

Lab course Image Analysis I ST 2023

Hubert Kanyamahanga
(kanyamahanga@ipi.uni-hannover.de)

*Lab 2: Discriminative Probabilistic Models
Logistic Regression & Pytorch Basics*



Content of the 2nd Lab

Goal:

- Application of Logistic regression for **classification** purposes
- Introduction to Neural Networks
- Fundamental concepts about how **Pytorch** works

Tasks:

1. Implement and train Logistic Regression classifiers by
 - ✓ first, using **Scikit-learn** and then **Pytorch** Frameworks
2. Implement and train Neural Networks
 - ✓ apply them to synthetic datasets and existing benchmark
3. Discussion and evaluation (**Visually** and **Quantitatively**)

Recall from the 1st Lab

- **Generative classifiers:**

- Wanted to compute $p(C|\mathbf{x})$ with the help of the theorem of Bayes $p(C|\mathbf{x}) = \frac{p(\mathbf{x}|C) \cdot p(C)}{p(\mathbf{x})}$
- Determine the parameters of the likelihood $p(\mathbf{x}|C)$ (training)
- Determine the prior $p(C)$
- ‘**Generative**’: Generate synthetic data sets by sampling from the joint distribution

$$p(\mathbf{x}, C) = p(\mathbf{x}|C) \cdot p(C)$$

Now in the 2nd Lab

- **Discriminative classifiers:**

- Direct modelling of $p(C|\mathbf{x})$
- Focus on the separating surfaces in feature space
- In general, this leads to simpler models and, therefore, requires fewer training samples

Overview: Lab 2

- Discriminative probabilistic classifiers:
 - Logistic Regression
 - Training
- Neural networks:
 - Model
 - Training
 - Hyper-parameters
- Frameworks
 - Overview



Overview: Lab 2

- Discriminative probabilistic classifiers:
 - Logistic Regression
 - Training
- Neural networks:
 - Model
 - Training
 - Hyper-parameters
- Frameworks
 - Overview



1. Binary classification: Logistic Sigmoid Function

- **Binary classification** with two classes L^1, L^2 (*object, background*)
- We start with a **Bayesian view** and express the posterior for L^1 using the **Theorem of Bayes**:

$$p(C=L^1|\mathbf{x}) = \frac{p(\mathbf{x}|C=L^1) \cdot p(C=L^1)}{p(\mathbf{x}|C=L^1) \cdot p(C=L^1) + p(\mathbf{x}|C=L^2) \cdot p(C=L^2)} = \frac{1}{1 + \frac{p(\mathbf{x}|C=L^2) \cdot p(C=L^2)}{p(\mathbf{x}|C=L^1) \cdot p(C=L^1)}} = \frac{1}{1 + e^{-a}}$$

with $a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C = L^1) \cdot p(C = L^1)}{p(\mathbf{x} | C = L^2) \cdot p(C = L^2)} = \ln \frac{p(C = L^1 | \mathbf{x})}{p(C = L^2 | \mathbf{x})}$

- **Result:** Formally, the posterior $p(C=L^1|\mathbf{x})$ can be expressed using the **logistic sigmoid function** $\sigma(a) = \frac{1}{1 + e^{-a}}$
- **Assumptions:** $\Sigma_1 = \Sigma_2 = \Sigma$ (Lecture slide 8 ch.8), $a(\mathbf{x})$ is a linear function of the features!

$$p(C=L^1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \cdot \mathbf{x} + w_0)}} = \sigma(a(\mathbf{x})) = \sigma(\mathbf{w}^T \cdot \mathbf{x} + w_0)$$

```
#Pytorch: Binary task with 2 fts x0 and x1
torch.sigmoid(w_0*X[:,0] + w_1*X[:,1] + b)
```



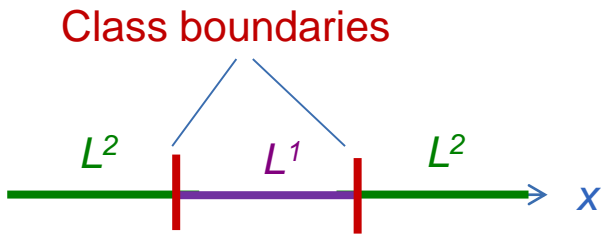
Complex models: Feature Space Mapping

- General equation for the posterior: $p(C=L^1|\mathbf{x}) = \sigma[a(\mathbf{x})] = \frac{1}{1 + e^{-a(\mathbf{x})}}$
- For **identical covariance matrices**, $a(\mathbf{x})$ was a linear function of the features \mathbf{x} :
$$a(\mathbf{x}) = \mathbf{w}^T \cdot \mathbf{x} + w_0$$
- **Problem**: Increasing complexity of the models for the probability densities $\rightarrow a(\mathbf{x})$ becomes more complex, too, e.g. a **quadratic form** for normal distributions
- Rather than using more complex models for probability distributions, **we expand the feature space**
 - \rightarrow Feature Space Mapping $\Phi(\mathbf{x}) = [\Phi_1(\mathbf{x}), \Phi_2(\mathbf{x}), \dots, \Phi_N(\mathbf{x})]^T$
- In the expanded feature space, we can still work with linear models for the exponent
 - Example for 2D feature space, i.e. $\mathbf{x} = (x_1, x_2)^T$: $\Phi(\mathbf{x}) = (1, x_1, x_2, x_1 \cdot x_2, x_1^2, x_2^2)^T$
 - \rightarrow “quadratic expansion” where $\Phi(\mathbf{x})$ is frequently a polynomial function

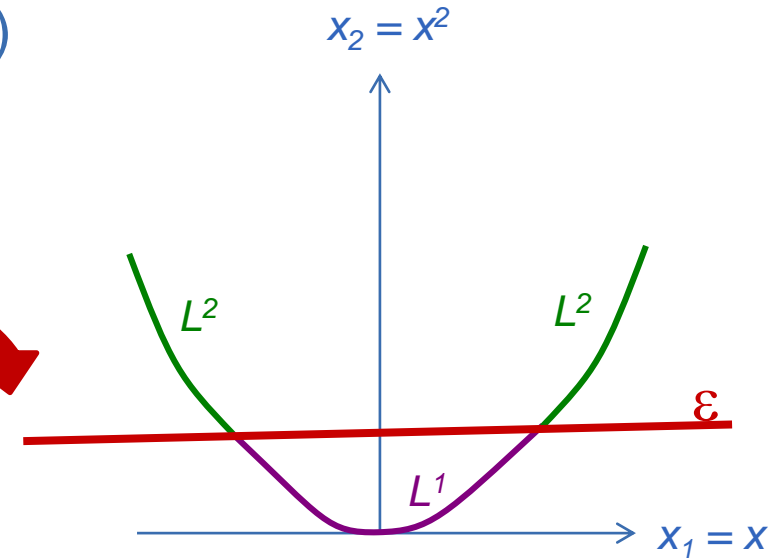
Example: Not linearly separable

- Transition to a higher-dimensional feature vector $\Phi(\mathbf{x})$
- Example

Feature space transformation



$$\Phi(\mathbf{x}) = \begin{pmatrix} x \\ x^2 \end{pmatrix}$$



1D feature space (x), 2 classes
Not linearly separable

After feature space transformation:
2D feature space (x, x^2)
Classes can be separated by a plane ϵ

Weights: $[[0.549 \quad 0.603] \leftarrow L^1$
 $[0.549 \quad 0.603] \leftarrow L^2$

Weights: $[[0.549 \quad 0.715 \quad 0.603] \leftarrow L^1$
 $[0.549 \quad 0.715 \quad 0.603] \leftarrow L^2$

2. Multi-class classification: Softmax Function

- The posterior $p(C = L^k | \mathbf{x})$ for each class L^k is modelled using the **softmax function**:

$$p(C=L^k|\mathbf{x}) = \frac{\exp [a_k(\mathbf{x})]}{\sum_j \exp [a_j(\mathbf{x})]}$$

with $a_k(\mathbf{x}) = \ln [p(\mathbf{x}|C=L^k)] + \ln [p(C=L^k)]$

- Assumptions about $p(\mathbf{x}|C=L^k)$ and $p(C=L^k)$ lead to models for $a_k(\mathbf{x})$
- Again, feature space mapping can help to obtain linear models:
 $a_k(\mathbf{x}) = a_k(\Phi(\mathbf{x})) = \mathbf{w}_k^\top \cdot \Phi(\mathbf{x})$
- In training, one parameter vector \mathbf{w}_k per class has to be determined



Overview: Lab 2

- Discriminative probabilistic classifiers:
 - Logistic Regression
 - Training
- Neural networks:
 - Model
 - Training
 - Hyper-parameters
- Frameworks
 - Overview

Multi-class classification: Training with ML

- Question: Which class label does a feature vector \mathbf{x} belong to?

- Wanted: How to compute posteriors:

$$p(C = L^k | \mathbf{x}_n) = \frac{\exp[\mathbf{w}_k^T \cdot \Phi(\mathbf{x}_n)]}{\sum_{j=1}^M \exp[\mathbf{w}_j^T \cdot \Phi(\mathbf{x}_n)]} = y_{nk}$$

- Given: Class labels C_n for each training sample \mathbf{x}_n

- Approach: Maximum Likelihood Training

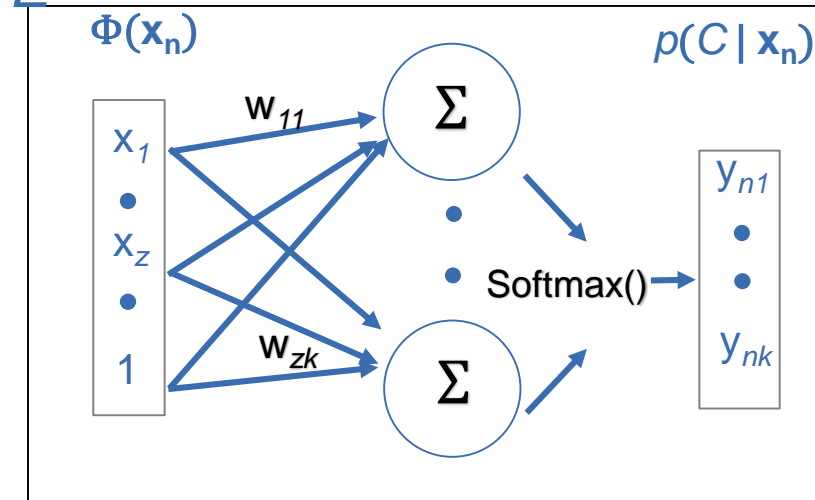
→ Similar to binary case minimize the negative log-likelihood E

$$E(\mathbf{w}_1, \dots, \mathbf{w}_M) = -\sum_{n=1}^N \sum_{k=1}^M t_{nk} \cdot \ln(y_{nk}) \rightarrow \min$$

with binary indicator variables $t_{nk} = \begin{cases} 1 & \text{if} \\ 0 & \text{otherwise} \end{cases} \quad C_n = L^k$

N: Number of samples

M: Number of classes



Training (cont'): Parameters Estimation

- **Procedure:** Compute the posteriors and take the class label to be the index of the maximum (`np.argmax(posteriors, axis=1)`)

1. Compute the Gradient $\nabla E(\mathbf{w}^{\tau-1})$:
$$\nabla_{\mathbf{w}_j} E(\mathbf{w}) = \sum_{n=1}^N (y_{nj} - t_{nj}) \cdot \Phi(\mathbf{x}_n)$$

2. Compute the Hessian Matrix $\nabla \nabla E(\mathbf{w}^{\tau-1}) = \mathbf{H}(\mathbf{w}^{\tau-1})$:

$$\mathbf{H}_{jk} = \nabla_{\mathbf{w}_j} \nabla_{\mathbf{w}_k} E(\mathbf{w}) = \sum_{n=1}^N y_{nk} \cdot (\mathbf{I}_{kj} - y_{nj}) \cdot \Phi(\mathbf{x}_n) \cdot \Phi(\mathbf{x}_n)^T$$

3. Update the weights: Newton-Raphson

\mathbf{I}_{jk} ... Elements of a unit matrix

$$\mathbf{w}^{\tau} = \mathbf{w}^{\tau-1} - \mathbf{H}^{-1} \cdot \nabla E(\mathbf{w}^{\tau-1}) \quad (\text{Ref. Ch. 8 Slide 11})$$

4. Increase iteration count τ and repeat until convergence and

5. Calculate posteriors:
$$p(C = L^k | \mathbf{x}_n) = \frac{\exp[\mathbf{w}_k^T \cdot \Phi(\mathbf{x}_n)]}{\sum_{j=1}^M \exp[\mathbf{w}_j^T \cdot \Phi(\mathbf{x}_n)]} = y_{nk}$$



Maximum Likelihood Training: Problem of Overfitting

- Maximum Likelihood training has the tendency to overfit the classifier to the training data
→ regularization with a prior for \mathbf{w}
- Solution to keep the numerical values of \mathbf{w} small:
 - Gaussian Prior $p(\mathbf{w})$ with expectation value $\mathbf{0}$ and Covariance Matrix $\sigma^2 \cdot \mathbf{I}$
- Requires hyper-parameter σ which is **either fixed by the user** or determined via a procedure such as **cross-validation**
- **Regularisation** (Gaussian prior) with exp. $\mathbf{0}$ and Covariance $\sigma \cdot \mathbf{I}$
- Negative logarithm of $p(\mathbf{w} \mid \mathbf{t}, \mathbf{x}_1, \dots, \mathbf{x}_N)$ (excluding constant terms):

$$E(\mathbf{w}) = -\sum_{n=1}^N \left[t_n \cdot \ln(y_n) + (1 - t_n) \cdot \ln(1 - y_n) \right] + \frac{\mathbf{w}^T \cdot \mathbf{w}}{2 \cdot \sigma^2} \rightarrow \min$$

Maximum Likelihood Training: Solution- Regularization

- Minimization of

$$E(\mathbf{w}) = -\sum_{n=1}^N \left[t_n \cdot \ln(y_n) + (1 - t_n) \cdot \ln(1 - y_n) \right] + \frac{\mathbf{w}^T \cdot \mathbf{w}}{2 \cdot \sigma^2} \rightarrow \min$$

leads to the numerical values of \mathbf{w} that are as small as possible (indicated by σ)

- Gradient has to be extended compared to the ML method:

$$\nabla E(\mathbf{w}_1, \dots, \mathbf{w}_M) = \left[\nabla_{\mathbf{w}_1} E(\mathbf{w}_1, \dots, \mathbf{w}_M)^T, \dots, \nabla_{\mathbf{w}_M} E(\mathbf{w}_1, \dots, \mathbf{w}_M)^T \right]^T + \frac{1}{\sigma^2} \cdot \mathbf{w}$$

- This is also true for the Hesse Matrix:

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \left[y_n \cdot (1 - y_n) \cdot \Phi(\mathbf{x}_n) \cdot \Phi(\mathbf{x}_n)^T \right] + \frac{1}{\sigma^2} \cdot \mathbf{I}$$

i.e. in the main diagonal, the weights of the direct observations for \mathbf{w} are added

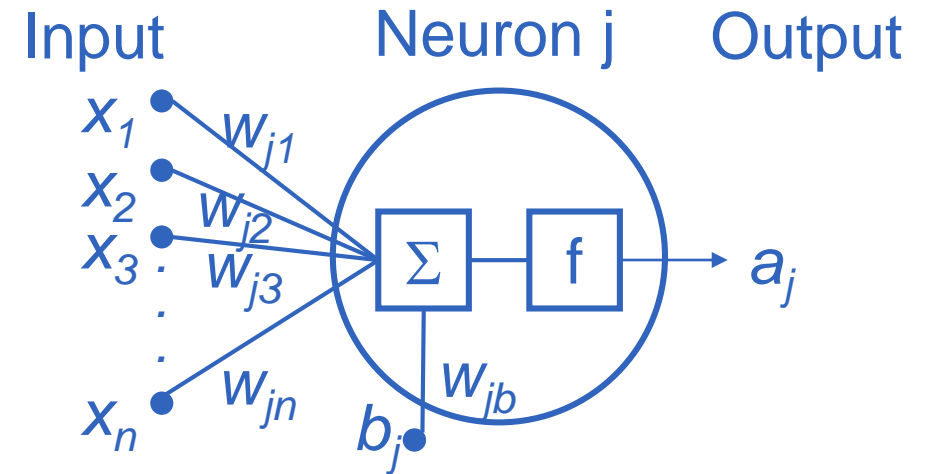
Overview: Lab 2

- Discriminative probabilistic classifiers:
 - Logistic Regression
 - Training
- Neural networks:
 - Model
 - Training
 - Hyper-parameters
- Frameworks
 - Overview



Artificial Model of a Neuron

- **Input variables** x_i : Components of the feature vector \underline{x}
- Each input variable is multiplied with a weight;
Determine **weighted sum** $z_j = \sum \underline{w}_{ji} \cdot \underline{x}_i + b_j = \underline{w}_j^T \cdot \underline{x} + b_j$
- b_j : **Bias**, considered to be a component of each feature vector
→ $\underline{x} = [\underline{x}^T \ 1]^T$ and $\underline{w}_j = [\underline{w}_j^T \ b_j]^T$
→ E.g. **For 1D feat. Space X**, $\underline{W} : [[0.549 \ 0.603]]$
→ Simplified notation : $l_j = \underline{w}_j^T \cdot \underline{x}$
- Output a_j of the neuron j :
 $a_j = f(l_j) = f(\underline{w}_j^T \cdot \underline{x})$ with $f(l_j)$... activation function



Linear binary classifier based on a **Single neuron** → **Perceptron, logistic regression**

Neural Networks

- Networks consisting of **several layers** of neurons
→ More complex decision boundaries possible
- Example: **two layers** and “**feed forward**” - architecture

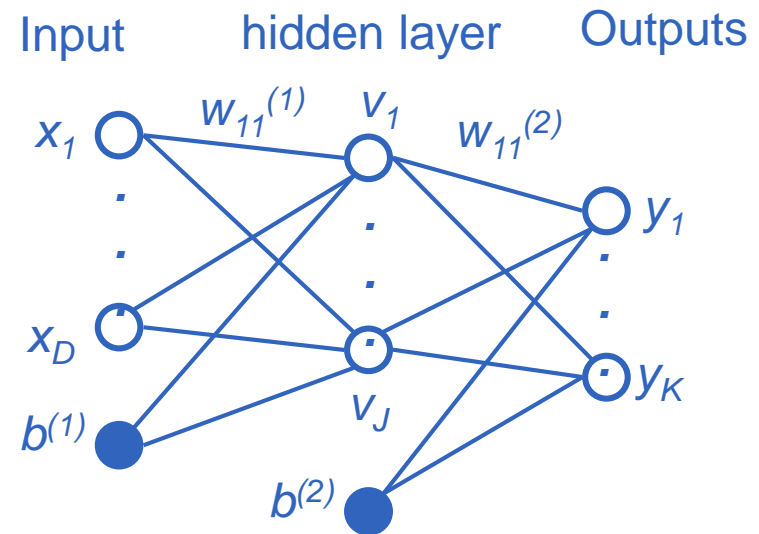
- **Input:** features x_i

- **Hidden layer** with neurons v_j :

$$v_j = f(\sum w_{ji}^{(1)} \cdot x_i)$$

- **Output layer:** degree of membership to class C^k

$$y_k = f(\sum w_{ki}^{(2)} \cdot v_i)$$



→ **Multilayer Perceptron (MLP)**

Overview: Lab 2

- Discriminative probabilistic classifiers:
 - Logistic Regression
 - Training
- Neural networks:
 - Model
 - Training
 - Hyper-parameters
- Frameworks
 - Overview



Neural Networks: Training with Gradient Descent

- **Goal:** Minimizing the Loss Function → Optimizer: Gradient Descent

- **Procedure:**

1. Compute the **posteriors** : $y_n = f(\mathbf{w}^T \cdot \mathbf{x} + w_b) \rightarrow$ Sigmoid function

2. Compute the **BCELoss** : $E(\mathbf{w}) = -\sum_{n=1}^N [t_n \cdot \ln(y_n) + (1 - t_n) \cdot \ln(1 - y_n)] \rightarrow \min$

3. Compute the **gradients** : $\nabla E(\mathbf{w})$, with respect to weights and bias

4. Update the weights: **Gradient Descent** [Bishop, 2006]

$$\mathbf{w}^\tau = \mathbf{w}^{\tau-1} - \mathbf{H}^{-1} \cdot \nabla E(\mathbf{w}^{\tau-1}) \longrightarrow \mathbf{w}^\tau = \mathbf{w}^{\tau-1} - \eta \cdot \nabla E(\mathbf{w}^{\tau-1}) \quad (\text{ref. slide 27})$$

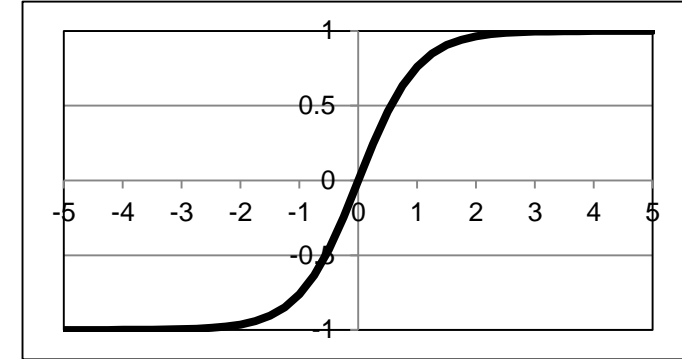
5. Run the training for a specific number of iterations/epochs until convergence

Activation Functions (non-linearity)

- Tangens Hyperbolicus (TanH):

- Output is centered at 0.0
- Small gradients for large / small inputs

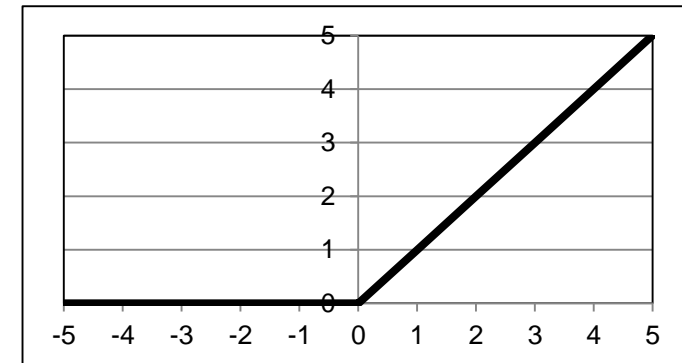
→ Slow convergence in training!



$$f(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

- Rectified Linear Unit (ReLU):

- Prevents vanishing gradients
- Neurons can 'die'
- Mostly used for Deep Networks

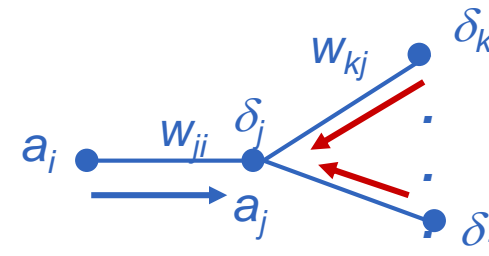


$$f(a) = \max(0, a)$$

Optimization: Gradient Descent

- Back-propagation for computing the gradients:
 - Forward step:
 - Calculate output y_{nk} from \mathbf{x}_n and the current values of \mathbf{w}
 - Save the output a_j as well as $f'(l_j)$ in every neuron j
 - The classification error and δ_k is calculated from y_{nk}
 - Actual back-propagation:
 - δ_j is calculated from δ_k and $f'(l_j)$ successively for each layer from δ_j and a_j :

$$\frac{\partial E_n}{\partial w_{ji}} = \delta_j \cdot a_j$$



Stochastic Gradient Descent

- Stochastic Gradient Descent (SGD):
 - Gradients are averaged over Mini-Batch with B samples:

$$w'_{ji} = \frac{1}{B} \cdot \sum_{b=1}^B \frac{\partial E_{nb}}{\partial w_{jib}}$$

- Weight update with learning rate η

$$w_{ji}^{(\tau+1)} = w_{ji}^{(\tau)} - \eta \cdot w'_{ji}$$

```
# Define the optimizer, here SGD: Stochastic Gradient Descent
optimizer = optim.SGD(net.parameters(), lr=Learning_rate)
```

Overview: Lab 2

- Discriminative probabilistic classifiers:
 - Logistic Regression
 - Training
- Neural networks:
 - Model
 - Training
 - Hyper-parameters
- Frameworks
 - Overview



Neural Networks: Hyper-parameters

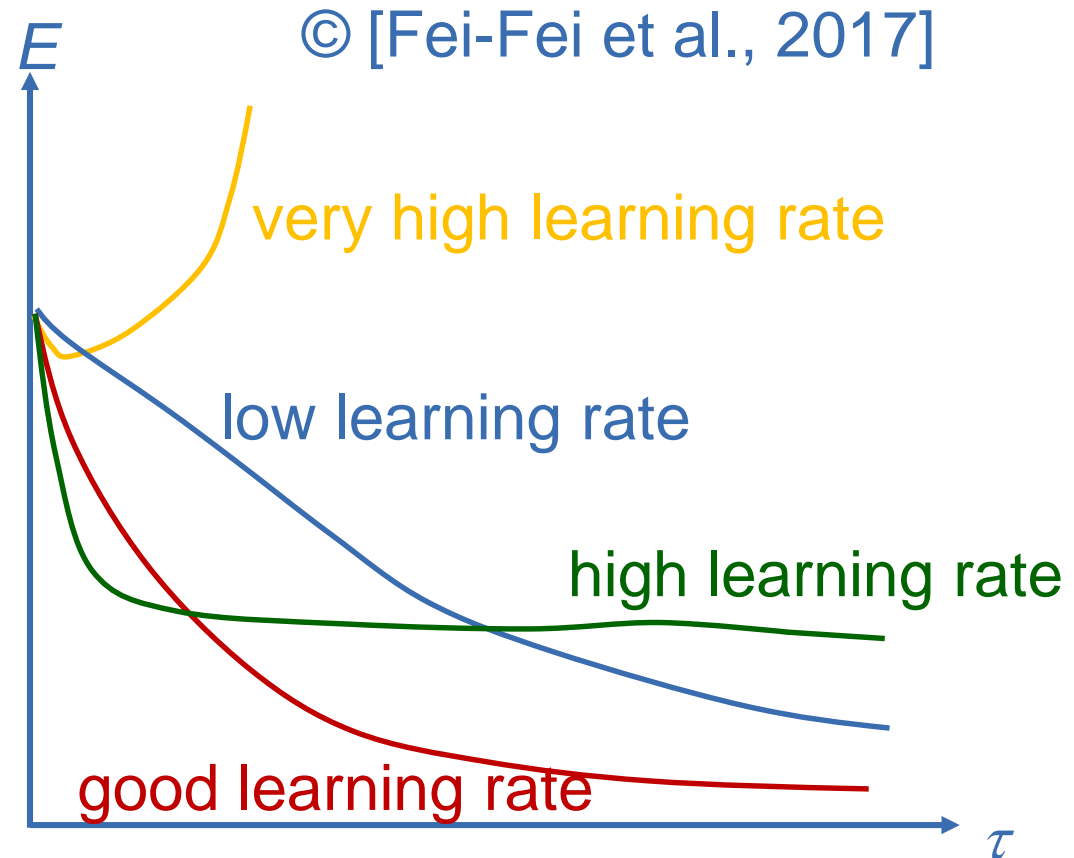
- Network architecture
 - Depth: Number of layers
 - Width: Number neurons per layer
 - Nonlinearity / activation function
 - ...
- Training
 - Optimizer (SGD / Momentum / Adam / RMSProp, ...)
 - Learning rate
 - Number of iterations
 - Mini-Batch size
 - ...

Training Neural Networks: Learning rate

- Learning rate η in gradient descent is an important hyper-parameter
- Needs to be tuned carefully!

- Good η leads to ...
 - Fast convergence
 - Strong minimum of E
- Adapt η in the iteration process
- Example:
exponential decay with small ε

$$\eta = \eta_0 \cdot (1 - \varepsilon)^{k \cdot \tau}$$



Overview: Lab 2

- Discriminative probabilistic classifiers:
 - Logistic Regression
 - Training
- Neural networks:
 - Model
 - Training
 - Hyper-parameters
- Frameworks
 - Overview



Frameworks overview: Pytorch

- Python package for machine learning, backed by Facebook
 - As of February 2022, PyTorch is the most used deep learning framework on Papers With Code
 - PyTorch also helps take care of many things such as GPU acceleration (making your code run faster) behind the scenes
- Pytorch has Five essential modules for building NN:

Module	What does it do?
<i>torch.Tensor</i>	Equivalent to Numpy arrays and matrices
<i>torch.autograd</i>	Automatic differentiation for all operations on tensors
<u><i>torch.nn</i></u>	Building Neural Networks
<u><i>torch.optim</i></u>	Various optimizations algorithms
<u><i>torch.utils.data(.Dataset, DataLoader)</i></u>	Data structures + transformation (augmentation)

Submission of Results

- Submission deadline: **July 3rd, 2023 before 11:00 am**
- **Assignment** → Jupyter Notebook (only digital)
 - Use meaningful variable names
 - Write comments if required
 - Answer concisely but completely
 - Consider acceptance rules
 - Consider [IAI_23_Lab_Technical_Details.pdf](#) about file naming scheme
- **Tutorial: Neural Networks: June 19th, 2023 at 11:30 am**
- Introduction to the 3rd Lab: **July 3rd, 2023 at 1:00 pm**

