

# Star–galaxy classification using deep convolutional neural networks

Edward J. Kim<sup>1★</sup> and Robert J. Brunner<sup>1,2,3,4</sup>

<sup>1</sup>*Department of Physics, University of Illinois, Urbana, IL 61801, USA*

<sup>2</sup>*Department of Astronomy, University of Illinois, Urbana, IL 61801, USA*

<sup>3</sup>*Department of Statistics, University of Illinois, Champaign, IL 61820, USA*

<sup>4</sup>*National Center for Supercomputing Applications, Urbana, IL 61801, USA*

Accepted 2016 October 12. Received 2016 October 4; in original form 2016 June 25

## ABSTRACT

Most existing star–galaxy classifiers use the reduced summary information from catalogues, requiring careful feature extraction and selection. The latest advances in machine learning that use deep convolutional neural networks (ConvNets) allow a machine to automatically learn the features directly from the data, minimizing the need for input from human experts. We present a star–galaxy classification framework that uses deep ConvNets directly on the reduced, calibrated pixel values. Using data from the Sloan Digital Sky Survey and the Canada–France–Hawaii Telescope Lensing Survey, we demonstrate that ConvNets are able to produce accurate and well-calibrated probabilistic classifications that are competitive with conventional machine learning techniques. Future advances in deep learning may bring more success with current and forthcoming photometric surveys, such as the Dark Energy Survey and the Large Synoptic Survey Telescope, because deep neural networks require very little, manual feature engineering.

**Key words:** methods: data analysis – methods: statistical – techniques: image processing – surveys – stars: statistics – galaxies: statistics.

## 1 INTRODUCTION

Currently ongoing and forthcoming large-scale photometric surveys, such as the Dark Energy Survey (DES) and the Large Synoptic Survey Telescope (LSST), aim to collect photometric data for hundreds of millions to billions of stars and galaxies. Due to the sheer volume of data, it is not possible for human experts to manually classify them, and the separation of photometric catalogues into stars and galaxies has to be automated. Furthermore, any classification approach must be probabilistic in nature. A fully probabilistic classifier enables a user to adopt probability cuts to obtain a pure sample for population studies, or to optimize the allocation of observing time by selecting objects for follow-up. Ideally, however, the probability estimates themselves would be retained for all sources and used in subsequent analyses to improve or enhance a particular measurement (Ross et al. 2011; Seo et al. 2012).

With machine learning, we can use algorithms to automatically create accurate source catalogues with well-calibrated posterior probabilities. Machine learning techniques have been a popular tool in many areas of astronomy (Ball et al. 2008; Banerji et al. 2010; Carrasco Kind & Brunner 2013; Ivezić et al. 2014; Kamdar, Turk & Brunner 2016). Artificial neural networks were first applied to the problem of star–galaxy classification in the work of Odewahn et al.

(1992), and they have become a core part of the astronomical image processing software SExtractor (Bertin & Arnouts 1996). Other successfully implemented examples of applying machine learning to the star–galaxy classification problem include decision trees (Weir, Fayyad & Djorgovski 1995; Suchkov, Hanisch & Margon 2005; Ball et al. 2006; Vasconcellos et al. 2011; Sevilla-Noarbe & Etayo-Sotos 2015), support vector machines (Fadely, Hogg & Willman 2012), and classifier combination strategies (Kim, Brunner & Kind 2015).

Almost all star–galaxy classifiers published in the literature use the reduced summary information available from astronomical catalogues. Constructing catalogues requires careful engineering and considerable domain expertise to transform the reduced, calibrated pixel values that comprise an image into suitable features, such as magnitudes or shape information of an object. In a branch of machine learning called deep learning (LeCun, Bengio & Hinton 2015), features are not designed by human experts; they are learned directly from data by deep neural networks. Deep learning methods learn multiple levels of features by transforming the feature at one level into a more abstract feature at a higher level. For example, when an array of pixel values is used as input to a deep learning method, the features in the first layer might represent locations and orientations of edges. The second layer could assemble particular arrangements of edges into more complex shapes, and subsequent layers would detect objects as combinations of low-level features. These multiple layers of abstraction progressively amplify aspects

★ E-mail: jkim575@illinois.edu

of the input that are important for classification tasks. Deep learning has been applied successfully to galaxy morphological classification in Sloan Digital Sky Survey (SDSS; Dieleman, Willett & Dambre 2015b) and Cosmic Assembly Near-infrared Deep Extragalactic Legacy Survey (Hurtas-Company et al. 2015) and to photometric redshift estimation (Hoyle 2016), but it has not yet been applied to the problem of source classification.

In this paper, we present a star–galaxy classification framework that uses a convolutional neural network (ConvNet) model directly on the images from the SDSS and the Canada–France–Hawaii Telescope Lensing Survey (CFHTLenS). We compare its performance with a standard machine learning technique that uses the reduced summary information from catalogues, and we demonstrate that our ConvNet model is able to produce accurate and well-calibrated probabilistic classifications with very little feature engineering by hand. In Section 2, we describe the data sets used in this paper and the pre-processing steps for preparing the image data sets. We provide a brief overview of deep learning and ConvNets in Section 3, and discuss our strategy for preventing overfitting in Section 4. In Section 5, we describe a state-of-the-art tree-based machine learning algorithm, to which the performance of our ConvNet model is compared. We present the main results of our ConvNet model in Section 6, and we outline our conclusions in Section 7.

## 2 DATA

To demonstrate the performance of our ConvNet model, we use photometric and spectroscopic data sets with different characteristics and compositions. In this section, we briefly describe these data sets and the image pre-processing steps for retrieving cutout images.

### 2.1 Sloan Digital Sky Survey

The SDSS (York et al. 2000) phases I–III obtained photometric data in five bands,  $u$ ,  $g$ ,  $r$ ,  $i$ , and  $z$ , covering  $14\,555\text{ deg}^2$ , more than one-third of the entire sky. The resulting catalogue contains photometry of over 300 million stars and galaxies with a limiting magnitude of  $r \approx 22$ , making the SDSS one of the largest sky surveys ever undertaken. The SDSS also conducted an expansive spectroscopic follow-up of more than three million stars and galaxies (Eisenstein et al. 2011). In this paper, we use a subset of the photometric and spectroscopic data contained within the Data Release 12 (DR12; Alam et al. 2015), which is publicly available through the online CasJobs server<sup>1</sup> (Li & Thakur 2008).

Using the CasJobs server, we randomly select a total of 65 000 sources, which are either stars or galaxies. In this work, we exclude objects that clearly are neither stars nor galaxies. Most of the excluded objects are QSOs or quasars. Quasars appear as point sources, rather than resolved sources similar to galaxies, and many of them have one or more saturated pixels in the images. However, unlike any known stars, their spectra show strong and broad emission lines. Quasars are also different from galaxies because of their intrinsic variability on a wide range of time-scales, which may be due to variation in the accretion rate or instabilities of the accretion disc around the black hole (Popović et al. 2012). Thus, many studies exclude quasars in the binary star–galaxy classification scheme (e.g. Vasconcellos et al. 2011; Fadelly et al. 2012). Expanding the historical star–galaxy classification problem to include additional classes,

e.g. nsng (neither star nor galaxy), may have advantages (Ball et al. 2006), and we plan to present the results of this multiclass problem in a future paper.

We also exclude some bad photometric observations as follows. We consider only objects with no warning flags in the spectroscopic measurement ( $z_{\text{Warning}} = 0$ ); the half-light radius in the  $r$  band is less than 30 arcsec as measured by the exponential and de Vaucouleurs light profiles; the error on the spectroscopic redshift measurement is less than 0.1; and the spectroscopic redshift is less than 2.

To create training images, we obtain the image FITS files for SDSS fields containing these objects in five photometric bands:  $u$ ,  $g$ ,  $r$ ,  $i$ , and  $z$ . We use the astrometry information in the FITS headers in the MONTAGE<sup>2</sup> software to align each image to the reference ( $r$ -band) image. We then use SExtractor to find the pixel positions of the 65 000 objects we have selected, and to centre each object in the cutout image. Magnitudes in the SDSS photometric catalogue are expressed as inverse hyperbolic sine magnitudes (also known as luptitudes; Lupton, Gunn & Szalay 1999), and we follow the SDSS convention and convert all flux values to luptitudes. Finally, in order to account for the effect of Galactic dust, extinction corrections in magnitudes are applied following Schlegel, Finkbeiner & Davis (1998). In the end, we have cutout images of size  $48 \times 48$  pixels with luptitude values in each pixel. We note that we have experimented with increasing the pixel dimensions to  $60 \times 60$  and  $72 \times 72$  pixels, but do not find noticeable improvement in the performance of our model.

In the end, we have 17 344 stars and 47 656 galaxies available for the training and testing processes. The apparent magnitudes range from  $10.7 < r < 23.1$ , and the galaxies in this sample have a mean redshift of  $z \sim 0.36$ . We randomly split the objects into training, held-out validation, and blind test sets of size 40 000, 10 000, and 15 000, respectively. We note that cross-validation is often avoided in deep learning in favour of hold-out validation, since cross-validation is computationally expensive. We also note that we perform a blind test, and the test set is not used in any way to train or calibrate the algorithms. The first two panels of Fig. 8 show the number of objects and the fraction of stars in the test set as functions of  $r$ -band magnitude. Similarly, Fig. 10 shows the number of objects and the fraction of stars in the test set as functions of  $g - r$  colour. The normalized kernel density estimate distributions for the training and validation sets are almost identical to those of the test set, and they are nearly indistinguishable when overlapped. We do not show the distributions for the training and validation sets in Figs 8 and 10 to avoid cluttering the plots.

### 2.2 Canada–France–Hawaii Telescope Lensing Survey

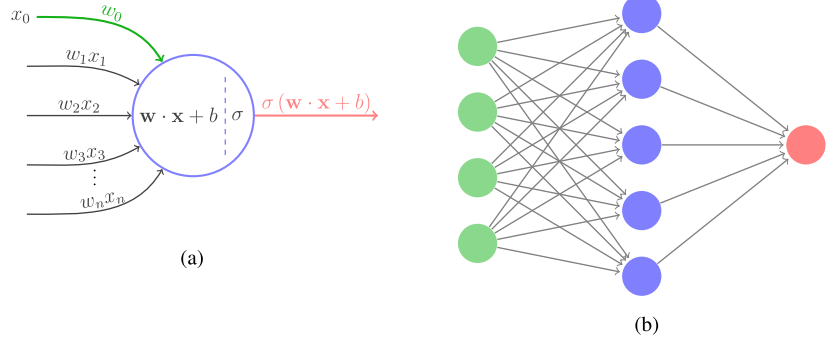
We also use photometric data from CFHTLenS<sup>3</sup> (Heymans et al. 2012; Hildebrandt et al. 2012; Erben et al. 2013). This catalogue consists of more than 25 million objects with a limiting magnitude of  $i_{\text{AB}} \approx 25.5$ . It covers a total of  $154\text{ deg}^2$  in the four fields (named W1, W2, W3, and W4) of the CFHT Legacy Survey (Gwyn 2012) observed in the five photometric bands:  $u$ ,  $g$ ,  $r$ ,  $i$ , and  $z$ .

We have cross-matched reliable spectroscopic galaxies from the Deep Extragalactic Evolutionary Probe Phase 2 (DEEP2; Davis et al. 2003; Newman et al. 2013), the SDSS Data Release 10 (SDSS-DR10; Alam et al. 2015), the Visible imaging Multi-Object

<sup>1</sup> <http://skyserver.sdss.org/casjobs/>

<sup>2</sup> <http://montage.ipac.caltech.edu/>

<sup>3</sup> <http://www.cfhtlens.org/>



**Figure 1.** (a) A mathematical model of a biological neuron. (b) A schematic diagram of a neural network with one hidden layer.

Spectrograph (VIMOS) Very Large Telescope (VLT) Deep Survey (VVDS; Le Fèvre et al. 2005; Garilli et al. 2008), and the VIMOS Public Extragalactic Redshift Survey (VIPERS; Garilli et al. 2014). We have selected only sources with very secure redshifts and no bad flags (quality flags  $-1$ ,  $3$ , and  $4$  for DEEP2; quality flag  $0$  for SDSS; quality flags  $3$ ,  $4$ ,  $23$ , and  $24$  for VIPERS and VVDS).

We obtain FITS images for each  $1 \text{ deg}^2$  CFHTLenS pointing that contains objects with spectroscopic labels. We create cutout images of size  $96 \times 96$  pixels by using a similar method to that described in Section 2.1. Finally, images are downsampled to  $48 \times 48$  pixels to reduce the computational cost.

In the end, we have 8545 stars and 57 843 galaxies available for the training and testing processes. The apparent magnitudes range from  $13.9 < r < 25.6$ , and the galaxies in this sample have a mean redshift of  $z \sim 0.59$ . We randomly split the objects into training, held-out validation, and blind test sets of size 40 000, 10 000, and 13 278, respectively. Figs 2 and 4 show the distribution of objects in the test set as functions of  $i$ -band magnitude and  $g - r$  colour. We do not show the distributions for the training and validation sets, since the normalized kernel density estimate distributions for the training and validation sets are almost identical to those of the test set.

### 3 DEEP LEARNING

Neural networks have many hyperparameters, including those that specify the network itself (e.g. the size and non-linearity of each layer) and those that specify how the network is trained (e.g. the mini-batch size or the learning rate). Furthermore, the architecture of a neural network can have a significant impact on its performance. In this section, we provide a brief description of key hyperparameters in our ConvNet model, and also present the network architecture.

#### 3.1 Neural networks

An artificial neuron in most artificial neural networks is represented as a mathematical function that models a biological neural structure (Aggarwal 2014). A schematic representation is shown in Fig. 1(a). Let  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  be a vector of inputs to a given neuron,  $\mathbf{w} = (w_1, w_2, \dots, w_n)$  be a vector of weights, and  $b$  be the bias. Then, the output of the neuron is (Rosenblatt 1961)

$$y = \sigma(\mathbf{w} \cdot \mathbf{x} + b), \quad (1)$$

where  $\sigma$  is the activation function (or non-linearity). The most popular non-linearity at present is the rectified linear unit (ReLU; Nair & Hinton 2010),  $\sigma(x) = \max(0, x)$ . ReLUs generally allow much faster training of deep neural networks with many layers. However, ReLU units can sometimes result in dead neurons whose output is always zero. To mitigate this problem, we use leaky ReLUs (Maas et al. 2013) that have a small, non-zero slope in the negative region,

$$\sigma(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0.01x & \text{if } x < 0. \end{cases} \quad (2)$$

Many deep learning models use feedforward neural network architectures with multiple layers, where each neuron in one layer is connected to the neurons of the subsequent layer (LeCun et al. 2015). A schematic representation is shown in Fig. 1(b). All layers except the input and output layers are conveniently called hidden layers.

We find a set of weights and biases such that, given  $N$  samples, the output from the network  $\mathbf{y} = (y_1, y_2, \dots, y_N)$  approximates the desired output  $\hat{\mathbf{y}} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N)$  as closely as possible for all input  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ . We can formulate this as the minimization of a loss function  $L(\mathbf{y}, \hat{\mathbf{y}})$  over the training data. In this work, we use cross-entropy (also called log loss; Murphy 2012) as the loss function. For binary classification, the cross-entropy per sample is given by

$$L(y_j, \hat{y}_j) = -\hat{y}_j \log_2 y_j - (1 - \hat{y}_j) \log_2 (1 - y_j), \quad (3)$$

where  $\hat{y}_j$  is the actual truth value (e.g. 0 or 1) of the  $j$ -th data, and  $y_j$  is the probability prediction made by the model. We compute the loss function by taking the average of all cross-entropies in the sample. Thus, the loss function becomes

$$L(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{N} \sum_{j=1}^N \hat{y}_j \log_2 y_j + (1 - \hat{y}_j) \log_2 (1 - \hat{y}_j). \quad (4)$$

To find the weights  $\mathbf{w}$  and biases  $\mathbf{b}$  which minimize the loss, we use a technique called gradient descent, where we use the following rules to update the parameters in each layer  $l$ :

$$\begin{aligned} \mathbf{w}_l &\rightarrow \mathbf{w}'_l = \mathbf{w}_l - \eta \frac{\partial L}{\partial \mathbf{w}_l} \\ \mathbf{b}_l &\rightarrow \mathbf{b}'_l = \mathbf{b}_l - \eta \frac{\partial L}{\partial \mathbf{b}_l}, \end{aligned} \quad (5)$$

where  $\eta$  is a small, positive number known as the learning rate. The gradients can be computed using the backpropagation procedure (Rumelhart, Hinton & Williams 1988). A common approach

**Table 1.** Summary of ConvNet architecture and hyperparameters. Note that pooling layers have no learnable parameters.

Type	Filters	Filter size	Padding	Non-linearity	Initial weights	Initial biases
Convolutional	32	$5 \times 5$	–	Leaky ReLU	Orthogonal	0.1
Convolutional	32	$3 \times 3$	1	Leaky ReLU	Orthogonal	0.1
Pooling	–	$2 \times 2$	–	–	–	–
Convolutional	64	$3 \times 3$	1	Leaky ReLU	Orthogonal	0.1
Convolutional	64	$3 \times 3$	1	Leaky ReLU	Orthogonal	0.1
Convolutional	64	$3 \times 3$	1	Leaky ReLU	Orthogonal	0.1
Pooling	–	$2 \times 2$	–	–	–	–
Convolutional	128	$3 \times 3$	1	Leaky ReLU	Orthogonal	0.1
Convolutional	128	$3 \times 3$	1	Leaky ReLU	Orthogonal	0.1
Convolutional	128	$3 \times 3$	1	Leaky ReLU	Orthogonal	0.1
Pooling	–	$2 \times 2$	–	–	–	–
Fully connected	2048	–	–	Leaky ReLU	Orthogonal	0.01
Fully connected	2048	–	–	Leaky ReLU	Orthogonal	0.01
Fully connected	2	–	–	Softmax	Orthogonal	0.01

to speed up training is to split the training data into mini-batches (LeCun et al. 1998a). In mini-batch gradient descent, instead of computing the gradients in equation (5) for the entire training data, we only compute the gradient of randomly chosen training examples at each step. As training examples are usually correlated, the gradient computed from each mini-batch is a good approximation of the overall gradient (Bottou 1998). As a result, mini-batch gradient descent results in much faster convergence. However, there is a trade-off: the lower the batch size is, the lower the convergence rate will be; the higher the batch size is, the longer it will take to compute the gradient at each step (Bousquet & Bottou 2008). Thus, a moderate batch size, combined with a decaying learning rate, is generally used in practice. We use a batch size of 128 in this work.

We define an epoch as a single, complete pass through the training data, and full training usually requires many epochs. At the end of each epoch, we evaluate the loss function on the validation set, and the model that minimizes the validation loss is chosen as the best model.

### 3.2 Convolutional neural networks

ConvNet (Fukushima 1980; LeCun et al. 1998b) is a type of deep, feedforward neural network that has recently become a popular approach in the computer vision community. In a typical ConvNet, the first few stages are composed of two types of layers: convolutional layers and pooling layers.

The input to a convolutional layer is an image, and the output channels of each layer are called feature maps. To produce output feature maps, we convolve each feature map with a set of weights called filters, and apply a non-linearity such as ReLU to the weighted sum of these convolutions. Different feature maps use different sets of filters, but all neurons in a feature map share the same set of filters. Mathematically, we replace the dot product in equation (1) with a sum of convolutions. Thus, the  $k$ -th feature map is given by

$$y^k = \sigma \left( \sum_m w_m^k * x_m + b^k \right), \quad (6)$$

where we sum over the set of input feature maps,  $*$  is the convolution operator, and  $w_m^k$  represents the filters.

Typically, a pooling layer computes the maximum of a local  $2 \times 2$  patch in a feature map (Krizhevsky, Sutskever & Hinton 2012). Since the pooling layer aggregates the activations of neighbouring units from the previous layer, it reduces the dimensionality

of the feature maps and makes the model invariant to small shifts and distortions (Boureau, Ponce & LeCun 2010). Two or more layers of convolution and pooling are stacked, followed by more convolutional and fully connected layers.

### 3.3 Neural network architecture

We present the overall architecture of our ConvNet model in Table 1. The network consists of 11 trainable layers. The first convolutional layer filters the  $5 \times 44 \times 44$  input image (i.e.  $44 \times 44$  images in five bands *ugriz*) with 32 square filters of size  $5 \times 5 \times 5$ . We have also experimented with using only three bands *gri* (for three channels of RGB) and four bands *ugri* and *griz* (corresponding to RGBA), and using only colours, e.g.  $u - g$ ,  $g - r$ ,  $r - i$ , and/or  $i - z$ , but we find that using magnitudes in all five bands *ugriz* yields the best performance.

The leaky ReLU non-linearity is applied to the output of the first convolutional layer (and all subsequent layers), and the second convolutional layer filters it with 32 filters of  $32 \times 3 \times 3$ . In the second convolutional layer (and all subsequent convolutional layers), we pad the input with zeros spatially on the border (i.e. the zero-padding is 1 pixel for  $3 \times 3$  convolutional layers) such that the spatial resolution is preserved after convolution. Max-pooling with filters of size  $2 \times 2$  follows the second convolutional layer. A stack of six additional convolutional layers, all with filters of size  $3 \times 3$ , is followed by three fully connected layers. The first two fully connected layers have 2048 channels each, and the third performs binary classification.

The output of the final fully connected layer is fed to a softmax function. The softmax function is given by

$$P(G | \mathbf{x}) = \frac{e^{\mathbf{x} \cdot \mathbf{w}_G}}{\sum_i e^{\mathbf{x} \cdot \mathbf{w}_i}}, \quad (7)$$

where we sum over the different possible values of the class label (i.e. star or galaxy), and interpret its output as the posterior probability that an object is a galaxy (or a star). We note that we could also try to solve a regression problem, e.g. by normalizing the output values that the network produces for each class. However, we find that solving a regression problem instead of using the softmax function results in significantly worse performance.

We have performed a manual search to explore more than 200 combinations of different architectures and hyperparameters to find an architecture that minimizes the loss function (equation 4) on the



validation set of the SDSS data. The architecture described in this section provides the best performance on the SDSS validation set. To test how this model performs across different, related data, we use the same architecture on the CFHTLenS data set.

The architecture of Krizhevsky et al. (2012) uses relatively large receptive fields ( $11 \times 11$ ) in the first convolutional layers. Zeiler & Fergus (2014) and Dieleman et al. (2015b) also use large receptive fields of  $7 \times 7$  and  $6 \times 6$  in the first convolution layer, respectively. However, we find that using a receptive field larger than  $5 \times 5$  in the first convolutional layer leads to worse performance. This result is in agreement with the network of Simonyan & Zisserman (2014), which has become known as ‘VGGNet’. VGGNet features an extremely homogeneous architecture that only performs  $3 \times 3$  convolutions. Using a large receptive field instead of a stack of multiple  $3 \times 3$  convolutions leads to a shallower network, and it is often preferable to increase the depth by using smaller receptive fields. However, we find that replacing the first layer with a stack of two  $3 \times 3$  convolutional layers increases the validation error, and thus use a  $5 \times 5$  convolution in the first layer.

In the remaining layers, we still follow VGGNet and add many  $3 \times 3$  convolutions (with zero-padding of size 1 pixel). Note that with the padding of 1 pixel for  $3 \times 3$  convolutional layers, the spatial resolution will be preserved after convolution. Such preservation of spatial resolution allows us to build relatively deep networks, as shown in Table 1. The main contribution of VGGNet is in showing that the depth plays an important role in good performance. In our case, we start with four convolutional layers and progressively add more layers, while monitoring the validation loss; we stop at eight convolutional layers after we find that adding more layers leads to worse performance. We conjecture that a greater depth and hence larger number of parameters lead to overfitting in our case.

The choice of momentum, learning rate, and initial weights is crucial for achieving high predictive performance and speeding up the learning process (Sutskever et al. 2013). To train our models, we use mini-batch gradient descent with a batch size of 128 and Nesterov momentum (Bengio, Boulanger-Lewandowski & Pascanu 2013) of 0.9. We initialize the learning rate  $\eta$  at 0.003 for all layers and decrease it linearly with the number of epochs from 0.003 to 0.0001 over 750 epochs. We also initialize the weights in each layer with random orthogonal initial conditions (Saxe, McClelland & Ganguli 2013). We use slightly positive values ( $b = 0.01$  or  $0.1$ ) for all biases. We find initializing biases to a small constant value helps eliminate dead neurons by ensuring that all ReLU neurons fire in the beginning.

To implement our model, we use PYTHON and the Lasagne library (Dieleman et al. 2015a), which is built on top of Theano (Theano Development Team 2016). The Theano library simplifies the use of GPU for computation, and using the GPU allows about an order of magnitude faster training than using just the CPU. We note that training our network takes about 40 h on an NVIDIA Tesla K40 GPU. In the interest of scientific reproducibility, we make all our code available at <https://github.com/EdwardJKim/dl4astro>.

## 4 REDUCING OVERFITTING

Our convolution neural network has  $11 \times 10^6$  learnable parameters, while there are only  $4 \times 10^4$  images in the training set. As a result, the model is very likely to overfit without regularization. In this section, we describe the techniques we used to minimize overfitting.

### 4.1 Data augmentation

One common method to combat overfitting is to artificially increase the number of training data by using label-preserving transformations (Krizhevsky et al. 2012; Dieleman, Willett & Dambre 2015b; Dieleman, De Fauw & Kavukcuoglu 2016). Each image is transformed as follows.

- (i) Rotation: rotating an image does not change whether the object is a star or a galaxy. We exploit this rotational symmetry and randomly rotate each image by a multiple of  $90^\circ$ .
- (ii) Reflection: we flip each image horizontally with a probability of 0.5 to exploit mirror symmetry.
- (iii) Translation: we also have translational symmetry in the images. Given an image of size  $48 \times 48$  pixels, we extract a random contiguous crop of size  $44 \times 44$ . Each cropping is equivalent to randomly shifting a  $44 \times 44$  image by up to 4 pixels vertically and/or horizontally.
- (iv) Gaussian noise: we introduce random Gaussian noise to each pixel values by using a similar method to Krizhevsky et al. (2012).

In addition to artificially increasing the size of the data set, these data augmentation schemes make the resulting model more invariant to rotation, reflection, translation, and small noise in the pixel values. We also note that the data augmentation steps add almost no computational cost, as they are performed on the CPU while the GPU is training the ConvNets on images.

### 4.2 Dropout

We use a regularization technique called dropout (Hinton et al. 2012) in the fully connected layers. Dropout consists of randomly setting to zero the output of each hidden neuron of the previous layer with probability 0.5. The weights of the remaining neurons are multiplied by 0.5 to preserve the scale of input values to the next layer. Since a neuron can be removed at any time, it cannot rely on the presence of other neurons in the same layer and is forced to learn more robust features.

### 4.3 Model combination

To make final classifications, we use our ConvNet model to make 64 sets of predictions for 64 transformations of the input images: 4 rotations, 4 horizontal translations, and 4 vertical translations (with random horizontal reflections). Although we use an identical network architecture for all transformations, we consider each set of predictions as separate results from different models. Finally, we use a model combination technique known as Bayesian model combination (BMC; Monteith et al. 2011), which uses Bayesian principles to generate an ensemble combination of different models. Although the data augmentation step in Section 4.1 should make our ConvNet model invariant to these types of transformations, we find that applying BMC still results in a significant increase in performance. For a thorough discussion of BMC, we refer the reader to Monteith et al. (2011, see also Carrasco Kind & Brunner 2014) for application of BMC to photometric redshift estimation, and Kim et al. (2015) for combining star–galaxy classifiers).

## 5 TREES FOR PROBABILISTIC CLASSIFICATIONS

To compare the performance of ConvNets with machine learning algorithms that use standard photometric features, we use a machine

learning framework called Trees for Probabilistic Classifications (TPC). TPC is a parallel, supervised machine learning algorithm that uses prediction trees and random forest techniques (Breiman et al. 1984; Breiman 2001) to produce a star–galaxy classification. A complete description of TPC is beyond the scope of this paper, and we refer the reader to Carrasco Kind & Brunner (2013) and Kim et al. (2015) for more details. While other random forest implementations exist, we have chosen TPC, because it is tailored specifically for astronomical use (Carrasco Kind & Brunner 2013); it has been tested for astronomical use cases, including photometric redshift estimation (Sánchez et al. 2014) and star–galaxy classification (Kim et al. 2015); and it uses parallelism to handle large data sets on distributed memory systems.

We train two TPC models on the SDSS data set by using different sets of attributes. The first model, which we denote  $\text{TPC}_{\text{phot}}$ , is trained with only nine photometric attributes: the extinction-corrected model magnitudes in five bands ( $u, g, r, i, z$ ) and their corresponding colours ( $u - g, g - r, r - i, i - z$ ). The second model, which we denote  $\text{TPC}_{\text{morph}}$ , is trained with the concentration parameter in each band in addition to the magnitudes and colours, for a total of 14 dimensions. The concentration is defined as the difference between the point spread function magnitude (psfMag) and the composite model magnitude (cModelMag), i.e. concentration  $\equiv \text{psfMag} - \text{cModelMag}$ . The SDSS pipeline uses a parametric method based on the concentration, an object is classified as a galaxy if concentration  $> 0.145$ . We find that the concentration is an excellent morphological feature for star–galaxy separation, and including more morphological features does not show noticeable improvement in performance. The concentration is a good example of carefully handcrafted feature extraction; we show in Section 6 that ConvNets do not require such feature engineering.

We also train two models on the CFHTLenS data set.  $\text{TPC}_{\text{phot}}$  is trained with the five magnitudes and their corresponding colours:  $u, g, r, i, z, u - g, g - r, r - i, i - z$ . Since the CFHTLenS catalogue does not provide the concentration parameter,  $\text{TPC}_{\text{morph}}$  uses SECTRACTOR’s FLUX\_RADIUS (the half-light radius), A\_WORLD (the semimajor axis), and B\_WORLD (the semiminor axis) for morphological features, in addition to the five magnitudes and their corresponding colours, for a total of 12 dimensions.

## 6 RESULTS AND DISCUSSION

In this section, we first describe the performance metrics that were used for evaluating the models. We then present the classification performance of our ConvNet model on the CFHTLenS and SDSS data sets, and compare it with the performance of TPC.

### 6.1 Classification metrics

Probabilistic classifiers, rather than only assigning discrete labels to each source, produce a continuous probability distribution of whether each source is a star or a galaxy. To evaluate the performance of probabilistic classifiers, many studies (e.g. Henrion et al. 2011; Fadely et al. 2012) convert probability estimates into class labels by choosing a probability threshold, e.g. a source is classified as a star if  $P_{\text{class}} < 0.5$ , and a galaxy if  $P_{\text{class}} > 0.5$ . However, using a fixed threshold ignores the model’s operating conditions, such as science requirements, misclassification costs, and stellar distribution. Furthermore, the probability threshold of 0.5 is not necessarily optimal for an unbalanced data set, where galaxies outnumber stars.

Following Kim et al. (2015), we use performance metrics that are well suited for probabilistic classifiers: area under the curve

**Table 2.** The definition of the classification performance metrics.

Metric	Meaning
AUC	Area under the receiver operating curve
MSE	Mean squared error
$c_g$	Galaxy completeness
$p_g$	Galaxy purity
$c_s$	Star completeness
$p_s$	Star purity
$p_g(c_g = x)$	Galaxy purity at $x$ galaxy completeness
$c_s(p_s = x)$	Star completeness at $x$ star purity
CAL	Calibration error with overlapping binning
$ \Delta N_g /N_g$	Absolute error in number of galaxies
log loss	Cross-entropy

(AUC) for the receiver operating characteristic (ROC) curve, completeness and purity, and the mean squared error (MSE). A good probabilistic classifier should also provide well-calibrated posterior probabilities. Thus, to evaluate calibration performance, we also use the calibration error and the absolute error in the estimation of number of galaxies. The definition of the metrics is summarized in Table 2.

#### 6.1.1 Receiver operating characteristic Curve

An ROC curve is the most commonly used method for evaluating the overall performance of a binary classifier (Swets, Dawes & Monahan 2000). In an ROC curve, we plot the true positive rate as a function of the false positive rate by varying the classification threshold. The AUC quantifies the overall performance in a single number.

#### 6.1.2 Completeness and purity

Let  $N_g$  be the number of true galaxies classified as galaxies, and  $M_g$  be the number of true galaxies classified as stars. Then the galaxy completeness  $c_g$  (also called recall or sensitivity) is given by

$$c_g = \frac{N_g}{N_g + M_g}. \quad (8)$$

Let  $M_s$  be the number of true stars classified as galaxies. Then the galaxy purity  $p_g$  (also called precision or positive predictive value) is given by

$$p_g = \frac{N_g}{N_g + M_s}. \quad (9)$$

We define the star completeness and purity in a similar way. As discussed in our previous work (Kim et al. 2015), we adopt weak lensing science requirements of the DES (Soumagnac et al. 2015), and compute  $p_g$  at  $c_g = 0.960$  and  $c_s$  at  $p_s = 0.970$ .

#### 6.1.3 Mean Squared Error

We also use MSE (also known as the Brier score Brier 1950) as a performance metric. We define MSE as

$$\text{MSE} = \frac{1}{N} \sum_{j=1}^N (y_j - \hat{y}_j)^2, \quad (10)$$

The MSE can be considered as both a score function that quantifies how well a set of probabilistic predictions is calibrated and a loss function.

**Table 3.** A summary of the classification performance metrics as applied to the CFHTLenS data. The definition of the metrics is summarized in Table 2. The bold entries highlight the best performance values within each column. Note that  $p_g(c_g = 0.96)$  and  $c_s(p_s = 0.97)$  require adjusting threshold values (i.e. probability cuts), while other metrics do not. To obtain a galaxy completeness of  $c_g = 0.96$ , we choose the threshold values 0.9972, 0.9963, and 0.9927 for ConvNet, TPC<sub>morph</sub>, and TPC<sub>phot</sub>, respectively; for star purity  $p_s = 0.97$ , we choose 0.6990, 0.5297, and 0.8570 for ConvNet, TPC<sub>morph</sub>, and TPC<sub>phot</sub>, respectively.

Classifier	AUC	MSE	$p_g(c_g = 0.96)$	$c_s(p_s = 0.97)$	CAL	$ \Delta N_g /N_g$	log loss
ConvNet	<b>0.9948</b>	0.0112	<b>0.9972</b>	0.8971	<b>0.0197</b>	<b>0.0029</b>	<b>0.0441</b>
TPC <sub>morph</sub>	0.9924	<b>0.0109</b>	0.9963	<b>0.9268</b>	0.0245	0.0056	0.0809
TPC <sub>phot</sub>	0.9876	0.0189	0.9927	0.8044	0.0266	0.0101	0.1085

#### 6.1.4 Calibration error

A fully probabilistic classifier predicts not only the class label, but also its confidence level on the prediction. In a well-calibrated classifier, the posterior class probability estimates should coincide with the proportion of objects that truly belong to a certain class. Probability calibration curves (or reliability curves; DeGroot & Fienberg 1983) are often used to display this relationship, where we bin the probability estimates and plot the fraction of positive examples versus the predicted probability in each bin (see Figs 5 and 11).

The problem with a binning approach is either too few or too many bins can distort the evaluation of calibration performance. Thus, we adopt a calibration measure based on overlapping binning (Caruana & Niculescu-Mizil 2004). We order the predicted values  $P_{\text{class}}$  and put the first 1000 elements in the first bin. We calculate the true probability  $P_{\text{gal}}$  by counting the true galaxies in this bin. The calibration error for this bin is  $|P_{\text{gal}} - P_{\text{class}}|$ . We then repeat this for the second bin (2 to 1001), the third bin (3 to 1002), and so on, and average the binned calibration errors. Thus, the overall calibration error is given by

$$\text{CAL} = \frac{1}{N-s} \sum_{b=1}^{N-s} \sum_{j=b}^{b+s-1} \left| P_{\text{class},j} - \frac{\sum_{j=b}^{b+s-1} P_{\text{gal},j}}{s} \right|, \quad (11)$$

where  $s = 1000$  is the bin length, which is chosen approximately equal to the number of objects in the testing set divided by the number of bins used in the calibration curve, i.e.  $s \approx N/10$ .

#### 6.1.5 Number of galaxies

Ideally, the probabilistic output of a classifier would be used in subsequent scientific analyses. For example, one can weight each object by the probability that it is a galaxy when measuring autocorrelation functions of luminous galaxies (Ross et al. 2011). In other words, given a well-calibrated classifier, instead of counting each galaxy equally, a galaxy could be counted as  $P_{\text{class}}$ , the probability estimate. This should in principle remove the contamination effect of stars. For a perfect classifier, we can count the total number of galaxies in the sample by summing the values of  $P_{\text{class}}$ . Thus, we measure the reliability of classifier output by the absolute error in the estimation of number of galaxies,

$$\frac{|\Delta N_g|}{N_g} = \frac{|N_g - \sum_{j=1}^N P_{\text{class},j}|}{N_g}. \quad (12)$$

## 6.2 CFHTLenS

As described in Section 3.2, we train our ConvNet model by monitoring its performance on the validation set. Once we have finished

training the model, we evaluate its performance on the blind test set. We also use the same training and validation sets to train and tune the hyperparameters of TPC<sub>morph</sub> and TPC<sub>phot</sub>, and perform classifications on the same test set to compare their performance with that of ConvNet.

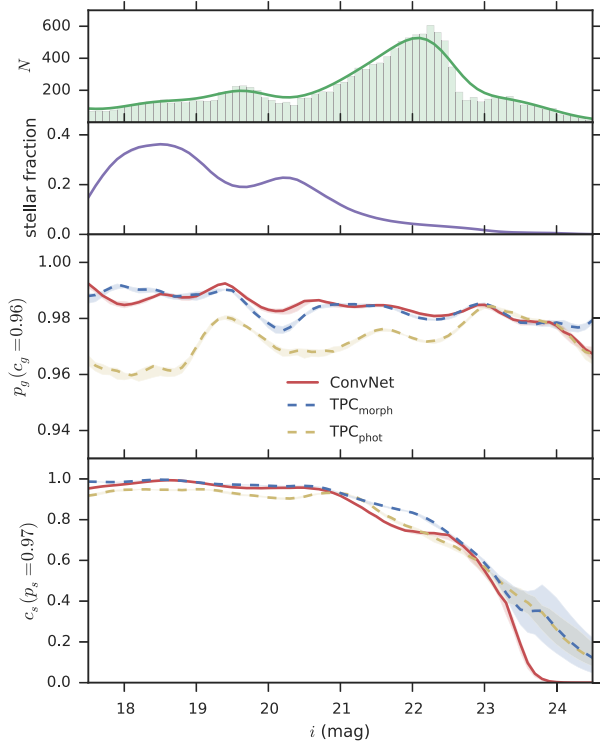
We present in Table 3 a summary of the results obtained by applying ConvNet, TPC<sub>morph</sub>, and TPC<sub>phot</sub> on the test set of the CFHTLenS data. The bold entries highlight the best technique for any particular metric. ConvNet outperforms TPC<sub>morph</sub> in five metrics (AUC,  $p_g$ , CAL,  $|\Delta N_g|/N_g$ , and log loss), while TPC<sub>morph</sub> performs better in two metrics (MSE and  $c_g$ ). It is not surprising that TPC<sub>phot</sub>, which is trained on only magnitudes and colours, performs significantly worse than both ConvNet and TPC<sub>morph</sub>. Thus, magnitudes and colours alone are not sufficient to separate stars from galaxies, and morphology is critical in separating stars from galaxies. ConvNet is able to learn the morphological features automatically from the images, and the performance of ConvNet is therefore comparable to that of TPC<sub>morph</sub>, which is trained on both morphological and photometric attributes.

In Fig. 2, we compare the galaxy purity and star completeness values for ConvNet, TPC<sub>morph</sub>, and TPC<sub>phot</sub>, as a function of  $i$ -band magnitude for the differential counts. We use kernel density estimation (KDE; Silverman 1986) with a Gaussian kernel. As the first panel shows, KDE is able to smooth the fluctuations in the distribution without binning. While ConvNet shows a slightly better performance than TPC<sub>morph</sub> in galaxy purity, ConvNet performs slightly worse than TPC<sub>morph</sub> in star completeness. Again, TPC<sub>phot</sub> performs significantly worse than both ConvNet and TPC<sub>morph</sub>, and this suggests that ConvNets are able to learn the shape information automatically from the images. We note that, at these operating conditions ( $c_g = 0.96$  or  $p_s = 0.97$ ), both ConvNet and TPC<sub>morph</sub> outperform the star–galaxy classification provided by the CFHTLenS pipeline (Hildebrandt et al. 2012) over all magnitudes.

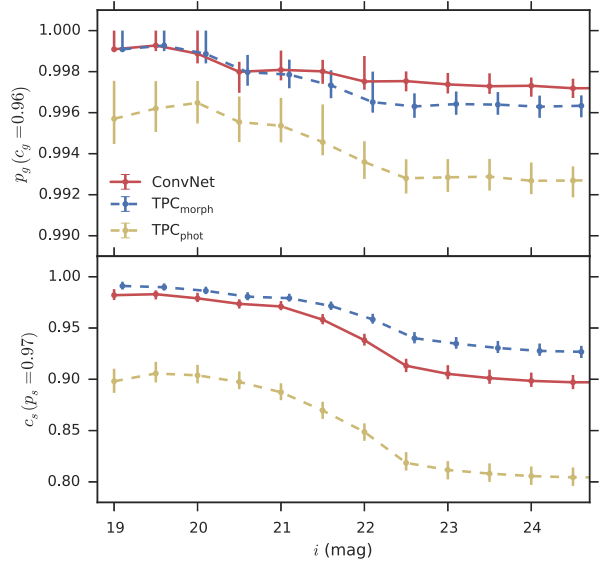
In Fig. 3, we show the overall galaxy purity and star completeness values as a function of  $i$ -band magnitude for the integrated counts. ConvNet is able to maintain a galaxy purity of 0.9972 up to  $i \sim 24.5$ , while the galaxy purity of TPC<sub>morph</sub> drops to 0.9963. However, TPC<sub>morph</sub> performs better than ConvNet in terms of star completeness, maintaining a purity of 0.9252 up to  $i \sim 24.5$ , while ConvNet drops to 0.8966.

We also show the galaxy purity and star completeness values as functions of  $g-r$  colour in Fig. 4. TPC<sub>morph</sub> provides slightly better completeness and purity than ConvNet between  $0.8 \lesssim g-r \lesssim 1.6$  while ConvNet outperforms TPC<sub>morph</sub> in the remaining regions.

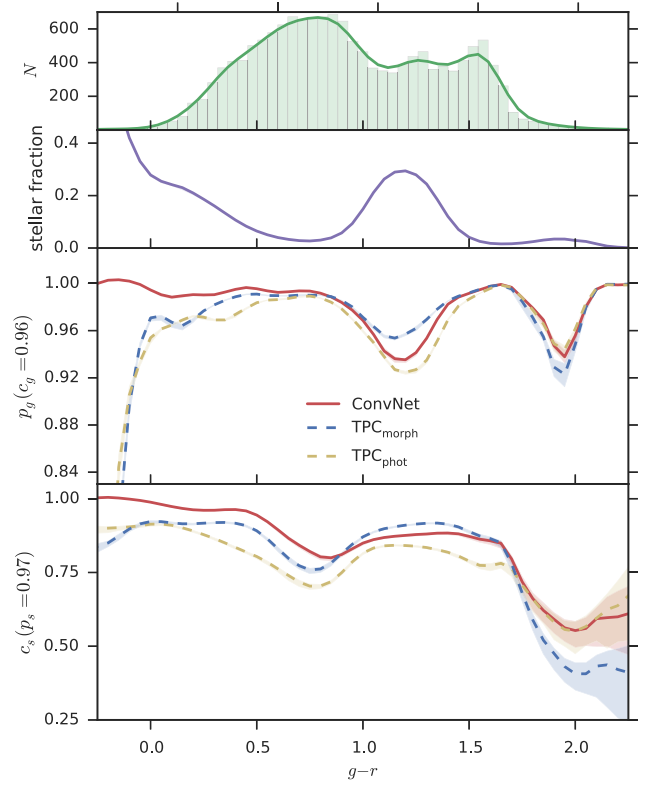
Fig. 5 shows the calibration curves that compare  $P_{\text{gal}}$ , the fraction of objects that are galaxies (as determined from their spectra), to  $P_{\text{class}}$ , the probabilistic outputs produced by ConvNet and TPC<sub>morph</sub>. The calibration curve for our ConvNet model is nearly diagonal, which implies that ConvNet is well calibrated and we can treat its probabilistic output as the probability that an object is a galaxy.



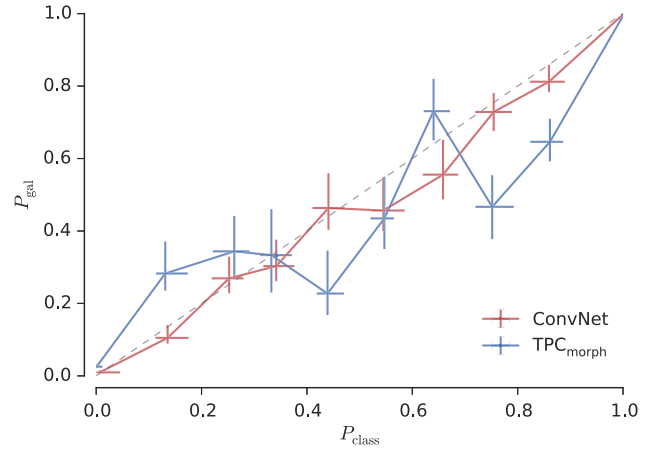
**Figure 2.** Galaxy purity and star completeness values as functions of the  $i$ -band magnitude (differential counts) as estimated by KDE in the CFHTLenS data set. The top panel shows the histogram with a bin size of 0.1 mag and the KDE for objects in the test set. The second panel shows the fraction of stars estimated by KDE as a function of magnitude. The bottom two panels compare the galaxy purity and star completeness values for ConvNet (red solid line),  $\text{TPC}_{\text{morph}}$  (blue dashed line), and  $\text{TPC}_{\text{phot}}$  (yellow dashed line) as functions of magnitude. The  $1\sigma$  confidence bands are estimated by bootstrap sampling.



**Figure 3.** Galaxy purity and star completeness as functions of the  $i$ -band magnitude (integrated counts) in the CFHTLenS data set. The upper panel compares the galaxy purity values for ConvNet (red solid line),  $\text{TPC}_{\text{morph}}$  (blue dashed line), and  $\text{TPC}_{\text{phot}}$  (yellow dashed line). The lower panel compares the star completeness values. The  $1\sigma$  error bars are computed following the method of Paterno (2004) to avoid the unphysical errors of binomial or Poisson statistics.



**Figure 4.** Similar to Fig. 2 but as a function of  $g - r$  colour. The bin size of histogram in the top panel is 0.05.

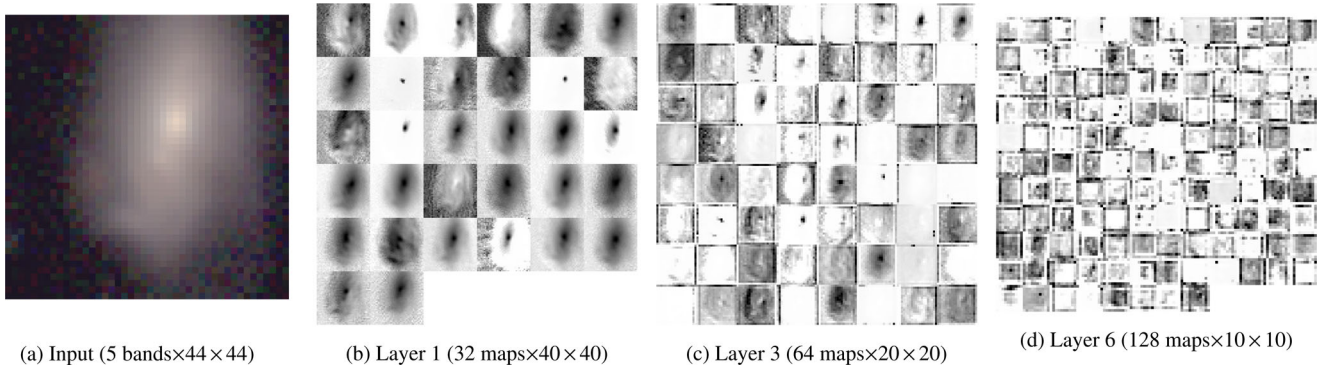


**Figure 5.** The calibration curves for ConvNet (red) and  $\text{TPC}_{\text{morph}}$  (blue) as applied to the CFHTLenS data set.  $P_{\text{gal}}$  is the fraction of objects that are galaxies, and  $P_{\text{class}}$  is the probabilistic outputs generated by the classifiers. The dashed line displays the relationship  $P_{\text{gal}} = P_{\text{conv}}$ . The  $1\sigma$  error bars are computed following the method of Paterno (2004).

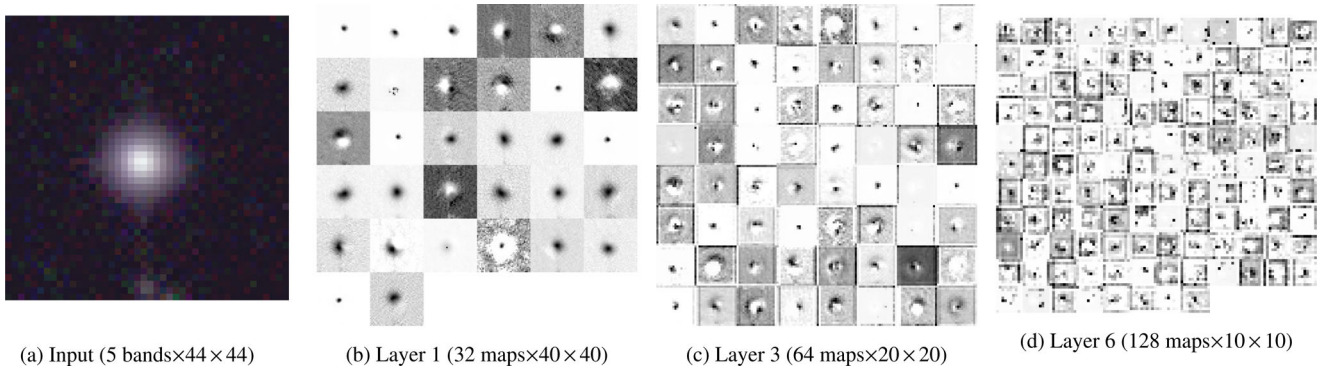
In contrast, the calibration curve for the probabilistic output of  $\text{TPC}_{\text{morph}}$  is apparently not as well-calibrated as ConvNet. These calibration curves visually confirm the results in Table 3 that the calibration error of ConvNet is about 20 per cent lower than that of  $\text{TPC}_{\text{morph}}$ . While probabilistic predictions can be further calibrated by using, e.g. isotonic calibration (Zadrozny & Elkan 2001), we do not consider additional probability calibration in this work.

It is informative to visualize how an input image activates the neurons in the convolutional layers. Figs 6 and 7 show the activations





**Figure 6.** (a) A sample  $44 \times 44$  RGB image of a galaxy in the CFHTLenS data set. The RGB image is created by mapping R  $\rightarrow$   $i$ -band magnitude, G  $\rightarrow$   $r$ -band magnitude, and B  $\rightarrow$   $g$ -band magnitude. (b) Activations on the first convolutional layer when a  $5 \times 44 \times 44$  image is fed into the network. (c) Activations on the third convolutional layer. (d) Activations on the sixth convolutional layer. Each image in (b), (c), and (d) is a feature map corresponding to the output for one of the learned features.



**Figure 7.** Similar to Fig. 6 but for a star in the CFHTLenS data set.

of the network when images of a galaxy and a star are fed into the network. The size of feature maps decreases with depth, and layers near the input layer have fewer filters while the later layers have more. The low-level features, e.g. edges or blobs, of the input images are still recognizable in the first convolutional layer. Subsequent layers use these low-level features to detect higher level features, and the final layer is a classifier that uses these high-level features. Thus, by performing hierarchical abstraction from low-level to high-level features, ConvNets are able to utilize shape information in the classification process.

### 6.3 SDSS

We have also trained and tested our ConvNet model on the SDSS data set, and we present in Table 4 the same six metrics for ConvNet, TPC<sub>morph</sub>, and TPC<sub>phot</sub>. The bold entries highlight the best technique for any particular metrics. In contrast with the CFHTLenS data set in Section 6.2, it is apparent that TPC<sub>morph</sub> outperforms ConvNet in all metrics except CAL and cross-entropy. Both ConvNet and TPC<sub>morph</sub> still outperform TPC<sub>phot</sub> in all six metrics by a significant amount, as magnitudes and colours alone are not sufficient to separate stars from galaxies. Although ConvNet performs worse than TPC<sub>morph</sub> on the SDSS data, its performance is much closer to TPC<sub>morph</sub>, as ConvNet is able to learn the shape information automatically from the images.

In Fig. 8, we compare the galaxy purity and star completeness values for ConvNet, TPC<sub>morph</sub>, and TPC<sub>phot</sub> as a function of  $r$ -band magnitude for the differential counts in the SDSS data. We note

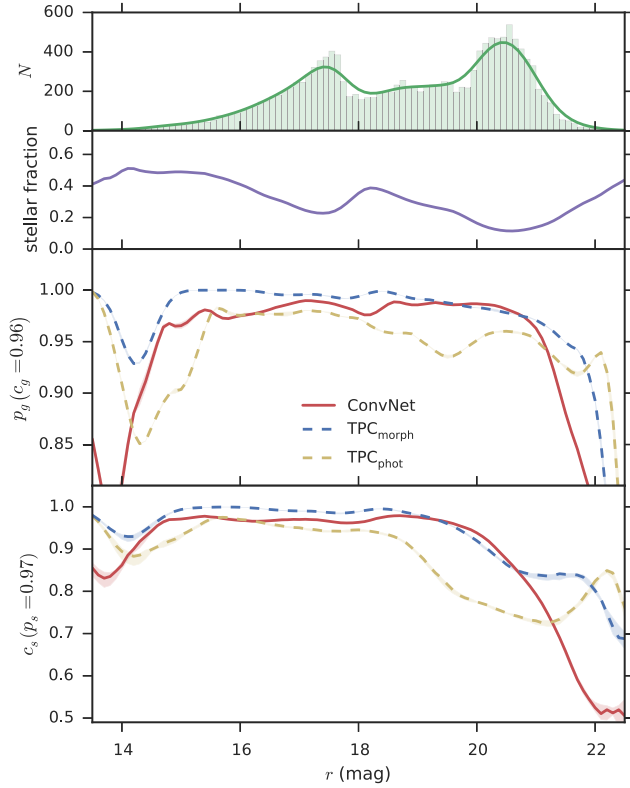
that TPC<sub>morph</sub> outperforms the star–galaxy classifier used by the SDSS pipeline (i.e. an object is classified as a galaxy if concentration  $> 0.145$ ) over all magnitudes. We do not show the SDSS classifications to avoid cluttering the plots. While ConvNet shows a similar but slightly worse performance than TPC<sub>morph</sub>, the galaxy purity and star completeness values of ConvNet begin to drop at faint magnitudes  $i \lesssim 21$ . Again, TPC<sub>phot</sub> performs significantly worse than both ConvNet and TPC<sub>morph</sub> at bright magnitudes. One reason that ConvNet fails to outperform TPC<sub>phot</sub>, especially at faint magnitudes, might be its overreliance on morphological features. Near a survey’s limit, the measurement uncertainties generally increase, and morphology is not a reliable metric for star–galaxy classification. Another possibility is that data augmentation has a negative effect at faint magnitudes, as the network may get confused by additional examples of faint galaxies that look like point sources. Data augmentation however is indispensable, since it improves the overall performance greatly.

In Fig. 9, we show the overall galaxy purity and star completeness values as a function of magnitude for the integrated counts. ConvNet is able to maintain a galaxy purity of 0.9915 up to  $i \sim 22.5$ , while TPC<sub>morph</sub> provides a galaxy purity of 0.9977. TPC<sub>morph</sub> also outperforms ConvNet in terms of star completeness, maintaining a purity of 0.9810 up to  $i \sim 22.5$ , while the star completeness of ConvNet drops to 0.9500.

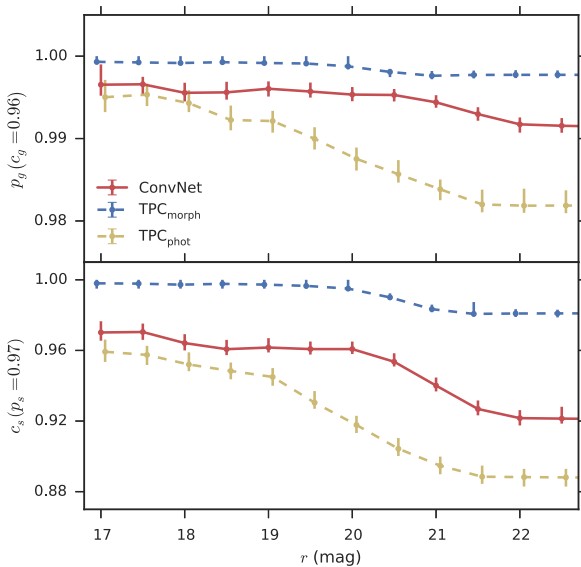
We also show the galaxy purity and star completeness values as a function of  $g - r$  colour in Fig. 10. ConvNet performs slightly better than TPC<sub>morph</sub> in both galaxy completeness and star purity between  $0.7 \lesssim g - r \lesssim 2.0$ , where the stellar fraction is relatively low. On

**Table 4.** A summary of the classification performance metrics as applied to the SDSS data. To obtain a galaxy completeness of  $c_g = 0.96$ , we choose the threshold values 0.7558, 0.9989, and 0.9360 for ConvNet, TPC<sub>morph</sub>, and TPC<sub>phot</sub>, respectively; for star purity  $p_s = 0.97$ , we choose 0.6046, 0.0547, and 0.7449 for ConvNet, TPC<sub>morph</sub>, and TPC<sub>phot</sub>, respectively.

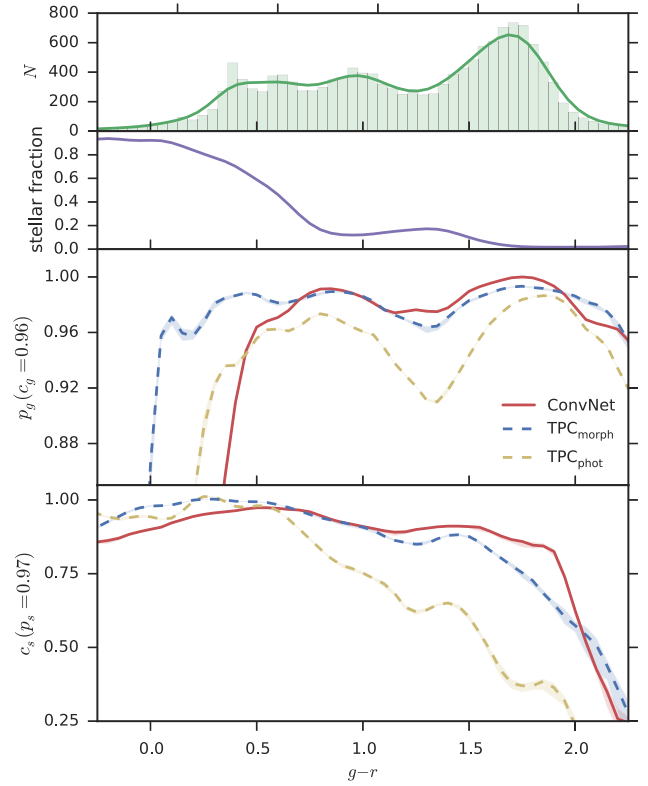
Classifier	AUC	MSE	$p_g(c_g = 0.96)$	$c_s(p_s = 0.97)$	CAL	$ \Delta N_g /N_g$	log loss
ConvNet	0.9952	0.0182	0.9915	0.9500	<b>0.0243</b>	0.0157	<b>0.0731</b>
TPC <sub>morph</sub>	<b>0.9967</b>	<b>0.0099</b>	<b>0.9977</b>	<b>0.9810</b>	0.0254	<b>0.0044</b>	0.0914
TPC <sub>phot</sub>	0.9886	0.0283	0.9819	0.8879	0.0316	0.0160	0.1372



**Figure 8.** Galaxy purity and star completeness as function of the  $r$ -band magnitude for the differential counts in the SDSS data set.



**Figure 9.** Galaxy purity and star completeness as functions of the  $r$ -band magnitude for the integrated counts in the SDSS data set.



**Figure 10.** Similar to Fig. 8 but as a function of  $g - r$  colour. The bin size of histogram in the top panel is 0.05.

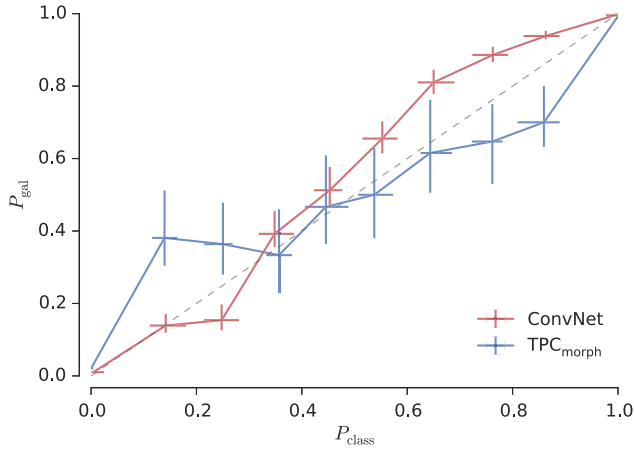
the other hand, both TPC<sub>morph</sub> and TPC<sub>phot</sub> outperform ConvNet in the region  $g - r \lesssim 0.8$  where the stellar fraction is higher.

Fig. 11 shows the calibration curves of ConvNet and TPC<sub>morph</sub>. The calibration curve of ConvNet in Fig. 11 is not as well calibrated as the calibration curve in Fig. 5, where the same ConvNet model was applied to the CFHTLenS data set. However, ConvNet may still be better calibrated than TPC<sub>morph</sub>, even when it is applied to the SDSS data set. Although it is not straightforward to compare the two calibration curves by visual inspection, Table 4 shows that the CAL metric of ConvNet is lower than that of TPC<sub>morph</sub>.

Figs 12 and 13 show the activations when images of a galaxy and a star are fed into the network. Similarly to Figs 6 and 7 in Section 6.2, the feature maps show hierarchical abstraction from low-level features in the first convolutional layer to high-level features in the subsequent layers. This hierarchical abstraction is what enables ConvNets to learn morphological features automatically from images.

## 7 CONCLUSIONS

We have presented a ConvNet for classifying stars and galaxies in the SDSS and CFHTLenS photometric images. For the CFHTLenS



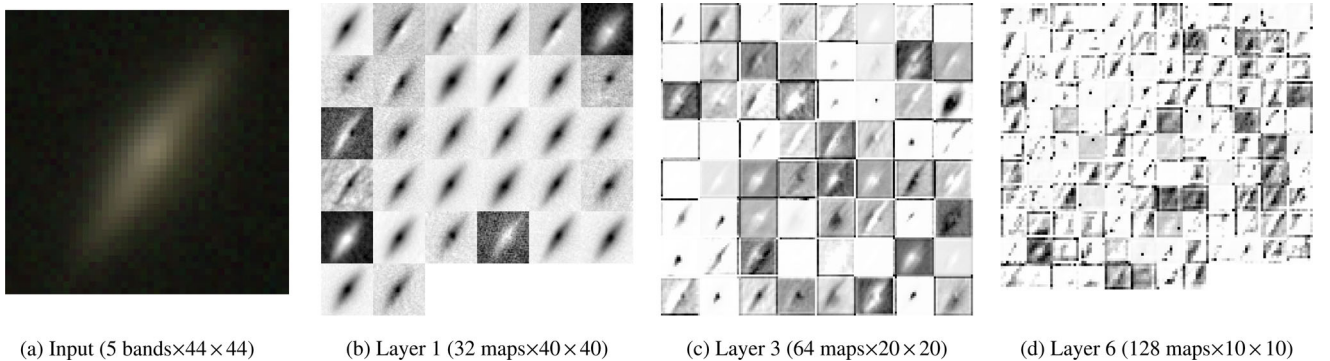
**Figure 11.** Calibration curves for ConvNet (red) and  $\text{TPC}_{\text{morph}}$  (blue) as applied to the SDSS data set.

data set, the network is able to provide a classification that is as accurate as a random forest algorithm (TPC), while the probability estimates of our ConvNet model appear to be better calibrated. When the same network architecture is applied to the SDSS data set, the network fails to outperform TPC, but the probabilities are still slightly better calibrated. The major advantage of ConvNets is that useful features are learned automatically from images, while traditional machine learning algorithms require feature engineering as a separate process to produce accurate classifications.

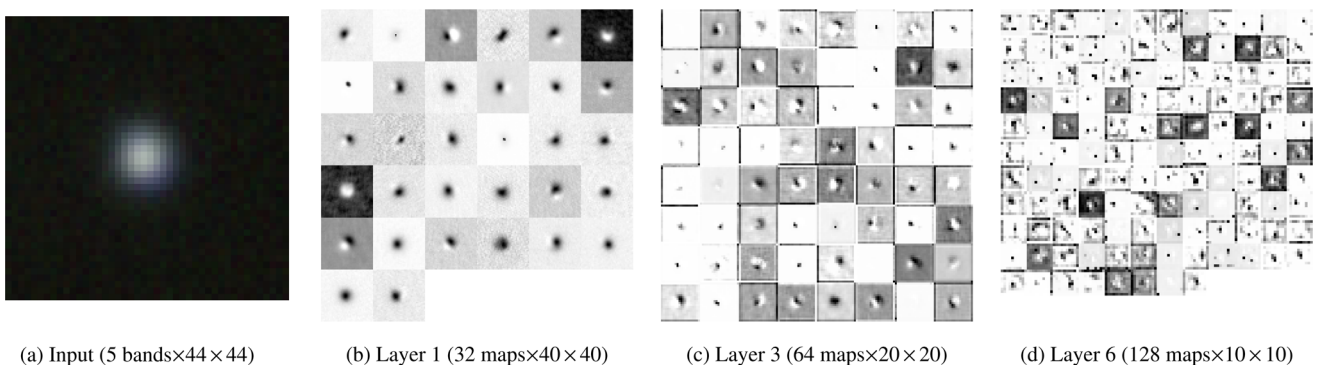
ConvNets have recently achieved record-breaking results in many image classification tasks (LeCun et al. 2015) and have been quickly and widely adopted by the computer vision community. One of the

main reasons for the success is that ConvNets are general-purpose algorithms that are applicable to a variety of problems without the need for designing a feature extractor. The lack of requirement for feature extraction is a huge advantage, e.g. when the task is to classify 1000 classes in the ImageNet data set (Russakovsky et al. 2015), as a good feature extractor for identifying images of cats would be of little use for classifying sailboats, and it is impractical to design a separate feature extractor for each class. However, when there already exists a good feature extractor for the problem at hand, e.g. the concentration parameter, the weight-averaged `spread_model` parameter from DES (Desai et al. 2012; Crocce et al. 2016), or even the SExtractor software, conventional machine learning algorithms that have been shown to be effective, such as TPC (Carrasco Kind & Brunner 2013; Kim et al. 2015), remain a viable option. As the ‘no free lunch’ theorem (Wolpert 1996) states, there is no model that works for every problem. For the CFHTLenS data set, our ConvNet model outperforms TPC. Since the SDSS catalogue provides the concentration parameter that is highly optimized for star–galaxy classification, TPC works better for SDSS.

Although we used various techniques to combat overfitting, it is possible that our ConvNet model has overfit the data. Overfitting could explain why our ConvNet model with maximal information fails to significantly outperform a standard machine learning algorithm that uses the reduced summary information from catalogues. The most effective way to prevent overfitting would be to simply collect more training images with spectroscopic follow-up, as the performance of ConvNets generally improves with more training data. However, spectroscopic observations are expensive and time-consuming, and it is unclear if sufficient training data will be available in future photometric surveys. If enough training data become



**Figure 12.** Similar to Fig. 6 but for a galaxy in the SDSS data set.



**Figure 13.** Similar to Fig. 12 but for a star in the SDSS data set.



available in DES or LSST, ConvNets become an attractive option because it can be applied immediately on reduced, calibrated images to produce well-calibrated posterior probabilities. We also note that using more training images will require the use of multi-GPU systems, which was beyond the scope of this work.

Deep learning is a rapidly developing field, and recent developments include improved network architectures. For future work, we plan to train more ConvNet variants, such as the Inception Module (Szegedy et al. 2015) and Residual Network (He, Zhang, Ren & Sun 2015). To improve the predictive performance, we have combined the predictions of different models across multiple transformations of the input images (Section 4.3). To further improve the performance, we could also train several networks with different architectures and combine the models. For example, the winning solution of Dieleman et al. (2015b) for the Galaxy Zoo challenge was based on a ConvNet model, and it required averaging many sets of predictions from models with different neural network architectures. Furthermore, future work could compare the performance of other deep learning variants, such as deep belief networks (Hinton, Osindero & Teh 2006), deep Boltzmann machines (Salakhutdinov & Hinton 2009), or multilayer perceptrons (Wasserman & Schwartz 1988).

It is also likely that the performance will be improved not only by training multiple network architectures, but also by combining them with different star–galaxy classifiers. In Kim et al. (2015), we combined a purely morphological classifier, a supervised machine learning method (TPC), an unsupervised machine learning method based on self-organizing maps, and a hierarchical Bayesian template fitting method, and demonstrated that our combination technique improves the overall performance over any individual classification method. ConvNets could be included as a different machine learning paradigm in this classifier combination framework to produce further improvements in predictive performance.

Our ConvNet model is a supervised algorithm, and one of the criticisms of supervised techniques is their difficulty in extrapolating past the limits of available spectroscopic training data. Since it is difficult to assess the classification performance without a deeper spectroscopic sample, we evaluated the performance using a test set that is drawn from the same underlying distribution as the spectroscopic sample. However, when our ConvNet model – trained on sources from a spectroscopic sample – is applied to a photometric sample – which is often fainter than our training set – the performance of ConvNet will be less reliable. Combining our ConvNet model with unsupervised methods (e.g. a template fitting method) in the aforementioned meta-classification framework will help improve the efficacy of star–galaxy classification beyond the limits of spectroscopic training data.

Finally, we are currently exploring different strategies for including objects that are neither stars nor galaxies (e.g. quasars). The results of this multiclass problem will be presented in subsequent papers.

## ACKNOWLEDGEMENTS

The authors thank the referee for a careful reading of the manuscript and comments that improved this work. The authors acknowledge support from the National Science Foundation Grant No. AST-1313415. RJB acknowledges support as an Associate within the Center for Advanced Study at the University of Illinois.

This work used the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant number ACI-1053575.

This work has made use of the Theano library (<http://deeplearning.net/software/theano/>), the Lasagne library (<http://lasagne.readthedocs.io/>), and the Nolearn library (<http://pythonhosted.org/nolearn/>).

This work is based on observations obtained with MegaPrime/MegaCam, a joint project of CFHT and CEA/DAPNIA, at the Canada–France–Hawaii Telescope (CFHT) which is operated by the National Research Council (NRC) of Canada, the Institut National des Sciences de l’Univers of the Centre National de la Recherche Scientifique (CNRS) of France, and the University of Hawaii. This research used the facilities of the Canadian Astronomy Data Centre operated by the National Research Council of Canada with the support of the Canadian Space Agency. CFHTLenS data processing was made possible thanks to significant computing support from the NSERC Research Tools and Instruments grant programme.

Funding for SDSS-III has been provided by the Alfred P. Sloan Foundation, the Participating Institutions, the National Science Foundation, and the U.S. Department of Energy Office of Science. The SDSS-III website is <http://www.sdss3.org/>.

SDSS-III is managed by the Astrophysical Research Consortium for the Participating Institutions of the SDSS-III Collaboration including the University of Arizona, the Brazilian Participation Group, Brookhaven National Laboratory, Carnegie Mellon University, University of Florida, the French Participation Group, the German Participation Group, Harvard University, the Instituto de Astrofísica de Canarias, the Michigan State/Notre Dame/JINA Participation Group, Johns Hopkins University, Lawrence Berkeley National Laboratory, Max Planck Institute for Astrophysics, Max Planck Institute for Extraterrestrial Physics, New Mexico State University, New York University, Ohio State University, Pennsylvania State University, University of Portsmouth, Princeton University, the Spanish Participation Group, University of Tokyo, University of Utah, Vanderbilt University, University of Virginia, University of Washington, and Yale University.

## REFERENCES

- Aggarwal C. C., 2014, *Data Classification: Algorithms and Applications*. CRC Press, Boca Raton, FL
- Alam S. et al., 2015, *ApJS*, 219, 12
- Ball N. M., Brunner R. J., Myers A. D., Tchong D., 2006, *ApJ*, 650, 497
- Ball N. M., Brunner R. J., Myers A. D., Strand N. E., Alberts S. L., Tchong D., 2008, *ApJ*, 683, 12
- Banerji M. et al., 2010, *MNRAS*, 406, 342
- Bengio Y., Boulanger-Lewandowski N., Pascanu R., 2013, in Ward R., Deng L., eds, *IEEE Int. Conf., Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, Piscataway, NJ, p. 8624
- Bertin E., Arnouts S., 1996, *A&AS*, 117, 393
- Bottou L., 1998, *Online Learn. Neural Netw.*, 17, 142
- Boureau Y.-L., Ponce J., LeCun Y., 2010, in Fürnkranz J., Joachims T., eds, *Proc. 27th Int. Conf. Mach. Learn. (ICML-10)*, p. 111. Available at: <http://www.icml2010.org/papers/638.pdf>
- Bousquet O., Bottou L., 2008, in Platt J. C., Koller D., Singer Y., Roweis S., eds, *Advances in Neural Information Processing Systems*. Cambridge, MA, p. 161
- Breiman L., 2001, *Mach. Learn.*, 45, 5
- Breiman L., Friedman J. H., Olshen R. A., Stone C. J., 1984, *Classification and Regression Trees*. CRC Press, New York
- Brier G. W., 1950, *Mon. Weather Rev.*, 78, 1
- Carrasco Kind M., Brunner R. J., 2013, *MNRAS*, 432, 1483
- Carrasco Kind M., Brunner R. J., 2014, *MNRAS*, 442, 3380
- Caruana R., Niculescu-Mizil A., 2004, in Kim W., Kohavi R., eds, *Proc. Tenth ACM SIGKDD Int. Conf. Knowl. Discovery Data Min.* ACM, New York, NY, p. 69



- Crocce M. et al., 2016, MNRAS, 455, 4301
- Davis M. et al., 2003, in Guhathakurta P., ed., Proc. SPIE Vol. 4834, Discoveries and Research Prospects from 6- to 10-Meter-Class Telescopes II. SPIE, Bellingham, p. 161
- DeGroot M. H., Fienberg S. E., 1983, The Statistician. Wiley R. Stat. Soc., Hoboken, NJ, p. 12
- Desai S. et al., 2012, ApJ, 757, 83
- Dieleman S. et al., 2015a, Lasagne: First release [Data set]. Zenodo. Available at: <http://doi.org/10.5281/zenodo.27878>
- Dieleman S., Willett K. W., Dambre J., 2015b, MNRAS, 450, 1441
- Dieleman S., De Fauw J., Kavukcuoglu K., 2016, in Balcan M., Weinberger K., eds, Proc. 33rd Int. Conf. Mach. Learn. JMLR, New York, p. 1889
- Eisenstein D. J. et al., 2011, AJ, 142, 72
- Erben T. et al., 2013, MNRAS, 433, 2545
- Fadely R., Hogg D. W., Willman B., 2012, ApJ, 760, 15
- Fukushima K., 1980, Biol. Cybern., 36, 193
- Garilli B. et al., 2008, A&A, 486, 683
- Garilli B. et al., 2014, A&A, 562, A23
- Gwyn S. D., 2012, AJ, 143, 38
- He K., Zhang X., Ren S., Sun J., 2015, preprint ([arXiv:1512.03385](https://arxiv.org/abs/1512.03385))
- Henriem M., Mortlock D. J., Hand D. J., Gandy A., 2011, MNRAS, 412, 2286
- Heymans C. et al., 2012, MNRAS, 427, 146
- Hildebrandt H. et al., 2012, MNRAS, 421, 2355
- Hinton G. E., Osindero S., Teh Y.-W., 2006, Neural Comput., 18, 1527
- Hinton G. E., Srivastava N., Krizhevsky A., Sutskever I., Salakhutdinov R. R., 2012, preprint ([arXiv:1207.0580](https://arxiv.org/abs/1207.0580))
- Hoyle B., 2016, Astron. Comput., 16, 34
- Huertas-Company M. et al., 2015, ApJS, 221, 8
- Ivezić Ž., Connolly A. J., VanderPlas J. T., Gray A., 2014, Statistics, Data Mining, and Machine Learning in Astronomy: A Practical Python Guide for the Analysis of Survey Data. Princeton Univ. Press, Princeton, NJ
- Kamdar H. M., Turk M. J., Brunner R. J., 2016, MNRAS, 455, 642
- Kim E. J., Brunner R. J., Kind M. C., 2015, MNRAS, 453, 507
- Krizhevsky A., Sutskever I., Hinton G. E., 2012, in Pereira F., Burges C. J. C., Bottou L., Weinberger K. Q., eds, Advances in Neural Information Processing Systems. Curran Assoc., Inc., Red Hook, NY, p. 1097
- Le Fèvre O. et al., 2005, A&A, 439, 845
- LeCun Y. A., Bottou L., Orr G. B., Müller K.-R., 1998a, in Neural networks: Tricks of the trade. Springer-Verlag, Berlin, p. 9
- LeCun Y. et al., 1998b, Proc. IEEE, 86, 2278
- LeCun Y., Bengio Y., Hinton G., 2015, Nature, 521, 436
- Li N., Thakar A. R., 2008, Comput. Sci. Eng., 10, 18
- Lupton R. H., Gunn J. E., Szalay A. S., 1999, AJ, 118, 1406
- Maas A. L., Hannun A. Y., Ng A. Y., 2013, in Dasgupta S., McAllester D., eds, The 30th Int. Conf. Mach. Learn. 2013 Workshop on Deep Learning for Audio, Speech, and Language Processing. J. Mach. Learn. Res.
- Monteith K. et al., 2011, in Marko K., Werbos P., eds, Int. Joint Conf. Neural Networks (IJCNN). IEEE, Piscataway, NJ, p. 2657
- Murphy K. P., 2012, Machine Learning: a Probabilistic Perspective. MIT Press, Cambridge, MA
- Nair V., Hinton G. E., 2010, in Fürnkranz J., Joachims T., eds, Proc. 27th Int. Conf. Mach. Learn. (ICML-10). p. 807. Available at: <http://www.icml2010.org/papers/432.pdf>
- Newman J. A. et al., 2013, ApJS, 208, 5
- Odehahn S. C., Stockwell E. B., Pennington R. L., Humphreys R. M., Zumach W. A., 1992, AJ, 103, 318
- Paterno M., 2004, Technical Report, Calculating Efficiencies and their Uncertainties. Department of Energy, Washington D.C.
- Popović L., Jovanović P., Stalevski M., Anton S., Andrei A. H., Kovačević J., Baes M., 2012, A&A, 538, A107
- Rosenblatt F., 1961, Technical report AD0256582, Principles of Neurodynamics. Perceptrons and the Theory of Brain Mechanisms. DTIC Document. Cornell Aeronaut. Lab Inc., Buffalo, NY
- Ross A. J. et al., 2011, MNRAS, 417, 1350
- Rumelhart D. E., Hinton G. E., Williams R. J., 1988, Cogn. Model., 5, 1
- Russakovsky O. et al., 2015, Int. J. Comput. Vis., 115, 211
- Salakhutdinov R., Hinton G. E., 2009, in van Dyk D., Welling M., eds, AISTATS. J. Mach. Learn. Res. p. 3
- Sánchez C. et al., 2014, MNRAS, 445, 1482
- Saxe A. M., McClelland J. L., Ganguli S., 2013, preprint ([arXiv:1312.6120](https://arxiv.org/abs/1312.6120))
- Schlegel D. J., Finkbeiner D. P., Davis M., 1998, ApJ, 500, 525
- Seo H.-J. et al., 2012, ApJ, 761, 13
- Sevilla-Noarbe I., Etayo-Sotos P., 2015, Astron. Comput., 11, 64
- Silverman B. W., 1986, Density Estimation for Statistics and Data Analysis, Vol. 26. CRC press, New York
- Simonyan K., Zisserman A., 2014, preprint ([arXiv:1409.1556](https://arxiv.org/abs/1409.1556))
- Soumagnac M. T. et al., 2015, MNRAS, 450, 666
- Suchkov A. A., Hanisch R. J., Margon B., 2005, AJ, 130, 2439
- Sutskever I., Martens J., Dahl G., Hinton G., 2013, in Dasgupta S., McAllester D., eds, Proc. 30th Int. Conf. Mach. Learn. (ICML-13), J. Mach. Learn. Res. p. 1139
- Swets J. A., Dawes R. M., Monahan J., 2000, Sci. Am., p. 83
- Szegedy C. et al., 2015, in Bischof H., Forsyth D., Schmid C., Sclaroff S., eds, Proc. IEEE Conf. Comput. Vis. Pattern Recognit. IEEE, Piscataway, NJ, p. 1
- Theano Development Team 2016, preprint ([arXiv:1605.02688](https://arxiv.org/abs/1605.02688))
- Vasconcellos E. C., de Carvalho R. R., Gal R. R., LaBarbera F. L., Capelato H. V., Frago Campos Velho H., Trevisan M., Ruiz R. S. R., 2011, AJ, 141, 189
- Wasserman P. D., Schwartz T., 1988, IEEE Expert, 3, 10
- Weir N., Fayyad U. M., Djorgovski S., 1995, AJ, 109, 2401
- Wolpert D. H., 1996, Neural Comput., 8, 1341
- York D. G. et al., 2000, AJ, 120, 1579
- Zadrozny B., Elkan C., 2001, in Brodley C. E., Danyluk A. P., eds, Int. Conf. Mach. Learn. Morgan Kaufmann, Burlington, MA, p. 609
- Zeiler M. D., Fergus R., 2014, in Van Gool L., Pollefeys M., eds, Computer vision–ECCV 2014. Springer-Verlag, Berlin, p. 818

This paper has been typeset from a  $\text{\LaTeX}$  file prepared by the author.