# **Introduction to classification**

Gianluca Campanella

Classification

k-nearest neighbours classifier

**Metrics** 

Cost-benefit analysis

# Classification

#### Regression

Aim Predict a continuous value Loss How 'off' (numerically) our predictions are

Classification

Aim Predict a class Loss How 'inaccurate' the predicted classes are

# *k*-nearest neighbours classifier

Given a new observation...

- Find the *k* 'most similar' training sample(s)
- Use the most common class among them as prediction

#### Questions

- How do we define similarity?
- How many neighbours do we use?

# *k*-nearest neighbours classifier



From Burton DeWilde's blog

# **Choice of** *k*

- Larger  $k \rightarrow$  smoother boundaries, less 'noisy'
- If k = N, we always predict the majority class



From CS231n: Convolutional Neural Networks for Visual Recognition

## Minkowski distance

$$\left(\sum_{i}|x_{i}-y_{i}|^{p}\right)^{1/p}$$

$$p=1$$
 Manhattan distance  $\sum_i |x_i - y_i|$   
 $p=2$  Euclidean distance  $\sqrt{\sum_i (x_i - y_i)^2}$ 

#### **Uniform weights**

- All *k* neighbours contribute equally to the prediction
- Actual distance to each is ignored

#### **Distance weights**

- Contributions are weighted by 1/distance
- Closer neighbours influence the prediction more

As the number of variables (coordinates) increases...

- The volume of the space increases
- Pairwise distances become more similar  $\rightarrow$  sparsity
- Some samples have huge neighbourhoods  $\rightarrow$  'hubs'

# **Metrics**

# **Classification accuracy**

#### **Classification accuracy**

- Percentage of correct predictions
- Higher is better

#### **Classification error**

- Percentage of incorrect predictions (inverse of accuracy)
- Lower is better

# **Confusion matrix**



- Gives a better understanding of behaviour
- Can be used to define multiple performance metrics

Sensitivity (a.k.a. true positive rate)

 $\frac{\sum \text{True positive}}{\sum \text{Actual} = 1}$ 

**Specificity** (a.k.a. true negative rate)

 $\frac{\sum \text{True negative}}{\sum \text{Actual} = 0}$ 

**Sensitivity** (a.k.a. true positive rate)

 $\frac{\sum \text{True positive}}{\sum \text{Actual} = 1}$ 

#### **Perfect sensitivity**

- All sick identified as sick
- Negative test result definitely rules out disease

Specificity (a.k.a. true negative rate)  $\frac{\sum \text{True negative}}{\sum \text{Actual} = 0}$ 

**Sensitivity** (a.k.a. true positive rate)

 $\frac{\sum \text{ True positive}}{\sum \text{ Actual} = 1}$ 

#### **Perfect sensitivity**

- All sick identified as sick
- Negative test result definitely rules out disease

**Specificity** (a.k.a. true negative rate)

 $\frac{\sum True \ negative}{\sum Actual = 0}$ 

# Perfect specificity

- No healthy identified as sick
- Positive test result useful for ruling in disease

**Sensitivity** (a.k.a. true positive rate)

 $\frac{\sum \text{ True positive}}{\sum \text{ Actual} = 1}$ 

#### **Perfect sensitivity**

- All sick identified as sick
- Negative test result definitely rules out disease

# **Specificity** (a.k.a. true negative rate)

 $\frac{\sum True \ negative}{\sum Actual = 0}$ 

## Perfect specificity

- No healthy identified as sick
- Positive test result useful for ruling in disease

Can we maximise both at the same time?

#### 100% sensitivity

• 'Everyone is a terrorist!'

#### 100% sensitivity

- 'Everyone is a terrorist!'
- $\bullet\,$  All terrorists are stopped  $\rightarrow$  100% sensitivity

#### 100% sensitivity

- 'Everyone is a terrorist!'
- $\bullet\,$  All terrorists are stopped  $\rightarrow\,100\%$  sensitivity
- $\bullet\,$  No one can enter the country!  $\rightarrow$  0% specificity

#### 100% sensitivity

- 'Everyone is a terrorist!'
- $\bullet\,$  All terrorists are stopped  $\rightarrow\,100\%$  sensitivity
- No one can enter the country!  $\rightarrow$  0% specificity

#### 100% specificity

• 'No one is a terrorist!'

#### 100% sensitivity

- 'Everyone is a terrorist!'
- $\bullet\,$  All terrorists are stopped  $\rightarrow\,100\%$  sensitivity
- No one can enter the country!  $\rightarrow$  0% specificity

#### 100% specificity

- 'No one is a terrorist!'
- $\bullet\,$  All non-terrorists are allowed in  $\rightarrow\,100\%$  specificity

#### 100% sensitivity

- 'Everyone is a terrorist!'
- $\bullet\,$  All terrorists are stopped  $\rightarrow\,100\%$  sensitivity
- No one can enter the country!  $\rightarrow$  0% specificity

#### 100% specificity

- 'No one is a terrorist!'
- $\bullet\,$  All non-terrorists are allowed in  $\rightarrow\,100\%$  specificity
- $\bullet\,$  All terrorists are also allowed into the country!  $\rightarrow$  0% sensitivity

# ROC and AUC

#### Receiver Operating Characteristic (ROC) curve



Sensitivity vs (1 – specificity)  $\rightarrow$  TP rate vs FP rate

# **ROC and AUC**

## Receiver Operating Characteristic (ROC) curve



#### Area Under the Curve (AUC)

- Probability that Prediction(actual 1) > Prediction(actual 0)
- Random guess  $\rightarrow$  AUC = 50% (diagonal)
- Higher is better

Sensitivity vs (1 – specificity)  $\rightarrow$  TP rate vs FP rate

# **Cost-benefit analysis**

- Assume that the four possible outcomes of a classification problem have (numerical) benefits and costs
  - > 0 desirable (e.g. profit)
  - $= 0 \ neutral$
  - < 0 undesirable (e.g. loss)
- These 'benefits and costs' needn't be symmetrical

You have 20 people enrolled in an outdoor activity paying £30 each

- Before the activity, you check the weather forecast and either:
  - Go ahead, which costs you £5 per participant
  - Cancel and refund the participants in full
- If you decide to go ahead, the day of the activity it will either:
  - Be sunny, in which case you get to keep the profit
  - Rain, in which case you'll have to refund the participants in full

What is the 'benefits and costs' matrix?



- What are the sensitivity and specificity?
- What is the expected profit?