$

Furthermore, we need to think about in which order we generate the pixels of the image. Intuitively, it would make sense to do left-right top-down or top-down left-right. However, in practice, randomly ordering the pixels works just as well. One just needs to make sure that the ordering remains the same for all data points.

6.1 Fully visible sigmoid belief network

In a *fully visible sigmoid belief network* (FVSBN) [Frey, 1998], each timestep has its own function f_i that is modeled by logistic regression,

$$f_i(\mathbf{x}_{1:i-1}) = \sigma\left(\alpha_0^{(i)} + \alpha_1^{(i)}x_1 + \cdots + \alpha_{i-1}^{(i)}x_{i-1}\right).$$

At the i -th timestep, we have i parameters denoted by $\theta_i = [\alpha_0, \dots, \alpha_{i-1}]$. Thus, the total number of parameters is

$$\sum_{i=1}^n |\theta_i| = \sum_{i=1}^n i = \frac{n^2 + n}{2} \in \mathcal{O}(n^2),$$

¹¹ A naive solution would be to assume that all points are independent, which would result in

$$p(\mathbf{x}) = p(x_1) \cdot p(x_n).$$

Then, we only require n parameters. However, in practice, this would result in a random sampling of pixels, making the generations incoherent.

which is much better than the exponential number of parameters we had before.

6.2 Neural autoregressive density estimator

The problem with FVSBN is that they are likely not expressive enough for any meaningful tasks, since they only consist of a single linear layer for each timestep. The *neural autoregressive density estimator* (NADE) [Uribe et al., 2016] offers an alternative parametrization based on MLPs, where we have hidden layers, increasing expressivity, and the weights are shared between timesteps, decreasing the number of parameters. Specifically, the hidden layer activations can be computed by the following,

$$\begin{aligned} \mathbf{h}_i &= \sigma(\mathbf{W}_{:,1:i-1} \mathbf{x}_{1:i-1} + \mathbf{b}) \\ \hat{x}_i &= \sigma(\mathbf{V}_{i,:} \mathbf{h}_i + c_i). \end{aligned}$$

The advantage of shared parameters is that the total number of parameters gets reduced from $\mathcal{O}(n^2d)$ to $\mathcal{O}(nd)$, and the hidden unit activations can be evaluated in $\mathcal{O}(nd)$ by using an alternative recursive definition,

$$\begin{aligned} \mathbf{a}_0 &= \mathbf{0} \\ \mathbf{a}_{i+1} &= \mathbf{a}_i + \mathbf{W}_{:,i} x_i \\ \mathbf{h}_i &= \sigma(\mathbf{a}_i + \mathbf{b}). \end{aligned}$$

Since NADE is a model for binary data, $\hat{x} \in [0, 1]$ is the probability $p(x_i \mid \mathbf{x}_{1:i-1})$ at each timestep i . NADE is trained by maximizing the average log-likelihood,

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \log p(\mathbf{x}^t) &= \frac{1}{T} \sum_{t=1}^T \log \prod_{i=1}^n p(x_i^{(t)} \mid \mathbf{x}_{1:i-1}^{(t)}) \\ &= \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^n \log p(x_i^{(t)} \mid \mathbf{x}_{1:i-1}^{(t)}). \end{aligned}$$

Thus, we optimize the log-likelihood exactly, which was not possible for VAEs.

During training, a *teacher forcing* approach is used, where the ground truth values of pixels are used for conditioning, rather than the predicted ones. This leads to more stable training. However, at inference time we use the predicted value, since we do not have access to the ground truth.

6.3 Masked autoencoder distribution estimation

The idea behind *masked autoencoder distribution estimation* (MADE) [Germain et al., 2015] is to construct an autoencoder which fulfills the autoregressive property, such that its outputs can be used as conditionals. In order to achieve this, no computational path between output unit \hat{x}_{k+1} and any of the input units x_{k+1}, \dots, x_n may exist. This has the result that \hat{x}_{k+1} is only conditioned on $\mathbf{x}_{1:k}$, relative to the ordering.

The way to do this is to first assign a number from 1 to n to each input unit x_i , which we call the ordering. Then, for each hidden unit, uniformly

sample an integer m , between 1 and $n - 1$. In the hidden layers, we then only allow units in layer ℓ to propagate to units in layer $\ell + 1$ with higher or equal value. Finally, we allow connections between the last hidden layer and the output only to units with value that is strictly greater. Let $m^{(\ell)}(k)$ be the value assigned to the k -th element of layer ℓ , then these constraints can be encoded in mask matrices,

$$\begin{aligned} \mathbf{M}_{ij}^{\mathbf{W}^{(\ell)}} &= \mathbb{1}\{m^{(\ell-1)}(j) \leq m^{(\ell)}(i)\} \\ \mathbf{M}_{ij}^{\mathbf{V}} &= \mathbb{1}\{m^{(\ell-1)}(j) < m^{(\ell)}(i)\}. \end{aligned}$$

Then, we alter the weight matrices by

$$\begin{aligned} \bar{\mathbf{W}}^{(\ell)} &= \mathbf{W}^{(\ell)} \odot \mathbf{M}^{\mathbf{W}^{(\ell)}} \\ \bar{\mathbf{V}} &= \mathbf{V} \odot \mathbf{M}^{\mathbf{V}}, \end{aligned}$$

and use those instead.

However, a problem with this approach is that it requires very large hidden layers to retain expressivity.

6.4 Generating images

Pixel-RNN. The idea behind *Pixel-RNN* [Van Den Oord et al., 2016] is to generate image pixels starting from the corner and modeling the dependency on previous pixels using an RNN. In particular, a pixel value is dependent on its top-left neighboring pixels and the RNNs hidden state; see Figure 16. However, the problem with this approach is that the generation of new pixels depends on the hidden state, making the generation process sequential, and thus slow.

Pixel-CNN. We can solve the problem of *Pixel-RNN* by assuming that pixel values only depend on a context region around the pixel. This is exactly what the *Pixel-CNN* does [Van Den Oord et al., 2016]. This allows for parallelization during training, because the context region values are known during training.¹² Just as *Pixel-RNN*, it starts from the top-left corner, but it models the dependencies with a CNN. During training, we need to make sure that only the previously generated pixels are used for prediction of the next, thus we need a mask that masks out all unknown pixel values. However, stacking layers of masked convolutions creates a blind spot in the convolution; see Figure 17. The solution to this is to combine horizontal and vertical stacks of convolutions [Van den Oord et al., 2016]. The former conditions on the row so far, while the latter conditions on all rows above. The final output is obtained by summing the two outputs.

To enforce the autoregressive property, the model also needs to go over the color channels in an autoregressive manner. Thus the conditional probability is expressed as

$$p(\mathbf{x}_i | \mathbf{x}_{1:i-1}) = p(x_{i,R} | \mathbf{x}_{1:i-1})p(x_{i,G} | \mathbf{x}_{1:i-1}, x_{i,R})p(x_{i,B} | \mathbf{x}_{1:i-1}, x_{i,R}, x_{i,G}).$$

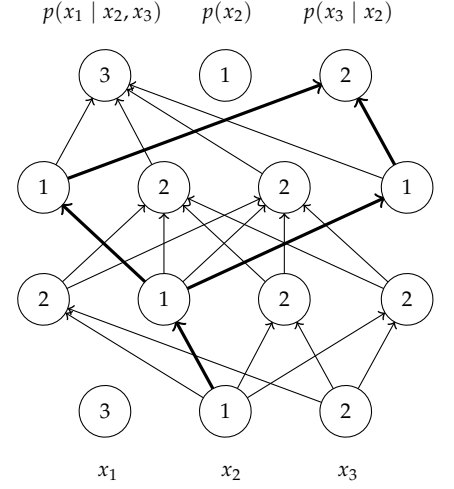


Figure 15. MADE masking with $n = 3$.

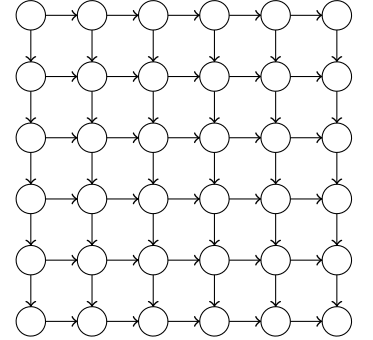


Figure 16. Pixel-RNN generation process.

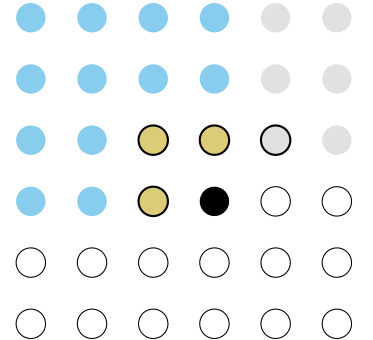


Figure 17. Pixel-CNN generation process. The black pixel depends explicitly on the yellow pixels, where the thick-lined pixels denote the masked convolutional layer. The black pixel should also depend on the gray pixels, but it does not due to the way the stacked masked convolutions work; a blind spot.

¹² However, we still have to sequentially predict every token during inference, but this is an issue with all autoregressive models.

6.5 Generating audio

WaveNet [van den Oord et al., 2016] adapts the Pixel-CNN framework for audio data. However, audio typically has much longer time horizons, since they consist of 16000 samples per second. To be able to capture these long-term dependencies efficiently, *WaveNet* incorporates *dilated convolutions* [Yu and Koltun, 2016]. This allows for an exponential increase in the receptive field. In dilated convolutions, the filter is applied with gaps between the filter elements. *WaveNet* increases the dilation factor as we go up in the layers to attain an exponential receptive field; see Figure 18.

6.6 Variational RNN

RNNs can also be used to generate sequences by sampling $\mathbf{h}^{(0)}$ and then sequentially predicting the next element, and updating the hidden state. However, the generation of new sequences is very slow, because of the sequential nature of the generation process. Another limitation of vanilla RNNs is that their structure is entirely deterministic and thus limited in expressive power.

The *variational RNN* [Chung et al., 2015] (VRNN) is a recurrent version of the VAE. A different VAE is used at each timestep. However, the VAEs are not the same as we have seen. In particular, the prior distribution is no longer a standard Gaussian, but follows the distribution $\mathcal{N}(\boldsymbol{\mu}_0^{(t)}, \text{diag}(\boldsymbol{\sigma}_0^{(t)}))$, where $(\boldsymbol{\mu}_0^{(t)}, \boldsymbol{\sigma}_0^{(t)})$ is predicted by the network $\phi_{\text{prior}}(\mathbf{h}^{(t)})$ based on the hidden state of the RNN at each timestep. This introduces randomness to the generation.

The *conditional VRNN* (C-VRNN) is an extension to the VRNN, where we can condition on *e.g.* style or, in the case of MNIST, which digit should be generated. During training, the model predicts the conditional variable, called the posterior, and hidden state from the input $\mathbf{x}^{(t)}$. From the hidden state it also predicts the conditional variable, called the prior. It then minimizes the reconstruction error from the style variable, and the KL divergence between prior and posterior. During inference, it only uses the priors for next token prediction, predicted from the hidden state, since it does not have access to $\mathbf{x}^{(t)}$. Then, given the predicted token, it predicts the posterior conditional variables, which are used to predict the next hidden state $\mathbf{h}^{(t+1)}$.

6.7 Transformers

Transformers [Vaswani et al., 2017] are used in nearly all the state-of-the-art models at the moment. At the basis of transformers lies the self-attention mechanism, where we extract key, value, and query representations from the input. These are then used in a soft-lookup mechanism,

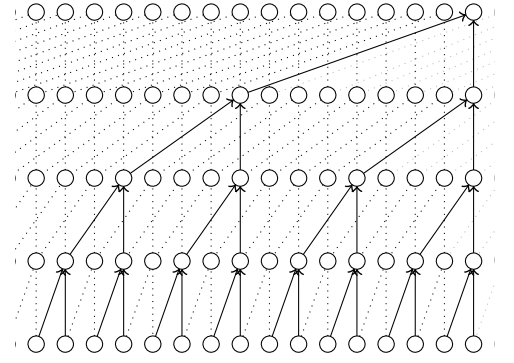


Figure 18. Stacked dilated convolutional layers in *WaveNet*.

where the query and key values decide how much attention is paid to the values.¹³

It computes the keys, values, and queries by

$$\mathbf{K} = \mathbf{X}\mathbf{W}_K$$

$$\mathbf{V} = \mathbf{X}\mathbf{W}_V$$

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q,$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{W}_K, \mathbf{W}_V, \mathbf{W}_Q \in \mathbb{R}^{d \times d}$. Now, we have key, value, and query representations of the input that are all in $\mathbb{R}^{n \times d}$. Then, we can compute the output by

$$\mathbf{Y} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right)\mathbf{V}.$$

Intuitively, the softmax computes how much “attention” should be given to each of the previous timesteps by outputting a probability distribution over timesteps.

Since $\mathbf{Y} \in \mathbb{R}^{n \times d}$ is of the same dimensionality as \mathbf{X} , we can stack self-attention layers. This is the basis of the transformer architecture. However, this does not respect the autoregressive property, thus we need to use a mask \mathbf{M} to prevent the model from accessing future timesteps, where

$$\mathbf{M} = \begin{bmatrix} -\infty & -\infty & \cdots & -\infty \\ 0 & -\infty & \cdots & -\infty \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & -\infty \end{bmatrix}.$$

Then, the attention computation becomes

$$\mathbf{Y} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} + \mathbf{M}\right)\mathbf{V}.$$

The computational complexity of self-attention is $\mathcal{O}(n^2d)$, which is quite high for large inputs, such as audio and images. But, we gain a maximum path length of $\mathcal{O}(1)$ to any of the previous timesteps, making sure that no information gets “forgotten”, and it allows for easy error propagation during training, in contrast to RNNs.

¹³ Note the parallel with dictionaries in programming languages, which use a hard-lookup.

References

- Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. A recurrent latent variable model for sequential data. *Advances in neural information processing systems*, 28, 2015.
- Brendan J Frey. *Graphical models for machine learning and digital communication*. MIT press, 1998.
- Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. Made: Masked autoencoder for distribution estimation. In *International conference on machine learning*, pages 881–889. PMLR, 2015.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
- Benigno Uribe, Marc-Alexandre Côté, Karol Gregor, Iain Murray, and Hugo Larochelle. Neural autoregressive distribution estimation. *Journal of Machine Learning Research*, 17(205):1–37, 2016.
- Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio, 2016.
- Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional image generation with pixelcnn decoders. *Advances in neural information processing systems*, 29, 2016.
- Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *International conference on machine learning*, pages 1747–1756. PMLR, 2016.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions, 2016.