

# 1 Methodological framework

This Section will outline the theoretical background of convolutions from a deep learning point of view, as well as provide a brief overview of image segmentation as a computer vision task. Furthermore, the Methodological framework required to develop the U-Net architecture will be highlighted, followed by the description of a training loop and the central algorithms utilized. Finally, validation metrics used to assess the performance of the developed deep learning system will be derived.

## 1.1 Convolutional layers

Convolutional layers incorporated into a deep neural network which utilized the backpropagation algorithm (Rumelhart et al., 1986) was initially proposed by (LeCun et al., 1989) to classify handwritten numbers. The layer LeCun et al. (1989) presented consists of an arbitrary amount of filters, which are small two dimensional matrices (e.g.  $(3 \times 3)$  pixels) designed to capture a certain structure in the image such as lines or edges . Each filter contains trainable weights, which are learned from the data during backpropagation (LeCun et al., 1989) and gradient descent. When a filter is convolved with all possible local neighborhoods from the input, it outputs a feature map which represent where the input image triggered a response from the filter (Zeiler et al., 2010). Moreover, inputting feature maps to a convolutional layer allows for the filters to respond to combinations of lower level structures, which trains the layer to detect more complicated patterns (Fukushima, 1980). Additionally, stacking convolutional layers in a network-architecture structure increases the field of view for each subsequent layer, which makes each layer observe an increasingly complex pattern of higher order feature maps at increasingly larger spatial scales (Fukushima, 1980). As a result, convolutional layers are able to discern between object and background as they perceive only a limited view of the scene. The convolutional layer is also invariant to translation of the object, since the filter is constant when creating the feature map (LeCun et al., 1989). I.e. the filter is detecting the same feature at all locations in the image, known as weight sharing (LeCun et al., 1989).

The number of trainable parameters for a convolutional layer is equal to the size of a filter times the number of filters. As a result, the number of trainable parameters is invariant to the spatial extent of the input images. On the other hand, fitting a fully connected layer to spatial gridded data consists of associating a separate trainable parameter to each pixel. As such, the size of a fully connected layer scales with the size of the image, which increases the risk of overfitting the network. In the case of the convolutional layer, LeCun et al. (1989) notes that reducing the number of trainable parameters through weight sharing constrains the solution space such that overfitting is avoided while still having enough trainable parameters to fit the layer to the data. Furthermore, the fully

Sitere  
boka  
til  
Good-  
fel-  
low2016?

connected layer is not invariant to translation as each trainable parameter is exclusive to their respective pixel, hence no weight sharing. As such, the layer is unable to detect a similar object at a different position, reducing their usefulness for image-based prediction tasks.

Finally, Ciresan et al. (2012b) showed that the processing time of a convolutional layer is significantly shortened by utilizing a graphics processing unit (GPU), due to their large amount of compute cores compared to traditional Central Processing Units (CPUs). Furthermore, the authors of Krizhevsky et al. (2012) provided the first publicly available implementation of a CNN running on a GPU by utilizing the Nvidia Compute Unified Device Architecture api. Krizhevsky et al. (2012) also demonstrated that their results are tied to the performance of the GPU in terms of available memory as well as floating point operations, with the implications that a better GPU as well as larger datasets would improve their results. As such, the modern implementation of convolutional layers, and by extension deep learning architectures which stacks many convolutional layers, can be initiated with more trainable parameters since they can process greater datasets consisting of larger samples due to their GPU implementation.

The convolutional layer can be described mathematically by utilizing the previously described principle of allowing the filter to only perceive a local neighborhood of the input. Consider the value of a single point  $y_{i,j} \in Y \in \mathbb{R}^2$  where  $i, j$  denote the position in the x and y direction as a single output from a convolution. Let  $X \in \mathbb{R}^3$  be an input image of size  $(A \times B \times D)$  consisting of a single channel, and  $W \in \mathbb{R}^3$  be a symmetric filter of size  $(r \times r \times D)$ . Then, the value at a single point  $y_{i,j}$  is given as follows,

$$y_{i,j} = \sum_{a=1-\frac{r}{2}}^{\frac{r}{2}} \sum_{b=1-\frac{r}{2}}^{\frac{r}{2}} \sum_{d=1}^D W_{a+\frac{r}{2}, b+\frac{r}{2}, d} X_{i+a, j+b, d} \quad (1)$$

Where the subscript notation is used in  $W$  and  $X$  to denote indexes similarly to  $Y$ . Repeating Equation (1) across all points  $x \in X$  by applying a sliding window technique returns the convolution of  $X$  with filter  $W$ , which results in an output  $Y$  with size  $(A - r + 1) \times (B - r + 1)$ . Note that the above definition only applies for  $X_{1 \leq i+a \leq A, 1 \leq j+b \leq B}$ . The size of the output can be adjusted by padding the input  $X$  by a size  $P$  in each direction or increasing the stride  $S$  of the sliding window, which reformulates the output size of  $Y$  in a single dimension as a function

$$Y_{\text{dim}} = \lfloor \frac{A - r + 2P}{S} + 1 \rfloor \quad (2)$$

The convolutional layer adds the convolution described in Equation (1) with a bias term  $B \in \mathbb{R}^2$  of same spatial shape as  $Y$ , as well as applying an activation function  $g$  to each

$y_{i,j}$  which introduces nonlinearity. In summary, the output of a convolutional layer can be described as

$$Y' = g(Y + B) = g(W^T X + B) \quad (3)$$

If the number of filters increases from 1 to  $N$ , Equation (3) is repeated for all filters, resulting in an output  $Y \in \mathbb{R}^3$  of size  $(Y_{\text{dim1}}, Y_{\text{dim2}}, N)$ .

## 1.2 Image segmentation

Image segmentation is a computer vision task where pixels are assigned labels according to some predetermined rules. It is common to define an image segmentation task either as a study of countable *things* (Instance segmentation), or recognizing similarly textured *stuff* (Semantic segmentation) (Kirillov et al., 2018). The task for this thesis, which is labelling sea ice concentration according to its predicted concentration class, falls into the latter category following the definition of *stuff* in Adelson (2001). I.e. the current task is to assign each pixel in a predicted scene a single class label.

Network architectures based on the Convolutional Neural Network (CNN) (LeCun et al., 1989; Ciresan et al., 2012b; Krizhevsky et al., 2012; Simonyan and Zisserman, 2014; Szegedy et al., 2014; He et al., 2015b; Huang et al., 2016) can be used to perform pixelwise semantic segmentation, however the CNN architectures listed have been developed for image classification i.e. predicting a single label for the entire image. Ciresan et al. (2012a) presented an approach where a CNN (see the architecture of Ciresan et al. (2012b)) was used to predict a label for all pixels in an image. Instead of processing the entire image at once, Ciresan et al. (2012a) applied a sliding window technique which predicted each pixel by using their surrounding neighborhood as input. However, due to only processing parts of the image at once, the segmentation algorithm in Ciresan et al. (2012a) is computationally expensive as the CNN must be run for all neighborhoods as well as limiting the context in which each pixel is classified in to a surrounding local context (Ronneberger et al., 2015).

To capture the global context of a scene, network architectures such as Long et al. (2015); Noh et al. (2015); Ronneberger et al. (2015); Badrinarayanan et al. (2017); Chen et al. (2018) implement the Encoder-Decoder architecture, where the entire input scene is first processed by a CNN-like architecture referred to as the Encoder to produce a signal. The signal is then used as input to a subsequent network which reconstructs the encoded signal to match the resolution of the original image through upsampling. Long et al. (2015); Ronneberger et al. (2015); Badrinarayanan et al. (2017) all apply the Deconvolution architecture proposed by Zeiler et al. (2010) to upsample the encoded signal through the use

of a trainable deconvolutional layer, which will be described in greater detail in Subsection (1.3.3). However, other upsampling techniques exist, such as unpooling used in Noh et al. (2015) which undoes a maxpool (downsampling) layer (maxpooling is described in Subsection (1.3.2)).

This thesis will utilize the U-Net architecture proposed by Ronneberger et al. (2015). The U-Net was originally developed for medical image segmentation, however the architecture has shown promising results for both pan-arctic seasonal (Andersson et al., 2021) and regional short term (Grigoryev et al., 2022) sea ice concentration forecasting amongst other applications. Another aspect which makes the U-Net more suitable to the current task, compared to other previously mentioned image-to-image architectures is that the network converges quickly, which is ideal when working with a dataset consisting of few samples (Ronneberger et al., 2015).

### 1.3 Describing the U-Net architecture

Figure (1) shows the U-Net architecture. This section intends to describe the different components constituting the architecture from a technical point of view.

#### 1.3.1 The convolutional block

A single convolutional block consists of two repeat convolutional layers, each followed by the Rectified Linear Unit (ReLU) nonlinear activation function, which is defined as  $f(x) = \max(0, x)$  (Nair and Hinton, 2010). Each convolution is performed using a  $3 \times 3$  window. The original formulation of the U-Net also does not apply padding to the input, resulting the convolutional filter only being applied to the entries of the input where the filter is never out of bounds. With a stride of  $S = 1$ , this results in each convolutional layer reducing the spatial extent by two pixels in each direction following Equation (2). It is also noted that the number of feature maps is doubled after each downsampling step, which is performed by the pooling layers.

#### 1.3.2 Maxpooling

Pooling operations are used to reduce the spatial extent of the current feature maps, by downsampling the data in the spatial dimensions. As seen in Figure (1), the U-Net downsamples the data in the contracting path through  $2 \times 2$  maximum pool layers with a stride of 2. This specific configuration causes the spatial resolution to be halved. In the max-pool layer, a filter runs through each input channel and chooses the maximum value inside the neighborhood of the filter. As such, the extreme values in each feature

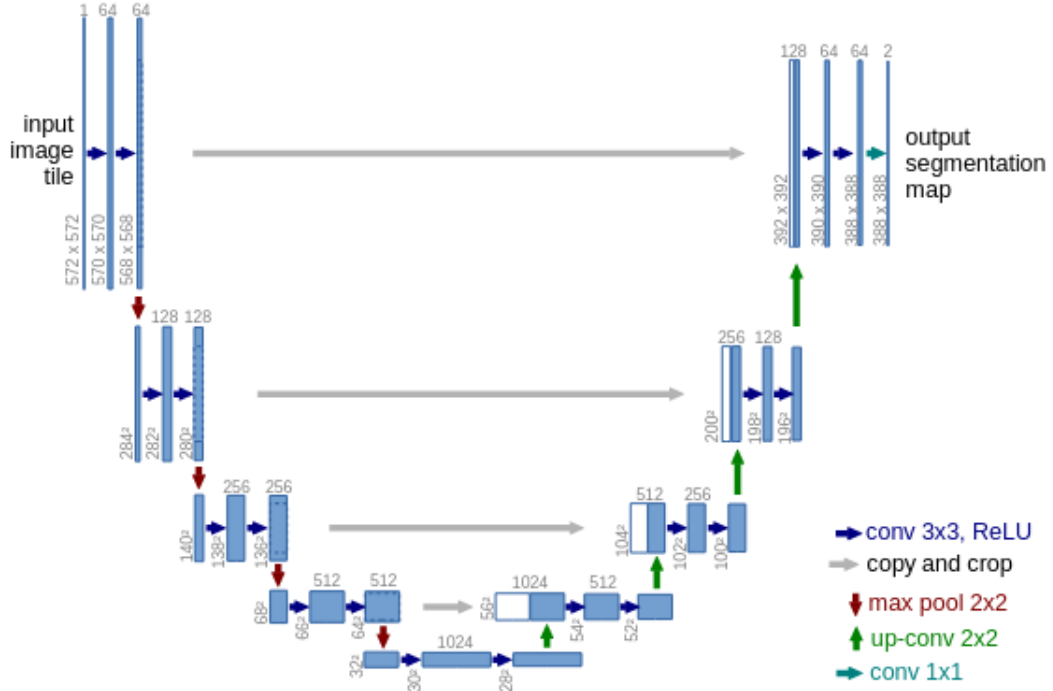


Figure 1: The U-Net architecture. The blue boxes represent feature maps, with the lower left numbers determining the spatial resolution and the top number the amount of feature maps. White boxes in the expansive path (right side / decoder) are the copied feature maps from the contractive path (left side / encoder). Arrows denote the different operations. Note that the original U-Net only performs *valid* convolutions, i.e. convolution without padding to match the input. This causes a convolutional layer to slightly decrease the spatial extent. As a result, the copied features from the contracting path are also cropped to match the dimensionality in the expansive path. Figure extracted from Ronneberger et al. (2015).

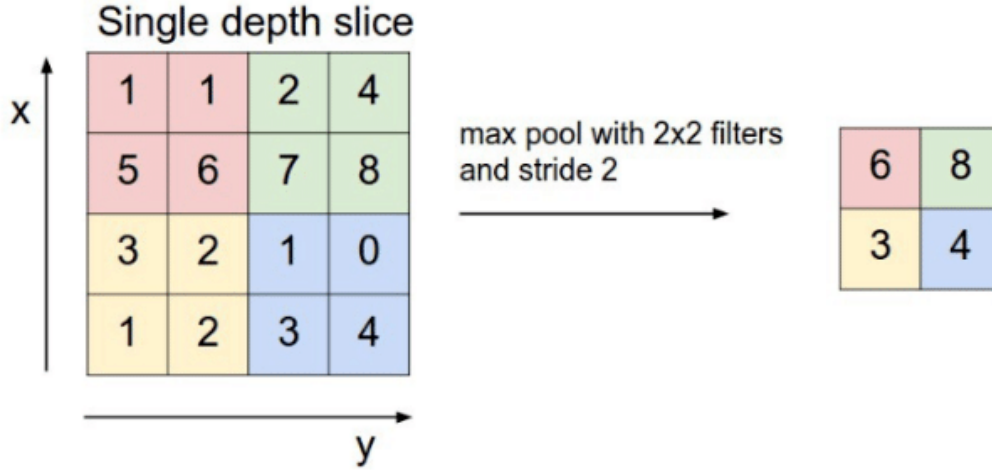


Figure 2: The max-pool operation for a  $2 \times 2$  filter with a stride of 2. Figure taken from (RADU et al., 2020)

map is retained at the expense of rejecting the rest of the data. Since the maxpooling operation is rejecting some parts of the data, it may be regarded as a regularizer for the network which aid in avoiding overfitting. See Figure (2) for a graphical description.

### 1.3.3 Transposed convolutions

Transposed convolution was proposed by Zeiler et al. (2010) (note the incorrect use of deconvolution, this is not the mathematical inverse of a convolution) to increase the resolution of a feature map. The method was first utilized by Long et al. (2015) to connect the coarse output of an encoder with the image resolution of the target (it is referred to as both *backwards convolution* and *deconvolution* in the proceedings paper). Similar to the convolutional layer, the transposed convolutional layer involves striding a convolutional filter with trainable parameters across a feature map. However, the transposed convolutional layer projects a singular entry from the input through the convolutional kernel to produce an output that is larger than the input.

In the Encoder architecture, lower level feature maps provide spatial information regarding where stuff is located in a scene, whereas higher level feature maps contain information regarding what is in the scene at the expense of losing spatial information (Long et al., 2015). To circumvent this, Ronneberger et al. (2015) concatenate the features from the contracting path with the output from the transposed convolution at the same level of depth, i.e. where the number of feature maps are equal at the end of the convolutional block. The concatenation operation is possible in Ronneberger et al. (2015) since they

crop the feature maps in the encoder in their spatial dimensions to match the spatial dimensionality of the feature maps in the decode. The operation can be seen in Figure(1) denoted by the gray arrow. The resulting convolutional layer is then trained to make a more precise prediction due to the concatenated input (Ronneberger et al., 2015).

### 1.3.4 Outputs

The output layer of the U-Net is denoted by the turquoise arrow at the right side of Figure (1). The arrow denote that the input processed by a convolutional layer with a number of filters equal to the number of output classes. Each filter is of size  $(1 \times 1)$  and maps each layer in the input feature map to their respective class probability map of equal spatial shape (Ronneberger et al., 2015). By inspecting Figure (1), the U-Net outputs two feature maps, and from each feature map the pixelwise probability of belonging to the associated class can be computed.

## 1.4 Training procedure for the U-Net

This subsection aims to demonstrate how Ronneberger et al. (2015) trained the U-Net, and will consequently highlight some different hyperparameters and exemplify some functions and operations which are used performed training. Hyperparameters refer to model parameters which are not updated during training (Yu and Zhu, 2020), and may directly influence the model architecture or the training procedure. This section will not describe how samples were preprocessed and loaded, and modifications made which reflects concerns regarding medical images may be noted but not explained.

Training the U-Net start by assigning random values to the weights of the network. Since the U-Net utilizes the ReLU activation function after each convolutional layer in the convolutional blocks (Ronneberger et al., 2015), it is standard for each layer to draw the weights from a normal distribution with  $\mu = 0$  and standard deviation  $\sigma = \sqrt{\frac{2}{n_l}}$ , where  $n_l$  is the number of inputs to the layer He et al. (2015a). This weight initialization scheme ensures that variance of the feature maps are approximately equal, i.e. avoids varying the activation of input signals between layers He et al. (2015a); Ronneberger et al. (2015).

The process of training the U-Net involves making predictions on all training data. For each sample, what is predicted is compared against a ground truth label associated with the sample. For the U-Net, a pixelwise prediction map is created by computing the pixelwise softmax which is an extension of the softmax function (Bridle, 1990) defined as

$$p_{k,i,j} = \frac{e^{a_{k,i,j}}}{\sum_{k'=1}^K e^{a_{k',i,j}}} \quad (4)$$

where  $a_k$  is the feature map for feature channel  $k$  of input  $x$  and  $K$  is the number of output classes.  $i, j$  is the spatial coordinates. Similarly to the standard softmax function (Bridle, 1990), Equation (4) is  $\approx 1$  for the class that has maximum  $a_{k,i,j}$  and  $\approx 0$  for all other classes, albeit depthwise in the channel dimension for all pixels (Ronneberger et al., 2015).

To quantify the prediction error, a loss function is defined. The overall goal of training a neural network is to minimize the loss function with respect to the trainable weights. For the U-Net, a weighted variation of the cross entropy loss function is proposed (Ronneberger et al., 2015).

$$L(p) = \sum_{i,j \in \mathbb{Z}^2} w_{i,j} \log(p_{l,i,j}) \quad (5)$$

where  $w$  is a predefined weight map and  $p_{l,i,j}$  is the prediction made at pixel  $i, j$  at the true label  $l$ .

The error computed by the loss function is then sent backwards throughout the network according to the backpropagation algorithm (Rumelhart et al., 1986), which effectively computes the gradient of the loss function with regards to the trainable parameters

$$\frac{\partial L}{\partial w_l} = \frac{\partial L}{\partial p} \frac{\partial p}{\partial w_l} \quad (6)$$

where  $w_l$  is the trainable parameters associated with the  $l$ -th layer. The gradient of the loss for a weight at a given layer shown in Equation (6) is used by an optimizer to adjust the weights such that the loss is minimized with respect to the weights (gradient descent).

Ronneberger et al. (2015) uses the stochastic gradient descent with momentum optimizer implemented in the Machine Learning library Caffe (Jia et al., 2014), where the optimizer is defined as follows,

$$w_l^{t+1} = \gamma(w_l^t - w_l^{t-1}) - \mu \frac{\partial L}{\partial w_l^t} \quad (7)$$

In Equation (7), the superscript  $t$  was added to  $w$  and refers to training step, which is defined as a prediction and subsequent backpropagation of a batch of samples, where the size of a batch is a pre-determined hyperparameter.  $\gamma$  and  $\mu$  are the momentum



and learning rate hyperparameters respectively. Note that  $\gamma$  is introduced by momentum stochastic gradient descent, whereas the learning rate  $\mu$  is a hyperparameter common for all deep learning models and determines the rate of weight adjustment as seen in Equation (7).

When all training samples have been inspected once by the U-Net, the training data is shuffled and the above outlined training procedure is repeated. The process of going through all the training data once is known as an epoch, which is a hyperparameter which determines the duration of training. Too low of an epoch can result in an underfitted model, although too many epochs may cause the model to overfit to the data.

## 1.5 Forecast verification metrics

Verification schemes provide insight into how a forecasting system performs. For this thesis, verification metrics serve a dual purpose. From a model development point of view, verification metrics will be used to increase the skill of the model. However, the same metrics will also be utilized to assess the quality of a prediction as well as explain the physical interpretation of the model (Casati et al., 2008). The model developed for this thesis predicts a scene consisting of labelled pixels, as described in Section (1.2). It was mentioned in Section (1) that the developed model is aimed towards operational end users, which is partly achieved by validating the model against metrics of end user relevance. Furthermore, it can be assumed that the model and target observation won't differ much outside of the marginal ice zone (Fritzner et al., 2020). Thus, this section will introduce metrics which are aware of the sea ice edge, as the sea ice edge is a relevant quantity for maritime end users operating in the Arctic (Melsom et al., 2019) as well as providing skill scores stemming from where the change in sea ice concentration is occurring. The following subsections will describe how to determine the position of the sea ice edge, as well as its length according to Melsom et al. (2019), and derive the Integrated Ice Edge Error (Goessling et al., 2016), with regards to a spatially gridded dataset of deterministic sea ice concentration values.

The Integrated Ice Edge Error is chosen among similar sea ice edge metrics (Melsom et al., 2019; Dukhovskoy et al., 2015) as it has been shown to be less sensitive to isolated ice patches (Palerm et al., 2019). Furthermore, the work of Melsom et al. (2019) recommend the Integrated Ice Edge Error amongst other metrics for its intuitive interpretation as well as for the possibility to provide the spatial distribution of IIEE areas.

### 1.5.1 Defining the Ice Edge

The sea ice edge for a given spatial distribution of sea ice concentration values is derived on a per pixel basis, and defined as the grid cells which meets the following condition,

$$c_{i,j} \geq c_e \wedge \min(c_{i-1,j}, c_{i+1,j}, c_{i,j-1}, c_{i,j+1}) < c_e \quad (8)$$

In Condition (8),  $c \subset C$  are gridded sea ice concentration values, with  $i, j$  denoting indexes.  $c_e$  is a given concentration threshold. The entries  $c_{i,j}$  which adhere to the condition in Condition (8) form the set  $E$  where  $\exists e_{i,j}, \forall c_{i,j}$  where Condition (8) holds (Melsom et al., 2019).

Moreover, all the entries in  $E$  each contribute to the total length of the sea ice edge, with each entries' length contribution determined based on that entries' 4-connected neighborhood. Using this formulation, the different combination of entry neighborhood in  $E$  can result in three different length contributions. For the following contributions,  $s$  is the spatial resolution of the grid.

- A neighborless pixel is assumed to yield a contribution the length of the diagonal of a grid cell ( $l = \sqrt{2}s$ ). Here it is assumed that the grid cell only have diagonal neighbors.
- A pixel with two or more 4-neighbors contributes with its spatial resolution (length of the grid cell)  $l = s$ .
- A pixel with one 4-neighbor contributes the mean value between the length of the grid cell and length of the diagonal of the grid cell  $l = \frac{s+\sqrt{2}s}{2}$ . It is assumed that the grid cell also have a diagonal neighbor.

The final length of the sea ice edge length then becomes

$$L = \sum_{e \text{ in } E} l^e \quad (9)$$

where the superscript  $l^e$  denotes the length associated with the entry  $e$  according to the algorithm listed above. I.e. the sum of all contributions.

### 1.5.2 Integrated Ice Edge Error

The IIEE is an error metric which compares a forecast  $f$  to a predefined ground truth target  $t$  Goessling et al. (2016). The metric is defined as

$$\text{IIEE} = \text{O} + \text{U} \quad (10)$$

where

$$\text{O} = \int_A \max(C_f - C_t, 0) dA \quad (11)$$

and

$$\text{U} = \int_A \max(C_t - C_f, 0) dA \quad (12)$$

with  $A \in \mathbb{R}^2$  being the area of interest, and is of similar size as  $C$ . Subscript  $f, t$  denotes whether  $C$  contains forecasted or target sea ice concentration values. In Equations (11 and 12),  $C$  is binary and is 1 if its concentration value is above some predefined threshold, and 0 elsewhere (Goessling et al., 2016). From the definition of the metric, it can be seen that the IIEE is a sum of the forecast overshoot and undershoot compared to the ground truth target.

Additionally, the IIEE can also be represented as a spatial metric by removing the integral with respect to  $A$  in Equation (11 and 12). In this way, the metric is used to define the set of pixels which constitutes its area. To clearly distinguish between the area  $\text{O}$  and the set of pixels used to compute  $\text{O}$ ,  $A^+$  will be used to note the latter. Similarly,  $A^-$  will represent the set of pixels constituting  $\text{U}$ . Finally, it can be seen that  $A^+$  and  $A^-$  represent the spatial distribution of False Positive and False Negatives of the forecast respectively.

The length of the ice edge has a strong influence on the IIEE (Goessling and Jung, 2018; Palerme et al., 2019). Hence, to ensure that forecast errors are comparable across seasons, IIEE is normalized with the length of the ice edge, as mentioned in Section (1.2.2)

## References

- Adelson, E. H.: On seeing stuff: the perception of materials by humans and machines, in: SPIE Proceedings, edited by Rogowitz, B. E. and Pappas, T. N., SPIE, <https://doi.org/10.1117/12.429489>, 2001.
- Andersson, T. R., Hosking, J. S., Pérez-Ortiz, M., Paige, B., Elliott, A., Russell, C., Law, S., Jones, D. C., Wilkinson, J., Phillips, T., Byrne, J., Tietsche, S., Sarojini, B. B., Blanchard-Wrigglesworth, E., Aksenov, Y., Downie, R., and Shuckburgh, E.: Seasonal

- Arctic sea ice forecasting with probabilistic deep learning, *Nature Communications*, 12, <https://doi.org/10.1038/s41467-021-25257-4>, 2021.
- Badrinarayanan, V., Kendall, A., and Cipolla, R.: SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39, 2481–2495, <https://doi.org/10.1109/tpami.2016.2644615>, 2017.
- Bridle, J. S.: Probabilistic Interpretation of Feedforward Classification Network Outputs, with Relationships to Statistical Pattern Recognition, in: *Neurocomputing*, pp. 227–236, Springer Berlin Heidelberg, [https://doi.org/10.1007/978-3-642-76153-9\\_28](https://doi.org/10.1007/978-3-642-76153-9_28), 1990.
- Casati, B., Wilson, L. J., Stephenson, D. B., Nurmi, P., Ghelli, A., Pocerich, M., Damrath, U., Ebert, E. E., Brown, B. G., and Mason, S.: Forecast verification: current status and future directions, *Meteorological Applications*, 15, 3–18, <https://doi.org/10.1002/met.52>, 2008.
- Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., and Yuille, A. L.: DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40, 834–848, <https://doi.org/10.1109/tpami.2017.2699184>, 2018.
- Ciresan, D., Giusti, A., Gambardella, L., and Schmidhuber, J.: Deep Neural Networks Segment Neuronal Membranes in Electron Microscopy Images, in: *Advances in Neural Information Processing Systems*, edited by Pereira, F., Burges, C., Bottou, L., and Weinberger, K., vol. 25, Curran Associates, Inc., URL <https://proceedings.neurips.cc/paper/2012/file/459a4ddcb586f24efd9395aa7662bc7c-Paper.pdf>, 2012a.
- Ciresan, D., Meier, U., and Schmidhuber, J.: Multi-column deep neural networks for image classification, in: *2012 IEEE Conference on Computer Vision and Pattern Recognition*, IEEE, <https://doi.org/10.1109/cvpr.2012.6248110>, 2012b.
- Dukhovskoy, D. S., Ubnoske, J., Blanchard-Wrigglesworth, E., Hiester, H. R., and Proshutinsky, A.: Skill metrics for evaluation and comparison of sea ice models, *Journal of Geophysical Research: Oceans*, 120, 5910–5931, <https://doi.org/10.1002/2015jc010989>, 2015.
- Fritzner, S., Graversen, R., and Christensen, K. H.: Assessment of High-Resolution Dynamical and Machine Learning Models for Prediction of Sea Ice Concentration in a Regional Application, 125, <https://doi.org/10.1029/2020jc016277>, neural Networks for predicting Sea-Ice concentration are only slightly more accurate than persistence forecasting for short-term predictions., 2020.
- Fukushima, K.: Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position, *Biological Cybernetics*, 36, 193–202, <https://doi.org/10.1007/bf00344251>, 1980.
- Goessling, H. F. and Jung, T.: A probabilistic verification score for contours: Methodology and application to Arctic ice-edge forecasts, *Quarterly Journal of the Royal Meteorological Society*, 144, 735–743, <https://doi.org/10.1002/qj.3242>, 2018.

- Goessling, H. F., Tietsche, S., Day, J. J., Hawkins, E., and Jung, T.: Predictability of the Arctic sea ice edge, *Geophysical Research Letters*, 43, 1642–1650, <https://doi.org/10.1002/2015gl067232>, 2016.
- Grigoryev, T., Verezhenskaya, P., Krinitskiy, M., Anikin, N., Gavrikov, A., Trofimov, I., Balabin, N., Shpilman, A., Eremchenko, A., Gulev, S., Burnaev, E., and Vanovski, V.: Data-Driven Short-Term Daily Operational Sea Ice Regional Forecasting, *Remote Sensing*, 14, <https://doi.org/10.3390/rs14225837>, URL <https://www.mdpi.com/2072-4292/14/22/5837>, 2022.
- He, K., Zhang, X., Ren, S., and Sun, J.: Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification, <https://doi.org/10.48550/ARXIV.1502.01852>, 2015a.
- He, K., Zhang, X., Ren, S., and Sun, J.: Deep Residual Learning for Image Recognition, 2015b.
- Huang, G., Liu, Z., van der Maaten, L., and Weinberger, K. Q.: Densely Connected Convolutional Networks, <https://doi.org/10.48550/ARXIV.1608.06993>, 2016.
- Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., and Darrell, T.: Caffe: Convolutional Architecture for Fast Feature Embedding, <https://doi.org/10.48550/ARXIV.1408.5093>, 2014.
- Kirillov, A., He, K., Girshick, R., Rother, C., and Dollár, P.: Panoptic Segmentation, <https://doi.org/10.48550/ARXIV.1801.00868>, 2018.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E.: ImageNet Classification with Deep Convolutional Neural Networks, in: *Advances in Neural Information Processing Systems*, edited by Pereira, F., Burges, C., Bottou, L., and Weinberger, K., vol. 25, Curran Associates, Inc., URL <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>, 2012.
- LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., and Jackel, L. D.: Backpropagation Applied to Handwritten Zip Code Recognition, *Neural Computation*, 1, 541–551, <https://doi.org/10.1162/neco.1989.1.4.541>, 1989.
- Long, J., Shelhamer, E., and Darrell, T.: Fully convolutional networks for semantic segmentation, in: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE, <https://doi.org/10.1109/cvpr.2015.7298965>, 2015.
- Melsom, A., Palerme, C., and Müller, M.: Validation metrics for ice edge position forecasts, *Ocean Science*, 15, 615–630, <https://doi.org/10.5194/os-15-615-2019>, 2019.
- Nair, V. and Hinton, G.: Rectified Linear Units Improve Restricted Boltzmann Machines Vinod Nair, vol. 27, pp. 807–814, 2010.
- Noh, H., Hong, S., and Han, B.: Learning Deconvolution Network for Semantic Segmentation, in: *2015 IEEE International Conference on Computer Vision (ICCV)*, IEEE, <https://doi.org/10.1109/iccv.2015.178>, 2015.
- Palerme, C., Müller, M., and Melsom, A.: An Intercomparison of Verification Scores

- for Evaluating the Sea Ice Edge Position in Seasonal Forecasts, *Geophysical Research Letters*, 46, 4757–4763, <https://doi.org/10.1029/2019gl082482>, 2019.
- RADU, M. D., COSTEA, I. M., and STAN, V. A.: Automatic Traffic Sign Recognition Artificial Intelligence - Deep Learning Algorithm, in: 2020 12th International Conference on Electronics, Computers and Artificial Intelligence (ECAI), IEEE, <https://doi.org/10.1109/ecai50035.2020.9223186>, 2020.
- Ronneberger, O., Fischer, P., and Brox, T.: U-Net: Convolutional Networks for Biomedical Image Segmentation, in: *Lecture Notes in Computer Science*, pp. 234–241, Springer International Publishing, [https://doi.org/10.1007/978-3-319-24574-4\\_28](https://doi.org/10.1007/978-3-319-24574-4_28), 2015.
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J.: Learning representations by back-propagating errors, *Nature*, 323, 533–536, <https://doi.org/10.1038/323533a0>, 1986.
- Simonyan, K. and Zisserman, A.: Very Deep Convolutional Networks for Large-Scale Image Recognition, <https://doi.org/10.48550/ARXIV.1409.1556>, 2014.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A.: Going Deeper with Convolutions, <https://doi.org/10.48550/ARXIV.1409.4842>, 2014.
- Yu, T. and Zhu, H.: Hyper-Parameter Optimization: A Review of Algorithms and Applications, <https://doi.org/10.48550/ARXIV.2003.05689>, 2020.
- Zeiler, M. D., Krishnan, D., Taylor, G. W., and Fergus, R.: Deconvolutional networks, in: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, IEEE, <https://doi.org/10.1109/cvpr.2010.5539957>, 2010.