

CS 360: Machine Learning

Sara Mathieson, Sorelle Friedler

Spring 2024



HAVERFORD
COLLEGE

Admin

- **Project meetings today** in lab with all groups
 - Try to come to the same lab as your group
- **Midterm April 25** in class (this Thursday)
 - Can still do handout videos for extra credit! Up to 24 hours before exam
 - Create your own “study sheet” (front and back) to use during the exam
 - No other notes or resources

Faculty Interview TODAY!

- The candidate joining us is Dakota Lambert. Dakota received their PHD from Stony Brook University in the Department of Linguistics and Institute for Advanced Computational Science, with a research focus on understanding communication structures between humans and computers using algebraic and logical analysis.
- Date: **23rd April 2024** (Today)
- Research talk: **4:30 pm** onwards in **H109**. The research talk will be preceded by a tea at around **4:15 pm**
- Title: **Finite automata: interpretation and normal forms**
- Abstract: In this lesson we will explore finite-state automata, a standard representation of abstract machines that implement regular expressions. These are commonly used, perhaps with extensions, in compiler design, language modeling, and machine learning. We will see how to read standard depictions of these abstract machines, to understand which languages they represent, and to determine whether two machines implement the same functionality.

Outline for April 23

- Review temperature, cross-entropy, loss functions, and ML pipeline
- Convolutional filters, backpropagation, and CNNs
- Gaussian mixture model handout
- Perceptron and SVM

Outline for April 23

- Review temperature, cross-entropy, loss functions, and ML pipeline
- Convolutional filters, backpropagation, and CNNs
- Gaussian mixture model handout
- Perceptron and SVM

Review temperature in language generation

Handout 21, Q1

- Temperature = 1
 - Scores (logits) $[-0.6931472, -0.9162907, -2.3025851]$
 - Probabilities $[0.5, 0.4, 0.1]$
- Temperature = ~~0.01~~ **100 (corrected!)**
 - Scores (logits) $[-0.00693147, -0.00916291, -0.02302585]$
 - Probabilities $[0.33536728108, 0.3346197634, 0.3300129554]$
- Temperature = ~~100~~ **0.01 (corrected!)**
 - Scores (logits) $[-69.31472, -91.629074, -230.25851]$
 - Probabilities $[0.9999999, 2.0370382e-10, 1.26765211e-70]$

Temperature: don't always pick the letter with maximum probability

```
>>> print(extend_text("To be or not to be", temperature=0.01))  
To be or not to be the duke  
as it is a proper strange death,  
and the
```

```
>>> print(extend_text("To be or not to be", temperature=1))  
To be or not to behold?  
  
second push:  
gremio, lord all, a sistermen,
```

```
>>> print(extend_text("To be or not to be", temperature=100))  
To be or not to bef ,mt'&o3fpadm!$  
wh!nse?bws3est--vgerdjw?c-y-ewznq
```

Q: is `tf.argmax` different from `temp=1`?

Core Machine Learning Pipeline

```
def train_step(model, images, labels, loss_func, optimizer):  
    with tf.GradientTape() as tape:  
1        predictions = model(images)  
2        loss = loss_func(labels, predictions)  
3    gradients = tape.gradient(loss, model.trainable_variables)  
4    optimizer.apply_gradients(zip(gradients, model.trainable_variables))  
  
    return loss, predictions
```


linear regression

① $\hat{y} = \vec{w} \cdot \vec{x}$ features \rightarrow squared error (MSE)

② $J(\vec{w}) = \frac{1}{2}(\hat{y} - y)^2 = \frac{1}{2}(\vec{w} \cdot \vec{x} - y)^2$

③ $\nabla J(\vec{w}) = (\vec{w} \cdot \vec{x} - y)\vec{x}$

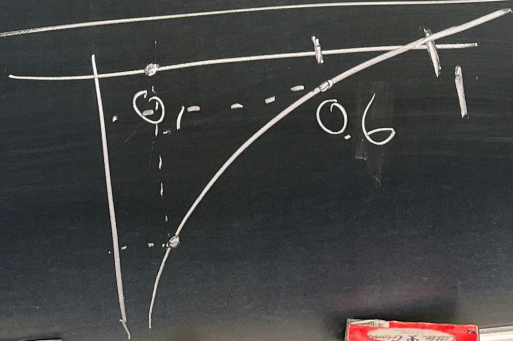
$y=1$ dog

$y=2$ cat

$y=3$ bird

④ $\vec{w} \leftarrow \vec{w} - \alpha \nabla J(\vec{w})$

cat $\Rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$



Stochastic Gradient Descent (high-level)

```
set w = 0 vector
while cost J(w) still changing (or max iter reached):
    shuffle data points
    for i = 1...n:
        w <- w - alpha(derivative of J(w) wrt xi)
    store J(w)
```


Gradient descent variations

- **Batch gradient descent**
 - Go over *all* training data before making a weight update
- **Stochastic gradient descent**
 - Shuffle data and make a weight update after *each training example*
- **Mini-batch gradient descent**
 - Update after a “*mini-batch*” of training examples (i.e. 50)
 - Vocab word: **epoch** (one pass through the data)

3 important pieces to SGD

(with logistic regression concrete example)

- Hypothesis function (prediction)

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- Cost/loss function (want to minimize)

$$J(\mathbf{w}) = - \sum_{i=1}^n y_i \log h_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))$$

- Gradient of cost wrt single data point \mathbf{x}_i

$$\nabla J_{\mathbf{x}_i}(\mathbf{w}) = (h_{\mathbf{w}}(\mathbf{x}_i) - y_i) \mathbf{x}_i$$

binary cross
entropy

(2 classes)

$$H(y, \hat{y}) = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

$$y \in \{0, 1\}$$

$$\hat{y} \in [0, 1] \text{ "prob pos"}$$

categorical cross
entropy

(K classes)

ci-far $K=10$

$$H(\vec{y}, \hat{\vec{y}}) = - \sum_{k=1}^K y_k \log_2 \hat{y}_k$$

$$-1 \cdot \log_2(0.6)$$

$$\vec{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\hat{\vec{y}} = \begin{bmatrix} 0.2 \\ 0.6 \\ 0.2 \end{bmatrix}$$

$$\hat{\vec{y}} = \begin{bmatrix} 0.8 \\ 0.1 \\ 0.1 \end{bmatrix}$$

Sparse categorical
cross entropy

$$y=2$$

example

$$y \in \{1, 2, \dots, K\}$$

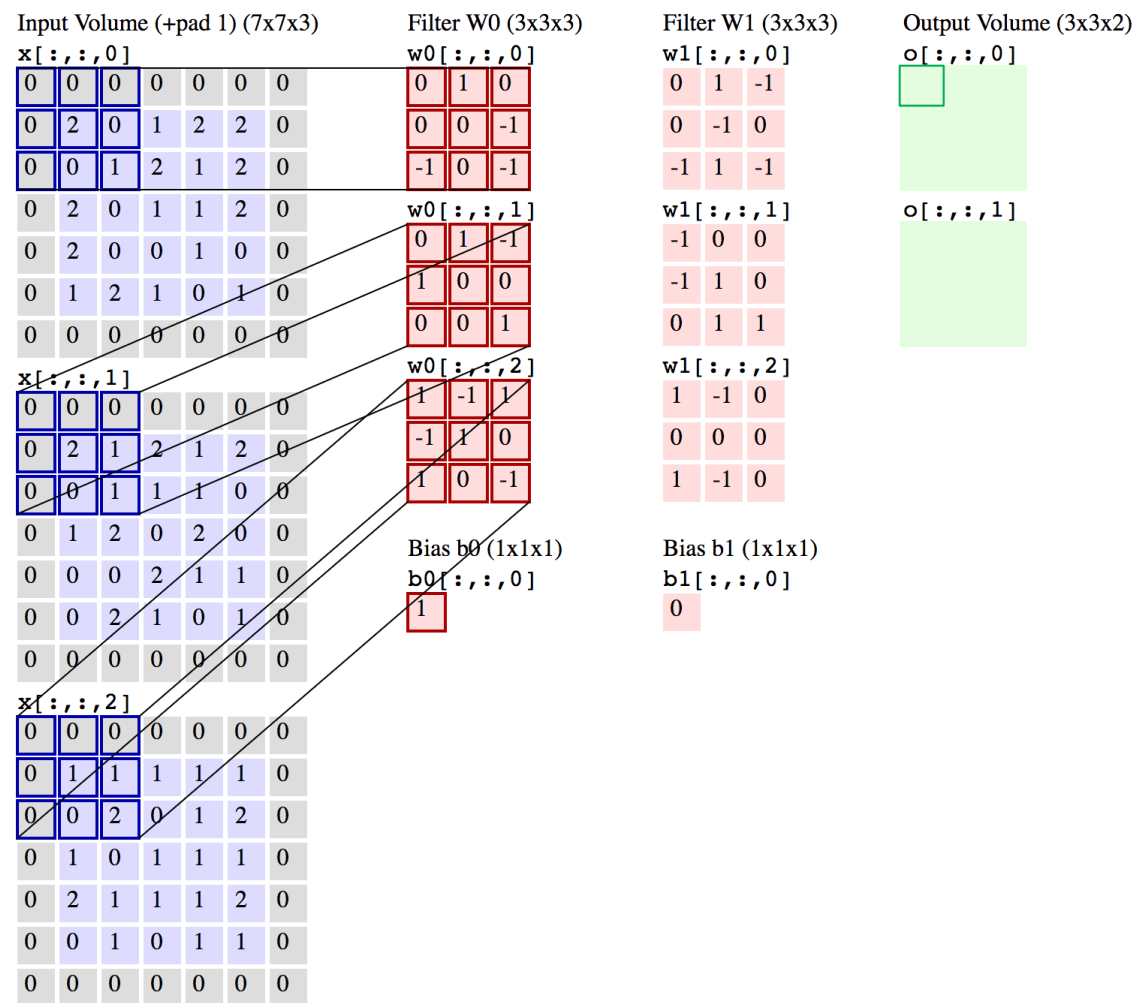
general

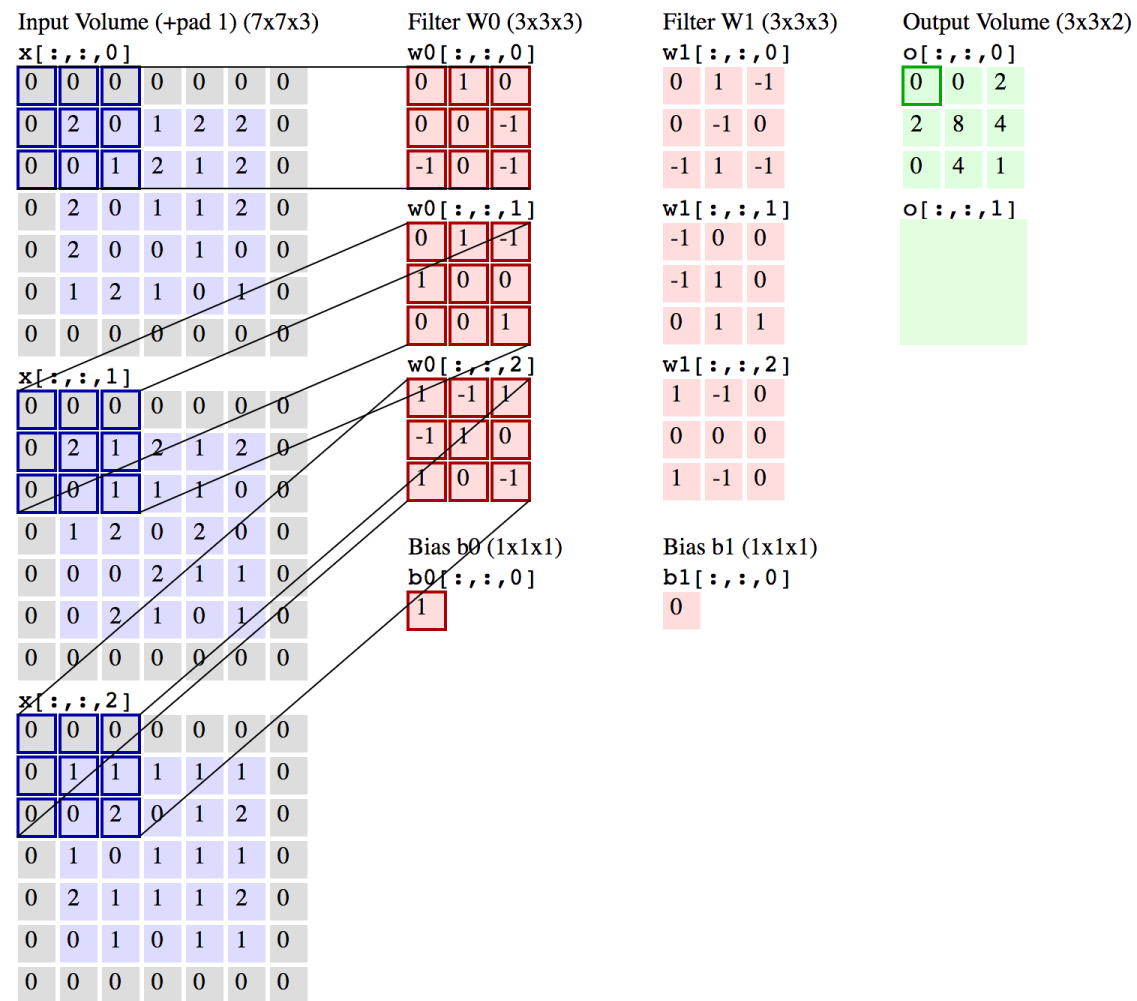
$$-\log \hat{y}_k \text{ if } y=k$$

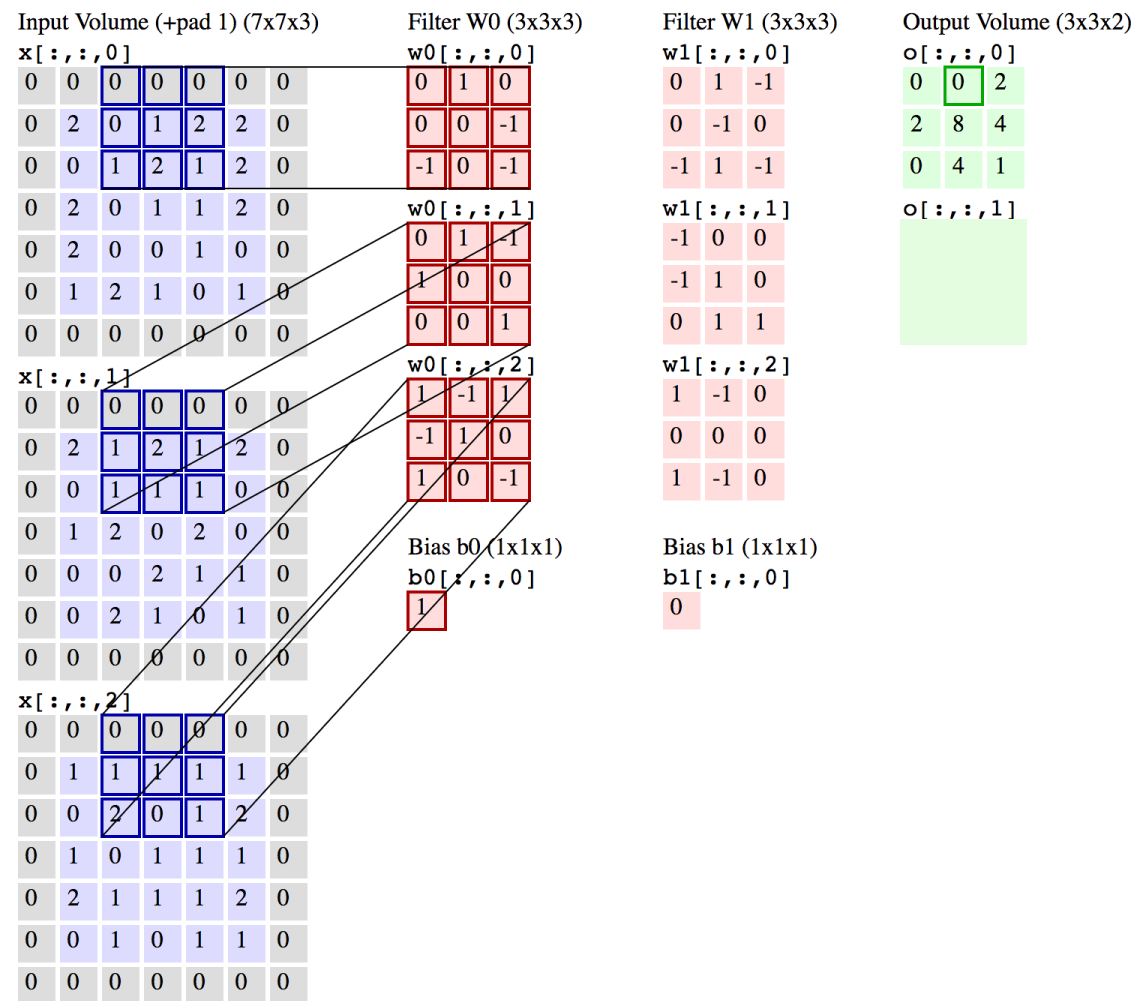
Outline for April 23

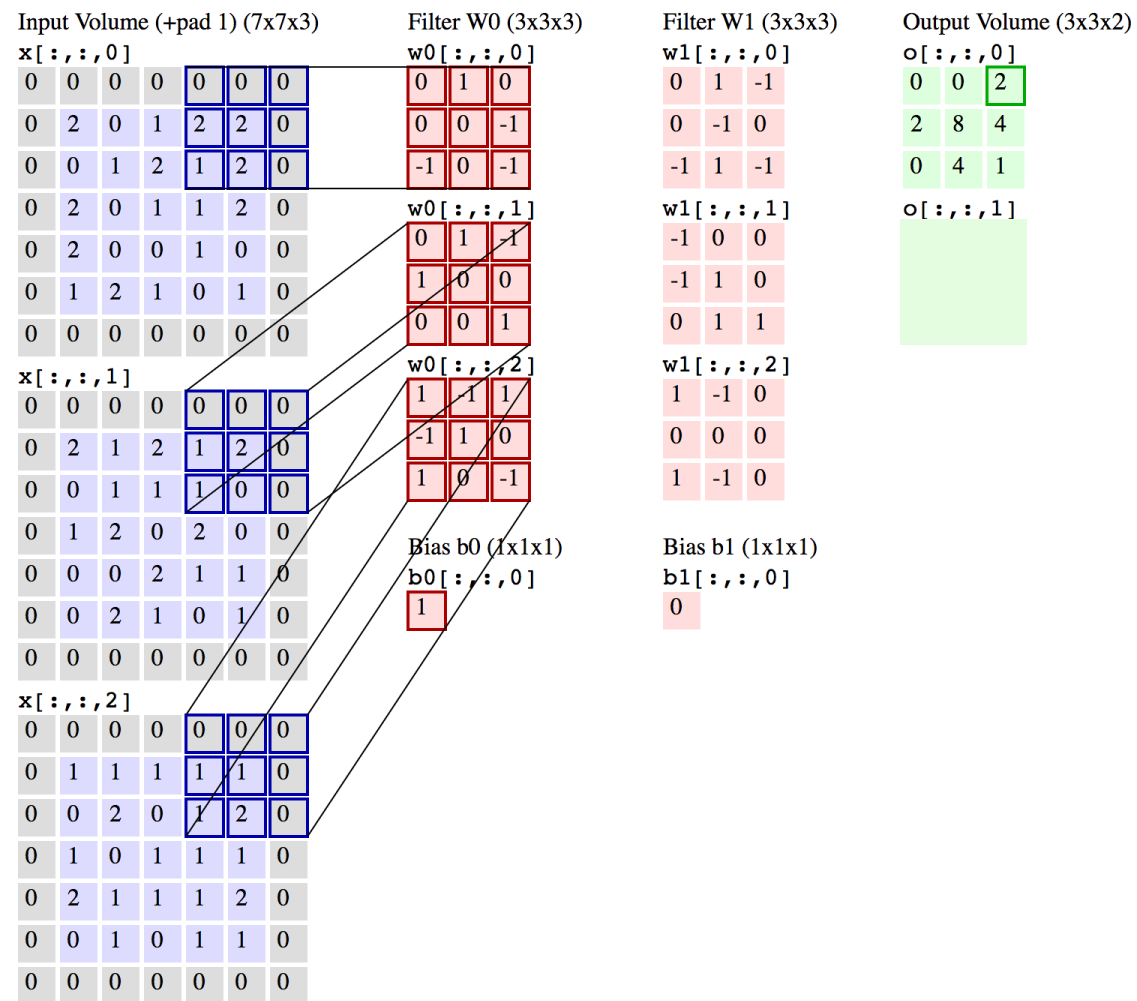
- Review temperature, cross-entropy, loss functions, and ML pipeline
- Convolutional filters, backpropagation, and CNNs
- Gaussian mixture model handout
- Perceptron and SVM

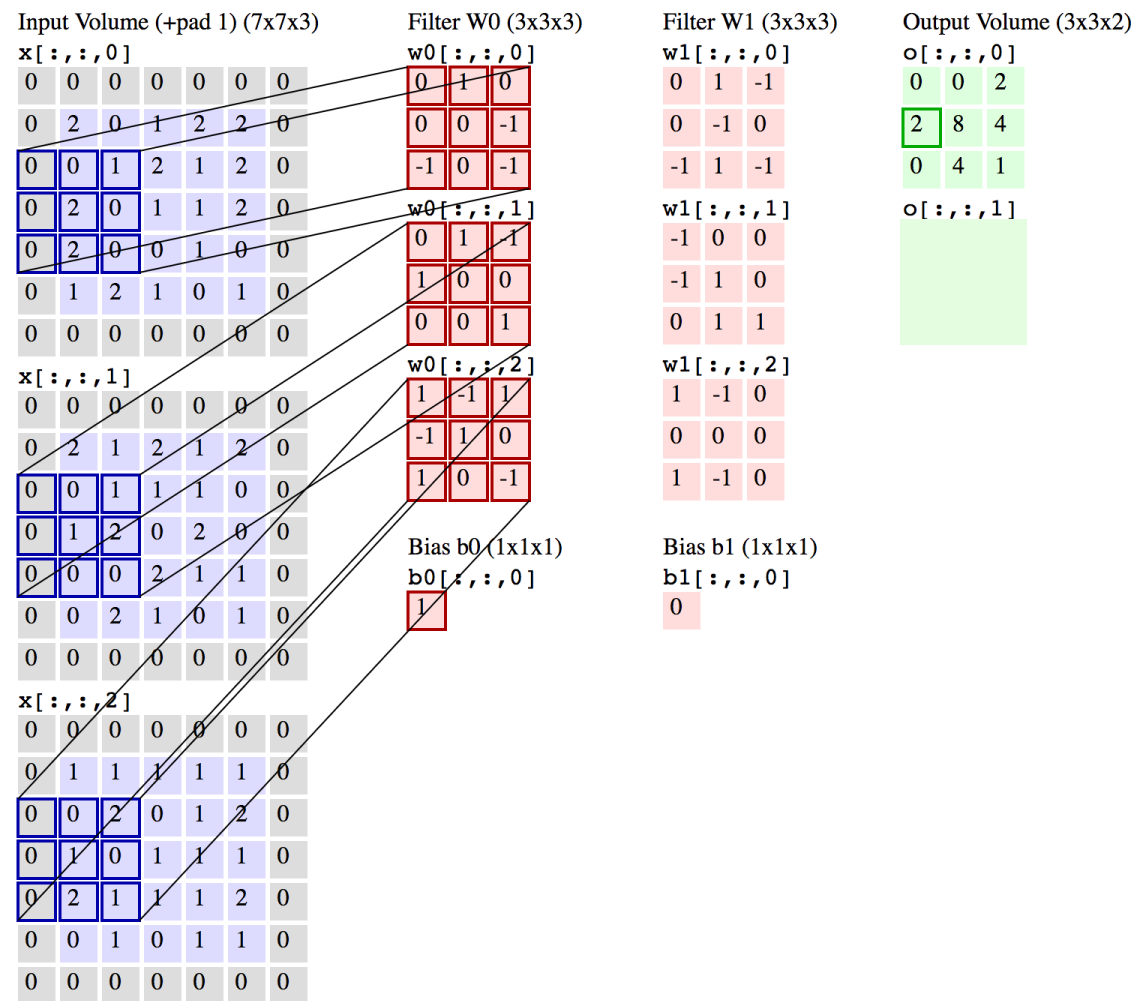
Handout 25, page 1 and 2!

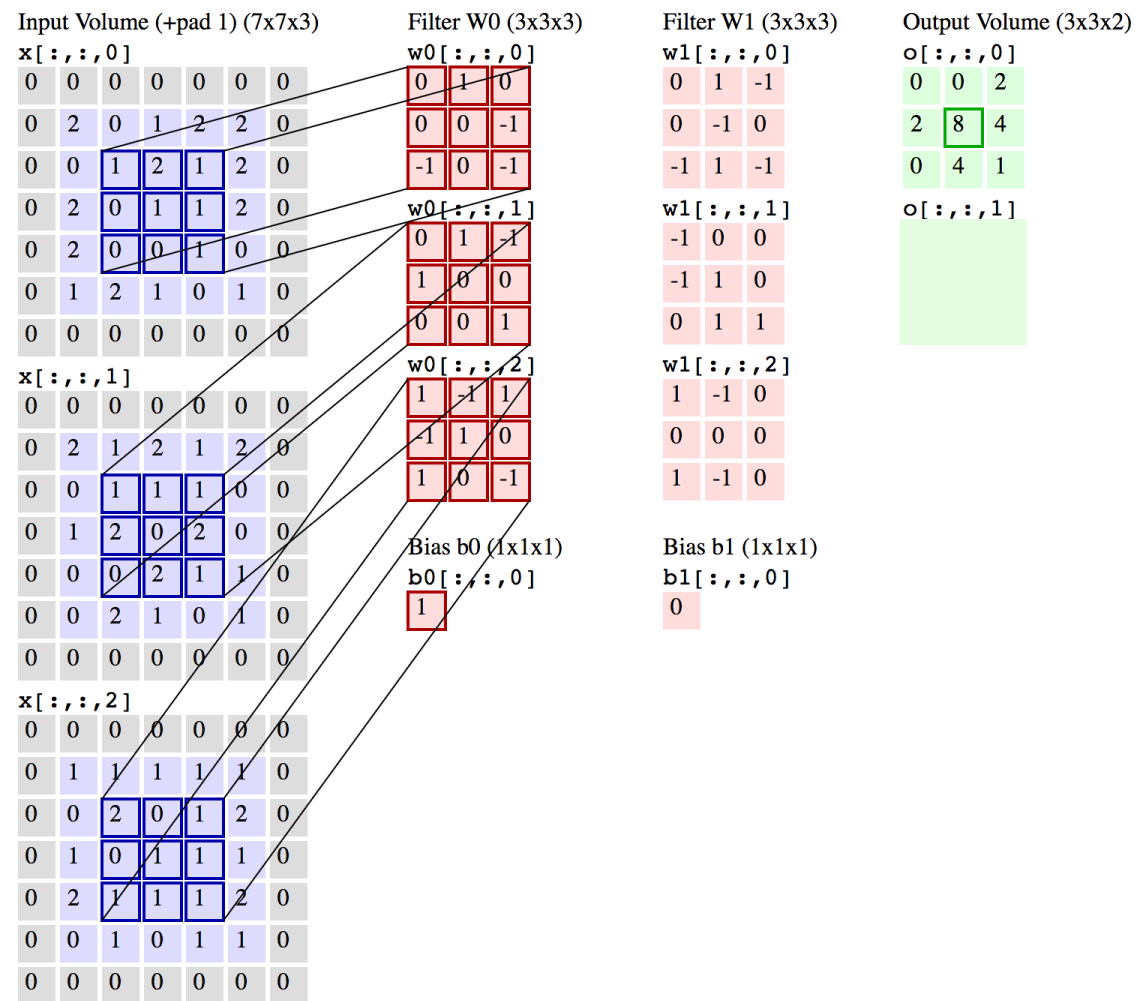


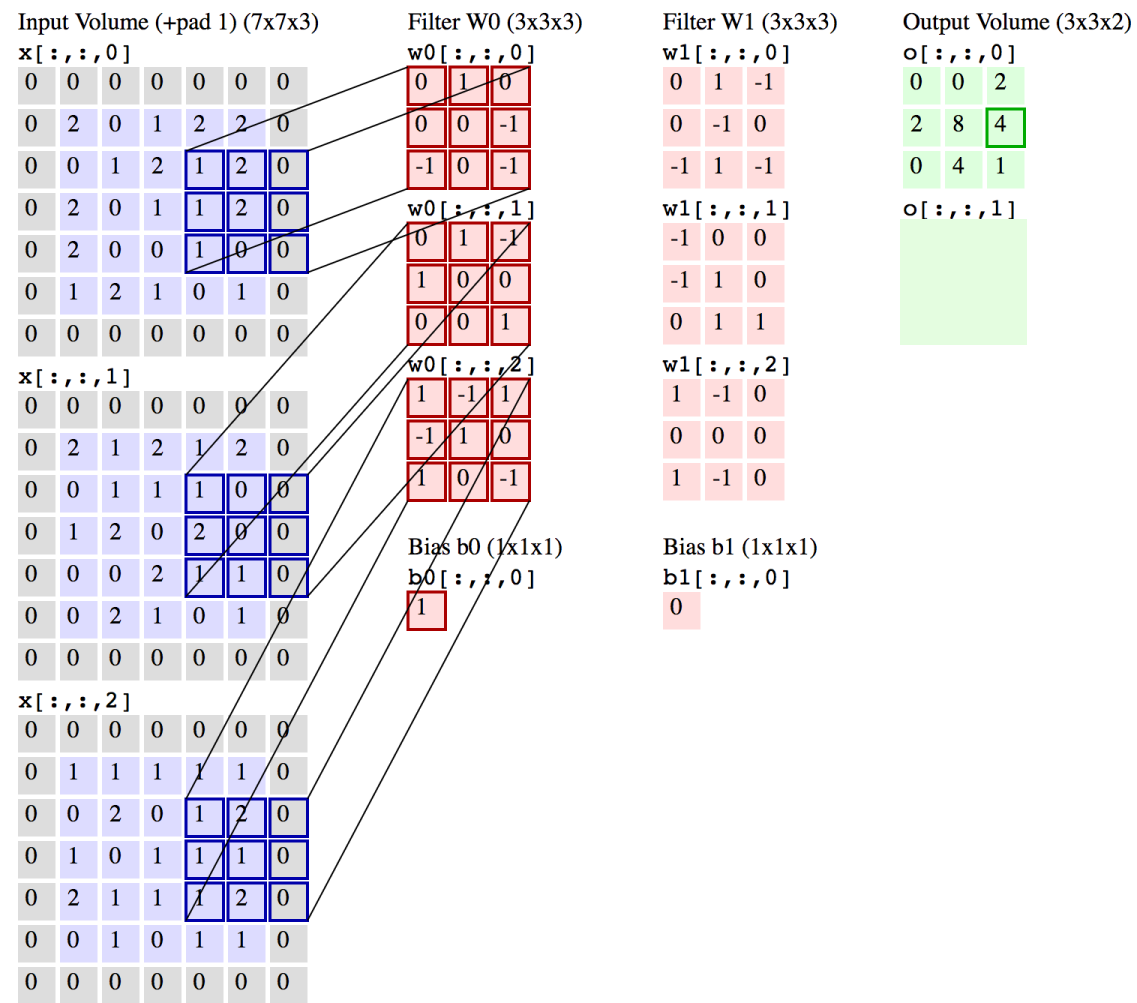


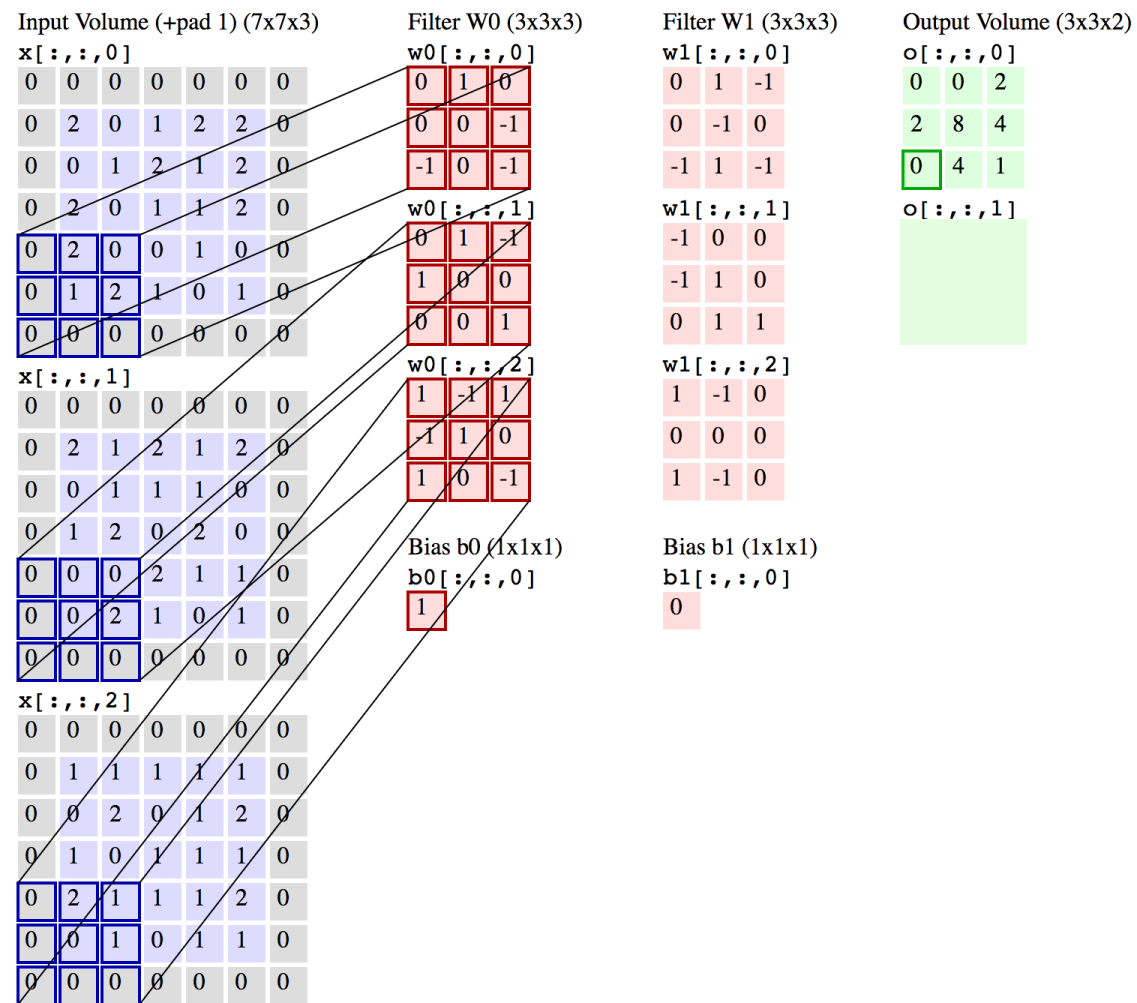


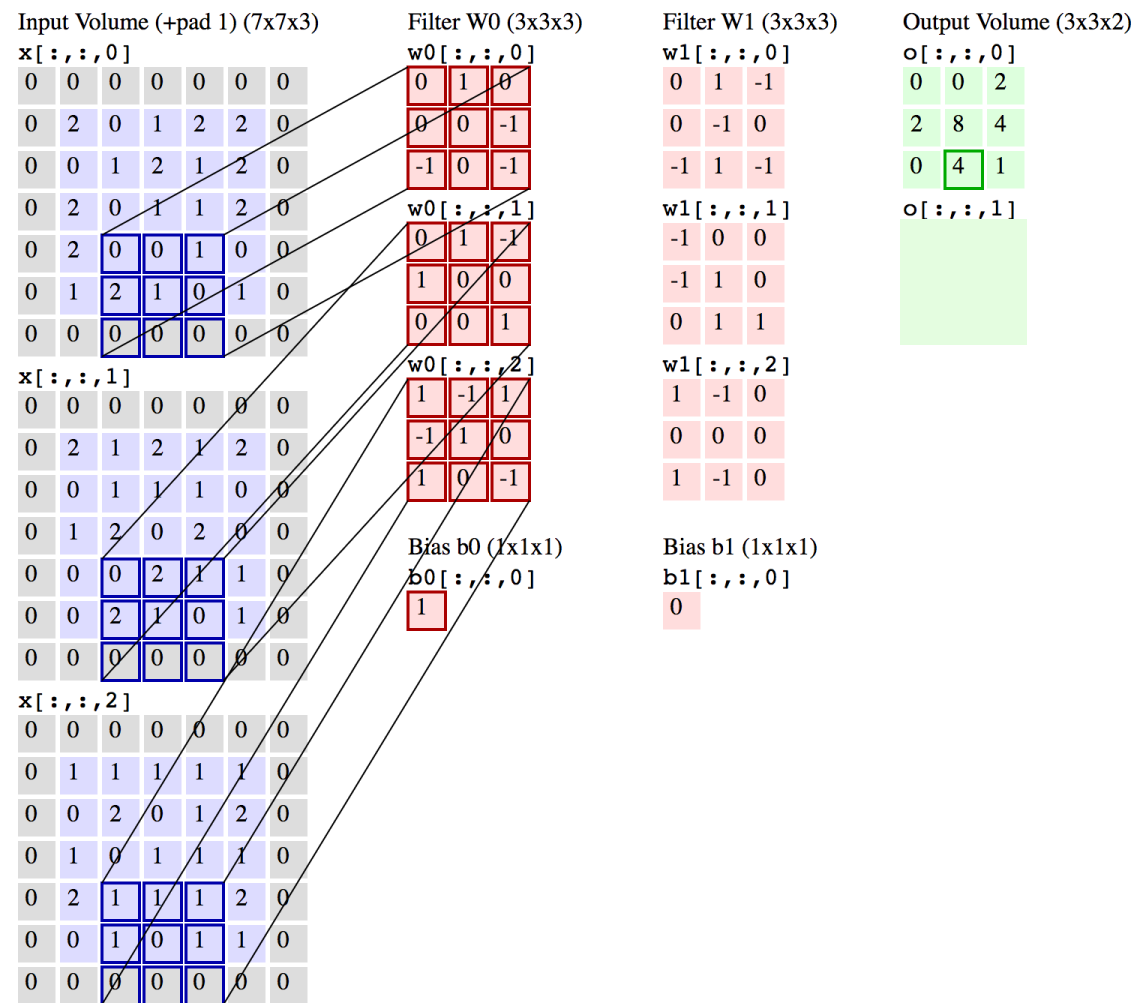


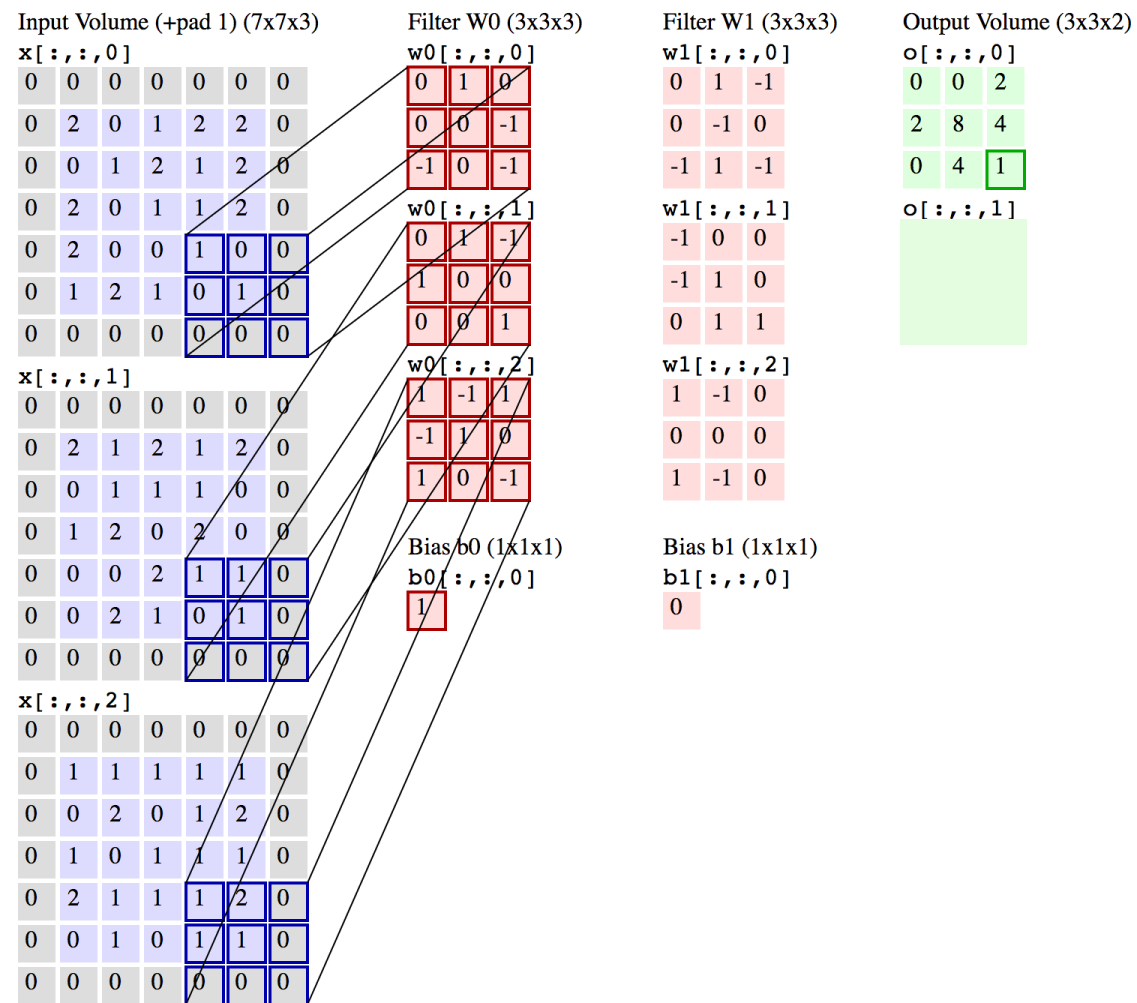


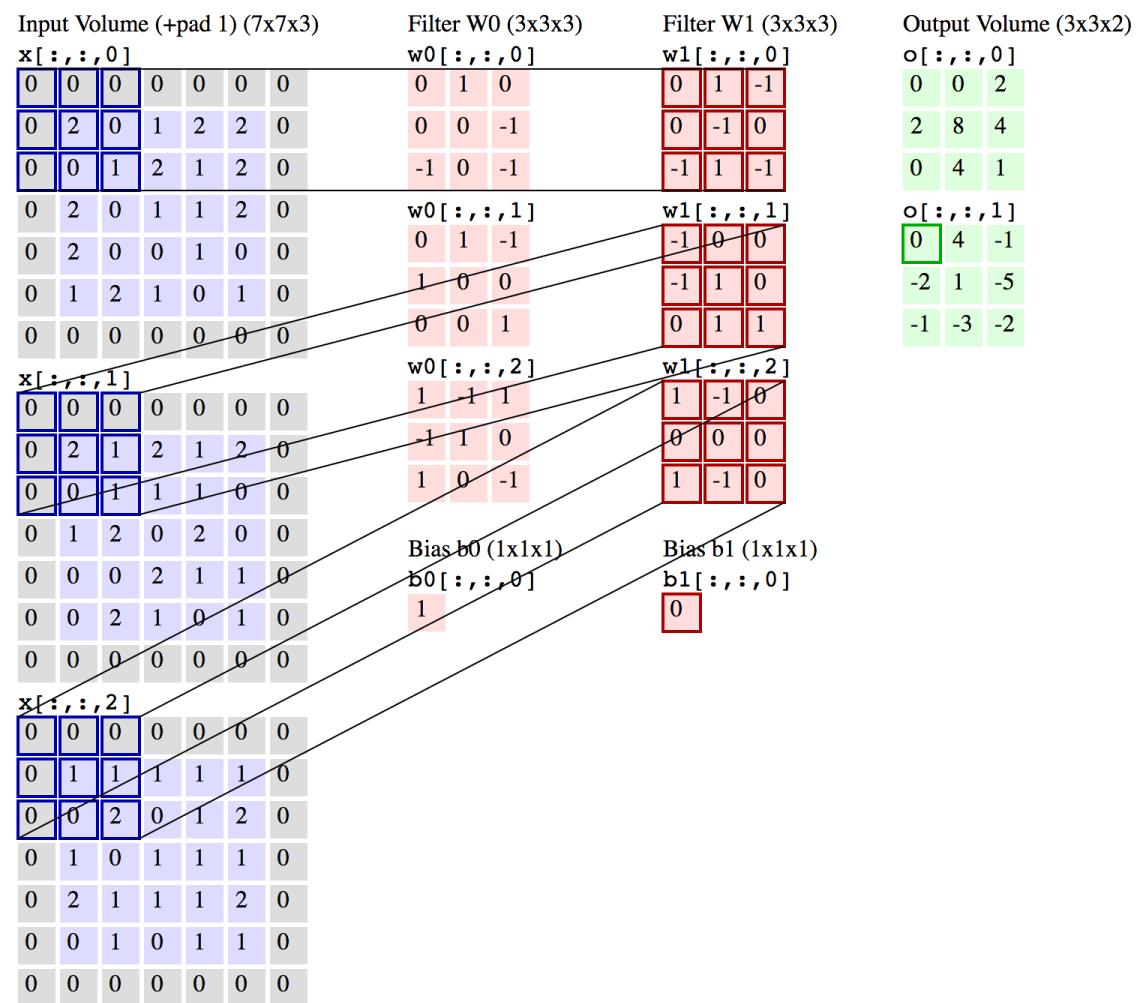


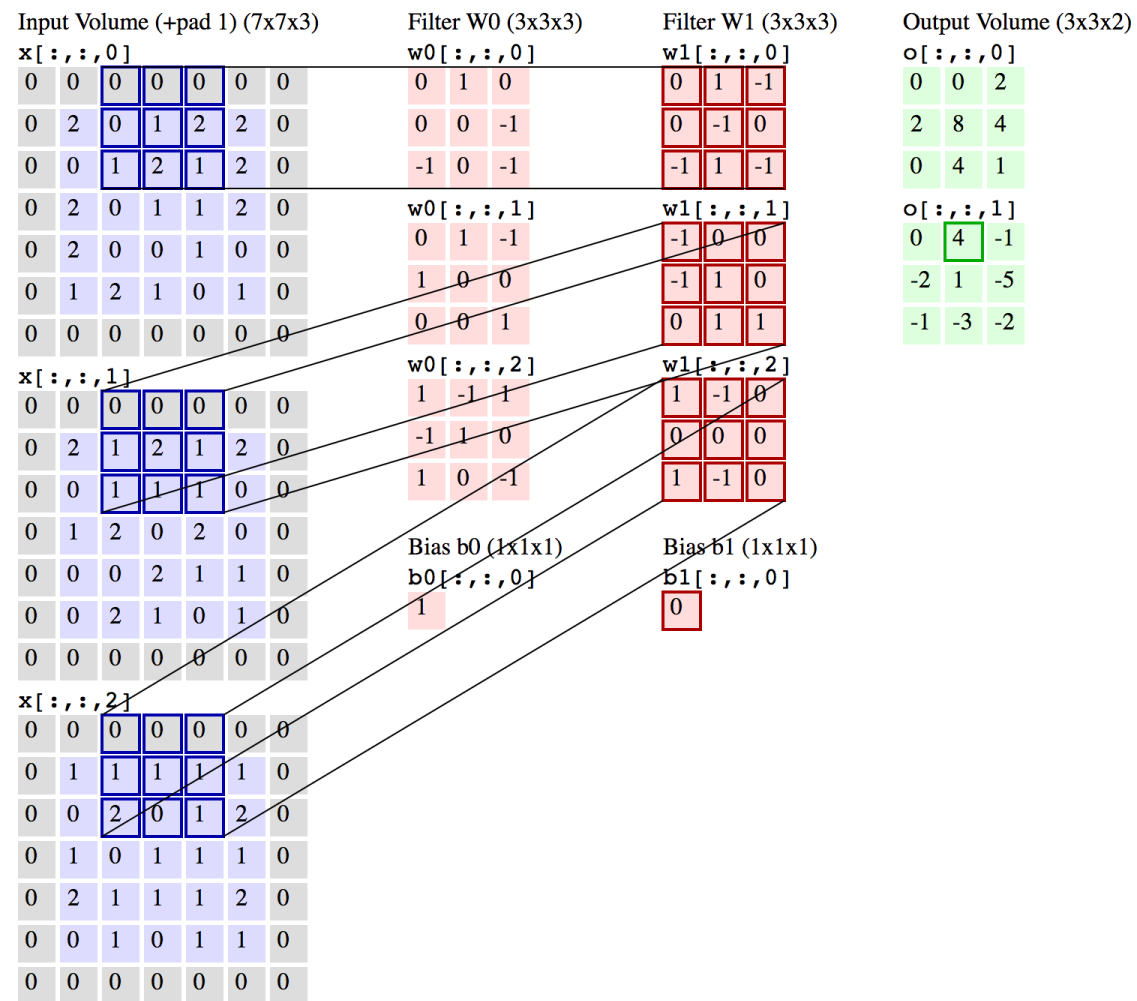


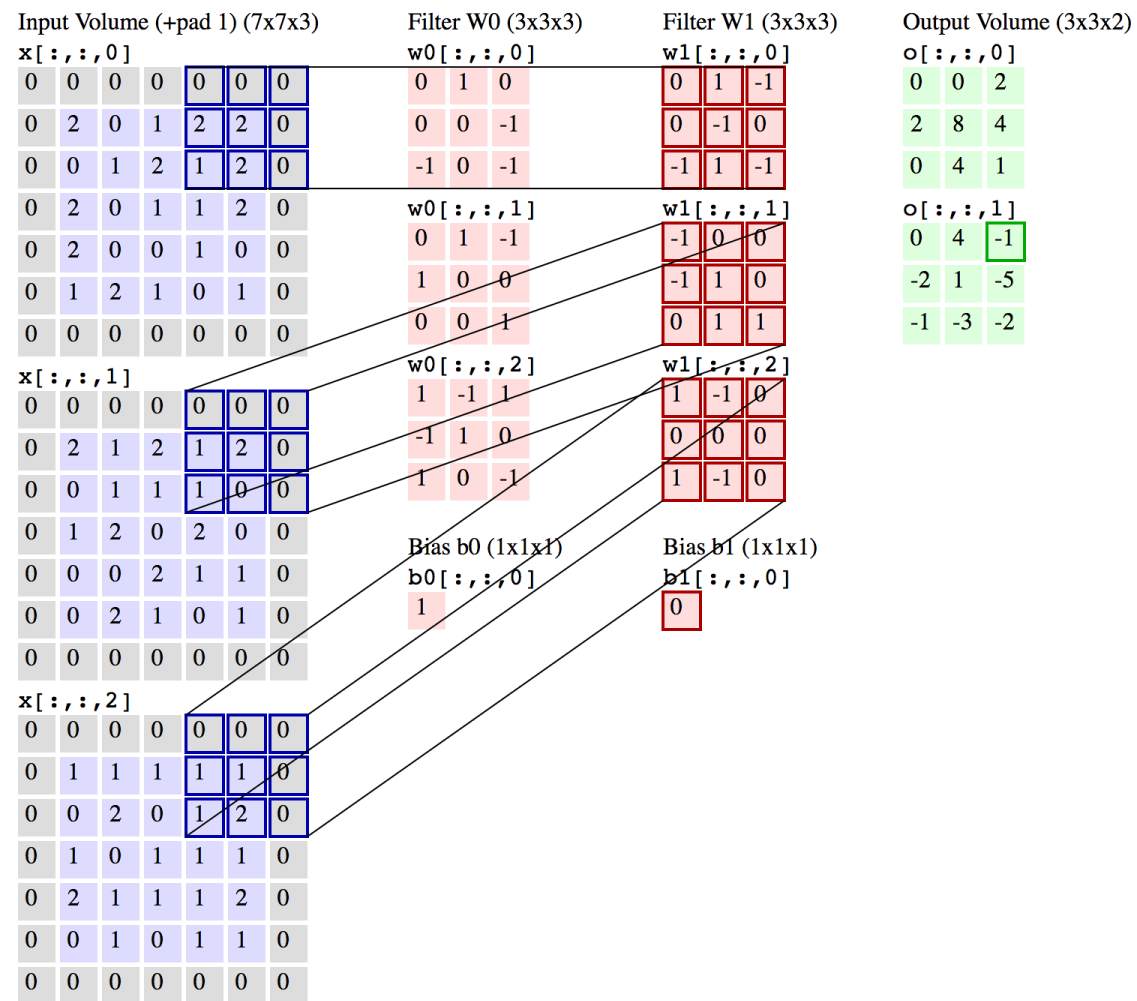


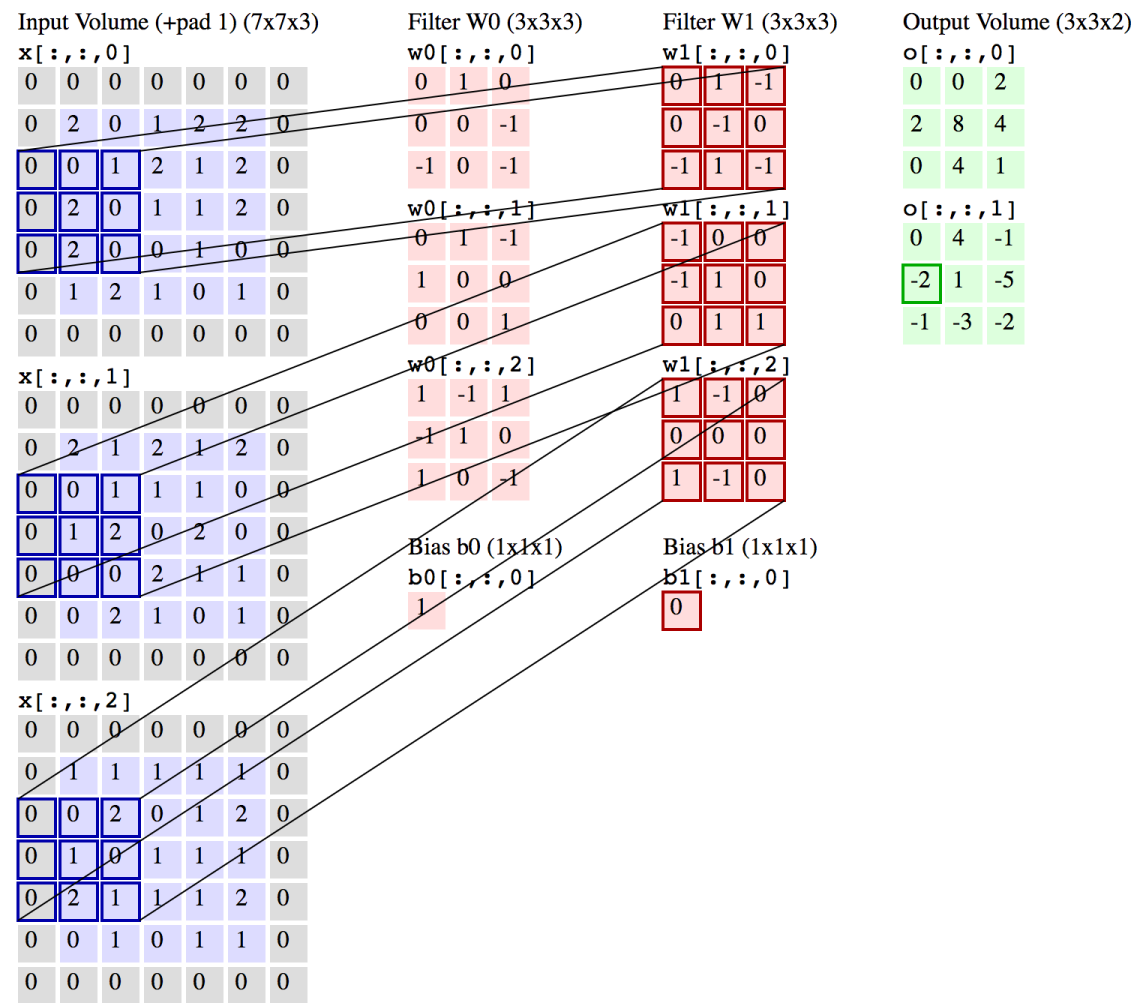


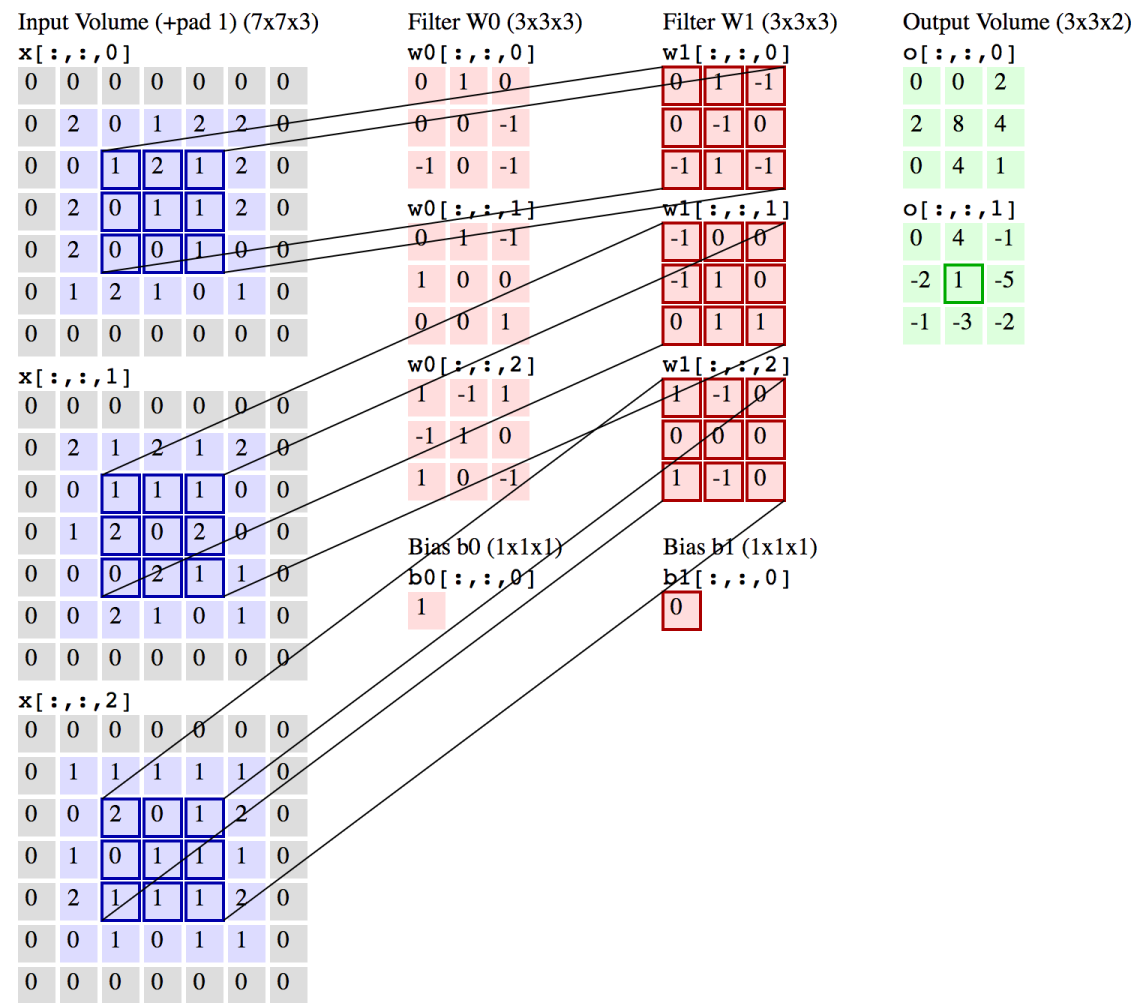


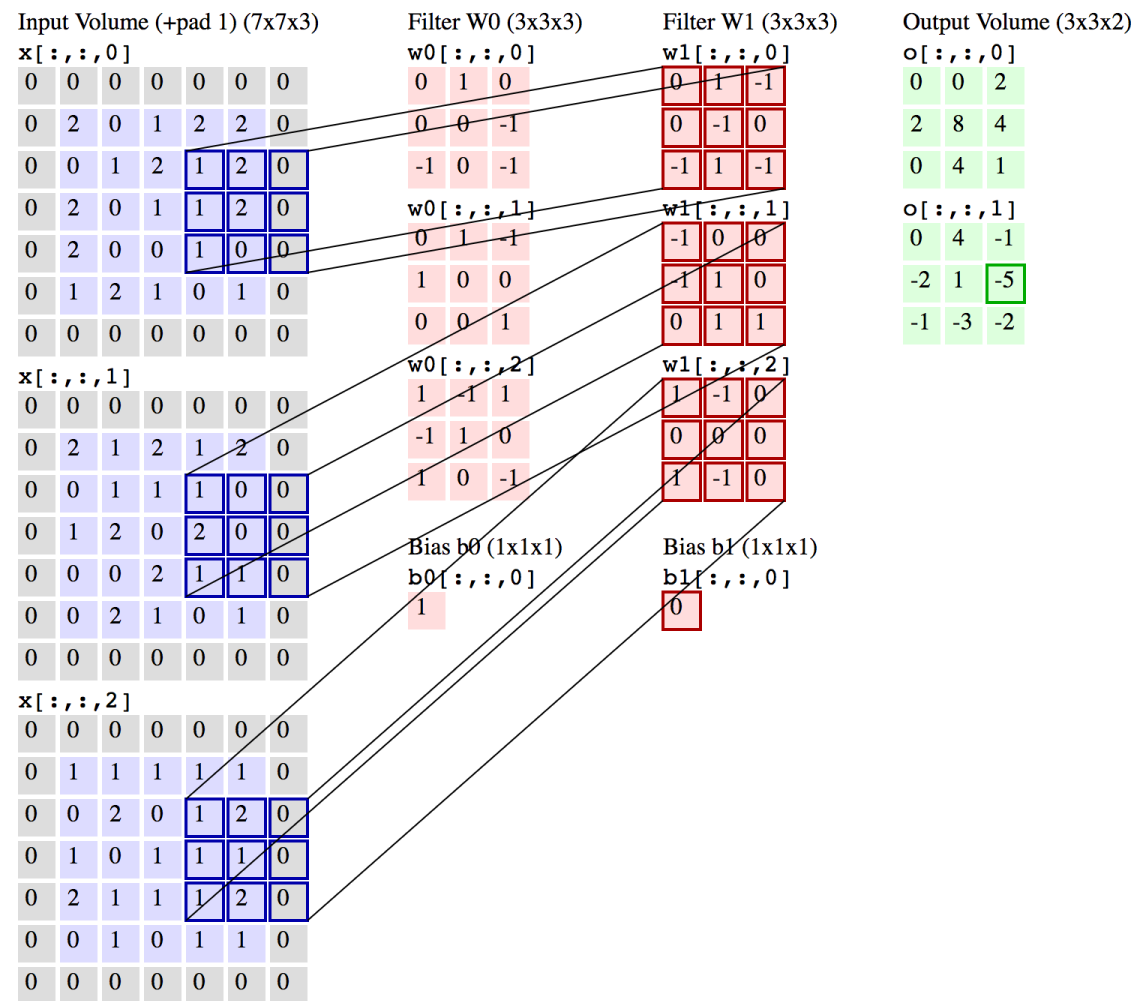


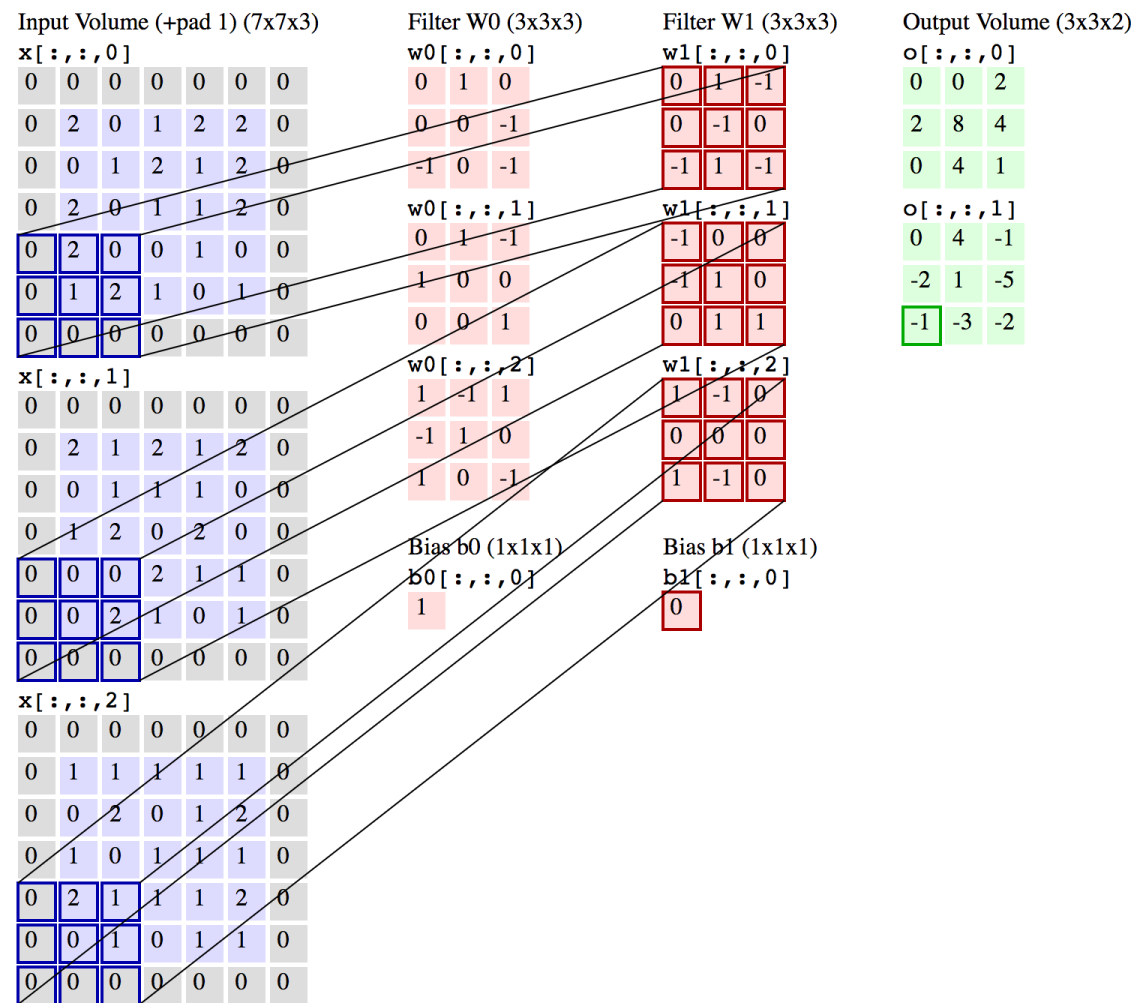


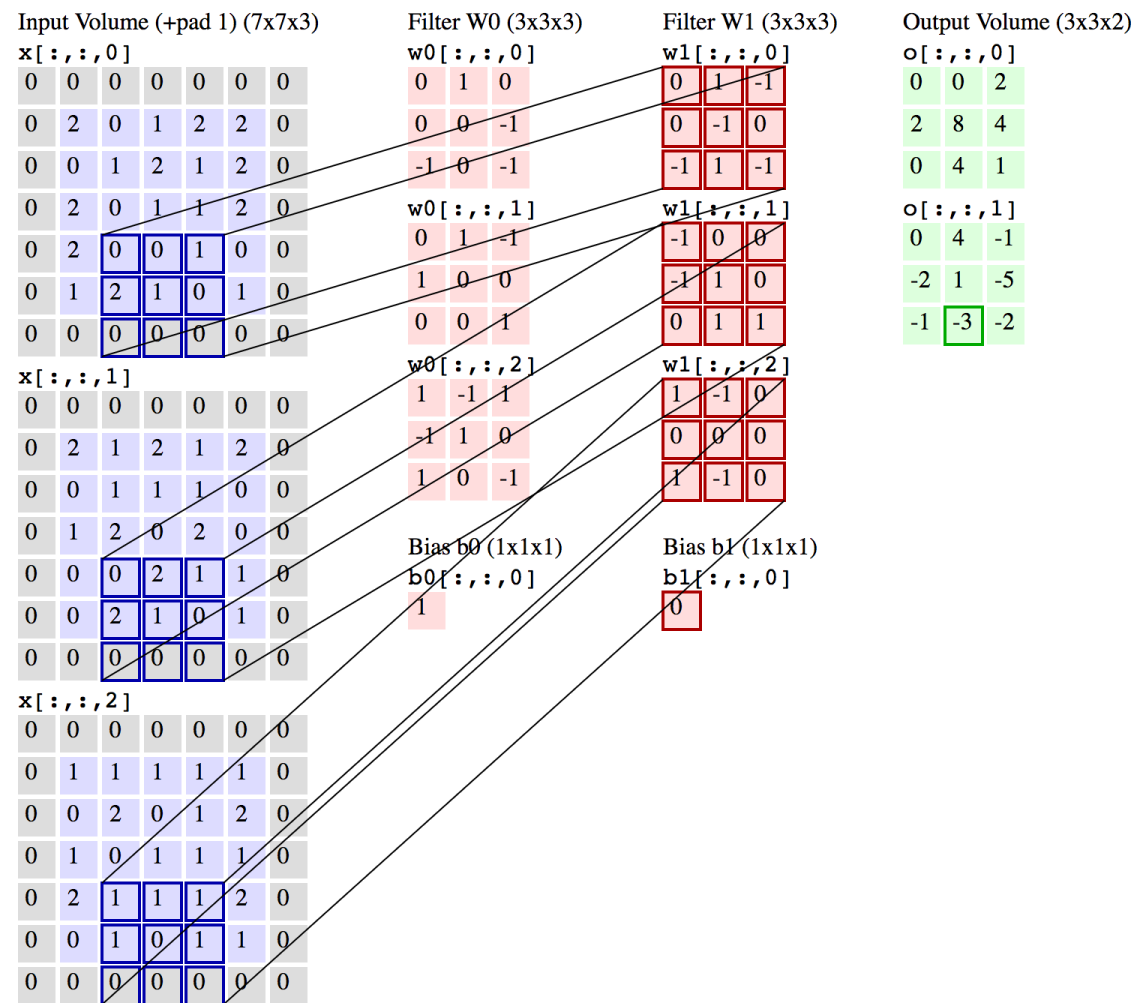


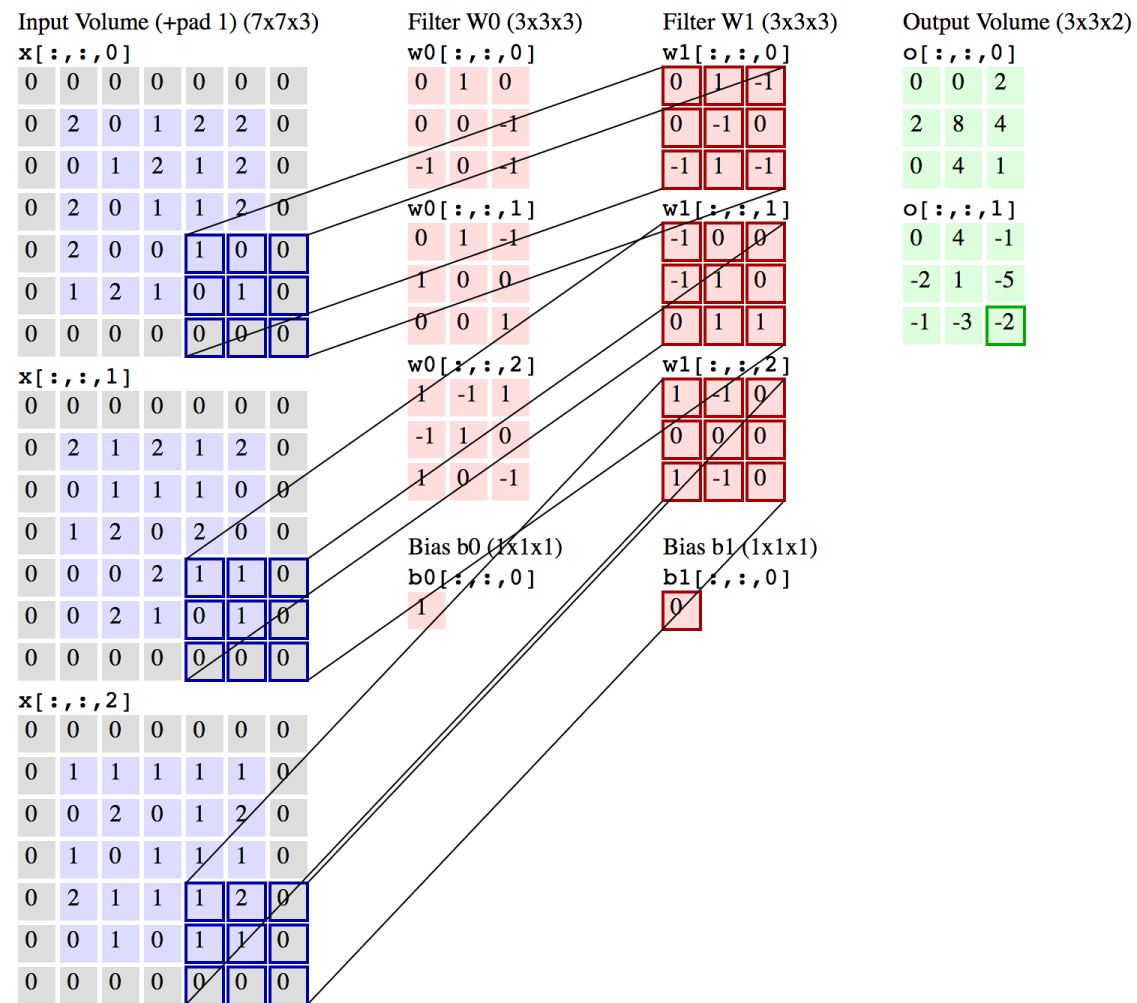












$$w_1^{(new)} = 1 - \left(-\frac{3}{4}\right)$$

$$w_2^{(new)} = 5 - \left(\frac{1}{4}\right)$$

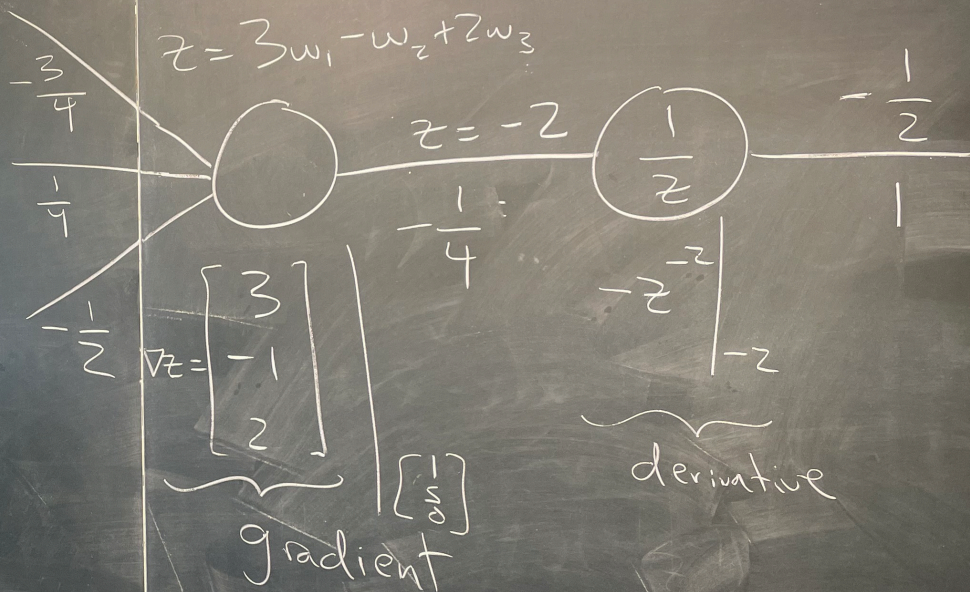
$$w_3^{(new)} = 0 - \left(-\frac{1}{2}\right)$$

$$w_1 = 1$$

$$w_2 = 5$$

$$w_3 = 0$$

$$z = 3w_1 - w_2 + 2w_3$$



#3

CONV →

0	-2	-4	-8
1	7	-3	-1
-1	5	2	0
-7	3	2	1

max
POOL →

7	-1
5	2

RELU →

7	0
5	2

flatten

7
0
5
2

fc

w_{ij}

(Scores
logits)

-1
3
0

Softmax

0.017
0.936
0.046

3 classes

want small

7 - 2 - 1 - 1

$$-1 + 1 - 1 + 1 = 0$$

$$-4 + 2 + 1 + 1 = 0$$

$$\frac{5 - 3 + 2 \cdot 1}{2} + 1$$

3

$$\frac{e^{-1}}{e^{-1} + e^3 + e^0}$$

all zeros
⇒

0.5
0.5
0.5

params $3 \cdot 3 = 9$

w/ bias = 10

Output size formula in one dimension

- W = input width (same for height though, but not for depth!)
- F = filter size
- P = padding (on one side)
- S = strides

formula for output size:

$$\frac{W - F + 2P}{S} + 1$$

Handout 18

(a) Which steps require parameter learning? (out of CONV, RELU, POOL, FLATTEN, FC)

CONV, FC

(b) First layer params $5*5*3*20 + 20 = 1520$

(c) Second layer params $3*3*20*10 + 10 = 1810$

(d) Third layer params $8*8*10*10 + 10 = 6410$

(e) Total # params 9740

If we had a FC with $p_1=100$ and $p_2=50$, we would have 312,860 params to learn (check this after class). CNN is much better!

Outline for April 23

- Review temperature, cross-entropy, loss functions, and ML pipeline
- Convolutional filters, backpropagation, and CNNs
- Gaussian mixture model handout
- Perceptron and SVM

Handout 24

$$W = \begin{bmatrix} 0.2 & 0.5 & 0.3 \\ 0.7 & 0.1 & 0.2 \\ 0.1 & 0.4 & 0.5 \\ 0.8 & 0.1 & 0.1 \end{bmatrix} \begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{matrix}$$

cluster
1

cluster
2

cluster
3

$$\begin{matrix} n=4 \\ K=3 \end{matrix}$$

$$M_k = \sum_{i=1}^n w_{ik}$$

$$M_1 = 1.8$$

$$\Rightarrow \pi_1 = \frac{1.8}{4}$$

$$\pi_2 = \frac{1.1}{4}$$

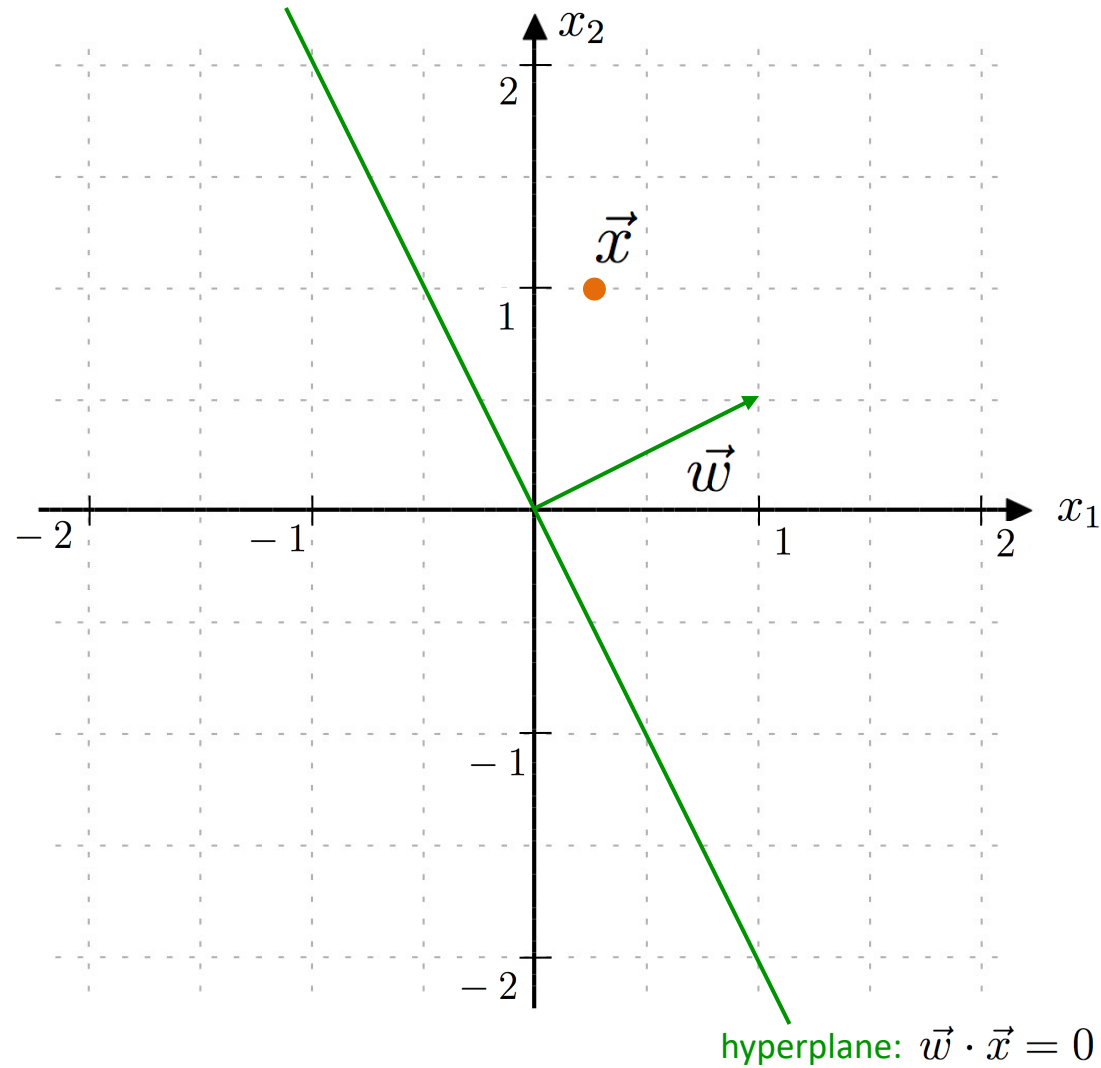
$$\pi_3 = \frac{1.1}{4}$$

Outline for April 23

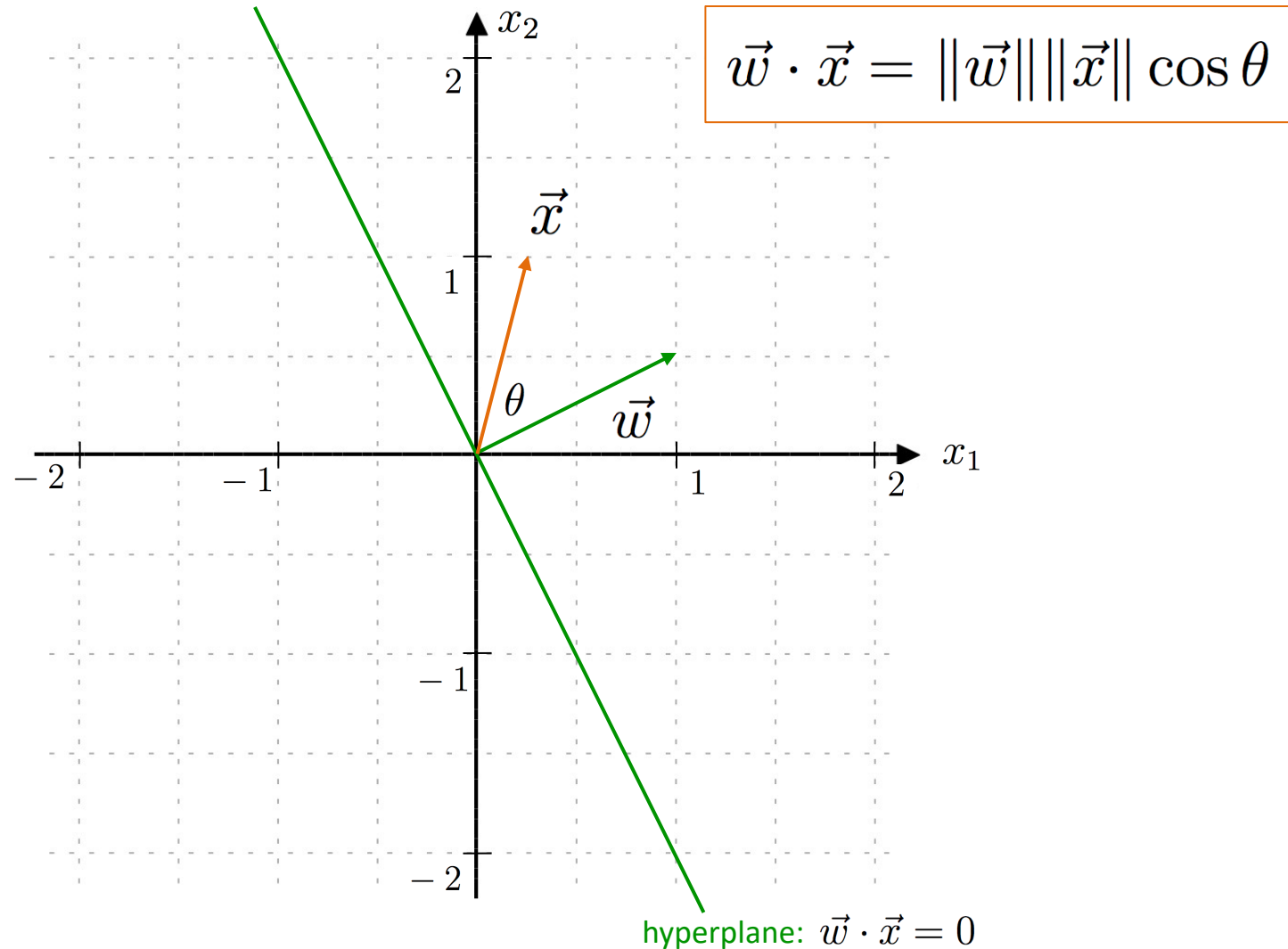
- Review temperature, cross-entropy, loss functions, and ML pipeline
- Convolutional filters, backpropagation, and CNNs
- Gaussian mixture model handout
- Perceptron and SVM

Handout 25, page 3!

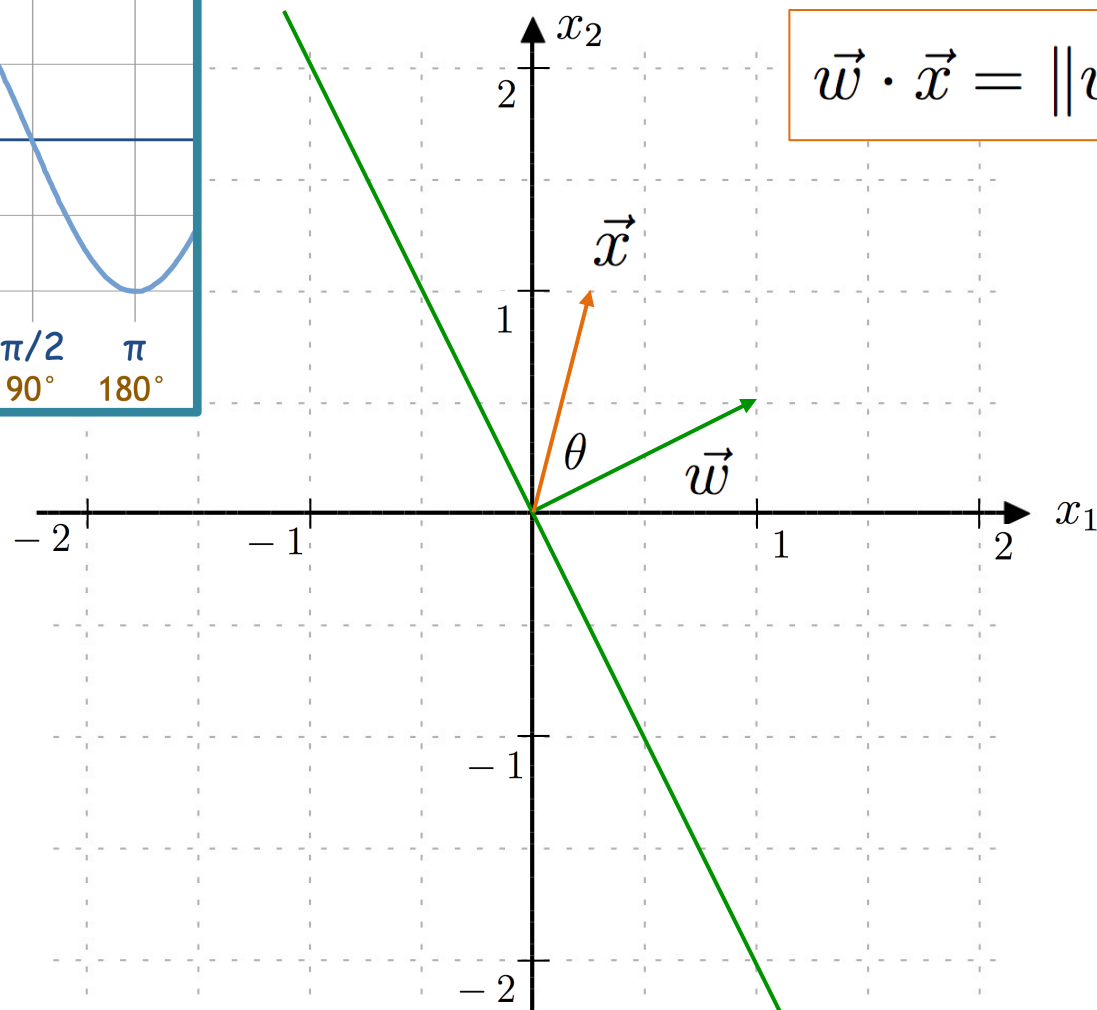
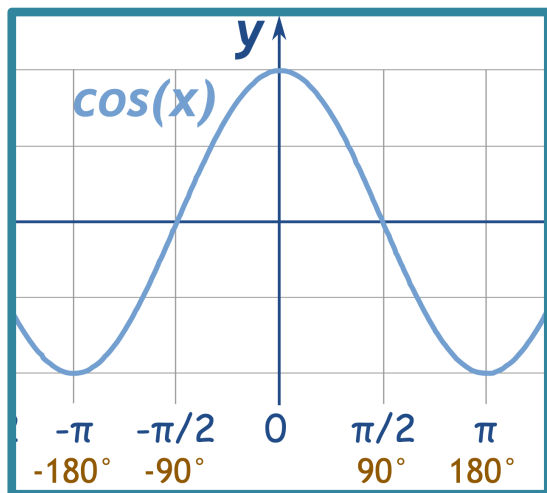
Intuition behind the dot product



Intuition behind the dot product



Intuition behind the dot product



$$\vec{w} \cdot \vec{x} = \|\vec{w}\| \|\vec{x}\| \cos \theta$$

hyperplane: $\vec{w} \cdot \vec{x} = 0$

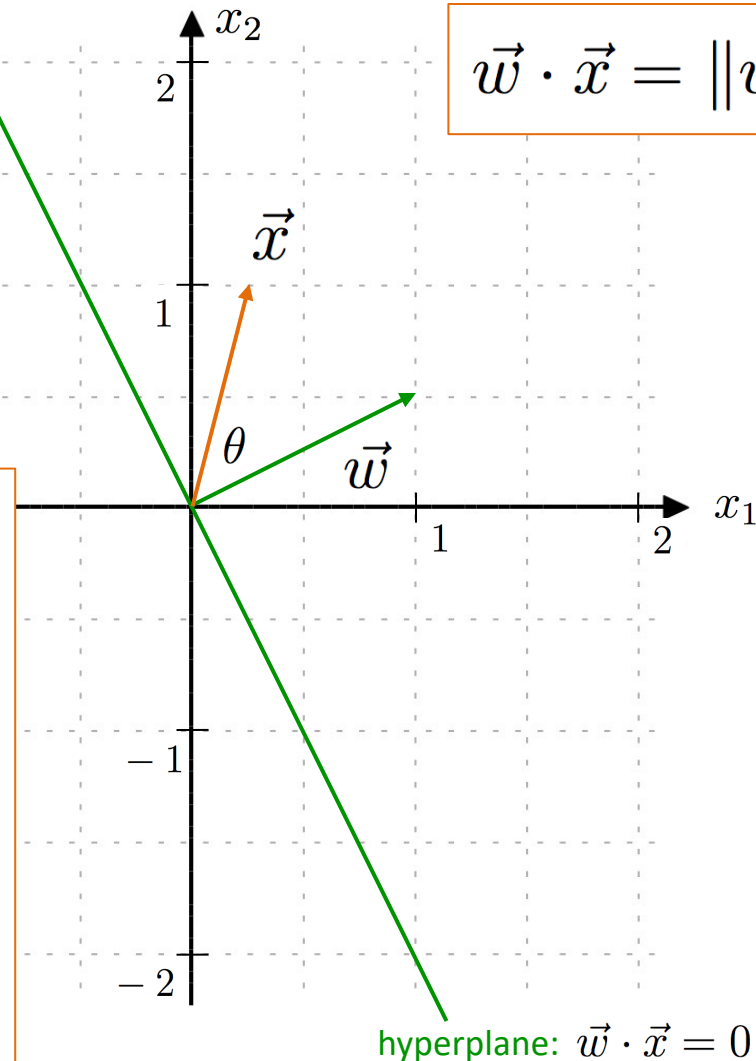
Intuition behind the dot product

Takeaway: we only care about the sign of the angle between \mathbf{x} and \mathbf{w}

- If $\cos \theta > 0$, \mathbf{x} is on the same side of the hyperplane as \mathbf{w} , so we classify it as positive
- If $\cos \theta < 0$, \mathbf{x} is on the opposite side from \mathbf{w} , so we classify it as negative

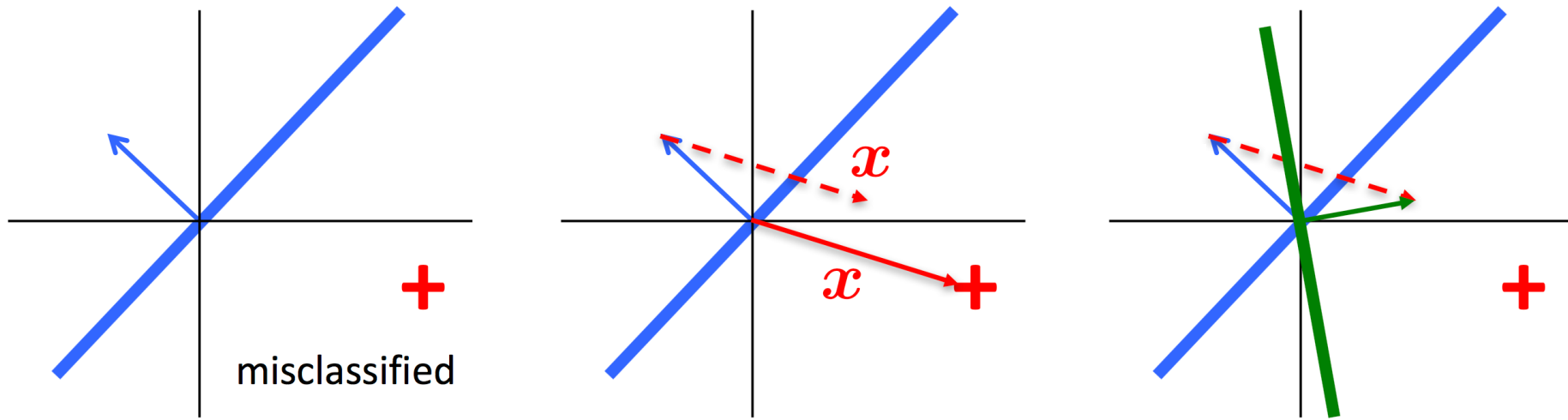
Takeaway: we only care about the sign of the angle between \mathbf{x} and \mathbf{w}

- If $\cos \theta > 0$, \mathbf{x} is on the same side of the hyperplane as \mathbf{w} , so we classify it as positive
- If $\cos \theta < 0$, \mathbf{x} is on the opposite side from \mathbf{w} , so we classify it as negative

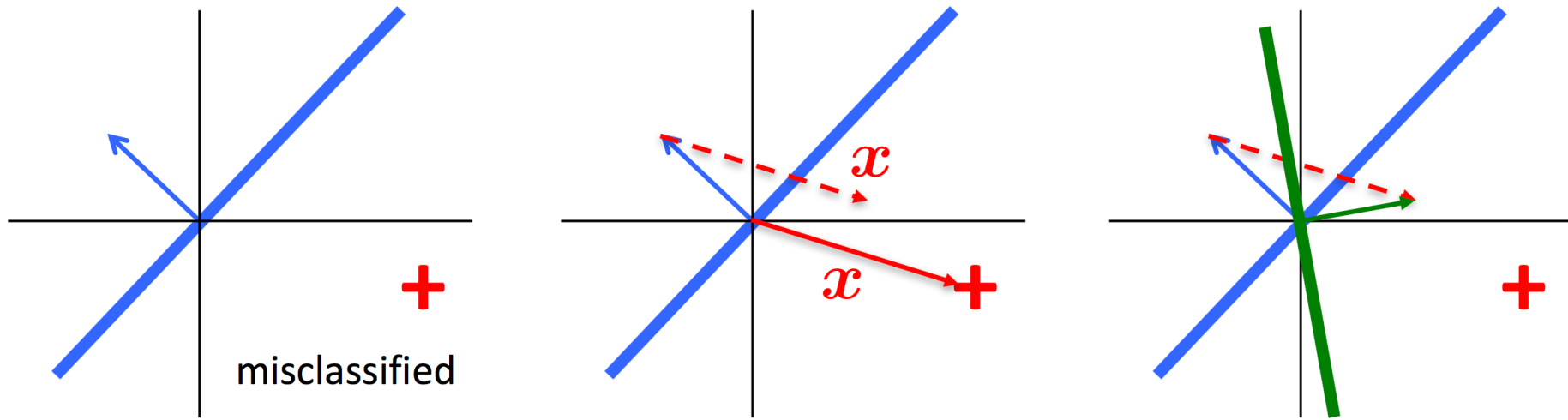


$$\vec{w} \cdot \vec{x} = \|\vec{w}\| \|\vec{x}\| \cos \theta$$

Perceptron algorithm and intuition



Perceptron algorithm and intuition



Let $\vec{w} = [0, 0, \dots, 0]^T$

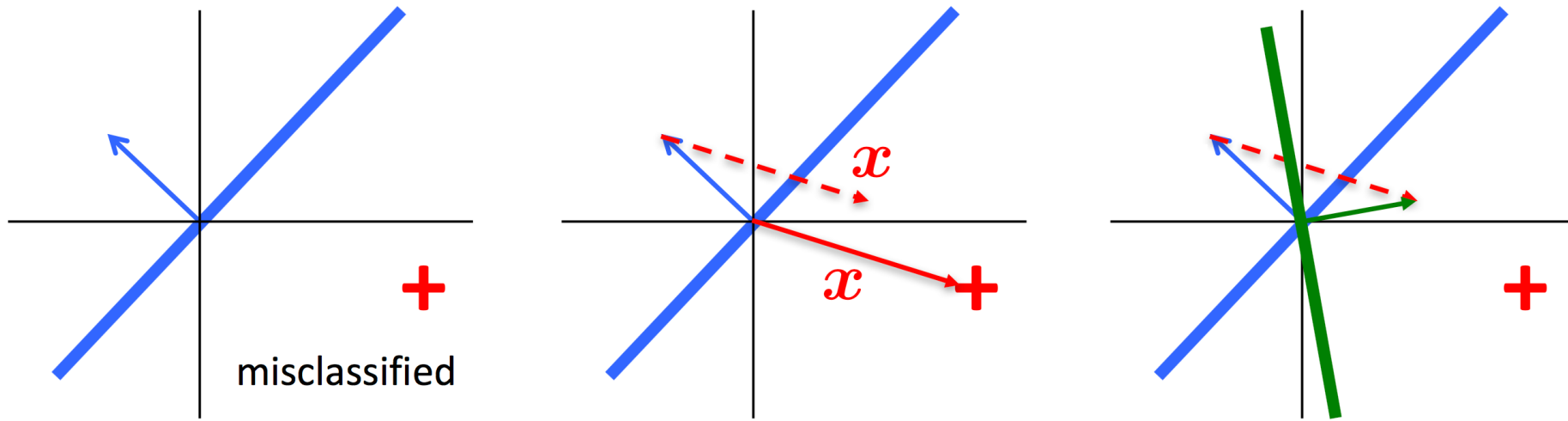
Repeat until convergence:

Receive training example (\vec{x}_i, y_i)

If $y_i(\vec{w}^T \vec{x}_i) \leq 0$ (incorrectly classified)

$$\vec{w} \leftarrow \vec{w} + \alpha y_i \vec{x}_i$$

Perceptron algorithm and intuition



Let $\vec{w} = [0, 0, \dots, 0]^T$

Repeat until convergence:

Receive training example (\vec{x}_i, y_i)

If $y_i(\vec{w}^T \vec{x}_i) \leq 0$ (incorrectly classified)

$$\vec{w} \leftarrow \vec{w} + \alpha y_i \vec{x}_i$$

Convergence:

- All data points correctly classified
- Fixed number of iterations passed

Often: $\alpha = 1$ (only changes magnitude of weight vector)

Functional and Geometric Margins

SVM classifier:
(same as perceptron)

$$h(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$$

Functional and Geometric Margins

SVM classifier:
(same as perceptron)

$$h(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$$

Functional Margin:

$$\hat{\gamma}_i = y_i(\vec{w} \cdot \vec{x}_i + b)$$

Functional and Geometric Margins

SVM classifier:
(same as perceptron)

$$h(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$$

Functional Margin:

$$\hat{\gamma}_i = y_i(\vec{w} \cdot \vec{x}_i + b)$$

Geometric Margin:
(distance between
example and hyperplane)

$$\gamma_i = y_i \left(\frac{\vec{w}}{\|\vec{w}\|} \cdot \vec{x}_i + \frac{b}{\|\vec{w}\|} \right)$$

Functional and Geometric Margins

SVM classifier:
(same as perceptron)

$$h(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$$

Functional Margin:

$$\hat{\gamma}_i = y_i(\vec{w} \cdot \vec{x}_i + b)$$

Geometric Margin:
(distance between
example and hyperplane)

$$\gamma_i = y_i \left(\frac{\vec{w}}{\|\vec{w}\|} \cdot \vec{x}_i + \frac{b}{\|\vec{w}\|} \right)$$

Note:

$$\gamma_i = \frac{\hat{\gamma}_i}{\|\vec{w}\|}$$

Optimization Problem: try 1

Goal: maximize the minimum distance
between example and hyperplane

$$\gamma = \min_{i=1, \dots, n} \gamma_i$$

Optimization Problem: try 1

Goal: maximize the minimum distance
between example and hyperplane

$$\gamma = \min_{i=1, \dots, n} \gamma_i$$

Formulation: optimize a function with
respect to a constraint

$$\begin{aligned} \max_{\gamma, \vec{w}, b} \quad & \gamma \\ \text{s.t.} \quad & y_i(\vec{w} \cdot \vec{x}_i + b) \geq \gamma, \quad i = 1, \dots, n \\ \text{and} \quad & \|\vec{w}\| = 1 \end{aligned}$$

(force functional and geometric
margin to be equal)

Optimization Problem: try 2

Idea: substitute functional margin
divided by magnitude of weight vector

$$\begin{aligned} \max_{\hat{\gamma}, \vec{w}, b} \quad & \frac{\hat{\gamma}}{\|\vec{w}\|} \\ \text{s.t.} \quad & y_i(\vec{w} \cdot \vec{x}_i + b) \geq \hat{\gamma}, \quad i = 1, \dots, n \end{aligned}$$

(gets rid of non-convex constraint)

Optimization Problem: try 3

Idea: put arbitrary constraint on functional margin

$$\hat{\gamma} = 1$$

$$\begin{aligned} \min_{\vec{w}, b} \quad & \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.t.} \quad & y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

Optimization Problem: try 3

Idea: put arbitrary constraint on functional margin

$$\hat{\gamma} = 1$$

$$\begin{aligned} \min_{\vec{w}, b} \quad & \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.t.} \quad & y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

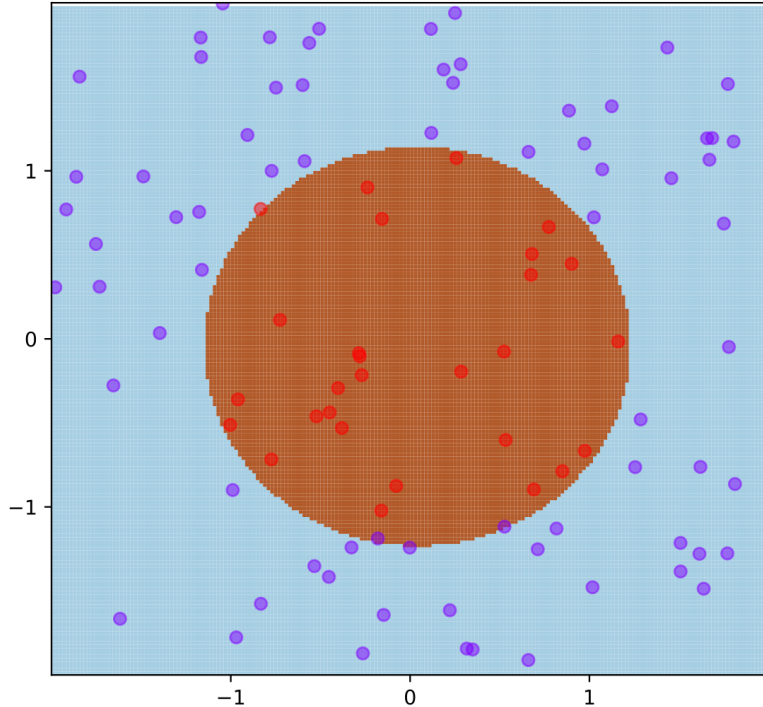
$$\begin{aligned} \min_{\vec{w}, b} \quad & \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.t.} \quad & -y_i(\vec{w} \cdot \vec{x}_i + b) + 1 \leq 0, \quad i = 1, \dots, n \end{aligned}$$

Kernel Idea

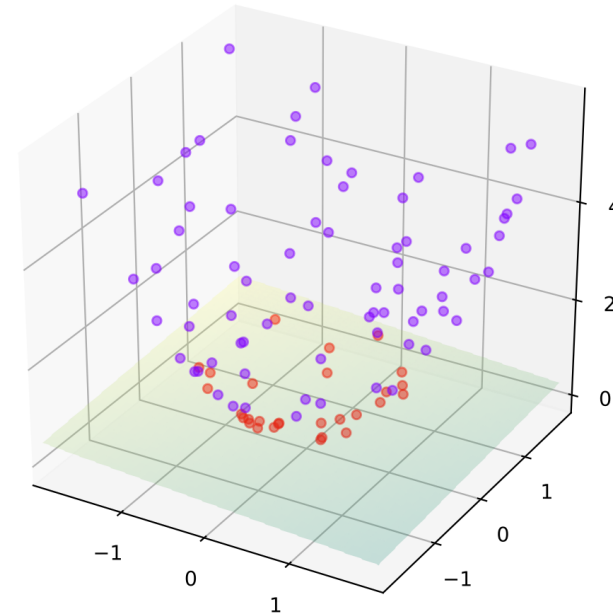
- By solving the dual form of the problem, we have seen how all computations can be done in terms of inner products between examples
- One example of an inner product is the dot product, which is the linear version of SVMs
- But there are many others!
- Intuition: if points are close together, their kernel function will have a large value (measure of similarity)

Kernel Trick example

Feature mapping: $\varphi(\mathbf{x}) = (x_1, x_2, x_1^2 + x_2^2)$



Original feature space



Mapping after applying kernel
(can now find a hyperplane)

Kernel function: $K(\mathbf{x}, \mathbf{z}) = \mathbf{x} \cdot \mathbf{z} + \|\mathbf{x}\|^2 \|\mathbf{z}\|^2$

Gaussian Kernel

- Gaussian kernel is near 0 when points are far apart and near 1 when they are similar
- Also called Radial Basis Function (RBF) kernel

$$K(\vec{x}, \vec{z}) = \exp\left(-\frac{\|\vec{x} - \vec{z}\|^2}{2\sigma^2}\right)$$

Gaussian Kernel

- Gaussian kernel is near 0 when points are far apart and near 1 when they are similar
- Also called Radial Basis Function (RBF) kernel

$$K(\vec{x}, \vec{z}) = \exp \left(-\frac{\|\vec{x} - \vec{z}\|^2}{2\sigma^2} \right)$$

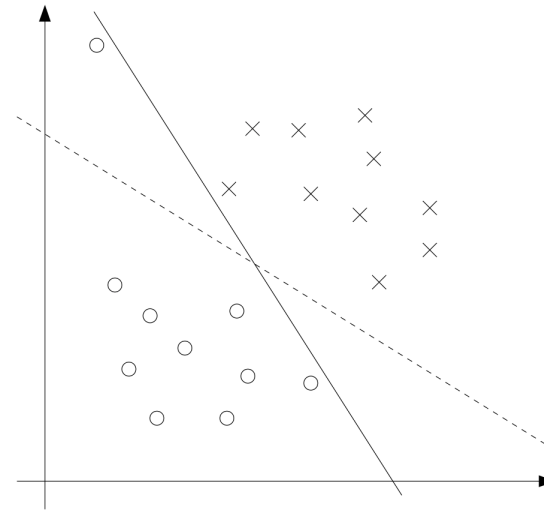
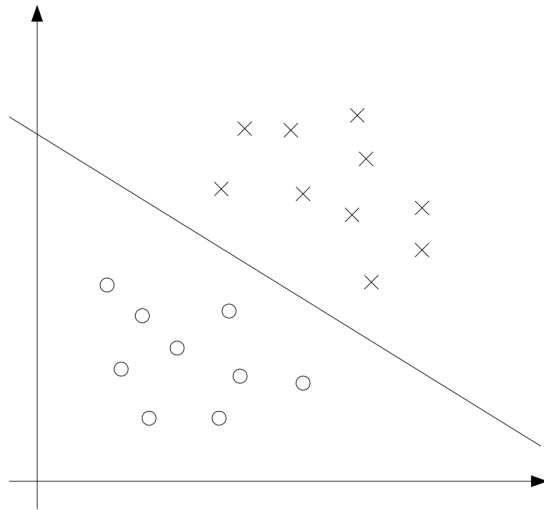
Often re-parametrized by
gamma

$$\gamma = \frac{1}{2\sigma^2}$$

$$K(\vec{x}, \vec{z}) = \exp \left(-\gamma \|\vec{x} - \vec{z}\|^2 \right)$$

Soft-margin SVMs (non-separable case)

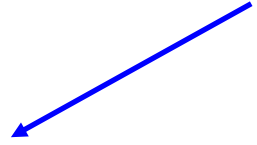
- Idea: we will use regularization to add a cost for each point being incorrectly classified by the hyperplane
- Hopefully many costs will be 0, but we can accommodate a few outliers



Soft-margin SVMs (non-separable case)

- New optimization problem with regularization

$$\begin{aligned} \min_{\xi, \vec{w}, b} \quad & \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\ \text{and} \quad & \xi_i \geq 0, \quad i = 1, \dots, n \end{aligned}$$

"flexible margin" 

Meta-optimization process

- Incremental SVM optimization algorithm

Meta-optimization process

- Incremental SVM optimization algorithm
- Choose a subset S of examples and run optimization to get alpha values

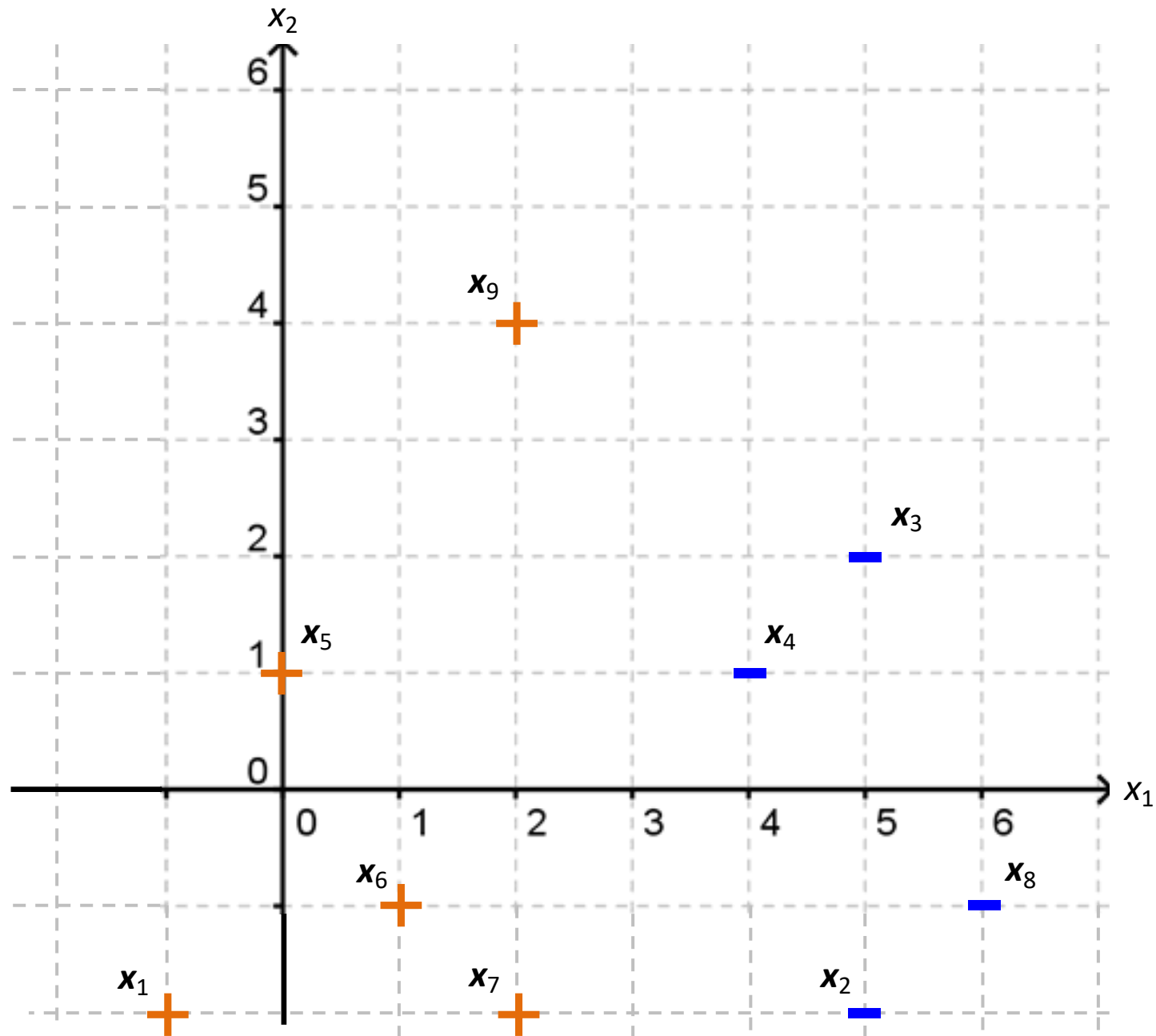
Meta-optimization process

- Incremental SVM optimization algorithm
- Choose a subset S of examples and run optimization to get alpha values
- Identify which alpha values are 0 \Rightarrow these cannot be support vectors in final solution!

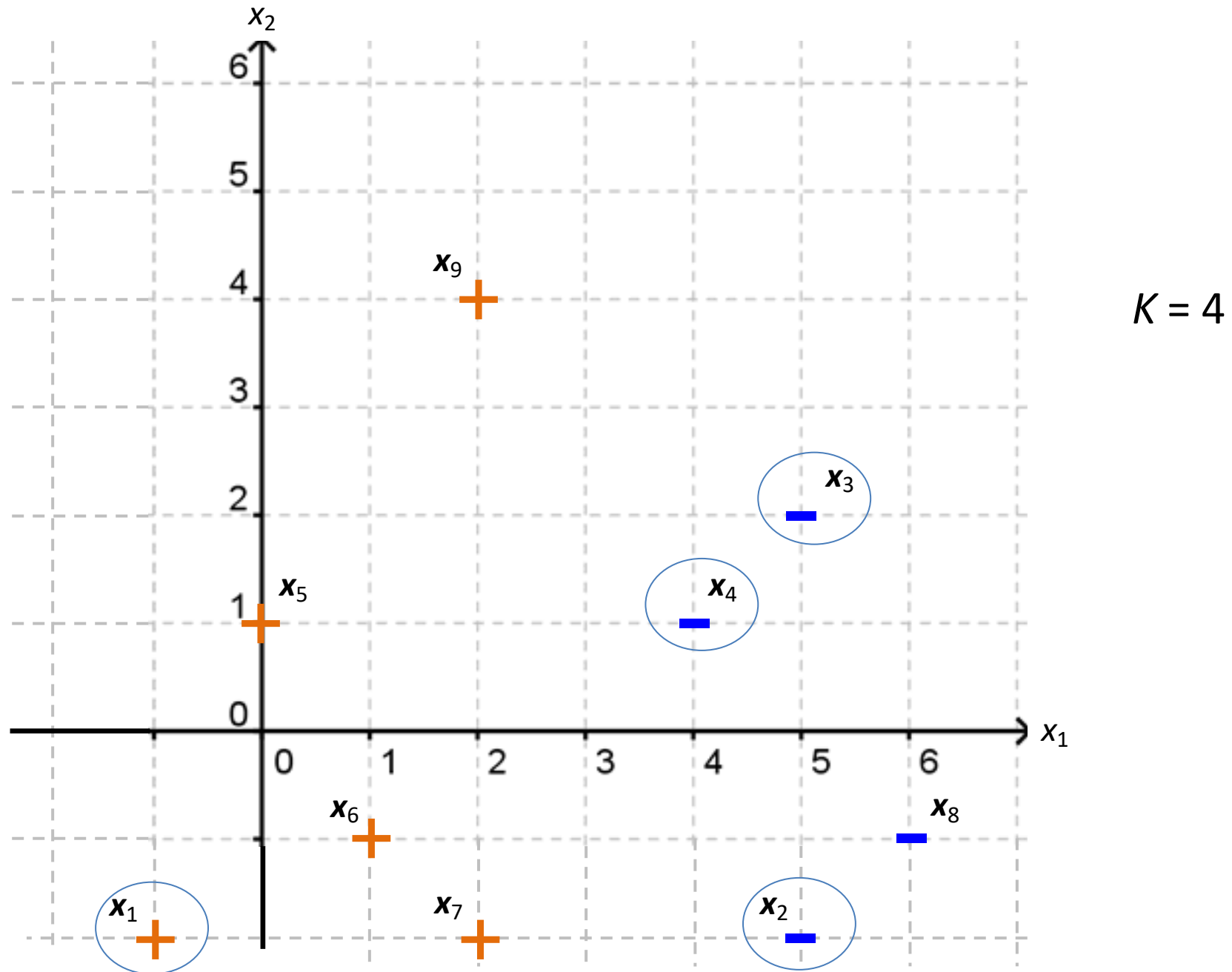
Meta-optimization process

- Incremental SVM optimization algorithm
- Choose a subset S of examples and run optimization to get alpha values
- Identify which alpha values are 0 \Rightarrow these cannot be support vectors in final solution!
- Discard these points and add new ones; repeat

Meta-optimization: example



Meta-optimization: example



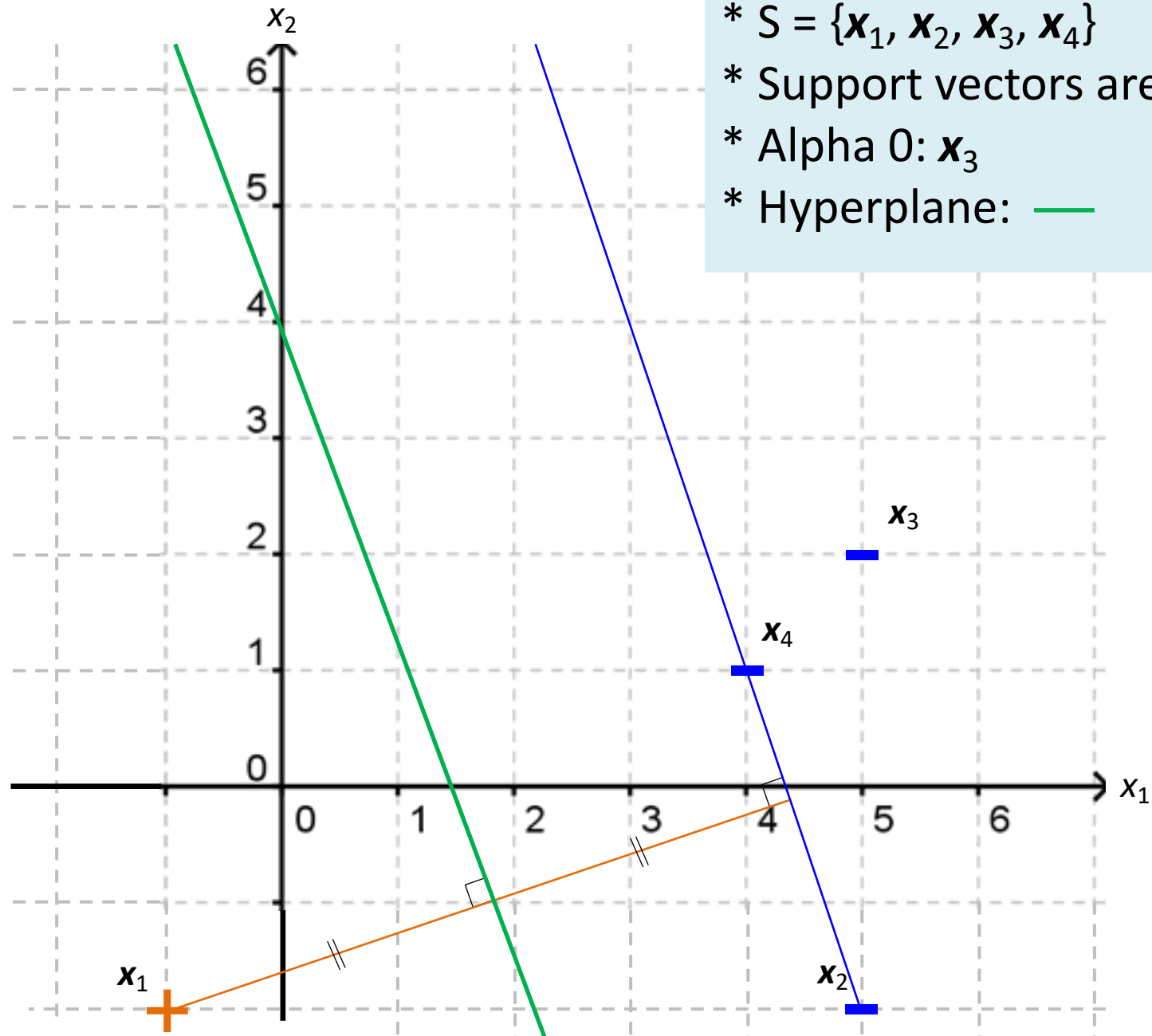
Round 1:

* $S = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4\}$

* Support vectors are: $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_4$

* Alpha 0: \mathbf{x}_3

* Hyperplane: —



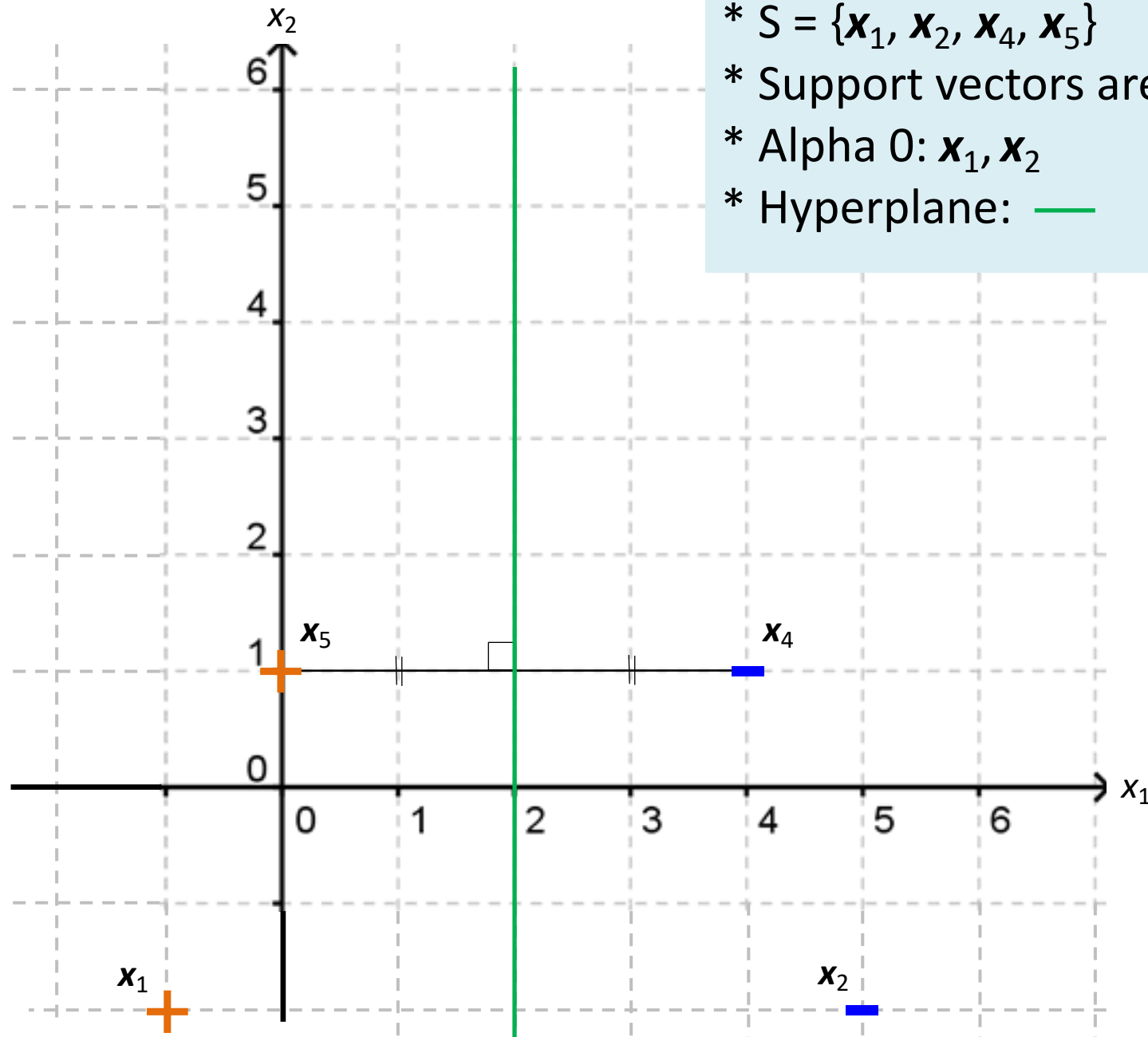
Round 1:

* $S = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_4, \mathbf{x}_5\}$

* Support vectors are: $\mathbf{x}_4, \mathbf{x}_5$

* Alpha 0: $\mathbf{x}_1, \mathbf{x}_2$

* Hyperplane: —



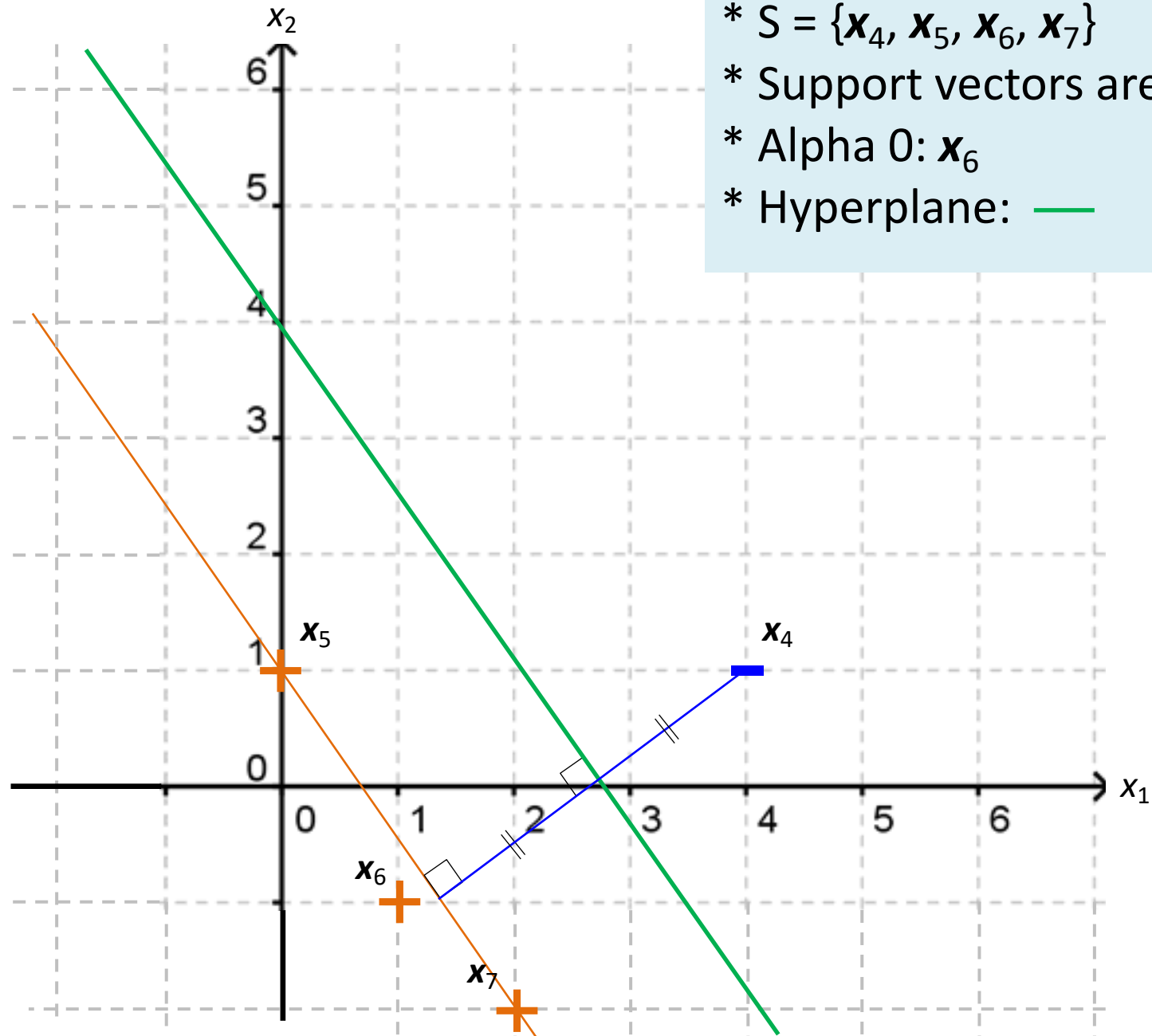
Round 3:

* $S = \{\mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6, \mathbf{x}_7\}$

* Support vectors are: $\mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_7$

* Alpha 0: \mathbf{x}_6

* Hyperplane: —



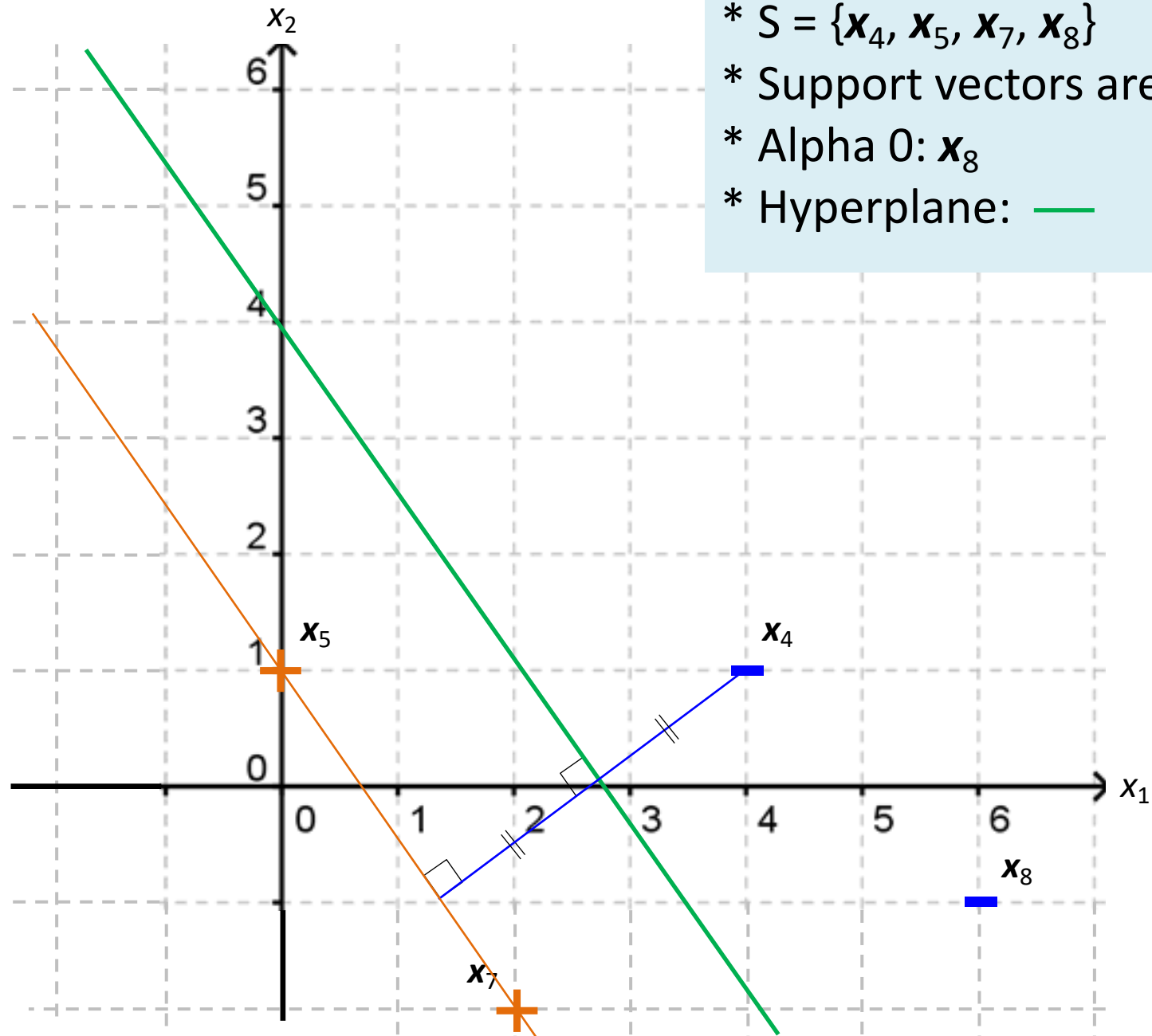
Round 4:

* $S = \{\mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_7, \mathbf{x}_8\}$

* Support vectors are: $\mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_7$

* Alpha 0: \mathbf{x}_8

* Hyperplane: —



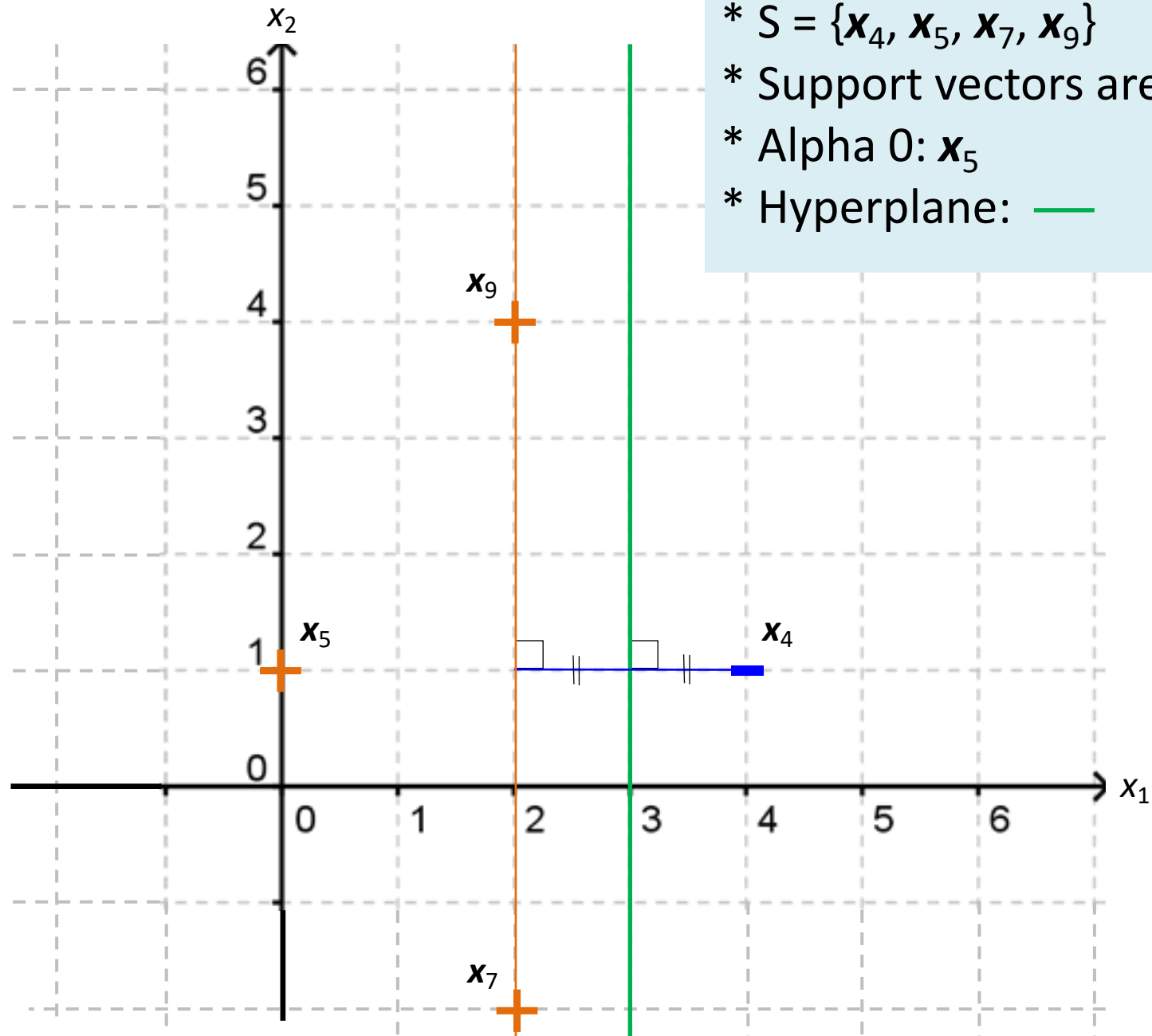
Round 5:

* $S = \{\mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_7, \mathbf{x}_9\}$

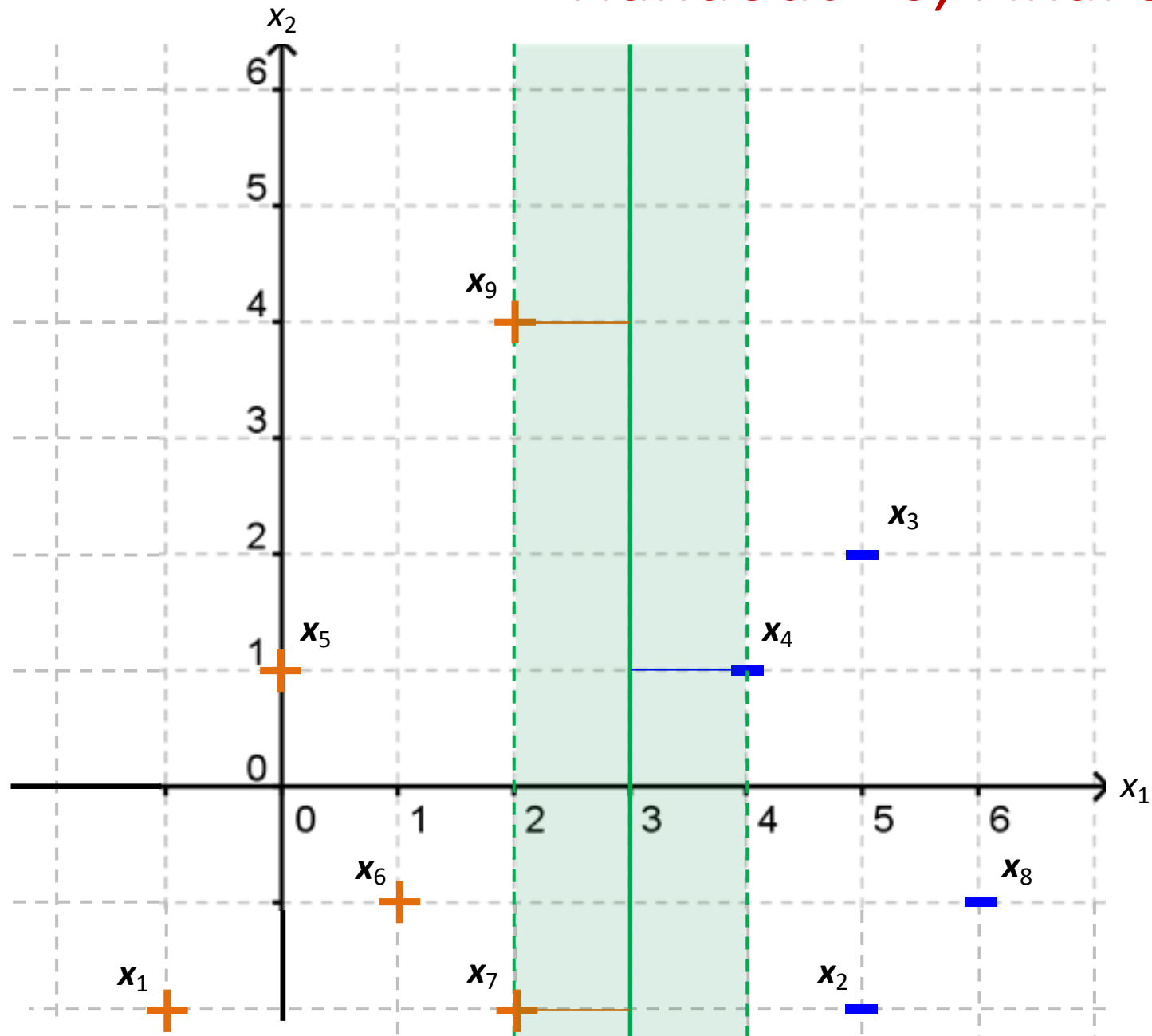
* Support vectors are: $\mathbf{x}_4, \mathbf{x}_7, \mathbf{x}_9$

* Alpha 0: \mathbf{x}_5

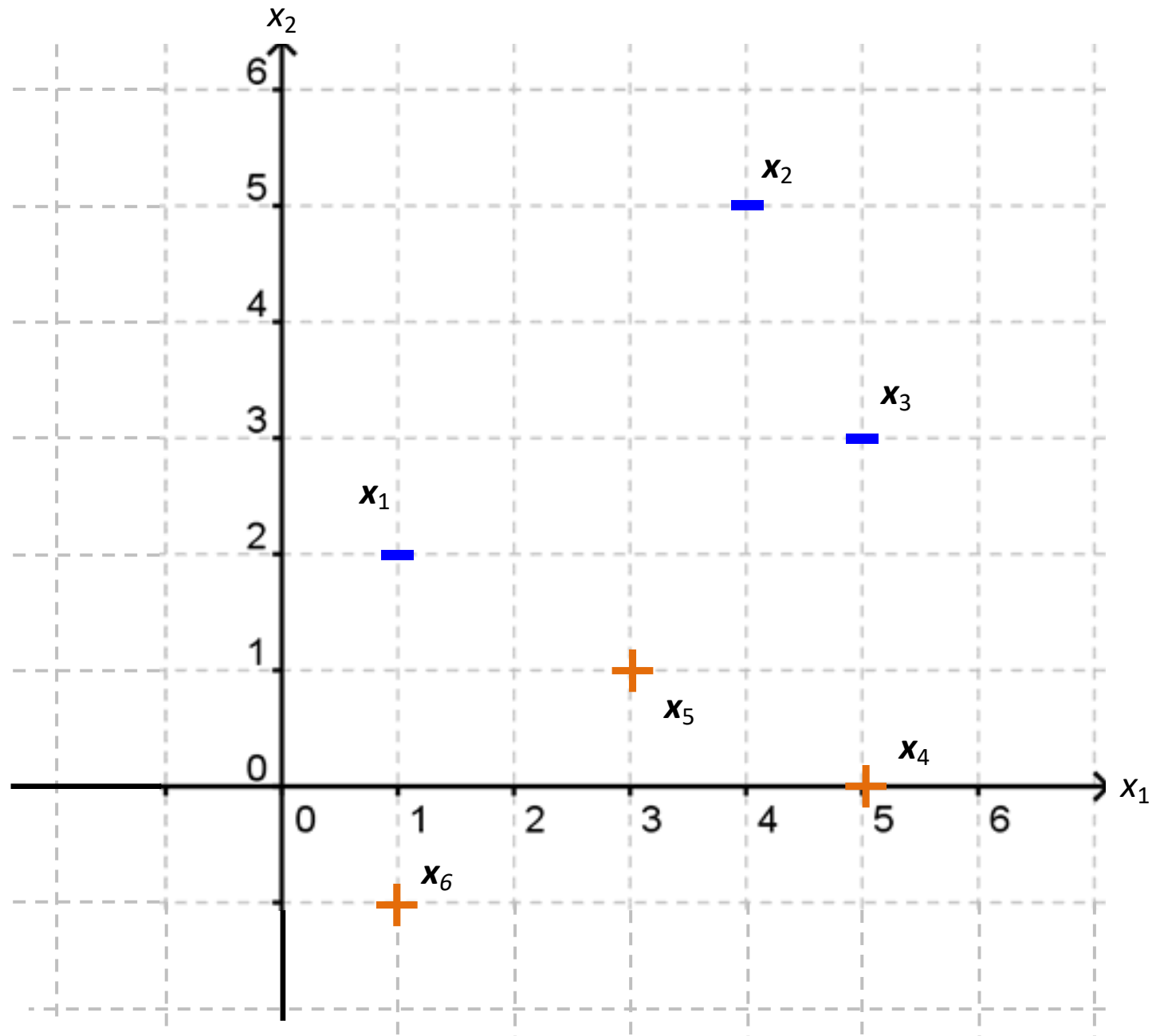
* Hyperplane: —



Handout 16, Final Solution



Discuss with a partner: what are the support vectors?



Discuss with a partner: what are the support vectors?

