# Convolutional Neural Networks (CNNs) and Support Vector Machines (SVMs)

**CSCI-P556 Applied Machine Learning**

**Lecture 18**

**D.S. Williamson**
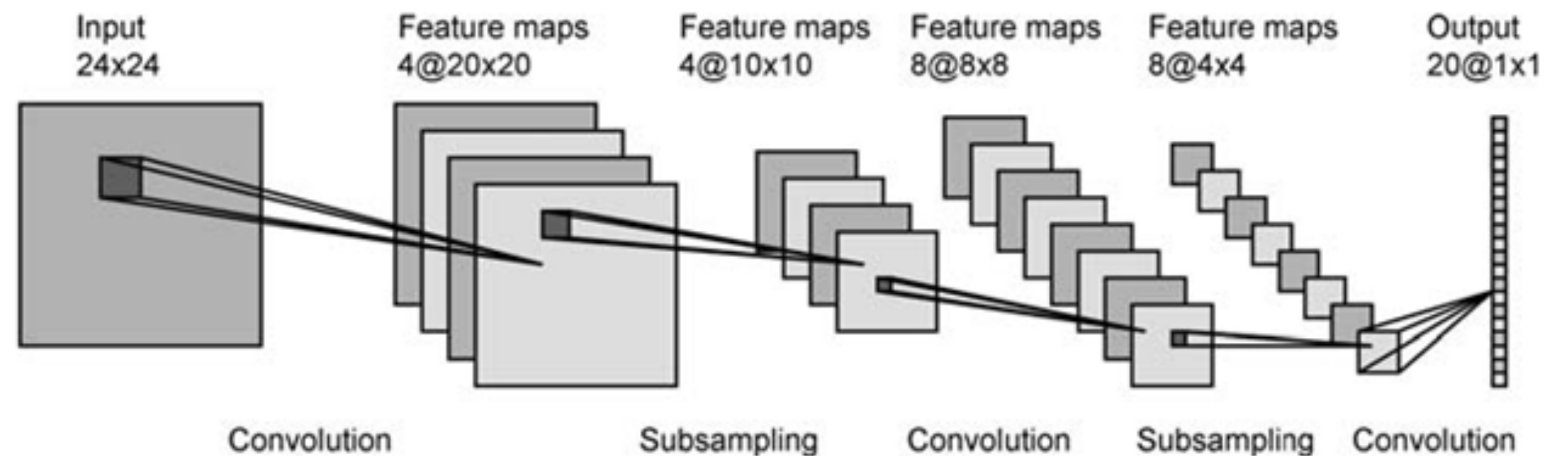
# Agenda and Learning Outcomes

## Today's Topics

- **Topics**:

  - Convolutional Neural Networks (CNNs)

  - Support Vector Machines, Part I


- **Announcements**

  - Quiz #2 on Thursday

  - Project proposal comments (tomorrow)

  - hw#3 coming soon (today or tomorrow)

# Convolutional Neural Networks

# Convolutional Neural Networks (CNNs)

- Used for processing data with grid-like topology (i.e. images). Networks use convolution in place of general matrix multiplication

- Good at capturing local (short-term) dependencies and correlations (e.g. correlations amongst adjacent pixels in an image, or dependencies across nearby frequencies of an audio signal)

- There are four main operations in CNNs

  - Convolution

  - Nonlinear Activation Function (i.e. ReLU)

  - Pooling (or Subsampling)

  - Classification

# Convolution

- Convolution is a linear mathematical operation on two functions

- Given functions *x(t)* and *w(t),* the convolution of *x(t)* and *w(t)* is as follows

$$s(t) = x(t) * w(t) = \sum_{a=-\infty}^{\infty} x(a)w(t-a)$$

- This step is used to extract "features" from an input, so it is also referred to as the ***feature map*** stage

# Example: Convolution on an Image

- Suppose you are given the following binary image, X



- You want to convolve this image with matrix, W as shown below





Image

Convolved
Feature

Images courtesy of  *the data science blog*
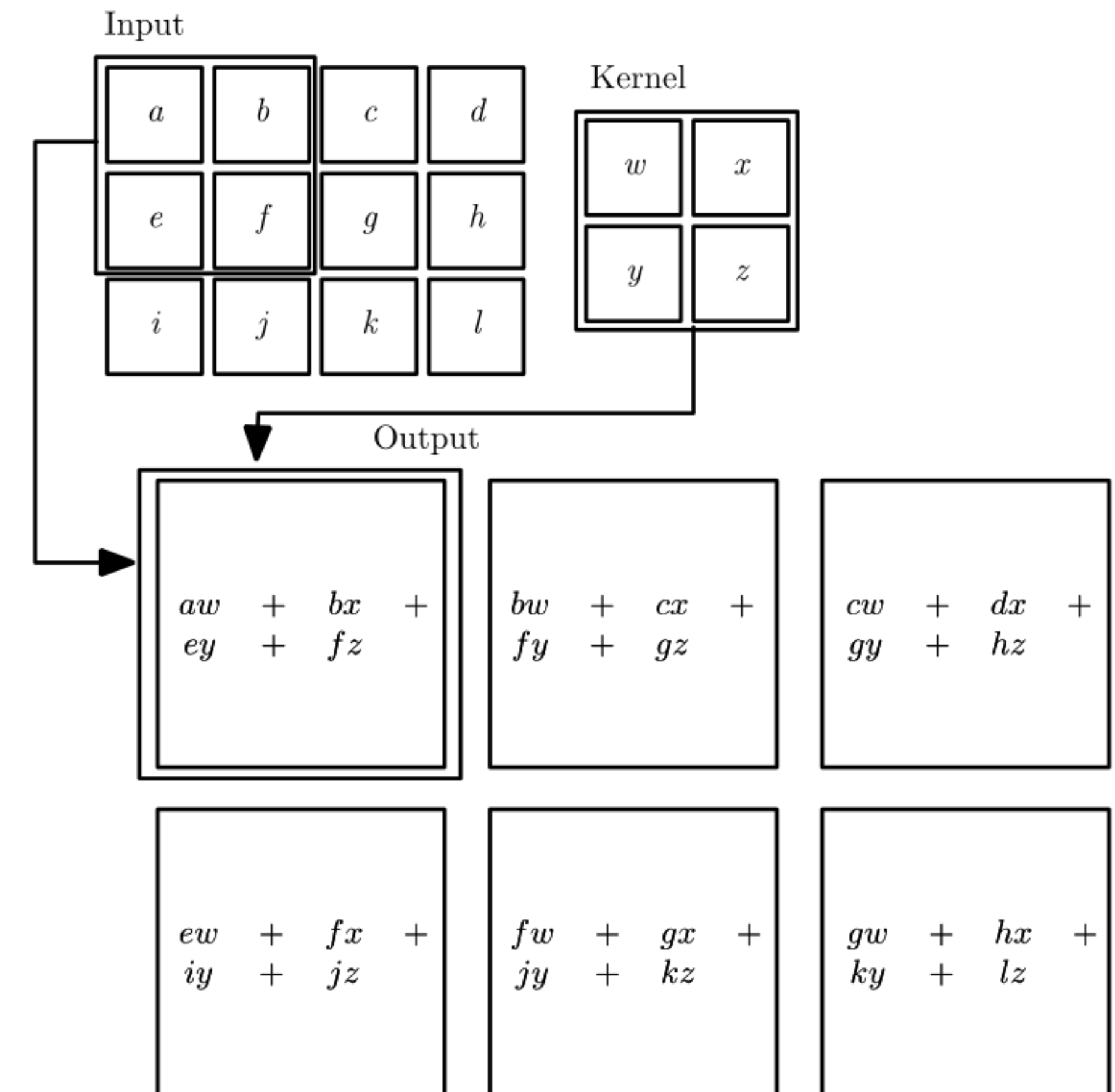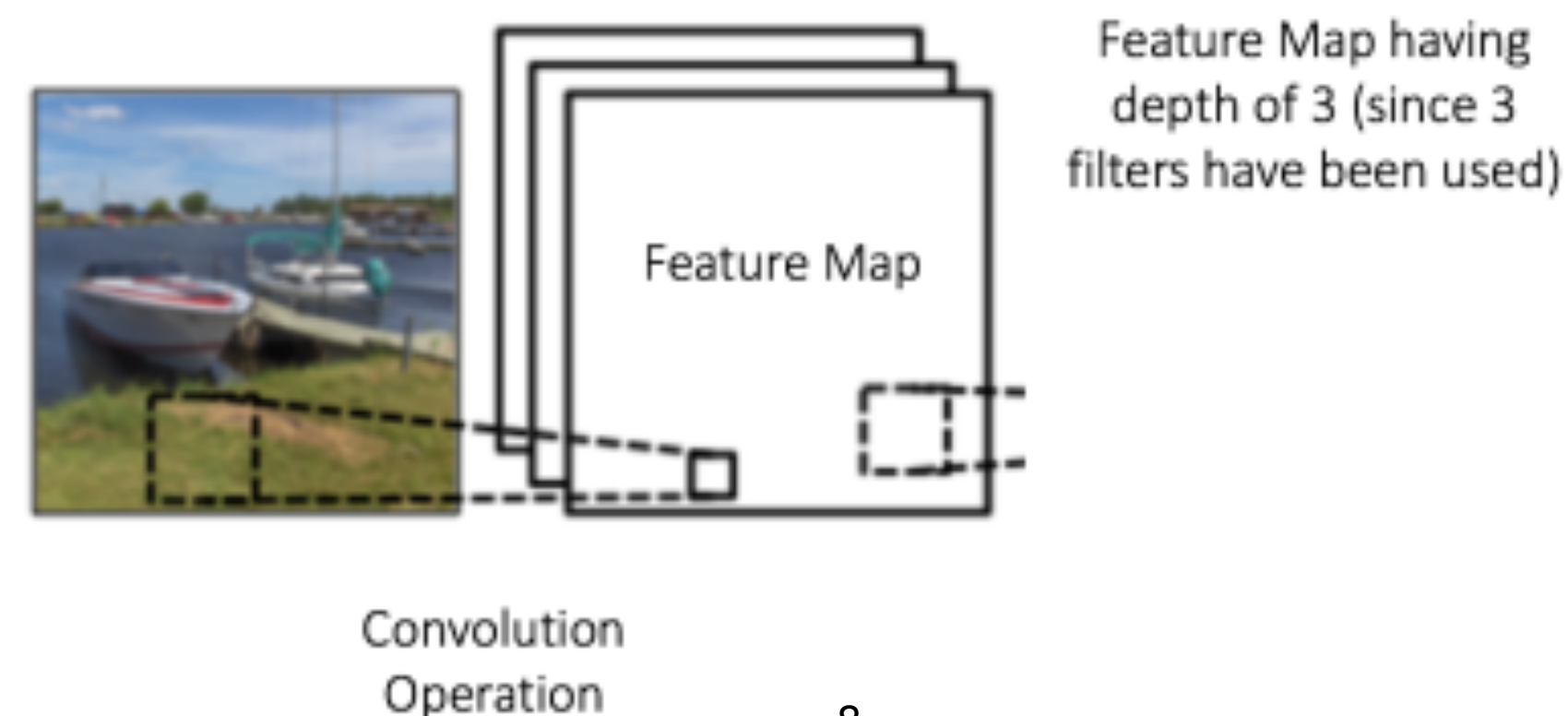
# Convolution Operation

- A mathematical depiction of the convolution operation on an input image

- A CNN learns the values of the filter (or kernel) on its own during the training process



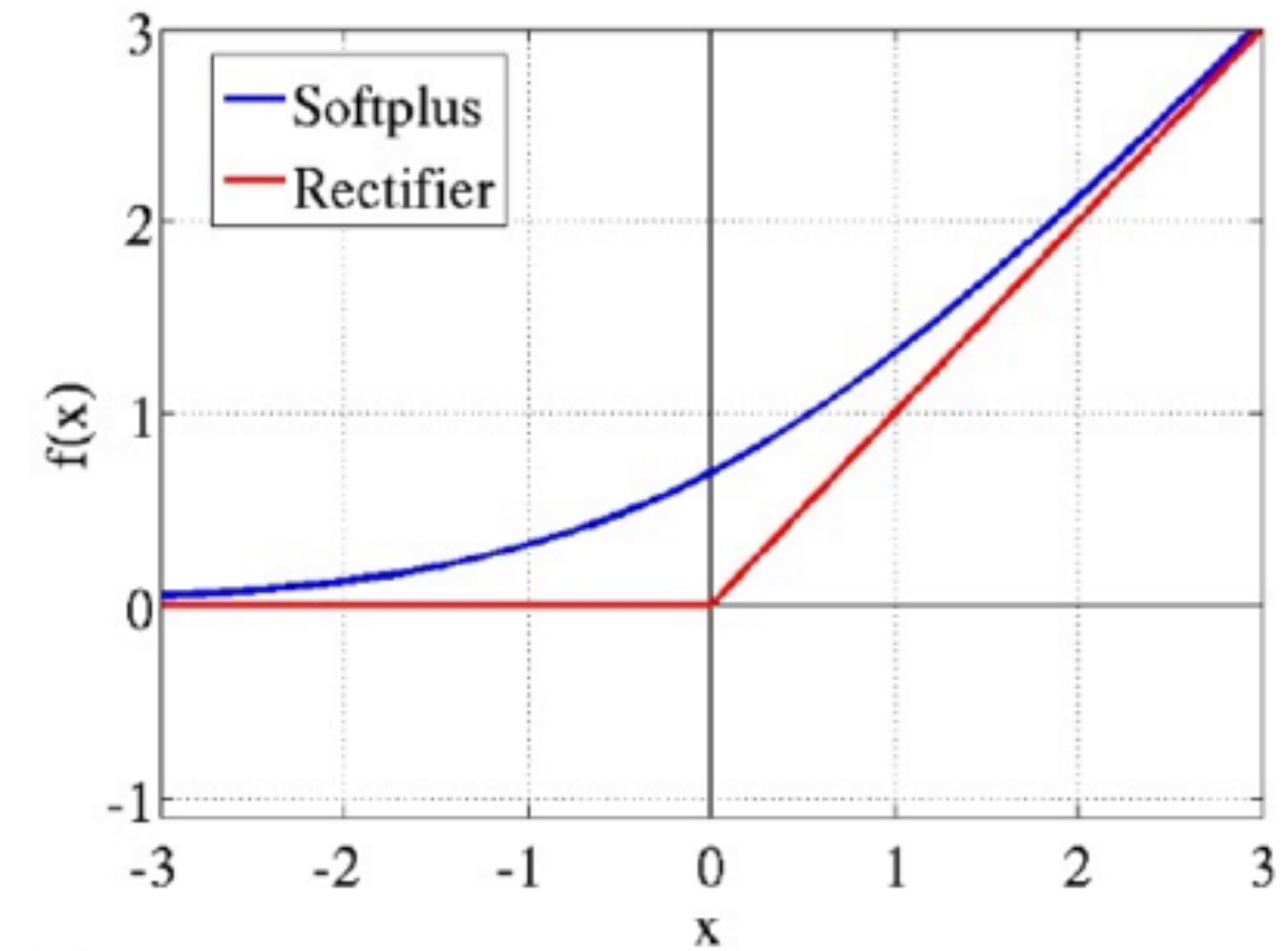Images courtesy of  *(Goodfellow 2016)*

# Feature Map

- The size of the Feature Map (resulting image after convolution) is controlled by three parameters

  - **Depth**: Number of different filters to use for the convolution operation

  - **Stride**: Number of pixels used to slide the filter across the input

  - **Zero-padding**: May pad the input with zeros around the border



Feature Map having depth of 3 (since 3 filters have been used)

Feature Map

Convolution Operation
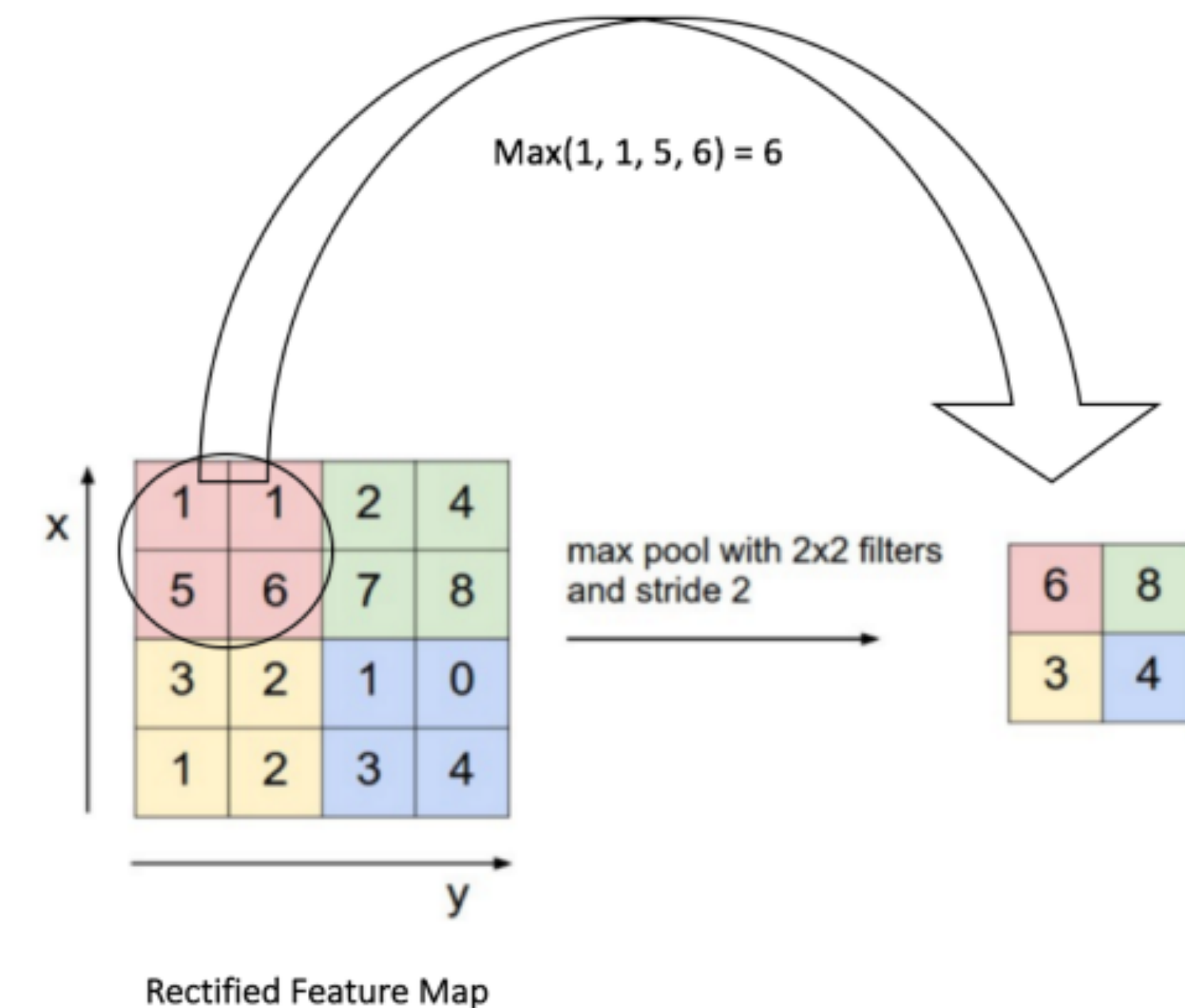
Images courtesy of *the data science blog*

# Rectified Linear (ReLU) Activation

- A rectified linear (ReLU) activation operation may be applied after the convolution operation

  - Introduces nonlinearity to the network

  - ReLU(x) = max(0, x)

  - Applied to every element (pixel)

  - Negative values are replaced by 0

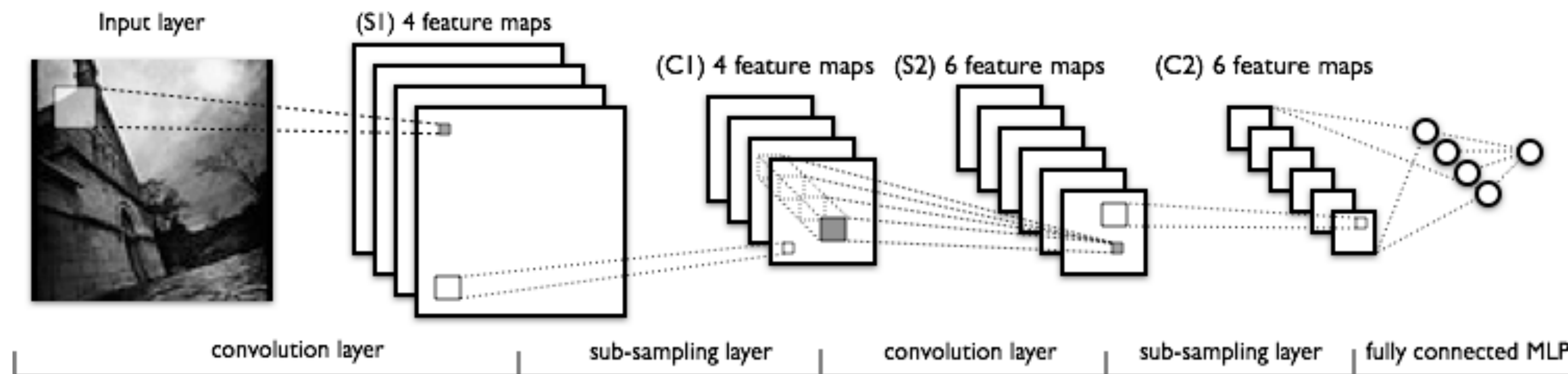- Other nonlinear activation functions may be used instead

# Pooling

- Pooling (aka spatial pooling, subsampling, or downsampling) is used to reduce the dimensionality of the feature map

- Different types of pooling include: Max, Average, Sum, etc.

- A window is defined, and the pooling operation is performed over the elements within that window

- The pooling window slides over the feature map by the stride amount

- It is applied to each feature map



Max(1, 1, 5, 6) = 6

max pool with 2x2 filters and stride 2

Rectified Feature Map

# CNN

- Multiple layers of Convolution, Activation, and Pooling may be used in a CNN

- These layers act as feature extraction, to find useful features from the input

- Generally, a final Fully Connected layer is added via a DNN for classification or regression purposes
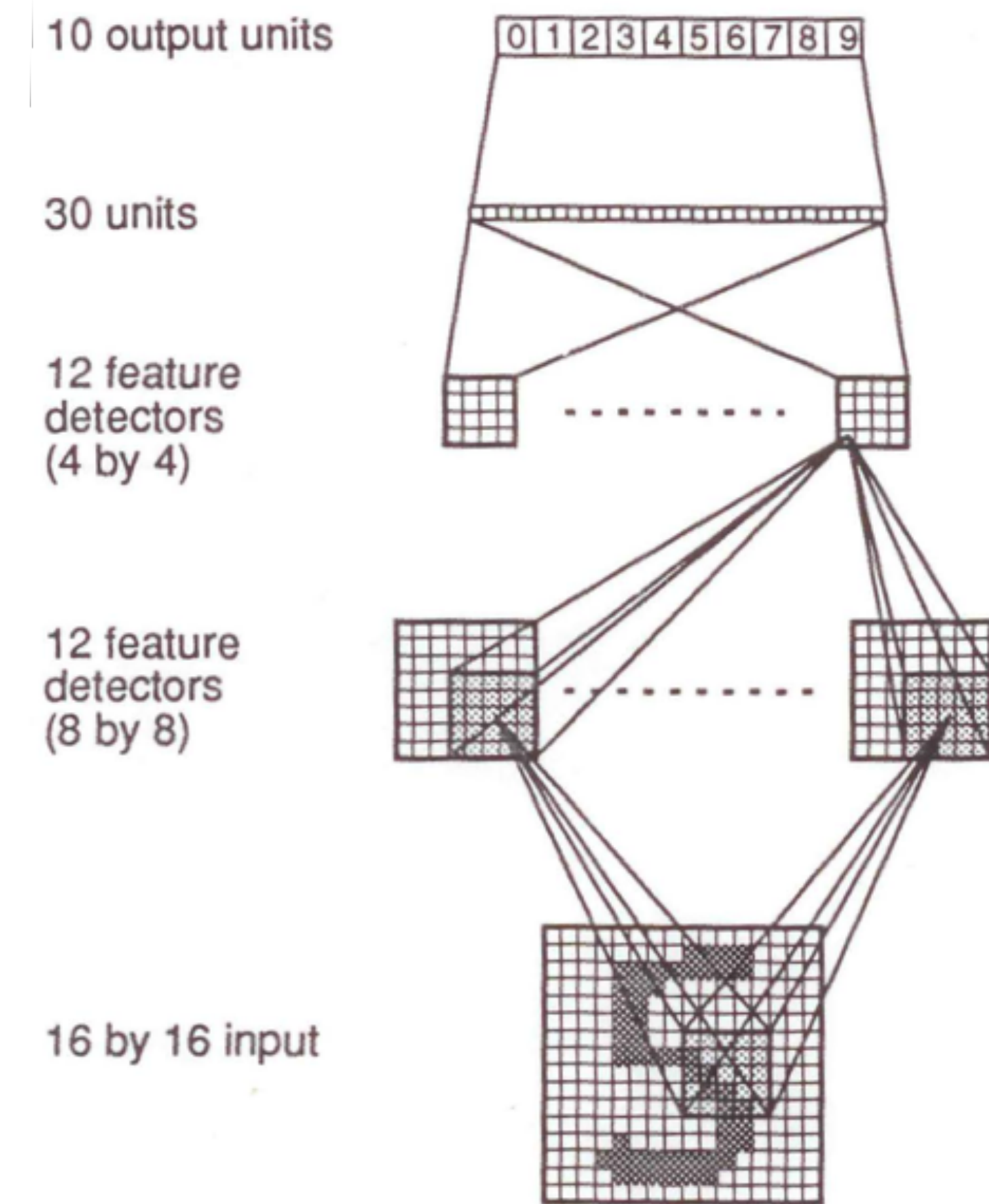
# CNN Training

- The Backpropagation algorithm is used to train the parameters of a CNN

- Basic steps:

  - Randomly initialize all filters (or kernels) and weights

  - Propagate the input forward through the layers of the CNN (convolution, activation, pooling, DNN) to get an output(s)

  - Calculate the error between the actual output(s) and the desired output(s)

  - Use Backpropagation to calculate the gradients and deltas, and then update the filters and weights accordingly

# An Early CNN Application

## Handwritten Zip code Recognition (1989)

- **Input**: binary pixels for each digit (16 x 16 dimensional)

- Padded with -1.

- **Output**: 10 digits, using sigmoid activations

- Scaled hyperbolic tanget activation function

- **Architecture**: 4 layers (12x8x8 – 12x4x4 -30 -10)

  - Layer 1: kernel of size 5x5, stride of 2, depth of 12

  - Layer 2: kernel of size 5x5, stride of 2, depth of 12

  - Layer 3: Fully-connected layer with 30 units

  - Layer 4: 10 unit output layer

- **Performance**: Trained on 7300 digits and tested on 2000 new ones

  - Achieved 1% error on the training set and 5% error on the test set

  - If allowing rejection (no decision), 1% error on the test set

  - This task is not easy



10 output units `0 1 2 3 4 5 6 7 8 9`

30 units

12 feature detectors (4 by 4)

12 feature detectors (8 by 8)

16 by 16 input

# CNN in PyTorch

- The following classes can be used within a defined class that implements the desired CNN architecture. (e.g. in the *__init__()* and *forward()* functions)

  - **Convolution**: Conv2d (or Conv1d) to perform convolution operation

  - **Pooling**: Maxpool2D, Avgpool2D,..

# Support Vector Machines

# Warning!

- This is just a first introduction to SVMs – very basic version of SVM introduced
- We will discuss more on the advanced techniques later
  - Hence, Strong assumptions in this case

# The Classification Problem

- ***Binary classification*** can be viewed as the task of ***separating classes in the feature space***

- This is accomplished by formulating a decision boundary
  - **w** is the slope of the line
  - **b** is the intercept

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

$$\mathbf{w} \cdot \mathbf{x} + b > 0$$

$$\mathbf{w} \cdot \mathbf{x} + b < 0$$

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

# Recall: Perceptrons

- **Perceptron Convergence Theorem**: If a classification problem is linearly separable, a perceptron will reach a solution in a finite number of iterations

- The solution weight vector is **<u>not</u>** unique. There are **_infinite possible solutions_** and decision boundaries.

  - Perceptrons find any separating hyperplane

  - The hyperplane depends on initialization and ordering of training points

# Recall: Perceptrons

- **Perceptron Convergence Theorem**: If a classification problem is linearly separable, a perceptron will reach a solution in a finite number of iterations

- The solution weight vector is **not** unique. There are **_infinite possible solutions_** and decision boundaries.

  - Perceptrons find any separating hyperplane

  - The hyperplane depends on initialization and ordering of training points

- **_If done differently_**

# Recall: Perceptrons

- **Perceptron Convergence Theorem**: If a classification problem is linearly separable, a perceptron will reach a solution in a finite number of iterations

- The solution weight vector is **not** unique. There are ***infinite possible solutions*** and decision boundaries.

  - Perceptrons find any separating hyperplane

  - The hyperplane depends on initialization and ordering of training points
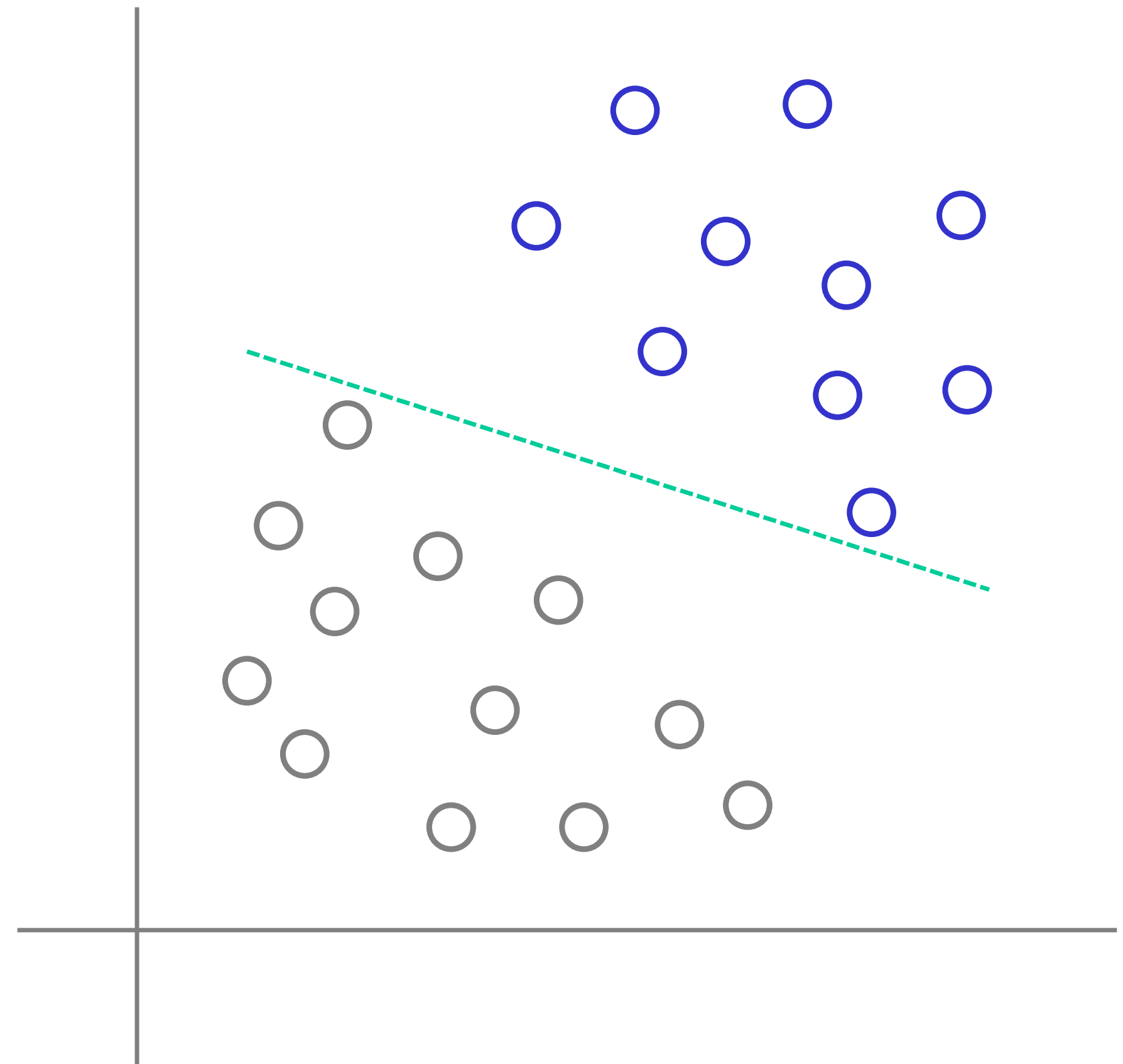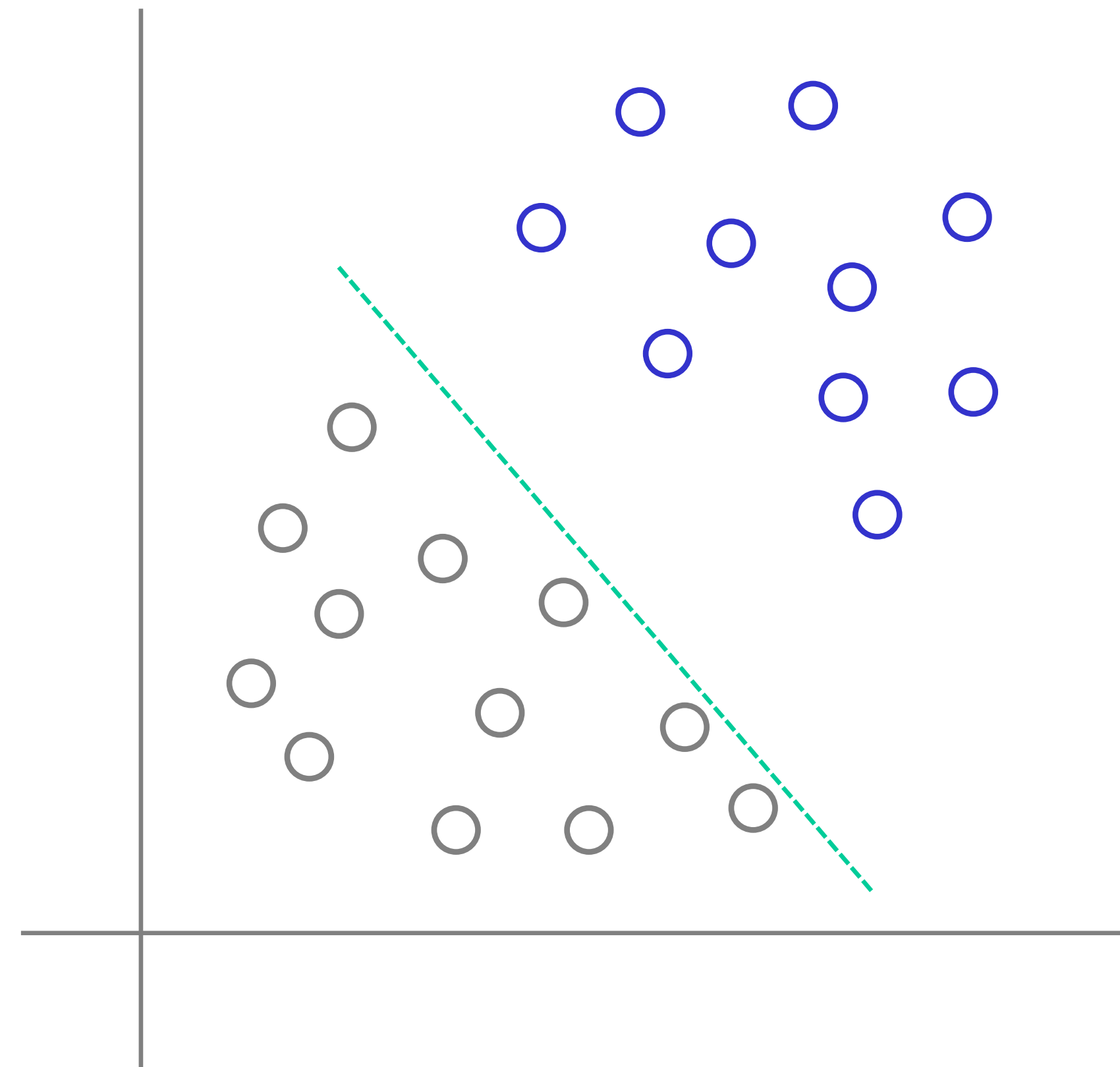
- ***If done differently….Again***

# Recall: Perceptrons

- **Perceptron Convergence Theorem**: If a classification problem is linearly separable, a perceptron will reach a solution in a finite number of iterations

- The solution weight vector is **not** unique. There are **_infinite possible solutions_** and decision boundaries.

  - Perceptrons find any separating hyperplane

  - The hyperplane depends on initialization and ordering of training points

- **_If done differently….Again...And Again_**

# Linear Separators

- Which is the best linear separator?
  - Depends on the goal
  - Goal is to classify **accurately** and **generalize** to new examples.

# Notion of Margins

- Many different hyperplanes can classify the data. Which one will work best?

- The **hyperplane** that **maximizes** the **separation between the two classes** (the margin)

# Intuition of a Margin

- Consider points A, B, and C

- We are quite confident in our prediction for A because it is far from the decision boundary.

- In contrast, we are not so confident in our prediction for C because a slight change in the decision boundary may flip the decision

- Given a training set, we would like to make all predictions correct and confident! This leads to the concept of margin

# Why Max Margin?

- Minimizes generalization error. Works well on **Future data**

- Minimizes Complexity. **Fewer support vectors**

- Minimizes the capacity of the classifier. **Eliminates overfitting**

# Decision Boundary

## Need to know distance from point to the decision boundary

- Given a decision boundary (e.g. linear discriminant function)

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$$

- To find its distance to a given pattern $\mathbf{x}$, project $\mathbf{x}$ onto the decision boundary

# Decision Boundary (cont.)

- **x** can be re-written as a function of the projection and the weights

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{||\mathbf{w}||}$$

- $\mathbf{x}_p$ is **x**'s projection

- The second term arises from the fact that the weight vector is perpendicular to the decision boundary

- The algebraic distance $r$ is positive if **x** is on the positive side of the boundary and negative if **x** is on the negative side

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$$

# Decision Boundary (cont.)

- Since **x** can be written in terms of its projection and $r$, then the decision boundary can as well:

$$g(\mathbf{x}) = g(\mathbf{x}_p + r\frac{\mathbf{w}}{||\mathbf{w}||})$$

$$= \mathbf{w}^T(\mathbf{x}_p + r\frac{\mathbf{w}}{||\mathbf{w}||}) + b$$

$$= \mathbf{w}^T\mathbf{x}_p + b + r||\mathbf{w}||$$

$$= r||\mathbf{w}||$$

$$g(\mathbf{x}) = \mathbf{w}^T\mathbf{x} + b = 0$$



- Thus, $r = \dfrac{g(\mathbf{x})}{||\mathbf{w}||}$

As a special case, for the origin, r = b/||w||

# Margins

- Distance from example $\mathbf{x}_i$ to the separator is $r = \dfrac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|}$

- Examples closest to the hyperplane are **support vectors**.

- **Margin** $\gamma$ of the separator is the distance between support vectors.

# Notation

- We denote the classifier, $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, for all $\mathbf{x} \in \mathbb{R}^n$

- **Assumptions**

  - Supporting hyperplanes $\mathbf{w}^T \mathbf{x} + b = \pm 1$

  - The distance between the two supporting hyperplanes is the margin, which is $\gamma = 2$

- To achieve **scale invariance** we divide the classifier by $||\mathbf{w}||_2$.
  Then the supporting hyperplanes are $\mathbf{w}^T \mathbf{x} + b = \pm \dfrac{1}{||\mathbf{w}||_2}$

  margin is $\gamma = \dfrac{2}{||\mathbf{w}||_2}$.

# Max margin Classifier

- Given a **linearly separable** training set $S = \{(\mathbf{x}^{(i)}, y^{(i)}) : i = 1, \cdots, N\}$, we would like to find a classifier with a maximum margin, $\gamma$

- This can be represented as an optimization problem.

$$\max_{\mathbf{w}, b, \gamma} \gamma \quad \text{subject to: } y^{(i)} \frac{(\mathbf{w}^T \mathbf{x}^{(i)} + b)}{||\mathbf{w}||} \geq \gamma, \quad i = 1, \cdots, N$$

Constraint ensures accurate classification

Nasty optimization problem! Let's make it look nicer!

- Let $\gamma' = \gamma ||\mathbf{w}||$, this is equivalent to

$$\max_{\mathbf{w}, b, \gamma'} \frac{\gamma'}{||\mathbf{w}||} \quad \text{subject to: } y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq \gamma', \quad i = 1, \cdots, N$$

# Max margin Classifier

- Note that rescaling $\mathbf{w}$ and $b$ by $(1/\gamma')$ will not change the classifier, we can thus further reformulate the optimization problem

$$\max_{\mathbf{w},b,\gamma'} \frac{\gamma'}{||\mathbf{w}||} \quad \text{subject to: } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq \gamma', \quad i = 1,\cdots,N$$

$$\max_{\mathbf{w},b} \frac{1}{||\mathbf{w}||} \quad \text{subject to: } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1, \quad i = 1,\cdots,N$$

- Note that maximizing the geometric margin is equivalent to minimizing the magnitude of $\mathbf{w}$ subject to maintaining a functional margin of at least 1

# Solving the problem

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|^2$$

$$\text{subject to}: \ y^i(\mathbf{w}\cdot\mathbf{x}^i + b) \geq 1, \quad i = 1,\cdots,N$$

- This results in a ***quadratic optimization problem*** with *linear inequality constraints*.

- This is a well-known class of mathematical programming problems for which several (non-trivial) algorithms exist.

  - One could solve for $\mathbf{w}$ using any of these methods

- We will see that it is useful to first formulate an equivalent dual optimization problem and solve it instead

  - This requires a bit of machinery.

# Constrained Optimization

- The general optimization problem can be written as such, for generic functions

$$\min_{x} f(x) \quad \text{subject to: } g_i(x) \leq 0, \quad i = 1, \cdots, m$$

- To solve the above optimization problem, **consider the following cost function known as the Lagrangian**

$$\mathscr{L}(x, \alpha) = f(x) + \sum_{i} \alpha_i g_i(x)$$

- Under certain conditions it can be shown that for a solution $x'$ to the above problem, we have

$$f(x') = \min_{x} \max_{\alpha} \mathscr{L}(x, \alpha) = \max_{\alpha} \min_{x} \mathscr{L}(x, \alpha)$$

**Primal Form**

**Dual Form subject to $\alpha_i \geq 0$**

# Dual Problem

- After simplifying the inequality, the problem becomes:

$$\min \frac{1}{2} ||\mathbf{w}||^2 \quad \text{subject to: } 1 - y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \leq 0, \quad i = 1, \cdots, N$$

- **The Lagrangian is then**

$$\mathscr{L}(\mathbf{w}, b, \alpha) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + \sum_{i=1}^{N} \alpha_i \left[1 - y^{(i)}(\mathbf{w}^T\mathbf{x} + b)\right], \text{ subject to } \alpha_i \geq 0$$

- We want to solve $\max_{\alpha} \min_{x} \mathscr{L}(\mathbf{w}, b, \alpha)$ s.t. $\alpha_i \geq 0$

- **Setting the gradient of $\mathscr{L}$ w.r.t. $\mathbf{w}$ and $b$ to zero, we have**

$$\mathbf{w} - \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i = 0 \qquad \mathbf{w} = \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i$$

$$\sum_{i=1}^{N} \alpha_i y^i = 0$$

# Dual Problem

- If we substitute $\mathbf{w} = \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i$ in $\mathscr{L}$, we have

$$
\begin{aligned}
L(\boldsymbol{\alpha}) &= \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_{i=1}^{N} \alpha_i \{ y^i (\mathbf{w} \cdot \mathbf{x}^i + b) - 1 \} \\
&= \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^i y^j < \mathbf{x}^i \cdot \mathbf{x}^j > - \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^i y^j < \mathbf{x}^i \cdot \mathbf{x}^j > - b \sum_{i=1}^{N} \alpha_i y^i + \sum_{i=1}^{N} \alpha_i \\
&= \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^i y^j < \mathbf{x}^i \cdot \mathbf{x}^j >
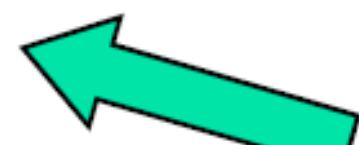\end{aligned}
$$

- Note that $\sum_{i=1}^{N} \alpha_i y^i = 0$

- This is a function of $\alpha_i$ only

# Dual Problem

- The new objective function is in terms of $\alpha_i$ only. It is known as the ***dual problem***: if we know all $\alpha_i$, then we know $\mathbf{w}$

  - The original problem is known as at the primal problem

- **The objective function of the dual problem needs to be maximized!**

$$\max L(\boldsymbol{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^i y^j < \mathbf{x}^i \cdot \mathbf{x}^j >$$

$$\text{subject to} \quad \alpha_i \geq 0, i = 1, \ldots, n, \qquad \sum_{i=1}^{N} \alpha_i y^i = 0$$

Properties of $\alpha_i$ when we introduce the Lagrange multipliers

The result when we differentiate the original Lagrangian w.r.t. b

# Dual Problem

$$\max L(\mathbf{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^i y^j < \mathbf{x}^i \cdot \mathbf{x}^j >$$

$$\text{subject to} \quad \alpha_i \geq 0, i = 1,...,n, \qquad \sum_{i=1}^{N} \alpha_i y^i = 0$$

- This is also a **_quadratic programming (QP) problem_**. A global maximum of $\alpha_i$ can always be found using a QP solver (beyond scope of this class). Luckily, SVMs have been implemented within many platforms.

- $\mathbf{w}$ can be recovered by $\mathbf{w} = \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i$

- $b$ can also be recovered as well (just a minute)

# Characteristics of the Solution

- Many of the $\alpha_i$ are zero. $\mathbf{w}$ is a linear combination of only a small number of data points

- In fact, optimization theory requires that the solution to satisfy the following KKT conditions:

$$\alpha_i \geq 0, i = 1,...,n,$$

$$y^i (\sum_{j=1}^{N} \alpha_i y^j < \mathbf{x}^j \cdot \mathbf{x}^i > + \mathrm{b}) \geq 1 \qquad \boxed{\text{Functional margin} \geq 1}$$

$$\alpha_i \{ y^i (\sum_{j=1}^{N} \alpha_i y^j < \mathbf{x}^j \cdot \mathbf{x}^i > + \mathrm{b}) - 1 \} = 0 \qquad \boxed{\begin{array}{l}\alpha_i \text{ is nonzero only when} \\ \text{functional margin} = 1\end{array}}$$

- **$\mathbf{x}_i$ with non-zero $\alpha_i$ are called support vectors (SV)**

  - Let $t_j, (j = 1, \cdots, s)$ be the indices of the $s$ support vectors. We can write $\mathbf{w} = \sum_{j=1}^{s} \alpha_{t_j} y^{t_j} \mathbf{x}^{t_j}$

  - The decision boundary is determined only by the SV

# Solving for b, the bias

- Note that we know that for support vectors the functional margin = 1

- We can use this information to solve for $b$

- We can use any support vector to achieve this

$$y^i \left( \sum_{j=1}^{s} \alpha_{t_j} y^{t_j} < \mathbf{x}^{t_j} \cdot \mathbf{x}^i > + \mathrm{b} \right) = 1$$

- A numerically more stable solution is to use all support vectors (see a ML textbook)

# Classifying new examples
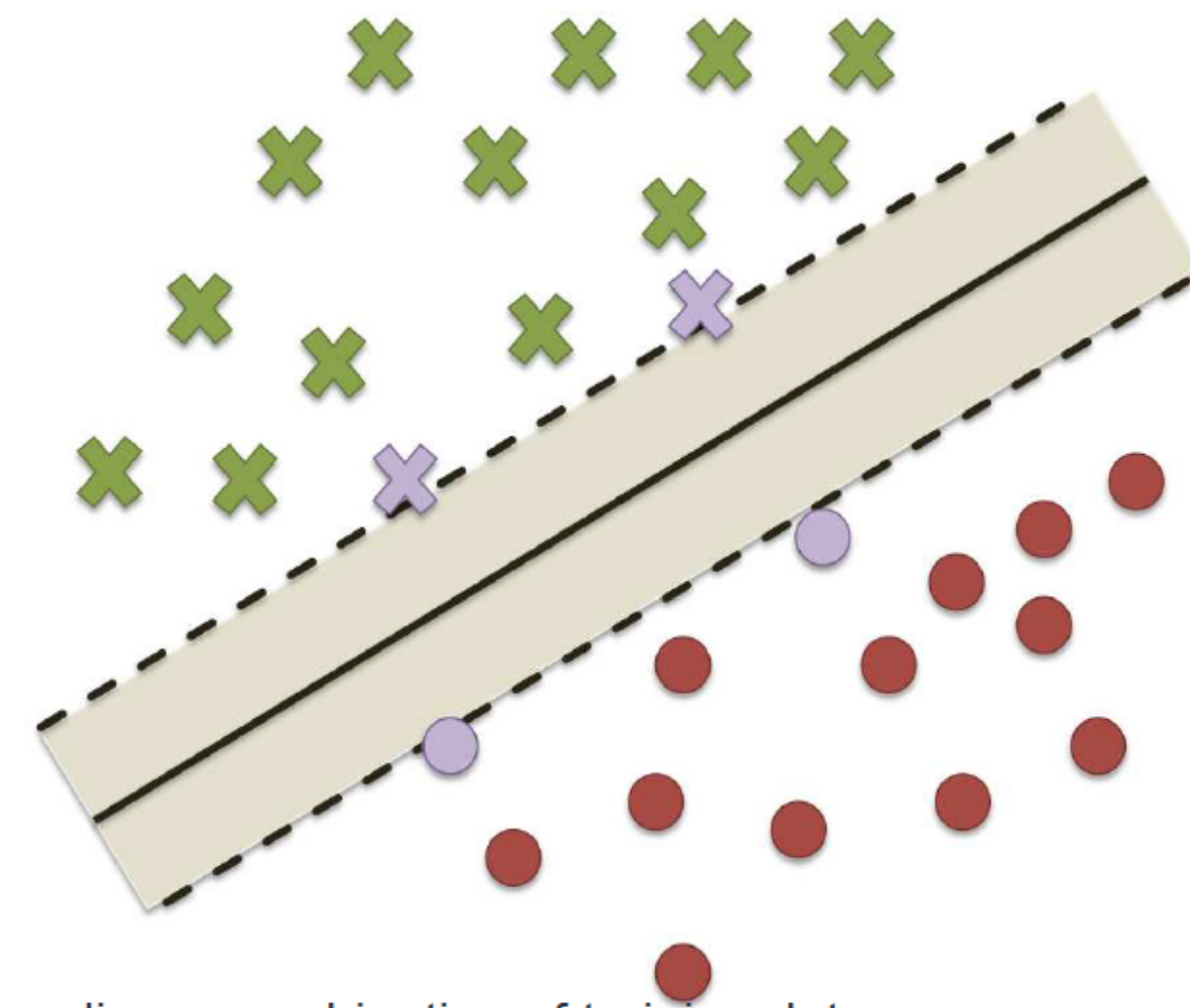
- For classifying with a new input $\mathbf{z}$

  Represents inner product between the two (e.g. support vector and input)

- Compute $\mathbf{w}^T\mathbf{z} + b = \sum_{j=1}^{2} \alpha_{t_j} y^{t_j} < \mathbf{x}^{t_j} \cdot \mathbf{z} > + b$

- Classify $\mathbf{z}$ as positive if the sum is positive, and negative otherwise

- Note: $\mathbf{w}$ need not be formed explicitly, rather we can classify $\mathbf{z}$ by taking a weighted sum of the inner products with the support vectors

  - This is useful when we generalize from inner product to kernel functions later

# Support vectors

- Only points, $\mathbf{x}_i$, that lie on the supporting hyperplanes have $\alpha_i > 0$. These are called the **_support vectors_**. Complexity of the solution only depends on the number of support vectors
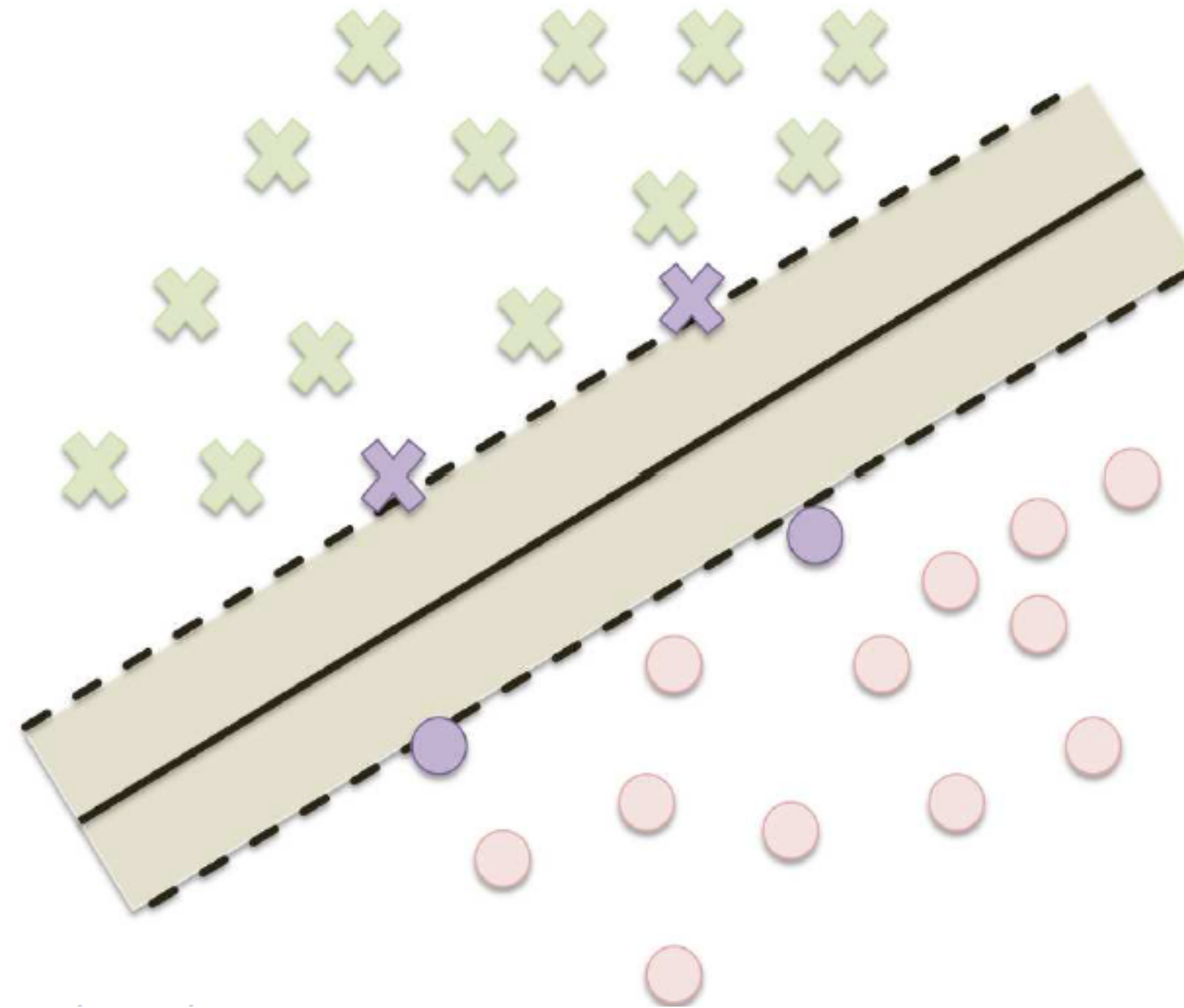
$$\mathbf{w} = \sum_{i=1}^{N} y_i \, \alpha_i \, \mathbf{x}_i = \sum_{\text{support vectors}} y_i \, \alpha_i \, \mathbf{x}_i$$



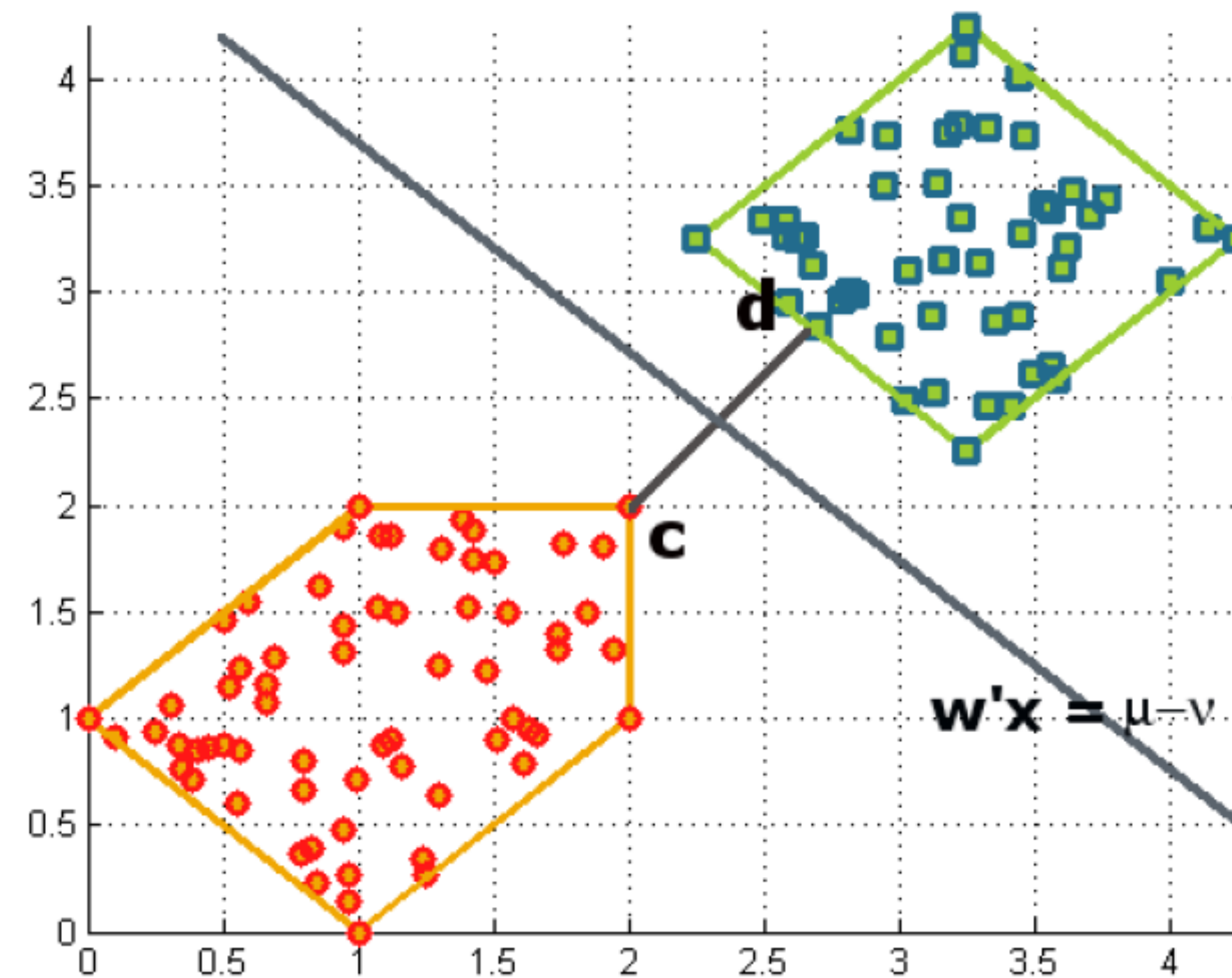Recall that $\mathbf{w}$ is a linear combination of training data

# Support Vectors

- Learned model will not change if we delete all the data (e.g., only the support vectors matter)
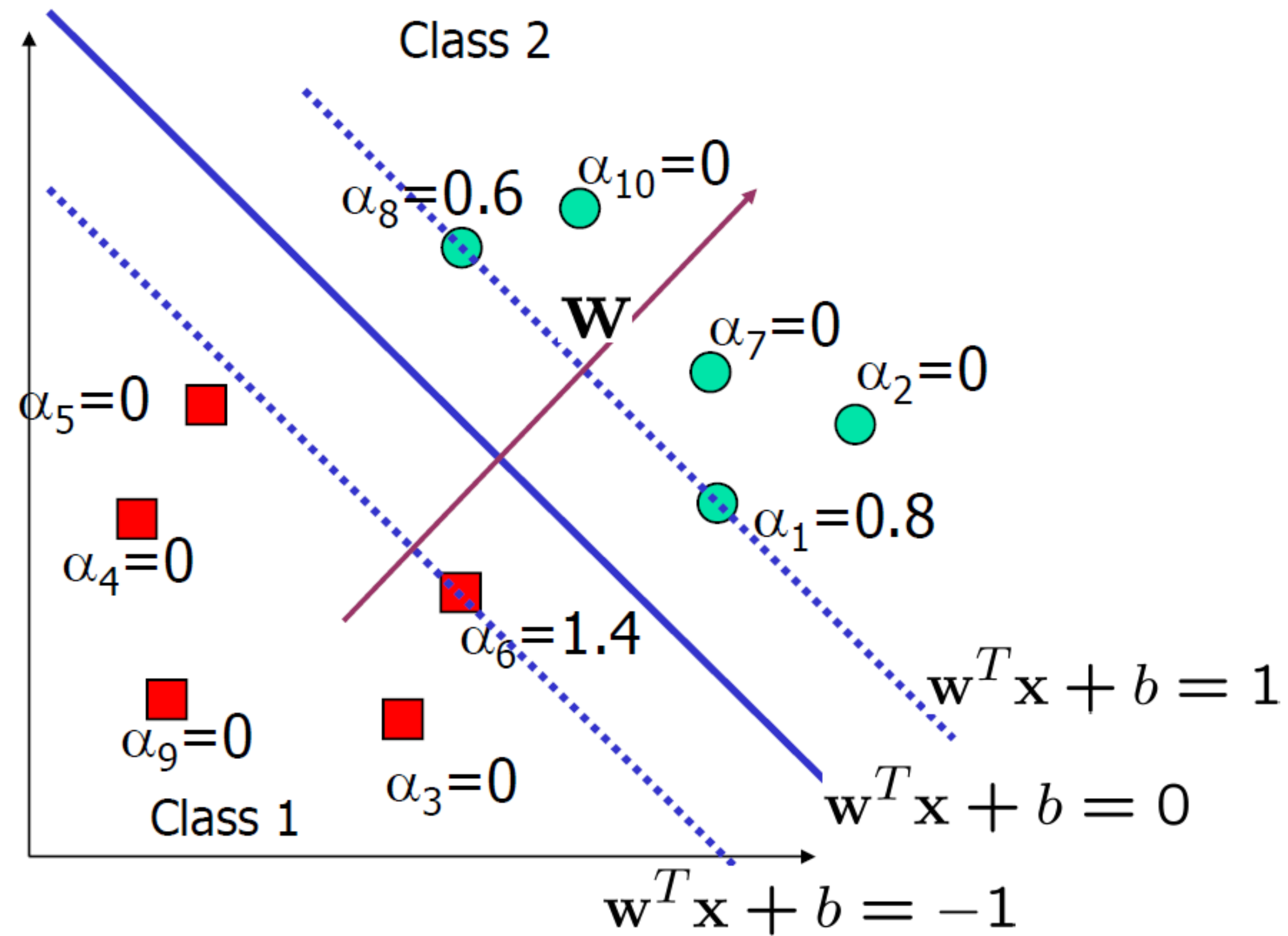
# Geometric Perspective

- Maximizing margin is equivalent to ***maximizing the distance between the two closet points on the convex hulls of the two sets***

# Geometric Perspective (2)

# Summary

- We demonstrated that we prefer to have linear classifiers with large margin.

- We formulated the problem of finding the maximum margin linear classifier as a quadratic optimization problem

- This problem can be solved by solving its dual problem, and efficient QP algorithms are available.

- Problem solved?

- How about non-linear data? – Kernels

- How about noise? – Soft Margin SVMs

# Next Class

**Support Vector Machines II:** soften assumptions