

# Deep Learning for NLP

Universität Bielefeld

## Lecture 2 – ML principles

**Dr. Steffen Eger**

**[seminar.dldh@gmail.com](mailto:seminar.dldh@gmail.com)**



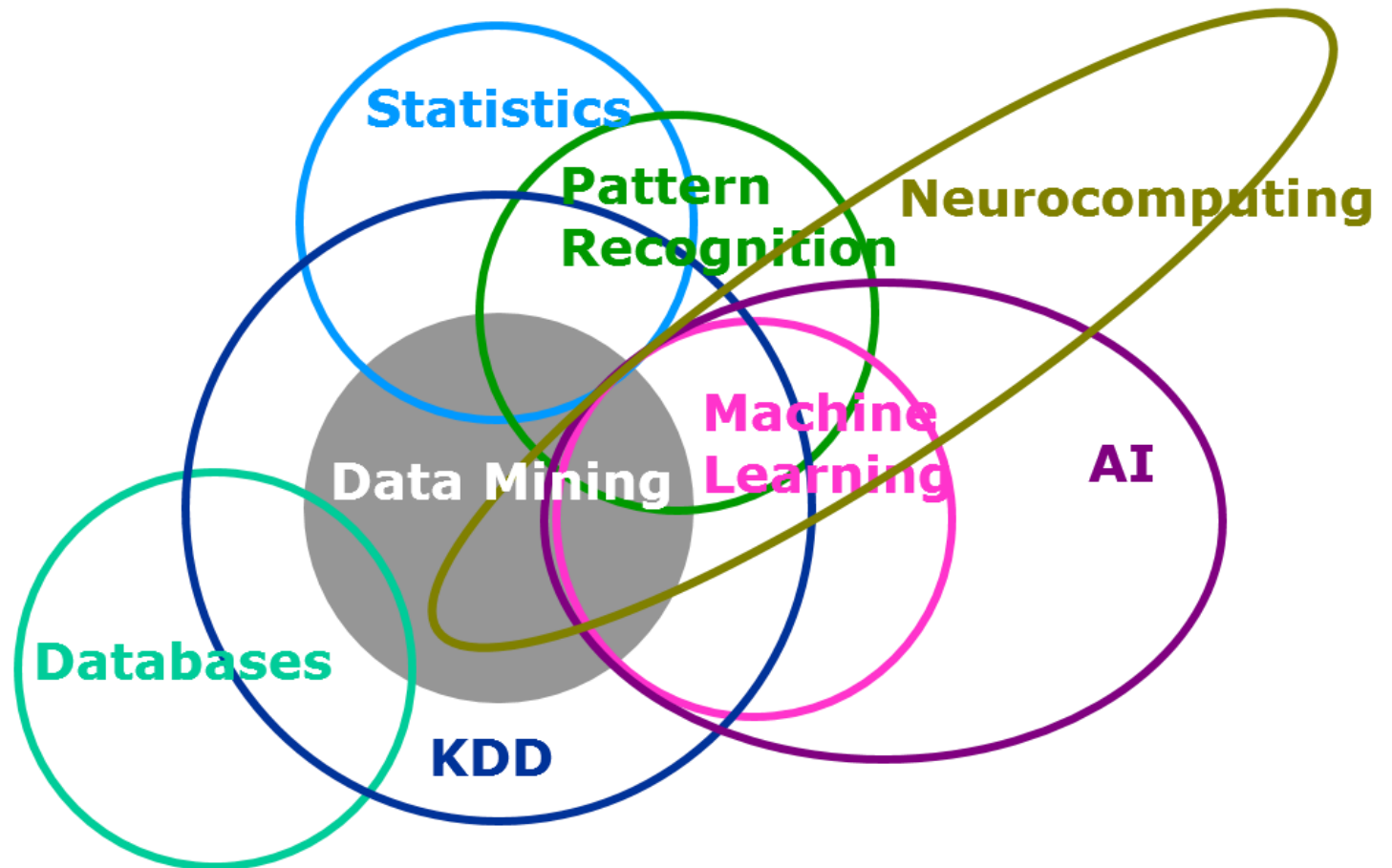
# Previous and upcoming lectures

- Previous lecture:
  - (Historical) Survey of Deep Learning
- Next lecture:
  - Training Networks with Hidden Layers (Backpropagation algorithm)
  - Language Model

# This lecture:

- Machine Learning Principles
  - Train/dev/test split
  - Loss functions
  - Evaluation
- **Learning goals:**
  - Understand ML & DL foundations

**ML principles**



# Standard Setup (Supervised setting)

- We have **data**
  - $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

Alternative Notation:  $(\mathbf{x}_i, t_i)$

$\mathbf{x}_i$  or  $(\mathbf{x}_i, t_i)$  is also called *instance*  
 $t_i$  is called truth / gold label

- And a statistical model  $f_{\theta}(\mathbf{x})$
- We also specify some *loss* function, e.g.,

Outputs/predictions of model  
are denoted as  $y$  or  $\hat{y}$

$$\bullet \frac{1}{N} \cdot \sum_{(\mathbf{x}, y)} (y - f_{\theta}(\mathbf{x}))^2$$

**Goal** is then to optimize/minimize our loss over  $\theta$ :

We are looking for parameters that bring our model close to the data

$$\text{Loss: } \frac{1}{N} \cdot \sum_{(x,y)} (y - f_{\theta}(x))^2$$

- Goal is then to optimize/minimize our loss over  $\theta$
- HOWEVER: our – real - goal is NOT to fit the **given** data well (which is called **overfitting**)



If data is non-pathological, we'll always find a model  $f_{\theta}$  that perfectly fits it



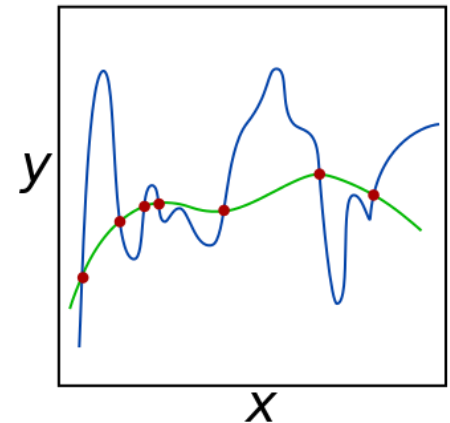
- The (real) goal in ML is to **generalize** particularly well to unseen data
- Hence, our (real) goal is more to minimize the expected loss rather than the actually observed loss:

$$E \left[ \sum (y - f_{\theta}(x))^2 \right] \quad (*)$$

where the expectation is over the distribution of datapoints:  $(x, y) \sim D$

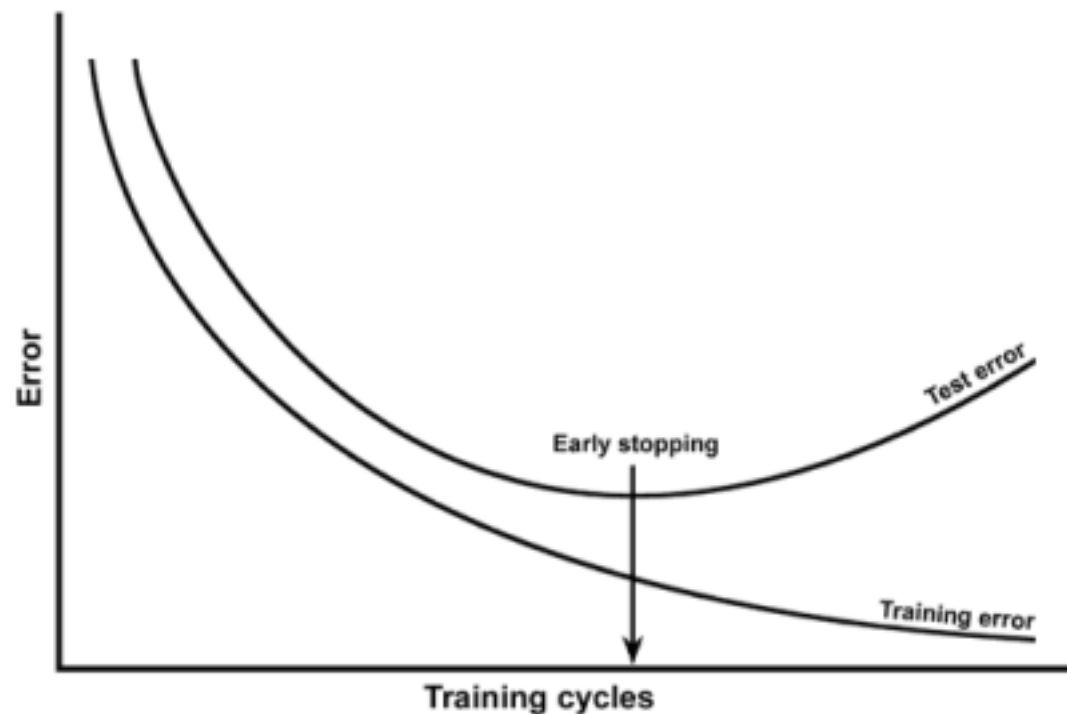
- Because we're interested in generalization performance, **we always split our data**:
  - Test data vs. Train data
  - Test data represents the true underlying distribution
- Model must perform well on the test data; performance on the training data is not of interest

- To achieve generalization, one usually introduces *regularization*
  - a term added to the loss function that constrains the parameters to be small, for example
  - Intuition: use a smaller polynomial to fit your data (less expressive function) if possible



- We further split it into development+(proper) train data
- On the proper training data, we optimize the **parameters  $\theta$**  of our model  $f_{\theta}$
- On the dev(elopment) data, we optimize the **hyperparameters** of our model
  - E.g. Learning rate, regularization terms/coefficients, number of epochs, etc.
  - In particular: can use dev data for **early stopping**, etc. (extremely common and popular)

In particular: can use dev data for **early stopping**, etc.  
(extremely common and popular)



- Could choose some fraction of your overall data for testing, e.g., 75%/25% split (train/test) or 90%/10%

*Why can a 90%-10% train/test split be dangerous?*

*A: Too little test data*

*B: Too much train data*

*C: May accidentally pick hard test cases*

*D: May accidentally pick easy test cases*

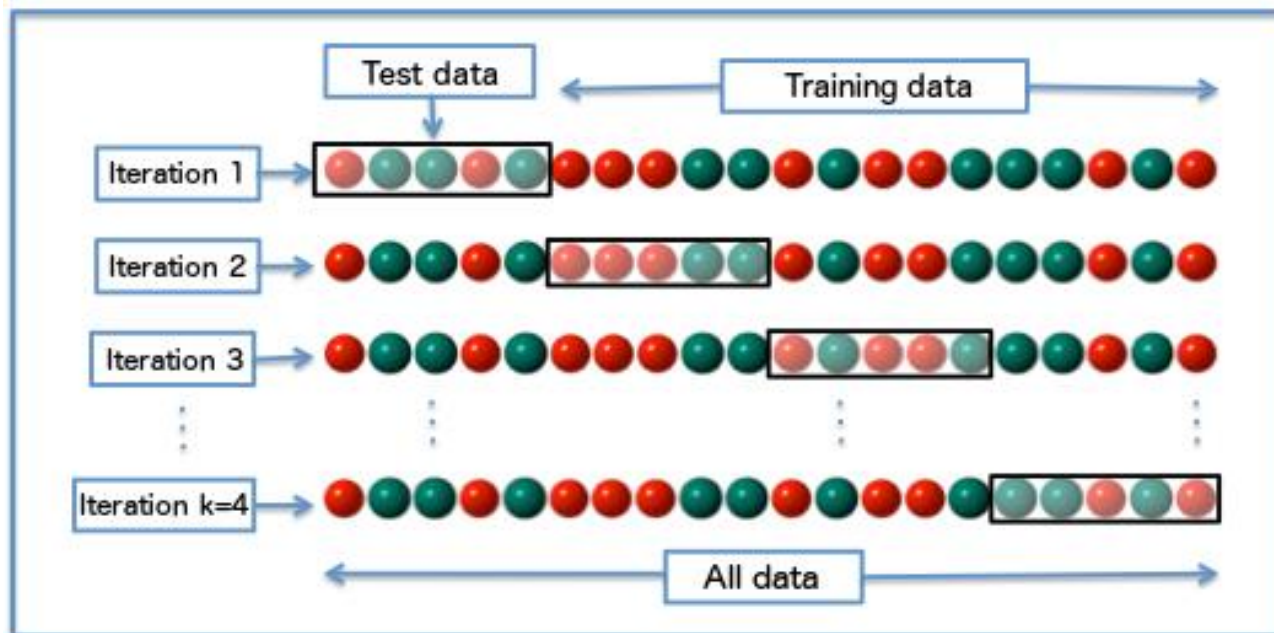
*E: Performance estimate is unreliable*



- Could choose some fraction of your overall data for testing, e.g., 75%/25% split (train/test) or 90%/10%
- Better is k-fold cross-validation:
  - Train on  $(k-1)$  equally sized folds, test on the remaining
  - Repeat  $k$  times



# Test data

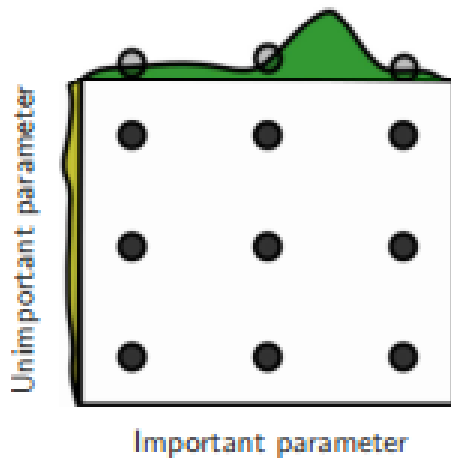


- This way, estimate of true performance is much more unbiased
- However, also computationally more costly
- Commonly used methods: **10-fold cross-validation**, **leave-one-out-validation**
- Other strategies:
  - e.g. *random subsample validation* in which test folds are randomly selected

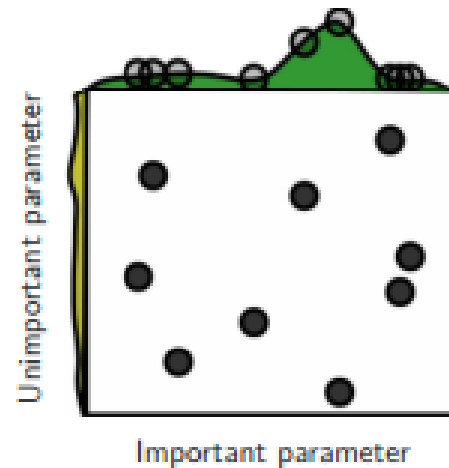
- How to search for good hyperparameters?
  - Grid search
  - Random search (Bergstra and Bengio 2012)
  - Bayesian Methods

# Hyperparameter search

Grid Layout



Random Layout



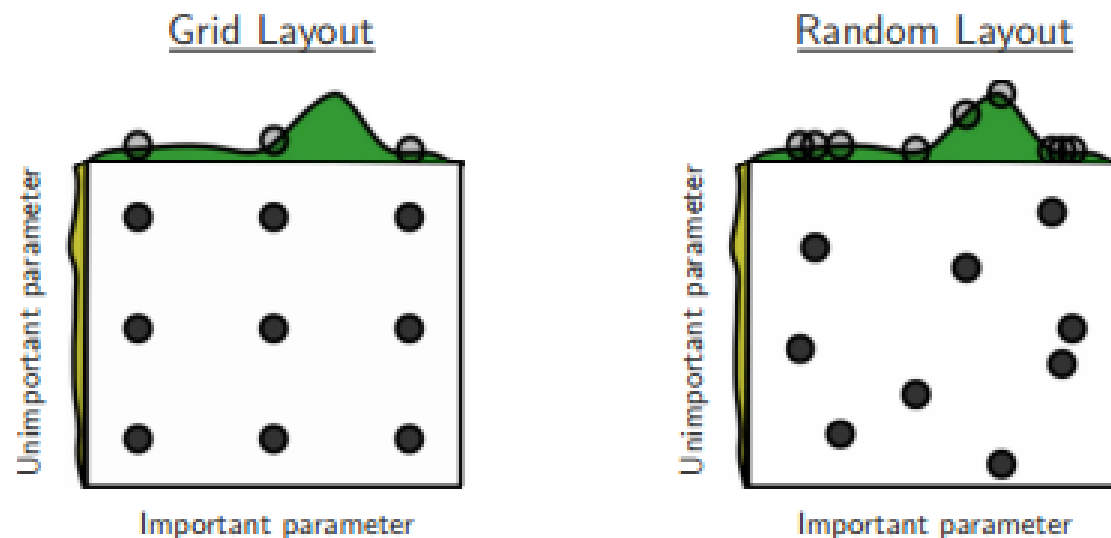


Figure 1: Grid and random search of nine trials for optimizing a function  $f(x,y) = g(x) + h(y) \approx g(x)$  with low effective dimensionality. Above each square  $g(x)$  is shown in green, and left of each square  $h(y)$  is shown in yellow. With grid search, nine trials only test  $g(x)$  in three distinct places. With random search, all nine trials explore distinct values of  $g$ . This failure of grid search is the rule rather than the exception in high dimensional hyper-parameter optimization.

## Loss functions

In the remainder, we'll call

- The true labels  $\mathbf{t} = (t_1, \dots, t_m)$  (“truth”)
- Our network's predictions  $\mathbf{y} = (y_1, \dots, y_m)$

- We had said before that our goal in ML is to solve (a variant of)

$$\sum_{(x,t) \in T} (f_{\theta}(x) - t)^2$$

where  $f_{\theta}$  was our model, parametrized by  $\theta$

- What we optimized there was the so-called **square loss**

$$\ell(t, y) = (y - t)^2$$

- Multi-dimensional square loss would look as follows:

$$\ell(\mathbf{t}, \mathbf{y}) = \sum_j (y_j - t_j)^2 = \|\mathbf{y} - \mathbf{t}\|^2$$



There are other loss functions commonly used in machine learning such as

- **0-1 loss:**  $\ell(\mathbf{t}, \mathbf{y}) = \begin{cases} 1, & \text{if } \mathbf{t} \neq \mathbf{y} \\ 0, & \text{if } \mathbf{t} = \mathbf{y} \end{cases}$
- **Multi-dim Hinge loss:**  $\ell(\mathbf{t}, \mathbf{y}) = \sum_j \max(0, y_j - y_t + 1)$ 
  - where  $t$  is the index where  $\mathbf{t}_t = 1$
- **Cross-entropy loss**
  - $\ell(\mathbf{t}, \mathbf{y}) = -\sum_j t_j \log(y_j)$  (Minimum value is achieved when  $\mathbf{t} = \mathbf{y}$ . In this case  $\ell(\mathbf{t}, \mathbf{y}) = H(\mathbf{y})$ , the entropy of  $\mathbf{y}$ )

Let's assume  $\mathbf{t}$  is a **one-hot vector**  
i.e.,  $\mathbf{t} = [0, 0, \dots, 1, 0, \dots, 0]$

# Examples

- Suppose that for an input  $x$ 
  - $t = (0,1,0,0)$
  - $y = (0.25,0.3,0.4,0.05)$

Let's assume  $t$  is a **one-hot vector**  
i.e.,  $t = [0,0,\dots,1,0,\dots,0]$   
→ multi-label classification

# Examples

- Suppose that for an input  $x$ 
  - $t = (0,1,0,0)$
  - $y = (0.25,0.3,0.4,0.05)$

# Examples

- Square loss is

$$\mathbf{t} = (0,1,0,0)$$

$$\mathbf{y} = (0.25,0.3,0.4,0.05)$$

# Examples

- Square loss is

$$\mathbf{t} = (0,1,0,0)$$

$$\mathbf{y} = (0.25,0.3,0.4,0.05)$$

- $0.25^2 + 0.7^2 + 0.4^2 + 0.05^2$

# Examples

- 0-1 loss is

$$\mathbf{t} = (0,1,0,0)$$

$$\mathbf{y} = (0.25,0.3,0.4,0.05)$$

# Examples

- 0-1 loss is

- 1

$$t = (0,1,0,0)$$

$$y = (0.25,0.3,0.4,0.05)$$

# Examples

- Multi-dim Hinge Loss is

$$\mathbf{t} = (0, 1, 0, 0)$$
$$\mathbf{y} = (0.25, 0.3, 0.4, 0.05)$$

$$\ell(\mathbf{t}, \mathbf{y}) = \sum_j \max(0, y_j - y_t + 1)$$



- Multi-dim Hinge Loss is

$$\mathbf{t} = (0,1,0,0)$$

$$\mathbf{y} = (0.25,0.3,0.4,0.05)$$

- $(0.25 - 0.3 + 1) + (0.3 - 0.3 + 1)$   
 $+ (0.4 - 0.3 + 1) + (0.05 - 0.3 + 1)$

$$\ell(\mathbf{t}, \mathbf{y}) = \sum_j \max(0, y_j - y_t + 1)$$

Whenever  $y_j - y_t + 1 \leq 0$

$$\Leftrightarrow y_j \leq y_t - 1$$

we occur no loss for class  $j$ . The constant 1 is the „margin“

- Cross-Entropy loss is

$$\mathbf{t} = (0,1,0,0)$$
$$\mathbf{y} = (0.25,0.3,0.4,0.05)$$

$$\ell(\mathbf{t}, \mathbf{y}) = - \sum_j t_j \log(y_j)$$

# Examples

- Cross-Entropy loss is

$$\mathbf{t} = (0,1,0,0)$$

$$\mathbf{y} = (0.25,0.3,0.4,0.05)$$

- $-\log(0.3)$

$$\ell(\mathbf{t}, \mathbf{y}) = - \sum_j t_j \log(y_j)$$

Loss over whole training data is the sum over loss for each example

- $L = \sum \ell_i$
- where  $\ell_i = \ell(\mathbf{t}_i, \mathbf{y}_i)$

# Cross-entropy loss

The 'natural' loss for softmax is cross-entropy

- $\ell(\mathbf{t}, \mathbf{y}) = -\sum_j t_j \log(y_j)$ 
  - $\mathbf{t} = (t_1, \dots, t_m)$  is the target output (distribution)
  - $\mathbf{y} = (y_1, \dots, y_m)$  is the network prediction (distribution)

# Cross-entropy loss

- Cross-entropy  $H(\mathbf{p}, \mathbf{q}) = -\sum_x p(x) \log q(x)$ 
  - is (related to) a measure of distance between two (discrete) probability distributions
  - CE is minimized when  $\mathbf{p} = \mathbf{q}$ 
    - But it's not zero then, but  $H(\mathbf{p})$ , the *entropy* of distribution  $\mathbf{p}$
- It's not symmetric,  $H(\mathbf{p}, \mathbf{q}) \neq H(\mathbf{q}, \mathbf{p})$
- $H(\mathbf{p}, \mathbf{q}) = H(\mathbf{p}) + \text{KL}(\mathbf{p}, \mathbf{q})$ 
  - where KL is the Kullback-Leibler divergence

# Square loss vs. CE loss

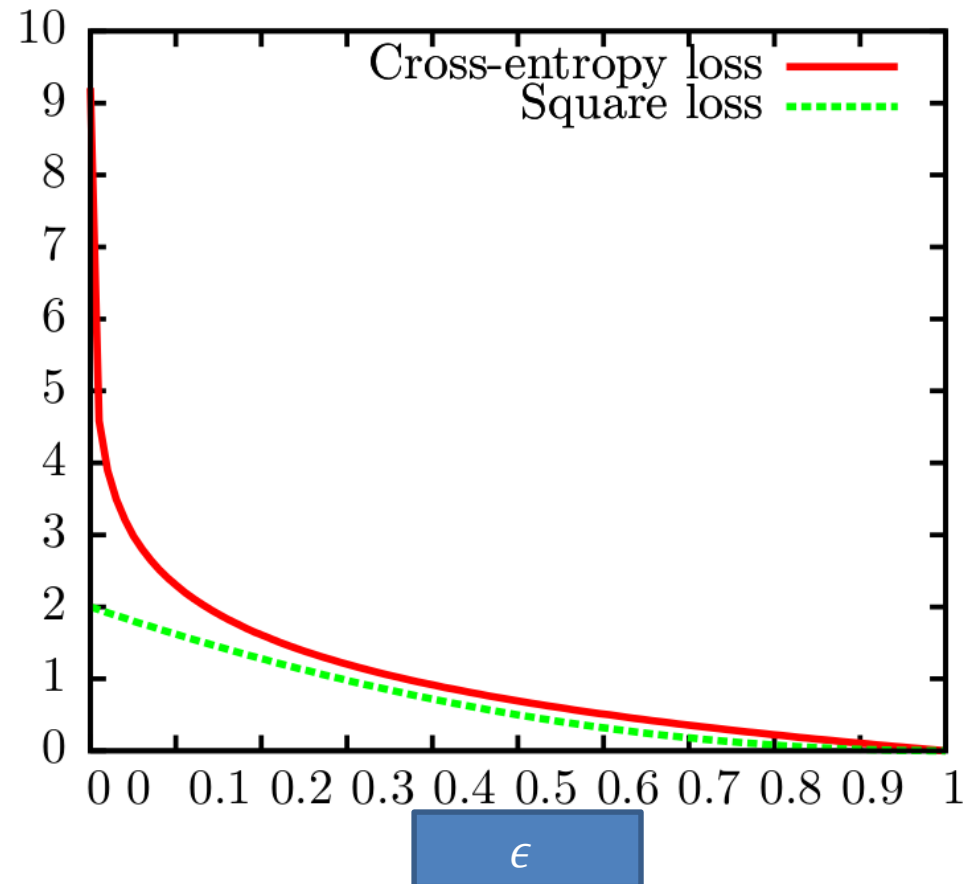
Is square loss a good loss function for multi-class prediction?

- Consider
  - $\mathbf{t} = (1, 0)$
  - $\mathbf{y} = (\epsilon, 1 - \epsilon)$
- Square loss  $2(1 - \epsilon)^2$
- Cross entropy loss  $-\log \epsilon$

# Square loss vs. CE loss

Is square loss a good loss function for multi-class prediction?

- Consider
  - $\mathbf{t} = (1, 0)$
  - $\mathbf{y} = (\epsilon, 1 - \epsilon)$
- Square loss  $2(1 - \epsilon)^2$
- Cross entropy loss  $-\log \epsilon$







## Evaluation

- We often use **accuracy**:
  - How many instances are correctly classified divided by the number of instances in the test set
- Btw. Assume our classifier is a DL model with softmax in the last layer to perform multi-class classification
  - How to compute the accuracy? i.e.,
  - How do we get a decision from the probability distribution?

Why are we not evaluating models using loss functions directly?



- When data is imbalanced:
  - F1
- When outputs are continuous:
  - E.g., **Squared distance (MSE)**, **correlation**, etc.
- When output is a sequence:
  - E.g., **edit distance**, **BLEU**, etc.

- Even when your outputs are discrete classes, accuracy is sometimes not a good evaluation measure
- Say, you want to predict whether a patient has a rare disease Q
- Which system is better?

- System A

	Prediction is Q	Prediction is not Q
Patient has Q	0	10
Patient has not Q	1	1004

- System B

	Prediction is Q	Prediction is not Q
Patient has Q	4	6
Patient has not Q	5	1000

- Both systems have an acc. of  $1004/1015 = 99\%$

## For Class „Disease“

- System B has a *precision* of  $4/9 = 44\%$
- and a *recall* of  $4/10 = 40\%$
- System A has precision of 0% and recall of 0%

## For Class „No Disease“, both systems are very close:

- Precision A: 1004/1014, Pr B: 1000/1006
- Recall A: 1004/1005, Rec B: 1000/1005

- When there are more than two classes, precision and recall for class  $i$  are defined as

- $\text{Precision}_i = C(i,i)/\text{sum}(C(:,i))$

- $\text{Recall}_i = C(i,i)/\text{sum}(C(i,:))$

where  $C$  is a *confusion matrix* as above

- From precision and recall, one can compute the „ $F1$  score“:

The harmonic mean of  $P$  and  $R$ :

$$2PR/(P + R)$$



# Quizz

Pred

True

	Politics	Sports	Religion
Politics	2	4	3
Sports	1	9	5
Religion	6	8	9

*What is the recall of class „Sports“?*

- A: 9/17
- B: 9/21
- C: 9/15
- D: 20/42



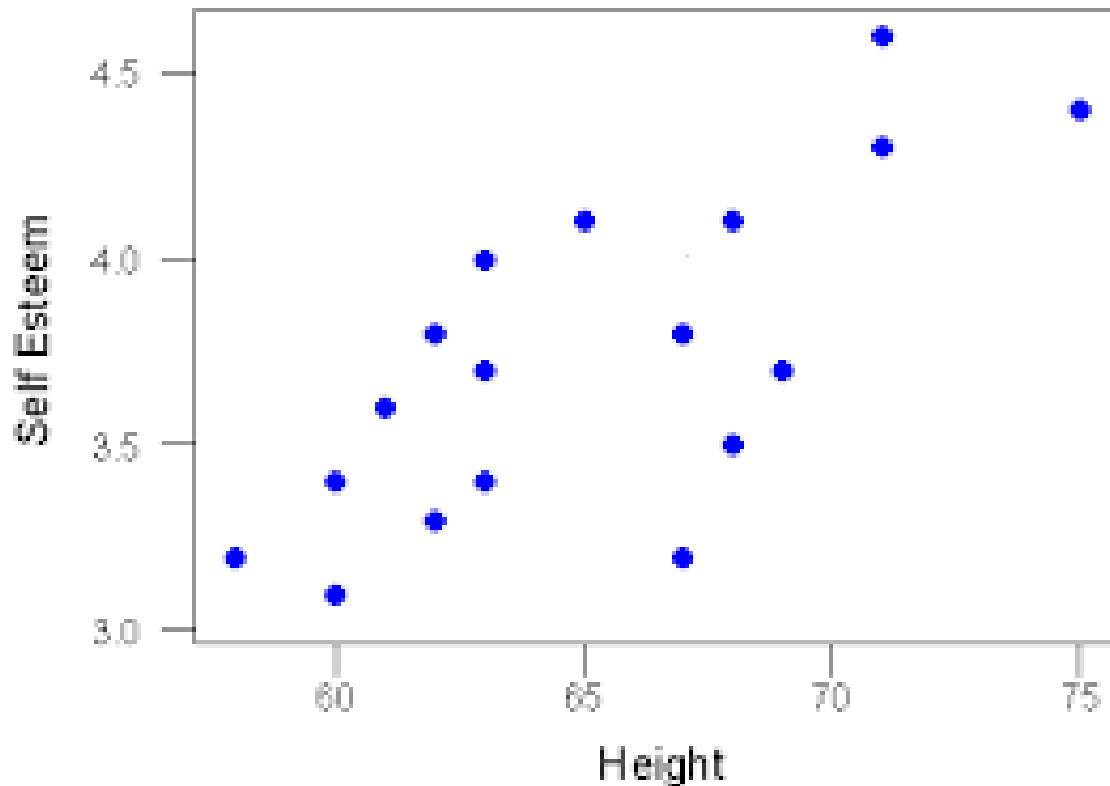
# Evaluation – F1 measure

- For two or more classes, one typically computes the F1-score of each class and then combines this in an overall score:
  - For example, averaging all the F1 scores
  - This is called *macro F1*



Can use when your output are **discrete classes** which are **imbalanced**

# Evaluation - Correlation



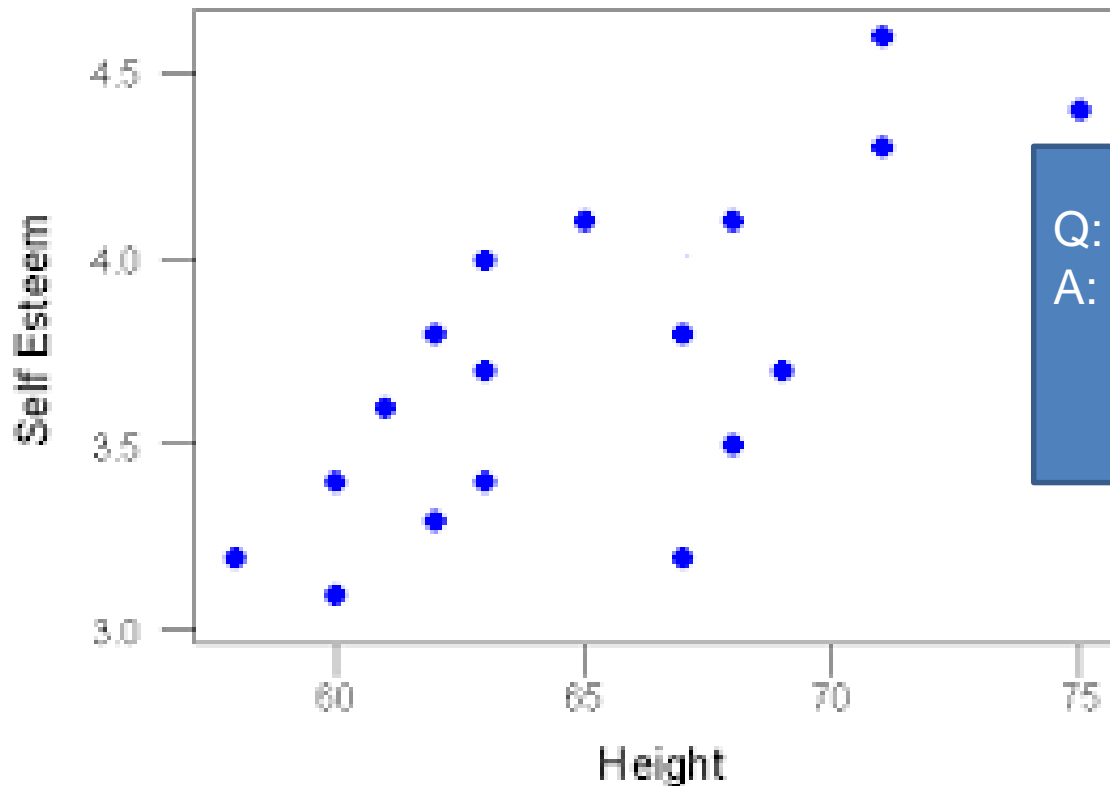
```
>>> from scipy.stats import pearsonr  
>>> # try also spearmanr  
>>> x=[1,2,3]; y=[-1,-4,6]  
>>> pearsonr(x,y)  
0.682
```

Measures linear relationship between two variables: if one increases, does the other also increase (in expectation)?



Can use when your output is **continuous** (real numbers) rather than discrete

# Evaluation - Correlation



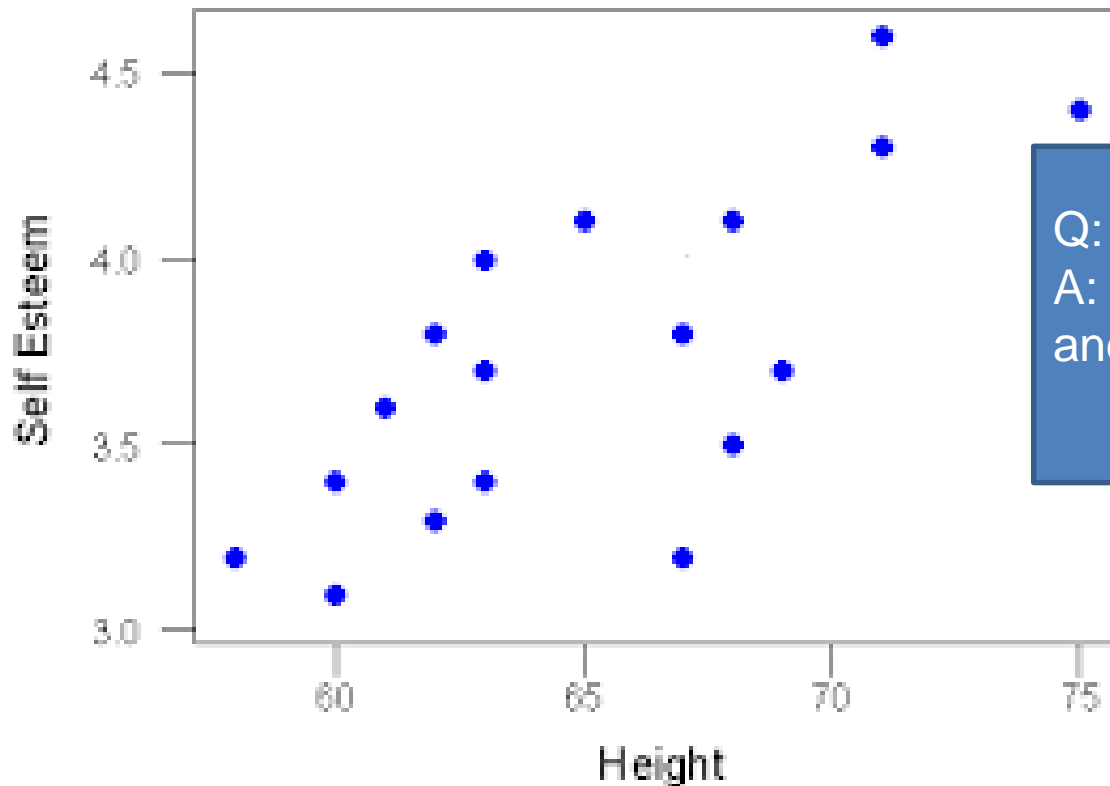
Q: What are we correlating?  
A:

Measures linear relationship between two variables: if one increases, does the other also increase (in expectation)?



Can use when your output is **continuous** (real numbers) rather than discrete

# Evaluation - Correlation



Q: What are we correlating?  
A: (Vector of) model predictions  
and truth labels

Measures linear relationship  
between two variables:  
if one increases, does the  
other also increase  
(in expectation)?



Can use when your output is **continuous** (real numbers) rather than discrete

# DISCUSS

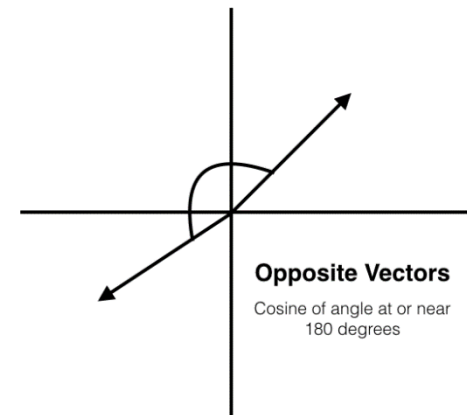
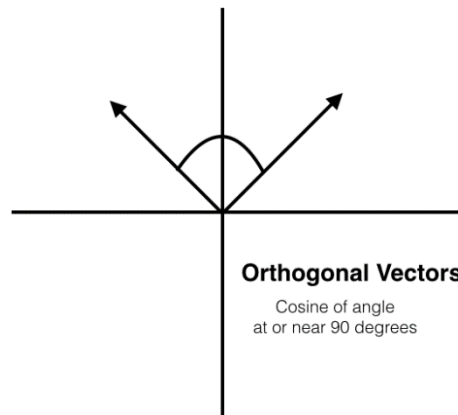
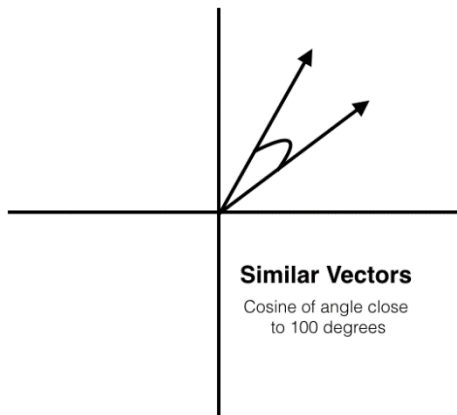
When would you prefer correlation over MSE?  
(and vice versa)



# Evaluation – Cosine Similarity

For two vectors  $x, y$ , their *cosine similarity* is:

$$\frac{x \cdot y}{||x|| \cdot ||y||} \in [-1, 1]$$



Can use when your output is a **vector**

# Evaluation – Edit distance

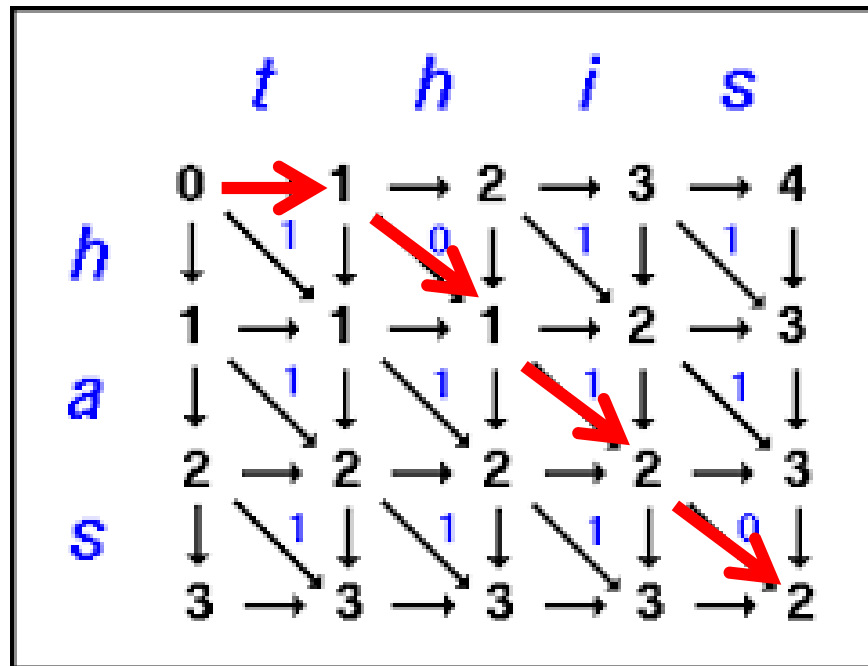
- Minimal number of insertions, deletions and substitutions to transform one sequence into another

		<i>t</i>	<i>h</i>	<i>i</i>	<i>s</i>	
	0	→ 1	→ 2	→ 3	→ 4	
<i>h</i>	↓	↘ 1	↓	↘ 0	↓	↘ 1
	1	→ 1	→ 1	→ 2	→ 3	
<i>a</i>	↓	↘ 1	↓	↘ 1	↓	↘ 1
	2	→ 2	→ 2	→ 2	→ 3	
<i>s</i>	↓	↘ 1	↓	↘ 1	↓	↘ 0
	3	→ 3	→ 3	→ 3	→ 2	



# Evaluation – Edit distance

- Minimal number of insertions, deletions and substitutions to transform one sequence into another



```
>>> import editdistance as ed
>>> ed.eval("this","has")
2
```



Can use when your output are **strings** (words, text, documents)

- There are many more evaluation measures
  - For Machine Translation (MT): BLEU scores,...
  - For Summarization: ROUGE (n-gram overlap),....
  - Etc.
- Choosing which evaluation measure is an important field of research
- For higher level NLP tasks, some automatic measures may correlate poorly with human evaluation → need for new measures



How would you evaluate machine translation and why is it difficult?



- Foundations of ML
  - Train vs. Dev vs. Test Set
  - (Expected) Loss function optimization
  - Evaluation Measures
- Cross-Entropy vs Square Loss for Neural Nets

## General/Advice:

- In class, we use square loss or (more often) cross-entropy loss
- Output layer typically has softmax activation function (multi-class classification) → in this case, always choose CE loss
- Loss function and evaluation measure go together, i.e., one may change loss when another evaluation metric is used