# KNN Imputation

## Christian Werner
*(Quantitative geneticist and biostatistician)* **EiB, CIMMYT**, Texcoco (Mexico)

## Filippo Biscarini
*(Biostatistician, bioinformatician and quantitative geneticist)* **CNR-IBBA**, Milan (Italy)
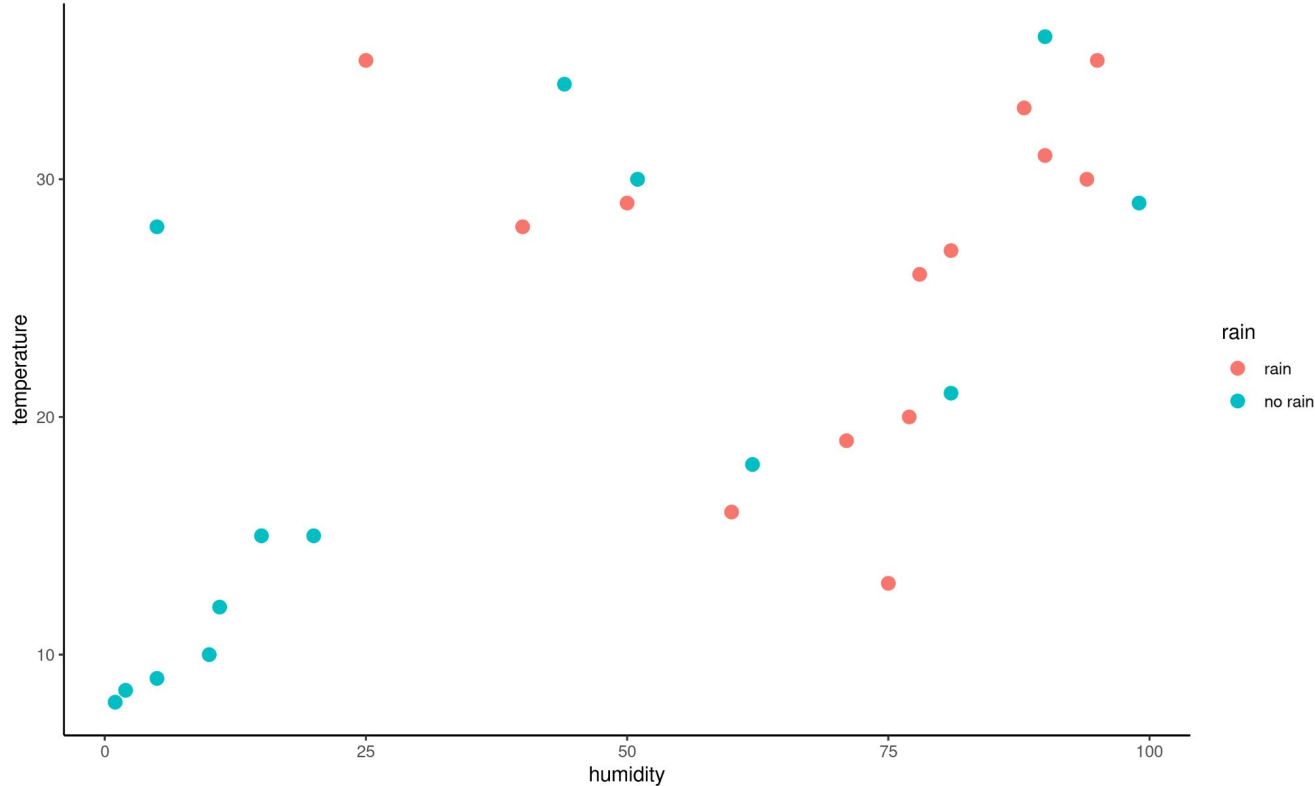
HerrFalloppio

## Oscar González-Recio
*(Computational biologist and quantitative geneticist)* **INIA-UPM**, Madrid (Spain)
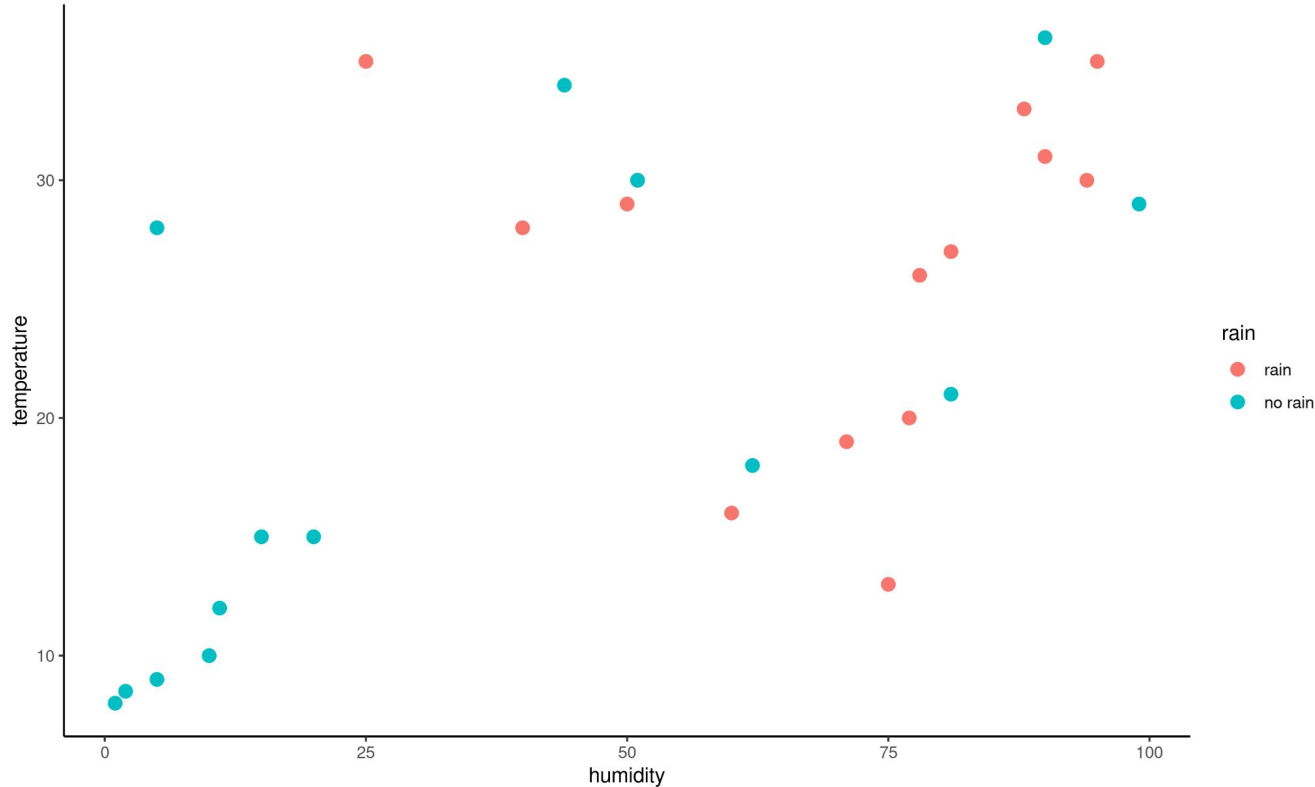
OscarGenomics

# One-minute k-nearest neighbors (KNN)!



1. We collect data on temperature, humidity and rain, for a number of days

2. New day: will it rain?

# k-nearest neighbors (KNN): a bit of math



KNN ALGORITHM

1. the values for humidity and temperature define the new datapoint in input

2. measure the Euclidean distance (Pythagoras theorem inside!) with all other points

3. find the closest point (neighbor) to our new point

4. assign to the new point the label (colour) of this nearest neighbor (k=1)

# One-minute k-nearest neighbors (KNN)!

If it looks like a duck, swims like a duck, and quacks like a duck, then it probably is a duck.
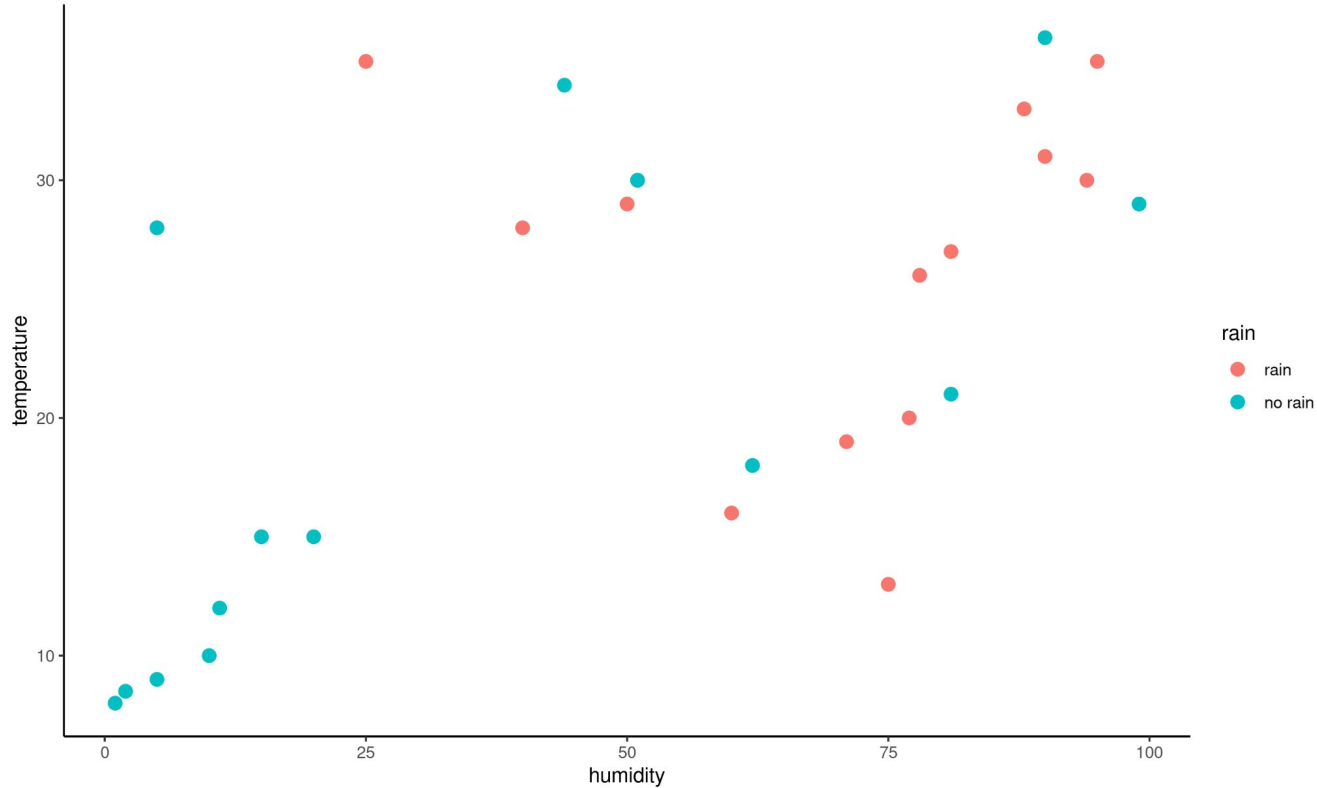
$$Pr(Y = j | X = x_0)$$

# KNN: more than 3-D

- We saw a simple example with few datapoints and only 2 dimensions
- When we have (typically) many dimensions (e.g. many SNPs) and many datapoints (e.g. many samples) we need the machine to do it! (and we can no longer visualize it)
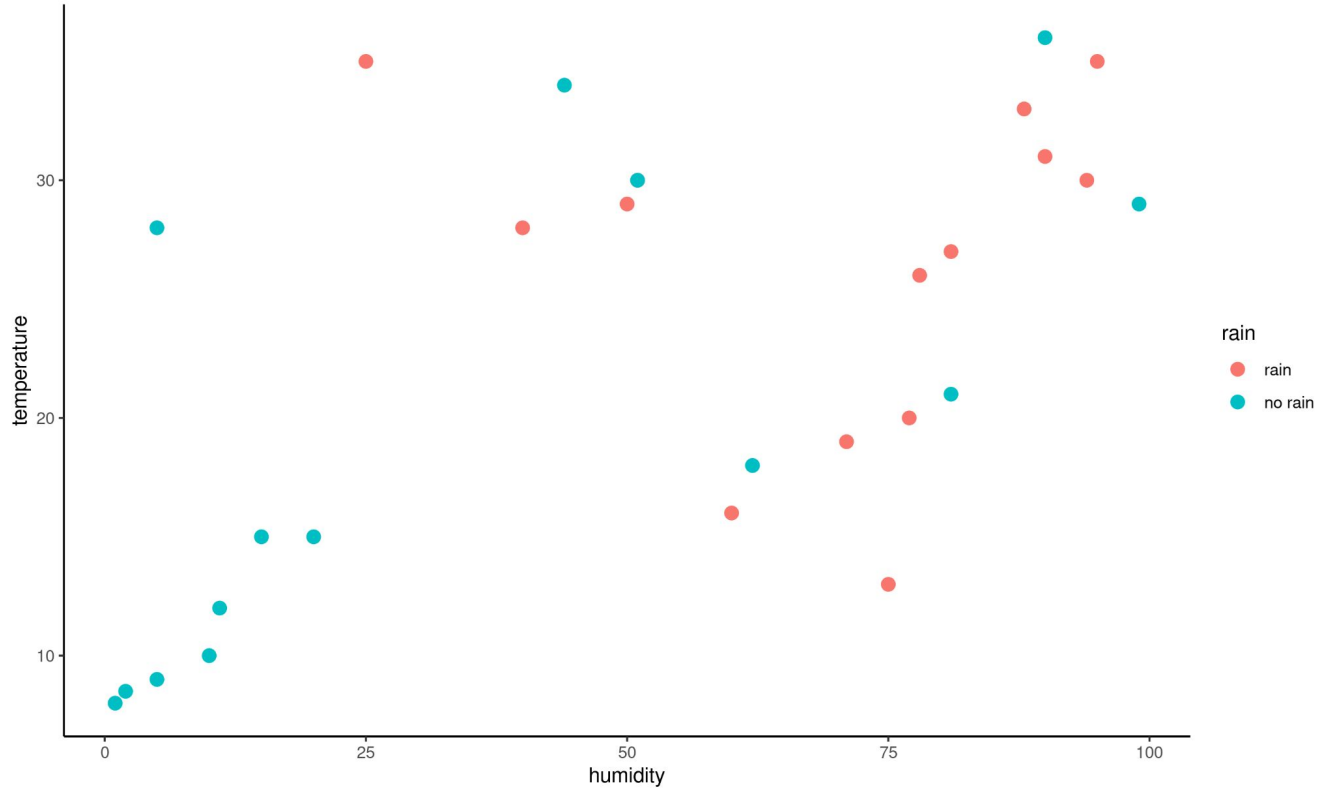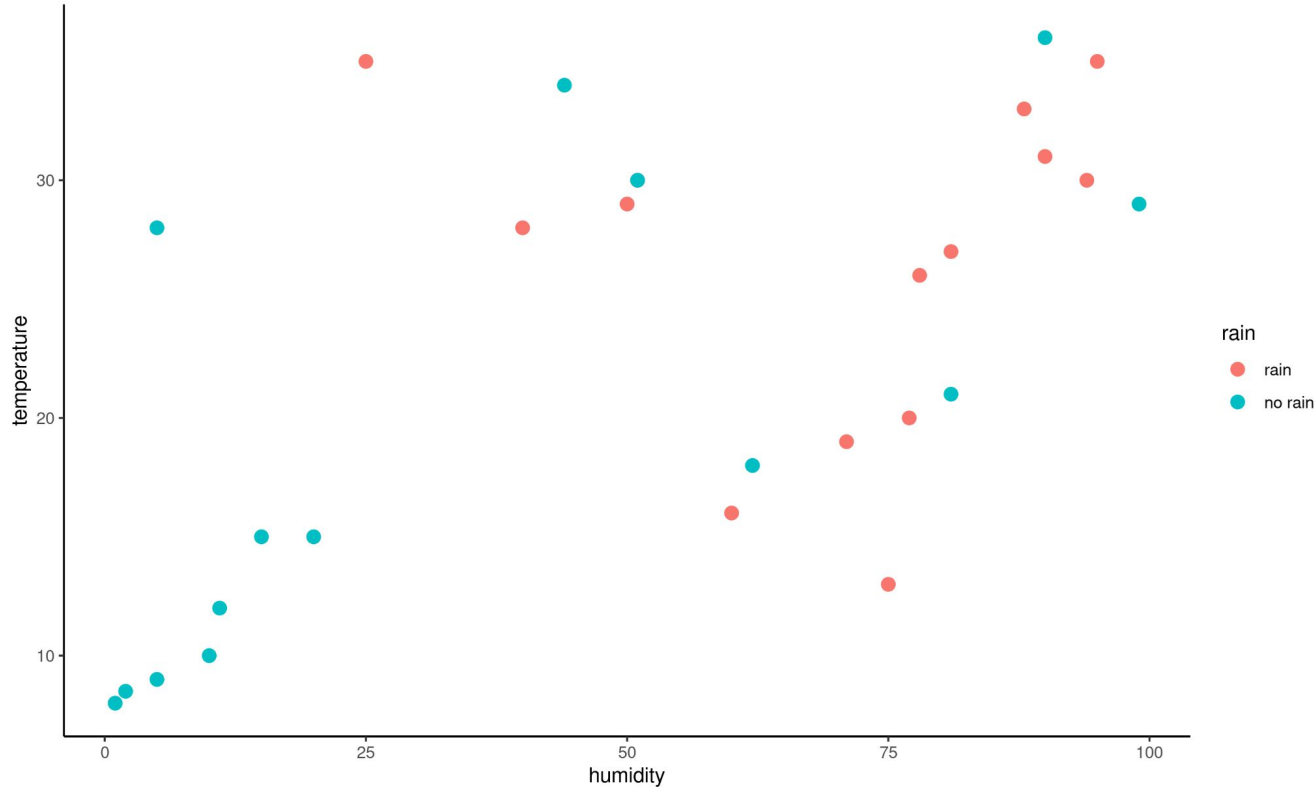
# KNN: K?



- K-neighborhood → majority!

# KNN: majority?



- K-neighborhood → majority!

    - average
    - weighted majority
    - weighted avg
    - etc.

# KNN: Euclidean distance?



- similarity↔dissimilarity
- [0, +∞]
- many possible distance metrics (Hamming, Chebyshev, Jaccard etc.) → see here

# KNN: a few things to tweak

KNN is said to be non-parametric, still:

- size of neighborhood (K)

- type of distance

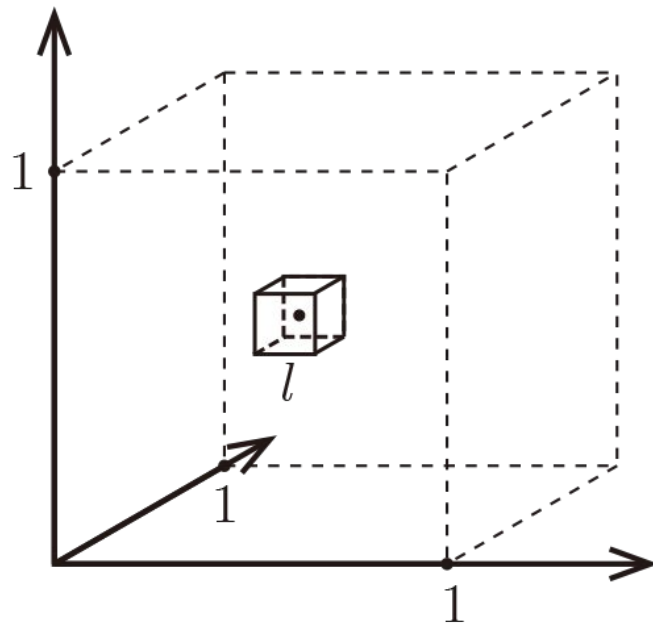- type of assignment metric (majority, average, weighted metrics, etc.)

# The curse of dimensionality

KNN assumes that **similar (close) points share similar labels/target**

unfortunately, **in high dimensional spaces points tend to never be close together**

- increasing the number of dimensions (parameters) of the problem increases and complicates the identification of k neighbors which are close enough to the data point to be classified/predicted
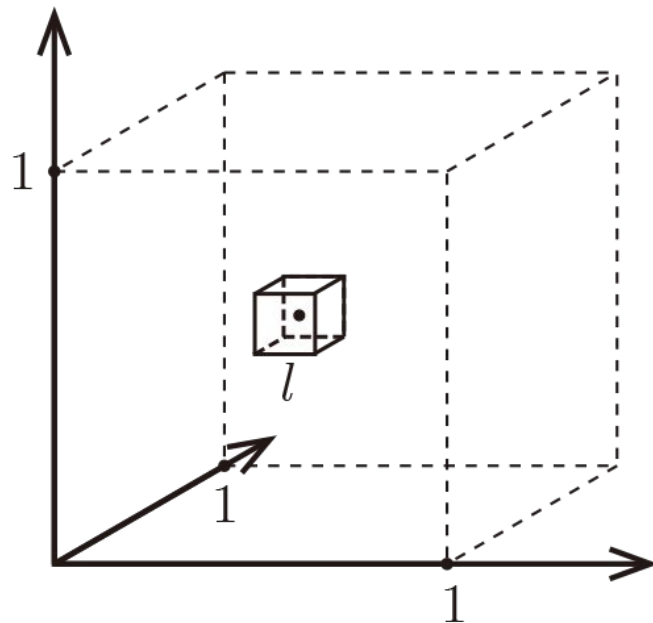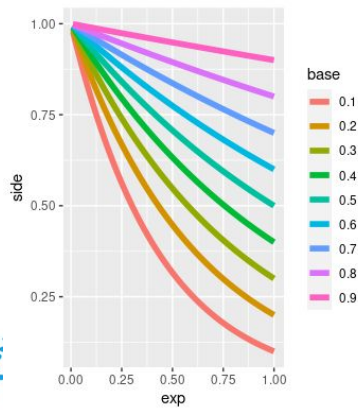


From: https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote02_kNN.html

# The curse of dimensionality

$$l = \left(\frac{k}{n}\right)^{1/d}$$

- **l**: side of the hypercube that include the k neighbours
- **n**: sample size
- **d**: n. of dimensions

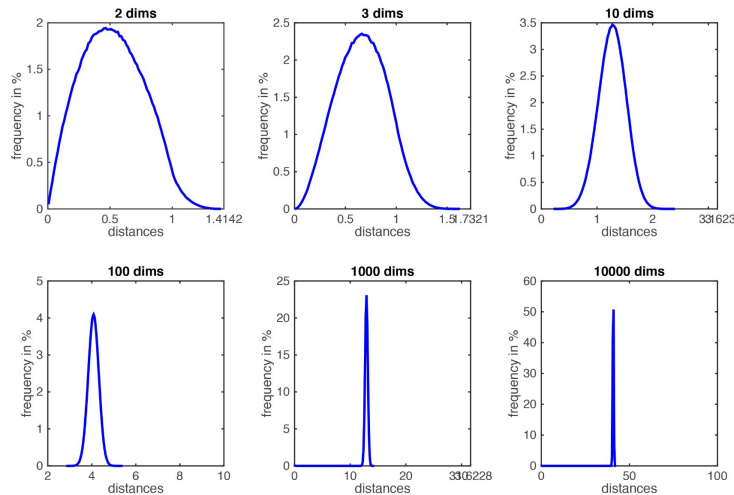- with n constant (data size), **the hypercube in which the *k* neighbors lie gets bigger as *d* increases**

# The curse of dimensionality

distributions of all pairwise distances in the neighbourhood between randomly drawn points within d-dimensional unit hypercubes: as the number of dimensions *d* grows, all distances concentrate within a very small range ("*the night where all cows are black*")



From: https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote02_kNN.html

# KNN: as lazy as it gets

- KNN is a **lazy algorithm**: each time new datapoints are added (e.g. to be predicted) pairwise distances with all existing datapoints (over all dimensions) must be calculated (there is no predictive function to be learnt)

- However:
    - when new data are available, there's no need to retrain the model (no parameters to estimate or fine-tune) → excellent for applications where data are added incrementally (e.g. on-line learning, update predictions)

# What about imputation?

- Ok, we learnt about KNN, but imputation?

- The **imputation of missing SNP genotypes** is a type of prediction (slightly "*sui generis*"), where non-missing values can be considered as training observations and missing values as the test observations

- KNNI is a more advanced imputation method compared to mean/median imputation

- ⚠️ Yet, more specialised method for genomic data have been developed

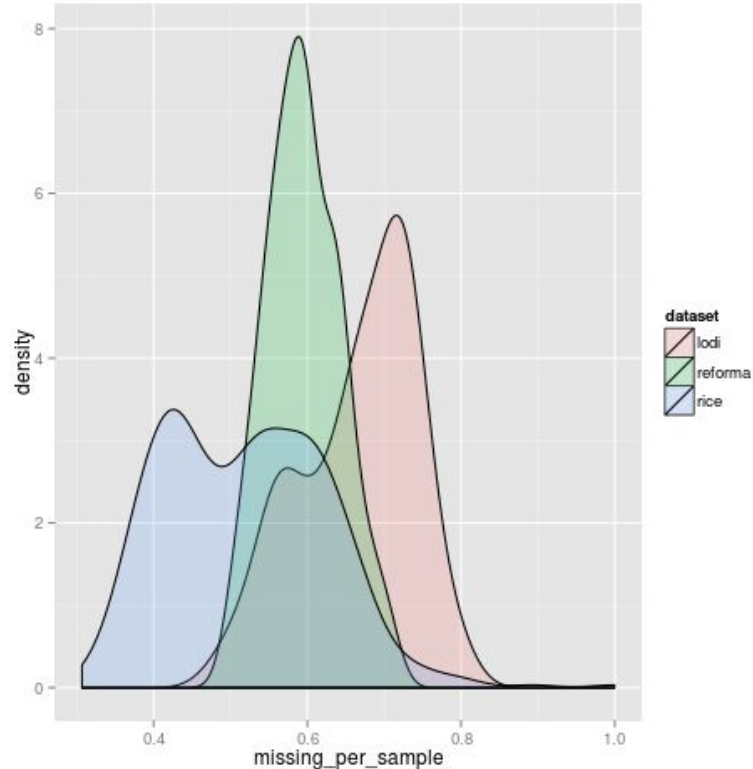# Why <span style="color:red">KNN</span> imputation?

CrossMark

## Marker imputation efficiency for genotyping-by-sequencing data in rice (*Oryza sativa*) and alfalfa (*Medicago sativa*)

Nelson Nazzicari · Filippo Biscarini ·
Paolo Cozzi · E. Charles Brummer ·
Paolo Annicchiarico
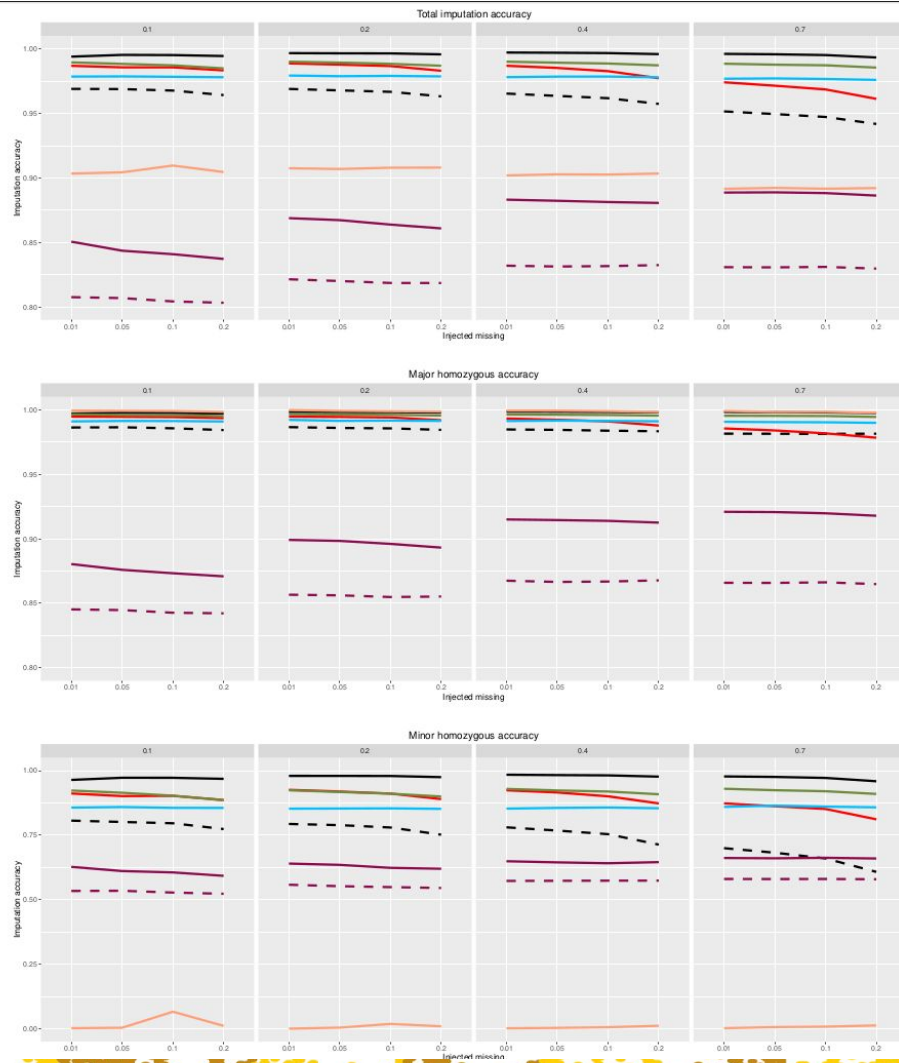
# Why KNN imputation?



- <u>GBS data</u>: > 50% missing data

- <u>rice</u>, *Oryza sativa*: ordered markers (known reference genome)

- <u>alfalfa</u>, *Medicago sativa* (lodi & reforma datasets): unordered markers (no reference genome available - at the time)
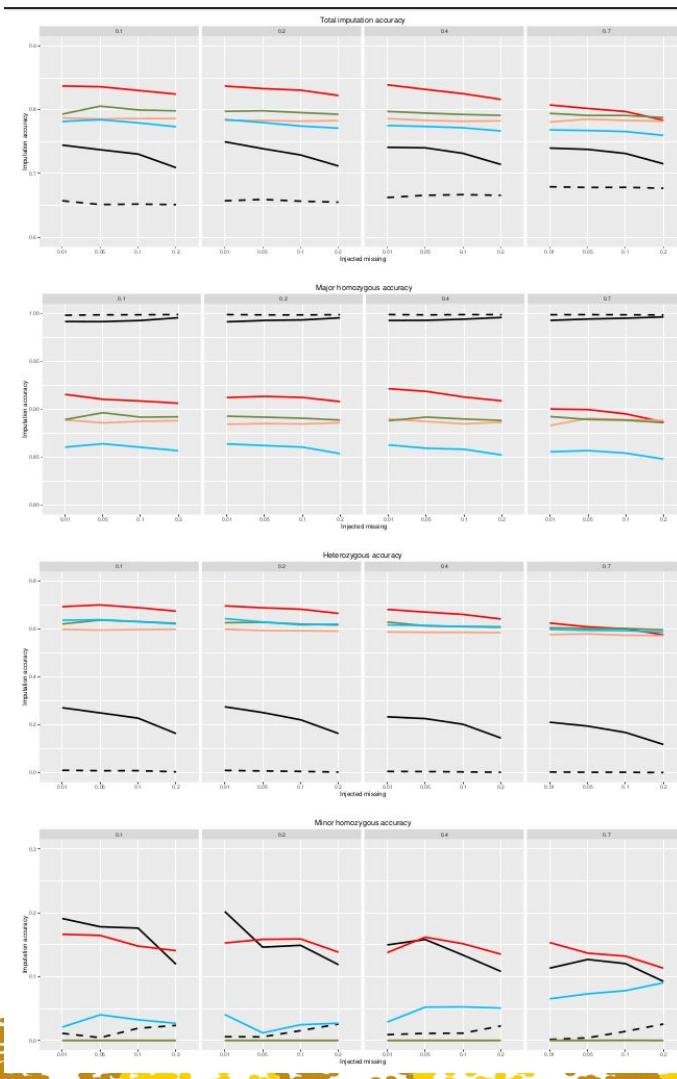
# KNN imputation - rice

- x-axis: % injected missing [1%-20%]
- panels: % allowed missing [10%-70%]
- total accuracy, AA, BB
- MNI (salmon), KNNI (red), SVDI (blue), RFI
- (green), Beagle with ordered markers (solid black), Beagle with unordered markers (dashed black), FILLIN with ordered markers (purple) FILLIN with unordered markers (dashed purple)

# KNN imputation - alfalfa

- x-axis: % injected missing [1%-20%]
- panels: % allowed missing [10%-70%]
- total accuracy, AA, AB, BB
- MNI (salmon), KNNI (red), SVDI (blue), RFI
- (green), Beagle with ordered markers (solid
  black), Beagle with unordered markers (dashed
  black), FILLIN with ordered markers (purple)
  FILLIN with unordered markers (dashed purple)

# Not only genotypes

## Integrating heterogeneous across-country data for proxy-based random forest prediction of enteric methane in dairy cattle

Enyew Negussie,[1]* Oscar González-Recio,[2] Mara Battagin,[3] Ali-Reza Bayat,[4] Tommy Boland,[5]
Yvette de Haas,[6] Aser Garcia-Rodriguez,[7] Philip C. Garnsworthy,[8] Nicolas Gengler,[9]
Michael Kreuzer,[10] Björn Kuhla,[11] Jan Lassen,[12] Nico Peiren,[13] Marcin Pszczola,[14]
Angela Schwarm,[15] Hélène Soyeurt,[9] Amélie Vanlierde,[16] Tianhai Yan,[17] and Filippo Biscarini[18]

**NEXT LECTURE**

Genotype imputation with KNN: a demonstration
(R code)