



Machine Learning- Exercise 3

AdaBoost

Stephanie Kaes & George Lydakis
<lastname>@vision.rwth-aachen.de

RWTH Aachen University - Computer Vision Group
<http://www.vision.rwth-aachen.de/>

2023-12-06

Content



1. Cross-Validation
2. AdaBoost Recap
3. Exercise 3: AdaBoost
 - a) Simple Classifier
 - b) Simple AdaBoost
 - c) Cross-Validation AdaBoost
 - d) Least Squares AdaBoost
 - e) Application

Cross Validation



Two variants of CV are very popular:

Hold-out CV

Typically 80/20 i.e.

- Use 80% of the data for training
- Use 20% for testing

For ex.: *sklearn train-test-split function*

- No validation data
- Requires quite large datasets
- Problematic if data is imbalanced

K-folds CV

1. Split the training dataset into k equally sized training data and validation data subsets.
2. Train the model with the training data.
3. Validate the performance of the model using the validation data.
4. Repeat steps 1-3 or break

...but there are many more!

Cross Validation



Val	Train	Train	Train	Train
Train	Val	Train	Train	Train
Train	Train	Val	Train	Train
Train	Train	Train	Val	Train
Train	Train	Train	Train	Val

How do we apply k-folds CV ?

1. Split the training dataset into k equally sized training data and validation data subsets.
2. Train the model with the training data.
3. Validate the performance of the model using the validation data.
4. Repeat steps 1-3 (5-10 times) or break

Cross Validation



Benefits of k-folds CV?

- **Reduces Overfitting:** If a model consistently performs well across various subsets, it is more likely to generalize well.
- **Maximizes Data Utilization:** Instead of using only one fixed split, we use multiple combinations, ensuring that each data point is used for both training and testing at some point. Can be useful if test dataset is really small or data shows class imbalances.
- **Enhances Robustness:** Averaging across multiple folds reduces the impact of randomness in a single train-test split -> more stable performance estimate.

Cross Validation



Further Benefits and Applications of k-folds CV:

- **Model Performance Estimation:** Better assessment of model's performance on unknown data than single-split
- **Hyperparameter Tuning:** Selecting the best set of hyperparameters by assessing the model's performance across different parameter configurations using random/grid search etc.
- **Model Selection:** If multiple models shall be compared: Select the best-performing model by comparing the average performance across folds

Cross Validation



Warning

- k-folds CV might require many resources
- Be careful when datasets show an inherent order (time-series data etc.) – k-folds CV might not be applicable here
- k-folds CV assumes that the data points are independent and equally distributed which might not be true

How to use CV



1. Data Preparation

1. Split your dataset into features and target variable
2. Apply further preprocessing if necessary

2. Select a Cross-Validation Technique

3. Model Training and Evaluation Loop

For each fold or iteration of the cross-validation:

- Split the data into training and testing sets.
- Train your model on the training set.
- Evaluate the model on the testing set (apply your desired performance metrics)

4. Aggregate Results

Collect the performance metrics obtained from each fold or iteration.

5. Compute Average Performance

Calculate the average of the performance metrics obtained from all folds.

6. Analyze Results

AdaBoost: Main Idea

Components

- ▶ $c_k(\mathbf{x})$: “weak” or base classifier
 - ▶ Condition: $< 50\%$ training error over any distribution
- ▶ $C(\mathbf{x})$: “strong or final classifier”



AdaBoost: Main Idea

- Combine a weighted set of weak classifiers
- Resulting classifier should be better than the weak classifiers („Ensemble Learning“)



AdaBoost: Main Idea

Construct a strong classifier as a thresholded linear combination of the weighted classifiers:

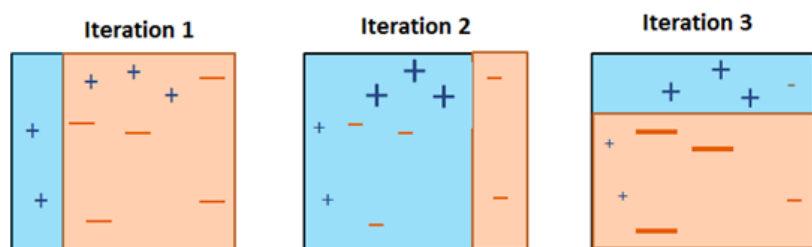
$$C(\mathbf{x}) = \text{sign} \left(\sum_{k=1}^K \alpha_k c_k(\mathbf{x}) \right)$$



AdaBoost: Main Idea

Main idea

- ▶ Instead of resampling, reweight misclassified training examples.
 - ▶ Increase the chance of being selected in a sampled training set.
 - ▶ Or increase the misclassification cost when training on the full set.



Data points weighted with $w_i^{(k)}$



$$H = \text{sign} \left(0.38 \times \begin{array}{|c|} \hline \text{orange} \\ \hline \end{array} + 0.58 \times \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} + 0.87 \times \begin{array}{|c|} \hline \text{orange} \\ \hline \end{array} \right)$$

$$= \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} \quad \text{Final Classifier / Strong Classifier}$$

$$C(\mathbf{x}) = \text{sign} \left(\sum_{k=1}^K \alpha_k C_k(\mathbf{x}) \right)$$

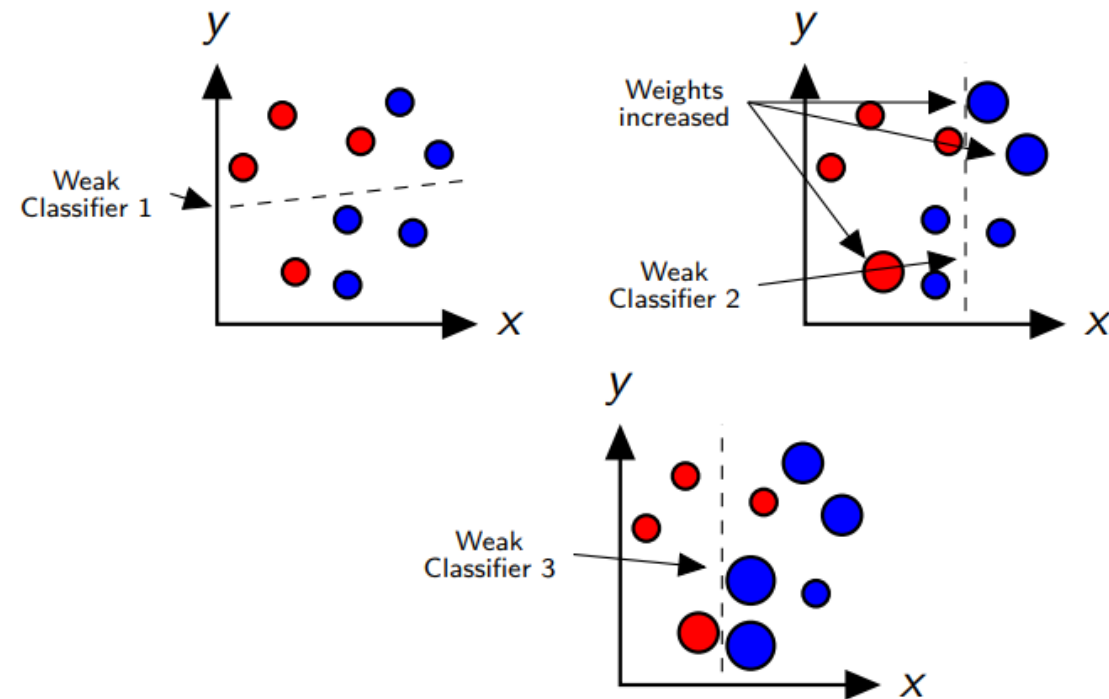
AdaBoost: Boosting Part



Consider a 2D feature space with **positive** and **negative** examples.

Each weak classifier splits the training examples with at least 50% accuracy.

Examples misclassified by a previous weak learner are given more emphasis at future rounds.



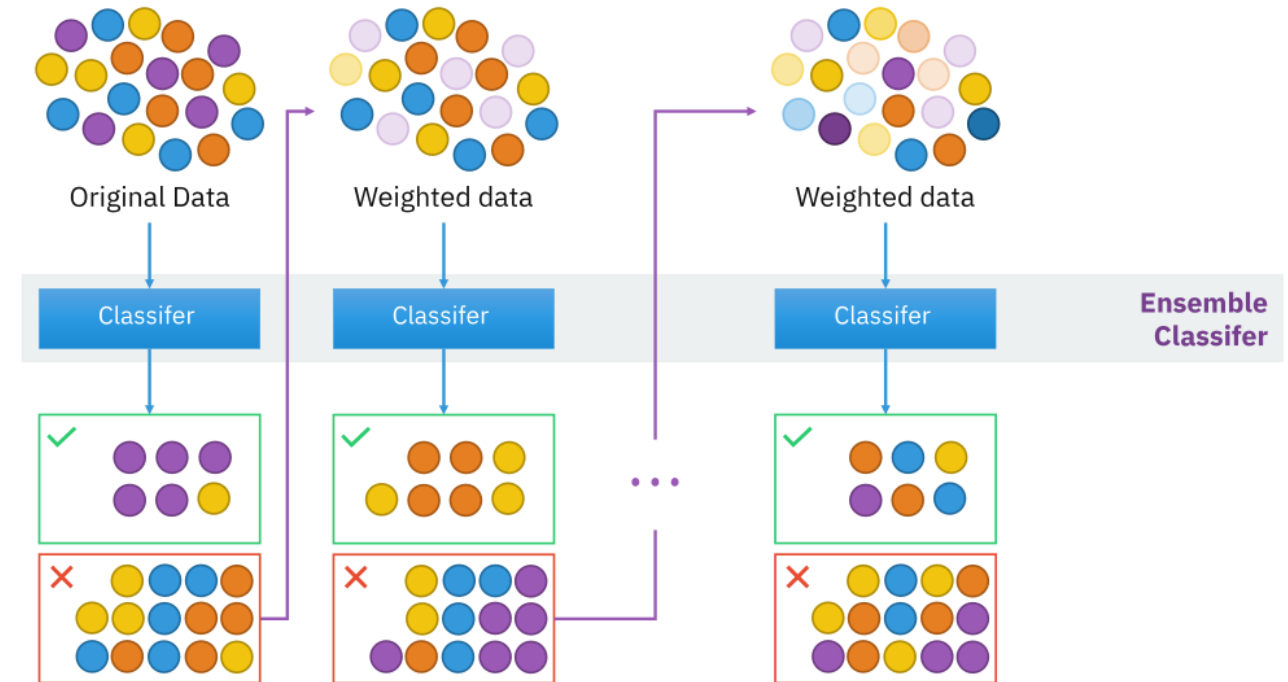
Final classifier is combination of the weak classifiers.

AdaBoost: Boosting Part



How do we emphasize wrongly classified data points?

If previously misclassified: Increase probability to be chosen as training data for upcoming weak classifier



AdaBoost: Algorithm



- ▶ Initialization: Set $w_n^{(1)} = \frac{1}{N}$ for $n = 1, \dots, N$.
- ▶ For $k = 1, \dots, k$ iterations
 - ▶ Train a new weak classifier $c_k(\mathbf{X})$ using current weights $\mathbf{W}^{(k)}$ by minimizing the weighted error function
 - ▶ estimate the weighted error of this classifier on \mathbf{X} :

$$\epsilon_k = \frac{\sum_{n=1}^N w_n^{(k)} I(c_k(\mathbf{X}) \neq y_n)}{\sum_{n=1}^N w_n^{(k)}}$$

- ▶ Calculate a weighting coefficient for $c_k(\mathbf{X})$:

$$\alpha_k = \ln \left\{ \frac{1 - \epsilon_k}{\epsilon_k} \right\}$$

- ▶ Update the weighting coefficients:

$$w_n^{(k+1)} = w_n^{(k)} \exp\{\alpha_k I(c_k(\mathbf{X}) \neq y_n)\}$$

Exercise 3

Solution



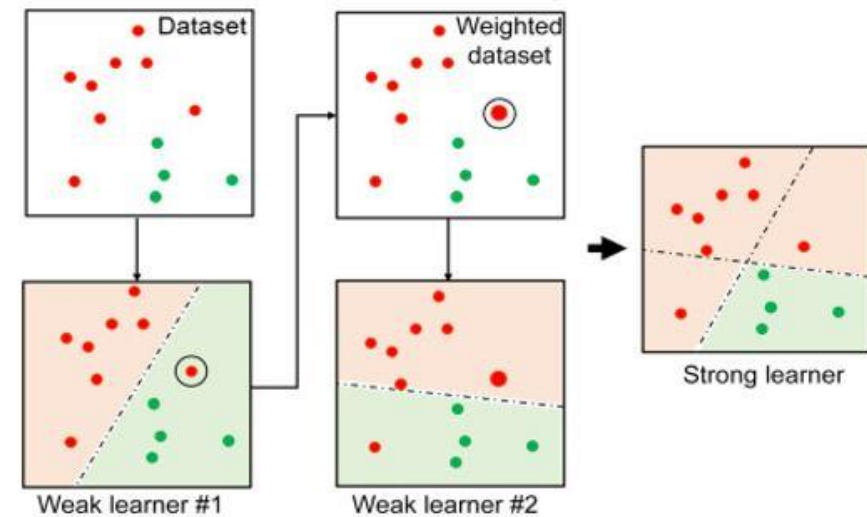
Solution - Exercise 3a



- (a) When AdaBoost is applied to a weak learner with a linear decision boundary, what is the form of the resulting decision boundary of the boosted classifier?

All single classifiers create linear boundaries

→ Final classifier will be a piecewise linear function



Solution - Exercise 3b

Code: Simple Classifier Function

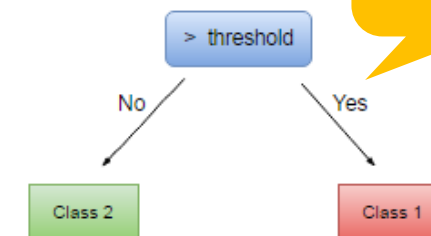


- (b) As you know from the lecture, the foundation of the Adaboost algorithm are the so-called weak classifiers. For a start we are going to use extremely simple classifiers as weak classifiers.

$$c(\mathbf{x}|j, \theta, p) := \begin{cases} 1 & \text{if } px_j > p\theta \\ -1 & \text{otherwise,} \end{cases}$$

dimension (pointing to j)

threshold value (pointing to θ)



Decision stump

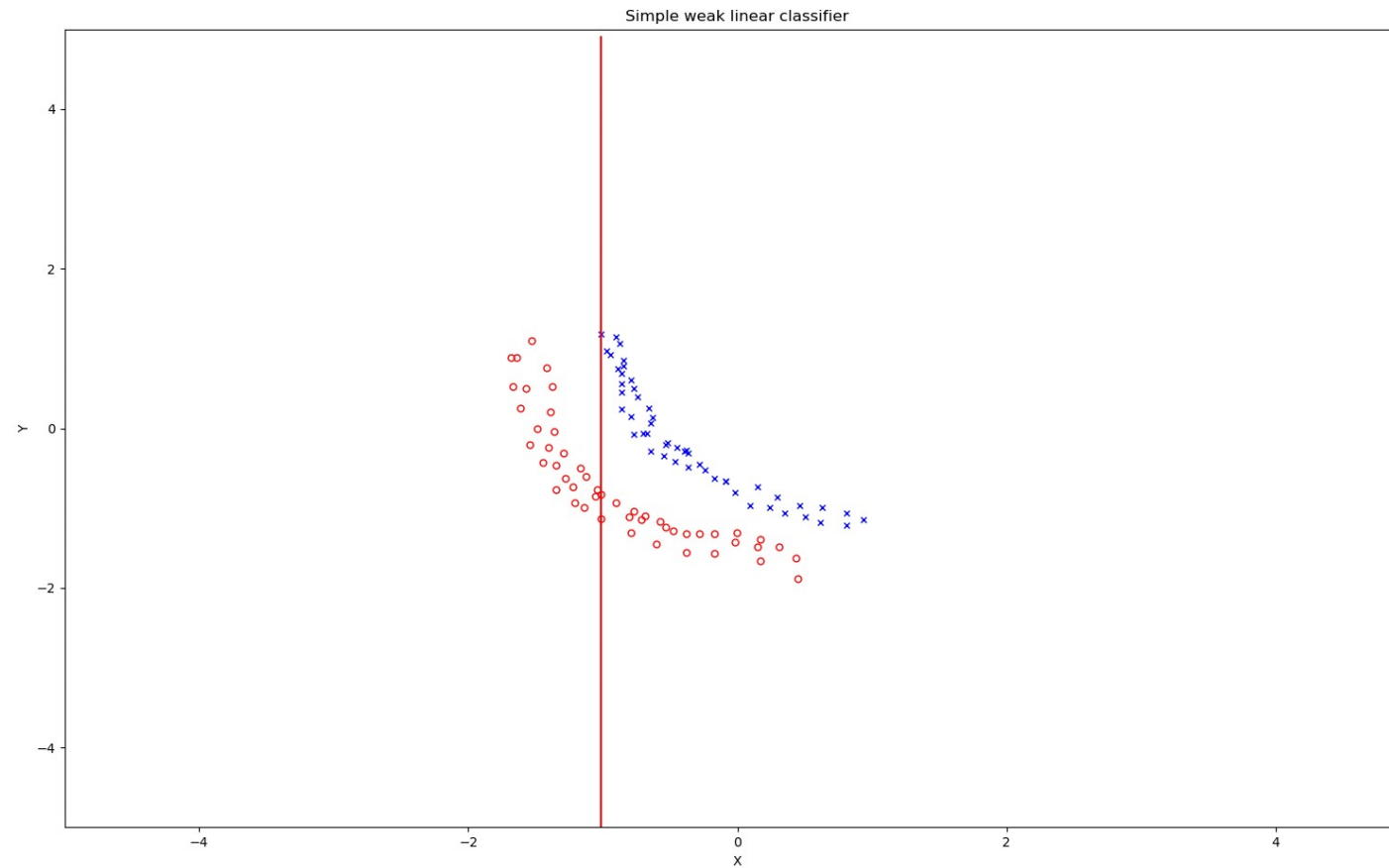
where $p \in \{-1, 1\}$ is a parity indicating the direction of the inequality sign. This classifier ignores all entries of \mathbf{x} ^{except} those in dimension j , x_j . To train this weak classifier on weighted data, we use the following learning rule:

$$(\hat{j}, \hat{\theta}) := \operatorname{argmin}_{j, \theta} \frac{\sum_{i=1}^n w_i \mathbf{I}\{y_i \neq c(\mathbf{x}|j, \theta, p)\}}{\sum_{i=1}^n w_i}$$

We count how often this is true.

For each dimension, there is an optimal decision boundary θ . We aim to minimize this function, i.e. we try to find a classifier where the number of misclassified data points is minimal.

Solution - Exercise 3b



Solution - Exercise 3c

Code: AdaBoost Simple + Eval



```

4 def simpleClassifier(X, Y):
5     # Select a simple classifier
6     #
7     # INPUT:
8     # X      : training examples (numSamples x numDim)
9     # Y      : training labels (numSamples x 1)
10    #
11    # OUTPUT:
12    # theta   : threshold value for the decision (scalar)
13    # j       : the dimension to "look at" (scalar)
14
15    #####Start Subtask 1b#####
16    N, D = X.shape
17
18    # Initialize least error
19    le = 1
20    j = 1 # dimension
21    theta = 0 # decision value
22
23    # Iterate over dimensions, which j to choose
24    for jj in range(D):
25
26        # Find interval to choose theta
27        val = X[:, jj] # shape: (100 x 1)
28
29        sVal = np.sort(val) # TODO: returns unique sorted values, shape: (100 x 1)
30        idx = np.argsort(val)
31        change = np.where(np.roll(Y[idx], -1) + Y[idx] == 0)[0] # shape: (36 x 1)
32
33        # Calculate thresholds for which we want to check classifier error.
34        # Candidates for theta are always between two points of different classes.
35
36        th = (sVal[change[change < len(X)-1]] + sVal[change[change < len(X)-1]+1])/2 # shape: (35 x 1)
37
38        error = np.zeros(len(th)) # error-placeholder for each value of threshold th
39

```

where class labels change,
we have potential θ values

potential θ values

Solution - Exercise 3c

Code: AdaBoost Simple + Eval



```

40 # Iterate over candidates for theta
41 for t in range(len(th)):
42     # Initialize temporary labels for given j and theta
43     cY = np.ones(N)*(-1) # shape: (100 x 1)
44
45     # Classify
46     cY[X[:, jj] > th[t]] = 1 # Set all values to one, which are bigger then current threshold
47
48     # Calculate error for given classifier
49     error[t] = sum([Y[i] != cY[i] for i in range(N)])
50
51     # Visualize potential threshold values
52     print('J = {0} \t Theta = {1} \t Error = {2}\n'.format(jj, th[t], error[t]))
53
54     le1 = min(error/N)
55     ind1 = np.argmin(error/N)
56     le2 = min(1-error/N)
57     ind2 = np.argmin(1-error/N)
58     le0 = min([le1, le2, le])
59
60     if le == le0:
61         continue
62     else:
63         le = le0
64         j=jj+1 # Set theta to current value of threshold
65         # Choose theta and parity for minimum error
66         if le1 == le:
67             theta = th[ind1]
68         else:
69             theta = th[ind2]
70
71 #####End Subtask#####
72 return j, theta

```

Solution - Exercise 3c

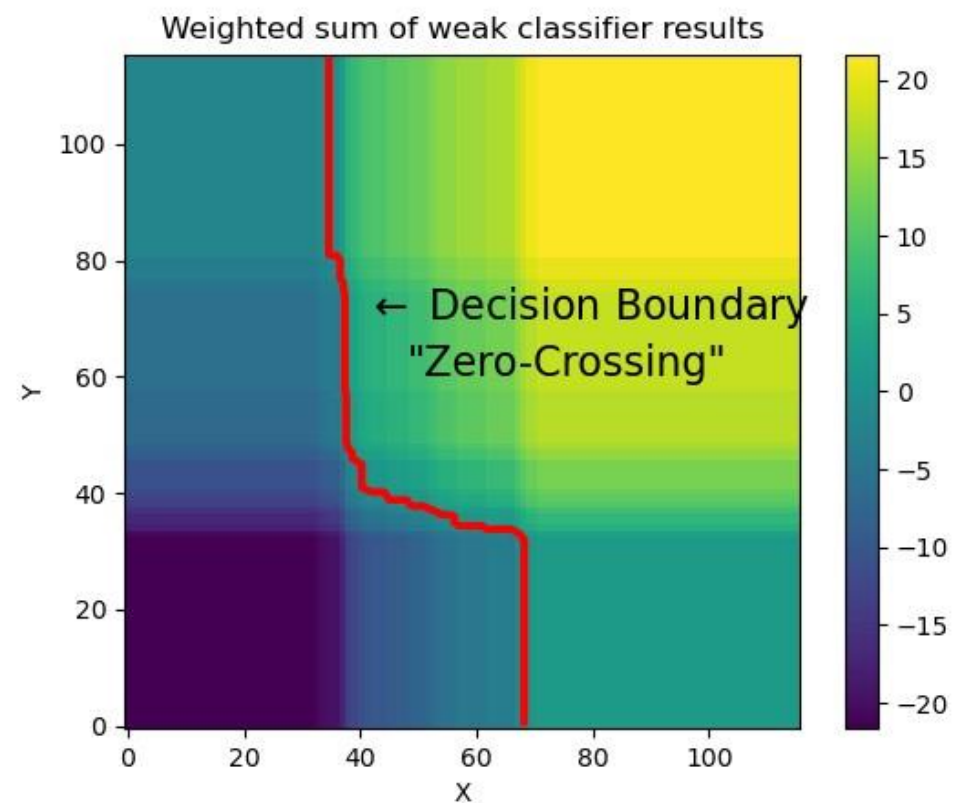
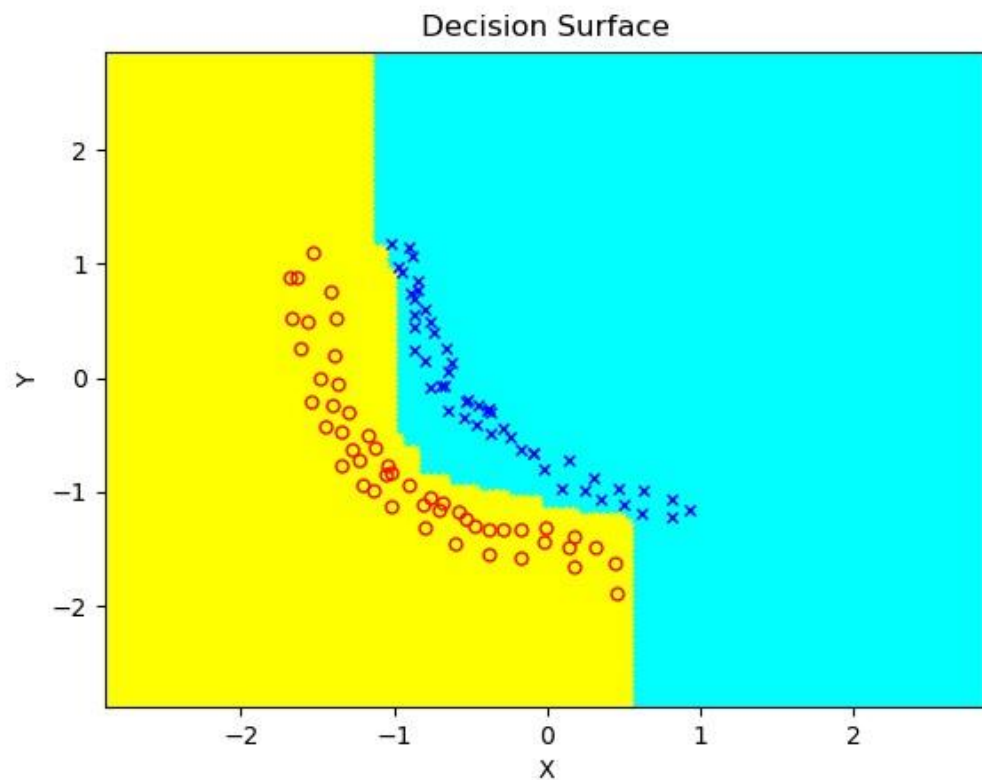
Code: AdaBoost Simple + Eval



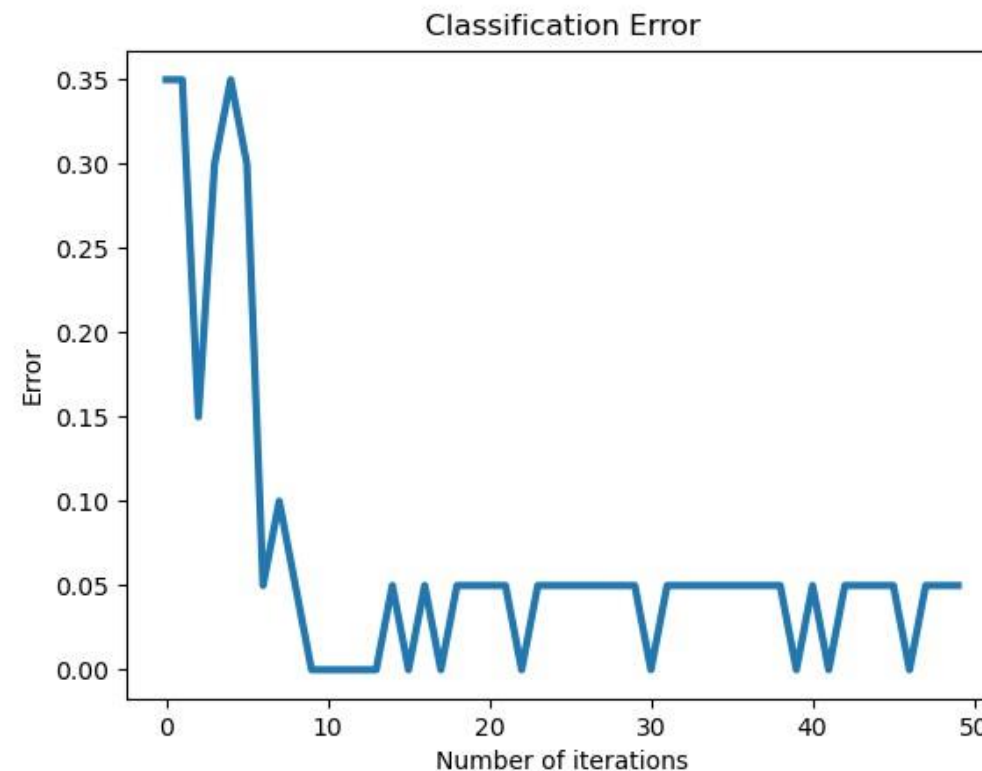
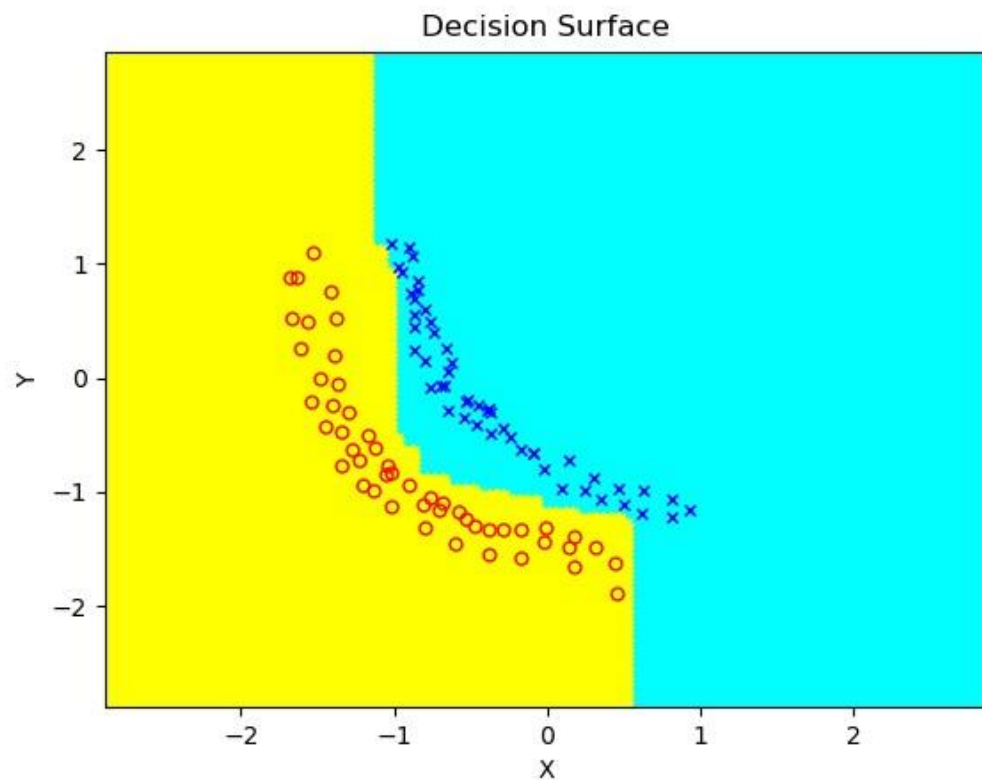
A quick recap of the AdaBoost algorithm: Given a weak classifier and training samples (x_i, y_i) , $i = 1 \dots N$ (where $x_i \in \mathbb{R}^d$ are data points in d dimensional space and $y_i \in \{-1, +1\}$ are class labels), the AdaBoost algorithm for the two-class problem is:

1. Initialize the weights $w_i^1 = \frac{1}{N}$
2. For $k = 1 \dots K$
 - (a) Train the weak classifier c_k on the weighted data i.e. select the optimal parameter combination
 - (b) Compute the error: $e_k = \frac{\sum_{i=1}^N w_i^k \mathbf{I}\{y_i \neq c_k(x_i)\}}{\sum_{i=1}^N w_i^k}$
 - (c) Compute the voting weight for the weak classifier: $\alpha_k = \frac{1}{2} \ln \left(\frac{1-e_k}{e_k} \right)$
 - (d) Recalculate the weights: $w_i^{k+1} = w_i^k \exp(-\alpha_k y_i c_k(x_i))$ and normalize w such that it is a probability distribution
3. Resulting classifier: $c_K(\mathbf{x}) = \text{sgn} \left(\sum_{k=1}^K \alpha_k c_k(\mathbf{x}) \right);$

Solution - Exercise 3c



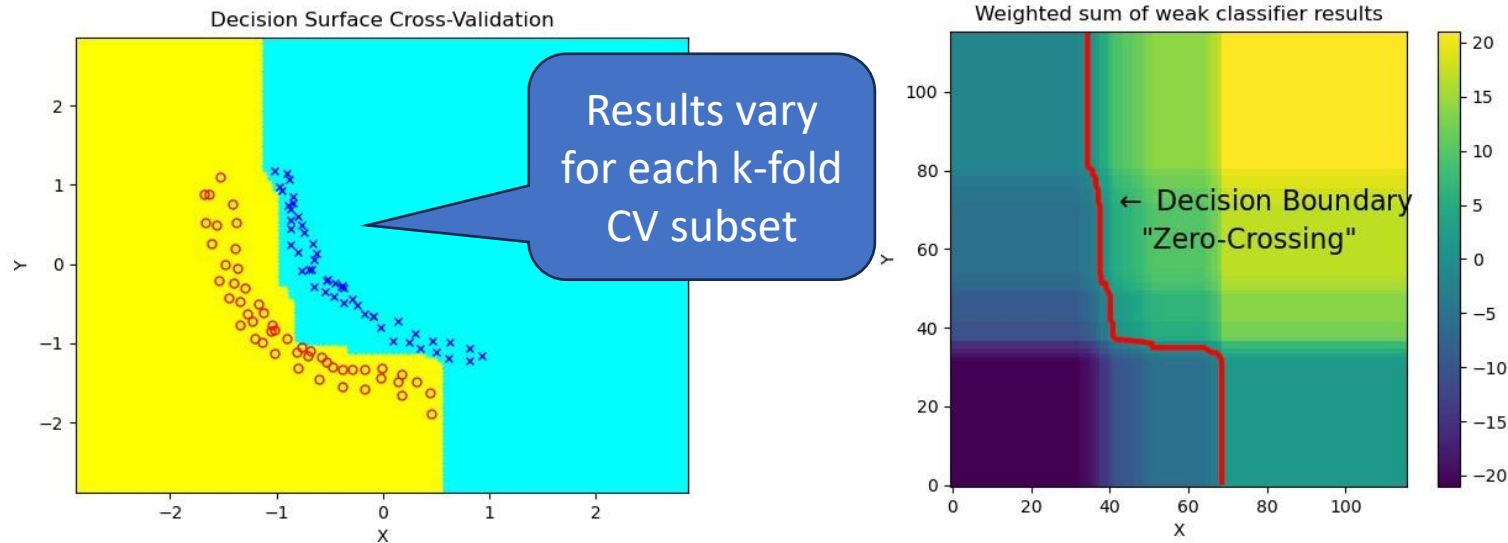
Solution - Exercise 3c



Solution - Exercise 3d

Code: AdaBoost Cross

- (d) Next, we are going to add a cross-validation step to the training procedure. This means that only a part of the available training data is used for actual training. The remaining part is used to estimate the generalization performance of the learned classifier. To this end, split the training data according to the given percentage (percent = 0.4 means 40% is used for **validation**). After each iteration k of the algorithm, estimate the classification error of the current boosting classifier (not the base classifier) by cross-validation.



Solution - Exercise 3d



Plot the cross validation error estimates vs. the number k of iteration. What do you observe?

You can see that the decision boundaries vary due to the differences in the training datasets

Why do we apply Cross Validation?

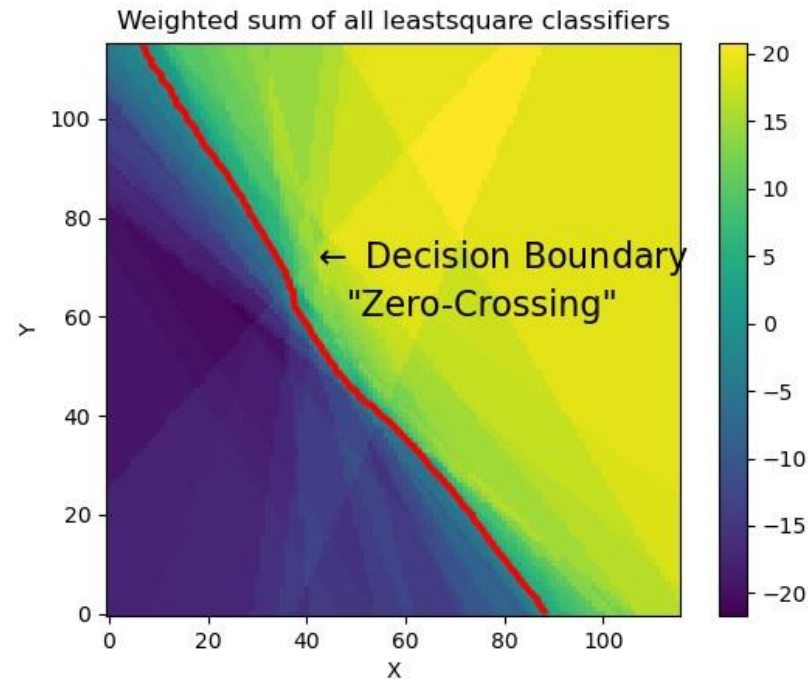
- Prevent overfitting
- Find optimal parameter set

Solution - Exercise 3e

Code: AdaBoost LSLC + Eval



- (e) Now, we are going to use a more complex weak classifier in the AdaBoost framework, namely the least squares classifier (already implemented in `leastSquare`). Implement the AdaBoost algorithm in



Solution - Exercise 3f

Code: AdaBoost USPS



- (f) Now, you are going to use the least-squares based AdaBoost on real data, i.e. the USPS data (provided in `usps.mat`). The dataset consists of a matrix `X` and a label vector `Y`. Each row of the matrix `X` is an image of size 20×14 and can be viewed with `matplotlib.pyplot.imshow(X[0, :].reshape(20, 14, order='F').copy())`. The first 5000 rows of `X` contain the images of the digit 2, and the rest contains the images of the digit 9. Perform a random split of the 10000 data points into two equally sized subsets, one for training and one for validation. Run this at least three times and plot the cross validation error estimates vs. the number k of iterations.

"2" vs. "9"

