

Intro to ML for Input Recognition

CHI Course Computational Interaction 2017

Otmar Hilliges

8 May 2017

Explicit interaction



Implicit interaction



Designing Interactive Systems

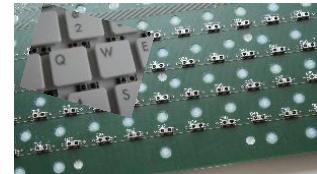
Interaction
Task



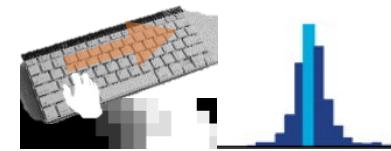
Sensing
Technology



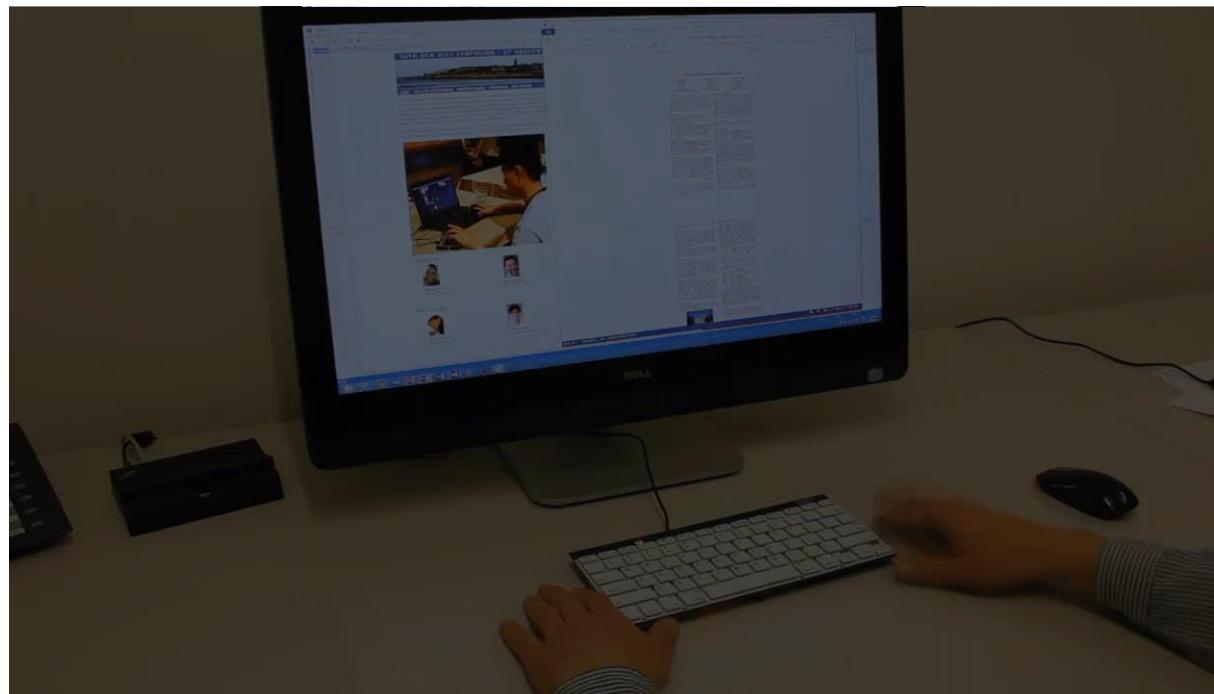
Sensor
Placement



Input
Recognition



System
Output



Today's topic

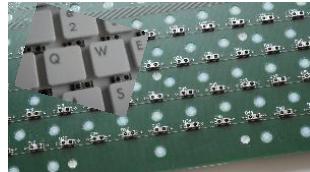
Interaction
Task



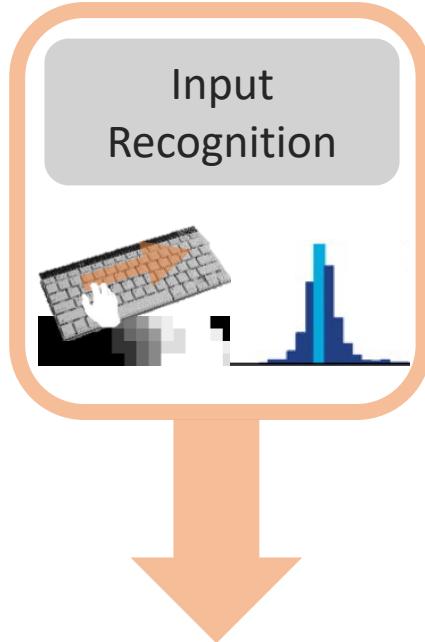
Sensing
Technology



Sensor
Configuration



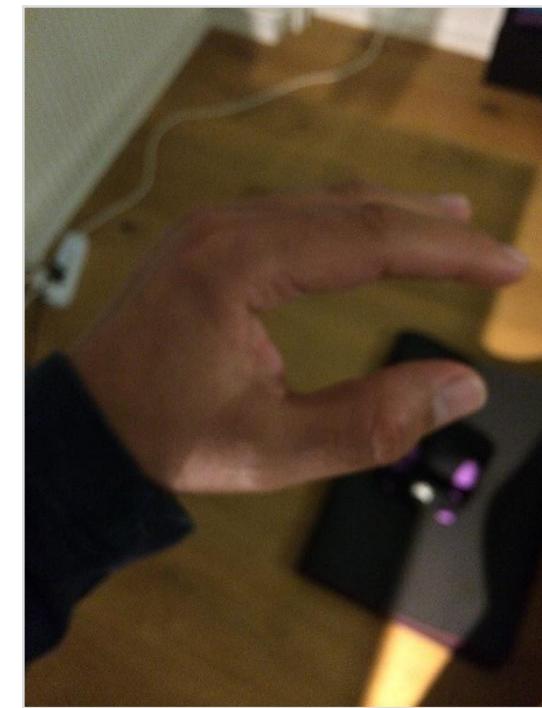
Input
Recognition



System
Output



Today: Machine Learning for Input Recognition



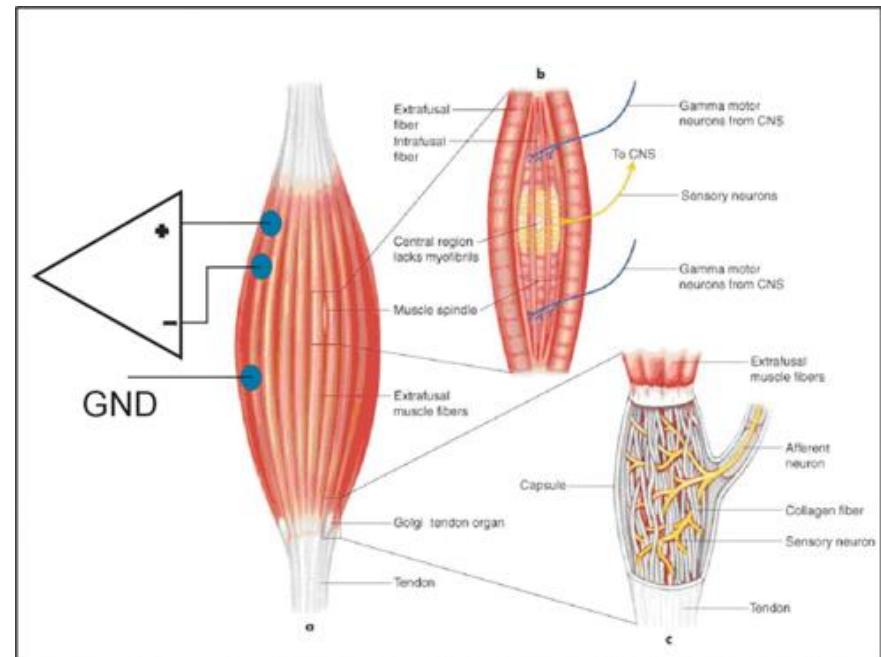
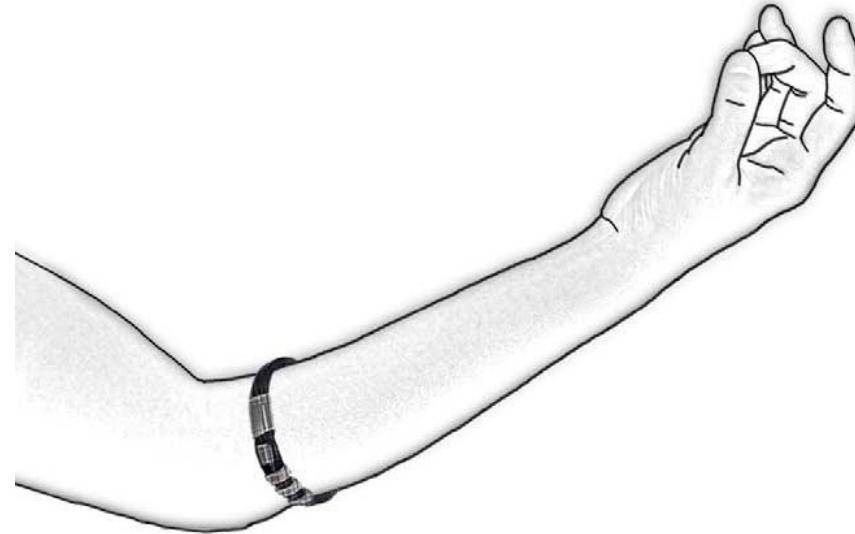
Muscle-Computer Interfaces

Hands free interaction – even when holding an object

EMG or Electromyography
(primarily used in medical therapy)

Action potential generated by muscle when signal arrives from motor neuron

Can be sensed non-invasively by placing electrodes on the skin



Video



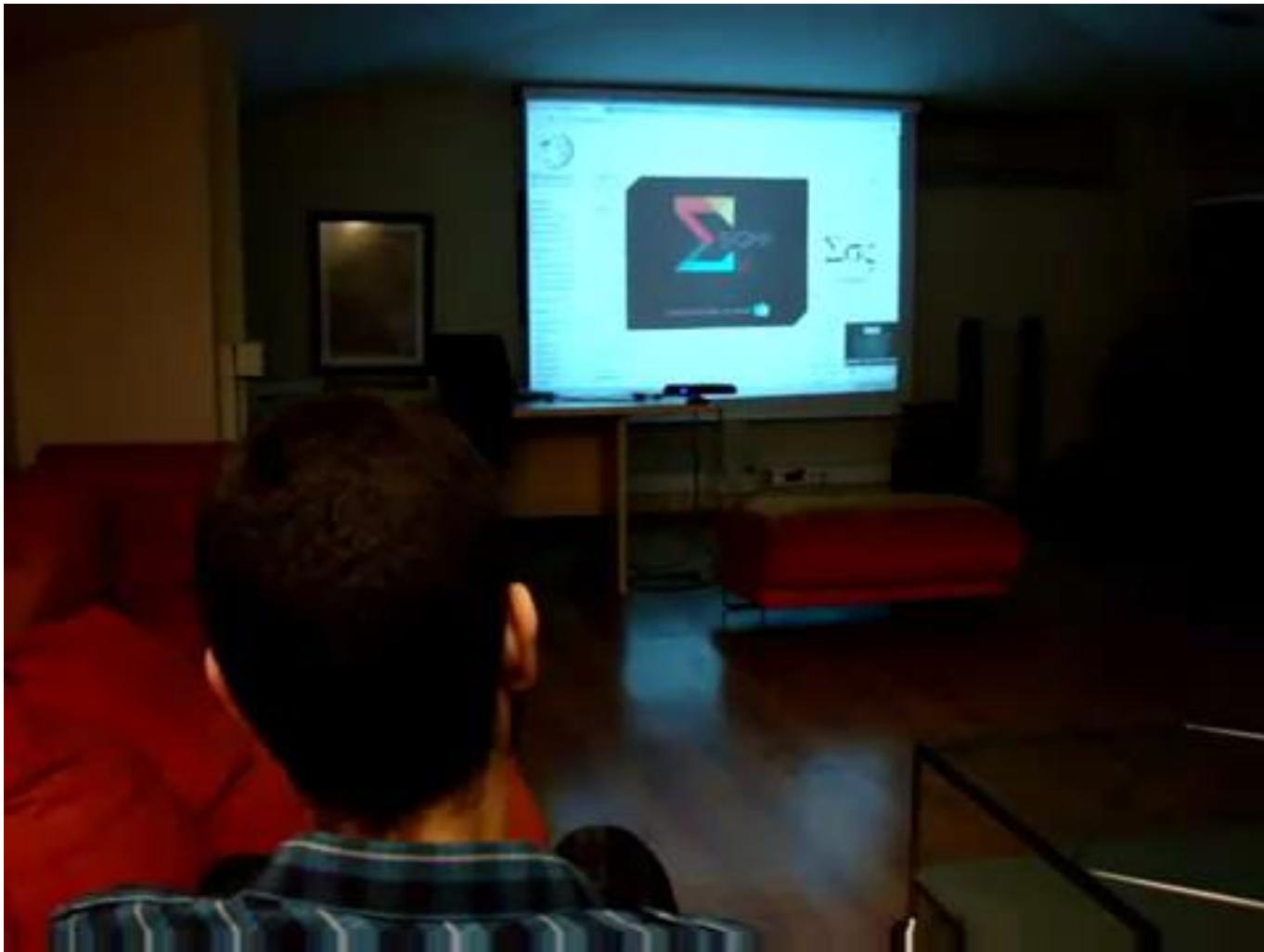
[Saponas et al., *Proceedings of ACM UIST '09*]

Touch & Activate: Video

Touch & Activate: Adding Interactivity to Existing Objects
using Active Acoustic Sensing

Makoto Ono, Buntarou shizuki, and Jiro Tanaka
University of Tsukuba

Gesture Recognition



<http://www.sigmar.d.com/portfolio/sigmanil-framework/>

Take Home Message

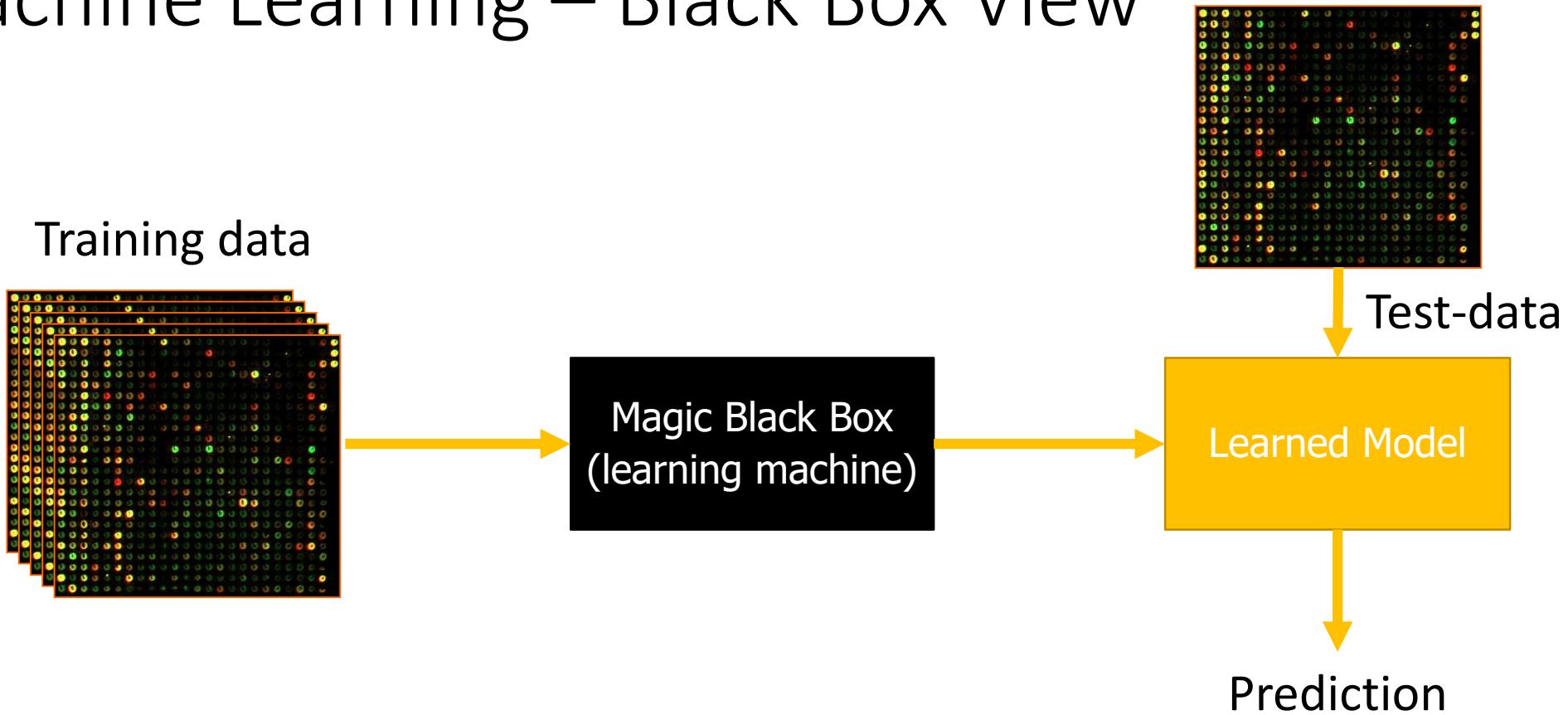
Humans use a variety of channels to interact & communicate

- Human behavior is highly complex
- Human behavior is highly ambiguous
- To create we intuitive interfaces we want to use:
 - Small, low-effort natural motion
 - Avoid heavy user instrumentation

In consequence,

- Recognizing human input is challenging
- Most heuristic / traditional approaches are not powerful enough to model the variance in (gesture) execution
- Many different sensor modalities, often noisy data

Machine Learning – Black Box View



Training data:

- Positive *and* negative examples. Data has to be labelled (ground truth).

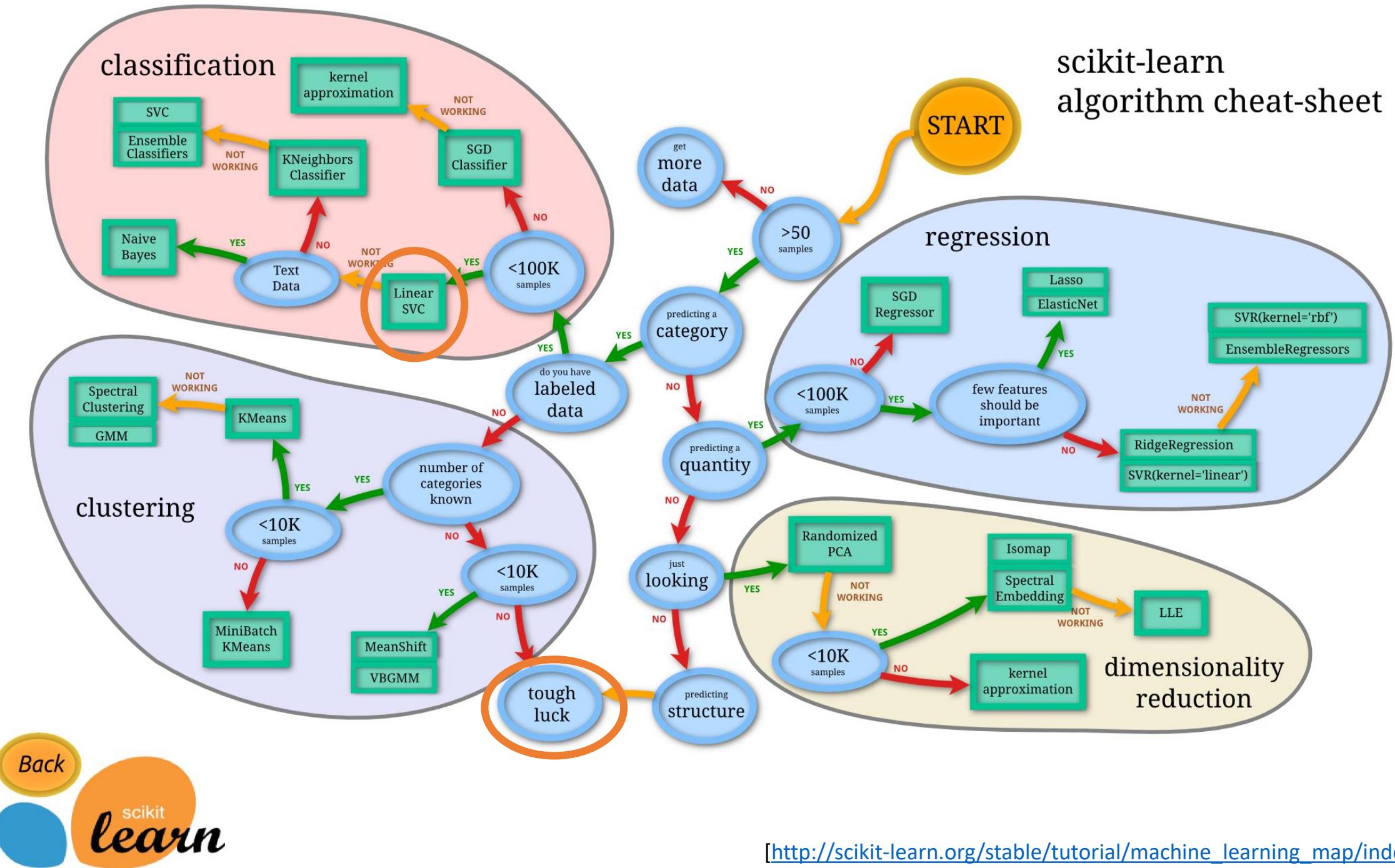
Model:

- The model can distinguish between classes or estimate real valued outcome.

Categorization of Machine Learning Tasks (+HCI examples)

- ***Classification***: inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one or more of these classes (**recognize gestures**).
- ***Regression***: outputs are continuous rather than discrete (**regress x,y,z positions from electro-magnetic field**).
- ***Clustering***: divide a set of inputs into (a priori unknown groups).
(Find similarity in users in some data domain)
- ***Dimensionality reduction***: simplifies inputs by mapping them into a lower-dimensional space.
(Visualize high-dimensional data for human consumption)

scikit-learn algorithm cheat-sheet

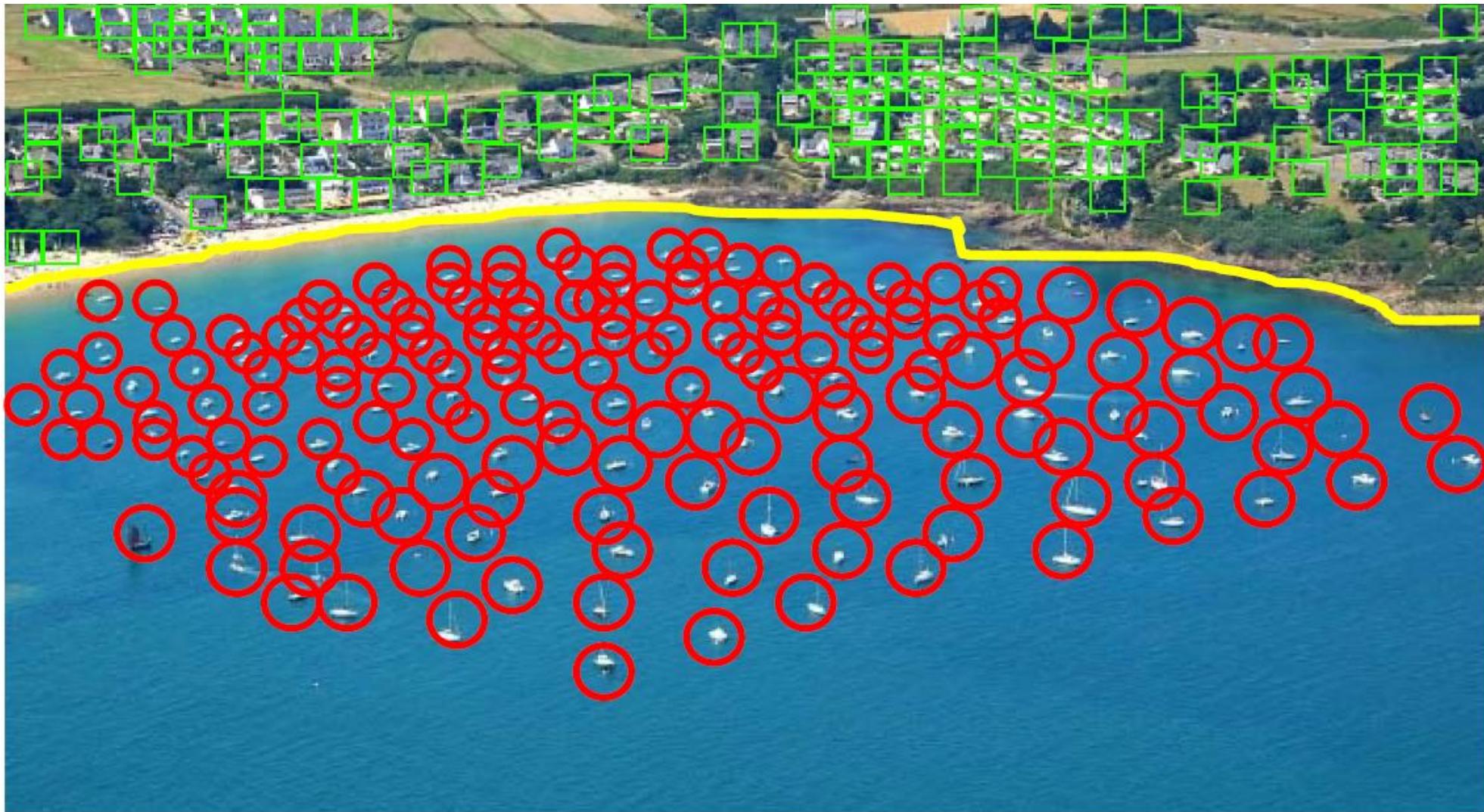


Let's get our hands dirty!

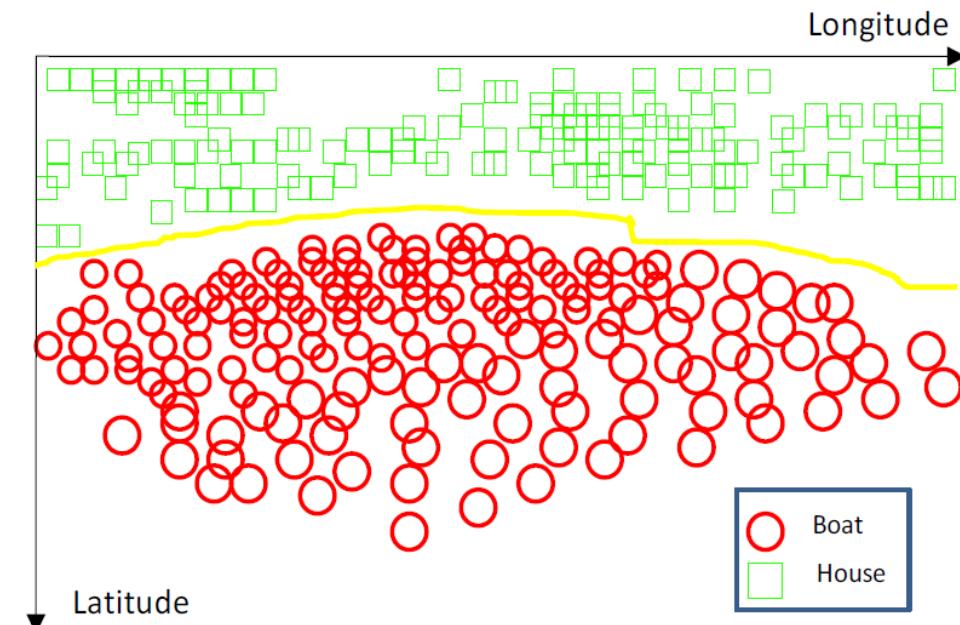
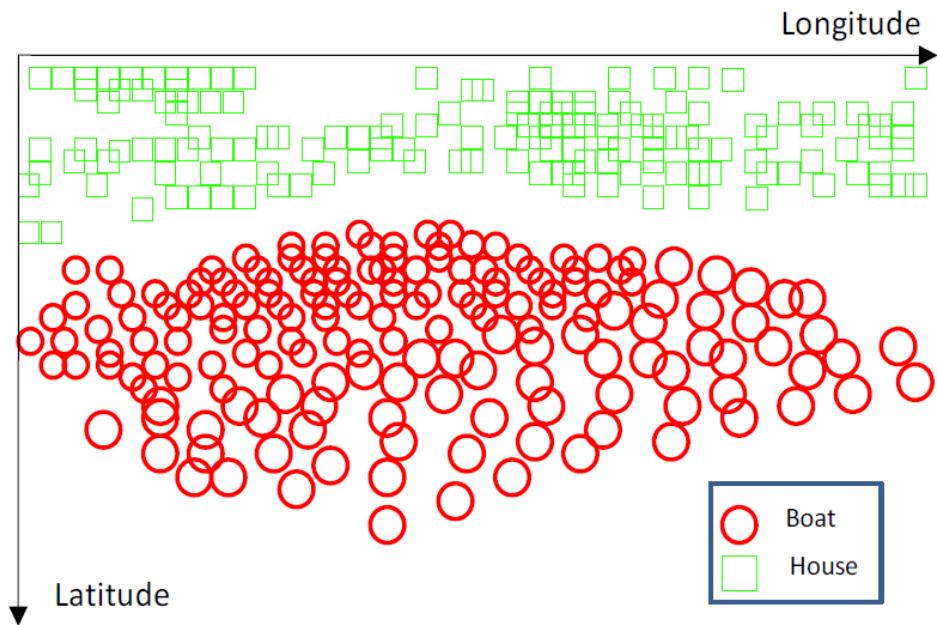
Classification: The Task



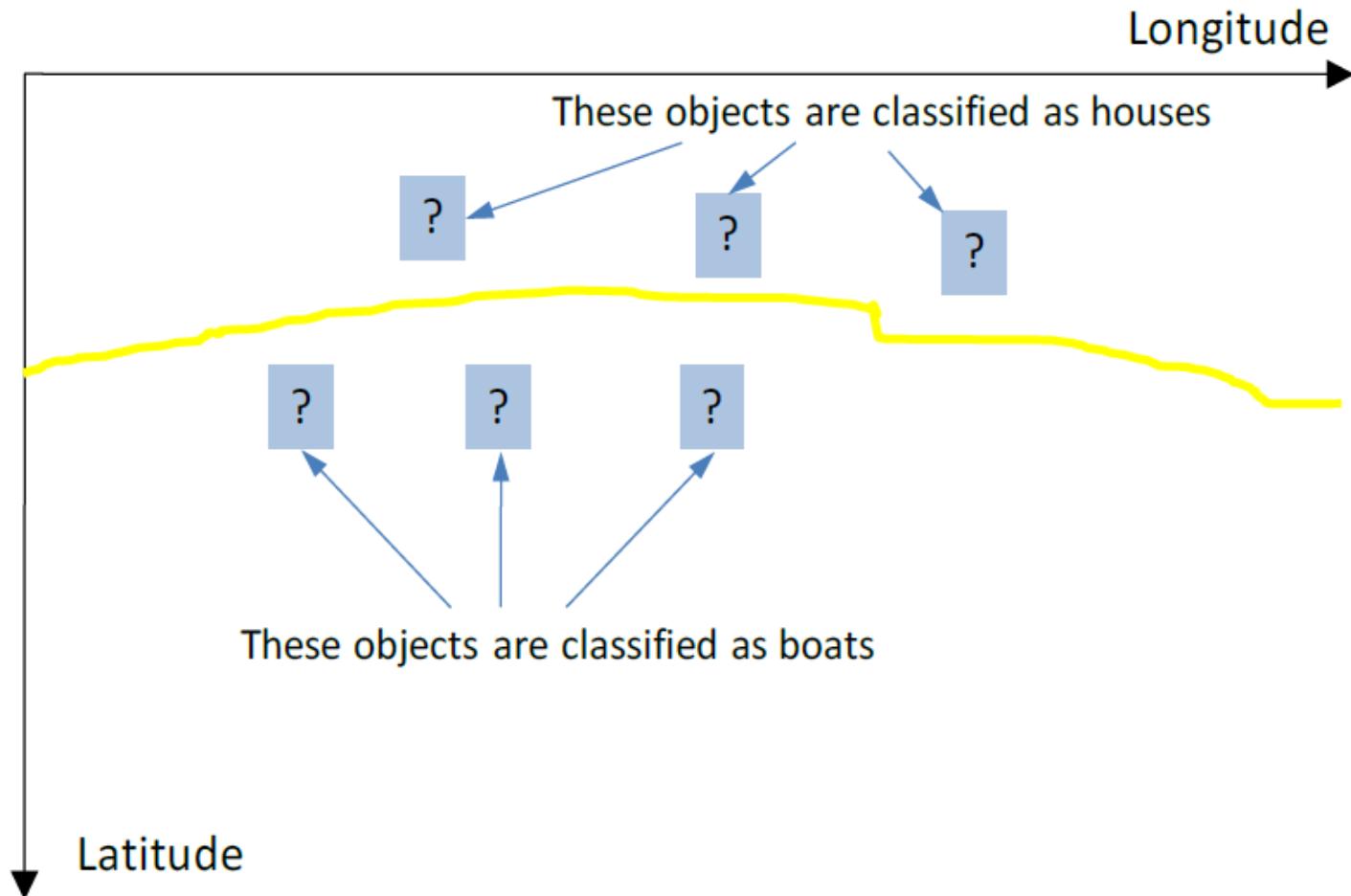
Classification: Decision Surface



Classification Algorithms – Training time



Classification Algorithms – Test time

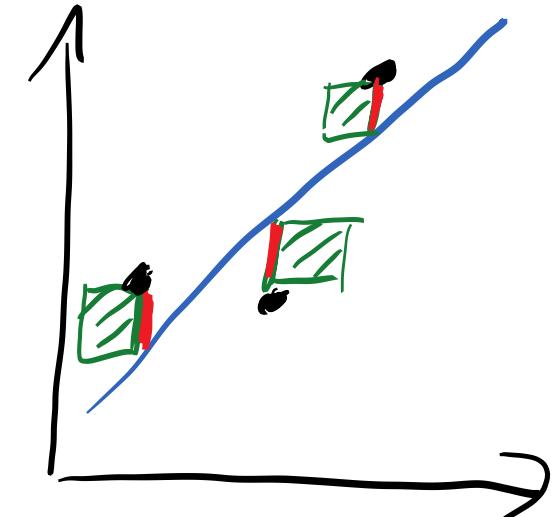


Classification: Errors



Supervised learning - Ingredients

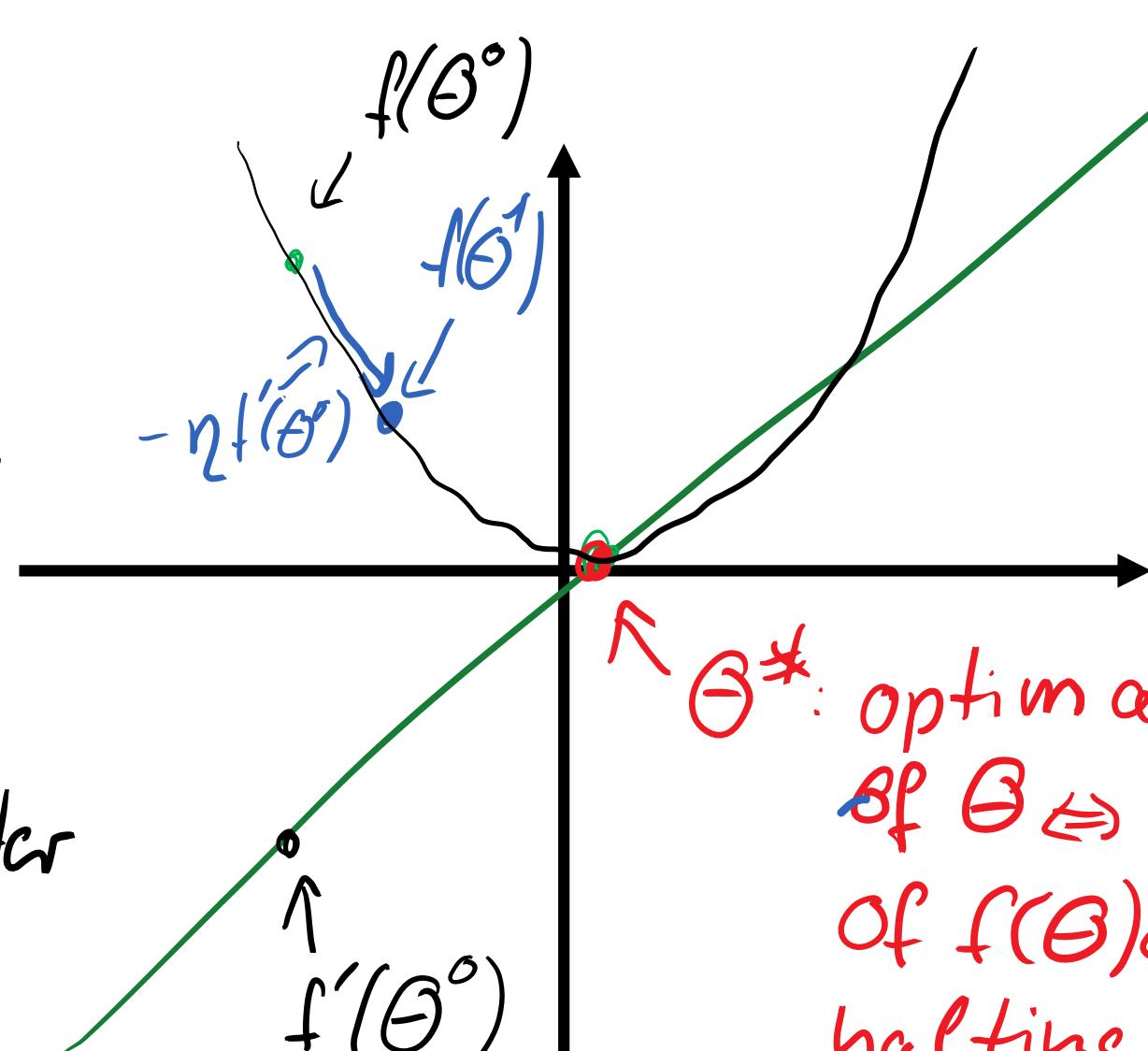
- **Training data** $D = \{x_1, y_1, \dots, x_n, y_n\}$, where $y_i = \begin{cases} 1 & \text{if } x \text{ positive} \\ -1 & \text{otherwise} \end{cases}$
- **Model** $g: X \rightarrow Y$ -> g is element of the hypothesis space G and has parameters Θ . We write its predictions as $g(x|\Theta)$
- **Loss function** $E(\Theta|X) = \sum_i L(y_i, g(x_i|\Theta))$
- **Optimization procedure** $\Theta^* = \operatorname{argmin}_{\Theta} E(\Theta|X)$



Gradient Descent in One Picture

update rule:

- randomly pick start value for parameter θ
- compute negative gradient at θ^0 :
 $-f'(\theta^0)$
- update parameter
 $\theta^1 = \theta^0 - \eta f'(\theta^0)$
↑ learning rate



$\hat{\theta}^*$: optimal value
if $\theta \Leftarrow$ minimum
of $f(\theta) \Leftarrow$
halting point of
gradient descent

Support Vector Machine (SVM)

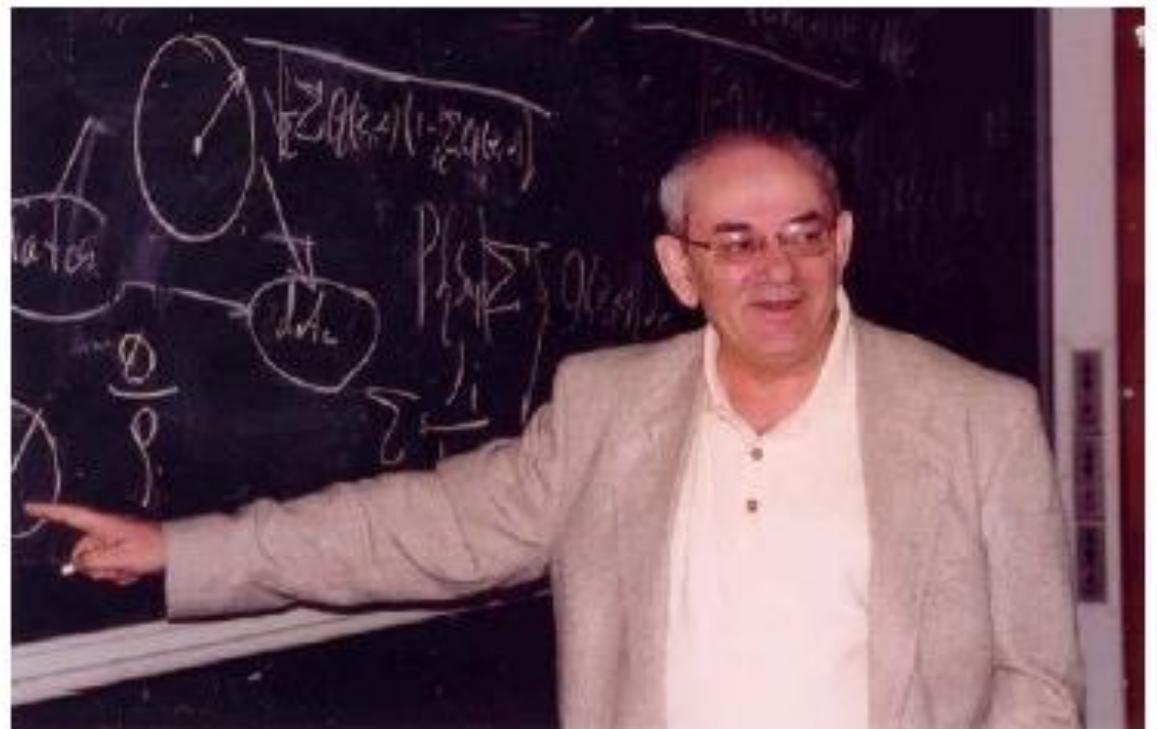
Derived from statistical learning theory by Vapnik, et al. in 1992

Simple SVM achieved accuracy comparable to neural-network with hand-designed features in a handwriting recognition task

SVM widely used in object detection & recognition, content-based image retrieval, text recognition, biometrics, speech recognition, etc.

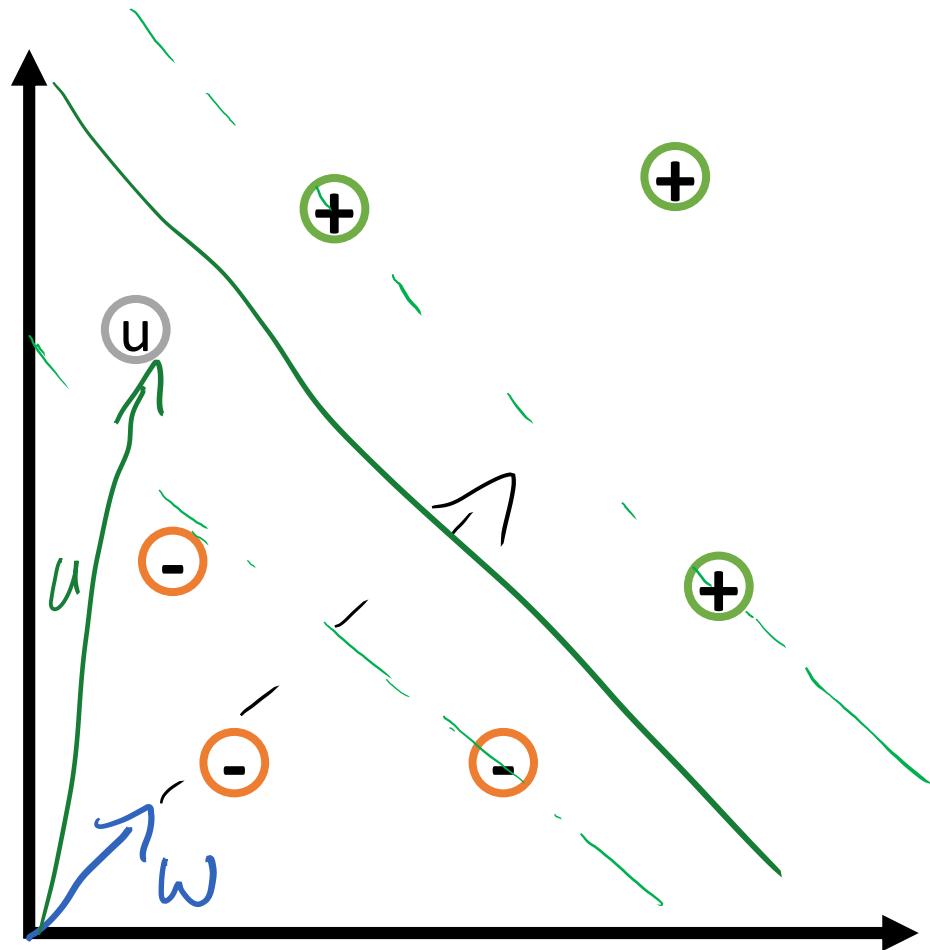
Also used for regression (will not cover today)

Vladimir Vapnik





SVM Derivation (Intuition with some slight of hand)



$$\bar{w} \cdot \bar{u} \geq c$$

with $b = -c$

$\bar{w} \cdot \bar{u} + b \geq 0$

then \dagger

how to find w & b ?

SVM Derivation (Intuition with some slight of hand)

We need to introduce constraints:

$$(\bar{w} \cdot \bar{x}_+ + b) \geq 1$$

$$(\bar{w} \cdot \bar{x}_- + b) \leq -1$$

$$y_i = \begin{cases} 1 & \text{if } x \text{ positive} \\ -1 & \text{otherwise} \end{cases}$$

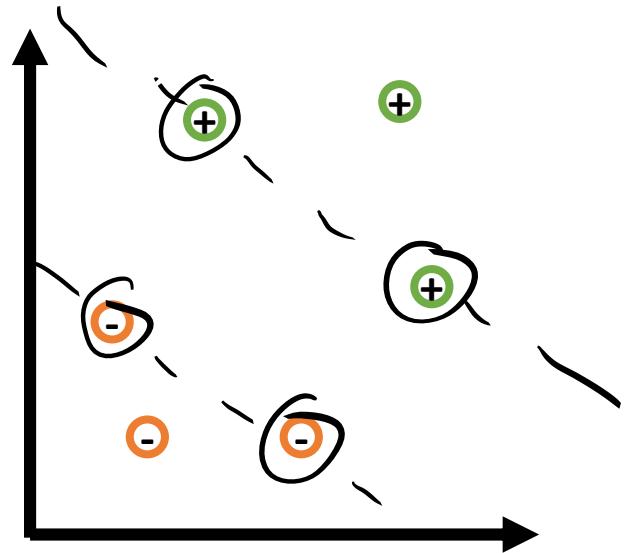
Multiplying with labels y_i gives:

$$y_i(\bar{w} \cdot \bar{x}_+ + b) \geq 1$$

$$y_i(\bar{w} \cdot \bar{x}_- + b) \leq 1$$

now we are going
to ask that for
points on the margin:

$$y_i(\bar{w} \cdot \bar{x}_i + b) - 1 = 0$$



SVM Derivation (Intuition with some slight of hand)

$$y_i(\bar{w} \cdot \bar{x}_i + b) - 1 = 0$$

$$\text{width} = \frac{(x_+ - x_-)}{\|\omega\|} = \frac{2}{\|\omega\|} \quad *$$

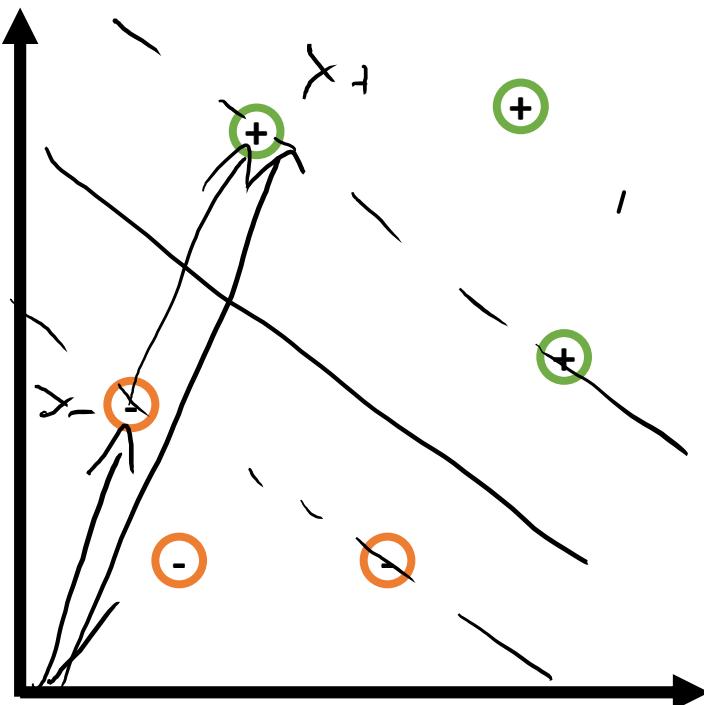
$$\cancel{1+b} + \cancel{1+b} = 2$$

$$x_+: \|\omega \cdot x_+ + b\| - 1 = 0$$

$$\omega \cdot x_+ = 1 - b$$

$$x_-: -1(\omega \cdot x_- + b) - 1 = 0$$

$$\omega \cdot x_- = -1 - b$$



$$\max \frac{2}{\|\omega\|} = \max \frac{1}{\|\omega\|} = \min \|\omega\|$$

$\approx \frac{1}{2} \|\omega\|^2$

primal problem

Now Some Practical Issues (iPython Notebook)

SVM Derivation (Intuition with some slight of hand)

$$y_i(\bar{w} \cdot \bar{x}_i + b) - 1 = 0$$

$$\min \frac{1}{2} \|w\|^2$$

// dual problem

$$L = \frac{1}{2} \|\bar{w}\|^2 - \sum \alpha_i [y_i(\bar{w} \cdot \bar{x}_i + b) - 1]$$

$$\frac{\partial L}{\partial w} = \bar{w} - \sum \alpha_i y_i x_i = 0 \Rightarrow \boxed{\bar{w} = \sum \alpha_i y_i \bar{x}_i}$$

$$\frac{\partial L}{\partial b} = -\sum \alpha_i y_i = 0 \Rightarrow \boxed{\sum \alpha_i y_i = 0}$$

$$L = \frac{1}{2} (\sum \alpha_i y_i \bar{x}_i)(\sum \alpha_j y_j \bar{x}_j) - (\sum \alpha_i y_i \bar{x}_i)(\sum \alpha_j y_j \bar{x}_j) - \underbrace{\sum \alpha_i y_i b}_{= 0} + \sum \alpha_i$$

$$= \sum \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i \cdot x_j$$

SVM Derivation (Intuition with some slight of hand)

$$L = \sum \alpha_i - \frac{1}{2} \sum \sum \underbrace{\alpha_i \alpha_j}_{\text{Output}} \underbrace{y_i y_j}_{\text{coeffs from training data}} \boxed{\bar{x}_i \cdot \bar{x}_j}$$

optimization
only depends
on dot product
of sample
vecs

DECISION RULE:

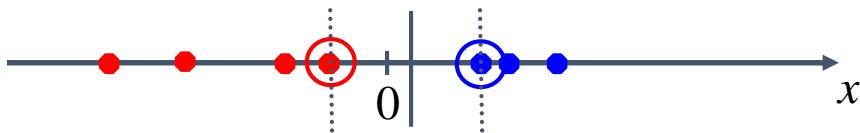
$$\sum \alpha_i y_i \boxed{\bar{x}_i \cdot \bar{u}} + b \geq 0 \quad \text{THEN } +$$

decision also only depends on
dot product

Non-linear SVMs

Non-linear SVMs

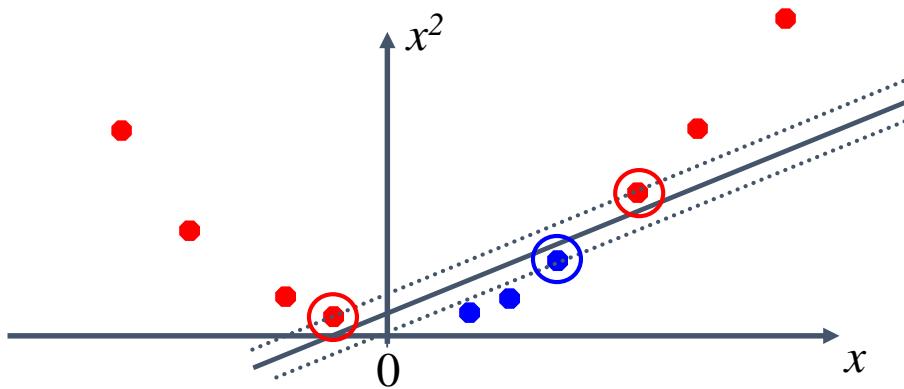
Datasets that are linearly separable (with some noise) work out great:



But what are we going to do if the dataset is just too hard?

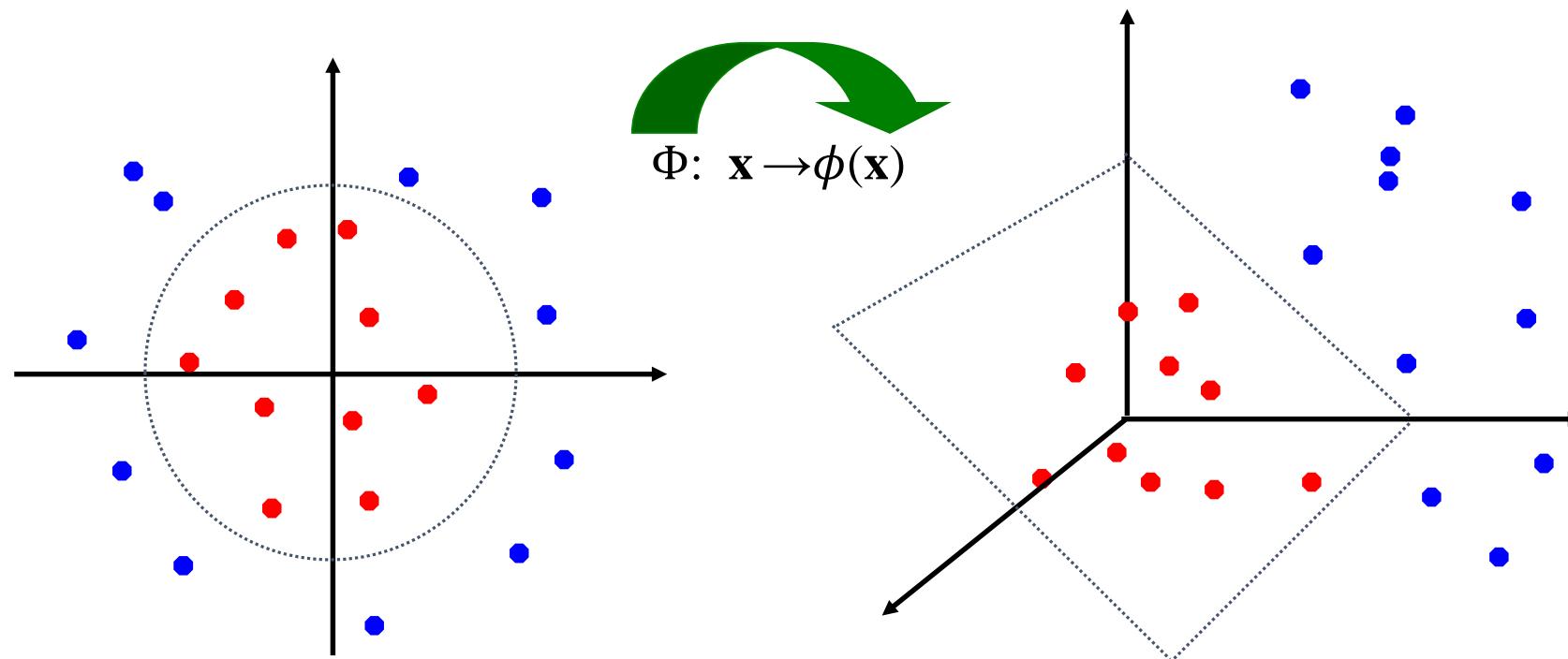


How about ... mapping data to a higher-dimensional space:



Non-linear SVMs: Feature spaces

General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is linearly separable:



Non-linear SVMs: The Kernel Trick

With this mapping, our discriminant function is now:

$$y(x) = w^T \phi(x) + b = \sum_{i \in SV} \alpha_i \boxed{\phi(x_i)^T \phi(x)} + b$$

No need to know mapping function ϕ explicitly, because we only use the **dot product** of feature vectors in both training and test.

A *kernel function* is defined as a function that corresponds to a inner product of two feature vectors in some expanded feature space:

$$K(\mathbf{x}_i, \mathbf{x}_j) \equiv \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

Non-linear SVMs: The Kernel Trick

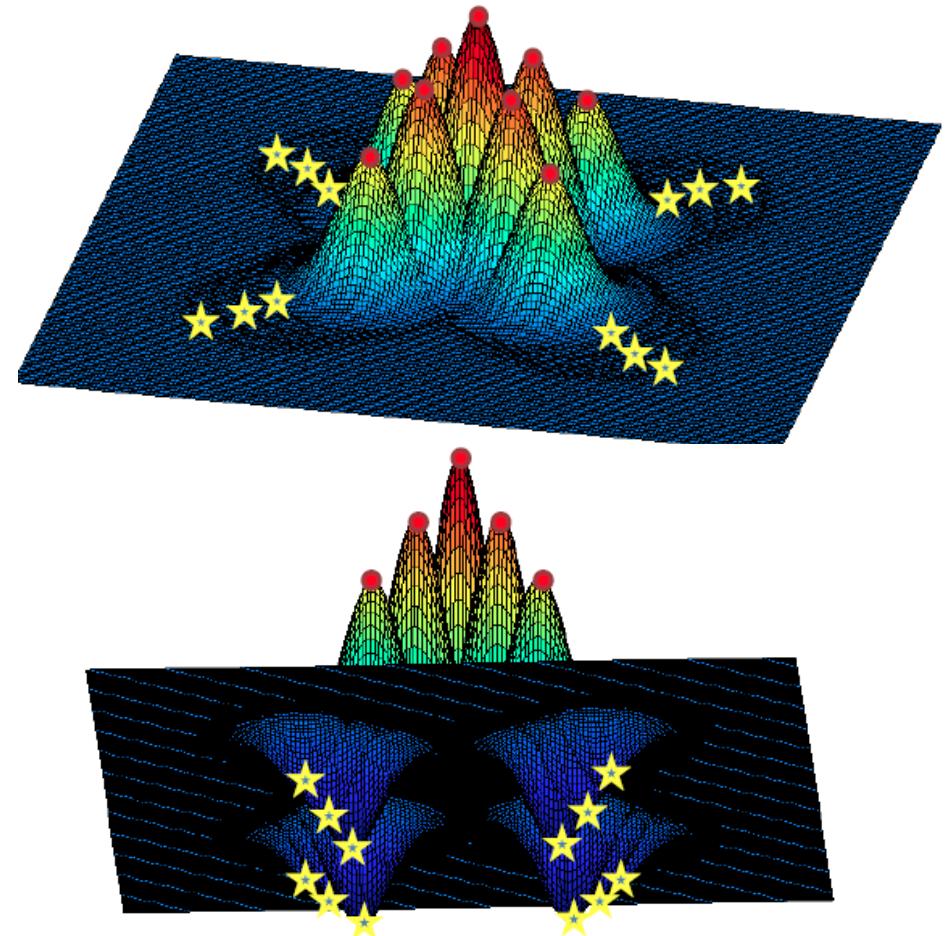
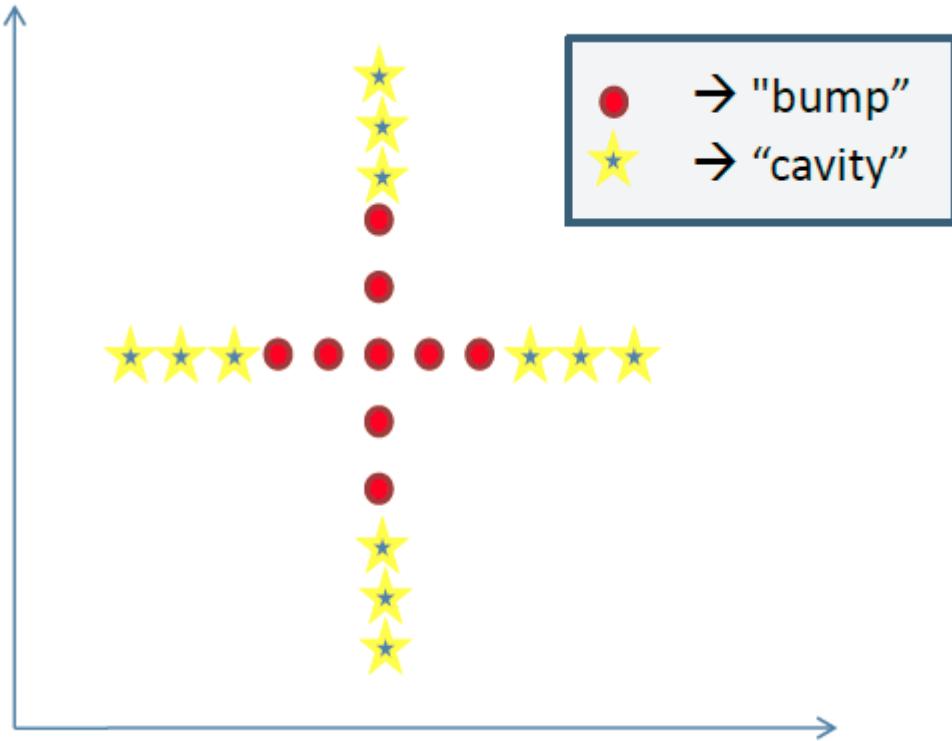
Commonly used kernel functions:

- Linear kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- Polynomial kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$
- Gaussian (Radial-Basis Function Kernel (RBF)) kernel:
$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$
- Sigmoid: $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta_0 \mathbf{x}_i^T + \beta_1)$

In general, functions that satisfy *Mercer's condition* can be kernel functions.

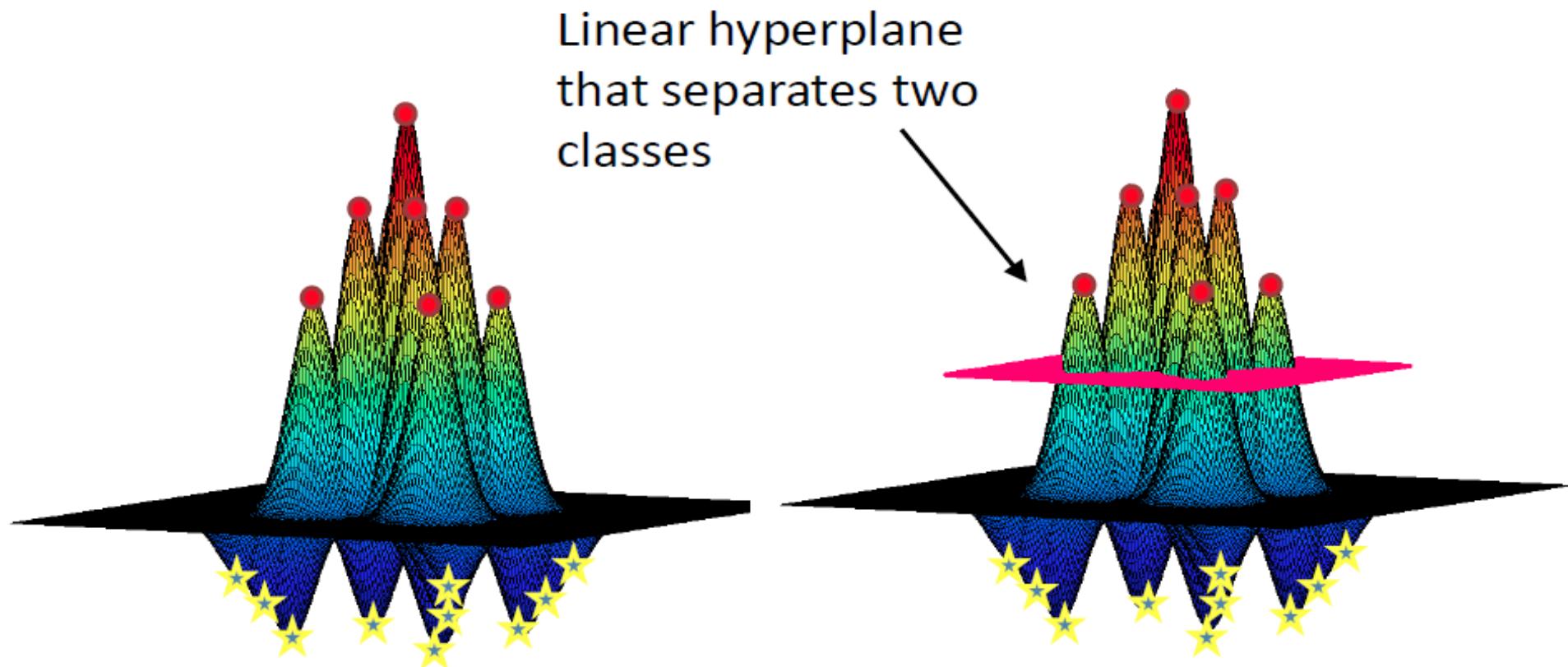
Example: Gaussian Kernel

Geometrically, this is a “bump” or “cavity” centered at the training data point x :



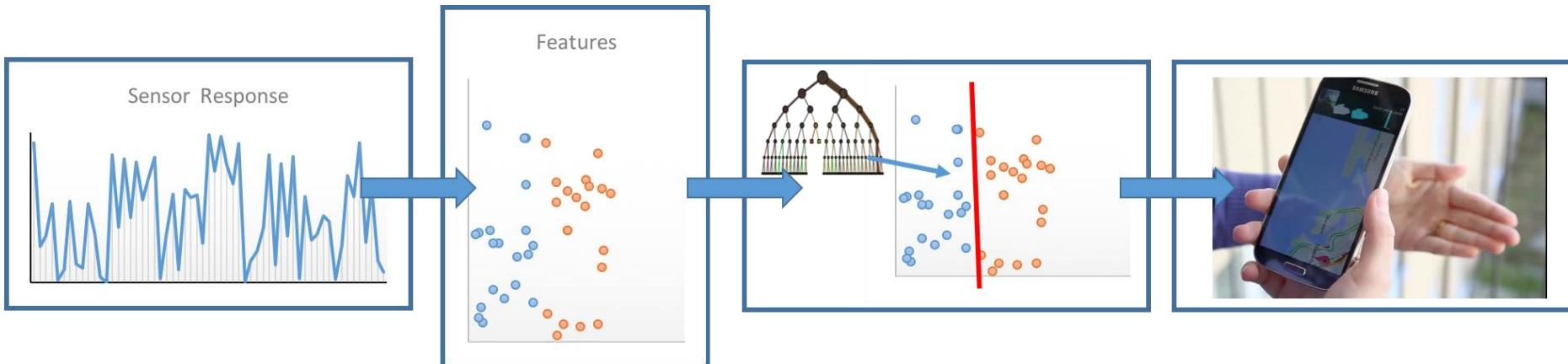
Example: Gaussian Kernel - Classification

Geometrically, this is a “bump” or “cavity” centered at the training data point x :



General Approach to supervised learning

- 1) Collect and label data.
- 2) Extract features.
- 3) Load the dataset.
- 4) Summarize the dataset.
- 5) Visualize the dataset.
- 6) Evaluate some ML algorithms.
- 7) Make predictions.
- 8) Integrate best performing model into your pipeline



Typical Issues in Computational Input Recognition

- Data is sometimes extremely **noisy**
- A lot of **variance** between subjects
- Collecting **data** is **expensive**
- Need models that can be evaluated at **interactive rates**
- **Feature engineering** can be an “art”
- *Spotting* versus *recognition*

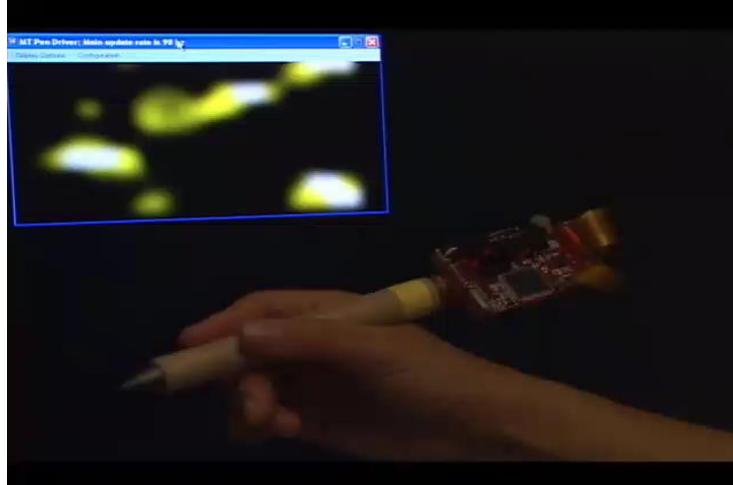
Outlook

Traditional Machine Learning Pipeline

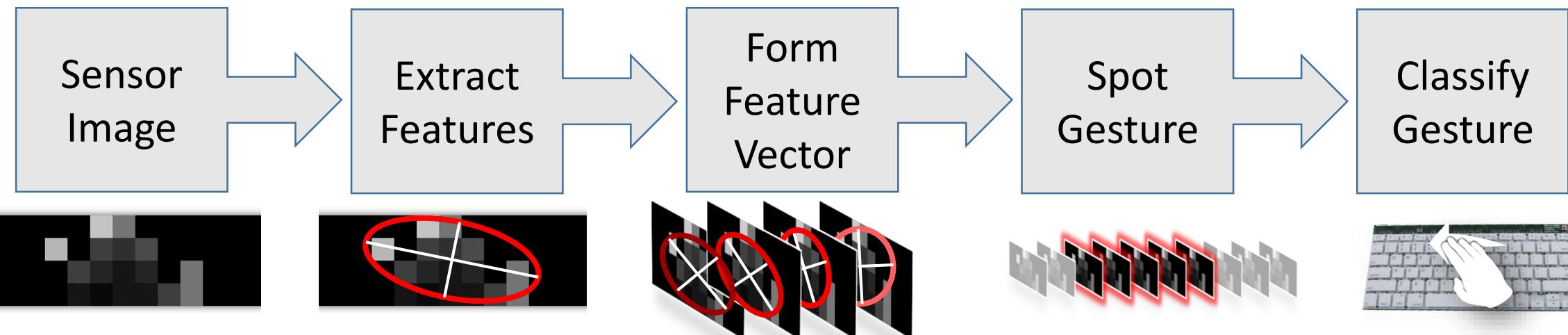
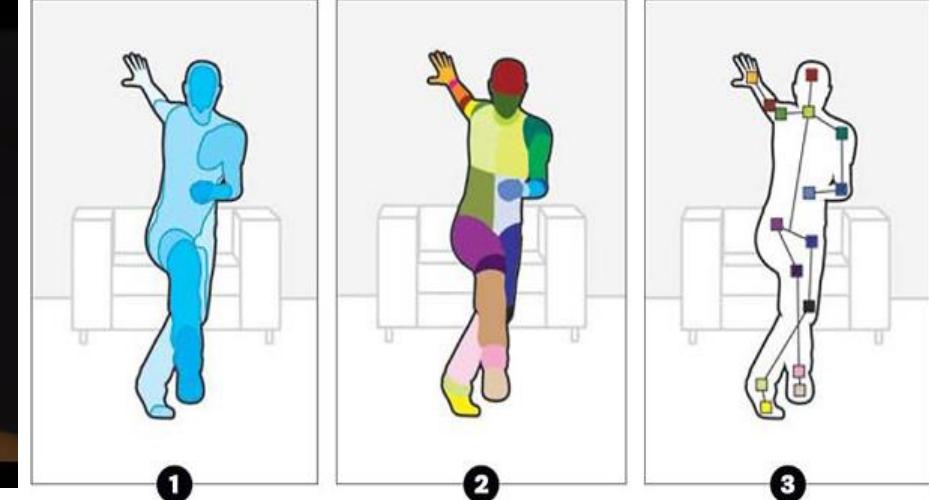
[Zhang *et al.*, CHI'14]



[Song *et al.*, CHI'11]



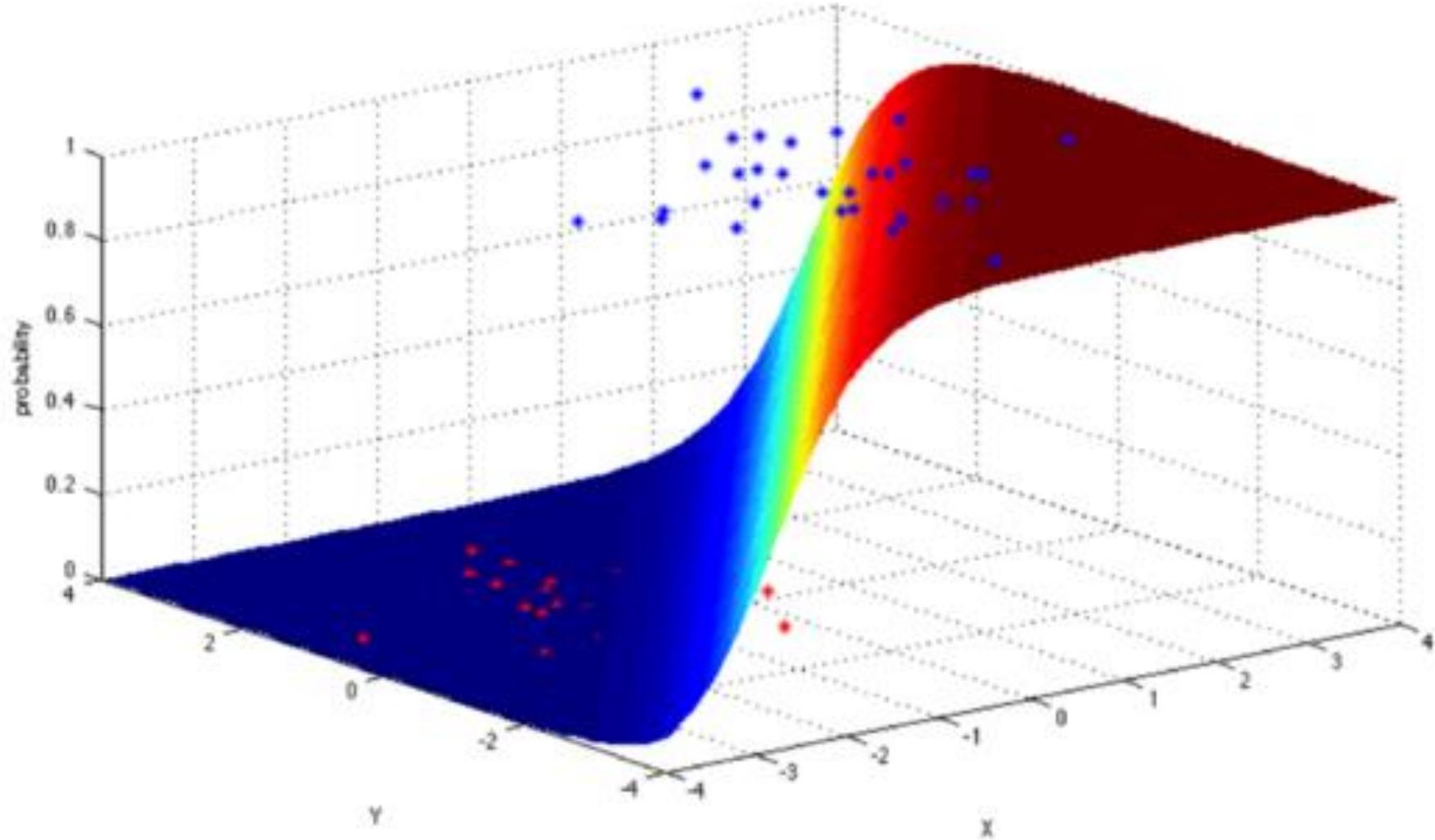
[Nowozin *et al.*, MSR TR '11]



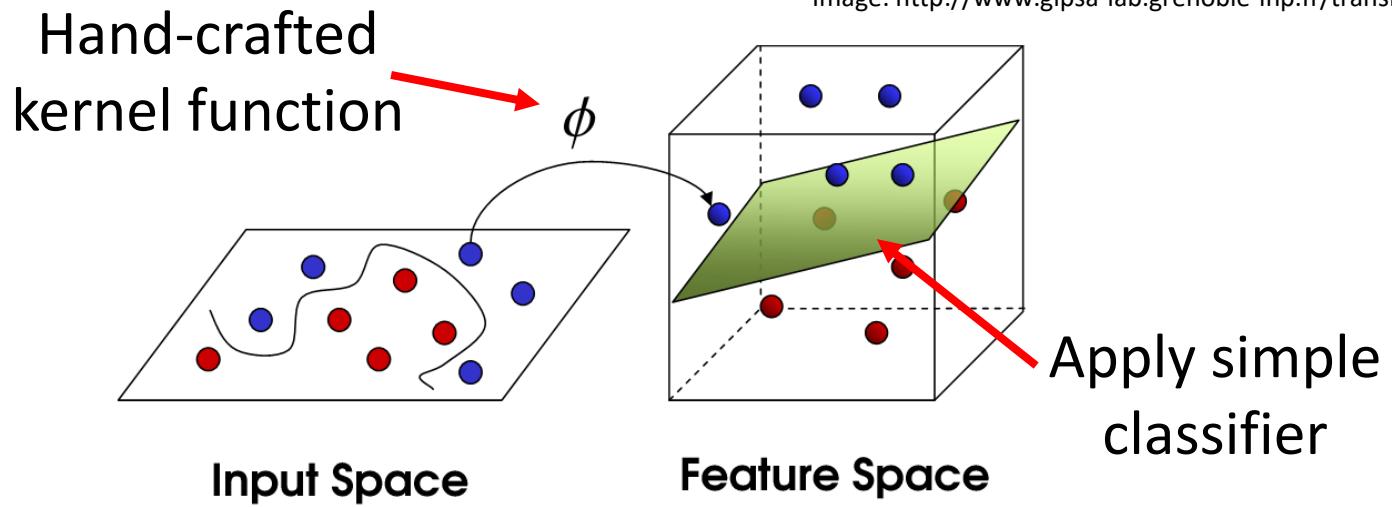
Feature (parameter) learning: Random Forests

COARSE
DEPTH ESTIMATION

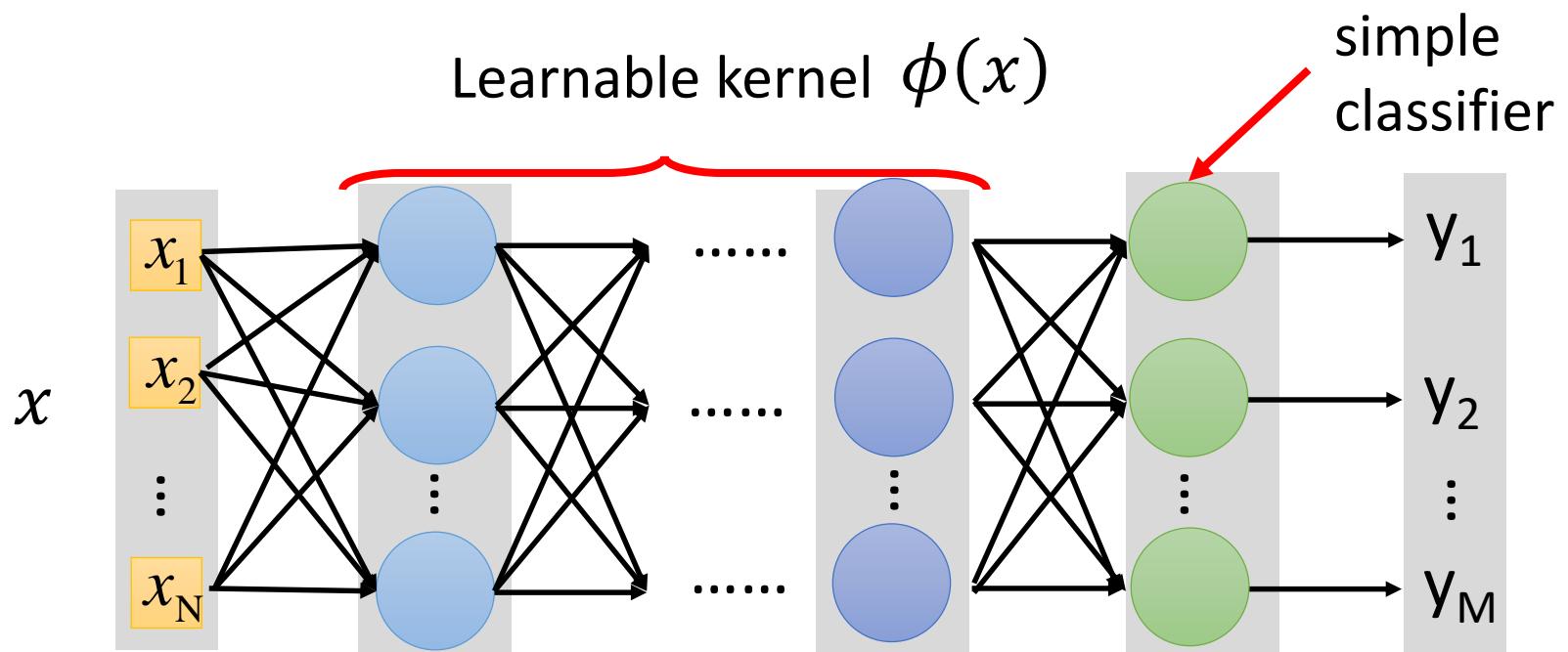
Representation learning



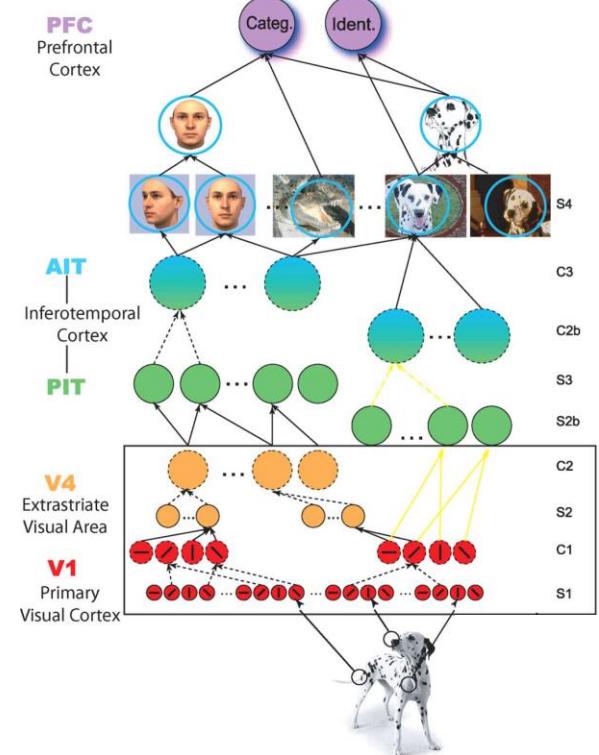
Traditional ML



Deep Learning

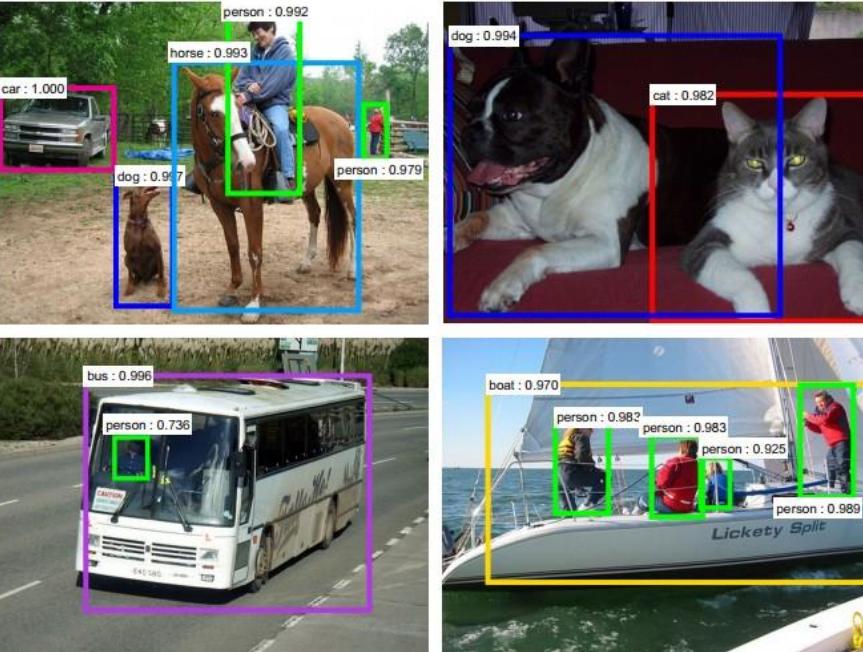


Deep-learning



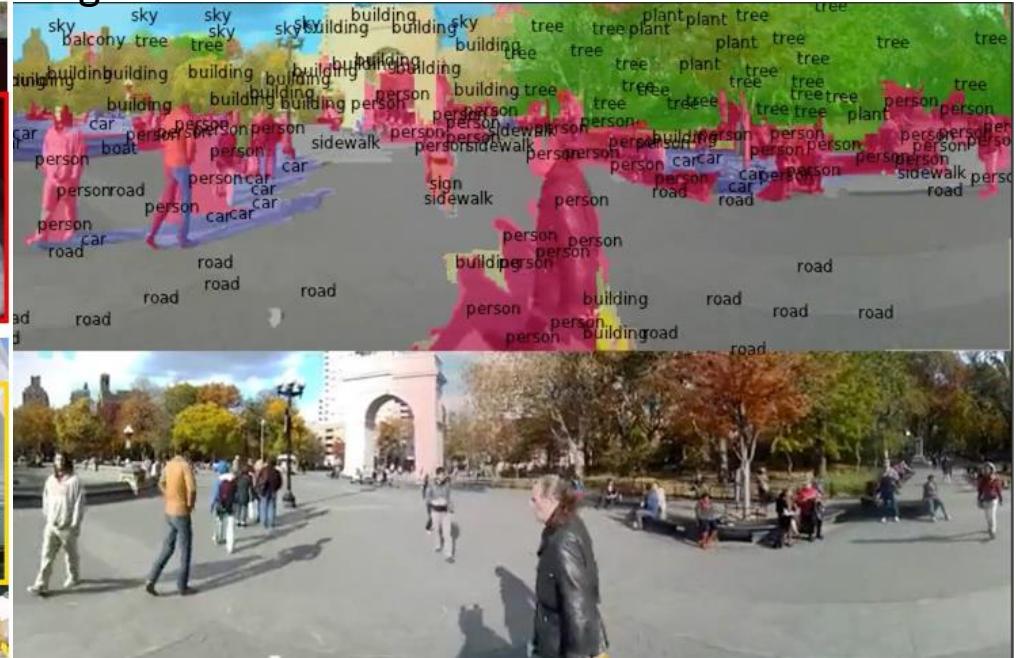
[Serre et al. PNAS 2007]

Detection



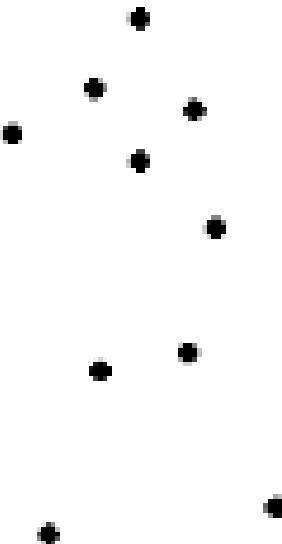
[Faster R-CNN: Ren, He, Girshick, Sun 2015]

Segmentation



[Farabet et al., 2012]

What about motion?



G. Johansson, "Visual Perception of Biological Motion and a Model For Its Analysis", *Perception and Psychophysics* 14, 201-211, 1973.