

tldr: We will introduce the concepts relevant to so called “deep learning” — our fundamental processes are based on computations performed over differentiable graphs, where nodes correspond to operations and edges correspond to operands.

**Instructor** Chris Curro, EE '15, MEE '16; professor@curro.cc

**Reference Textbook** Ian Goodfellow and Yoshua Bengio and Aaron Courville. 2016. *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>

**Assignments** There will be a handful of programming assignments — I recommend using Python and TensorFlow for these — each will be due either 1 or 2 weeks after the assigned date. There will be 2 larger projects: a midterm and final.

**Citations** Plagiarism will not be tolerated. All cases of suspected plagiarism will be submitted to the Dean's office for investigation. Feel free to ask questions of your peers, but please cite them for any help you receive. Cite resources you may utilize from the web and elsewhere.

**Quizzes** There will be quizzes most weeks. These quizzes will test understanding of assigned research papers. Expect 1-3 papers on most weeks. If you must miss a quiz, please let me know before hand and we will arrange appropriate accommodations, otherwise you receive a zero for that quiz.

**Grading** Grading breakdown in table at bottom of page. If you fail to submit an assignment you will fail the course. Unexcused late assignments will have a single letter grade deducted per 2 days late. The maximum grade for any tardy assignment is a B.

**Attendance** We will not take attendance, but it may factor into your participation score. Participation score is multifaceted. We will discuss this during the first class.

**Office hours** We will arrive at an appropriate schedule during the first class. Expect 1 or 2 hours per week. Additional hours by appointment. Office hours will be conducted remotely on Microsoft Teams.

Grading	
Assignments	30%
Projects	30%
Quizzes	30%
Participation	10%

## Approximate list of topics

**Introduction** Linear regression. Regression with basis functions. Gradient descent. Automatic differentiation; reverse mode and forward mode. Affine projection. Multi-layer perceptrons. Activation functions. Cross validation.  $L_1$  and  $L_2$  regularization. Dropout, batch normalization, and friends. Logistic regression. Binary cross entropy, and other entropy based loss functions. Weight initialization.

**Convolutions and friends** Convolutional layers. Strided convolutions. Pooling. Residual connections. Transposed convolutions.

**Applications and other techniques** Autoencoders. Super-resolution. Image inpainting. Speech generation. Speech recognition. Music generation. Image generation. Recommender systems. Text classification. Natural language generation. Reinforcement learning. Style transfer. Content transfer.

## Midterm project - due Oct 27.

The goal of the midterm project is to reproduce results from a contemporary research paper.

Procedure:

1. Find a paper of interest.
2. Pick a reasonable subset of the results to reproduce in the time allotted with present resource constraints.
3. Submit a proposal to me. If approved, continue. Else, go back to step 1 or 2, according to feedback.
4. Write code to reproduce the experiment. Document any necessary assumptions or changes from the paper.
5. Submit code and proof/evidence of reproduction. Submit a ~1-page document explaining your engagement with the work.

## Final project - due Dec 15.

The goal of the final project is to attempt to produce original research. We define success as a well-demonstrated engagement with the topic of the work.

Procedure:

1. Familiarize yourself with a topic of interest.
2. Propose amendment(s) to or suggest a novel application of an extant methodology.
3. Present the proposal to the class community for feedback. Iterate as necessary.
4. Write code to perform the experiment and produce the results.
5. Write a paper in contemporary conference-style describing your experiments and engagement with the work.
6. Produce a presentation (need not be slides) to present to your peers and guests.

## Assignment 1

tldr: Perform linear regression of a noisy sinewave using a set of gaussian basis functions with learned location and scale parameters. Model parameters are learned with stochastic gradient descent. Use of automatic differentiation is required. Hint: note your limits!

**Problem Statement** Consider a set of scalars  $\{x_1, x_2, \dots, x_N\}$  drawn from  $\mathcal{U}(0, 1)$  and a corresponding set  $\{y_1, y_2, \dots, y_N\}$  where:

$$y_i = \sin(2\pi x_i) + \epsilon_i \quad (1)$$

and  $\epsilon_i$  is drawn from  $\mathcal{N}(0, \sigma_{\text{noise}})$ . Given the following functional form:

$$\hat{y}_i = \sum_{j=1}^M w_j \phi_j(x_i | \mu_j, \sigma_j) + b \quad (2)$$

with:

$$\phi(x | \mu, \sigma) = \exp \frac{-(x - \mu)^2}{\sigma^2} \quad (3)$$

find estimates  $\hat{b}$ ,  $\{\hat{\mu}_j\}$ ,  $\{\hat{\sigma}_j\}$ , and  $\{\hat{w}_j\}$  that minimize the loss function:

$$J(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2 \quad (4)$$

for all  $(x_i, y_i)$  pairs. Estimates for the parameters must be found using stochastic gradient descent. A framework that supports automatic differentiation must be used. Set  $N = 50, \sigma_{\text{noise}} = 0.1$ . Select  $M$  as appropriate. Produce two plots. First, show the data-points, a noiseless sinewave, and the manifold produced by the regression model. Second, show each of the  $M$  basis functions. Plots must be of suitable visual quality.

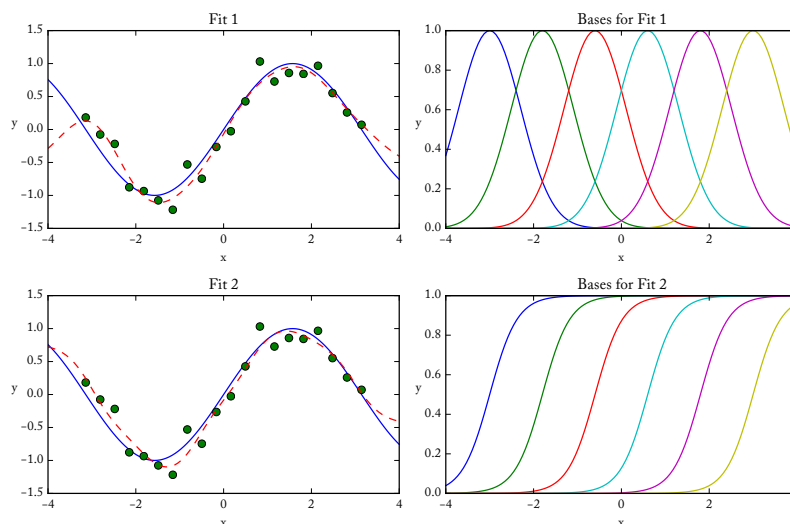


Figure 1: Example plots for models with equally spaced sigmoid and gaussian basis functions.

## Assignment 2

tldr: Perform binary classification on the spirals dataset using a multi-layer perceptron. You must generate the data yourself.

**Problem Statement** Consider a set of examples with two classes and distributions as in Figure 2. Given the vector  $x \in \mathbb{R}^2$  infer its target class  $t \in \{0, 1\}$ . As a model use a multi-layer perceptron  $f$  which returns an estimate for the conditional density  $p(t = 1 | x)$ :

$$f: \mathbb{R}^2 \rightarrow [0, 1] \quad (5)$$

parametrized by some set of values  $\theta$ . All of the examples in the training set should be classified correctly (i.e.  $p(t = 1 | x) > 0.5$  if and only if  $t = 1$ ). Impose an  $L^2$  penalty on the set of parameters. Produce one plot. Show the examples and the boundary corresponding to  $p(t = 1 | x) = 0.5$ . The plot must be of suitable visual quality. It may be difficult to find an appropriate functional form for  $f$ , write a few sentences discussing your various attempts.

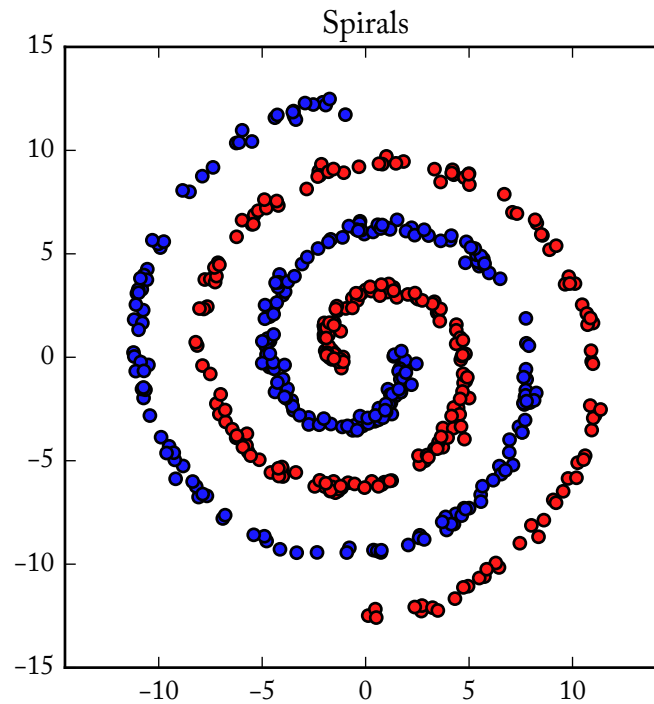


Figure 2: Sample spiral data.

### Assignment 3

tldr: Classify MNIST digits with a (optionally convolitional) neural network. Get at least 95.5% accuracy on the test test.

**Problem Statement** Consider the MNIST dataset consisting of 50,000 training images, and 10,000 test images. Each instance is a  $28 \times 28$  pixel handwritten digit zero through nine. Train a (optionally convolutional) neural network for classification using the training set that achieves at least 95.5% accuracy on the test set. Do not explicitly tune hyperparameters based on the test set performance, use a validation set taken from the training set as discussed in class. Use dropout and an  $L^2$  penalty for regularization. Note: if you write a sufficiently general program the next assignment will be very easy.

Do not use the built in MNIST data class from TensorFlow.

**Extra challenge (optional)** In addition to the above, the student with the fewest number of parameters for a network that gets at least 80% accuracy on the test set will receive a prize. There will be an extra prize if any one can achieve 80% on the test set with a single digit number of parameters. For this extra challenge you can make your network have any crazy kind of topology you'd like, it just needs to be optimized by a gradient based algorithm.

## Assignment 4

tldr: Classify CIFAR10. Achieve performance similar to the state of the art. Classify CIFAR100. Achieve a top-5 accuracy of 90%.

**Problem Statement** Consider the CIFAR10 and CIFAR100 datasets which contain  $32 \times 32$  pixel color images. Train a classifier for each of these with performance similar to the state of the art (for CIFAR10). It is your task to figure out what is state of the art. Feel free to adapt any techniques from papers you read. I encourage you to experiment with normalization techniques and optimization algorithms in this assignment. Write a paragraph or two summarizing your experiments. Hopefully you'll be able to reuse your MNIST program.

## Assignment 5

tldr: Classify the AG News dataset posted on the course website. Achieve an accuracy similar to the state of the art as of 2017 [1].

**Problem Statement** Consider the AG News dataset [2] which contains headlines and descriptions for a large set of news articles. Perform proper cross validation and achieve a classification accuracy similar to those listed in Table 2 of [1]. Use methods similar to those discussed in class.

## References

- [1] L. Wu, A. Fisch, S. Chopra, K. Adams, A. Bordes, and J. Weston, "Starspace: Embed all the things!" *CoRR*, vol. abs/1709.03856, 2017. [Online]. Available: <http://arxiv.org/abs/1709.03856>
- [2] "AG News Dataset." [Online]. Available: <http://ee.cooper.edu/~curro/cgml/week5/ag-news-csv.tar.gz>