

Machine Learning – Take home messages

Week 11

Recap of the unit ...

Machine Learning Models

Unsupervised

Supervised

Clusterin
g

Dimensi-
onality
Reductio
n

Linear

Non-Linear

Linear
regress-
ion

Logistic
regressio
n/classifi
cation

K nearest
neighbor

SVM

Decision
Tree

Random
Forest

Neural
Network

Deep
Learning

Recap of the unit ...

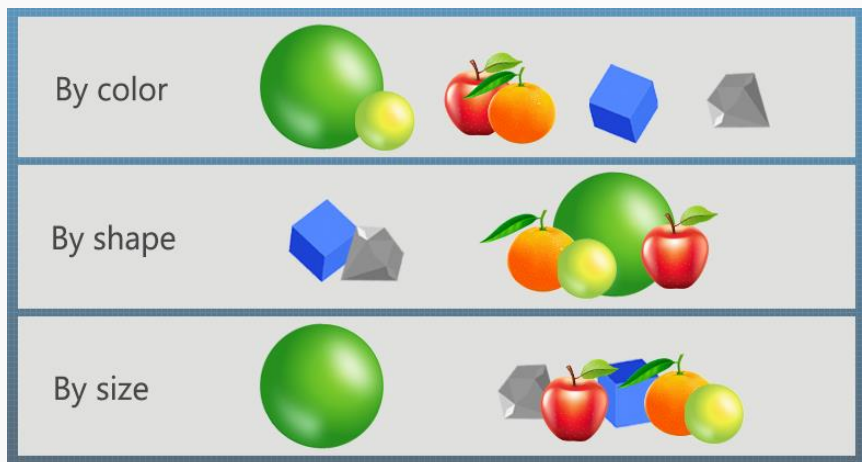
Unsupervised

Clustering

Dimensionality reduction



Cluster can vary based on the properties



Distance measures

Euclidean
Mahalanobis
Jaccard
Cosine
Manhattan

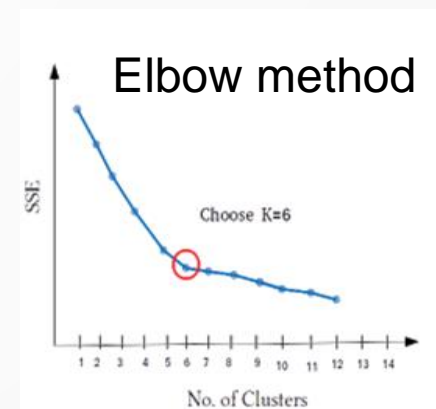
Algorithms

Kmeans++

KMeans

Intra-cluster distances are minimized

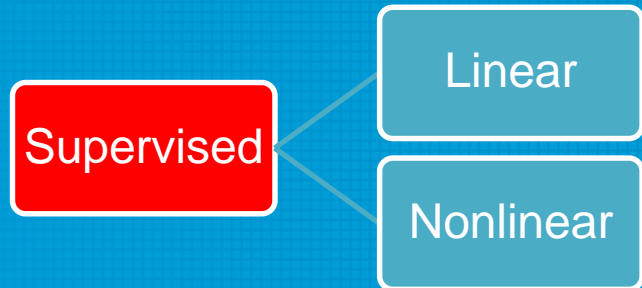
Inter-cluster distances are maximized



Performance measures

Rand Index
Purity
Mutual Information
Silhouette Coefficient

Recap of the unit ...



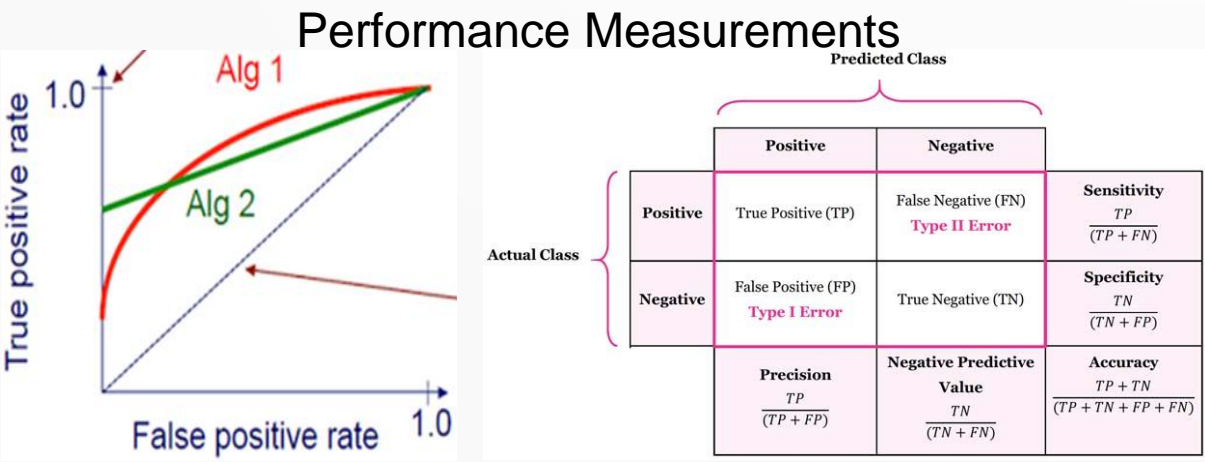
What we want to learn?

$$h : X \rightarrow Y$$



How the machine learns?

$$\min_{h \in H} \frac{1}{n} \sum_{i=1}^n L(y_i, h(x_i))$$



Model complexity

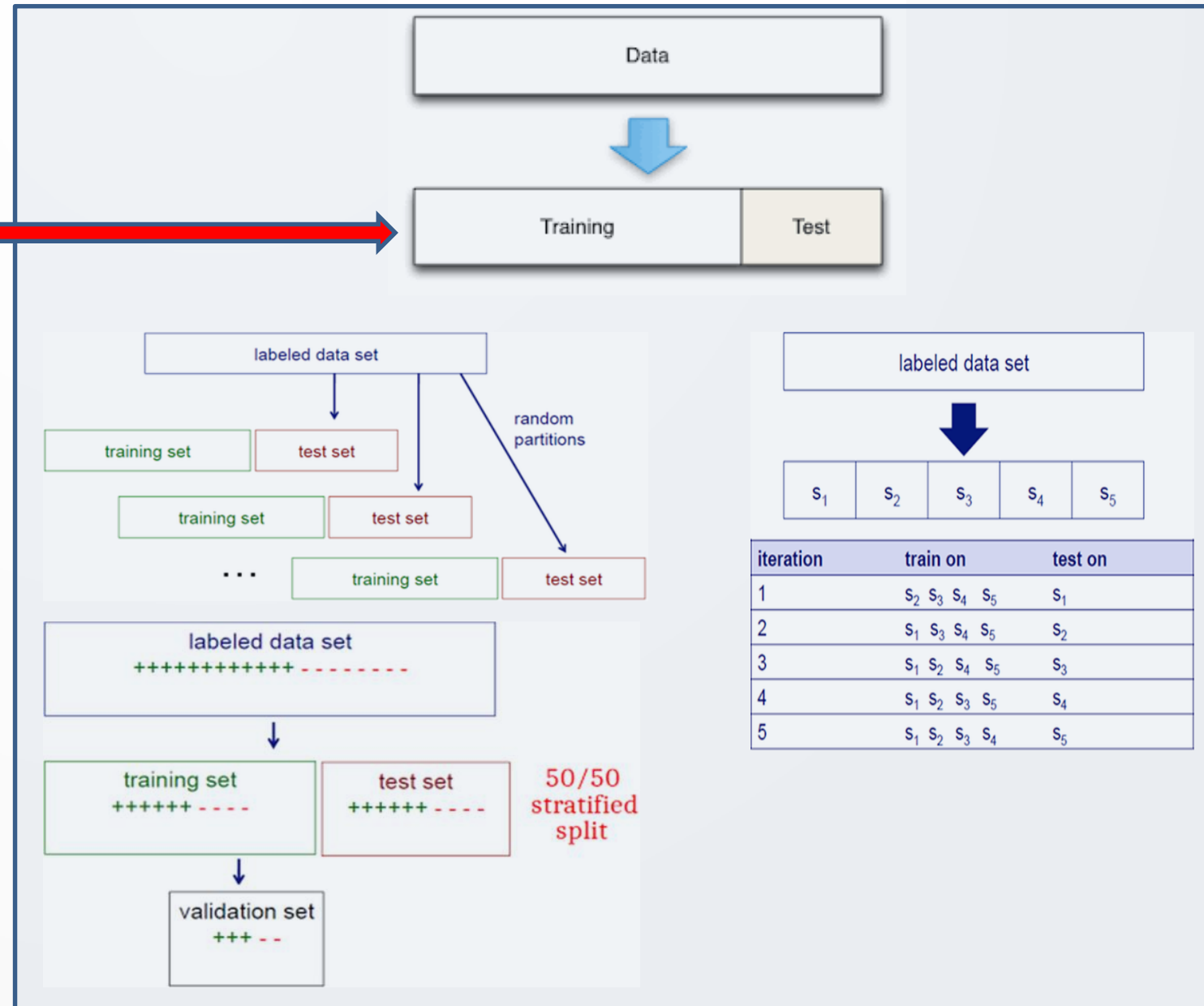
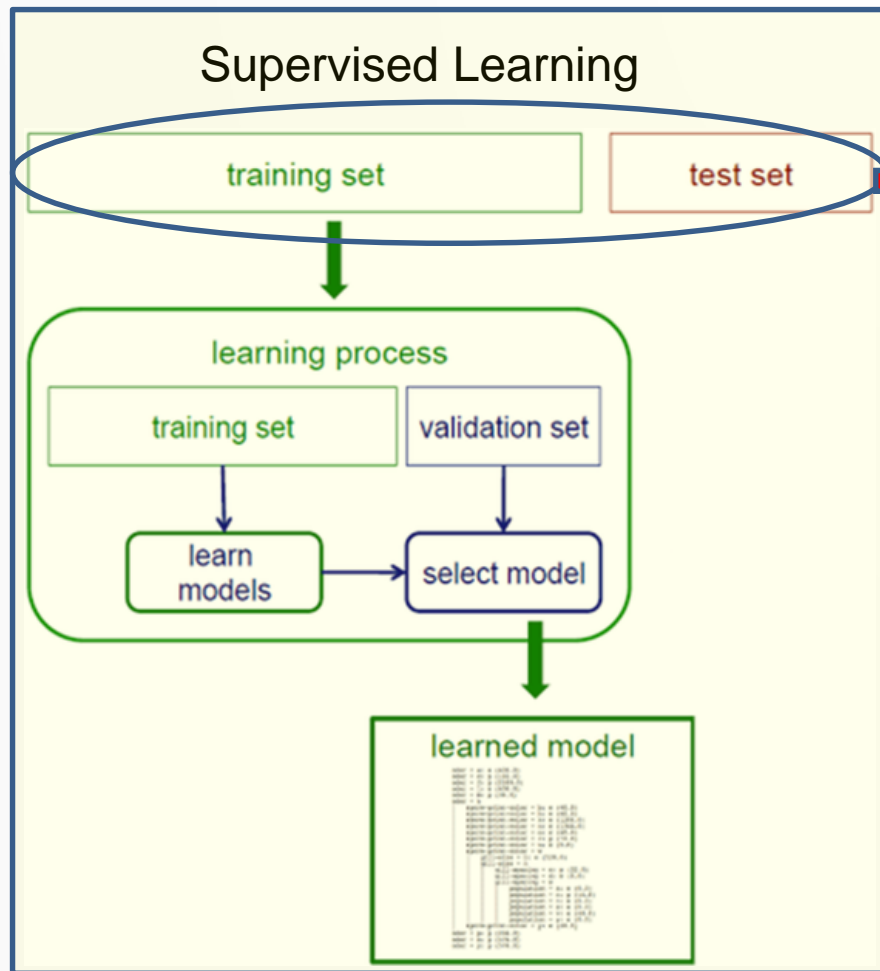
$$R_{str}(h) = R_{emp}(h) + \lambda C(h)$$

Recap of the unit ...

Supervised

Linear

Nonlinear

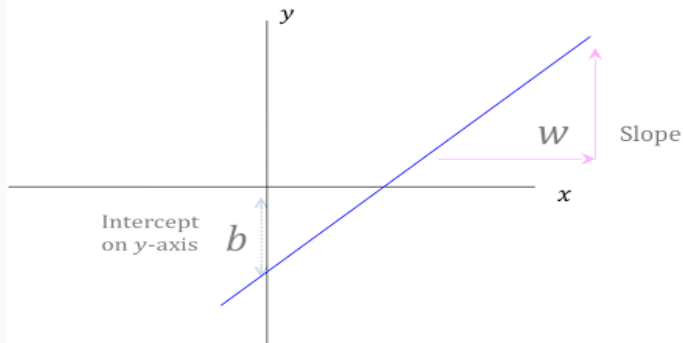


Recap of the unit ...

Supervised

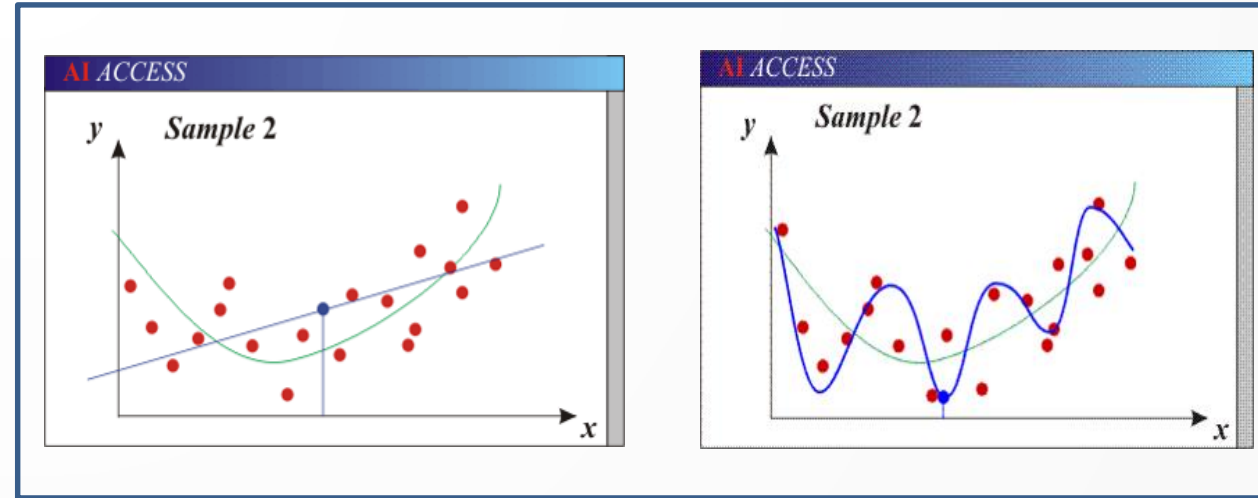
Linear

Nonlinear



$$L(y_i, \mathbf{x}_i^T \mathbf{w}) = (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

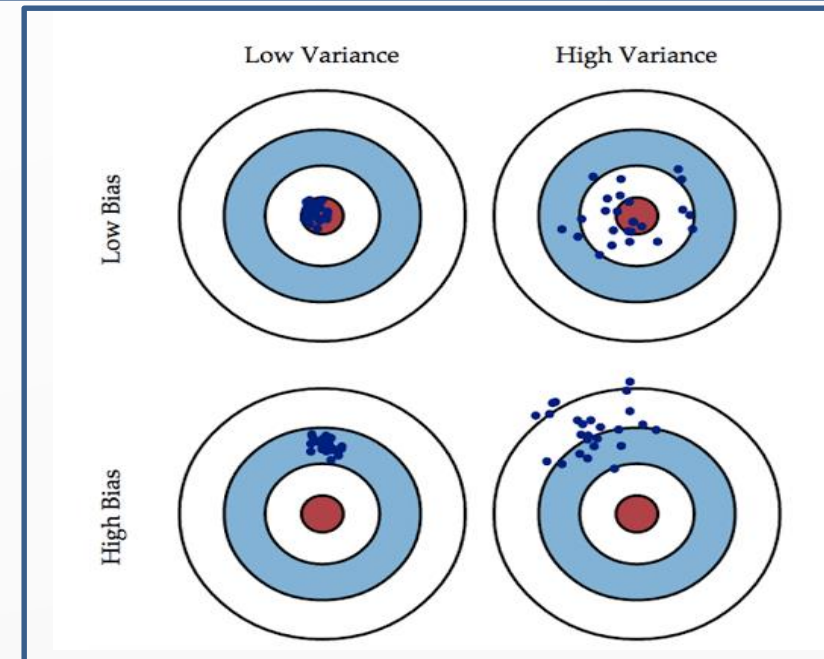
$$L(y_i, \mathbf{x}_i^T \mathbf{w}) = \log(1 + \exp(-y_i \mathbf{x}_i^T \mathbf{w}))$$



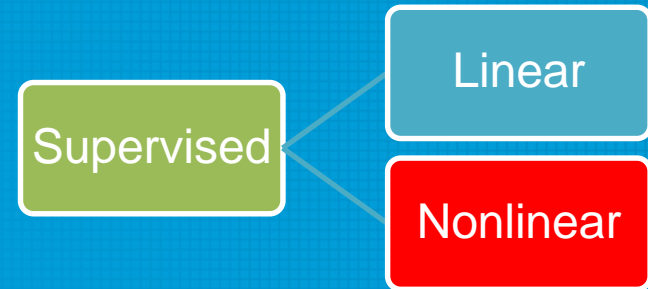
$$\min_w \frac{1}{n} \sum_i L(y_i, \mathbf{x}_i^T \mathbf{w}) + \lambda_1 ||\mathbf{w}||_1$$

$$\min_w \frac{1}{n} \sum_i L(y_i, \mathbf{x}_i^T \mathbf{w}) + \lambda_2 ||\mathbf{w}||_2^2$$

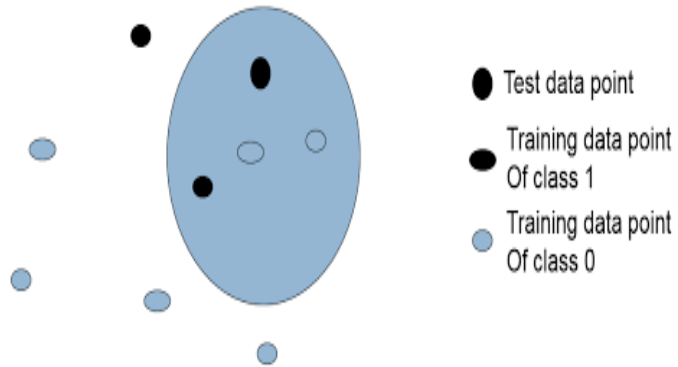
$$\min_w \frac{1}{n} \sum_i L(y_i, \mathbf{x}_i^T \mathbf{w}) + \lambda_1 ||\mathbf{w}||_1 + \lambda_2 ||\mathbf{w}||_2^2$$



Recap of the unit ...

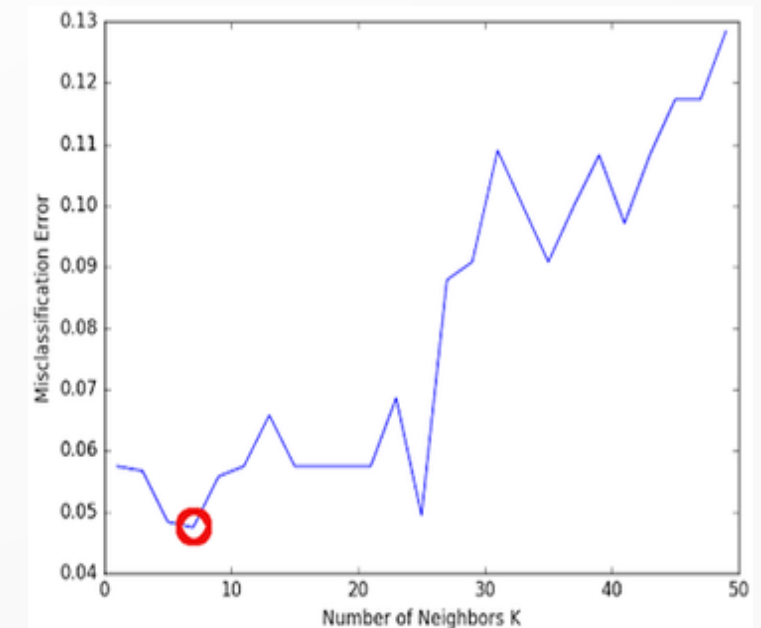


K Nearest Neighbor



- Small values of K
 - Restrains the region of a given prediction
 - Forces classifier to be more **focused on the close regions and neighbours**
 - This will result in a **low bias** and **high variance**
- Higher values of K
 - Asking for more **information from distant training points**
 - Smoother decision boundaries
 - **Lower variance but increases bias**

- Finding the best K
 - There is **no rule of thumb** in selecting K_{\max} since it depends on your desired rate of exploration for K
 - A simple and handy method
 - **Cross-validation** to partition your data into test and training samples
 - **Evaluate model** with different ranges of K values
 - The misclassification error can be used as a measurement of performance



Recap of the unit ...

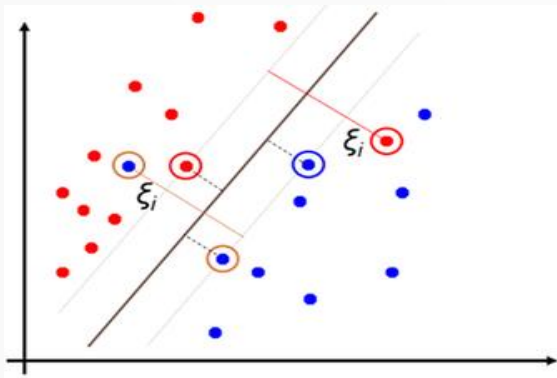
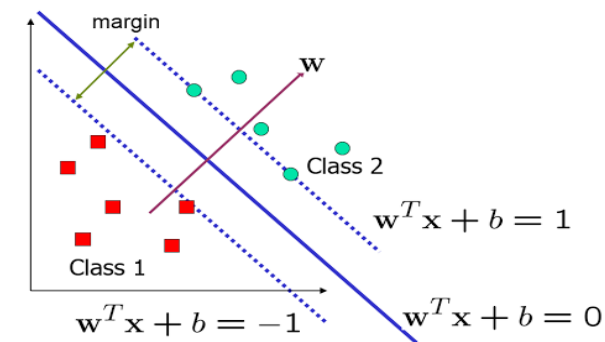
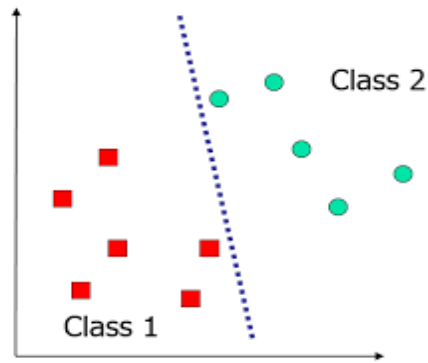
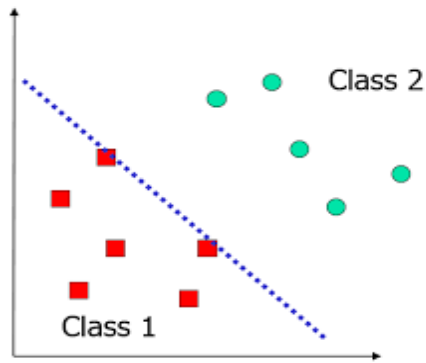
Supervised

Linear

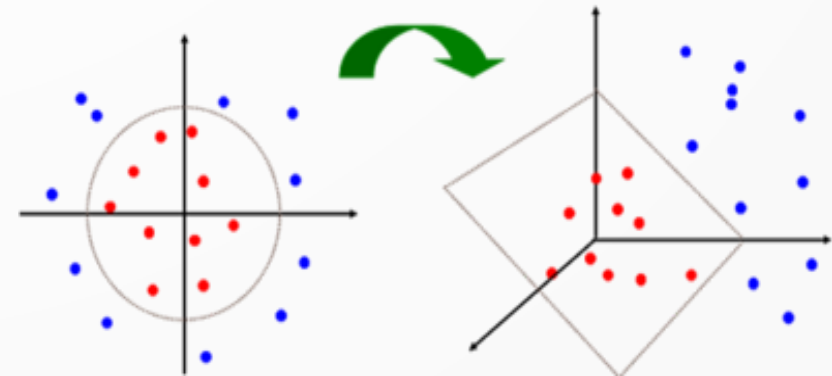
Nonlinear

$$\min_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

subject to: $\sum_{i=1}^n \alpha_i y_i = 0$ and $\alpha_i \geq 0$ for all i



- Nonlinearly separable data
- Kernel



$$\min_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \longleftrightarrow K(X_i, X_j)$$

Subject to : $\sum_{i=1}^n \alpha_i y_i = 0$ and $0 \leq \alpha_i \leq C \quad \forall i$

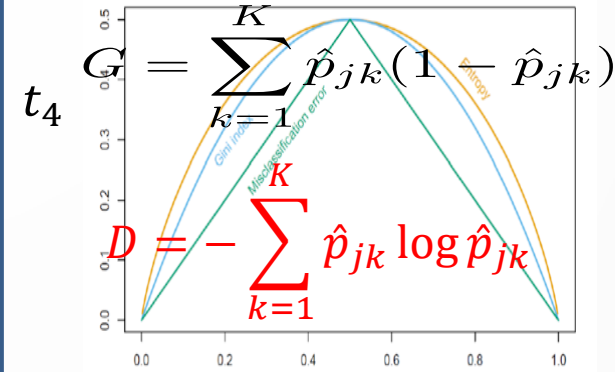
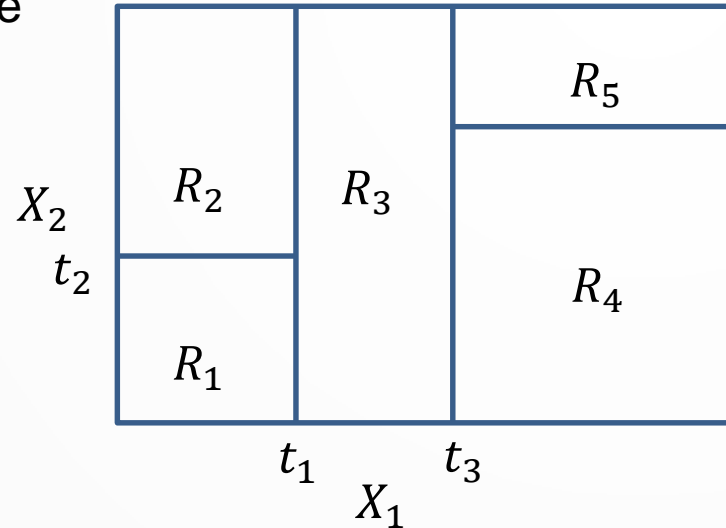
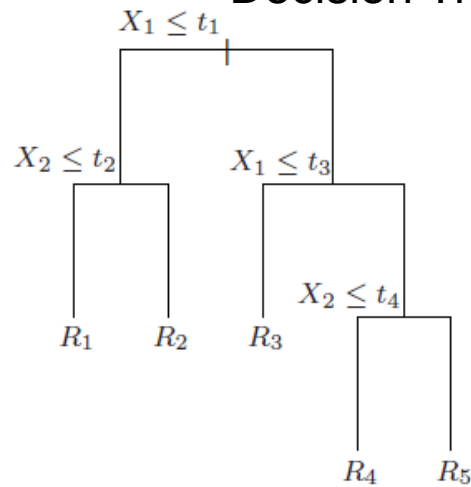
Recap of the unit ...

Supervised

Linear

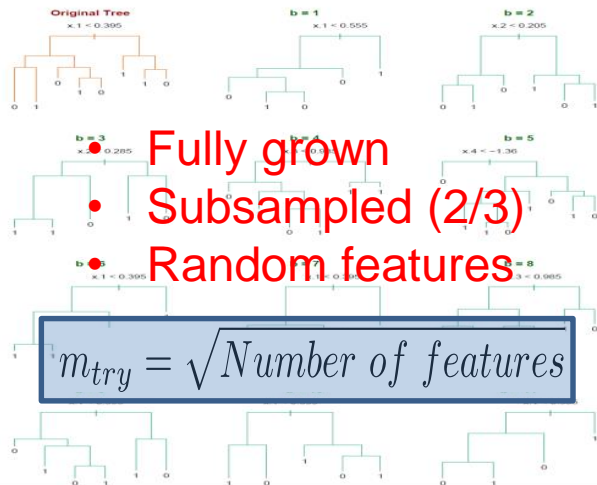
Nonlinear

Decision Tree



- Simple
- Explainable
- **Overfitting**

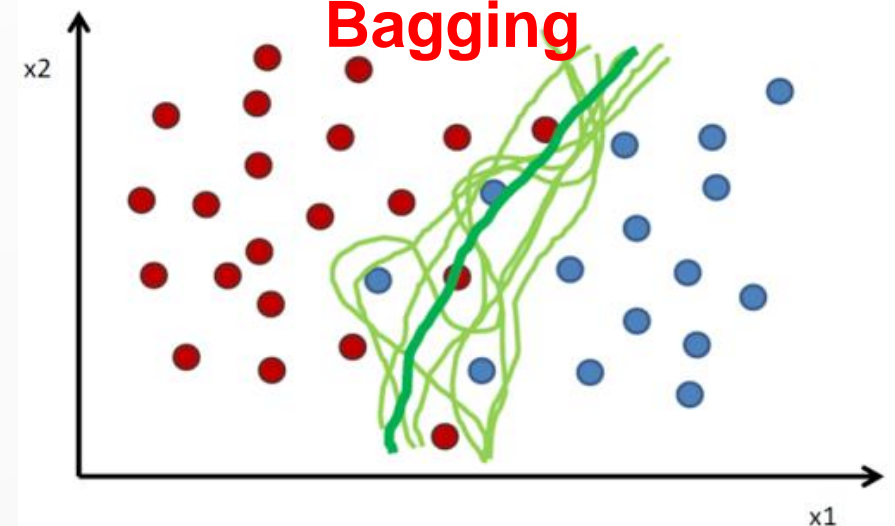
- Pruning
 - Pre-pruning
 - Post-pruning



Power of collective knowledge



Bagging



Recap of the unit ...

Supervised

Linear

Nonlinear

Hyperparameter

Parameter

Performance
measurement

Data splitting

Overfitting

Model
generalisation

Recap of the unit ...

Supervised

Linear

Nonlinear

Design questions

Which ML model to use?

- Sample size
- Class balance
- Dimensionality of the data
- Application environment
- Realtime or asynchronous

How to optimize the model?

- Identify hyperparameters (especially sensitive ones)
- Do not try to use any rule of thumb
- Do not use the test data in the model selection process
- Select the best model based on cross-validation results
- Generate multiple test performance and report the aggregated results

eVALUate

- Completing the survey is an important way that Deakin listens to student voices about their experience (especially during this time)
- Responses are confidential, students are not identified, and eVALUate reports are only sent out after student results are released
- You must only give respectful and polite feedback to better represent your views
 - There is a useful set of guidelines on constructive feedback for eVALUate developed by the University of Tasmania – [here](#).
- You can complete your eVALUate surveys on any device, from anywhere using this link <https://deakin.is/evaluating-us>

Thank You.