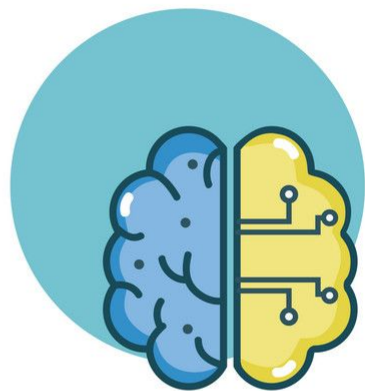


INTRODUCTION TO MACHINE LEARNING

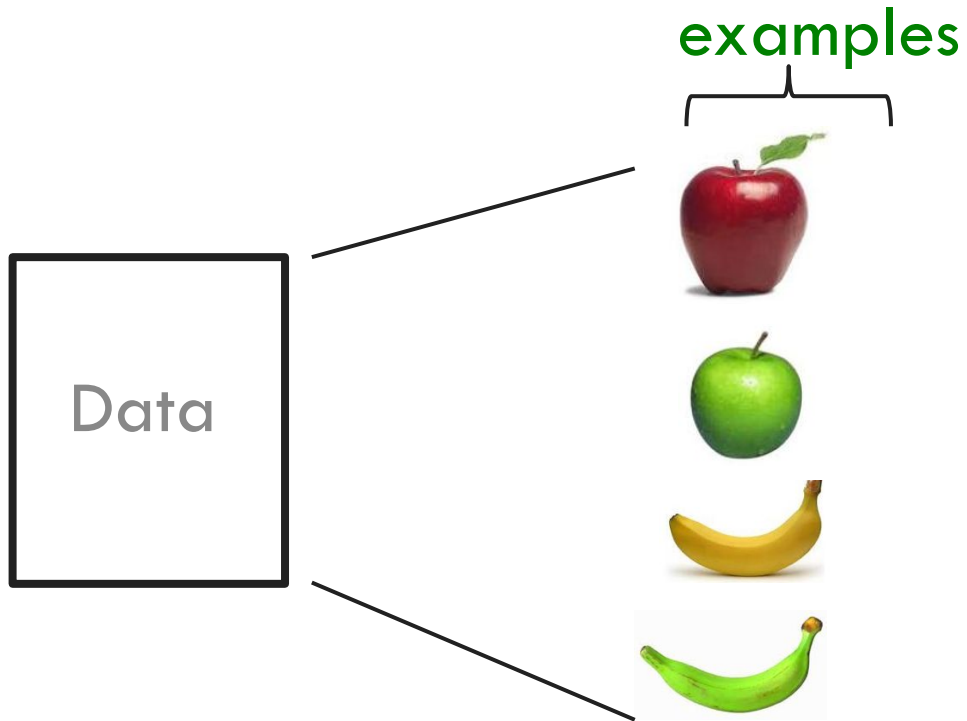
K-NEAREST NEIGHBOR



Elisa Ricci



RECAP: DATA



RECAP: REPRESENTING EXAMPLES

examples



What is an example?

How is it represented?

RECAP: FEATURES

examples



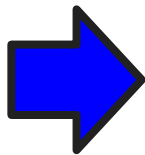
features

$f_1, f_2, f_3, \dots, f_n$

$f_1, f_2, f_3, \dots, f_n$

$f_1, f_2, f_3, \dots, f_n$

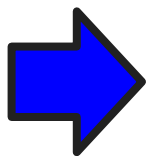
$f_1, f_2, f_3, \dots, f_n$



How our algorithms
actually “view” the data

RECAP: FEATURES

examples



features

red, round, leaf, 3oz, ...

green, round, no leaf, 4oz, ...

yellow, curved, no leaf, 8oz, ...

green, curved, no leaf, 7oz, ...

How our algorithms actually “view” the data

Features are the questions we can ask about the examples

Features in general are represented with vectors

APPLES VS. BANANAS

Weight	Color	Label
4	Red	Apple
5	Yellow	Apple
6	Yellow	Banana
3	Red	Apple
7	Yellow	Banana
8	Yellow	Banana
6	Yellow	Apple

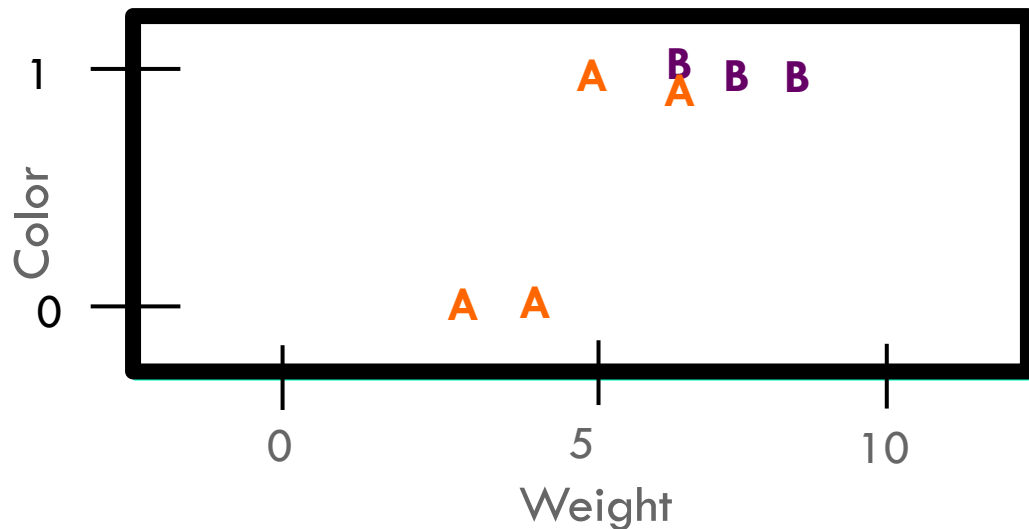


Can we visualize this data?

APPLES VS. BANANAS

We can turn features into numerical values

Weight	Color	Label
4	0	Apple
5	1	Apple
6	1	Banana
3	0	Apple
7	1	Banana
8	1	Banana
6	1	Apple



Examples are points in an n -dimensional space where n is the number of features

EXAMPLES IN A FEATURE SPACE

We can imagine to have 3 categories of interest

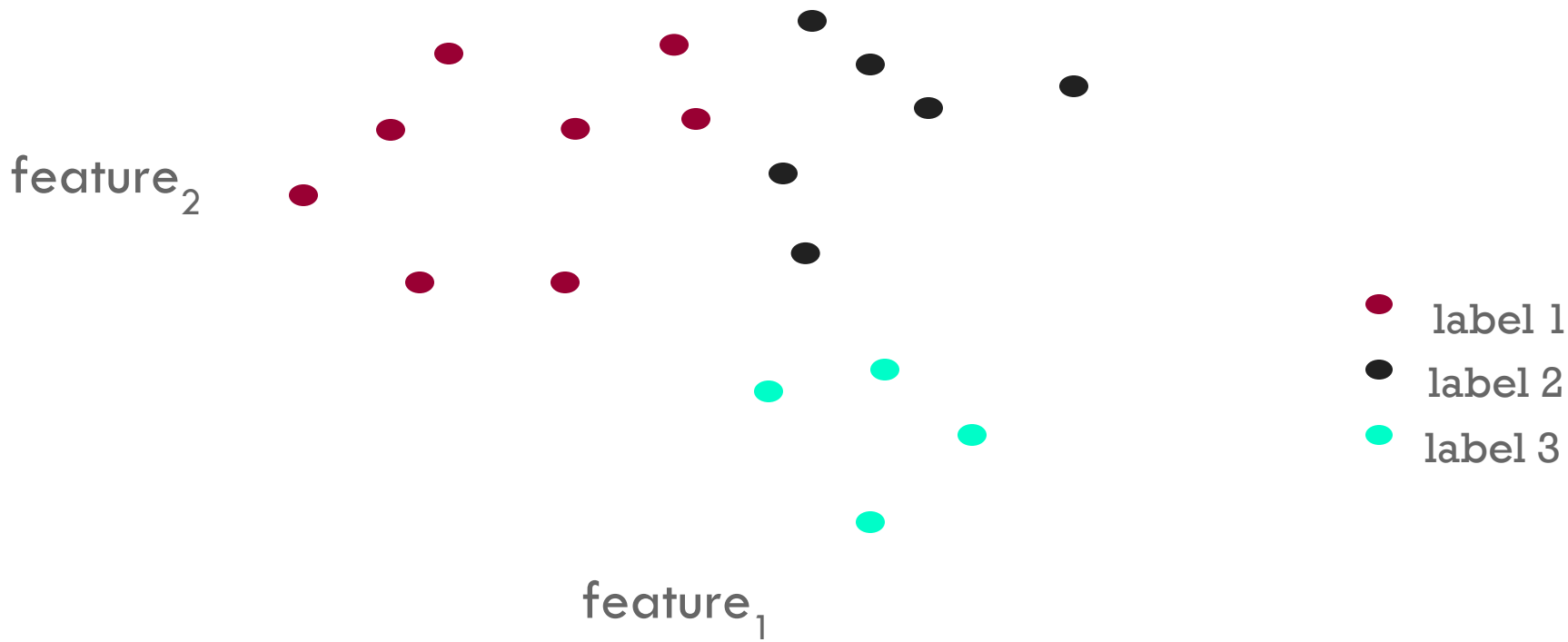


feature₂

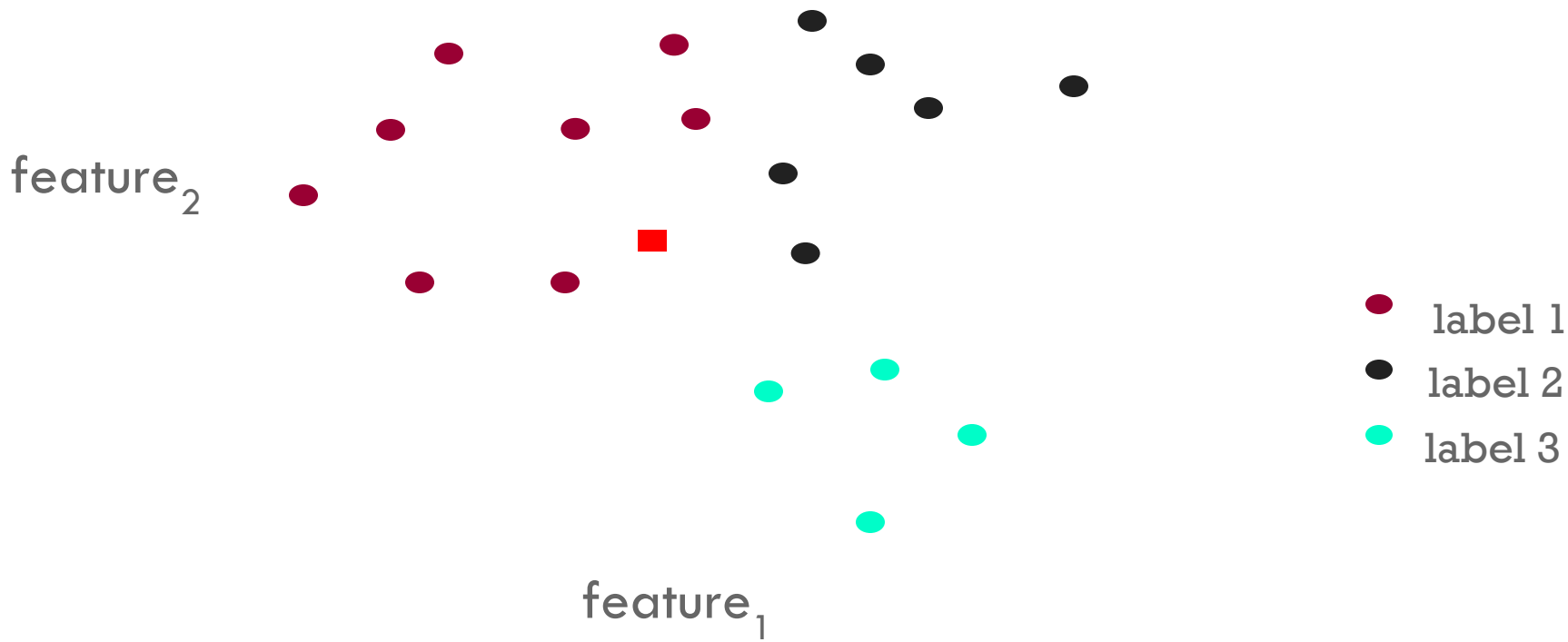
feature₁

- label 1
- label 2
- label 3

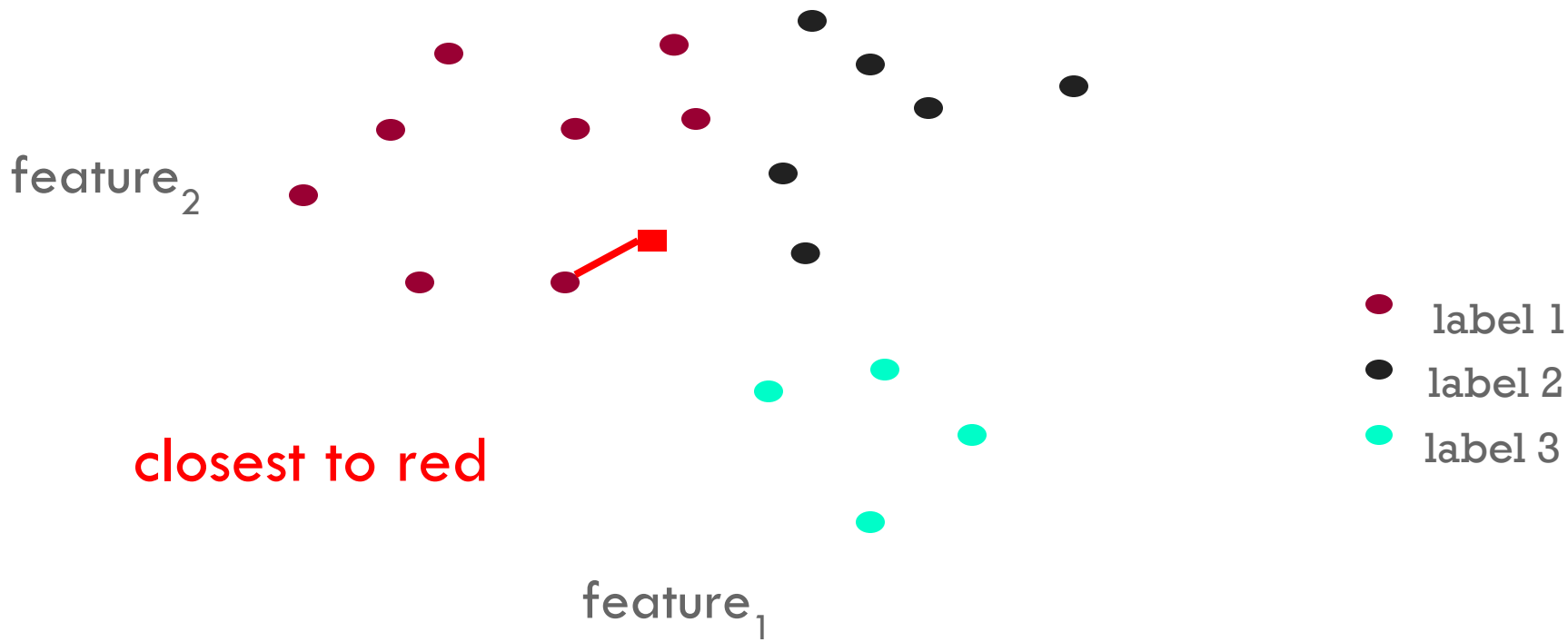
TRAINING SET



TEST EXAMPLE: WHICH CLASS?



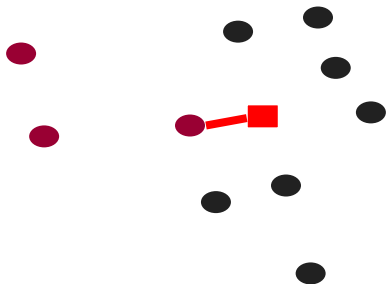
TEST EXAMPLE: WHICH CLASS?



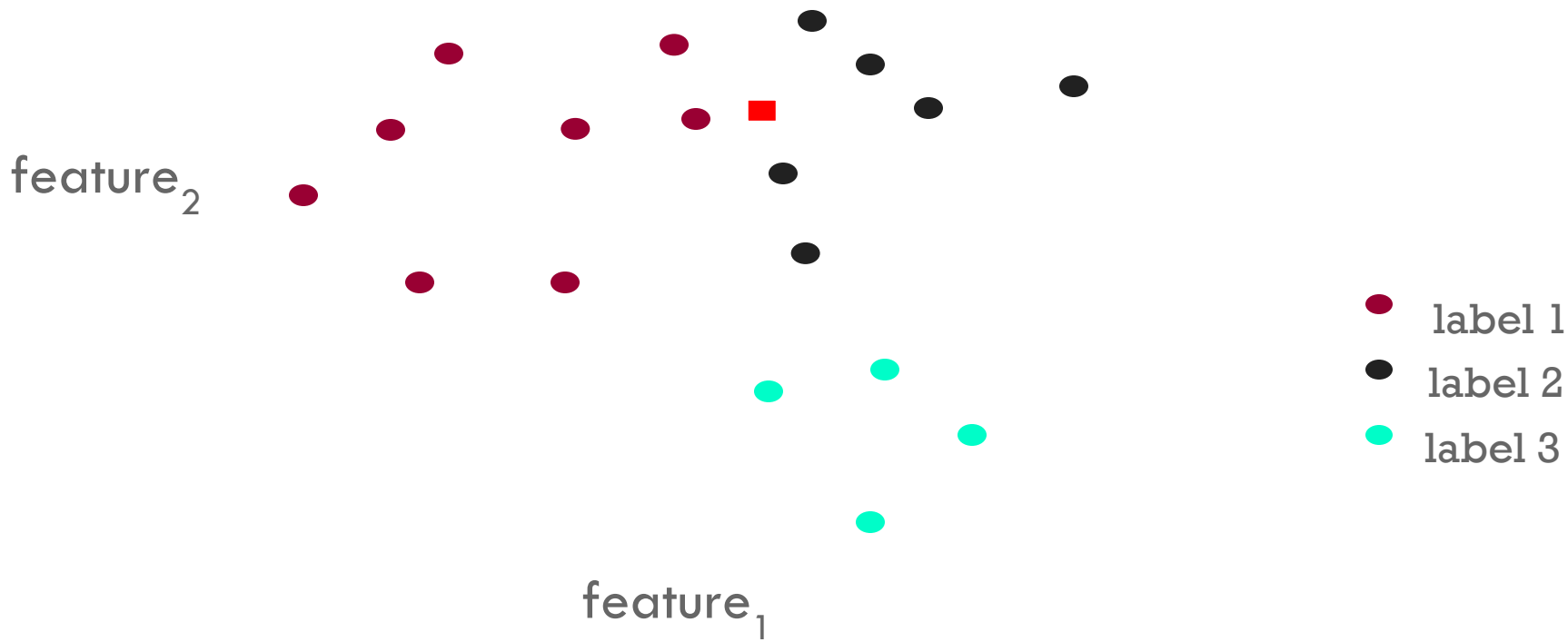
A CLASSIFICATION ALGORITHM?

To classify an example d :

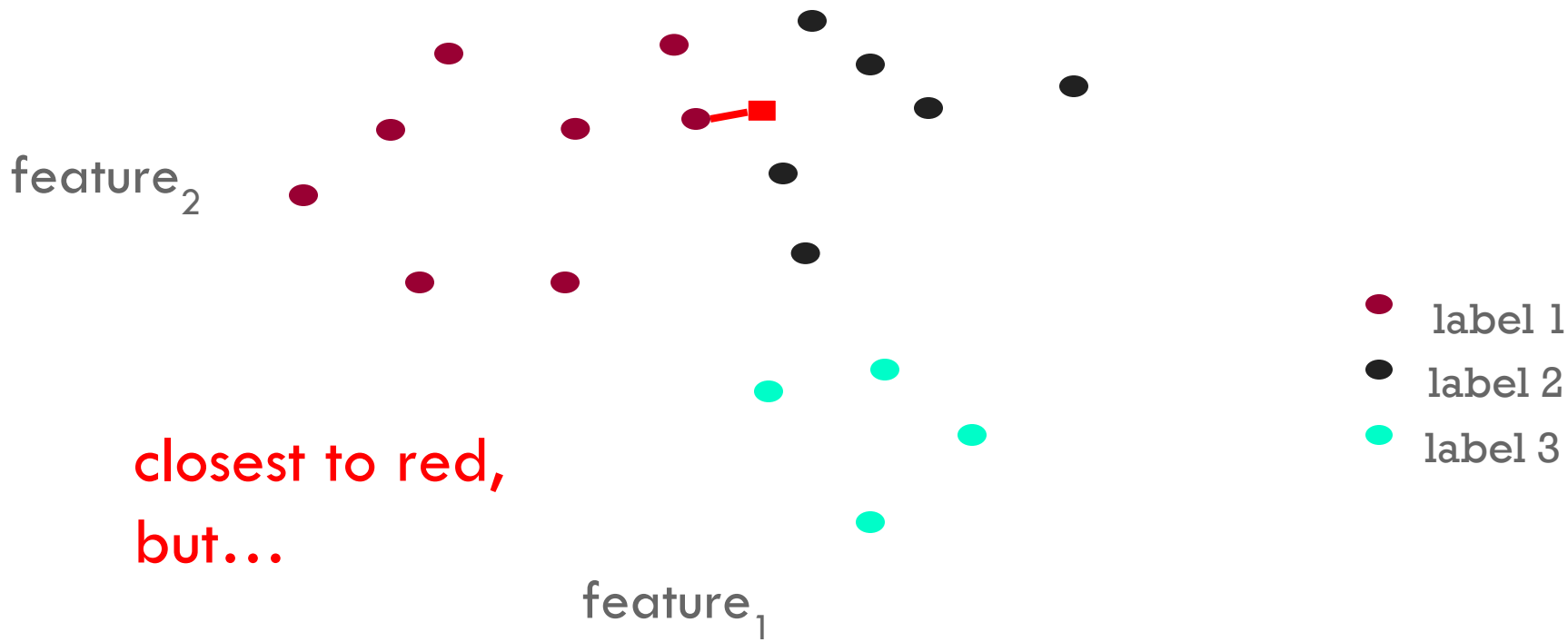
Label d with the label of the closest example to d in the training set



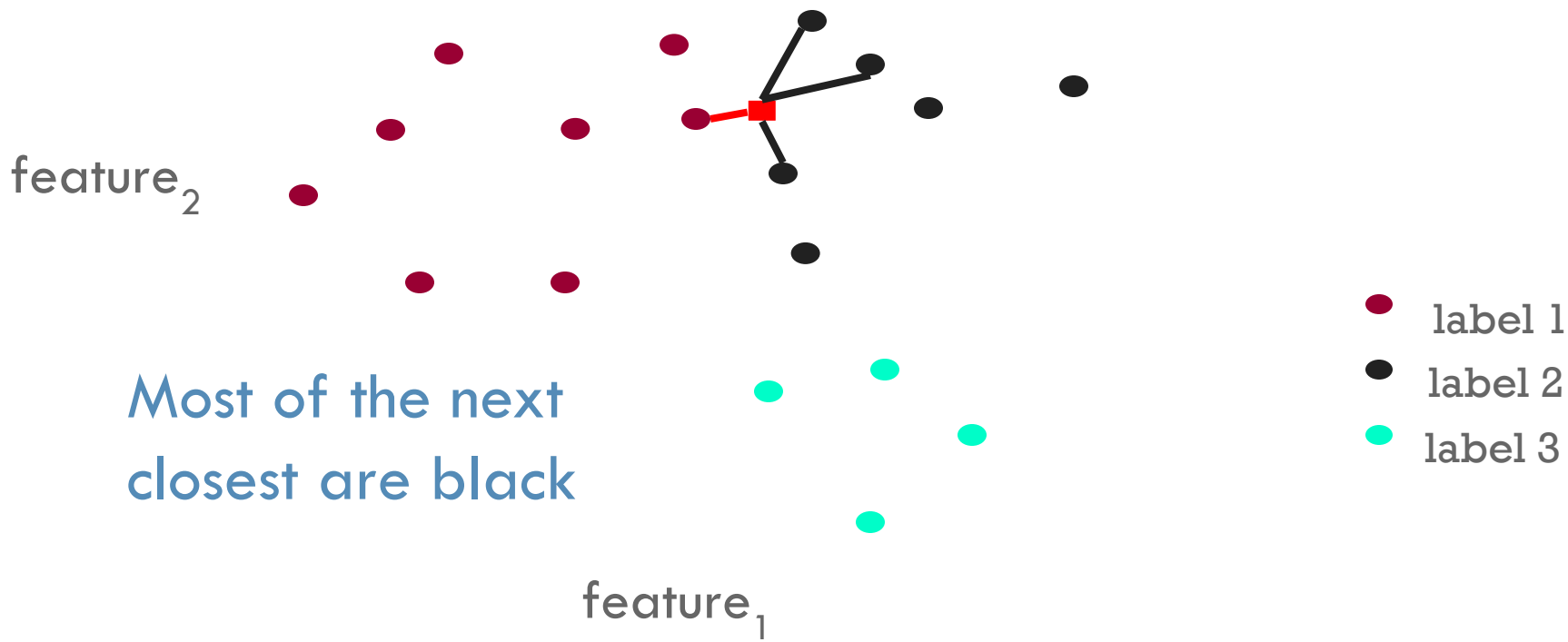
WHAT ABOUT THIS EXAMPLE?



WHAT ABOUT THIS EXAMPLE?



WHAT ABOUT THIS EXAMPLE?



K-NEAREST NEIGHBOR (K-NN)

To classify an example d :

- Find k nearest neighbors of d
- Choose as the label the **majority label** within the k nearest neighbors

Robust to noisy data by considering k -nearest neighbors

K-NEAREST NEIGHBOR (K-NN)

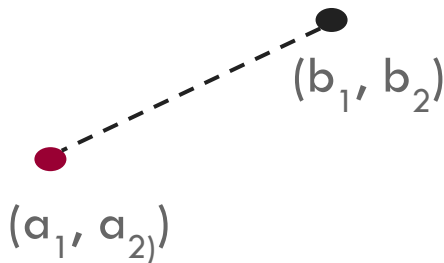
To classify an example d :

- Find k *nearest* neighbors of d
- Choose as the label the *majority label* within the k nearest neighbors

How do we measure “nearest”?

EUCLIDEAN DISTANCE

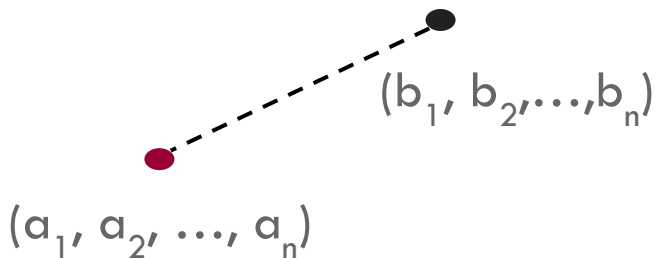
In two dimensions, how do we compute the distance?



$$D(a, b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$$

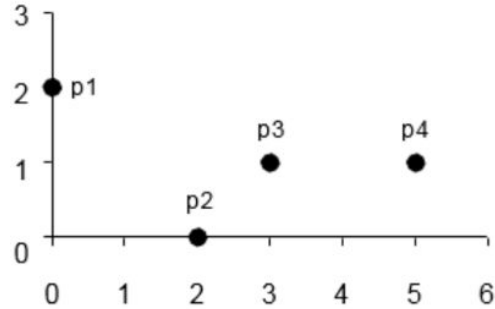
EUCLIDEAN DISTANCE

In n-dimensions, how do we compute the distance?



$$D(a, b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots + (a_n - b_n)^2}$$

EUCLIDEAN DISTANCE: EXAMPLE



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

EUCLIDEAN DISTANCE

- Euclidean Distance makes sense when different features are comparable (e.g. each is variable measured in the same units).
 - For instance, if the measurements are different, say length and weight, it is not clear.
- **Standardization:** When features are not comparable we can standardize them by dividing by their corresponding standard deviation. This makes them all equally important.

STANDARDIZATION & SCALING

- **Standardization or Z-score normalization**

- Rescale the data so that the mean is 0 and the standard deviation from the mean (standard scores) is 1

$$x_{norm} = \frac{x - \mu}{\sigma}$$

μ is mean, σ is a standard deviation from the mean
(standard score)

- **Min-Max scaling**

- Scale the data to a fixed range – between 0 and 1

$$x_{morm} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

DISTANCE AND SIMILARITY

- Measuring distance/similarity is a domain-specific problem and there are many, many different variations
- Similarity
 - Numerical measure of how alike two data objects are.
 - Is higher when objects are more alike.
 - Often falls in the range $[0,1]$
- Distance
 - Numerical measure of how different are two data objects
 - Lower when objects are more alike
 - Minimum dissimilarity is often 0
 - Upper limit varies

MINKOWSKI DISTANCE

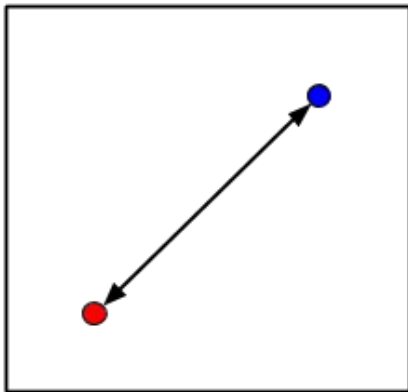
- Minkowski Distance is a generalization of Euclidean Distance

$$D(a, b) = \sum_{k=1}^p |a_k - b_k|^r$$

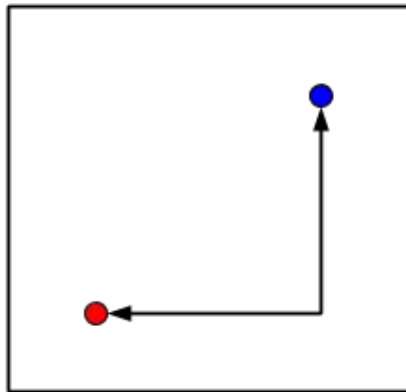
- r is a parameter, p is the number of dimensions
- $r = 1$. City block (Manhattan, L1 norm) distance.
- $r = 2$. Euclidean distance
- $r \rightarrow \infty$. “supremum” (L_∞ norm) distance.
 - This is the maximum difference between any component of the vectors

DISTANCES

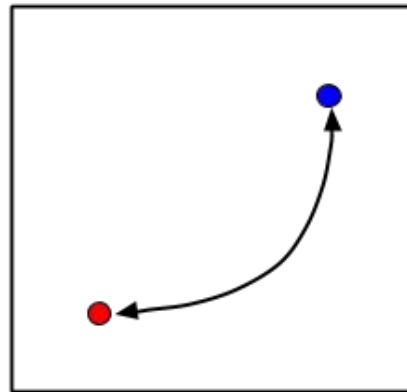
Euclidean



Manhattan

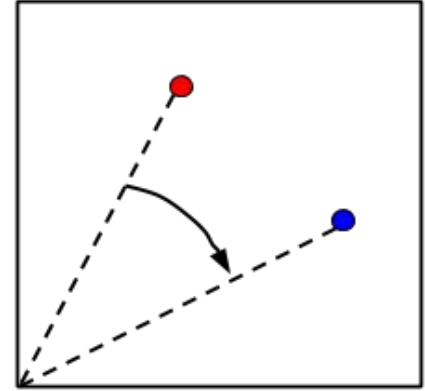


Minkowski



COSINE SIMILARITY

Cosine Similarity



Given two document vectors:

$$\cos(d_1, d_2) = (d_1 \cdot d_2) / (||d_1|| ||d_2||),$$

Where \cdot indicates vector dot product and $||d||$ is the length of vector d .

Example:

$$d_1 = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0$$

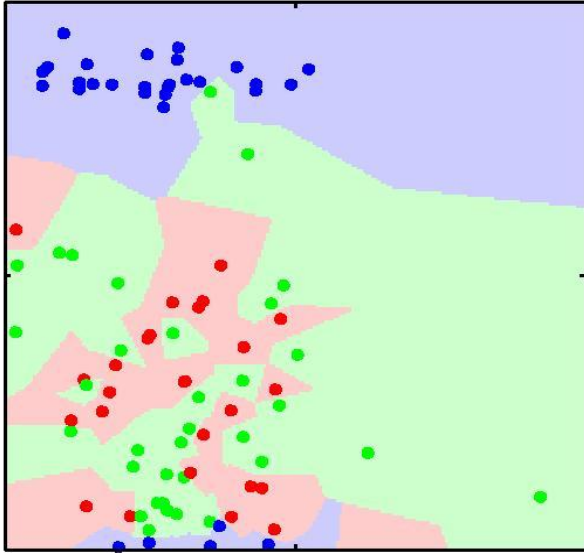
$$d_2 = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 2$$

$$d_1 \cdot d_2 = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$$

$$||d_1|| = (3*3 + 2*2 + 0*0 + 5*5 + 0*0 + 0*0 + 0*0 + 2*2 + 0*0 + 0*0)^{0.5} = (42)^{0.5} = 6.481$$

$$||d_2|| = (1*1 + 0*0 + 0*0 + 0*0 + 0*0 + 0*0 + 0*0 + 1*1 + 0*0 + 2*2)^{0.5} = (6)^{0.5} = 2.245$$

$$\cos(d_1, d_2) = .3150$$

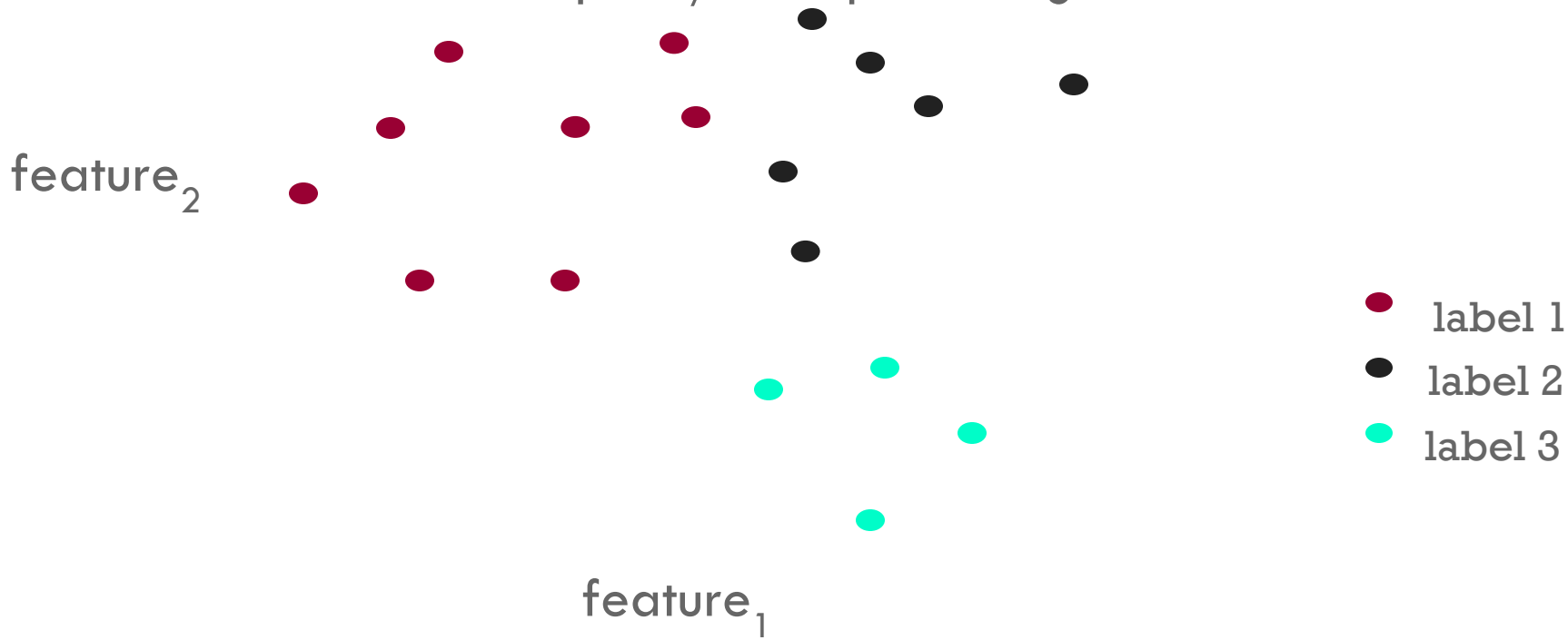


DECISION BOUNDARIES

HYPER-PARAMETERS

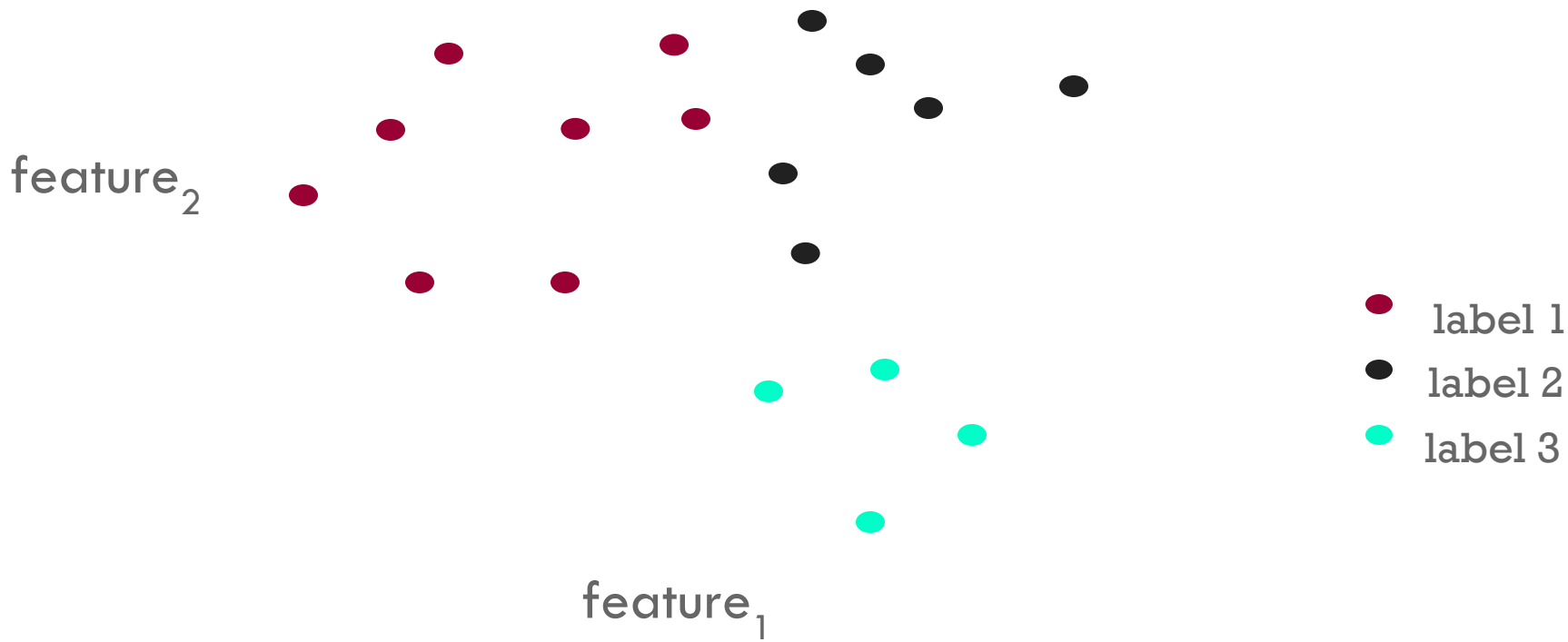
DECISION BOUNDARIES

The **decision boundaries** are places in the features space where the classification of a point/example changes



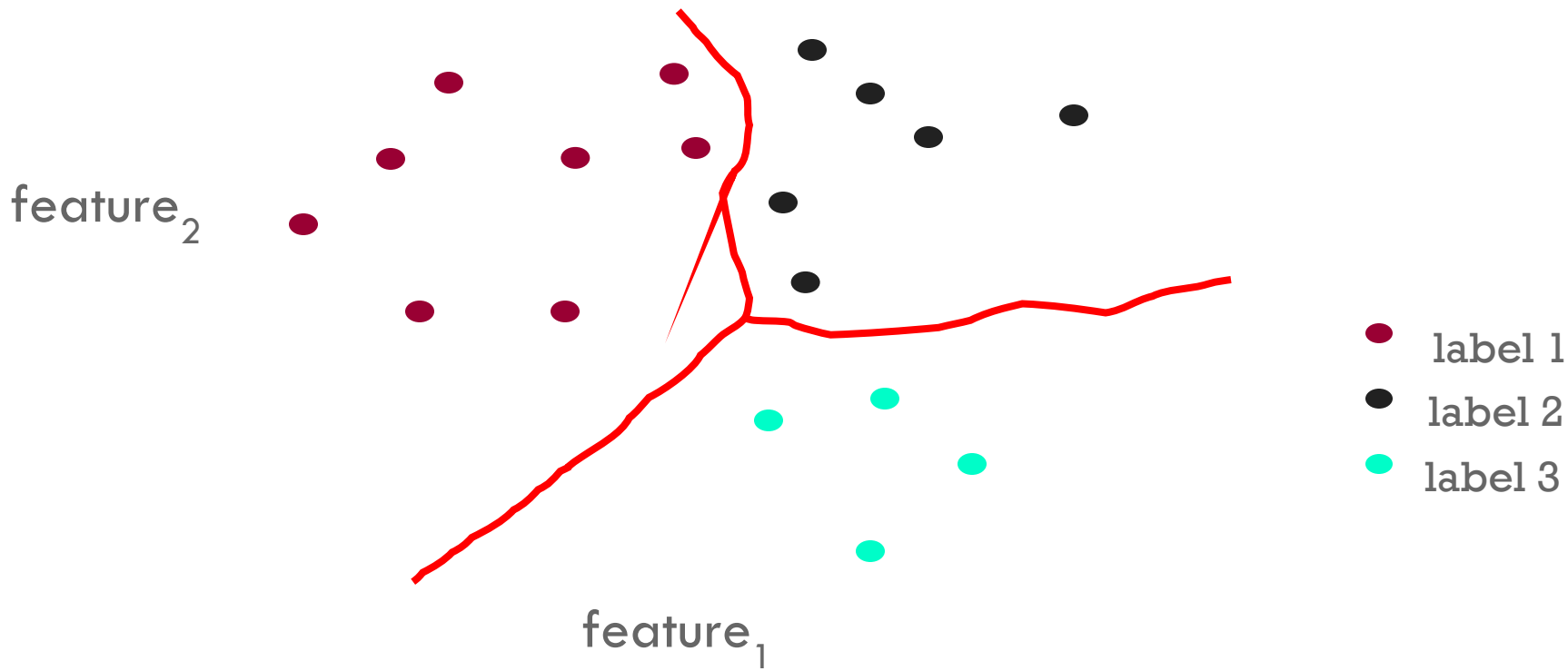
DECISION BOUNDARIES

Where are the decision boundaries for k-NN?



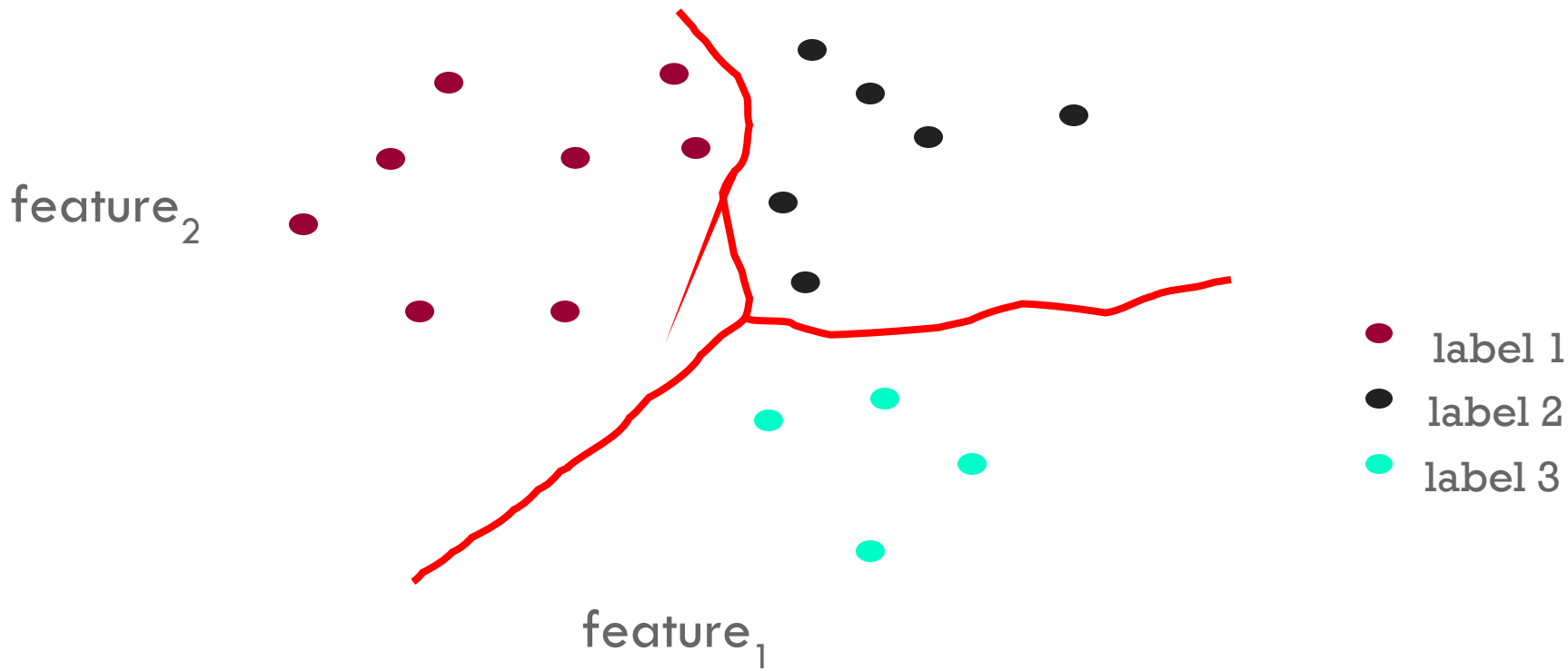
DECISION BOUNDARIES

Where are the decision boundaries for k-NN?



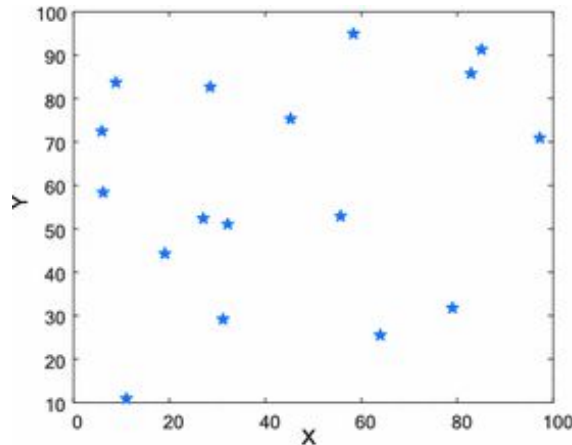
DECISION BOUNDARIES

k-NN gives **locally** defined decision boundaries between classes

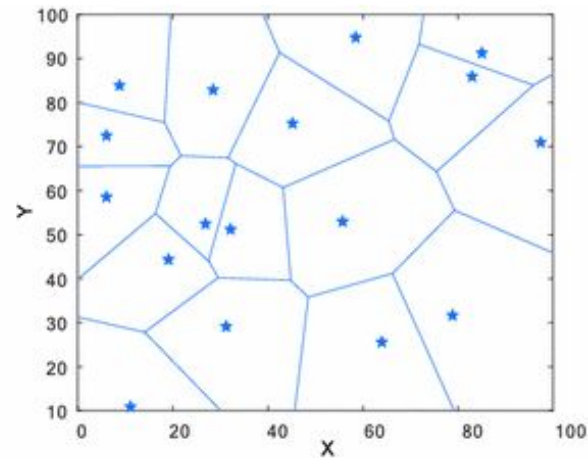


VORONOI DIAGRAM

- Describes the areas that are nearest to any given point, given a set of data.
- Each line segment is equidistant between two points



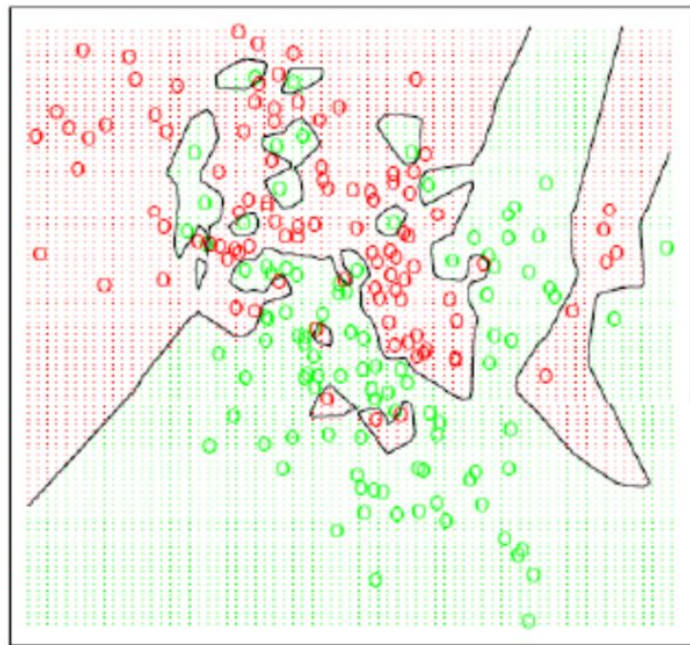
(a) Initial points distribution



(b) Voronoi diagram by initial points

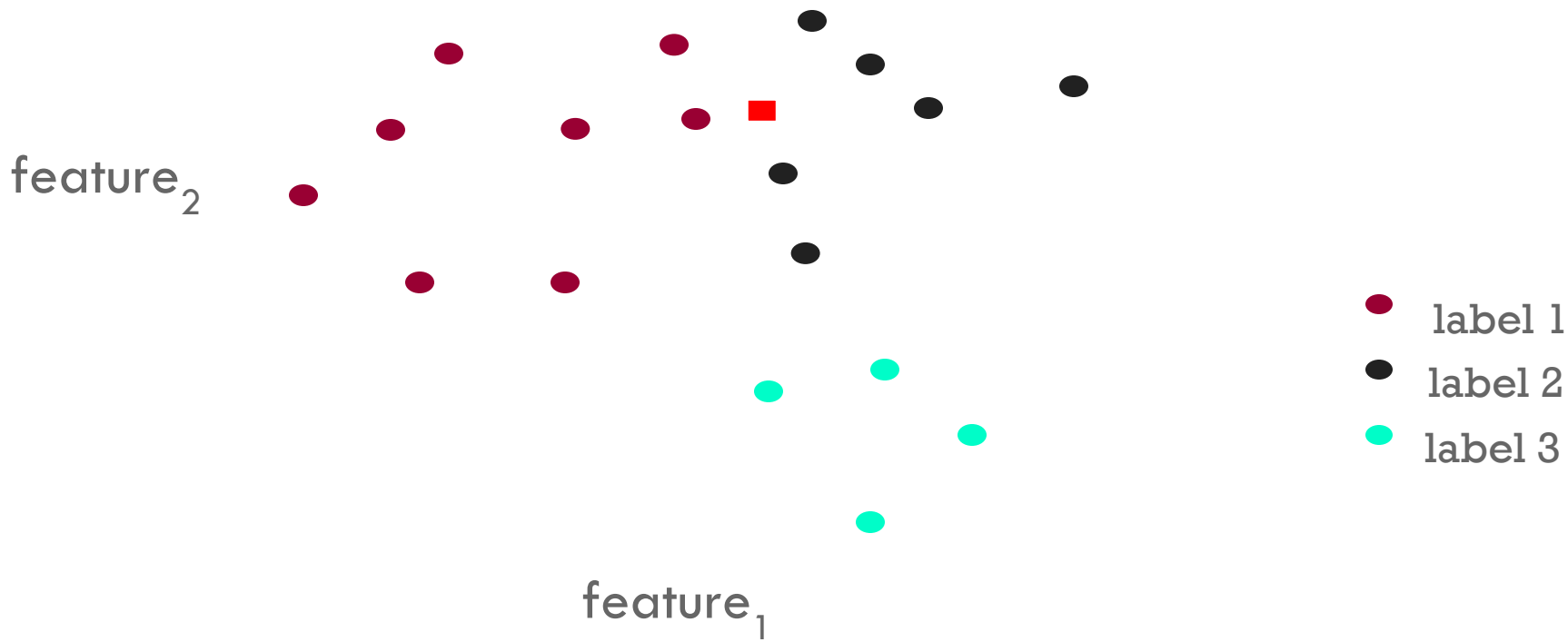
DECISION BOUNDARIES

- k-NN algorithm does not explicitly compute decision boundaries.
- The decision boundaries form a subset of the Voronoi diagram for the training data.
- The more examples that are stored, the more complex the decision boundaries can become



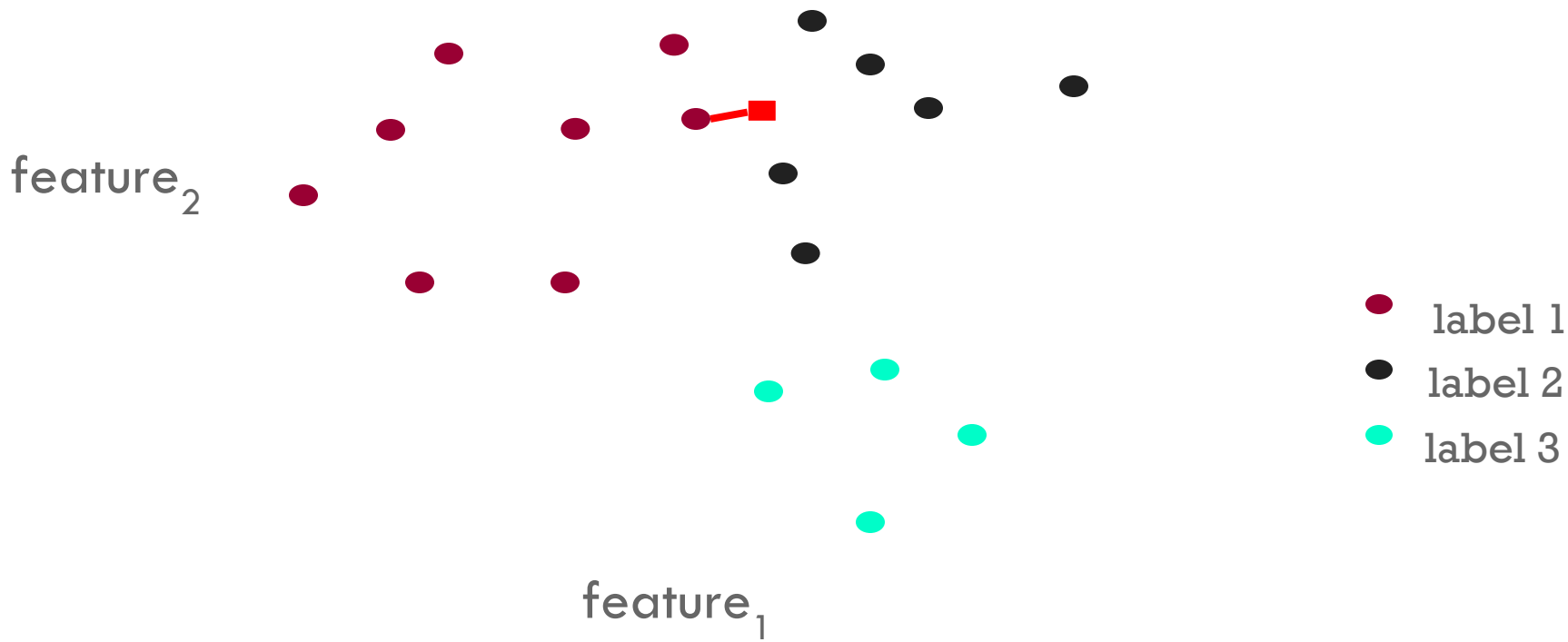
CHOOSING K

What is the label with $k = 1$?



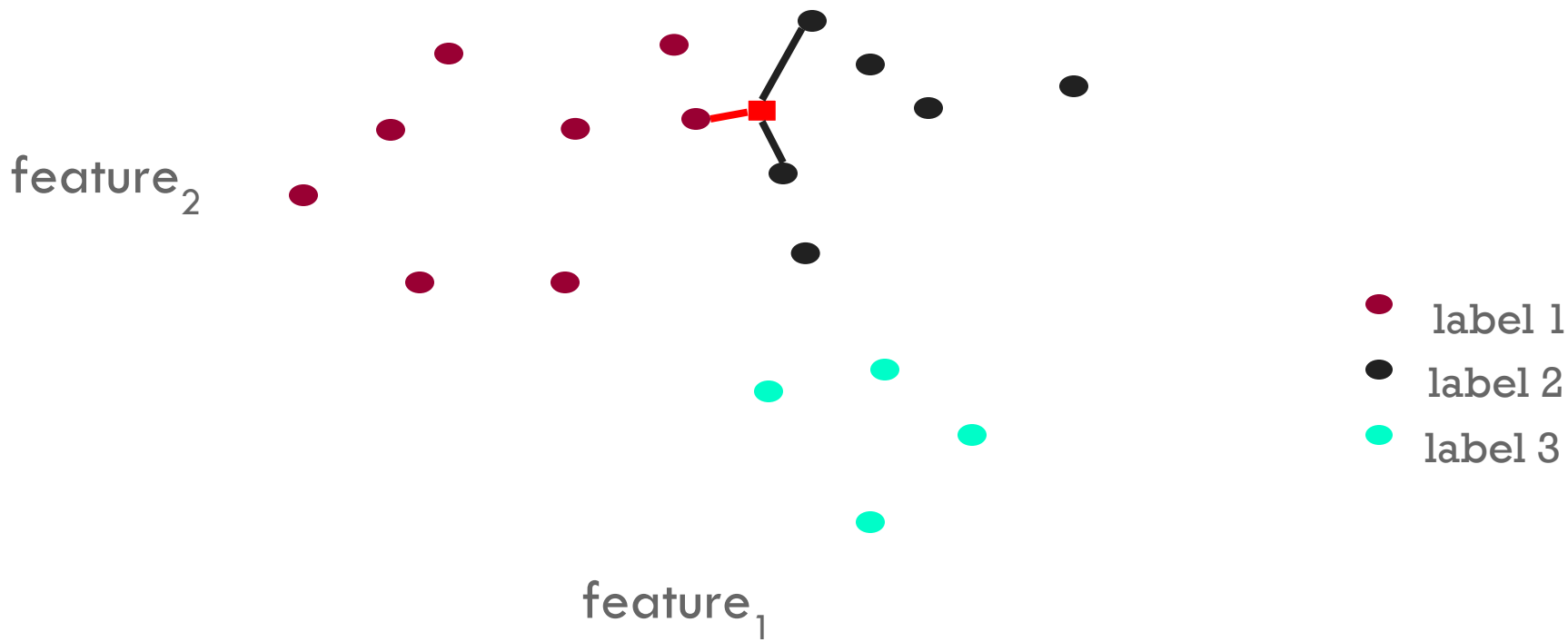
CHOOSING K

We choose red. Do you agree?



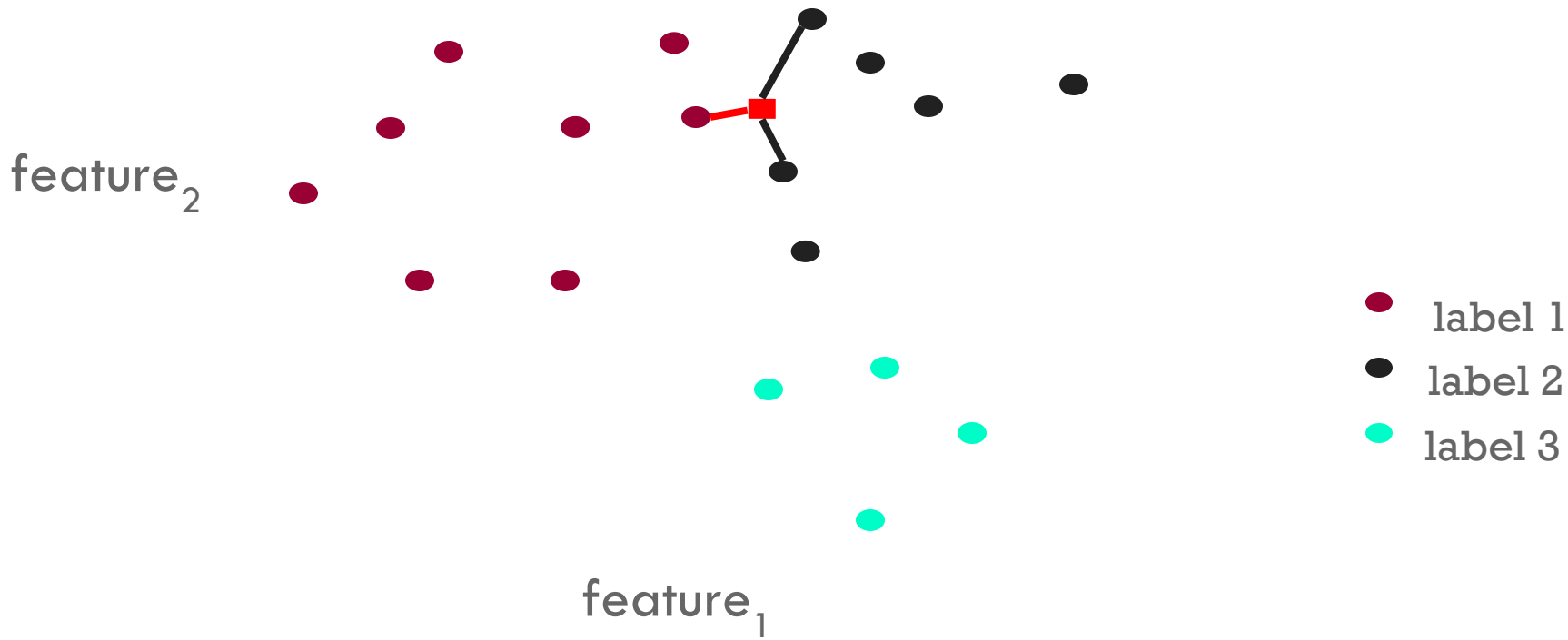
CHOOSING K

What is the label with $k = 3$?



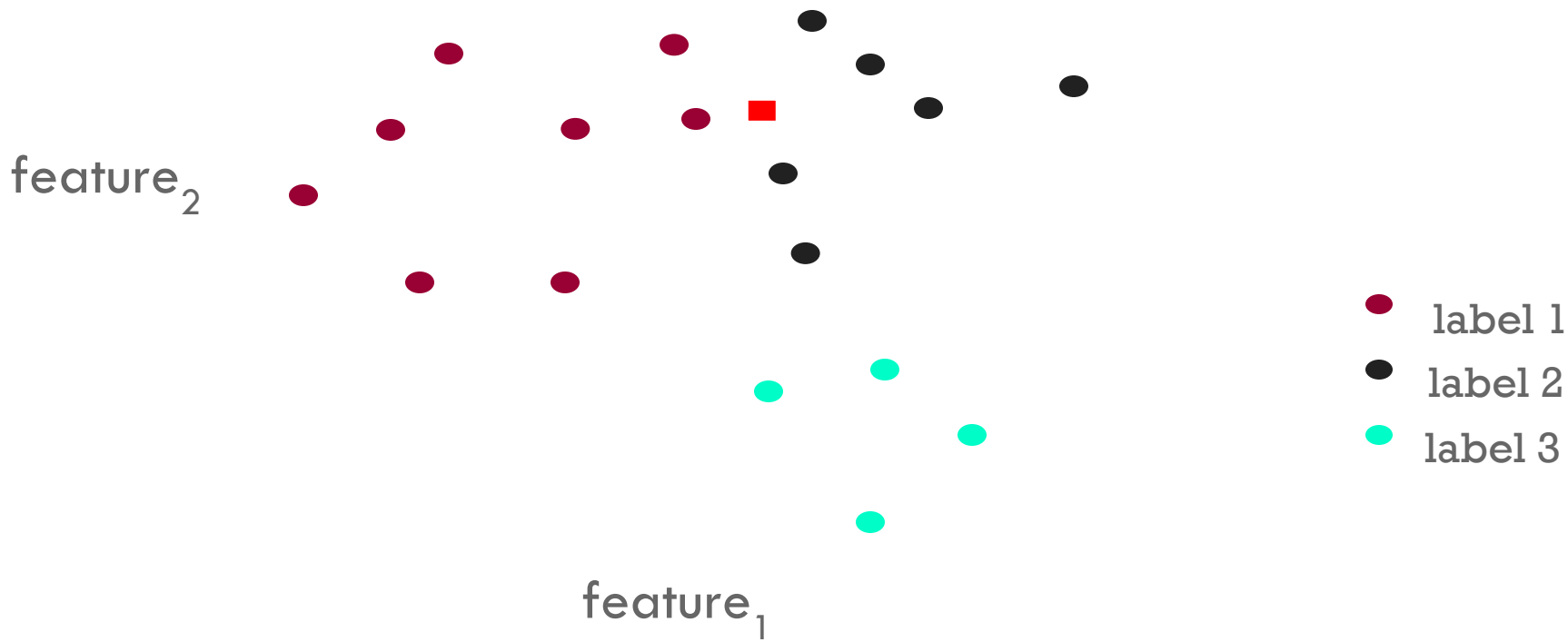
CHOOSING K

We choose black. Do you agree?



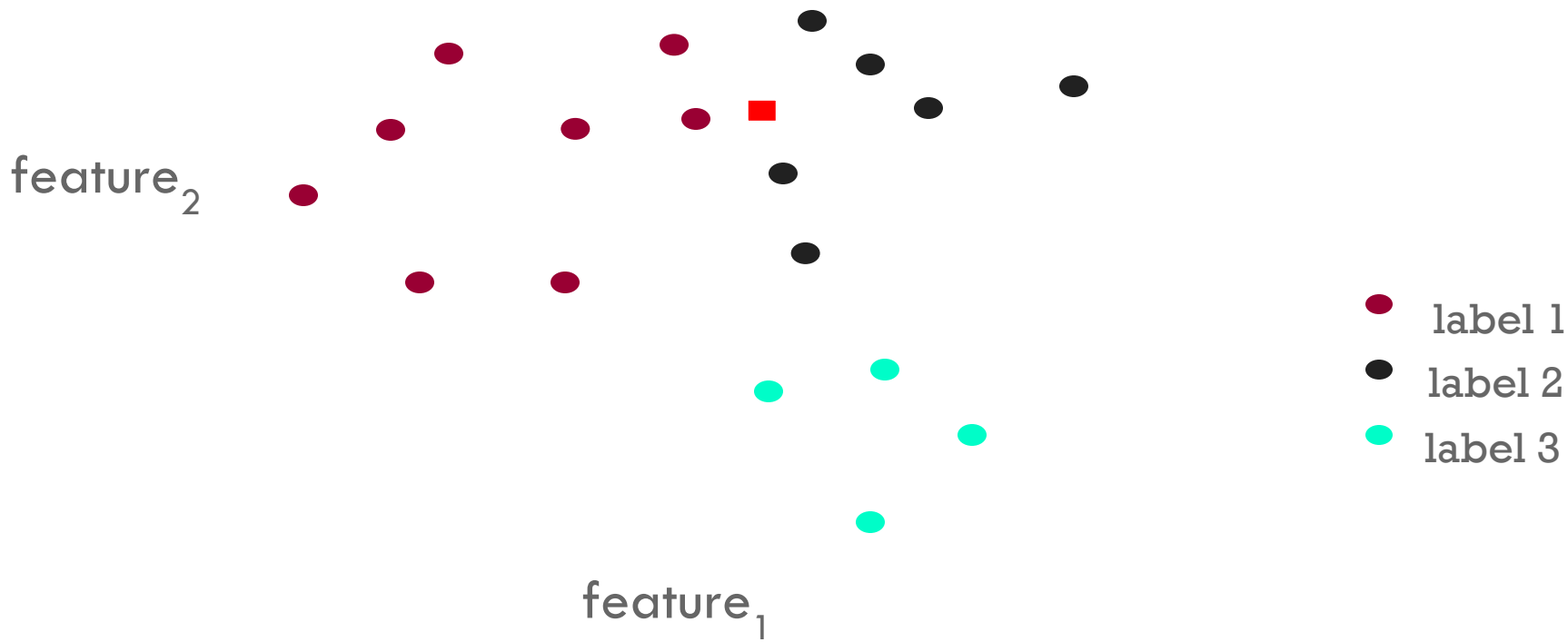
CHOOSING K

What is the label with $k = 20$?

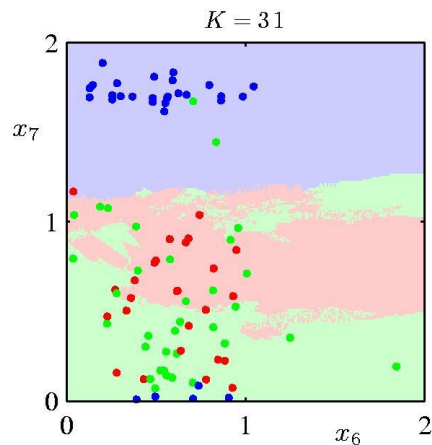
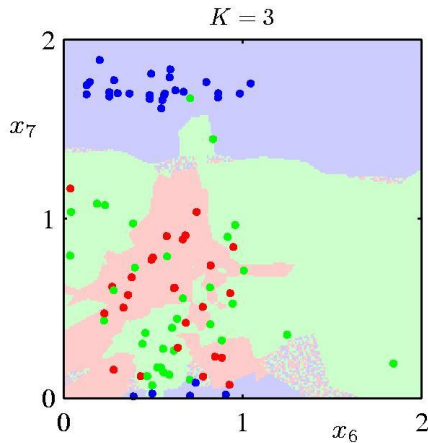
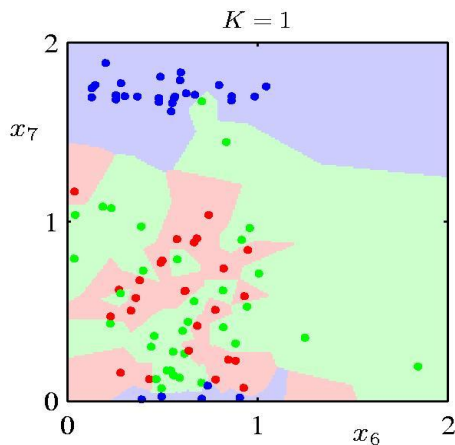


CHOOSING K

What is the label with $k = 20$? Red! (class with max number of datapoints, weird)



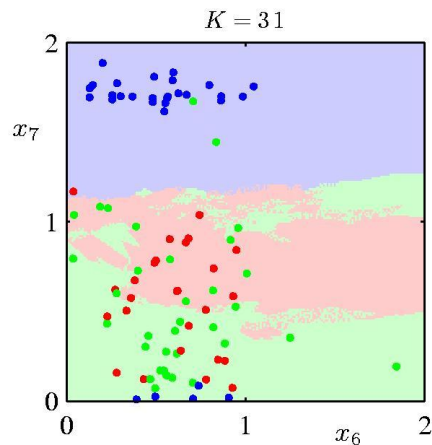
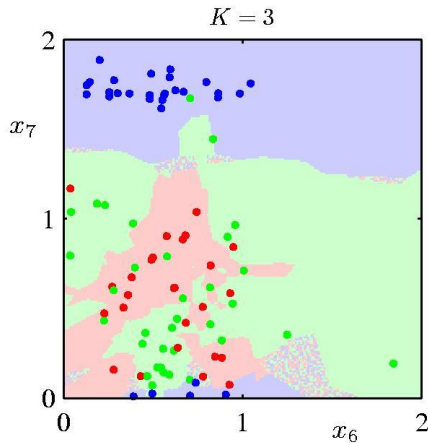
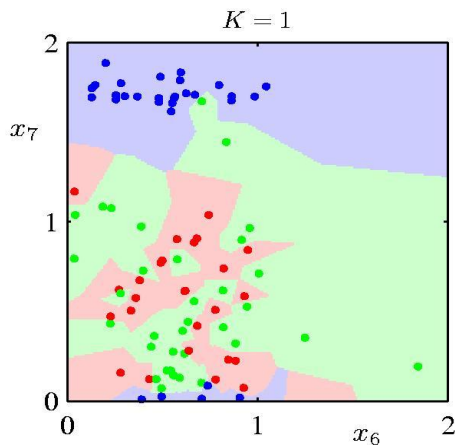
THE IMPACT OF K



What is the role of k ?

How does it relate to overfitting and underfitting?

THE IMPACT OF K



Larger k produces smoother boundary effect
When $k=N$, always predict the majority class

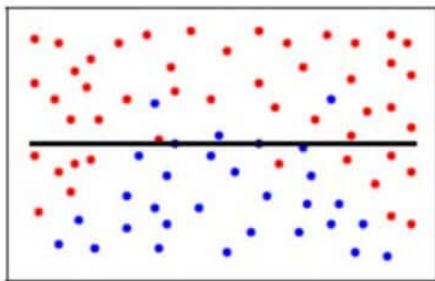
RECAP: UNDERFITTING OVERFITTING

- The factors determining how well a machine learning algorithm will perform are its ability to:
 - Make the **training error small**.
 - Make the **gap between training and test error small**.
- These two factors correspond to the cases of **underfitting** and **overfitting**.
- **Underfitting** occurs when the model is not able to obtain a sufficiently low error value on the training set.
- **Overfitting** occurs when the gap between the training error and test error is too large.

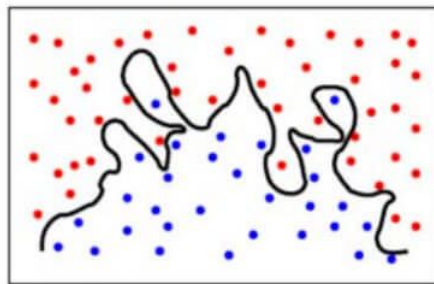
RECAP: UNDERFITTING OVERFITTING

- Classification

Underfitting



Overfitting



K-NEAREST NEIGHBOR (K-NN)

To classify an example d :

- Find k nearest neighbors of d
- Choose as the class the majority class within the k nearest neighbors

How do we choose k ?

HOW TO PICK K

- Common heuristics:
 - often 3, 5, 7
 - choose an odd number to avoid ties
- Use validation set
- Use cross-validation
- Rule of thumb is $k < \sqrt{n}$, where n is the number of training examples

K-NEAREST NEIGHBOR (K-NN)

To classify an example d :

- Find k nearest neighbors of d
- Choose as the class the majority class within the k nearest neighbors

Any other idea?

K-NN VARIATIONS

Instead of k nearest neighbors, count majority from all examples within a fixed distance

Weighted k -NN:

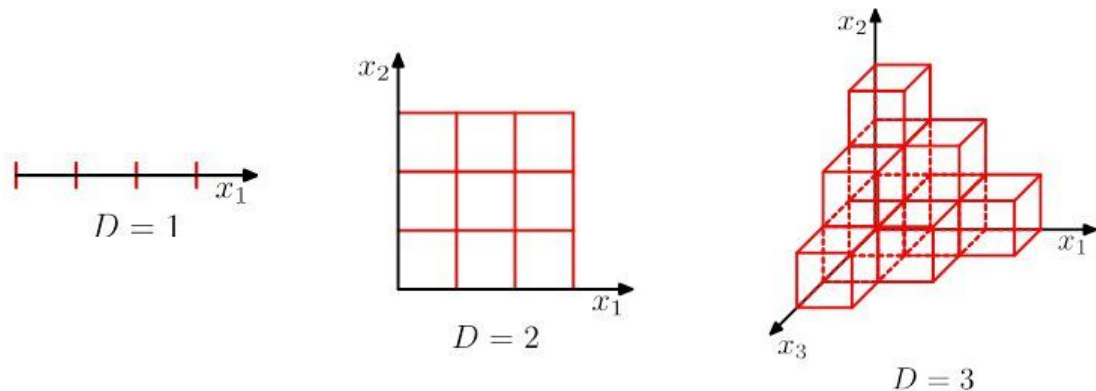
- Right now, all examples are treated equally
- Weight the “vote” of the examples, so that closer examples have more vote/weight
- Often use some sort of exponential decay

LAZY LEARNER VS EAGER LEARNER

- k-NN is belongs to the class of **lazy learning** algorithms
- **Lazy learning**: Simply stores training data (or perform limited processing) and operates when it is given a test example
- **Eager learning** (e.g. Decision Trees, SVMs): Given a training set, constructs a classification model before receiving new test data to classify
- **Lazy learners**: less time in training but more time in predicting

CURSE OF DIMENSIONALITY

- **Curse of dimensionality:** In high dimensions almost all points are far away from each other
- The size of the data space grows exponentially with the number of dimensions.
- This means that the size of the data set must also grow exponentially in order to keep the same density.



CURSE OF DIMENSIONALITY

- The success of k-NN is very dependent on having a **dense data set**.
- Every machine learning algorithm needs a dense data set for accurate prediction
- What makes k-nearest neighbors special?
- k-NN requires a point to be close in **every single dimension**.
- Some algorithms can create models based on single dimensions, and only need points to be close together along that axis.
- k-NN doesn't work that way. It needs all points to be close along every axis in the data space.

CURSE OF DIMENSIONALITY: SUMMARY

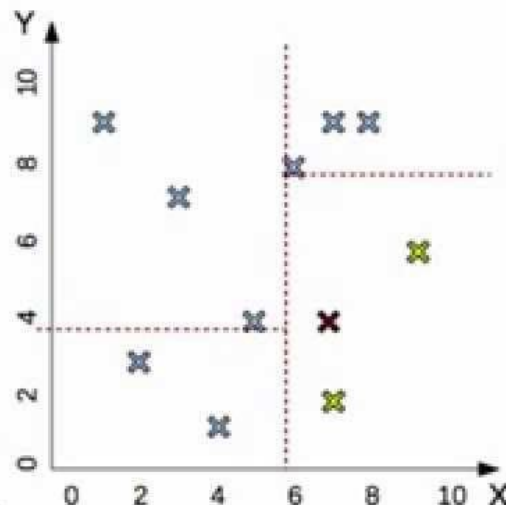
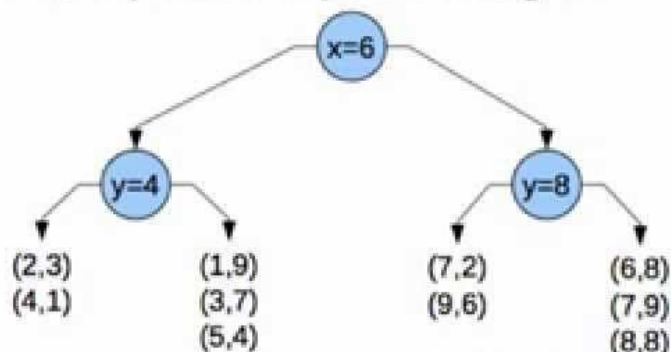
- The problem is fundamentally that there isn't enough data available for the number of dimensions.
- As the number of dimensions increases the size of the data space increases, and the amount of data needed to maintain density also increases.
- Without dramatic increases in the size of the data set, k-NN loses all predictive power.

COMPUTATIONAL COST

- Linear algorithm (no pre-processing):
 - Compute distance for all N datapoints
 - Complexity of distance computation: $O(kN)$
 - No additional space needed
- Tree-based **data structures**: pre-processing
 - Often used in applications: k-d trees (k-dimensional trees).

K-D TREE

- Building a K-D tree from training data:
 - $\{(1,9), (2,3), (4,1), (3,7), (5,4), (6,8), (7,2), (8,8), (7,9), (9,6)\}$
 - pick random dimension, find median, split data, repeat
- Find NNs for new point $(7,4)$
 - find region containing $(7,4)$
 - compare to all points in region



K-NEAREST NEIGHBOR: SUMMARY

- Non parametric model
- When to consider
 - Instance map to points in R^N
 - Few (e.g. less than 20) features per instance
 - Lots of training data
- Advantages
 - Training is very fast (no training)
 - Learn complex target functions
- Disadvantages
 - Slow at query time
 - Easily fooled by irrelevant attributes

K-NEAREST NEIGHBOR: ADVANTAGES

- Easy to program
- No optimization or training required
- Classification accuracy can be very good, it can outperform more complex models

K-NEAREST NEIGHBOR: ISSUES

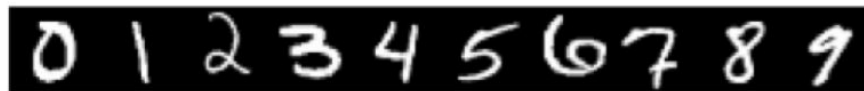
- Choose distance measure
 - Most common: Euclidean
- Choosing k
- Curse of Dimensionality
- Memory based technique.
 - Must make a pass through the data for each classification. This can be prohibitive for large data sets.

AN ASIDE: PARAMETRIC VS NON-PARAMETRIC MODELS

- Parametric models we have a finite number of parameters
 - Linear regression, logistic regression, and linear Support Vector Machines
- Nonparametric models: the number of parameters is (potentially) infinite. The complexity of the model grows with the number of training data.
 - K-nearest neighbor, decision trees, or RBF kernel SVMs are considered as non-parametric learning algorithms since the number of parameters grows with the size of the training set.

EXAMPLE

Decent performance when large training set



- Yann LeCunn – MNIST Digit Recognition
 - Handwritten digits
 - 28x28 pixel images: $d = 784$
 - 60,000 training samples
 - 10,000 test samples
- Nearest neighbour is competitive

	Test Error Rate (%)
Linear classifier (1-layer NN)	12.0
K-nearest-neighbors, Euclidean	5.0
K-nearest-neighbors, Euclidean, deskewed	2.4
K-NN, Tangent Distance, 16x16	1.1
K-NN, shape context matching	0.67
1000 RBF + linear classifier	3.6
SVM deg 4 polynomial	1.1
2-layer NN, 300 hidden units	4.7
2-layer NN, 300 HU, [deskewing]	1.6
LeNet-5, [distortions]	0.8
Boosted LeNet-4, [distortions]	0.7

EXAMPLE

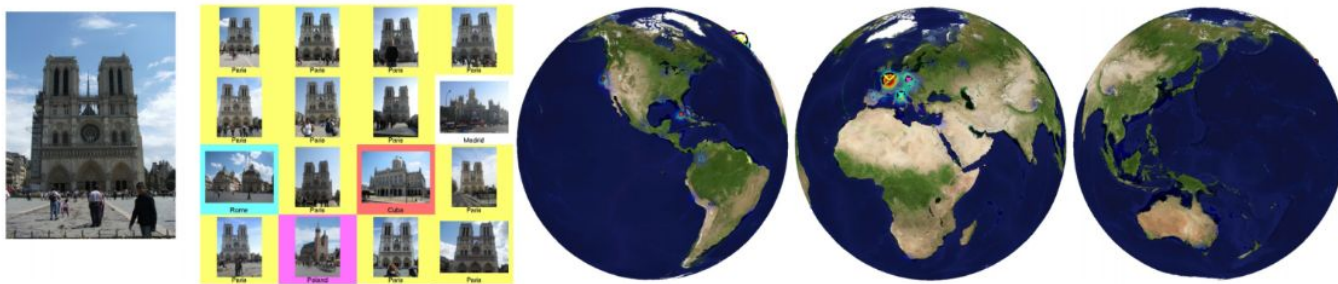
Problem: Where (e.g., which country or GPS location) was this picture taken?



[Paper: James Hays, Alexei A. Efros. im2gps: estimating geographic information from a single image. CVPR'08. Project page: <http://graphics.cs.cmu.edu/projects/im2gps/>]

EXAMPLE

- Get 6M images from Flickr with gps info (dense sampling across world)
- Represent each image with meaningful features
- Do kNN (large k better, they use $k = 120$)



<http://graphics.cs.cmu.edu/projects/im2gps/>

<http://graphics.cs.cmu.edu/projects/im2gps/results.html>

<https://www.youtube.com/watch?v=dgif39IKT9A>

QUESTIONS?

