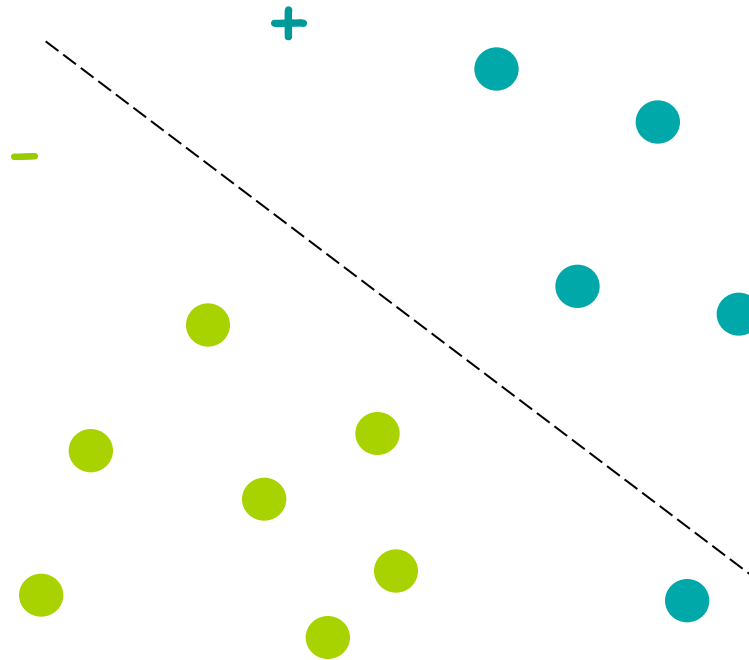


# Optimization Methods in Machine Learning

Lecture 12,  
IPMs for SVMs

# An Interior Point Method for SVM

# Support Vector Machines



## Optimization Problem

Two sets of points:  $X_+ \subset \mathbf{R}^p$  and  $X_- \subset \mathbf{R}^p$

Total number of points:  $n$

$$\begin{aligned} \min_{\xi, w, \beta} \quad & \frac{1}{2} w^\top w + c \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & -y_i(w^\top x_i + \beta) \leq -1 + \xi_i, \quad i = 1, \dots, n \\ & \xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

$$y_i = 1 \quad \text{if } x_i \in X_+$$

$$y_i = -1 \quad \text{if } x_i \in X_-$$

## Optimization Problem

At optimality  $w^* = \sum_{i=1}^n \alpha_i y_i x_i$ ,  $0 \leq \alpha_i \leq c$

$$\begin{aligned} \min_{\alpha, \beta, \xi} \quad & \frac{1}{2} \alpha^\top Q \alpha + c \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & -Q\alpha + y\beta + s_i - \xi_i = -1, \quad i = 1, \dots, n \\ & s_i \geq 0, \xi \geq 0, 0 \leq \alpha_i \leq c, \quad i = 1, \dots, n, \end{aligned}$$

$$Q := D_y X X^\top D_y \quad \Leftrightarrow \quad Q_{ij} = y_i y_j x_i^\top x_j$$

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^\top Q \alpha - e^\top \alpha \\ \text{s.t.} \quad & y^\top \alpha = 0, \\ & 0 \leq \alpha \leq c, \end{aligned}$$

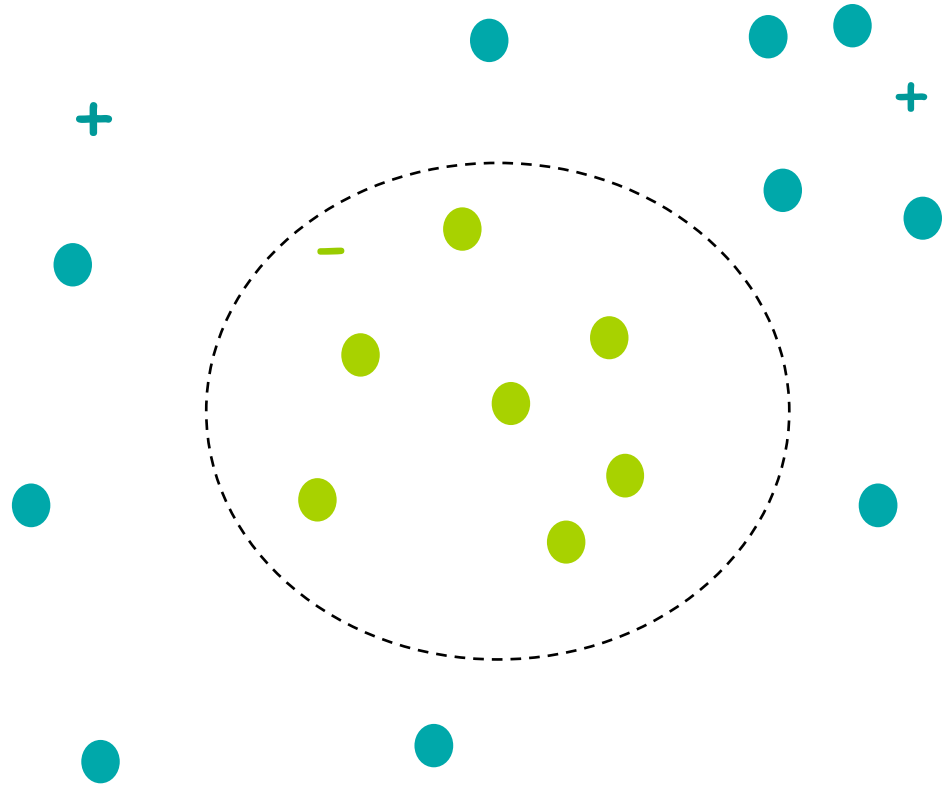
# Kernel SVM

$$Q_{ij} = y_i y_j x_i^\top x_j \rightarrow Q_{ij} = y_i y_j \phi(x_i)^\top \phi(x_j) = y_i y_j K(x_i, x_j)$$

Kernel operation:  $K(x_i, x_j) = \phi(x_i)^\top \phi(x_j)$

Examples:

- $K(x_i, x_j) = \exp^{-\|x_i - x_j\|^2 / 2\sigma^2}$
- $K(x_i, x_j) = (x_i^\top x_j / a_1 + a_2)^d$



## Optimality Conditions

Dual problem

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^{\top} Q \alpha - e^{\top} \alpha \\ \text{s.t.} \quad & y^{\top} \alpha = 0, \\ & 0 \leq \alpha \leq c, \end{aligned}$$

KKT conditions

$$\begin{aligned} \alpha_i s_i &= 0, \quad i = 1, \dots, n, \\ (c - \alpha_i) \xi_i &= 0, \quad i = 1, \dots, n, \\ y^{\top} \alpha &= 0, \\ -Q \alpha + y \beta + s - \xi &= -e, \\ 0 \leq \alpha \leq c, \quad s &\geq 0, \quad \xi \geq 0. \end{aligned}$$

## Interior Point methods

### Relaxed KKT conditions

$$\alpha_i s_i = \mu, \quad i = 1, \dots, n,$$

$$(c - \alpha_i) \xi_i = \mu \quad i = 1, \dots, n,$$

$$y^\top \alpha = 0,$$

$$-Q\alpha + y\beta + s - \xi = -e,$$

$$0 < \alpha < c, \quad s > 0, \quad \xi > 0.$$



## A Newton step of IPM

Linearize perturbed KKT conditions

$$\alpha_i \Delta s_i + s_i \Delta \alpha_i = \mu - s_i \alpha_i, \quad i = 1, \dots, n,$$

$$(c - \alpha_i) \Delta \xi_i - \xi \Delta \alpha_i = \mu - (c - \alpha_i) \xi_i \quad i = 1, \dots, n,$$

$$y^\top \Delta \alpha = -y^\top \alpha,$$

$$-Q \Delta \alpha + y \Delta \beta + \Delta s - \Delta \xi = -e + Q \alpha - y \beta - s + \xi,$$

Let  $\mathcal{A} = \text{diag}(\alpha)$ ,  $\mathcal{S} = \text{diag}(S)$  and  $\Xi = \text{diag}(\xi)$

$$\Delta s = \mathcal{A}^{-1} \mu e - s - \mathcal{A}^{-1} \mathcal{S} \Delta \alpha,$$

$$\Delta \xi = (C - \mathcal{A})^{-1} \mu e - \xi + (C - \mathcal{A})^{-1} \Xi \Delta \alpha,$$

$$y^\top \Delta \alpha = -y^\top \alpha,$$

$$-Q \Delta \alpha + y \Delta \beta + \Delta s - \Delta \xi = -e + Q \alpha - y \beta - s + \xi,$$

## Solving the linear system

$$y^\top \Delta\alpha = -y^\top \alpha,$$

$$\begin{aligned} & -Q\Delta\alpha + y\Delta\beta + \mathcal{A}^{-1}\mu e - s - \mathcal{A}^{-1}S\Delta\alpha - (C - \mathcal{A})^{-1}\mu e - \xi - (C - \mathcal{A})^{-1}\Xi\Delta\alpha \\ & = -e + Q\alpha - y\beta - s + \xi, \end{aligned}$$



$$y^\top \Delta\alpha = -y^\top \alpha,$$

$$\begin{aligned} & -(Q + \mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi)\Delta\alpha + y\Delta\beta \\ & = -e + Q\alpha - y\beta - \mathcal{A}^{-1}\mu e + (C - \mathcal{A})^{-1}\mu e, \end{aligned}$$

## Solving the linear system

$$\begin{bmatrix} y^\top & 0 \\ -(Q + \mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi) & y \end{bmatrix} \begin{pmatrix} \Delta\alpha \\ \Delta\beta \end{pmatrix} = \begin{pmatrix} -y^\top \alpha \\ -e + Q\alpha - y\beta - \mathcal{A}^{-1}\mu e + (C - \mathcal{A})^{-1}\mu e \end{pmatrix}$$



$$[y^\top (Q + \mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi)^{-1} y] \Delta\beta = -y^\top \alpha + y^\top (Q + \mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi)^{-1} (-e + Q\alpha - y\beta - \mathcal{A}^{-1}\mu e - (C - \mathcal{A})^{-1}\mu e)$$

## Forming storing inverting matrices

Need to compute  $Q+D$  and solve a system of lin equation with it (factorize)

$$Q_{ij} = y_i y_j K(x_i, x_j)$$

$Q$  is  $n \times n$ , typically dense matrix

- $K(x_i, x_j) = \exp^{-\|x_i - x_j\|^2 / 2\sigma^2}$
- $K(x_i, x_j) = (x_i^\top x_j / a_1 + a_2)^d$

$Q+D$  has to be factorized at each iteration –  $O(n^3)$  flops

## Scherman-Morrison-Woodbury formula

$$(M + UV^T)^{-1} = M^{-1} - M^{-1}U(I + V^T M^{-1}U)^{-1}V^T M^{-1}$$

$$\text{Let } Q_{ij} = y_i y_j \phi(x_i)^T \phi(x_j) \Rightarrow Q = VV^T$$

k = the number of columns in V is the dimension of feature space -  $\phi(x)$

$$(D + VV^T)^{-1} = D^{-1} - D^{-1}V(I + V^T D^{-1}V)^{-1}V^T D^{-1}$$

$O(n)$              $O(nk)$              $O(k^2n)$              $O(nk)$

Per iteration complexity is  $O(nk^2)$  and storage is  $O(nk)$

## Return to the linear formulation

$$\begin{aligned} \min_{\xi, w, \beta} \quad & \frac{1}{2} w^\top w + c \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & -y_i(w^\top x_i + \beta) \leq -1 + \xi_i, \quad i = 1, \dots, n \\ & \xi \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

$$L(w, b, \alpha, \xi) = \frac{1}{2} w^\top w + c \sum_{i=1}^n \xi_i - \sum_i \alpha_i (y_i(w^\top x_i + \beta) - 1 + \xi_i) - \nu^\top \xi$$

$$\nabla_w L = w - \sum_{i=1}^n \alpha_i y_i x_i = 0$$

$$\nabla_\xi L = c - \alpha - \nu = 0$$

$$\nabla_\beta L = y^\top \alpha = 0$$

$$\alpha \geq 0, \quad \nu \geq 0$$

## Optimality Conditions

### Perturbed KKT conditions

$$\alpha_i s_i = \mu, \quad i = 1, \dots, n,$$

$$(c - \alpha_i) \xi_i = \mu, \quad i = 1, \dots, n,$$

$$y^\top \alpha = 0,$$

$$w - \sum_{i=1}^n \alpha_i y_i x_i = 0,$$

$$-y_i (w^\top x_i + \beta) + s_i - \xi_i = -1, \quad i = 1, \dots, n,$$

$$0 < \alpha < c, \quad s > 0, \quad \xi < 0.$$

Or in vector matrix terms...

## Perturbed KKT conditions

$$As = \mu e,$$

$$(C - A)\xi = \mu e$$

$$y^\top \alpha = 0,$$

$$w - (YX)^\top \alpha = 0,$$

$$-YXw - y\beta + s - \xi = -e$$



## The Newton system

$$\mathcal{A}\Delta s + \mathcal{S}\Delta\alpha = \mu e - \mathcal{A}s$$

$$(\mathcal{C} - \mathcal{A})\Delta\xi - \Xi\Delta\alpha = \mu e - (\mathcal{C} - \mathcal{A})\xi$$

$$y^\top \Delta\alpha = -y^\top \alpha,$$

$$\Delta w - (\mathcal{Y}\mathcal{X})^\top \Delta\alpha = -w + (\mathcal{Y}\mathcal{X})^\top \alpha,$$

$$-\mathcal{Y}\mathcal{X}\Delta w - y\Delta\beta + \Delta s - \Delta\xi = -e - \mathcal{Y}\mathcal{X}w - y\beta - s + \xi$$



$$y^\top \Delta\alpha = -y^\top \alpha,$$

$$\Delta w - (\mathcal{Y}\mathcal{X})^\top \Delta\alpha = -w + (\mathcal{Y}\mathcal{X})^\top \alpha,$$

$$-\mathcal{Y}\mathcal{X}\Delta w - y\Delta\beta - (\mathcal{A}^{-1}\mathcal{S} + (\mathcal{C} - \mathcal{A})^{-1}\Xi)\Delta\alpha =$$

$$-e + \mathcal{Y}\mathcal{X}w + y\beta - s + \xi - \mathcal{A}^{-1}\mu e + s + (\mathcal{C} - \mathcal{A})^{-1}\mu e - \xi$$

## The Newton system

$$\begin{bmatrix} (\mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi) & YX & y \\ (YX)^\top & -I & 0 \\ y^\top & 0 & 0 \end{bmatrix} \begin{pmatrix} \Delta\alpha \\ \Delta w \\ \Delta\beta \end{pmatrix} = \begin{pmatrix} r_1 \\ r_2 \\ r_3 \end{pmatrix}$$

## The Newton system

$$\begin{bmatrix} -(\mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi) & YX & y \\ (YX)^\top & -I & 0 \\ y^\top & 0 & 0 \end{bmatrix} \begin{pmatrix} \Delta\alpha \\ \Delta w \\ \Delta\beta \end{pmatrix} = \begin{pmatrix} r_1 \\ r_2 \\ r_3 \end{pmatrix}$$

$$I + (YX)^\top (\mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi)^{-1} YX \Delta w - y \Delta\beta = r$$

or

$$(YX(YX)^\top + (\mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi)^{-1}) \Delta\alpha - (YX)y \Delta\beta = r$$

## The Newton system

$$\begin{bmatrix} (\mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi) & YX & y \\ (YX)^{\top} & -I & 0 \\ y^{\top} & 0 & 0 \end{bmatrix} \begin{pmatrix} \Delta\alpha \\ \Delta w \\ \Delta\beta \end{pmatrix} = \begin{pmatrix} r_1 \\ r_2 \\ r_3 \end{pmatrix}$$

$O(nk^2)$

$$I + (YX)^{\top} (\mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi)^{-1} YX \Delta w - y \Delta\beta = r$$

or

$O(n^3)$

$$(YX(YX)^{\top} + (\mathcal{A}^{-1}S + (C - \mathcal{A})^{-1}\Xi)^{-1}) \Delta\alpha - (YX)y \Delta\beta = r$$

# Complexity

## Interior point method for nonlinear SVMs:

- Need to solve  $(Q + D)p = r$ .  $Q$  is completely dense, with rank  $k_d$ .
- If  $k_d \sim n$ , then each IPM iteration is  $O(n^3)$  operations and  $O(n^2)$  memory.
- If  $k_d \ll n$ , then each IPM iteration is  $O(nk_d^2)$  operations and  $O(nk_d)$  memory.

## Interior point method for linear SVMs:

- Can solve  $(VV^T + D)p = r$  or  $(I + VD^{-1}V^T)p = r$ .
- $V = XY$  and is as sparse as the data. In large scale cases  $V$  can be sparse and the complexity per step is similar to linear programming.

# Optimization methods for convex problems

- Interior Point methods
  - Best iteration complexity  $O(\log(1/\epsilon))$ , in practice  $<50$ .
  - Worst per-iteration complexity (sometimes prohibitive)
- Active set methods
  - Exponential complexity in theory, often linear in practice.
  - Better per iteration complexity.
- Gradient based methods
  - $O(1/\sqrt{\epsilon})$  or  $O(1/\epsilon)$  iterations
  - Matrix/vector multiplication per iteration
- Nonsmooth gradient based methods
  - $O(1/\epsilon)$  or  $O(1/\epsilon^2)$  iterations
  - Matrix/vector multiplication per iteration
- Block coordinate descent
  - Iteration complexity ranges from unknown to similar to FOMs.
  - Per iteration complexity can be constant.