

Strengths and weaknesses of different machine learning models

APSC 8280: Machine learning applied to plant science

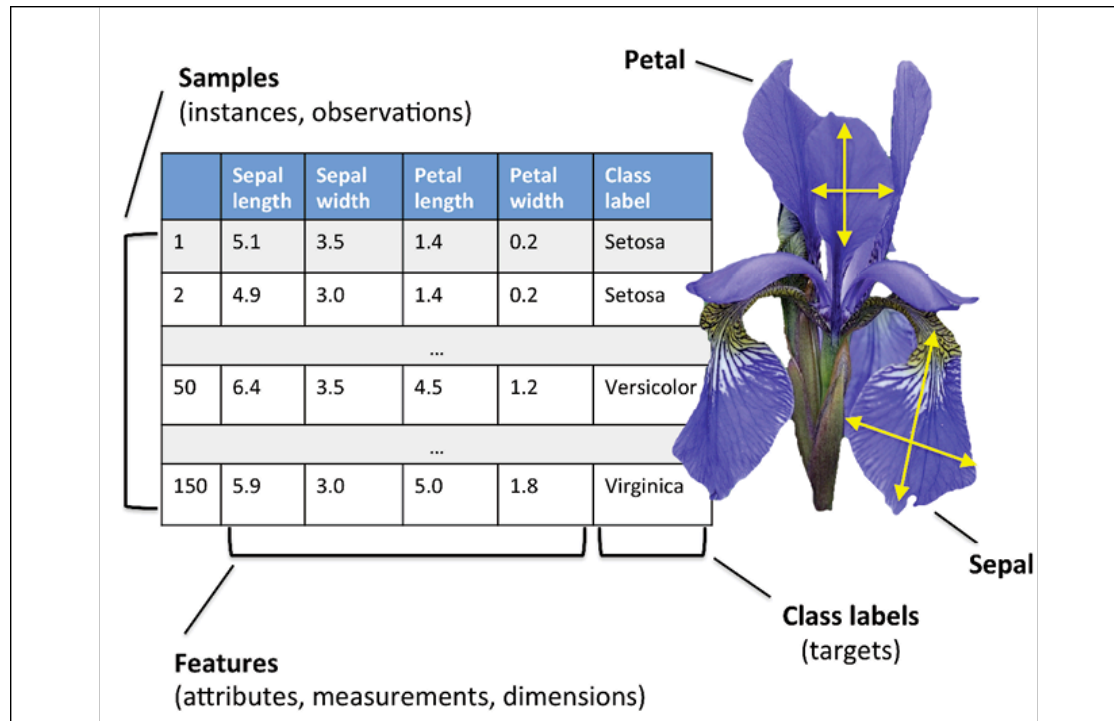
Outline

- Review lab 1
- Demo 1: End-to-end machine learning
- Quiz on machine learning fundamentals
- Algorithms
 - Naïve Bayes
 - Regression techniques
 - Support vector machines
 - Nearest neighbors
 - Decision trees
 - Ensemble methods

Demo: End-to-End Machine Learning

Datasets

Classification: flower identification



Regression: House prices prediction



Nearest neighbors

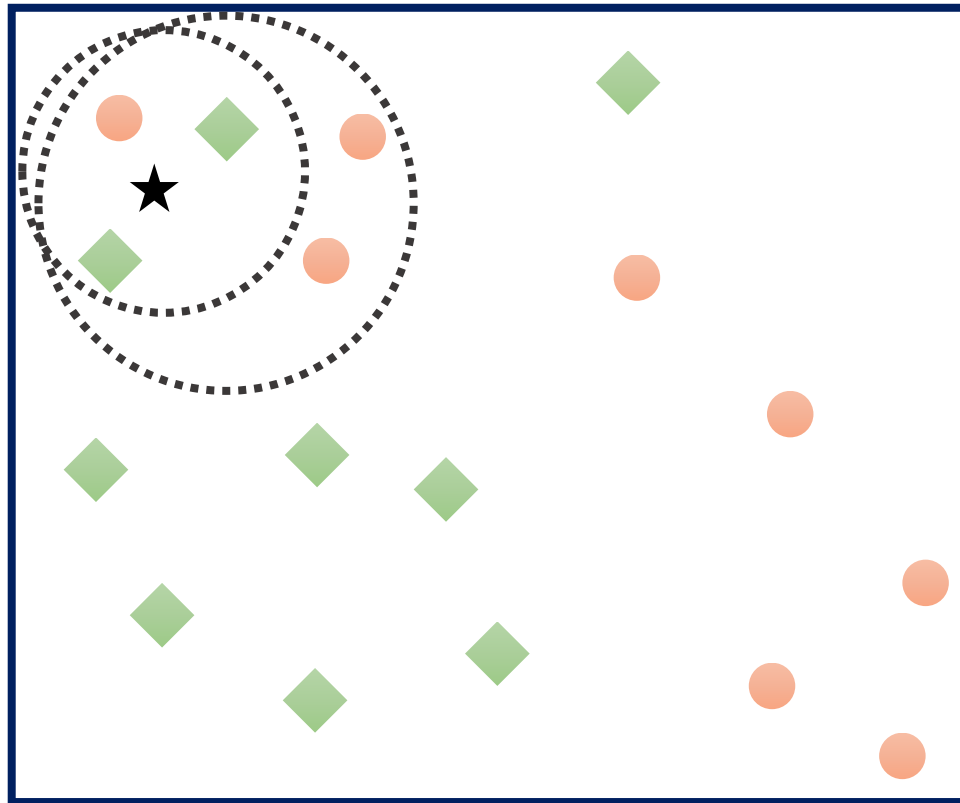
Plant disease prediction

◆ Healthy

● Diseased

$k = 3$

$k = 5$



- Classifies unlabeled examples by assigning them the class of similar labeled examples.
- Best for data with numerous complex relationships among features and target variables
- Rule of thumb: start with k equal to the square root of the number of training examples
- Large number of training examples makes the choice of k less important
- Heavily affected by the scale of the data

Nearest neighbors

Strengths

- ✓ Simple and effective
- ✓ Easy to train
- ✓ Makes no assumptions

Weaknesses

- × Does not produce a model (instance-based)
- × Hard to test
- × Struggles nominal features and missing data

Naïve Bayes

- Naïve Bayesian classifiers compute probability of each class using feature values from training data
- Best applied to problems in which information from many features should be considered simultaneously in order to estimate the probability of an outcome.
- Assumes all features are equally important and independent

$$P(\text{Healthy}|\text{data}) = \frac{P(\text{data}|\text{Healthy})P(\text{Healthy})}{P(\text{data})}$$

Diagram illustrating the Naïve Bayes formula with labels:

- Posterior probability** points to $P(\text{Healthy}|\text{data})$
- likelihood** points to $P(\text{data}|\text{Healthy})$
- Prior probability** points to $P(\text{Healthy})$
- Marginal likelihood** points to $P(\text{data})$

Naïve Bayes

Strengths

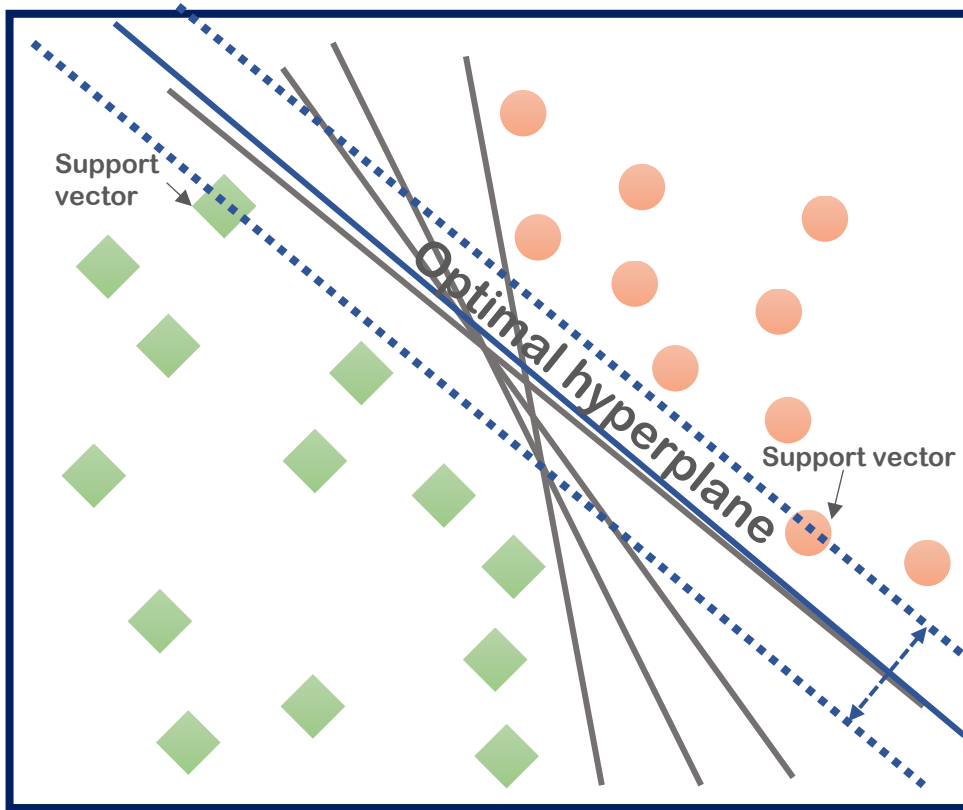
- ✓ Simple, fast and effective
- ✓ Handles noisy and missing data well
- ✓ Requires relatively few examples for training
- ✓ Easy to obtain the estimated probability for a prediction

Weaknesses

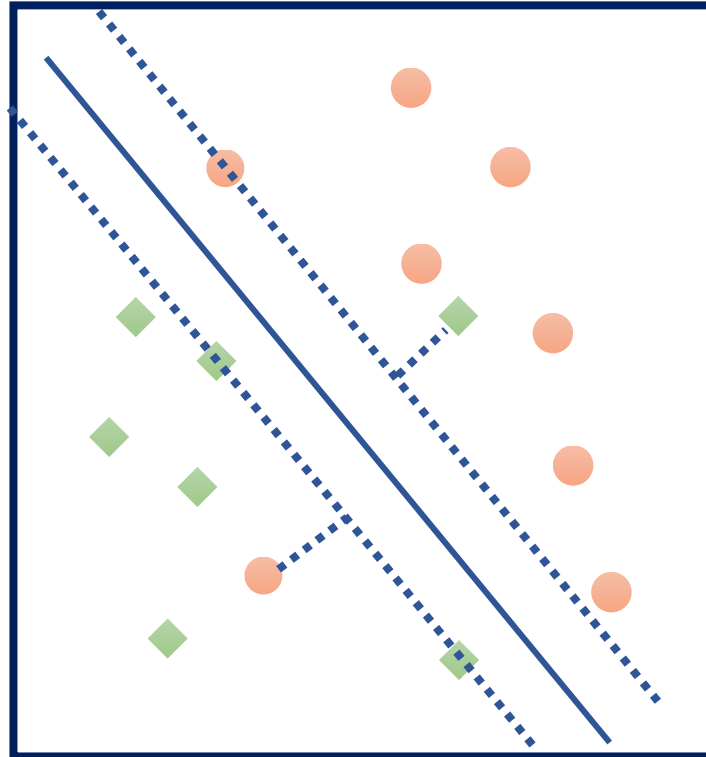
- × Relies on the assumption of equally important independent features
- × Not ideal for datasets with many numeric features
- × Estimated probabilities are less reliable than the predicted classes

Support vector machines

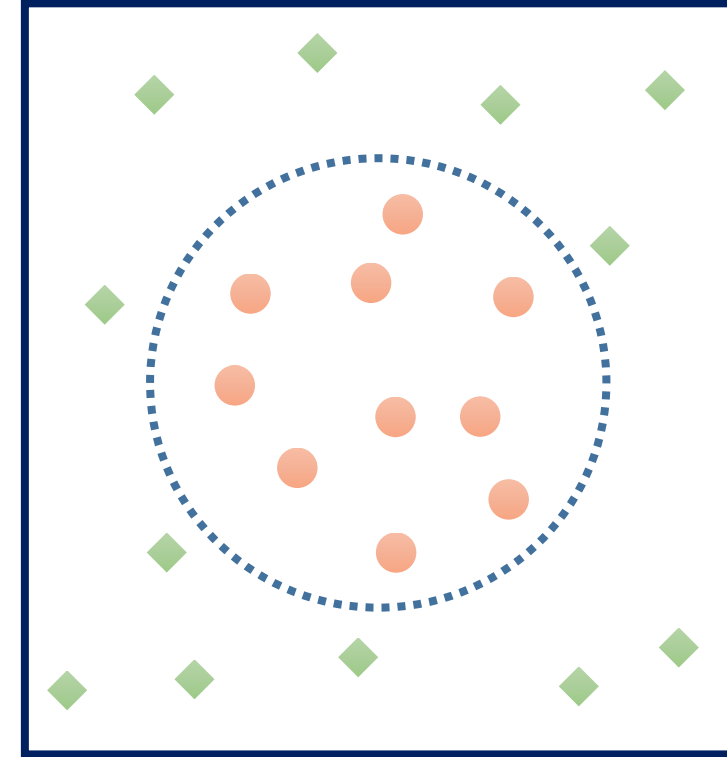
Maximum margin hyperplane



Not linearly separable



Nonlinear kernels



- SVMs use multidimensional surfaces to define relationships between features and outcomes
- Support vectors are data points from each class that are closest to the MMH.
- They combine aspects of instance-based methods and linear regression
- Uses the concept of maximum margin hyperplane: greatest separation between the classes

Support vector machines

Strengths

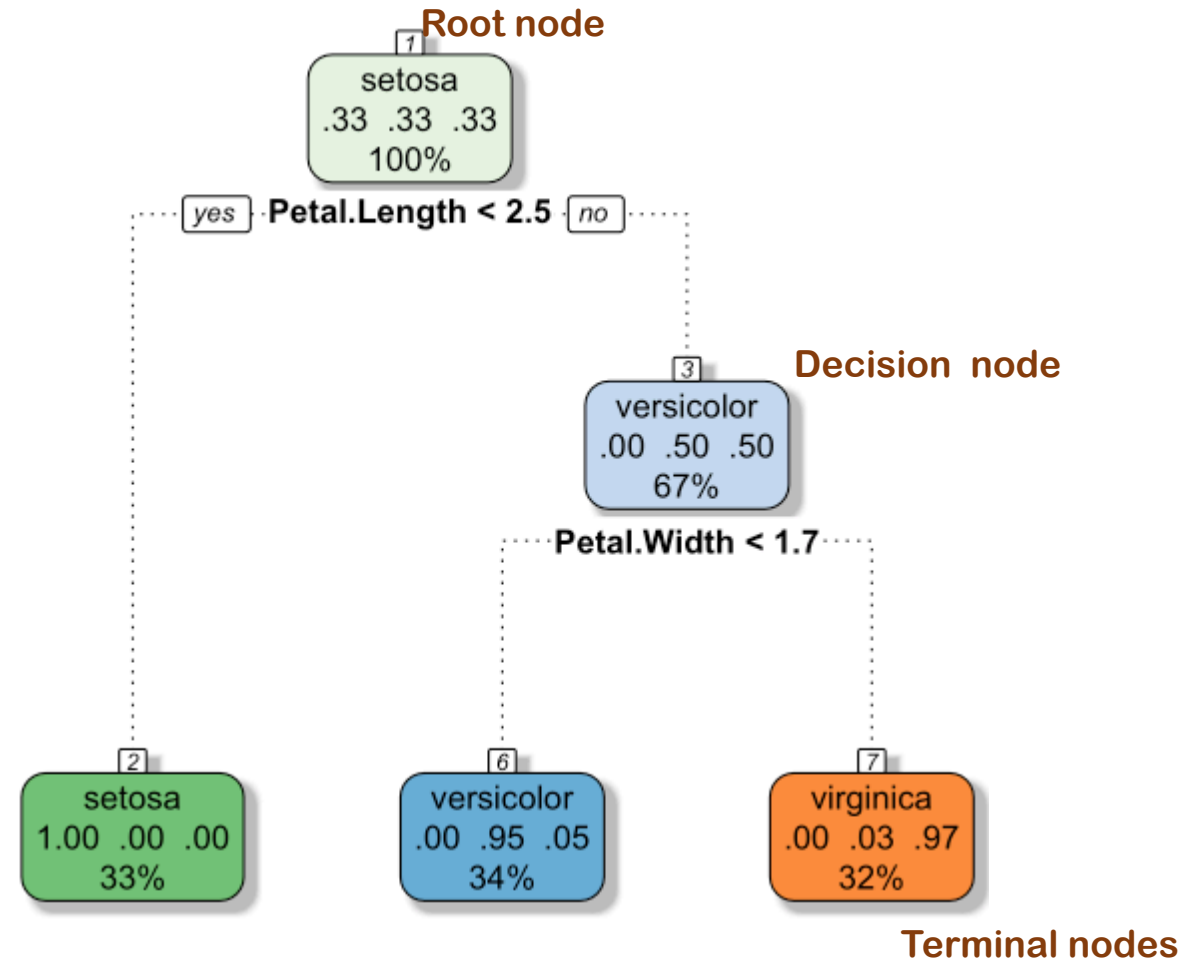
- ✓ Can be used for classification or regression
- ✓ Not overly influenced by noisy data and not very prone to overfitting
- ✓ Popularity due to high profile wins in data mining competitions

Weaknesses

- × So many options (kernels and model parameters)
- × Slow to train
- × Complex black box model

Decision trees

- DTs use a tree structure to model the relationships among the features and target variable
- Based on recursive partitioning (divide and conquer)
- Stopping criterion:
 - Homogenous
 - No more features to distinguish examples
 - The tree has grown to a predefined size limit
- Human readable format
- Not an ideal fit for large number of nominal features with many levels or many numeric features. Why?



Decision trees

Strengths

- ✓ Wide applicability
- ✓ Can handle numeric or nominal features, missing data
- ✓ Excludes unimportant features
- ✓ Works well with both small and large datasets
- ✓ Readily interpretable
- ✓ More efficient than complex models

Weaknesses

- × Typically biased towards splits on features with large number of levels
- × Very easy to overfit or underfit
- × Can have trouble modeling complex relationships
- × Small changes in training data can result in large changes in logic
- × Large trees are difficult to interpret

Regression techniques

Regression involves specifying a relationship between a dependent variable and one or more independent variables

$$\hat{y} = \beta_0 x_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \varepsilon$$

Parameters

Hyperparameters

$$RSS = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \cdots + (y_n - \hat{y}_n)^2$$

$$+ \lambda \left[\frac{(1 - \alpha)}{2} \sum_{i=1}^p \beta_i^2 + \alpha \sum_{i=1}^p |\beta_i| \right] \text{ Elastic net}$$

Classification: logistic regression

$$z = \beta_0 x_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \varepsilon$$

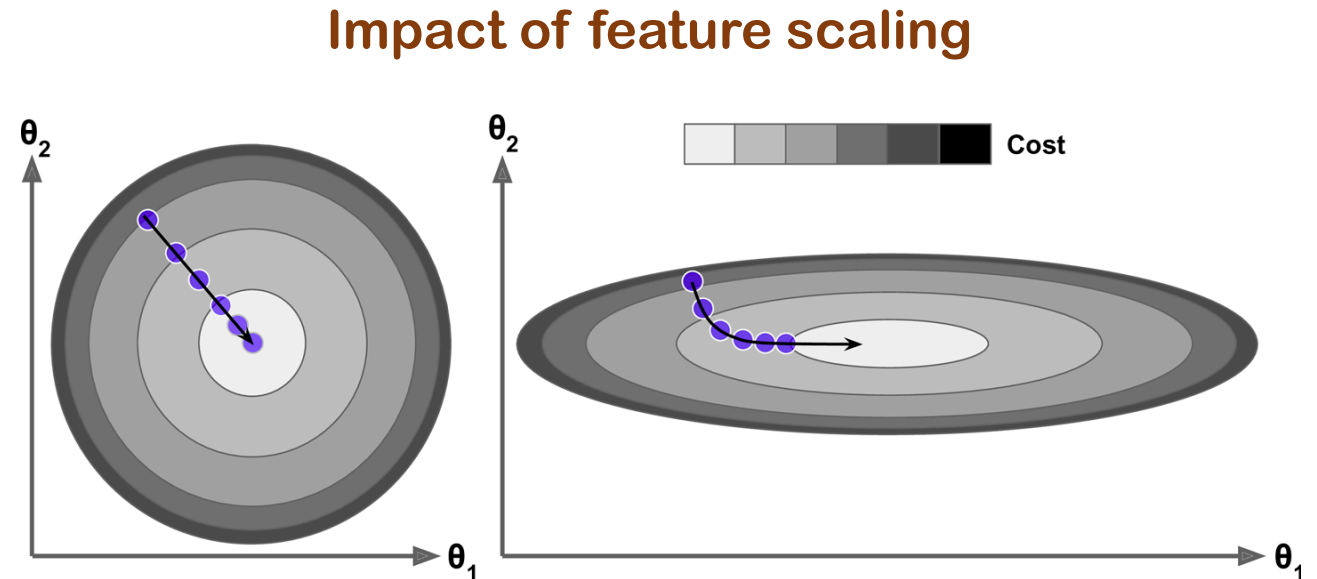
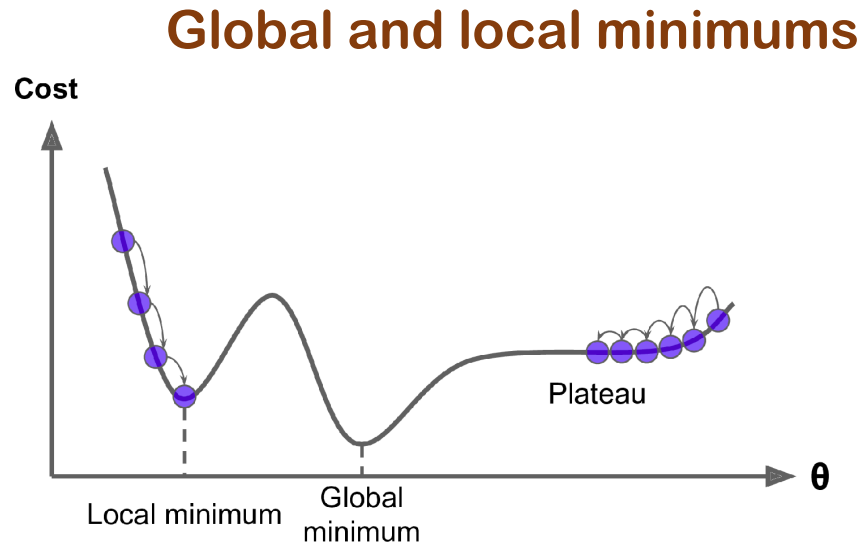
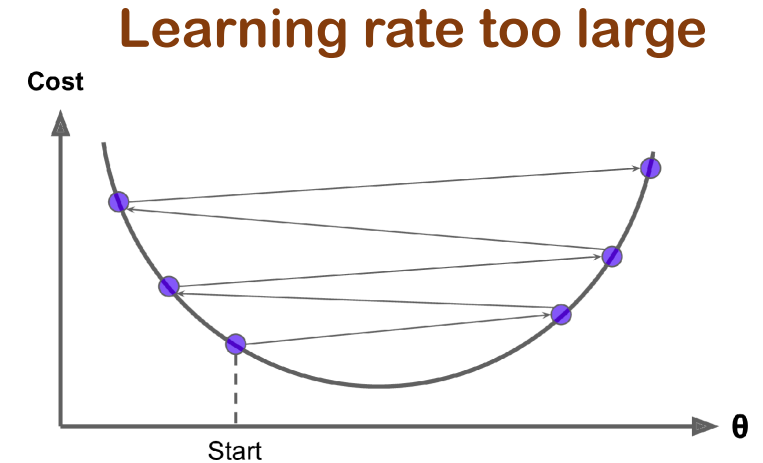
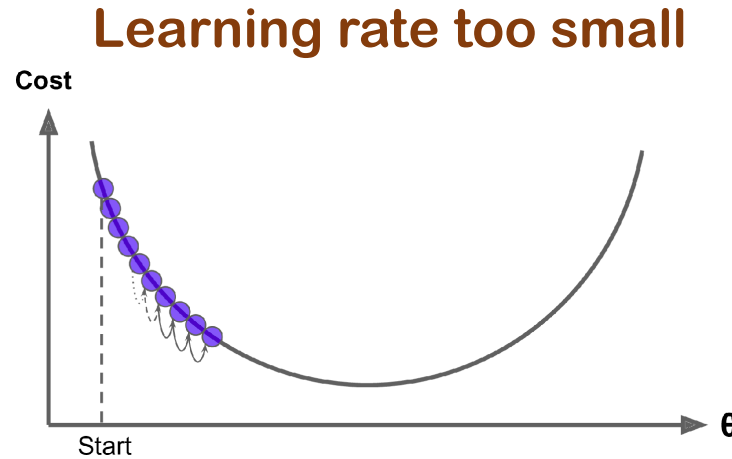
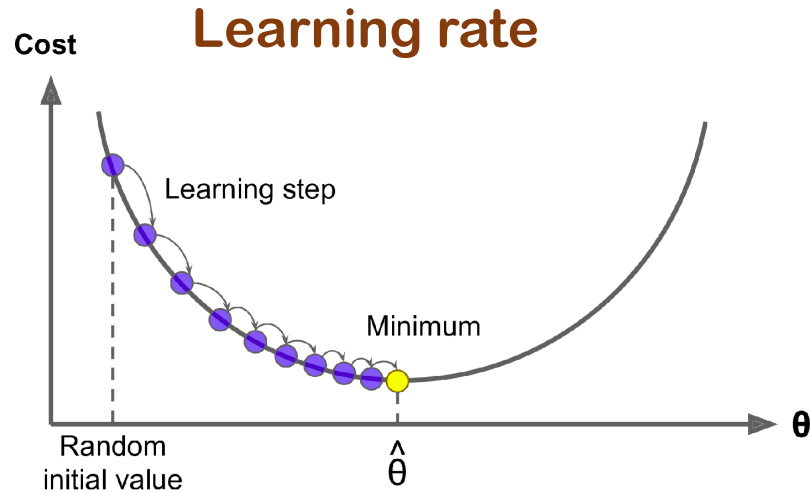
$$p = \frac{1}{1 + e^{-z}}$$

$$\hat{y} = \begin{cases} 0 & \text{if } p < 0.5 \\ 1 & \text{if } p \geq 0.5 \end{cases}$$

Multi-class classification

$$p_k = \frac{\exp(z_k)}{\sum_{k=1}^K \exp(z_k)}$$

Gradient descent and learning rates



Regression techniques

Strengths

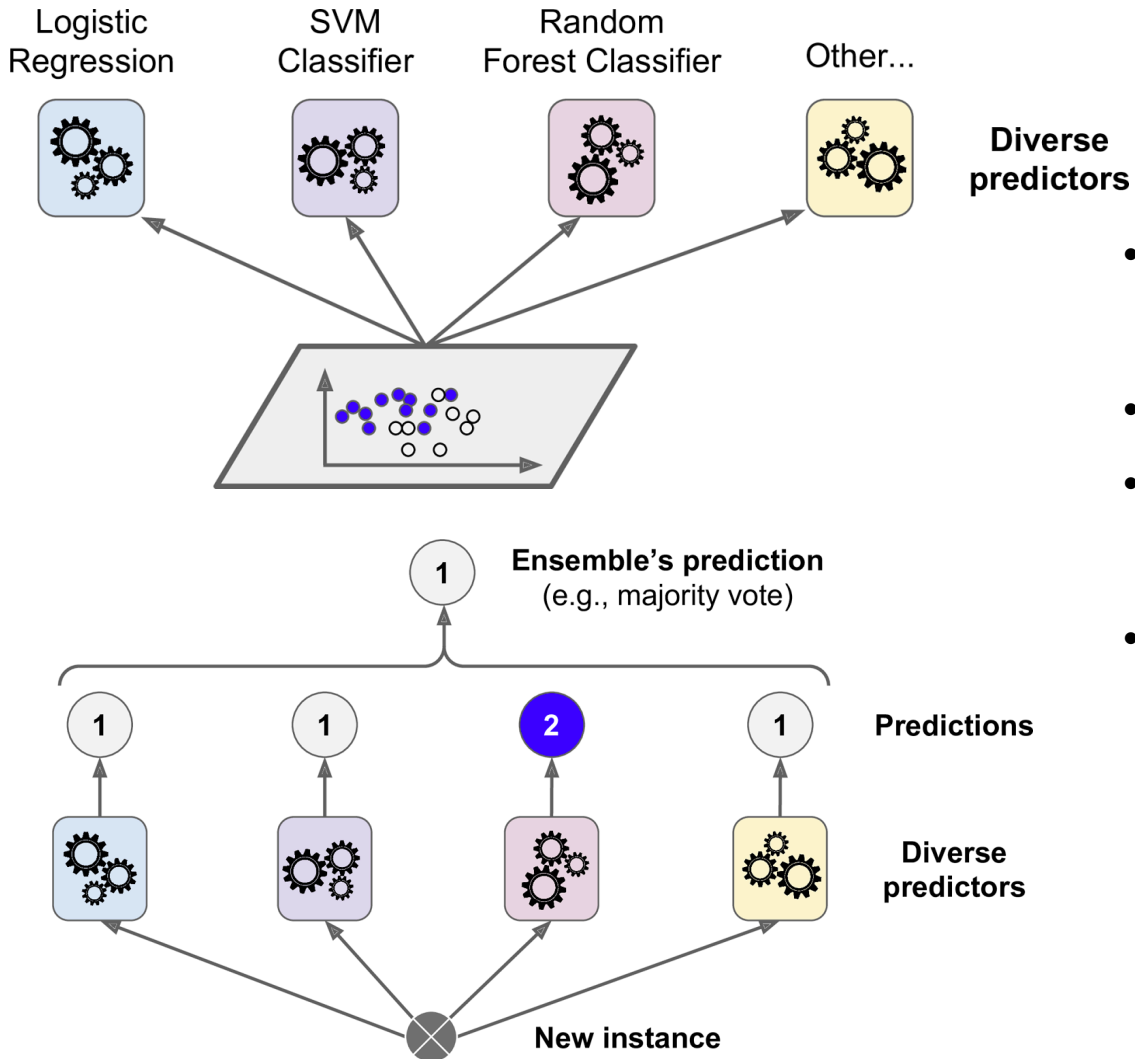
- ✓ By far the most common approach for modeling numeric data
- ✓ Can be adapted to model any task
- ✓ Provides estimates of the size and strength of relationships among features and the outcome

Weaknesses

- × Makes strong assumptions about the data
- × Model's form must be specified in advance
- × Cannot handle missing data
- × Works with only numeric features
- × Requires good knowledge of statistics to use/understand the models properly

Ensemble methods

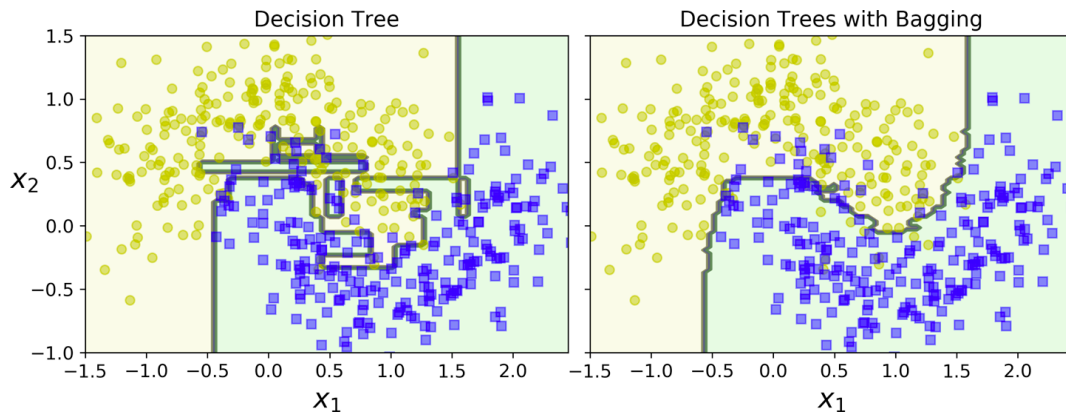
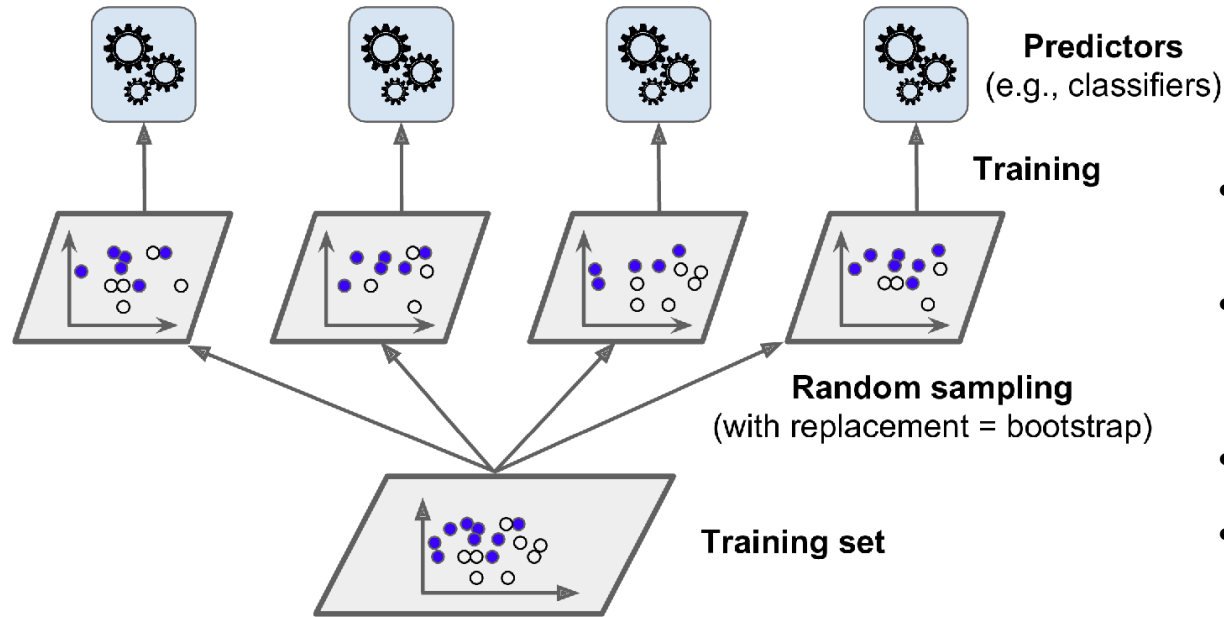
Voting classifiers



- Voting often achieves higher accuracy than the best classifier in the ensemble.
- Learners need to be diverse
- Hard voting uses class labels, soft voting using class probabilities
- Different kinds of voting schemes can be devised

Ensemble methods

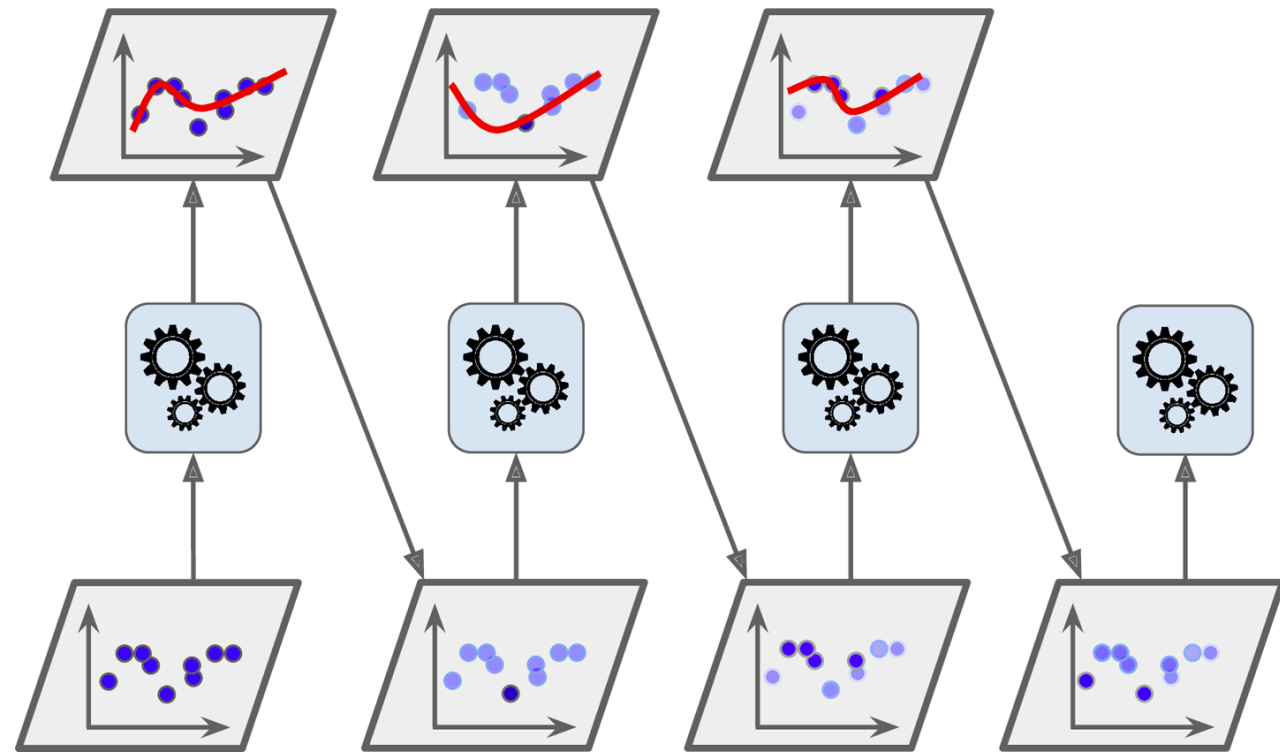
Bootstrap aggregating (Bagging)



- Uses the same training model but different random subsets of the data.
- Sampling with replacement: bagging, without replacement: pasting
- Out of bag instances can be used as a validation set
- Random forest is a bagging of decision trees with random subset of features used at each node of the tree
- Extremely randomized trees (extra-trees) use random thresholds at each node

Ensemble methods

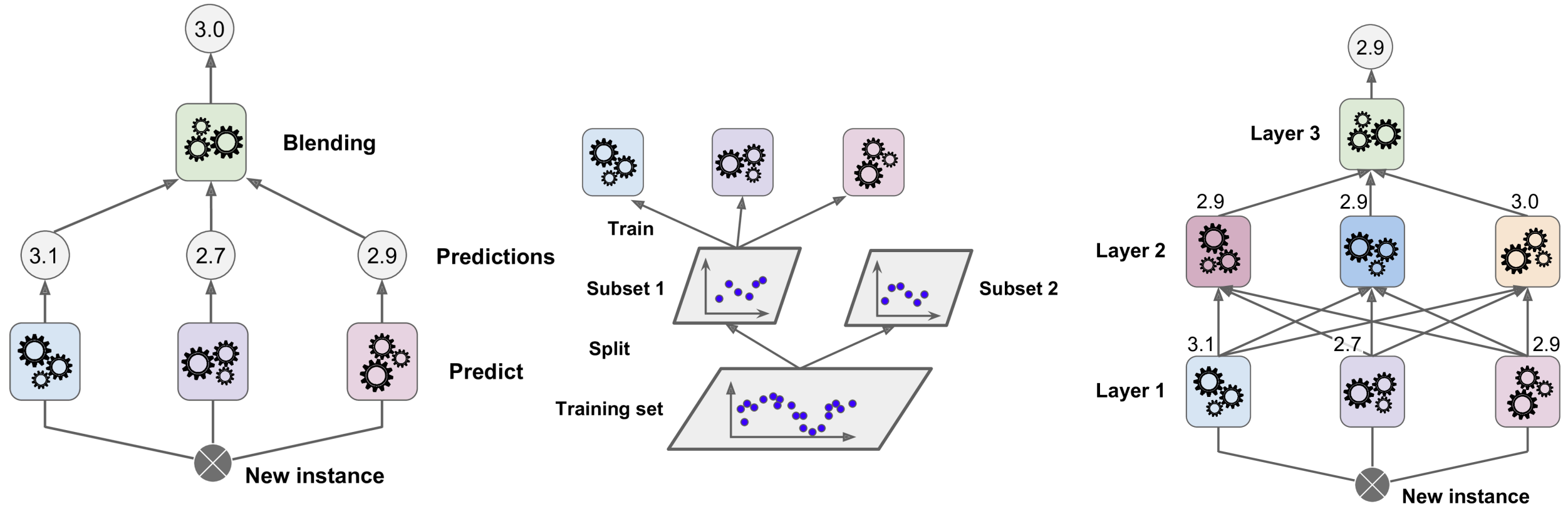
Boosting



- Training learners sequentially
- The new predictor tries to correct the errors its predecessor
- Pay more attention to the training instances that the predecessor underfitted (AdaBoost)
- In gradient boosting, the new predictor is fitted to the residual errors made by the previous predictor

Ensemble methods

Stacking



Ensemble methods

Strengths

- ✓ Improved performance
- ✓ Models typically have lower variances
and are more robust
- ✓ Can capture simple and complex
relationships in the data

Weaknesses

- × Can be difficult to train
- × Much larger parameter space
- × Reduced model interpretability

Discussion questions

1. How would you define Machine Learning?
2. What are the two most common supervised tasks?
3. What type of Machine Learning algorithm would you use to allow a robot to walk in various unknown terrains?
4. What type of algorithm would you use to segment your customers into multiple groups?
5. What type of learning algorithm relies on a similarity measure to make predictions?
6. What is the difference between a model parameter and a learning algorithm's hyperparameter?
7. What do model-based learning algorithms search for? What is the most common strategy they use to succeed? How do they make predictions?
8. If your model performs great on the training data but generalizes poorly to new instances, what is happening? Can you provide three possible solutions?
9. What is a test set, and why would you want to use it?
10. What is the purpose of a validation set?
11. What can go wrong if you tune hyperparameters using the test set?
12. What is an online learning system?