

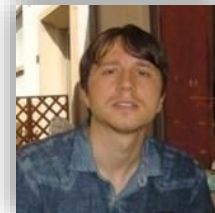
Deep Learning for NLP



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Lecture 2 – ML principles

Dr. Steffen Eger
Wei Zhao
Niraj Dev Pandey



Natural Language Learning Group (NLLG)
Technische Universität Darmstadt

This lecture:

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



ML principles

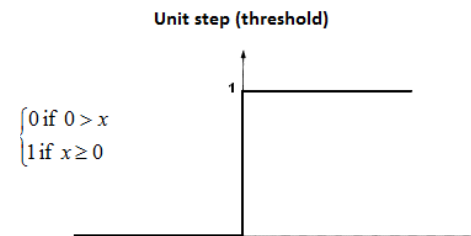
Before we start: Notation

- Throughout, we write
 - Vectors (elements of \mathbf{R}^d) as bold face letters
 - E.g., $\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{x}$
 - Other conventions used in literature:
 1. \vec{u}, \vec{x}, \dots (adopted from physics)
 2. Nothing: u, v, w, x

This can be very confusing if in addition you don't specify the range of variables („meaningless formulae“)
- Scalars as ordinary letters: a, b, c, d
- Matrices are bold and uppercase: $\mathbf{U}, \mathbf{V}, \mathbf{W}, \mathbf{X}, \mathbf{Y}$

Before we start: Notation

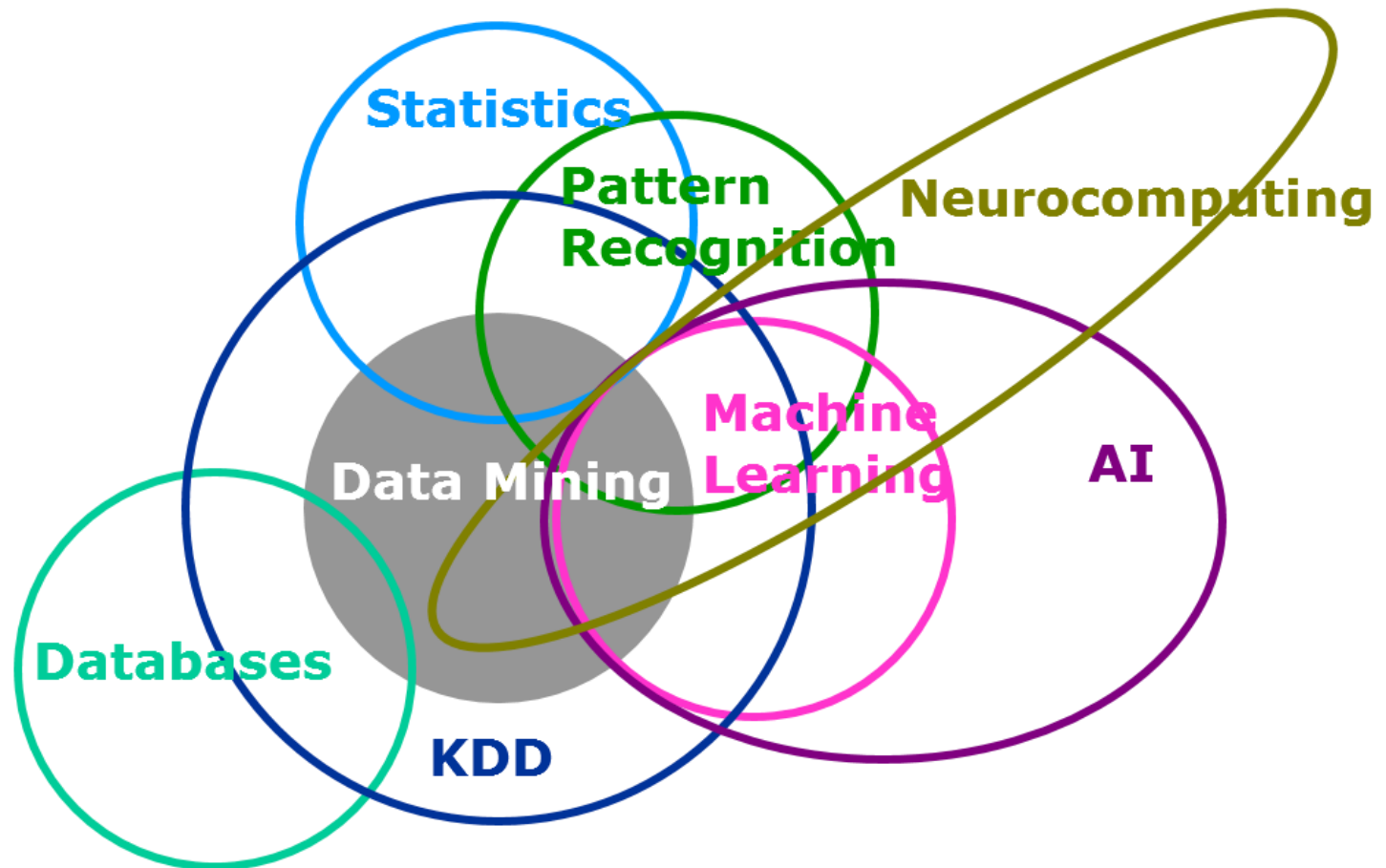
- Often, we assume that vectors are **row vectors**, i.e., lie in $\mathbf{R}^{1 \times n}$
 - I'll try to be explicit on this
- A dot can mean different things, depending on context:
 - Multiplication: $a \cdot b$
 - Dot product: $\mathbf{x} \cdot \mathbf{y} = \sum x_i \cdot y_i$
 - Matrix-vector or matrix-matrix multiplication: $\mathbf{x} \cdot \mathbf{W}$, $\mathbf{U} \cdot \mathbf{V}$
(note that dimensions must fit here - same is for dot product!)
- A note on differentiability:
 - What is $\sigma'(x)$ for this function?



Before we start: Notation

- Elementary vector-matrix multiplication:
 - What is $x \cdot E$, where
 - $x \in \mathbf{R}^{1 \times n}$ is a 1-hot vector and $E \in \mathbf{R}^{n \times d}$?
- Cosine similarity:
 - For two vectors x, y , their *cosine similarity* is defined as:
 - $\frac{x \cdot y}{||x|| \cdot ||y||} \in [-1, 1]$

- A fancy (mystifying) name for a rather „ordinary“ field of study



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}

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

Standard Setup

- We have **data**
 - $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
- And a statistical model $f_{\theta}(\mathbf{x})$
- We also specify some *loss* function, e.g.,
 - $\frac{1}{N} \cdot \sum_{(x,y)} (y - f_{\theta}(\mathbf{x}))^2$
- **Goal** is then to optimize/minimize our loss over θ :
 - We are looking for parameters that bring our model close to the data

- Loss: $\frac{1}{N} \cdot \sum_{(x,y)} (y - f_{\theta}(x))^2$
- Goal is then to optimize/minimize our loss over θ
- HOWEVER: our (real) goal is NOT to fit the **given** data well (which is called **overfitting**)
 - Every sufficiently rich model f_{θ} can do this
 - Moreover, if our data is non-pathological, we'll always find a model f_{θ} that perfectly fits it

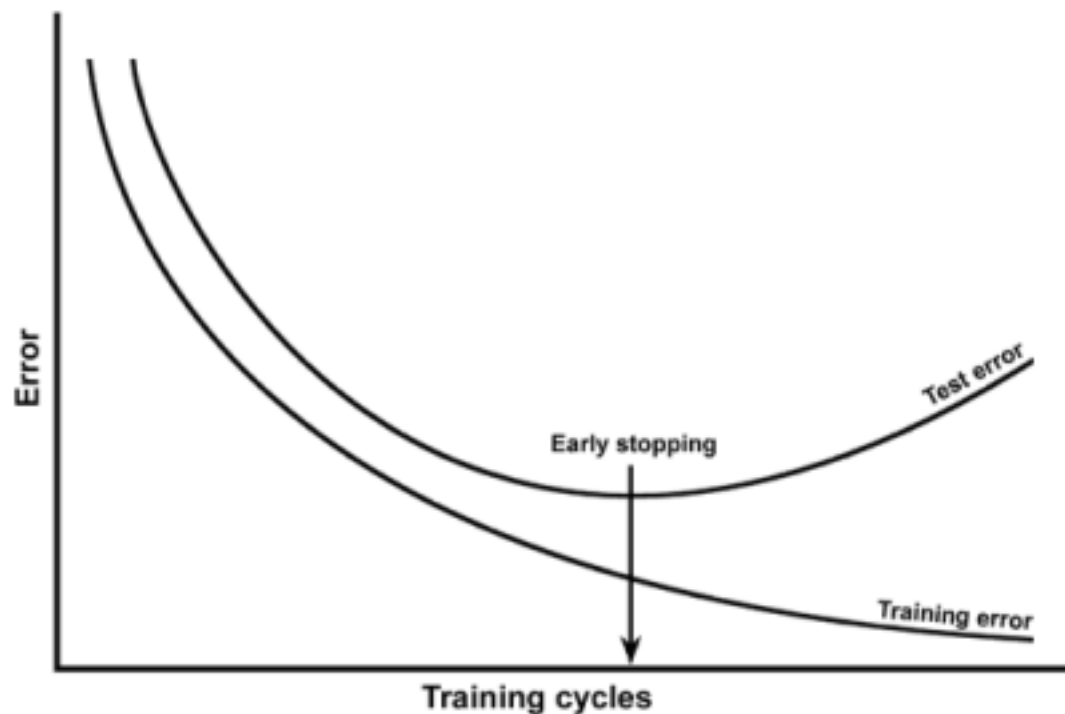


- Can assume our previous optimization problem is a „proxy“ of our true optimization problem (*); but must keep the latter in mind, too
- Because we're interested in generalization performance, **we always split our data:**
 - Test data vs. Train data
 - Test data represents the true underlying distribution
- Your model must perform well on the test data; its performance on the training data is not of interest to us

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

Training data

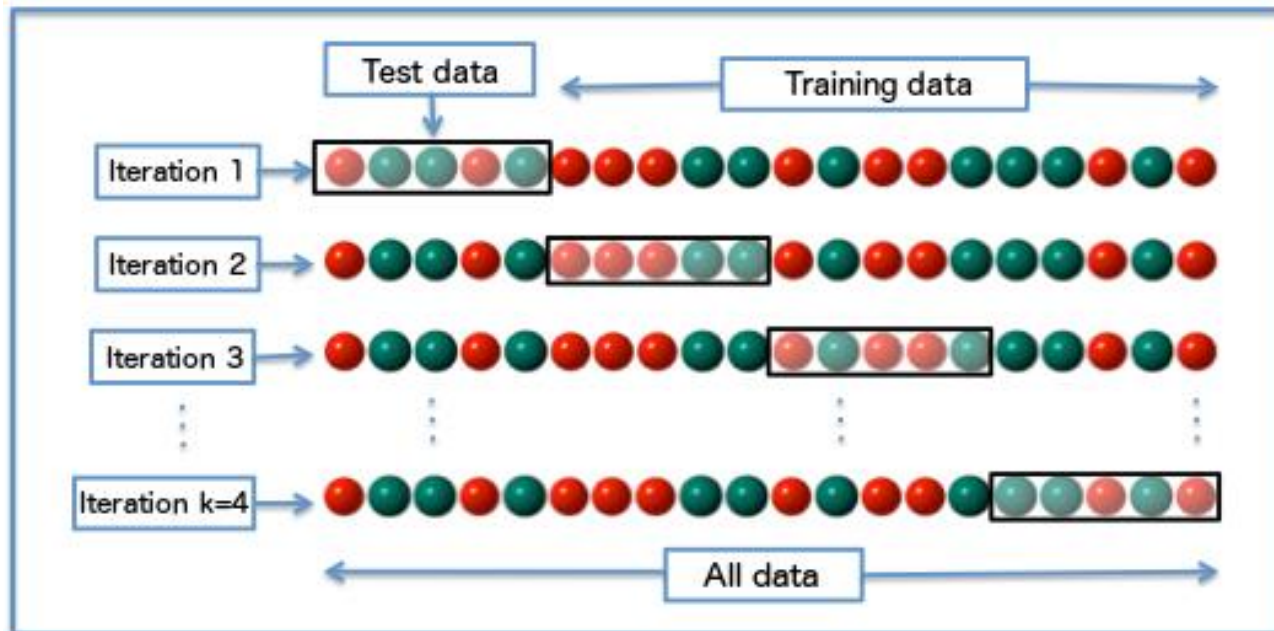
In particular: can use dev data for **early stopping**, etc.
(extremely common and popular)



- Could choose some fraction of your overall data for testing, e.g., 75%/25% split (train/test) or 90%/10%
- That's not always a good idea: by chance, you could pick 10% that are not representative of the overall problem
- Better is k-fold cross-validation:
 - Train on $(k-1)$ equally sized folds, test on the remaining
 - Repeat k times

Test data

- Better is k-fold cross-validation:
 - Train on $(k-1)$ equally sized folds, test on the remaining
 - Repeat k times



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

Hyperparameter search

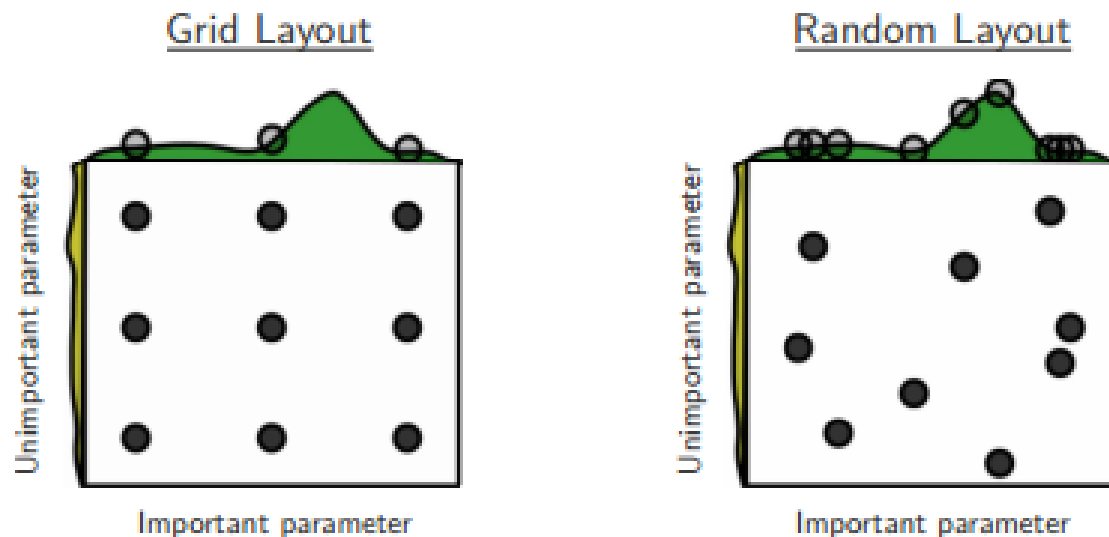
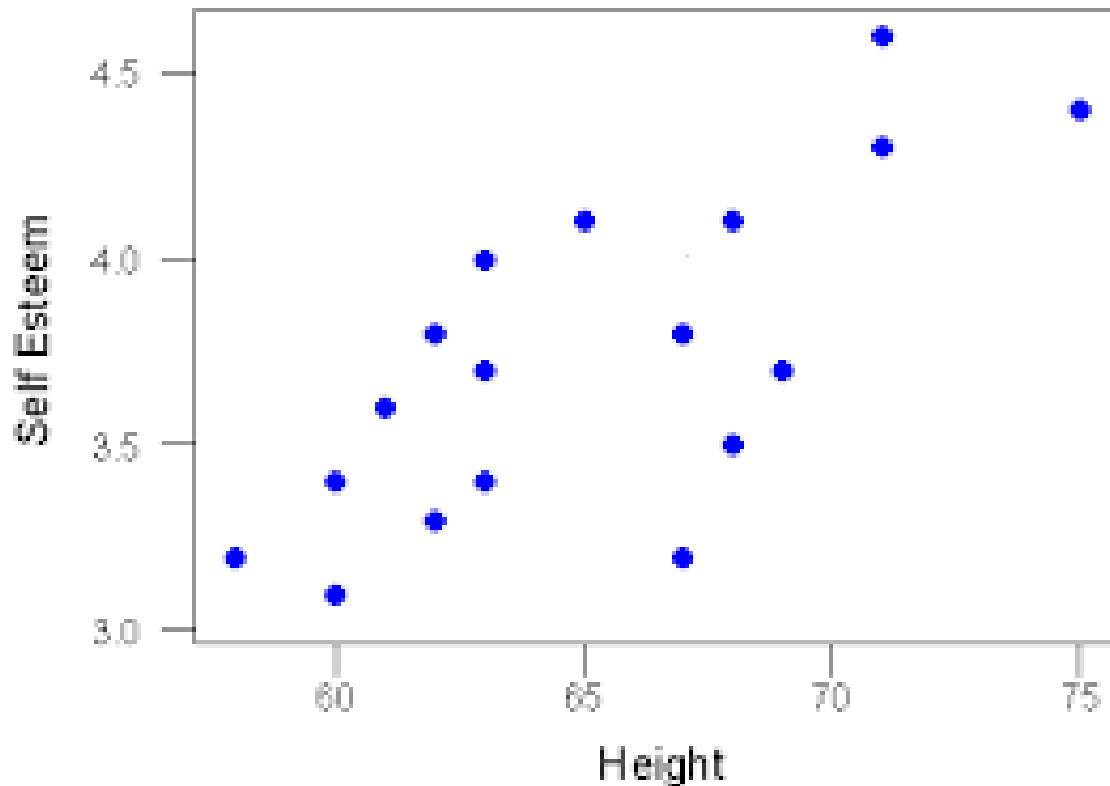


Figure 1: Grid and random search of nine trials for optimizing a function $f(x,y) = g(x) + h(y) \approx g(x)$ with low effective dimensionality. Above each square $g(x)$ is shown in green, and left of each square $h(y)$ is shown in yellow. With grid search, nine trials only test $g(x)$ in three distinct 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.

- We often use **accuracy** on the test data to evaluate our model performance:
 - How many instances are correctly classified divided by the number of instances in the test set
- However, there are other plausible evaluation measures:
 - E.g., when your outputs are continuous:
 - **Squared distance (MSE)**, **cosine**, **correlation**, etc.
 - Say, your output is a sequence:
 - Could use **edit distance**, for instance

Evaluation - Correlation



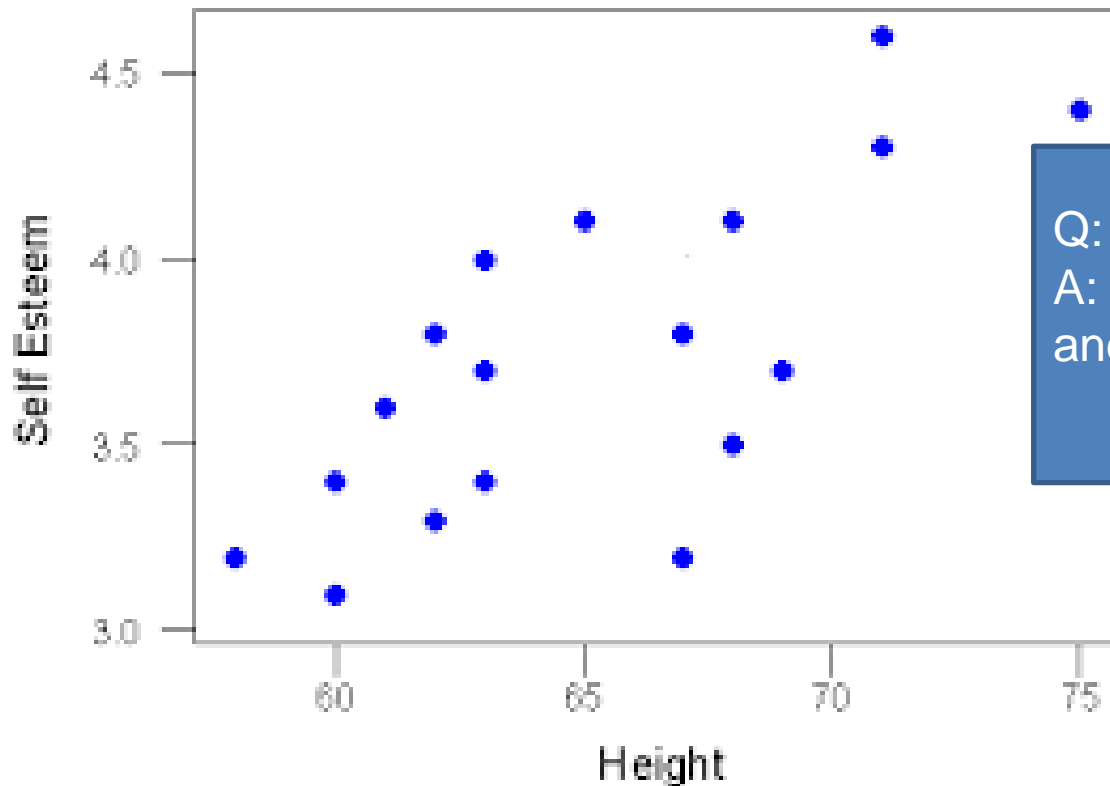
```
>>> from scipy.stats import pearsonr  
>>> # try also spearmanr  
>>> x=[1,2,3]; y=[-1,-4,6]  
>>> pearsonr(x,y)  
0.682
```

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: (Vector of) model predictions
and truth labels

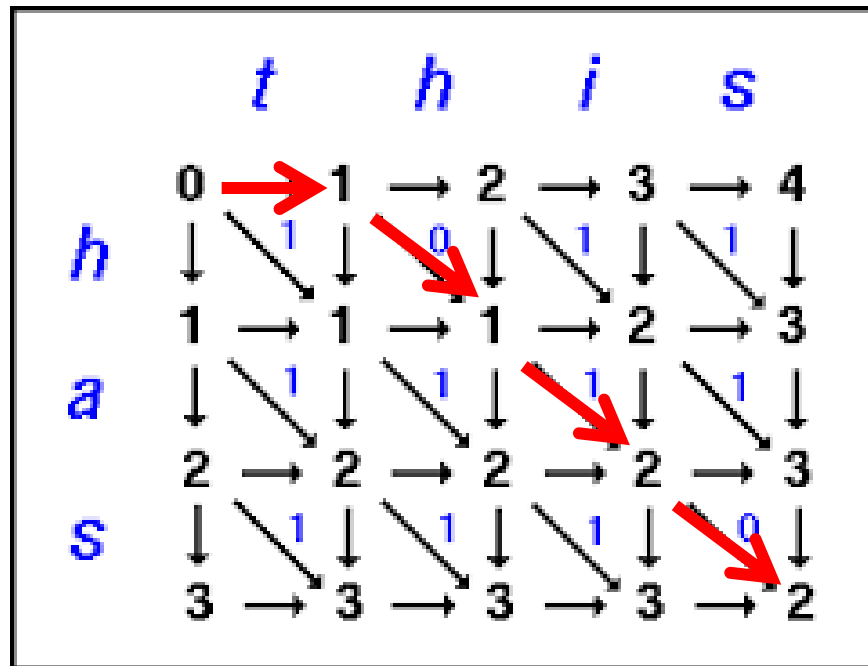
Measures linear relationship
between two variables:
if one increases, does the
other also increase
(in expectation)?



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

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)

- 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?

Terminology (precision/recall)

	Prediction is Q	Prediction is not Q
Patient has Q	True positive	False negative
Patient has not Q	False positive	True negative

- 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**, defined as $2PR/(P + R)$

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
- There are different manners in which this can be done with different names
 - e.g. *micro F1* vs. *macro F1*



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

- 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 import field of research
- Importantly: For higher level NLP tasks, automatic measures may correlate poorly with human evaluation
→ need for new measures





Loss functions

- In the remainder, we'll call
 - The true labels $\mathbf{t} = (t_1, \dots, t_m)$ (“truth”)
 - Our network's predictions $\mathbf{y} = (y_1, \dots, y_m)$
- Note that we have m output units – indexed by j
- As before the number of samples is N :
$$\{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\}$$
 - We index them by i

Square loss

- 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_{θ} was our model, parametrized by θ

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

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

- Multi-dimensional square loss would look as follows:

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

Types of loss functions

- 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:** $\ell(\mathbf{t}, \mathbf{y}) = \sum_j \max(0, y_j - y_{t_j} + 1)$
 - where t is the index where $t_{t_j} = 1$
 - **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})

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

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

Examples

- Multi-dim Hinge Loss is

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

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

Examples

- 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“

Examples

- 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

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

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

- $-\log(0.3)$

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

Training data loss

- 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 $H(\mathbf{p}, \mathbf{q}) = -\sum_x p(x) \log q(x)$
 - 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 $-\log \epsilon$

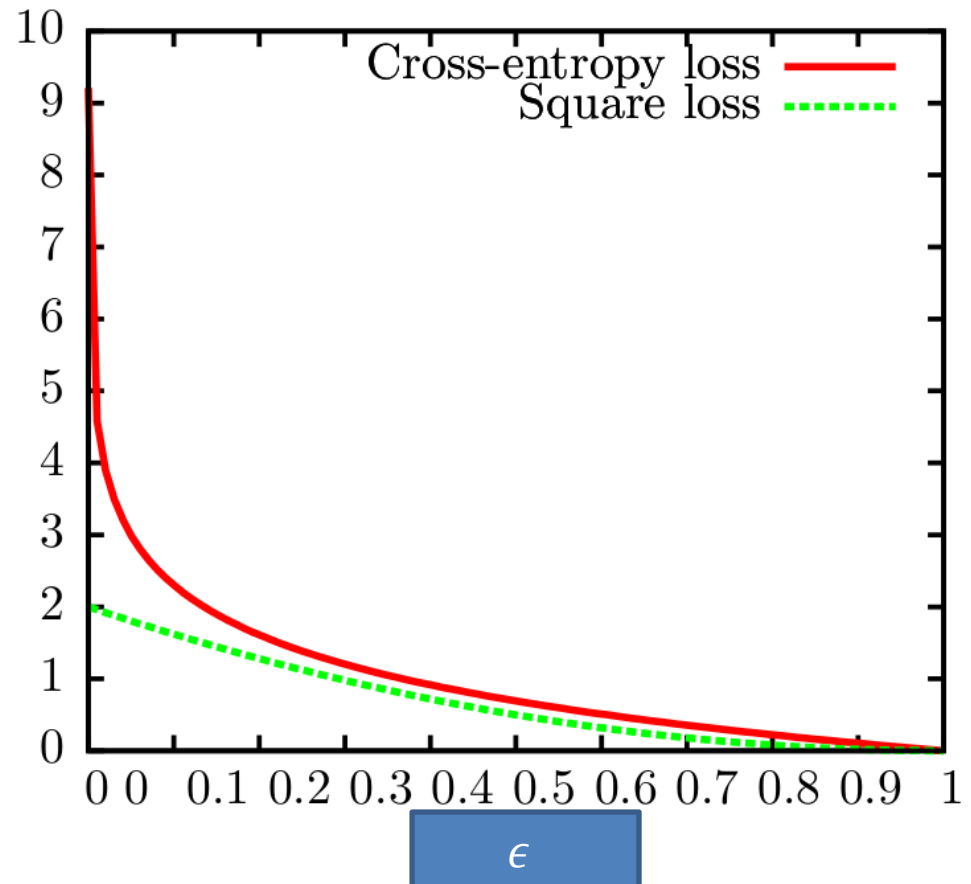
Square loss vs. CE loss

- 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 $-\log \epsilon$



- 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

References

- Maas et al. 2013, Rectifier non-linearities improve neural network acoustic models. In ICML, vol. 30
- He et al., 2015, Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification
- Clevert et al., 2015, Fast and Accurate Deep Network Learning by Exponential Linear Units
- Goodfellow et al., 2013, Maxout Networks
- Glorot and Bengio, 2010, Understanding the difficulty of training deep feedforward neural networks
- Ioffe and Szegedy, 2015, Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift