# COMPUTATIONAL INTELLIGENCE
## DEEP LEARNING

## Performance Improvement of Deep Neural Networks

**Adrian Horzyk**

horzyk@agh.edu.pl

**AGH University of Science and Technology Krakow, Poland**

**When training DNN we usually struggle with the improvement of hyperparameters, structures and training models to achieve better training speed and final performance. We can try (some ideas):**

- **Collect more training data (and label them for supervised training).**

- **Diversify training data to represent a computational task better.**

- **Use different network architectures and different numbers of layers and neurons.**

- **Use different activation functions and different sequences of various layers.**

- **Experiment with various hyperparameters and try different combinations of them.**

- **User regularization, dropout, optimization methods (e.g. Adam optimizer).**

- **Train a chosen network longer with different or changing learning rates.**

**How to quickly and smarter choose between various training strategies?**

**We always have limited resources (time and computational power) to solve a given problem and must cut costs in the commercial implementations!**

# Orthogonalization

## Orthogonalization:

- Is a clear-eyed process about what to tune and how to achieve a supposed effect.

- Is the process that let us refer to individual hyperparameters in such a way that we can fix a selected training problem by tuning on a limited subset of hyperparameters.



- Why do we prefer to use drones over helicopters?

- Which one is easier to control and why?

- Is it easier to control a single knob changing a single parameter or a compound joystick changing many parameters at the time?

- Have you tried to fly a helicopter and/or a drone in the past? What is your experience?

**What about the car controllers like a weal, pedals, knobs, shifts and buttons?**

**Is it easier to control it (e.g. speed) when each parameter is controlled separately?**

**What do you prefer to control the car:**

- **set of controllers** (like weal, pedals, knobs, shifts and buttons) that control individual parameters of the car (speed, direction, etc.) or

- **an integrated controller** (like joystick) that can control a combination of parameters (like speed and direction) by the same move?

**When developing and training the model we usually follow the following chain of goals:**

1. **Fit a training set well** on a cost function trying to achieve the human-level of performance.

2. **Fit a dev set** (validation set) well on a cost function to get good generalization properties verified during the training process.

3. **Fit a test set** well on a cost function to be sure that the generalization is good enough and validated on the data that were not used during the training process.

4. Next, we hope that **the model will perform well in real world**.

- If the model does not fit well in any of the first three steps, we need to know what we can do with the model, its hyperparameters and the training to achieve the goal!

- Therefore, we want to define knobs (hyperparameters, optimization strategies that can help us) for each step to control the training process to fit the model well.

**In certain steps of the training process, we can use different knobs:**

| Fit a training set | Fit a dev set | Fit a test set | Fit on real-word data |
|---|---|---|---|
| Bigger network | Smaller network | Bigger dev set | Training data were not representative |
| Adam optimizer | Regularization | | Cost function not enough well defined |
| Xavier initialization | Bigger training set | | |
| Early stopping | Early stopping | | |
| Add new features | Add new features | Add new features | |
| | | | |
| | | | |

**When adapting the network we usually train it with different hyperparameters and comparing achieved precision and recall:**

- **Precision** – defines the percentage of correct classifications, e.g. if the achieved precision is 98% after the training is finished, and the network says that the input is a car, there is a 98% that it really is a car.

- **Recall** – is the percentage of correctly classified objects (inputs) for training classes, e.g. how many cars of all the cars from training data were correctly classified?

| Classifier | Precision | Recall | $F_1$ Score |
|------------|-----------|--------|-------------|
| Classifier A | 96% | 90% | 92,90% |
| Classifier B | **98%** | 88% | 92,73% |
| Classifier C | 94% | **93%** | **93,50%** |

$$F_1 = \frac{2}{Precision^{-1} + Recall^{-1}}$$

$$F_\beta = \frac{(1 + \beta^2) \cdot Precision \cdot Recall}{\beta^2 \cdot Precision + Recall} =$$

$$= \frac{(1 + \beta^2) \cdot TP}{(1 + \beta^2) \cdot TP + \beta^2 \cdot FN + FP}$$

$\beta$ - how many times recall is more important than precision
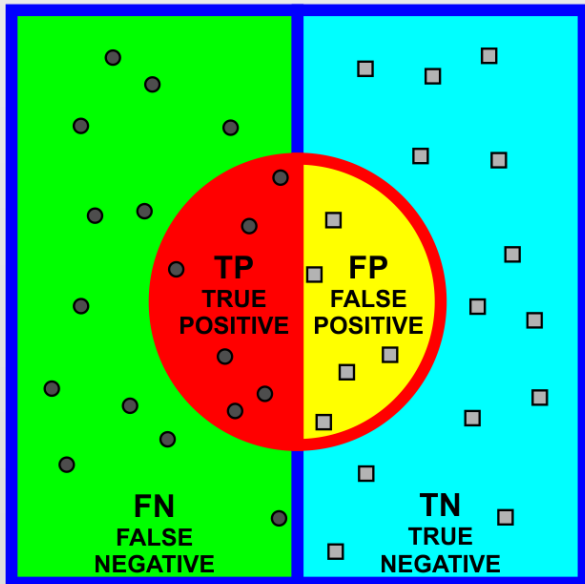TP – true positive, FP – false positive, FN – false negative

- **Which classifier from the above three is the best one?**

- **It turns out that there is often a trade-off between precision and recall, but we want to care about both of them!**

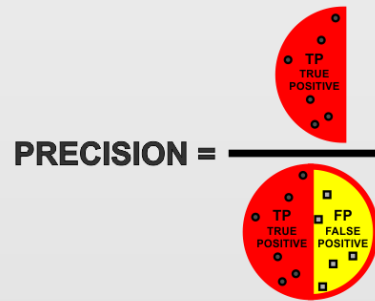- **We sometimes use F1 Score that is a harmonic mean of the precision and recall.**

# Confusion Matrix groups the results of classification:

- **TP (true positive)** – is the number of examples correctly classified as positive (class A).

- **FP (false positive)** – is the number of examples incorrectly classified as positive (class A).

- **TN (true negative)** – is the number of examples correctly classified as negative (class B).

- **FN (false negative)** – is the number of examples incorrectly classified as negative (class B).
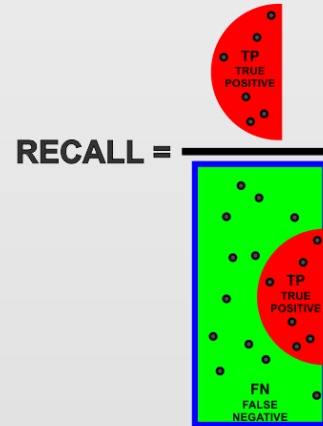
● GROUND-TRUTH POSITIVE
□ GROUND-TRUTH NEGATIVE



**TP** TRUE POSITIVE

**FP** FALSE POSITIVE

**FN** FALSE NEGATIVE

**TN** TRUE NEGATIVE

CLASSIFIED AS POSITIVE
CLASSIFIED AS NEGATIVE

PRECISION =

RECALL =

ACCURACY=

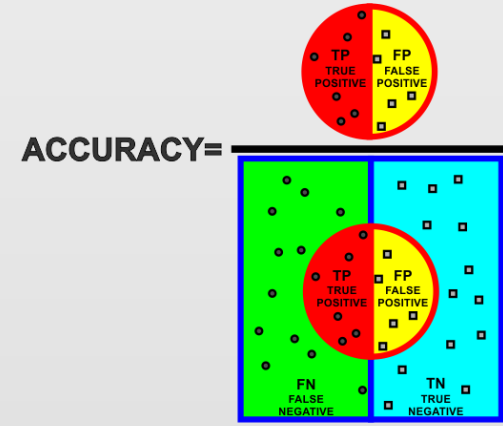**Precision** = TP / (TP + FP)

**Precision** – a ratio of how many examples were correctly classified as positive (class A) to all examples classified as positive (while not all are really positive, i.e. of class A).

**Recall** = TP / (TP + FN)

**Recall** – a ratio of how many examples were correctly classified as positive (class A) to all positive (class A) examples in the training set.

**Accuracy** = (TP + TN) / ALL

**Accuracy** – a ratio of how many examples were correctly classified to all examples in the training set.

## The most popular measures of results are:

| CONFUSION MATRIX | | | Prevalence = PP / ALL |
|---|---|---|---|
| All Examples (ALL) | Defined as Positive (P) = TP + FN | Defined as Negative (N) = FP + TN | Accuracy (ACC) = (TP + TN) / ALL |
| Predicted as Positive (PP) = TP + FP | True positive (TP) | False Positive (FP) | Precision = Positive Predictive Value (PPV) = TP / PP |
| | | | False Discovery Rate (FDR) = FP / PP |
| Predicted as Negative (PN) = FN + TN | False Negative (FN) | True Negative (TN) | False Omission Rate (FOR) = FN / PN |
| | | | Negative Predictive Value (NPV) = TN / PN |
| | True Positive Rate (TPR) = | False Positive Rate (FPR) = | $F_1 = 2 \cdot$ Precision $\cdot$ Recall / (Precision + Recall) |
| | Recall = Sensitivity = TP / P | Fall-out = FP / N | Positive Likelihood Ratio (LR+) = TPR / FPR |
| | False Negative Rate (FNR) = | True Negative Rate (TNR) = | Negative Likelihood Ratio (LR−) = FNR / TNR |
| | Miss Rate = FN / P | Selectivity = Specificity = TN / N | Diagnostic Odds Ratio (DOR) = LR+ / LR− |

**When we have results collected by many classifiers, we need to choose the best one, preferably using a single criterion that takes into account e.g. various positive or negative classifications for all classes separately:**

- **Compute the average error or harmonic mean for all classes and classifiers to compare them:**

| Classifier | Class A | Class B | Class C | Class D | Average | Harmonic Mean |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| A | 95% | 90% | 94% | 99% | 94.5% | 94.39% |
| B | 96% | 93% | 97% | 94% | 95.0% | 94.97% |
| C | 92% | 93% | 95% | 97% | 94.3% | 94.21% |
| D | 94% | 95% | 99% | 94% | 95.5% | 95.46% |
| E | 97% | 98% | 95% | 97% | 96.8% | 96.74% |
| F | 99% | 91% | 96% | 92% | 94.5% | 94.39% |

- **Thanks to such measures we can more easily point out the best classifier taking into account results collected for all classes.**

**Sometimes an application must run in real-time, so we cannot simply choose the classifier with the best accuracy, precision, or recall, but we must take into account the classification time:**

- **The accuracy must be the highest but available at the acceptable time, e.g. < 100 ms**

| Classifier | Accuracy | Classification Time |
|:---:|:---:|:---:|
| A | 94.5% | 70 ms |
| B | 95.0% | 95 ms |
| C | 94.3% | 35 ms |
| D | 95.5% | 240 ms |
| E | 96.8% | 980 ms |
| F | 94.5% | 60 ms |

- **Sometimes we must take into account various criteria to find out the suitable classifier, e.g. we choose this one with the highest accuracy if its classification time is lower than 100 ms.**

- **The accuracy is optimized, while the classification time must be satisfied.**

- **So we have to do with multi-criteria optimization.**

**Train, dev (validation), and test set should be set up in such a way that they share data of all distributions in the same way:**

- When we would like to create a classifier (or another predictor) for data coming from the various data distributions, e.g. the whole world or different countries of the company, we should take care about way how the data are distributed to the train, dev, and test sets. On the other hand, we can train the model almost perfectly on the train data and validate it on the dev set, but it will not work on a test set!

**Examples:**

- If we train and validate the model on data coming from rich people, it will rather not work for people with low incomes and vice versa.

- If we train and validate the model for men, it will rather not work for women and vice versa.

- If we train and validate the model for data coming from Europe, it will rather not work data coming from China or US and vice versa.

**The train, validation and test target must be the same, i.e. train, dev, and test data must be taken from the same data distributions, i.e. they must be representative for the solved problem.**

## Old way of splitting data (for small datasets < 100 thousands), e.g.:

* Train set : dev set : test set = 60% : 20% : 20%

* Train set : dev set : test set = 70% : 15% : 15%

* Train set : dev set : test set = 80% : 10% : 10%

## New way of splitting data (for large datasets used in deep learning):

* Train set : dev set : test set = 98% : 1% : 1%

* Because training data today have huge amount of training samples (> 1.000.000), so 1% is enough for validation or testing (1% from 1.000.000 is 10.000 of validation or testing examples), and thanks to it we can use more examples (data) for training!

* The test set should be big enough to give high confidence in the overall performance of the trained system or solved task.

**Training examples might be treated in the same or different way, i.e. they can influence the training process with different strength:**

- We can modify the definition of the error function in such a way to add the **strengthening factor $s^{(i)}$** for each training example to let it influence on (impact) the training process with a different strength:

- $J(w, b) = \frac{1}{\sum_{i=1}^{m} s^{(i)}} \sum_{i=1}^{m} s^{(i)} \cdot L\left(a^{(i)}, y^{(i)}\right)$

- In this way, we can avoid some unwanted classifications.


**Data attributes** can also have **different values for the training process** (e.g. we want the gender not to influence the training process much), so we can bind different strengths with different attributes, weakening these which should have reduced influence on the classification process and strengthening those which are especially important from the classification point of view.

**Sometimes we even avoid using some attributes** like race, sex, age, disabilities, health condition, political or religious affiliation etc. not to discriminate some groups of people or other objects because of the law or equal opportunities.

**Suppose that we try to classify some medical images:**



A chest X-ray

**The classification can be made by different humans or their teams:**

for some groups of training examples

| Classification made by: | | Produces the error |
|---|---|---|
| (a) | Typical human | 25% |
| (b) | Typical doctor (expert) | 4% |
| (c) | Experienced doctor (expert) | 1% |
| (d) | Team of experienced doctors (experts) | 0.4% |

**Human-level error (d) is defined as the lowest possible error that might be achieved by the human team of the best experienced experts. We assume that nobody contemporarily can do it better!**
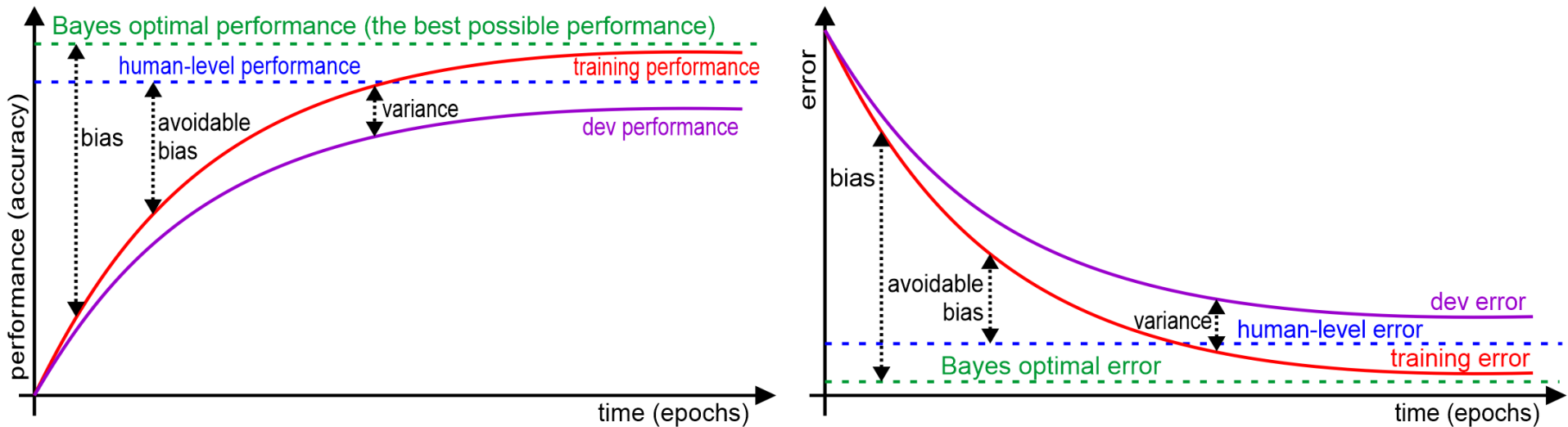
# Human-level Performance

**Avoidable bias** is an error defined by the difference between the training error and the human-level error:

avoidable bias = training error – human level error

**Bias** is an error defined by the difference between the training error and the Bayes error: bias = training error – Bayes level error

**Variance** is an error defined by the difference between the dev error and the training error: variance = dev error – training error



**Human level-error and performance** are defined by the human team of the world's best experts:

- Then cannot be contemporary surpassed by any human or the human team.

- If they would be surpassed in the future, then they automatically set a new human-level performance and a new human-level error.

**Bayes optimal (the best possible) error and performance** are defined by the blurred and noisy training examples that nobody and nothing can recognize or differentiate them due to the low quality:

- They can be never surpassed: Bayes optimal error ≤ Human-level error   and   Bayes optimal performance ≥ Human-level performance

- The Bayes optimal performance can be higher than a human-level performance.

- The Bayes optimal error can be lower than a human-lever error.

- It is many times very close to the human-level performance, where humans are very good at.

- Sometimes it is equal to the human-level performance when data are labelled by humans, so we cannot surpass this level in these cases.

**Analysis of the collected results and error levels allows us to look for a better solution by the implementation of some tips and tricks to improve the model performance.**

**Let's analyse a few examples:**

| Error Type | Model A | Model B | Model C | Model D |
|---|---|---|---|---|
| Bayes error* | 0.5% | 0.5% | 0.5% | 0.5% |
| Human-level error | 0.7% | 0.7% | 0.7% | 0.7% |
| Bias / Avoidable Bias | 3.0% / 2.8% | 3.0% / 2.8% | 0.5% / 0.3% | 0.3% / 0.1% |
| Training error | 3.5% | 3.5% | 1.0% | 0.8% |
| Variance | 4.3% | 0.8% | 4.0% | 0.2% |
| Dev error | 7.8% | 4.0% | 5.0% | 1.0% |
| Conclusion: | High Bias High Variance First focus on Bias | High Bias Low Variance Focus on Bias | Low Bias High Variance Focus on Variance | Low Bias & Variance Quite well-trained model |

\* Bayes error is not always known because the human-level abilities sometimes disallow to determine it. It might be experimentally determined or known due to the constructed training data set.

**Tips and methodology:**

**We should not try to decrease variance if the avoidable bias is still high.**
**First, we should always decrease bias and when it is low enough, start decreasing variance.**

During the model development and tuning, we try to minimize bias and variance.

Bias tells us about the ability of the model to adapt to the training data.

- Bias is the basic evaluation of the quality of the constructed model. If bias is poor, the model should be reconstructed and/or we should use bias-rising methods to minimize it.

Variance defines the generalization property of the model and tells us how well it can generalize about training data, dealing with a dev set and probably also with a test set and real-world data.

- If variance is poor, the constructed model is useless, because the main goal of machine learning is the generalization that allows us to use the trained models to real-world data.

- We should reconstruct the model and/or adapt variance-rising methods to achieve smaller variance.

Human-level performance reflects the quality of training data and the wisdom and experience of the humanity to solve the problem of a given kind.

- If the training data quality is poor or training data are contradictory, we cannot achieve better performance because nobody (event the team of human experts) can do the given task better.

- We can try to rise the quality of the training data (train set) to rise the human-level performance.

- Sometimes training data are described/defined by too small subset of attributes that do not allow to differentiate them enough (ambiguity producing contradictions) and to discriminate them in the classification process. In this case, we should redefine the train set adding new attributes (features) that describe the train objects in such a way that the diversity of them allow us to discriminate them.

**Generally, it is not easy to surpass the human-level performance, especially for various perceptions, speech and image recognition etc.**

**Surpassing human-level performance is possible for many problems:**

- **Product recommendations**

- **Online advertising**

- **Predicting transit time in logistics**

- **Loan approvals**

- **Many big data problems where humans cannot analyse them**

- **Non-natural perception tasks that were not evolved in humans over millions of years**

- **Various structural data requiring complex comparisons and analysis**

- **etc.**

**It occurs when the achieved training error (e.g. 0.3%) is less than the human-level error (e.g. 0.5%).**

**It may be difficult to establish how much the human-level performance might be overpassed and what would be the final performance and the bias of training that could be still avoided.**

# Guidelines for Minimizing Errors

**What can we do when the quality indicators of the model are low?**

**What are the guidelines for minimizing errors and improving performance?**

| Bayes error or Human-level error is high | Training error is high Bias is high | Dev (validation) error is high Variance is high |
|---|---|---|
| Clean data (remove possible bugs, inaccuracy, blurred data, outliers etc.) | Better or bigger network architecture (more parameters to enrich the too simple model) | Better or smaller network architecture (less parameters avoiding overshooting) |
| Add discriminating attributes | Xavier initialization | Regularization (L2, dropout) |
| Remove ambiguity (contradictory training data) | Use better optimization algorithms (e.g. Adam, RMSprop, momentum) | Bigger training set of the same distributions as the dev set |
| Reconstruct training data | Train longer (more epochs) | Early stopping |
|  | Transfer learning | Transfer learning |
|  | Better hyperparameter (network structure) search | Better hyperparameter (network structure) search |
|  |  | Data augmentation |

# Error Analysis and Overgoing Troubles

When you train the network, trying to implement various tips and tricks, but you are still unsatisfied of the achieved results, you can try to analyse results, e.g. incorrectly classified examples, and overgo these troubles implementing special routines into them:

- Check to which classes belong incorrectly classified examples? Are they of one or more classes? Do one class patterns prevail in them or not?

- Focus your effort on the most numerous incorrectly classified examples of one class because it can help you to decrease the error the most (ceiling) if you succeed.

- Are trained classes represented evenly in the training set? If not, try to balance the size of all classes, e.g. using augmentation to the less numerous classes or to reduce unevenly the learning rates implemented to various classes taking into account the number of examples which represent them.

- You can try to strengthen the training process for the incorrectly classified examples, e.g. use different learning rates for various training examples, i.e. the bigger learning rates for examples that are difficult to train.

- Check what the neurons of the network represent and whether the classification is not based on the object surrounding instead of the classified object self.

- Finally, try to find out all possible categories of errors and count up their occurrences:

| Example | Too big | Blurry | Mislabeled | Cars | Data Distribution 1 | Weak representation of this class | Comments |
|---------|---------|--------|------------|------|---------------------|-----------------------------------|----------|
| 1 | | ✓ | ✓ | | | | |
| 2 | | | | ✓ | | ✓ | |
| … | | | | | ✓ | | |
| % of total: | 15% | 42% | 18% | 32% | 12% | 18% | |

# Cleaning and Correcting Mislabeled Data

**Deep learning algorithms are usually robust, so the random errors and mislabeled training data should not spoil much the training process, but if there is a lot of incorrectly labeled data, they should be corrected:**
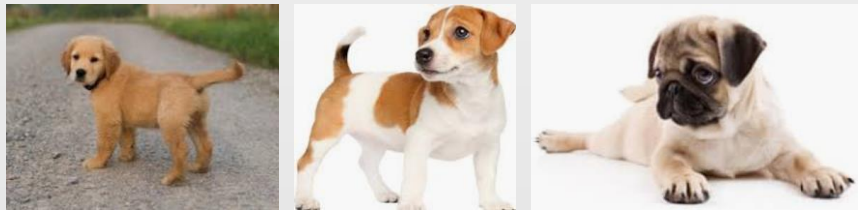
- **How to correct the training set when it consists of thousands/millions of examples?**

- **If the number of mislabeled examples is not too big (> 10%), we can try to learn the model using all correctly and incorrectly labeled examples, then filter out all misclassified examples and correct or remove those which are mislabeled, next, continue or start the training process from scratch again and again until we correct enough mislabeled data and achieve satisfying results of training the model.**

- **We can also use unsupervised training method to cluster training data, next, in each cluster, filter out all differently labeled examples to the most numerous class(es) represented be each cluster, and correct the mislabeled examples.**

- **If training data contain blurry or misleading examples, we can also remove them from the training set (cleaning it). Such examples are removed during the error analysis of the filtered out incorrectly classified examples. After removing of such examples, we start the training process again and again until we remove enough poor-quality examples and achieve satisfying results of training the model.**

**When training the model using data from different distributions, we should construct training, dev and test sets from all distributions!**

**EXAMPLE (a possible data mismatch problem):**

**High-quality data distribution**

**Low-quality data distribution**



Distribution 1: 100000 examples

Distribution 2: 20000 examples

**Take the data from both distributions together and shuffle them:**

Distribution 1+2: 120000 examples

**Next, split them into training, dev and test sets:**
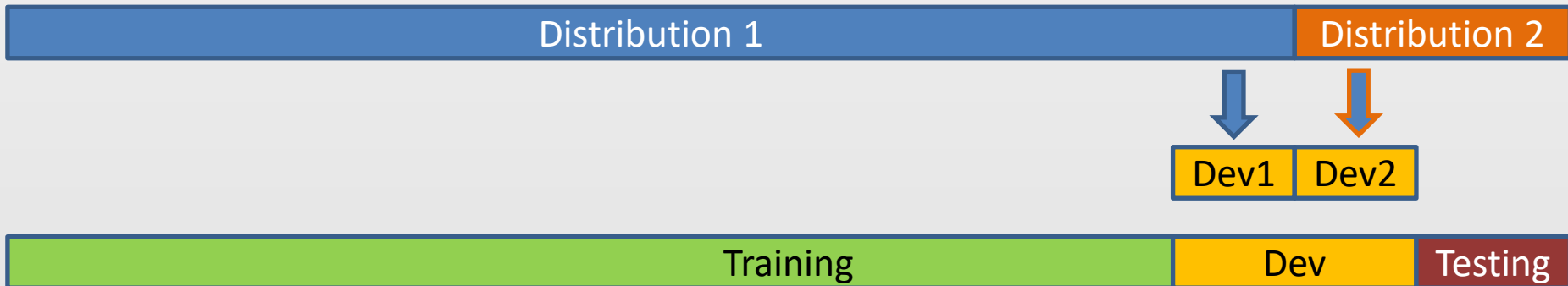
Training | Dev | Testing

**If you don't put the data together and shuffle them, they might be trained, validated, and tested on different distributions and dev and test results might be very poor:**

Distribution 1 | Distribution 2

Training | Dev | Testing

**If the dev set is composed from various distributions, but the training set is taken only from one distribution (or a subset of distributions), then dev error on the subsets of dev sets will differ!**

| Distribution 1 | Distribution 2 |
|---|---|

| Dev1 | Dev2 |
|---|---|

| Training | Dev | Testing |
|---|---|---|

**In this case, we usually achieve the following:**

human-level error < training error < **dev1 error < dev error < dev2 error** < testing error

This indicates the **data mismatch problem**, i.e. the model has been trained on a limited subset of distributions (not all distributions).

The difference between testing error and dev error indicates **the overfitting problem.**

When dev data, testing data or real-world data differ from training data, we can try to **artificially synthetize new training data** which will be more similar to real-world data (noise data augmentation), e.g.:

- We can add typical noise to training data.

- Blur training data.

- Add some distortions to training data.

# Let's start to change hyperparameters!

- ✓ **Improving performance of the training**
- ✓ **Speeding up the training process**
- ✓ **Not stacking in local minima**
- ✓ **Using less computational resources to get the model**
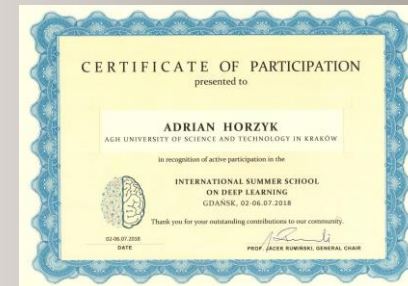
# Bibliography and Literature

1. Nikola K. Kasabov, *Time-Space, Spiking Neural Networks and Brain-Inspired Artificial Intelligence*, In Springer Series on Bio- and Neurosystems, Vol 7., Springer, 2019.
2. Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning*, MIT Press, 2016, ISBN 978-1-59327-741-3 or PWN 2018.
3. Holk Cruse, *Neural Networks as Cybernetic Systems*, 2nd and revised edition
4. R. Rojas, *Neural Networks*, Springer-Verlag, Berlin, 1996.
5. *Convolutional Neural Network* (Stanford)
6. *Visualizing and Understanding Convolutional Networks*, Zeiler, Fergus, ECCV 2014
7. IBM: https://www.ibm.com/developerworks/library/ba-data-becomes-knowledge-1/index.html
8. NVIDIA: https://developer.nvidia.com/discover/convolutional-neural-network
9. JUPYTER: https://jupyter.org/

**Adrian Horzyk**

**horzyk@agh.edu.pl**

**Google: Horzyk**

CERTIFICATE OF PARTICIPATION
presented to

ADRIAN HORZYK
AGH UNIVERSITY OF SCIENCE AND TECHNOLOGY IN KRAKÓW

in recognition of active participation in the

INTERNATIONAL SUMMER SCHOOL
ON DEEP LEARNING
GDAŃSK, 02-06.07.2018

**University of Science and Technology in Krakow, Poland**

AGH

# XXX

## XXXXX:

- xxxxx

# XXXXX:

- xxxxx