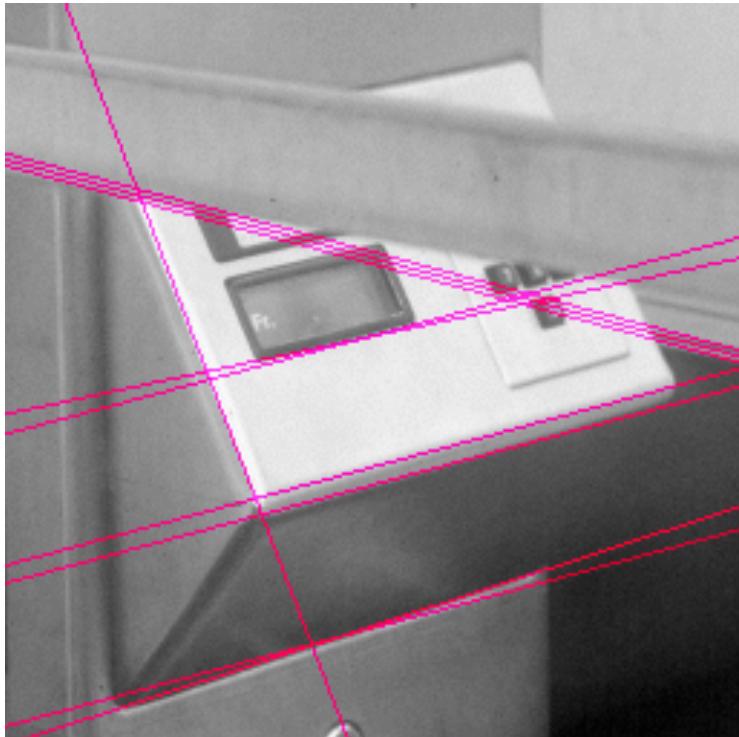


Classification



Dan Witzner Hansen

Today

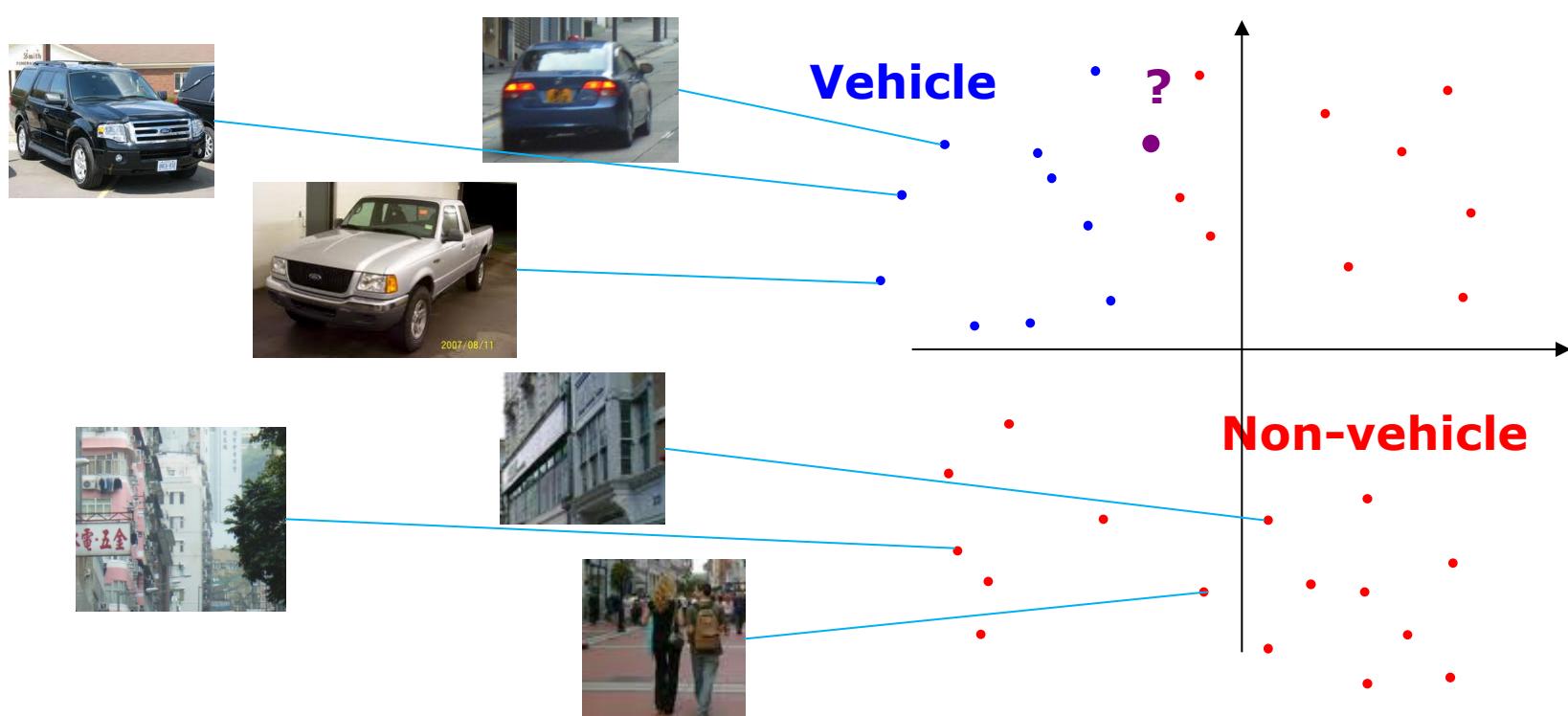
- Linear Classification (binary)
- HOG feature
- Key Concepts:
 - Classification as regression
 - Decision boundary
 - Loss function
 - Metrics to evaluate classification

SUMMARY OF VIDEO LECTURE

Image Classification

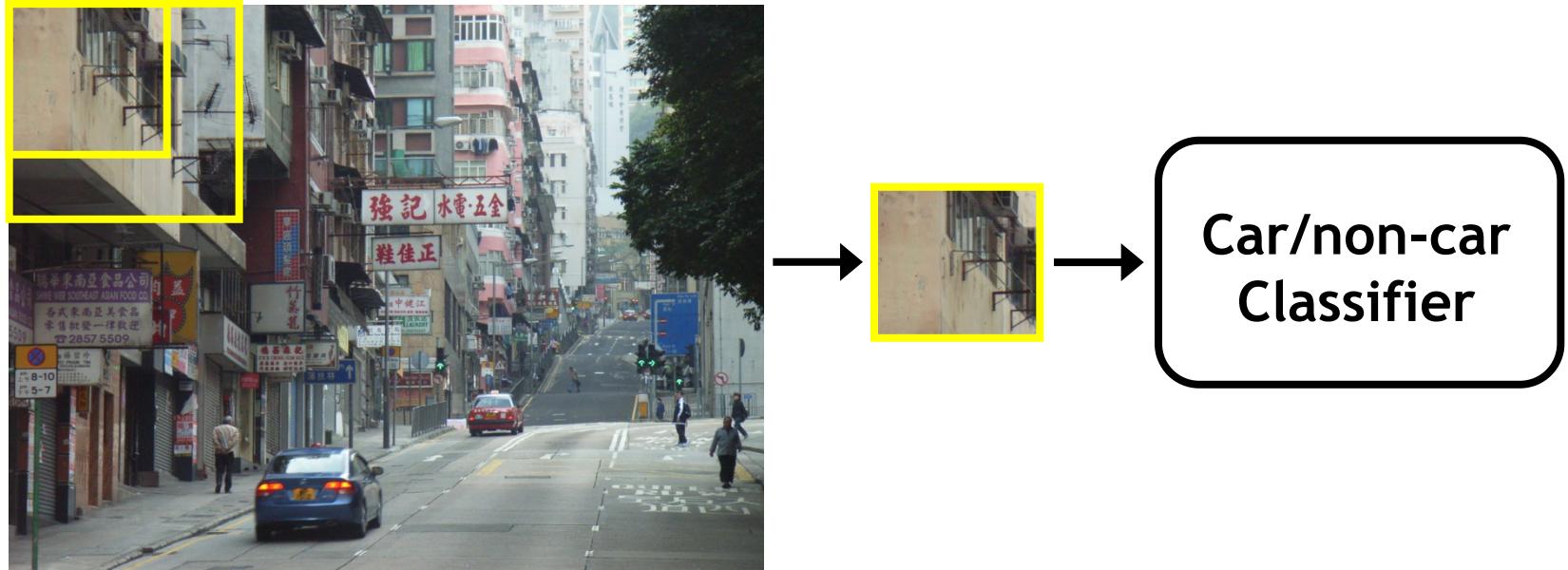
Supervised learning

- Provide training images to learn what are vehicles and what are not.



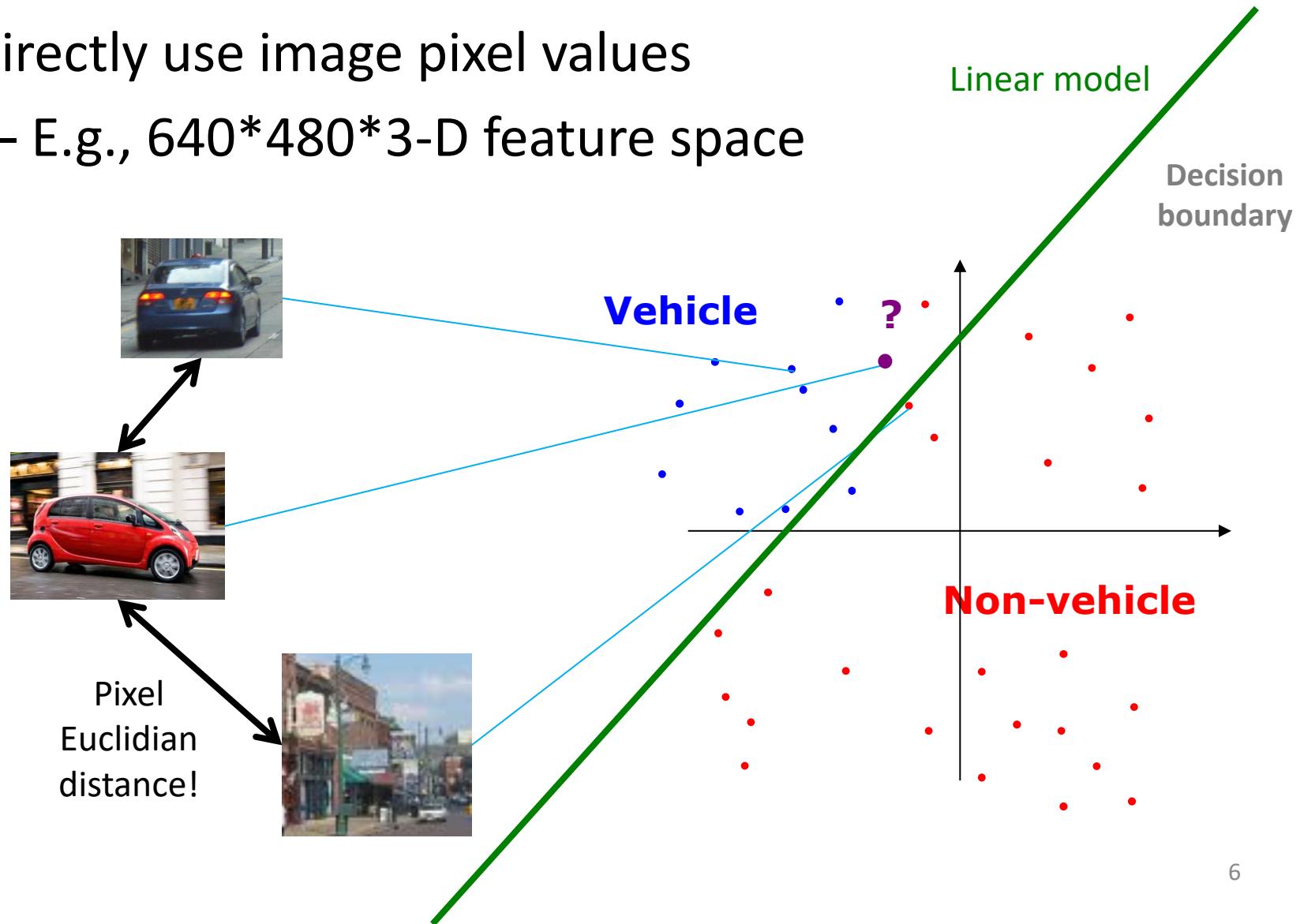
Sliding windows Object detection from object classification

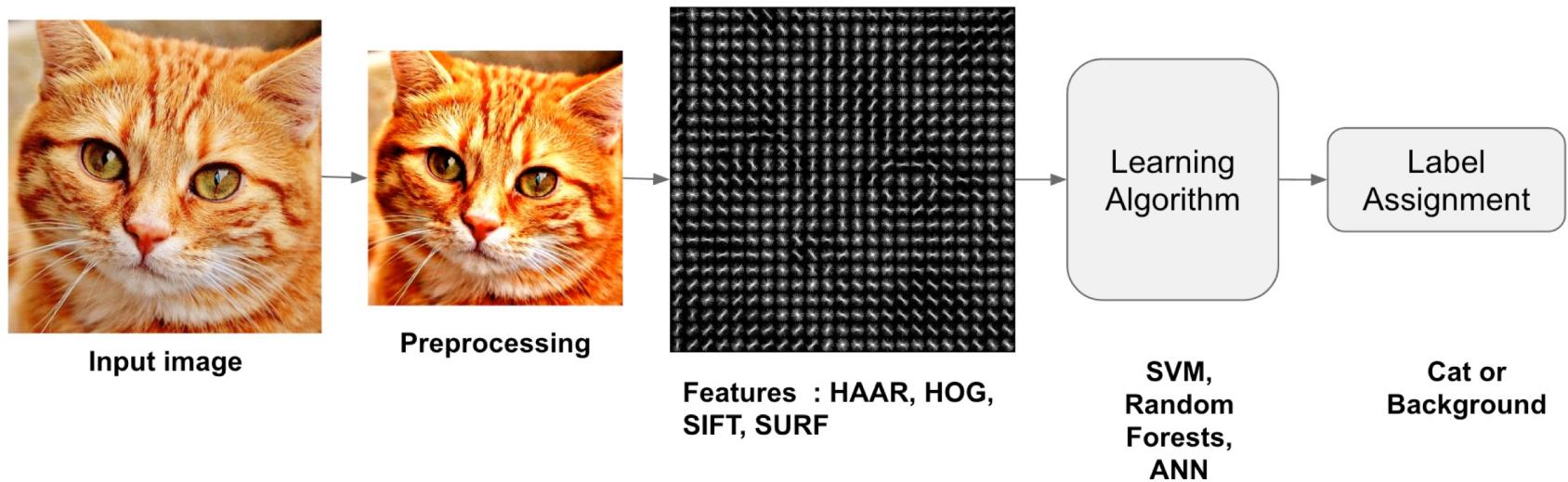
- Test all possible image segments (also different sizes)



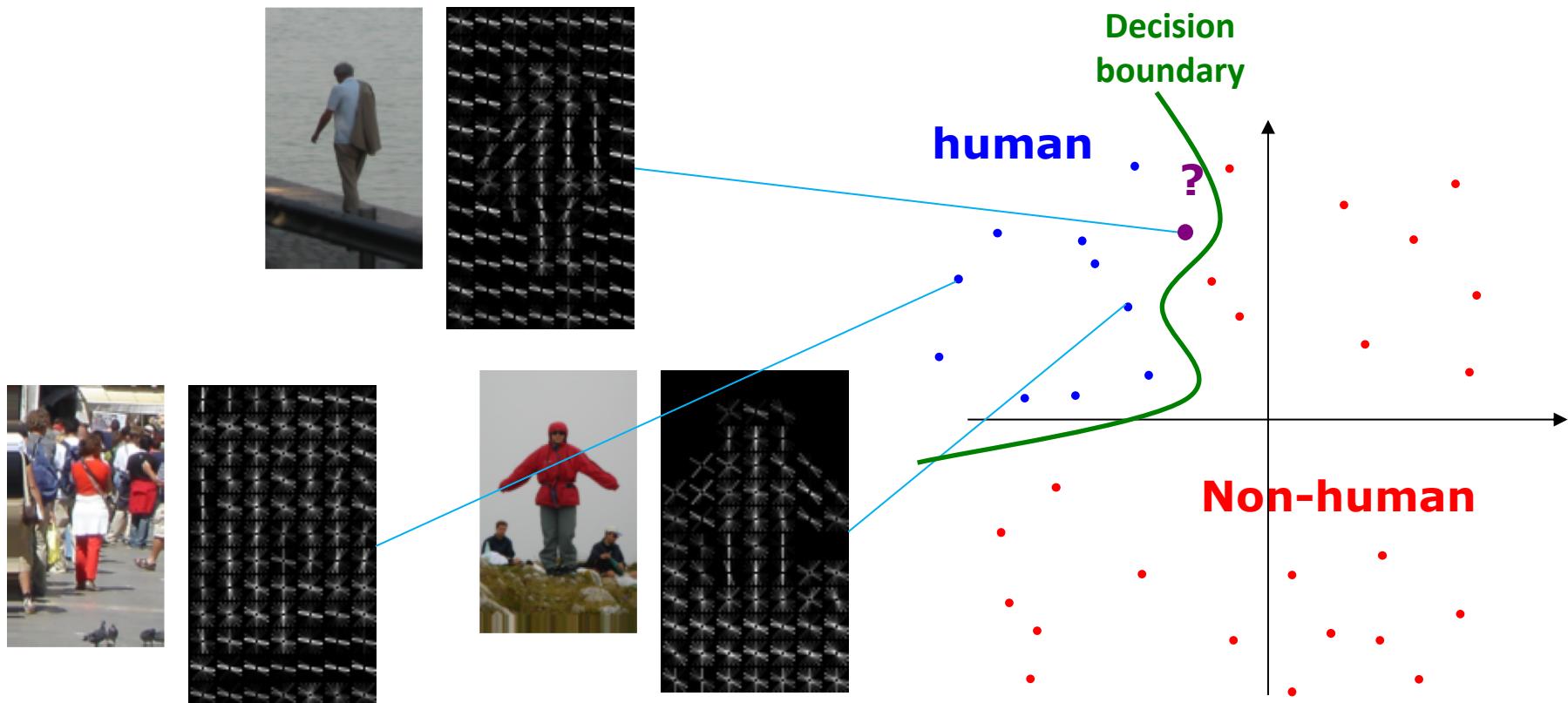
Simple solution – image pixels

- Directly use image pixel values
 - E.g., $640 \times 480 \times 3$ -D feature space



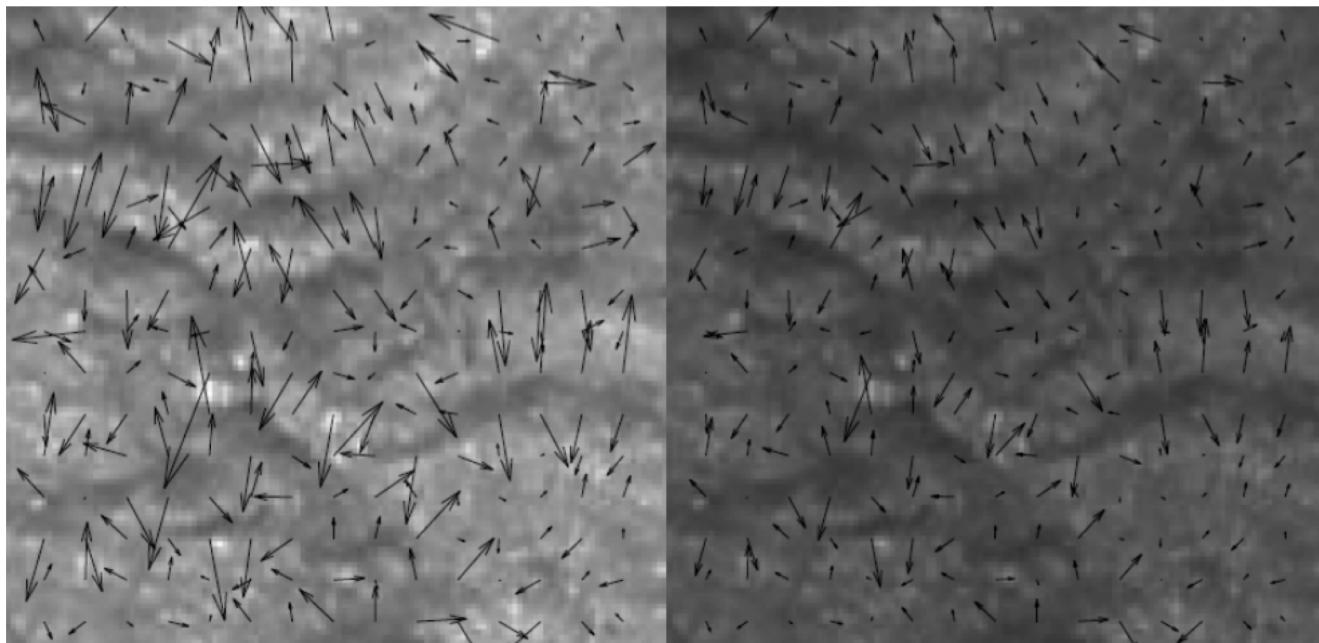


Extracting features and Non-linear descision boundaries



Describe image patches by orientation

- Gradient Magnitude is affected by illumination changes – orientation is not



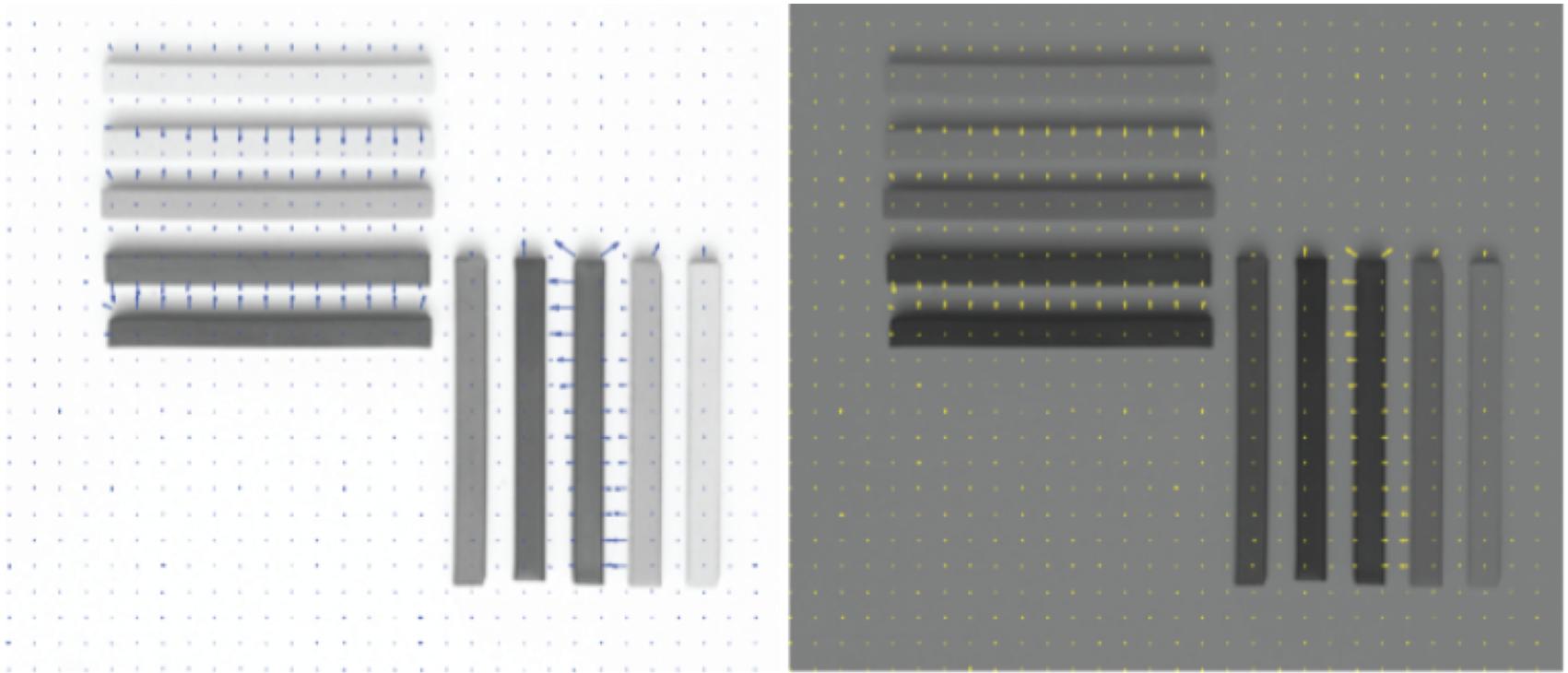
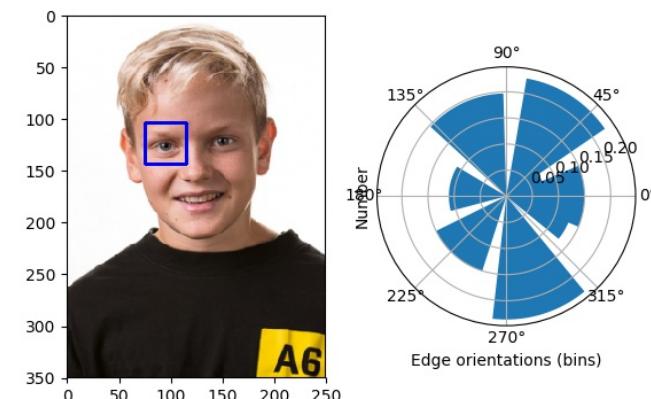
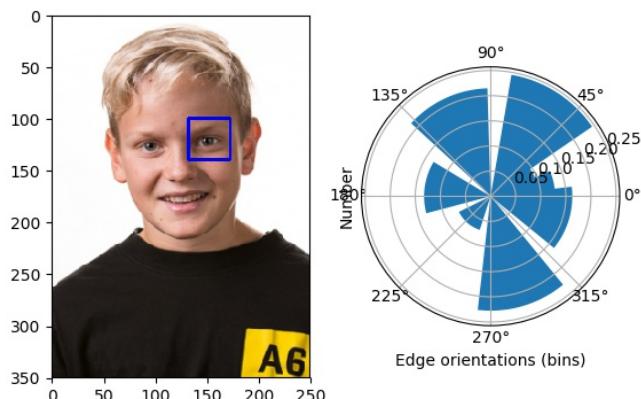
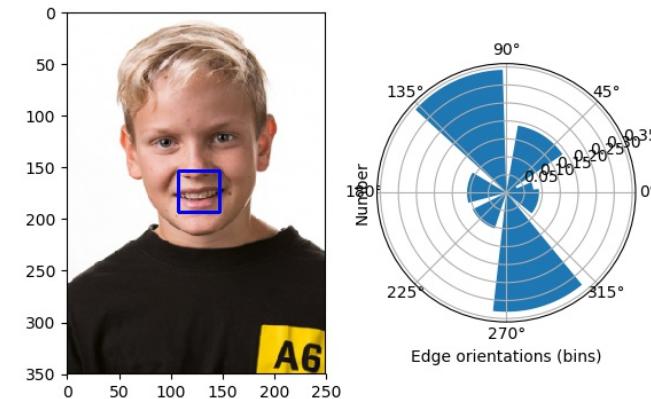
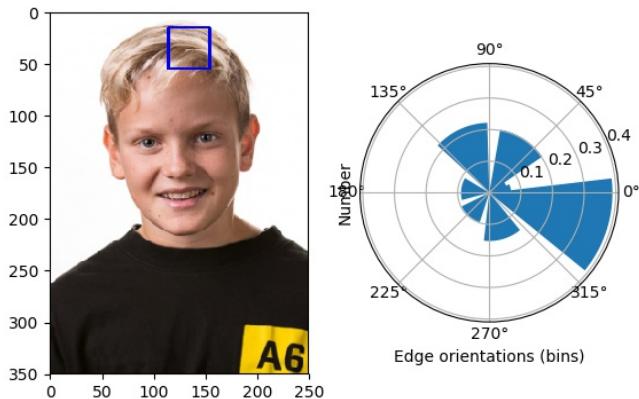
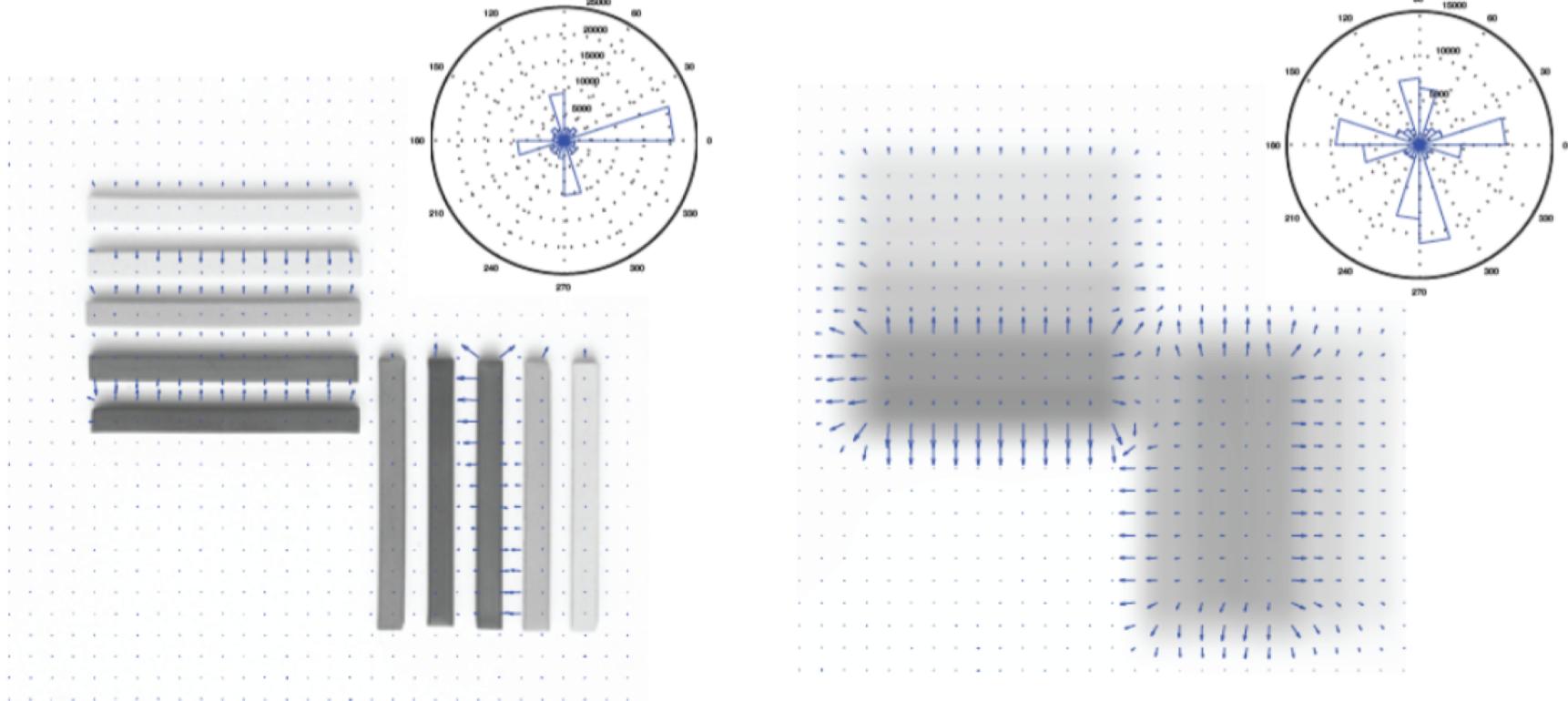


FIGURE 5.7: The magnitude of the image gradient changes when one increases or decreases the intensity. The orientation of the image gradient does not change; we have plotted every 10th orientation arrow, to make the figure easier to read. Note how the directions of the

Comparing orientation histograms

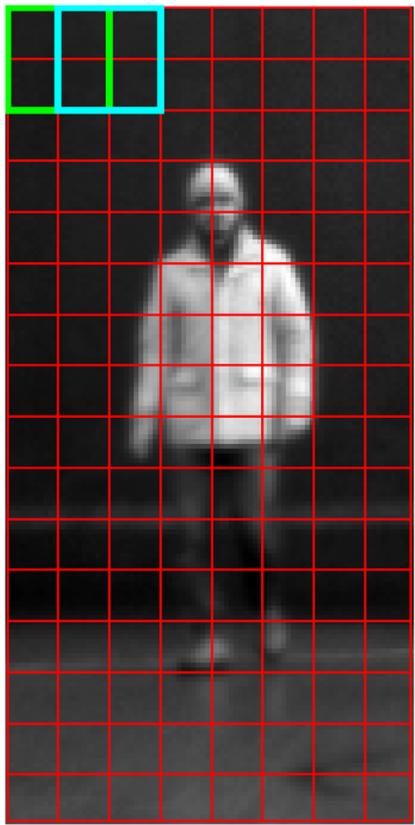


Orientation at different scales

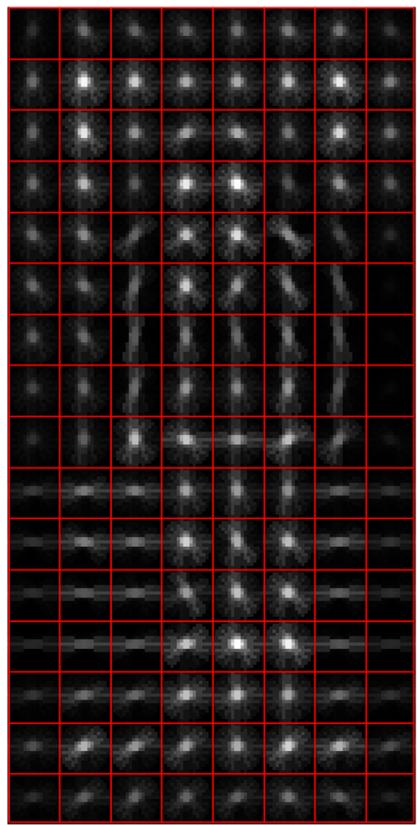


Idea of HOG features

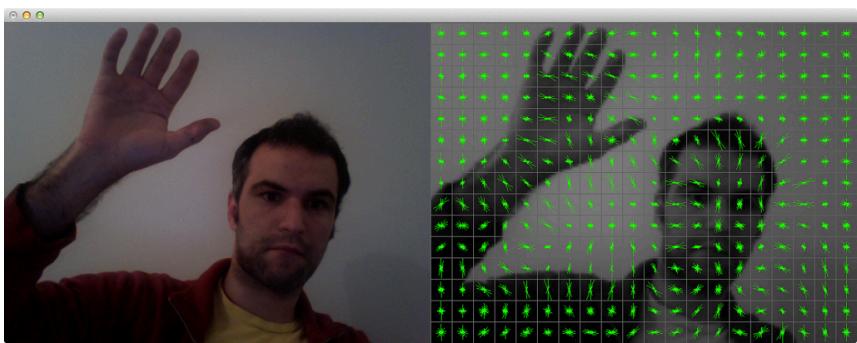
- Divide image into small sub-images: “cells”
- Cells can be rectangular (R-HOG) or circular (c-HOG)
- Accumulate a histogram of edge orientations within that cell
- The combined histogram entries are used as the feature describing the object
- To provide better illumination invariance (light, shadows etc) normalize the cells across larger regions incorporating multiple cells: “blocks”



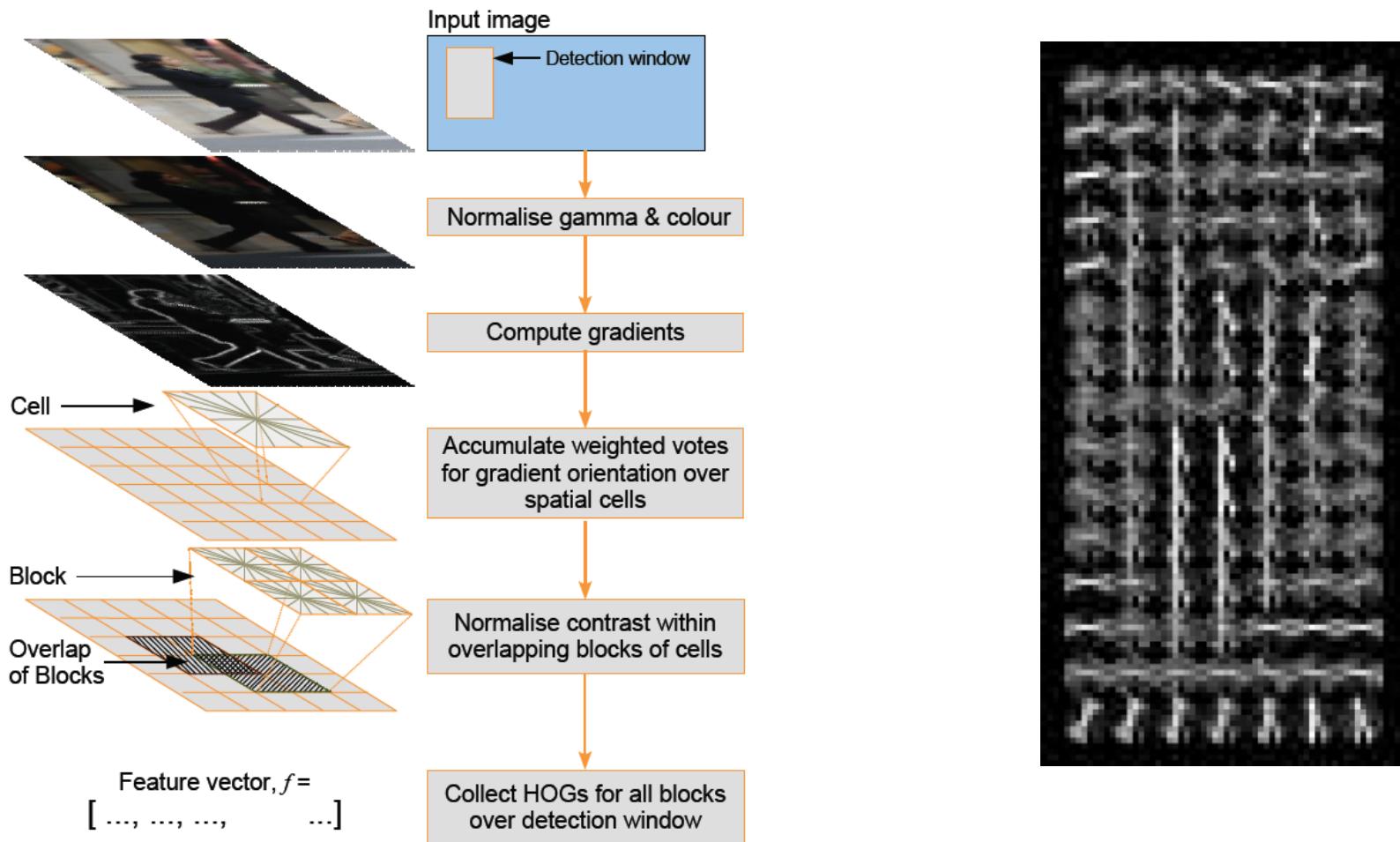
(a)



(b)

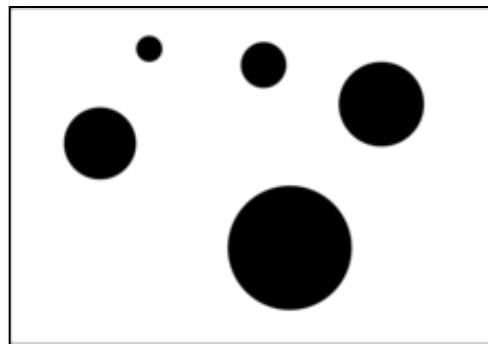


HOG

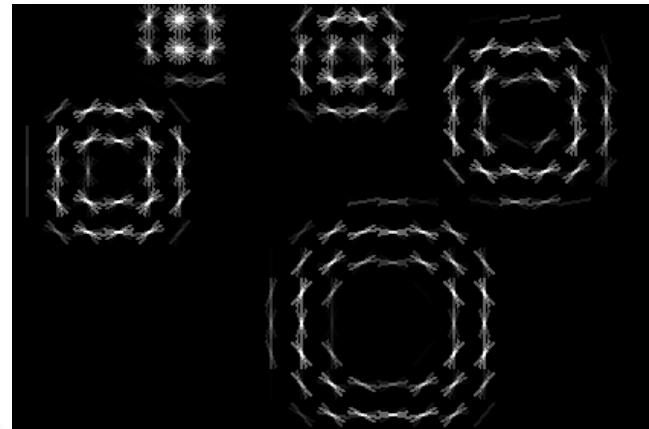


*Navneet Dalal and Bill Triggs. Histograms of Oriented Gradients for Human Detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, SanDiego, USA, June 2005. Vol. II, pp. 886-893.

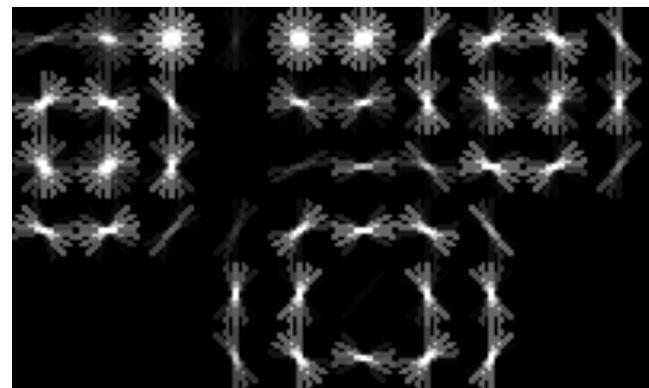
Histogram of Oriented Gradients (HoG)



HoGify

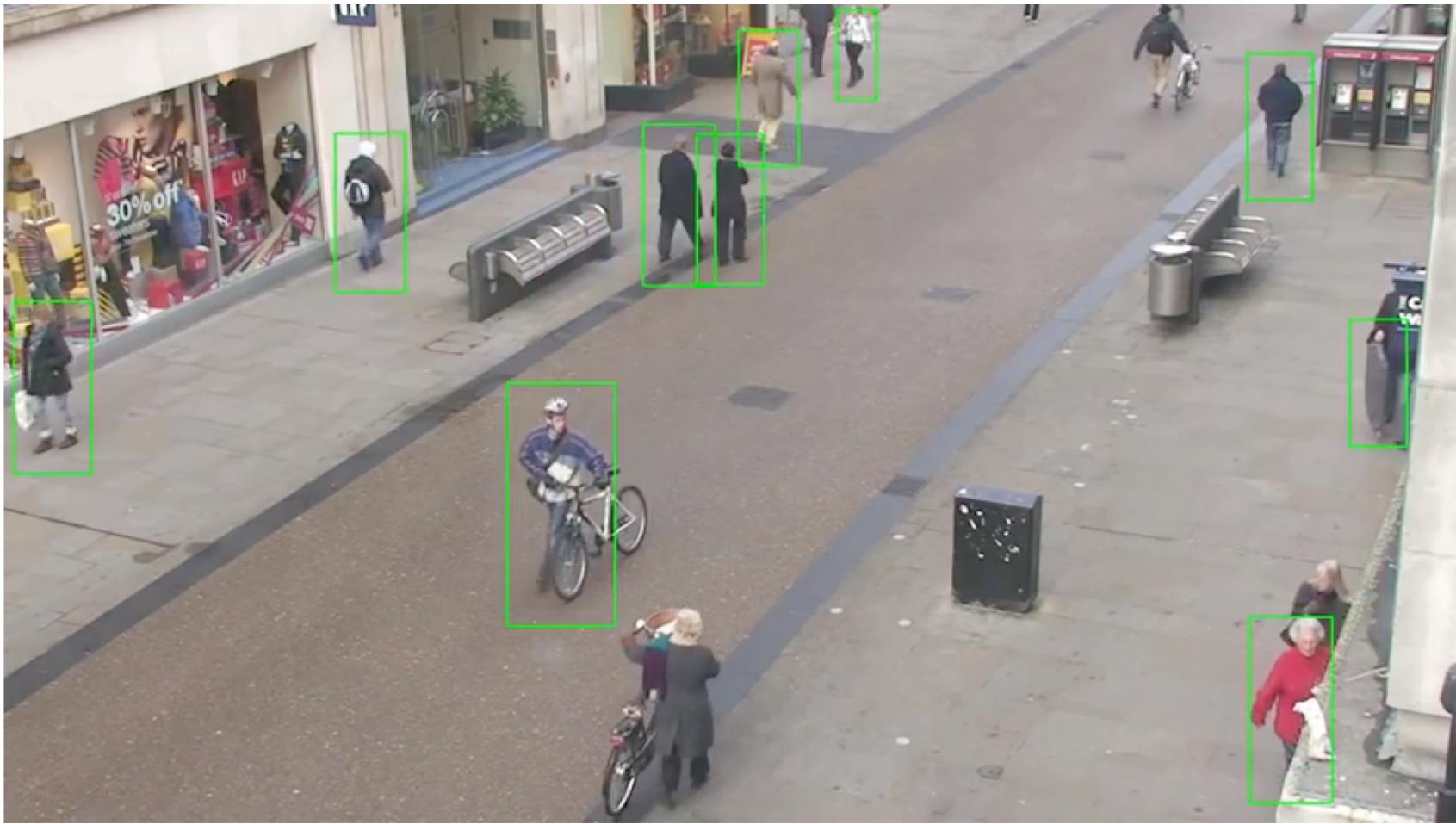


10x10 cells



20x20 cells

Results



Classification vs regression

- We are interested in mapping the input $\mathbf{x} \in \mathcal{X}$ to a *label* $t \in \mathcal{Y}$
- In regression typically $\mathcal{Y} = \mathbb{R}$
- Now \mathcal{Y} is categorical

Classification as regression

- Can we do this task using what we have learned in previous lectures?
- Simple hack: Ignore that the output is categorical!
- Suppose we have a binary problem, $t \in \{-1, 1\}$
- Assuming the standard model used for regression

$$y = f(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

- How can we obtain \mathbf{w} ?
- Use least squares, $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$. How is \mathbf{X} computed? and \mathbf{t} ?
- Which loss are we minimizing? Does it make sense?

$$\ell_{square}(\mathbf{w}, t) = \frac{1}{N} \sum_{n=1}^N (t^{(n)} - \mathbf{w}^T \mathbf{x}^{(n)})^2$$

- How do I compute a label for a new example? Let's see an example

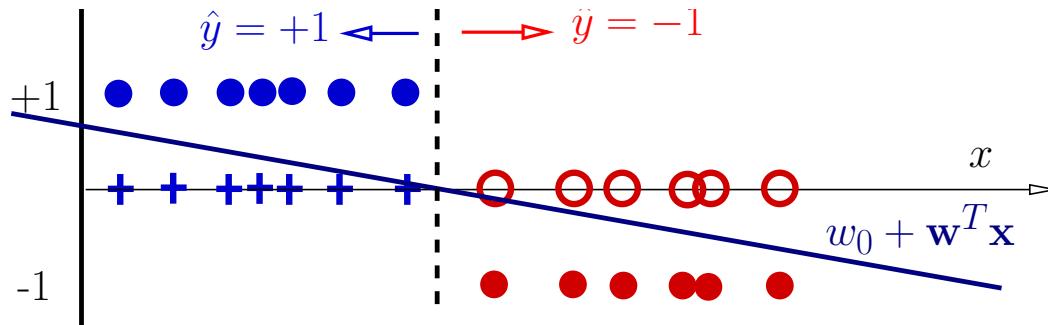
Classification as regression

- One dimensional example (input x is 1-dim)



- The colors indicate labels (a blue plus denotes that $t^{(i)}$ is from class -1 , red circle that $t^{(i)}$ is class 1)

Decision Rule



- How can I mathematically write this rule?

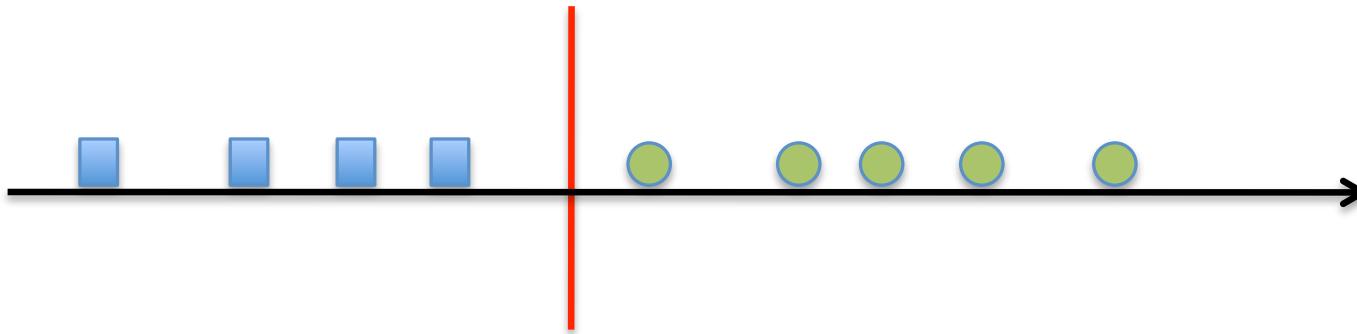
$$y = \text{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

- This specifies a **linear classifier**: it has a **linear boundary (hyperplane)**

$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"

Example in 1D



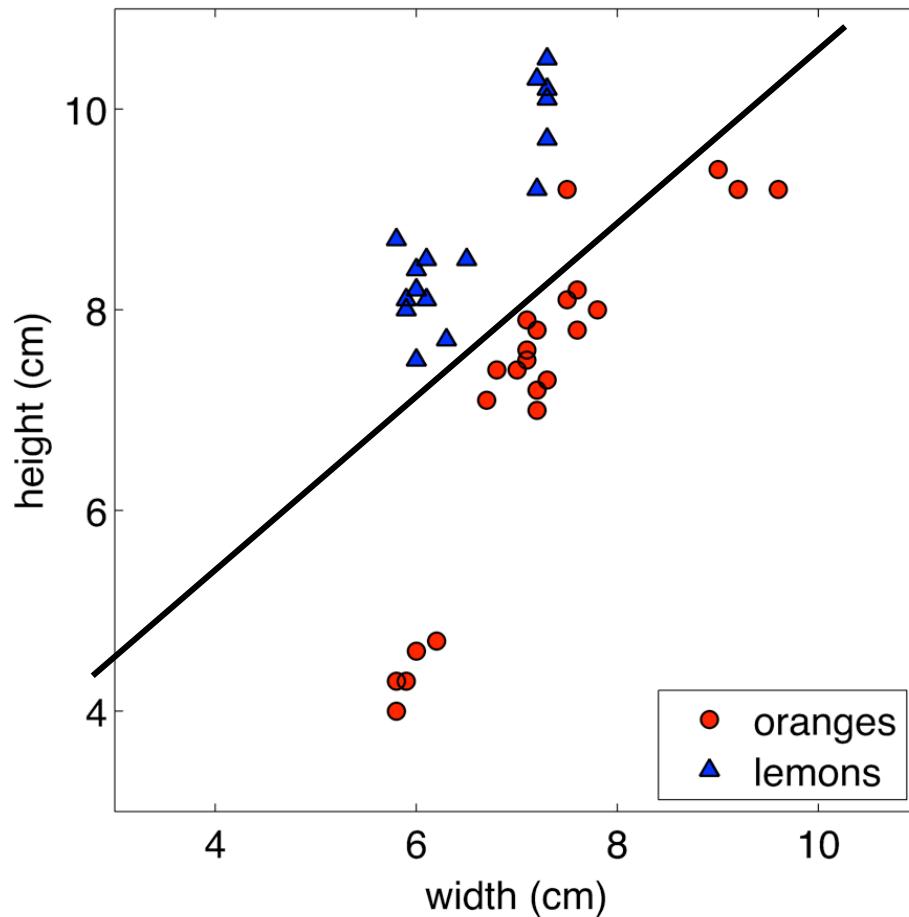
- The linear classifier has a linear boundary (hyperplane)

$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

which separates the space into two "half-spaces"

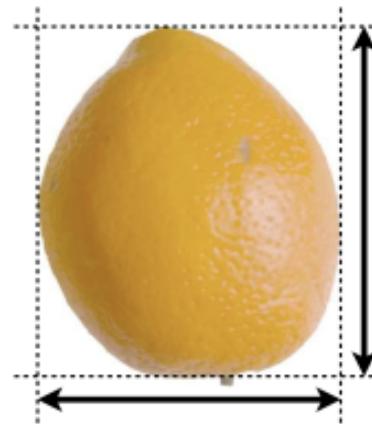
- In 1D this is simply a threshold

Classification Lemons and oranges

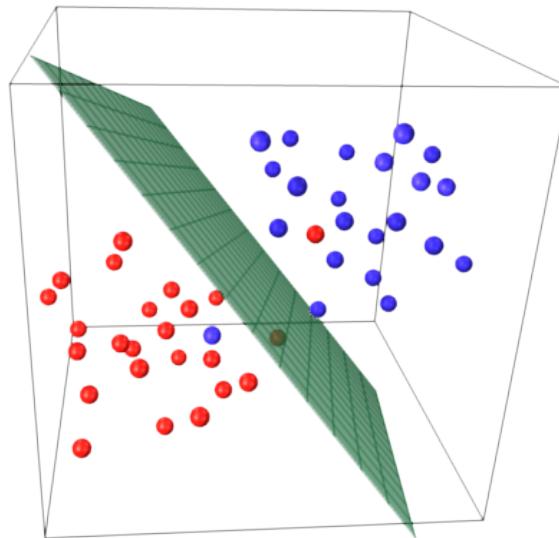


Can construct simple linear decision boundary:

$$y = \text{sign}(w_0 + w_1 x_1 + w_2 x_2)$$



Example in 3D



- The linear classifier has a linear boundary (hyperplane)

$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

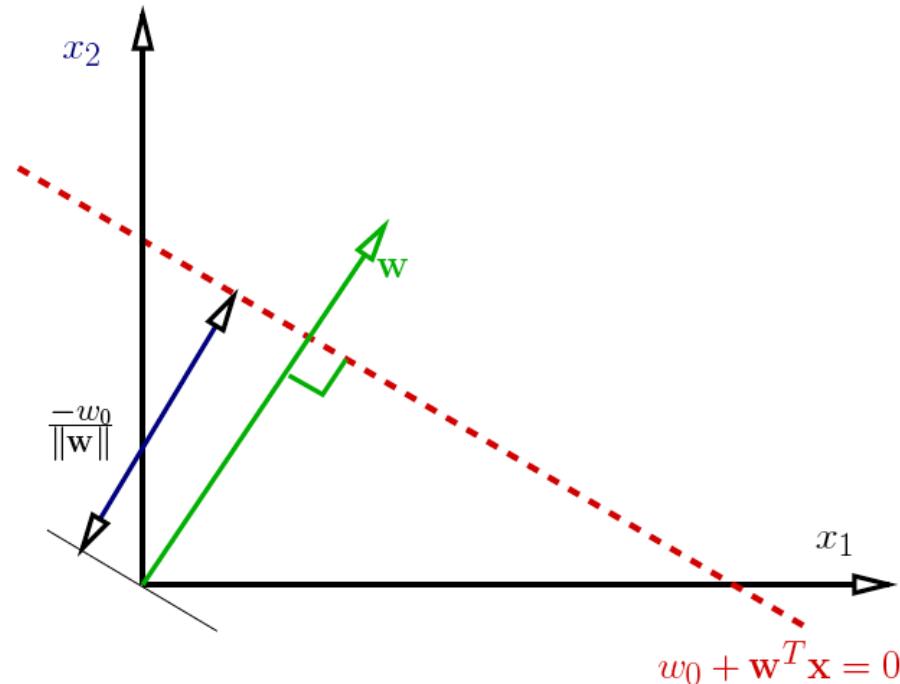
which separates the space into two "half-spaces"

- In 3D this is a plane
- What about higher-dimensional spaces?

Sometimes written in
homogeneous form
 $[w w_0] \begin{bmatrix} x \\ 1 \end{bmatrix} = 0$

Geometry

$\mathbf{w}^T \mathbf{x} = 0$ a line passing though the origin and orthogonal to \mathbf{w}
 $\mathbf{w}^T \mathbf{x} + w_0 = 0$ shifts it by w_0



Loss Functions

- Classifying using a linear decision boundary reduces the data dimension to 1

$$y(\mathbf{x}) = \text{sign}(w_0 + \mathbf{w}^T \mathbf{x})$$

- What is the cost of being wrong?
- **Loss function:** $L(y, t)$ is the loss incurred for predicting y when correct answer is t
- For medical diagnosis: For a diabetes screening test is it better to have false positives or false negatives?
- For movie ratings: The "truth" is that Alice thinks E.T. is worthy of a 4. How bad is it to predict a 5? How about a 2?

Other Loss Functions

- Zero/one loss for a classifier

$$L_{0-1}(y(\mathbf{x}), t) = \begin{cases} 0 & \text{if } y(\mathbf{x}) = t \\ 1 & \text{if } y(\mathbf{x}) \neq t \end{cases}$$

- Asymmetric Binary Loss

$$L_{ABL}(y(\mathbf{x}), t) = \begin{cases} \alpha & \text{if } y(\mathbf{x}) = 1 \wedge t = 0 \\ \beta & \text{if } y(\mathbf{x}) = 0 \wedge t = 1 \\ 0 & \text{if } y(\mathbf{x}) = t \end{cases}$$

- Squared (quadratic) loss

$$L_{\text{squared}}(y(\mathbf{x}), t) = (t - y(\mathbf{x}))^2$$

- Absolute Error

$$L_{\text{absolute}}(y(\mathbf{x}), t) = |t - y(\mathbf{x})|$$

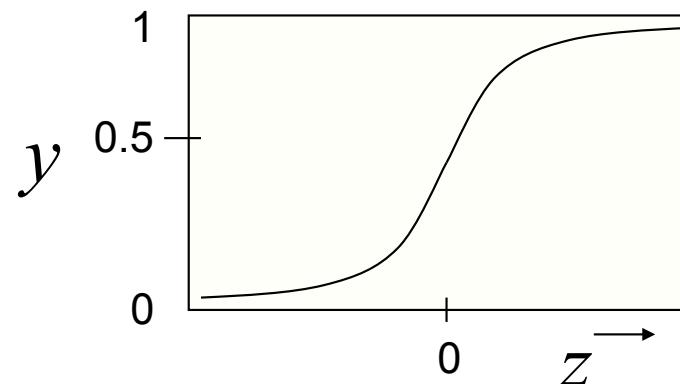
Logistic Regression

- An alternative: replace the $\text{sign}(\cdot)$ with the **sigmoid** or **logistic function**
- We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

where the sigmoid is defined as

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



- The output is a smooth function of the inputs and the weights. It can be seen as a smoothed and differentiable alternative to $\text{sign}(\cdot)$

Logistic regression

- We assumed a particular functional form: sigmoid applied to a linear function of the data

$$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

where the sigmoid is defined as

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

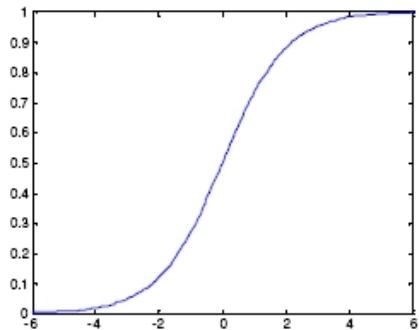
- ▶ One parameter per data dimension (feature) and the bias
- ▶ Features can be discrete or continuous
- ▶ Output of the model: value $y \in [0, 1]$
- ▶ Allows for gradient-based learning of the parameters

Shape of the logistic function

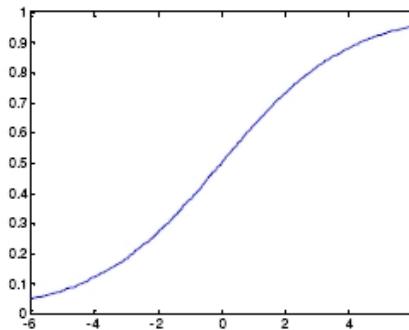
- Let's look at how modifying w changes the shape of the function
- 1D example:

$$y = \sigma(w_1 x + w_0)$$

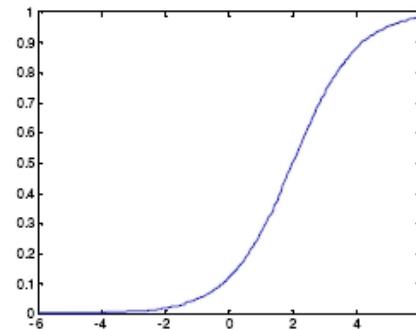
$$w_0 = 0, w_1 = 1$$



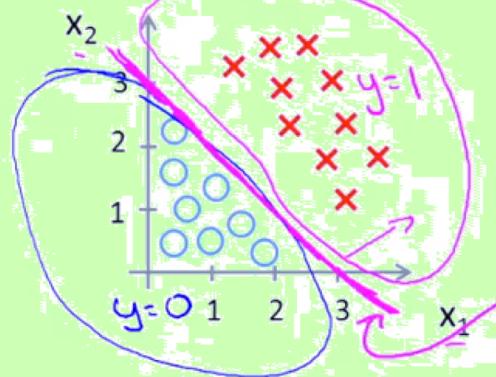
$$w_0 = 0, w_1 = 0.5$$



$$w_0 = -2, w_1 = 1$$



Decision Boundary



Predict " $y = 1$ " if $\underline{-3 + x_1 + x_2 \geq 0}$

$$\Theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \leftarrow$$

$$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\frac{-3}{||} + \frac{1}{||} x_1 + \frac{1}{||} x_2$$

Decision boundary

$$x_1, x_2$$

$$\rightarrow h_{\theta}(x) = 0.5$$

$$\boxed{x_1 + x_2 = 3}$$

$$\Theta^T x$$

$$\rightarrow \underline{x_1 + x_2 \geq 3}$$

$$\left(\begin{array}{l} x_1 + x_2 < 3 \\ y = 0 \end{array} \right)$$

Higher order terms??

Non-linear decision boundaries

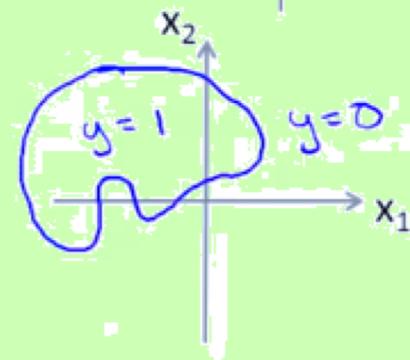


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$\Theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

Predict " $y = 1$ " if $x_1^2 + x_2^2 \geq 1$

$x_1^2 + x_2^2 \geq 1$



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 + \dots)$$

Change feature vector
Andrew N

Probabilistic interpretation

- If we have a value between 0 and 1, let's use it to model class probability

$$p(C = 0|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) \quad \text{with} \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

- Substituting we have

$$p(C = 0|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - w_0)}$$

- Suppose we have two classes, how can I compute $p(C = 1|\mathbf{x})$?
- Use the marginalization property of probability

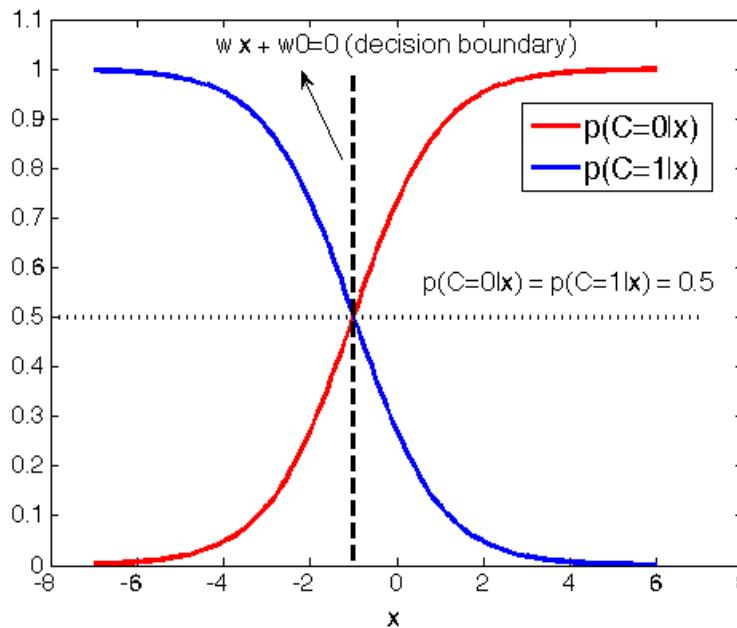
$$p(C = 1|\mathbf{x}) + p(C = 0|\mathbf{x}) = 1$$

- Thus

$$p(C = 1|\mathbf{x}) = 1 - \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - w_0)} = \frac{\exp(-\mathbf{w}^T \mathbf{x} - w_0)}{1 + \exp(-\mathbf{w}^T \mathbf{x} - w_0)}$$

Decision Boundary for Logistic Regression

- What is the decision boundary for logistic regression?
- $p(C = 1|\mathbf{x}, \mathbf{w}) = p(C = 0|\mathbf{x}, \mathbf{w}) = 0.5$
- $p(C = 0|\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) = 0.5$, where $\sigma(z) = \frac{1}{1+\exp(-z)}$
- Decision boundary: $\mathbf{w}^T \mathbf{x} + w_0 = 0$
- Logistic regression has a **linear decision boundary**

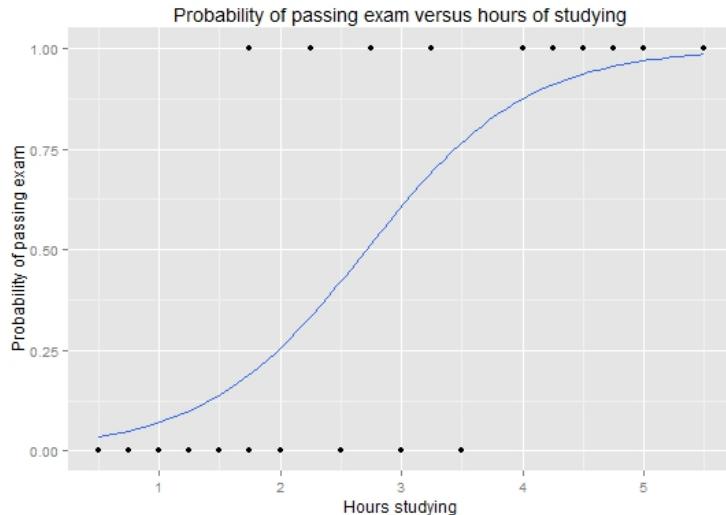


Example

- **Problem:** Given the number of hours a student spent learning, will (s)he pass the exam?
- Training data (top row: $\mathbf{x}^{(i)}$, bottom row: $t^{(i)}$)

Hours	0.50	0.75	1.00	1.25	1.50	1.75	1.75	2.00	2.25	2.50	2.75	3.00	3.25	3.50	4.00	4.25	4.50	4.75	5.00	5.50
Pass	0	0	0	0	0	0	1	0	1	0	1	0	1	0	1	1	1	1	1	1

- Learn \mathbf{w} for our model, i.e. logistic regression (coming up)
- Make predictions:



Hours of study	Probability of passing exam
1	0.07
2	0.26
3	0.61
4	0.87
5	0.97

Fitting parameters to logistic model

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

m examples

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix} \quad x_0 = 1, y \in \{0, 1\}$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose parameters θ ?

Cost function

→ Linear regression:

logistic

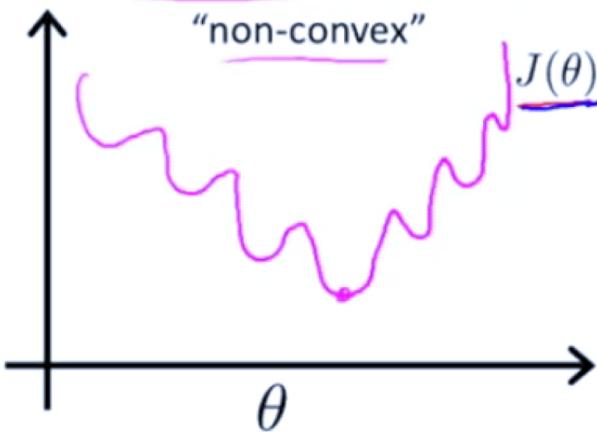
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

cost($h_\theta(x^{(i)})$, $y^{(i)}$)

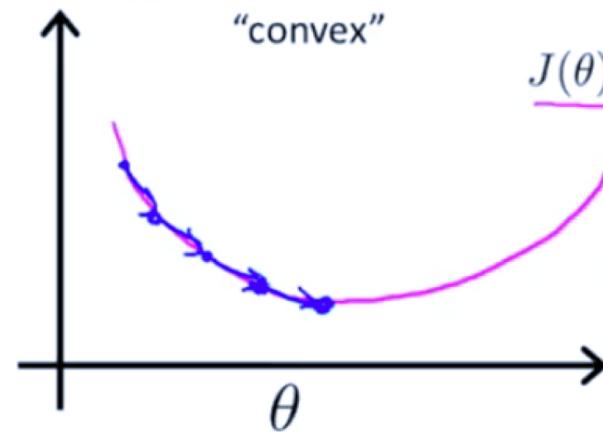
→ Cost($h_\theta(x^{(i)})$, $y^{(i)}$)

$$\frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\frac{1}{2} h_\theta(x^{(i)})^2$$



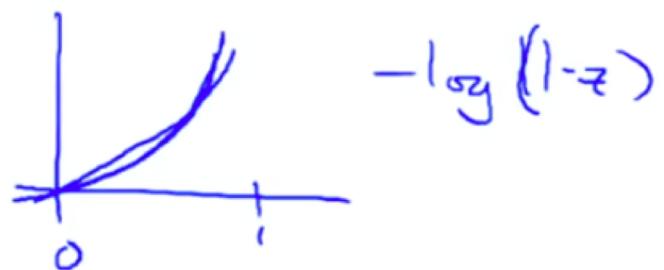
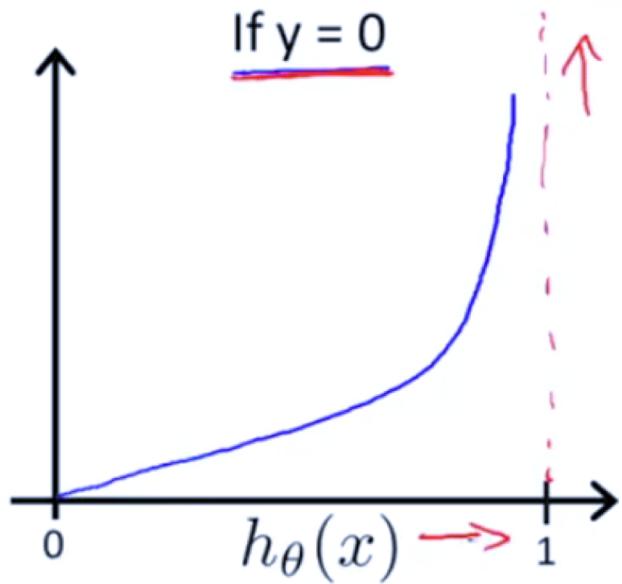
"non-convex"



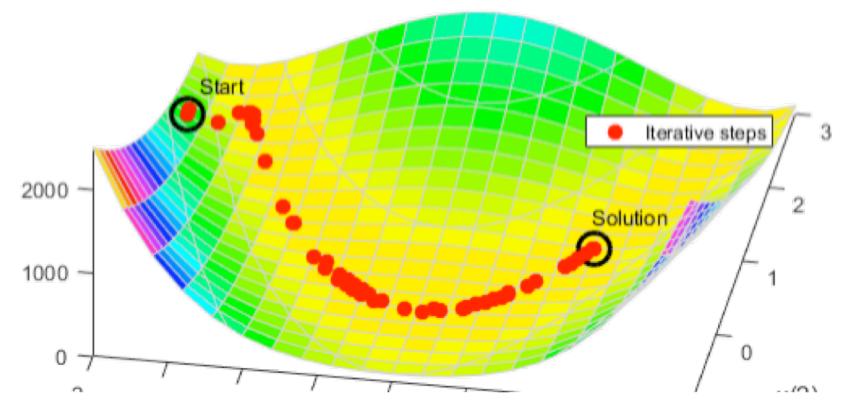
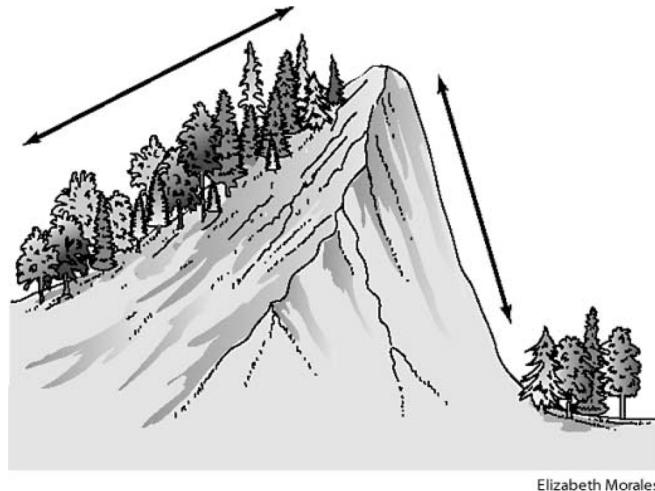
"convex"

Logistic regression cost function

$$\text{Cost}(h_\theta(x^{(i)}, y^{(i)})) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ \boxed{-\log(1 - h_\theta(x))} & \text{if } y = 0 \end{cases}$$



A few notes on non-linear parameter learning

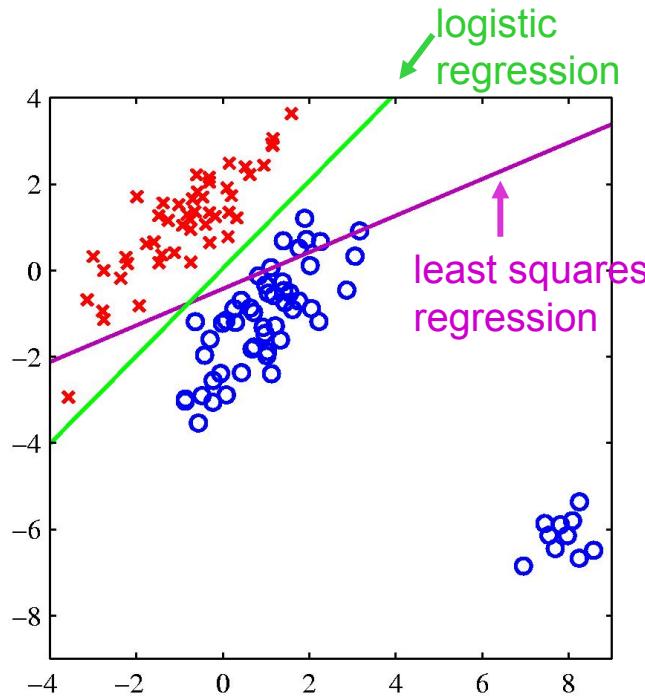
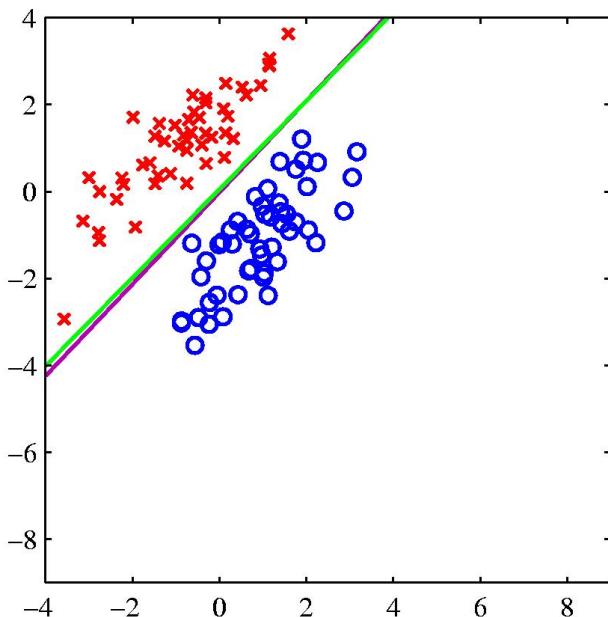


Repeat until convergence {

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

Logistic Regression vs Least Squares regression



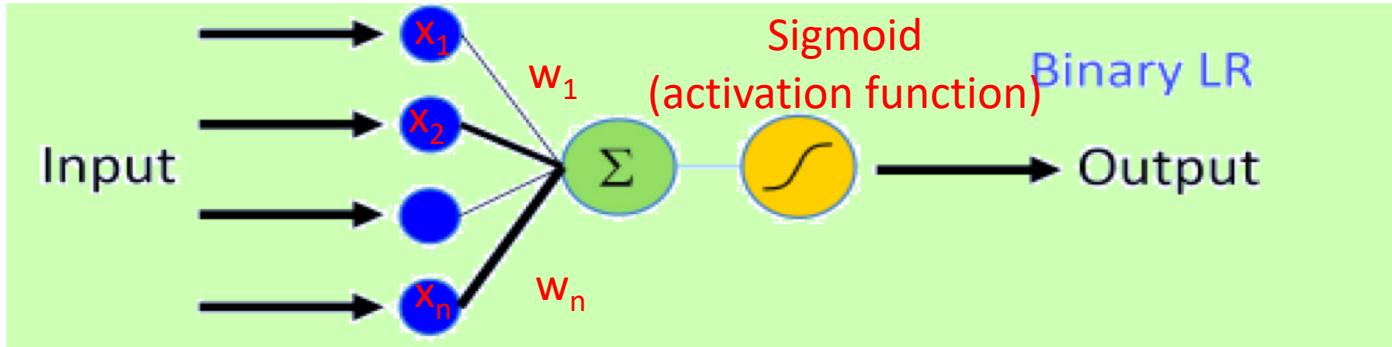
If the right answer is 1 and the model says 1.5, it loses, so it changes the boundary to avoid being “too correct” (tilts away³³ from outliers)

One way to think about it...

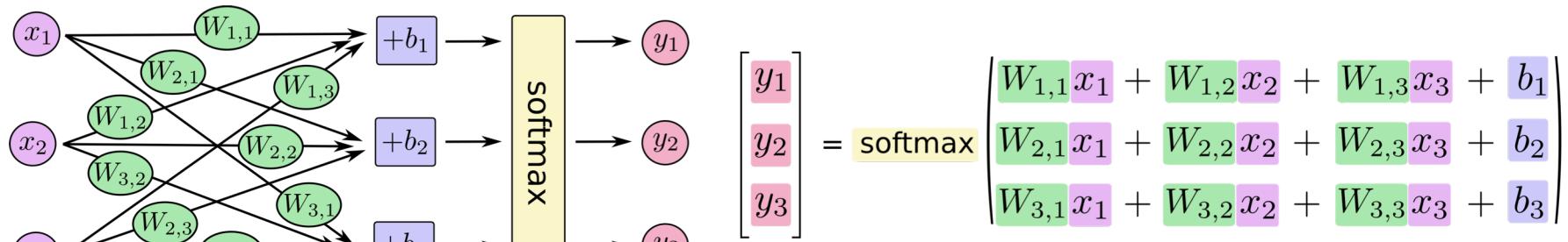
- Training labels dictate that two examples are the same or different (in some sense).
- Features and distance measures define visual similarity
- Classifiers try to learn weights or parameters for features and distance measures so that visual similarity predicts label similarity

Visualizing logistic regression

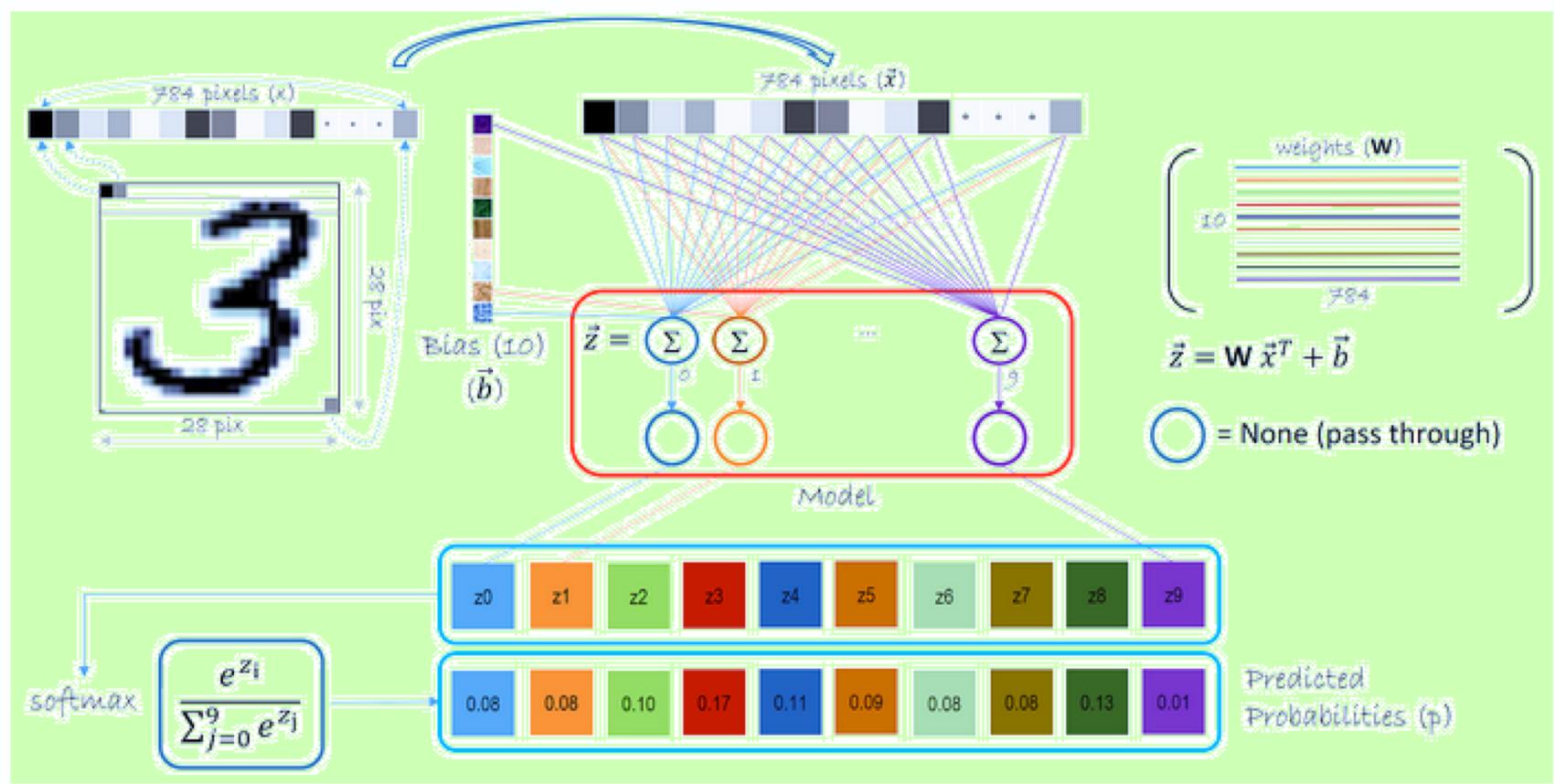
Binary classification



Multi class classification

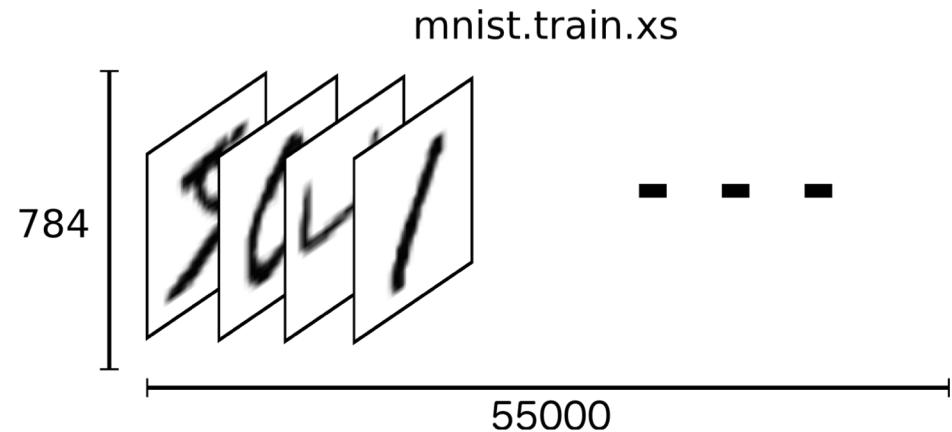


$$P(y = j | \mathbf{x}) = \frac{e^{\mathbf{x}^\top \mathbf{w}_j}}{\sum_{k=1}^K e^{\mathbf{x}^\top \mathbf{w}_k}}$$



MNIST

label = 5	5	label = 0	0	label = 4	4	label = 1	1	label = 9	9
label = 2	2	label = 1	1	label = 3	3	label = 1	1	label = 4	4
label = 3	3	label = 5	5	label = 3	3	label = 6	6	label = 1	1
label = 7	7	label = 2	2	label = 8	8	label = 6	6	label = 9	9



Logistic regression using keras

```
def build_logistic_model(input_dim, output_dim):
    model = Sequential()
    model.add(Dense(output_dim, input_dim=input_dim, activation='softmax'))

    return model
```

Can Classes Always be Separated

Causes of non perfect separation:

- Model is too simple
- Noise in the inputs (i.e., data attributes)
- Simple features that do not account for all variations
- Errors in data targets (miss labelings)

Should we make the model complex enough to have perfect separation in the training data?

Metrics

How to evaluate how good my classifier is? How is it doing on dog vs no-dog?



TP (True Positive)

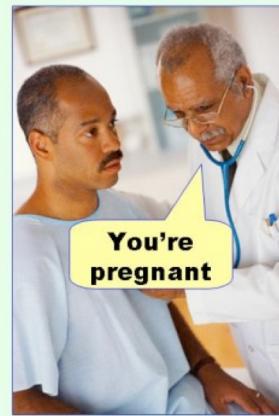
FP (False Positive)

FN (False Negative)

Measured /
Perceived

		Reality	
		True	False
True	True	Correct 	Type I False Positive
	False	Type II False Negative	Correct

Type I error
(false positive)

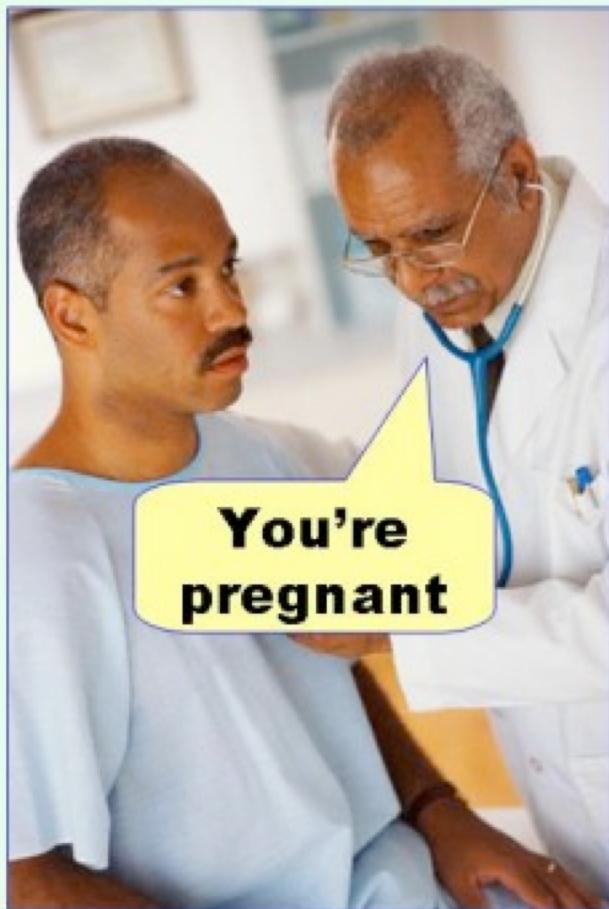


Type II error
(false negative)



Type I error

(false positive)

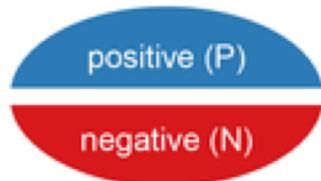


Type II error

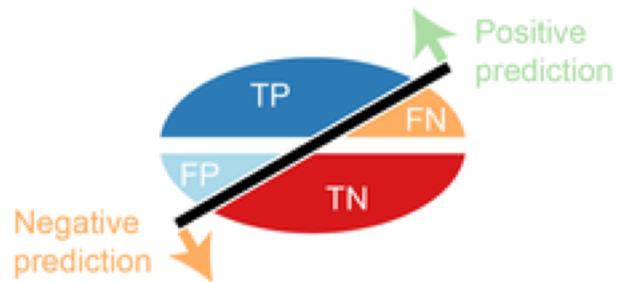
(false negative)



Observed labels

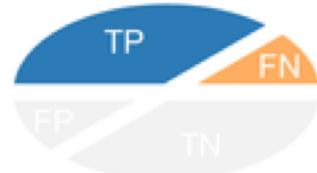


Four outcomes of a classifier



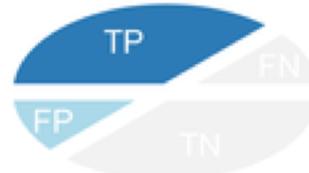
x-axis Recall

Sensitivity



y-axis Precision

Positive predictive value



Metrics

How to evaluate how good my classifier is?

- Recall: is the fraction of relevant instances that are retrieved

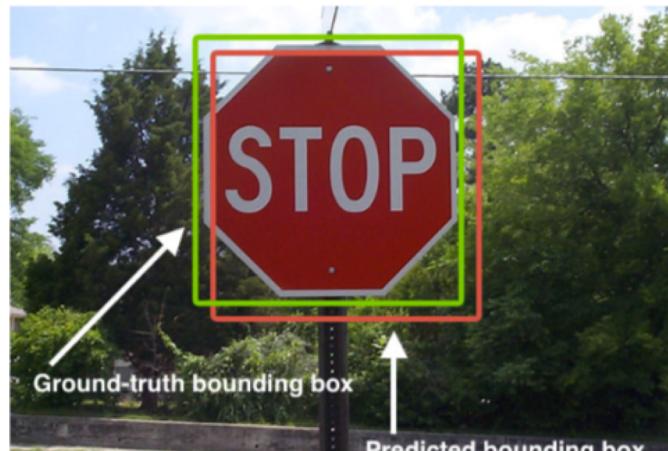
$$R = \frac{TP}{TP + FN} = \frac{TP}{\text{all groundtruth instances}}$$

- Precision: is the fraction of retrieved instances that are relevant

$$P = \frac{TP}{TP + FP} = \frac{TP}{\text{all predicted}}$$

- F1 score: harmonic mean of precision and recall

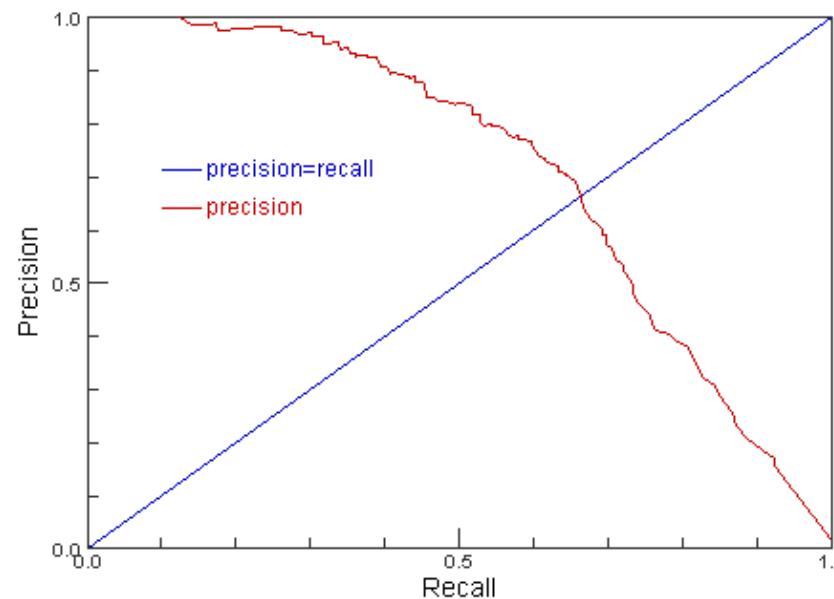
$$\text{IoU} = \frac{\text{area}(B_p \cap B_{gt})}{\text{area}(B_p \cup B_{gt})}$$



More on metrics

How to evaluate how good my classifier is?

- **Precision:** is the fraction of retrieved instances that are relevant
- **Recall:** is the fraction of relevant instances that are retrieved
- **Precision Recall Curve**



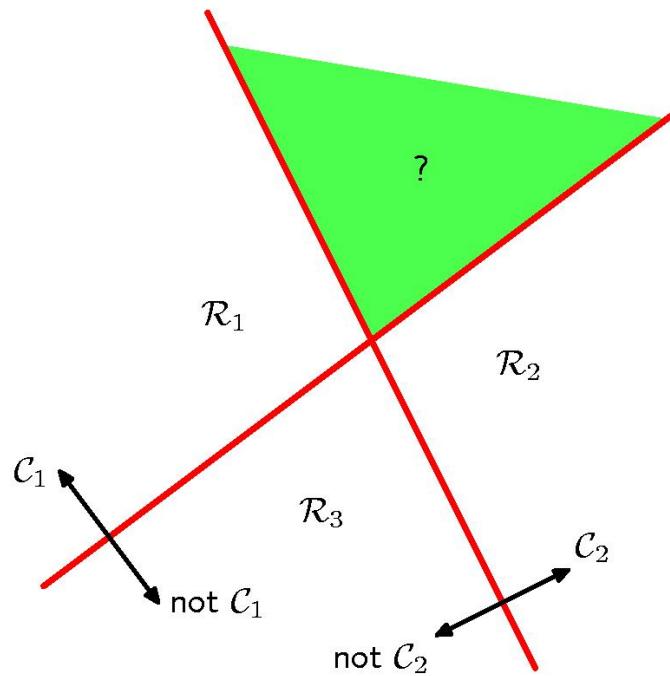
- **Average Precision (AP):** mean under the curve

Multi-class classification



Idea of discriminant functions $K > 2$

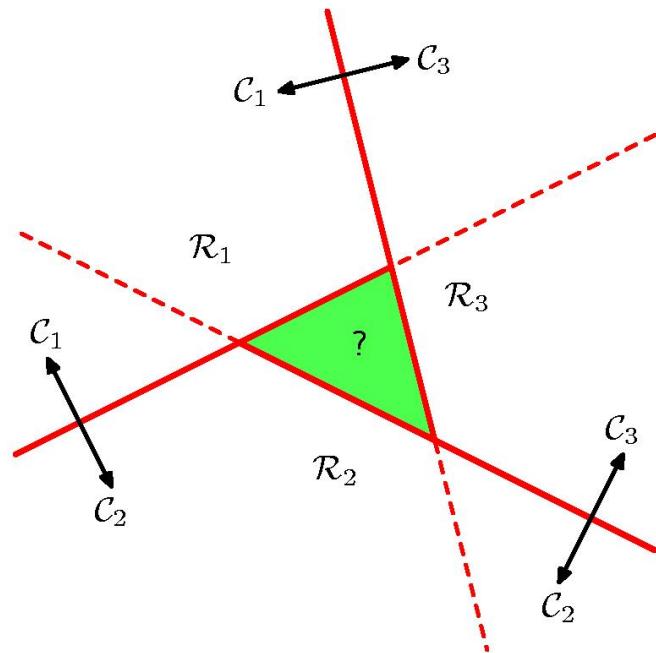
- First idea: Use $K - 1$ classifiers, each solving a two class problem of separating point in a class C_k from points not in the class.
- Known as **1 vs all** or **1 vs the rest** classifier



- PROBLEM: More than one good answer for green region!

Idea of discriminant functions $K > 2$

- Another simple idea: Introduce $K(K - 1)/2$ two-way classifiers, one for each possible pair of classes
- Each point is classified according to majority vote amongst the disc. func.
- Known as the **1 vs 1 classifier**



- PROBLEM: Two-way preferences need not be transitive

Next week

- Non-linear models
 - Neural networks
 - SVM
 - Kernels
 - WATCH THE VIDEOS BEFORE CLASS

FOR NN

- https://github.com/Hvass-Labs/TensorFlow-Tutorials/blob/master/02_Convolutional_Neural_Network.ipynb