# CS 360: Machine Learning

Sara Mathieson, Sorelle Friedler

Spring 2024

# Admin

- **Lab 4** and **Lab 5** graded
  - Any regrade requests (including midterm) must be brought within 1 week of receiving your grade

- **Lab 6** was due last night (see Piazza for runtime issues if you're taking a late day)

- **Lab 7** posted, due Thurs April 4
  - Last lab with required partners
  - We can form partners during lab today

- **Project proposal** due April 8 (short)

# Outline for March 26

- SVM extensions

- Introduction to neural networks

- Fully connected (FC) neural networks

- Image data format and intro to Lab 7

# Outline for March 26

- SVM extensions

- Introduction to neural networks

- Fully connected (FC) neural networks

- Image data format and intro to Lab 7

# SVM dual optimization problem

$\boxed{\text{SVM}}$ optimization problem

$$\min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 \quad s.t.$$

$$y_i (\vec{w} \cdot \vec{x}_i + b) - 1 \geq 0 \qquad i = 1 \cdots n$$

$$L(\vec{w}, b, \vec{\alpha}) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^{n} \alpha_i [y_i (\vec{w} \cdot \vec{x}_i + b) - 1]$$
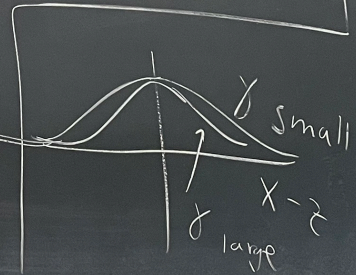
↑ ↑ ↑

take gradient
for all ⇒
Set to 0

really ≈ far
far

$\boxed{\vec{w}}$

$$\nabla_{\vec{w}} L(\vec{w}, b, \vec{\alpha}) = \vec{w} - \sum_{i=1}^{n} \alpha_i y_i \vec{x}_i = \vec{0}$$

$$\Rightarrow \boxed{\vec{w} = \sum_{i=1}^{n} \alpha_i y_i \vec{x}_i}$$

$\alpha_i > 0$ if $\vec{x}_i$
support vector

$\alpha_i = 0$ o.w.

γ small
X - ẑ
γ large

$$\boxed{b} \quad \frac{\partial h(\vec{w}, b, \vec{\alpha})}{\partial b} = \sum_{i=1}^{n} \alpha_i y_i = 0 \qquad\qquad y_i = \{-1, +1\}$$

$$\Rightarrow \quad \boxed{\sum_{i: y_i = -1} \alpha_i = \sum_{i: y_i = +1} \alpha_i}$$

$$\boxed{dual}$$

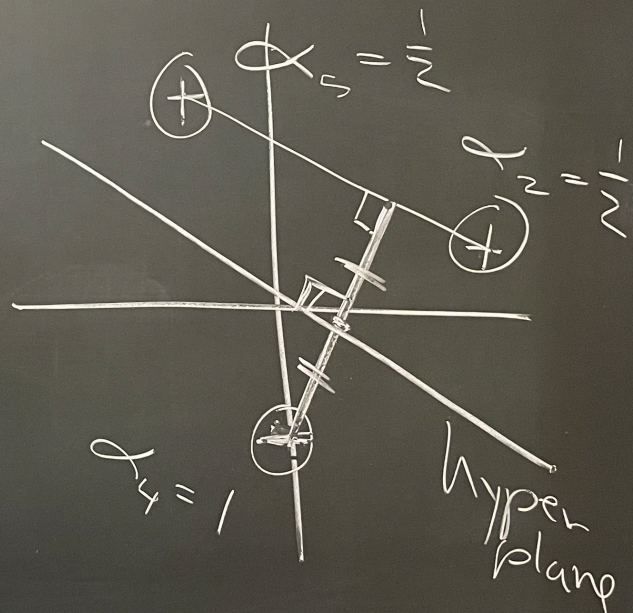$$\max_{\vec{\alpha}} W(\vec{\alpha}) = \sum_{i=1}^{n} \alpha_i$$

$$- \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} y_i y_j \alpha_i \alpha_j \boxed{\vec{x_i} \cdot \vec{x_j}}$$

dot product
flexibility !

$$\text{s.t.} \quad \boxed{\alpha_i \geq 0}$$

$$\& \quad \sum_{i=1}^{n} \alpha_i y_i = 0$$

$$\alpha_s = \frac{1}{2}$$

$$\ell_2 = \frac{1}{2}$$
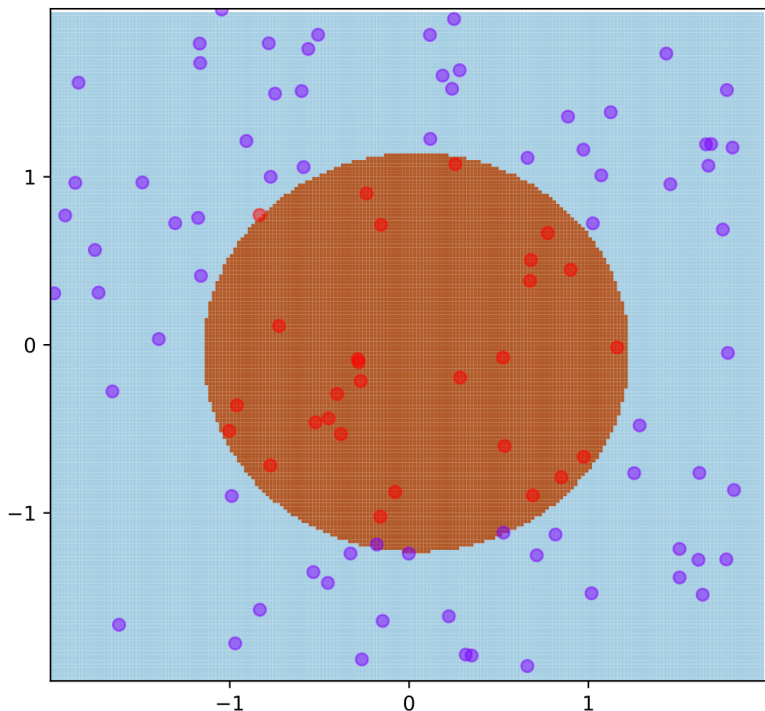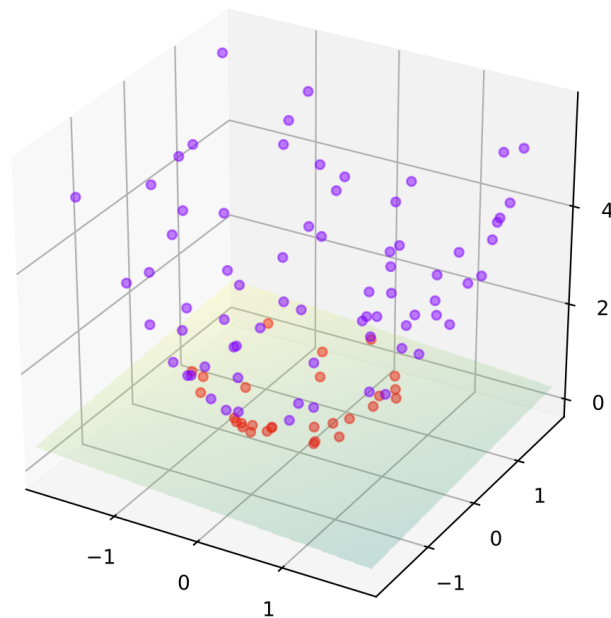
$$\alpha_4 = 1$$

hyper
plane

# 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: $\varphi(\boldsymbol{x}) = (x_1, \ x_2, \ x_1^2 + x_2^2)$



Original feature space



Mapping after applying kernel
(can now find a hyperplane)

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

# 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

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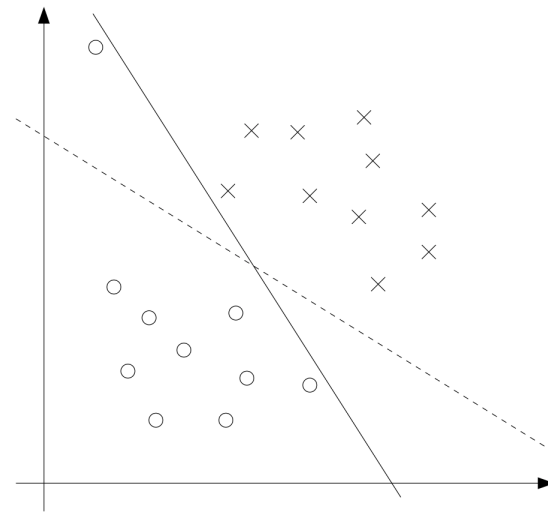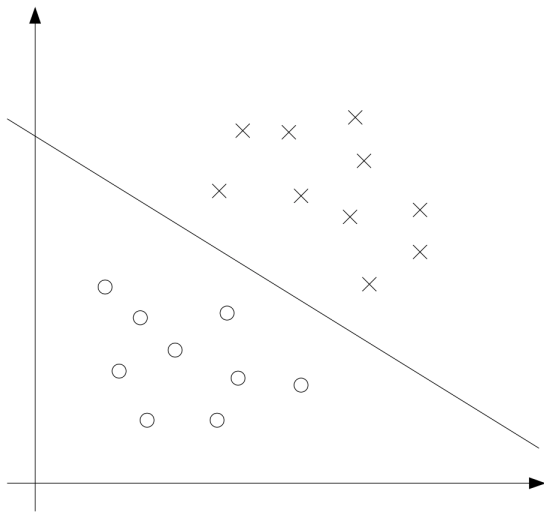
# Gaussian Kernel

- 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

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

$$K(\vec{x}, \vec{z}) = \exp\left(-\gamma\|\vec{x} - \vec{z}\|^2\right)$$
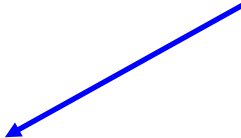
# Soft-margin SVMs (non-separable case)

- 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

Figure: Andrew Ng

# Soft-margin SVMs (non-separable case)

- New optimization problem with regularization

$$\min_{\xi, \vec{w}, b} \quad \frac{1}{2}\|\vec{w}\|^2 + C\sum_{i=1}^{n}\xi_i$$

"flexible margin"

$$\text{s.t.} \quad y_i(\vec{w}\cdot\vec{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \cdots, n$$

$$\text{and} \quad \xi_i \geq 0, \quad i = 1, \cdots, n$$

# 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
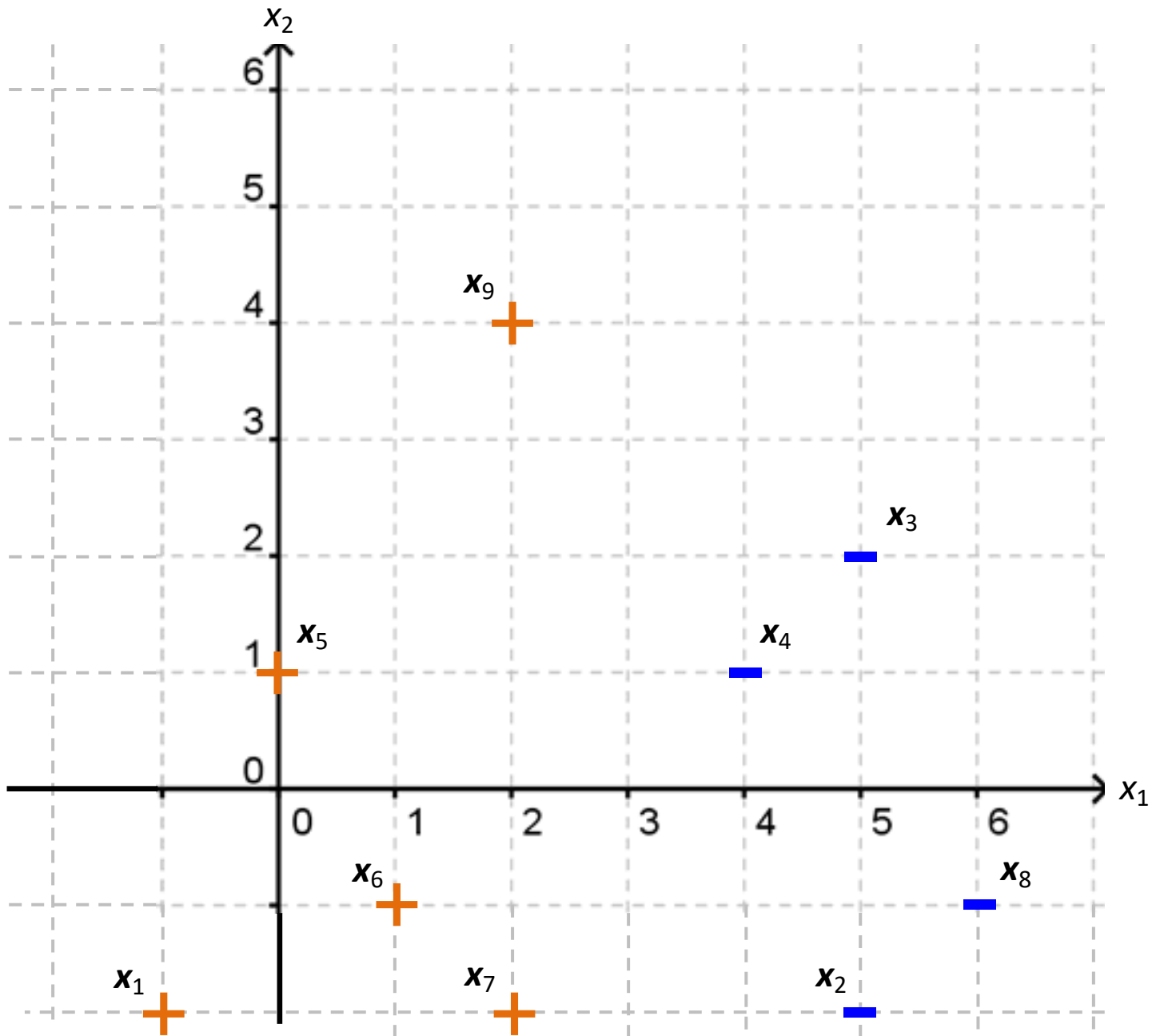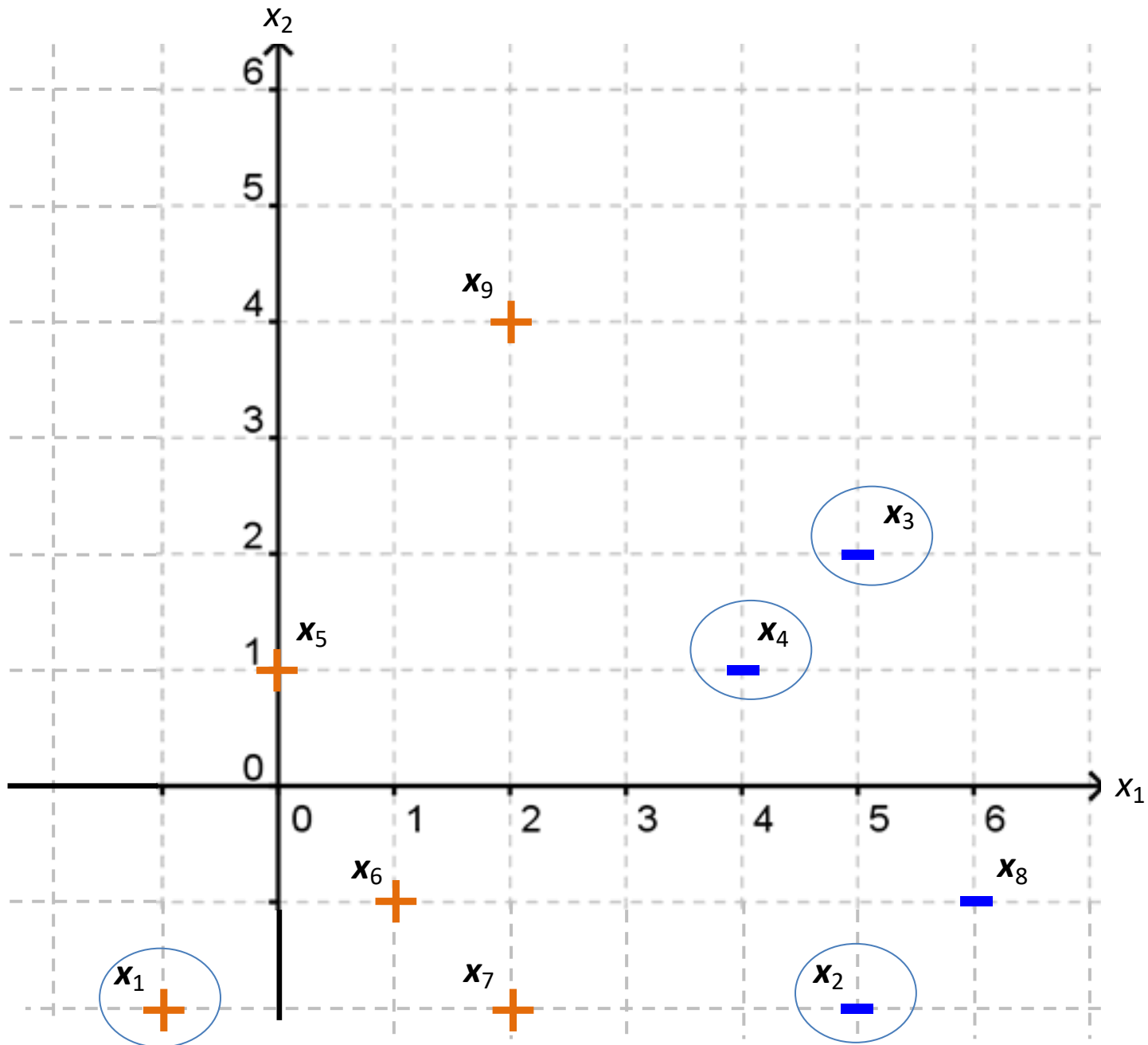
# Meta-optimization process

- Incremental SVM optimization algorithm

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

- Identify which alpha values are 0 => 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 => these cannot be support vectors in final solution!

- Discard these points and add new ones; repeat

$K = 4$

Round 1:
* S = {$x_1$, $x_2$, $x_3$, $x_4$}
* Support vectors are: $x_1$, $x_2$, $x_4$
* Alpha 0: $x_3$
* Hyperplane: ——

Round 1:
* S = {$x_1$, $x_2$, $x_4$, $x_5$}
* Support vectors are: $x_4$, $x_5$
* Alpha 0: $x_1$, $x_2$
* Hyperplane: —

Round 3:
* S = {$x_4$, $x_5$, $x_6$, $x_7$}
* Support vectors are: $x_4$, $x_5$, $x_7$
* Alpha 0: $x_6$
* Hyperplane: —

Round 4:
* S = {$x_4$, $x_5$, $x_7$, $x_8$}
* Support vectors are: $x_4$, $x_5$, $x_7$
* Alpha 0: $x_8$
* Hyperplane: —

Round 5:
* S = {$x_4$, $x_5$, $x_7$, $x_9$}
* Support vectors are: $x_4$, $x_7$, $x_9$
* Alpha 0: $x_5$
* Hyperplane: ——

Handout 16, Final Solution

# Discuss with a partner

1. If $\vec{x}_i$ is a support vector, what can we say about it? Circle all that apply:

    (a) its Lagrange multiplier $\alpha_i > 0$

    (b) its Lagrange multiplier $\alpha_i = 0$

    (c) $y_i(\vec{w} \cdot \vec{x}_i + b) = 0$

    (d) $y_i(\vec{w} \cdot \vec{x}_i + b) = 1$

    (e) $\vec{x}_i$ lies on the margin

# Discuss with a partner

1. If $\vec{x}_i$ is a support vector, what can we say about it? Circle all that apply:
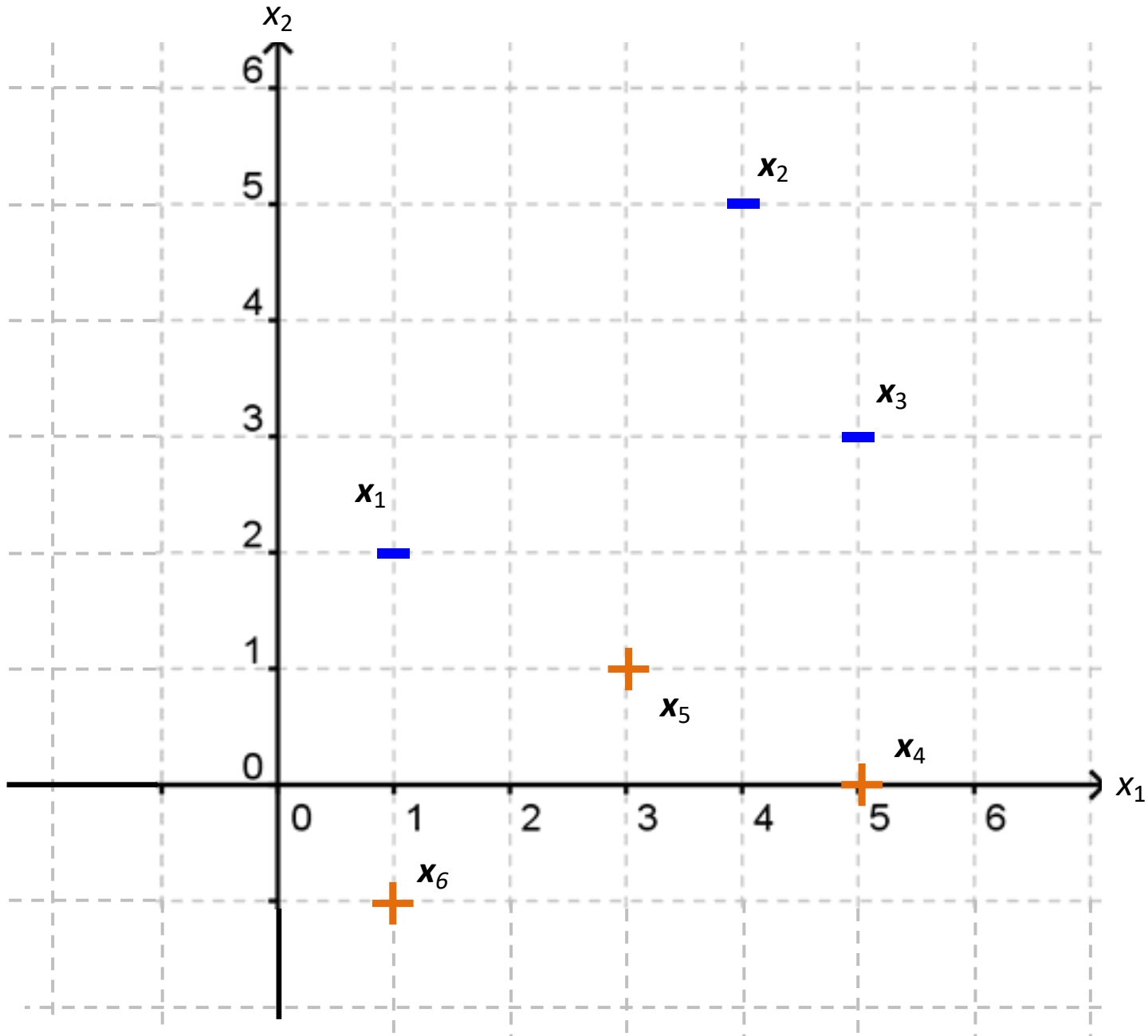
(a) its Lagrange multiplier $\alpha_i > 0$

(b) its Lagrange multiplier $\alpha_i = 0$
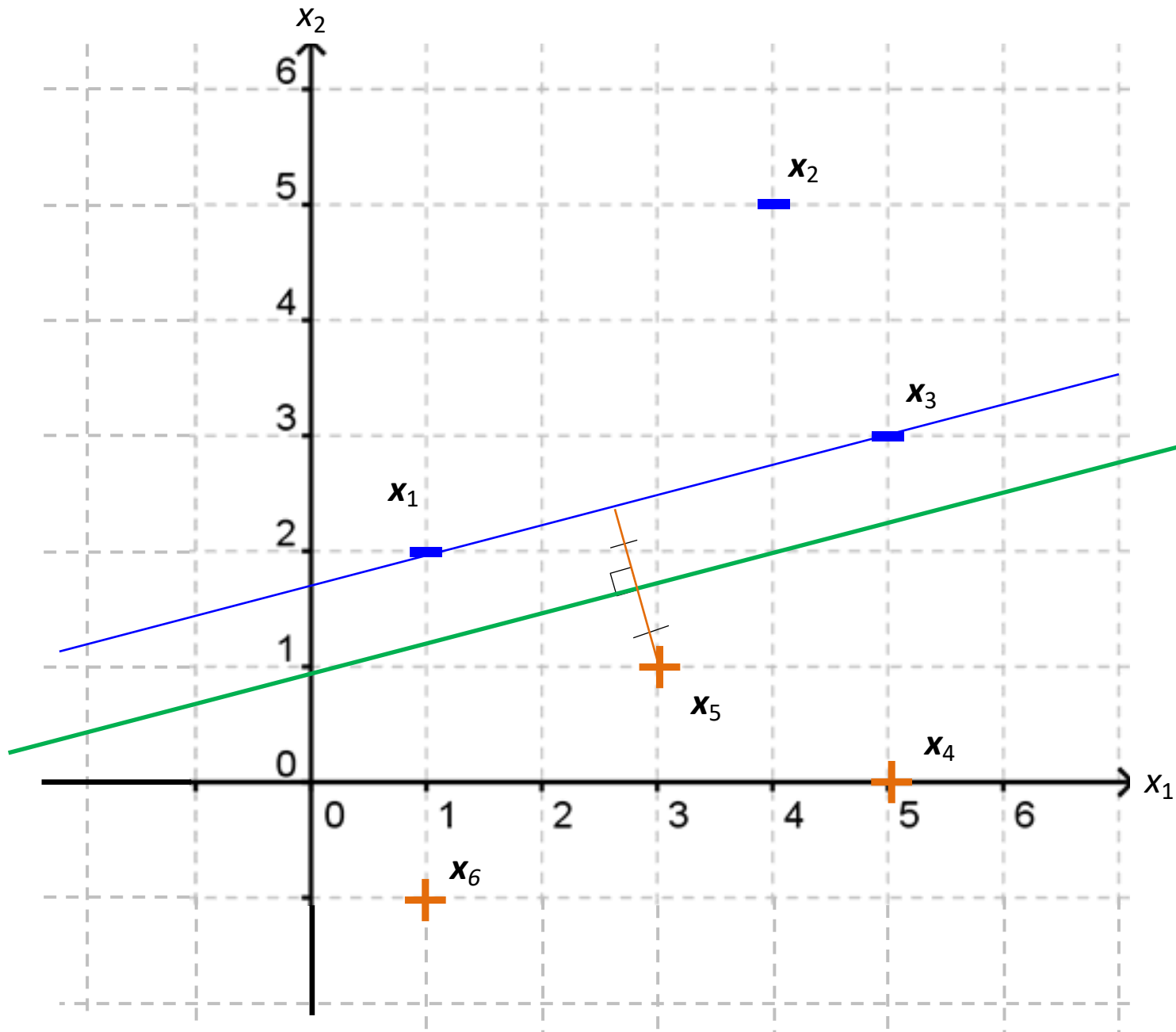
(c) $y_i(\vec{w} \cdot \vec{x}_i + b) = 0$

(d) $y_i(\vec{w} \cdot \vec{x}_i + b) = 1$

(e) $\vec{x}_i$ lies on the margin

# Discuss with a partner: what are the support vectors?

# Discuss with a partner: what are the support vectors?

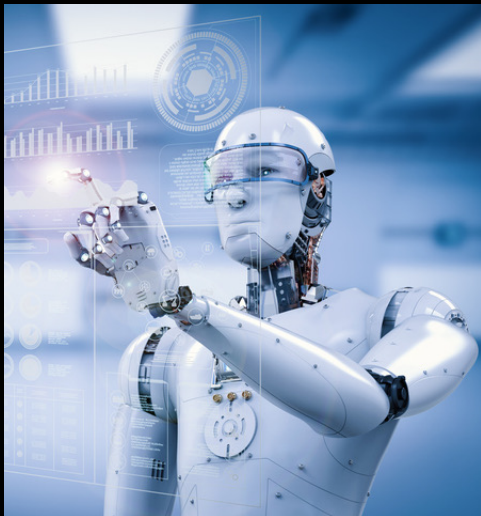# Disadvantages of SVMs

- Difficult to choose a kernel function

- Does not naturally take into account the correlations between features

- Hard to understand and interpret what the model has learned

# Outline for March 26

- SVM extensions

- Introduction to neural networks

- Fully connected (FC) neural networks

- Image data format and intro to Lab 7

# MACHINE LEARNING

# Biological Inspiration

# Goal: learn from complicated inputs



$X_1$

$X_2$

$X_3$

$X_4$

$X_5$

$X_6$

input data

**?**

$Y_1$ (glasses)

$Y_2$ (smiling)

$Y_3$ (eye size)

parameters

# Idea: transform data into lower dimension



input data     hidden layer     parameters

$X_1$ $X_2$ $X_3$ $X_4$ $X_5$ $X_6$

$Y_1$ (glasses)
$Y_2$ (smiling)
$Y_3$ (eye size)

# Multi-layer networks = "deep learning"



input data

hidden layer 1

hidden layer 2

$X_1$
$X_2$
$X_3$
$X_4$
$X_5$
$X_6$

$Y_1$ (glasses)
$Y_2$ (smiling)
$Y_3$ (eye size)

parameters

# Example from my research:
# learning about evolution from genetic data



statistics

hidden layer 1

hidden layer 2

population sizes

selection

# History of Neural Networks

- Perceptron can be interpreted as a simple neural network

- Misconceptions about the weaknesses of perceptrons contributed to declining funding for NN research

- Difficulty of training multi-layer NNs contributed to second setback

- Mid 2000's: breakthroughs in NN training contribute to rise of "deep learning"

# Number of papers that mention "deep learning" over time



2006: Hinton and Salakhutdinov
make a break-through in
initializing deep learning networks

# Big picture for today

- Neural networks can approximate any function!

# Big picture for today

- Neural networks can approximate any function!

- For our purposes in ML, we want to use them to approximate a function from our inputs to our outputs

# Big picture for today

- Neural networks can approximate any function!

- For our purposes in ML, we want to use them to approximate a function from our inputs to our outputs

- We will train our network by asking it to minimize the loss between its output and the true output

# Big picture for today

- Neural networks can approximate any function!

- For our purposes in ML, we want to use them to approximate a function from our inputs to our outputs

- We will train our network by asking it to minimize the loss between its output and the true output

- We will use SGD-like approaches to minimize loss

# Outline for March 26

- SVM extensions

- Introduction to neural networks

- Fully connected (FC) neural networks

- Image data format and intro to Lab 7

# Fully Connected Neural Network Architecture

$$X = \begin{bmatrix} -\overline{x}_i - \end{bmatrix}$$

$n \times p$

(first hidden layer)

$$b^{(1)} = \begin{bmatrix} -7 \\ 2 \\ 0.1 \end{bmatrix}$$

$$h_i^{(1)} = a\left(\vec{w}_i^{(1)} \cdot \vec{x} + b_i^{(1)}\right)$$

↑ activation function    linear function!

$$H^{(1)} = a\left(\underbrace{X}_{n \times p} \underbrace{W^{(1)}}_{p \times p_1} + \underbrace{b^{(1)}}_{p_1 \times 1}\right)$$

$\underbrace{\phantom{XW}}_{n \times p_1}$

broadcast to $n \times p_1$

applied element-wise

one example

$p$ features

$$H^{(2)} = a\left(\underbrace{H^{(1)}}_{n \times p_1} \underbrace{W^{(2)}}_{p_1 \times p_2} + \underbrace{b^{(2)}}_{p_2 \times 1}\right)$$

$n \times 1$

prediction for each example

$$\hat{y} = a\left(\underbrace{H^{(2)}}_{n \times p_2} \underbrace{W^{(3)}}_{p_2 \times 1} + \underbrace{b^{(3)}}_{_{1}}\right)$$

fake ones

## K classes

## 5 classes

$$\text{pred } \hat{y} = \left[ \frac{1}{10}, \frac{1}{10}, \frac{3}{10}, \frac{4}{10}, \frac{1}{10} \right]$$

$$y = [0, 0, 0, 1, 0]$$

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

$$H(y, \hat{y}) = -1 \cdot \log_2 \left( \frac{4}{10} \right)$$

# Option 1: sigmoid function

- Input: all real numbers, output: [0, 1]

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- Derivative is convenient

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

# Option 2: hyperbolic tangent

- Input: all real numbers, output: [-1, 1]

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

# Option 3: Rectified Linear Unit (ReLU)

- Return *x* if *x* is positive (i.e. threshold at 0)

$$f(x) = \max(0, x)$$

# Pros and Cons of Activation Functions

## 1) Sigmoid

- (-) When input becomes very positive or very negative, gradient approaches 0 (saturates and stops gradient descent)
- (-) Not zero-centered, so gradient on weights can end up all positive or all negative (zig-zag in gradient descent)
- (+) Derivative is easy to compute given function value!

More info:
http://cs231n.github.io/neural-networks-1/

# Pros and Cons of Activation Functions

## 1) Sigmoid

- (-) When input becomes very positive or very negative, gradient approaches 0 (saturates and stops gradient descent)
- (-) Not zero-centered, so gradient on weights can end up all positive or all negative (zig-zag in gradient descent)
- (+) Derivative is easy to compute given function value!

## 2) Tanh

- (-) Still has a tendency to prematurely kill the gradient
- (+) Zero-centered so we get a range of gradients
- (+) Rescaling of sigmoid function so derivative is also not too difficult

# Pros and Cons of Activation Functions

## 1) Sigmoid

- (-) When input becomes very positive or very negative, gradient approaches 0 (saturates and stops gradient descent)
- (-) Not zero-centered, so gradient on weights can end up all positive or all negative (zig-zag in gradient descent)
- (+) Derivative is easy to compute given function value!

## 2) Tanh

- (-) Still has a tendency to prematurely kill the gradient
- (+) Zero-centered so we get a range of gradients
- (+) Rescaling of sigmoid function so derivative is also not too difficult

## 3) ReLU

- (+) Works well in practice (accelerates convergence)
- (+) Function value very easy to compute! (no exponentials)
- (-) Units can "die" (no signal) if input becomes too negative throughout gradient descent

More info:
http://cs231n.github.io/neural-networks-1/

# Mini-batches

- So far in this class, we have considered *stochastic gradient descent*, where one data point is used to compute the gradient and update the weights

- On the flipside is *batch gradient descent*, where we compute the gradient with respect to all the data, and then update the weights

- A middle ground uses *mini-batches* of examples before updating the weights. This is the approach we will use in Lab 7.

# Notes about scores and softmax

- The output of the final fully connected layer is a vector of length *K* (number of classes)

- The raw scores are transformed into probabilities using the *softmax function*: (let $s_k$ be the score for class *k*)

$$\hat{y}_k = \frac{e^{s_k}}{\sum_{j=1}^{K} e^{s_j}}$$

- Then we apply *cross-entropy loss* to these probabilities

# Notes about scores and softmax

- The output of the final fully connected layer is a vector of length *K* (number of classes)

- The raw scores are transformed into probabilities using the *softmax function*: (let $s_k$ be the score for class *k*)

$$\hat{y}_k = \frac{e^{s_k}}{\sum_{j=1}^{K} e^{s_j}}$$

Think about outside of class:
- Why do we use exp?
- Why don't we just take the max score?

- Then we apply *cross-entropy loss* to these probabilities

① $\sigma(0) = \frac{1}{2}$

$\tanh(0) = 0$

$ReLU(0) = 0$

② $(3+1)\,4 + (4+1)\cdot4 + (4+1)1$

    ↑   ↑    ↑  ↑     ↑   ↑

    $P$   $P_1$  $P_1$  $P_2$   $P_2$  out

$= \boxed{41}$

③ $H(y, \hat{y}) = 1$   ✗

   $H(y, \hat{y}) = 3$

One example

$\left( 32, 32, 3 \right)$

$32 \cdot 32 \cdot 3 = \boxed{3072}$

$P$

# Outline for March 26

- SVM extensions

- Introduction to neural networks

- Fully connected (FC) neural networks

- Image data format and intro to Lab 7

# Lab 7 data pre-processing

- It is helpful to have our data be zero-centered, so we will subtract off the mean

- It is also helpful to have the features be on the same scale, so we will divide by the standard deviation

- We will compute the mean and std with respect to the *training data*, then apply the same transformation to all datasets

# Lab 7 data pre-processing

- Input is now itself a multi-dimensional array
  - Also known as a **tensor**!

- For images, often the shape of each image will be (width, height, 3) for RGB channels

- Need to "**flatten**" or "unravel" for fully connected networks