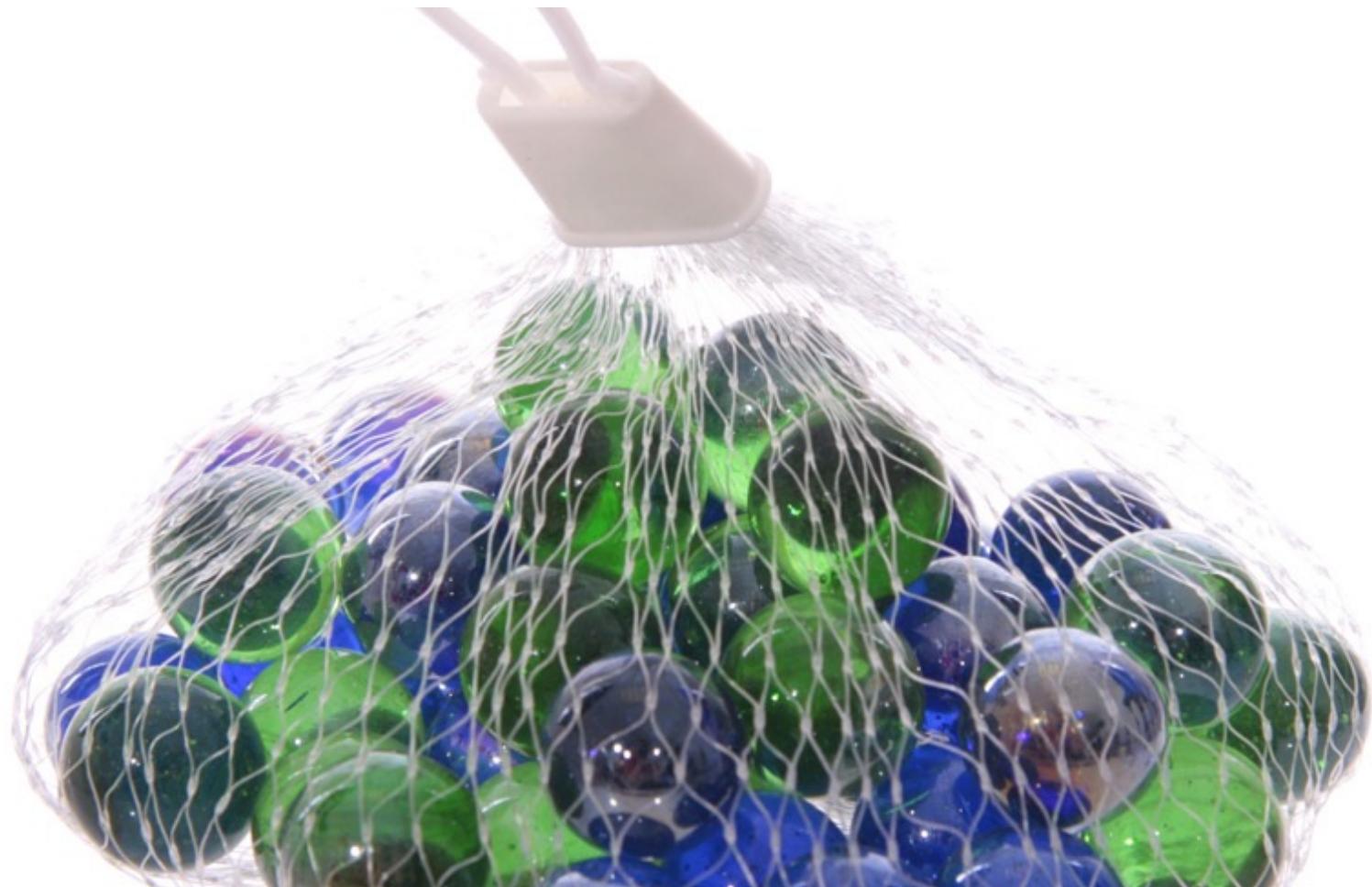


Image classification (II)



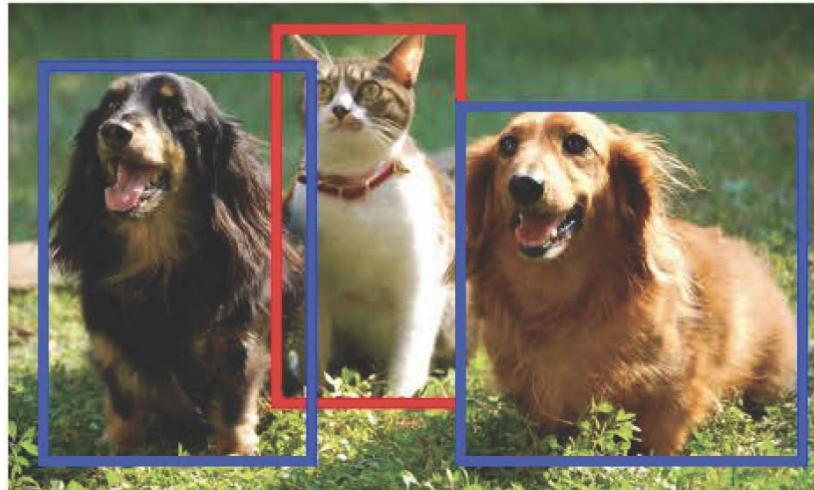
Computer Vision
Fall 2022, Lecture 11

Classification



CAT

Object Detection



CAT, DOG

Instance Segmentation



CAT, DOG

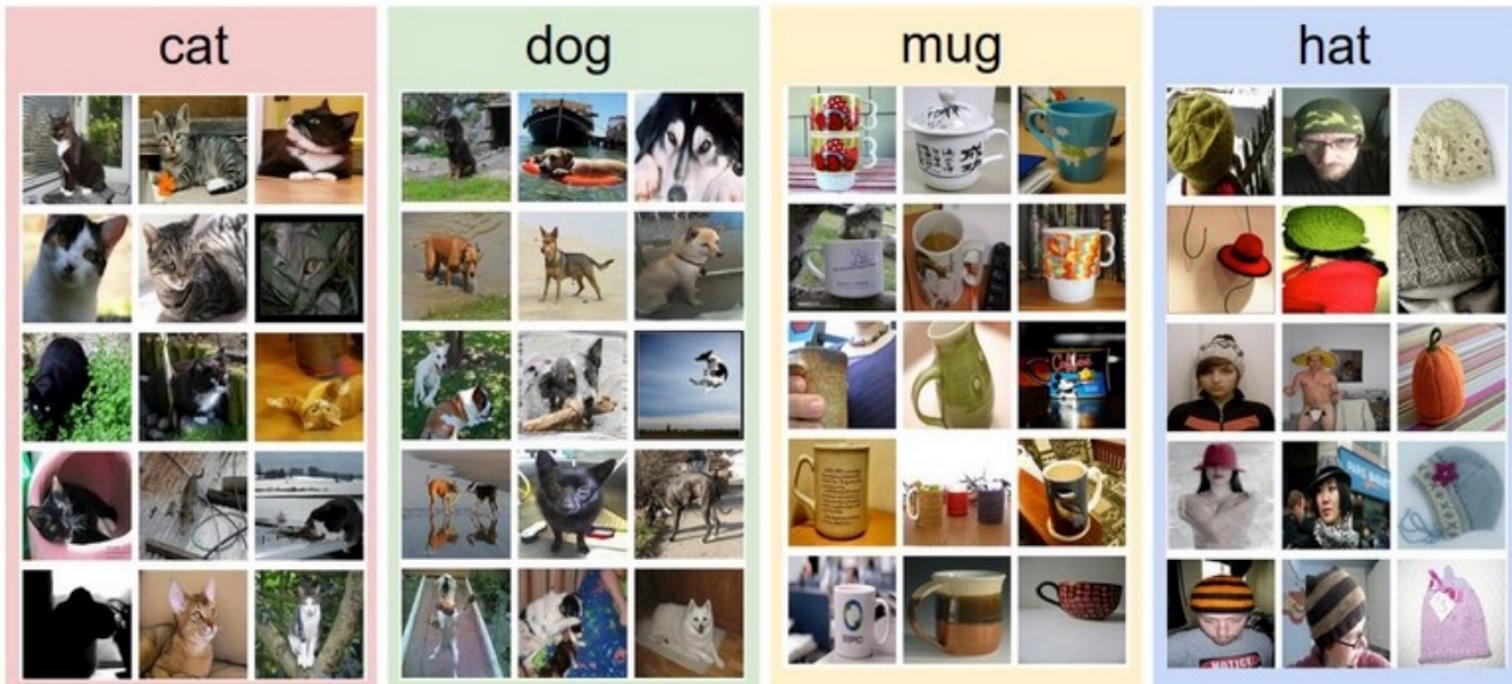
Data-driven approach

- Collect a database of images with labels
- Use ML to train an image classifier
- Evaluate the classifier on test images

```
def train(images, labels):  
    # Machine learning!  
    return model
```

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

Example training set

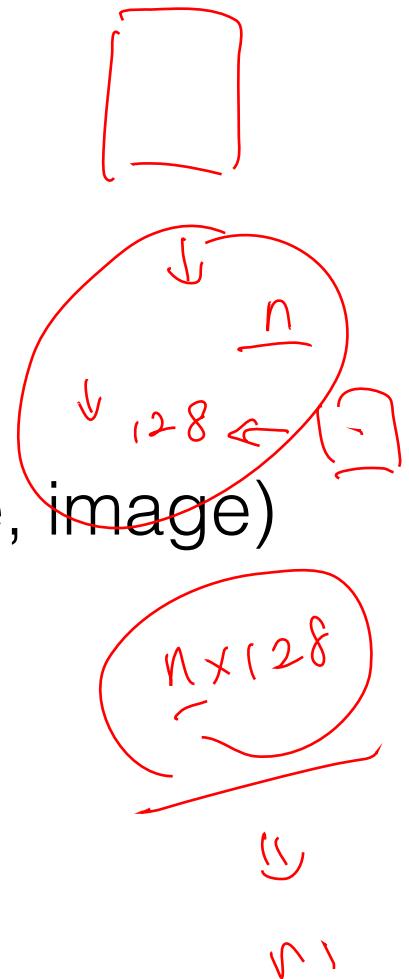


Bag-of-Words

represent a data item (document, texture, image)
as a histogram over words

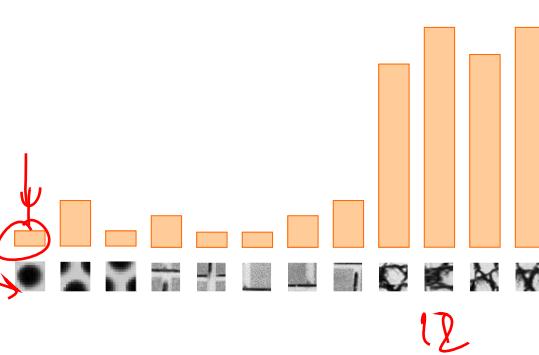
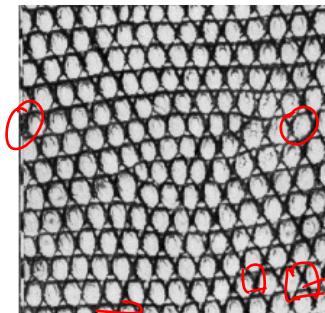
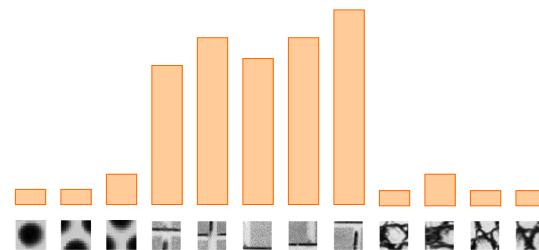
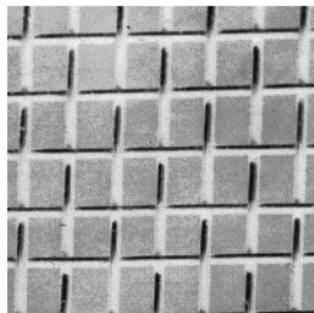
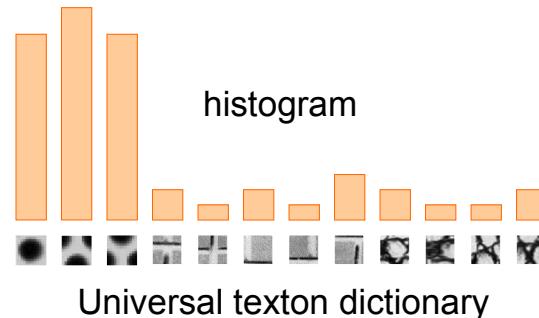
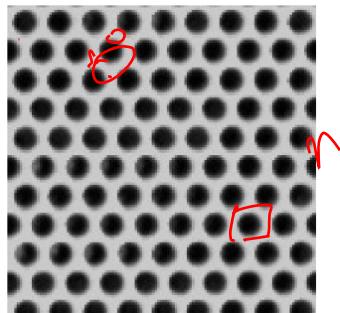
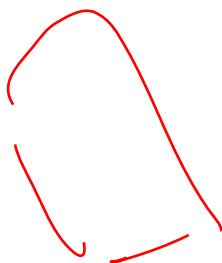
an old idea

(e.g., texture recognition and information retrieval)



Texture recognition

same dictionary



Standard BOW pipeline

(for image classification)

Dictionary Learning:

Learn Visual Words using clustering

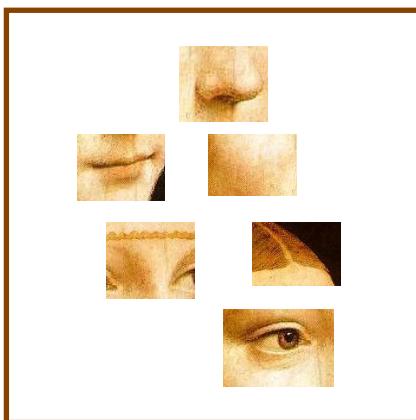
Encode:
build Bags-of-Words (BOW) vectors
for each image

Classify:
Train and test data using BOWs

Dictionary Learning:

Learn Visual Words using clustering

1. extract features (e.g., SIFT) from images



how many per image?

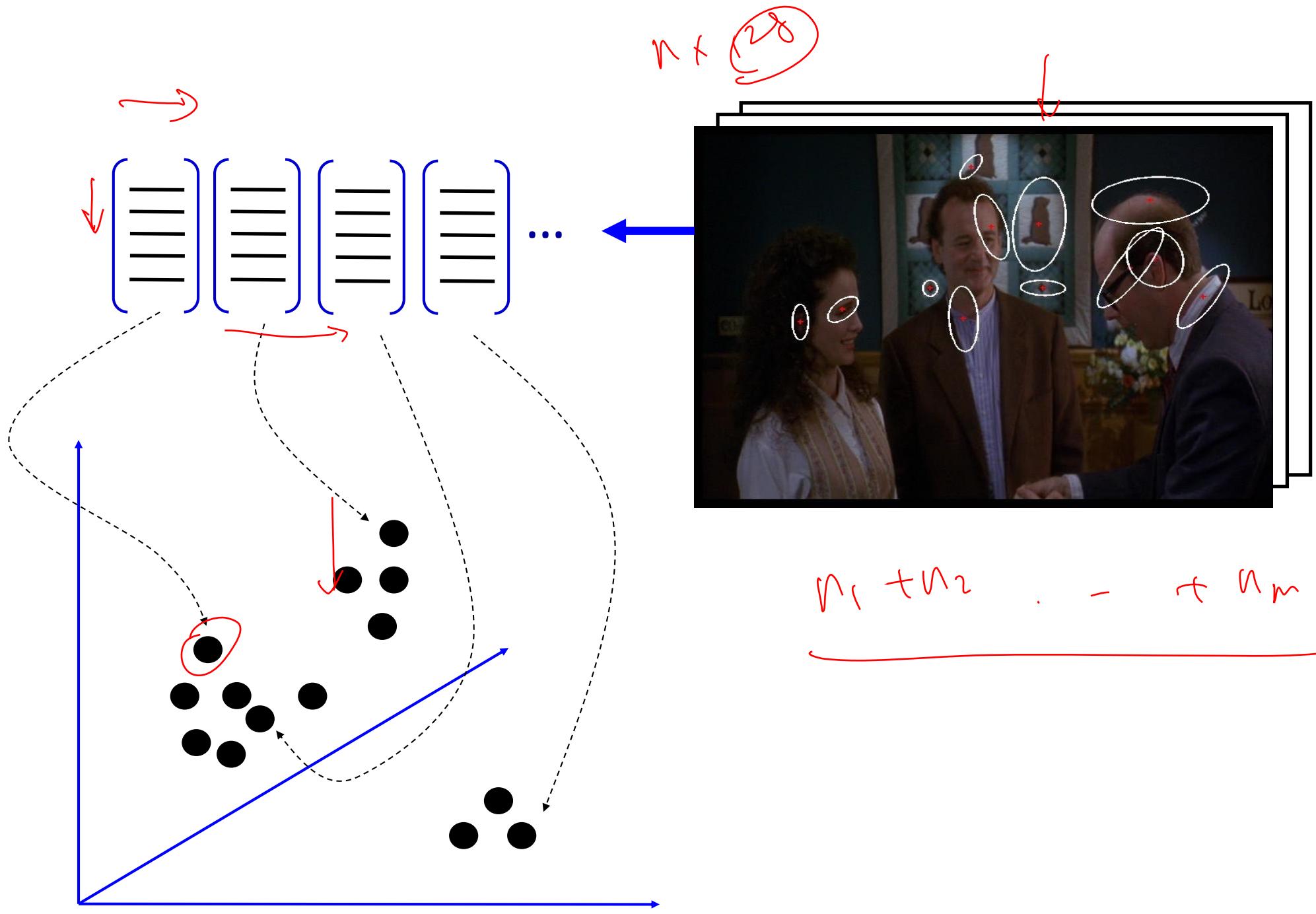
Dictionary Learning:

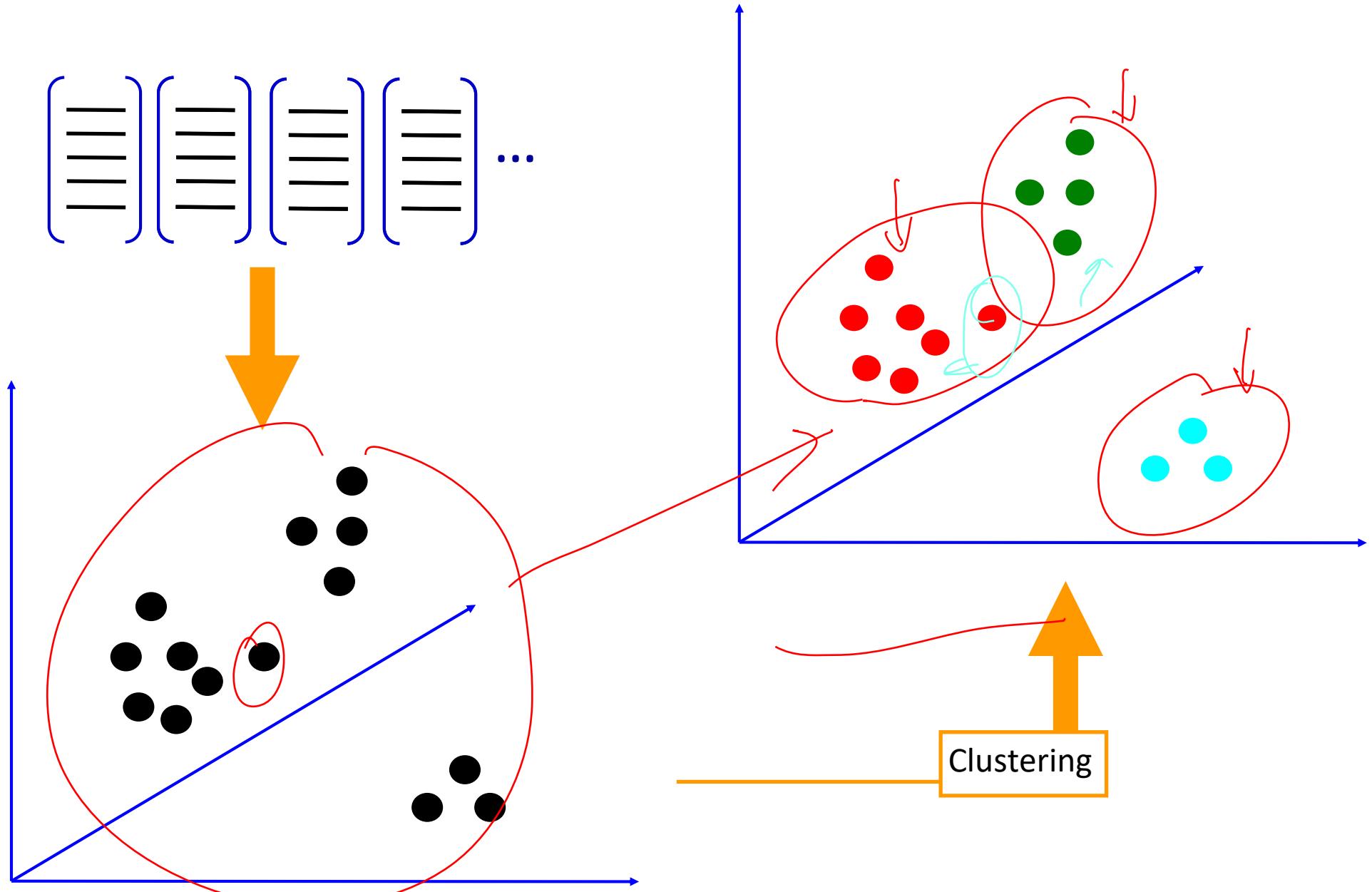
Learn Visual Words using clustering

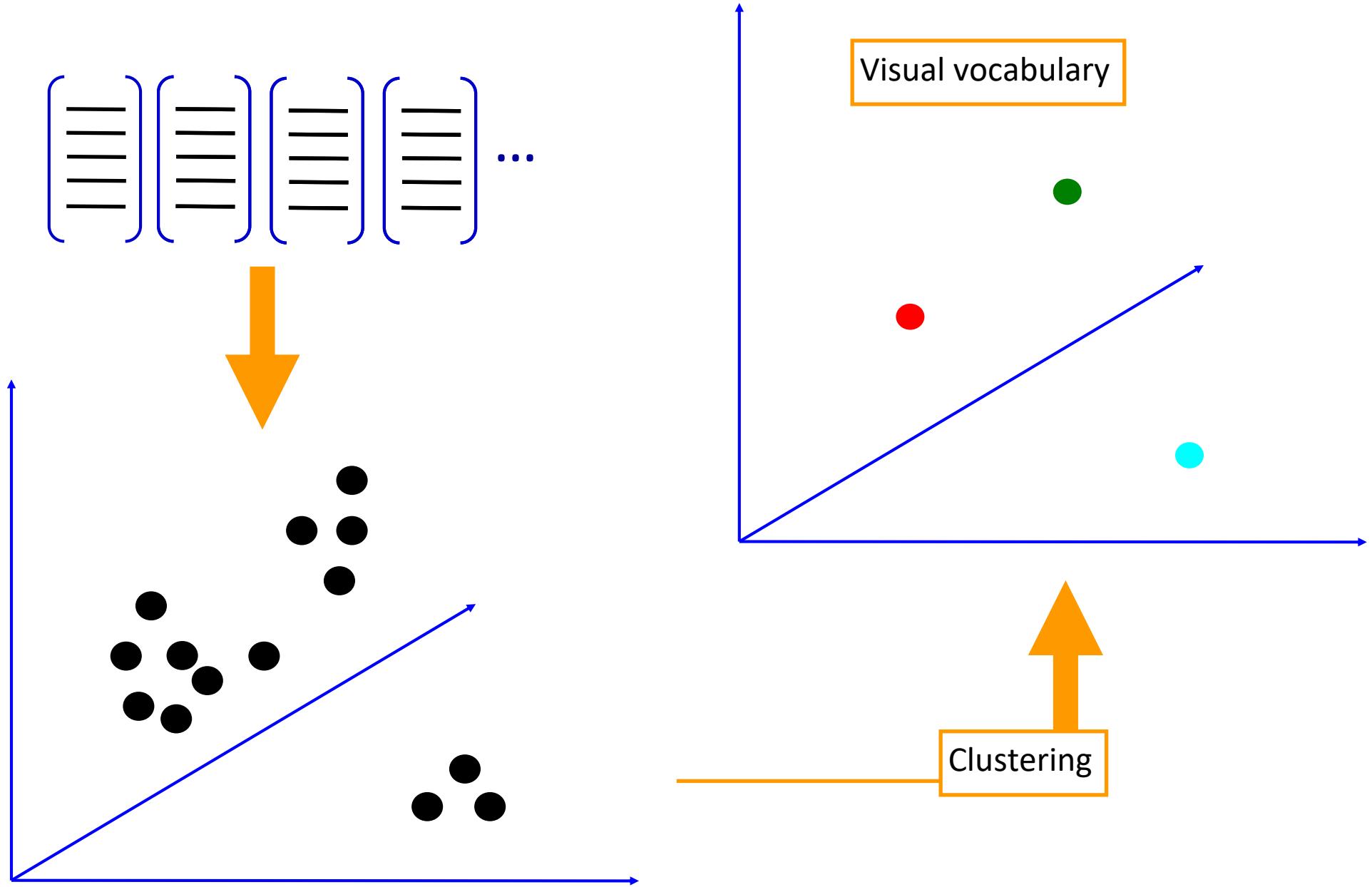
2. Learn visual dictionary (e.g., K-means clustering)



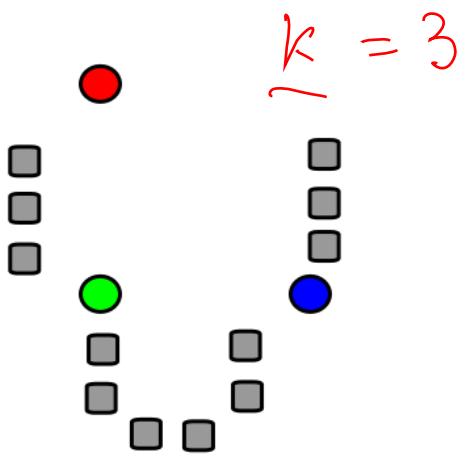
How many words we need?



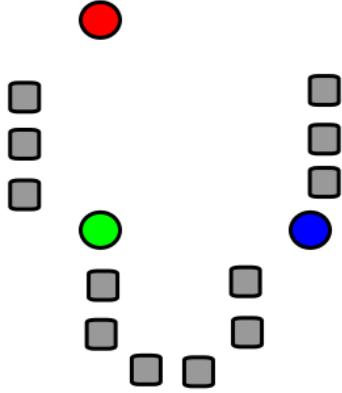




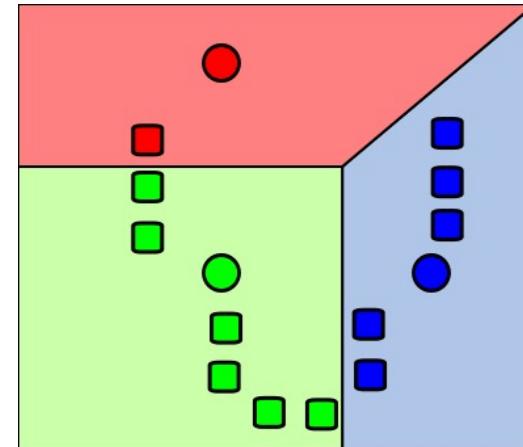
K-means clustering



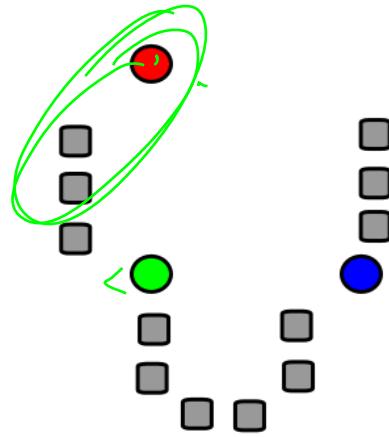
1. Select initial centroids at random



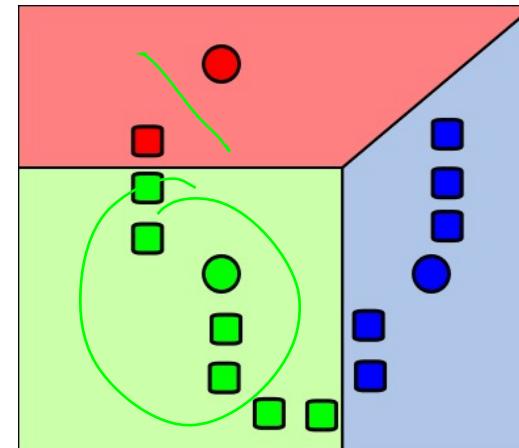
1. Select initial centroids at random



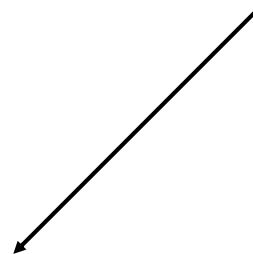
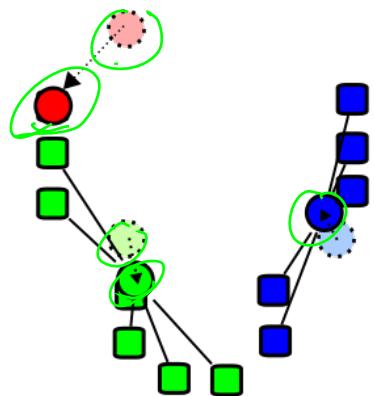
2. Assign each object to the cluster with the nearest centroid.



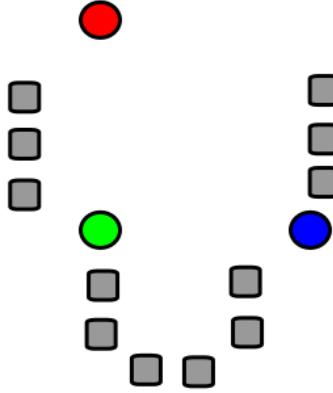
1. Select initial centroids at random



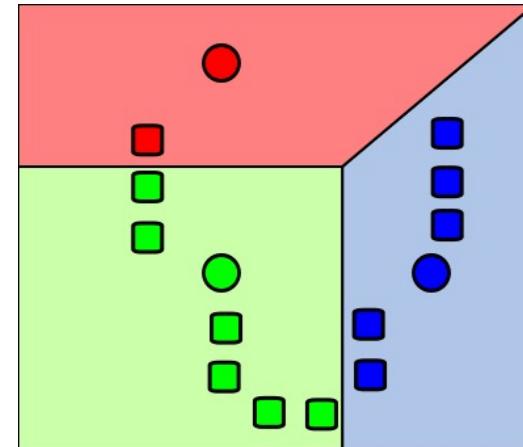
2. Assign each object to the cluster with the nearest centroid.



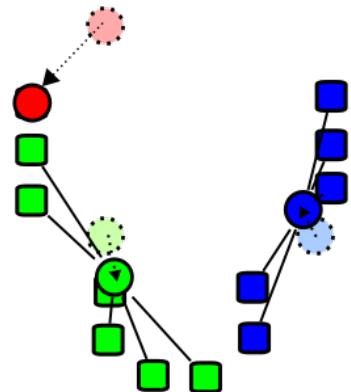
3. Compute each centroid as the mean of the objects assigned to it (go to 2)



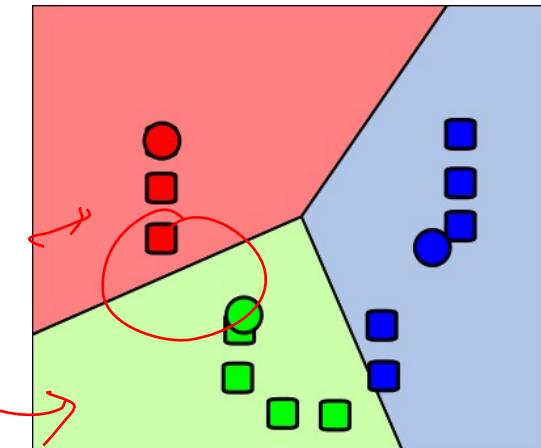
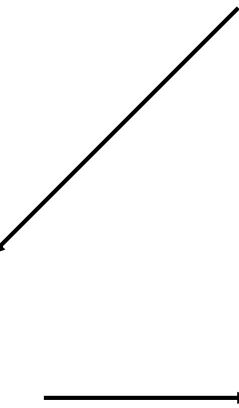
1. Select initial centroids at random



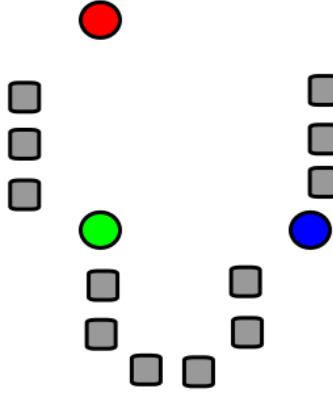
2. Assign each object to the cluster with the nearest centroid.



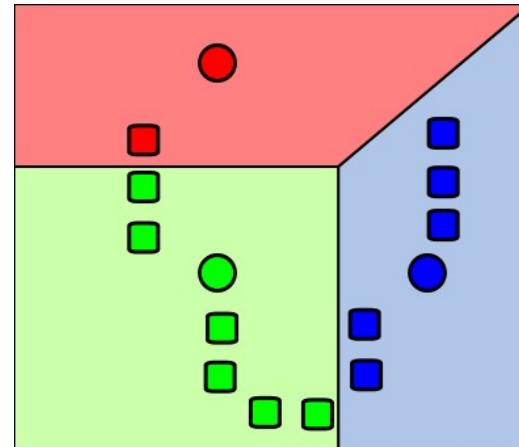
3. Compute each centroid as the mean of the objects assigned to it (go to 2)



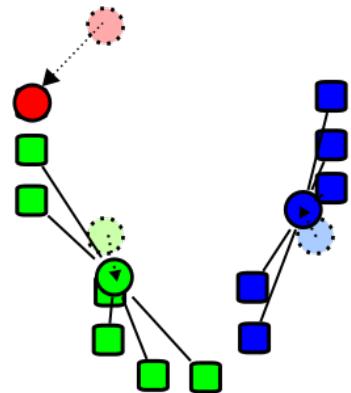
2. Assign each object to the cluster with the nearest centroid.



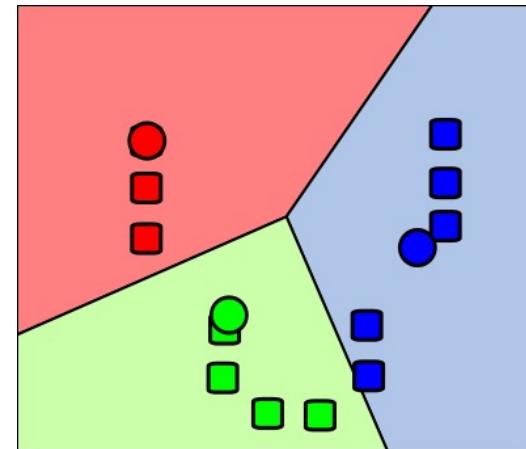
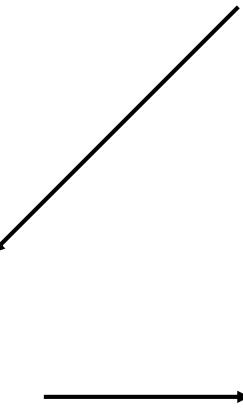
1. Select initial centroids at random



2. Assign each object to the cluster with the nearest centroid.



3. Compute each centroid as the mean of the objects assigned to it (go to 2)



2. Assign each object to the cluster with the nearest centroid.

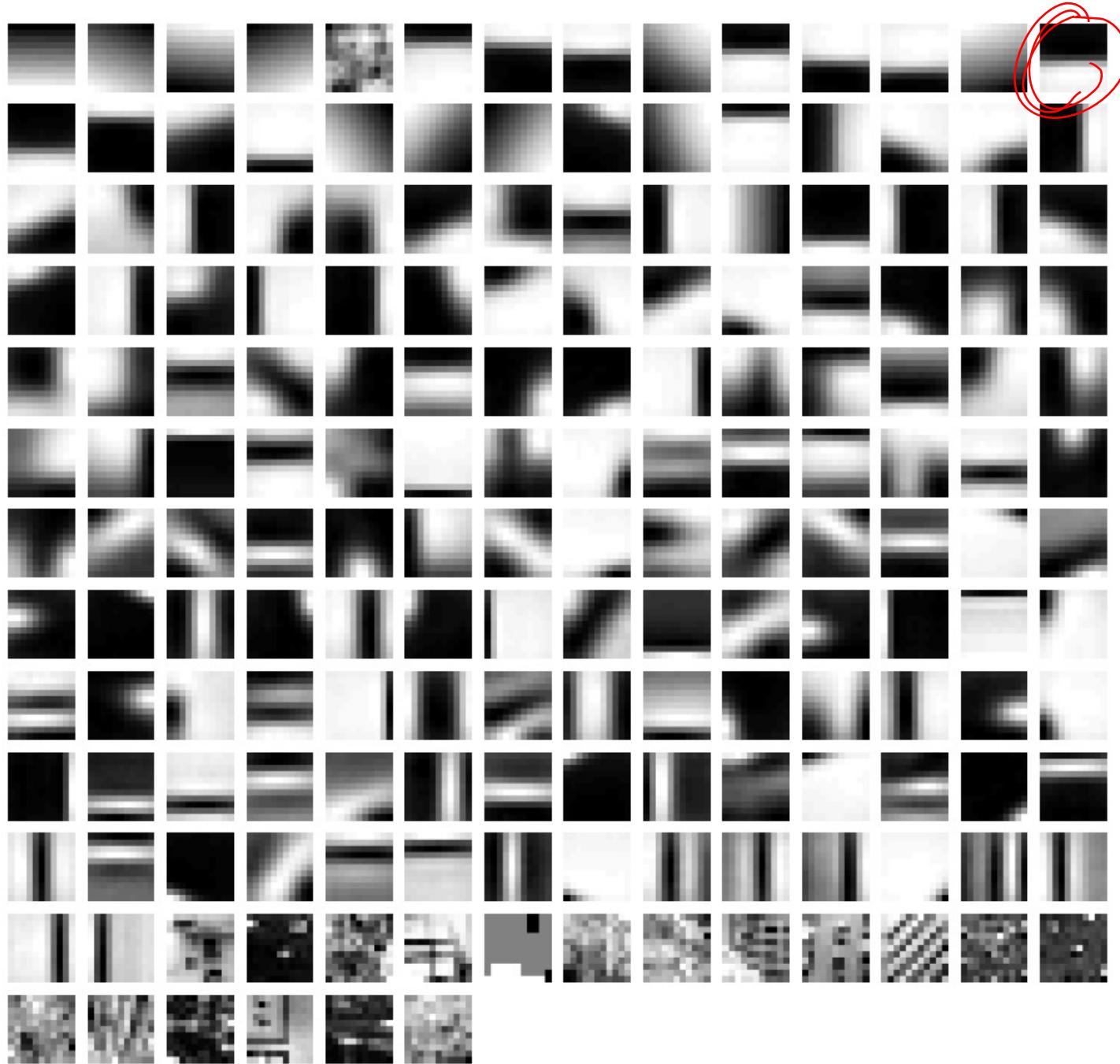
Repeat previous 2 steps until no change

$$D : \underbrace{k \times 128}$$

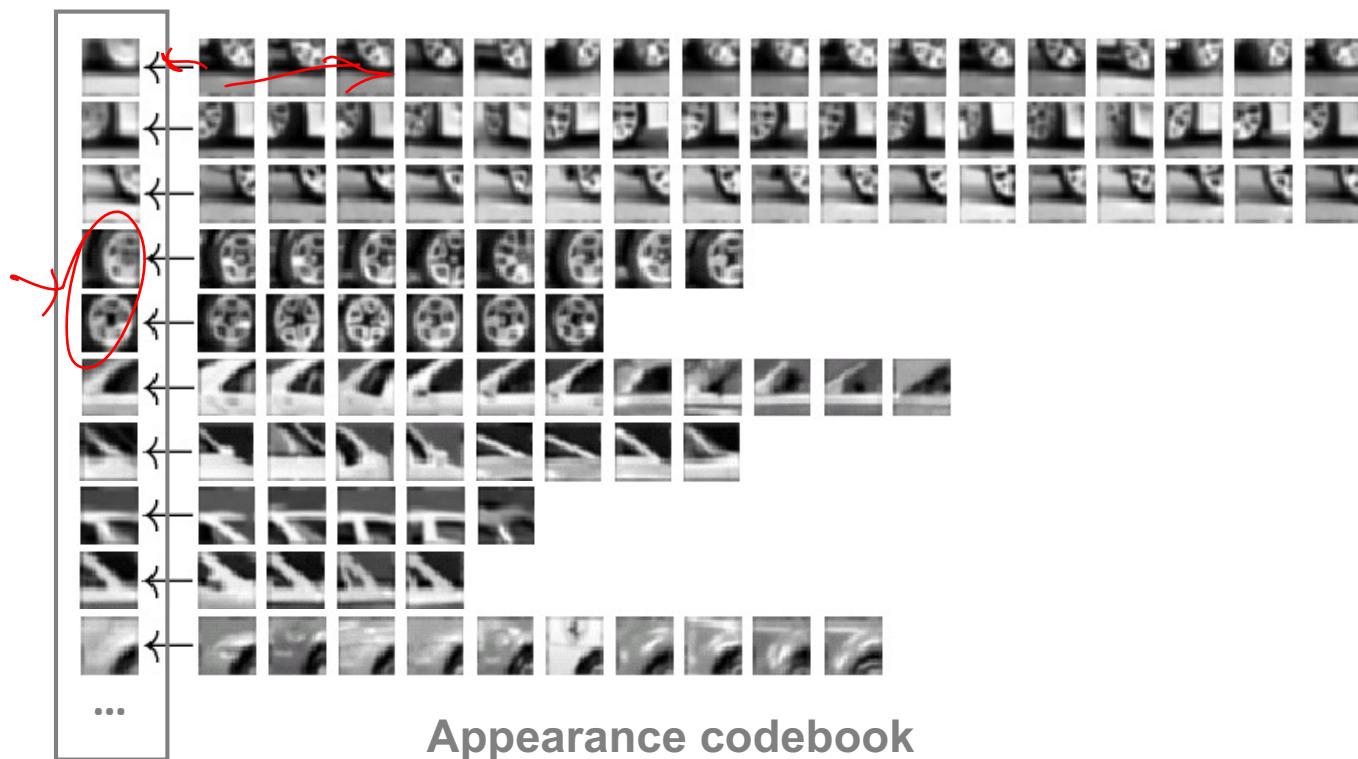
*From what **data** should I learn the dictionary?*

- Dictionary can be learned on separate training set
- Provided the training set is sufficiently representative, the dictionary will be “universal”

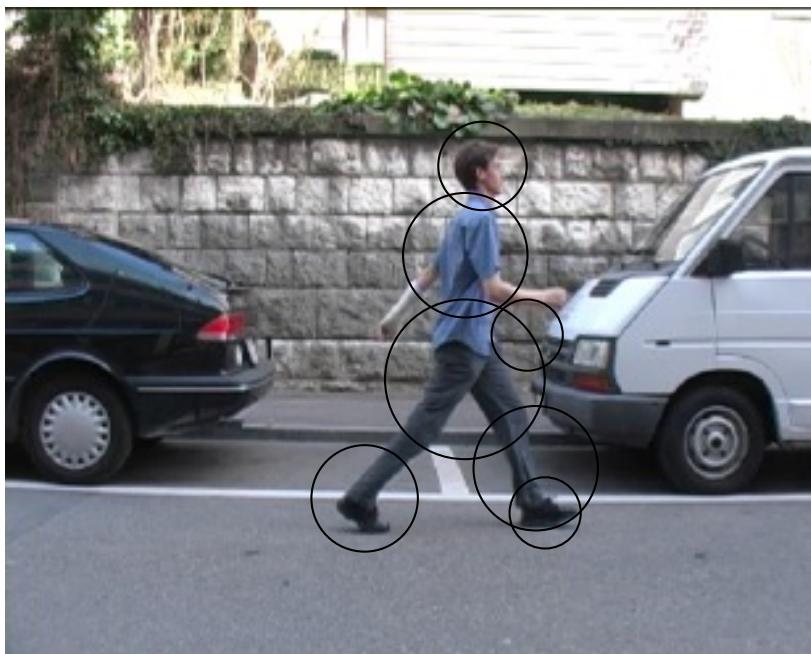
Example visual dictionary



Example dictionary



Another dictionary



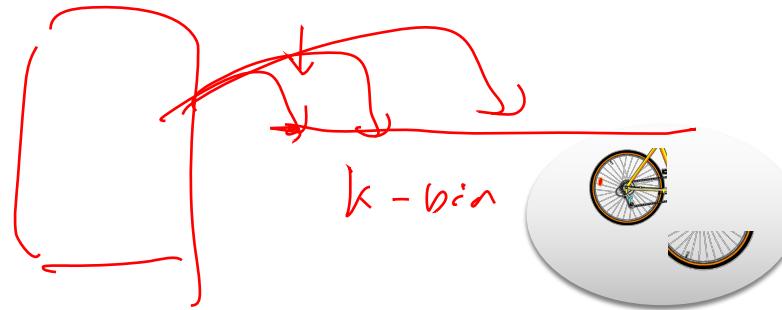
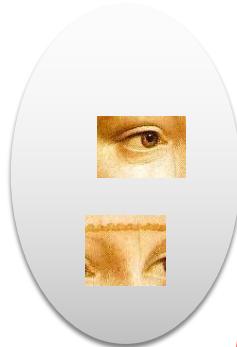
Appearance codebook

Dictionary Learning:

Learn Visual Words using clustering

Encode:
build Bags-of-Words (BOW) vectors
for each image

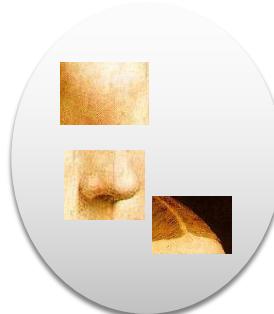
Classify:
Train and test data using BOWs



1. Quantization: image features get associated to a visual word (nearest cluster center)

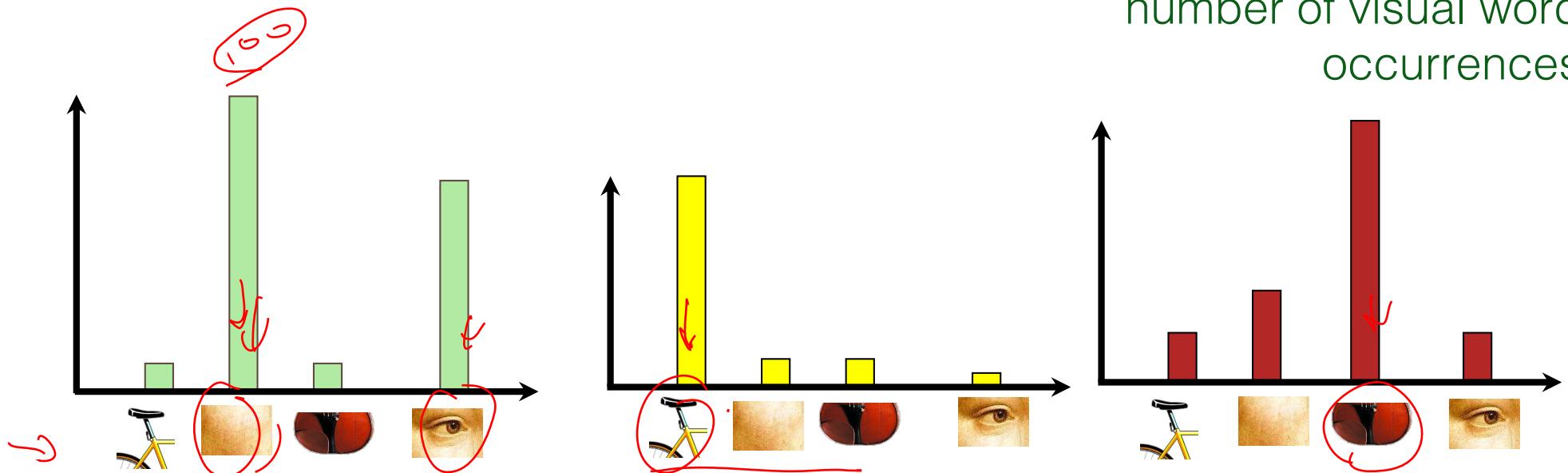
Encode:

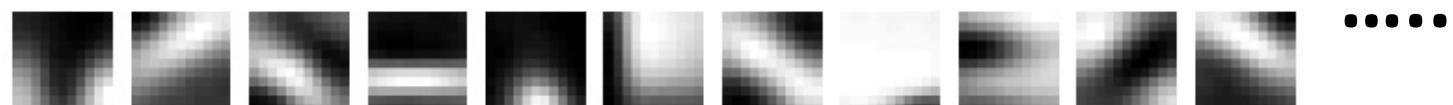
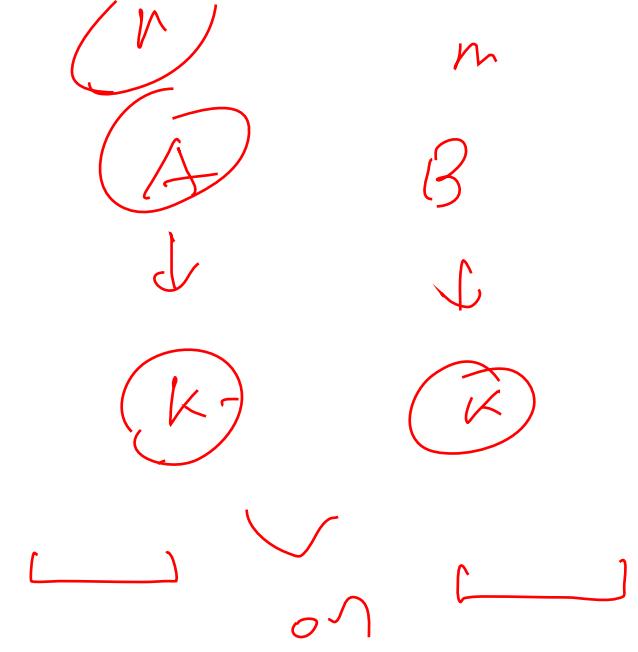
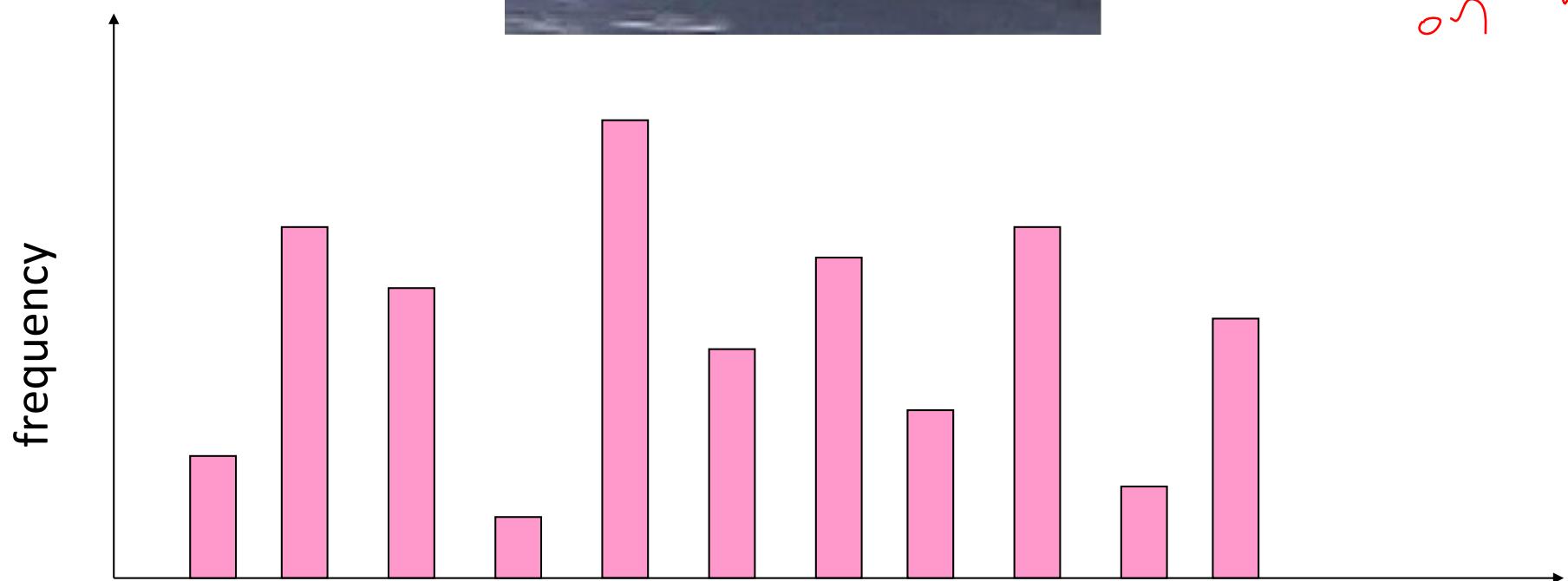
build Bags-of-Words (BOW) vectors
for each image



Encode:
build Bags-of-Words (BOW) vectors
for each image

2. Histogram: count the
number of visual word
occurrences





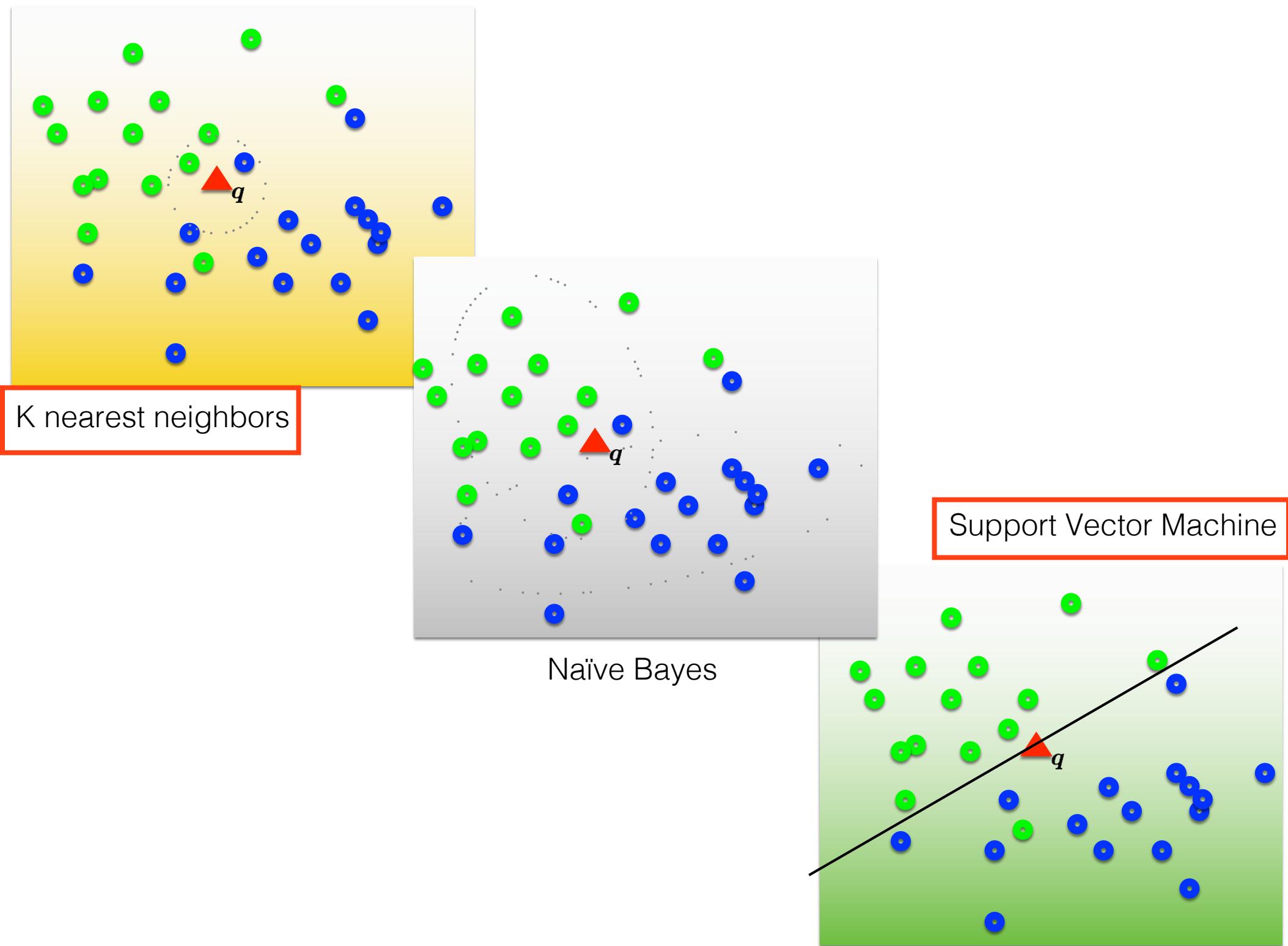
codewords

Dictionary Learning:

Learn Visual Words using clustering

Encode:
build Bags-of-Words (BOW) vectors
for each image

Classify:
Train and test data using BOWs



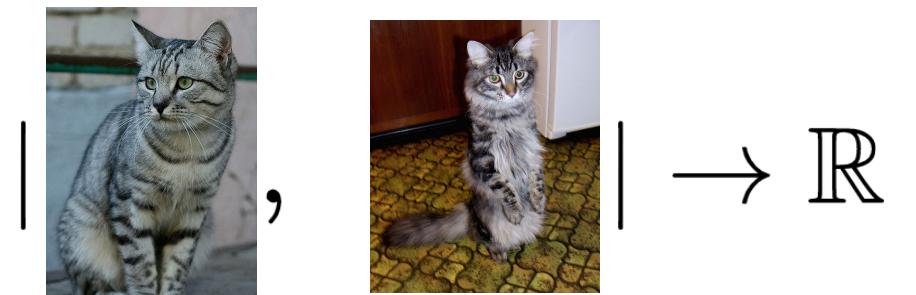
K nearest neighbors

Nearest neighbors

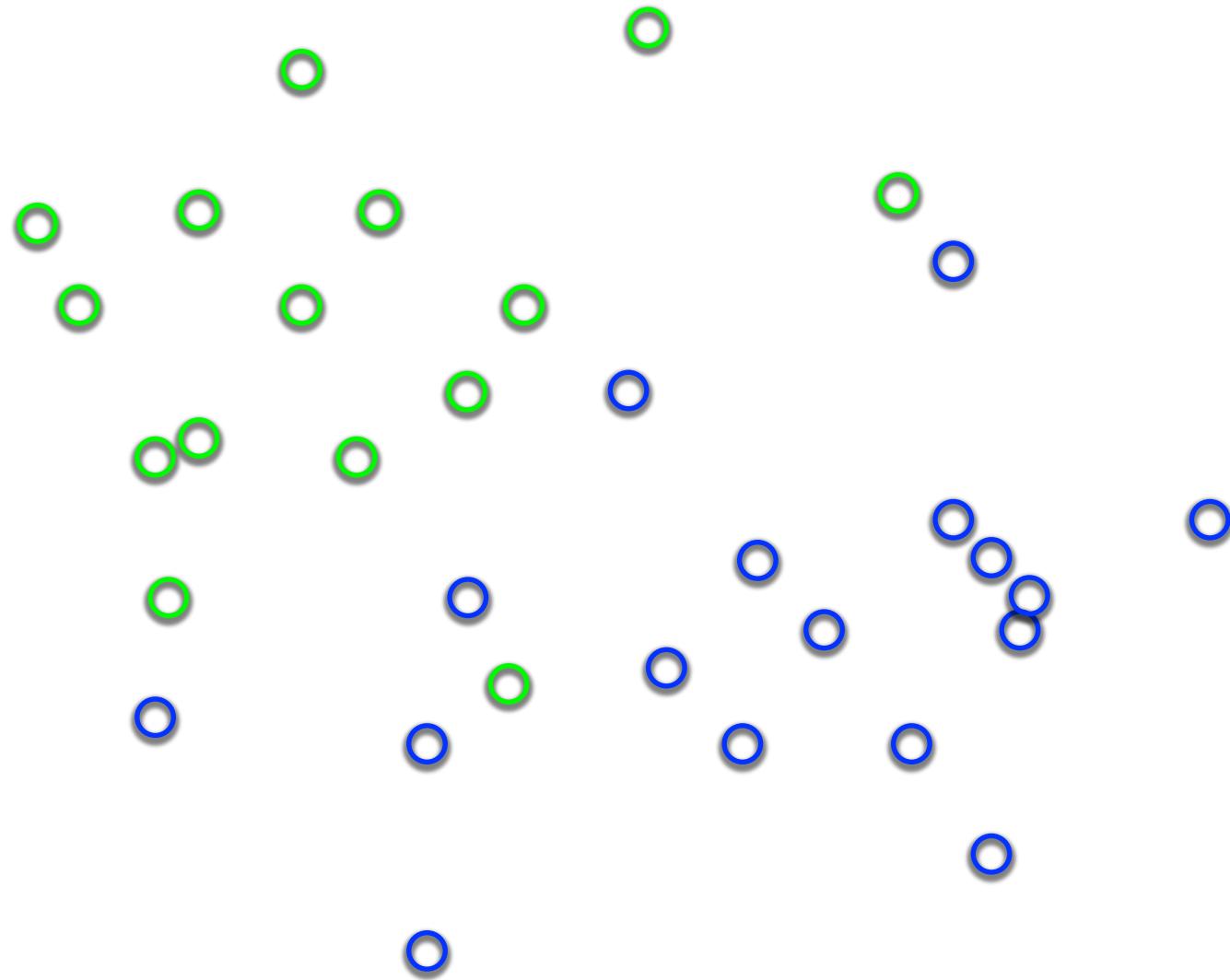


Distance Metric

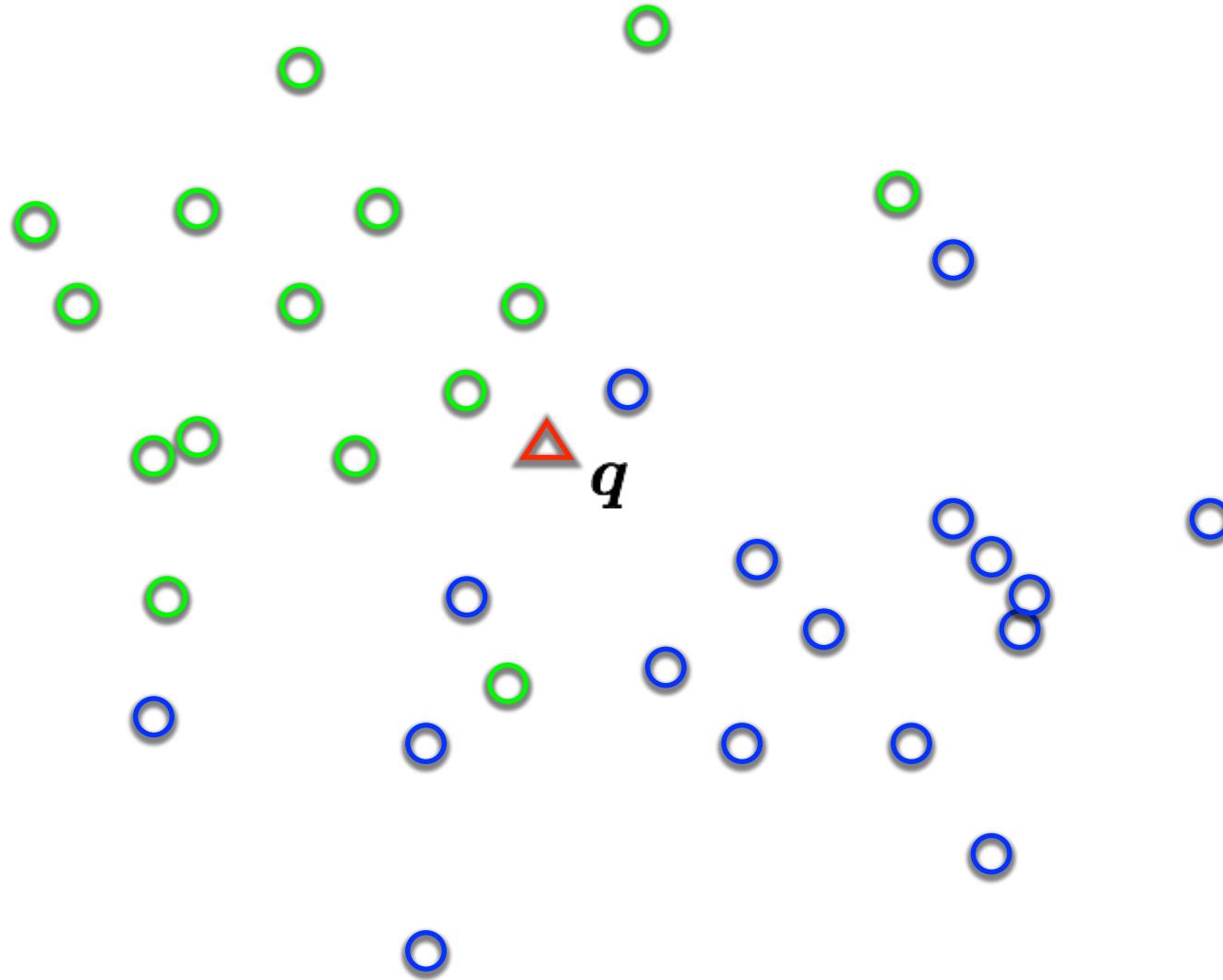
K=1



Distribution of data from two classes

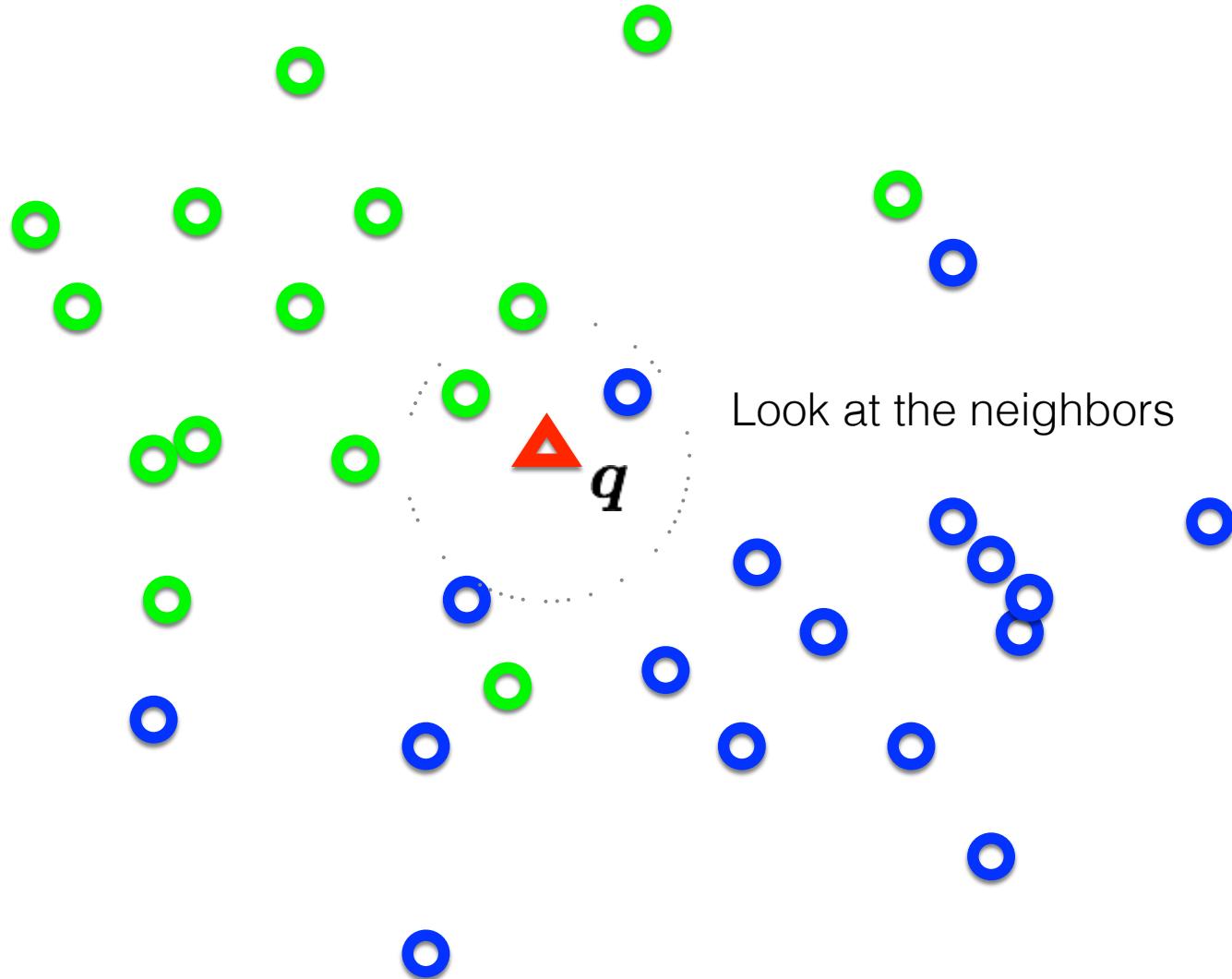


Distribution of data from two classes

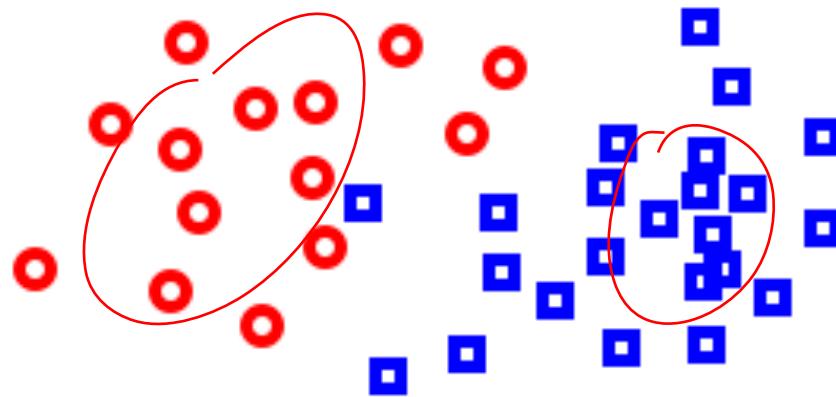


Which class does q belong to?

Distribution of data from two classes



K-Nearest Neighbor (KNN) Classifier



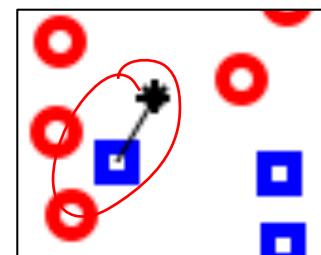
For a given query point q , assign the class of the nearest neighbor

Compute the k nearest neighbors and assign the class by majority vote.

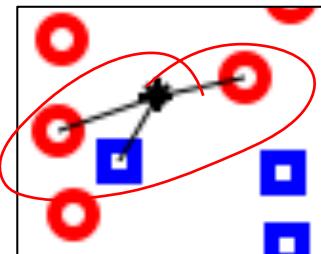
Non-parametric pattern classification approach

Consider a two-class problem where each sample consists of two measurements (x, y) .

$k = 1$



$k = 3$



What is the best distance metric between data points?

- Typically Euclidean distance
- Locality sensitive distance metrics
- Important to normalize.
Dimensions have different scales

How many K?

- Typically $k=1$ is good
- Cross-validation (try different k !)

Distance metrics

$$D(\mathbf{x}, \mathbf{y}) = \sqrt{(x_1 - y_1)^2 + \cdots + (x_N - y_N)^2} \quad \text{Euclidean}$$

$$D(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{x_1 y_1 + \cdots + x_N y_N}{\sqrt{\sum_n x_n^2} \sqrt{\sum_n y_n^2}} \quad \text{Cosine}$$

$$D(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \sum_n \frac{(x_n - y_n)^2}{(x_n + y_n)} \quad \text{Chi-squared}$$

Choice of distance metric

- Hyperparameter

L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

- Two most commonly used special cases of p-norm

$$\|x\|_p = \left(|x_1|^p + \cdots + |x_n|^p \right)^{\frac{1}{p}} \quad p \geq 1, x \in \mathbb{R}^n$$

CIFAR-10 and NN results

Example dataset: **CIFAR-10**

10 labels

50,000 training images

10,000 test images.

airplane



automobile



bird



cat



deer

dog

frog

horse

ship

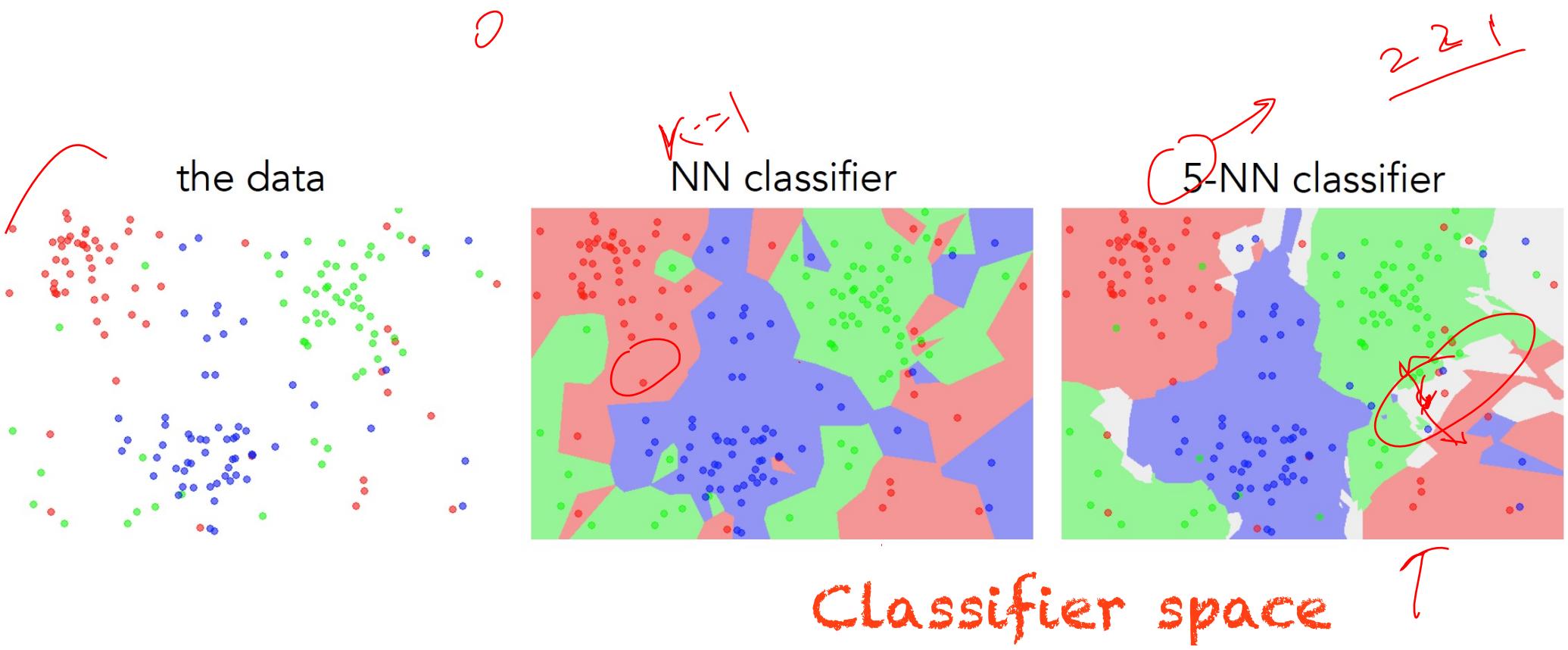
truck

For every test image (first column),
examples of nearest neighbors in rows



k-nearest neighbor

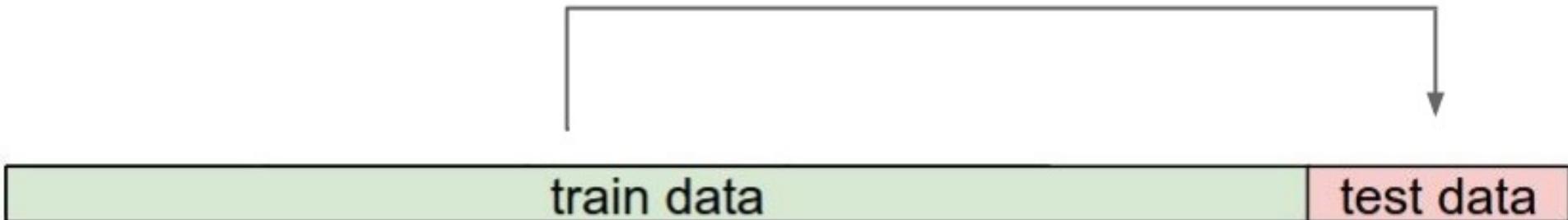
- Find the k closest points from training data
- Labels of the k points “vote” to classify



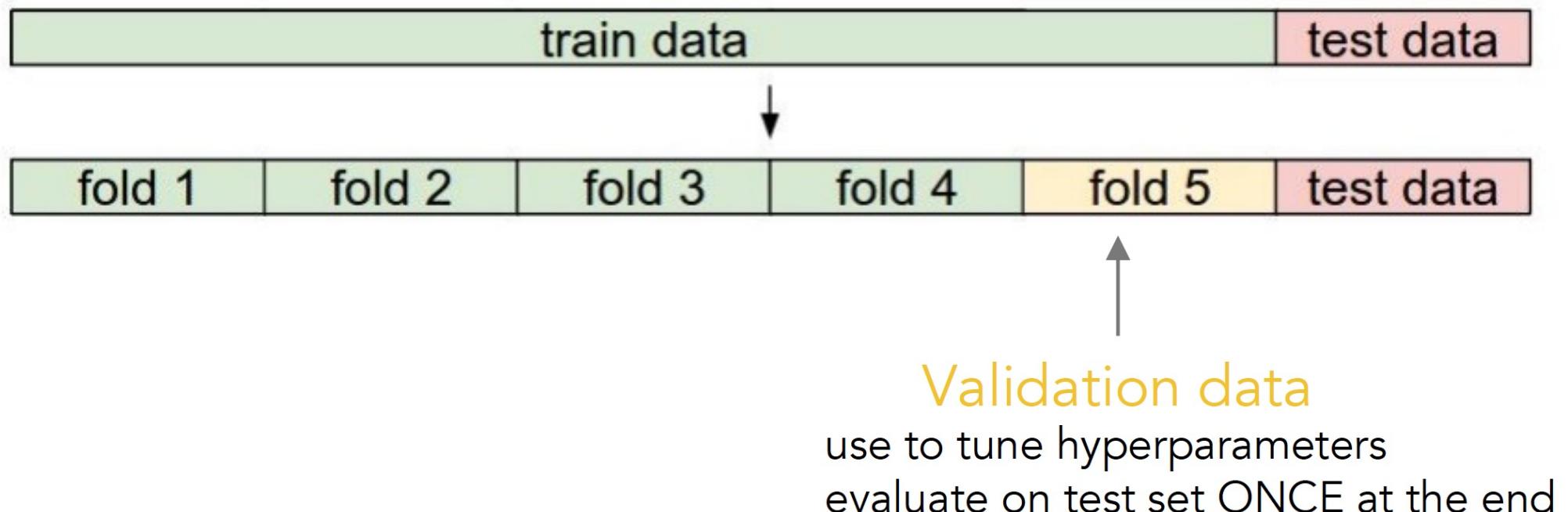
Hyperparameters

- What is the best distance to use?
- What is the best value of k to use?
- i.e., how do we set the hyperparameters?
- Very problem-dependent
- Must try them all and see what works best

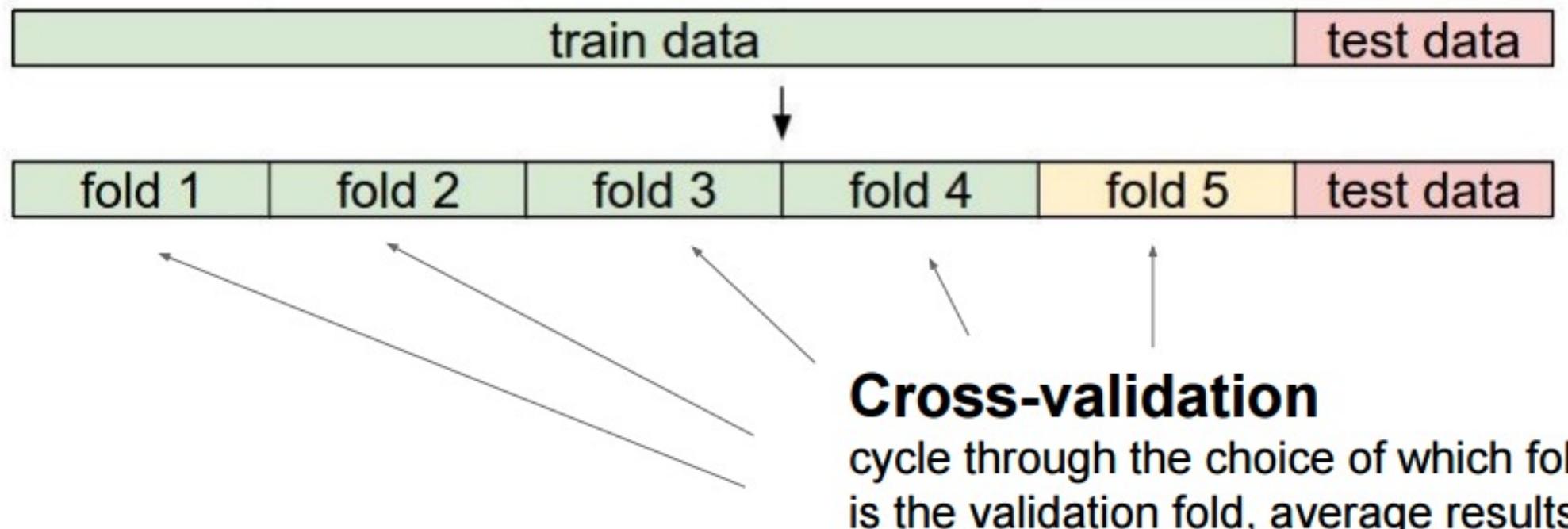
Try out what hyperparameters work best on test set.



Validation

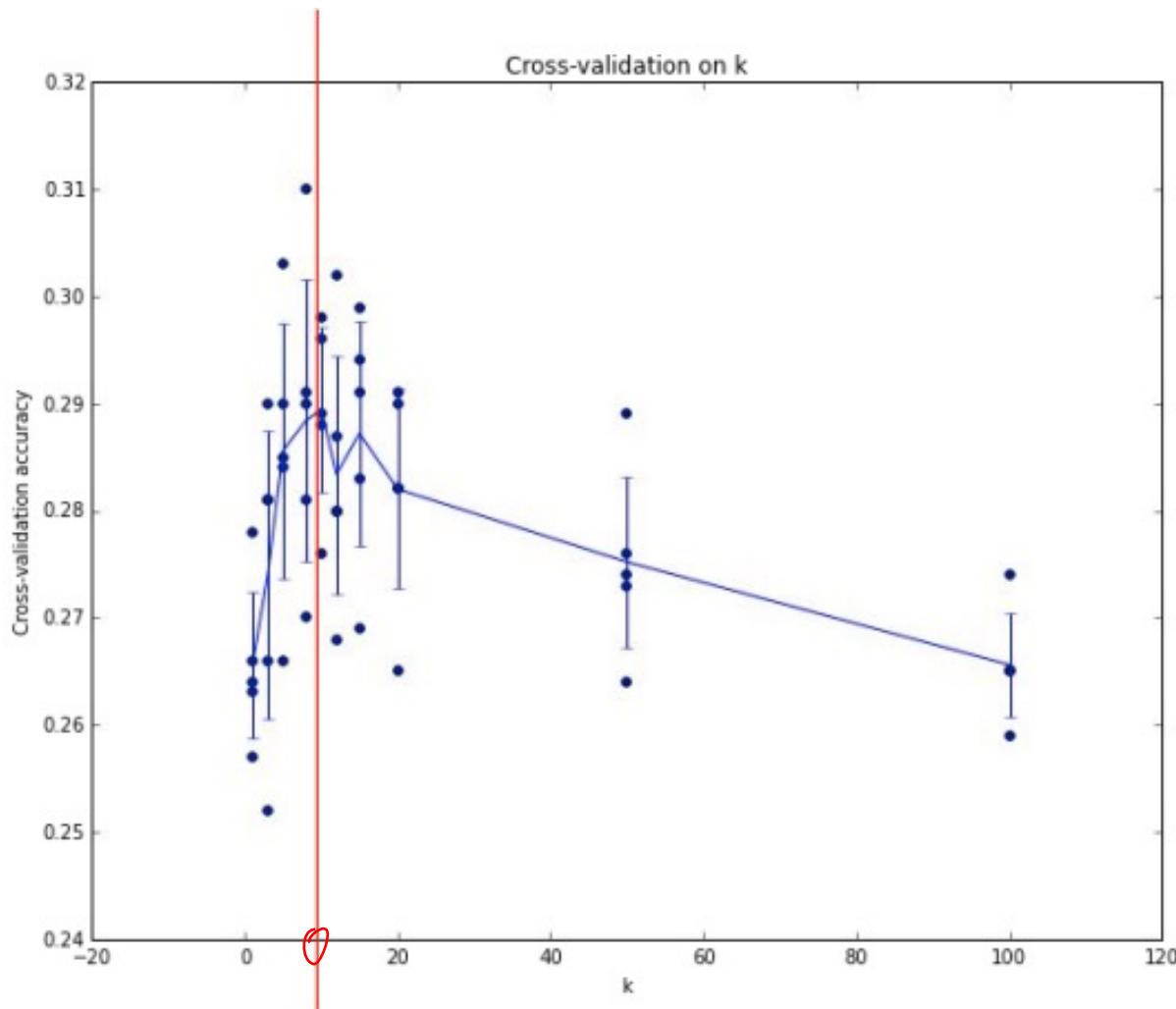


Cross-validation



4

→ ↴



Example of
5-fold cross-validation
for the value of **k**.

Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

(Seems that $k \approx 7$ works best
for this data)

How to pick hyperparameters?

- Methodology
 - Train, validate, test
- Train for original model
- Validate to find hyperparameters
- Test to understand generalizability

Pros

- simple yet effective

Cons

- search is expensive (can be sped-up)
- storage requirements
- difficulties with high-dimensional data

kNN -- Complexity and Storage

- N training images, M test images

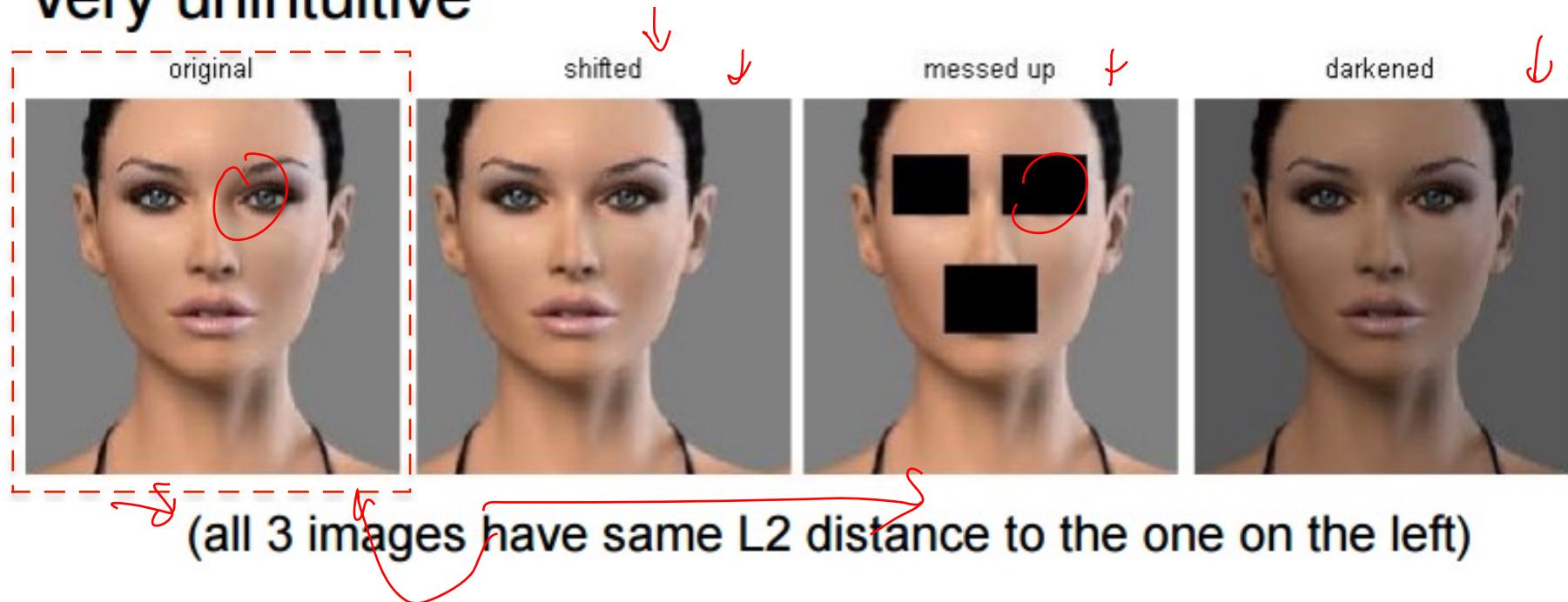
- Training: $O(1)$
- Testing: $O(MN)$

why?

- Hmm...
 - Normally need the opposite
 - Slow training (ok), fast testing (necessary)

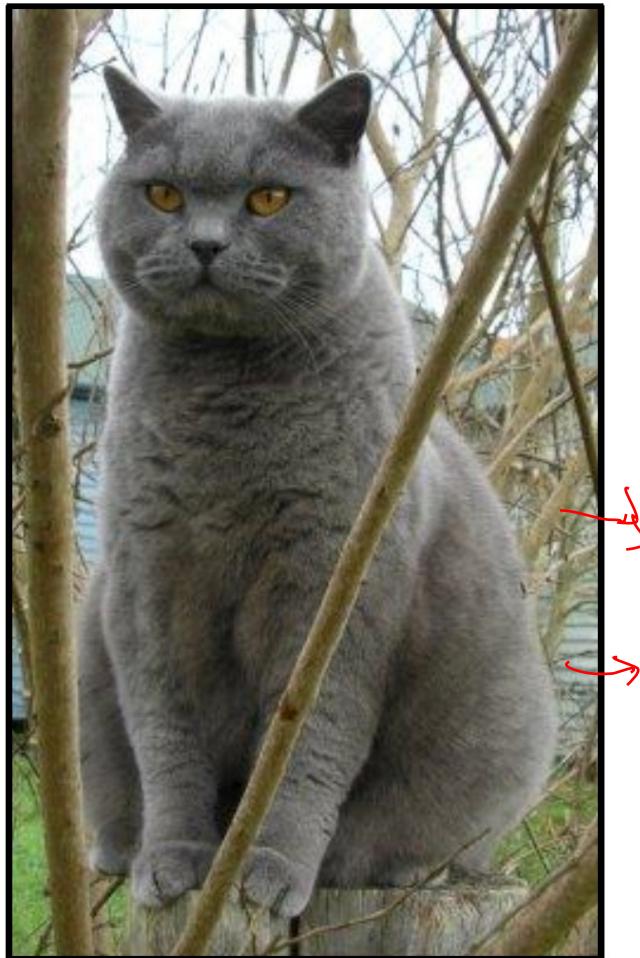
k-Nearest Neighbor on images **never used**.

- terrible performance at test time
- distance metrics on level of whole images can be very unintuitive



Support Vector Machine

Image Classification



(assume given set of discrete labels)
 $\overbrace{\{dog, cat, truck, plane, \dots\}}$

→ cat

Score function



class scores

Linear Classifier

define a **score function**

data (histogram)

$$f(x_i, W, b) = \underbrace{Wx_i}_\text{"weights"} + \underbrace{b}_\text{"bias vector"}$$

class scores

$5 \times K$

K

“parameters”

“weights”

“bias vector”

“parameters”

Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)

Convert image to histogram representation



input image

0.2	-0.5	0.1	2.0
1.5	1.3	2.1	0.0
0	0.25	0.2	-0.3

W

↓

56
231
24
2

x_i

+

1.1
3.2
-1.2

→

-96.8
437.9
61.95

$f(x_i; W, b)$

cat score

dog score

ship score

0 -1

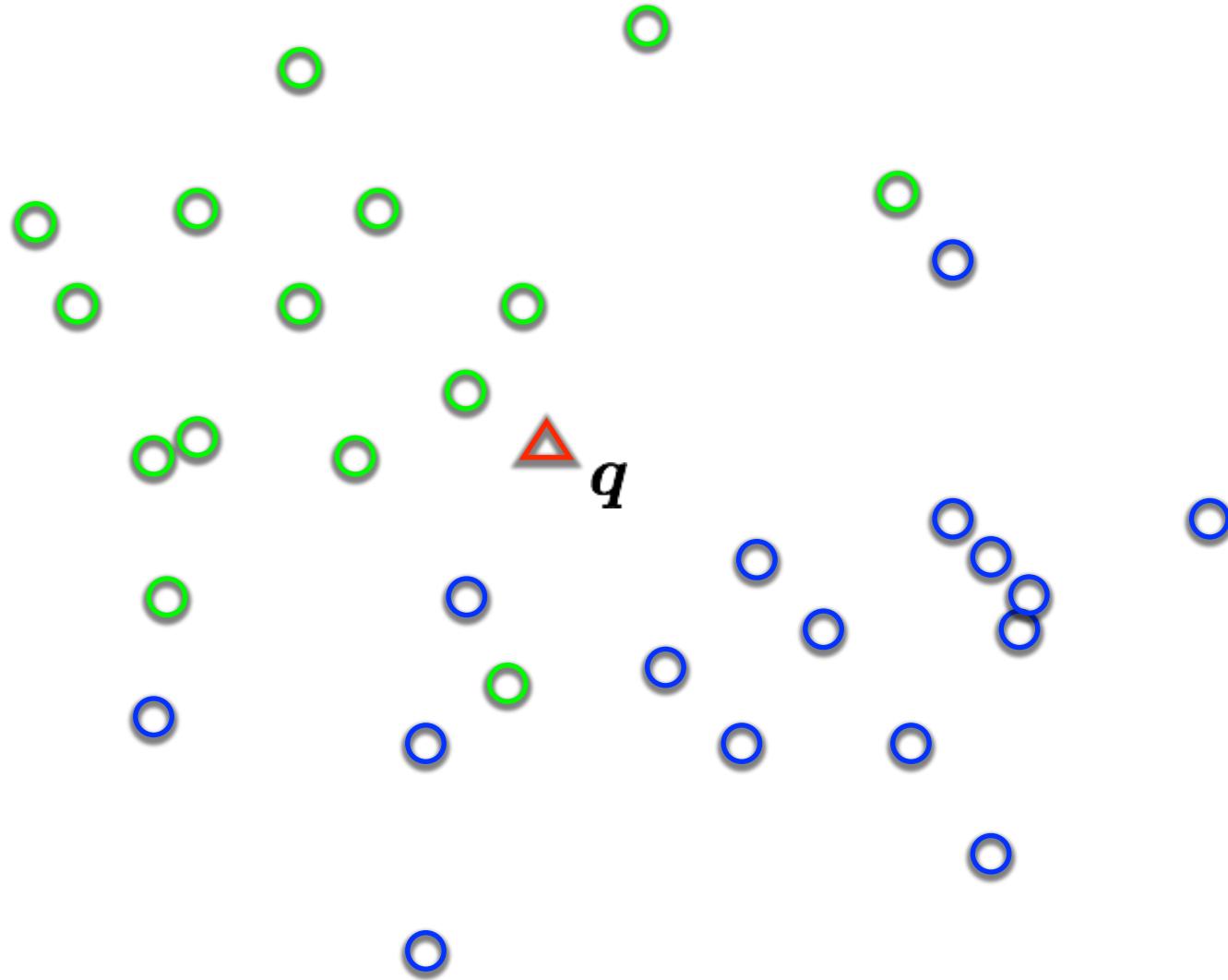
-96.8
437.9
61.95

cat score

dog score

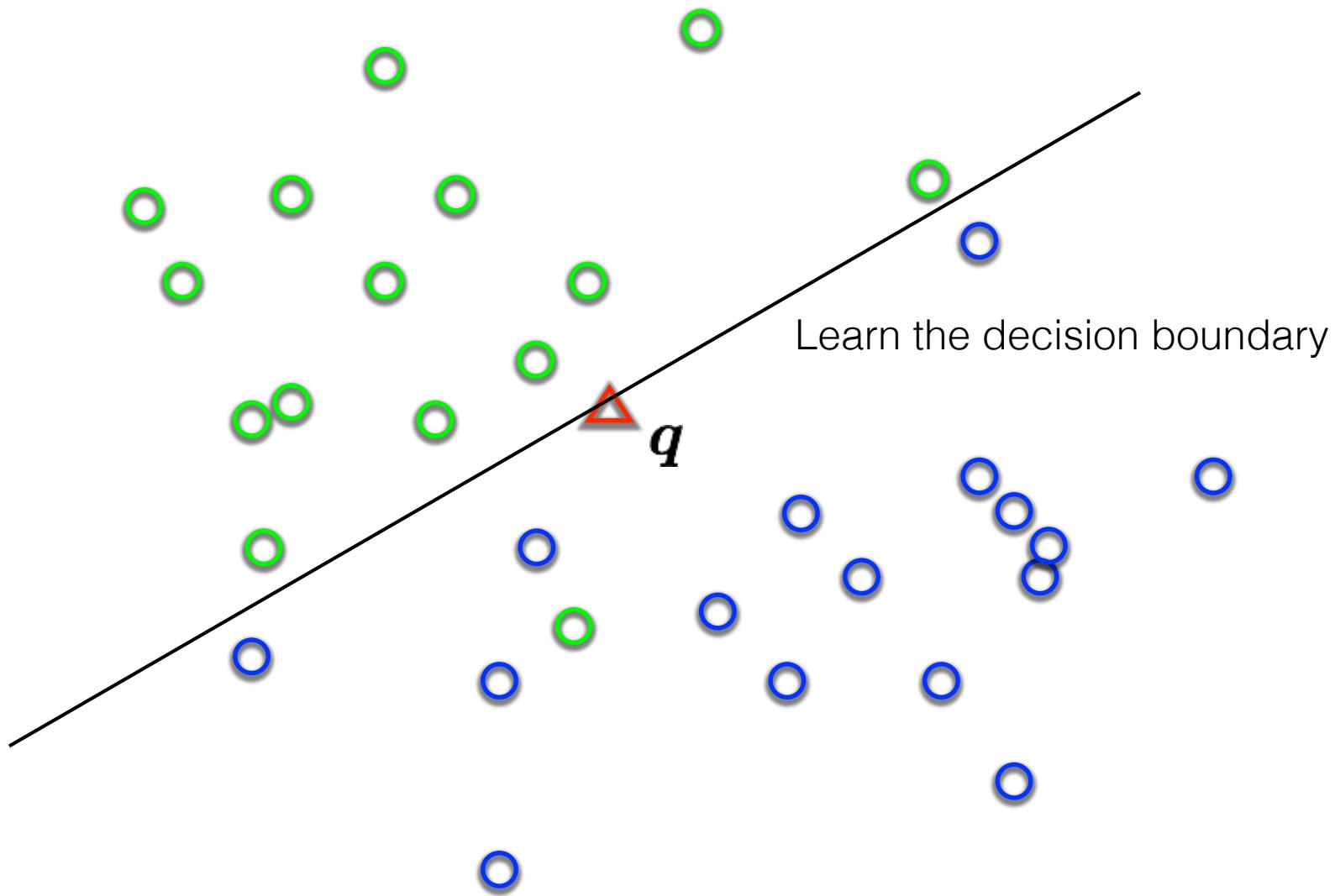
ship score

Distribution of data from two classes



Which class does q belong to . ?

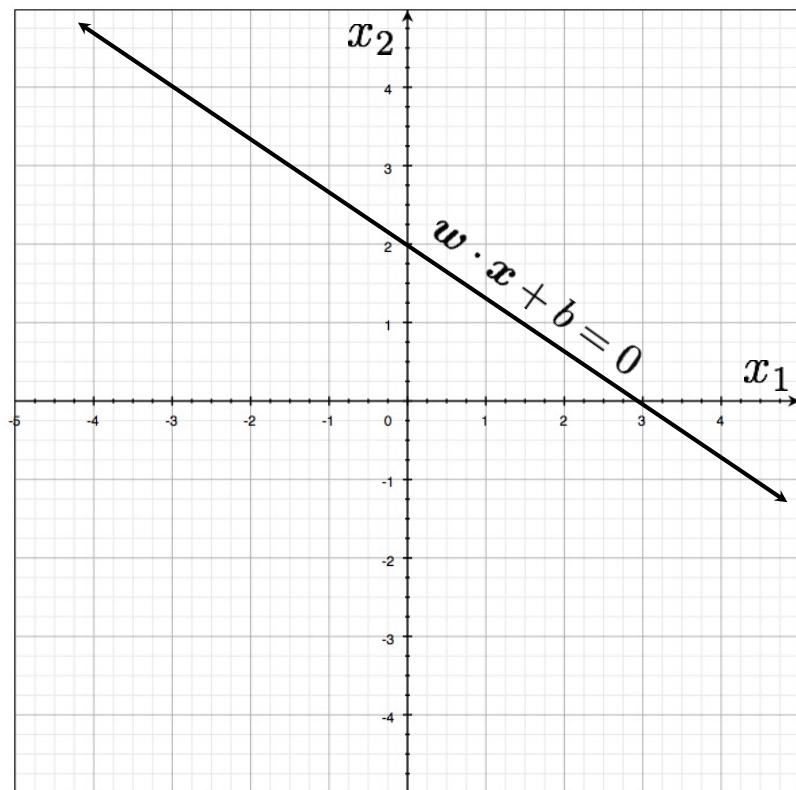
Distribution of data from two classes



First we need to understand hyperplanes...

Hyperplanes (lines) in 2D

$$w_1x_1 + w_2x_2 + b = 0$$



a line can be written as
dot product plus a bias

$$\cancel{w \cdot x + b} = 0$$
$$w \in \mathcal{R}^2$$

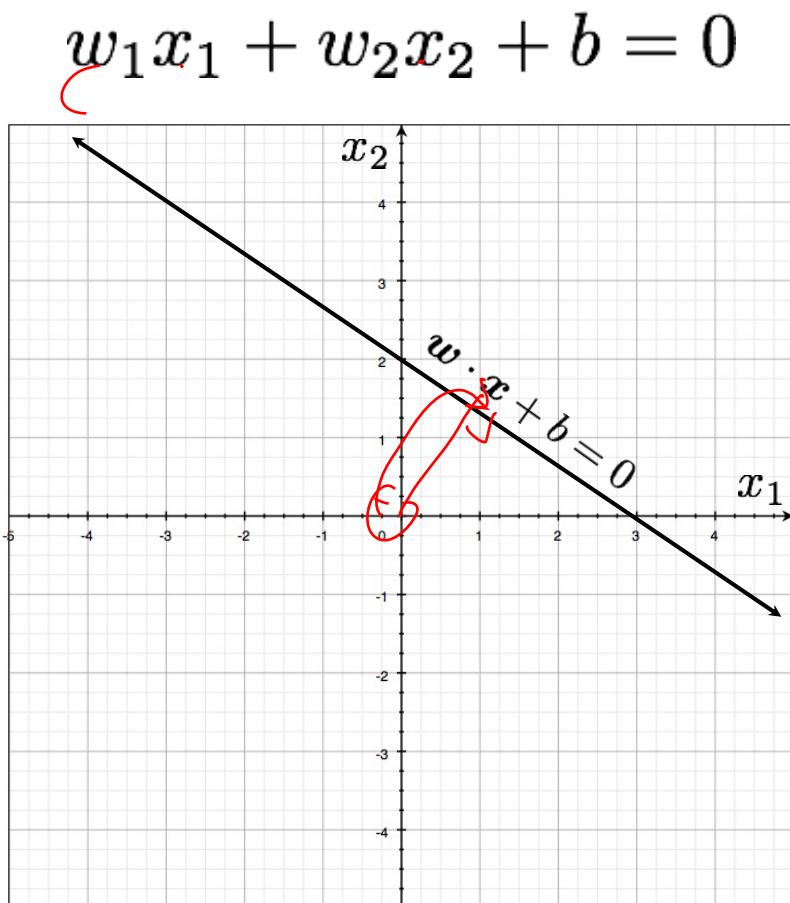
another version, add a weight 1
and push the bias inside

$$\underline{w \cdot x} = 0$$
$$\cancel{w} \in \mathcal{R}^3$$

$$(w, u, b)$$

Hyperplanes (lines) in 2D

$$\mathbf{w} \cdot \mathbf{x} + b = 0 \quad (\text{offset/bias outside}) \quad \mathbf{w} \cdot \mathbf{x} = 0 \quad (\text{offset/bias inside})$$

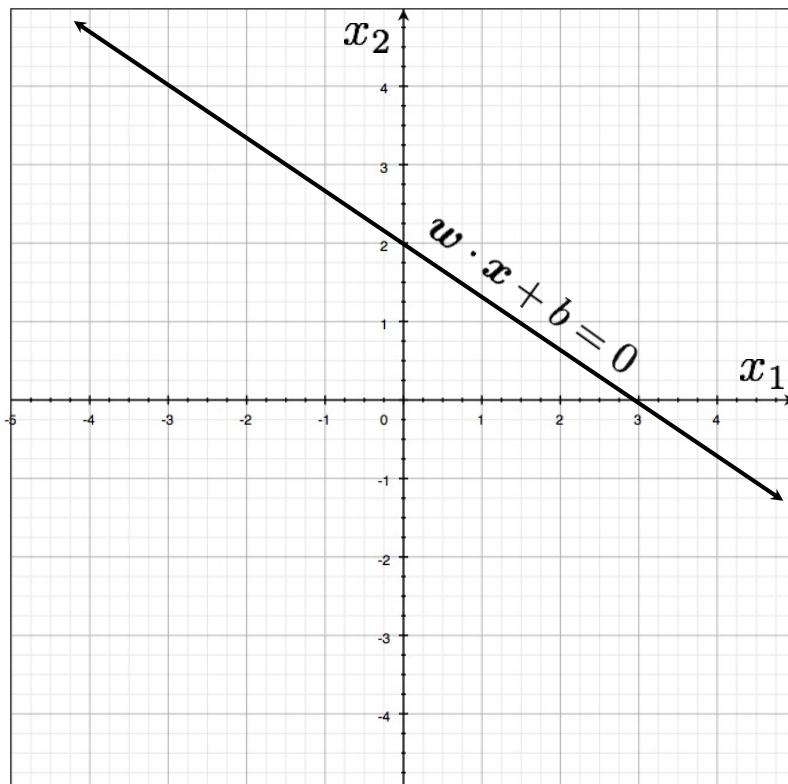


$$\frac{b}{\|\mathbf{w}\|}$$

Hyperplanes (lines) in 2D

$$\mathbf{w} \cdot \mathbf{x} + b = 0 \quad (\text{offset/bias outside}) \quad \mathbf{w} \cdot \mathbf{x} = 0 \quad (\text{offset/bias inside})$$

$$w_1x_1 + w_2x_2 + b = 0$$



Important property:
Free to choose any normalization of w

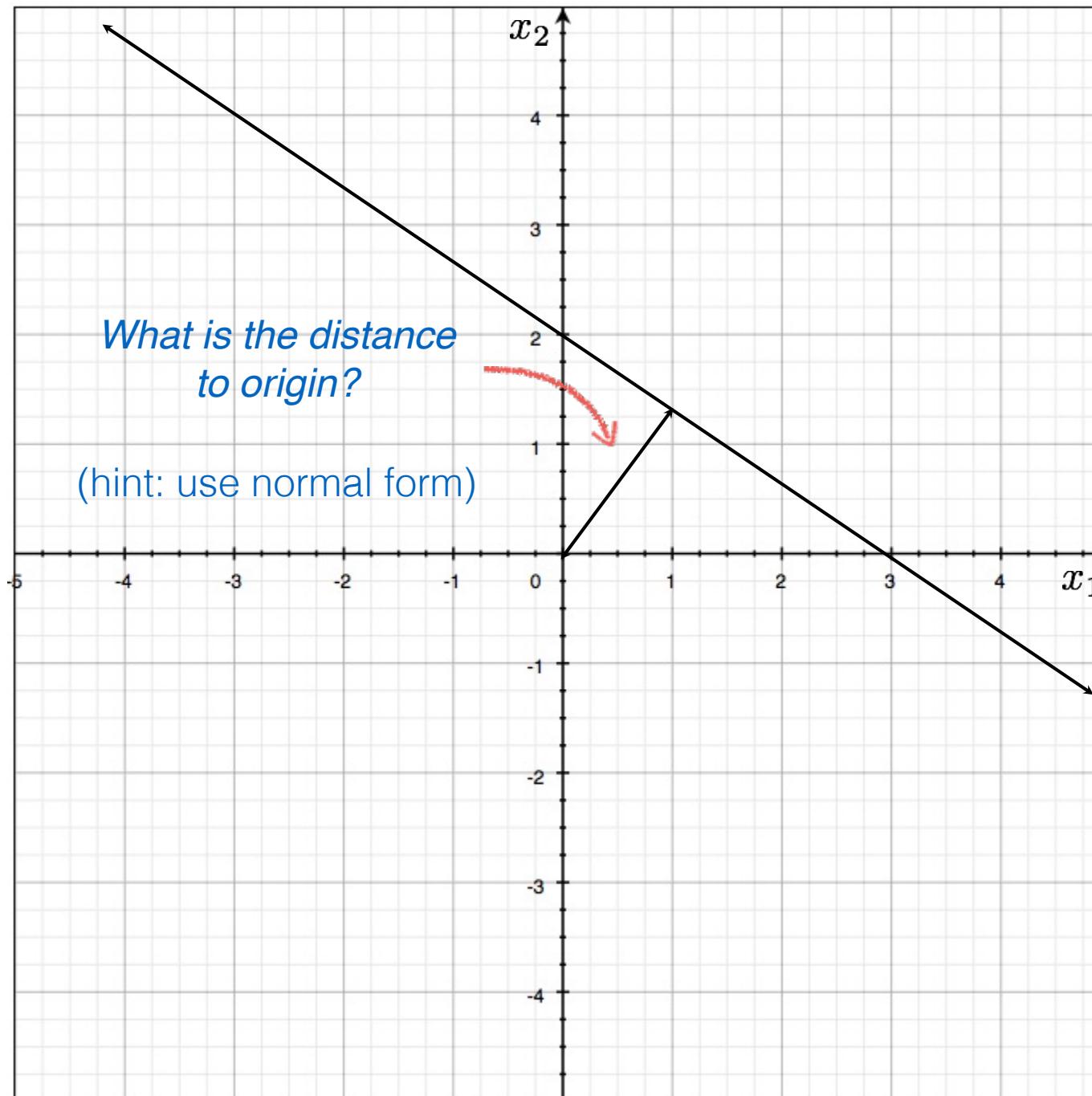
The line

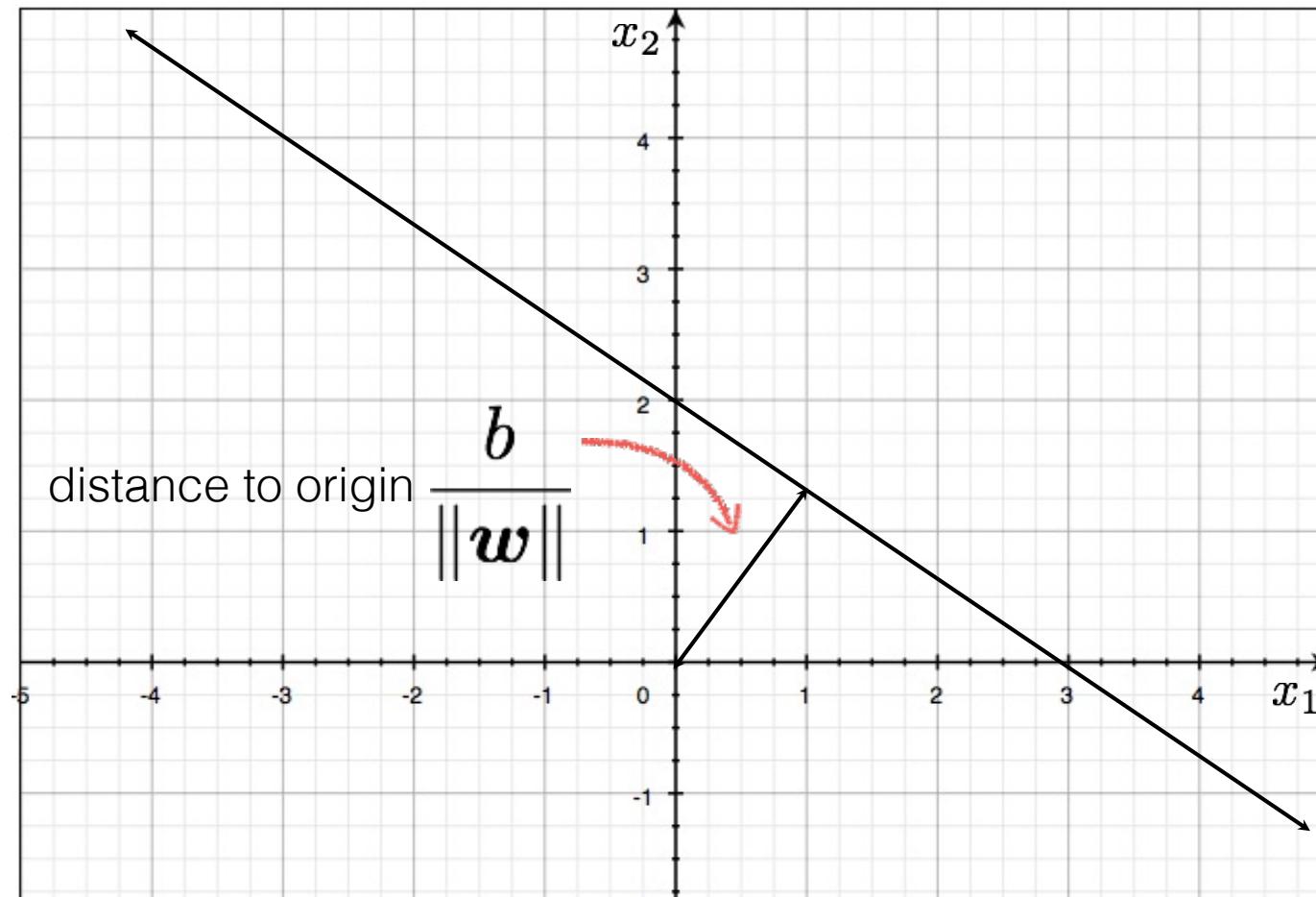
$$w_1x_1 + w_2x_2 + b = 0$$

and the line

$$\lambda(w_1x_1 + w_2x_2 + b) = 0$$

define the same line

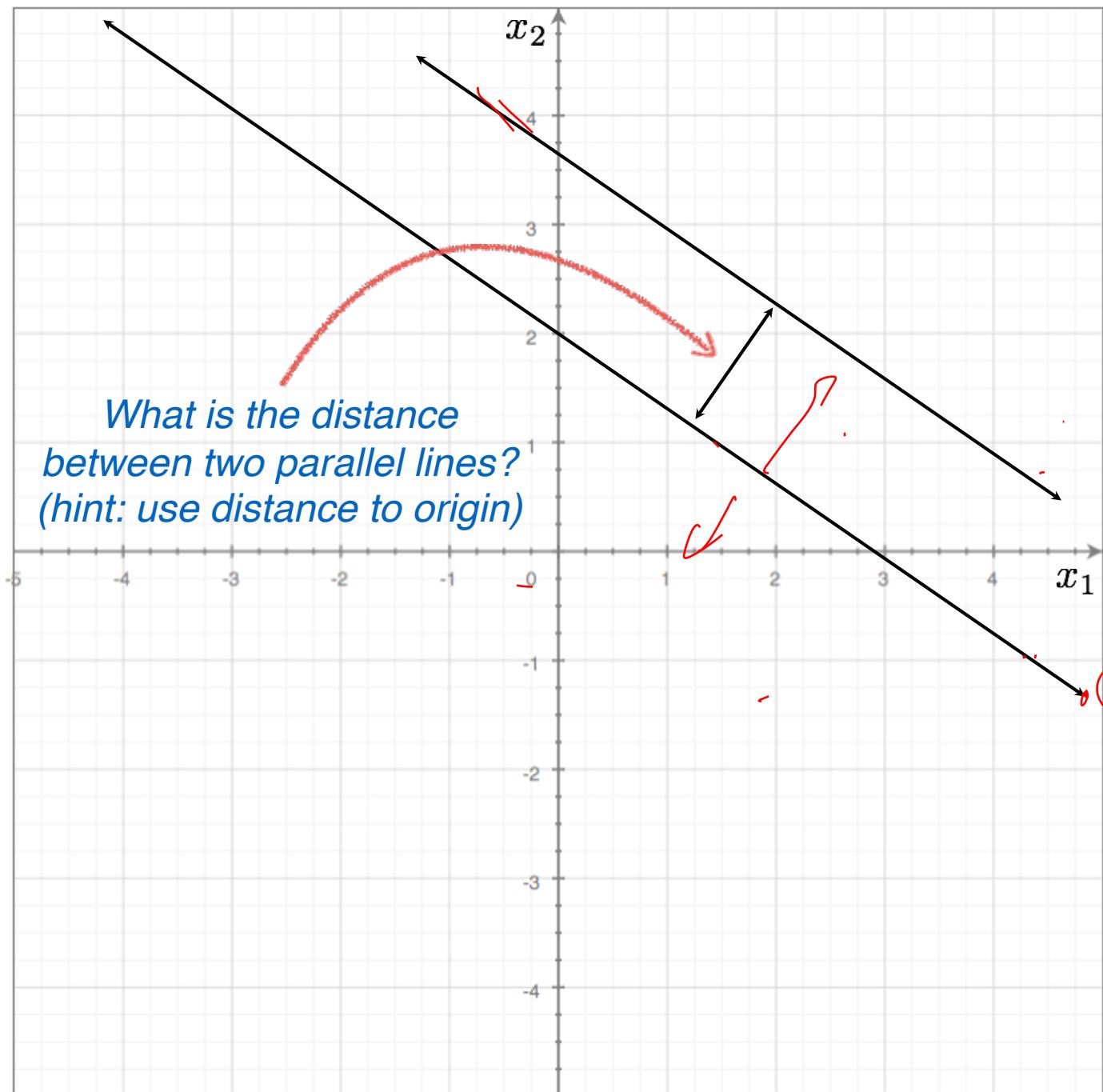




scale $w \cdot x + b = 0$ by $\frac{1}{\|w\|}$

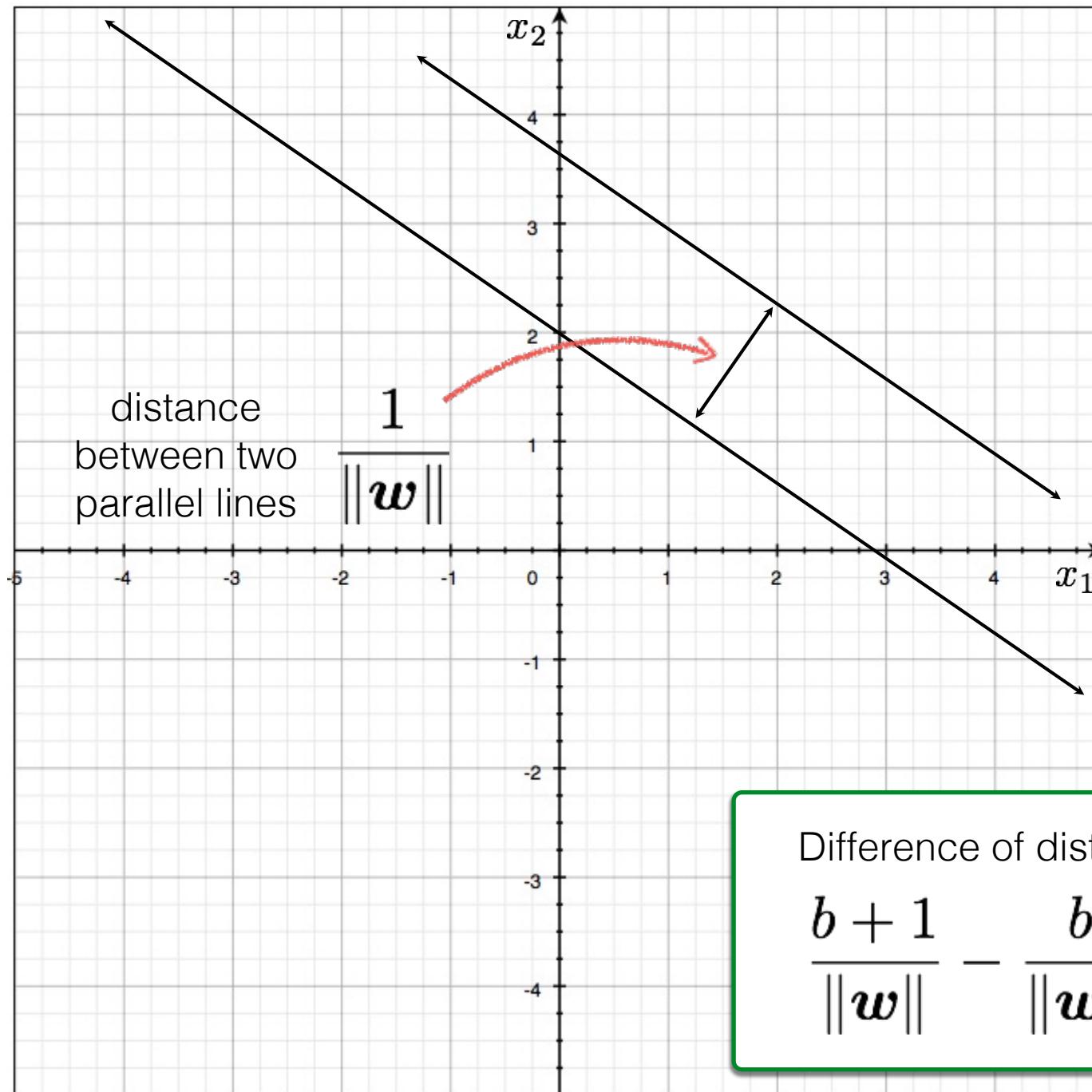
you get the normal form

$$x \cos \theta + y \sin \theta = \rho$$



$$\underline{w \cdot x + b} = -1$$

$$w \cdot x + b = 0$$

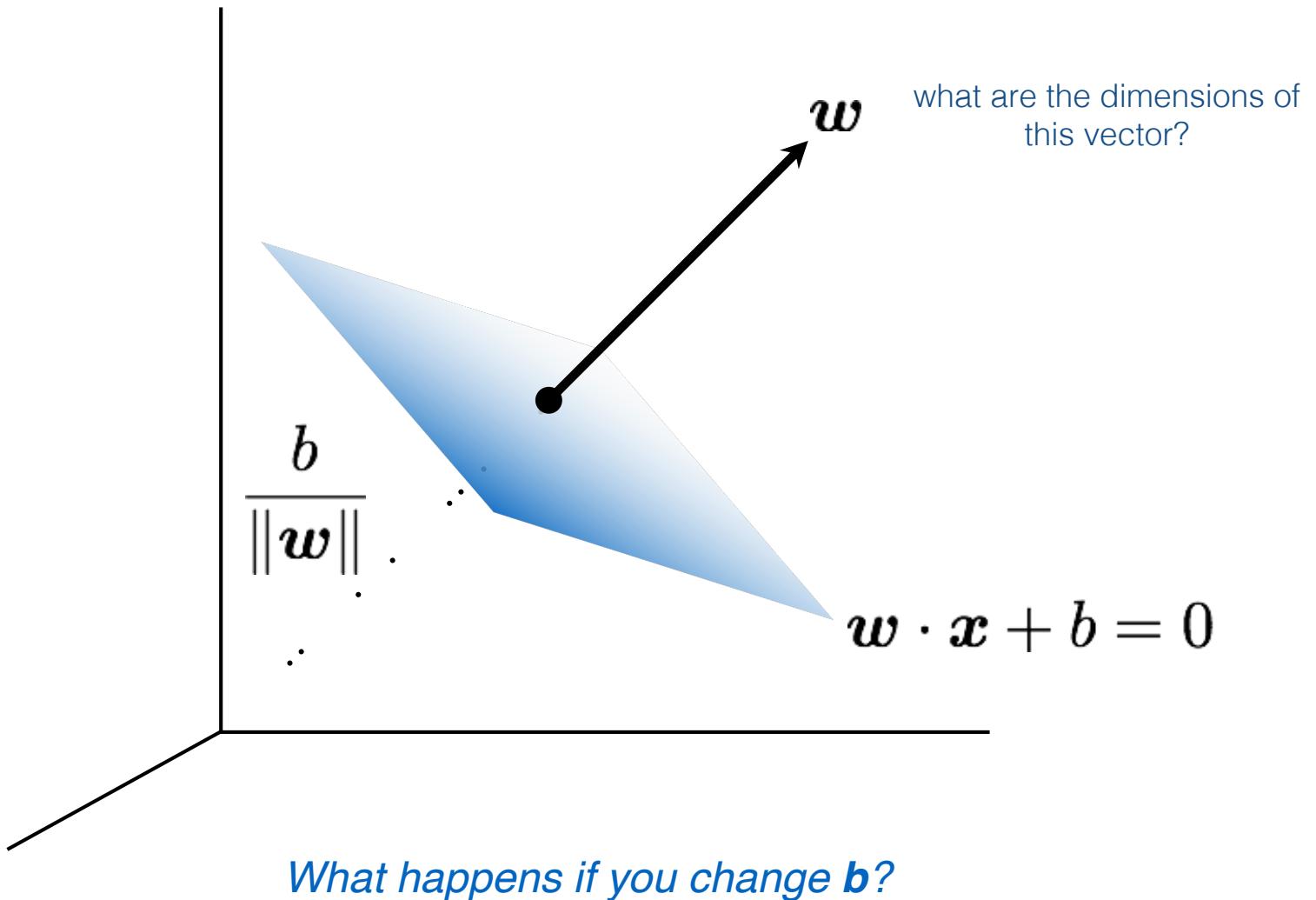


Difference of distance to origin

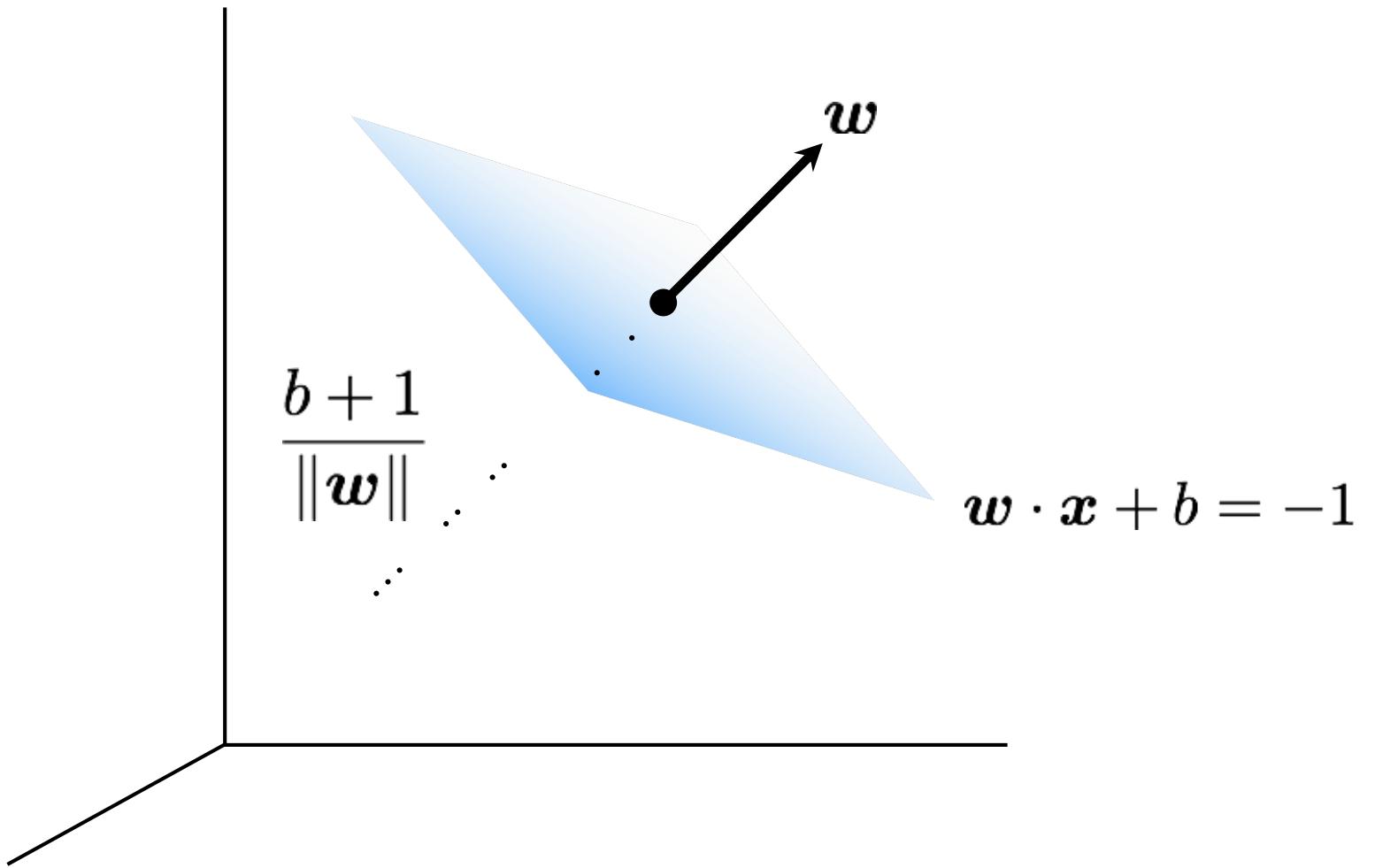
$$\frac{b+1}{\|w\|} - \frac{b}{\|w\|} = \frac{1}{\|w\|}$$

Now we can go to 3D ...

Hyperplanes (planes) in 3D

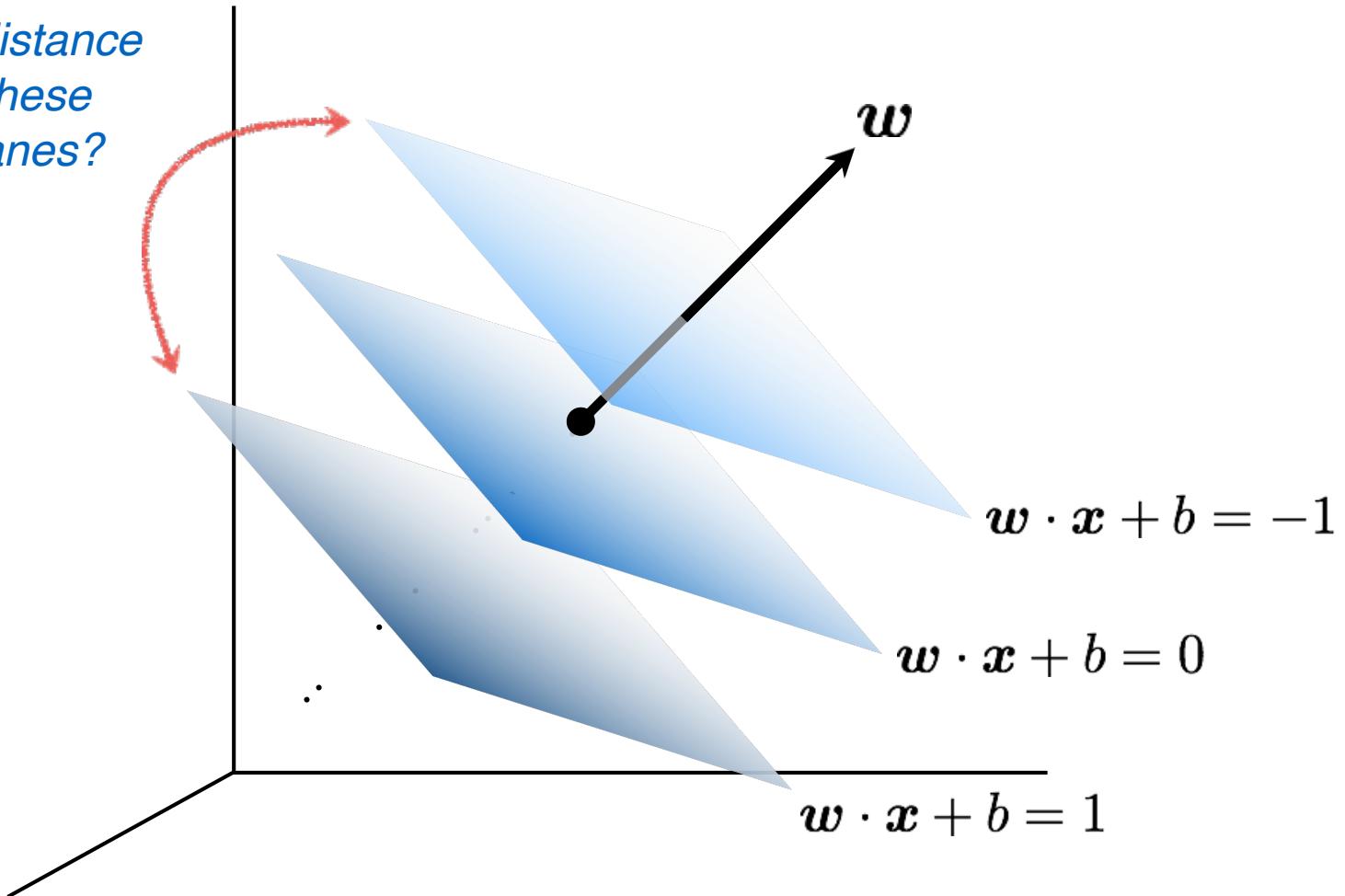


Hyperplanes (planes) in 3D

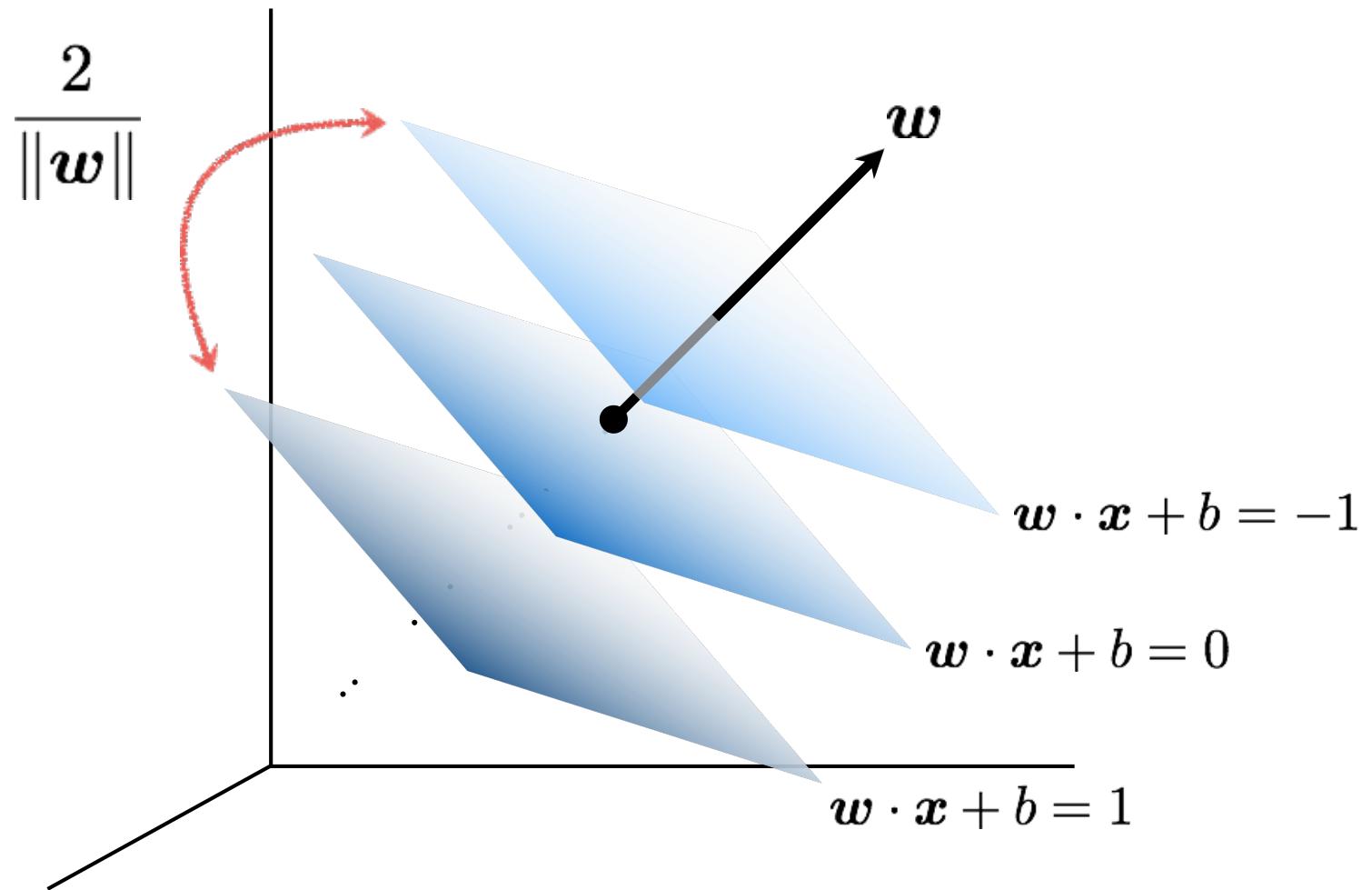


Hyperplanes (planes) in 3D

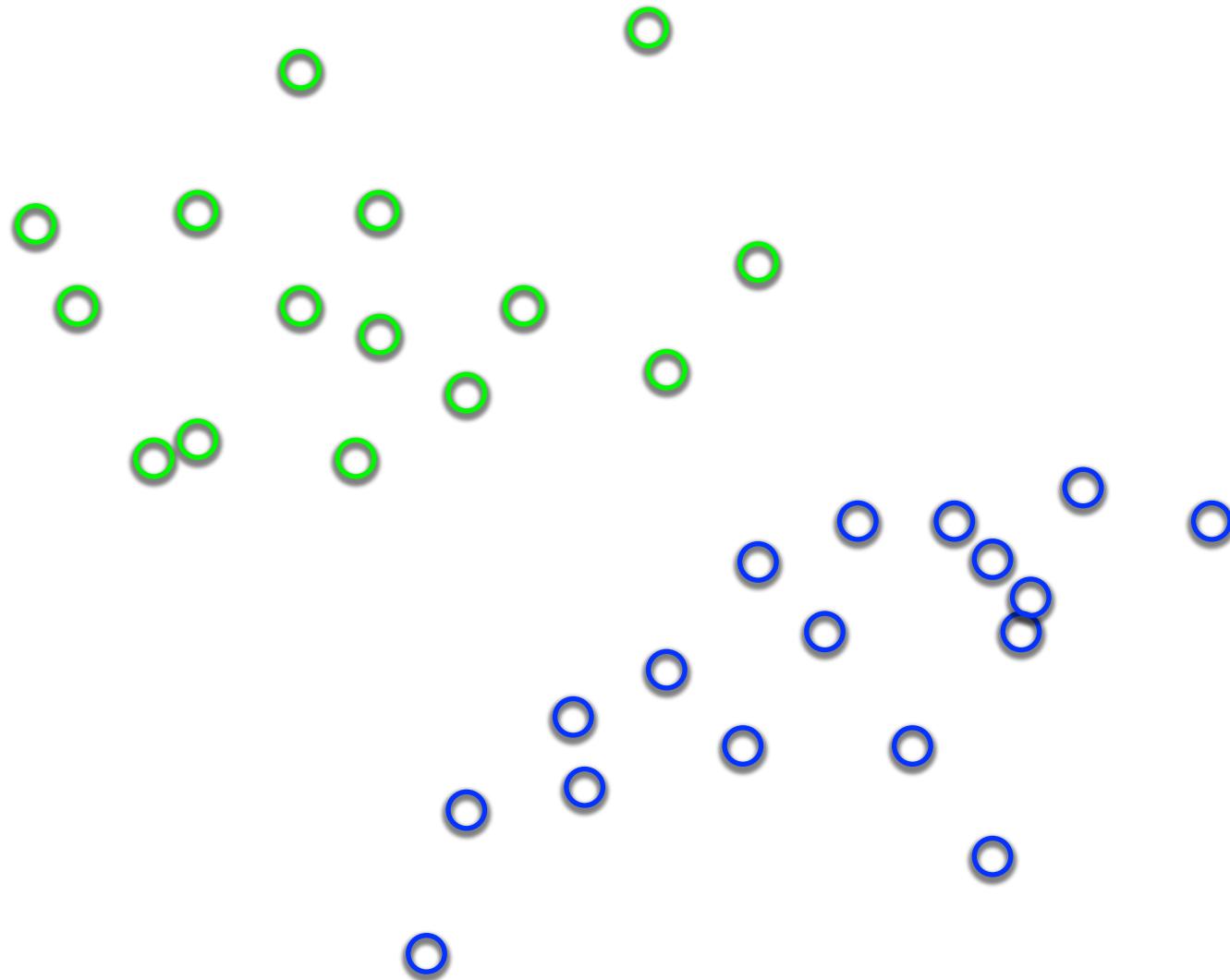
*What's the distance
between these
parallel planes?*



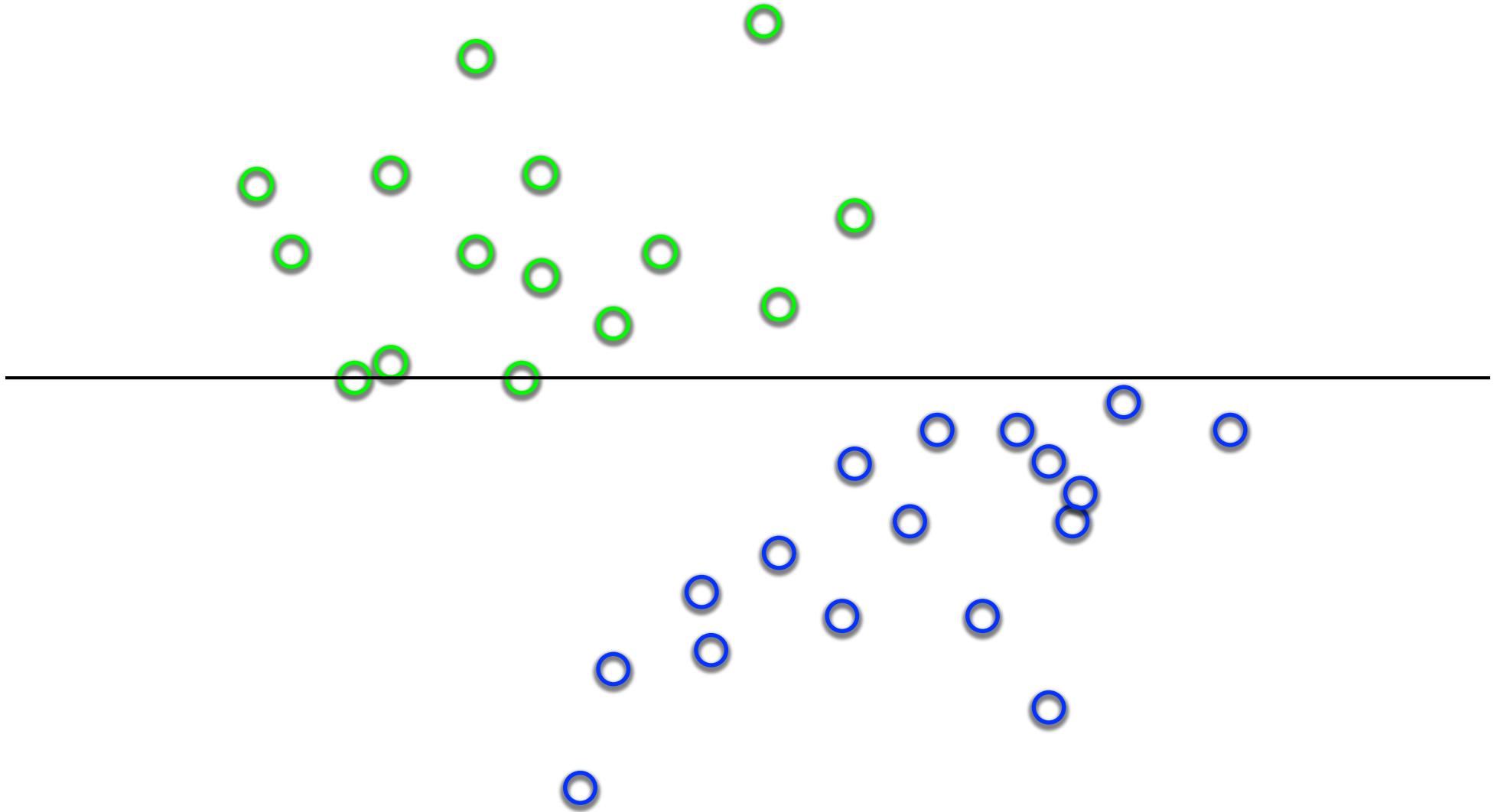
Hyperplanes (planes) in 3D



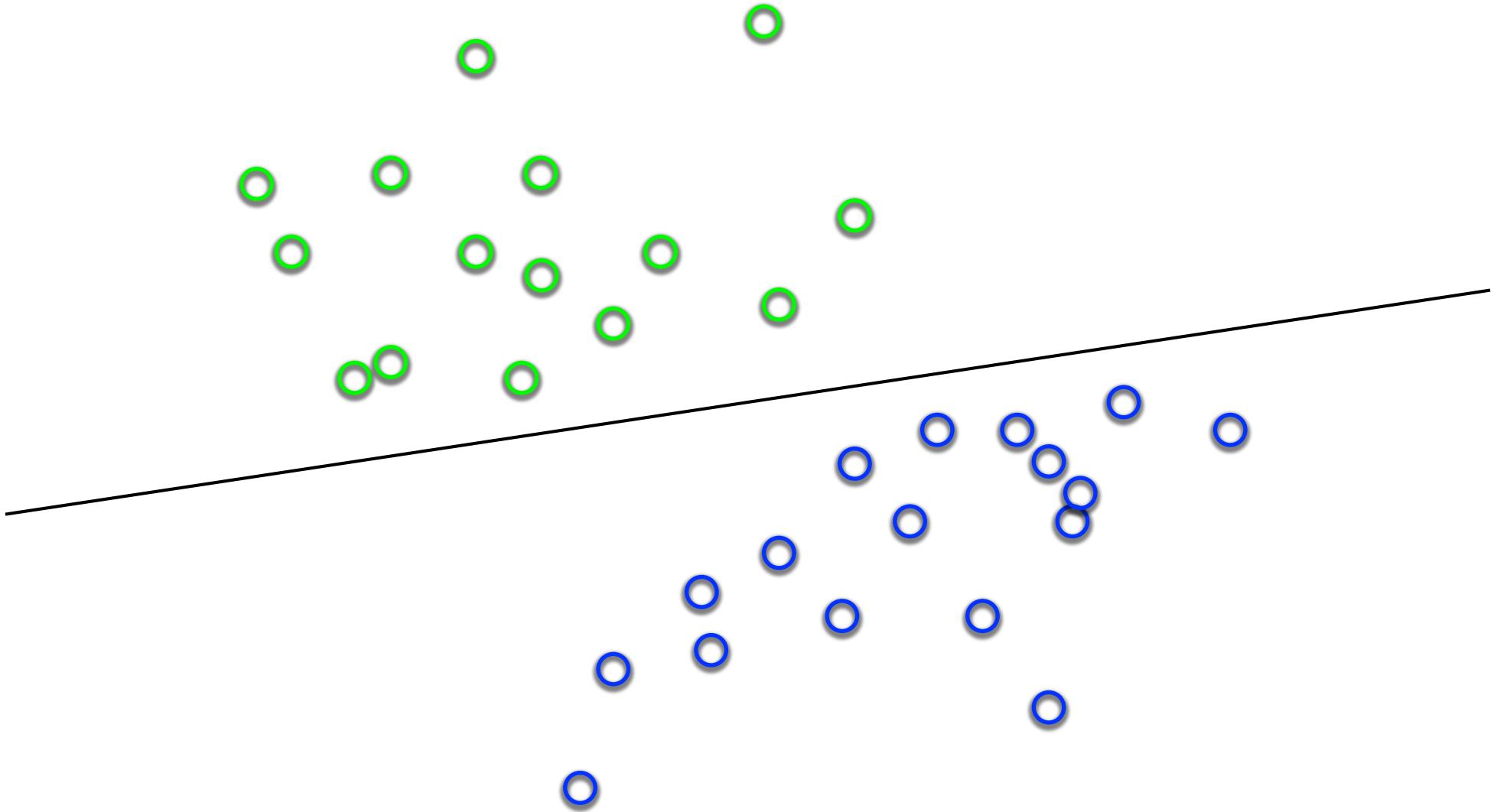
What's the best w ?



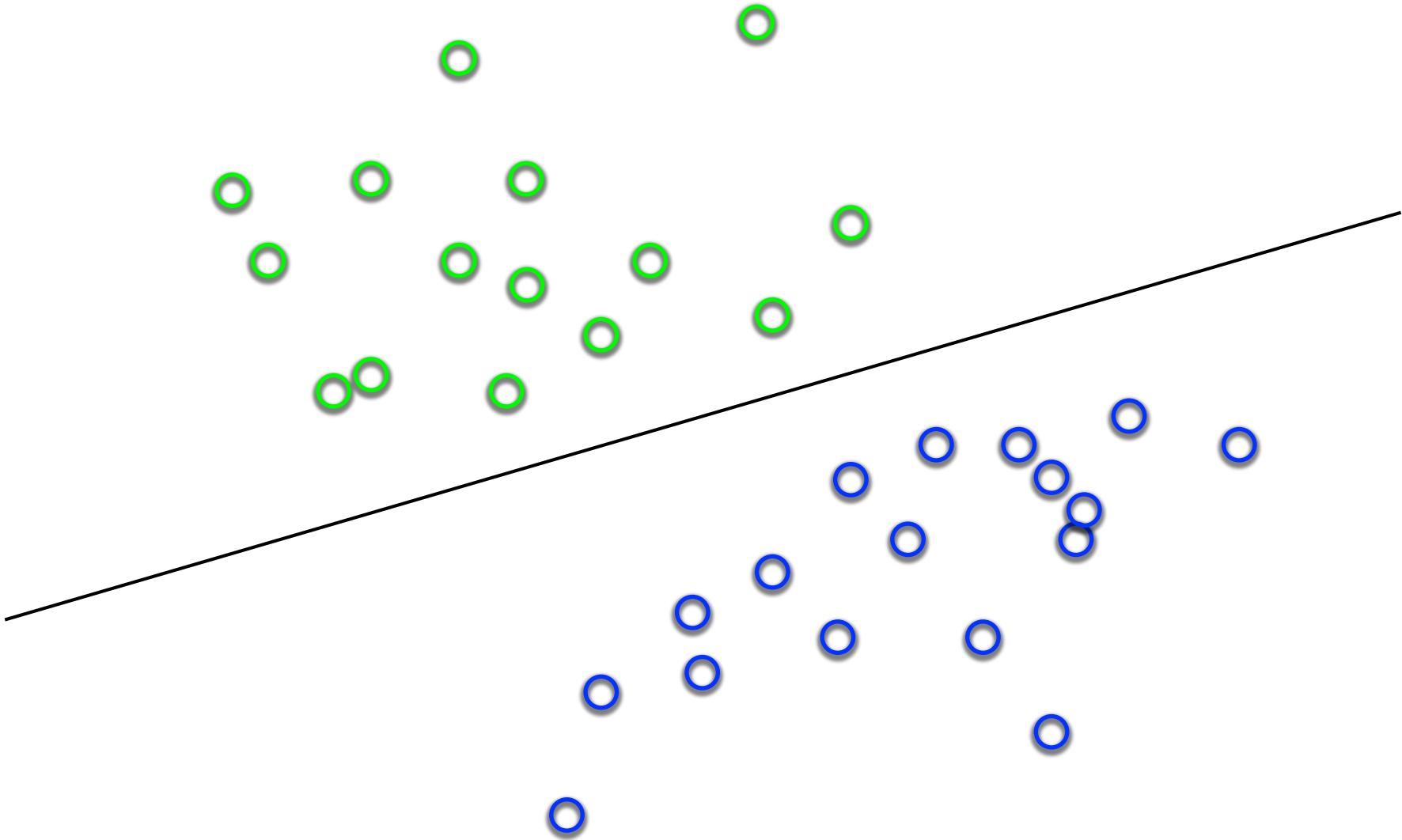
What's the best w ?



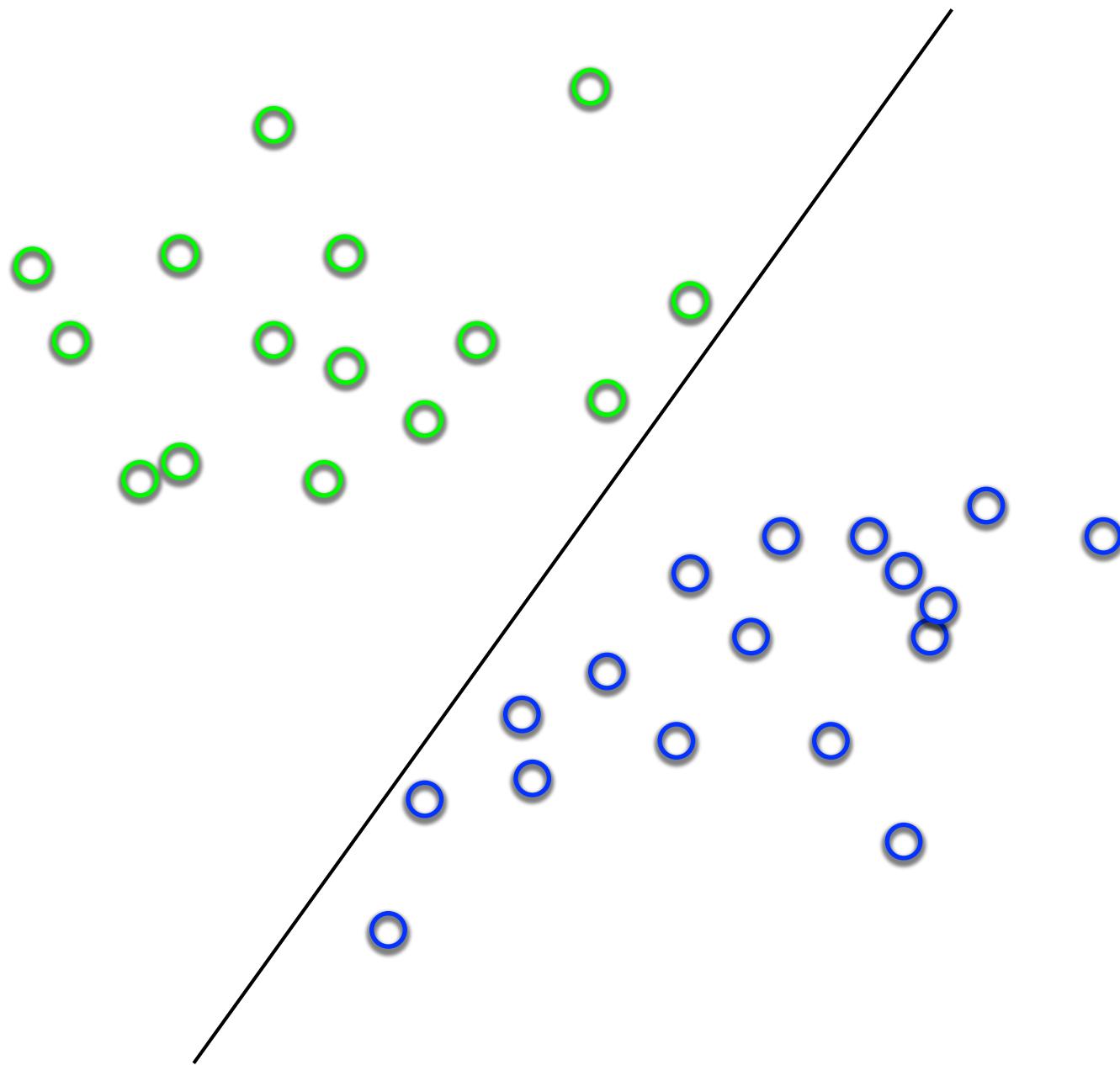
What's the best w ?



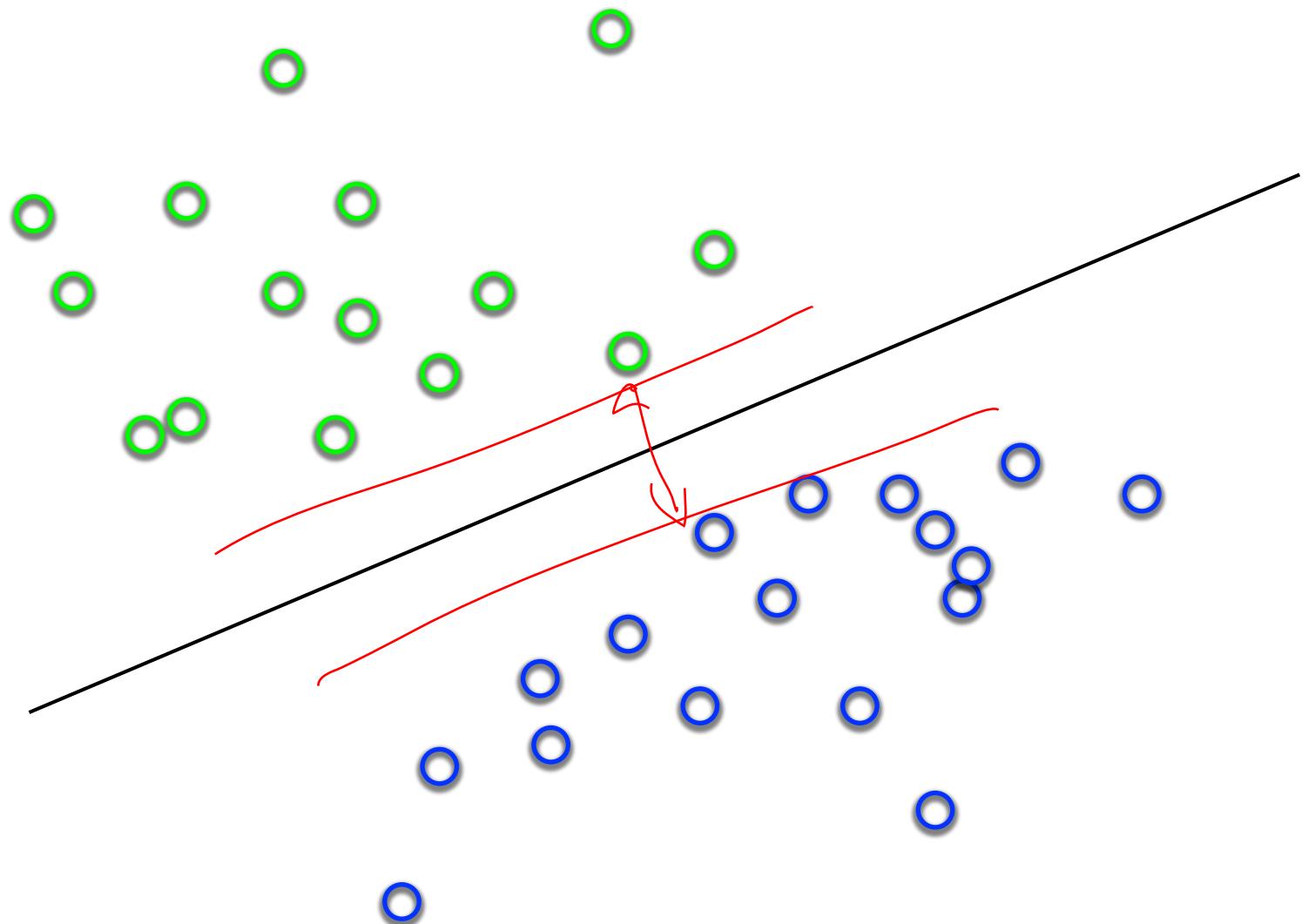
What's the best w ?



What's the best w ?

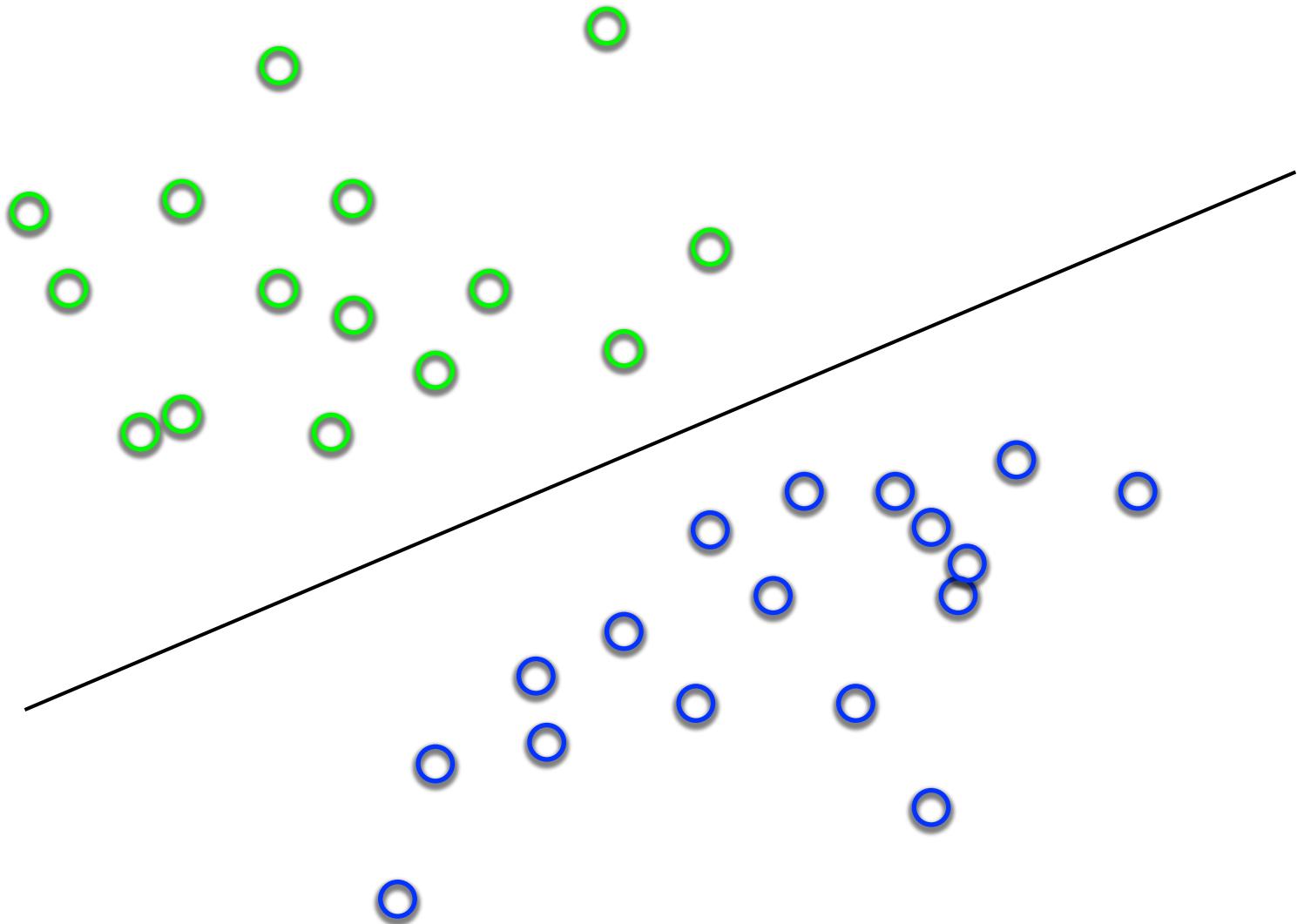


What's the best w ?



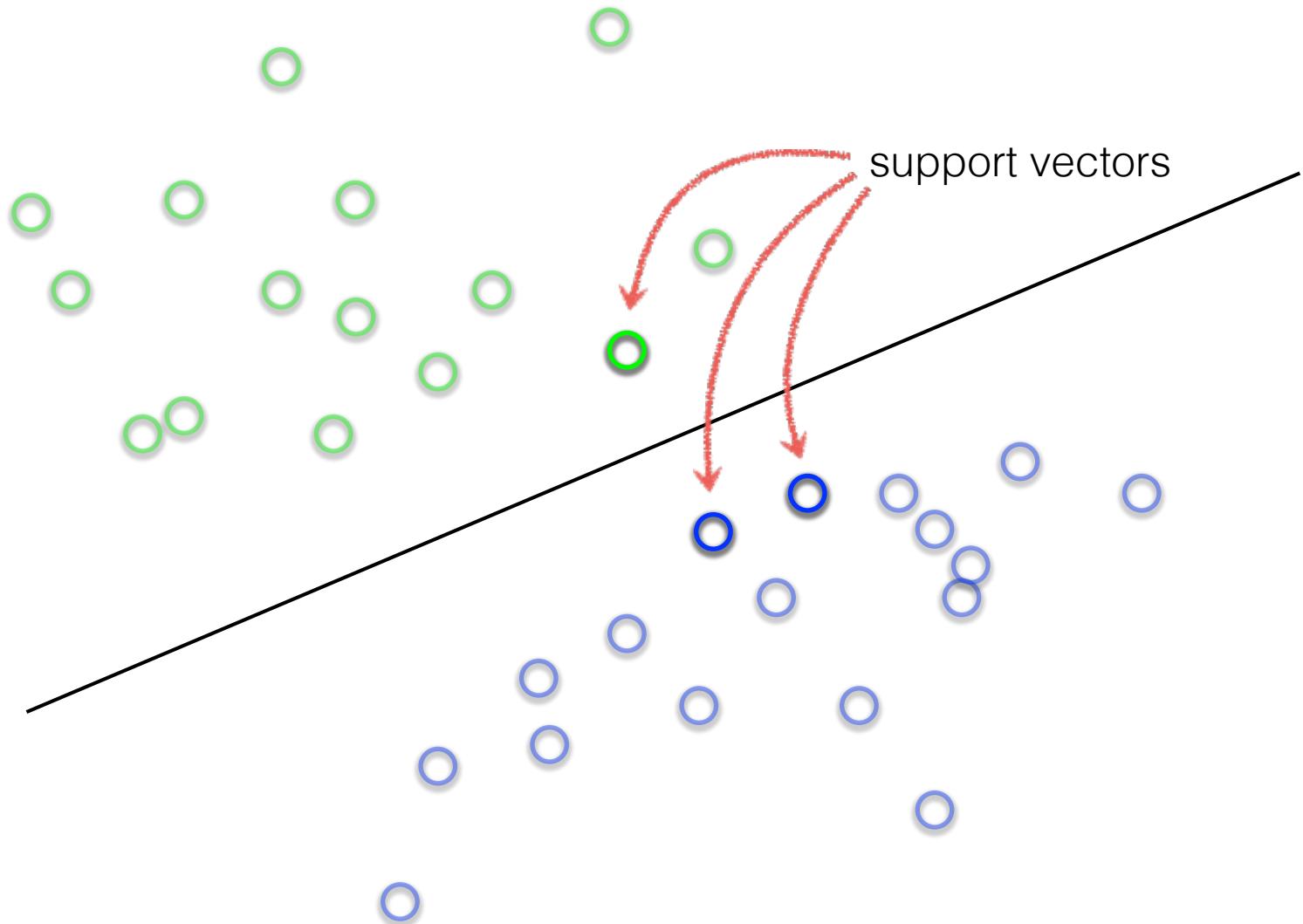
Intuitively, the line that is the
farthest from all interior points

What's the best w ?



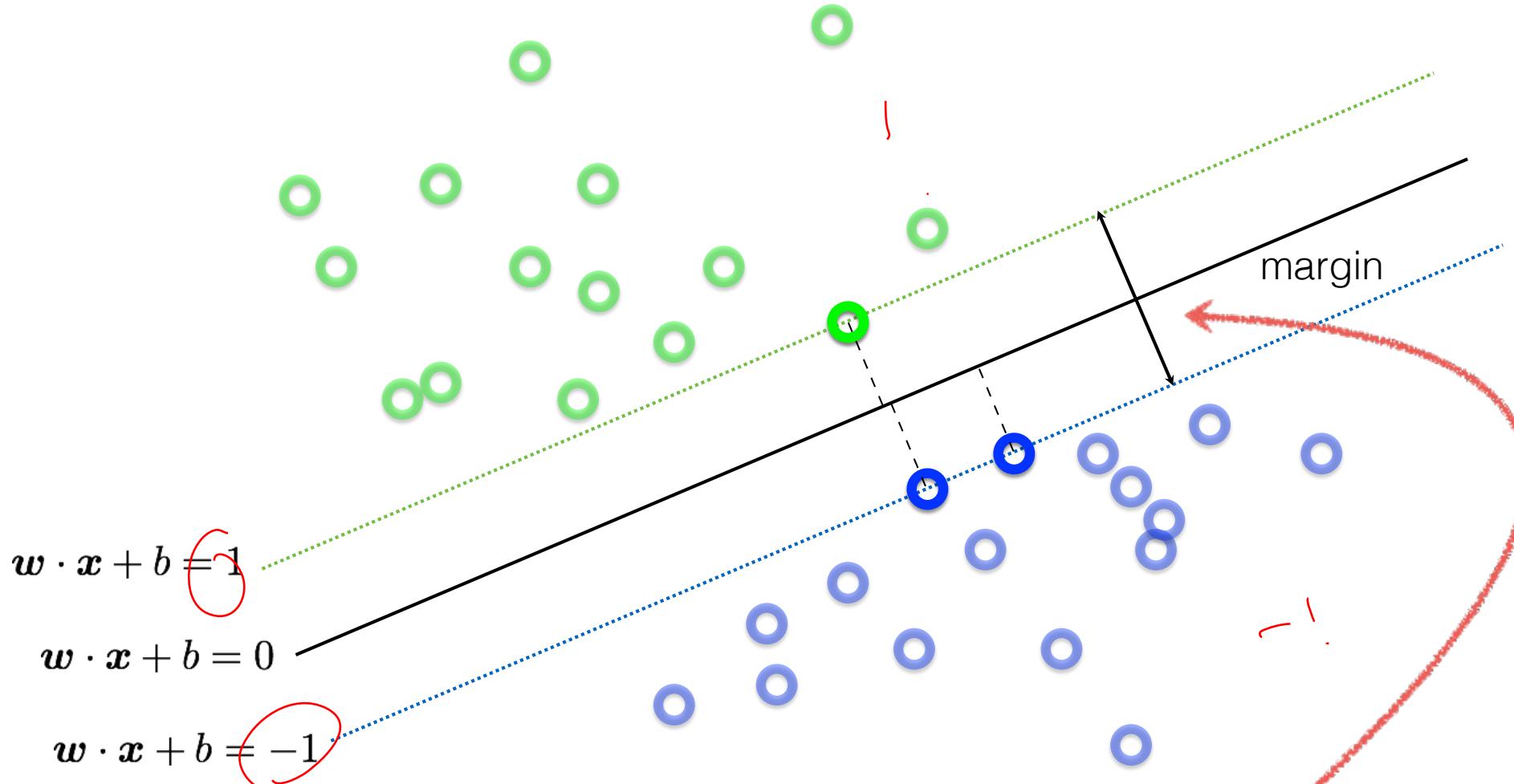
Maximum Margin solution:
most stable to perturbations of data

What's the best \mathbf{w} ?



Want a hyperplane that is far away from 'inner points'

Find hyperplane \mathbf{w} such that ...



the gap between parallel hyperplanes $\frac{2}{\|\mathbf{w}\|}$ is maximized

Can be formulated as a maximization problem

$$\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$$

subject to $\mathbf{w} \cdot \mathbf{x}_i + b \geq +1$ if $y_i = +1$
 $\mathbf{w} \cdot \mathbf{x}_i + b \leq -1$ if $y_i = -1$ for $i = 1, \dots, N$

What does this constraint mean?



label of the data point

Why is it +1 and -1?

Can be formulated as a maximization problem

$$\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$$

subject to $\mathbf{w} \cdot \mathbf{x}_i + b \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases} \text{ for } i = 1, \dots, N$

Equivalently,

Where did the 2 go?

$$\min_{\mathbf{w}} \|\mathbf{w}\|$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \text{ for } i = 1, \dots, N$

What happened to the labels?

'Primal formulation' of a linear SVM

$$\min_{\mathbf{w}} \|\mathbf{w}\|$$

Objective Function

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$ for $i = 1, \dots, N$

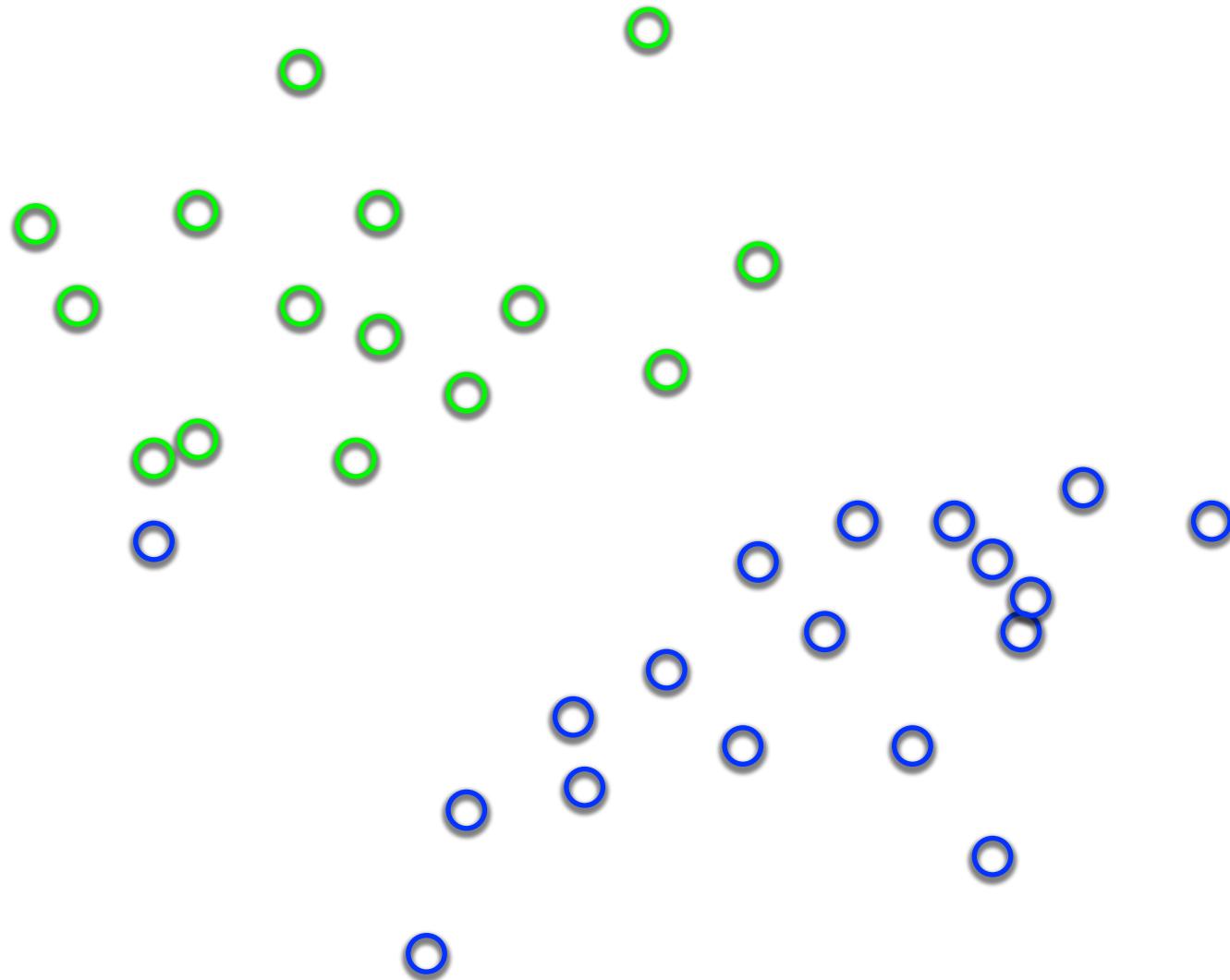
Constraints

This is a convex quadratic programming (QP) problem
(a unique solution exists)

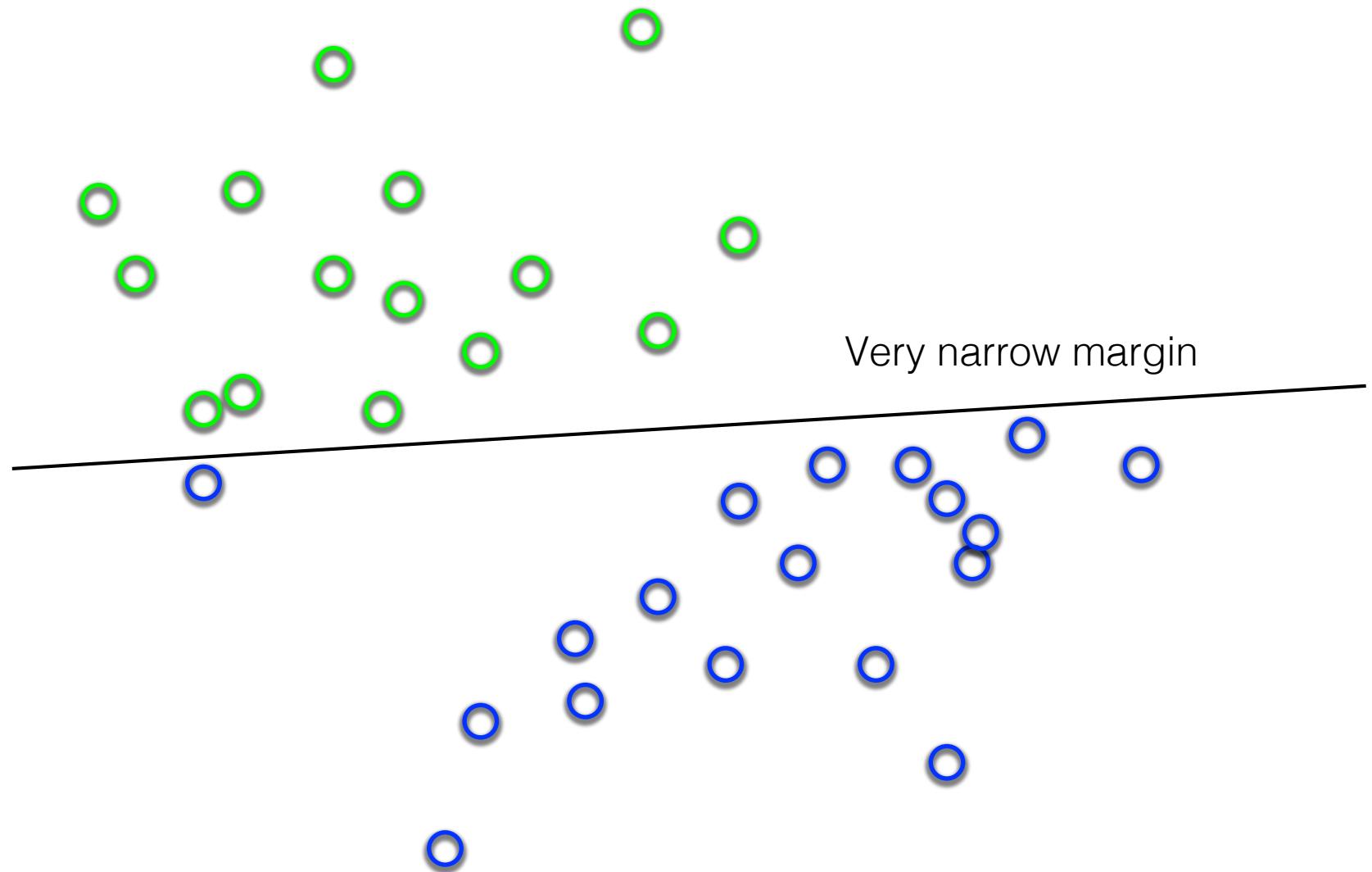
<https://scikit-learn.org/stable/modules/svm.html>

‘soft’ margin

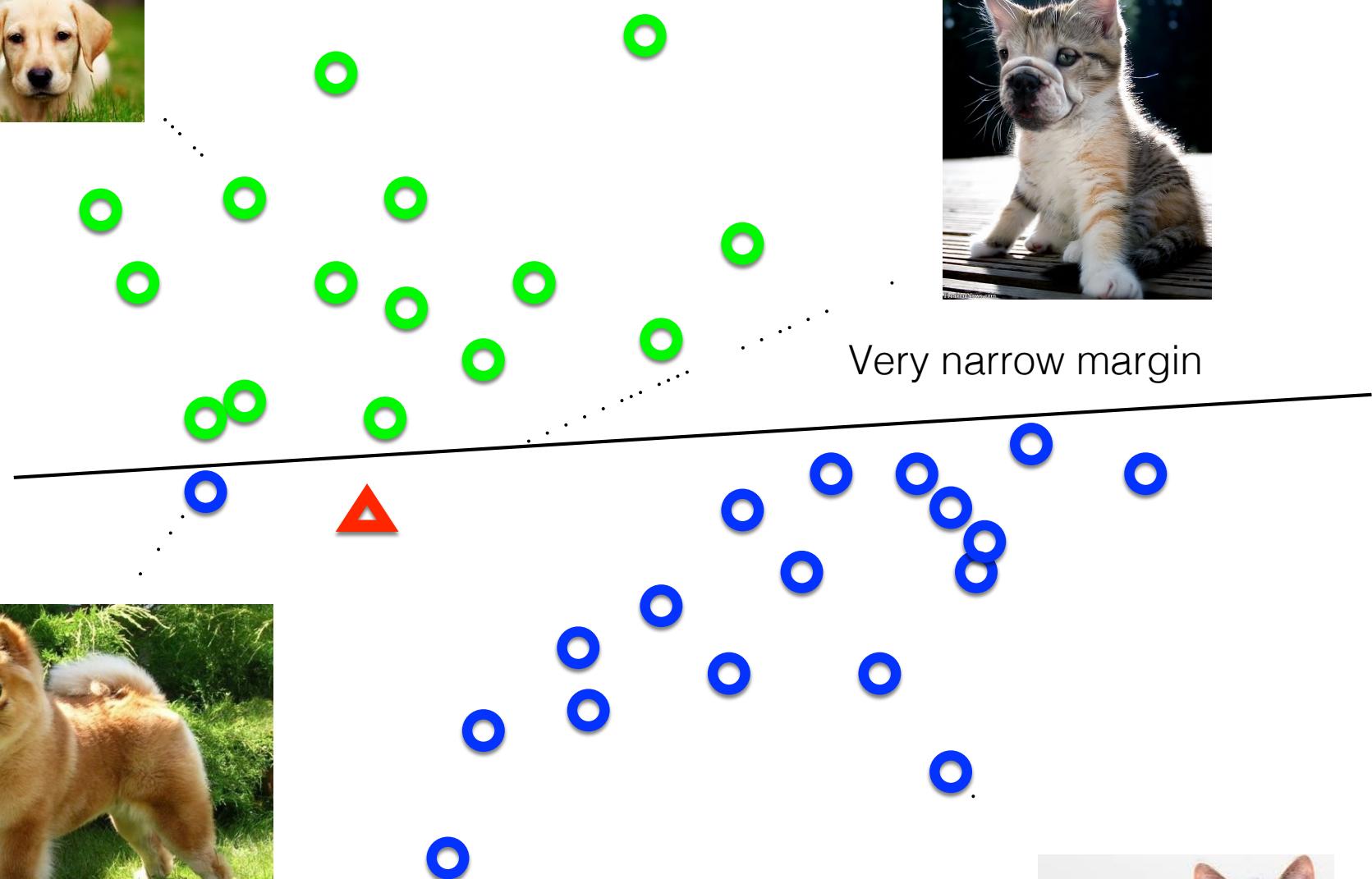
What's the best w ?



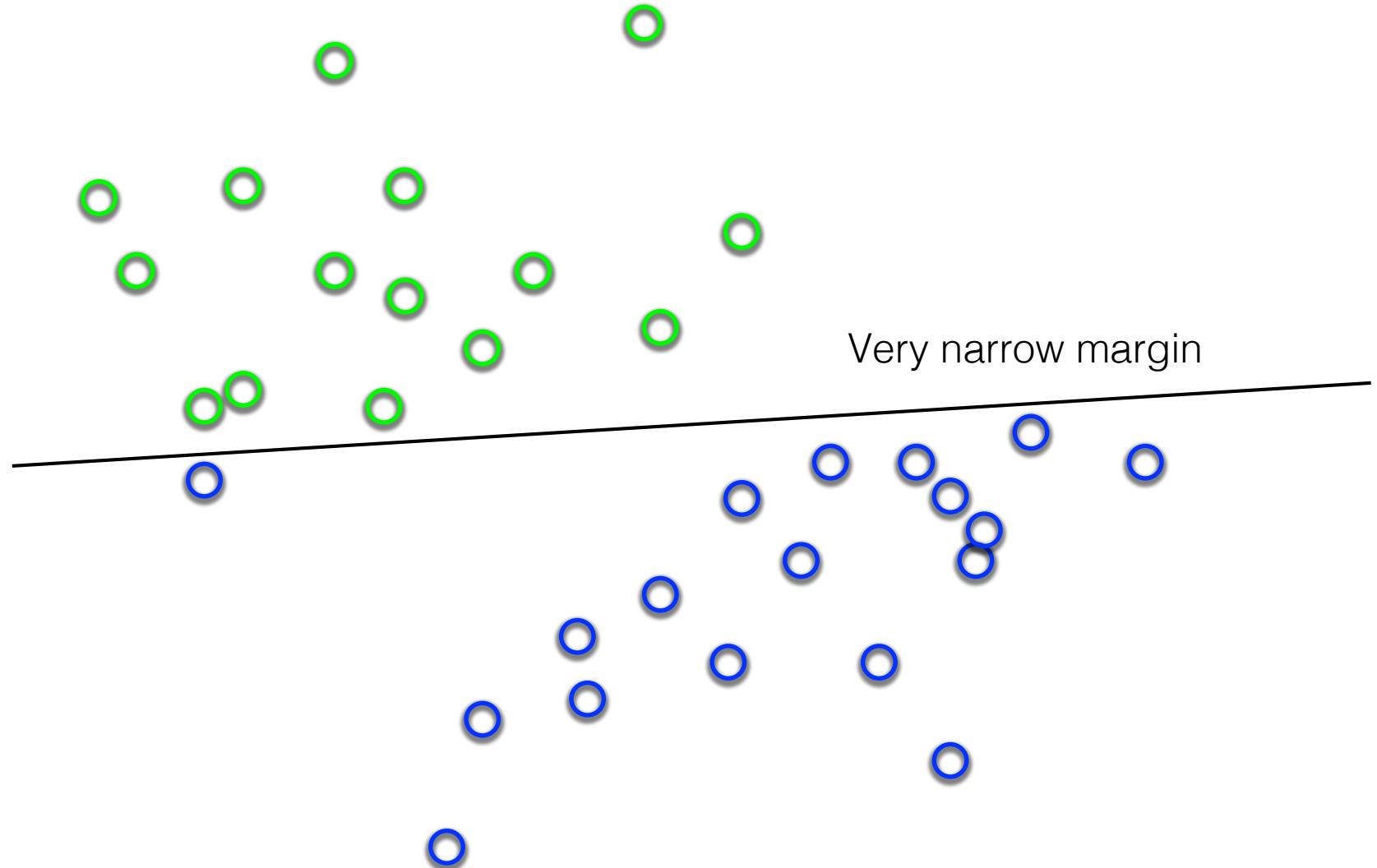
What's the best w ?



Separating cats and dogs

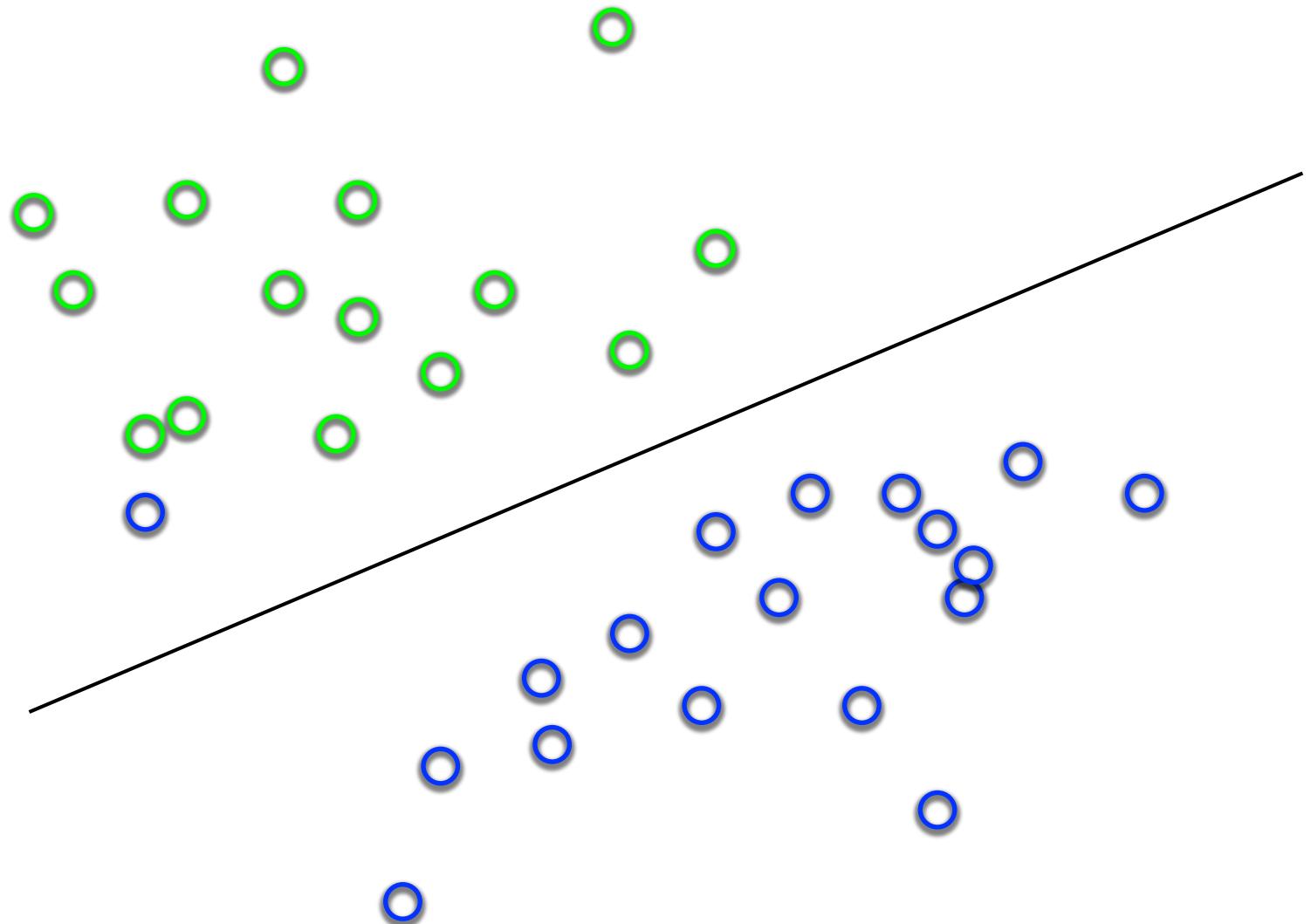


What's the best w ?



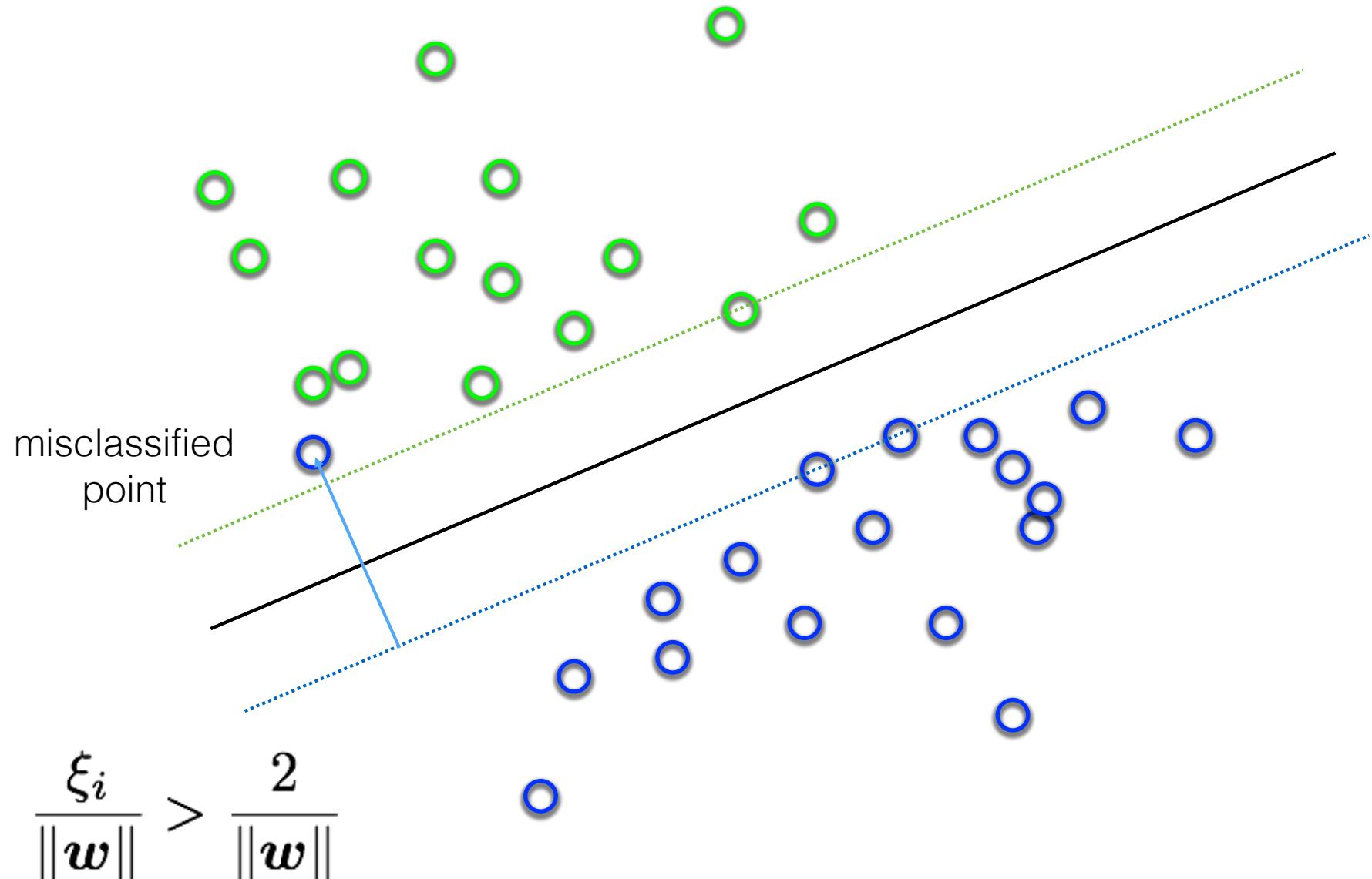
Intuitively, we should allow for some misclassification if we can get more robust classification

What's the best w ?



Trade-off between the MARGIN and the MISTAKES
(might be a better solution)

Adding slack variables $\xi_i \geq 0$



‘soft’ margin

objective

$$\min_{\boldsymbol{w}, \boldsymbol{\xi}} \|\boldsymbol{w}\|^2 + C \sum_i \xi_i$$

subject to

$$y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) \geq 1 - \xi_i \quad \text{for } i = 1, \dots, N$$

'soft' margin

objective

$$\min_{\mathbf{w}, \xi} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

subject to

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{for } i = 1, \dots, N$$

The slack variable allows for mistakes,
as long as the inverse margin is minimized.

‘soft’ margin

objective

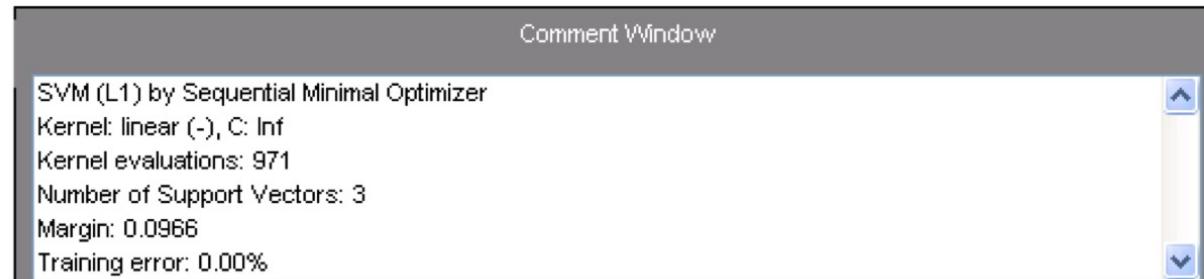
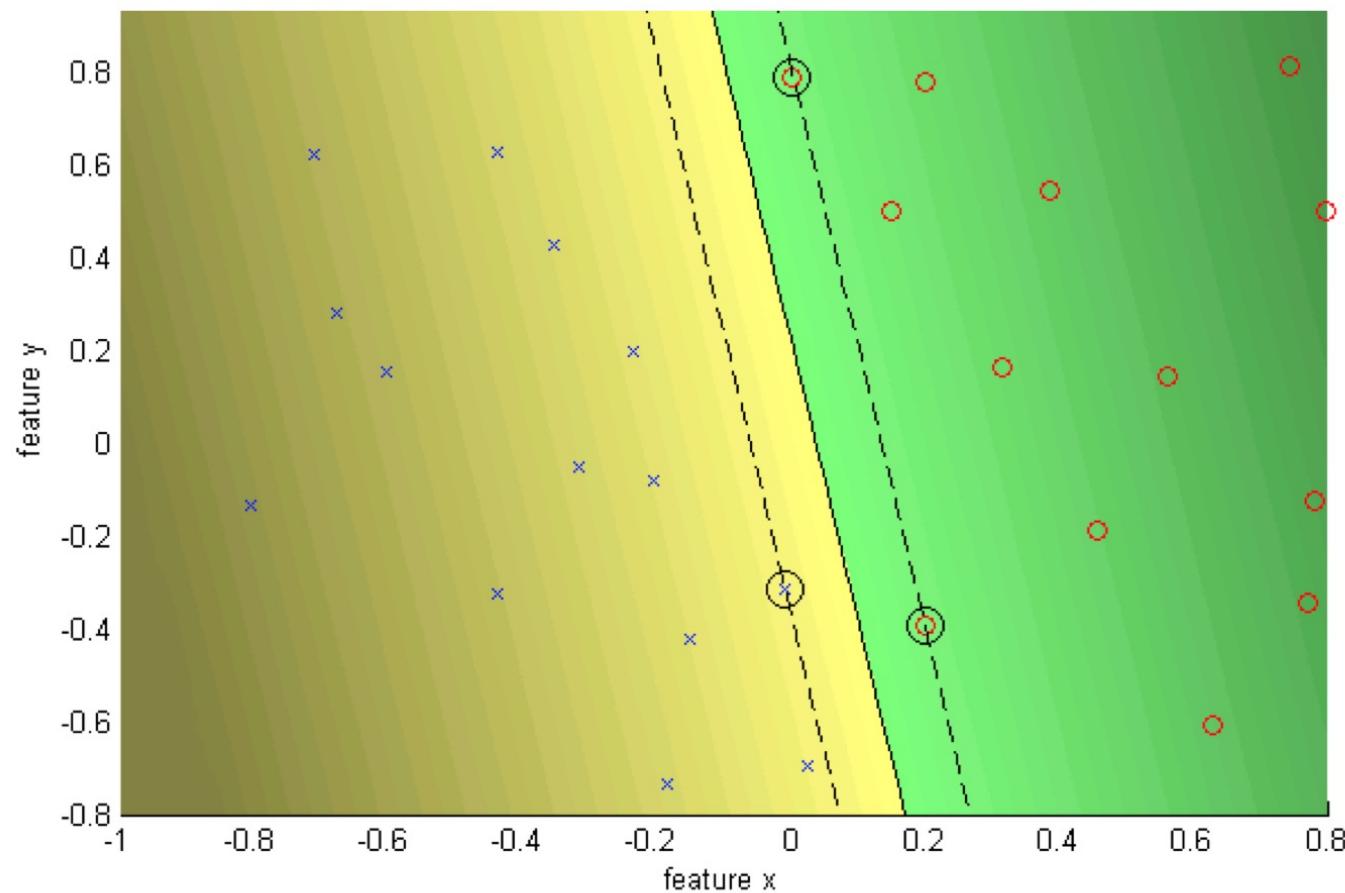
$$\min_{\boldsymbol{w}, \boldsymbol{\xi}} \|\boldsymbol{w}\|^2 + C \sum_i \xi_i$$

subject to

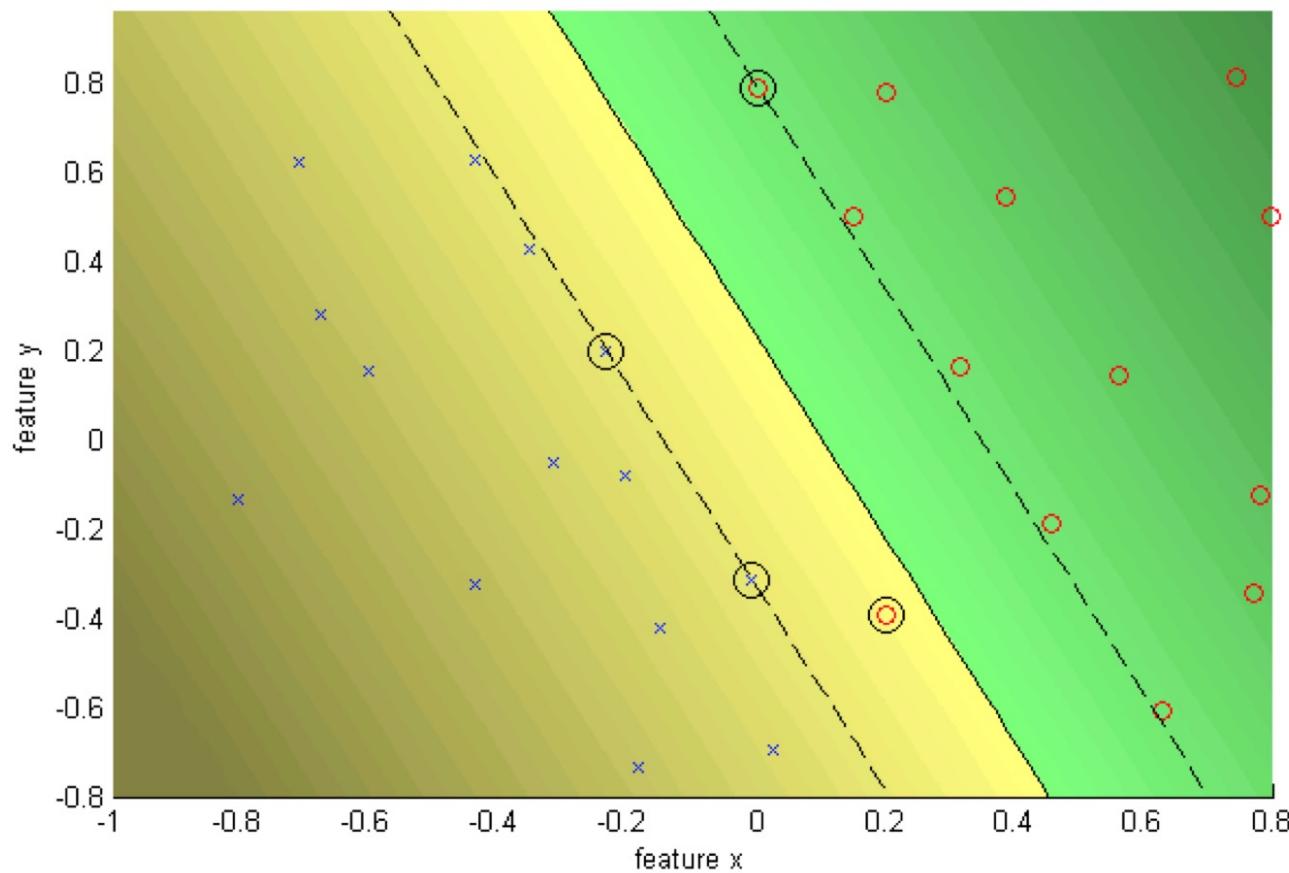
$$y_i(\boldsymbol{w}^\top \boldsymbol{x}_i + b) \geq 1 - \xi_i \quad \text{for } i = 1, \dots, N$$

- Every constraint can be satisfied if slack is large
- C is a regularization parameter
 - Small C: ignore constraints (larger margin)
 - Big C: constraints (small margin)
- Still QP problem (unique solution)

$C = \text{Infinity}$ hard margin



$C = 10$ soft margin



Comment Window

SVM (L1) by Sequential Minimal Optimizer
Kernel: linear (-), C: 10.0000
Kernel evaluations: 2645
Number of Support Vectors: 4
Margin: 0.2265
Training error: 3.70%

Nearest Neighbor is competitive

40281508803277064755529284686500876/71127400776386420140578274711366
50711167679664143112241087634006330217113109975414895351982739901029
8468682467933943144705960444612336459685608641865284554770782237018
7695346501828357808571101378507110114527623028596972136418240510226
93477149069842728100783331376131605747595849918501320348220251514889
82049962335648092836457294912860704116759914592504108408989425198980
3551721691955162286714604033223689853854520563283995719467131366090/
94368160413174951001162198403649071654525185470670258104571851900607
88573898868239756292881688791801720751902098623938021111429725/12199
14853434175074881539597690363982212868553949251514414435912233029009
9319097549201081493361525220266012030255795808950325908846884548549
692854579921834078393465623926061287982047750564674307507420897404
12845278113035703193631773084826529739099642972116747598821445161325
90662367728608302983253880019513960141712379749939282718091017796999
21010452828351781129784030788497858498138031785516574935471208160734
28308784084458566309376893495891288681379011470817457121130621280766
41992780136134111560701232522949812/61274000822292799275134941856283

MNIST Digit Recognition

- Handwritten digits
- 28x28 pixel images: $d = 784$
- 60,000 training samples
- 10,000 test samples

Yann LeCunn

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