

Title: Lecture - Machine Learning, PHYS 777

Speakers: Mohamed Hibat Allah

Collection/Series: Machine Learning (Elective), PHYS 777, February 24 - March 28, 2025

Subject: Condensed Matter, Other

Date: March 04, 2025 - 9:00 AM

URL: <https://pirsa.org/25030035>

Lecture 4

Last time:

- ↳ Architecture of feed-forward NNs for supervised learning (Weights $W_{ij}^{(l)}$, bias $b_i^{(l)}$ and activations)
- ↳ Expressivity of NNs.
- ↳ Cost/Loss functions.

Today:

- ↳ Learning algorithms.
- ↳ backpropagation
- ↳ Overfitting

Learning Algorithms

We use a learning algorithm (such as gradient descent) to minimize the cost function C with respect to all parameters.

①

$$\left. \begin{aligned} \frac{\partial C}{\partial w_{ij}^{(l)}} &\approx 0 \\ \frac{\partial C}{\partial b_j^{(l)}} &\approx 0 \end{aligned} \right\} \forall i, j, l$$

→ Unlike linear regression, as the # of parameters (weights & biases) increases, it quickly becomes impossible to solve ① analytically.

→ Recall gradient descent (GD)

$$\vec{p} \leftarrow \vec{p} - \eta \nabla_{\vec{p}} C \quad (\vec{p} = \text{vector of parameters}).$$

→ Today, we'll examine some alternative/improvements to GD.

↳ Cost/loss functions

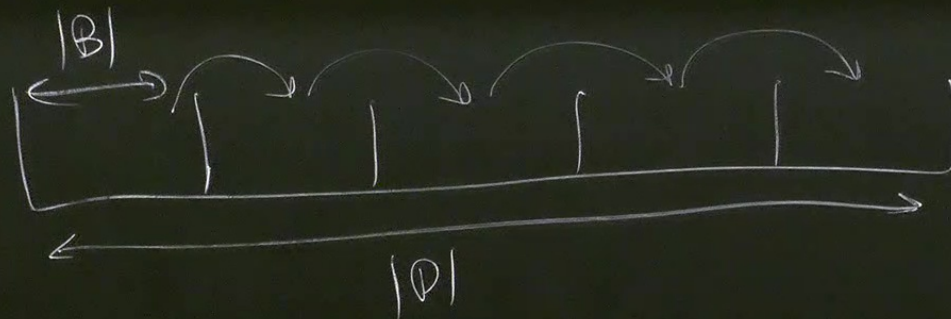
Mini-batch gradient descent:

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

↘

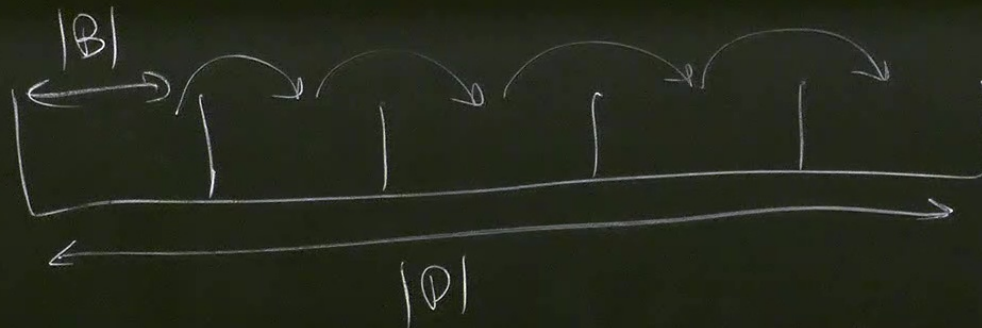
$$C \approx \frac{1}{|B|} \sum_{\vec{x} \in B} c(\vec{x})$$

$B \subset D$



Idea: approximate C and its gradient $\nabla_{\theta} C$ by summing over a
randomly chosen mini-batch B

$$\frac{\partial C}{\partial b_j} \approx 0$$

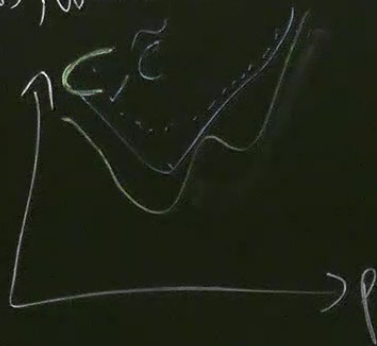


Idea: approximate C and its gradient $\nabla_{\theta} C$ by summing over a
randomly chosen mini-batch B

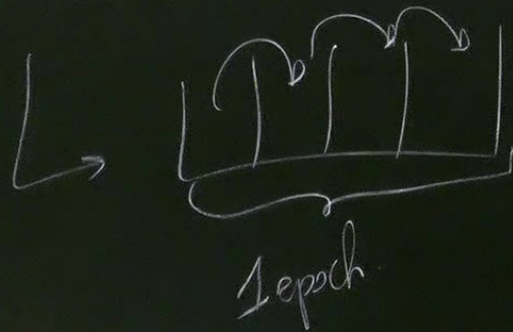
Benefits:

↳ Speeds up the calculation of the gradient / optimization / convergence.

↳ Introduces randomness, which can decrease the chance of getting stuck in a local min.



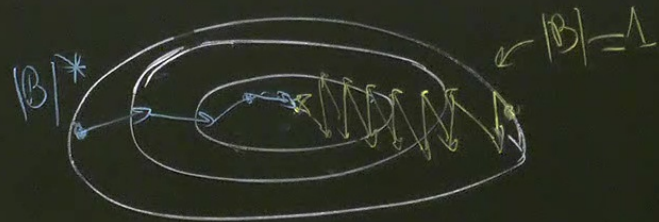
$|B|=1 \rightarrow$ Stochastic gradient descent (SGD)



B
 BCD

$\hookrightarrow |B|=1$ (a lot of noise)

$\hookrightarrow |B|^*$ in the middle for optimal optimization.





Epoch
000,000

Learning rate
0.03

Activation
Tanh

Regularization
None

Regularization rate
0

Problem type
Classification

DATA

Which dataset do you want to use?



Ratio of training to test data: 50%



Noise: 0



Batch size: 1



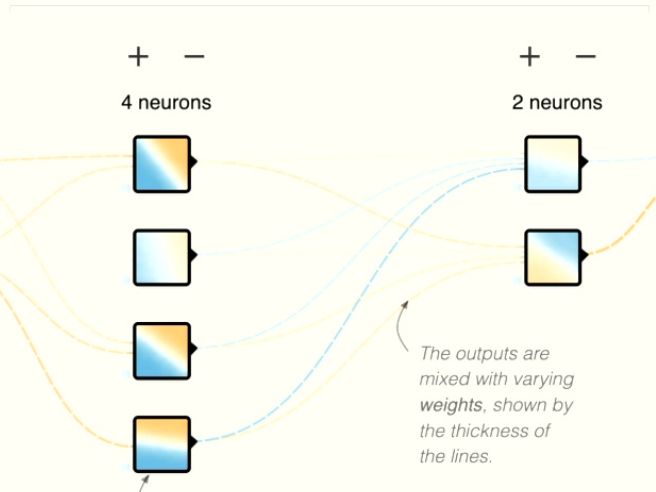
REGENERATE

FEATURES

Which properties do you want to feed in?

- X1
- X2
- X12
- X22
- X1X2
- sin(X1)
- sin(X2)

2 HIDDEN LAYERS

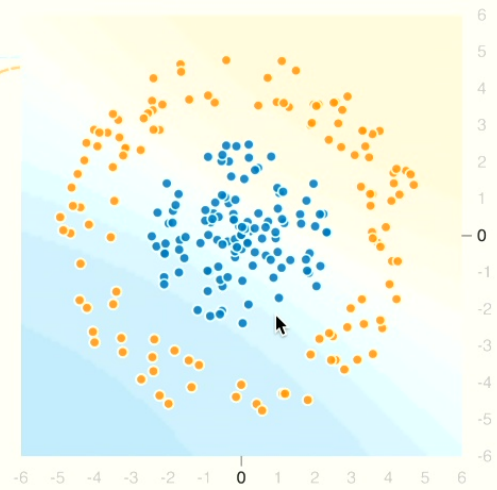


This is the output from one neuron. Hover to see it larger.

The outputs are mixed with varying weights, shown by the thickness of the lines.

OUTPUT

Test loss 0.492
Training loss 0.497



Colors shows data, neuron and weight values.



Epoch: 000,000
 Learning rate: 0.03
 Activation: Tanh
 Regularization: None
 Regularization rate: 0
 Problem type: Classification

DATA

Which dataset do you want to use?



Ratio of training to test data: 50%



Noise: 0



Batch size: 15



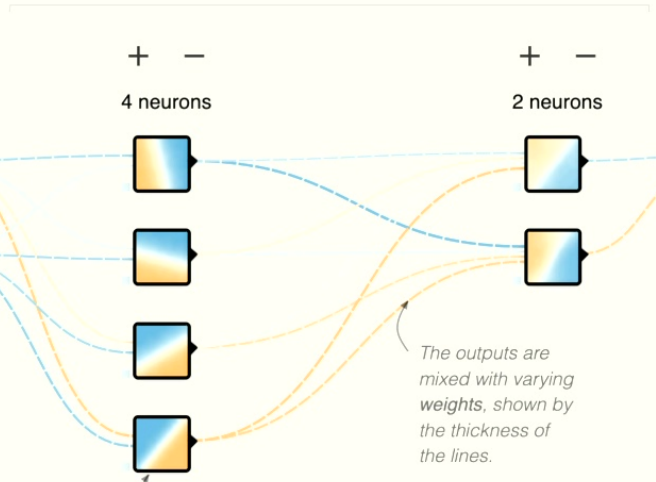
REGENERATE

FEATURES

Which properties do you want to feed in?

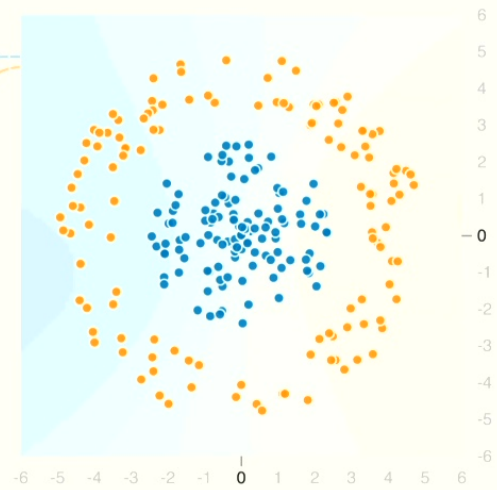
- X1
- X2
- X1²
- X2²
- X1X2
- sin(X1)
- sin(X2)

2 HIDDEN LAYERS



OUTPUT

Test loss 0.511
 Training loss 0.508



randomly chosen mini-batch B

Adding momentum:

$$\underbrace{\frac{d\vec{p}}{dt}}_{\text{Velocity}} = -\eta \nabla_{\vec{p}} C$$

randomly chosen mini-batch B

Adding momentum:

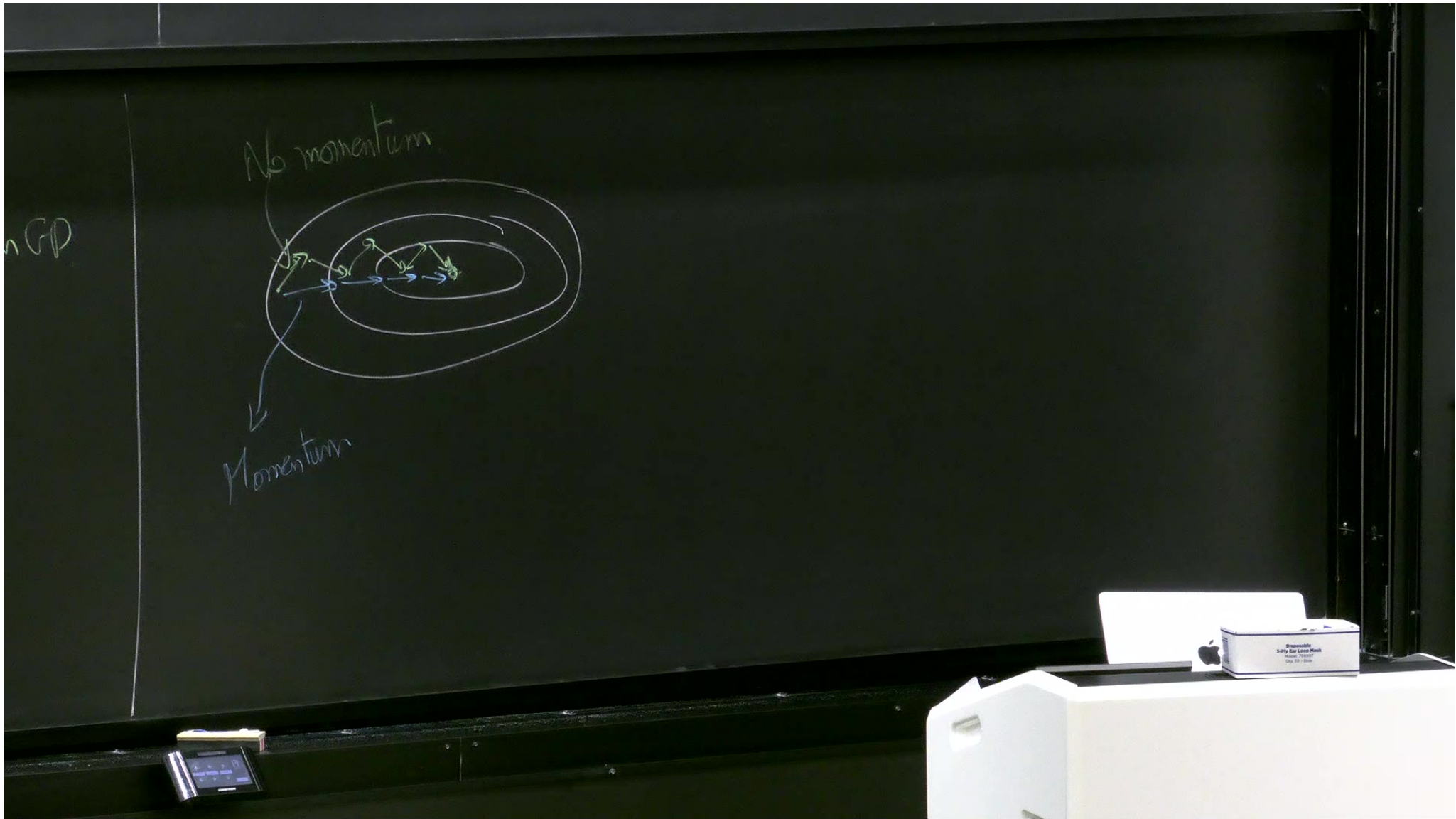
$$\underbrace{\frac{d\vec{p}}{dt}}_{\text{Velocity}} = -\eta \nabla_{\vec{p}} C$$

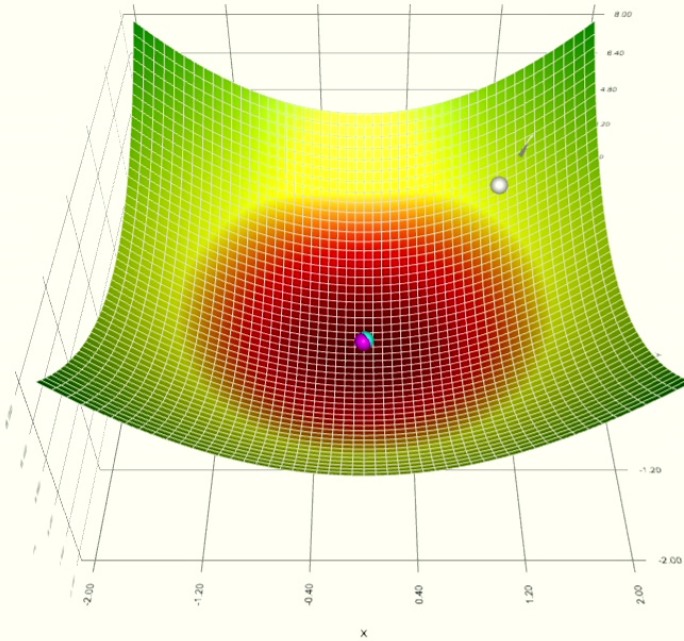
→

$$\frac{d\vec{v}}{dt} = -\nabla_{\vec{p}} C - \lambda \vec{v}$$
$$\frac{d\vec{p}}{dt} = v$$

$$\boxed{\Delta t = 1} \quad \vec{v}_{i+1} = \vec{v}_i + \nabla_{\vec{p}} C - \lambda \vec{v}_i$$
$$\vec{v}_{i+1} = \underbrace{(1 - \lambda)}_{\beta \text{ "momentum"}} \vec{v}_i + \nabla_{\vec{p}} C$$
$$\vec{p}_{i+1} = \vec{p}_i - \eta \vec{v}_i$$

$\boxed{\beta = 0}$ back to regular GP





Global Minimum

Overview Step-by-Step

- Gradient Arrows
- Adjusted Gradient Arrows
- Momentum Arrows
- Sum of Gradient Squared
- Path

Gradient Descent

Learning Rate: $1e^{-2}$

Momentum

Learning Rate: $1e^{-2}$

Decay rate: 0.800

Adagrad

Learning Rate: $1e^{-2}$

RMSprop

Learning Rate: $1e^{-2}$

Decay rate: 0.990

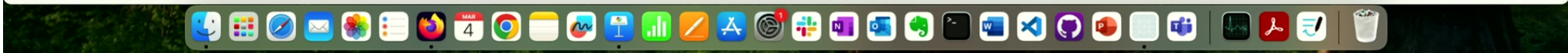
Adam

Learning Rate: $1e^{-2}$

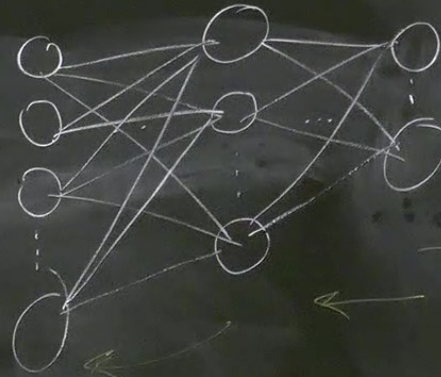
Beta1: 0.900

Beta2: 0.999

Pause Restart Playback speed: 0.2x



Lecture 4



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

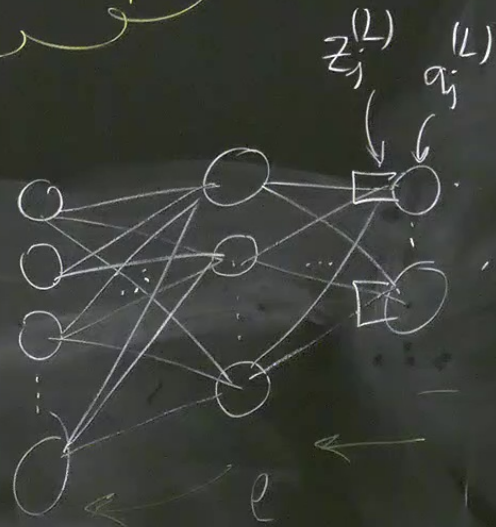
$$\frac{\partial c(\vec{x})}{\partial w_{ij}^{(l)}}, \frac{\partial c(\vec{x})}{\partial b_j^{(l)}}$$

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

$$z_j^{(l)} = \sum_{i=1}^{n_{l-1}} a_i^{(l-1)} w_{ij}^{(l)} + b_j^{(l)}$$

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

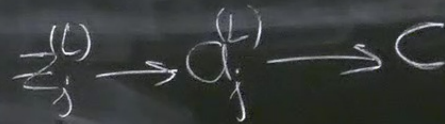
Lecture 4



$$\delta_j^{(L)} = \frac{\partial C}{\partial z_j^{(L)}}$$

$$= \frac{\partial C}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}}$$

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



→ for $l < L$

$$\delta_j^{(l)} = \frac{\partial C}{\partial z_j^{(l)}} = \sum_k \delta_k^{(l+1)} \left(W_{kj}^{(l+1)} \right)^T g'_j \left(z_j^{(l)} \right) \quad (2)$$

$$(3) \quad \frac{\partial C}{\partial W_{ij}^{(l)}} = a_i^{(l-1)} \delta_j^{(l)} \quad \frac{\partial C}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad (4)$$

(2) + (3) + (4) to be shown in HW #1.

Overfitting

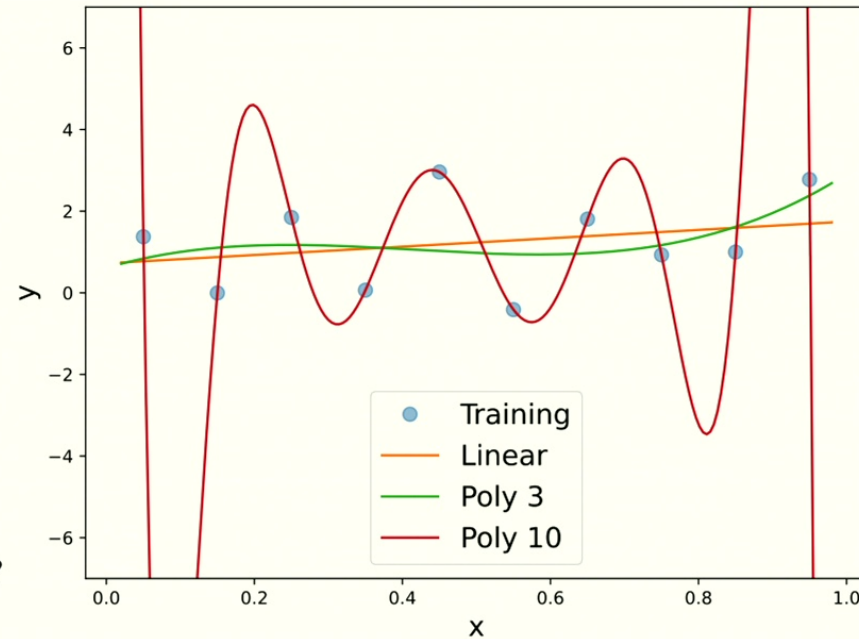
Training data:

```
N_train = 10
x = np.linspace(0.05,0.95,N_train)
s = np.random.randn(N_train)
y = 2*x+s
```

Models for fitting:

Polynomials of degree 1, 3, and 10

Question: Which model will give the best fit (lowest error) to the training data?



See Section II of arXiv:1803.08823 and the corresponding Notebook 1

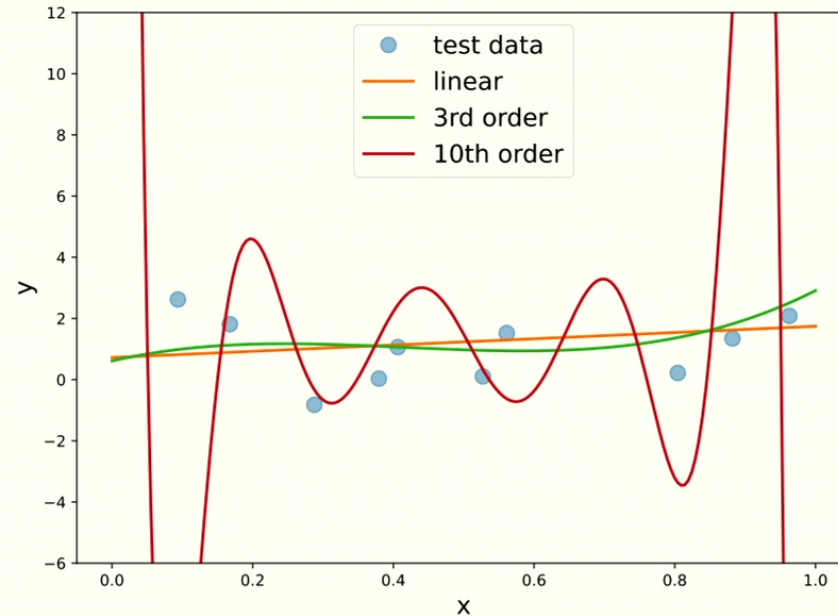
Overfitting

Testing data (from the same distribution as the training data):

```
N_test = 10
x_test = np.random.random(N_test)
s_test = np.random.randn(N_test)
y_test = 2*x_test+s_test
```

Questions:

- Which model from the previous slide will make the best predictions on the testing data?
- What would happen if we increased the number of data points in the training dataset?



See Section II of arXiv:1803.08823 and the corresponding Notebook 1

Don't Worry, You Can't Break It. We Promise.

Epoch 000,000 Learning rate 0.03 Activation Tanh Regularization None Regularization rate 0 Problem type Classification

DATA

Which dataset do you want to use?



Ratio of training to test data: 20%

Noise: 35

Batch size: 10

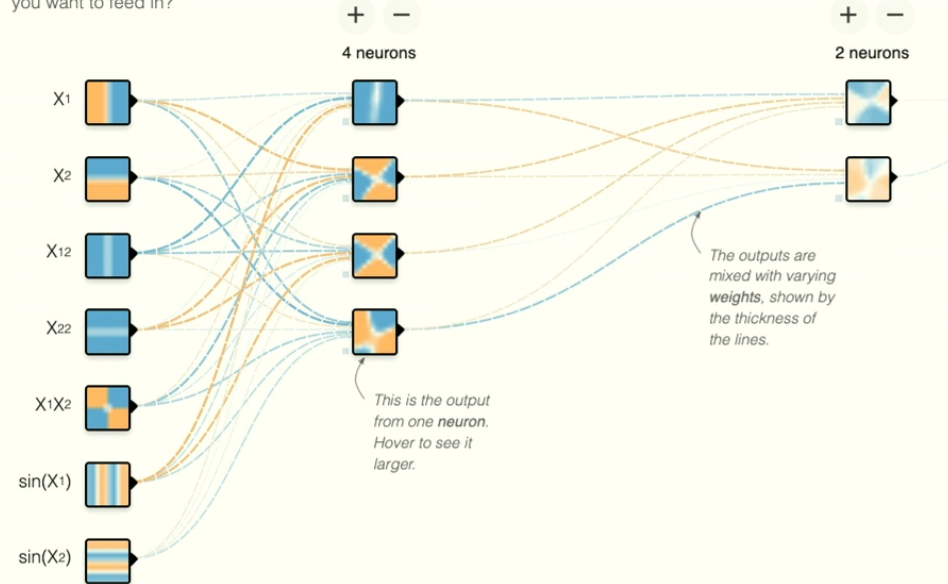
REGENERATE

FEATURES

Which properties do you want to feed in?

- X1
- X2
- X1²
- X2²
- X1X2
- sin(X1)
- sin(X2)

2 HIDDEN LAYERS



OUTPUT

Test loss 0.489
Training loss 0.486



Colors shows data, neuron and weight values. -1 0 1

Don't Worry, You Can't Break It. We Promise.

Epoch 003,116
Learning rate 0.03
Activation Tanh
Regularization None
Regularization rate 0
Problem type Classification

DATA
Which dataset do you want to use?
Ratio of training to test data: 20%
Noise: 35
Batch size: 10
REGENERATE

FEATURES
Which properties do you want to feed in?
X1
X2
X12
X22
X1X2
sin(X1)
sin(X2)

