

CSCE 633: Machine Learning

Lecture 24: SVM Optimization

Texas A&M University
Bobak Mortazavi

Support Vector Machines: Non-separable case

4. Maximize the Lagrangian with respect to dual variables (dual problem)

$$\begin{aligned} \max_{\alpha_i} L = \max_{\alpha_i} & \left\{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N \alpha_i \alpha_{i'} y_i y_{i'} \mathbf{x}_i^T \mathbf{x}_{i'} \right\} \\ \text{s.t. } & \sum_{i=1}^N \alpha_i y_i = 0 \text{ and } 0 \leq \alpha_i \leq C, \text{ for } i = 1, \dots, N \end{aligned}$$

Coordinate Ascent

$$w(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N \alpha_i \alpha_{i'} y_i y_{i'} x_i^T x_{i'}$$

- Unconstrained optimization problem
 - We have already seen gradient descent (or ascent, if we negate the optimization function), now we consider another optimization method called coordinate ascent

Loop until convergence

1 For $i = 1, \dots, m$

1a $\alpha_i = \arg \max_{\hat{\alpha}_i} \omega(\alpha_1, \dots, \hat{\alpha}_i, \dots, \alpha_m)$

- In the innermost loop, hold all variables constant except for some fixed α_i
- Re-optimize w with respect to just the parameter α_i
- When argmax of the inner loop can be performed efficiently, coordinate ascent can be a fairly efficient algorithm

Sequential Minimal Optimization (SMO)

$$w(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N \alpha_i \alpha_{i'} y_i y_{i'} \mathbf{x}_i^T \mathbf{x}_{i'}$$
$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0 \text{ and } 0 \leq \alpha_i \leq C, \text{ for } i = 1, \dots, N$$

Loop until convergence

1 For $i = 1, \dots, m$

1a Select some α_i and α_j to update

1b Re-optimize w wrt α_i and α_j while holding all other α_k 's fixed

- Let's assume that we optimize wrt α_1, α_2 , while $\alpha_3, \dots, \alpha_N$ are constant

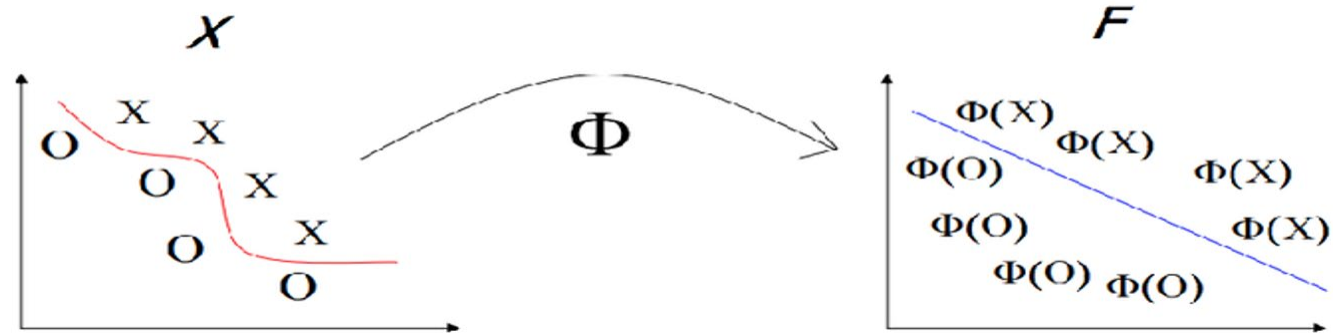
– From the constraint:

$$\alpha_1 y_1 + \alpha_2 y_2 = \sum_{i=3}^N \alpha_i y_i = \zeta \rightarrow \alpha_1 = (\zeta - \alpha_2 y_2) / y_1$$

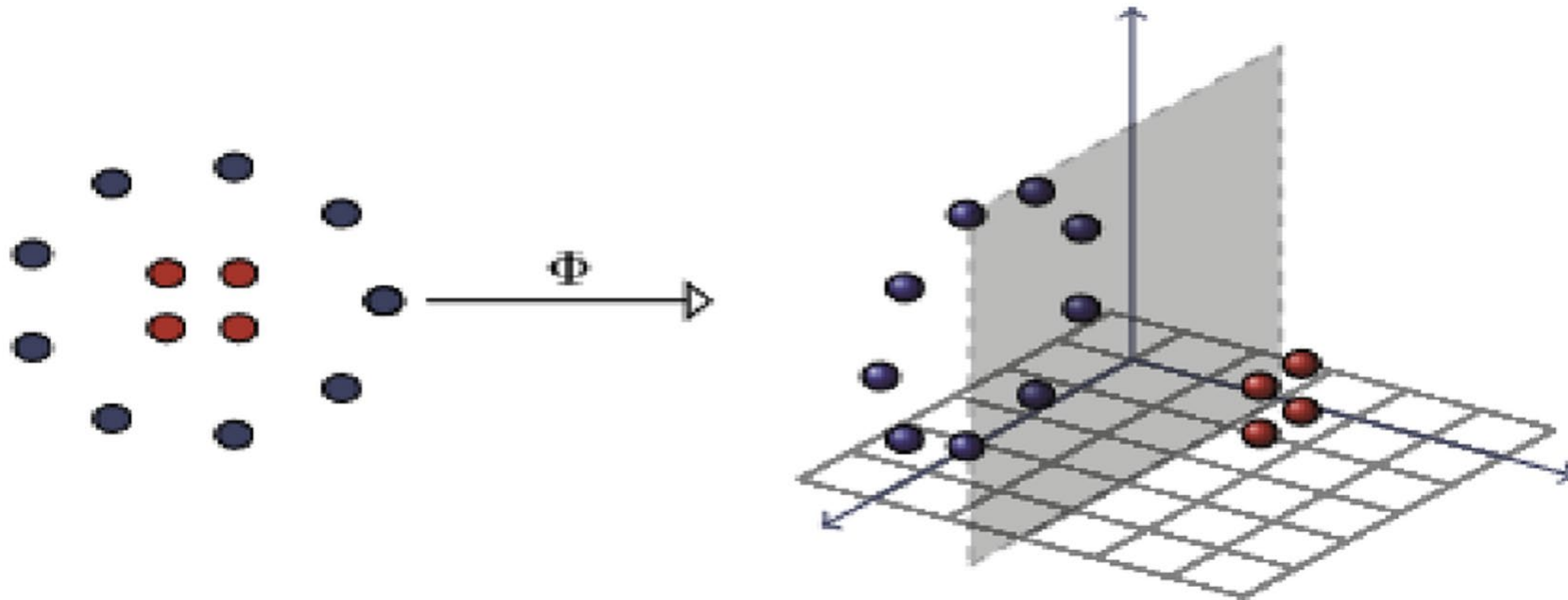
- The objective is $w(\alpha) = w((\zeta - \alpha_2 y_2) / y_1, \alpha_2, \dots, \alpha_N)$ (some quadratic function of $\alpha_2 \rightarrow$ easily solved by setting its derivative to zero)

Kernel Functions

- Motivation
 - Given a set of vectors, there are many tools available for one to use to detect linear relations among the data
 - But what if the relations are non-linear in the original space?
 - Solution: Map the data into a (possibly high dimensional) vector space where linear relations exist among the data, then apply a linear algorithm in this space



Higher Dimensions



Kernel Functions

- A function that takes as its input vectors in the original space and returns the dot product of the vectors in the feature space is called a kernel

- Definition

- A (positive semi-definite) kernel function $L(\cdot, \cdot)$ is a bivariate function for which any \mathbf{x}_m and \mathbf{x}_n

$$K(\mathbf{x}_m, \mathbf{x}_n) = K(\mathbf{x}_n, \mathbf{x}_m) \text{ and } K(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n) \phi(\mathbf{x}_m)$$

- Examples

$$K(\mathbf{x}_n, \mathbf{x}_m) = (\mathbf{x}_n^T \mathbf{x}_m)^2, K(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{\|\mathbf{x}_n - \mathbf{x}_m\|_2^2}{2\sigma^2}\right)$$

- Using kernels, we do not need to embed the data into the space explicitly, because a number of algorithms only require the inner products between input vector. We never need the coordinates of the data in the feature space
 - Kernels can be perceived as similarity measures

Kernel Functions

- Example

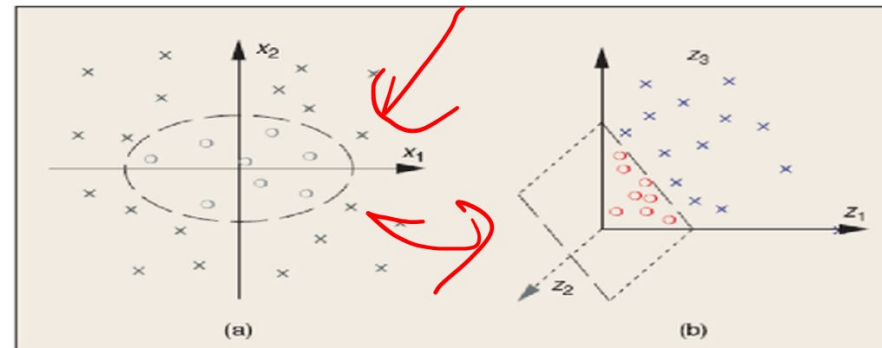
- Consider a two-dimensional input space $\mathbf{X} \in \mathbb{R}^2$ with the feature map:

$$\phi(\mathbf{x}): \mathbf{x} = [x_1, x_2]^T \rightarrow [x_1^2, x_2^2, \sqrt{2}x_1x_2]^T \in \mathbb{R}^3$$

- The inner product in the feature space is

$$\begin{aligned} \phi(\mathbf{x})\phi(\mathbf{z}) &= \langle (x_1^2, x_2^2, \sqrt{2}x_1x_2), (z_1^2, z_2^2, \sqrt{2}z_1z_2) \rangle \\ &= x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 = (x_1z_1 + x_2z_2)^2 = \langle \mathbf{x}, \mathbf{y} \rangle^2 \end{aligned}$$

- Then $K(\mathbf{x}, \mathbf{z}) = \langle \mathbf{x}, \mathbf{y} \rangle^2$



▲ 1. Effect of the map $\phi(x_1, x_2) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$. (a) Input space \mathcal{X} and (b) feature space \mathcal{H} .

Kernel Functions

- Mercer's condition
 - A bivariate function $K(\cdot, \cdot)$ is a positive semidefinite kernel function, if and only if, for any N and any $\mathbf{x}_1, \dots, \mathbf{x}_N$ the following matrix, called the Gram matrix, is positive semi-definite.

$$\mathbf{K} = \begin{pmatrix} \phi(\mathbf{x}_1)\phi(\mathbf{x}_1) & \cdots & \phi(\mathbf{x}_1)\phi(\mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_N)\phi(\mathbf{x}_1) & \cdots & \phi(\mathbf{x}_N)\phi(\mathbf{x}_N) \end{pmatrix} = \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

- Mercer's condition tells us whether or not a prospective kernel is actually a dot product in some space

Examples of Kernel Functions

- Polynomial kernel function with degree of d

$$K(\mathbf{x}_n, \mathbf{x}_m) = (\mathbf{x}_n^T \mathbf{x}_m + c)^d$$

– For $c \geq 0$ and d a positive integer

- Gaussian kernel, RBF (radial basis function) kernel, or Gaussian RBF kernel

$$K(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{\|\mathbf{x}_n - \mathbf{x}_m\|_2^2}{2\sigma^2}\right)$$

- Sigmoid Kernel

$$K(\mathbf{x}_n, \mathbf{x}_m) = \tanh(a(\mathbf{x}_n^T \mathbf{x}_m) + b)$$

- Most of those kernels have parameters to be tuned: d , c , σ^2 , etc. They are hyperparameters and are often tuned on holdout data or with cross-validation

Examples of Kernel Functions

- Document Similarity
 - Let x_{ij} be the # times a word j occurs in a document i (bag-of-words)
 - Cosine similarity between documents i and i' counts the number of shared words

$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = \tanh(a(\mathbf{x}_i^T \mathbf{x}_{i'}) + b)$$

- Edit distance (e.g. gene alignment)
 - # insertions, deletions, substitutions it takes to convert one gene to another

Rules for composing kernels

- There are infinite number of kernels to use
 - If $K(\mathbf{x}_i, \mathbf{x}_{i'})$ is a kernel, then $cK(\mathbf{x}_i, \mathbf{x}_{i'})$, $c > 0$, is a kernel
 - If $K(\mathbf{x}_i, \mathbf{x}_{i'})$ is a kernel, then $e^{K(\mathbf{x}_i, \mathbf{x}_{i'})}$ is a kernel
 - If $K_1(\mathbf{x}_i, \mathbf{x}_{i'})$ and $K_2(\mathbf{x}_i, \mathbf{x}_{i'})$ are kernels, then $\alpha K_1(\mathbf{x}_i, \mathbf{x}_{i'}) + \omega K_2(\mathbf{x}_i, \mathbf{x}_{i'})$, $\alpha, \omega > 0$ is a kernel
 - If $K_1(\mathbf{x}_i, \mathbf{x}_{i'})$ and $K_2(\mathbf{x}_i, \mathbf{x}_{i'})$ are kernels, then $K_1(\mathbf{x}_i, \mathbf{x}_{i'})K_2(\mathbf{x}_i, \mathbf{x}_{i'})$ is a kernel
- In practice, using which kernel, or which kernels to compose a new kernel, remains somewhat of a “black art”, though most people will start with polynomial and Gaussian RBF kernels.
- Example: Audio-visual speech recognition, where notion of similarity is differently defined for speech and image

Kernel Trick

- Many learning methods depend on computing inner products between features, e.g. linear regression
- For those methods, we can use a kernel function in the place of the inner products, i.e. “kernelizing” the methods, thus, introducing nonlinear features/basis
- Instead of first transforming the original features into the new feature space and then computing the inner product, we can compute the inner product in the new feature space directly through the kernel function.

Support Vector Machines: Kernel Trick

- Set of basis functions $\mathbf{z} = \phi(\mathbf{x})$
- Map the problem into the new space and solve as before
- For SVM, this leads to certain simplifications
- Setup for two classes
 - Input: $\mathbf{x} \in \mathbb{R}^N$
 - Output: $y \in \{-1, 1\}$
 - Training data: $\mathcal{D}^{train} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$
 - Non-separable model:

$$f(\mathbf{x}_i) = \begin{cases} 1 & \text{if } w^T \phi(\mathbf{x}_i) + w_0 > -(1 - \epsilon_i) \\ -1 & \text{if } w^T \phi(\mathbf{x}_i) + w_0 \leq -(1 - \epsilon_i) \end{cases}$$

Support Vector Machines: Linearly separable case

$$\min_w \frac{1}{2} w^T w + C \sum_{i=1}^N \epsilon_i, \text{ such that (s.t.) } y_i(w^T \phi(\mathbf{x}_i) + w_0) \geq 1 - \epsilon_i \text{ and } \epsilon_i > 0 \text{ for } i = 1, 2, \dots, N$$

1. Formulate Lagrangian function (primal problem)

$$L = \frac{1}{2} \|\mathbf{w}\|_2^2 - C \sum_{i=1}^N \epsilon_i - \sum_{i=1}^N \alpha_i (y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + w_0) - 1 + \epsilon_i) - \sum_{i=1}^N \mu_i \epsilon_i$$

2. Minimize Lagrangian to solve for primal variables w and w_0

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \phi(\mathbf{x}_i)$$

$$\frac{\partial L}{\partial w_0} = 0 \rightarrow 0 = \sum_{i=1}^N \alpha_i y_i$$

$$\frac{\partial L}{\partial \epsilon_i} = 0 \rightarrow 0 = C - \alpha_i - \mu_i$$

Support Vector Machines: Linearly separable case

$$\min_w \frac{1}{2} w^T w + C \sum_{i=1}^N \epsilon_i, \text{ such that (s.t.) } y_i(w^T x_i + w_0) \geq 1 - \epsilon_i \text{ and } \epsilon_i > 0 \text{ for } i = 1, 2, \dots, N$$

3. Substitute the primal variables w and w_0 into the Lagrangian and express in terms of dual variables α_i

$$\begin{aligned} L &= \frac{1}{2} \|\mathbf{w}\|_2^2 - C \sum_{i=1}^N \epsilon_i - \mathbf{w}^T \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i - w_0 \sum_{i=1}^N \alpha_i y_i + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i \epsilon_i - \sum_{i=1}^N \mu_i \epsilon_i \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N \alpha_i \alpha_{i'} y_i y_{i'} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_{i'}) \end{aligned}$$

Support Vector Machines: Kernel Trick

- Instead of transforming \mathbf{x}_i and $\mathbf{x}_{i'}$ through ϕ and computing their inner product, we directly apply the kernel function to the original space
- Lagrangian for the dual problem

$$L = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N \alpha_i \alpha_{i'} y_i y_{i'} K(\mathbf{x}_i, \mathbf{x}_{i'})$$

$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_{i'})$$

- For taking a decision in the test set

$$\sum_{i'=1}^N \alpha_i y_i K(\mathbf{x}_i, \mathbf{x})$$

- The cost of computing the primal variables is $O(N^3)$, while for the dual variable is $O(D^3)$. A kernel method can be useful in high dimensional settings

Support Vector Machines: Multiple Kernel Learning

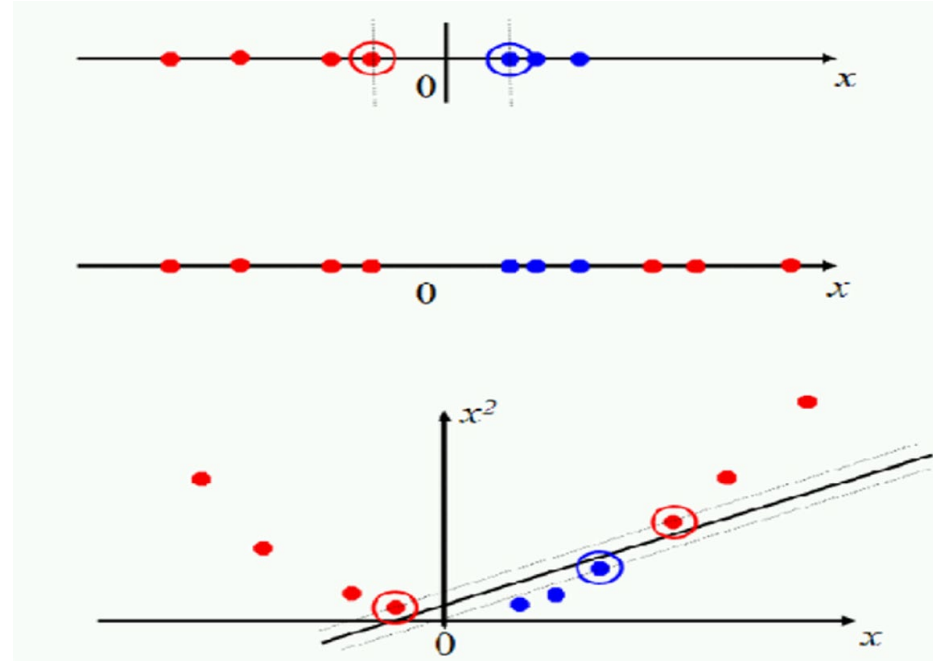
- Weighted sum of the kernels: $K(\mathbf{x}_i, \mathbf{x}_{i'}) = \sum_{l=1}^L v_l K_l(\mathbf{x}_i, \mathbf{x}_{i'}), v_l \geq 0$
- Objective function

$$L = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N \alpha_i \alpha_{i'} y_i y_{i'} \sum_{l=1}^L v_l K_l(\mathbf{x}_i, \mathbf{x}_{i'})$$

- Solved for both the support vectors, α_i and the weights
- For taking a decision in the test set

$$\sum_{i'=1}^N \alpha_{i'} y_{i'} \sum_{l=1}^L v_l K(\mathbf{x}_{i'}, \mathbf{x})$$

Support Vector Machines: Multiple Kernel Learning



- Projecting data that is not linearly separable into a higher dimensional space can make it linearly separable

Multi-class kernel machines

- One-vs-all approach
 - Define K two-class problems, each separating one class from all others
 - Learn K binary support vector machines $f_i(x), i = 1, \dots, K$
 - Class 1: Examples from class i
 - Class -1: Examples from all classes besides class p
 - Decision during testing

$$f_i(x) = \operatorname{argmax}_i f_i(x)$$

Multi-class kernel machines

- One-vs-one approach
 - Define $K(K-1)$ two-class problems, each separating class i from class j
 - Learn $K(K-1)$ binary SVM $f_{ij}(x)$
 - Class 1: Examples from class i
 - Class -1: Examples from class j
 - Note that $f_{ij} = -f_{ji}$
 - Decision during testing

$$f_i(x) = \operatorname{argmax}_i \left(\sum_j f_j(x) \right)$$

Multi-class kernel machines

- Comparison of one-vs-all and one-vs-one approach
 - One-vs-one
 - Requires $O(K^2)$ classifiers instead of $O(K)$
 - But each classifier is on average smaller $O\left(\frac{2N}{K}\right)$
 - One-vs-all approach solve $O(K)$ separate problems, each of size $O(N)$

Multi-class kernel machines

- Multi-class formulation
 - Define K weights for each class w_1, \dots, w_K and K bias terms w_{01}, \dots, w_{0K}
 - Training data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, y_i \in \{1, \dots, K\}$
 - Optimization criterion

$$\min_{w_1, \dots, w_K} \frac{1}{2} \sum_{k=1}^K \|w_k\|_2^2 + C \sum_{k=1}^K \sum_{i=1}^N \epsilon_{nk}$$
$$s. t. \mathbf{w}_{y_n}^T \mathbf{x}_n + w_0 \geq \mathbf{w}_k^T \mathbf{x}_i + w_{0k} + 2 - \epsilon_{nk}, \forall k \neq y_n$$

(ie. So that the weight for each class yields a sufficient margin form the other classes)

What have we learnt so far

- Kernel Functions $K(\mathbf{x}_i, \mathbf{x}_{i'})$
 - $K(\mathbf{x}_m, \mathbf{x}_n) = K(\mathbf{x}_n, \mathbf{x}_m)$ and $K(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n) \phi(\mathbf{x}_m)$
 - Positive definite kernel \rightarrow positive definite Gram matrix
- Kernel trick
 - Instead of first transforming the original features into the new features space and then computing the inner product, we can compute the inner product in the new features space directly through the kernel function
 - Kernel SVM: The cost of computing the primal variables is $O(N^3)$, while for the dual variable is $O(D^3)$. A kernel method can be useful in high dimensional settings
- Multi-class SVM
 - One-vs-one, one-vs-all, multiclass formulation

Takeaways and Next Time

- Primal-Dual formulation of SVM
- Kernelization
- Next time: Midterm review