

23

25

17

Crossing time

2 can cross at a time  $\rightarrow$  true is for slowest  
only with "flashlight"

Fibonacci Number

$$F_0 = 0$$
$$F_1 = 1$$

$$F_n = F_{n-1} + F_{n-2}$$
$$\forall n \geq 2$$

$F_0$	$F_1$	$F_2$	$F_3$	$F_4$	$F_5$	$F_6$	$F_7$	$F_8$	$F_9$	$F_{10}$
0	1	1	2	3	5	8	13	21	34	55

Task: Compute  $F_n$  = function of  $n$  = input

① Math - following recursive

```
Fib(n)
  if n == 0 return 0
  if n == 1 return 1
  (else) return Fib(n-1) + Fib(n-2)
```

④ correct? yes

③ how fast?  
exponential  
runtime

$$\Theta(2^n) \quad ??$$

tight  
asymptote

lower bound

upper bound

$\Omega$

$\Theta$

$O$

asympt

exact  
asymptote

asympt

ex  $\Theta(n^2) = f(n)$

$$C_2 \cdot n^2 \leq f(n) \leq C_1 \cdot n^2$$

low bound

upper  
bound



② Array computation  $Fib(n)$   
 $F = \text{array}$   $F[0] = 0$ ,  $F[1] = 1$

for  $i = 2 : n$   
 $F[i] = F[i-1] + F[i-2]$

return  $F[n]$

C. Storage  
A correct  
B how fast

LINEAR

$\Theta(n)$

Q: can we store  
only last 3 values?

### ③ Matrix-Multiplication-based $\text{Fib}(n)$

$$M = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \quad \text{claim} \quad M^n = \begin{bmatrix} F_{n+1} & F_n \\ F_n & F_{n-1} \end{bmatrix}$$

$2 \times 2$

induction proof base case  $M^1 = M = \begin{bmatrix} F_2 & F_1 \\ F_1 & F_0 \end{bmatrix}$ ?

inductive step  $n \rightarrow n+1$

$$\boxed{M^n = \begin{bmatrix} F_{n+1} & F_n \\ F_n & F_{n-1} \end{bmatrix}} \quad \text{ind hyp} \quad \Rightarrow \quad \boxed{M^{n+1} = \begin{bmatrix} F_{n+2} & F_{n+1} \\ F_{n+1} & F_n \end{bmatrix}} \quad \text{ind conclusion}$$

proof:  $M^{n+1} = M^n \times M =$

$$= \begin{bmatrix} F_{n+1} + F_n & F_{n+1} \\ F_{n+1} + F_n & F_n \end{bmatrix} = \begin{bmatrix} F_{n+2} & F_{n+1} \\ F_{n+1} & F_n \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$$

✓

② Matrix  $B_{bb}(n)$  Run Time?

$M = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$  ✓

compute  $M^n = A$

exponential of matrix

return  $M[1, 2]$  ✓

naive  $M^n =$   
 $M \cdot M \cdot M \cdot \dots \cdot M$   
n times

$\Theta(n)$  time

$\Theta(\log n)$

Fast exponentiation (example)

$M$

$M^2 = M \cdot M$

$M^4 = M^2 \cdot M^2$

$M^8 = M^4 \cdot M^4$

$M^{16} = M^8 \cdot M^8$

$M^{32} = M^{16} \cdot M^{16}$

$M^{64} = M^{32} \cdot M^{32}$

want  $M^{100} = M^{64} \cdot M^{32} \cdot M^4$

$M^{171} = M^{128} \cdot M^{32} \cdot M^8 \cdot M^2 \cdot M^1$

④ Generating Function

$$F(n) = \frac{\phi^n - \psi^n}{\sqrt{5}}$$

$\phi$  real number  $\in \mathbb{R}$   
 math-processor estimate  
 $\approx$  constant time.  
 $\approx \Theta(1)$

$$\phi = \frac{1 + \sqrt{5}}{2}$$

$$\psi = \frac{1 - \sqrt{5}}{2}$$

idea for  $\phi, \psi$   
 guess  
 $F(n) \approx a^n$  ? true  
 rec. Fib:

$$F_{n+1} = F_n + F_{n-1}$$

$$a^{n+1} = a^n + a^{n-1} \quad | \cdot a^{-n}$$

$$a^2 = a + 1$$

quad eq  $\Rightarrow \phi, \psi$ .

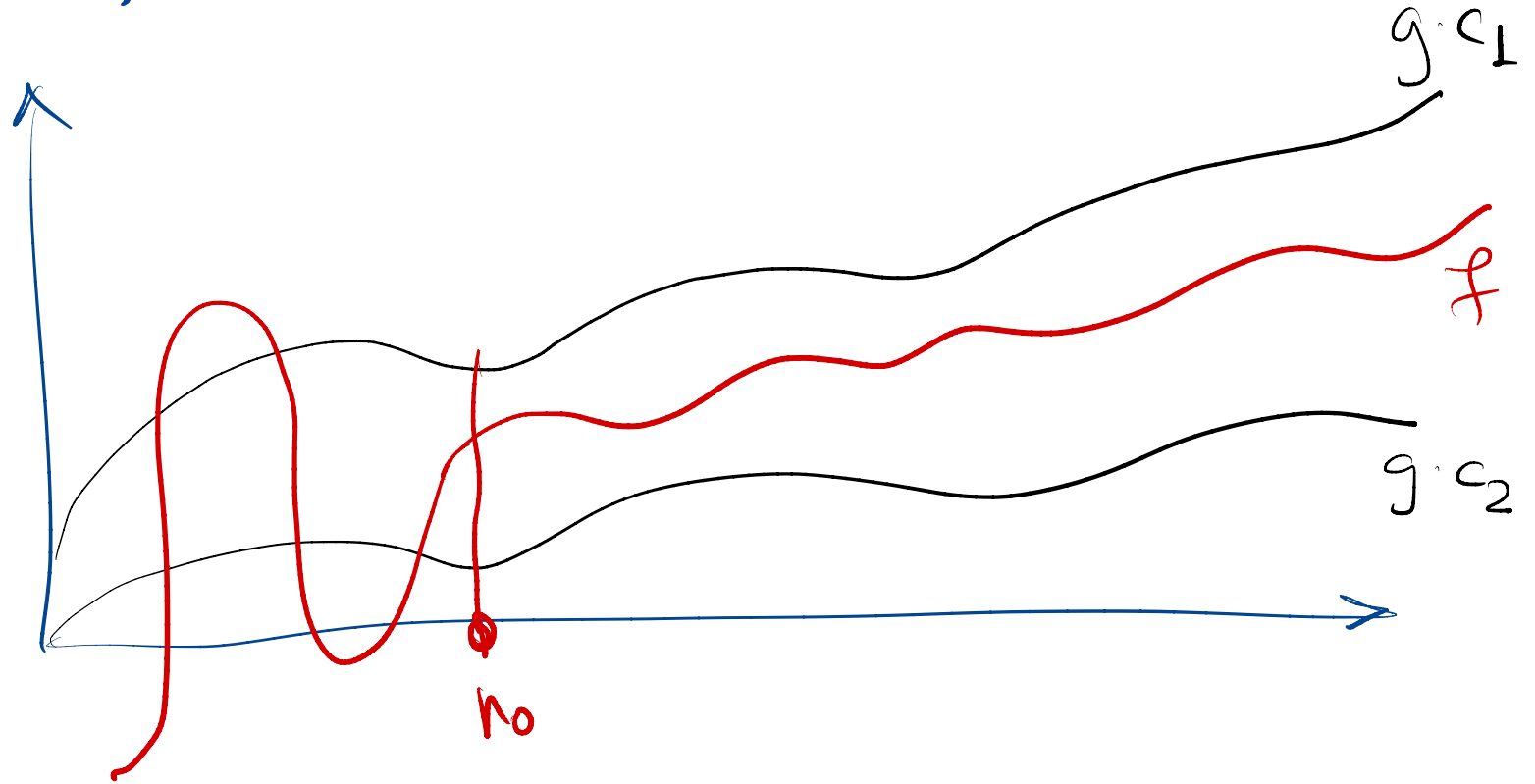
big O notation

$$f = \Theta(g)$$

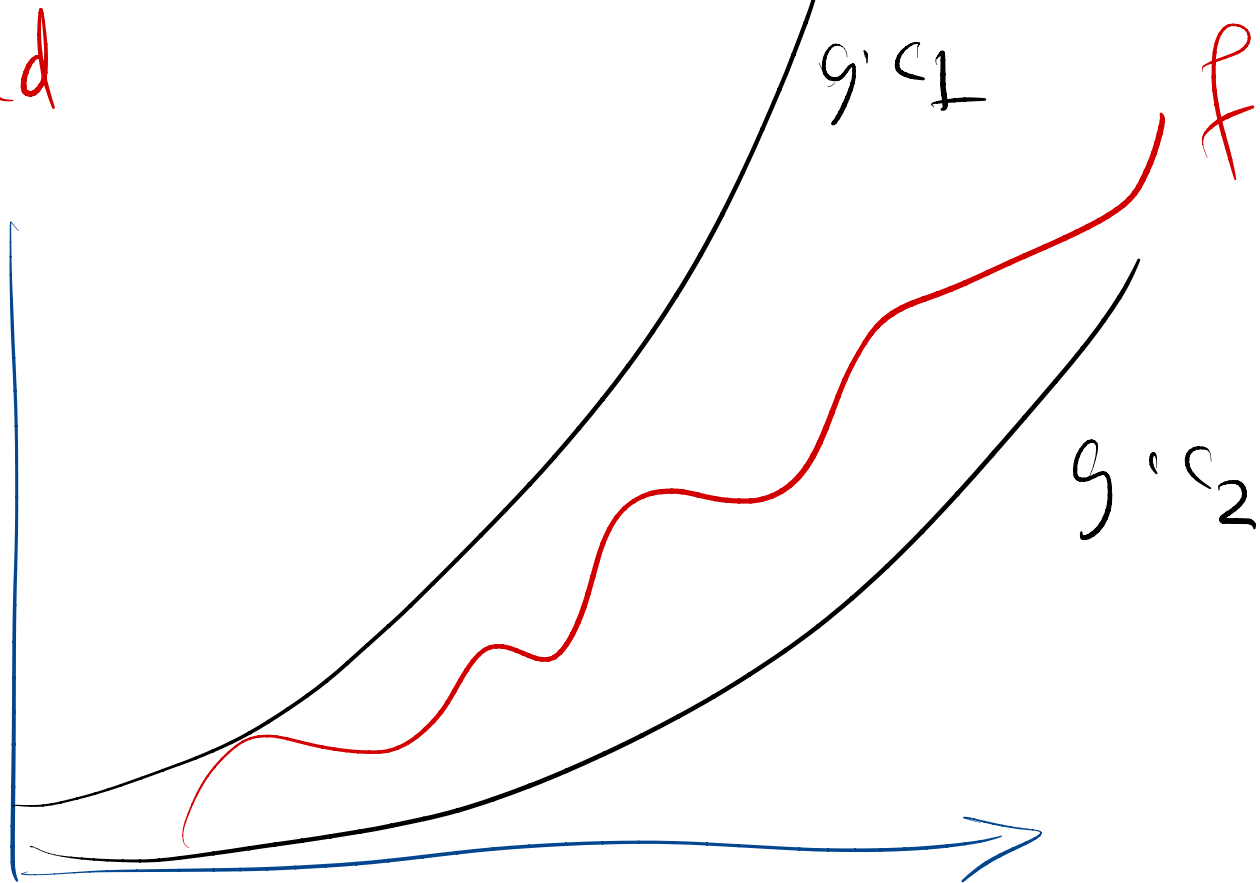
$$f(n) = \Theta(g(n))$$

$c_1, c_2 > 0$  constants  $c_1 > c_2$

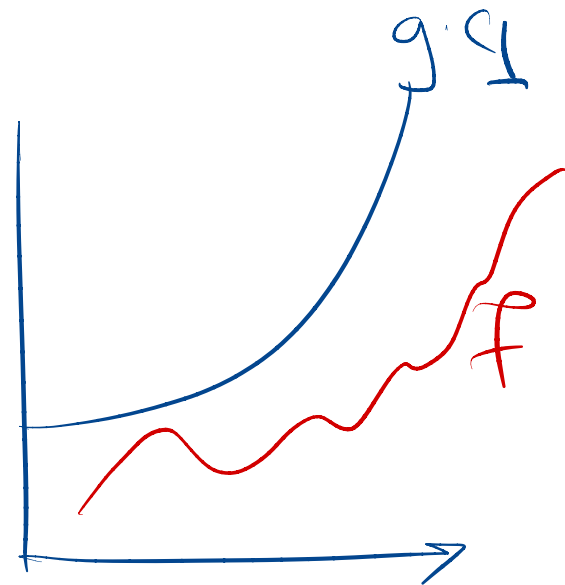
$$g \cdot c_2 \leq f \leq g \cdot c_1$$



$g = \text{quad}$



$O(g) = \text{upper bound } f = O(g)$   
 $\exists c_1 > 0 \quad f \leq g \cdot c_1$



$f(n)$

$2n^2 + 3n + 2$

$\Theta(n^2)$

$n^2 \log n + n^3$

$\Theta(n^3)$

$2^n + 5n^2 - 3$

$\Theta(2^n)$

$3^n + 2^n$

$\Theta(3^n)$

$3^n = \Theta(2^n)?$

$c_2 \cdot 2^n \leq 3^n \leq c_1 \cdot 2^n$

true

$\frac{3^n}{2^n} \leq c_1$

$(\frac{3}{2})^n \leq c_1$

$\lim_{n \rightarrow \infty} (\frac{3}{2})^n = \infty$   
**No**

$\log_a x = b \iff a^b = x$

$f = \log_2 x$

$g = \log_3 x$

$f = \Theta(g)$

$c_2 \cdot \log_3(n) \leq \log_2(n) \leq c_1 \cdot \log_3(n)$

$\log_3(x) = ? \log_2(x) \cdot \log_3 2?$

$\log_3 x = \log_2 x \cdot \log_3 2$   
 $3 = 3^{(\log_3 2) \log_2 x}$   
**X**

$$C_2 = 1$$
$$\log_3(n) \leq \log_2(n) \leq C_1 \cdot \log_3(n) \quad ??$$

2  $^{\log_2 x}$   
⊗

$$\log_2(n) \leq C_1 \cdot \log_2(n) \cdot \log_3 2$$

$$1 \leq C_1 \cdot \log_3 2 \quad \text{constant}$$





$$k=1 \quad T(n) = 1 + T(n/2) \quad \forall n$$

$$k=2 \quad = 1 + \left[ 1 + T(n/4) \right] = 2 + T(n/4)$$

$$k=3 \quad = 2 + \left[ 1 + T(n/8) \right] = 3 + T(n/8)$$

$$k=4 \quad = 3 + \left[ 1 + T(n/16) \right] = \boxed{4} + T(n/16)$$

pattern

$k$ :

$$k + T(n/2^k)$$

optional (if messy)

step  $k$

$$k + T(n/2^k)$$

induction proof

step  $k+1$

$$k+1 + T(n/2^{k+1})$$

$$T(n) = k + T\left(\frac{n}{2^k}\right)$$

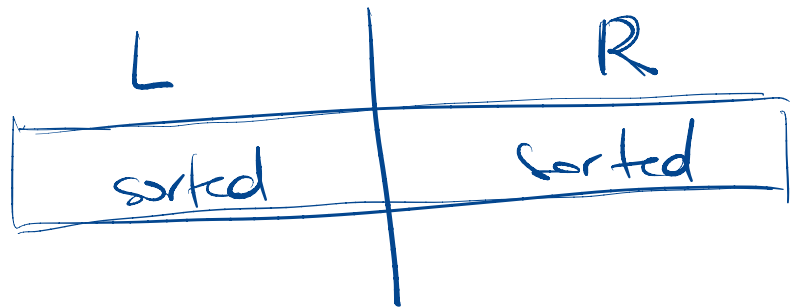
$$\text{last } k : \frac{n}{2^k} \approx 1 \Rightarrow T\left(\frac{n}{2^k}\right) \approx T(1) = \text{base}$$

$$n \approx 2^k$$

$$k \approx \log_2(n)$$

$$T(n) = \log_2(n) + \underbrace{T(1)}_{\text{const}} = \Theta(\log n)$$

Merge Sort  $A(b:e)$   
: sort Array A



Split in  $\frac{1}{2}$   $A(b:m)$   
 $A(m+1:e)$

MergeSort ( $A(b:m)$ )  $T(n/2)$

MergeSort ( $A(m+1:e)$ )  $T(n/2)$

Merge/Combine sorted  $A(b:m)$  & sorted  $A(m+1:e)$   
 $i = 1 (L)$   $j = 1 (R)$   $t = 1$

if  $L[i] < R[j]$

$C[t] = L[i],$

$i = i + 1$

$\Theta(n)$

else

$C[t] = R[j]$

$j = j+1$   
 $t = t+1$   
 until both arrays finished.

R.T  $T(n) = T(n/2) + T(n/2) + n$   
L-sort          R-sort          Mergeup

$k=1$

$2T(n/2) + n$

$k=2$

$\Rightarrow 2[2T(n/4) + n/2] + n = 4T(n/4) + 2n$

$k=3$

$= 4[2T(n/8) + n/4] + 2n = 8T(n/8) + 3n$

$k=4$

$= 8[2T(n/16) + n/8] + 3n = 16T(n/16) + 4n$

$k$  pattern

$= 2^k T(n/2^k) + kn$

$$\text{last } k: \frac{n}{2^k} \approx 1 \text{ last } \Leftrightarrow k \approx \log_2 n$$

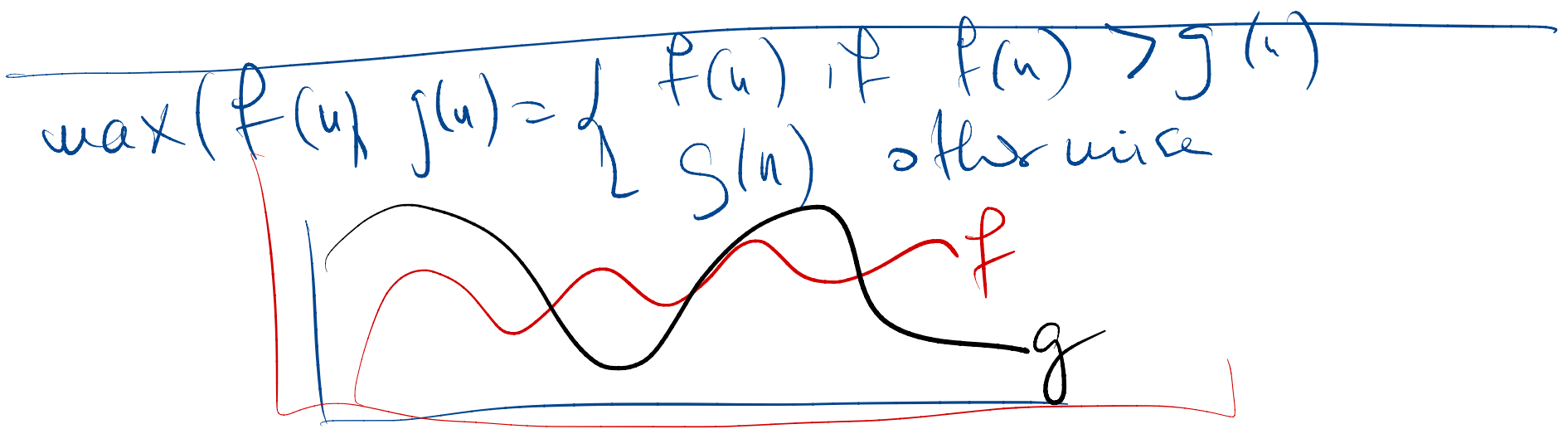
$$\begin{aligned} T(n) &= 2^{\log_2 n} T(1) + (\log_2 n) \cdot n \\ &= n \cdot \text{const} + n \log_2 n \\ &= \Theta(n \log n). \end{aligned}$$

HW1 3.1-1  $\max(f, g) = \Theta(f + g)$

Def  $\Theta$   $c_2(f + g) \leq \max(f, g) \leq c_1(f + g)$   
 $c_1, c_2 > 0$

guess  $c_1 = 2$   $c_2 = 1/3$   $\leftrightarrow \frac{1}{3}(f + g) \leq \max(f, g) \leq 2(f + g)$

$c_2 < c_1$



any function  $\min(f, g) \leq f \leq \max(f, g)$

$\ln \ln(n^2)$

$2^{\ln^2(n)}$

$\ln(\ln(n^2))$

$2^{\ln(n) \cdot \ln(n)}$

$\ln(\ln(n^2))$

asymptote  $\ll$

$n \ln(n)$

bigger

$\ln(\ln(n^2)) \ll$

explanation

$\log(n)$  not flaters  
 $n^2$  easy  $\ll n$   $\ln(n)$





## K-nearest neighbours

Javier Béjar © 1 2 3 4

LSI - FIB

Term 2012/2013

supervised  
(not this class)

$f$  (datapoint-representation)  $\approx$  label  
 $x_i$   $y_i$

classifier

1 Non Parametric Learning

- 2 K-nearest neighbours
  - K-nn Algorithm
  - K-nn Regression
  - Advantages and drawbacks
  - Application

## 1 Non Parametric Learning

- 2 K-nearest neighbours
  - K-nn Algorithm
  - K-nn Regression
  - Advantages and drawbacks
  - Application

## Parametric vs Non parametric Models

- In the models that we have seen, we select a hypothesis space and adjust a fixed set of parameters with the training data ( $h_\alpha(x)$ )
- We assume that the parameters  $\alpha$  summarize the training and we can forget about it
- This methods are called **parametric** models
- When we have a small amount of data it makes sense to have a small set of parameters and to constraint the complexity of the model (avoiding overfitting)

## Parametric vs Non parametric Models

- When we have a large quantity of data, overfitting is less an issue
- If data shows that the hypothesis has to be complex, we can try to adjust to that complexity
- A **non parametric** model is one that can not be characterized by a fixed set of parameters
- A family of non parametric models is **Instance Based Learning**

## Instance Based Learning

- Instance based learning is based on the memorization of the dataset
- The number of parameters is unbounded and grows with the size of the data
- There is not a model associated to the learned concepts
- The classification is obtained by looking into the memorized examples
- The cost of the learning process is 0, all the cost is in the computation of the prediction
- This kind learning is also known as **lazy learning**

1 Non Parametric Learning

- 2 K-nearest neighbours
  - K-nn Algorithm
  - K-nn Regression
  - Advantages and drawbacks
  - Application

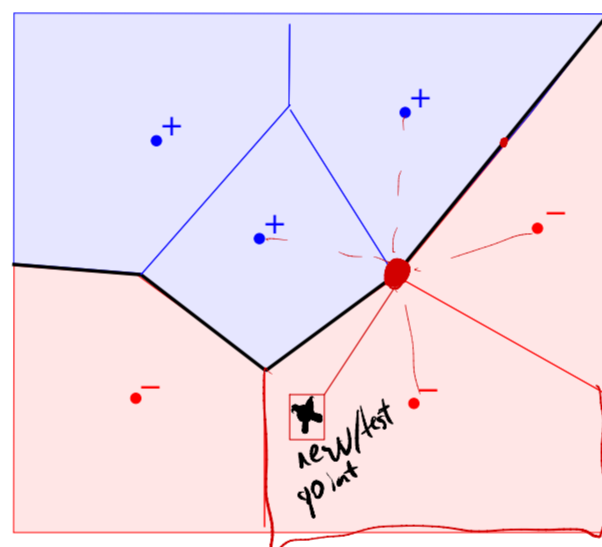
## K-nearest neighbours

- **K-nearest neighbours** uses the local neighborhood to obtain a prediction
- The  $K$  memorized examples more similar to the one that is being classified are retrieved
- A distance function is needed to compare the examples similarity
  - Euclidean distance ( $d(x_j, x_k) = \sqrt{\sum_i (x_{j,i} - x_{k,i})^2}$ )
  - Mahnattan distance ( $d(x_j, x_k) = \sum_i |x_{j,i} - x_{k,i}|$ )
- This means that if we change the distance function, we change how examples are classified



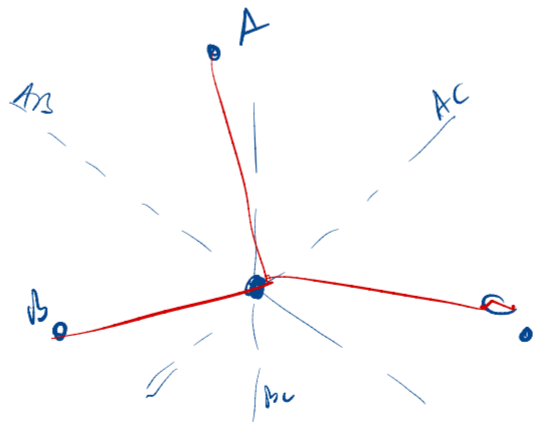
### K-nearest neighbours - hypothesis space (1 neighbour)

Binary  
Labels = {+, -}



{+, -}  
training set  
(I know labels)

unique partition  
of plane (the space)  
s.t.  $x$  is in  
the partition with  
closest centroid.



## K-nearest neighbours - Algorithm

- Training: Store all the examples
- Prediction:  $h(x_{new})$ 
  - Let be  $x_1, \dots, x_k$  the  $k$  more similar examples to  $x_{new}$
  - $h(x_{new}) = \text{combine\_predictions}(x_1, \dots, x_k)$
- The parameters of the algorithm are the number  $k$  of neighbours and the procedure for combining the predictions of the  $k$  examples
- The value of  $k$  has to be adjusted (crossvalidation)
  - We can overfit ( $k$  too low)
  - We can underfit ( $k$  too high)

Neighbourhood- $k$   
 $N(x_{new}) =$   
 the closest  
 $k$  points in  
 training

$k=3$

Good - close

Bad - far away

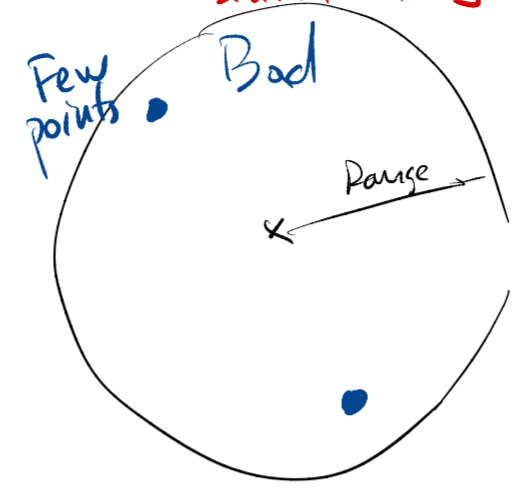
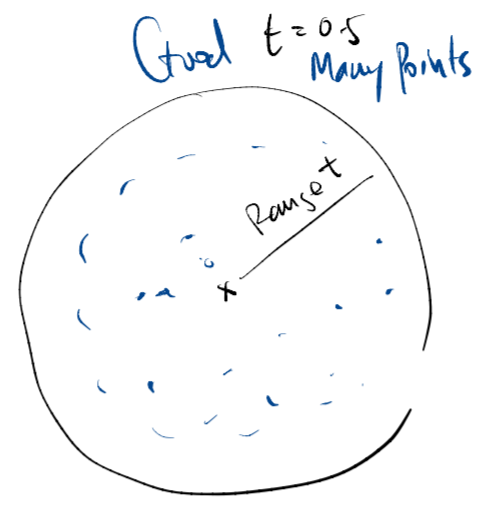


- data density - (not uniform)
  - easy to find close neighbors (dense / typical)
  - hard to find close neighbors (sparse / anomaly)

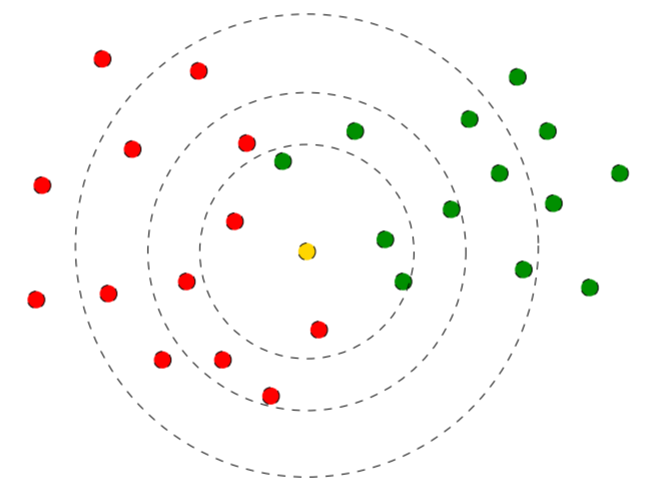


K NN  
top/closest k  
 $N_x = \{ \text{closest } k \text{ points to } x \}$

Range-NN  
"t"  
 $N_x = \{ \text{all points within } t \text{ similarity} \}$   
=  $\{ \text{Z-training} \mid \text{sim}(x, z) \geq t \}$



K-nearest neighbours - Prediction



Weighted version

$\psi(x, z) = \text{similarity}(x, z)$

use all points

$$\text{predict}(x, L) = \text{avg} \sum_{i=1}^N L(z_i) \cdot \psi(x, z_i)$$

$\downarrow$  test  $\downarrow$  label  $\downarrow$  label L for point  $z_i$   $\downarrow$  weighted avg

## Looking for neighbours

- Looking for the K-nearest examples for a new example can be expensive
- The straightforward algorithm has a cost  $O(n \log(k))$ , not good if the dataset is large
- We can use indexing with *k-d trees* (multidimensional binary search trees)
  - They are good only if we have around  $2^{dim}$  examples, so not good for high dimensionality
- We can use *locality sensitive hashing* (approximate k-nn)
  - Examples are inserted in multiple hash tables that use hash functions that with high probability put together examples that are close
  - We retrieve from all the hash tables the examples that are in the bin of the query example
  - We compute the k-nn only with these examples

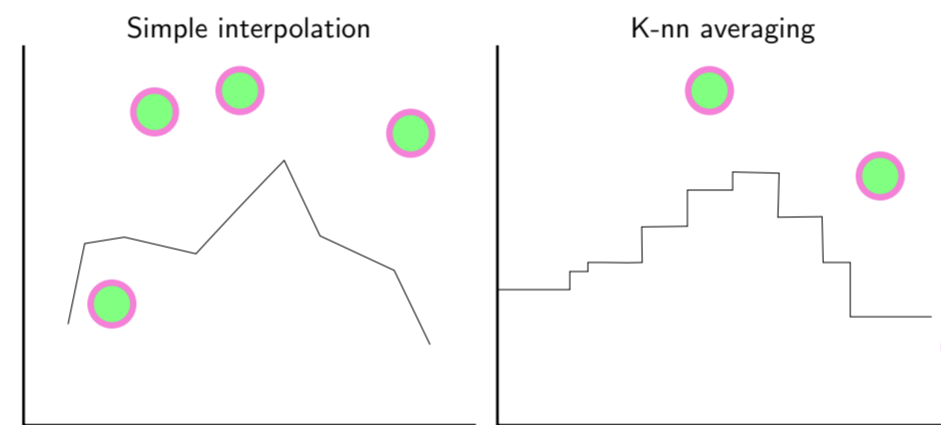
## K-nearest neighbours - Variants

- There are different possibilities for computing the class from the  $k$  nearest neighbours
  - Majority vote
  - Distance weighted vote
    - Inverse of the distance
    - Inverse of the square of the distance
    - Kernel functions (gaussian kernel, tricube kernel, ...)
- Once we use weights for the prediction we can relax the constraint of using only  $k$  neighbours
  - We can use  $k$  examples (local model)
  - We can use all examples (global model)

## K-nearest neighbours - Regression

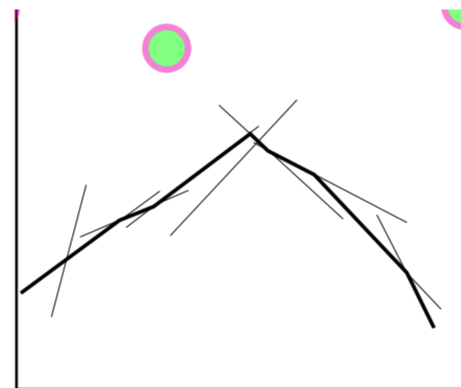
- We can extend this method from classification to regression
- Instead of combining the discrete predictions of k-neighbours we have to combine continuous predictions
- This predictions can be obtained in different ways:
  - Simple interpolation
  - Averaging
  - Local linear regression
  - Local weighted regression
- The time complexity of the prediction will depend on the method

## K-nearest neighbours - Regression



## K-nearest neighbours - Regression (linear)

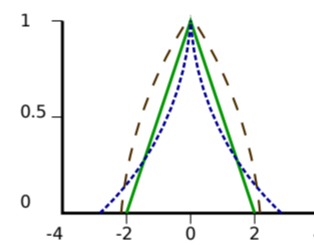
- K-nn linear regression fits the best line between the neighbors
- A linear regression problem has to be solved for each query (least squares regression)





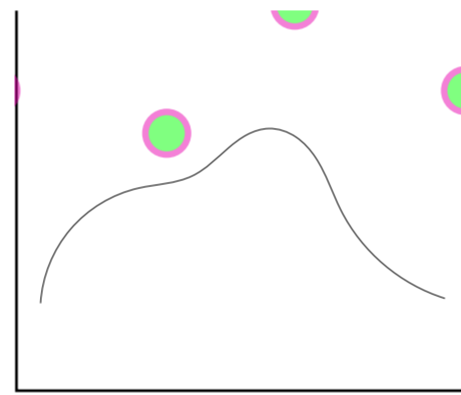
## K-nearest neighbours - Regression (LWR)

- Local weighted regression uses a function to weight the contribution of the neighbours depending on the distance, this is done using a **kernel function**



- Kernel functions have a width parameter that determines the decay of the weight (it has to be adjusted)
  - Too narrow  $\implies$  overfitting
  - Too wide  $\implies$  underfitting
- A weighted linear regression problem has to be solved for each query (gradient descent search)

### K-nearest neighbours - Regression (LWR)



## K-nearest neighbours - Advantages

- The cost of the learning process is zero
- No assumptions about the characteristics of the concepts to learn have to be done
- Complex concepts can be learned by local approximation using simple procedures

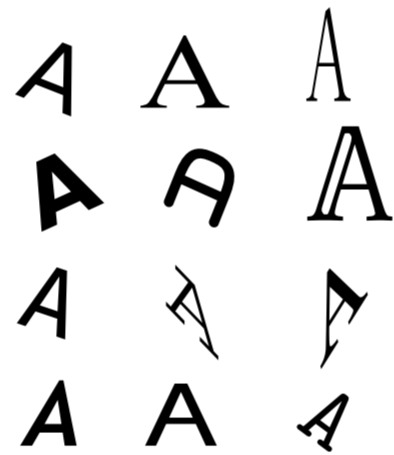
## K-nearest neighbours - Drawbacks

- The model can not be interpreted (there is no description of the learned concepts)
- It is computationally expensive to find the  $k$  nearest neighbours when the dataset is very large
- Performance depends on the number of dimensions that we have (*curse of dimensionality*)  $\implies$  *Attribute Selection*

## The Curse of dimensionality

- The more dimensions we have, the more examples we need to approximate a hypothesis
- The number of examples that we have in a volume of space decreases exponentially with the number of dimensions
- This is specially bad for k-nearest neighbors
  - If the number of dimensions is very high the nearest neighbours can be very far away

## Optical Character Recognition



- OCR capital letters
- 14 Attributes (All continuous)
- Attributes: horizontal position of box, vertical position of box, width of box, height of box, total num on pixels, mean x of on pixels in box, . . .
- 20000 instances
- 26 classes (A-Z)
- Validation: 10 fold cross validation

## Optical Character Recognition: Models

- K-nn 1 (Euclidean distance, weighted): accuracy 96.0%
- K-nn 5 (Manhattan distance, weighted): accuracy 95.9%
- K-nn 1 (Correlation distance, weighted): accuracy 95.1%

idea: change of variable

$$2^{\ln^2 n}$$

$$(\ln(n))!$$

$$m = \ln(n)$$

$$2^{m \cdot m}$$

$$m!$$