# Machine Learning
# V07: ML System Design

System development: What to give priority?
Example: Learning to read checks end-to-end

With material from Andrew Y. Ng, Coursera
See also [LeCun et al, "Gradient-Based Learning…", 1998]

# Educational objectives

- **Remember error-** and **ceiling analysis** as well as the initial **24h hack** as tools to be successful in ML

- **Know how** to **design** and **prioritize** complete machine learning **system pipelines**

- **Appreciate** the **elegance** of the design **that enables ent-to-end learning** for the check reading application of LeCun et al.

# 1. SYSTEM DEVELOPMENT: WHAT TO GIVE PRIORITY?

# Example 1: Building a spam classifier

```
From: cheapsales@buystufffromme.com
To: stdm@zhaw.ch
Subject:  Buy now!

Deal of the week! Buy now!
Rolex w4tchs - $100
Med1cine (any kind) - $50
Also low cost M0rgages
available.
```

```
From: Renate Stadelmann
To: stdm@zhaw.ch
Subject:  Holiday plans

Hi Thilo,
was talking to Philipp about
plans for New Year. Sauna and
surfing in winter? ;-)
Love, Renate
```

Supervised learning task
- $x$: features of email ➔ $y$: 1 (spam) or 0 (non-spam)

Practical features
- Create **list of** 50'000 **most frequent words** in training set

$$x = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix} \begin{matrix} \#Thilo \\ \#buy \\ \#deal \\ \#discount \\ \vdots \\ \#now \\ \vdots \end{matrix}, \; x_j = \begin{cases} 1: word\ occurrs\ in\ mail \\ 0: otherwise \end{cases}$$

# Example 1: Building a spam classifier (contd.)
## How to prioritize *algorithmic* work?

How to best invest the time to make it work (i.e., have low error)?
- **Collect** lots of **data** (e.g., "honeypot" project)?
- **Develop** sophisticated **features**?
  - …based on email routing information from email **header**
  - …for message **body**
    - ➔ Treat "discount" and "discounts" as same word? "Deal" and "Dealer"?
    - ➔ Features about punctuation?
- **Develop** sophisticated **algorithm** to detect misspellings?
  - ➔ e.g. "m0rtgage", "med1cine", "w4tches"

> **Advice**
> - Take **24h** to implement (rather: **hack**) a **complete system** including scoring
> - Use diagnostics to decide where to improve

Recommendation
1. Start with a **simple algorithm** that can be **implement**ed **quickly**
   - ➔ Implement it and test it on cross-validation data
2. Plot **learning curves** to diagnose if more data, more features, etc. are likely to help
3. **Error analysis**: Manually examine the CV examples that where misclassified
   - ➔ Is there a **systematic trend** in what type of examples are **misclassified**?

# Example 1: Building a spam classifier (contd.)
## Error analysis

Assume the following experimental outcome
- $N_{CV} = 500$ emails in CV set
- 100 emails are misclassified
➔ **Manually examine** the 100 **errors**

**Categorization** based on
1. **Type** of email
2. **Cues** (feature candidates) that would have helped the algorithms to classify correctly

| Type | Number |
|---|---|
| Pharma | 12 |
| Replica / faked goods | 4 |
| Phishing | 53 |
| Other | 31 |

| Cues | Number |
|---|---|
| Deliberate misspellings ("m0rgage", "med1cine", etc.) | 5 |
| Unusual email routing | 16 |
| Unusual punctuation ("!!!!!!!!" etc.) | 32 |

quite rare

this might help

# Example 1: Building a spam classifier (contd.)
## The importance of numerical evaluation (error analysis 2)

**Should** a stemmer **be used** (e.g., free "Porter stemmer")?
- Treats "discount" / "discounts" / "discounted" / "discounting" as the same word
- Makes e.g. "universe" / "university" undistinguishable
➔ Error analysis doesn't help much in deciding

Solution
- **Try** with & without
- **Compare numerical results** ➔ need a **single performance metric** for that (e.g. CV error; F-score)

| Method | CV error |
|---|---|
| Original: without stemming | 5% |
| With stemming | 3% |
| Additional: distinguishing upper vs. lower case | 3.2% |

good idea!

doesn't help

- Attention: If classes are **skewed** (e.g., cancer prediction), **regard** recall-precision **trade-off**
  ➔ Use  for example the F-measure instead of pure error (➔ compare V03)
  ➔ Give the **rare class** the **label 1/true**

# Example 2: Check reading application

Problem description:
- Read amount of $6
- Easy for humans, but time-consuming
→ Automation wanted



other numbers

courtesy amount

legal amount

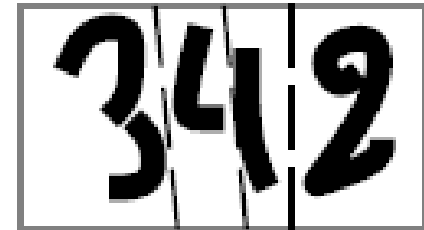Challenge:
- Which number?
- Diversity

# Example 2: Check reading application (contd.)
## What part of the *pipeline* to improve next?

Challenge
- Identify correct character string (*«342»*) on a piece of paper
- Therefore:
    - **Detect** all handwritten strings
    - [ **Identify** correct string (containing the amount) ]
    - Find correct **segmentation**
    - **Recognize** individual characters

System pipeline
- **What part** of the pipeline should you **spend** the **most time** trying to improve**?**
- Note: Identification of the correct string is omitted here (could be placed at the end)

| Image | → | Text detection | → | Character segmentation | → | Character recognition |

# Example 2: Check reading application (contd.)
## **Ceiling analysis**: Attributing errors to individual components

Ceiling analysis
1. Baseline ➔ measure the (CV) performance of the complete pipeline
2. Replace first component with ground truth (perfect results) ➔ measure performance
3. Replace next component with ground truth ➔ measure performance
4. …

| Image | ➔ | Text detection | ➔ | Character segmentation | ➔ | Character recognition |

1. Measure complete existing system

2a. Replace component by ground truth        2b. Use existing implementation

3a. Replace components with ground truth        3b. Use existing implementation

Baseline: Overall system performance    72%

...od idea!

...sn't help

...od idea!

# 2. SYSTEM EXAMPLE: LEARNING TO READ CHECKS END-TO-END

(Just the idea – details in appendix)

# A landmark work in Machine *Learning*

**LeCun et al., "Gradient-Based Learning Applied to Document Recognition", 1998**



## Outline

- <span style="color:green">Gradient-Based ML ✓</span>
- <span style="color:orange">Convolutional Neural Networks (→ later)</span>
- <span style="color:gray">Comparison with other Methods</span>
- Multi-Module Systems & Graph Transformer Networks
- Multiple Object Recognition & Heuristic Oversegmentation
- <span style="color:gray">Space Displacement Neural Networks</span>
- <span style="color:gray">GTN's as General Transducers</span>
- <span style="color:gray">On-Line Handwriting Recognition System</span>
- Check Reading System

# Standard and sequential supervised learning

## Supervised Learning



feature vectors     labels

Typical assumption on data:
- i.i.d.
- Surrounding tasks deemed simple(r)

## Sequential Supervised Learning



Typical assumptions on data:
- Sequence information matters
- Overall task has many challenging components (e.g., segmentation → recognition → sequence assembly)

# Approaches to classifying sequential data
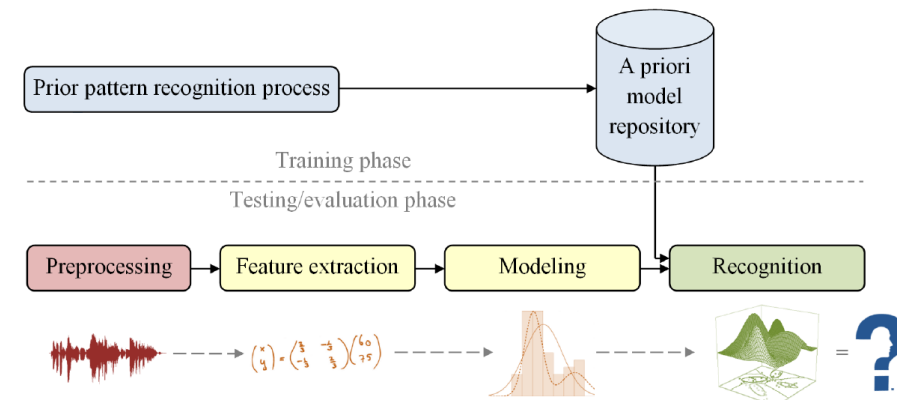
## «A bird in the hand…» approach

- Train standard classifier, extend it using a sliding window and post-processing (e.g., smoothing)

## Direct modeling approach

- Train a generative (statistical) model of the sequence generation process (e.g., HMM)

## «…two in the bush» approach

- Build a unified pattern recognition processing chain, optimize it globally with a unique criterion

See also: T.G. Dietterich, «Machine Learning for Sequential Data – A Review», 2002

# Proposed Solution: Global Learning
## Example: Reading handwritten strings

Challenge
- Identify correct character string (*«342»*) on a piece of paper
- Therefore: Find correct segmentation & recognize individual characters



Images sources for this section: → see references slide in appendix

# Global Learning
## Learning end-to-end

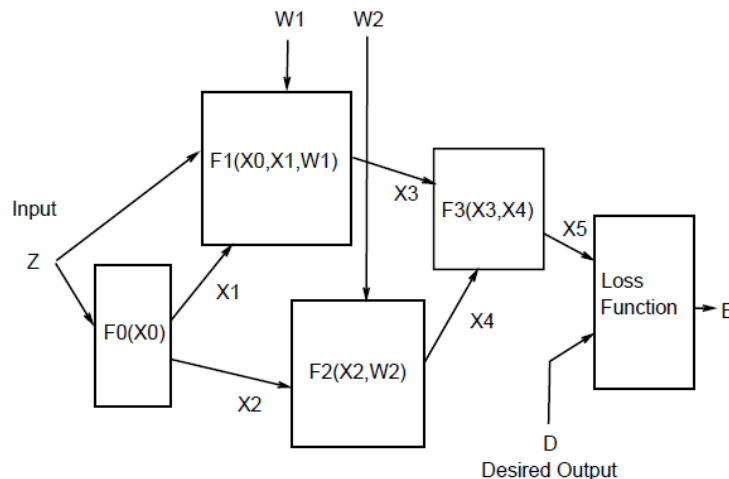What we know: Traditional pattern recognition system architecture



What we want: Train all parameters to optimize a global performance criterion

# Foundation: Gradient-based learning

i.e., gradient descent

A trainable system composed of heterogeneous modules:



**Backpropagation** (bp) can be used if...
- cost (loss) function is differentiable w.r.t. parameters
- modules are differentiable w.r.t. parameters

➔ Gradient-based learning is the unifying concept behind many machine learning methods (➔ see day 1 & 5)
➔ Object-oriented design approach: Each module is a class with a `fprop()` and `bprop()` method

**Graph transformer network** (GTN)
- General architecture for **individual components collectively trainable** via bp
- …using **graph structures** as input and output
➔ See appendix for more details

# Conclusions

Example: Learning to segment without intermediate labels

## Less need for manual labeling

- Ground truth only needed for final result (not for every intermediate result like e.g. segmentation)

## Early errors can be adjusted later due to…

- …unified training of all pattern recognition modules under one regime
- …postponing hard decisions until the very end

possible without crop marks?

## No call upon probability theory for modeling / justification

- Occam's razor: Choose easier discriminative model over generative one
  Vapnik: Don't solve a more complex problem than necessary
- No need for normalization when dealing with "penalties" instead of probabilities → no "other class" examples needed
- Less constrains on system architecture and module selection

# Review

- ML systems are **pipelines composed of** individual **components** that can be **develope**d **collaboratively** in a team

- Do **ceiling analysis** to **decide which component** in the pipeline is most likely to alter the result for good

- Do **qualitative analysis of wrong**ly predicted **examples** to get **insight** what is going wrong

- Do **numerical error analysis** (i.e., compare CV scores) to **prioritize algorithmic ideas**

- Have a **single performance metric** (e.g., error or F-measure) to **monitor** the **evolution** of your system **continuously**

- Consider **end-to-end training** (global optimization; e.g., GTNs)

# P06.2: Applying learning curves

Solve exercise 2 from P06:

- The given Python script draws a learning curve for (any) classifier and certain data sets

- Analyze the performance of the given setup → what insight can you get? What advice as how to proceed would you give?

- Can you also put validation curves to work for you in finding good parameters (see the given library)?
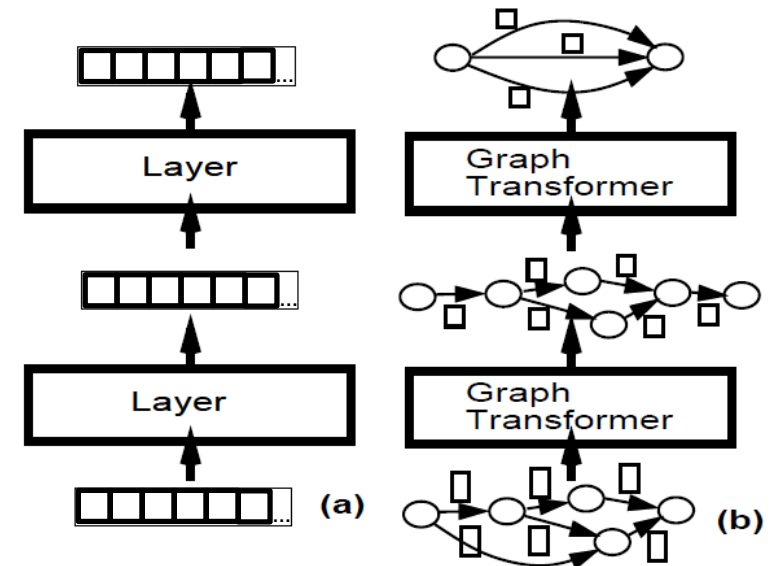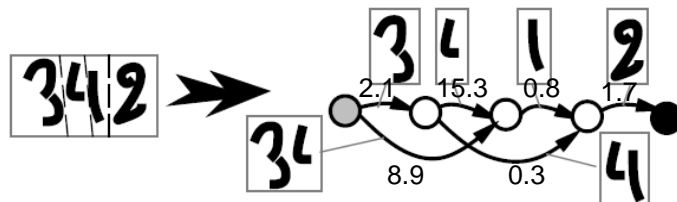
# APPENDIX

# Graph Transformer Networks

**Network of pattern recognition modules that successively refine graph representations of the input**

## GTNs

- Operate on graphs (b) of the input instead fixed-size vectors (a)

- Graph: DAG with numerical information ("penalties") at the arcs



➔ GTN takes gradients w.r.t. module parameters and numerical data at input arcs

# Example: Heuristic over-segmentation
## …for reading handwritten strings

**Viterbi Penalty**

Loss function: Average penalty (over all training data) of best (i.e., lowest penalty) correct path (i.e. associated with correct sequence)

Sum up all penalties

$\mathbf{G}_{vit}$ — Viterbi Path

Select path with least cum. penalty

$\mathbf{T}_{vit}$ — Viterbi Transformer

Combined penalties ➜ handled via integrated training

$\mathbf{G}_{int}$ — Interpretation Graph

Recognizes individual characters on single images (e.g. CNN)

$\mathbf{T}_{rec}$ — Recognition Transformer

NN NN NN NN NN NN

Contains all possible segmentations (arcs: penalties & images)

$\mathbf{G}_{seg}$ — Segmentation Graph

class label

character recognizer penalty for each class

PIECE OF THE INTERPRETATION GRAPH

| | |
|---|---|
| "0" 6.7 | "0" 7.9 |
| "1" 10.3 | "1" 11.2 |
| | "2" 6.8 |
| | "3" 0.2 |
| "8" 0.3 | "8" 13.5 |
| "9" 12.5 | "9" 8.4 |

Character Recognizer    Character Recognizer

W

8  0.1     3  0.5

candidate segment image

penalty given by the segmentor

PIECE OF THE SEGMENTATION GRAPH

For each arc in $\mathbf{G}_{seg}$: Create new one per (character) class, penalty and label attached

**Problems**:
1. Trivial solution possible (Recognizer ignores input & sets all outputs to small values)
2. Penalty does not take competing answers into account (i.e., ignores training signals)
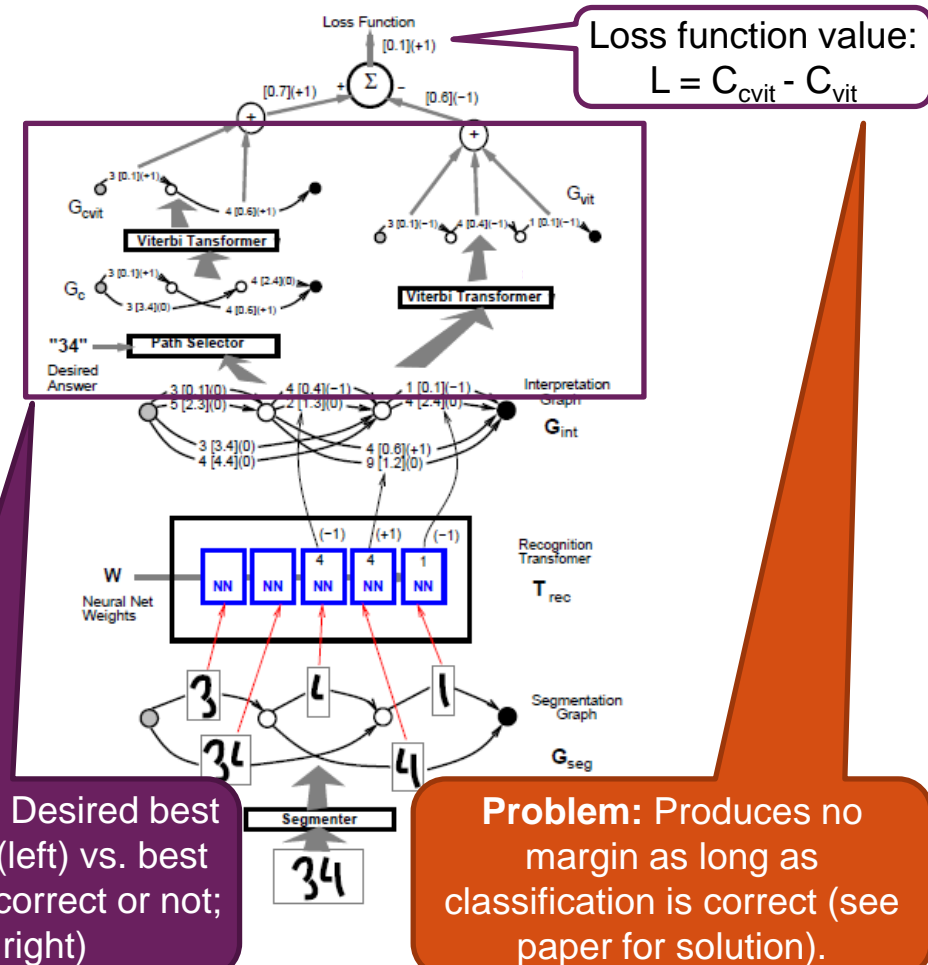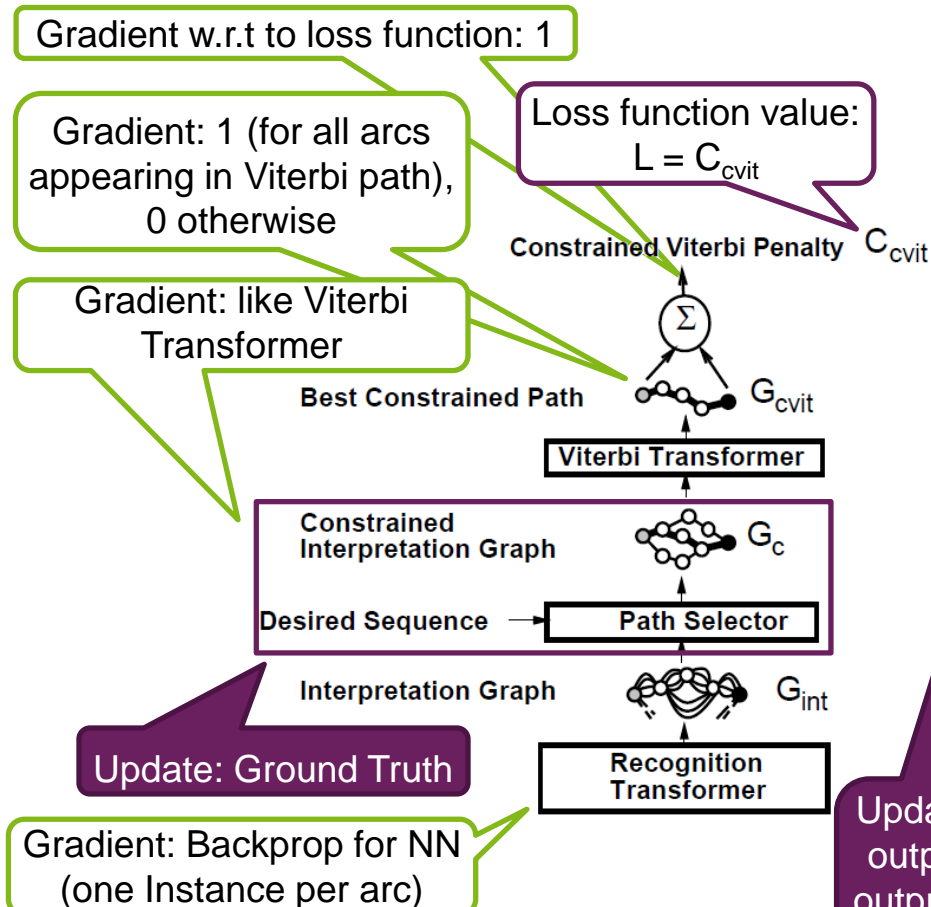
**Solved:** Discriminative training builds the class-"separating surfaces rather than modeling individual classes independently of each other" ➔ L=0 if best path is a correct path.

## «Viterbi» training

## Discriminative training

Gradient w.r.t to loss function: 1

Gradient: 1 (for all arcs appearing in Viterbi path), 0 otherwise

Gradient: like Viterbi Transformer

Loss function value:
$$L = C_{cvit}$$

Loss function value:
$$L = C_{cvit} - C_{vit}$$



Update: Ground Truth

Gradient: Backprop for NN (one Instance per arc)

Update: Desired best output (left) vs. best output (correct or not; right)

**Problem:** Produces no margin as long as classification is correct (see paper for solution).
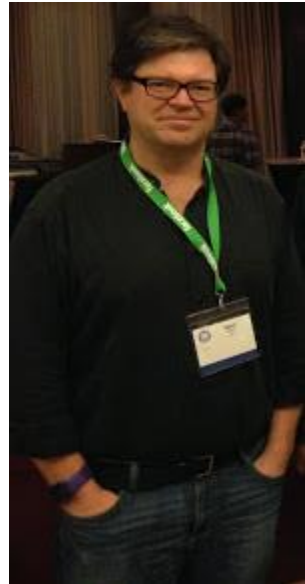
# Remarks

## Discriminative training

- Uses **all** available training **signals**
- Utilizes "penalties", **not probabilities**
  - ➔ **No** need for **normalization**
  - ➔ Enforcing normalization is *"complex, inefficient, time consuming, ill-conditions the loss function"* [according to paper]
- Is the *easiest/direct* **way** to achieve the objective of classification (**as opposed to Generative training**, that solves the more complex density estimation task as an intermediary result)

## List of possible GT modules

- All **building blocks of (C)NNs** (layers, nonlinearities etc.)
- **Multiplexer** (though not differentiable w.r.t. to switching input)
  ➔ can be used to dynamically rewire GTN architecture per input
- **min**-function (though not differentiable everywhere)
- **Loss** function

# Further Reading

- Original short paper: Bottou, Bengio & LeCun, "Global Training of Document Processing Systems using Graph Transformer Networks", 1997
http://www.iro.umontreal.ca/~lisa/pointeurs/bottou-lecun-bengio-97.pdf

- Landmark long paper: LeCun, Bottou, Bengio & Haffner, "Gradient-Based Learning Applied to Document Recognition", 1998
http://yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf

- Slide set by the original authors: Bottou, "Graph Transformer Networks", 2001
http://leon.bottou.org/talks/gtn

- Overview: Dietterich, "Machine Learning for Sequential Data: A Review", 2002
http://eecs.oregonstate.edu/~tgd/publications/mlsd-ssspr.pdf

- Recent work: Collobert, "Deep Learning for Efficient Discriminative Parsing", 2011
http://ronan.collobert.com/pub/matos/2011_parsing_aistats.pdf