

LECTURE 11

SPRING 2021

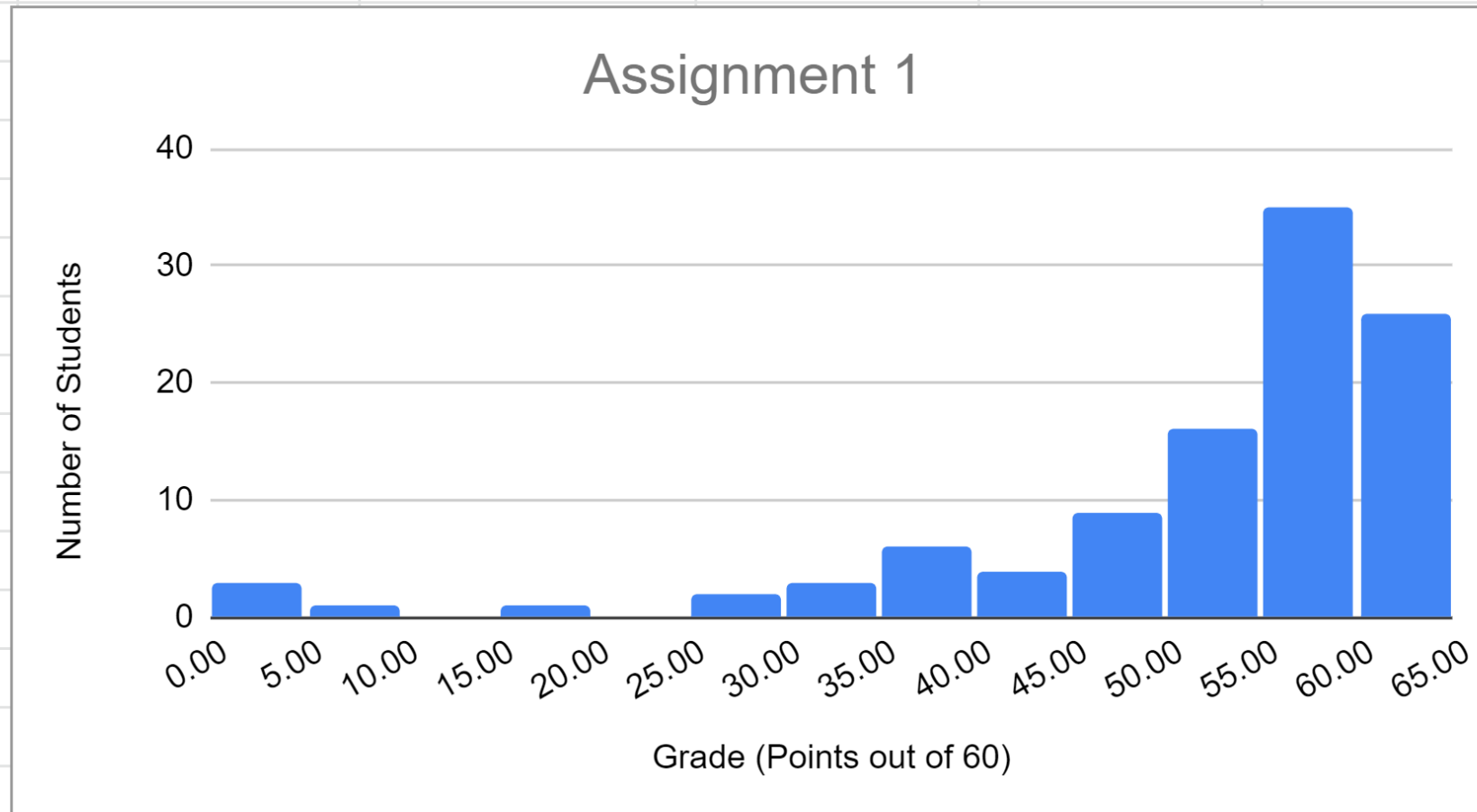
APPLIED MACHINE LEARNING
CIHANG XIE

SLIDE CREDIT:

DHRUV BATRA

ANDREW ZISSERMAN

HW1



Average: 50.54
Median: 56.00
STD: 13.28

Q2@HW2

```
def computeRegularizedCost(X, y, theta, lambdah):
    m = y.size
    J = (np.sum((X @ theta - y)**2))/2/m + lambdah * np.sum([i**2 for i in theta])
    return J

def gradientDescentWithRegularization(X, y, theta, alpha, num_iters, lambdah):
    m = y.shape[0]
    theta = theta.copy()
    J_history = []
    for i in range(num_iters):
        theta_tmp = []
        for j in range(len(theta)): # partial derivative
            gradient = (alpha/m) * np.sum((X @ theta - y) * X[:,j]) + 2 * lambdah * theta[j]
            new_theta = theta[j] - gradient
            theta_tmp.append(new_theta)
        theta = theta_tmp
        J_history.append(computeCost(X, y, theta))
    return theta, J_history
```

WEEK 5 GROUP ACTIVITIES

- Some groups have completed the SVM code
- There is one group has not submitted it

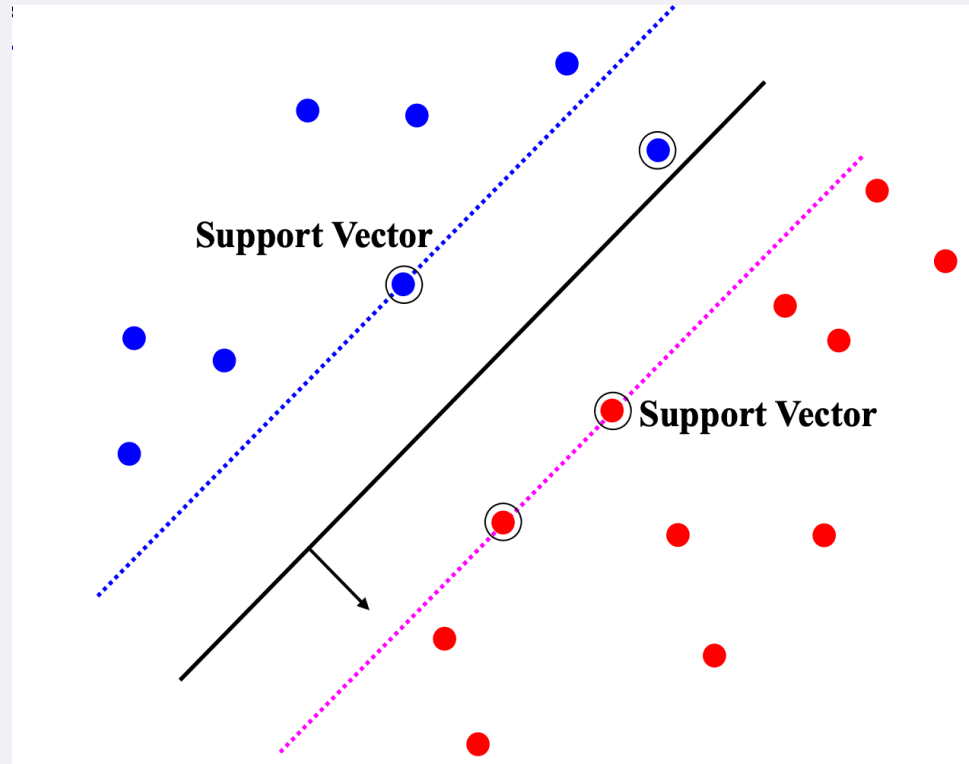
TODAY

- Support Vector Machine
 - review
 - Lagrangian duality
 - kernel trick

SUPPORT VECTOR MACHINE (SVM)

$$\min_{\theta} \frac{1}{2} \theta^T \theta$$

s.t. $y^{(i)} (\theta^T x^{(i)} + b) \geq 1, \forall i$



SOFT MARGIN SVM

$$\begin{aligned} \min_{\theta, \xi, b} \quad & \frac{1}{2} \theta^T \theta + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & y^{(i)} (\theta^T x^{(i)} + b) + \xi_i \geq 1, \\ & \xi_i \geq 0, \forall i \end{aligned}$$

ξ_i is the “slack” variable

- for $0 < \xi_i \leq 1$ point is between margin and correct side of hyperplane. This is a margin violation
- for $\xi_i > 1$ point is misclassified

SOFT MARGIN SVM

$$\begin{aligned} \min_{\theta, \xi, b} \quad & \frac{1}{2} \theta^T \theta + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & y^{(i)} (\theta^T x^{(i)} + b) + \xi_i \geq 1, \\ & \xi_i \geq 0, \forall i \end{aligned}$$

C is a regularization parameter:

- small C allows constraints to be easily ignored → large margin
- large C makes constraints hard to ignore → narrow margin
- $C = \infty$ enforces all constraints: hard margin

GRADIENT DESCENT FOR SVM

$$y^{(i)}(\theta^T x^{(i)} + b) + \xi_i \geq 1 \text{ \& } \xi_i \geq 0$$



$$\xi_i = \max \{0, 1 - y^{(i)}(\theta^T x^{(i)} + b)\}$$



$$\min_{\theta, b} \frac{1}{2} \theta^T \theta + C \sum_{i=1}^N \max \{0, 1 - y^{(i)}(\theta^T x^{(i)} + b)\}$$

GRADIENT DESCENT FOR SVM

$$\begin{aligned}\text{COST}(\theta, b) &= \frac{1}{2} \theta^T \theta + C \sum_{i=1}^N \max \{0, 1 - y^{(i)} (\theta^T x^{(i)} + b)\} \\ &= \sum_{i=1}^N \left(\frac{1}{2N} \theta^T \theta + C \max \{0, 1 - y^{(i)} (\theta^T x^{(i)} + b)\} \right)\end{aligned}$$

For each data point $x^{(i)}$

$$\frac{\partial \text{Cost}(\theta, b)}{\partial \theta_j} = \begin{cases} \frac{1}{N} \theta_j - C y^{(i)} x_j^{(i)} & , \text{ if } 1 - y^{(i)} (\theta^T x^{(i)} + b) > 0 \\ \frac{1}{N} \theta_j, & \text{ otherwise} \end{cases}$$

$$\frac{\partial \text{Cost}(\theta, b)}{\partial b} = \begin{cases} -C y^{(i)}, & \text{ if } 1 - y^{(i)} (\theta^T x^{(i)} + b) > 0 \\ 0, & \text{ otherwise} \end{cases}$$

RELATIONSHIP TO LOGISTIC REGRESSION

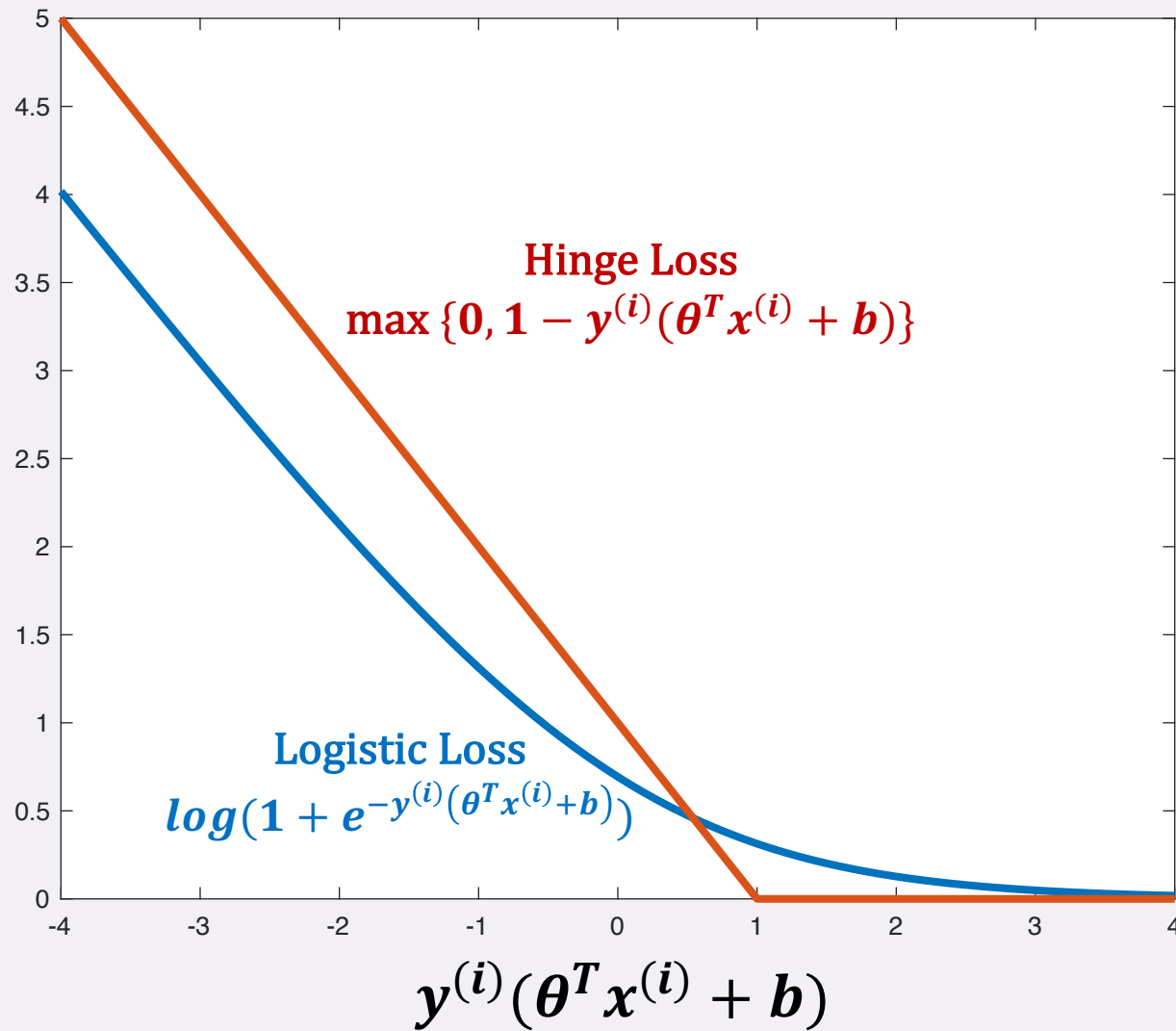
$$\min_{\theta, b} \lambda \theta^T \theta - \sum_i \log P(y^{(i)} | x^{(i)}; \theta, b)$$



$$\min_{\theta, b} \underbrace{\lambda \theta^T \theta}_{\text{Regularization}} + \sum_i \underbrace{\log(1 + e^{-y^{(i)}(\theta^T x^{(i)} + b)})}_{\text{Logistics Loss}}$$

$$\min_{\theta, b} \underbrace{\frac{1}{2} \theta^T \theta}_{\text{Regularization}} + C \sum_{i=1}^N \underbrace{\max\{0, 1 - y^{(i)}(\theta^T x^{(i)} + b)\}}_{\text{Hinge Loss}}$$

RELATIONSHIP TO LOGISTIC REGRESSION



Logistic loss is sometime viewed as the **smooth version** of the Hinge loss.



DUAL FORMULATION

LAGRANGIAN DUALITY

Primal

$$\begin{aligned} \min_{\theta, b} \quad & \frac{1}{2} \theta^T \theta \\ \text{s.t.} \quad & y^{(i)}(\theta^T x^{(i)} + b) \geq 1, \forall i \end{aligned}$$

The Lagrangian for the primal problem:

$$L = \frac{1}{2} \theta^T \theta + \sum_i \alpha_i [1 - y^{(i)}(\theta^T x^{(i)} + b)] \text{ where } \alpha_i \geq 0 \text{ are Lagrange multipliers.}$$

Now we want to solve: $\min_{\theta, b} \max_{\alpha} L(\theta, b, \alpha)$

$$\begin{aligned} \text{if } 1 > y^{(i)}(\theta^T x^{(i)} + b) &\rightarrow \text{min won't let it happen} \\ \text{if } 1 = y^{(i)}(\theta^T x^{(i)} + b) &\rightarrow \text{equivalent to } \min_{\theta, b} \frac{1}{2} \theta^T \theta \\ \text{if } 1 < y^{(i)}(\theta^T x^{(i)} + b) &\rightarrow \text{equivalent to } \min_{\theta, b} \frac{1}{2} \theta^T \theta \end{aligned}$$

LAGRANGIAN DUALITY

Primal

$$\begin{aligned} & \min_{\theta, b} \frac{1}{2} \theta^T \theta \\ \text{s.t. } & y^{(i)}(\theta^T x^{(i)} + b) \geq 1, \forall i \end{aligned}$$

The Lagrangian for the primal problem:

$$L = \frac{1}{2} \theta^T \theta + \sum_i \alpha_i [1 - y^{(i)}(\theta^T x^{(i)} + b)] \text{ where } \alpha_i \geq 0 \text{ are Lagrange multipliers.}$$

Now we want to solve:

$$\min_{\theta, b} \max_{\alpha} L(\theta, b, \alpha)$$



Slater's condition from convex optimization guarantees that these two optimization problems are equivalent!

This is the
dual problem



$$\max_{\alpha} \min_{\theta, b} L(\theta, b, \alpha)$$

LAGRANGIAN DUALITY

Primal

$$\begin{aligned} & \min_{\theta} \frac{1}{2} \theta^T \theta \\ \text{s.t. } & y^{(i)}(\theta^T x^{(i)} + b) \geq 1, \forall i \end{aligned}$$

$$\max_{\alpha} \min_{\theta, b} \frac{1}{2} \theta^T \theta + \sum_i \alpha_i [1 - y^{(i)}(\theta^T x^{(i)} + b)]$$

Dual

$$\frac{\partial L}{\partial \theta} = \theta - \sum_i \alpha_i y^{(i)} x^{(i)} = 0$$

$$\frac{\partial L}{\partial b} = - \sum_i \alpha_i y^{(i)} = 0$$

LAGRANGIAN DUALITY

Primal

$$\begin{aligned} & \min_{\boldsymbol{\theta}} \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} \\ \text{s.t. } & \mathbf{y}^{(i)} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} + \mathbf{b}) \geq 1, \forall i \end{aligned}$$

$$\begin{aligned} \max_{\boldsymbol{\alpha}} L(\boldsymbol{\alpha}) &= \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \mathbf{y}^{(i)} \mathbf{y}^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)} + \sum_i \alpha_i - \sum_i \sum_j \alpha_i \alpha_j \mathbf{y}^{(i)} \mathbf{y}^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)} \\ &= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \mathbf{y}^{(i)} \mathbf{y}^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)} \end{aligned}$$

Dual

$$\begin{aligned} \text{s.t. } & \alpha_i \geq 0, \forall i \\ & \sum_i \alpha_i \mathbf{y}^{(i)} = \mathbf{0} \end{aligned}$$

LAGRANGIAN DUALITY (SOFT MARGIN)

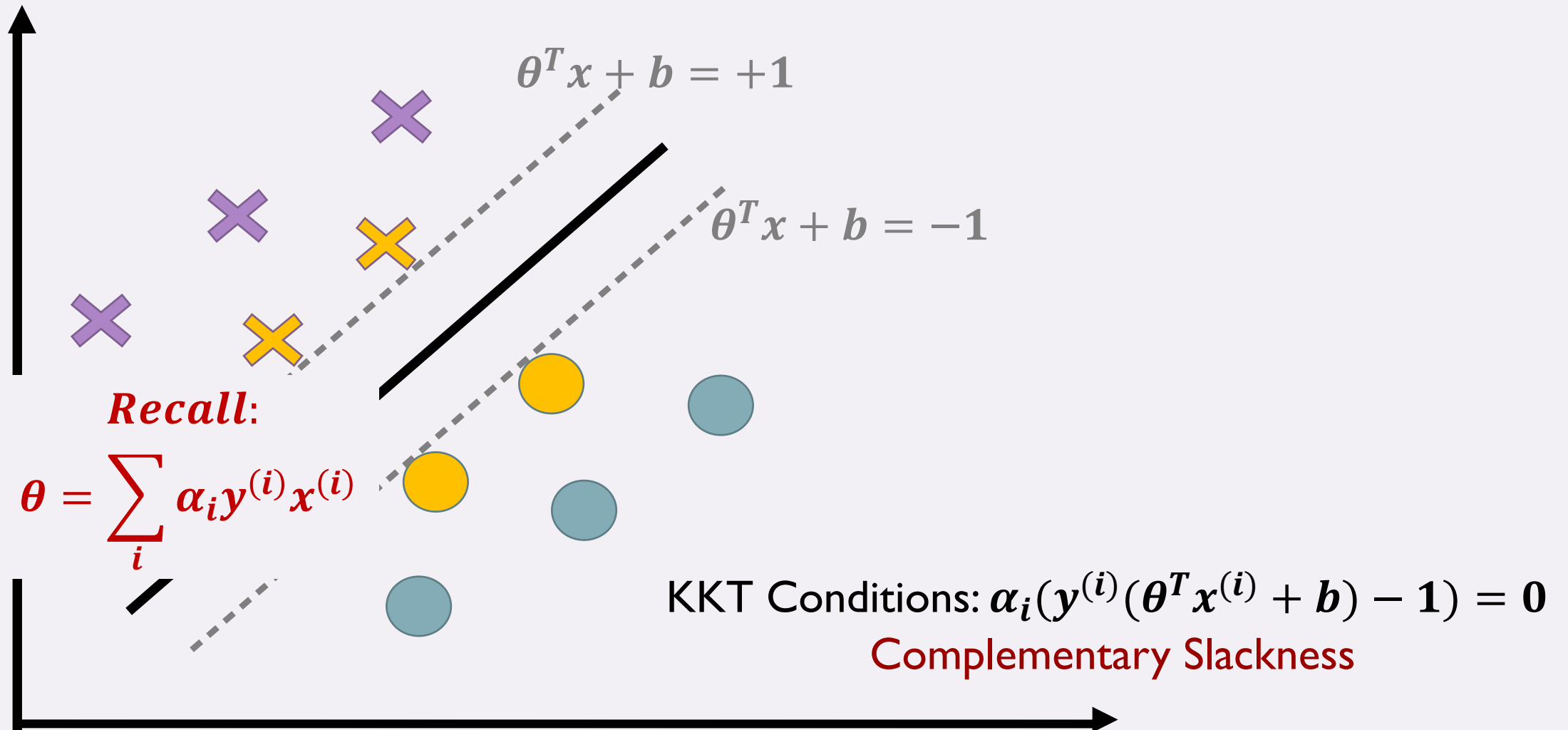
Primal

$$\begin{aligned} \min_{\theta, \xi, b} \quad & \frac{1}{2} \theta^T \theta + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & y^{(i)} (\theta^T x^{(i)} + b) + \xi_i \geq 1, \\ & \xi_i \geq 0, \forall i \end{aligned}$$

Dual

$$\begin{aligned} \max_{\alpha} \quad & L(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)T} x^{(j)} \\ \text{s.t.} \quad & C \geq \alpha_i \geq 0, \forall i \\ & \sum_i \alpha_i y^{(i)} = 0 \end{aligned}$$

SUPPORT VECTORS



KERNEL TRICK

$$x \rightarrow \phi(x) \quad R^d \rightarrow R^D$$

$$\begin{aligned} & \min_{\theta} \frac{1}{2} \theta^T \theta \\ \text{s.t.} \quad & y^{(i)} (\theta^T \phi(x^{(i)}) + b) \geq 1, \forall i \end{aligned}$$

If $D \gg d$ then there are many more parameters to learn for θ Can this be avoided?

KERNEL TRICK

$$x \rightarrow \phi(x) \quad R^d \rightarrow R^D$$

$$\max_{\alpha} L(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} \phi(x^{(i)})^T \phi(x^{(j)})$$

$$\text{s.t.} \quad \alpha_i \geq 0, \forall i \\ \sum_i \alpha_i y^{(i)} = 0$$

- $\phi(x)$ only occurs in pairs $\phi(x)^T \phi(x)$
- Once the scalar products are computed, only the N dimensional vector α needs to be learnt; it is not necessary to learn in the D dimensional space, as it is for the primal

KERNEL TRICK

Write $k(x_j, x_i) = \phi(x)^T \phi(x)$. This is known as **kernel**

$$\max_{\alpha} L(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} k(x^{(i)}, x^{(j)})$$

$$\text{s.t. } \alpha_i \geq 0, \forall i$$
$$\sum_i \alpha_i y^{(i)} = 0$$

KERNEL TRICK

$$\Phi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix} \quad \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

- Classifier can be learnt and applied without explicitly computing $\Phi(x)$
- All that is required is the kernel $k(x, z) = (x^T z)^2$

EXAMPLE KERNELS

- Linear $k(x, z) = x^T z$
- Polynomial of degree exactly d $k(x, z) = (x^T z)^d$
- Polynomial of degree up to d $k(x, z) = (x^T z + 1)^d$
- Gaussian $k(x, z) = \exp(-||x-z||^2)$
--- Infinite dimensional feature space

PROPERTIES

- Weight vector is a linear combination of data
- Only the points on the margins matter
 - Ignore the rest
- Only inner products between data points matter
 - Can use Kernels!
- Provides the widest possible separation between classes

QUIZ 3 (DUE TONIGHT)



Quiz 3

Not available until May 4 at 3:00pm | 4 pts | 5 Questions

HW2 (DUE MAY 9)



HW2

Available until May 9 at 11:59pm | Due May 9 at 11:59pm | 60 pts

A decorative graphic on the left side of the slide consisting of two parallel, wavy vertical lines. The inner line is a light purple color, and the outer line is a slightly darker shade of purple. They extend from the top to the bottom of the slide.

QUESTIONS?