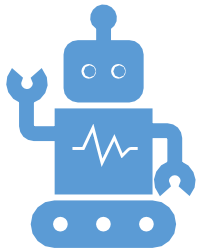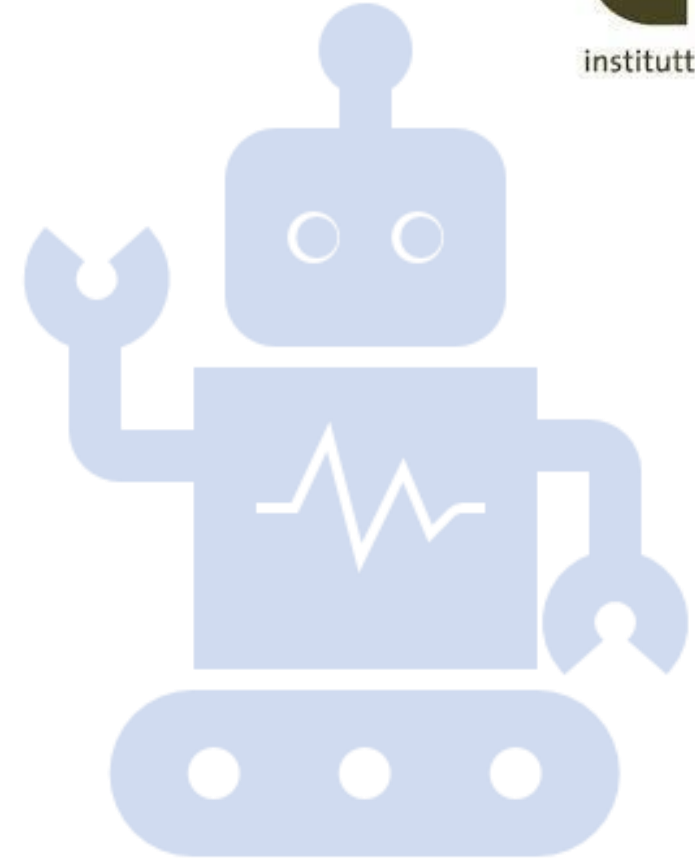# IN3050/IN4050 - Introduction to Artificial Intelligence and Machine Learning

Lecture 9

Topics in supervised learning

Jan Tore Lønning

overview

# Topics in supervised machine learning

- This presentation
- Videos
- Some jupyter notebooks
- scikit-learn

# Topics in supervised machine learning

1. This intro
2. ML in scikit-learn:
   - video
   - notebook: sklearn_assignment_2
3. Overfitting and regularization -
   - video
4. Regularization in scikit-learn
   - video
   - notebook: boston_regularization_cv

4. Bias-variance tradeoff
   - video
5. Cross-validation
   - video
   - notebook: boston_regularization_cv
6. Ensemble learning and random forest
   - video

- The  slides and videos was mostly made last year (2020)
- Included this slide to match the page numbering in the videos ☺

# scikit-learn

Jan Tore Lønning

# Why use toolkits

## Toolkits like scikit-learn

- More efficient, faster
- Tested code (error-free)
- Comprehensive
  - many different learners
- Flexible
  - many options
- Consistent interface
- Well-integrated with Python, NumPy, etc.

## Own implementation

Help us to

- Understand the inner workings of the algorithms
- Make informed choices regarding:
  - Learner
  - Parameter settings

# Other toolkits

## Why scikit?

- Why not jump directly to deep learning toolkits?
- Many problems can be solved with simpler learners
  - Don't use a jumbo jet when you are going to the store
- A base-line for more advanced learners
- Easy to use

## Toolkits for deep learning

- Tensor Flow
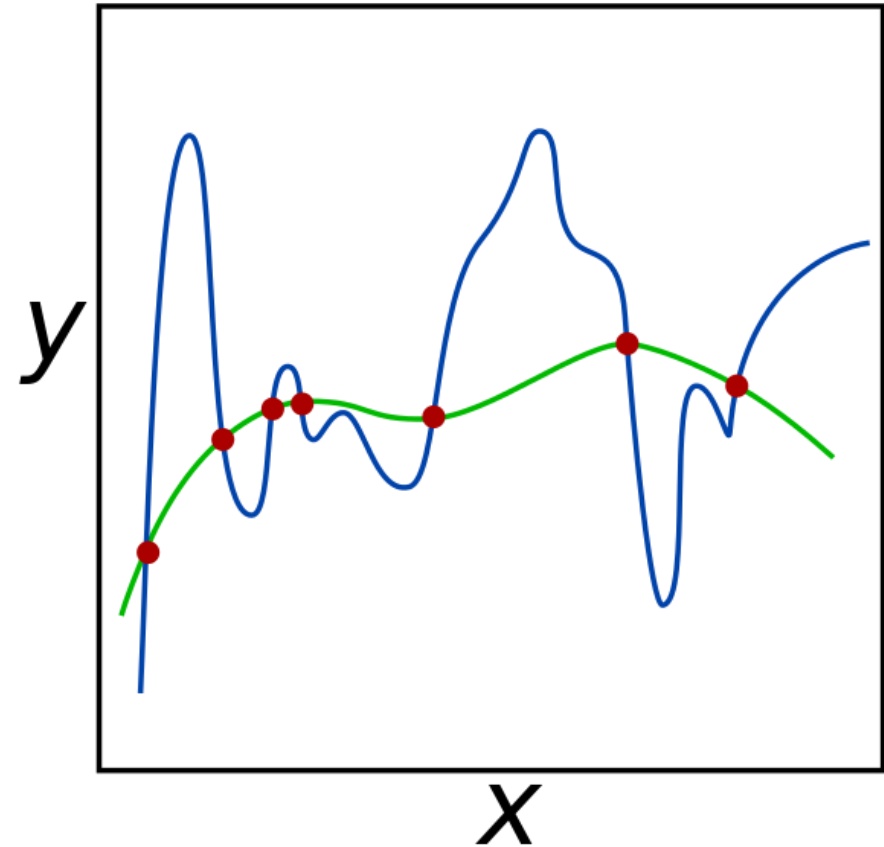  - with Keras
- PyTorch
- and more

# Overfitting and regularization

Jan Tore Lønning
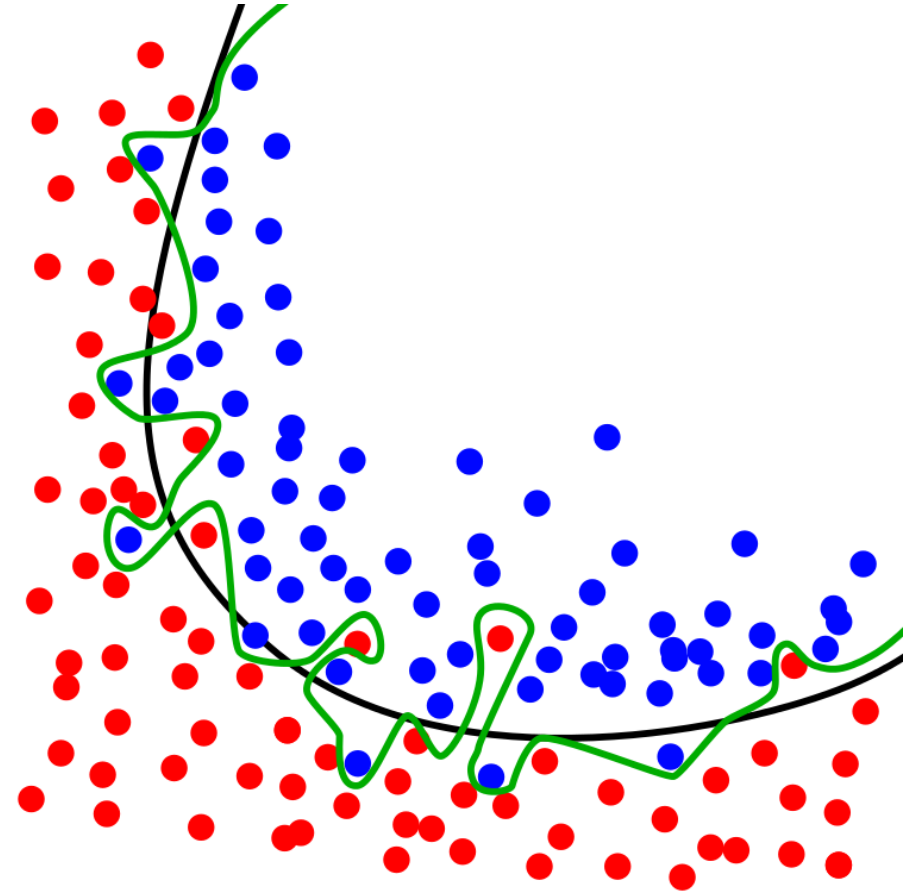
# Overfitting

- Goal in supervised ML:
    - Construct models that fit the training set
- Sometimes, we succeed too well:
    - Models fit training data very well
    - Does not generalize to other data



Source: Wikipedia

# Overfitting

- Classification

# How come?

- Whenever we choose a random sample of individuals from a larger population,
  there might be attributes that have a different distribution in the sample than in the population at large
- The learning algorithm may focus on these attributes

- Example:
- Pick 100 Swedes and 100 Norwegians at random
- We might find properties shared by the 100 Swedes and not the 100 Norwegians which are not representative for Swedes and Norwegians in general
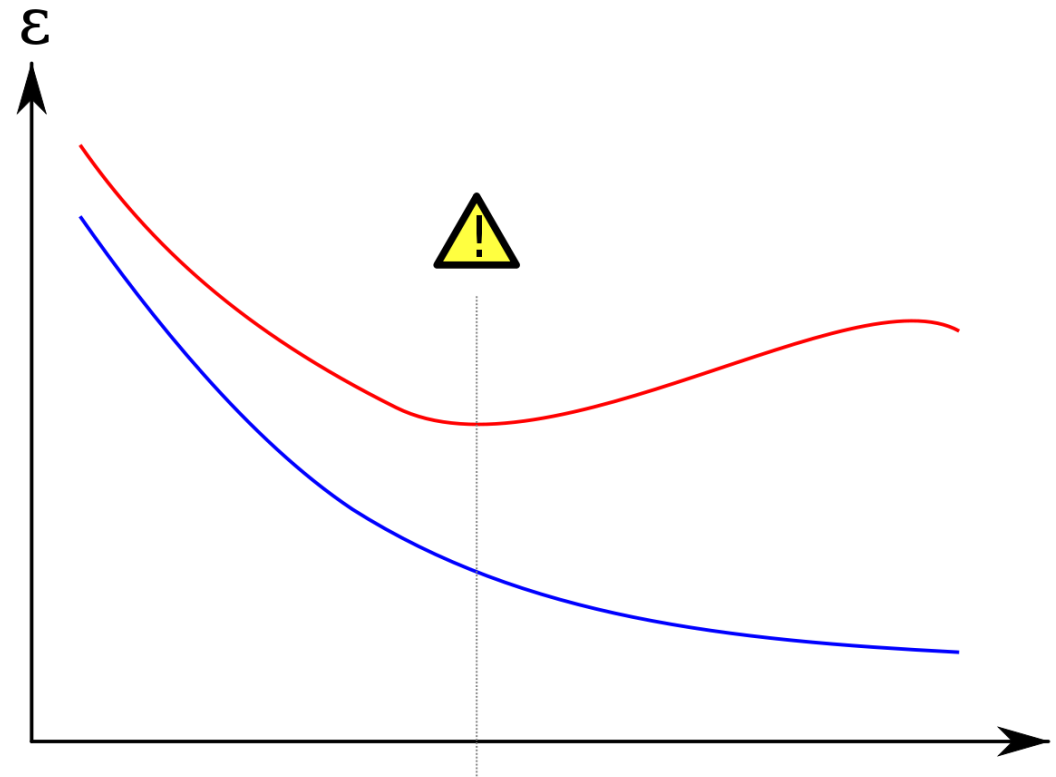
# Gradient descent-based algorithms

Classification:

- We replace the true objective:
  - accuracy, recall, precision, etc.
- with a loss function
- The loss function may continue to improve long after we have reached the best accuracy
- It moves weights to the best discriminators
- They are not necessarily representative

- Typically, overfitting with
  - many features
  - relative to the training data
- Advice:
  - At least ten times more training datapoints than features
- Not always possible
  - E.g., NLP:
    - Millions of features
    - A few thousand datapoints.

# Remedies

Alt. 1: Early stopping

Alt. 2: Regularization

$\varepsilon$

# Regularization

## So far

- Training set $(X, t)$, where
  - $X$ is the data (inputs)
  - $t$ the corresponding target values
- A model with weights $w$, s.t. it
  - from $X$ and $w$ predicts values $y$.
- A loss function which tells how well $y$ approximates $t$.
- Learning objective:
  Find the weights that minimize the loss:
- $\widehat{w} = \underset{w}{\operatorname{argmax}}(-Loss(X, t, w))$

## Regularization

- Replace
  - $-Loss(X, t, w)$ with
  - $-Loss(X, t, w) - \alpha R(w)$
- Where
  - $R(w)$ punishes large weights
  - $\alpha$ determines how hard they should be punished
  - $\alpha$ is to be tuned on a validation set

# Alternatives for regularization

| Ridge regularization | Other alternatives |
|---|---|

**Ridge regularization**

- The most common uses L2-distance
- $R(\boldsymbol{w}) = \boldsymbol{w} \cdot \boldsymbol{w} = \sum_{i=0}^{m} w_i^2$
- Example:
  - $\boldsymbol{w}_1 = (0,0,0,1), R(\boldsymbol{w}_1) = 1$
  - $\boldsymbol{w}_2 = (0.1, 0.2, 0.3, 0.4),$
    $R(\boldsymbol{w}_2) = 0.3$
  - Preference for $\boldsymbol{w}_2$

**Other alternatives**

- Lasso regularization:
  - Uses L1
  - $R(\boldsymbol{w}) = \sum_{i=0}^{m} |w_i|$

- Elastic nets:
  - uses both L2 and L1

# To the notebook

- In scikit-learn, you find regularization in
  - Ridge, for linear regression
  - RidgeClassifier, for linear regression classifier
  - LogisticRegression
  - and more

# Bias-variance tradeoff

# Meanings of "bias"

1. Bias term
   - Not today's topic

2. About a classifier:
   - *Unfairness, prejudice*
   - For ML-systems this might be because of the selection of
     - training examples, or
     - features

3. A statistics is biased
   - (this is what can be traded)

- Examples (of 2)
- Early speech recognition system for cars:
  - Performed much better on males than females
  - Reason: Mostly trained on men
- NLP, *word embeddings:*
  - *Man is to Computer Programmer as Woman is to Homemaker*

# Bias (meaning 2) contd.

- http://vectors.nlpl.eu/explore/embeddings/en/



https://www.shanelynn.ie/get-busy-with-word-embeddings-introduction/

# Bias-variance tradeoff

- Underlying idea
  - (from statistics, can be formalized)
- The training data $D = (X, t)$ can be described by an expression
- $t = f(x) + \varepsilon$, where
  - $\varepsilon$ is irreducible (noise) centered at 0
- The goal of the ML training is to find (an approximation) $g$ to $f$, where $(g(x) - f(x))^2$ is minimal for $x$ both from $D$, and from outside of D.

# Bias-variance tradeoff



- In general, we will not succeed in finding $g$ identical to $f$
- The expectation of the error can be split into
  - bias
  - variance

- Bias (meaning 3)
  - the classifiers systematically misses the target,
  - e.g. searching for linear classifiers to a non-linear problem
  - Underfitting

- Variance
  - Picking up some of the noise, being to sensible to small variation in the input data
  - cf. the green dec. boundary
  - Overfitting

# The tradeoff

- Traditional wisdom:
  - Lower bias gives higher variance,
  - and the other way around
- There is a sweet spot in the middle



http://scott.fortmann-roe.com/docs/BiasVariance.html

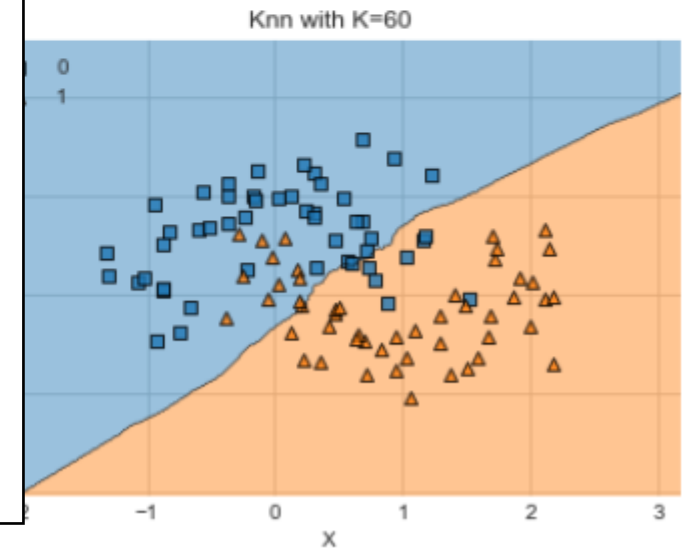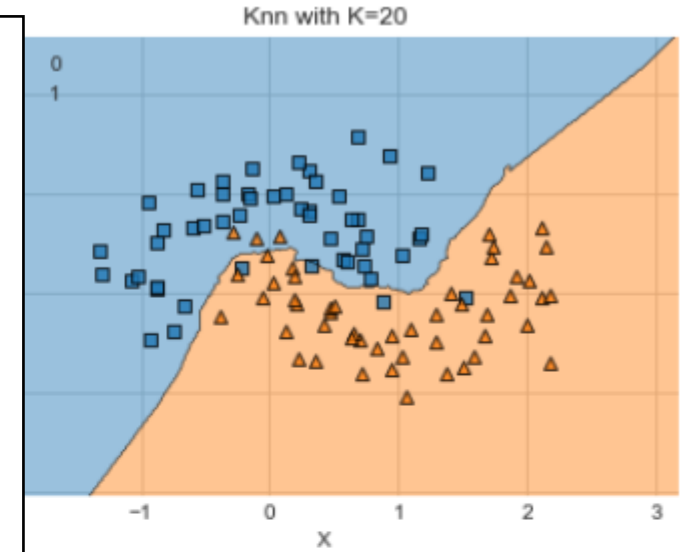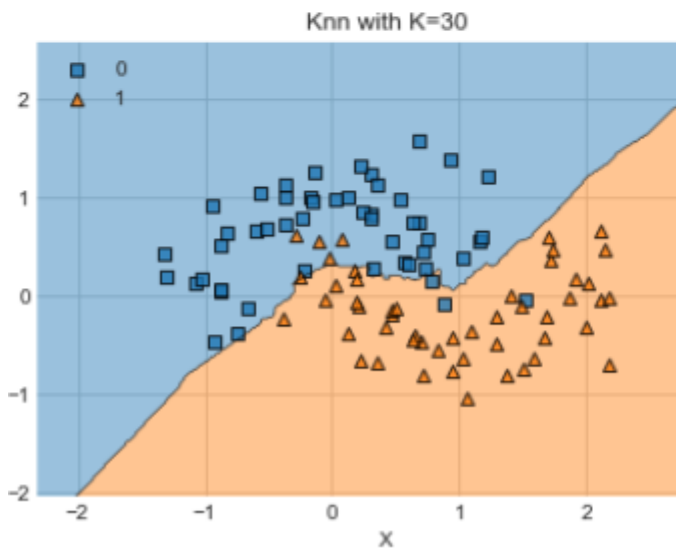# Metaphor of bias vs variance

1. *k*NN:
   - small *k*: high variance, low bias
   - large *k*: low variance, high bias

2. Polynomial:
   - small *n*, e.g., *n*=1 (linear):
     - high bias, low variance
   - large n:
     - low bias, high variance
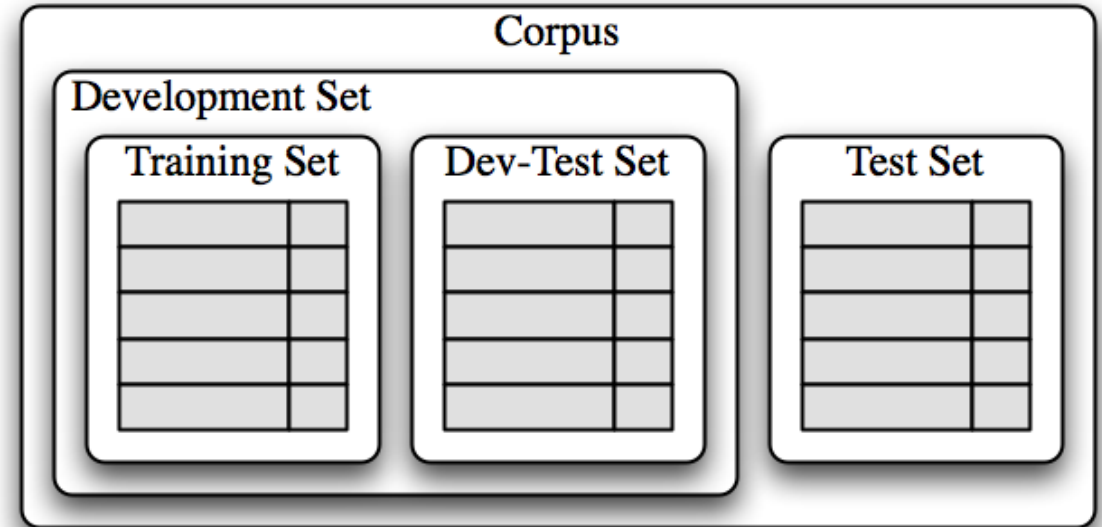
3. Regularization:
   - lower variance, increases bias
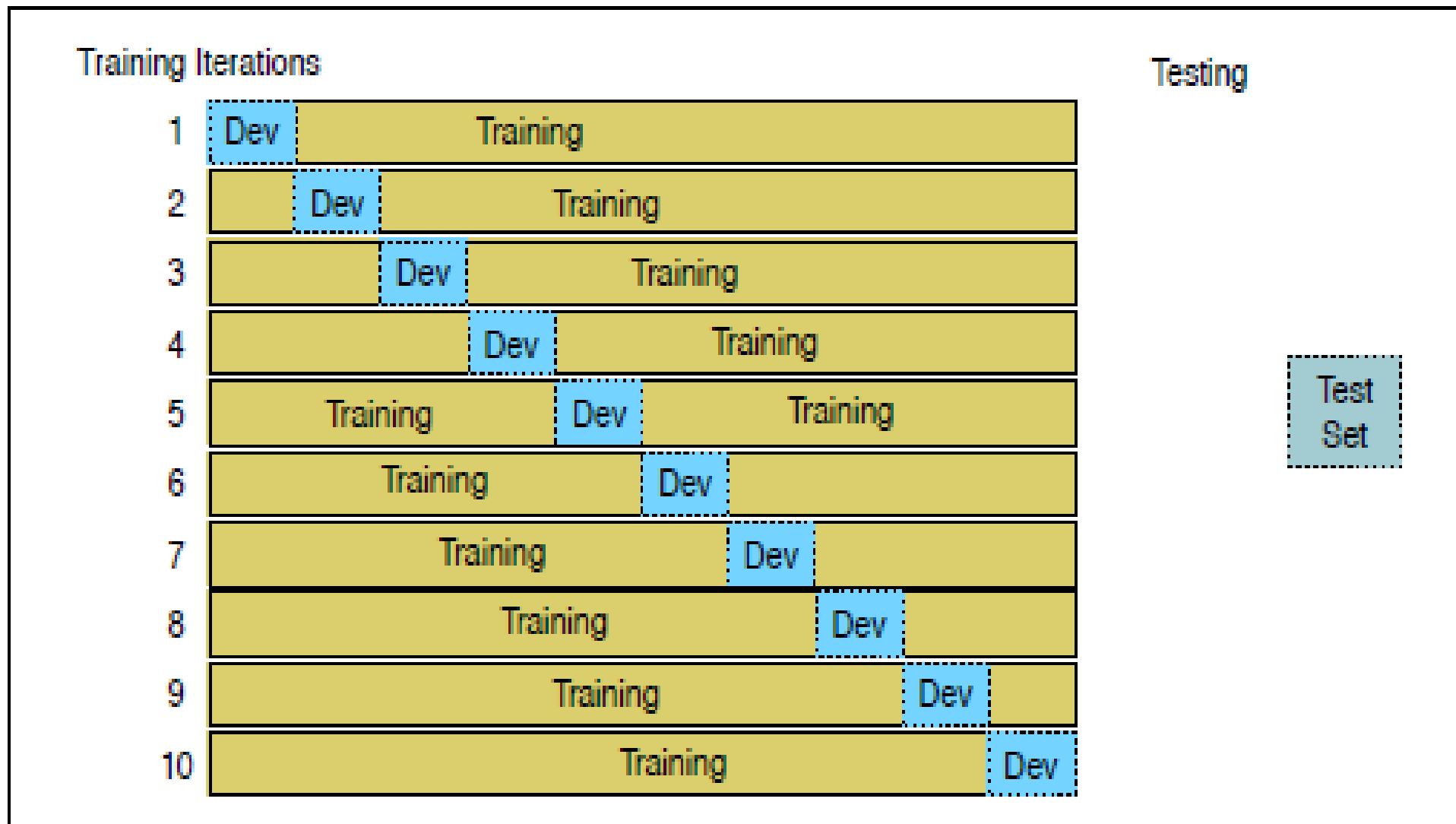
KNN visualization for the U-shaped dataset

https://towardsdatascience.com/knn-visualization-in-just-13-lines-of-code-32820d72c6b6    25

# Cross-validation

# Training and test sets (lecture 5)

- To measure improvement, we need (at least) two disjoint labeled sets:
  - Training set
  - Test set
- Train on the training set.
- Predict labels on the test set (after removing the labels)
- Compare the prediction to the given labels



https://www.nltk.org/book/ch06.html

- For repeated development we need (at least) two test sets.
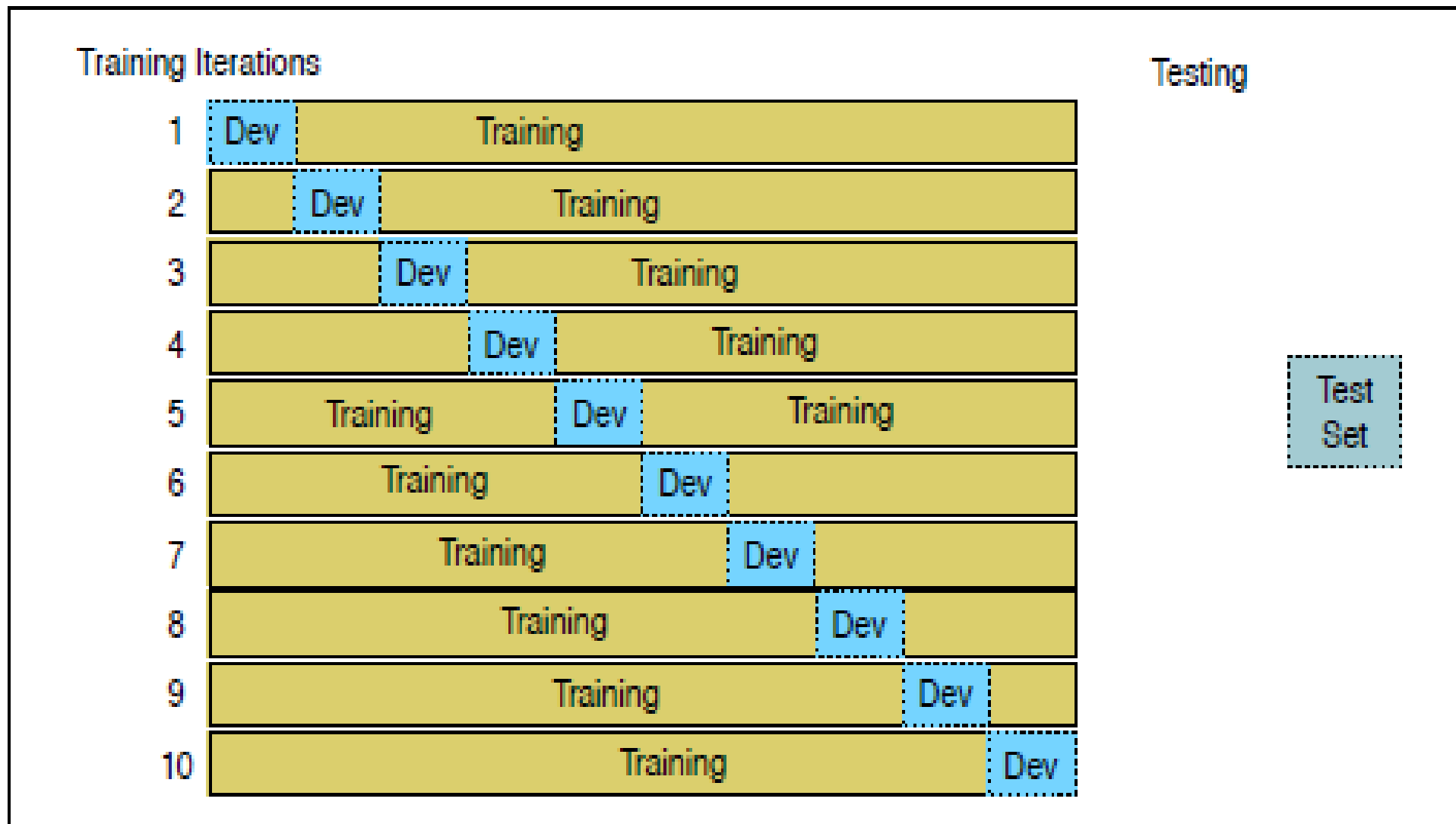  - One for repeated testing during development
  - One for final testing

**Figure 6.7** 10-fold crossvalidation

https://web.stanford.edu/~jurafsky/slp3/

# Cross-validation

- Small test sets ➔ Large variation in results
- N-fold cross-validation:
  - Split-off a final test set
  - Split the development set into $n$ equally sized bins (e.g. $n = 10$)
  - Conduct $n$ many experiments:
    - In experiment $m$, use part $m$ as test set and the $n$-1 other parts as training set.
  - This yields $n$ many results:
    - We can consider the mean of the results
    - We can consider the variation between the results.
      - Statistics!

**Figure 6.7** 10-fold crossvalidation

https://web.stanford.edu/~jurafsky/slp3/

# Marsland Ch. 2

- First edition 2009

- Second edition 2015:
  - Removed mistakes
  - Added chapter 2
    - With mistakes

- We use the same final test set.
  That is not changed with the folds.

- We might use several val-sets, say
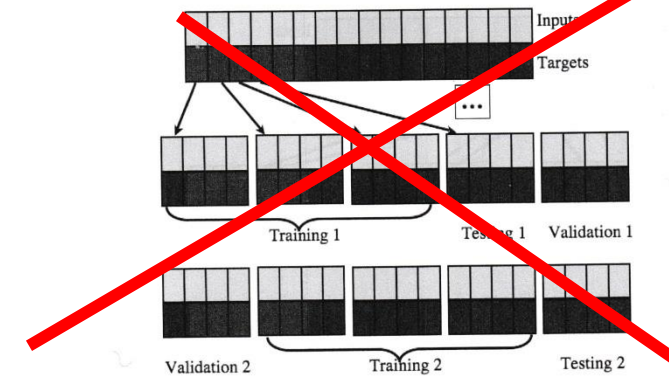  - One for testing and
  - one for tuning

# Precision and recall
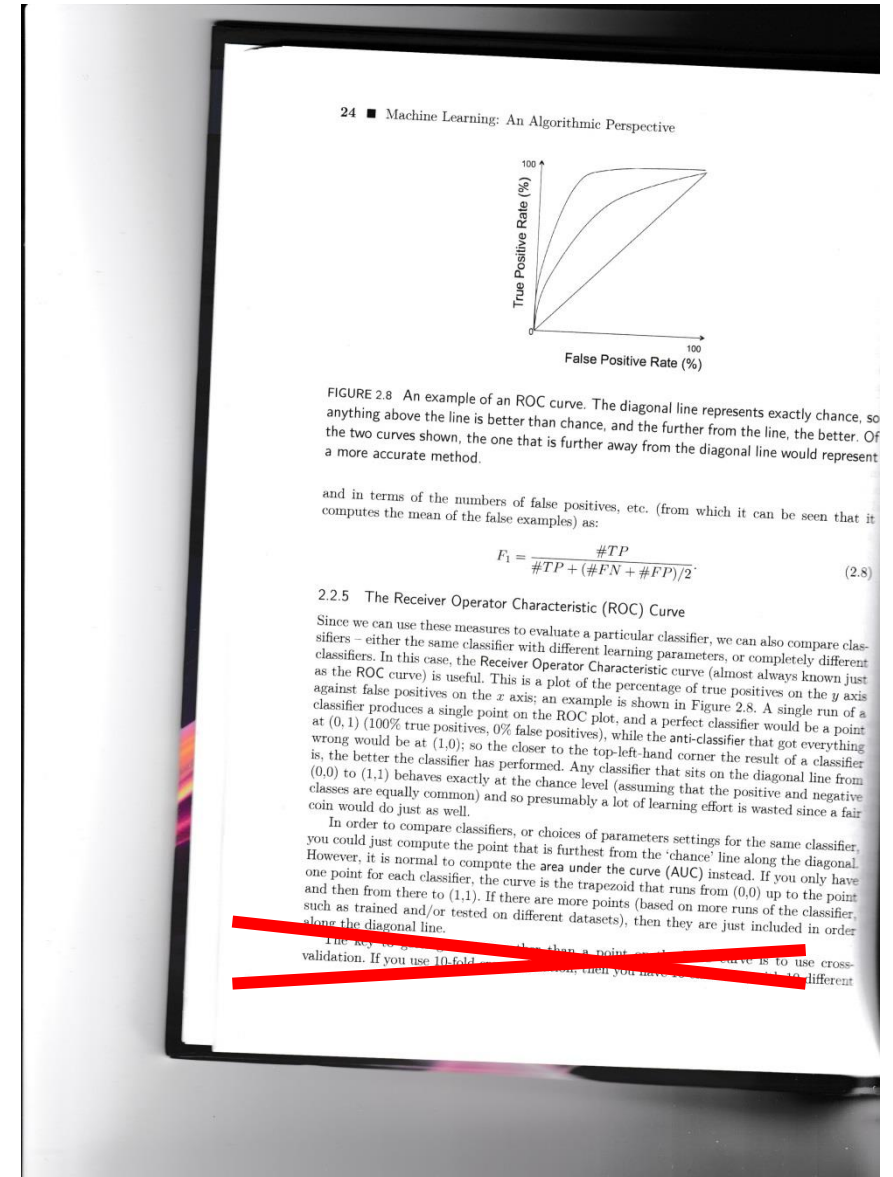
- The formulas for Precision and Recall are correct
- The explanation mixes things up

# 2.2.5 ROC curves

- We will not consider this section

- But in case you are interested in
  - Precision-recall curves, or
  - ROC-curves

- Look somewhere else.

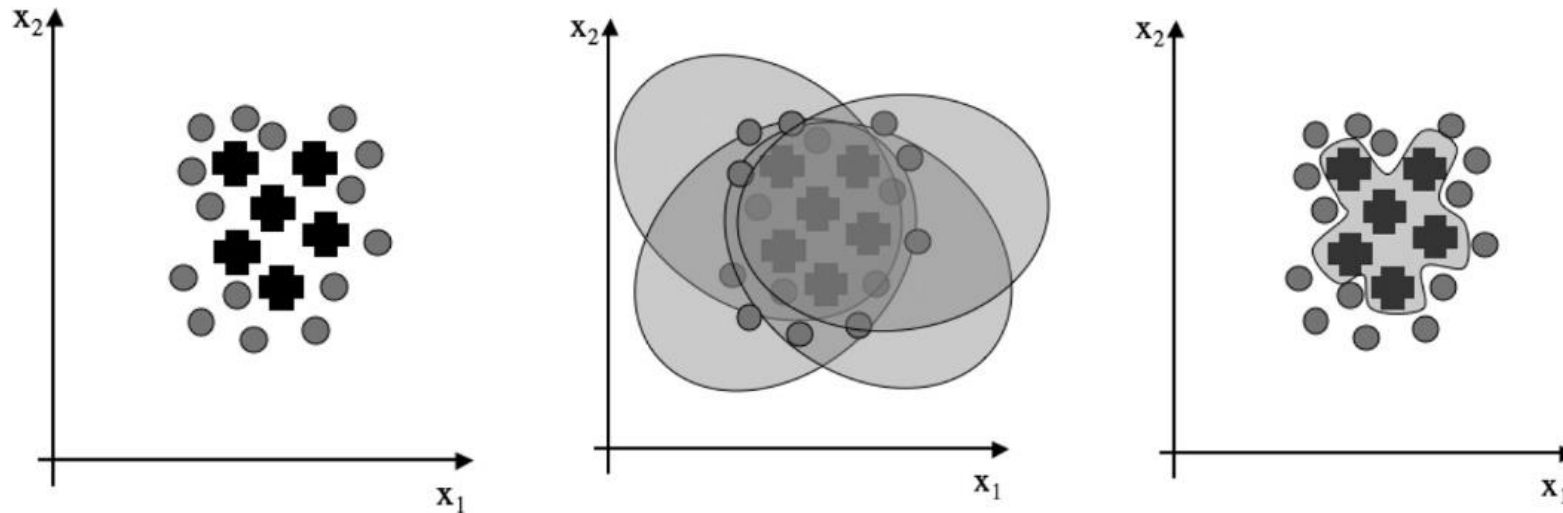- You do not use cross-validation for constructing them

# Example

- Notebook: boston_regularization_cv

# Ensemble Learning, Random Forest

# Voting classifiers



- Train several different classifiers
  - E.g. A LogReg
  - One or more *k*NNs
  - An SVM, etc.

- For prediction:
  - Run all the classifiers
  - Pick the majority vote (Hard vote)
  - Or the average if they provide probabilities (Soft vote)

# Observation

- The ensemble classifier may perform better than any of the individual classifiers.

- Even weak classifiers (accuracies just above 0.5) may perform well together, provided
  - They are independent
  - There are sufficiently many of them
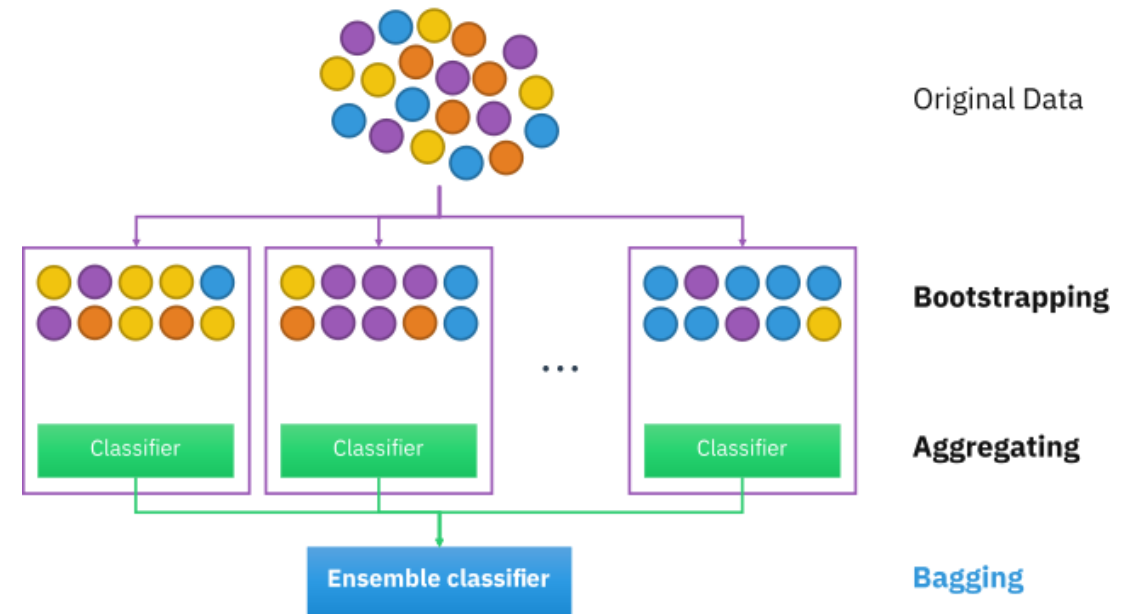
# Bagging (Bootstrap Aggregating)

- An alternative to different learners:

- Train the same learner on slightly different data

- With much training data, you may train on subsets of the data

- One way to produce different datasets is called *bootstrap.*

Bootstrap:

- Given a set of training data $D = (\boldsymbol{X}, \boldsymbol{t})$ of size $n$.

- Produce a new set by picking $n$ samples from $D$ with replacement.

- Produce several such sets, at least 50.

- Fit a model to each set.

- Prediction: use hard or soft vote.

- That's it!

# More on bagging

- Subagging:
  - Similar, but pick sets of smaller size than *D* (all of same size)
- Pasting:
  - Similar but sampling without replacement.

==========================

- By the way:
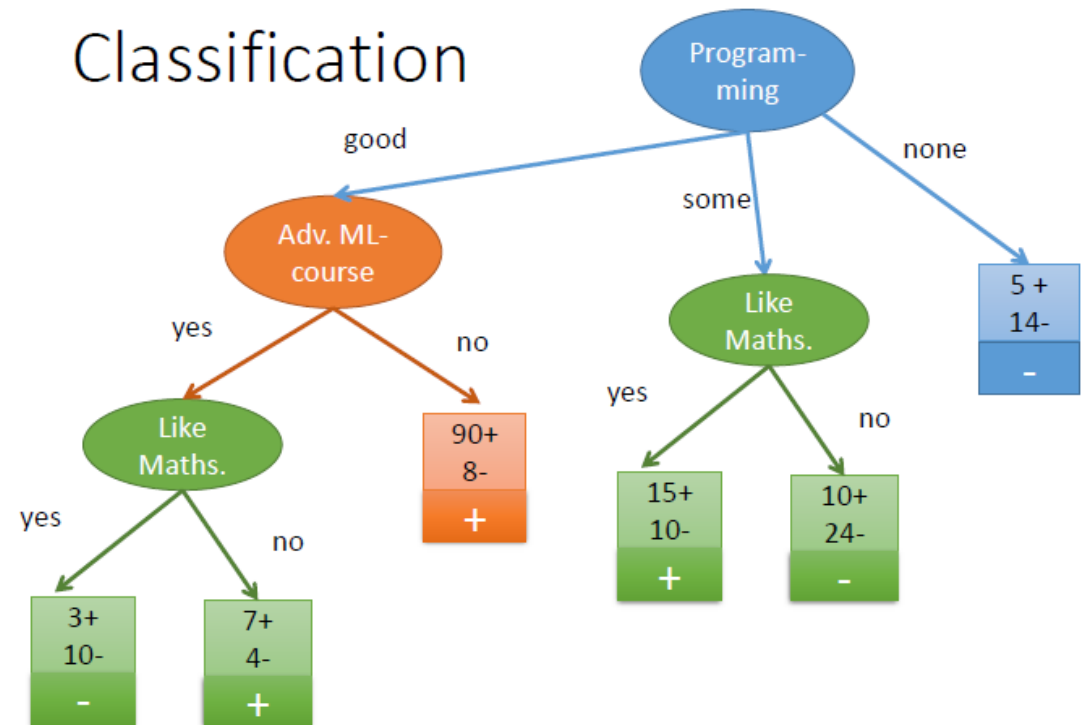  - Bootstrap may also be useful in evaluation.



https://en.wikipedia.org/wiki/Bootstrap_aggregating

# Back to decision trees

- Remember? (lecture 1)

- We used simplified criterion for selecting the stumps
  - Does not matter for the rest of the construction

==================

- Bias – variance:

- Stop the trees from growing:
  - More bias
  - Less variance

# Random forest

## Randomness, step 1

- Given training data $D = (\boldsymbol{X}, \boldsymbol{t})$:
- Use Bootstrap to construct sets of training data

========================

- Train a decision tree on each bootstrap sample

## Randomness, step 2

- At each step in the construction of the tree, consider only a subset of the available features
- In the algorithm:
  - a parameter *m*
  - *= the number of features to consider at each step*

# Random forest - properties

- Good results both on big and small datasets
- Embarrassingly parallel

# The 5+ lectures on supervised learning

## Understanding

Hopefully

- You understand basic algorithms including neural networks, s.t.
  - You can choose between them
  - You can choose hyper-parameter values and tune them

- You have got a foundation for further studies in deep learning

## Practical skills

- You can set up and run supervised ML experiments

- You can evaluate them