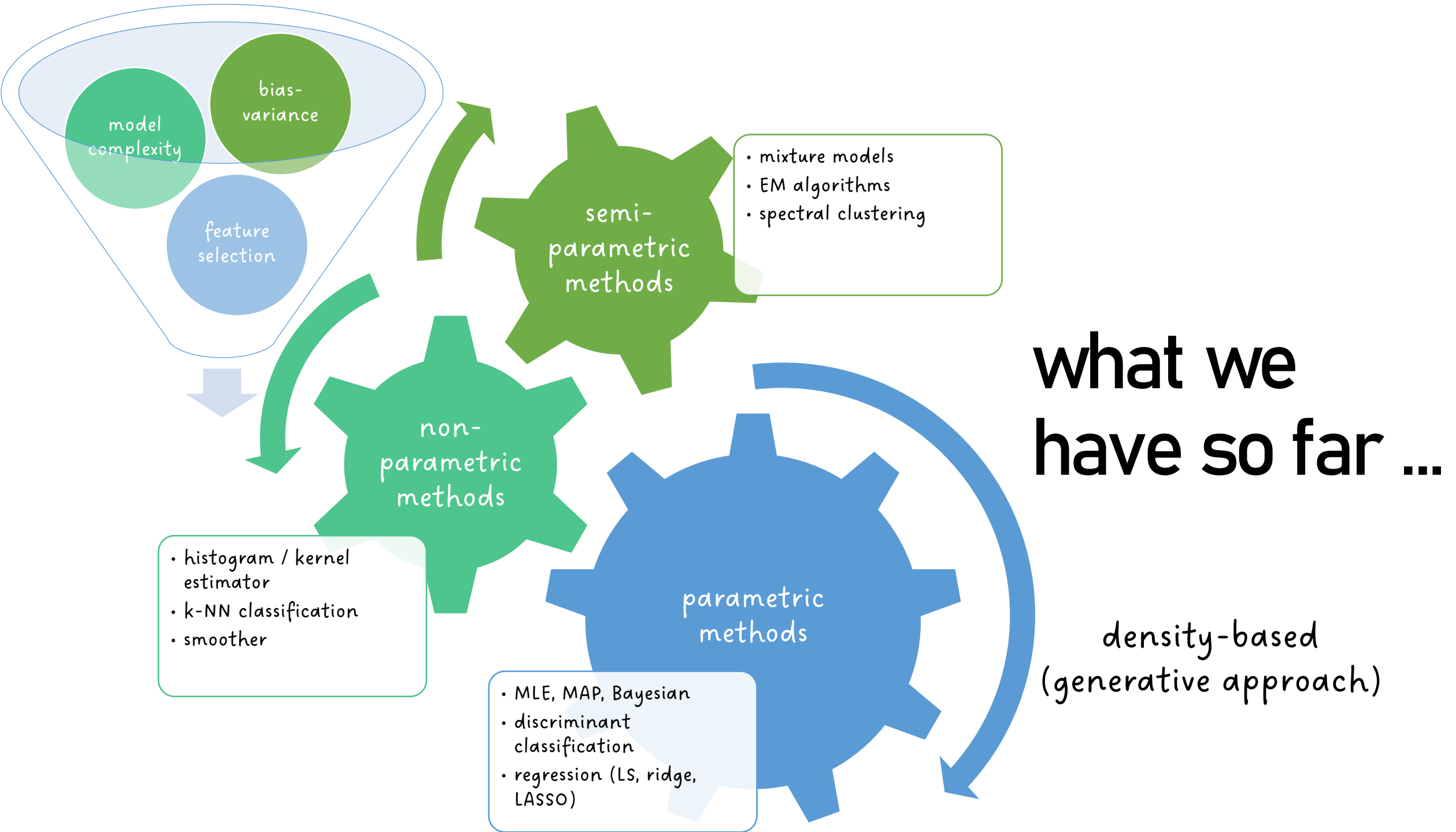


Kernel

STATS 303 Statistical Machine Learning

Spring 2022

Lecture 11



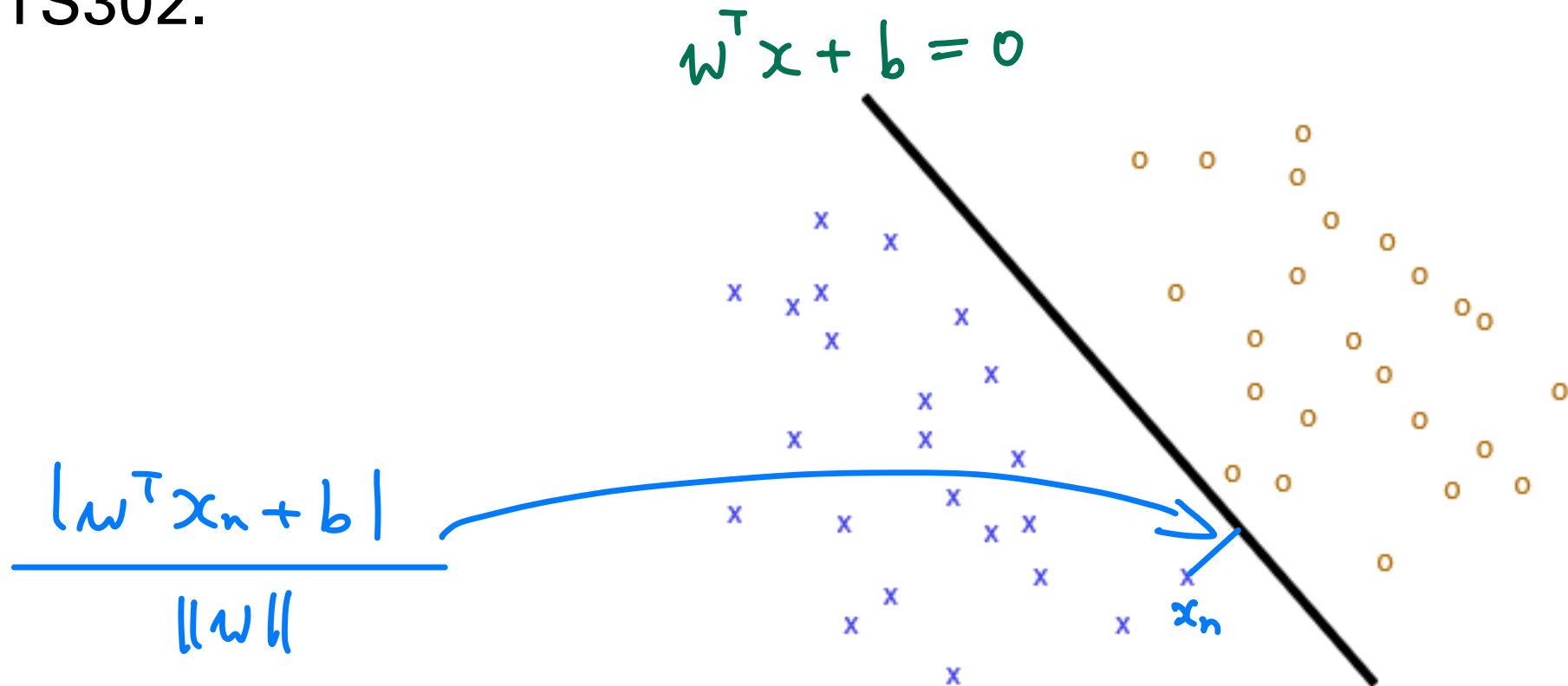
Vapnik's Principle

- If you possess a restricted amount of information for solving some problem, try to solve the problem directly and never solve a more general problem as an intermediate step.



linear SVM

- The best example that follows Vapnik's Principle is the support vector machine (SVM) that we learned in STATS302.



linear SVM

- Let $\{\mathbf{x}_n, t_n\}_{n=1}^N$ be the sample of training data for the classification problem, where $t_n \in \{+1, -1\}$ is the label.
- The linear SVM solves

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s. t.} \quad & t_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \text{for any } n \end{aligned}$$

- The dual problem is

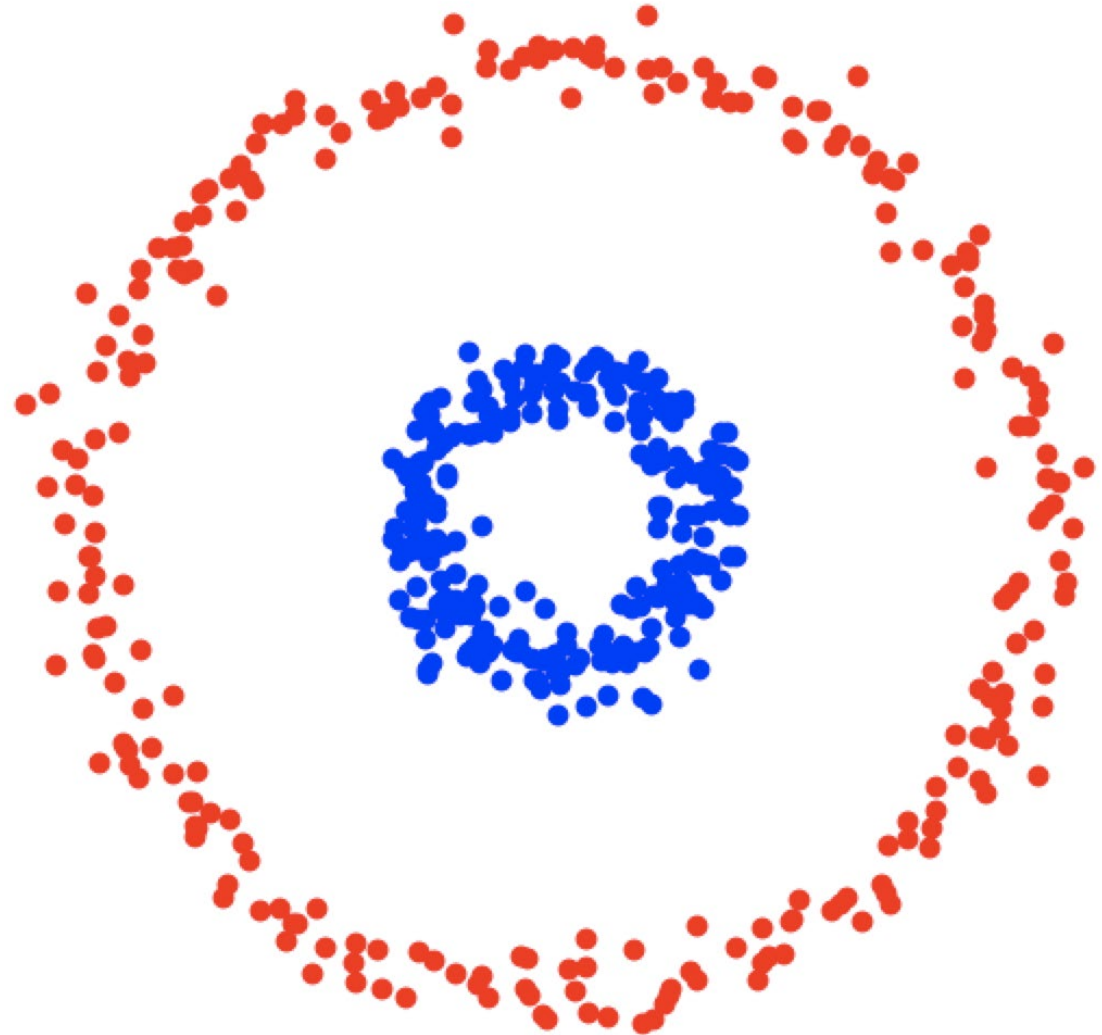
$$\begin{aligned} \max_{\mathbf{a}} \quad & \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \mathbf{x}_n^T \mathbf{x}_m \\ \text{s. t.} \quad & a_n \geq 0 \quad \text{for all } n; \quad \text{and} \quad \sum_{n=1}^N a_n t_n = 0 \end{aligned}$$

↑
Lagrange multipliers

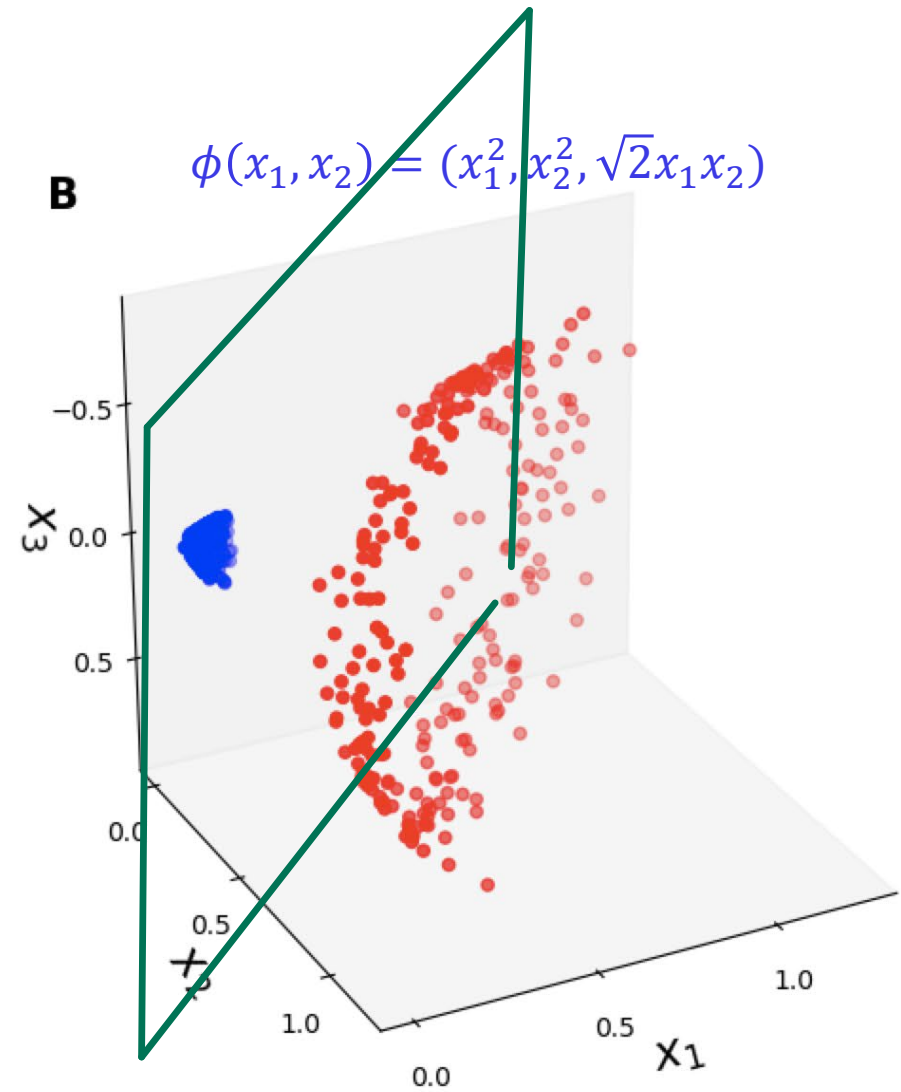
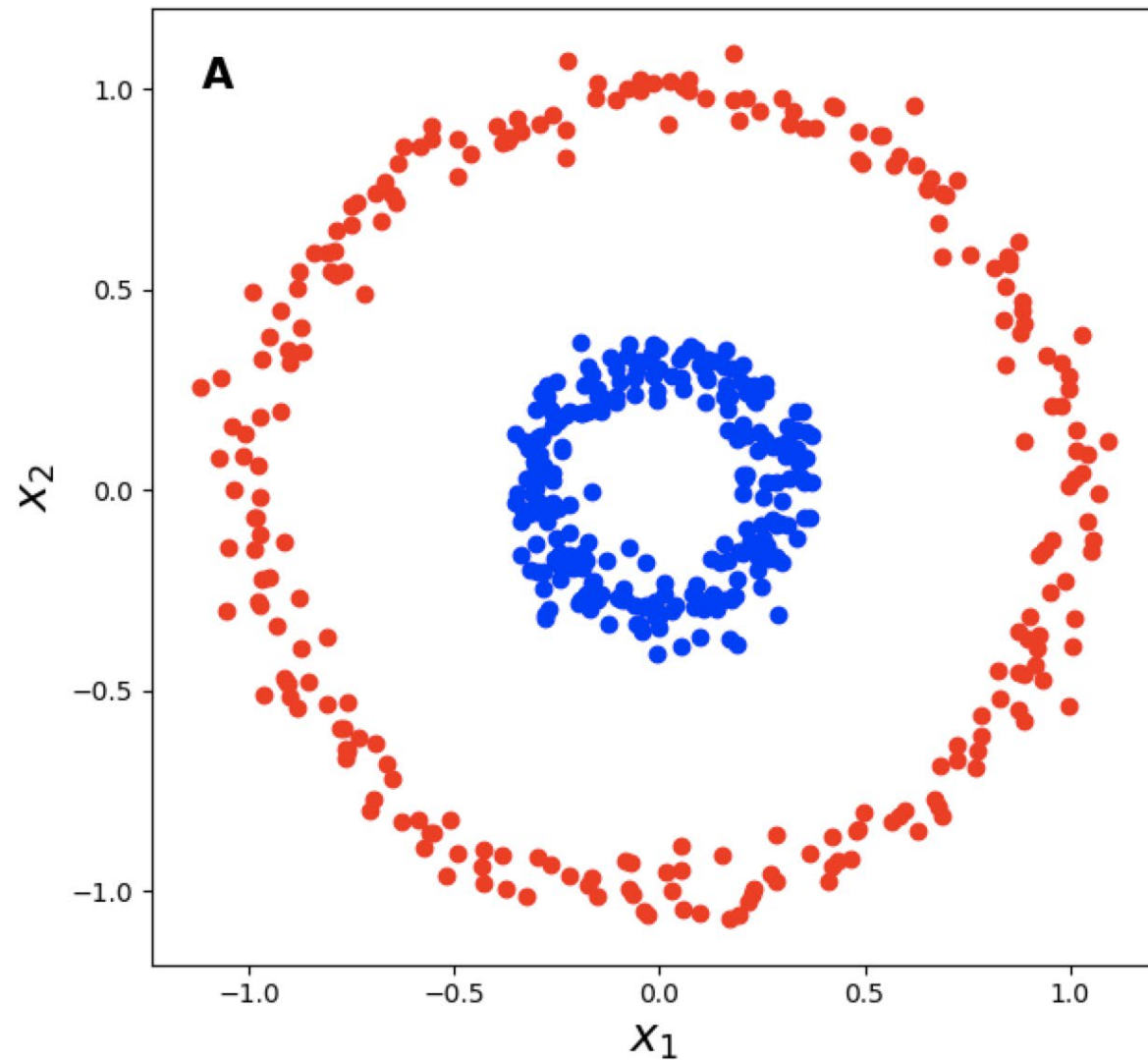
$$\mathbf{w} = \sum_{n=1}^N a_n t_n \mathbf{x}_n$$

kernel SVM

- If data are not linearly separable:



kernel SVM



kernel SVM

- Replacing each data point \mathbf{x}_n by $\phi(\mathbf{x}_n)$ in the formulation of SVM:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s. t.} \quad & t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 \quad \text{for any } n \end{aligned}$$

- Also, the dual problem becomes:

$$\begin{aligned} \max_{\mathbf{a}} \quad & \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) \\ \text{s. t.} \quad & a_n \geq 0 \quad \text{for all } n; \quad \text{and} \quad \sum_{n=1}^N a_n t_n = 0 \end{aligned}$$

kernel SVM

- Oftentimes, instead of finding a good mapping $\phi(\cdot)$, it is easier to define a “kernel” K , such that

$$K(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

- Given that, the dual form of the kernel SVM problem becomes

$$\begin{aligned} \max_{\mathbf{a}} \quad & \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m K(\mathbf{x}_n, \mathbf{x}_m) \\ \text{s. t.} \quad & a_n \geq 0 \quad \text{for all } n; \quad \text{and} \quad \sum_{n=1}^N a_n t_n = 0 \end{aligned}$$

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{n=1}^N a_n t_n K(\mathbf{x}, \mathbf{x}_n) + b$$

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$$

example of nonlinear map ϕ

- Let $\mathbf{x} = [x_1, x_2]^T, \mathbf{y} = [y_1, y_2]^T$
- Let $\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2]^T$
 $\phi(\mathbf{x})^T \phi(\mathbf{y}) = 1 + 2x_1y_1 + 2x_2y_2 + 2x_1x_2y_1y_2 + x_1^2y_1^2 + x_2^2y_2^2$
- $\phi(\mathbf{x})^T \phi(\mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^2 =: K(\mathbf{x}, \mathbf{y}) = (x_1y_1 + x_2y_2 + 1)^2$

use K instead of ϕ

- In practice, instead of starting with ϕ , we can directly start with K .
- Popular choices of K :
 - polynomial: $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^q$
 - radial basis function (RBF): $K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right)$
 - sigmoidal: $K(\mathbf{x}, \mathbf{y}) = \tanh(2\mathbf{x}^T \mathbf{y} + 1)$

food for thought

1. In general, what K can be chosen?
2. If e.g., $K(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$, what ϕ is behind it?

kernel PCA

kernel PCA

- Kernel methods are not just for SVMs.
- Suppose we want to do dimension reduction using principal component analysis (PCA) for some data that do not show a clear linear structure.
- Nevertheless, after applying a nonlinear function ϕ , the transformed data show a clear linear low-dimensional structure.
- We can apply PCA to $\{\phi(\mathbf{x}_n)\}_{n=1}^N$
- The question is how to avoid defining ϕ and use K instead.

$$\{x_n\}_{n=1}^N \xrightarrow{\text{nonlinear}} \{\phi(x_n)\}_{n=1}^N \xrightarrow{\text{apply PCA}} ?$$

Centralization: $\mu = \frac{1}{N} \sum_{n=1}^N \phi(x_n)$

Data matrix

$$\Phi = \begin{bmatrix} (\phi(x_1) - \mu)^T \\ (\phi(x_2) - \mu)^T \\ \vdots \\ (\phi(x_N) - \mu)^T \end{bmatrix}$$

PCA on $\{\phi(x_n)\}_{n=1}^N$

\Leftrightarrow Singular Value Decomposition (SVD) of Φ

\Leftrightarrow Eigenproblem of $\Phi\Phi^T$.

The (n,m) -th entry of $\Phi\Phi^T$ is given by

$$\begin{aligned} & (\phi(x_n) - \mu)^T (\phi(x_m) - \mu) \\ &= \phi(x_n)^T \phi(x_m) - \mu^T \phi(x_m) - \mu^T \phi(x_n) + \mu^T \mu \\ &= \phi(x_n)^T \phi(x_m) - \frac{1}{N} \sum_{\ell=1}^N \phi(x_\ell)^T \phi(x_m) - \\ & \quad \frac{1}{N} \sum_{\ell=1}^N \phi(x_\ell)^T \phi(x_n) + \frac{1}{N^2} \sum_{\ell=1}^N \sum_{k=1}^N \phi(x_\ell)^T \phi(x_k) \end{aligned}$$

Let K be the **Gram matrix** whose (n, m) -th entry is given by $K(n, m) = \phi(x_n)^T \phi(x_m)$

Let $\hat{\text{smiley}} := \Phi \Phi^T$. Then

$$\hat{\text{smiley}}(n, m)$$

$$\begin{aligned} &= \phi(x_n)^T \phi(x_m) - \frac{1}{N} \sum_{l=1}^N \phi(x_l)^T \phi(x_m) - \\ &\quad \frac{1}{N} \sum_{l=1}^N \phi(x_l)^T \phi(x_n) + \frac{1}{N^2} \sum_{l=1}^N \sum_{k=1}^N \phi(x_l)^T \phi(x_k) \\ &= K(n, m) - \frac{1}{N} \sum_{l=1}^N \overset{\text{"column sum"}}{K(l, m)} - \frac{1}{N} \sum_{l=1}^N \overset{\text{"row sum"}}{K(n, l)} + \\ &\quad \frac{1}{N^2} \sum_{l=1}^N \sum_{k=1}^N K(l, k) \\ &\quad \text{"column \& row sum"} \end{aligned}$$

Consider $\mathbf{1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$. Then $\mathbf{1} \mathbf{1}^T = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}$

$$\begin{aligned} \text{and } (\mathbf{1} \mathbf{1}^T K)(n, m) &= \sum_{l=1}^N \underbrace{(\mathbf{1} \mathbf{1}^T)(n, l)}_1 K(l, m) \\ &= \sum_{l=1}^N K(l, m) = \text{"column sum"}. \end{aligned}$$

Similarly, $(K \mathbb{1} \mathbb{1}^T)(n, m) = \sum_{l=1}^N K(n, l) = \text{"row sum"}$.

Therefore,

$$\Phi \Phi^T = K - \frac{1}{N} \mathbb{1} \mathbb{1}^T K - \frac{1}{N} K \mathbb{1} \mathbb{1}^T + \frac{1}{N^2} \mathbb{1} \mathbb{1}^T K \mathbb{1} \mathbb{1}^T$$

kernel regression

kernel regression

- In ridge regression, suppose our model is $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$:

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Taking gradient w.r.t. \mathbf{w} yields:

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n) \phi(\mathbf{x}_n) + \lambda \mathbf{w} \stackrel{\text{set}}{=} 0$$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \} \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a}$$

$$\text{Let } a_n = -\frac{1}{\lambda} (\mathbf{w}^T \phi(\mathbf{x}_n) - t_n)$$

kernel regression

- In ridge regression, suppose our model is $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$:

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Taking gradient w.r.t. \mathbf{w} yields: (cont'd from previous slide)

$$\text{Let } \Phi = \begin{bmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix}, \quad \mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{pmatrix}$$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \} \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a}$$

kernel regression

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- With $\mathbf{w} = \Phi^T \mathbf{a}$

$$\sum_{n=1}^N (\mathbf{w}^T \phi(x_n) - t_n)^2 = \begin{bmatrix} \mathbf{w}^T \phi(x_1) - t_1 \\ \mathbf{w}^T \phi(x_2) - t_2 \\ \vdots \\ \mathbf{w}^T \phi(x_N) - t_N \end{bmatrix}^T \begin{bmatrix} \mathbf{w}^T \phi(x_1) - t_1 \\ \mathbf{w}^T \phi(x_2) - t_2 \\ \vdots \\ \mathbf{w}^T \phi(x_N) - t_N \end{bmatrix}$$

$$= (\Phi \mathbf{w} - \mathbf{t})^T (\Phi \mathbf{w} - \mathbf{t}) \quad \text{where } \mathbf{t} = (t_1, t_2, \dots, t_N)^T$$

$$= (\Phi \Phi^T \mathbf{a} - \mathbf{t})^T (\Phi \Phi^T \mathbf{a} - \mathbf{t})$$

$$\text{Hence, } J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$

kernel regression

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$



$$K_{nm} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) \quad \text{"Gram matrix"}$$

$$\mathbf{K} = \Phi \Phi^T$$

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

kernel regression

- From $a_n = -\frac{1}{\lambda} \{ \mathbf{w}^T \phi(\mathbf{x}_n) - t_n \} \quad \mathbf{w} = \Phi^T \mathbf{a}$, we have

- Therefore,

$$\mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

Questions?

Reference

- *Kernel methods:*
 - [Al] Ch.13.5-13.7
 - [Bi] Ch.7.1, 6.1-6.2
- *Kernel PCA:*
 - [HaTF] Ch.14.5.4
- *Gaussian processes:*
 - [Bi] Ch.6.4.1-6.4.2, 6.4.5
 - [HaTF] Ch.5.8.1-5.8.2



Tasty & Healthy

Individual Pack

No Food Additives Added

Tuna
Floss