

**Mechatronic Engineering program:**  
**Python for machine learning and data science**

# Classification (revisited)

Ziemowit Dworakowski  
*AGH University of Krakow*

1

---

---

---

---

---

---

---

---

### What you should probably know already?

- Outliers**  
(their sources, ways to deal with them)
- Overlapping classes**  
(And how to deal with them)
- Overfitting**  
(And how to deal with it)
- Basic optimization methods**  
(Gradient, gradient with momentum, 1+1, grid search)
- Metaparameter configuration procedure**  
(2nd order optimization)
- Basic algorithms**  
(kNN, manually configured decision tree)

2

---

---

---

---

---

---

---

---

### Classification / regression analogy

The diagram illustrates the analogy between classification and regression in different dimensions. In 1D, classification involves a threshold, while regression involves a fitted function. In 2D, classification involves a decision boundary, and regression involves a fitted surface.

3

---

---

---

---

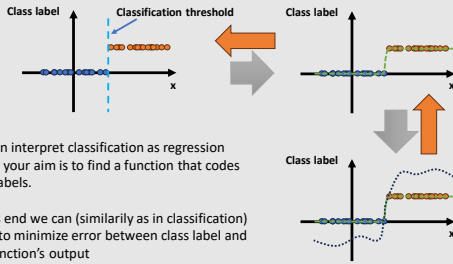
---

---

---

---

### Classification / regression analogy



You can interpret classification as regression where your aim is to find a function that codes class labels.

To this end we can (similarly as in classification) strive to minimize error between class label and this function's output

4

---

---

---

---

---

---

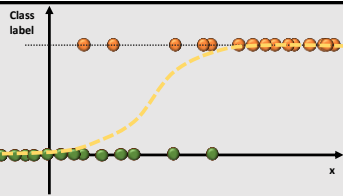
---

---

---

---

### Logistic regression



$$f(x) = \frac{1}{1 + e^{-(\beta_1 x + \beta_0)}}$$

This is again a sigmoid curve – here called a „logistic function“.  $\beta_1$  and  $\beta_0$  are parameters that need to be fit during training in such a way as to maximize **log likelihood** of correct label assignment.

We can use standard gradient-based optimization approach for that

5

---

---

---

---

---

---

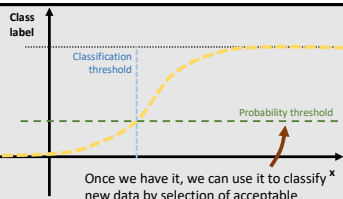
---

---

---

---

### Logistic regression



$$f(x) = \frac{1}{1 + e^{-(\beta_1 x + \beta_0)}}$$

This is again a sigmoid curve – here called a „logistic function“.  $\beta_1$  and  $\beta_0$  are parameters that need to be fit during training in such a way as to maximize **log likelihood** of correct label assignment.

We can use standard gradient-based optimization approach for that

6

---

---

---

---

---

---

---

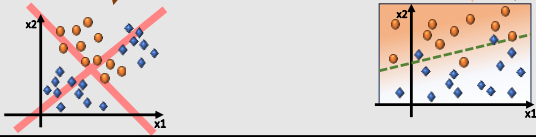
---

---

---

### Logistic regression

- LR is a classification equivalent of the linear regression
- It can work for multidimensional problems (we just have more parameters to learn)
- It can fit a „line classifier“ in the „best way possible“ (meaning that the classification has actually probabilistic interpretation)
- It cannot solve nonlinearly separable classification problems



7

---

---

---

---

---

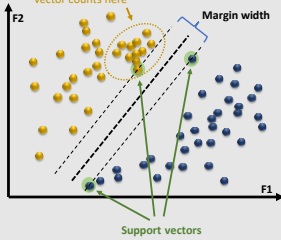
---

---

---

### SVM revisited: principle

It does not really matter that these points are so clumped together, only the support vector counts here



In regression SVM aimed to make margin as small as possible and fit all the data **inside** the margin.

In classification it does the opposite: it tries to make margin as wide as possible and keep all data **outside** of it

8

---

---

---

---

---

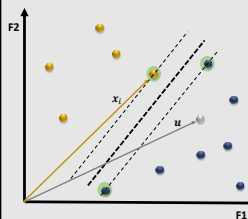
---

---

---

### SVM revisited: How does it do that?

- Let:
- $\omega$  - be a vector perpendicular to the margin
  - $x_i$  - be the  $i$ -th sample
  - $u$  - be an unknown sample
  - $y_i$  - be a class label (equal to either 1 or -1)
  - $\alpha_i$  - be a Lagrangian multiplier denoting whether a point should be treated as a support vector



In a nutshell, If we want to classify a new sample, we are using the following equation:

$$\sum_i \alpha_i y_i \vec{x}_i \circ \vec{u} + b \geq 0 \Rightarrow \text{THEN } \oplus$$

Note, that to classify a new point we calculate its dot product with all the support vectors – that will come in handy in a second...

9

---

---

---

---

---

---

---

---

If you want to know what is the source for this equation, here is a brief explanation. You don't need to learn that for the test though!

Decision rule:  
 $\bar{\omega} \circ \bar{u} + b \geq 0 \Rightarrow \text{THEN} \oplus$

For the training samples:  
 $\begin{cases} \bar{\omega} \circ \bar{x}_i + b \geq 1 \\ \bar{\omega} \circ \bar{x}_i + b \leq -1 \end{cases} \Rightarrow y_i(\bar{\omega} \circ \bar{x}_i + b) - 1 \geq 0$

For samples „on the margin“ (support vectors):  
 $y_i(\bar{\omega} \circ \bar{x}_i + b) - 1 = 0$

Margin width:  
 $(\bar{x}_i - \bar{x}_j) \circ \frac{\bar{\omega}}{\|\bar{\omega}\|} = \frac{2}{\|\bar{\omega}\|}$

We minimize  $\frac{1}{2} \|\bar{\omega}\|^2$   
 With respect to constraint  $\textcircled{1}$

$$L = \frac{1}{2} \|\bar{\omega}\|^2 - \sum_i \alpha_i [y_i(\bar{\omega} \circ \bar{x}_i + b) - 1]$$

$$\frac{\partial L}{\partial \bar{\omega}} = \bar{\omega} - \sum_i \alpha_i y_i \bar{x}_i = 0 \Rightarrow \bar{\omega} = \sum_i \alpha_i y_i \bar{x}_i$$

$$\frac{\partial L}{\partial b} = -\sum_i \alpha_i y_i = 0$$

After combination of the above equations, we have:  
 $L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \bar{x}_i \circ \bar{x}_j$

New decision rule:  
 $\sum_i \alpha_i y_i \bar{x}_i \circ \bar{u} + b \geq 0 \Rightarrow \text{THEN} \oplus$

10

---

---

---

---

---

---

---

---

---

---

### SVM revisited: Nonlinearly separable data

Nonlinearly separable data

„Soft margin“ idea

We can allow the algorithm to ignore a few points  
*(Sometimes ignoring just a few outliers makes the problem linearly separable)*

We minimize:  

$$\left[ \frac{1}{2} \sum_{i=1}^n \max(0, 1 - y_i(\bar{\omega} \cdot \bar{x}_i - b)) \right] + \lambda \|\bar{\omega}\|^2$$

Previous constraint now is just a weighted part of the objective function

$\lambda$  allows to balance importance of the margin width and influence of outliers

We classify data accepting small error rate

11

---

---

---

---

---

---

---

---

---

---

### SVM revisited: Nonlinearly separable data

„kernel trick“

Kernel  $K(x_1, x_2)$  Defines „similarity measure“ between two vectors of features. Examples of kernels include:

- Polynomial functions:  
 $K(x_1, x_2) = (1 + x_1^T x_2)^p$
- Radial-basis functions (RBF) (equivalent to weighted kNN model):  
 $K(x_1, x_2) = \exp(-\gamma \|x_1 - x_2\|^2)$
- Sigmoids: (equivalent to single-layer MLP):  
 $K(x_1, x_2) = \tanh(\beta_0 x_1^T x_2 + \beta_1)$

We can map the data nonlinearly into space with higher number of dimensions by using Kernel functions instead of dot products  
*(There is a chance that data that were not separable in euclidean space will be separable in our new kernel-based space)*

We classify data using linear SVM

12

---

---

---

---

---

---

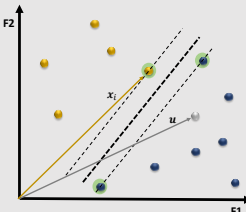
---

---

---

---

### SVM revisited: How does it do that?



Let:

- $\omega$  - be a vector perpendicular to the margin
- $x_i$  - be the  $i$ -th sample
- $\tilde{u}$  - be an unknown sample
- $y_i$  - be a class label (equal to either 1 or -1)
- $\alpha_i$  - be a Lagrangian multiplier denoting whether a point should be treated as a support vector
- $b$  - is a margin width parameter

In a nutshell, if we want to classify a new sample, we are using the following equation:

$$\sum_i \alpha_i y_i \tilde{x}_i \circ \tilde{u} + b \geq 0 \Rightarrow \text{THEN } \oplus$$

If we want to use nonlinear SVM, instead of dot product, we are using kernel function:

$$\sum_i \alpha_i y_i K(\tilde{x}_i, \tilde{u}) + b \geq 0 \Rightarrow \text{THEN } \oplus$$

A good visualization on how kernels affect data classification:  
<https://www.youtube.com/watch?v=Q7vTD-5VII>

13

---

---

---

---

---

---

---

---

---

---

### SVM revisited: How does it do that?

A good visualization on how kernels affect data classification:  
<https://www.youtube.com/watch?v=Q7vTD-5VII>

Very good lecture introduction to SVMs can be found here:  
[https://www.youtube.com/watch?v=\\_PwhiWxHK8o](https://www.youtube.com/watch?v=_PwhiWxHK8o)

Paper on SVM implementation in scikit is here (beware of challenge!):  
<https://www.csie.ntu.edu.tw/~cjlin/papers/libsvm.pdf>

14

---

---

---

---

---

---

---

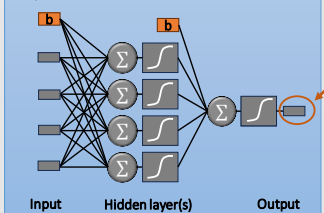
---

---

---

### ANN revisited: MLP for classification

All of this works exactly as in regression (same structure, training methods, configuration, etc.)



The only difference is here: we are using target data with labels (categorical binary variables) instead of numerical values

Input      Hidden layer(s)      Output

15

---

---

---

---

---

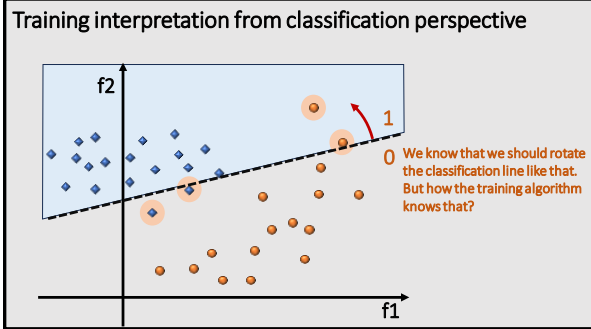
---

---

---

---

---



16

---

---

---

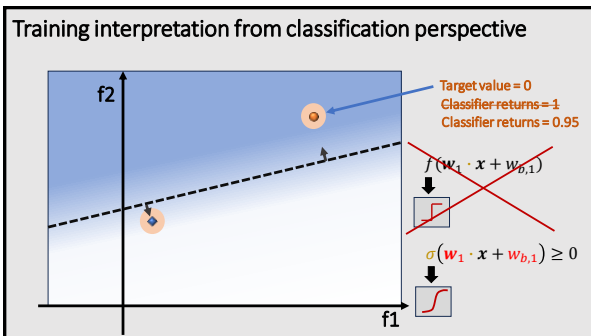
---

---

---

---

---



17

---

---

---

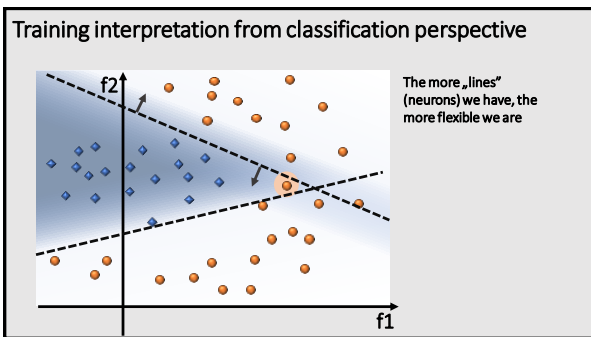
---

---

---

---

---



18

---

---

---

---

---

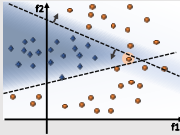
---

---

---

### ANN training: (Vanilla) gradient descent

*(standard, without additions)*



(Start from a random set of weights)

↓

Calculate gradient of error with respect to all the weights

↓

Modify the weights in the direction of the steepest gradient descent

19

---

---

---

---

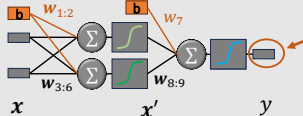
---

---

---

---

Consider a simple Multilayer Perceptron (MLP) network:



We also have a target  $t$  associated with each input vector  $x$ . Using that, we can calculate error  $d$ :

$$d = \frac{1}{2} (y - t)^2$$

Now we can calculate gradient of  $d$  with respect to weights  $w$ :

$$\nabla d = \left[ \frac{\partial d}{\partial w_1}, \frac{\partial d}{\partial w_2}, \frac{\partial d}{\partial w_3}, \dots, \frac{\partial d}{\partial w_9} \right]$$

The gradient now tells us how to adjust weights

$y = \sigma(w_8 \cdot x'_1 + w_9 \cdot x'_2 + w_7)$   
 $x'_1 = \sigma(w_3 \cdot x_1 + w_5 \cdot x_2 + w_1)$   
 $x'_2 = \sigma(w_2 \cdot x_1 + w_6 \cdot x_2 + w_2)$

You can combine this to get a  $f(w, x)$  function equivalent for a network

20

---

---

---

---

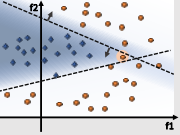
---

---

---

---

### ANN training: Adam optimizer



We now have a method that will be gradually dragging us towards a better classification accuracy.

What can we do to improve on it further?

- 1) We can accelerate optimization using averaged gradients (momentum): consecutive steps in the same direction cause the algorithm to speed up)
- 2) We can make the optimization more robust by including an average of square gradients of the past (oscillations tend to be stopped rapidly)
- 3) We can just use a sign of the derivative for each weight instead of its actual value (speedup for saddle points)

**Adaptive moment estimation (ADAM)**

21

---

---

---

---

---

---

---

---

### ANN training: Adam optimizer

Adam metaparameters:

- $\beta_1$  - Decay rate for averaged gradients (default: 0.9)
- $\beta_2$  - Decay rate for averaged gradient squares (default: 0.999)
- $\alpha$  - Learning rate (default: 0.001)
- $\epsilon$  - Small utility constant (default:  $10^{-9}$ , don't change)

\* If you want to know more (including math),

this is a good (basic) article on Adam optimizer here:

<https://www.geeksforgeeks.org/intuition-of-adam-optimizer/>

\* If you want to actually learn the method at the source, this article introduces Adam optimizer (it may be challenging for beginners though):

<https://arxiv.org/pdf/1412.6980.pdf>

22

---

---

---

---

---

---

---

---

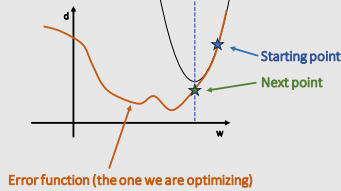
---

---

### ANN training: Limited memory BFGS (LM-BFGS)

(Broyden-Fletcher-Goldfarb-Shanno)

Instead of focusing on gradient, we can also use a Hessian (matrix of second derivatives) } So called: „Quasi-Newton“ or „Newton-based“ methods



23

---

---

---

---

---

---

---

---

---

---

### ANN training: Limited memory BFGS (LM-BFGS)

(Broyden-Fletcher-Goldfarb-Shanno)

Instead of focusing on gradient, we can also use a Hessian (matrix of second derivatives) } So called: „Quasi-Newton“ or „Newton-based“ methods



This allows often to outperform first-order-based methods (like Adam) provided that the input space is small (not many data points, not many dimensions)

LM-BFGS is one of the many algorithms from this family – available in popular libraries for ML in Python

A paper with a nice overview on second-order optimization algorithms : <https://cs.nyu.edu/~overton/mstheses/skajaa/msthesis.pdf>

24

---

---

---

---

---

---

---

---

---

---



ANN training: **Regularization**

*„How to prevent overfitting in such a way so we don't have to stop training early?“*

25

---

---

---

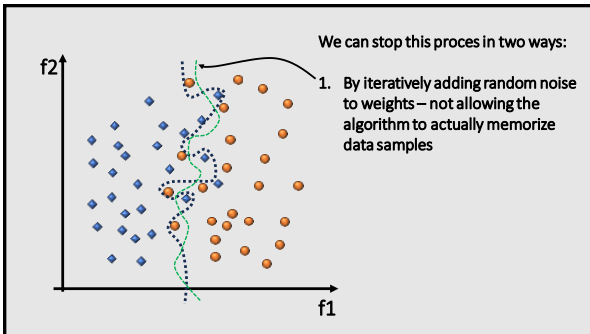
---

---

---

---

---



26

---

---

---

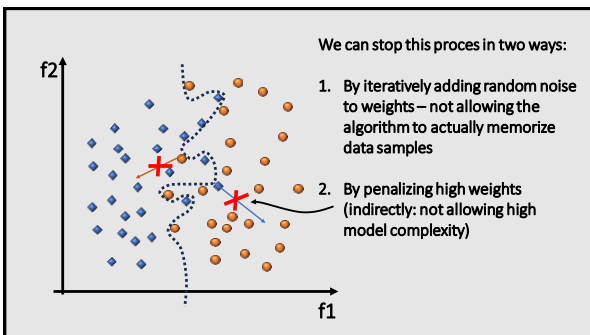
---

---

---

---

---



27

---

---

---

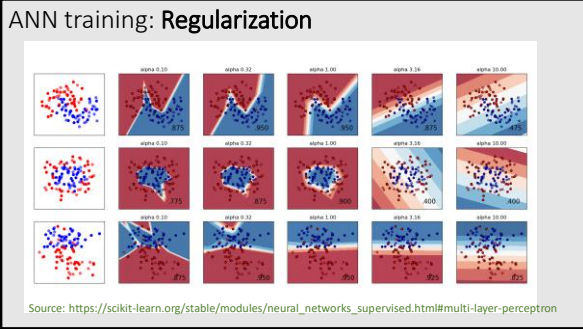
---

---

---

---

---



28

---

---

---

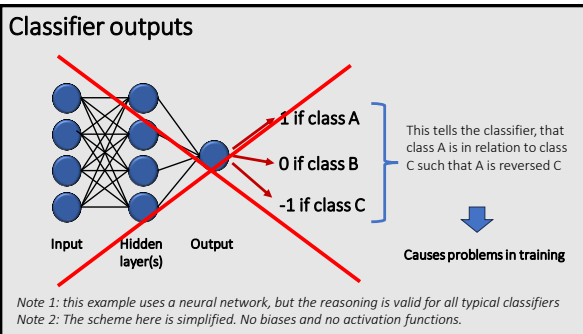
---

---

---

---

---



29

---

---

---

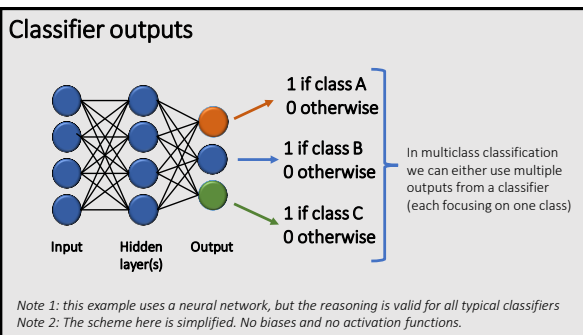
---

---

---

---

---



30

---

---

---

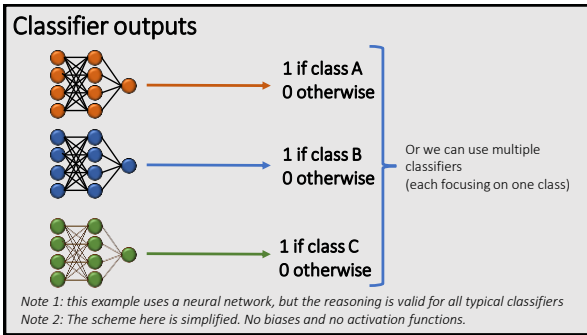
---

---

---

---

---



31

---

---

---

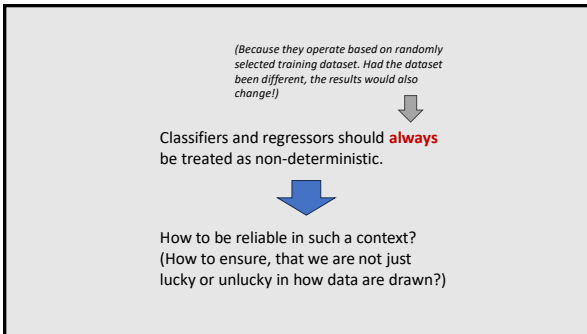
---

---

---

---

---



32

---

---

---

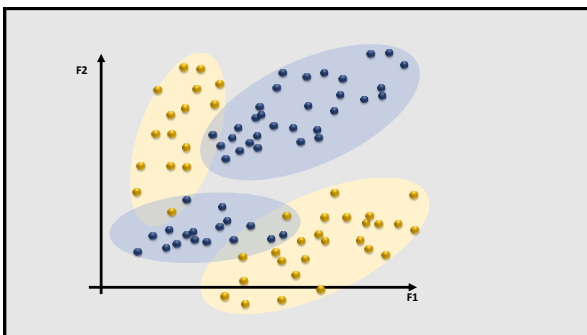
---

---

---

---

---



33

---

---

---

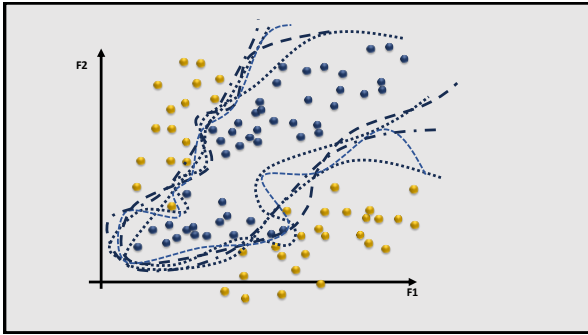
---

---

---

---

---



34

---

---

---

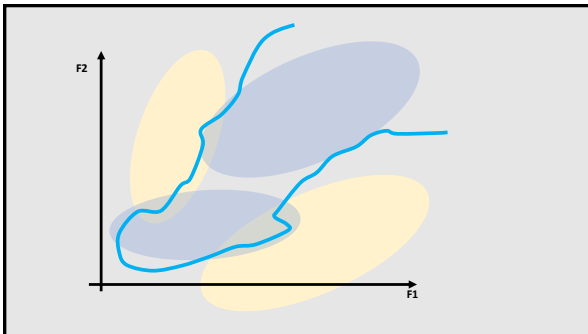
---

---

---

---

---



35

---

---

---

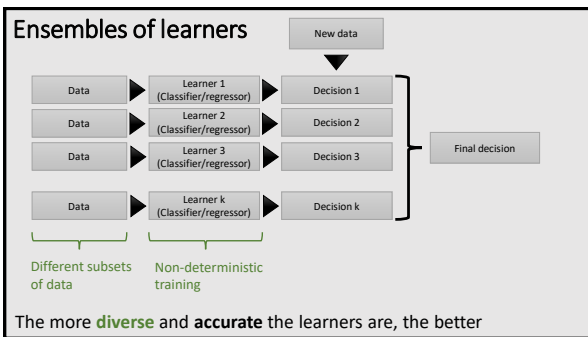
---

---

---

---

---



36

---

---

---

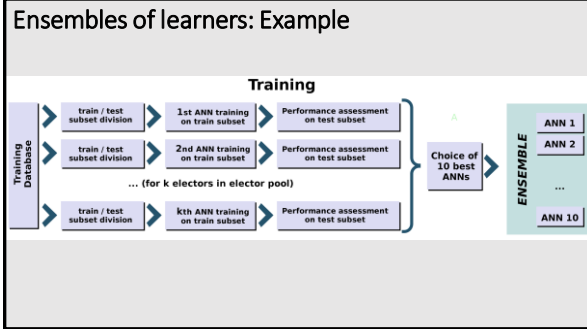
---

---

---

---

---



37

---

---

---

---

---

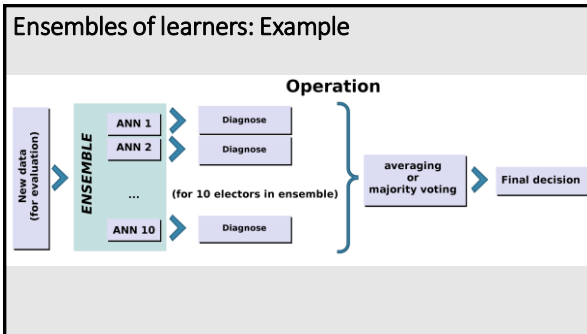
---

---

---

---

---



38

---

---

---

---

---

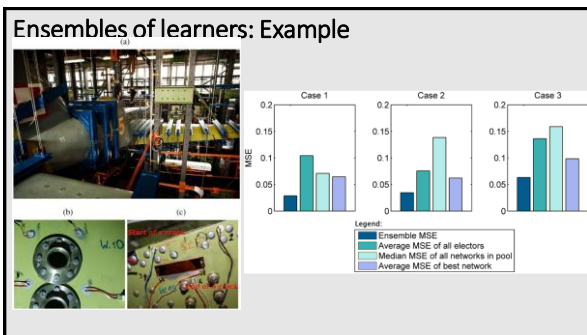
---

---

---

---

---



39

---

---

---

---

---

---

---

---

---

---

## Random forest



1. Do **Bagging** (Bootstrap Aggregating) by selecting repeatedly subsets of data using drawing with replacement
2. If dataset consists of many features, sample them also (let different subsets use also subsets of features!)
3. Train simple decision trees based on these subsets (each tree uses different subset)
4. Average responses from many trees

The more **diverse** and **accurate** the learners are, the better

40

---

---

---

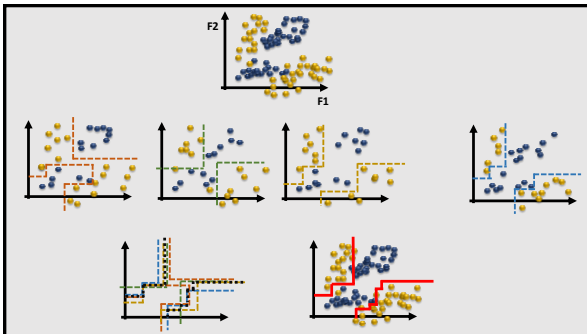
---

---

---

---

---



41

---

---

---

---

---

---

---

---

## Ensembles of learners

In general: Ensembles provide **increased reliability** at the cost of

- **lack of direct control over classification proces,**
- **almost no structural optimization possibilities,** and
- **a lot of required memory and time**

42

---

---

---

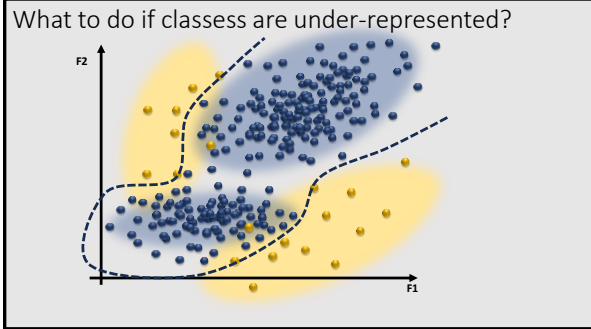
---

---

---

---

---



43

---

---

---

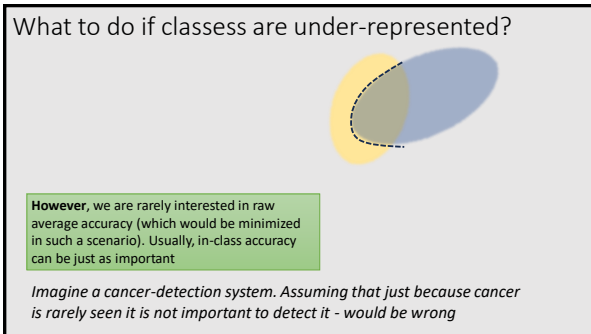
---

---

---

---

---



44

---

---

---

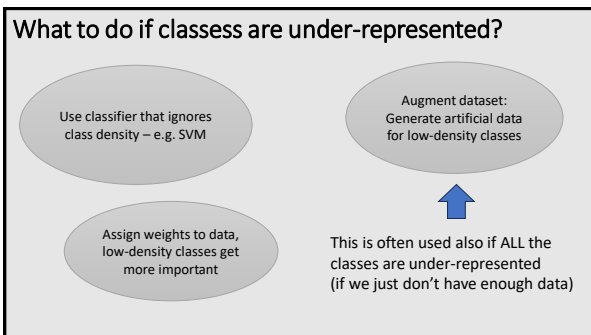
---

---

---

---

---



45

---

---

---

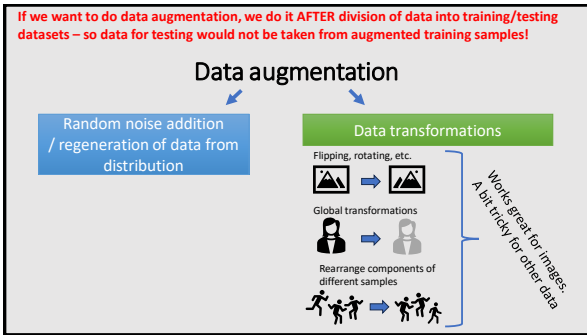
---

---

---

---

---



46

---

---

---

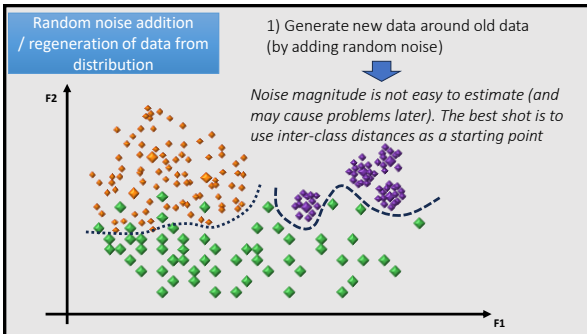
---

---

---

---

---



47

---

---

---

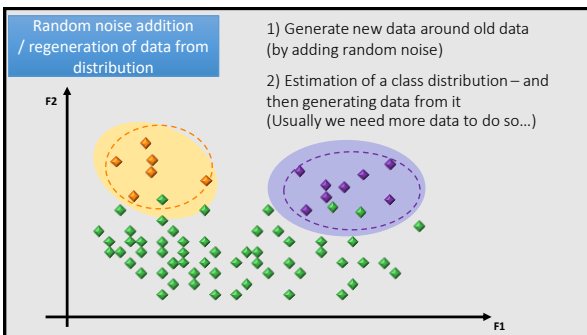
---

---

---

---

---



48

---

---

---

---

---

---

---

---



**Random noise addition / regeneration of data from distribution**

- 1) Generate new data around old data (by adding random noise)
- 2) Estimation of a class distribution – and then generating data from it (Usually we need more data to do so...)

↓

*If you have enough data to correctly derive its distribution – you usually have also enough data to do classification...*

*So it works mostly if you know additional context (what distribution can be assumed a priori)*

49

---

---

---

---

---

---

---

---

**One-class classification (novelty detection)**

50

---

---

---

---

---

---

---

---

**One-class classification (novelty detection)**

Sometimes we have a dataset with lots of examples belonging to just one class and none or almost none to the other classes:

- „**Healthy** people“
- „**Normal** behavior in subway“
- „**Normal** operation state of an assembly line
- „**Typical** weather conditions in September“

Then, our goal might be to learn this **normal** range of data, to detect any **anomalies (outliers, novelties)**

*In Novelty Detection we usually don't know what to expect (there is possibly infinite set of „norm breaches“)*

51

---

---

---

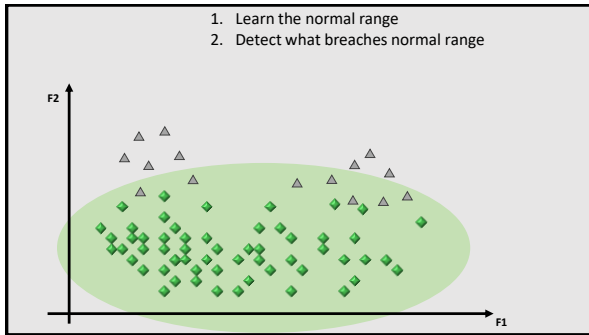
---

---

---

---

---



52

---

---

---

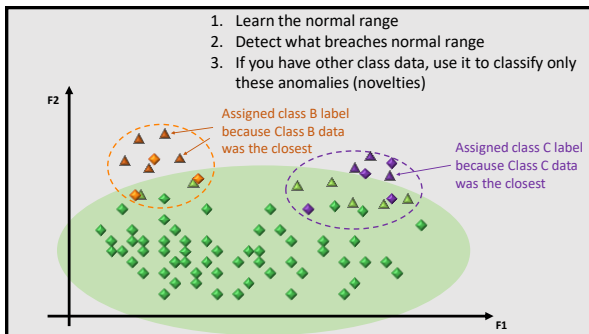
---

---

---

---

---



53

---

---

---

---

---

---

---

---

**One-class classification (novelty detection)**

How to learn „normal range“?

1. Estimate typical distances between samples within training set, detect anything that breaches it
2. Estimate underlying probability density for data – e.g. by fitting multidimensional gaussians into training data
3. Train a regression system to derive some of the features from others – in normal range (for known data) it should work very good, for novel samples it should produce large errors

54

---

---

---

---

---

---

---

---

### Things to remember\*:

1. Graphical interpretation of classification and regression
2. Logistic regression (idea, graphical interpretation, equation, features)
3. SVM: Principle of operation, Soft margin explanation, Kernel idea explanation
4. Graphical interpretation of MLP training (for one neuron)
5. MLP scheme
6. Basic ideas behind standard („vanilla“) gradient descent, ADAM and LM-BFGS
7. Explanation of two regularization methods
8. Configuration of ANNs for multiclass classification
9. Idea behind an ensemble approach for classification and its pros and cons
10. Steps and graphical interpretation of a random forest algorithm
11. Risk of having under-represented classes and three solutions to this problem
12. When can we do data augmentation (in relation to division of data into subsets)?
13. How can we augment data? Explain two methods to this end
14. What is a novelty detection problem? Explain three approaches to solve it

*\*You don't need to memorize equations if they are not explicitly mentioned in this list!*

---

---

---

---

---

---

---

---