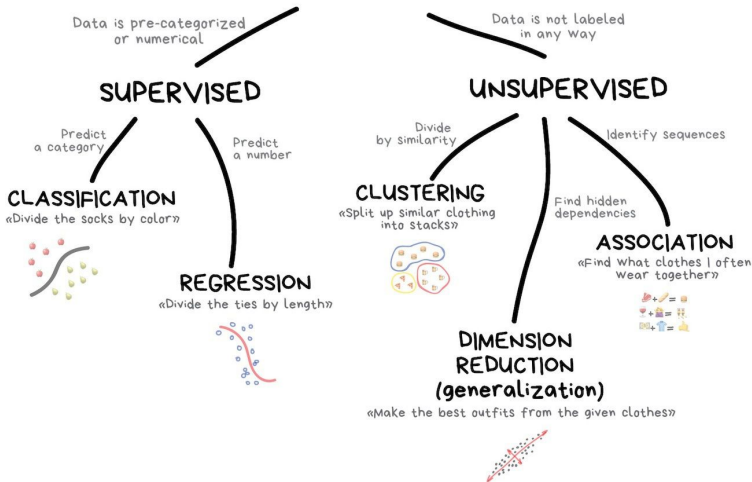


CLASSICAL MACHINE LEARNING



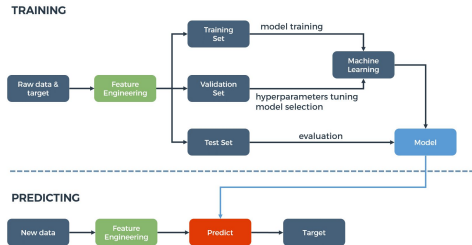
Data

Training set

- ☛ representative set of examples used for training, where the target value is known

Validation set

- ☛ representative set of examples used to tune the architecture of a learning algorithm and estimate prediction errors

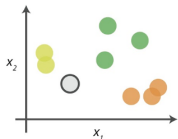


Independent test set (blind test)

- ☛ independently assess the performance of a predictive model
- ☛ never used during the training process
- ☛ error on the blind test provides an unbiased estimate of the generalization error

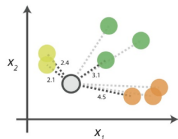
K-nearest neighbours (KNN)

0. Look at the data







Say you want to classify the grey point into a class. Here, there are three potential classes - lime green, green and orange.

1. Calculate distances









Start by calculating the distances between the grey point and all other points.

2. Find neighbours

Point	Distance	
	2.1	→ 1st NN
	2.4	→ 2nd NN
	3.1	→ 3rd NN
	4.5	→ 4th NN

Next, find the nearest neighbours by ranking points by increasing distance. The nearest neighbours (NNs) of the grey point are the ones closest in dataspace.

3. Vote on labels

Class	# of votes	
	2	→ Class  wins the vote! Point  is therefore predicted to be of class  .
	1	
	1	

Vote on the predicted class labels based on the classes of the k nearest neighbours. Here, the labels were predicted based on the $k=3$ nearest neighbours.

- 👑 given training data,
 $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$
and a test point x_u
- 👑 prediction rule:
look at the K most similar training examples to x_u
- 👑 for classification:
assign the majority class label (majority voting)
- 👑 for regression:
assign the average response
- 👑 The algorithm requires
 - parameter K : number of nearest neighbors to look for
 - distance function: to compute similarities between points

Naïve Bayes

Conditional probability:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

Probability of A and B

Probability of A given B

Probability of B

LIKELIHOOD
the probability of "B"
being TRUE given that "A" is TRUE

PRIOR
the probability of
"A" being TRUE

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

POSTERIOR
the probability of "A"
being TRUE given that "B" is TRUE

The probability
of "B" being
TRUE

@luminousmen.com

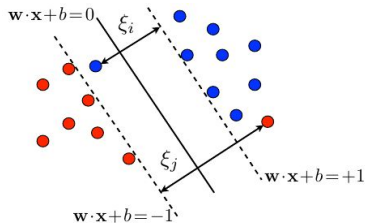
Maximum a posteriori (MAP):

$\arg \max P(A | B)$ for each class A

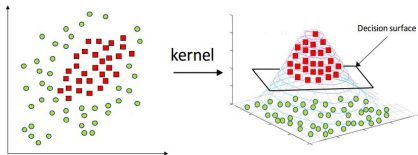
$$\arg \max P(A | B) = \arg \max P(B | A) \times \frac{P(A)}{P(B)}$$
$$= \arg \max P(B | A) \times P(A)$$

- uses Bayes's rule
- $P(B)$ will be the same for all classes
- for a set of features x_i we calculate the joint probability:
 $B = (x_1, x_2, \dots, x_n)$
- $\arg \max P(x_1, x_2, \dots, x_n | A) \times P(A)$
 $P(x_1, x_2, \dots, x_n | A) = \prod_i P(x_i | A)$
- simply the product of individual probabilities
- assumption: features are independent given a class

What if our data is not linearly separable?



- 👑 Use a soft SVM
 - Add terms to objective function to penalise points on the “wrong side” of boundary, but don’t forbid them
- 👑 Use a kernel function
 - Transform input data into a separable space
 - Kernel SVM allows non-linear boundaries in the original space and is a popular approach



Polynomials of degree exactly d : $K(u, v) = (u \cdot v)^d$

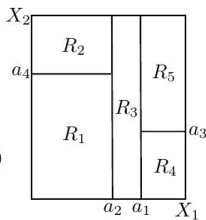
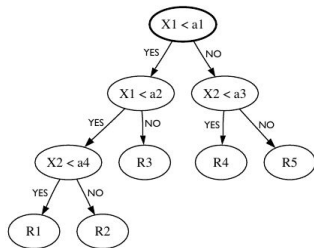
Polynomials of degree up to d : $K(u, v) = (u \cdot v + 1)^d$

Gaussian kernels: $K(\vec{u}, \vec{v}) = \exp\left(-\frac{\|\vec{u} - \vec{v}\|_2^2}{2\sigma^2}\right)$

Sigmoid: $K(u, v) = \tanh(\eta u \cdot v + \nu)$

How to build a decision tree from data?

- 1 Choose a decision point yielding best purity
- 2 Partition data into corresponding subsets
- 3 Reiterate with resulting subsets
- 4 Stop when regions are approximately pure



Impurity in classification

- 👑 misclassification
- 👑 Gini impurity: probability of incorrectly classifying a randomly chosen data point

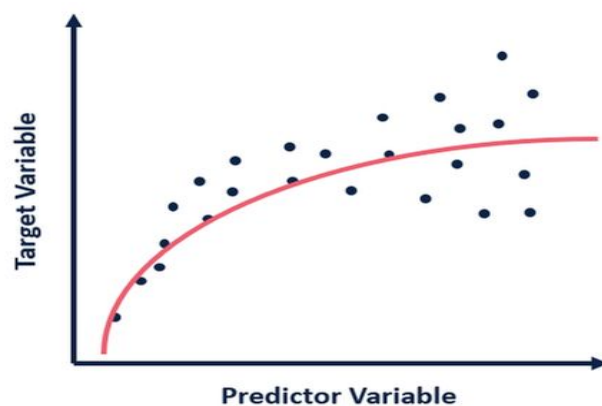
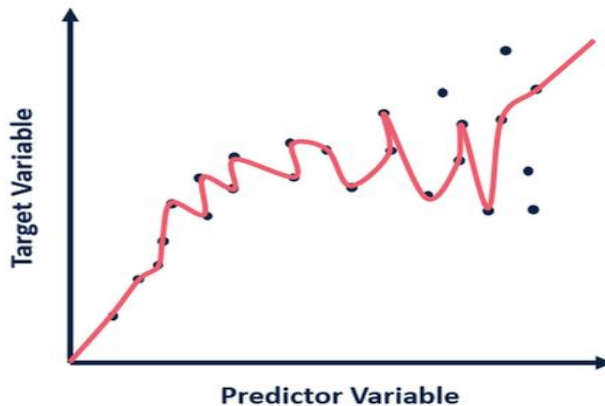
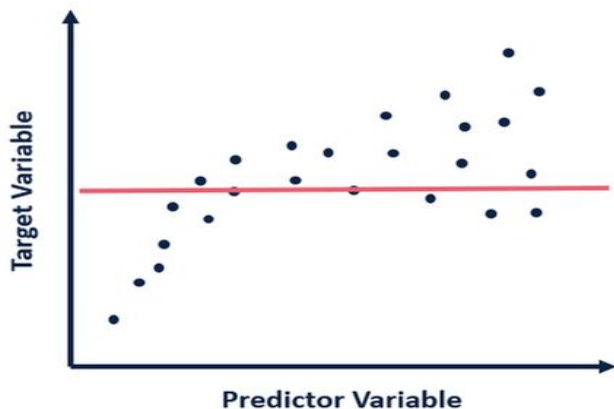
Impurity in regression

- 👑 mean squared error

$$F(R) = \sum_{x_i \in R} (y_i - \langle y \rangle)^2$$

Performance Estimation

- **Overfitting and underfitting**



- **Bias**

- Difference between prediction and real outcome

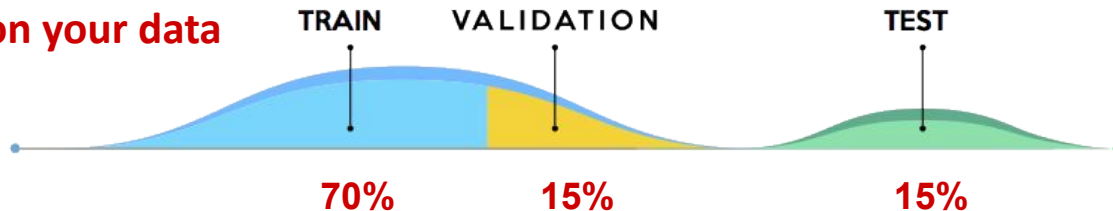
- **Variance**

- Variability of predictions

HOLDOUT STRATEGY

Highly depends on your data

1 Split your data into train / validation / test



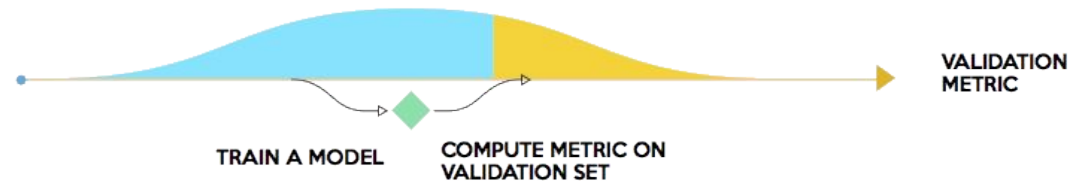
2 For each parameter combination

Parameter A (e.g., depth)

4	14
5	15
6	16
7	17

 Parameter B (e.g., n trees)

13	14
15	16
17	18



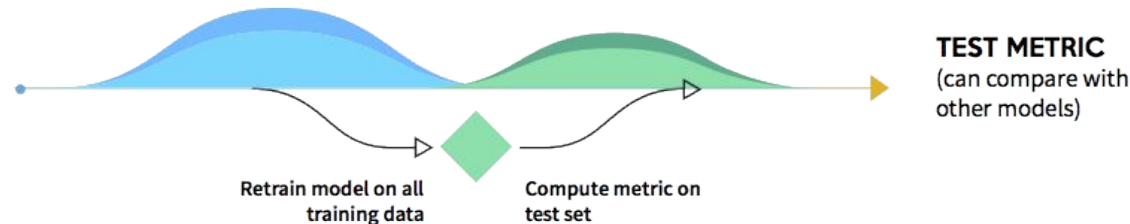
3 Choose the parameter combination with the best metric

Parameter A (e.g., depth)

6	14
---	----

 Parameter B (e.g., n trees)

14



Predictive Performance Metrics

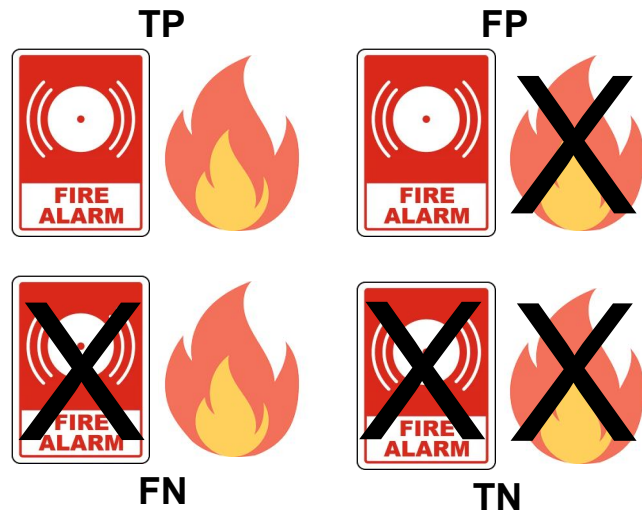
- To select the **best performing model** we need to be able to **compare them**
- **Predictive performance metrics**
 - On the **validation** set
 - On an **independent test set**
 - Make sure they are **consistent**
- There are multiple metrics for **classification** and **regression**
- **Classification**
 - We can derive several metrics from a **confusion matrix**

Statistical Classification Metrics

<div>Sensitivity Recall Power</div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>True Positive Rate</div>	TP	FP	FN	TN	TP	FP	FN	TN	<div>Precision</div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>Positive Predictive Value</div>	TP	FP	FN	TN	TP	FP	FN	TN	<div></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>False Discovery Rate</div>	TP	FP	FN	TN	TP	FP	FN	TN	<div>Type I Error α Fall Out</div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>False Positive Rate</div>	TP	FP	FN	TN	TP	FP	FN	TN	<div>Accuracy</div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div>F1 Score F Measure</div> <div><table><tr><td>2x TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>2x TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div></div>	2x TP	FP	FN	TN	2x TP	FP	FN	TN							
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
2x TP	FP																																																											
FN	TN																																																											
2x TP	FP																																																											
FN	TN																																																											
<div>Type II Error β</div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>False Negative Rate</div>	TP	FP	FN	TN	TP	FP	FN	TN	<div></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>True Discovery Rate</div>	TP	FP	FN	TN	TP	FP	FN	TN	<div></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>Negative Predictive Value</div>	TP	FP	FN	TN	TP	FP	FN	TN	<div>Specificity</div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>True Negative Rate</div>	TP	FP	FN	TN	TP	FP	FN	TN	<div>Confusion Matrix</div> <div><table><tr><td></td><td></td><td colspan="2">actual</td></tr><tr><td></td><td></td><td>T</td><td>F</td></tr><tr><td rowspan="2">predicted</td><td>p</td><td>TP</td><td>FP</td></tr><tr><td>N</td><td>FN</td><td>TN</td></tr></table></div> <div>TP: True Positive FP: False Positive FN: False Negative TN: True Negative</div> <div>actual = observed predicted = expected</div>			actual				T	F	predicted	p	TP	FP	N	FN	TN	<div>Matthews Correlation Coefficient</div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>difference of products</div> <div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div> <div>square root of product of sums</div>	TP	FP	FN	TN	TP	FP	FN	TN
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
		actual																																																										
		T	F																																																									
predicted	p	TP	FP																																																									
	N	FN	TN																																																									
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											

Type I and II errors

- Which Type Error is worse?
 - Fire alarm
- Type I:
 - Fire alarm **rings** when there is **no fire**
- Type II:
 - Fire alarm **fails to ring** when there **is fire**
- Which Type Error is worse in this case?



Which models is the best one?

Metric	MODEL1	MODEL2
Recall	0.6667	0.8333
Specificity	0.8333	0.6667
Precision	0.8000	0.7143
Accuracy	0.7500	0.7500
F1 Score	0.7273	0.7692
MCC	0.5071	0.5071

- **Model 1**
 - To minimise **Type I Error (better precision)**

- **Model 2**
 - To minimise **Type 2 Error (better recall)**

- **Benign vs. malignant cancer hypothetical cases**

MODEL 1	Actual disease	Actual healthy
Predicted disease	200	50
Predicted healthy	100	250

MODEL 2	Actual disease	Actual healthy
Predicted disease	250	100
Predicted healthy	50	200