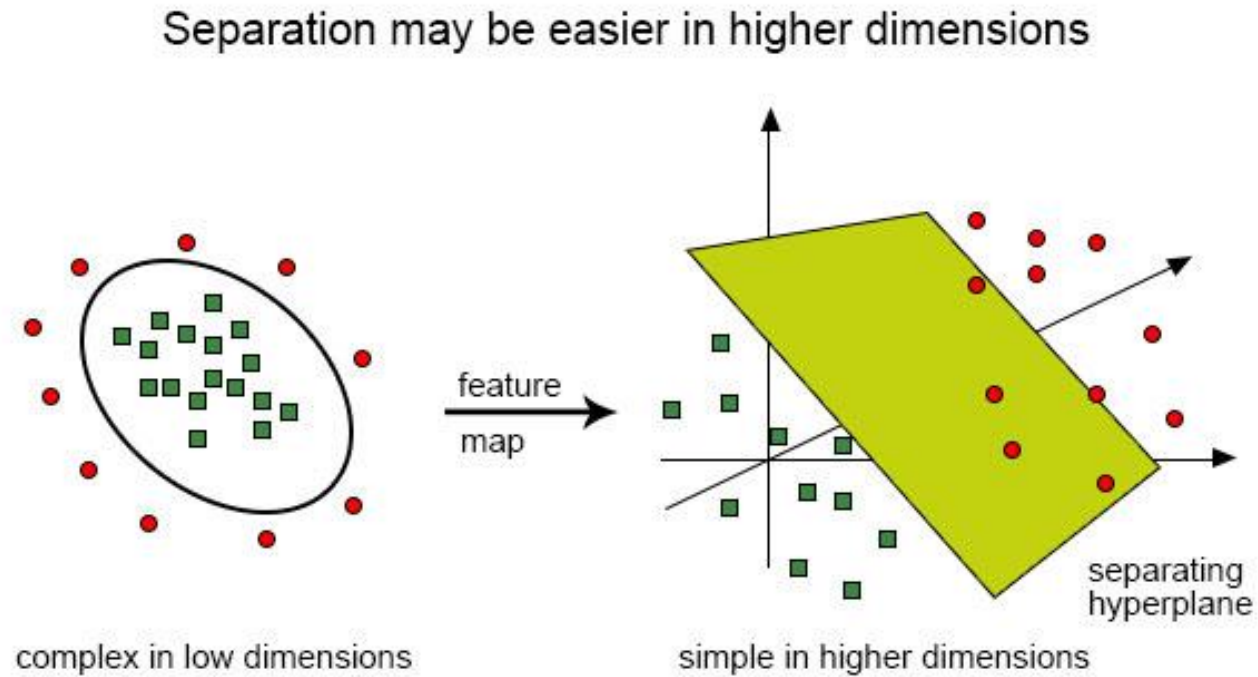ECE 57000

# Kernel method and Support vector machine

Chaoyue Liu

Fall 2024

# Linear models have limited performance (e.g., decision boundary is <u>linear</u>)

## Separation may be easier in higher dimensions

feature
map

complex in low dimensions

simple in higher dimensions

separating
hyperplane

# Feature map



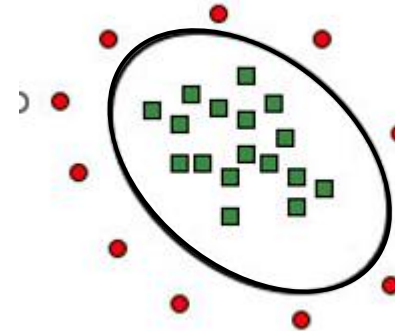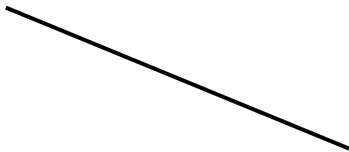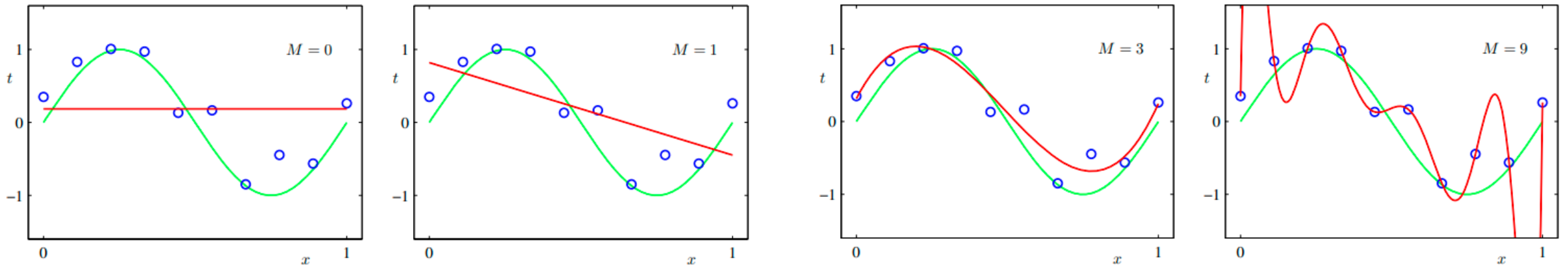| | Input space | Feature space |
|---|---|---|
| | $\mathbf{x} = (x_1, x_2)^T$ | $\phi(\mathbf{x}) = \left(1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2\right)^T$ |
| Model | $f_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$ | $f_{\mathbf{w}'}(\mathbf{x}) = \sigma(\mathbf{w}'^T \phi(\mathbf{x}))$ |
| Decision boundary | $\mathbf{w}^T \mathbf{x} = 0$ $(w_1 x_1 + w_2 x_2 = 0)$ | $\mathbf{w}'^T \phi(\mathbf{x}) = 0$ e.g. $x_1^2 + x_2^2 = 1$ |
| Decision boundary plot | | |

# Feature map

Input space: $x$

Feature space: $\phi(x) = (1, x, x^2, x^3, \cdots, x^M)^T$

Model: $f_{\mathbf{w}'}(x) = \mathbf{w}'^T \phi(x)$

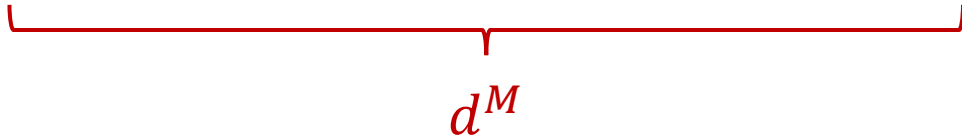The model is basically a $M$-degree polynomial of $x$



Larger $M$ gives better fit of the training data

# Downside of feature map

In practice, input $\mathbf{x}$ has a large dimension $d$. (for example, $d > 1000$)

What if we want the $M$-degree polynomial feature?

$$\phi(\mathbf{x}) = (1, x_1, \cdots, x_d, x_1^2, x_1 x_2, \cdots, x_1^M, \cdots, x_d^M)$$

$$\underbrace{\phantom{x_1, \cdots, x_d, x_1^2, x_1 x_2, \cdots, x_1^M, \cdots, x_d^M}}_{d^M}$$

Dimension of the feature vector $\phi(\mathbf{x})$, and weight $\mathbf{w}$: $\sim d^M$ (a huge number)

# Downside of feature map

In practice, input $\mathbf{x}$ has a large dimension $d$. (for example, $d > 1000$)

Dimension of the feature vector $\phi(\mathbf{x})$, and weight $\mathbf{w}$: $\sim d^M$ (a huge number)

**Q**: Do we really need to compute the whole vectors $\phi(\mathbf{x})$ and $\mathbf{w}$?

A: we only care about the <u>inner product</u> $\mathbf{w}^T \phi(\mathbf{x})$ which is a <u>scalar</u>.

recall: gradient descent: $\mathbf{w}^{t+1} - \mathbf{w}^t = \eta \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i) \cdot \phi(\mathbf{x}_i)$

$\mathbf{w} = \sum_{i=1}^n \alpha_i \cdot \phi(\mathbf{x}_i)$, for some coefficients $\alpha_i$

$$f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \sum_{i=1}^n \alpha_i \cdot \phi^T(\mathbf{x}_i) \phi(\mathbf{x})$$

We only need the inner products of feature vectors

# The "kernel trick"

$$\phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2)$$

$$\phi(\mathbf{z}) = (1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, \sqrt{2}z_1 z_2, z_2^2)$$

$$\phi^T(\mathbf{x})\phi(\mathbf{z}) = 1 + 2x_1 z_1 + 2x_2 z_2 + x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2$$

$$= (1 + x_1 z_1 + x_2 z_2)^2$$

$$= (1 + \mathbf{x}^T \mathbf{z})^2 \quad =: K(\mathbf{x}^T \mathbf{z})$$

More generally, for $M$-degree polynomials:

$$\phi(\mathbf{x}) = (1, x_1, \cdots, x_d, x_1^2, x_1 x_2, \cdots, x_1^M, \cdots, x_d^M) \qquad \sim d^M \text{ dimensional}$$

$$\phi^T(\mathbf{x})\phi(\mathbf{z}) = K(\mathbf{x}^T \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^M$$
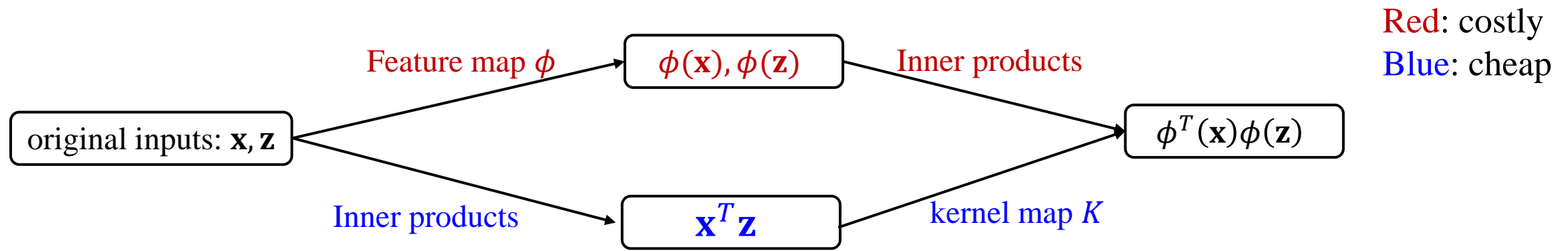
# The "kernel trick"

More generally, for $M$-degree polynomials:

$$\phi(\mathbf{x}) = (1, x_1, \cdots, x_d, x_1^2, x_1 x_2, \cdots, x_1^M, \cdots, x_d^M)$$

$$\phi(\mathbf{z}) = (1, z_1, \cdots, z_d, z_1^2, z_1 z_2, \cdots, z_1^M, \cdots, z_d^M)$$

$\sim d^M$ dimensional

$$\phi^T(\mathbf{x})\phi(\mathbf{z}) = K(\mathbf{x}^T \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^M$$

Red: costly
Blue: cheap

# The "kernel trick"

$$f(\mathbf{x}) = \mathbf{w}^\mathrm{T}\phi(\mathbf{x}) = \sum_{i=1}^{n} \alpha_i \cdot \phi^T(\mathbf{x}_i)\phi(\mathbf{x})$$

$$= \sum_{i=1}^{n} \alpha_i \cdot K(\mathbf{x_i}, \mathbf{x})$$

Instead of learning a high-dimensional weight vector $\mathbf{w}$, we just need to learn the coefficients $\alpha_i$.

# The "kernel trick"

One step further:
- we don't even need to know what feature maps are used
- we only need the kernel information

Kernel: $K(\cdot, \cdot): \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$    need a pair of inputs
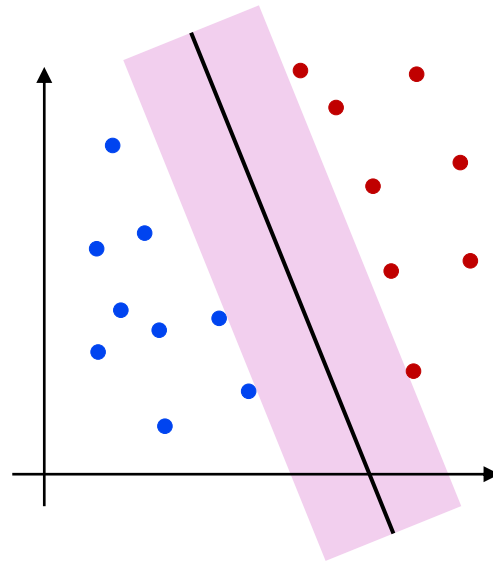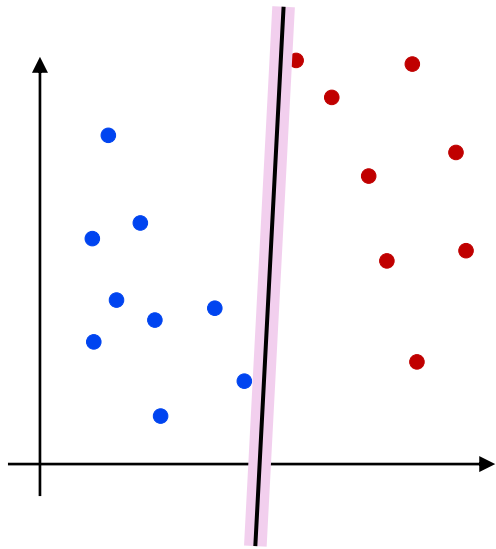
Polynomial kernel: $K(\mathbf{x}, \mathbf{z}) = (\mathbf{1} + \mathbf{x}^T \mathbf{z})^M$
Gaussian kernel: $K(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|^2 / \sigma^2)$
Laplace kernel: $K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|)$

# Support vector machine (SVM)

Setting: consider a <u>linearly separable</u> data with two classes



The latter is preferable, because it has a larger margin.

# Support vector machine (SVM)

Decision boundary: $\mathbf{w}^T\mathbf{x} + b = 0$
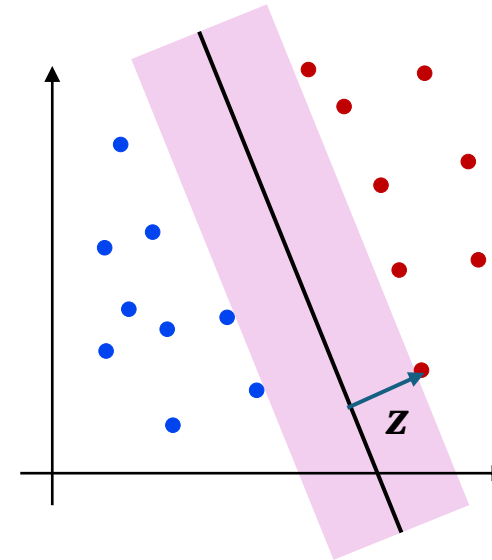
**Q**: does $(\mathbf{w}, b)$ and $\alpha(\mathbf{w}, b)$ give the same decision boundary?

Normalize $(\mathbf{w}, b)$:

let the closest data points $\mathbf{x}$ satisfy $y \cdot (\mathbf{w}^T\mathbf{x} + b) = 1$, and others with $y' \cdot (\mathbf{w}^T\mathbf{x}' + b) > 1$

**Compute the margin:**

$$\mathbf{w}^T\mathbf{x} + b = 1$$
$$\mathbf{w}^T(\mathbf{x} - \mathbf{z}) + b = 0$$

$\longrightarrow \quad \mathbf{w}^T\mathbf{z} = \|\mathbf{w}\|\|\mathbf{z}\| = 1$

margin: $\|\mathbf{z}\| = \dfrac{1}{\|\mathbf{w}\|}$

# Support vector machine (SVM)

**Goal**: maximize the margin $\frac{1}{\|\mathbf{w}\|}$, equivalent to minimize $\|\mathbf{w}\|^2$

minimize $\|\mathbf{w}\|^2$
subject to:    $y_i \cdot (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall i$

What if the data is <u>not</u> linearly separable?
Soft margin

minimize $\|\mathbf{w}\|^2 + C \sum_{i=1}^{n} \xi_i$
subject to:
$\quad y_i \cdot (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \forall i$
$\quad \xi_i \geq 0, \forall i$

$\Longrightarrow$

Minimize
$\|\mathbf{w}\|^2 + C \sum_{i=1}^{n} Hinge(y_i(\mathbf{w}^T \mathbf{x}_i + b))$

$Hinge(z) = \max(0, 1 - z)$