

Title: PSI 2017/2018 - Machine Learning for Many Body Physics - Lecture 3

Date: Apr 11, 2018 09:00 AM

URL: <http://pirsa.org/18040056>

Abstract:

Reminder:

MSE: Mean-Squared Error

CE: Cross-Entropy

GD: Gradient Descent

Outline for today:

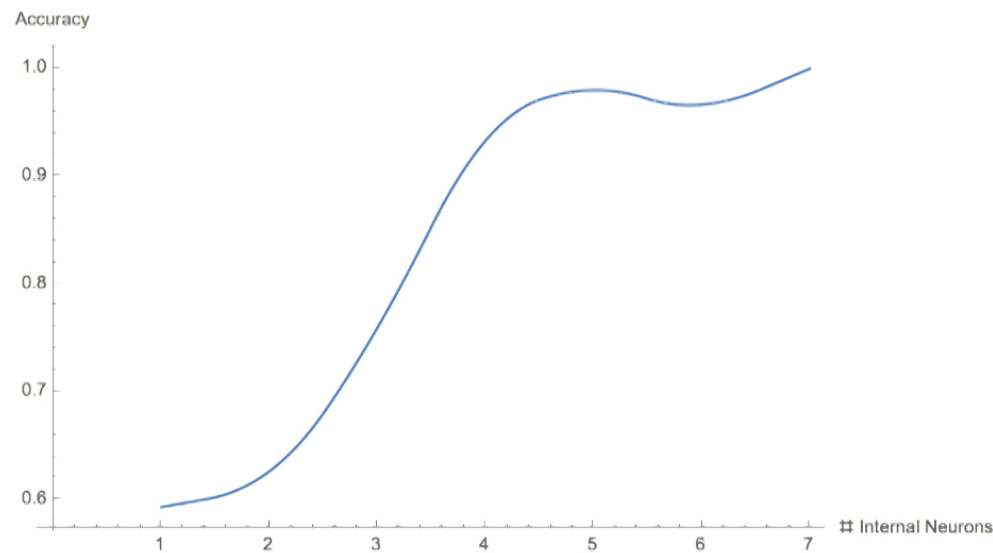
↳ alternatives to GD

↳ backpropagation

↳ overfitting

Group 1

Accuracy vs. the number of neurons in the hidden layer

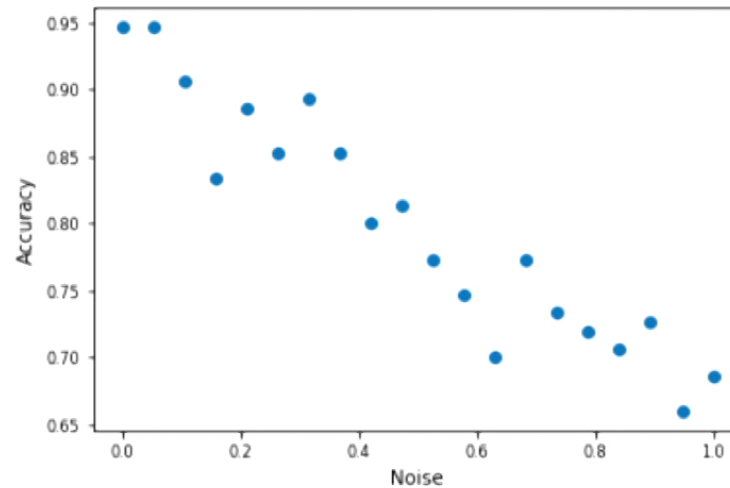


Carlos, Cheryne, Kfir, Vladimir



Group 2

Accuracy vs. the magnitude of noise in the data

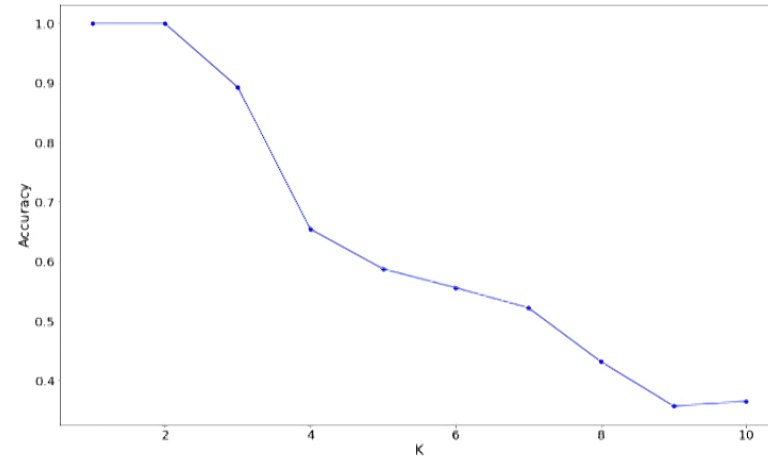
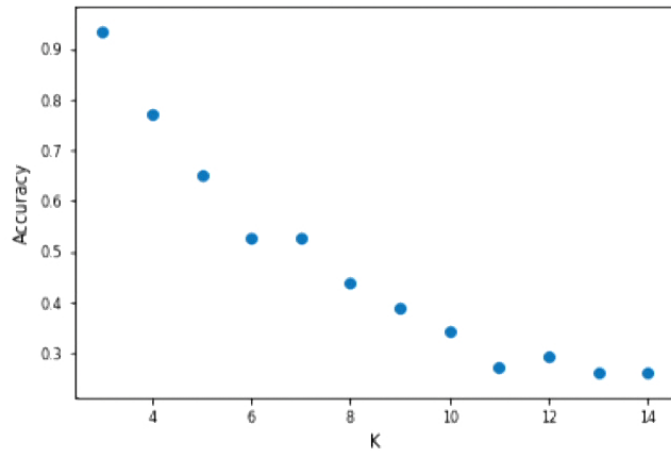


Ahmed, Gwyneth, Lei, Matt



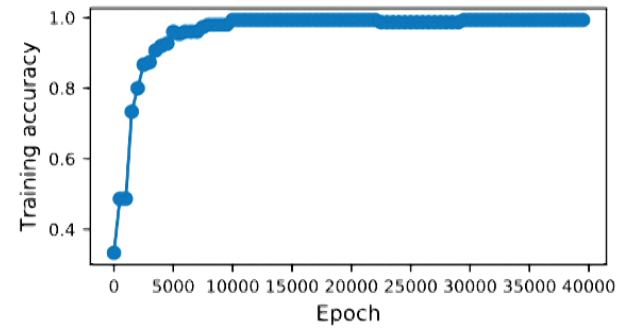
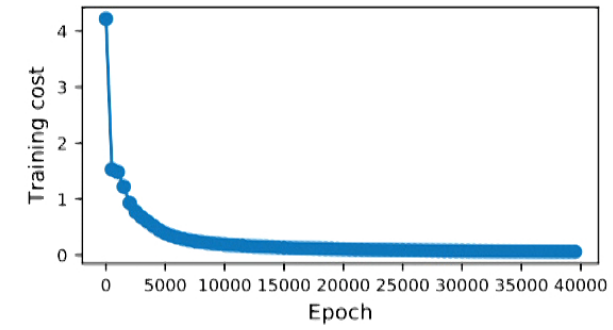
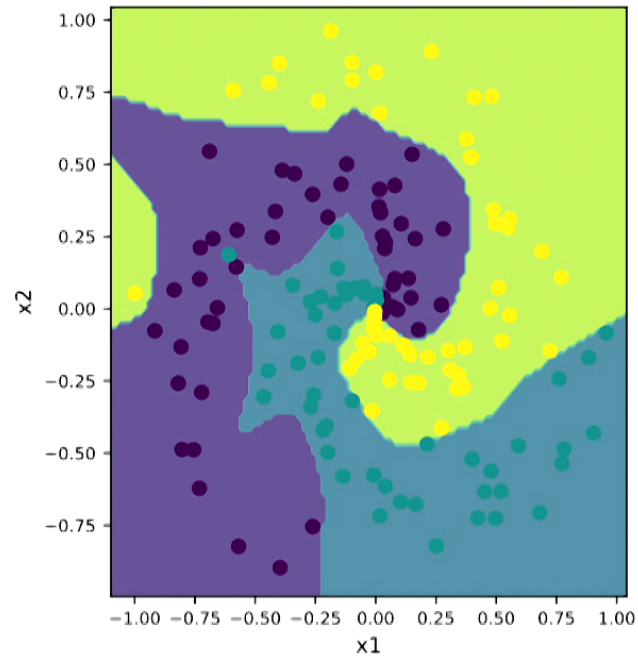
Group 3

Accuracy vs. the number of different labels (branches)



Lusine, Qi, Stavros, TC

mag_noise = 0.5
nH = 100



Alternatives to GD

① Stochastic GD (SGD)

Assume $C = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} c(\mathbf{x})$

Idea: Approximate the gradient $\vec{\nabla}_{\mathbf{p}} C$ by summing over a randomly-chosen "mini-batch" B of M samples (instead of summing over all of the dataset \mathcal{D})

Benefits

↳ speeds up calculation of gradient

↳ introduces randomness: decreases chance of getting stuck in a local min.

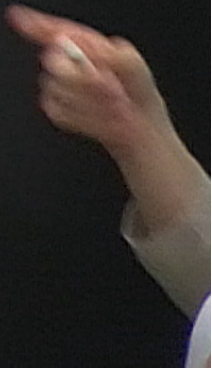
Adding momentum

↳ keep some "memory" \vec{v} of previous gradients

$$\vec{v} \rightarrow \gamma \vec{v} - \eta \vec{\nabla}_{\vec{p}} C$$

$$\vec{p} \rightarrow \vec{p} + \vec{v}$$

Benefit:
parameters P_k can keep changing
when the gradient gets small as
long as \vec{v} is not also small



Algorithms with Adaptive Learning Rates

↳ learning rates change, can be different for each of the weights and biases

↳ eg) ADAM, RMSProp

Many other training models
are available!

(see, for example, `tf.train`)

Backpropagation

↳ used to calculate the gradients needed for training

$$\text{Recall } C = \frac{1}{|\mathcal{D}|} \sum_{\vec{x} \in \mathcal{D}} c(\vec{x})$$

We will see how to calculate $\frac{\partial c(\vec{x})}{\partial W_{ij}^{(l)}}$ and $\frac{\partial c(\vec{x})}{\partial b_j^{(l)}}$ (corresponding to one training sample \vec{x}).

Then we average over many samples
to get gradients of C

Recall: $a_j^{(l)} = g_l(z_j^{(l)})$

$$z_j^{(l)} = \sum_i a_i^{(l-1)} W_{ij}^{(l)} + b_j^{(l)}$$

responding

Define $\delta_j^{(l)} = \frac{\partial c}{\partial z_j^{(l)}}$

c is written as a function of the last layer's output $a^{(L)}$: $c = c(a^{(L)})$

For the last layer: $\delta_j^{(L)} = \frac{\partial c}{\partial z_j^{(L)}} = \frac{\partial c}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}}$

$$\delta_j^{(L)} = \frac{\partial C}{\partial a_j^{(L)}} g_L'(z_j^{(L)}) \quad (1)$$

For layer $l < L$:

$$\delta_j^{(l)} = \frac{\partial c}{\partial z_j^{(l)}} = \sum_k \frac{\partial c}{\partial z_k^{(l+1)}} \frac{\partial z_k^{(l+1)}}{\partial z_j^{(l)}} = \dots$$

SHOW

$$\delta_j^{(l)} = \sum_k \delta_k^{(l+1)} W_{kj}^{(l+1)T} g_l'(z_j^{(l)}) \quad (2)$$

Now one can show that

$$\frac{\partial c}{\partial w_j^{(l)}} = a_i^{(l-1)} \delta_j^{(l)} \quad (3)$$

$$\frac{\partial c}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad (4)$$

Summary training algorithm (one "epoch")

- # forward pass

$$\vec{a}^{(0)} = \vec{x}$$

for $l=1, \dots, L$

$$\vec{a}^{(l)} = g_l(\vec{a}^{(l-1)} W^{(l)} + \vec{b}^{(l)})$$

• #backpropagation

get $\vec{\delta}^{(L)}$ from (1)

for $l = L-1, \dots, 1$:

get $\vec{\delta}^{(l)}$ from (2)

get $\frac{\partial C}{\partial W_{ij}^{(l)}}$ and $\frac{\partial C}{\partial b_j^{(l)}}$ from (3) & (4)

• #update

$$W_{ij}^{(l)} = W_{ij}^{(l)} - \eta \frac{\partial C}{\partial W_{ij}^{(l)}}$$

$$b_j^{(l)} = b_j^{(l)} - \eta \frac{\partial C}{\partial b_j^{(l)}}$$

} OR use an alternative to GD

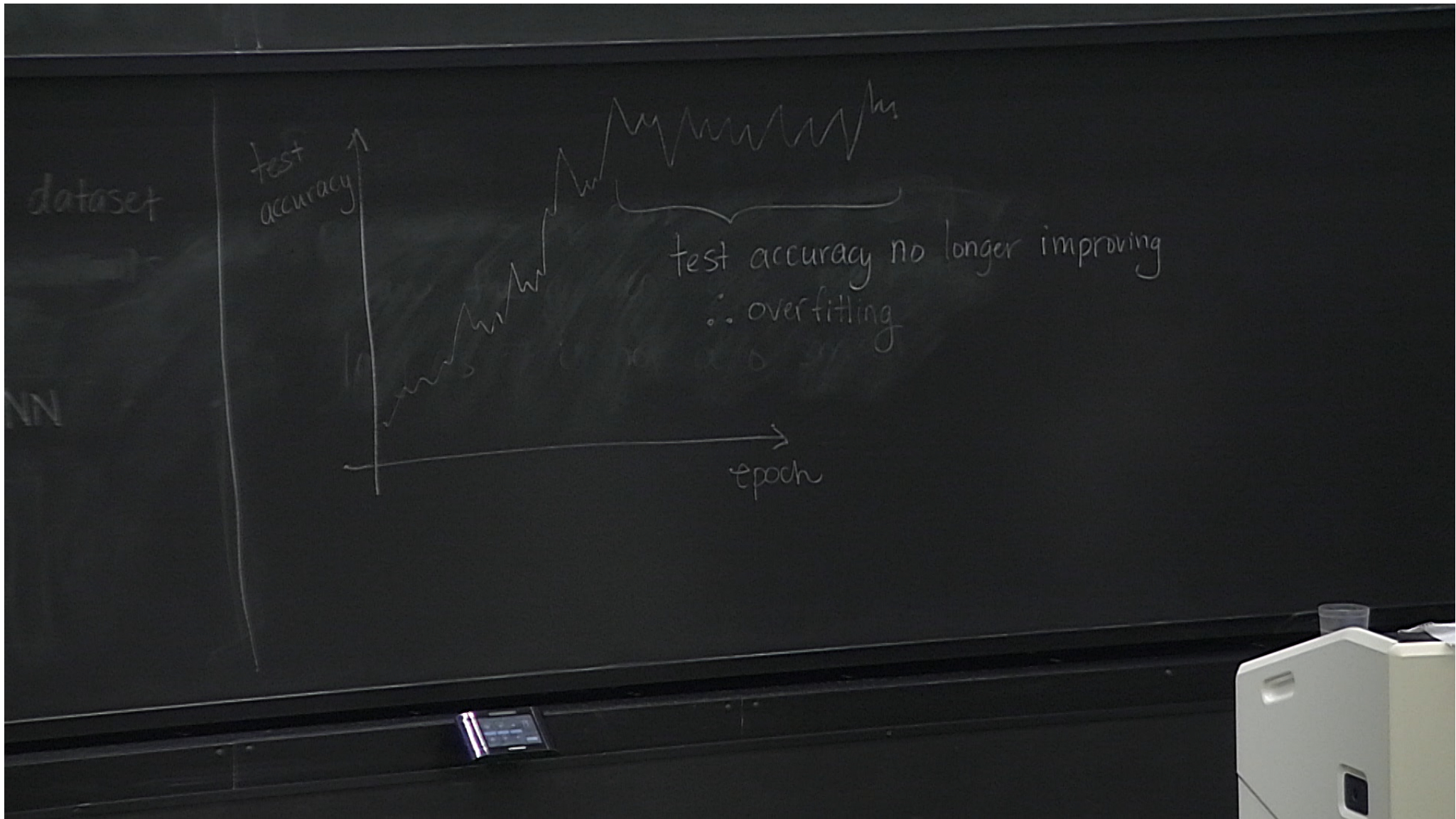
Overfitting

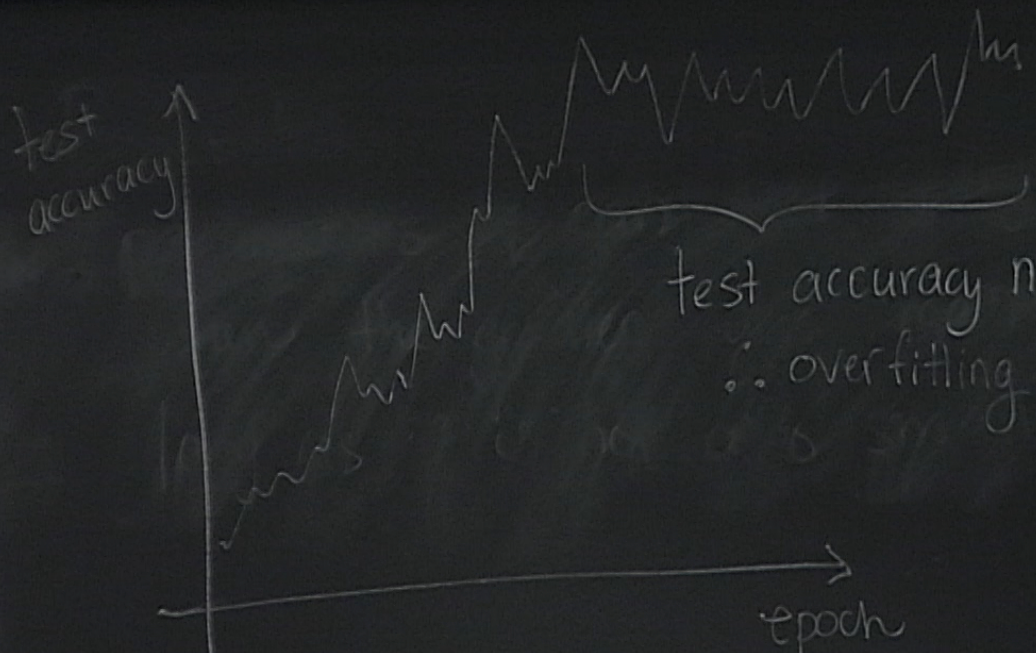
In Tutorial #1, we saw that C could be decreasing and (in some cases) the NN could get $\sim 100\%$ accuracy on the training data.
(% of label correctly predicted)

However, the NN was clearly overfitting!

We can introduce a separation of our dataset into training and testing data

testing data: never used to train the NN
(to adjust the weights and biases)





test accuracy no longer improving
∴ overfitting

Tomorrow = methods that prevent overfitting

