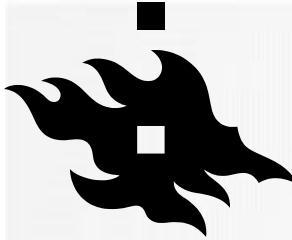


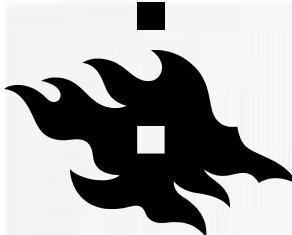
TODAY'S LECTURE

- Object detection and recognition
- Basic approaches with more traditional computer vision techniques
- Lecture 13: Deep learning in object detection and recognition
 - also discussion about the topical issues, e.g. how to fool the system
- Szeliski chapt. 14 has some material, but is quite out-dated



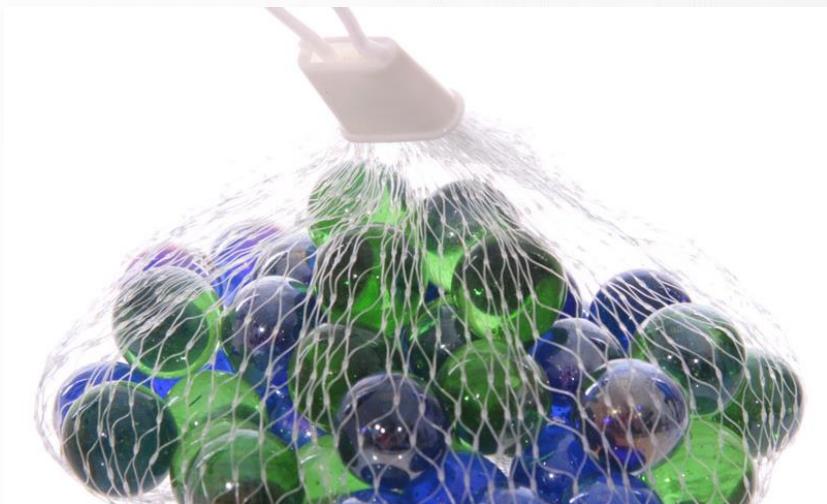
OBJECT DETECTION AND RECOGNITION

- **Object detection:** where is *this* object in the image?
 - *input:* a clear image of an object, or some kind of model of an object (e.g. duck) and an image (possibly) containing the object of interest
 - *output:* position, or a bounding box of the input object if it exists in the image
- **Object Recognition:** which object is depicted in the image?
 - *input:* an image containing unknown object(s)
 - *output:* position(s) and label(s) (names) of the objects in the image
 - The positions of objects are either acquired from the input, or determined based on the input image.
 - When labeling objects, there is usually a set of categories/labels which the system "knows" and between which the system can differentiate

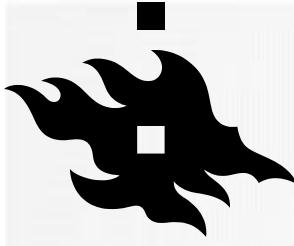


OBJECT DETECTION NOW

- Object detection is one of the most challenging computer vision tasks
- Aims at identifying the presence of various individual objects in an image
- Good results have been obtained when dealing with images with relatively simple image scenes and clear foreground objects
- The problem is not adequately addressed when dealing with the images and videos containing objects placed in arbitrary poses, with various shapes, and appearing in a cluttered and occluded environment

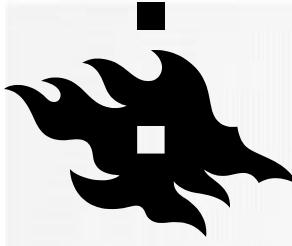


Han et al. (2018). Deep Learning for Visual Understanding: Part 2: Advanced Deep-Learning Techniques for Salient and Category-Specific Object Detection. IEEE Signal Processing Magazine.



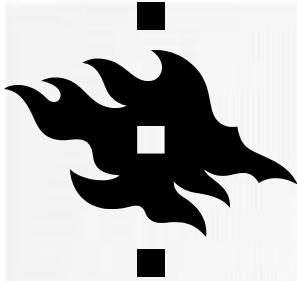
OBJECT DETECTION CATEGORIES

- Object detection can be divided into three directions: Objectness Detection (OD), Salient Object Detection (SOD), Category-specific Object Detection (COD)
- OD aims at detecting all possible objects appearing in each given image, does not mind about the object category
 - Quantifies how likely it is for an image window to contain an object of any class, as opposed to backgrounds
 - Challenge: different objects have large appearance variation
 - Outputs thousands of object proposals
- SOD highlights objects that draw attention
 - Mimics human's visual system
 - Challenge: how to associate desired visual stimulus with corresponding image regions

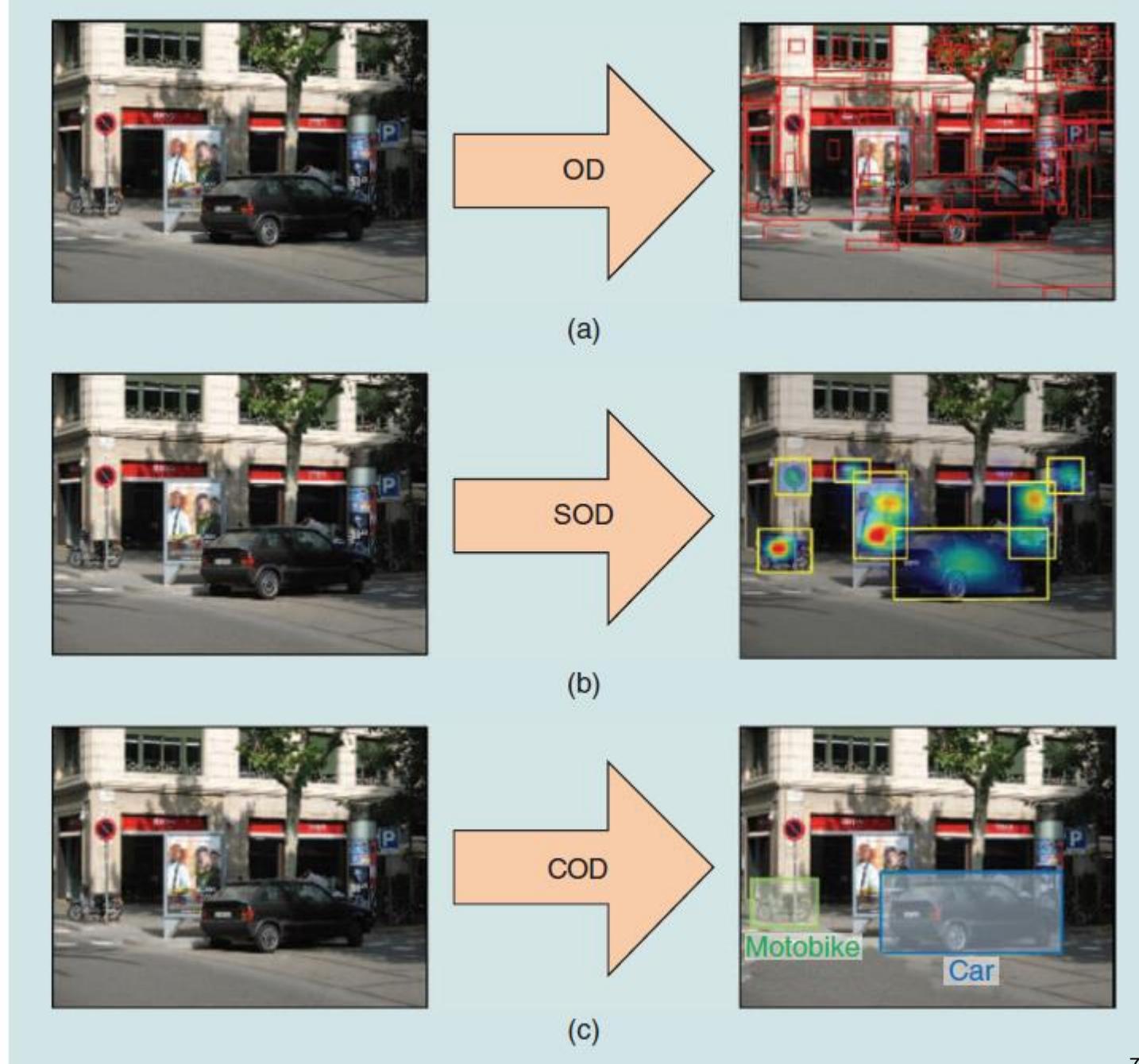


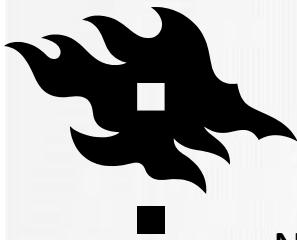
OBJECT DETECTION CATEGORIES

- COD detects predefined object categories
 - Challenge: in addition to recognizing the image areas where objects are and must recognize the object category
 - Solves a computational problem without trying to understand the visual attention
 - Multiclass classification problem
- SIFT and HOG (Histograms of Oriented Gradients) the most used features
- Object detection and recognition is always a machine learning task
 - the type and amount of training images is different for each
- Convolutional Neural Networks have been used since 2004, breakthrough 2014



Han et al. (2018). Deep Learning for Visual Understanding: Part 2: Advanced Deep-Learning Techniques for Salient and Category-Specific Object Detection. IEEE Signal Processing Magazine.

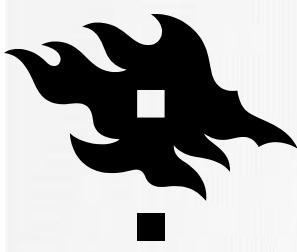




What do we mean by “object recognition”?

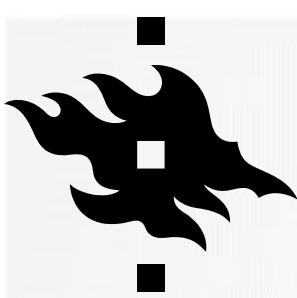
Next slides adapted from
Li, Fergus, & Torralba's excellent
[short course](#) on category and
object recognition





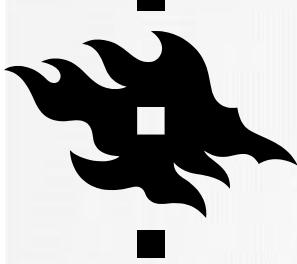
Verification: is that a lamp?





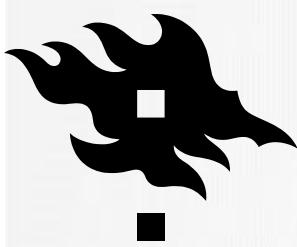
Detection: where are the people?





Identification: is that Potala Palace?

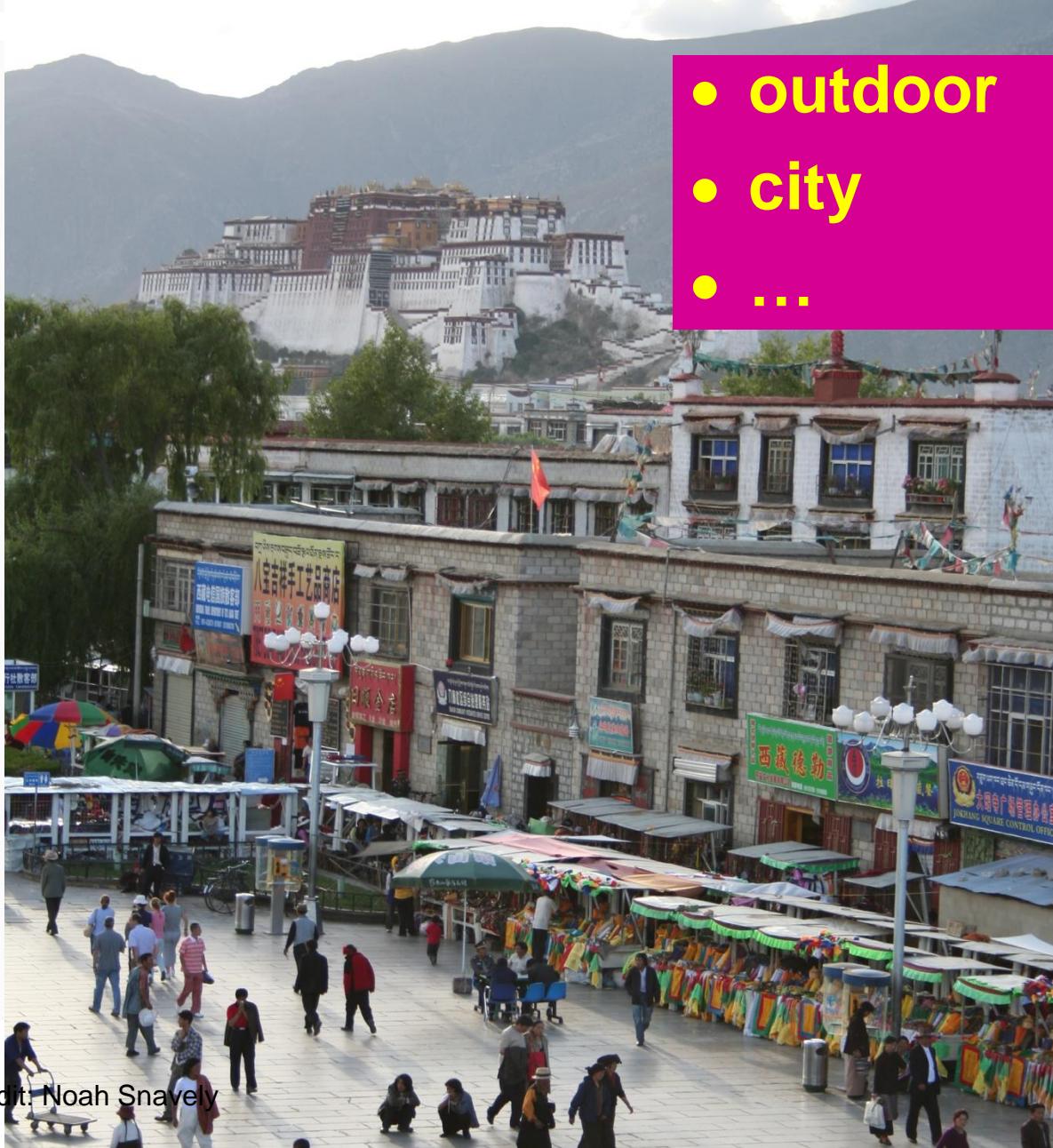
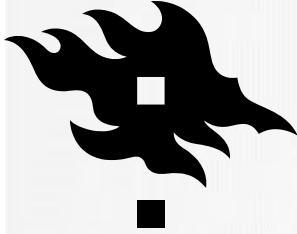




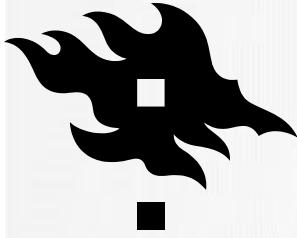
Object categorization

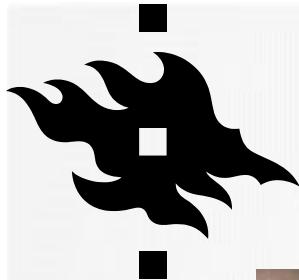


Scene and context categorization

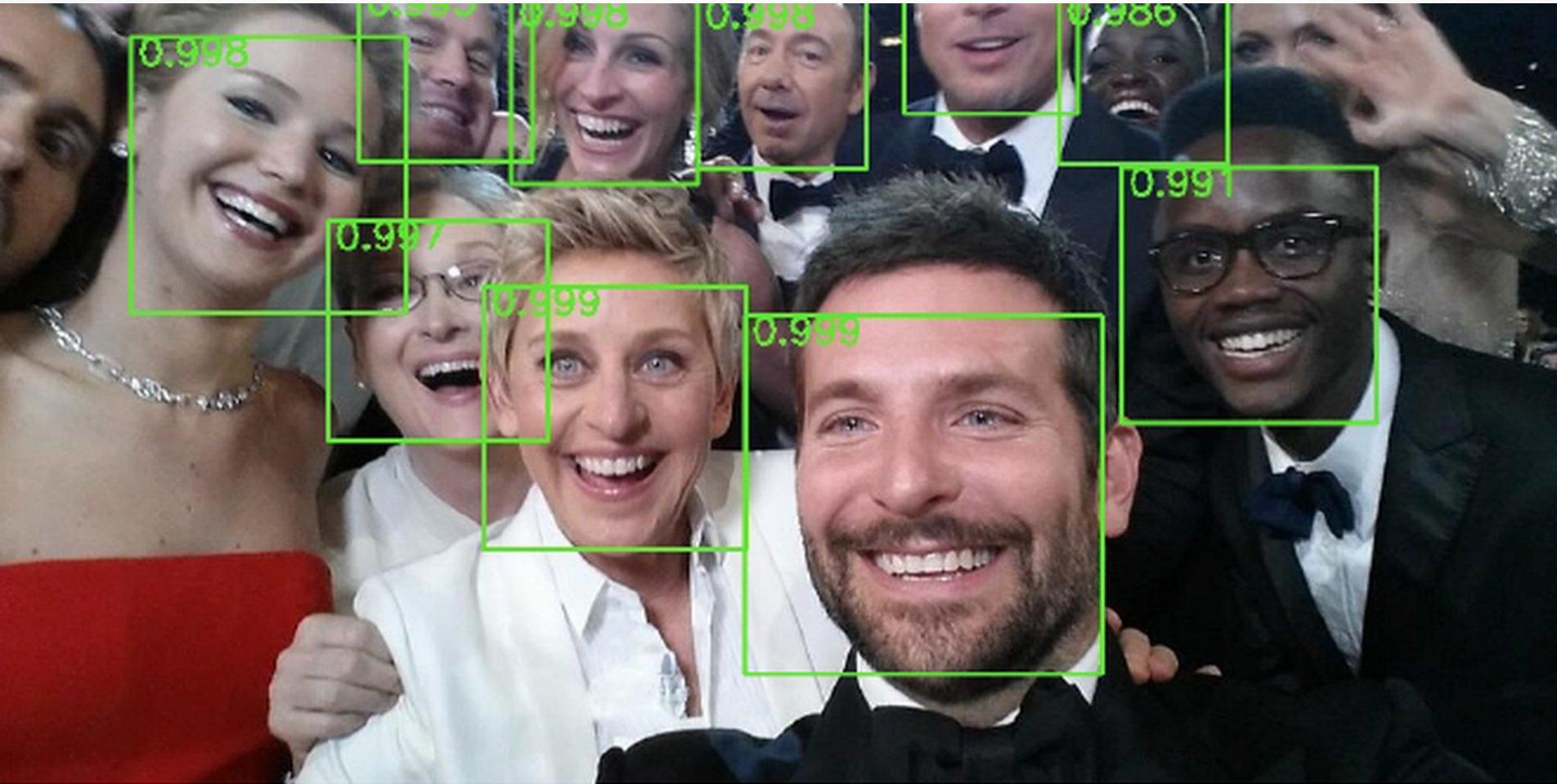


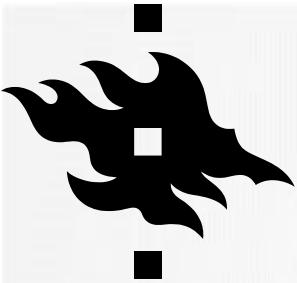
Activity / Event Recognition



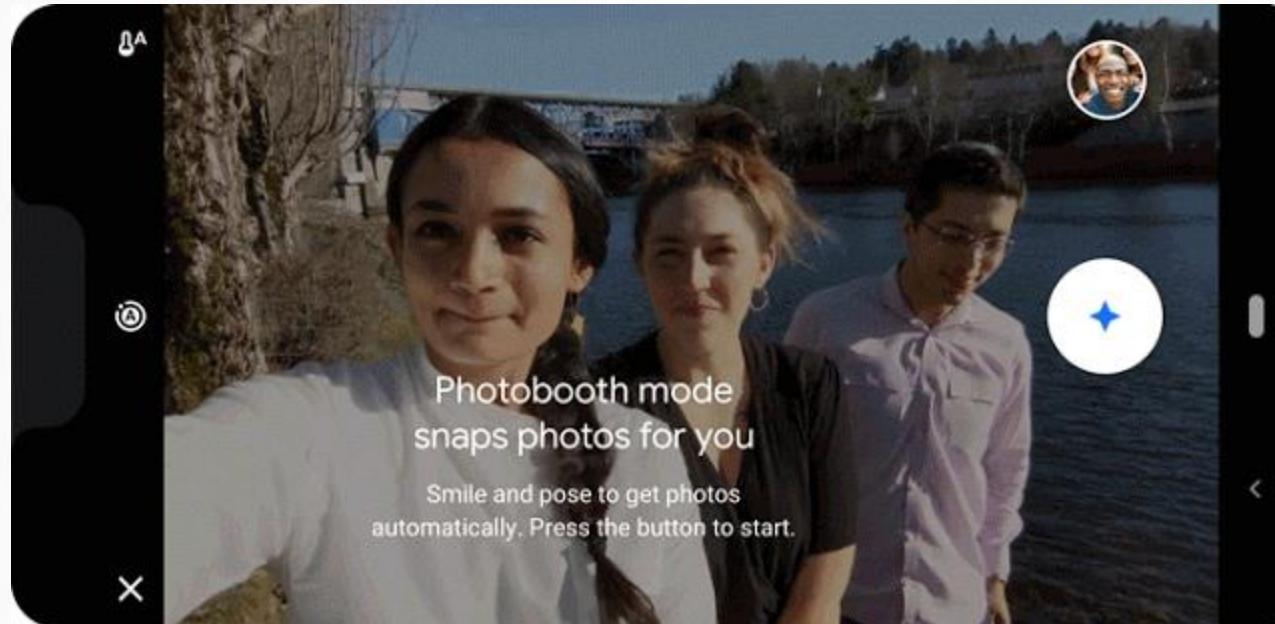


Applications: Photography



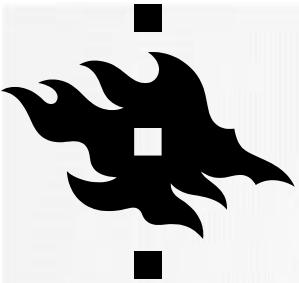


Applications: Shutter-free Photography



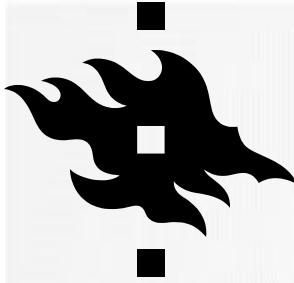
Take Your Best Selfie Automatically, with Photobooth on Pixel 3

<https://ai.googleblog.com/2019/04/take-your-best-selfie-automatically.html>



Applications: Assisted / autonomous driving

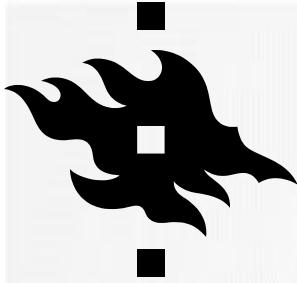




OLD SCHOOL OBJECT DETECTION

- Image segmentation
 - Segmentation and measuring the properties
 - Not for complicated images
- Template analysis
 - Take a small template image
 - Go through the image and look for a match by cross-correlation
 - Not robust for rotation, scaling, change in lighting





Object recognition Is it really so hard?

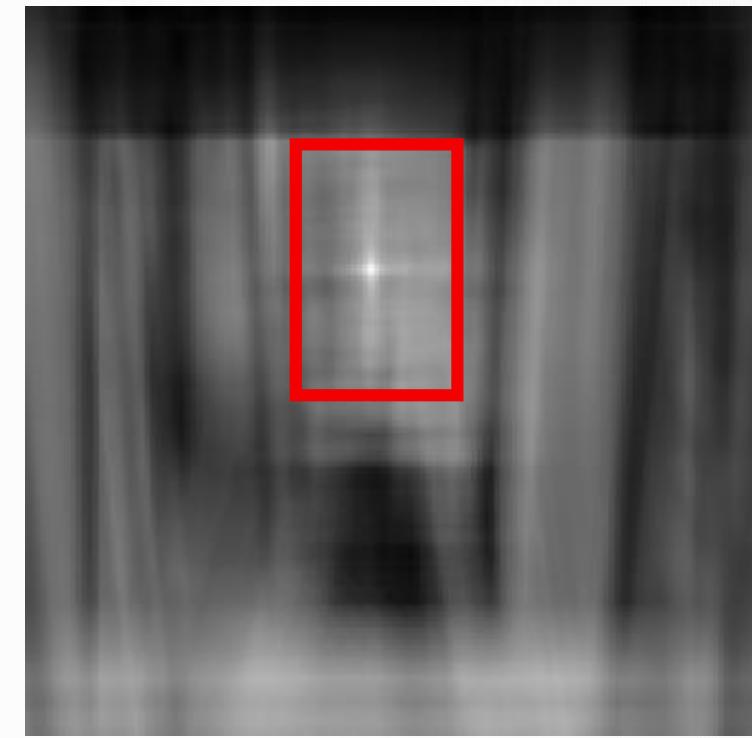
This is a chair

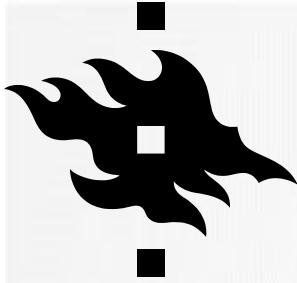


Find the chair in this image



Output of normalized correlation

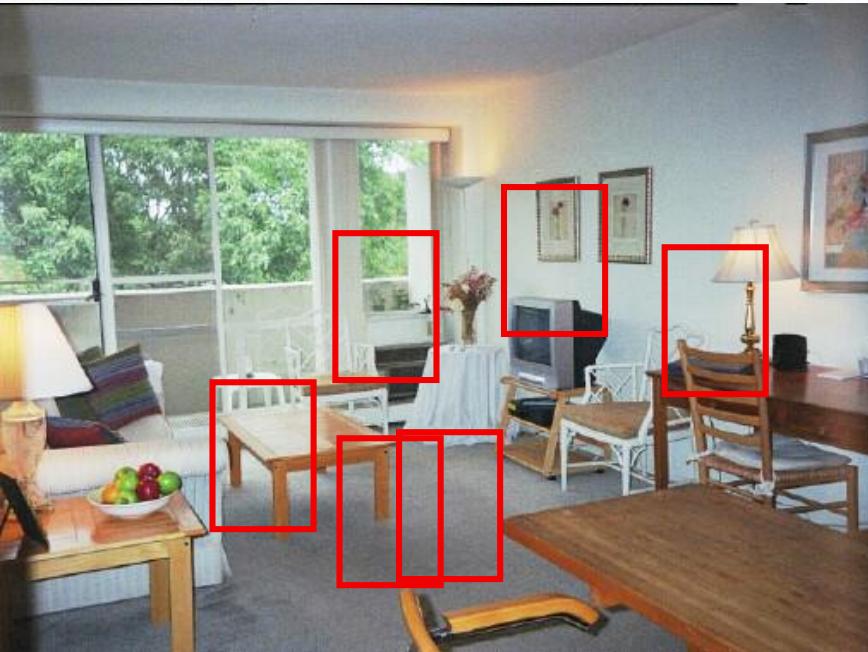




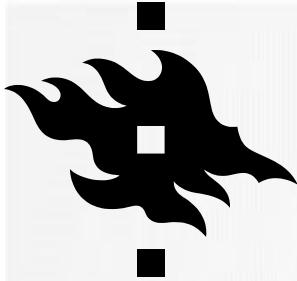
Object recognition Is it really so hard?



Find the chair in this image



Pretty much garbage
Simple template matching is not going to do the trick

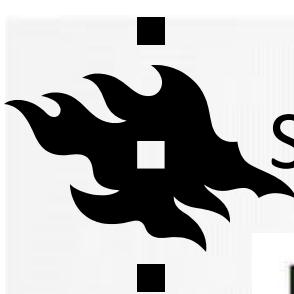


Object recognition Is it really so hard?

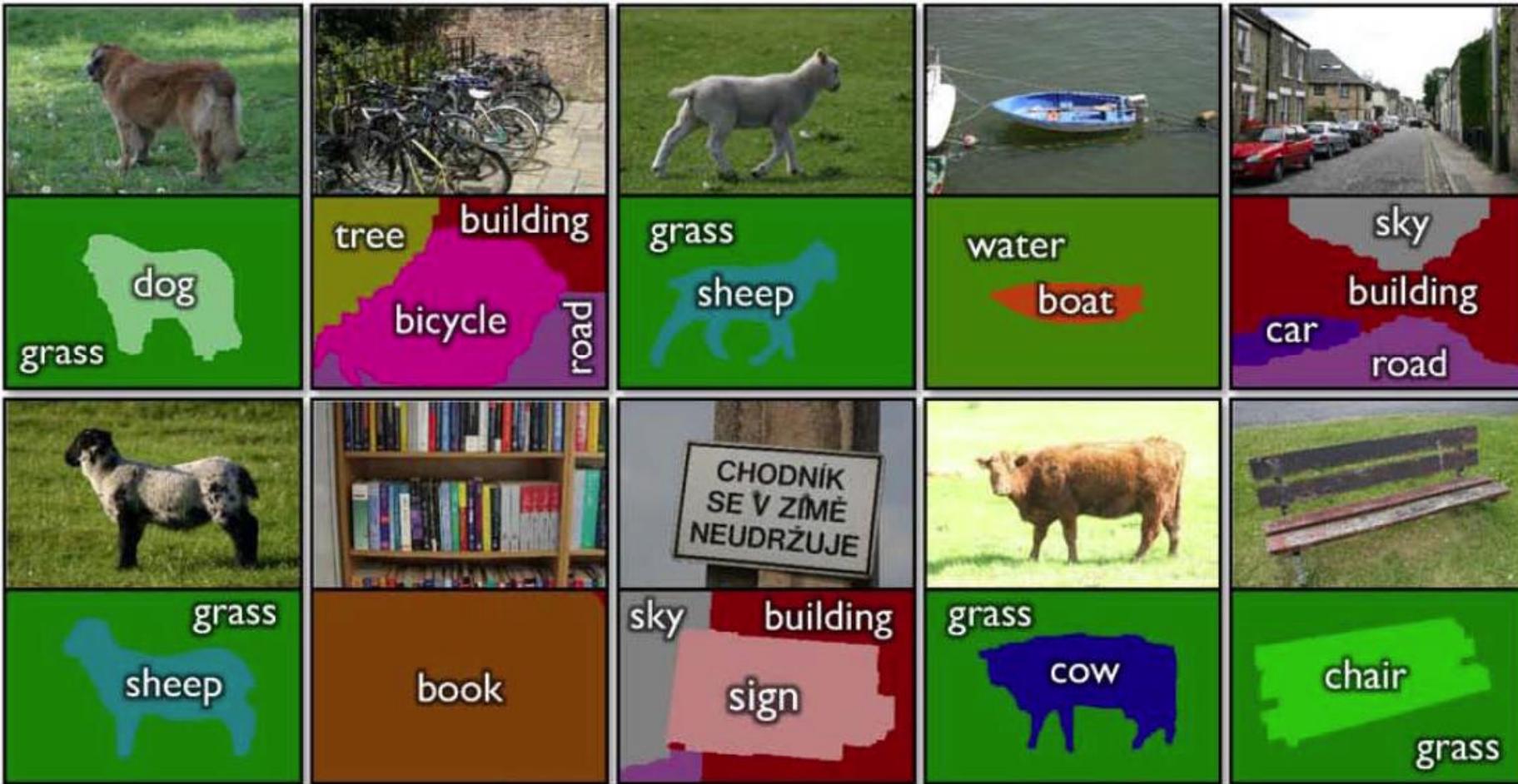
Find the chair in this image

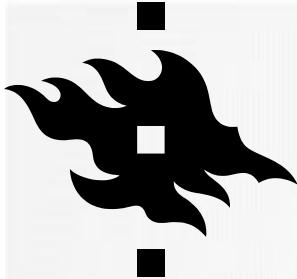


A “popular method is that of template matching, by point to point correlation of a model pattern with the image pattern. These techniques are inadequate for three-dimensional scene analysis for many reasons, such as occlusion, changes in viewing angle, and articulation of parts.” Nivatia & Binford, 1977.



Simultaneous recognition, detection, and segmentation



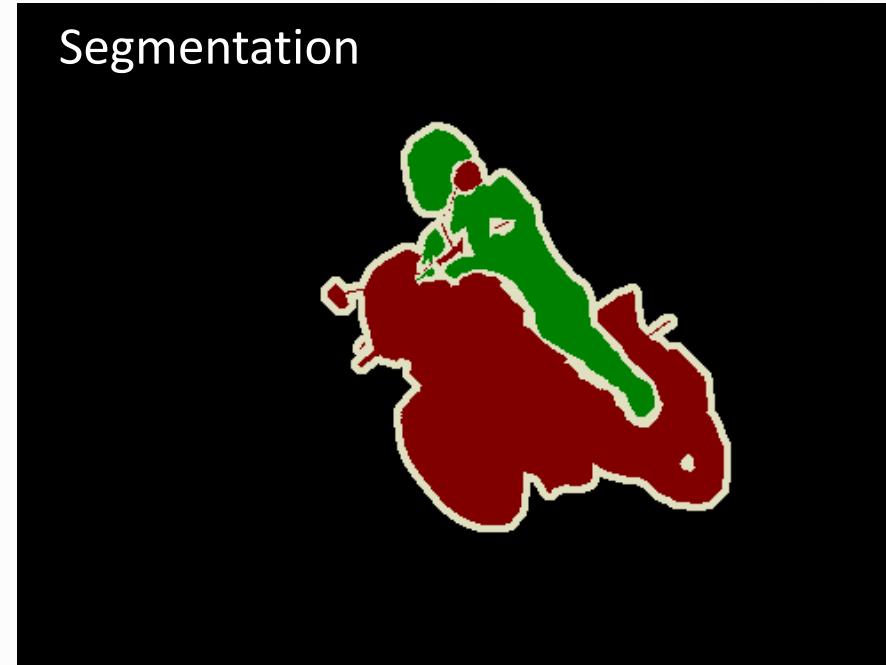
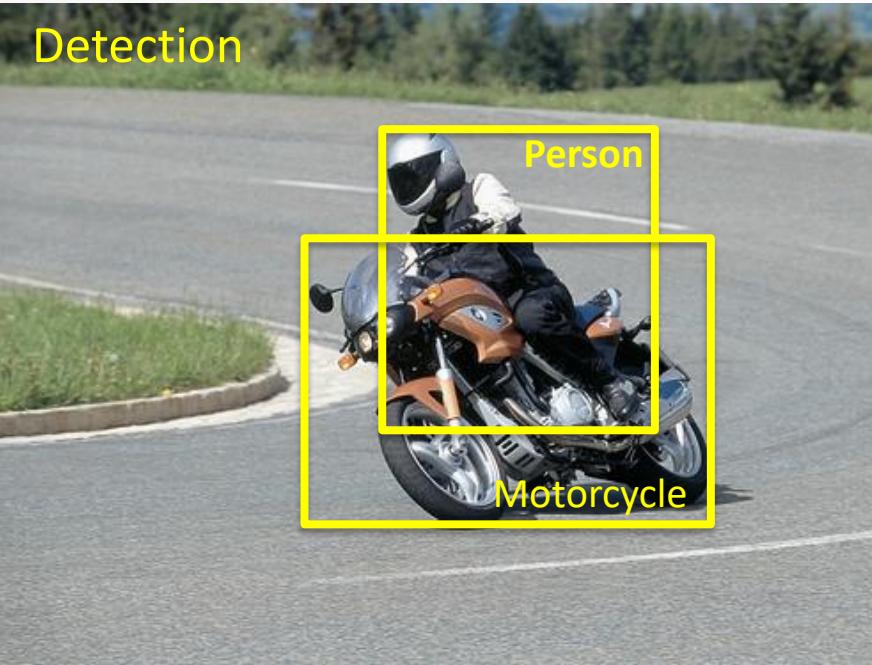


PASCAL VOC 2005-2012

20 object classes

22,591 images

Classification: person, motorcycle



Action: riding bicycle

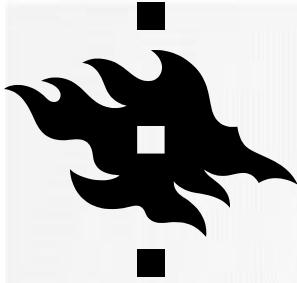


IMAGE CLASSIFICATION (IMAGE RECOGNITION)



(assume given set of discrete labels)
{dog, cat, truck, plane, ...}



cat

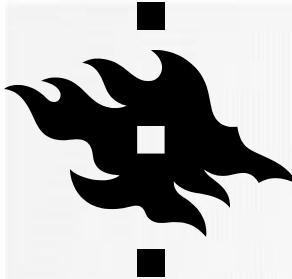


IMAGE CLASSIFICATION: PROBLEM



05	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	91	64
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	49	56	62	00	
81	49	31	73	55	79	14	29	93	71	40	67	55	58	30	03	49	13	36	65
92	70	95	23	04	60	11	42	62	51	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	63	39	41	92	36	54	22	40	40	28	66	33	13	80
24	47	19	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	69	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
03	46	68	87	57	62	20	72	03	46	33	67	46	55	12	52	63	93	53	69
04	42	16	73	58	35	39	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	34	65	99	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	40	88	81	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	71	48	

What the computer sees

image classification

82% cat
15% dog
2% hat
1% mug

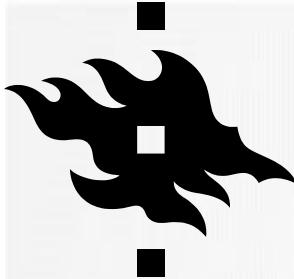
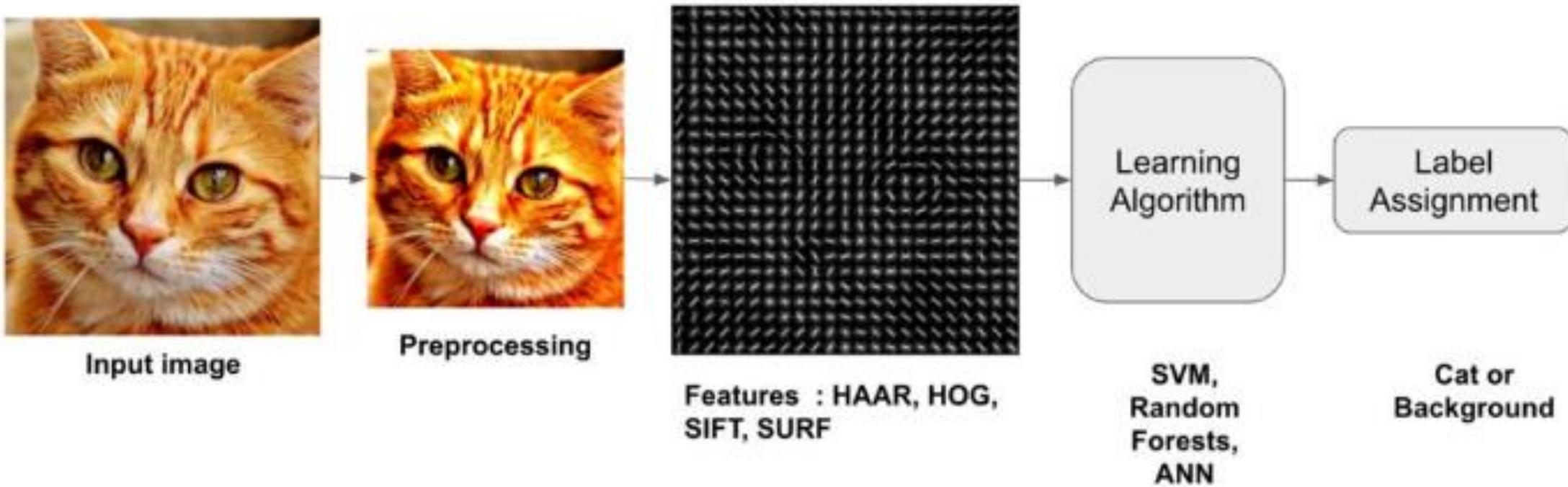
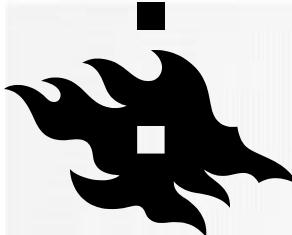


IMAGE CLASSIFIER



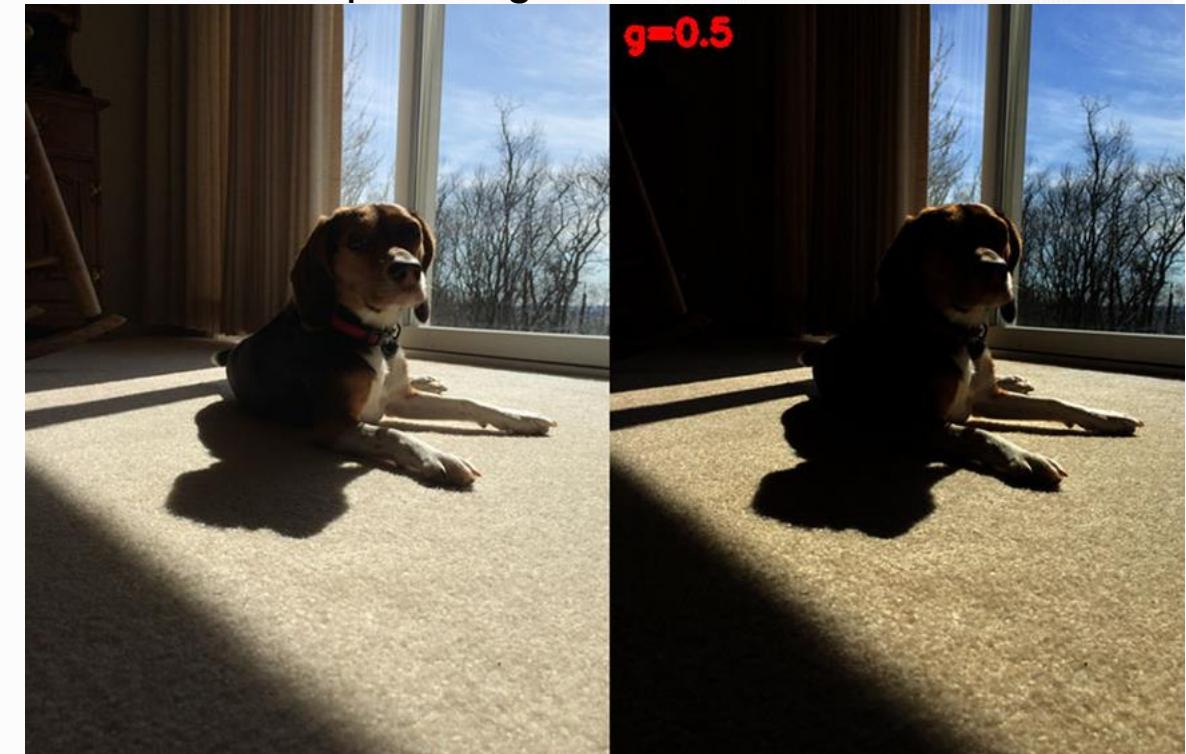
<https://www.learnopencv.com/image-recognition-and-object-detection-part1/>

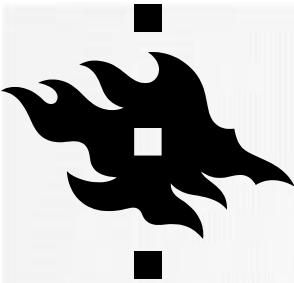


PREPROCESSING THE IMAGE

- Image is pre-processed to normalize contrast and brightness
- The suitable filetring method depends on the situation
- Image cropped and resized to a fixed size
- Feature detection
 - edge detection is used sometimes, but e.g. SIFT (and its variants), give much better performance
 - HOG is a regional descriptor (not local as SIFT) and most used in object detection

OpenCV gamma correction

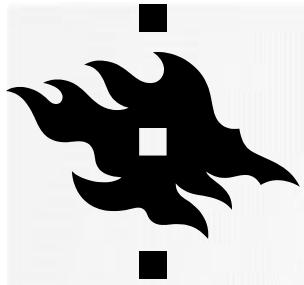




Why not use SIFT matching for everything?

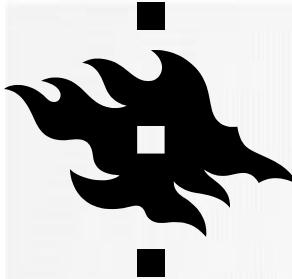
- Works well for object *instances* (or distinctive images such as logos)
- Not great for generic objects





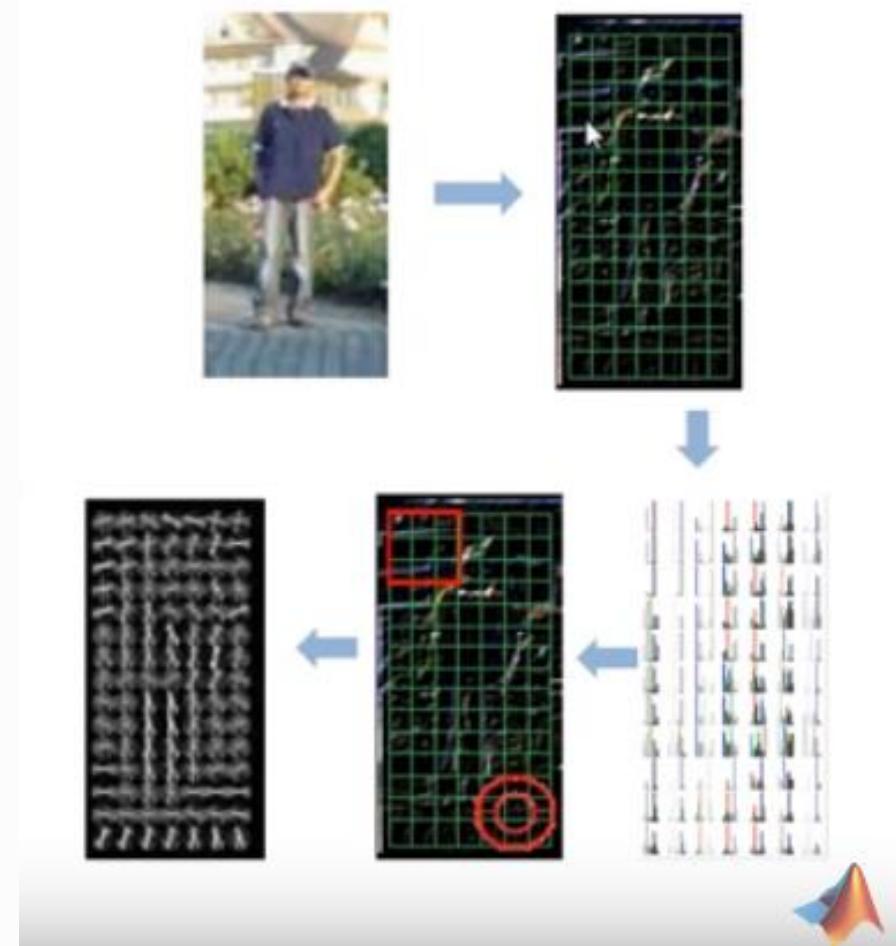
And it can get a lot harder

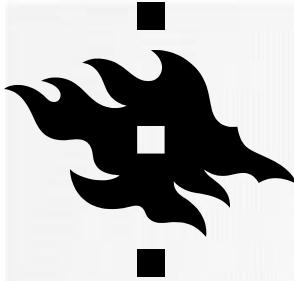




HISTOGRAM OF ORIENTED GRADIENTS (HOG)

- Compute gradients of the image, e.g. using $[-1, 0, 1]$ mask with no smoothing
- Divide into cells
- Compute a histogram of gradient orientations on each cell
- Group cells into overlapping blocks, normalize vector of histogram values
- Compute that for e.g. 1000 images and use those for training a classifier





DATA-DRIVEN APPROACH

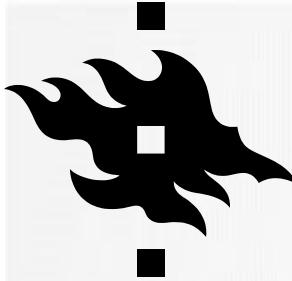
Collect a database of images with labels

Use ML to train an image classifier

Evaluate the classifier on test images

Example training set



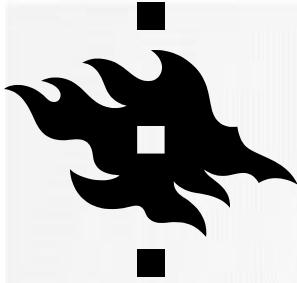


DATA-DRIVEN APPROACH

- Compute Bag-of-Visual-Words
- Support Vector Machine (SVM) was mainly used in traditional detection methods



<https://www.youtube.com/watch?v=nnGijZ9vok4>

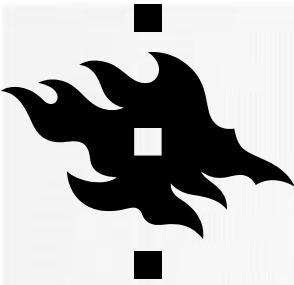


CalTech6 dataset



class	bag of features	bag of features	Parts-and-shape model
	Zhang et al. (2005)	Willamowski et al. (2004)	Fergus et al. (2003)
airplanes	98.8	97.1	90.2
cars (rear)	98.3	98.6	90.3
cars (side)	95.0	87.3	88.5
faces	100	99.3	96.4
motorbikes	98.5	98.0	92.5
spotted cats	97.0	—	90.0

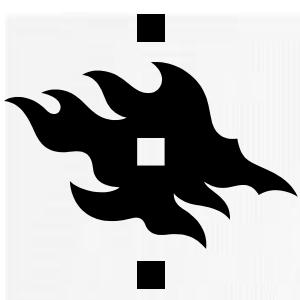
Works pretty well for image-level classification



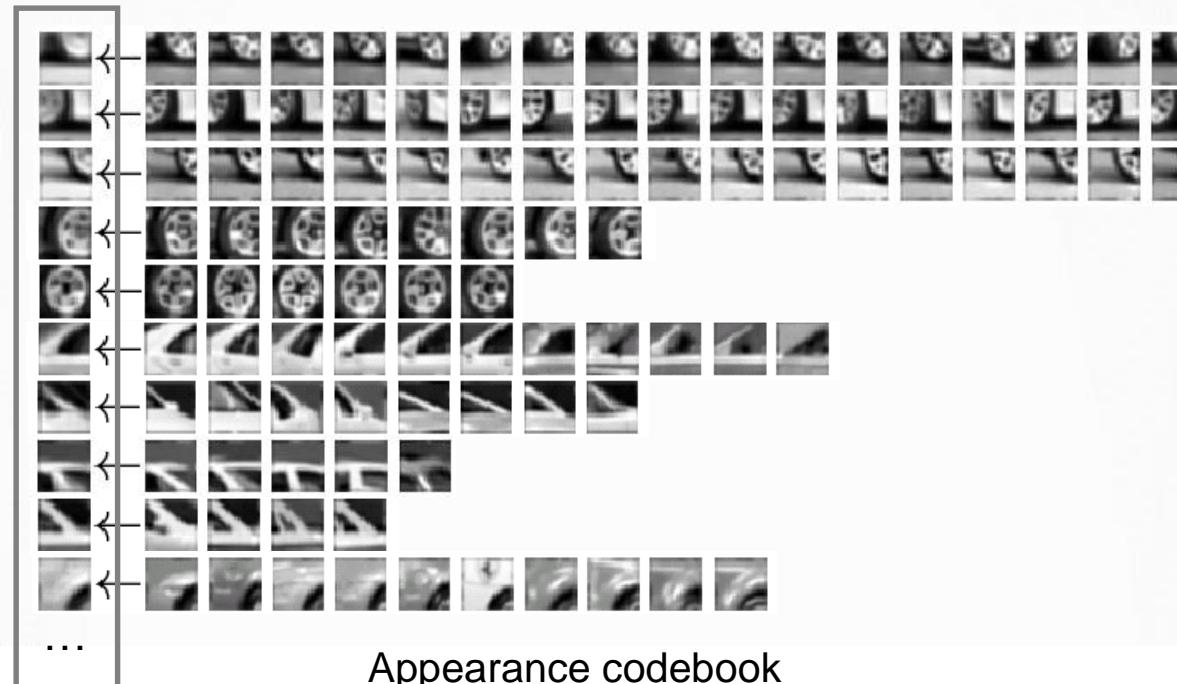
Dictionary Learning: Learn Visual Words using clustering

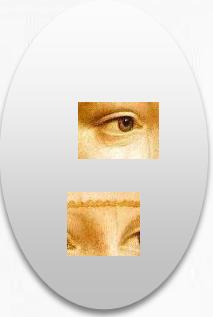
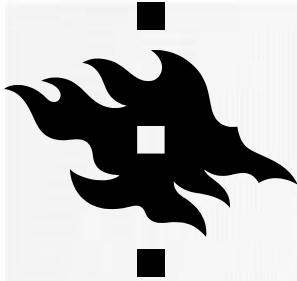
Encode:
build Bags-of-Words (BOW) vectors
for each image

Classify:
Train and test data using BOWs



EXAMPLE DICTIONARY

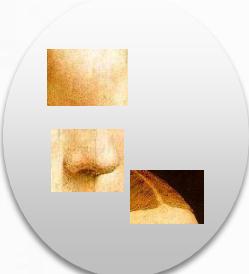


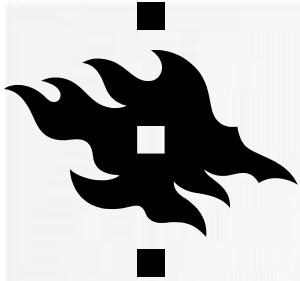


1. Quantization: image features gets associated to a visual word (nearest cluster center)



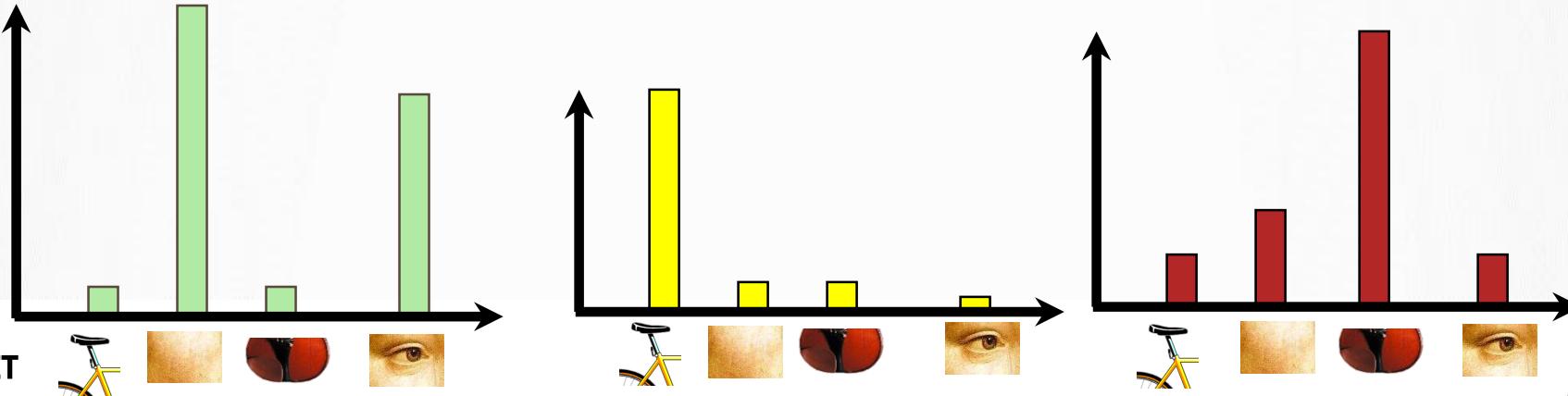
Encode:
build Bags-of-Words (BOW) vectors
for each image

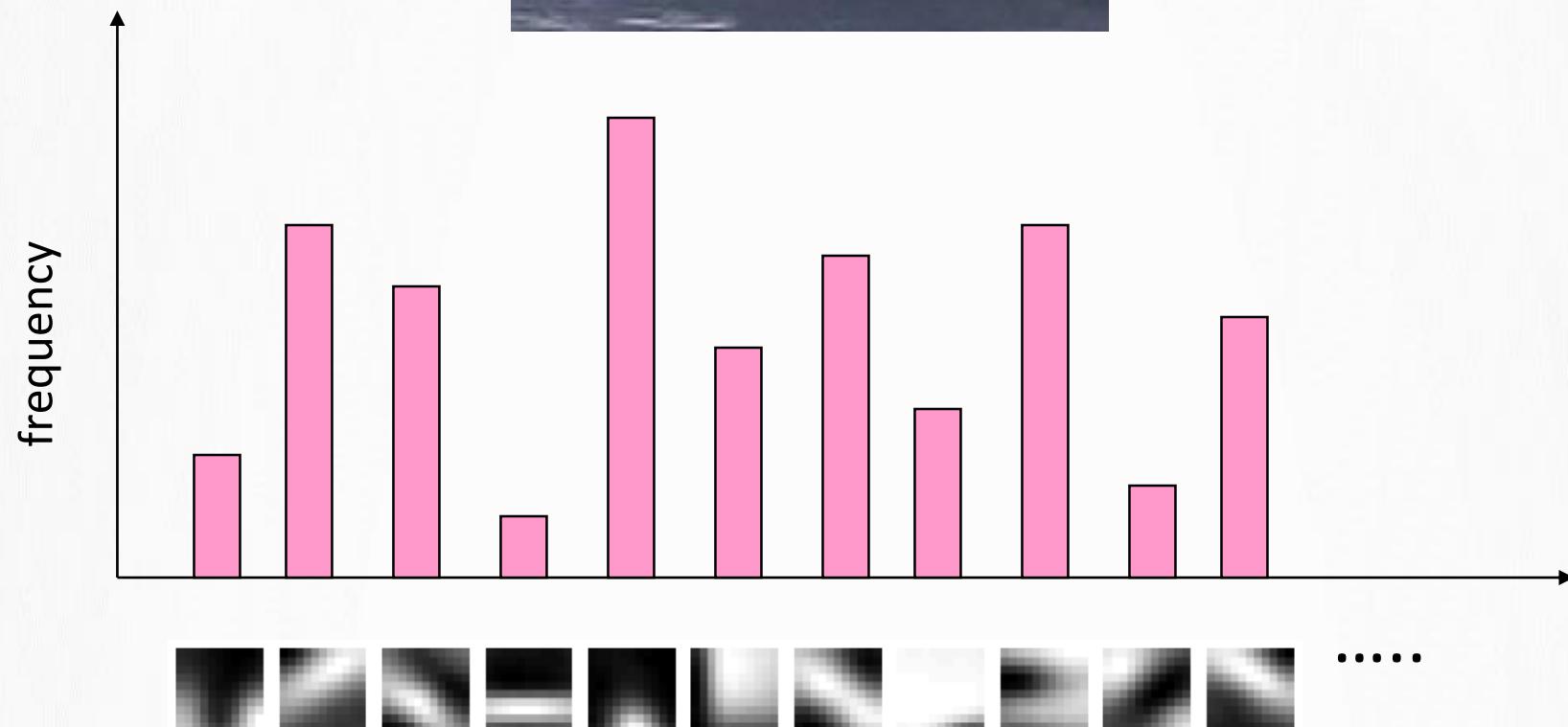
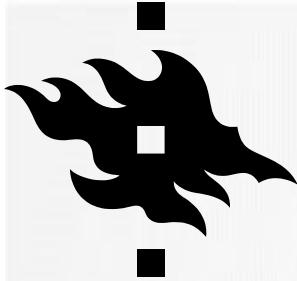


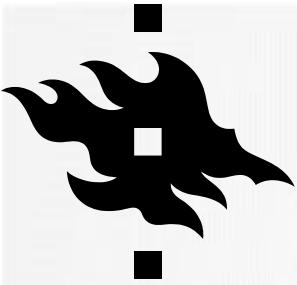


Encode:
build Bags-of-Words (BOW) vectors
for each image

2. Histogram: count the
number of visual word
occurrences



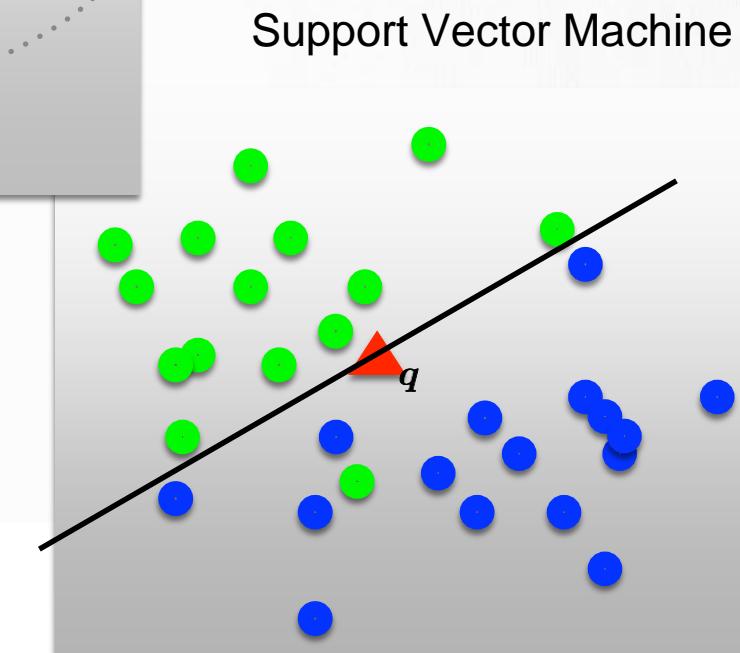
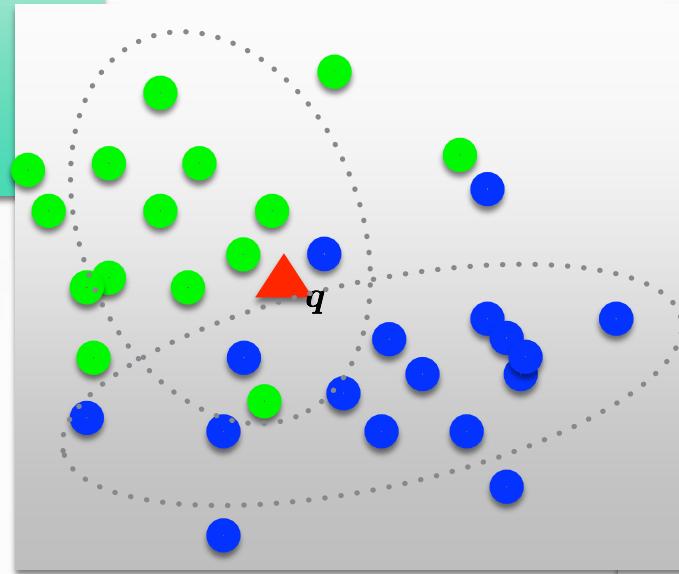
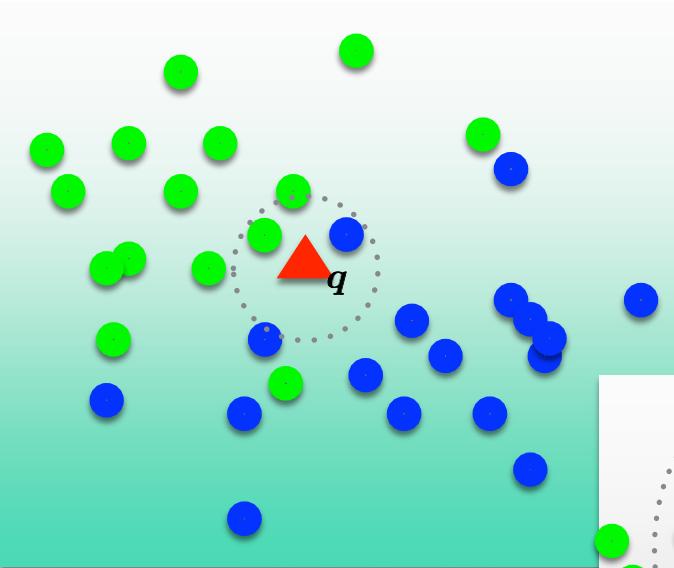


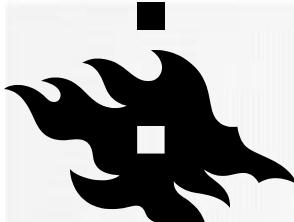


Dictionary Learning: Learn Visual Words using clustering

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

Classify:
Train and test data using BOWs





CIFAR-10 AND NN RESULTS

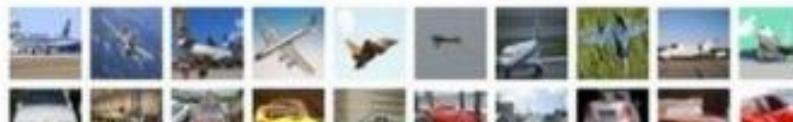
Example dataset: **CIFAR-10**

10 labels

50,000 training images

10,000 test images.

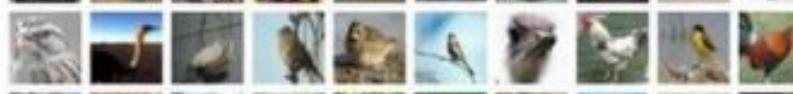
airplane



automobile



bird



cat



deer



dog



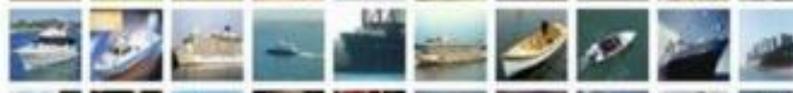
frog



horse



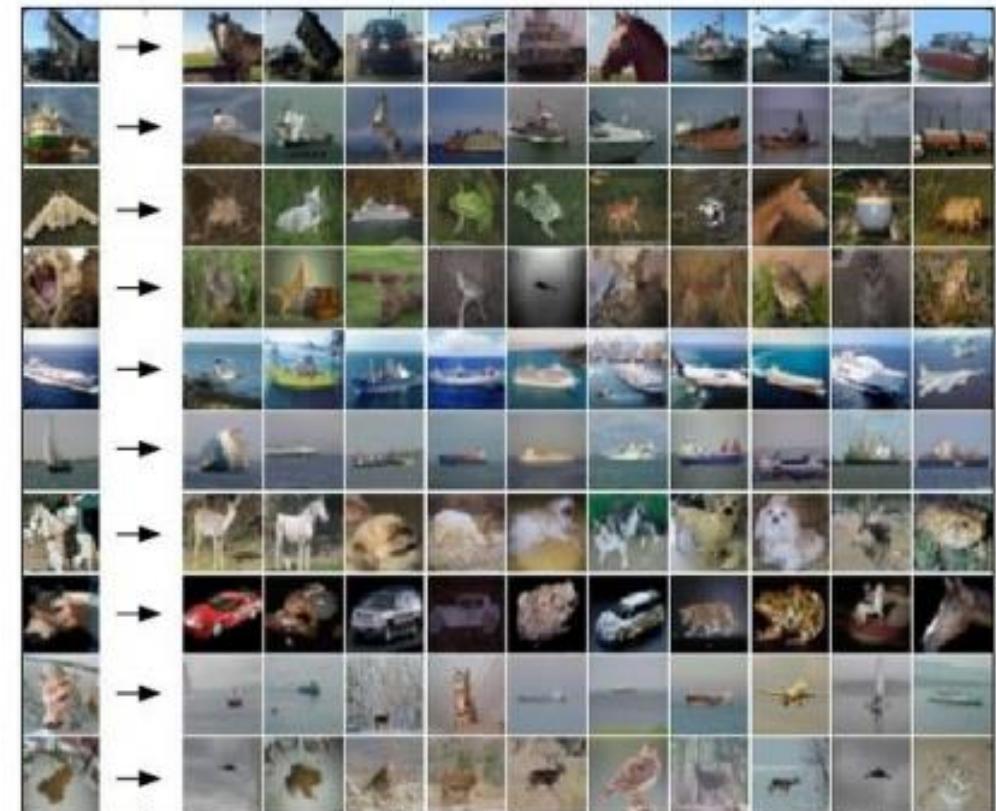
ship

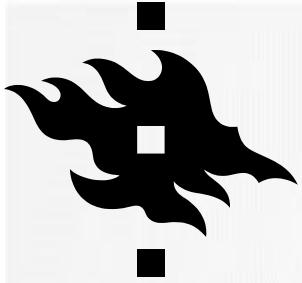


truck



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





SUPPORT VECTOR MACHINE

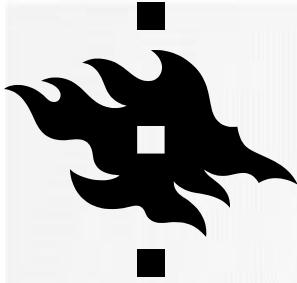


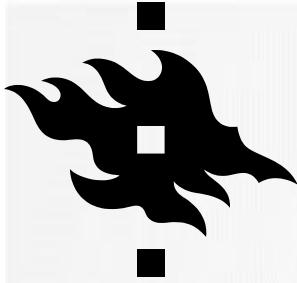
IMAGE CLASSIFICATION



(assume given set of discrete labels)
{dog, cat, truck, plane, ...}



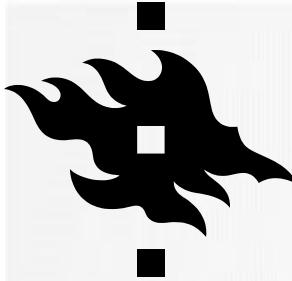
cat



SCORE FUNCTION



class scores



LINEAR CLASSIFIER

define a **score function**

$$f(x_i, W, b) = Wx_i + b$$

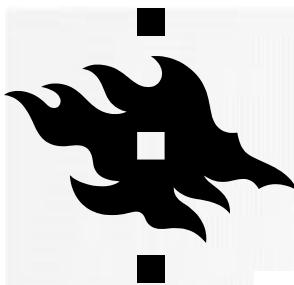
class scores

data (histogram)

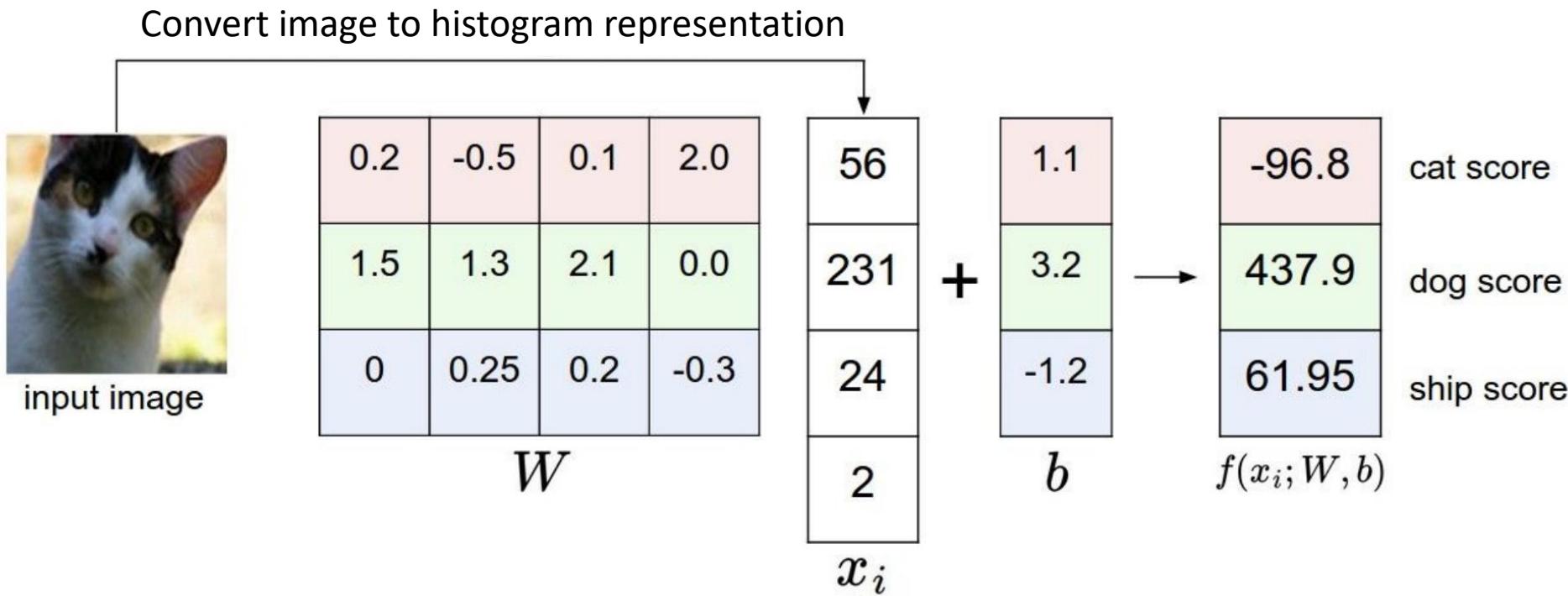
“weights”

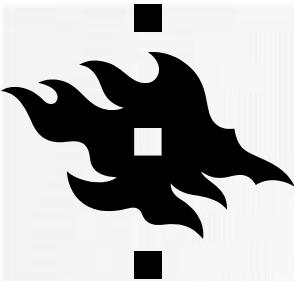
“bias vector”

“parameters”

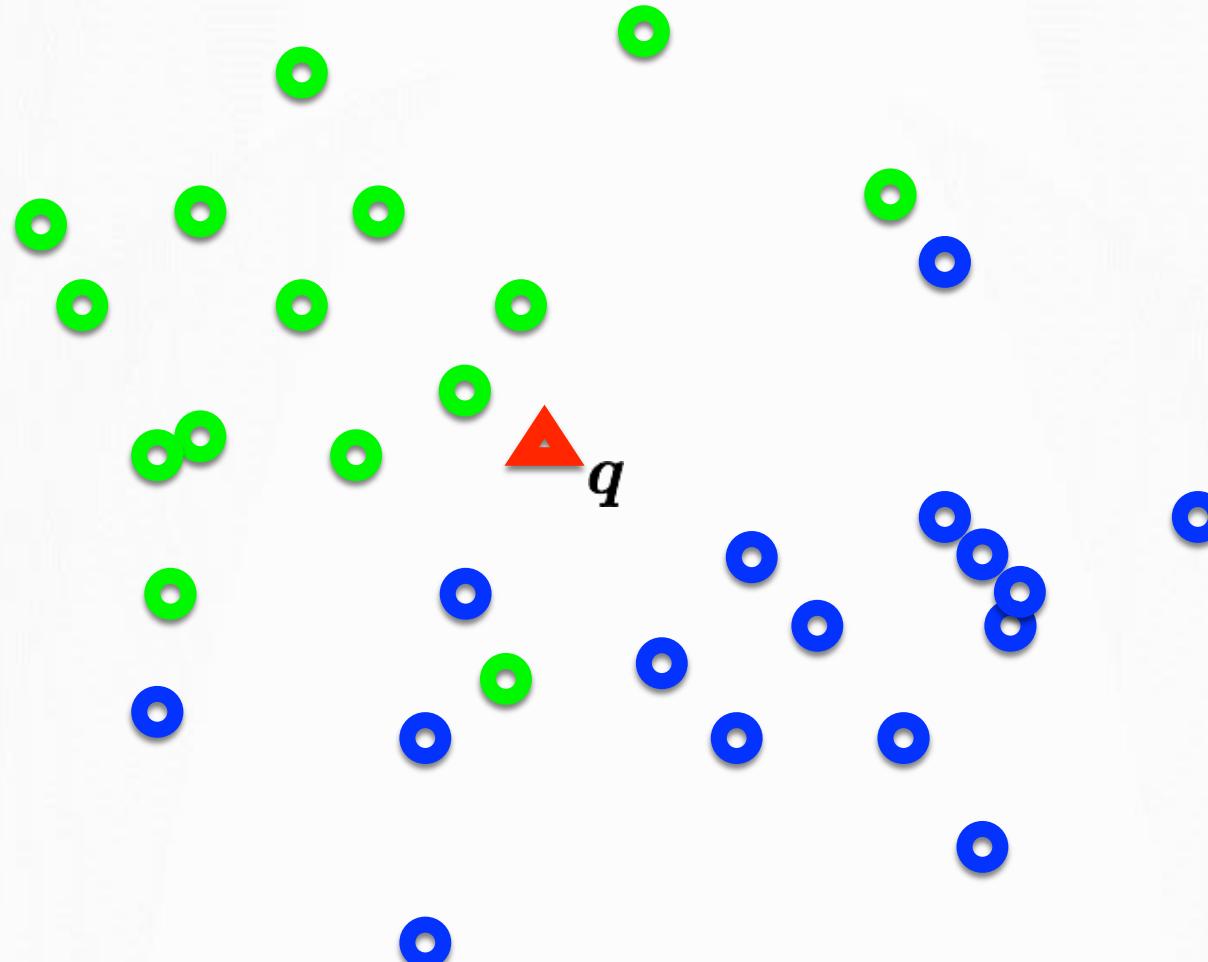


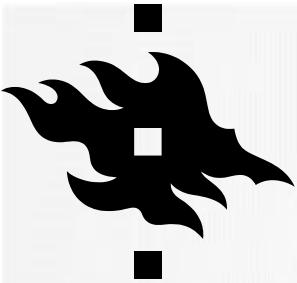
Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)



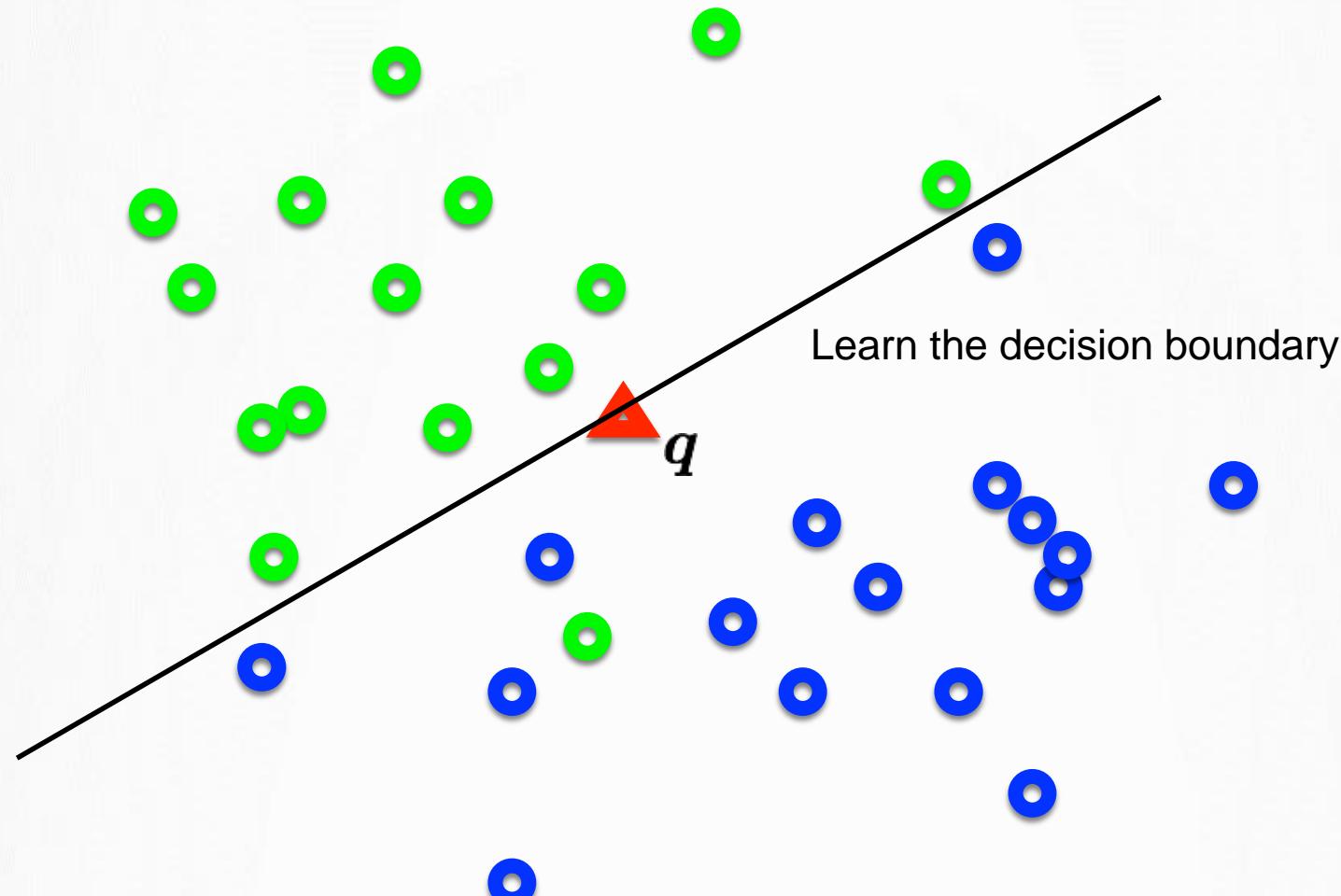


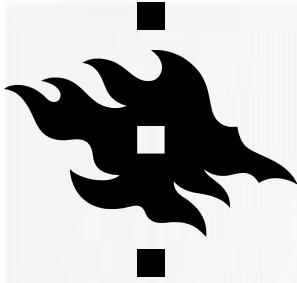
Distribution of data from two classes





Distribution of data from two classes

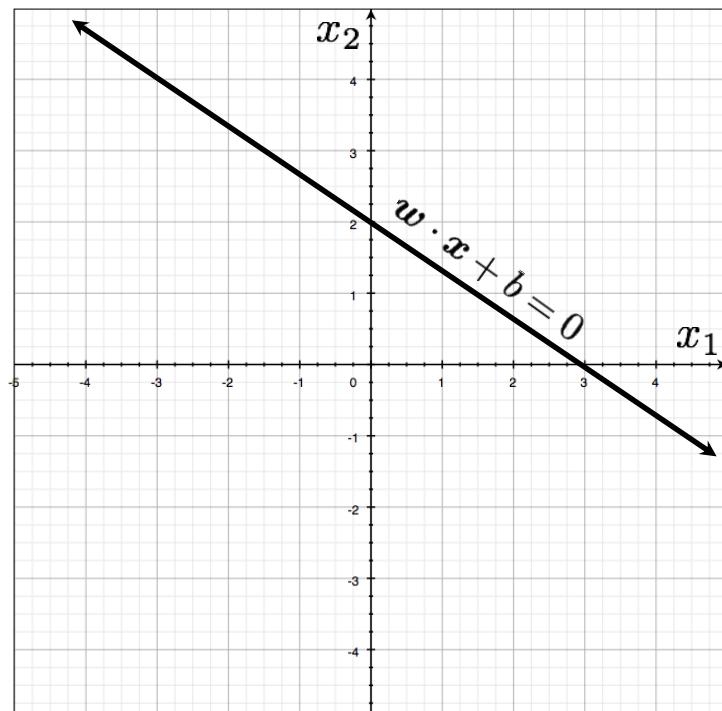




Hyperplanes (lines) in 2D

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

$$w_1x_1 + w_2x_2 + b = 0$$



Important property:
Free to choose any normalization of w

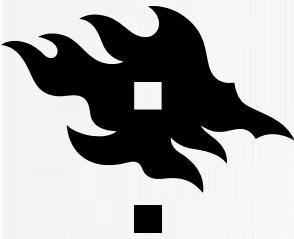
The line

$$w_1x_1 + w_2x_2 + b = 0$$

and the line

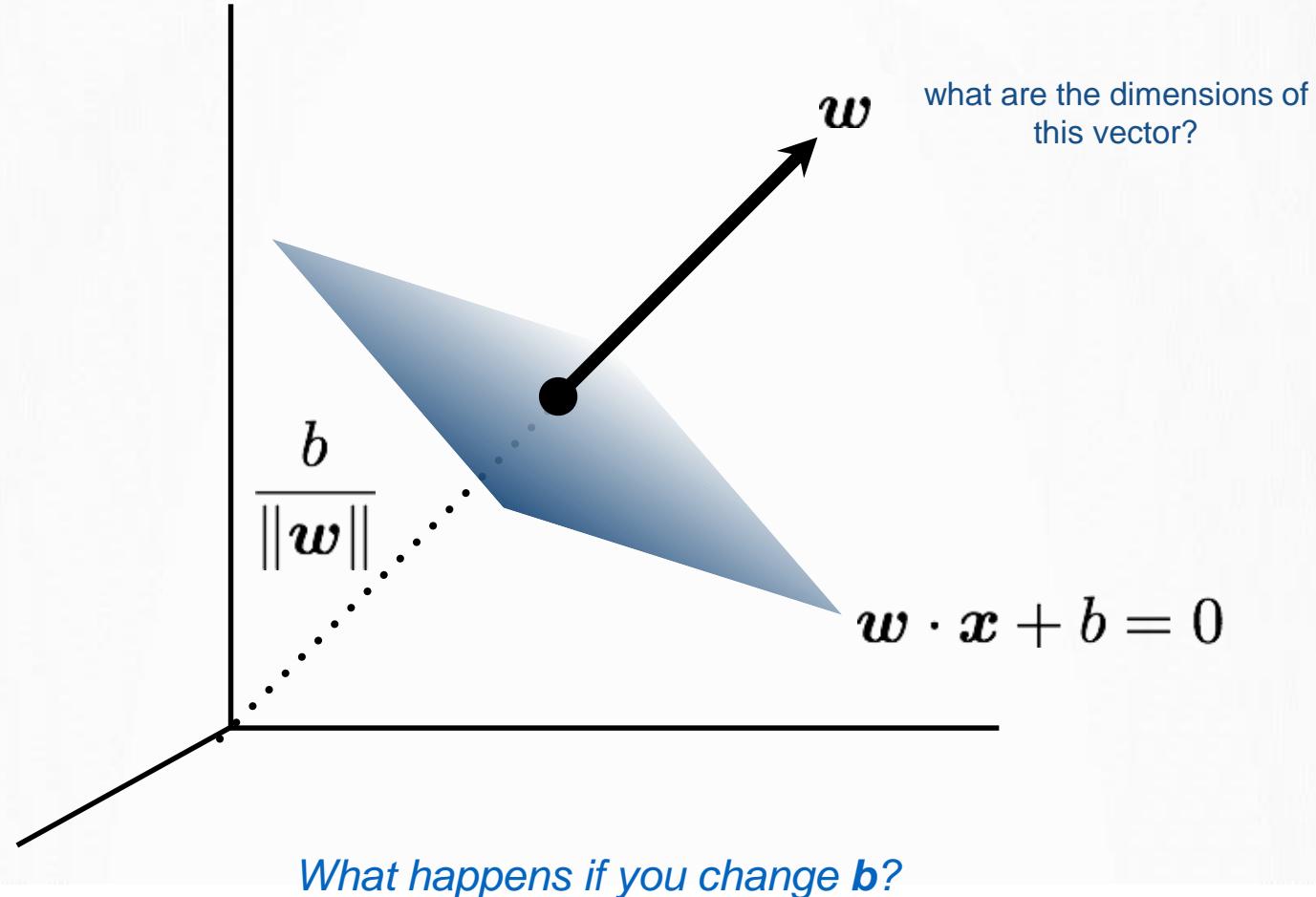
$$\lambda(w_1x_1 + w_2x_2 + b) = 0$$

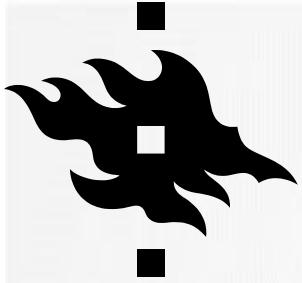
define the same line



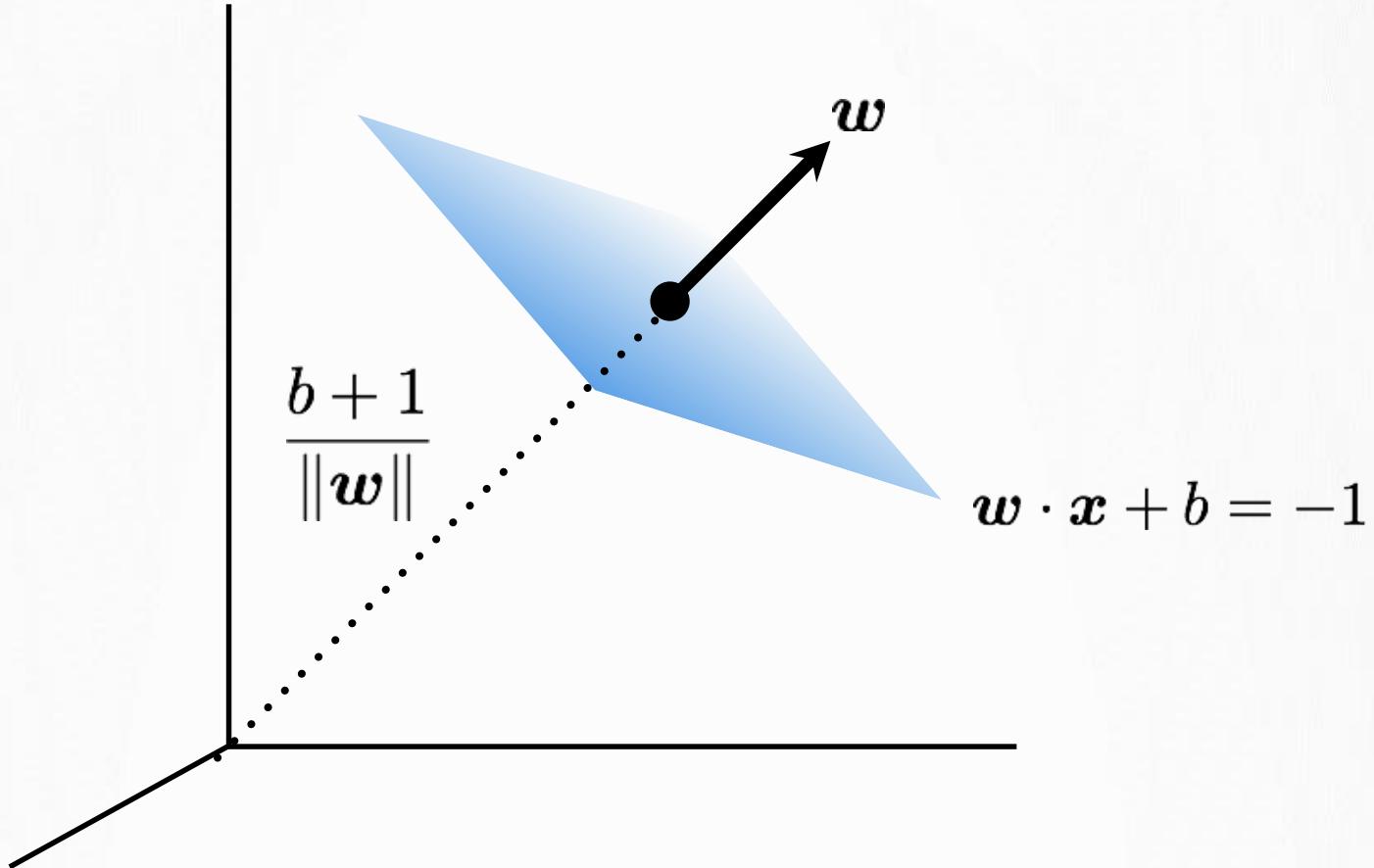
Now we can go to 3D ...

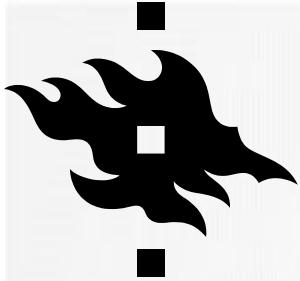
Hyperplanes (planes) in 3D





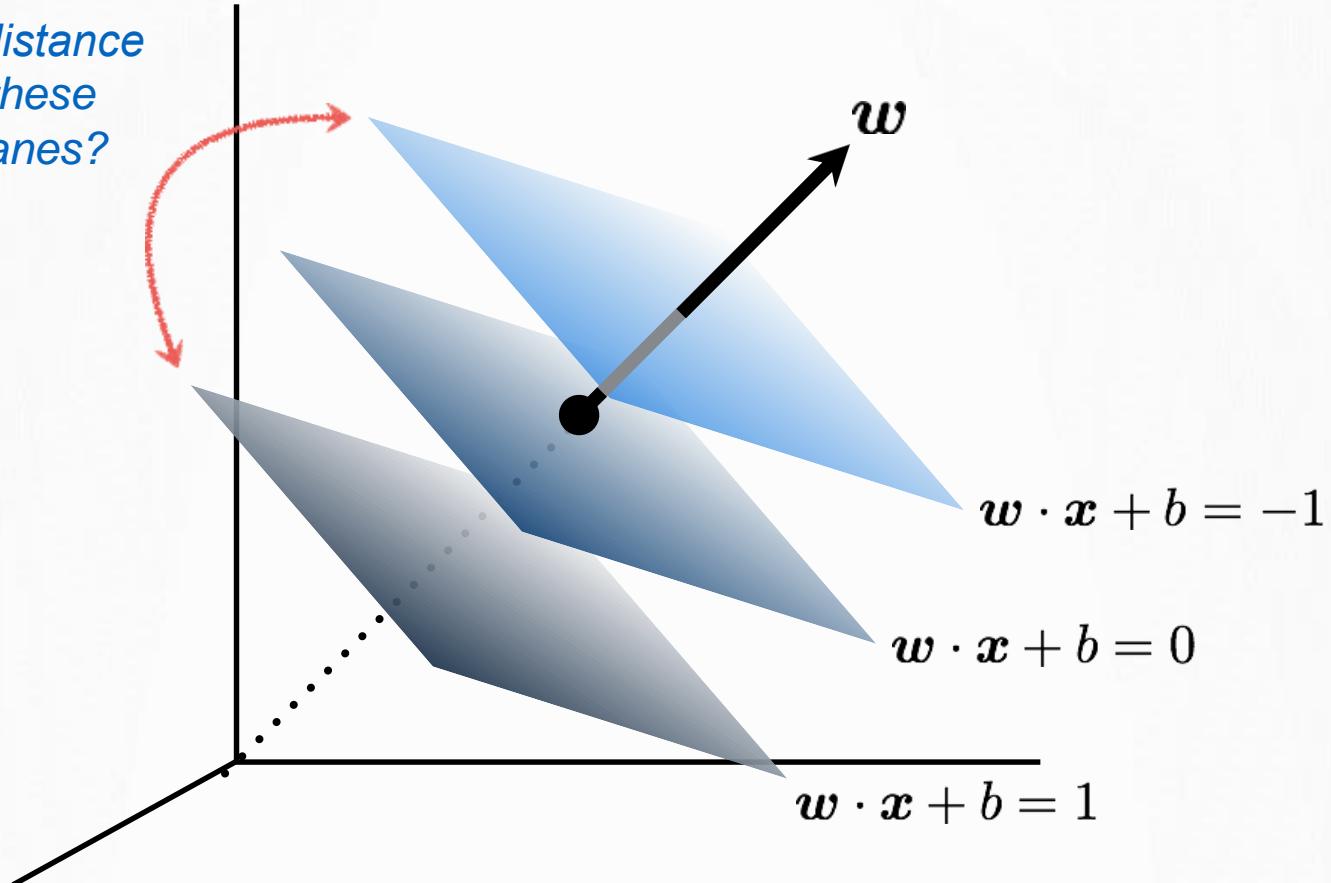
Hyperplanes (planes) in 3D

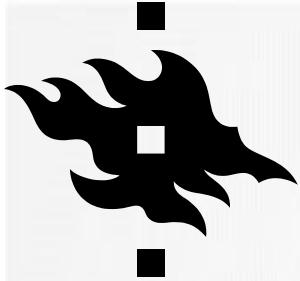




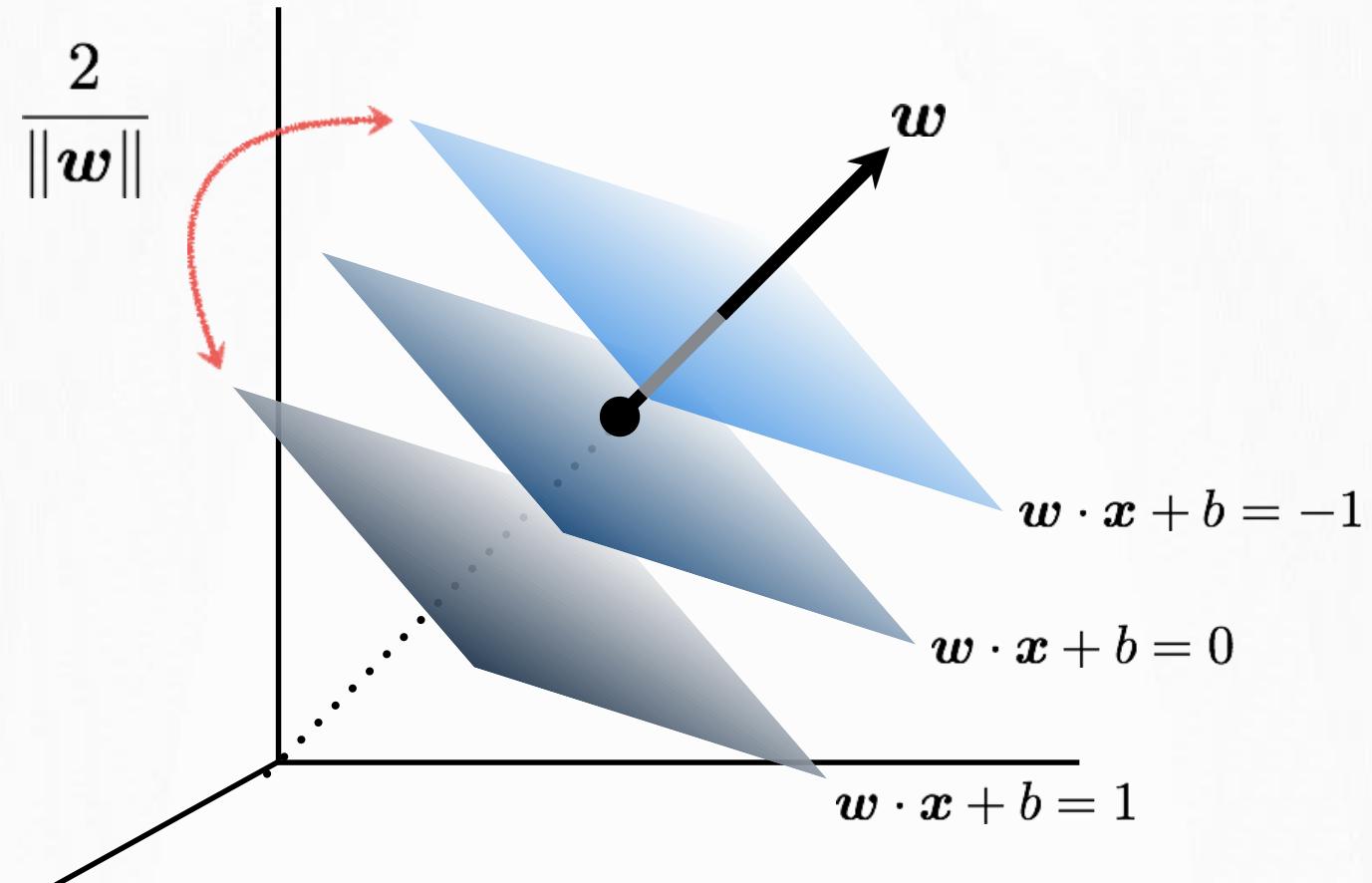
Hyperplanes (planes) in 3D

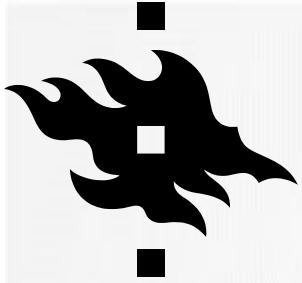
*What's the distance
between these
parallel planes?*



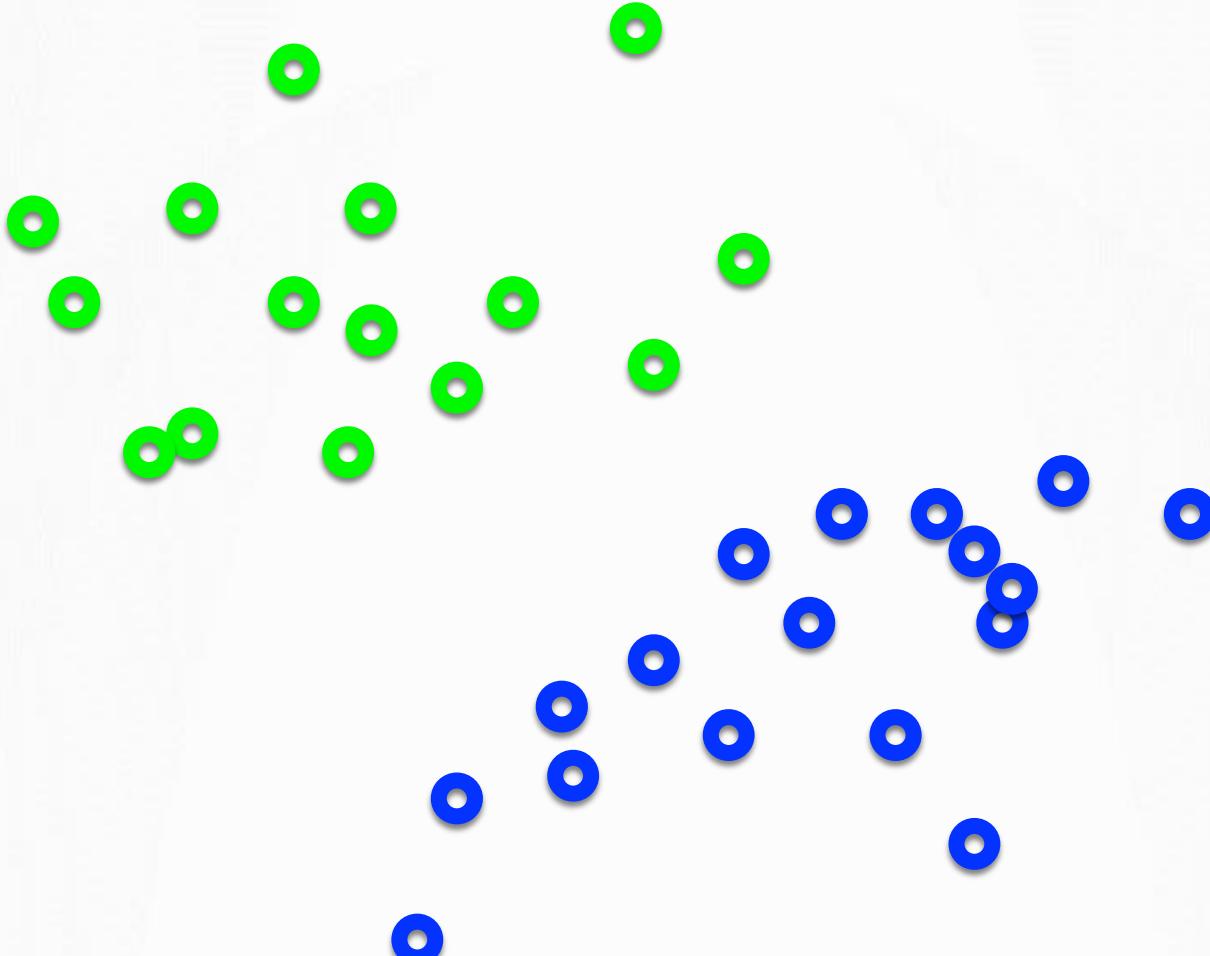


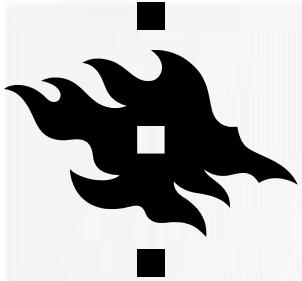
Hyperplanes (planes) in 3D



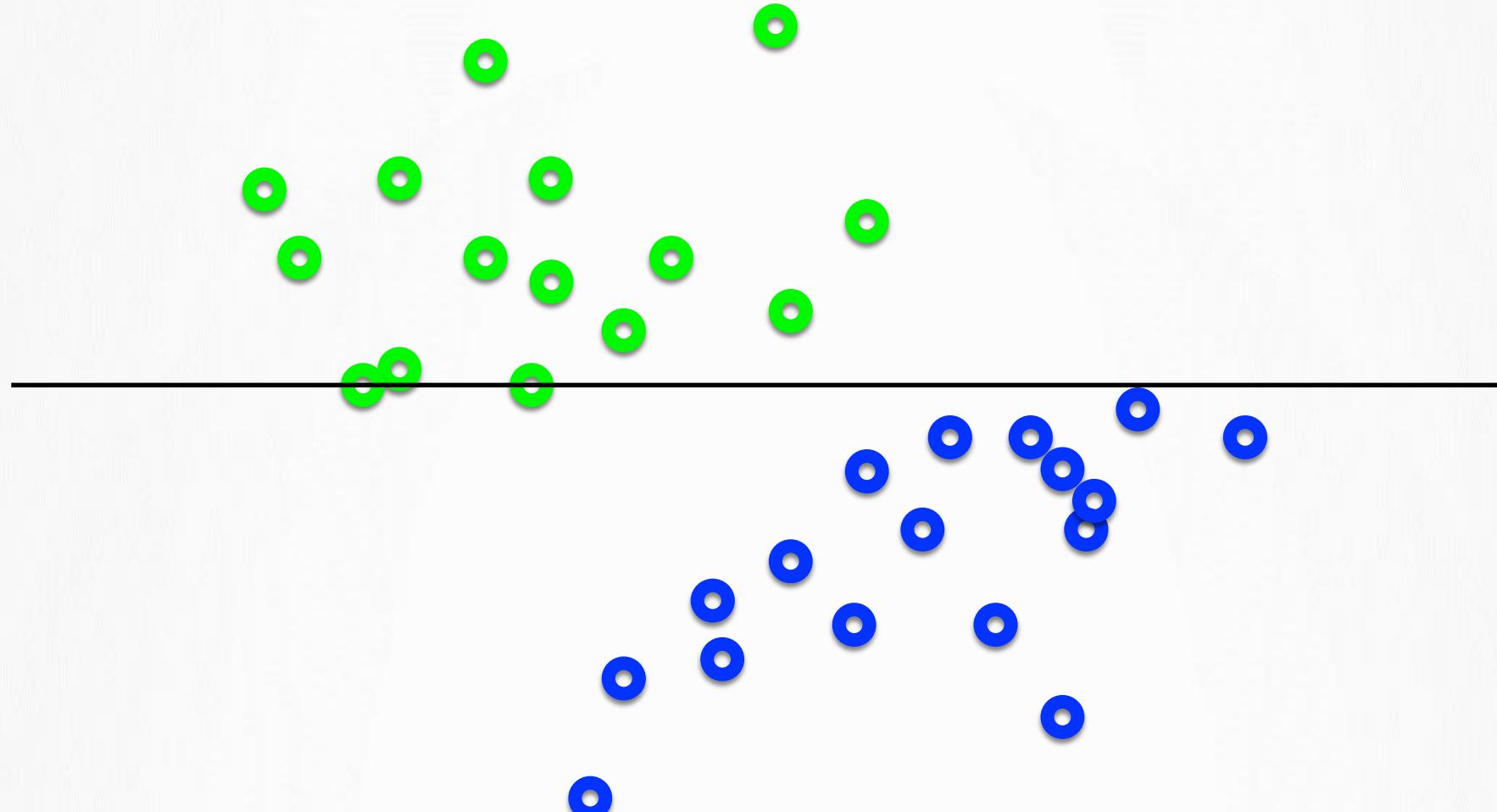


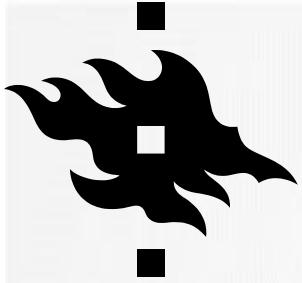
What's the best w?



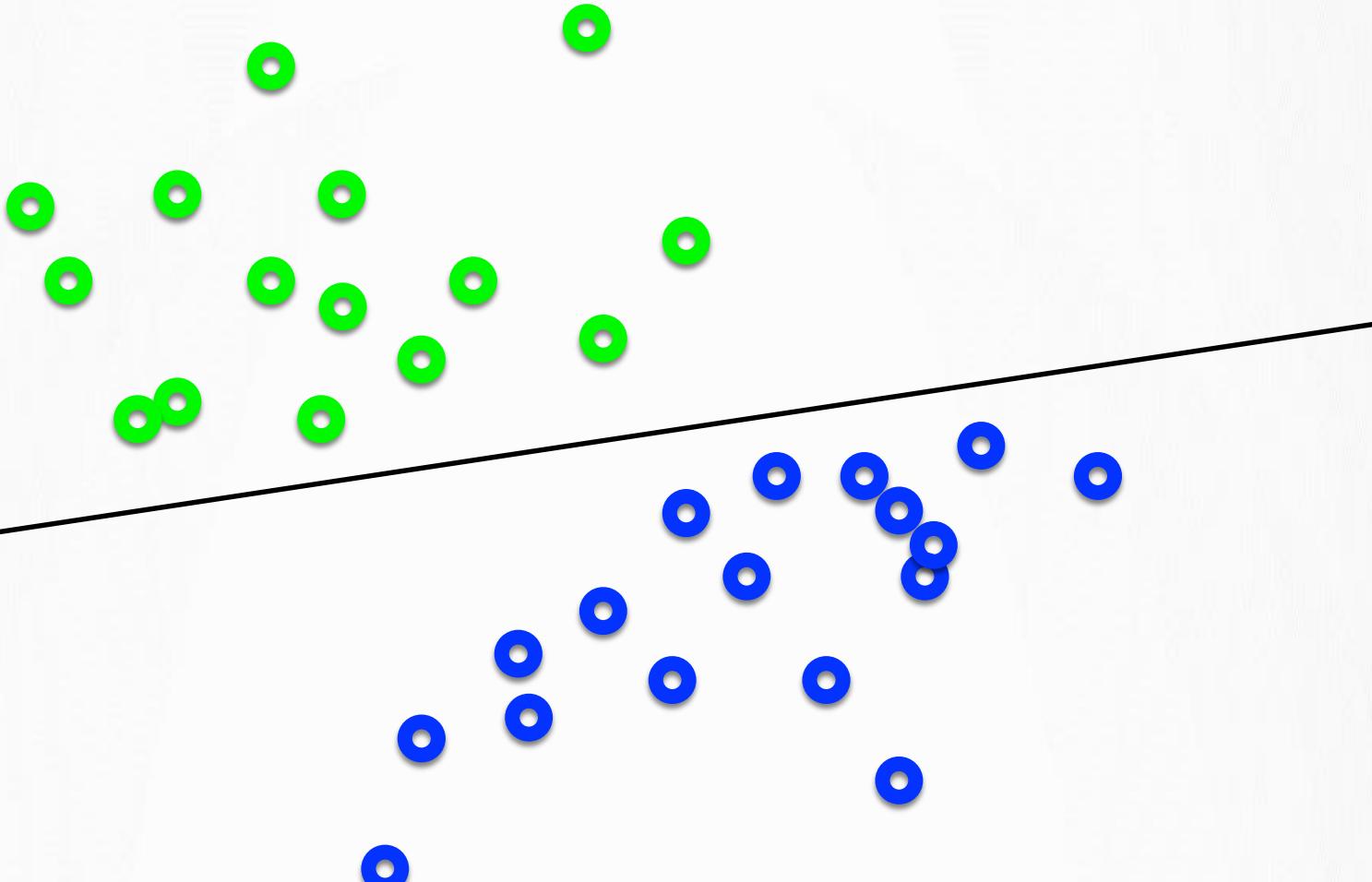


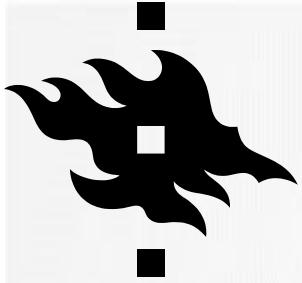
What's the best w?



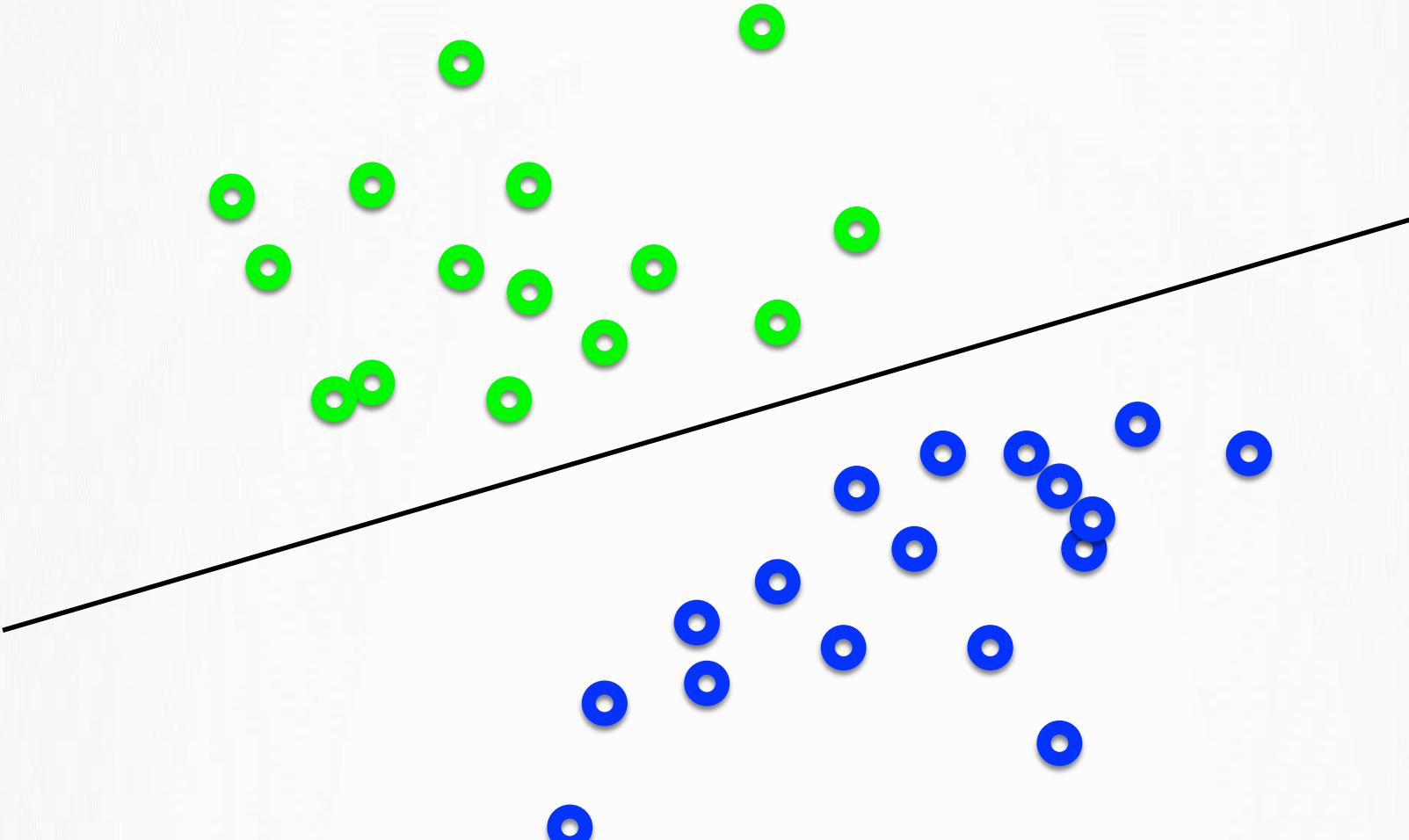


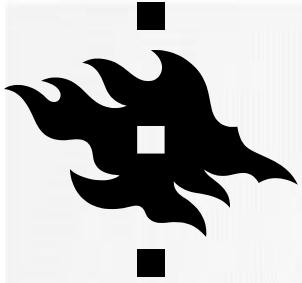
What's the best w?



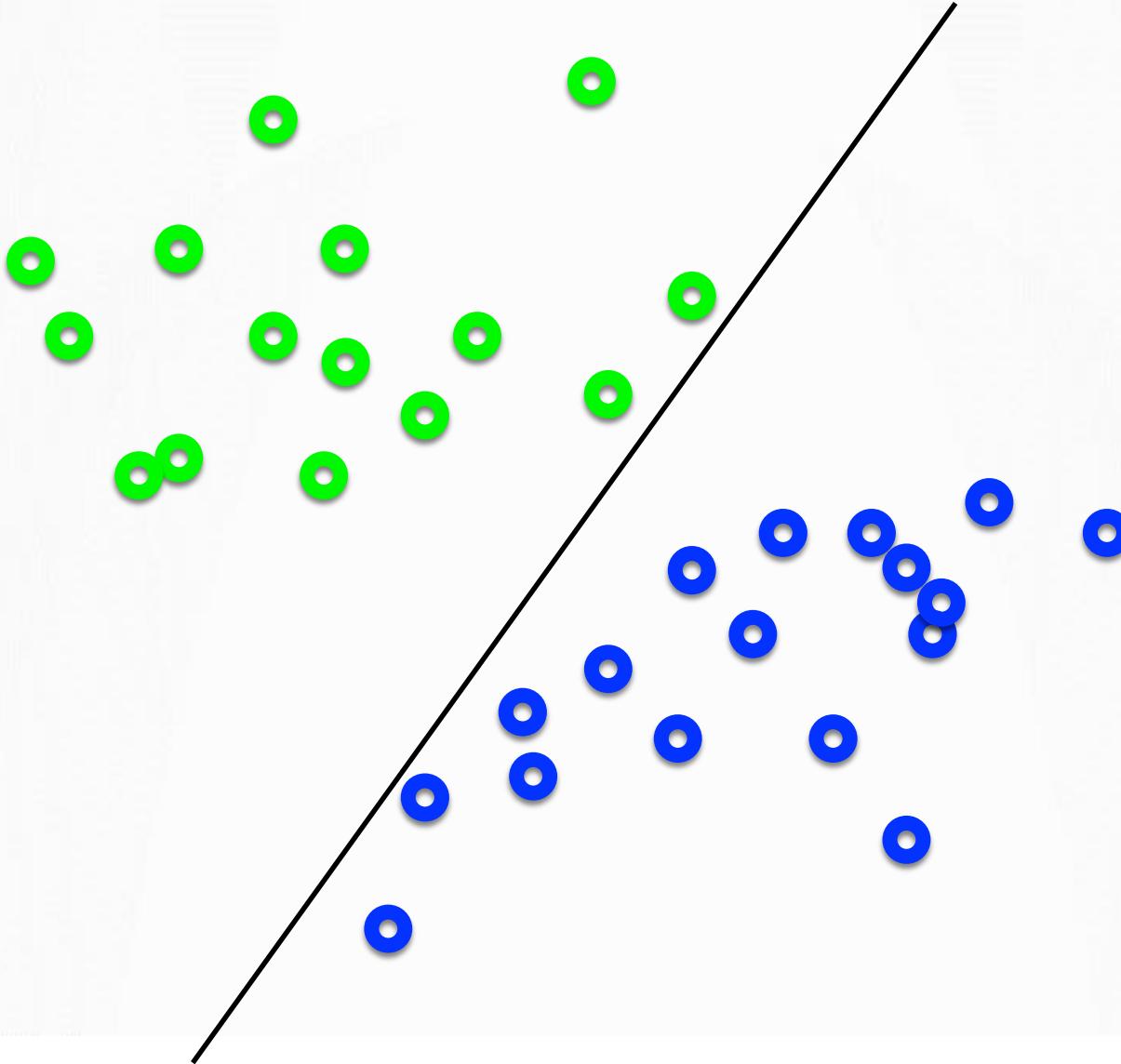


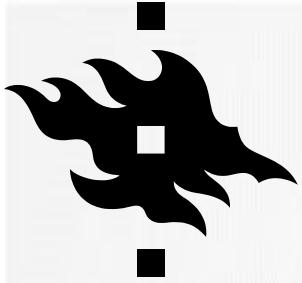
What's the best w ?



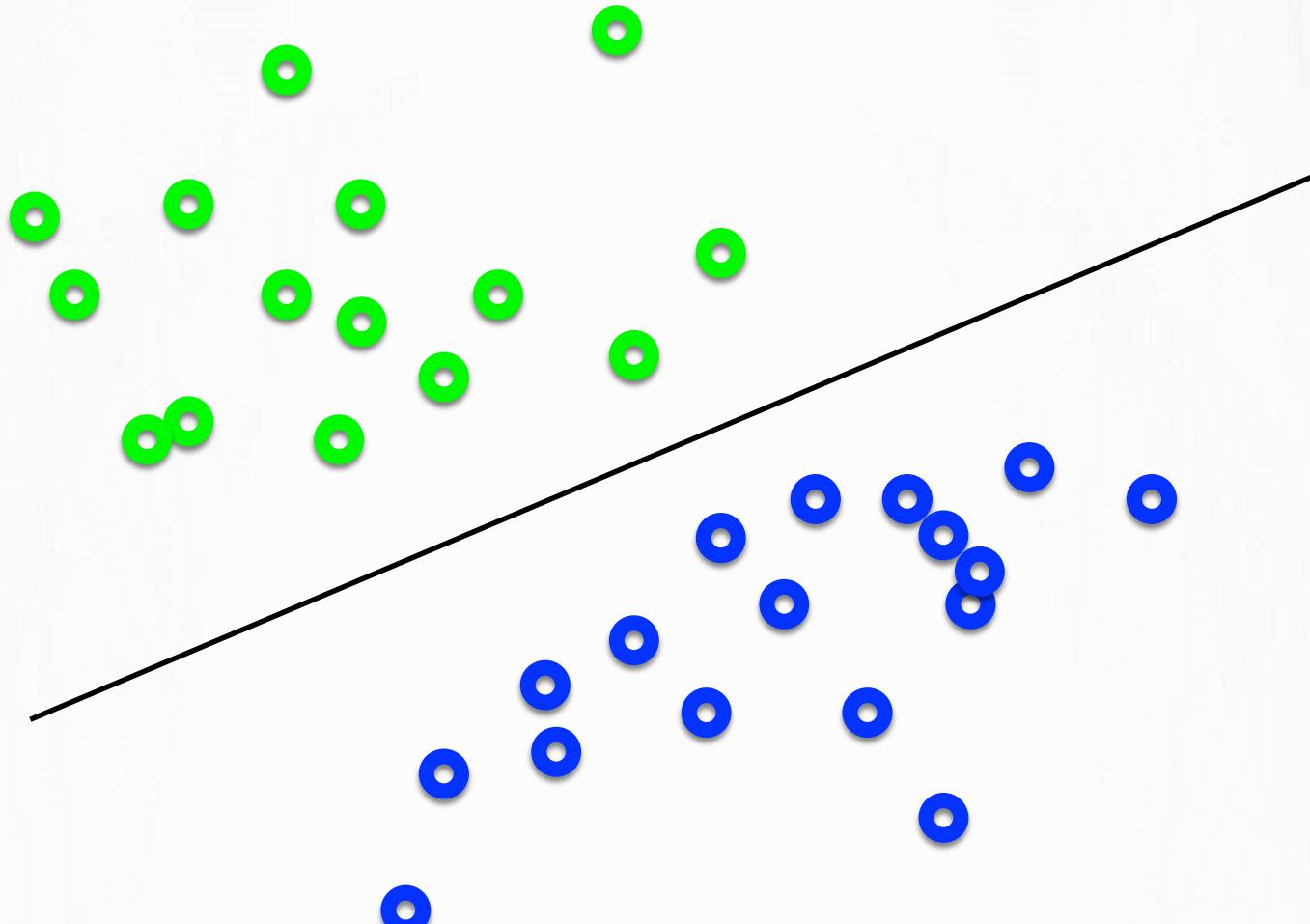


What's the best w ?

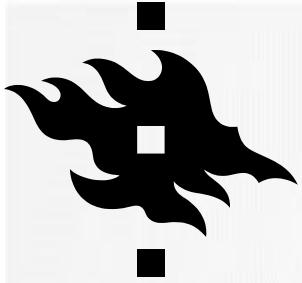




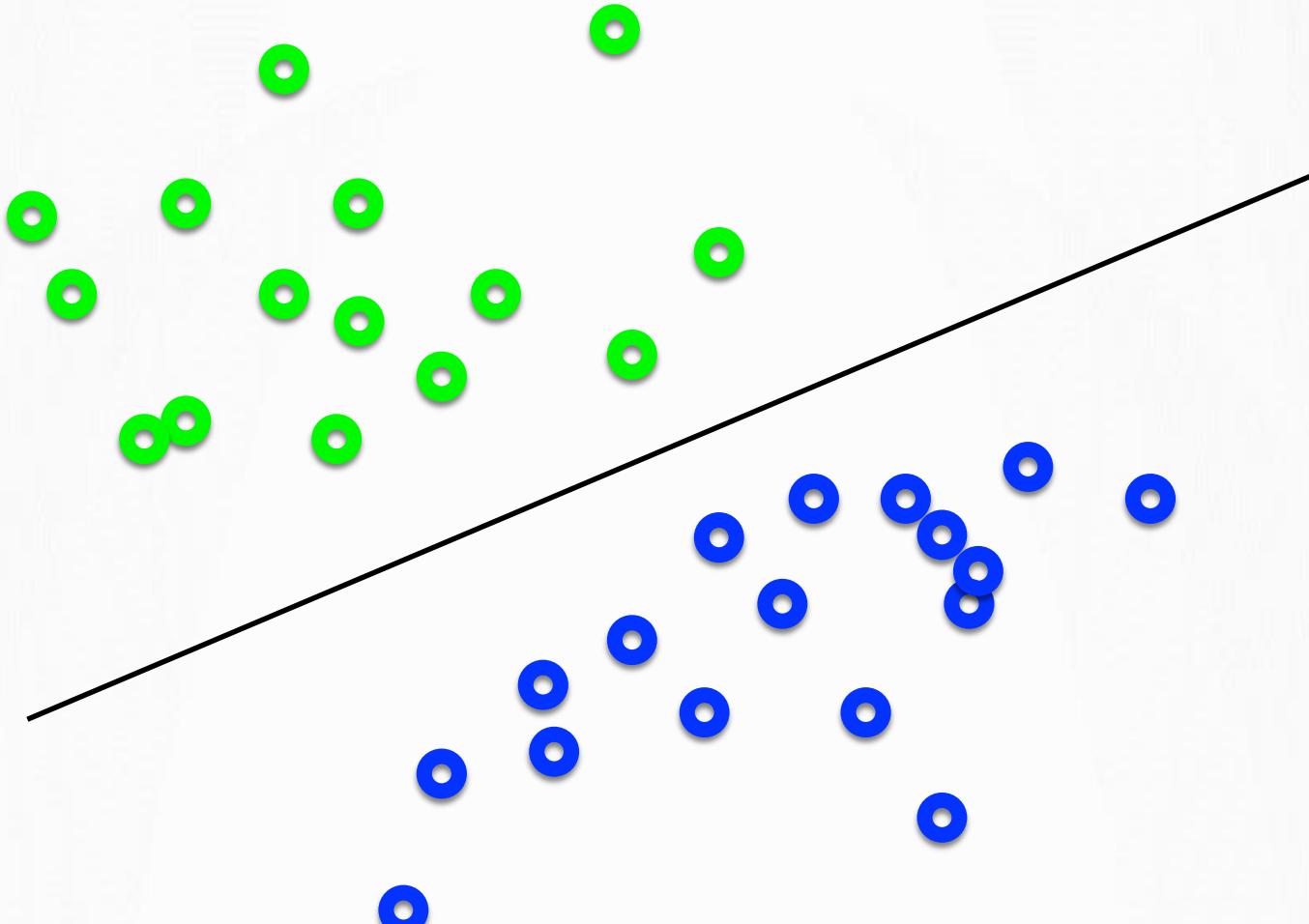
What's the best w ?



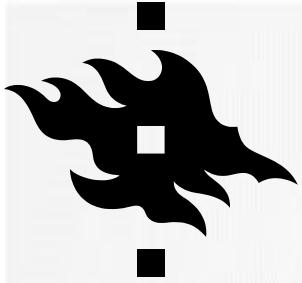
Intuitively, the line that is the
farthest from all interior points



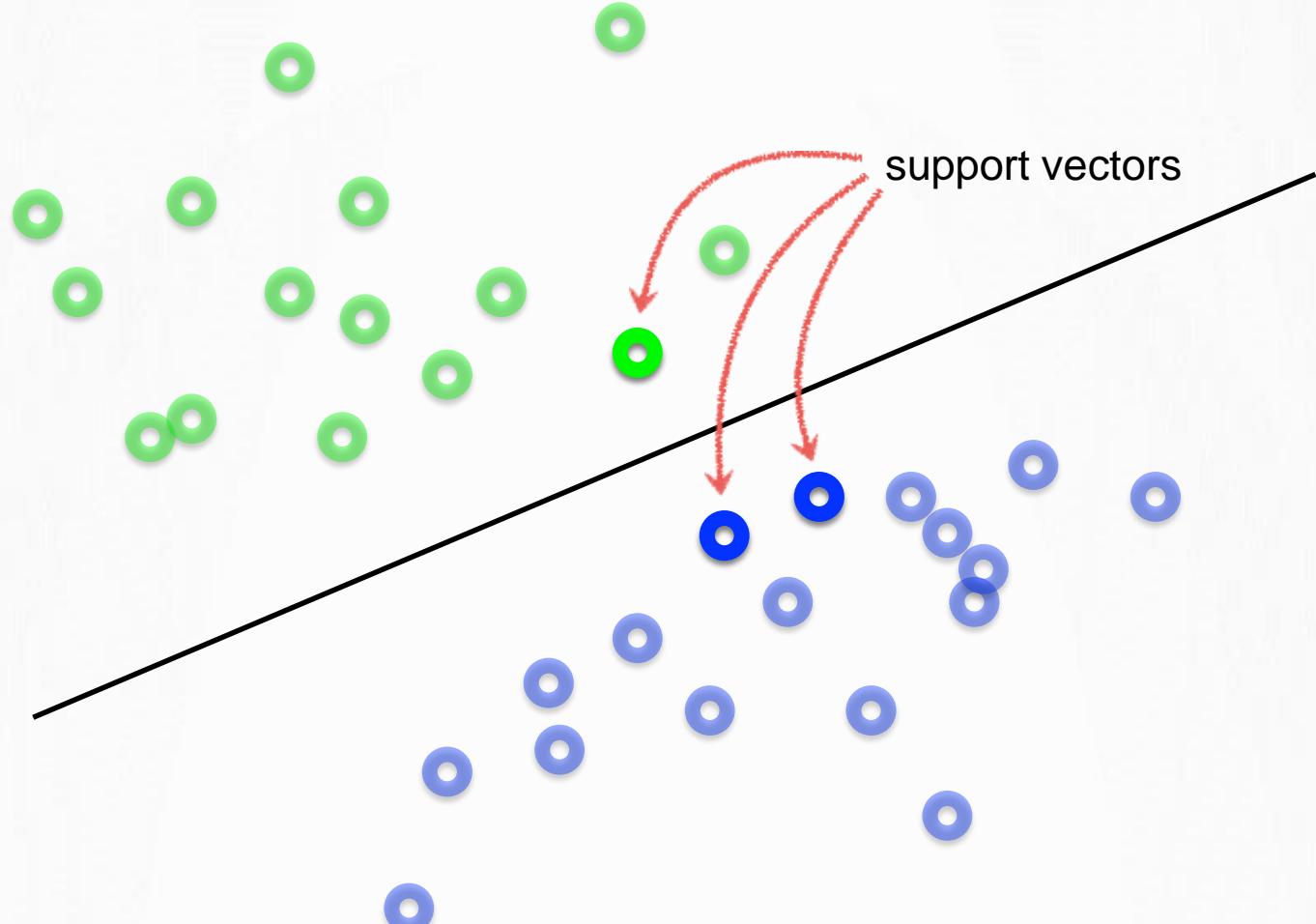
What's the best w ?

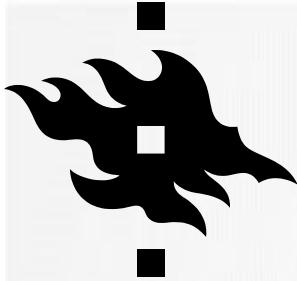


Maximum Margin solution:
most stable to perturbations of data

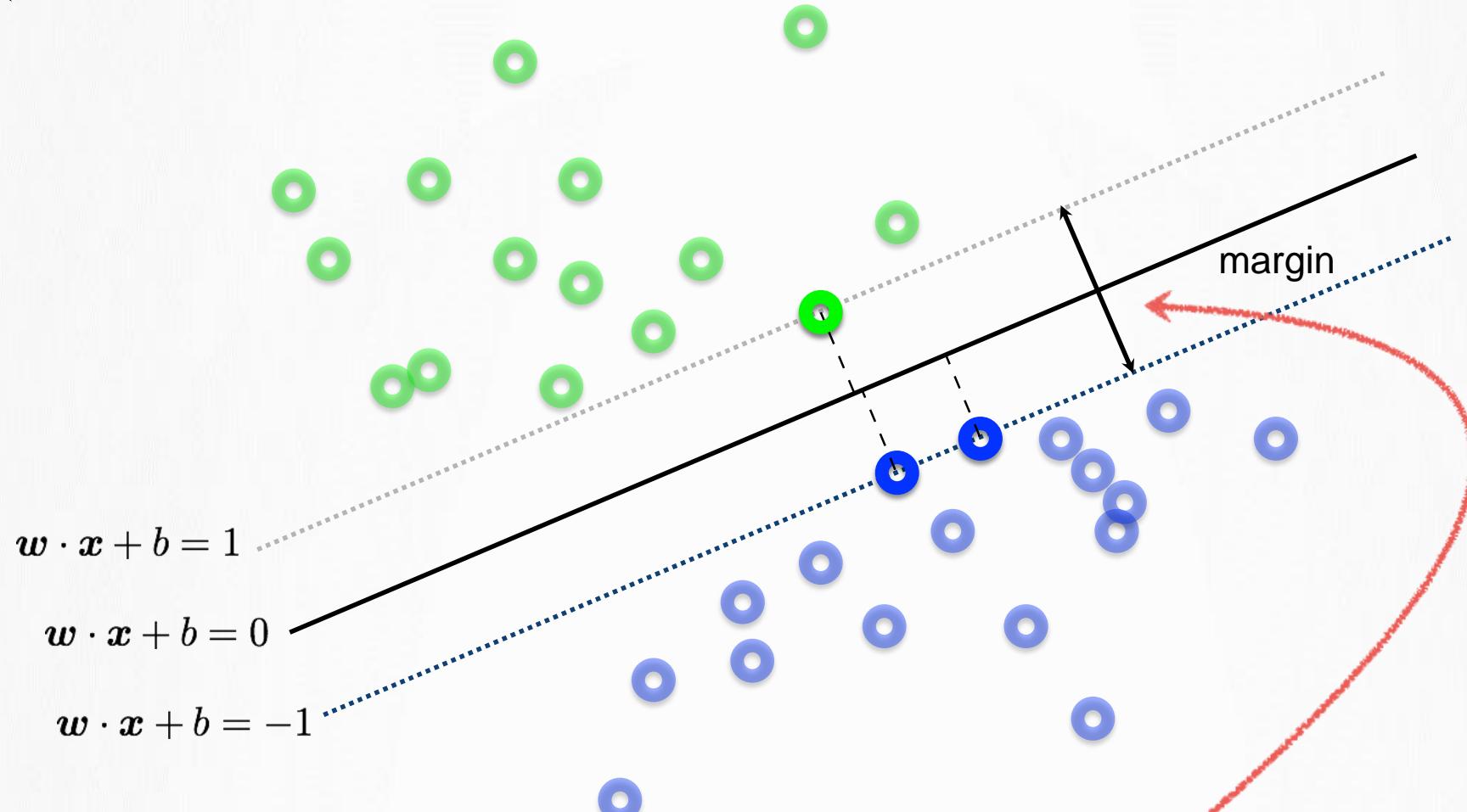


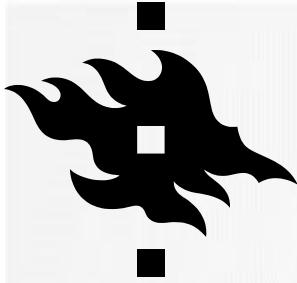
What's the best w ?





Find hyperplane w such that ...





Can be formulated as a maximization problem

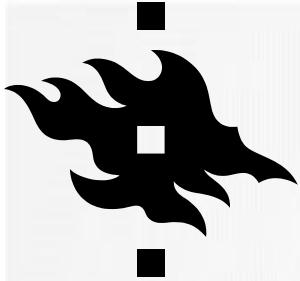
$$\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$$

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

What does this constraint mean?



label of the data point



Can be formulated as a maximization problem

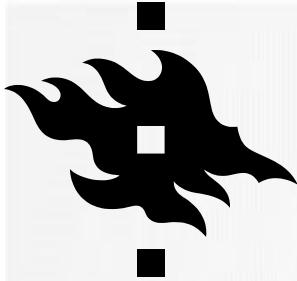
$$\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$$

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

Equivalently,

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

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



'Primal formulation' of a linear SVM

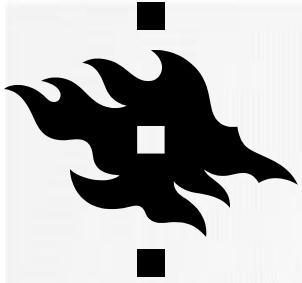
$$\min_{\mathbf{w}} \|\mathbf{w}\|$$

Objective Function

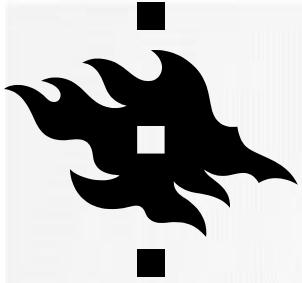
$$\text{subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \text{ for } i = 1, \dots, N$$

Constraints

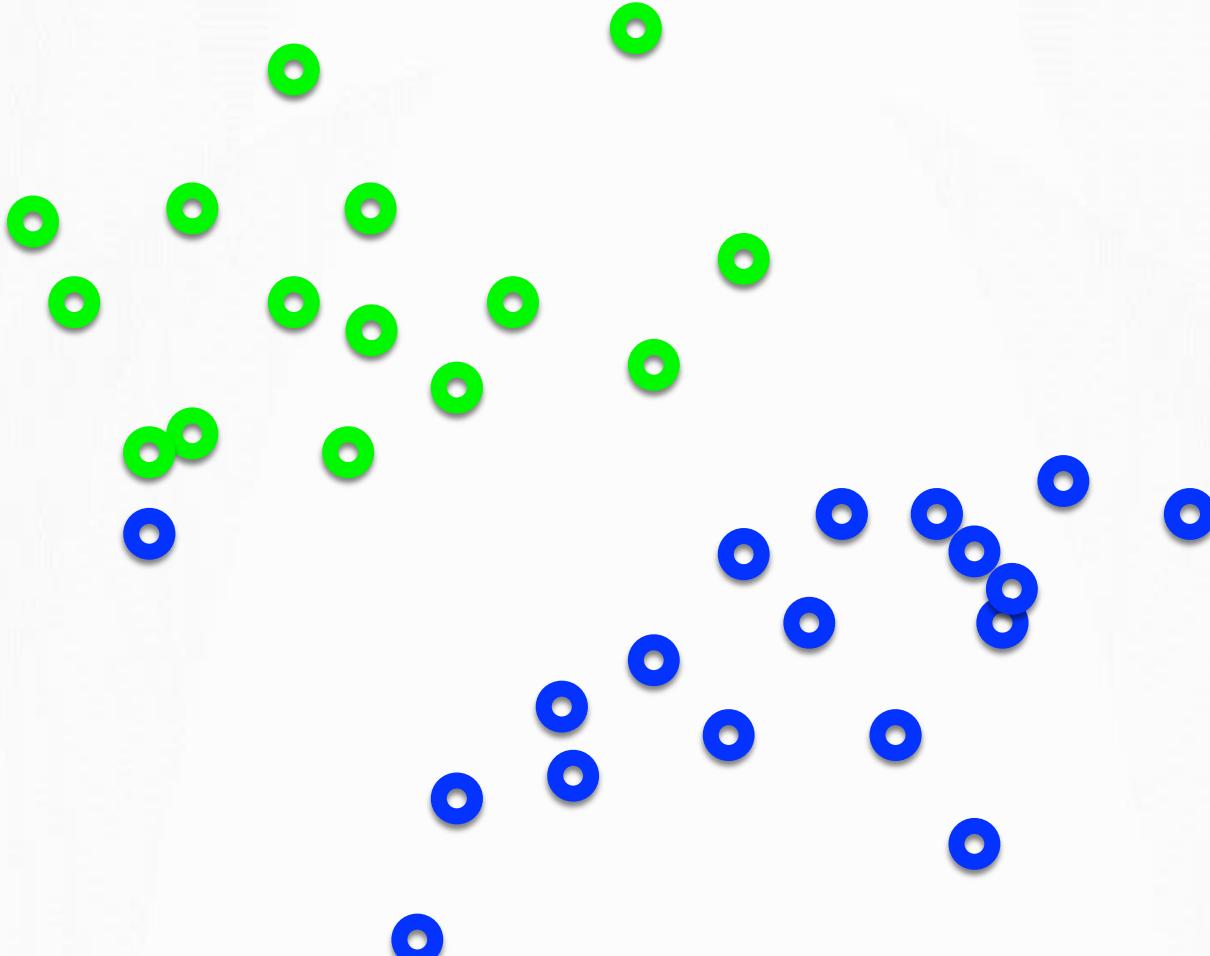
This is a convex quadratic programming (QP) problem
(a unique solution exists)

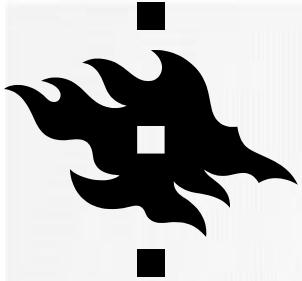


‘SOFT’ MARGIN

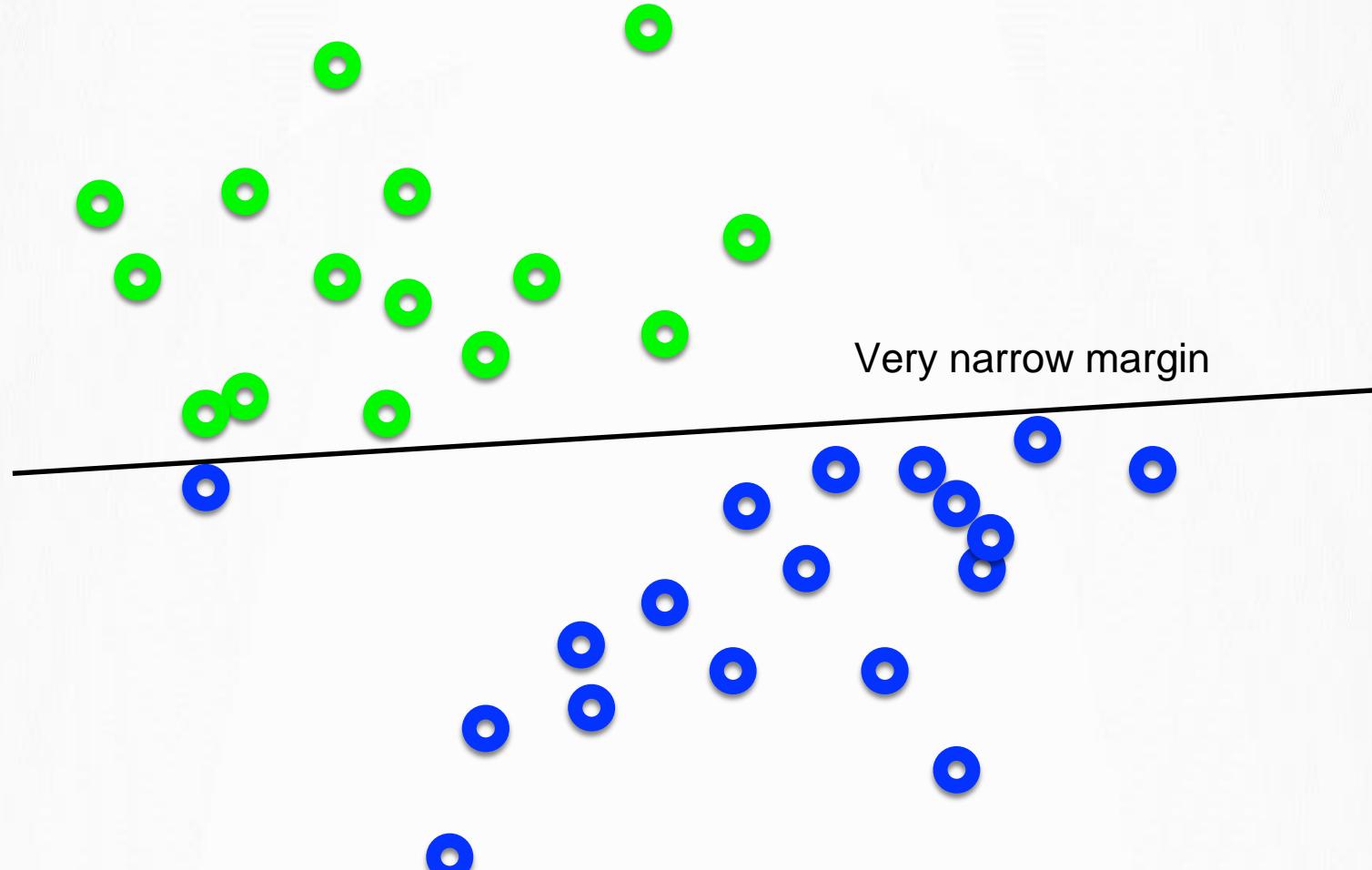


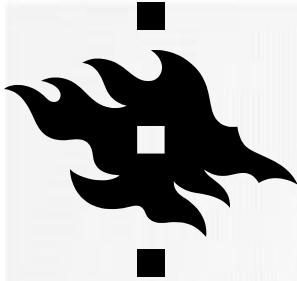
What's the best w?



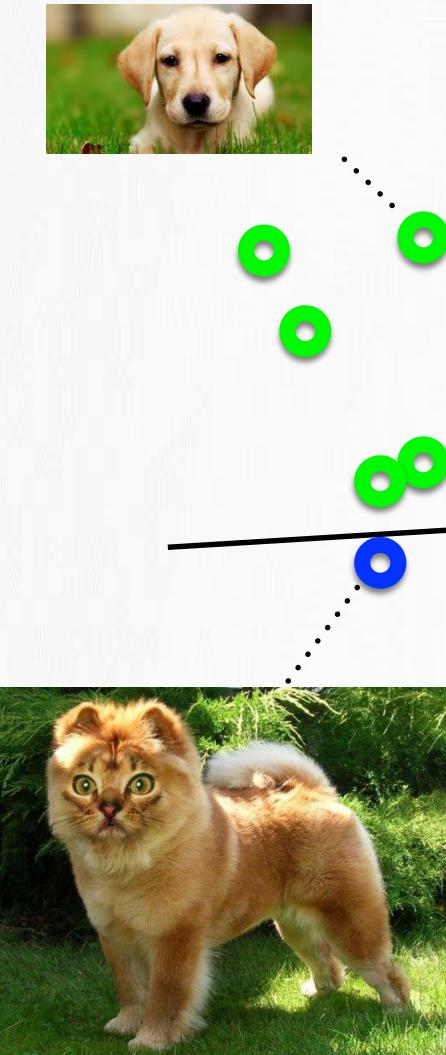


What's the best w?



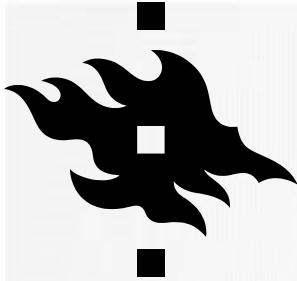


Separating cats and dogs

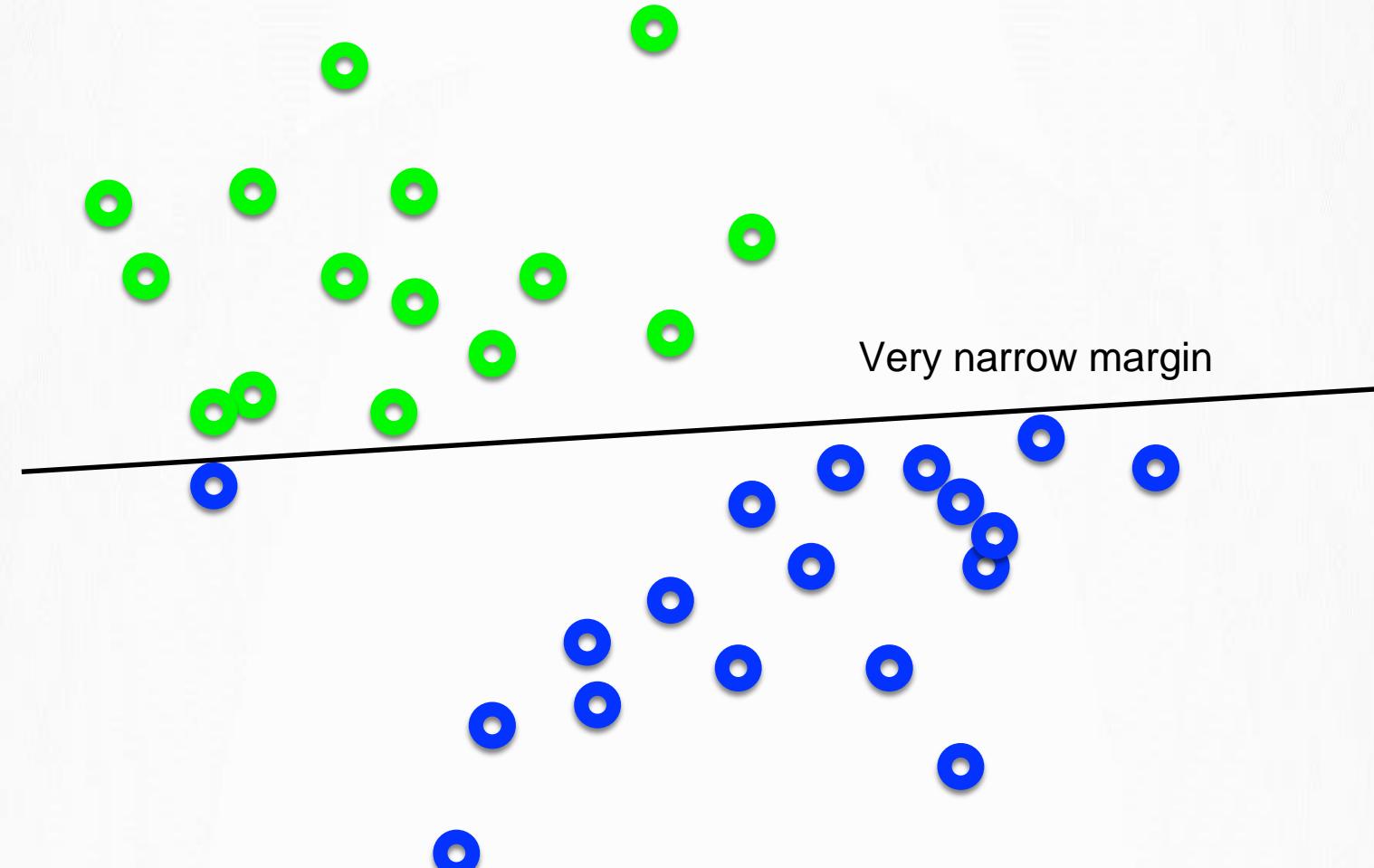


Very narrow margin

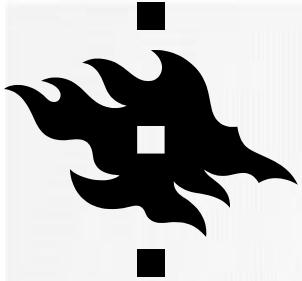




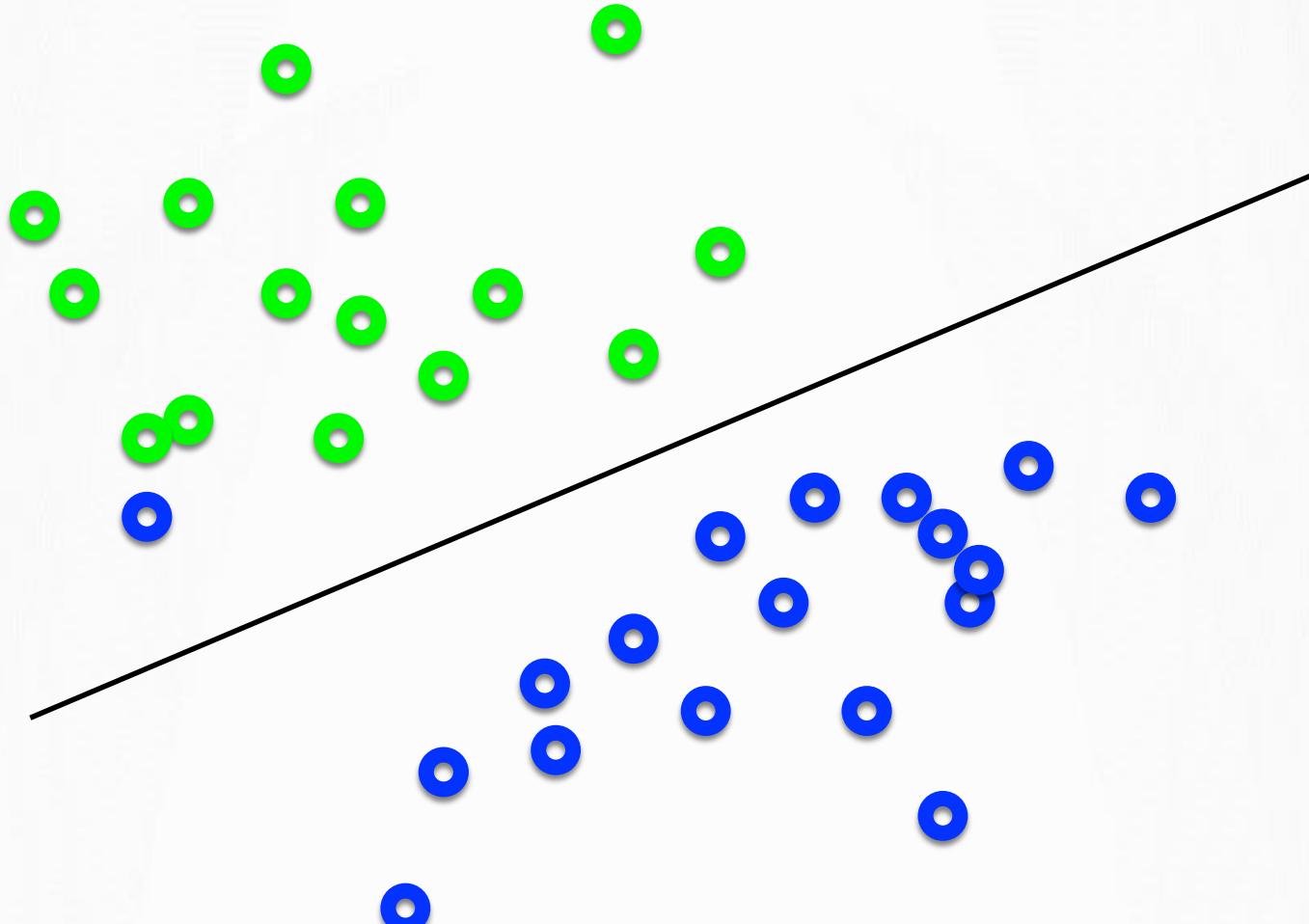
What's the best w ?



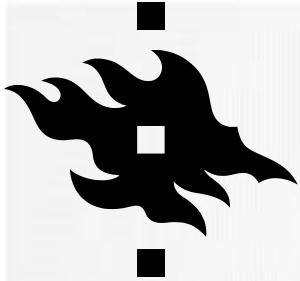
Intuitively, we should allow for some misclassification if we can get more robust classification



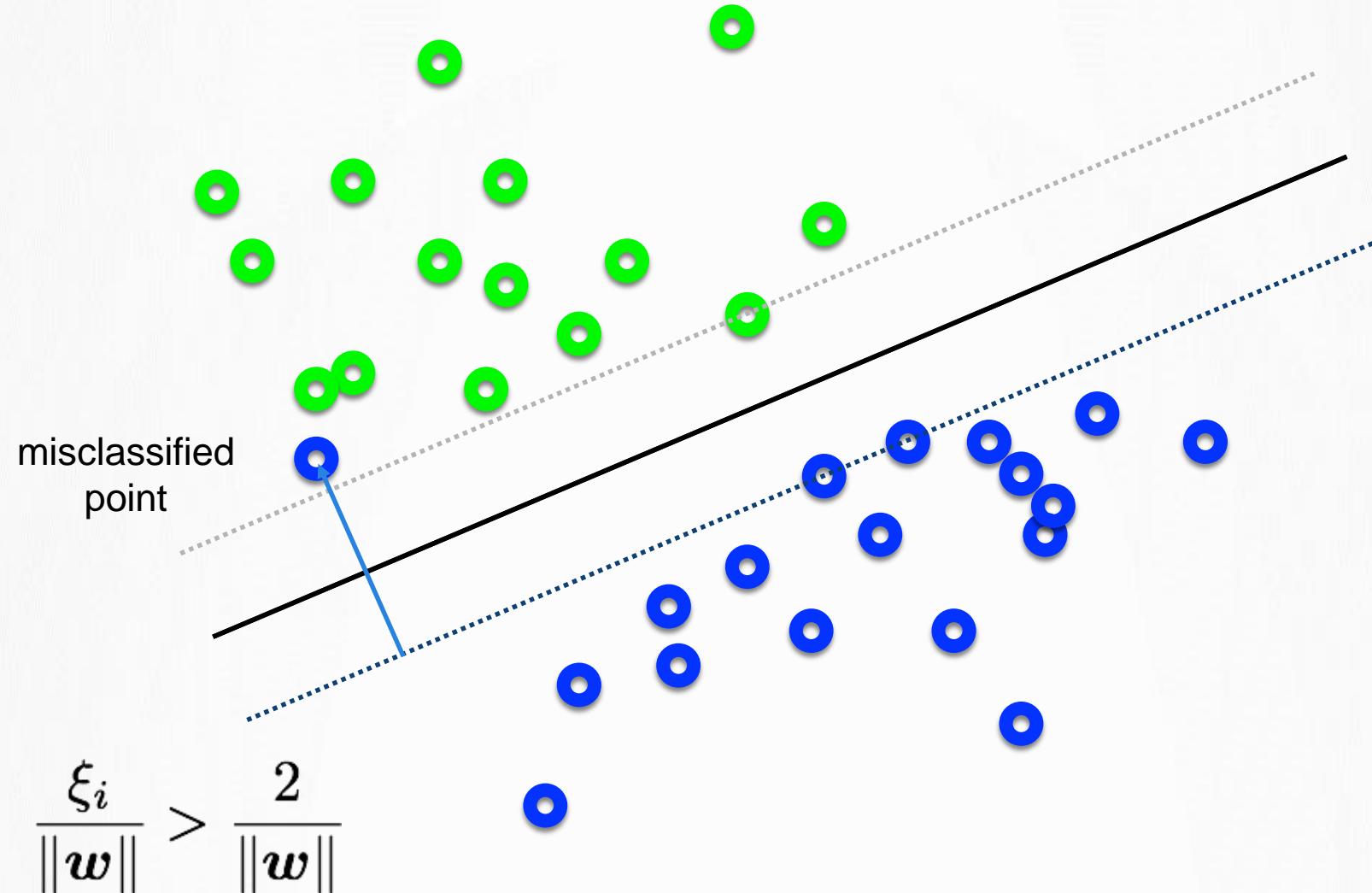
What's the best w?

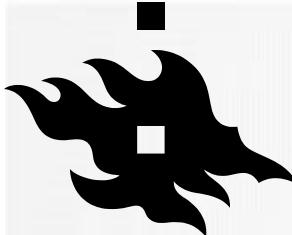


Trade-off between the MARGIN and the MISTAKES
(might be a better solution)



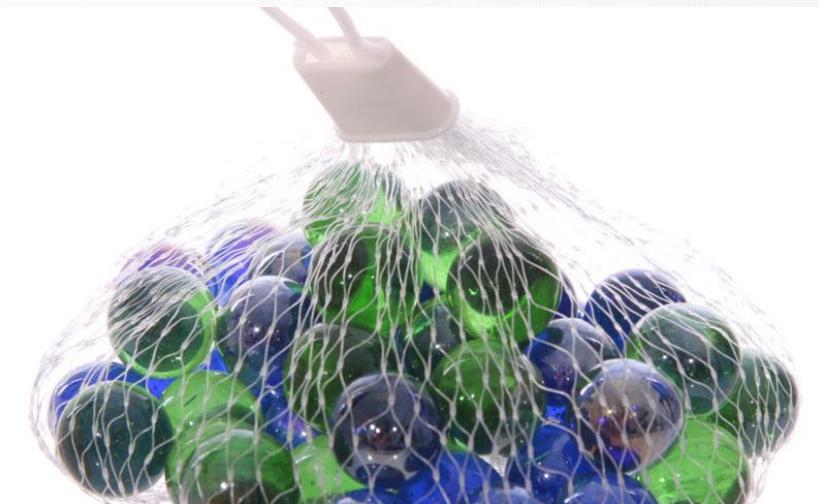
Adding slack variables $\xi_i \geq 0$





OBJECT DETECTION NOW

- Object detection is one of the most challenging computer vision tasks
- Aims at identifying the presence of various individual objects in an image
- Good results have been obtained when dealing with images with relatively simple image scenes and clear foreground objects
- The problem is not adequately addressed when dealing with the images and videos containing objects placed in arbitrary poses, with various shapes, and appearing in a cluttered and occluded environment



Han et al. (2018). Deep Learning for Visual Understanding: Part 2: Advanced Deep-Learning Techniques for Salient and Category-Specific Object Detection. IEEE Signal Processing Magazine.



60° 10 1.2 N, 24° 57 18 E