Support Vector Machines
○○○○○○○○○○○○○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○

# CS540 Introduction to Artificial Intelligence
## Lecture 5

Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

July 5, 2022

# Adverse Selection

Quiz

# Adverse Selection, ID

Quiz

# Course Rhythm

## Admin

# Course Grades

## Admin

Support Vector Machines

○○○○●○○○○○○○○○○○○○○

Subgradient Descent

○○○○○

Kernel Trick

○○○○○○○○○○○○

# Sharing Solutions

## Admin

# Sharing Solutions

- For solution sharing, please make sure it is Piazza note, not a Piazza question.

- For actual questions, please use a different name, *e.g.* "$M2Q1$ Question" or "Question about $M2Q1$".

- Make sure you tag the post correctly: $m2$, $m3$, or $d1$ in order to get the points.

- Please sign up before making the post and please do not sign up for more than 4 questions per week.

- I will either "good note" the post or leave a comment: if I leave a comment, please update your answers, reply to my comment, and remember to make the reply "unresolved" so I can see.

Support Vector Machines
0000000●00000000000

Subgradient Descent
00000

Kernel Trick
000000000000

# Maximum Margin Diagram

## Motivation

Support Vector Machines
○○○○○○○●○○○○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○

# SVM Weights
## Quiz

Support Vector Machines

○○○○○○○○●○○○○○○○○○

Subgradient Descent

○○○○○

Kernel Trick

○○○○○○○○○○○○

# SVM Weights Diagram

## Quiz

# SVM Weights

## Quiz

Support Vector Machines
○○○○○○○○○○●○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○

# SVM Weights Diagram
## Quiz

Support Vector Machines
○○○○○○○○○○○○●○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○

# Constrained Optimization Diagram

### Definition

Support Vector Machines
○○○○○○○○○○○○○●○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○

# Constrained Optimization Derivation

Definition

Support Vector Machines
○○○○○○○○○○○○○○●○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○

# Soft Margin Diagram

## Definition

Support Vector Machines
○○○○○○○○○○○○○○○●○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○

# Soft Margin Derivation

## Definition

Support Vector Machines
○○○○○○○○○○○○○○○●○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○

# SVM Formulations

## Definition

Support Vector Machines
○○○○○○○○○○○○○○○○○●○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○

# Soft Margin
## Quiz

Support Vector Machines

○○○○○○○○○○○○○○○○○●

Subgradient Descent

○○○○○

Kernel Trick

○○○○○○○○○○○○

# Soft Margin 2

## Quiz

Support Vector Machines
0000000000000000000

Subgradient Descent
●0000

Kernel Trick
00000000000

# Subgradient Descent
Definition

$$\min_w \frac{\lambda}{2} w^T w + \frac{1}{n} \sum_{i=1}^n \max\left\{0, 1 - (2y_i - 1)\left(w^T x_i + b\right)\right\}$$

- The gradient for the above expression is not defined at points with $1 - (2y_i - 1)\left(w^T x_i + b\right) = 0$.
- Subgradient can be used instead of a gradient.

# Subgradient 1

## Quiz

Support Vector Machines
○○○○○○○○○○○○○○○○○○○

Subgradient Descent
○○●○○

Kernel Trick
○○○○○○○○○○○○

# Subgradient 2

## Quiz

# Subgradient Descent Step

### Definition

- One possible set of subgradients with respect to $w$ and $b$ are the following.

$$\partial_w C \ni \lambda w - \sum_{i=1}^{n} (2y_i - 1) \, x_i \, \mathbb{1}_{\{(2y_i-1)(w^T x_i + b) \geqslant 1\}}$$

$$\partial_b C \ni - \sum_{i=1}^{n} (2y_i - 1)) \, \mathbb{1}_{\{(2y_i-1)(w^T x_i + b) \geqslant 1\}}$$

- The gradient descent step is the same as usual, using one of the subgradients in place of the gradient.

Support Vector Machines
0000000000000000000

Subgradient Descent
0000●

Kernel Trick
00000000000

# Regularization Parameter
### Definition

$$w = w - \alpha \sum_{i=1}^{n} z_i \mathbb{1}_{\{z_i w^T x_i \geqslant 1\}} x_i - \lambda w$$

$$z_i = 2y_i - 1, i = 1, 2, ..., n$$

- $\lambda$ is usually called the regularization parameter because it reduces the magnitude of $w$ the same way as the parameter $\lambda$ in $L2$ regularization.

- The stochastic subgradient descent algorithm for SVM is called PEGASOS: Primal Estimated sub-GrAdient SOlver for Svm.

Support Vector Machines
○○○○○○○○○○○○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
●○○○○○○○○○○○○

# Kernel Trick $1D$ Diagram

## Motivation

Support Vector Machines
○○○○○○○○○○○○○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○●○○○○○○○○○○○

# Kernelized SVM

Definition

- With a feature map $\varphi$, the SVM can be trained on new data points $\{(\varphi(x_1), y_1), (\varphi(x_2), y_2), ..., (\varphi(x_n), y_n)\}$.
- The weights $w$ correspond to the new features $\varphi(x_i)$.
- Therefore, test instances are transformed to have the same new features.

$$\hat{y}_i = \mathbb{1}_{\{w^T\varphi(x_i) \geqslant 0\}}$$

Support Vector Machines
○○○○○○○○○○○○○○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○●○○○○○○○○○○

# Kernel Trick for XOR

Quiz

Support Vector Machines
○○○○○○○○○○○○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○●○○○○○○○○

# Kernel Trick for XOR

Quiz

Support Vector Machines
oooooooooooooooooo

Subgradient Descent
ooooo

Kernel Trick
ooooo●oooooo

# Kernel Matrix
Definition

- The feature map is usually represented by a $n \times n$ matrix $K$ called the Gram matrix (or kernel matrix).

$$K_{ii'} = \varphi\left(x_i\right)^T \varphi\left(x_{i'}\right)$$

Support Vector Machines
OOOOOOOOOOOOOOOOOOO

Subgradient Descent
OOOOO

Kernel Trick
OOOOOO●OOOOOO

# Examples of Kernel Matrix
### Definition

- For example, if $\varphi(x) = \left(x_1^2, \sqrt{2}x_1 x_2, x_2^2\right)$, then the kernel matrix can be simplified.

$$K_{ii'} = \left(x_i^T x_{i'}\right)^2$$

- Another example is the quadratic kernel $K_{ii'} = \left(x_i^T x_{i'} + 1\right)^2$. It can be factored to have the following feature representations.

$$\varphi(x) = \left(x_1^2, x_2^2, \sqrt{2}x_1 x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1\right)$$

Support Vector Machines
○○○○○○○○○○○○○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○●○○○○○

# Examples of Kernel Matrix Derivation

Definition

Support Vector Machines
0000000000000000000

Subgradient Descent
00000

Kernel Trick
000000000000

# Popular Kernels
#### Discussion

- Other popular kernels include the following.

1. Linear kernel: $K_{ii'} = x_i^T x_{i'}$

2. Polynomial kernel: $K_{ii'} = \left(x_i^T x_{i'} + 1\right)^d$

3. Radial Basis Function (Gaussian) kernel:
   $$K_{ii'} = \exp\left(-\frac{1}{\sigma^2}\left(x_i - x_{i'}\right)^T\left(x_i - x_{i'}\right)\right)$$

- Gaussian kernel has infinite-dimensional feature representations. There are dual optimization techniques to find $w$ and $b$ for these kernels.

# Kernel Matrix

## Quiz

Support Vector Machines
ooooooooooooooooooo

Subgradient Descent
ooooo

Kernel Trick
oooooooooo●oo

# Kernel Matrix Math

Quiz

Support Vector Machines
○○○○○○○○○○○○○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○●○

# Kernel Matrix 2

## Quiz

Support Vector Machines
○○○○○○○○○○○○○○○○○○○

Subgradient Descent
○○○○○

Kernel Trick
○○○○○○○○○○○○●

# Kernel Matrix Math 2

Quiz