

INT104 knowledge domain

(Artificial Intelligence)

1 Data pre-processing

1.1 Data type

- Highly organized
- Usually with a label

1.2 Data Storage and Presentation

CSV, TSV, XML, JSON

1.3 Pearson's r correlation

$$r = \frac{N \sum xy - \sum x \sum y}{\sqrt{[N \sum x^2 - (\sum x)^2][N \sum y^2 - (\sum y)^2]}}$$

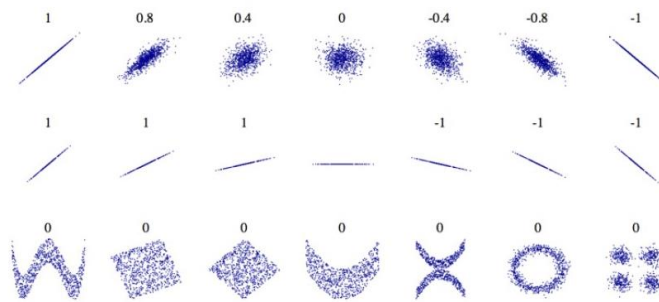


Figure 2-14. Standard correlation coefficient of various datasets (source: Wikipedia; public domain image)

1.4 Data Integration

Combine, Resolve conflicts, Remove redundant

1.5 Data Transformation

- Min-max normalization

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

- Z-score normalization

(Normalizing every value in a dataset such that the mean of all of the values is 0 and the standard deviation is 1)

$$x_{scaled} = \frac{x - mean}{sd}$$

- Normalization by decimal scaling

$$x_{scaled} = \frac{x}{10^j}$$

1.6 Feature Selection

- Filter methods – features are selected and ranked according to their relationships with the target
- Wrapper methods – it's a search for well-performing combinations of features
- Embedded methods – perform feature selection as part of the model training process

2 Classification & Model selection

2.1 Binary Classifier

- Classification algorithms find the mapping function to map the "x" input to "y" discrete output.

- Binary Classification refers to those classification tasks that have **two** class labels.

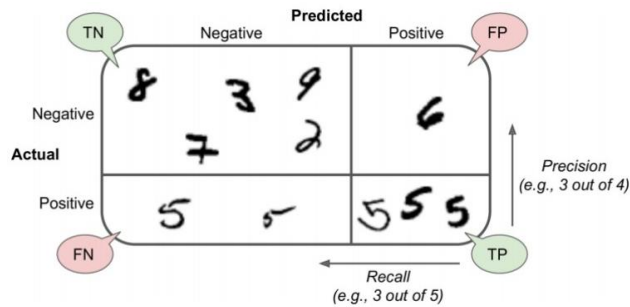
2.2 Performance Measures

Metrics to evaluate Classification models:

- **Accuracy**

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

- Confusion Matrix (not a metric but fundamental to others)



- True Positive (**TP**): Predict an observation belongs to a class and it actually does belong to that class.
- True Negative (**TN**): Predict an observation does not belong to a class and it actually does not belong to that class.
- False Positives (**FP**): Predict an observation belongs to a class and it actually does not belong to that class.
- False Negatives (**FN**): Predict an observation does not belong to a class and it actually does belong to that class.

- **Precision** and **Recall**

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

- **F1-score**

$$F_1 \text{ score} = \frac{2 \cdot TP}{2 \cdot TP + FP + FN} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

- F_β -score: a more flexible F score that combines precision and recall

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

- $\beta < 1$ focuses more on precision
- $\beta > 1$ focuses more on recall

- Specificity (Selectivity)

$$\text{Specificity} = \frac{TN}{TN + FP}$$

- Fall-out

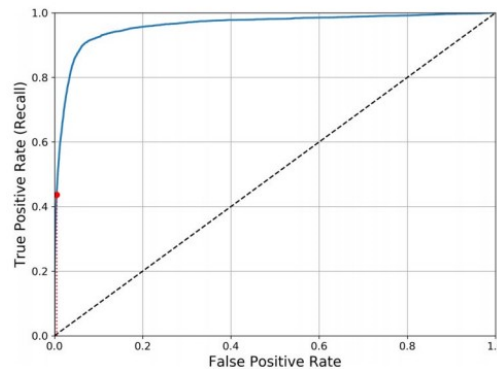
$$\text{Fall-out} = \frac{FP}{TN + FP}$$

- Miss

$$\text{Miss} = \frac{FN}{TP + FN}$$

• AUC&ROC

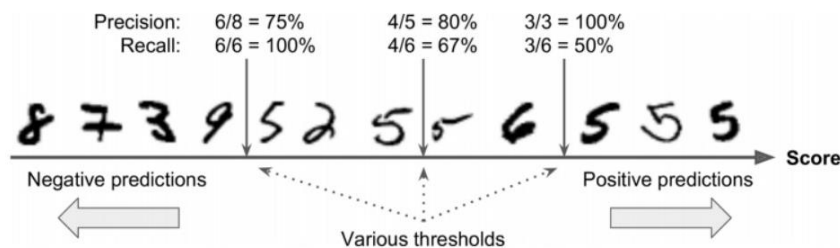
A ROC curve (receiver operating characteristic curve) is a graph showing the performance of a classification model at all classification thresholds. In another word, it presents **Recall (True Positive Rate) VS Fall-out (False Positive Rate)**



2.3 Precision/Recall Trade-off

The **threshold** of a decision function:

- If the threshold is small, FP is large → high recall, low precision (Strict, most real objects are detected)
- If the threshold is large, FN is large → low recall, high precision (Generous, most of the prediction are correct)

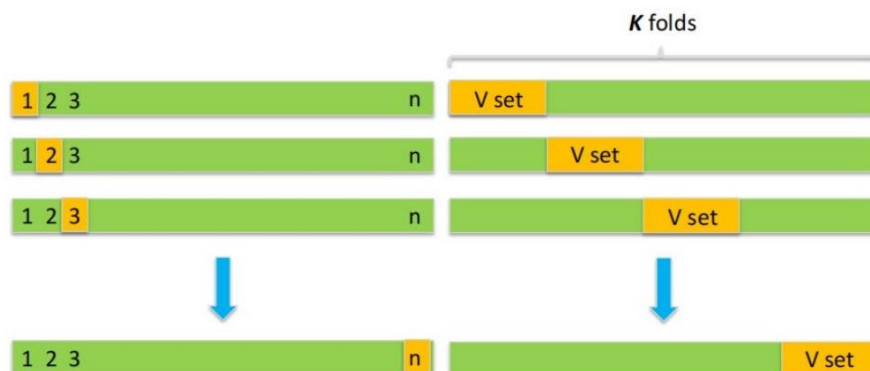


2.4 Cross Validation

To avoid selecting the parameters that perform best on the test data but maybe not the parameters that generalize best, we can further split the training set into training fold and validation fold:



- Training fold: used to fit the model
- Validation fold: used to estimate prediction error for model selection
- Test set: used for assessment of the prediction error of the final chosen model



Leave-one-out Cross Validation (LOOCV):

1. Split the training data set of size n into:
 - Training fold ($n-1$ elements)
 - Validation fold (1 element)

2. Fit the model using the training data set.
3. Evaluate the model using validation set and compute the corresponding mean squared error (MSE).
4. Repeat this process n times, producing n squared errors. The average of these n squared errors estimates the test MSE.

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n MSE_i$$

Merits:

- Less bias: it trains on almost all of the data means that the model is trained on as much data as possible

Demerits:

- Sensitive to outliers
- Due to n times model fitting, the computation of LOOCV is intensive

K-fold Cross Validation

1. Randomly divide the data set into K folds (they are mutually exclusive)
2. Treat one fold as a validation set (normally from the 1st to the Kth), and fit the model on the remaining K-1 folds. Then, computing the MSE on the observations in the validation set. (This operation is repeated K times)
3. Compute the mean MSE (or other performance measure) as following

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k MSE_i$$

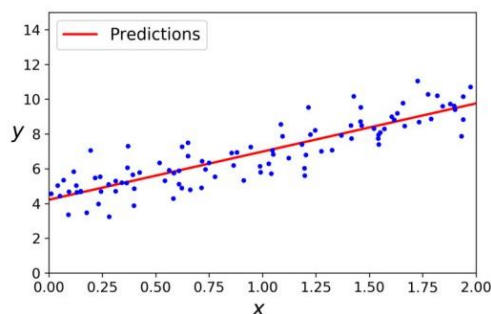
2.5 Multiclass Classification

Multiclass classification refers to classification tasks that can distinguish between more than two classes:

- One-versus-the-rest(OvR) strategy: train multiple binary classifiers for each class, select the class whose classifier outputs the highest score.
 - train N times
- One-versus-one (OvO) strategy: train a binary classifier for every pair of classes
 - train $N(N-1)/2$ times

3 Training models

3.1 Simple Linear Regression



- $Y = kx + b$
- Y: predicted value (output)
- x: feature value (input)
- k, b: model parameters (b: bias term, k: weight)

Usually the predicted value (fitted value) \hat{y} is not perfect. The difference between the fitted value and real value is known as residuals \hat{e}

$$\hat{y}^{(i)} = kx^{(i)} + b$$

$$\hat{e}_i = y^{(i)} - (kx^{(i)} + b) = y^{(i)} - \hat{y}^{(i)}$$

The regressed value usually pursues a minimum of **residual sum of square (RSS)**

$$RSS = \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

A linear model makes a prediction by simply computing a weighted sum of the input features, plus a constant called the bias term (also called the **intercept term**):

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

- \hat{y} is the predicted value.
- n is the number of features.
- x_i is the i^{th} feature value.
- θ_j is the j^{th} model parameter (including the bias term θ_0 and the feature weights $\theta_1, \theta_2, \dots, \theta_n$).

Linear Regression model prediction (**vectorized form**):

$$\hat{y} = h_{\theta}(\mathbf{x}) = \boldsymbol{\theta} \cdot \mathbf{x}$$

- $\boldsymbol{\theta}$ is the model's *parameter vector*, containing the bias term θ_0 and the feature weights θ_1 to θ_n .
- \mathbf{x} is the instance's *feature vector*, containing x_0 to x_n , with x_0 always equal to 1.
- $\boldsymbol{\theta} \cdot \mathbf{x}$ is the dot product of the vectors $\boldsymbol{\theta}$ and \mathbf{x} , which is of course equal to $\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$.
- h_{θ} is the hypothesis function, using the model parameters $\boldsymbol{\theta}$.

Evaluation:

- Cost function: **Mean Squared error (MSE)** for a Linear Regression model

$$\text{MSE}(\mathbf{X}, h_{\theta}) = \frac{1}{m} \sum_{i=1}^m \left(\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)} \right)^2$$

Training the model is the process to find the value of θ that minimizes the cost function.

- Normal Equation:

$$\hat{\boldsymbol{\theta}} = (X^T X)^{-1} X^T \mathbf{y}$$

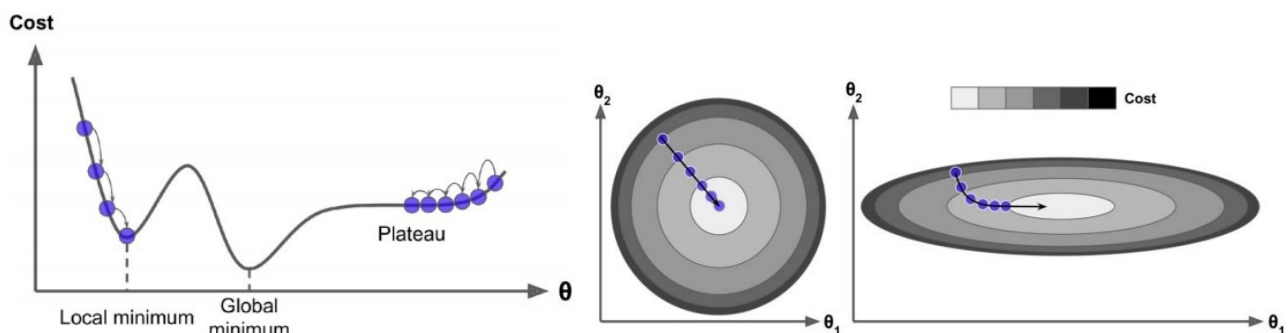
- $\hat{\theta}$ is the value of θ that minimizes the cost function
- \mathbf{y} is the vector of targeted values containing $y^{(1)}$ to $y^{(m)}$

3.2 Gradient Descent

The **MSE cost function** for a Linear Regression model is continuous and convex function.

Gradient Descent is guaranteed to approach arbitrarily close the **global minimum**.

Features with very different scales have different cost time.



3.2.1 Batch Gradient Descent

Use the whole training set to compute the gradients at every step

$$\nabla_{\theta} \text{MSE}(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} \text{MSE}(\theta) \\ \frac{\partial}{\partial \theta_1} \text{MSE}(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_n} \text{MSE}(\theta) \end{pmatrix} = \frac{2}{m} \mathbf{X}^T (\mathbf{X}\theta - \mathbf{y})$$

$$\theta^{(\text{next step})} = \theta - \eta \nabla_{\theta} \text{MSE}(\theta)$$

$$(\theta = [\theta_0, \theta_1, \dots, \theta_n]^T)$$

3.2.2 Stochastic Gradient Descent

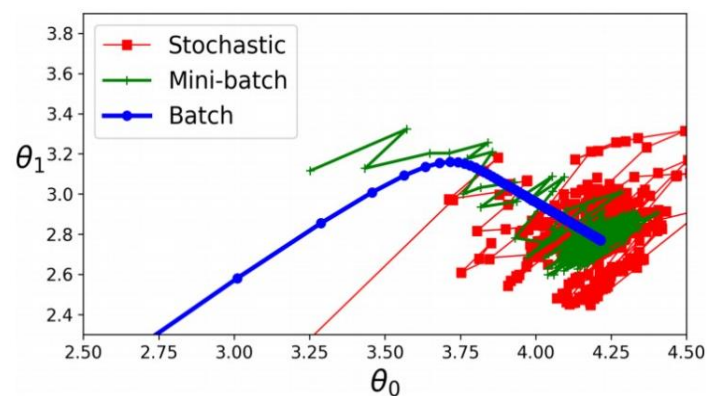
Picks a random instance in the training set at every step and compute based only on that single instance

$$\theta^{(\text{next step})} = \theta - \eta \nabla_{\theta} \text{MSE}(\theta; x^{(i)}, y^{(i)})$$

3.2.3 Mini-batch Gradient Descent

Computes the gradients on small random sets of instances called mini-batches

Compare these three methods

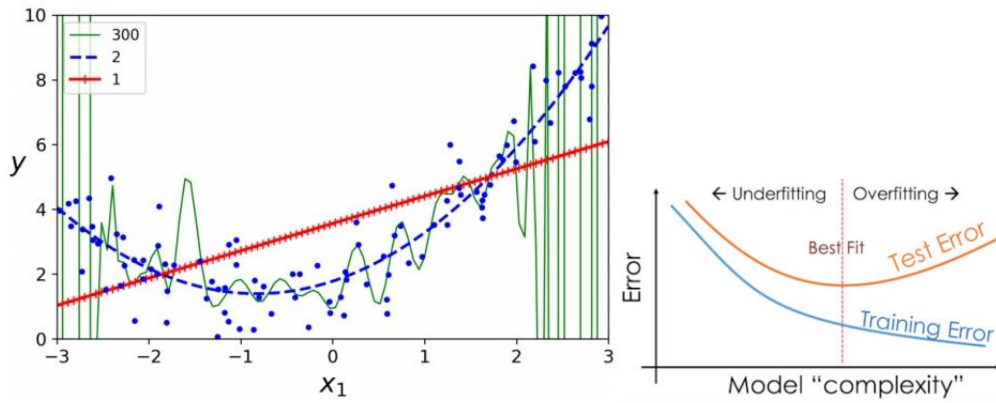


	Method	Pros	Cons
★	Batch Gradient Descent	Guaranteed convergence to the global minimum.	Computationally expensive for large datasets. Requires the entire dataset to be loaded into memory.
	Stochastic Gradient Descent	Computationally efficient for large datasets.	Convergence to the global minimum is not guaranteed. The noise in the updates can cause the loss function to oscillate.
★	Mini-Batch Gradient Descent	Better convergence than stochastic gradient descent. Computationally efficient for large datasets.	Convergence to the global minimum is not guaranteed.

3.3 Polynomial Regression

Better but problems:

- Bias: refers to the error from erroneous assumptions in the learning algorithm (inability to capture the underlying patterns in the data). [\[Generalization\]](#)
- Variance: refers the error from sensitivity to small fluctuations in the training data (difference in fits between data sets). [\[Robustness\]](#)



- Underfitting (1d): Poor performance on training data and poor performance on validation data
- Overfitting(300d): Good performance on training data but poor performance on validation data

3.4 Regularized Linear Models

- Ridge Regression(L2) cost function:

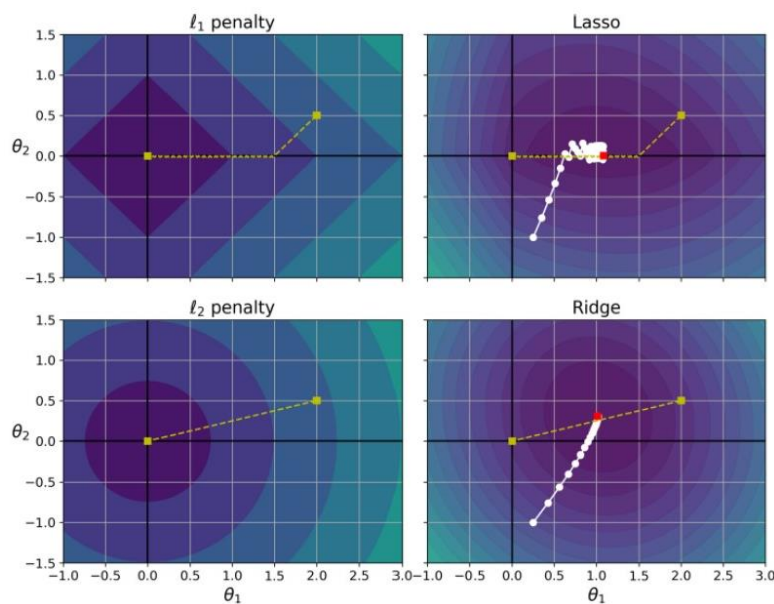
$$J(\theta) = \text{MSE}(\theta) + \alpha \frac{1}{2} \sum_{i=1}^n \theta_i^2$$

<Forces the learning algorithm to not only fit the data but also keep the model weights as small as possible>

- Lasso Regression(L1) cost function:

$$J(\theta) = \text{MSE}(\theta) + \alpha \sum_{i=1}^n |\theta_i|$$

<Lasso Regression automatically performs feature selection and outputs a sparse model>

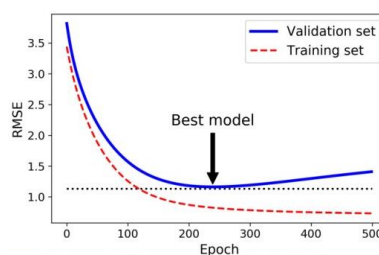


- Elastic Net cost function:

$$J(\theta) = \text{MSE}(\theta) + r\alpha \sum_{i=1}^n |\theta_i| + \frac{1-r}{2}\alpha \sum_{i=1}^n \theta_i^2$$

<It is a middle ground between Ridge Regression and Lasso Regression>

- Early stopping



<To stop training as soon as the validation error reaches a minimum>

3.5 Logistic Regression

Commonly used to estimate the probability that an instance belongs to a particular class

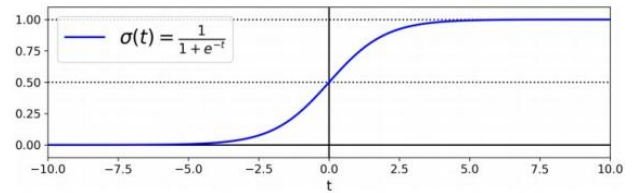
Logistic Regression model:

$$\hat{p} = h_{\theta}(\mathbf{x}) = \sigma(\mathbf{x}^T \boldsymbol{\theta})$$

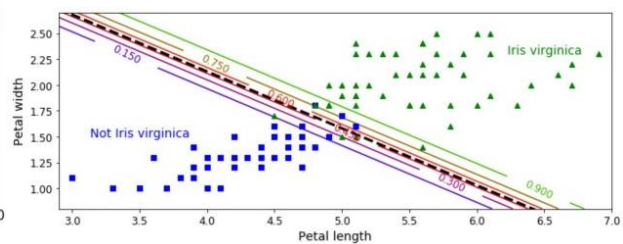
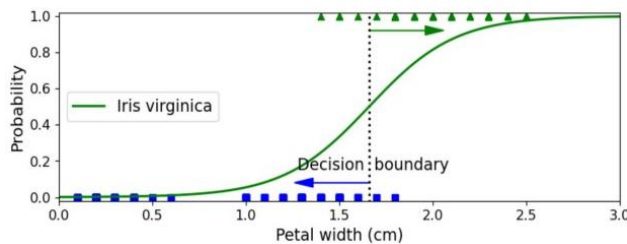
Sigmoid function:

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

$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$



Decision Boundaries:



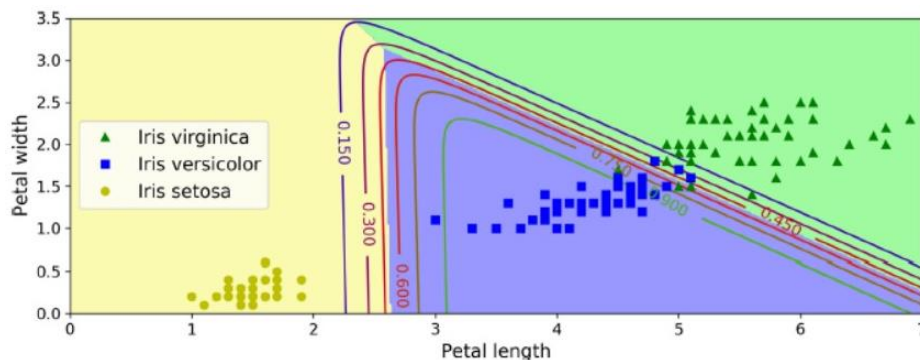
Softmax Regression: for Multinomial Logistic Regression

Softmax score for class k:

$$s_k(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\theta}^{(k)}$$

Softmax function:

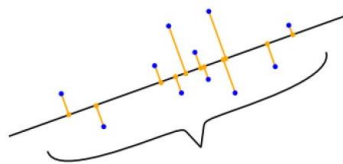
$$\hat{p}_k = \sigma(\mathbf{s}(\mathbf{x}))_k = \frac{\exp(s_k(\mathbf{x}))}{\sum_{j=1}^K \exp(s_j(\mathbf{x}))}$$



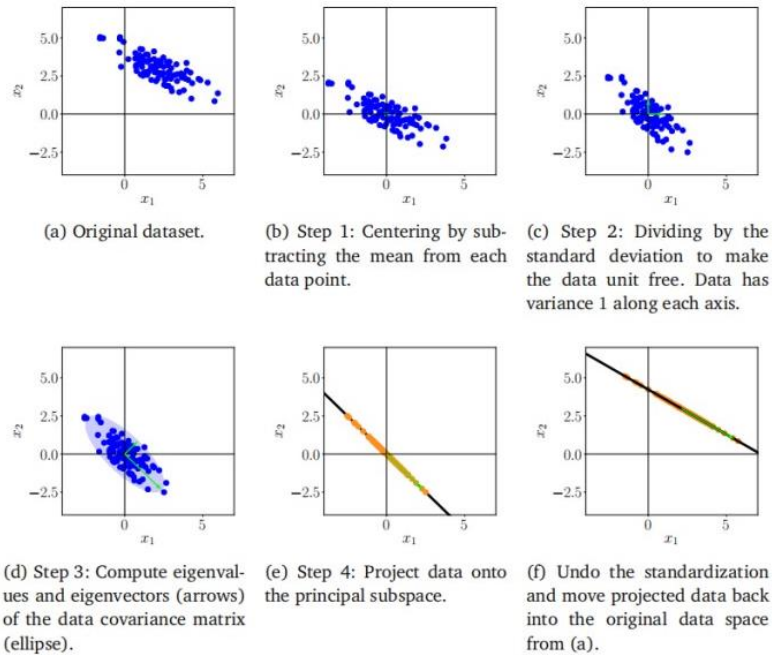
4 Dimensionality reduction

4.1 Principal Component Analysis (PCA)

PCA identifies the axis that accounts for the largest amount of variance in the training set. Focuses on the mutual independence of the components. Vectors are orthogonal.



Steps:



Variance on C1:

$$V_1 = \frac{1}{M} \sum_{i=1}^M (c_1^T x^{(i)})^2 = \frac{1}{M} \sum_{i=1}^M c_1^T x^{(i)} x^{(i)T} c_1 = c_1^T \left(\frac{1}{M} \sum_{i=1}^M x^{(i)} x^{(i)T} \right) c_1 = c_1^T S c_1$$

Data covariance matrix:

$$S = \frac{1}{M} \sum_{i=1}^M x^{(i)} x^{(i)T}$$

(S is an N*N matrix, N is the number of features, M is the total number of data points)

Practice: Given a dataset that consists of the following points below:

A=(2, 3), B=(5, 5), C=(6, 6), D=(8,9)

1. Calculate the covariance matrix for the dataset.
2. Calculate the eigenvalues and eigenvectors of the covariance matrix.

```
# Define the dataset
data <- matrix(c(2, 3, 5, 5, 6, 6, 8, 9), ncol = 2, byrow = TRUE)
colnames(data) <- c("x", "y")

# Calculate the covariance matrix
cov_matrix <- cov(data)

# Calculate the eigenvalues and eigenvectors
eigen_result <- eigen(cov_matrix)
eigenvalues <- eigen_result$values
eigenvectors <- eigen_result$vectors

# Print the results
cat("Covariance matrix:\n")
print(cov_matrix)
cat("\nEigenvalues:\n")
print(eigenvalues)
cat("\nEigenvectors:\n")
print(eigenvectors)
```

R

Singular Value Decomposition (SVD)

Theorem:

Let $A \in \mathbb{R}^{m \times n}$ be a rectangular matrix of rank $r \in [0, \min(m, n)]$. The SVD of A is a decomposition of the form:

$$\begin{matrix} n \\ \boxed{A} \\ m \end{matrix} = \begin{matrix} m \\ \boxed{U} \\ m \end{matrix} \begin{matrix} m \\ \boxed{\Sigma} \\ m \end{matrix} \begin{matrix} n \\ \boxed{V^T} \\ n \end{matrix}$$

Formula:

$$A = [x_1 \ \dots \ x_m]_{n \times m} = U \Sigma V^T = [u_1 \ \dots \ u_n]_{n \times n} \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_m \\ 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & 0 \end{bmatrix}_{n \times m} [v_1 \ \dots \ v_m]_{m \times m}^T$$

Principal components matrix:

$$V = \begin{pmatrix} | & | & & | \\ c_1 & c_2 & \dots & c_n \\ | & | & & | \end{pmatrix}$$

($c_1, c_2 \dots c_n$ are orthogonal)

Projecting Down to d Dimension:

$$X_{d-proj} = XW_d$$

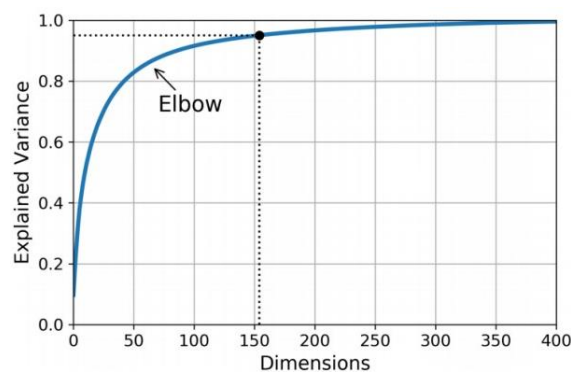
(W_d is the first d eigen vectors of data covariance matrix)

Explained **Variance Ratio**:

$$\frac{\lambda_1}{\lambda_1 + \lambda_2 + \dots + \lambda_n} \text{ (eigenvalue/ total eigenvalue)}$$

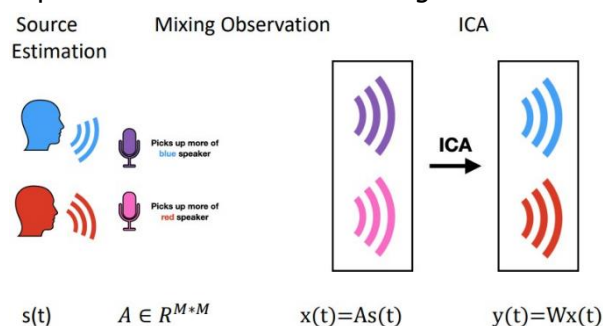
Choosing the Right Number of Dimensions:

- Choose the number of dimensions that add up to sufficiently large portion of the variance (e.g., 95%)
- Variance Ratio > 95%

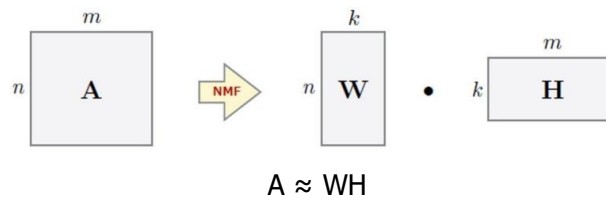


4.2 Independent Component Analysis (ICA)

ICA attempts to decompose (分解) a multivariate signal into independent non-Gaussian signals. Focuses on the mutual independence of the components. Vectors are not orthogonal.



4.3 Non-negative Matrix Factorization (NMF)



Optimization: minimize $D(A||WH)$ $s.t. W, H \geq 0$

A matrix factorization where everything is non-negative

- $A \in R^{n \times m}$ is original non-negative data
- $W \in R^{n \times k}$ is matrix of basis vectors, dictionary elements
- $H \in R^{k \times m}$ is matrix of activations, weights or gains
- NMF is not unique

4.4 Other Techniques

- Manifold Learning

Data lies on d-dimensional manifold is a part of an n-dimensional space (where $d < n$)

- Locally Linear Embedding (LLE)

LLE is a powerful nonlinear dimensionality reduction (NLDR) technique. It is a Manifold Learning technique that does not rely on projections.

- Multidimensional Scaling (MDS)

Trying to preserve the distances between the instances.

- Isomap

Trying to preserve the geodesic distances between the instances.

- t-Distributed Stochastic Neighbour Embedding (t-SNE)

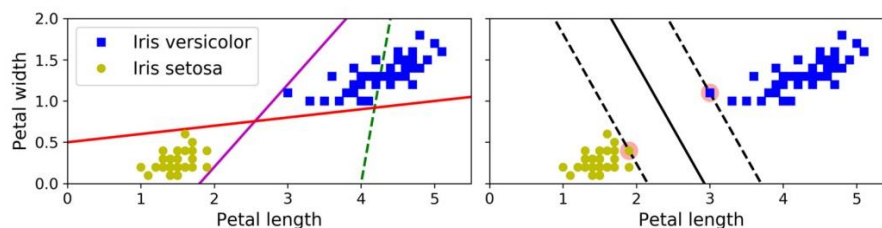
Trying to keep similar instances close and dissimilar instances apart.

5 Support Vector Machine (SVM)

5.1 Linear SVM Classification

5.1.1 Large Margin Classification

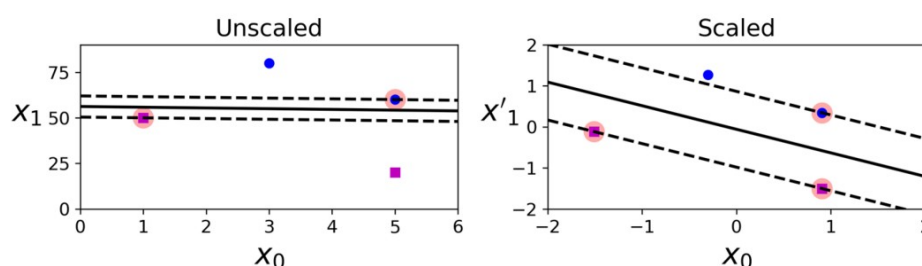
- If classes are linearly separable, two classes can be separated with a straight line
- SVM pursues the widest possible street between classes



- Decision boundary is not affected by more training instances
- It is determined by support vectors (instances located on the edge of street)

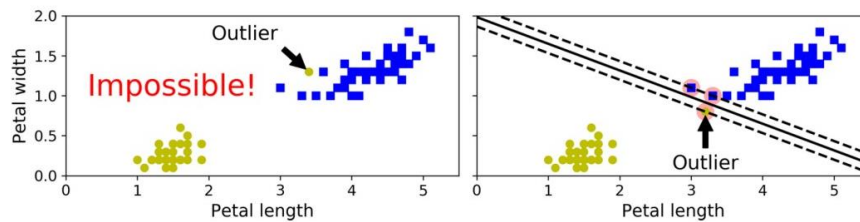
5.1.2 Feature Scales

The decision boundary could be much better if the feature is scaled:



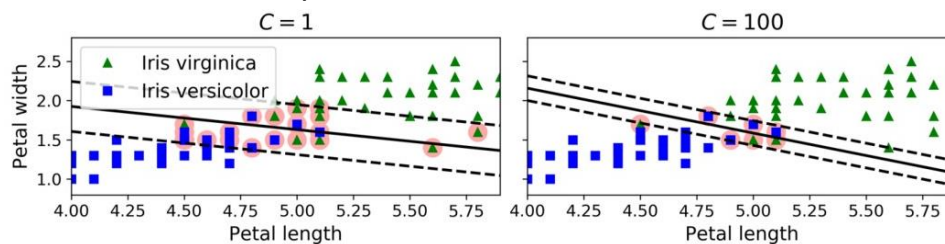
5.1.3 Hard Margin Classification

- All instances being off the street and on the right side is named "hard margin classification"
- The main limitation of hard margin classification is
 - The data must be **linearly separable**
 - There are **little outliers**



5.1.4 Soft Margin Classification

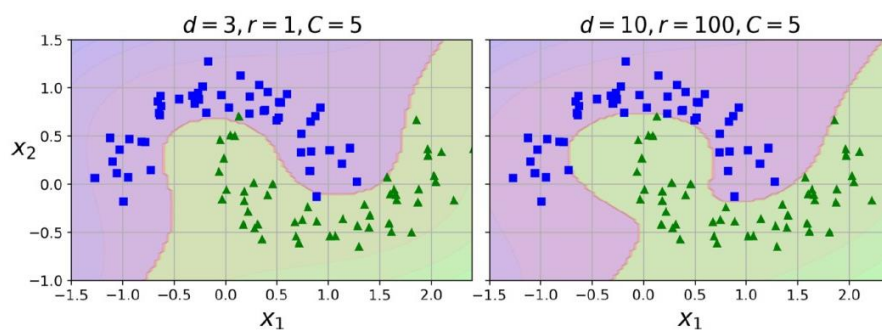
- Soft Margin Classification provides more flexibility
- The algorithm balances
 - The width of street
 - The amount of margin violations
- A hyper-parameter C is defined
 - A low value of C leads to more margin violation
 - A high value of C limits the flexibility



5.2 Nonlinear SVM Classification

5.2.1 Polynomial Kernel

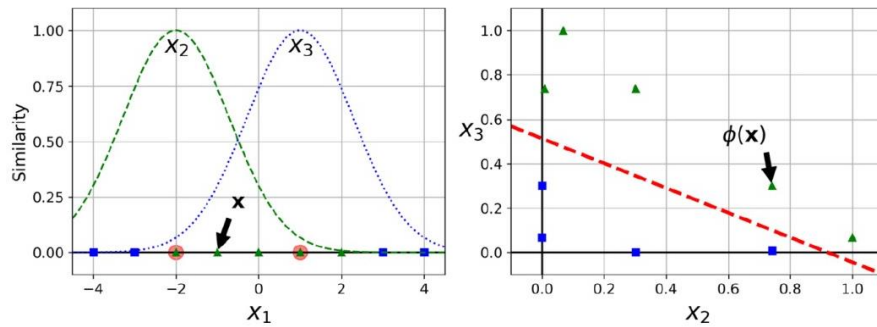
A kernel can be used instead of adding the polynomial features



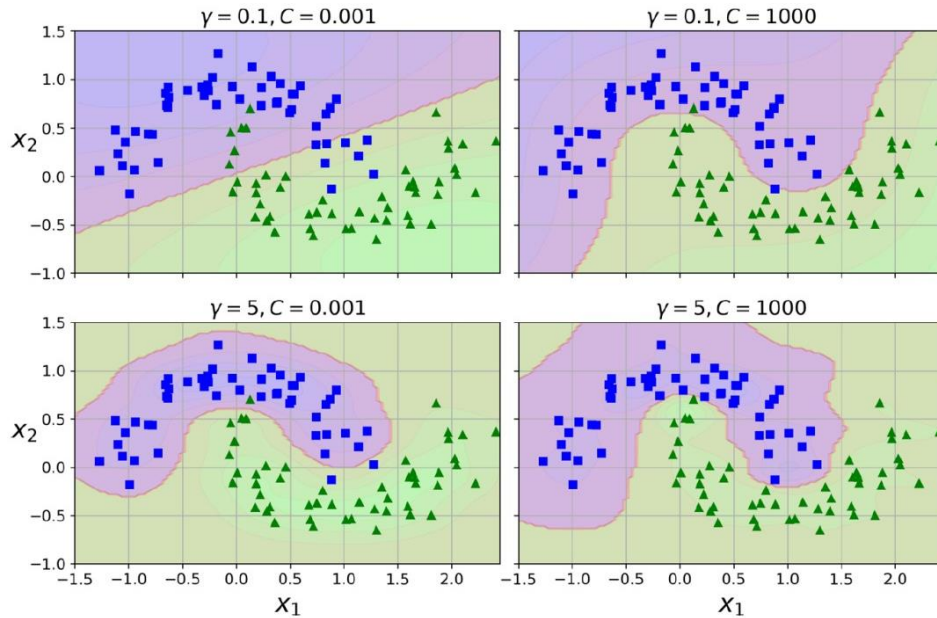
5.2.2 Similarity Feature

- Samples could be clustered with a landmark (reference point) I
- The similarity function could be defined as a Gaussian Radial Basis Function:

$$\Phi_{\gamma}(\mathbf{x}, l) = e^{-\gamma \|\mathbf{x} - l\|^2}$$

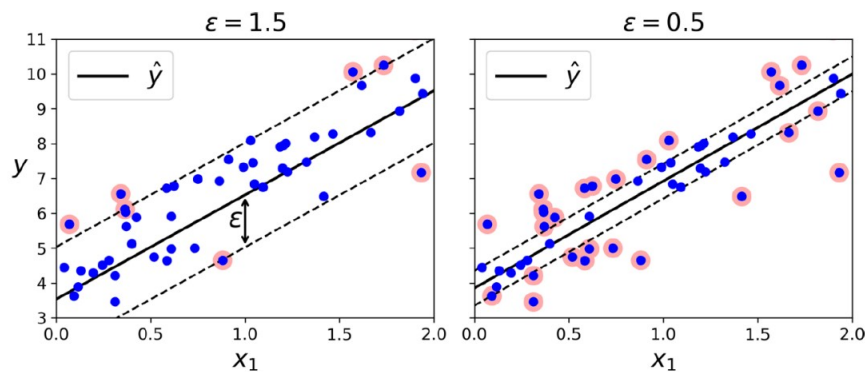


5.2.3 Gaussian RBF Kernel

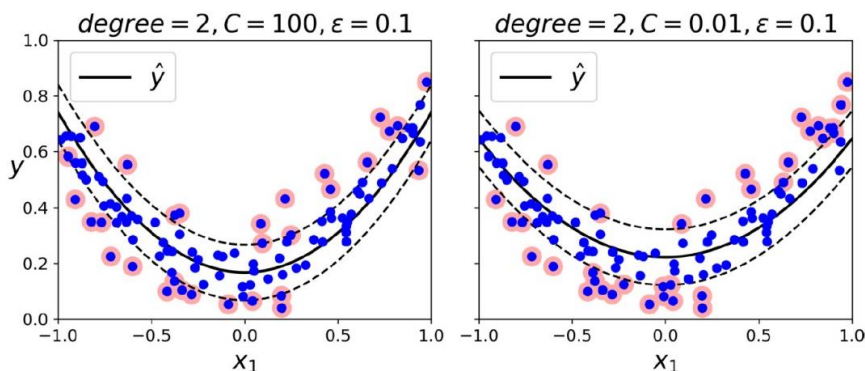


5.3 SVM Regression

- Fit as many instances as possible on the street
- Limit margin violations
- Width of street is controlled by a hyper-parameter ϵ



Use Kernel tricks:



6 Naïve Bayes

6.1 Bayes' Rule

$$\star P(c|\mathbf{x}) = \frac{P(c)P(\mathbf{x}|c)}{P(\mathbf{x})}$$

(Where c is considered as a class, \mathbf{x} is considered as a set of samples)

- $P(c)$ is named as prior probability
- $P(\mathbf{x}|c)$ is named as class-conditional probability (CCP, also known as "likelihood")
- $P(c|\mathbf{x})$ is named as posterior probability
- $P(\mathbf{x})$ is considered as evidence factor (observation)

6.2 Bayes' Rule for Classification

MAP estimation: maximise the posterior probability of observations

Presume that $\mathbf{x} \in D_c$ means that a sample \mathbf{x} belongs to class c (Dataset)

Recall Bayes' Rule:

$$PosteriorProbability = \frac{CCP \times PriorProbability}{Observation}$$

As the observation is same (the same training dataset), we have:

$$PosteriorProbability \propto CCP \times PriorProbability$$

According to Law of Large Numbers, the prior probability can be taken as the probability resulted from the frequency of observations:

$$p(D|\Theta) = p(c) \prod_c \prod_{\mathbf{x}_c \in D_c} p(\mathbf{x}_c|\Theta_c)$$

6.3 Naïve Bayes Classifier

A way to simplify the process is to assume the conditions / features of Θ_c are independent to each other

So, assume that $\Theta_c = (\theta_{c1}, \theta_{c2}, \dots, \theta_{cl})$, we have:

$$p(\mathbf{x}_c|\Theta_c) = \prod_{i=1}^l p(\theta_i)$$

Example:

<Dataframe>				
Outlook	Temperature	Humidity	Windy	Play
Overcast	Hot	High	False	Yes
Overcast	Cool	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Rainy	Mild	Normal	False	Yes
Rainy	Mild	High	True	No
Sunny	Hot	High	True	No
Sunny	Hot	High	False	No
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Sunny	Mild	Normal	True	Yes

Q1: Will you play on the day of Mild?

Temp.	Yes	No	p
Hot	2	2	0.28
Mild	4	2	0.43
Cool	3	1	0.28
p	0.64	0.36	

1. By this table we have $p(\text{Mild}|\text{Yes}) = 4/9 = 0.44$ and $p(\text{Mild}|\text{No}) = 2/5 = 0.4$
2. Posterior $p(\text{Yes}|\text{Mild}) = (p(\text{Mild}|\text{Yes}) \times p(\text{Yes})) / p(\text{Mild}) = (0.44 \times 0.64) / 0.43 = 0.65$
3. Posterior $p(\text{No}|\text{Mild}) = (p(\text{Mild}|\text{No}) \times p(\text{No})) / p(\text{Mild}) = (0.4 \times 0.36) / 0.43 = 0.33$
4. As $p(\text{Yes}|\text{Mild}) > p(\text{No}|\text{Mild})$, it is likely to play.

Q2: Check [INT104 NOTE \(zhtc.one\)](#) Chapter-3 to see multiply conditions Naïve Bayes Classifier by using R

7 Non-parametric Classification

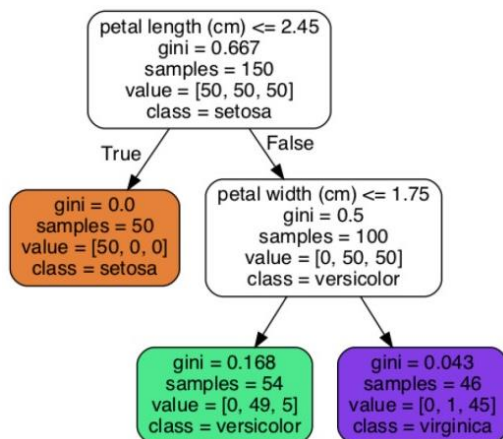
Parametric Methods limitations:

- We seldom obtain a “true” model due to the lack of prior knowledge
- For the same reason, we never know which model should be used

Non-parametric Methods merit:

- Can be used with arbitrary distributions and without the assumption that the forms of the underlying densities are known
- Can be used with multimodal distributions which are much more common in practice than unimodal distributions

7.1 Decision Tree



Content:

- Root
- Leaf
- Node
- Depth

7.1.1 Impurity

Gini is used to measure the impurity of node:

$$G_i = 1 - \sum_{k=1}^n p_{i,k}^2$$

($p_{i,k}$ is the ratio of class k instances among the training instances in the i -th node)

Entropy is also one of the possible ways to evaluate impurity:

$$H_i = - \sum_{k=1}^n p_{i,k} \log_2(p_{i,k})$$

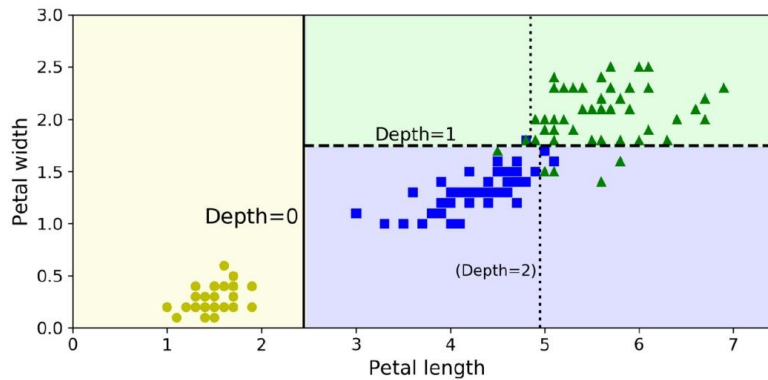
($p_{i,k}$ is the ratio of class k instances among the training instances in the i -th node ($p_{i,k}$ cannot be 0))

Question: A node applies to 0 Iris setosa, 49 Iris versicolor and 5 Iris virginica

$$G_i = 1 - (49/54)^2 = 0.1766118$$

$$H_i = -(49/54) * \log_2(49/54) = 0.1271982$$

7.1.2 Decision Boundaries



<White box method>

7.1.3 Instability

Decision tree is sensitive to both training set rotation and the samples in training set:

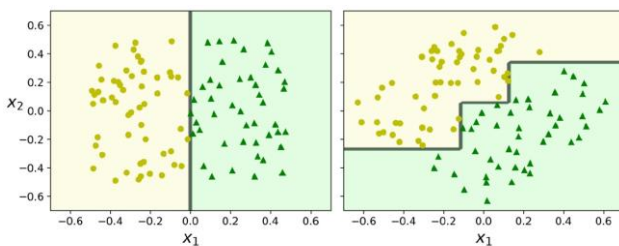


Figure 6-7. Sensitivity to training set rotation

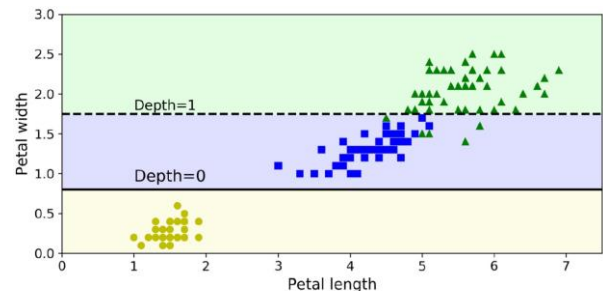


Figure 6-8. Sensitivity to training set details

7.1.4 Growing a tree

Classification and Regression Tree (CART)

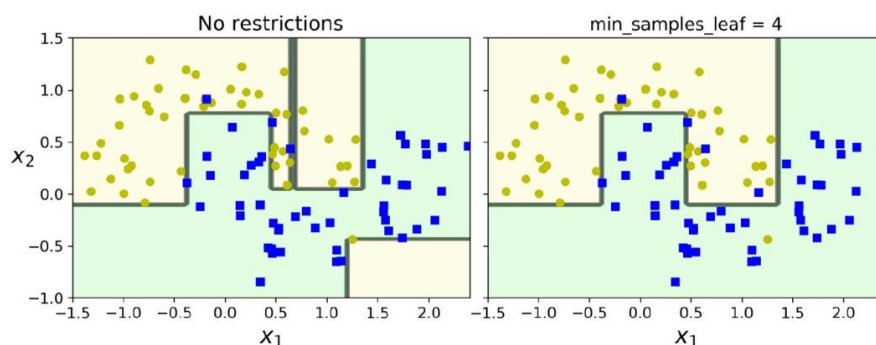
- Splitting a node into two subsets via single feature k and a threshold t_k

$$J(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}$$

- $G_{\text{left/right}}$ measures the impurity of the left/right subset
- $m_{\text{left/right}}$ is the number of instances in the left/right subset

Question:

- When the tree should stop grow?
- What problem is introduced if the tree grows until perfect impurity?
- If we allow the tree to grow too much, it may become too complex and memorize the training data instead of generalizing to new data. To avoid this, we can use a stopping criterion that tells the tree when to stop growing. Some commonly used to stop criteria include: maximum depth, minimum number of samples per leaf node, minimum impurity decrease, and maximum number of leaf nodes.
- If we let the tree grow until it perfectly separates the training data into groups with the same class, it will likely overfit the training data and perform poorly on new data. This is because the tree is too specific to the training data and cannot generalize well to new data. Therefore, it is important to use an appropriate stopping criterion to balance model complexity and generalization performance.



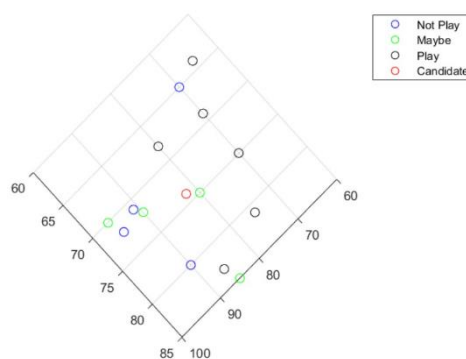
7.2 k-Nearest Neighbour (KNN)

1. **Load the data:** First, we need to load the training dataset into memory.
2. **Define the value of K:** We need to choose a value of K, which represents the number of nearest neighbors to consider when making a prediction.
3. **Calculate distance:** For each new instance to be classified, we calculate its distance to all instances in the training dataset using a distance metric such as Euclidean distance or Manhattan distance.
4. **Select the K-nearest neighbors:** We then select the K instances in the training dataset that are closest to the new instance based on the calculated distances.
5. **Assign the label:** We assign the most common class label among the K-nearest neighbors to the new instance.
6. **Repeat the process:** We repeat this process for each new instance we want to classify.
7. **Evaluate the model:** Finally, we evaluate the performance of the KNN model by testing it on a separate testing dataset and calculating metrics such as accuracy, precision, recall, and F1 score.

Example:

outlook	temperature	humidity	windy	play
sunny	85	85	1.0	maybe
sunny	80	90	3.0	no
overcast	83	86	0.8	yes
rainy	70	96	0.2	maybe
rainy	68	80	0.1	yes
rainy	65	70	2.8	no
overcast	64	65	2.6	yes
sunny	72	95	0.3	no
sunny	69	70	0.5	yes
rainy	75	80	0.4	maybe
sunny	75	70	2.2	yes
overcast	72	90	2.4	maybe
overcast	81	75	0.0	yes
rainy	71	91	2.9	no

Q1: Will we play golf in the case of 74 (temperature), 74 (humidity)? What value of k should be selected?

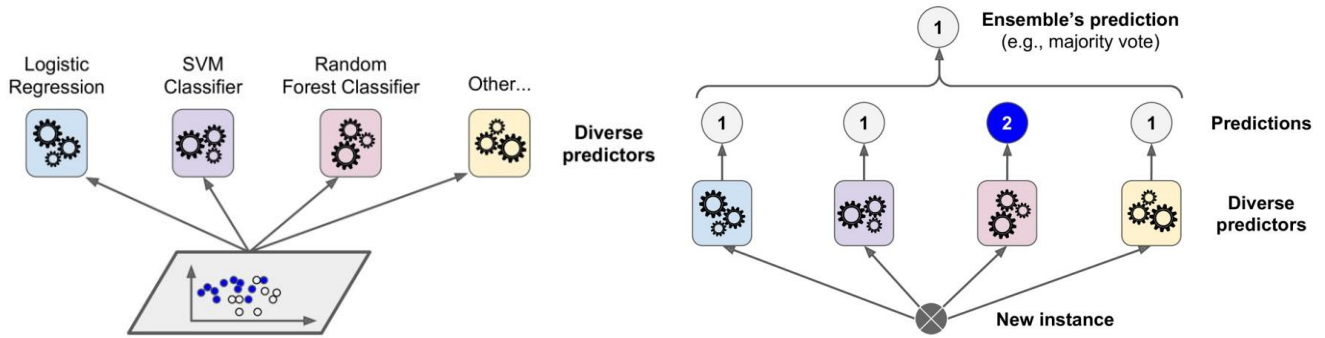


K=3

Q2: Check [INT104 NOTE \(zhtc.one\)](#) Chapter-1 to see the example of KNN by using R

8 Ensemble Learning & Random Forests

8.1 Voting Classifiers



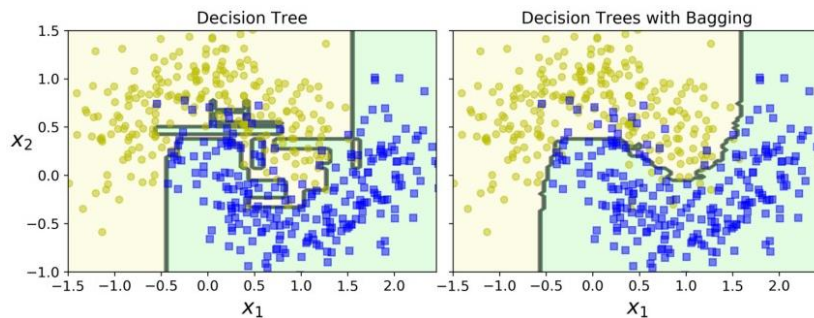
8.1.1 Ways to merge multiple classifiers:

Majority Voting, Weighted Voting, Bagging (Bootstrap Aggregating), Boosting, Stacking

8.1.2 Bagging and Pasting

Training with different subsets of data:

- **Bagging**: Sampling with Replacement. Since bootstrap sampling allows multiple instances to be selected multiple times in each subset, some instances may be included in multiple subsets, while others **may be left out**.
- **Pasting**: Sampling without Replacement. Each classifier is trained **on a different subset**, and their predictions are aggregated using averaging or voting.



8.2 Random Forest

8.2.1 Concepts

- Multiple decision trees with different samples and features can be grown
- Final decision could be made by voting

Random Forest achieves better performance than a single decision tree because it reduces overfitting, handles complex nonlinear relationships, handles high-dimensional data, and is robust to noisy data. It does this by building multiple trees on different subsets of the data and features and combining their predictions.

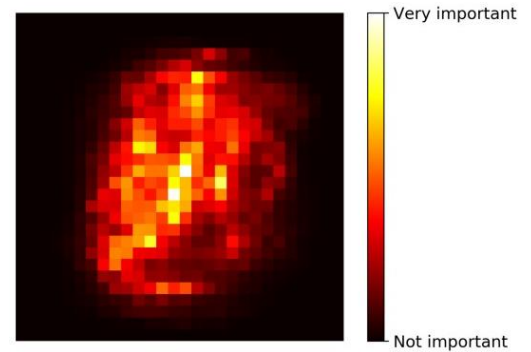
8.2.2 Random Patches & Random Subspaces

- **Random Subspace**: randomly select features for training
- **Random Patches**: randomly select features and randomly select samples

8.2.3 Feature Importance

Some common methods for measuring feature importance:

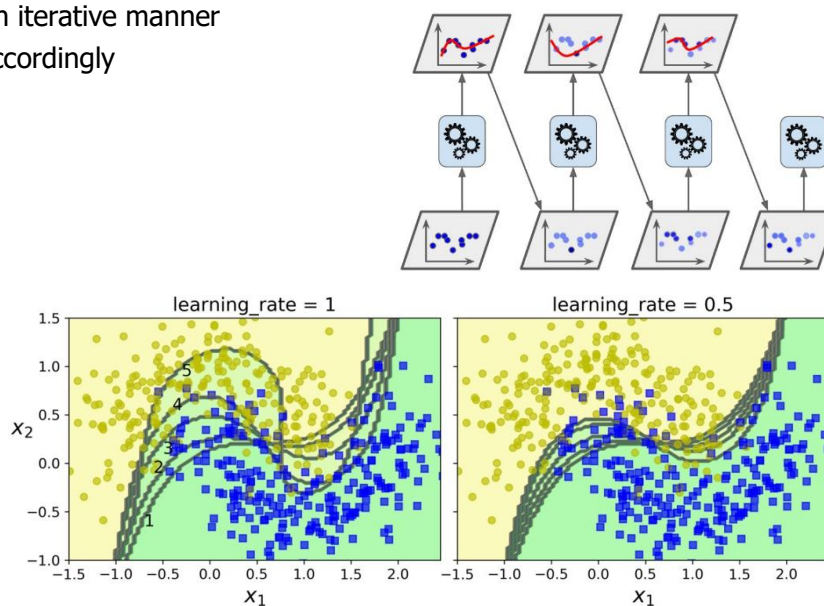
- Gini importance: For decision trees and random forests, the Gini importance is a measure of how often a feature is used to split the data across all trees in the ensemble. Features that are used more frequently to split the data have a higher importance.
- Permutation importance: This method measures the decrease in the model's performance when a feature's values are randomly permuted. Features that lead to a larger decrease in performance when permuted have a higher importance.
- Coefficient magnitude: For linear models, the magnitude of the coefficients reflects the contribution of each feature to the predicted outcome. Larger coefficients indicate higher feature importance.
- Recursive feature elimination: This method involves recursively removing features from the dataset and fitting a model on the remaining features until the optimal subset of features is found. The importance of a feature can be inferred from the order in which it is eliminated from the dataset.
- LASSO regularization: LASSO regularization can shrink some coefficients to zero, effectively removing some features from the model. The magnitude of the nonzero coefficients reflects the importance of each feature.



8.3 Boosting

8.3.1 AdaBoost (迭代推进算法)

- Weight samples in an iterative manner
- Train the classifier accordingly



8.3.2 Gradient Boosting (梯度推进算法)

- Residue is learned after each iteration

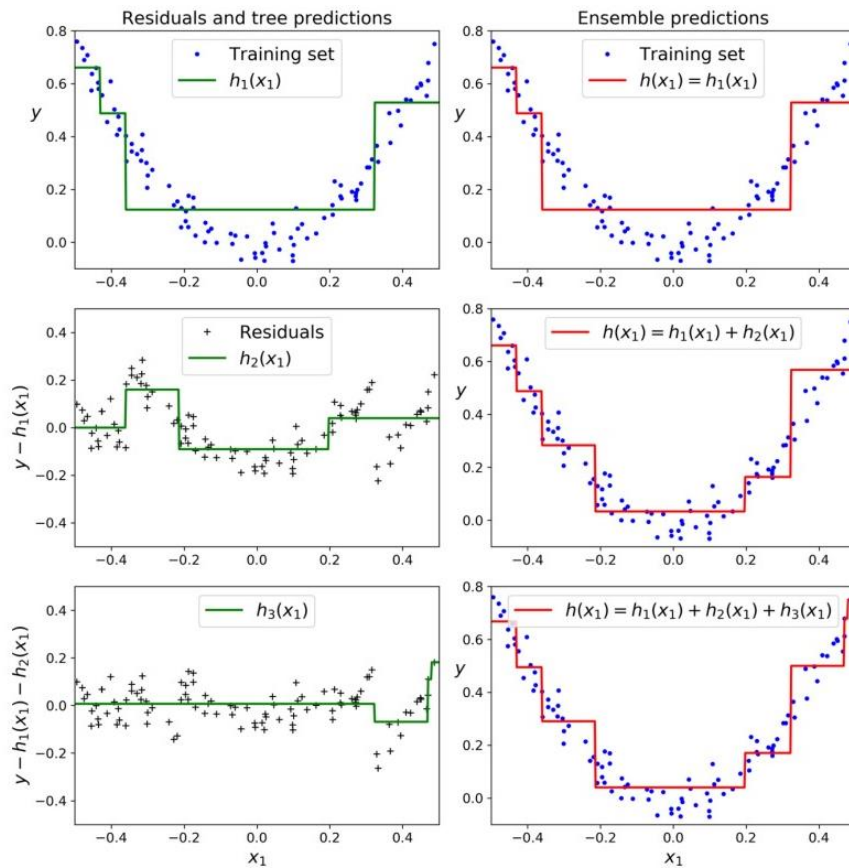


Figure 7-9. In this depiction of Gradient Boosting, the first predictor (top left) is trained normally, then each consecutive predictor (middle left and lower left) is trained on the previous predictor's residuals; the right column shows the resulting ensemble's predictions

8.3.3 Prevent Overfitting

By using early stop:

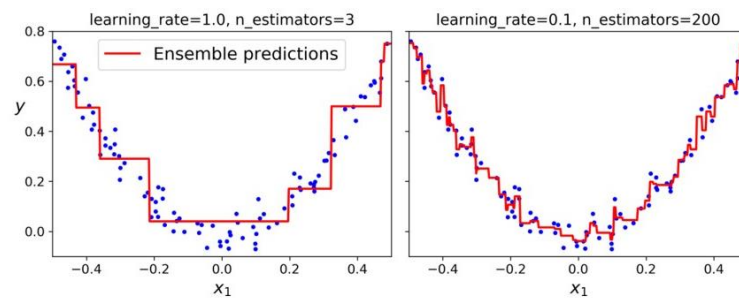


Figure 7-10. GBRT ensembles with not enough predictors (left) and too many (right)

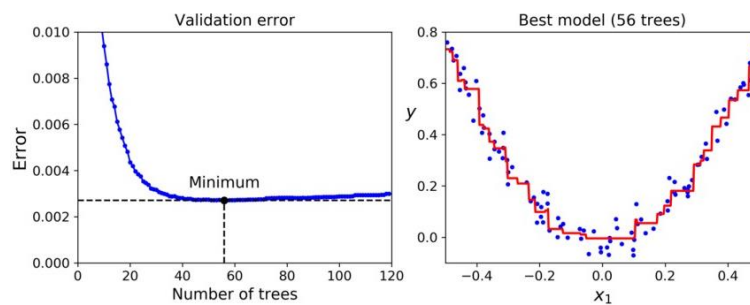
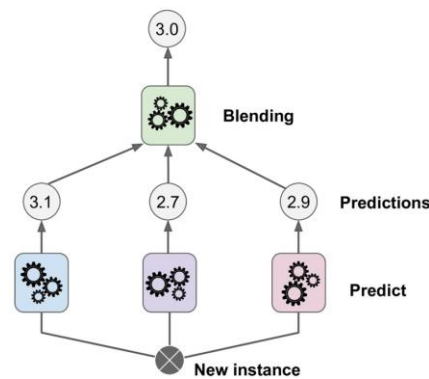


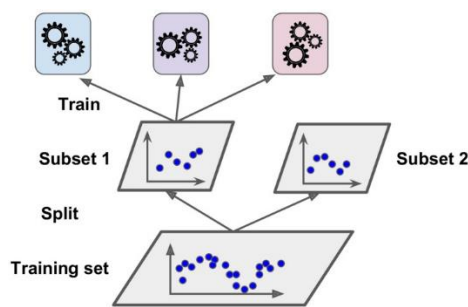
Figure 7-11. Tuning the number of trees using early stopping

8.3.4 Stacking

Combine the predictions of multiple base models by training a meta-model on the outputs of the base models:



1. Training the First Layer



[1] Split the data into training and testing sets.

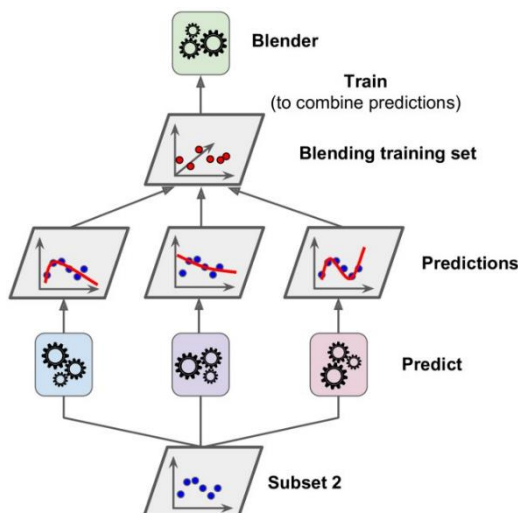
[2] Define the base models to use. These can be any machine learning models, such as decision trees, SVMs, neural networks, etc.

[3] Train each of the base models on the training data and record their predictions on the testing data.

[4] Use the predicted values as features for the meta-model.

[5] Train the meta-model on the predicted values and the actual target values to predict the final output.

2. Training the Blender



[1] Split the original training data into two parts: training and hold-out set.

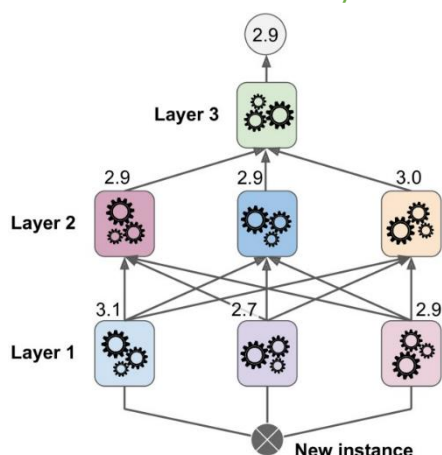
[2] Define the base models and train them on the training set.

[3] Use the trained base models to make predictions on the hold-out set.

[4] Concatenate the predictions made by the base models to create a new dataset.

[5] Train the blender on the new dataset and the target variable (or labels) of the hold-out set.

3. Predictions in a Multilayer Stacking Ensemble



[1] Split the original training data into three parts: training, validation, and hold-out set.

[2] Define the base models and train them on the training set.

[3] Use the trained base models to make predictions on the validation set.

[4] Concatenate the predictions made by the base models to create a new dataset, which becomes the input to the next layer.

[5] Define the second layer of models and train them on the predictions made by the base models in the first layer using the validation set.

[6] Use the trained second layer models to make predictions on the hold-out set.

[7] Concatenate the predictions made by the second layer models with the predictions made by the base models to create a new dataset, which becomes the input to the blender.

[8] Train the blender on new dataset and the target variable (or labels) of the hold-out set.

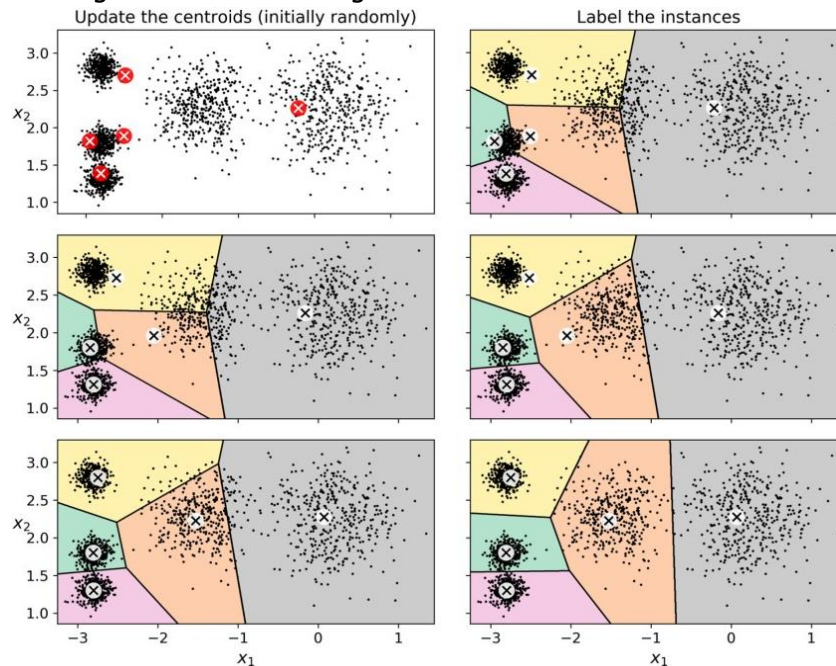
[9] Use the trained blender to make predictions on new data by passing the data through the base models to create the predicted features, and then passing those features through the second layer models to create a new set of features, and then finally passing those features through the blender to make the final prediction.

9 Unsupervised Learning

9.1 K-Means

9.1.1 Normal steps

- 1 The number of clusters k should be defined, where k is a positive integer
- 2 In the beginning, we determine the number of clusters (k) that we want and we assume the centroid or centre of these clusters
- 3 Then repeat the following steps until converge
 - Assign each training sample to the cluster with the nearest centroid
 - Calculate the new centroid for each cluster
- 4 Stop until no samples changes the cluster belonged after an iteration

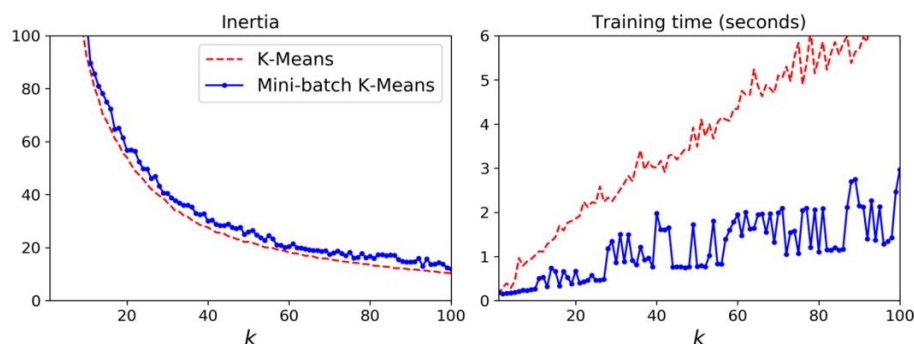


Define initial centroids:

- **Random selection**: Centroids are chosen randomly from the dataset, but this may not always produce good results.
- **K-Means++**: A variation of K-Means that selects centroids based on their distance from each other to improve results.
- **Manual selection**: Centroids are chosen manually based on prior knowledge of the dataset.
- **Hierarchical clustering**: Centroids are chosen using a hierarchical clustering algorithm that groups similar data points into clusters and then uses those clusters as initial centroids.
- **Density-based clustering**: Centroids are not explicitly defined in density-based clustering algorithms like DBSCAN, but instead the algorithm starts with a random point and expands the cluster by adding nearby points until a dense region is formed.

9.1.2 Mini-batch K-Means

Mini-batch K-Means has a higher inertia (惰性) than K-Means but it is much faster, especially as k increase:



How to choose k ?

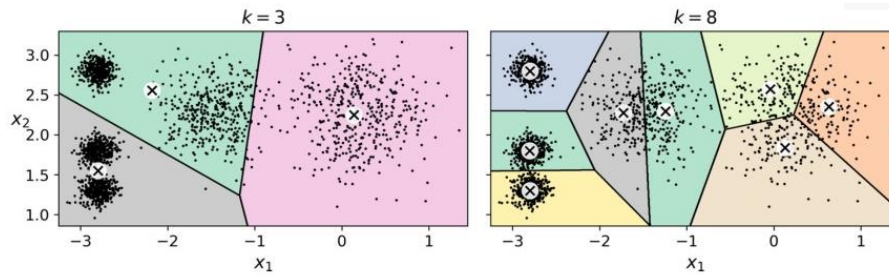


Figure 9-7. Bad choices for the number of clusters: when k is too small, separate clusters get merged (left), and when k is too large, some clusters get chopped into multiple pieces (right)

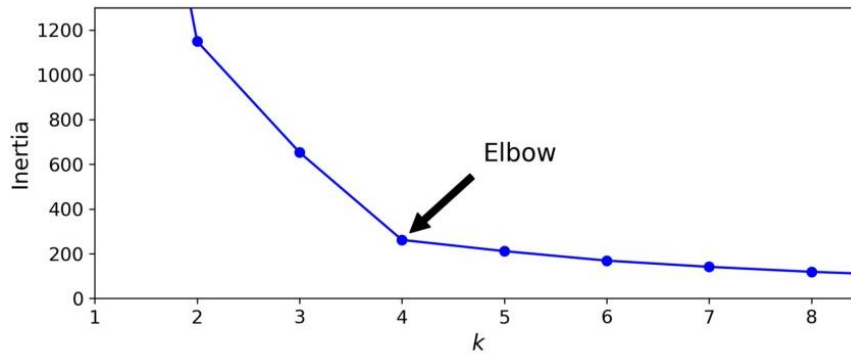
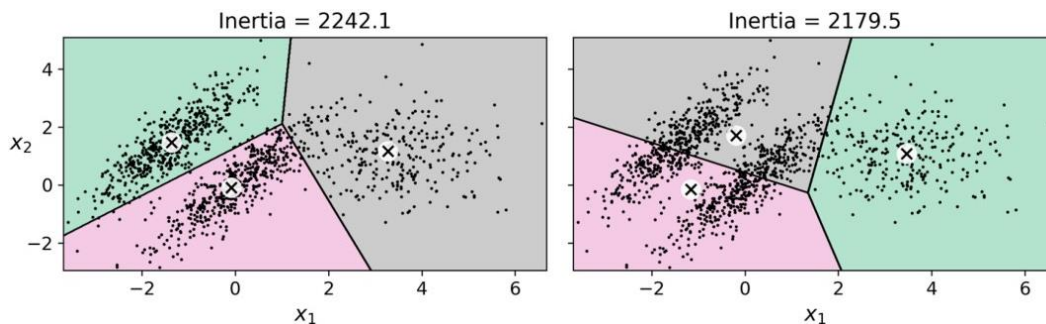


Figure 9-8. When plotting the inertia as a function of the number of clusters k , the curve often contains an inflexion point called the “elbow”

System Evaluation: by Silhouette score & Silhouette coefficient

K-Means suffers from imbalanced clusters, like fail to cluster these ellipsoidal blobs properly:



Two measures can be used to evaluate the quality of clustering results:

Adjusted Rand Index (ARI) and Normalized Mutual Information (NMI).

- ARI measures similarity between true and clustering results, with 1 indicating perfect agreement, 0 indicating random labelling, and negative values worse than random. If K-Means' ARI score is low for non-spherical or elongated data, it means the clusters are imbalanced.
- NMI also measures similarity between true and clustering results, but takes into account the size of clusters. If K-Means' NMI score is low for non-spherical or elongated data, it means the clusters are imbalanced.

9.1.3 Pre-processing with K-Means

Training Process:

1. K-Means Clustering
2. Replace Sample with Distances to Centroids (质心)
3. Train Classifier

Why such pre-processing method improves system performance?

- Reducing dimensionality: Instead of using the original features of the data, the distance to centroids reduces the dimensionality of the data. This can help to reduce noise and improve the signal-to-noise ratio, making it easier for the classifier to distinguish between different classes.
- Capturing underlying structure: The distances to centroids can capture the underlying structure of the

data more effectively than using the original features. This is because the centroids represent the centers of the clusters, and the distances to these centroids can provide information about how each sample is related to the different clusters.

- Improving clustering performance: K-Means clustering can be sensitive to the choice of initial centroids, and using the distances to centroids as inputs can improve the clustering performance by providing a better initialization of the centroids.
- Reducing computation: Using the distances to centroids instead of the original features can reduce the computational burden of the classifier, since there are fewer features to process. This can result in faster training and inference times, as well as lower memory usage.

9.1.4 Semi-Supervised Learning

Training Process:

1. Clustering with K-Means
2. Using Centroids to Approximate Samples
3. Train Classifier

9.2 DBSCAN

Training Process:

1. Count instances located within ϵ , which is called the instance's ϵ -neighbourhood
2. A instance that has a certain number of instances in ϵ -neighbourhood is considered as a core instance
3. All instances in the neighbourhood of a core instance belong to the same cluster
4. Any instance that is not a core instance and does not have one in its neighbourhood is considered as an anomaly

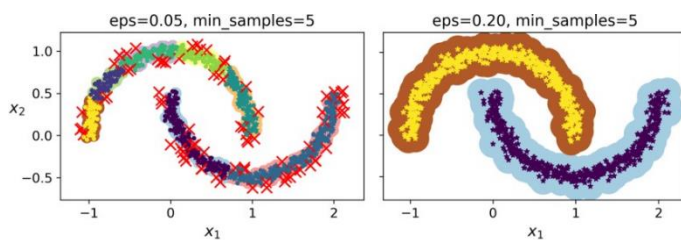


Figure 9-14. DBSCAN clustering using two different neighborhood radiuses

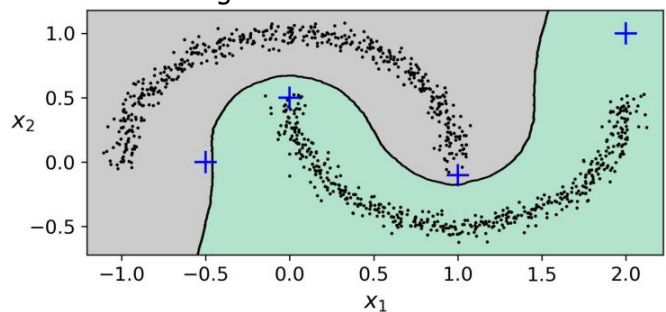


Figure 9-15. Decision boundary between two clusters

9.3 Hierarchical Clustering

9.3.1 Agglomerative Clustering

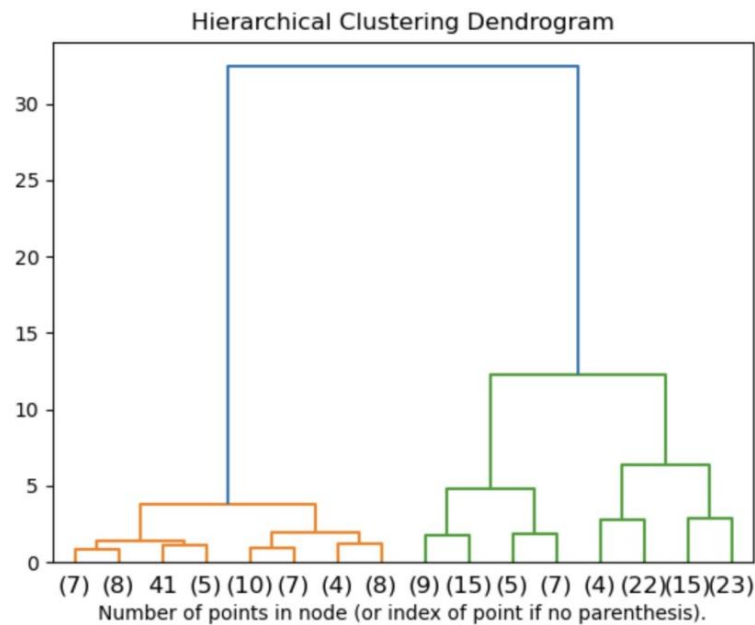
Training Process:

- 1 For n objects v_1, \dots, v_n , assign each to a singleton cluster $C_i = \{v_i\}$
- 2 Use any computable cluster similarity measure $D(C_i, C_j)$ e.g. Euclidean distance, cosine distance etc.
- 3 Repeat{
 - identify the two most similar clusters C_j and C_k (could be ties -> choose one pair)
 - delete C_j and C_k and add $(C_j \cup C_k)$ to the set of clusters
 } until just one cluster remains
- 4 Use a dendrogram diagram to show the sequence of cluster mergers

Some multiple strategies for distance measurement between joint clusters (D):

- single linkage: D is taken as the minimum distance between samples in sub-clusters
- complete linkage: D is taken as the maximum distance between samples in sub-clusters
- average linkage: D is taken as the average distance between each pair of samples in sub-clusters

Dendrogram:



10 Gaussian Mixture Model & Bayesian Statistics

10.1 Gaussian Mixture Model (GMM)

Mixture model could be used for clustering as well:

$$p(\mathbf{x}) = \sum_k \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)$$

- \mathcal{N} represents Gaussian distribution (known as Gaussian component in GMM)
- π is the weight of Gaussian component
- μ is the mean of Gaussian component
- Σ is the co-variance matrix of Gaussian component
- Each Gaussian component corresponds to a cluster

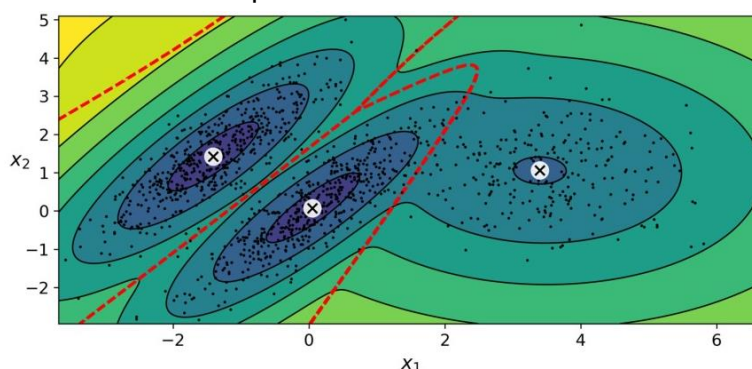
10.1.1 Expectation Maximisation

When there are missing values exist among the data or the model can be formulated more simply by assuming the existence of further unobserved data points, **Expectation Maximisation (EM)** can be used:

- **Expectation step:** find the posterior probability according to current model
- **Maximisation step:** calculate the new modal parameters

And there are several methods for determining the number of Gaussian components, including:

- **Visual inspection:** Plotting the data distribution and visually inspecting it can provide insights into the number of underlying components.
- **Information criteria:** Information criteria such as the Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC) can be used to compare models with different numbers of components and select the one with the best balance of fit and complexity.
- **Cross-validation:** Cross-validation can be used to evaluate the performance of models with different numbers of components and select the one with the best generalization performance.
- **Domain knowledge:** In some cases, domain knowledge or prior information about the data can be used to guide the selection of the number of components.



10.1.2 Central Limit Theorem

The sampling distribution of a sample mean is approximately normal if the sample size is large enough, even if the population distribution is not normal.

10.1.3 Covariance Matrix

- Covariance matrices shared across all Gaussian components
- Covariance matrices must be diagonal

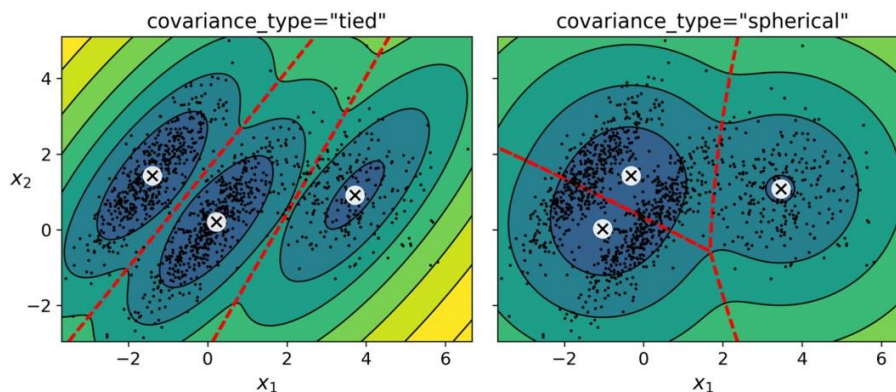


Figure 9-18. Gaussian mixtures for tied clusters (left) and spherical clusters (right)

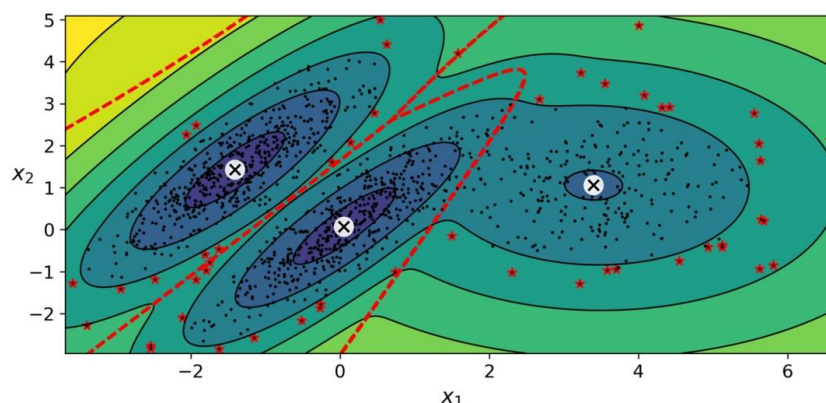
10.1.4 Anomaly Detection

How to **detect the outlier**?

One common approach is to use the Mahalanobis distance, which is a measure of the distance between a data point and the mean of the mixture component it is assigned to. A data point with a high Mahalanobis distance is considered an outlier because it is far from the center of the component and has a low probability of belonging to that component.

To detect outliers with GMMs using the Mahalanobis distance, the following steps can be taken:

1. Fit a GMM to the data using an appropriate number of mixture components.
2. Compute the Mahalanobis distance for each data point with respect to the mean and covariance matrix of the component it is assigned to.
3. Threshold the Mahalanobis distances using a predetermined cutoff value or a statistical test, such as the Q-Q plot or the Kolmogorov-Smirnov test.
4. Identify the data points with Mahalanobis distances above the threshold as outliers.



10.2 Model Selection

10.2.1 Concepts

In the context of comparing different models, the likelihood can be used as a basis for model selection. One approach is to compute the likelihood of the data under each model and compare them using statistical tests or information criteria such as the **Akaike Information Criterion (AIC)** or the **Bayesian Information Criterion (BIC)**. These criteria provide a way to balance the goodness-of-fit of the model to the data with the complexity of the model, by penalizing models with more parameters.

It's important to note that likelihood-based comparisons assume that the models being compared have the same form and differ only in their parameter values. If the models have different forms, such as different numbers of components in a mixture model, then additional methods may be needed to compare them, such as the [Bayesian model selection](#) or [cross-validation](#) techniques.

10.2.2 Formulas

Akaike Information Criterion:

$$AIC = 2p - 2 \log \hat{L}$$

Bayesian Information Criterion:

$$BIC = p \log(m) - 2 \log \hat{L}$$

- m is the number of instances
- p is the number of parameters learned by the model
- \hat{L} is the maximised value of the likelihood function of model

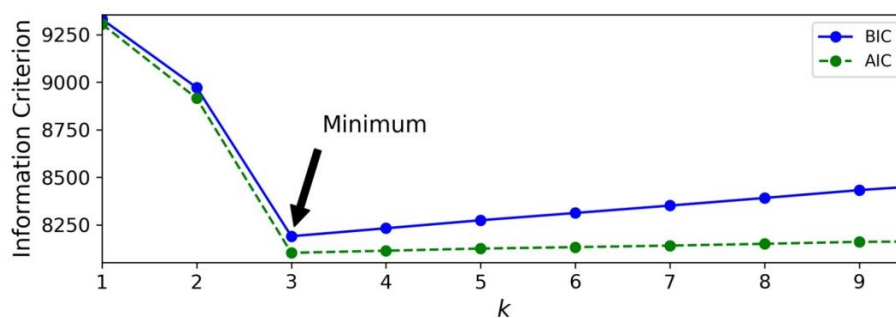


Figure 9-21. AIC and BIC for different numbers of clusters k

10.2.3 Exercise (binary distribution)

Question:

Suppose a set of experiments of rolling a coin is performed, the experiment results obtained are $t_1 = T$, $t_2 = H$, $t_3 = T$, $t_4 = T$, $t_5 = H$, $t_6 = H$, $t_7 = T$, $t_8 = T$, where t_i represents the result of experiment i and T , H means tail up and head up respectively.

- According to the experiments, build a model that assert the probability of each time appearing on the top of a dice. (Hint: in the form of $p(t = T) = p_0$ and $p(t = H) = p_0$)
- As a common knowledge, the dice rolling should be $p(t = T) = 1/2$ and $p(t = H) = 1/2$. Explain why the model resulted from experiments is not the same with the common knowledge.
- Calculate the model likelihood of the common knowledge model ($p(t = T) = 1/2$, $p(t = H) = 1/2$) and the model resulted in question before, given another set of validation experiments where 2 head-up and 1 tail-up obtained.

Answer:

- $p(t = H) = 3/8 = 0.375$, $p(t = T) = 5/8 = 0.625$
- The common knowledge can be taken as prior probability and the experiments can be taken as posterior probability. As there are very few experiments performed, it is very hard to get enough empirical data.
- $L_{\text{prior}} = p_{\text{prior}}(t = T) * p_{\text{prior}}(t = H) * p_{\text{prior}}(t = H) = 0.5 * 0.5 * 0.5 = 0.125$
 $L_{\text{post}} = p_{\text{post}}(t = T) * p_{\text{post}}(t = H) * p_{\text{post}}(t = H) = 0.625 * 0.375 * 0.375 = 0.088$

Usually, a common way to reduce the impact of an individual set of experiments is to combine the prior knowledge and the posterior probability distribution. A typical way is to introduce Bayesian estimation, where an extra attempt is introduced in the set of experiments with an evened-out results (i.e., for all possible n case, each case adds an extra "1" count). Please use a model that merges both prior knowledge and posterior probability to calculate the model likelihood for the given validation set.

- $p_{\text{Bayes}}(t = T) = (3+1)/(8+2) = 0.4$, $p_{\text{Bayes}}(t = H) = (5+1)/(8+2) = 0.6$
- $L_{\text{Bayes}} = p_{\text{Bayes}}(t = T) * p_{\text{Bayes}}(t = H) * p_{\text{Bayes}}(t = H) = 0.4 * 0.6 * 0.6 = 0.144$